

NEW DEVELOPMENTS in CIRCUITS, SYSTEMS, SIGNAL PROCESSING, COMMUNICATIONS and COMPUTERS

**Proceedings of the International Conference on Circuits, Systems,
Signal Processing, Communications and Computers (CSSCC 2015)**

**Vienna, Austria
March 15-17, 2015**

NEW DEVELOPMENTS in CIRCUITS, SYSTEMS, SIGNAL PROCESSING, COMMUNICATIONS and COMPUTERS

**Proceedings of the International Conference on Circuits, Systems,
Signal Processing, Communications and Computers (CSSCC 2015)**

**Vienna, Austria
March 15-17, 2015**

Copyright © 2015, by the editors

All the copyright of the present book belongs to the editors. All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, electronic, mechanical, photocopying, recording, or otherwise, without the prior written permission of the editors.

All papers of the present volume were peer reviewed by no less than two independent reviewers. Acceptance was granted when both reviewers' recommendations were positive.

Series: Recent Advances in Electrical Engineering Series | 45

ISSN: 1790-5117

ISBN: 978-1-61804-285-9

NEW DEVELOPMENTS in CIRCUITS, SYSTEMS, SIGNAL PROCESSING, COMMUNICATIONS and COMPUTERS

**Proceedings of the International Conference on Circuits, Systems,
Signal Processing, Communications and Computers (CSSCC 2015)**

**Vienna, Austria
March 15-17, 2015**

Organizing Committee

Editors:

Professor Nikos E. Mastorakis, Technical University of Sofia, Bulgaria
Professor Valeri Mladenov, Technical University of Sofia, Bulgaria
Professor Klimis Ntalianis, Technological Educational Institute of Athens, Greece

Program Committee:

Prof. Lotfi Zadeh (IEEE Fellow, University of Berkeley, USA)
Prof. Leon Chua (IEEE Fellow, University of Berkeley, USA)
Prof. Michio Sugeno (RIKEN Brain Science Institute (RIKEN BSI), Japan)
Prof. Dimitri Bertsekas (IEEE Fellow, MIT, USA)
Prof. Demetri Terzopoulos (IEEE Fellow, ACM Fellow, UCLA, USA)
Prof. Georgios B. Giannakis (IEEE Fellow, University of Minnesota, USA)
Prof. George Vachtsevanos (Georgia Institute of Technology, USA)
Prof. Abraham Bers (IEEE Fellow, MIT, USA)
Prof. David Staelin (IEEE Fellow, MIT, USA)
Prof. Brian Barsky (IEEE Fellow, University of Berkeley, USA)
Prof. Aggelos Katsaggelos (IEEE Fellow, Northwestern University, USA)
Prof. Josef Sifakis (Turing Award 2007, CNRS/Verimag, France)
Prof. Hisashi Kobayashi (Princeton University, USA)
Prof. Kinshuk (Fellow IEEE, Massey Univ. New Zealand),
Prof. Leonid Kazovsky (Stanford University, USA)
Prof. Narsingh Deo (IEEE Fellow, ACM Fellow, University of Central Florida, USA)
Prof. Kamisetty Rao (Fellow IEEE, Univ. of Texas at Arlington, USA)
Prof. Anastassios Venetsanopoulos (Fellow IEEE, University of Toronto, Canada)
Prof. Steven Collicott (Purdue University, West Lafayette, IN, USA)
Prof. Nikolaos Paragios (Ecole Centrale Paris, France)
Prof. Nikolaos G. Bourbakis (IEEE Fellow, Wright State University, USA)
Prof. Stamatios Kartalopoulos (IEEE Fellow, University of Oklahoma, USA)
Prof. Irwin Sandberg (IEEE Fellow, University of Texas at Austin, USA),
Prof. Michael Sebek (IEEE Fellow, Czech Technical University in Prague, Czech Republic)
Prof. Hashem Akbari (University of California, Berkeley, USA)
Prof. Yuriy S. Shmaliy, (IEEE Fellow, The University of Guanajuato, Mexico)
Prof. Lei Xu (IEEE Fellow, Chinese University of Hong Kong, Hong Kong)
Prof. Paul E. Dimotakis (California Institute of Technology Pasadena, USA)
Prof. M. Pelikan (UMSL, USA)
Prof. Patrick Wang (MIT, USA)
Prof. Wasfy B Mikhael (IEEE Fellow, University of Central Florida Orlando, USA)
Prof. Sunil Das (IEEE Fellow, University of Ottawa, Canada)
Prof. Panos Pardalos (University of Florida, USA)
Prof. Nikolaos D. Katopodes (University of Michigan, USA)
Prof. Bimal K. Bose (Life Fellow of IEEE, University of Tennessee, Knoxville, USA)
Prof. Janusz Kacprzyk (IEEE Fellow, Polish Academy of Sciences, Poland)
Prof. Sidney Burrus (IEEE Fellow, Rice University, USA)
Prof. Biswa N. Datta (IEEE Fellow, Northern Illinois University, USA)
Prof. Mihai Putinar (University of California at Santa Barbara, USA)
Prof. Wlodzislaw Duch (Nicolaus Copernicus University, Poland)
Prof. Tadeusz Kaczorek (IEEE Fellow, Warsaw University of Tehcnology, Poland)
Prof. Michael N. Katehakis (Rutgers, The State University of New Jersey, USA)
Prof. Pan Agathoklis (Univ. of Victoria, Canada)
Prof. P. Demokritou (Harvard University, USA)
Prof. P. Razelos (Columbia University, USA)
Dr. Subhas C. Misra (Harvard University, USA)

Prof. Martin van den Toorn (Delft University of Technology, The Netherlands)
Prof. Malcolm J. Crocker (Distinguished University Prof., Auburn University, USA)
Prof. S. Dafermos (Brown University, USA)
Prof. Urszula Ledzewicz, Southern Illinois University, USA.
Prof. Dimitri Kazakos, Dean, (Texas Southern University, USA)
Prof. Ronald Yager (Iona College, USA)
Prof. Athanassios Manikas (Imperial College, London, UK)
Prof. Keith L. Clark (Imperial College, London, UK)
Prof. Argyris Varonides (Univ. of Scranton, USA)
Prof. S. Furfari (Direction Generale Energie et Transports, Brussels, EU)
Prof. Constantin Udriste, University Politehnica of Bucharest, ROMANIA
Dr. Michelle Luke (Univ. Berkeley, USA)
Prof. Patrice Brault (Univ. Paris-sud, France)
Dr. Christos E. Vasios (MIT, USA)
Prof. Jim Cunningham (Imperial College London, UK)
Prof. Philippe Ben-Abdallah (Ecole Polytechnique de l'Universite de Nantes, France)
Prof. Photios Anninos (Medical School of Thrace, Greece)
Prof. Ichiro Hagiwara, (Tokyo Institute of Technology, Japan)
Prof. Metin Demiralp (Istanbul Technical University / Turkish Academy of Sciences, Istanbul, Turkey)
Prof. Andris Buikis (Latvian Academy of Science. Latvia)
Prof. Akshai Aggarwal (University of Windsor, Canada)
Prof. George Vachtsevanos (Georgia Institute of Technology, USA)
Prof. Ulrich Albrecht (Auburn University, USA)
Prof. Imre J. Rudas (Obuda University, Hungary)
Prof. Alexey L Sadovski (IEEE Fellow, Texas A&M University, USA)
Prof. Amedeo Andreotti (University of Naples, Italy)
Prof. Ryszard S. Choras (University of Technology and Life Sciences Bydgoszcz, Poland)
Prof. Remi Leandre (Universite de Bourgogne, Dijon, France)
Prof. Moustapha Diaby (University of Connecticut, USA)
Prof. Brian McCartin (New York University, USA)
Prof. Elias C. Aifantis (Aristotle Univ. of Thessaloniki, Greece)
Prof. Anastasios Lyrintzis (Purdue University, USA)
Prof. Charles Long (Prof. Emeritus University of Wisconsin, USA)
Prof. Marvin Goldstein (NASA Glenn Research Center, USA)
Prof. Costin Cepisca (University POLITEHNICA of Bucharest, Romania)
Prof. Kleanthis Psarris (University of Texas at San Antonio, USA)
Prof. Ron Goldman (Rice University, USA)
Prof. Ioannis A. Kakadiaris (University of Houston, USA)
Prof. Richard Tapia (Rice University, USA)
Prof. F.-K. Benra (University of Duisburg-Essen, Germany)
Prof. Milivoje M. Kostic (Northern Illinois University, USA)
Prof. Helmut Jaberg (University of Technology Graz, Austria)
Prof. Ardeshir Anjomani (The University of Texas at Arlington, USA)
Prof. Heinz Ulbrich (Technical University Munich, Germany)
Prof. Reinhard Leithner (Technical University Braunschweig, Germany)
Prof. Elbrous M. Jafarov (Istanbul Technical University, Turkey)
Prof. M. Ehsani (Texas A&M University, USA)
Prof. Sesh Commuri (University of Oklahoma, USA)
Prof. Nicolas Galanis (Universite de Sherbrooke, Canada)
Prof. S. H. Sohrab (Northwestern University, USA)
Prof. Rui J. P. de Figueiredo (University of California, USA)
Prof. Valeri Mladenov (Technical University of Sofia, Bulgaria)
Prof. Hiroshi Sakaki (Meisei University, Tokyo, Japan)
Prof. Zoran S. Bojkovic (Technical University of Belgrade, Serbia)

Prof. K. D. Klaes, (Head of the EPS Support Science Team in the MET Division at EUMETSAT, France)
Prof. Emira Maljevic (Technical University of Belgrade, Serbia)
Prof. Kazuhiko Tsuda (University of Tsukuba, Tokyo, Japan)
Prof. Milan Stork (University of West Bohemia , Czech Republic)
Prof. C. G. Helmis (University of Athens, Greece)
Prof. Lajos Barna (Budapest University of Technology and Economics, Hungary)
Prof. Nobuoki Mano (Meisei University, Tokyo, Japan)
Prof. Nobuo Nakajima (The University of Electro-Communications, Tokyo, Japan)
Prof. Victor-Emil Neagoe (Polytechnic University of Bucharest, Romania)
Prof. E. Protonotarios (National Technical University of Athens, Greece)
Prof. P. Vanderstraeten (Brussels Institute for Environmental Management, Belgium)
Prof. Annaliese Bischoff (University of Massachusetts, Amherst, USA)
Prof. Virgil Tiponut (Politehnica University of Timisoara, Romania)
Prof. Andrei Kolyshkin (Riga Technical University, Latvia)
Prof. Fumiaki Imado (Shinshu University, Japan)
Prof. Sotirios G. Ziavras (New Jersey Institute of Technology, USA)
Prof. Constantin Volosencu (Politehnica University of Timisoara, Romania)
Prof. Marc A. Rosen (University of Ontario Institute of Technology, Canada)
Prof. Alexander Zemliak (Puebla Autonomous University, Mexico)
Prof. Thomas M. Gatton (National University, San Diego, USA)
Prof. Leonardo Pagnotta (University of Calabria, Italy)
Prof. Yan Wu (Georgia Southern University, USA)
Prof. Daniel N. Riahi (University of Texas-Pan American, USA)
Prof. Alexander Grebennikov (Autonomous University of Puebla, Mexico)
Prof. Bennie F. L. Ward (Baylor University, TX, USA)
Prof. Guennadi A. Kouzaev (Norwegian University of Science and Technology, Norway)
Prof. Eugene Kindler (University of Ostrava, Czech Republic)
Prof. Geoff Skinner (The University of Newcastle, Australia)
Prof. Hamido Fujita (Iwate Prefectural University(IPU), Japan)
Prof. Francesco Muzi (University of L'Aquila, Italy)
Prof. Les M. Sztandera (Philadelphia University, USA)
Prof. Claudio Rossi (University of Siena, Italy)
Prof. Christopher J. Koroneos (Aristotle University of Thessaloniki, Greece)
Prof. Sergey B. Leonov (Joint Institute for High Temperature Russian Academy of Science, Russia)
Prof. Arpad A. Fay (University of Miskolc, Hungary)
Prof. Lili He (San Jose State University, USA)
Prof. M. Nasseh Tabrizi (East Carolina University, USA)
Prof. Alaa Eldin Fahmy (University Of Calgary, Canada)
Prof. Ion Carstea (University of Craiova, Romania)
Prof. Paul Dan Cristea (University "Politehnica" of Bucharest, Romania)
Prof. Gh. Pascovici (University of Koeln, Germany)
Prof. Pier Paolo Delsanto (Politecnico of Torino, Italy)
Prof. Radu Munteanu (Rector of the Technical University of Cluj-Napoca, Romania)
Prof. Ioan Dumitrache (Politehnica University of Bucharest, Romania)
Prof. Corneliu Lazar (Technical University Gh.Asachi Iasi, Romania)
Prof. Nicola Pitrone (Universita degli Studi Catania, Italia)
Prof. Miquel Salgot (University of Barcelona, Spain)
Prof. Amaury A. Caballero (Florida International University, USA)
Prof. Maria I. Garcia-Planas (Universitat Politecnica de Catalunya, Spain)
Prof. Petar Popivanov (Bulgarian Academy of Sciences, Bulgaria)
Prof. Alexander Gegov (University of Portsmouth, UK)
Prof. Lin Feng (Nanyang Technological University, Singapore)
Prof. Colin Fyfe (University of the West of Scotland, UK)
Prof. Zhaohui Luo (Univ of London, UK)

Prof. Mikhail Itskov (RWTH Aachen University, Germany)
Prof. George G. Tsytkin (Russian Academy of Sciences, Russia)
Prof. Wolfgang Wenzel (Institute for Nanotechnology, Germany)
Prof. Weilian Su (Naval Postgraduate School, USA)
Prof. Phillip G. Bradford (The University of Alabama, USA)
Prof. Ray Hefferlin (Southern Adventist University, TN, USA)
Prof. Gabriella Bognar (University of Miskolc, Hungary)
Prof. Hamid Abachi (Monash University, Australia)
Prof. Karlheinz Spindler (Fachhochschule Wiesbaden, Germany)
Prof. Josef Boercsoek (Universitat Kassel, Germany)
Prof. Eyad H. Abed (University of Maryland, Maryland, USA)
Prof. F. Castanie (TeSA, Toulouse, France)
Prof. Robert K. L. Gay (Nanyang Technological University, Singapore)
Prof. Andrzej Ordys (Kingston University, UK)
Prof. Harris Catrakis (Univ of California Irvine, USA)
Prof. T Bott (The University of Birmingham, UK)
Prof. Petr Filip (Institute of Hydrodynamics, Prague, Czech Republic)
Prof. T.-W. Lee (Arizona State University, AZ, USA)
Prof. Le Yi Wang (Wayne State University, Detroit, USA)
Prof. George Stavrakakis (Technical University of Crete, Greece)
Prof. John K. Galitos (Houston Community College, USA)
Prof. M. Petrakis (National Observatory of Athens, Greece)
Prof. Philippe Dondon (ENSEIRB, Talence, France)
Prof. Dalibor Bielek (Brno University of Technology, Czech Republic)
Prof. Oleksander Markovskyy (National Technical University of Ukraine, Ukraine)
Prof. Suresh P. Sethi (University of Texas at Dallas, USA)
Prof. Hartmut Hillmer (University of Kassel, Germany)
Prof. Bram Van Putten (Wageningen University, The Netherlands)
Prof. Alexander Iomin (Technion - Israel Institute of Technology, Israel)
Prof. Roberto San Jose (Technical University of Madrid, Spain)
Prof. Minvydas Ragulskis (Kaunas University of Technology, Lithuania)
Prof. Arun Kulkarni (The University of Texas at Tyler, USA)
Prof. Joydeep Mitra (New Mexico State University, USA)
Prof. Vincenzo Niola (University of Naples Federico II, Italy)
Prof. Ion Chrysosoverghi (National Technical University of Athens, Greece)
Prof. Dr. Aydin Akan (Istanbul University, Turkey)
Prof. Sarka Necasova (Academy of Sciences, Prague, Czech Republic)
Prof. C. D. Memos (National Technical University of Athens, Greece)
Prof. S. Y. Chen, (Zhejiang University of Technology, China and University of Hamburg, Germany)
Prof. Duc Nguyen (Old Dominion University, Norfolk, USA)
Prof. Tuan Pham (James Cook University, Townsville, Australia)
Prof. Jiri Klima (Technical Faculty of CZU in Prague, Czech Republic)
Prof. Rossella Cancelliere (University of Torino, Italy)
Prof. L.Kohout (Florida State University, Tallahassee, Florida, USA)
Prof. D' Attelis (Univ. Buenos Ayres, Argentina)
Prof. Dr-Eng. Christian Bouquegneau (Faculty Polytechnique de Mons, Belgium)
Prof. Wladyslaw Mielczarski (Technical University of Lodz, Poland)
Prof. Ibrahim Hassan (Concordia University, Montreal, Quebec, Canada)
Prof. Stavros J.Baloyannis (Medical School, Aristotle University of Thessaloniki, Greece)
Prof. James F. Frenzel (University of Idaho, USA)
Prof. Mirko Novak (Czech Technical University in Prague, Czech Republic)
Prof. Zdenek Votruba (Czech Technical University in Prague, Czech Republic)
Prof. Vilem Srovnal, (Technical University of Ostrava, Czech Republic)
Prof. J. M. Giron-Sierra (Universidad Complutense de Madrid, Spain)

Prof. Zeljko Panian (University of Zagreb, Croatia)
Prof. Walter Dosch (University of Luebeck, Germany)
Prof. Rudolf Freund (Vienna University of Technology, Austria)
Prof. Erich Schmidt (Vienna University of Technology, Austria)
Prof. Alessandro Genco (University of Palermo, Italy)
Prof. Martin Lopez Morales (Technical University of Monterey, Mexico)
Prof. Ralph W. Oberste-Vorth (Marshall University, USA)
Prof. Vladimir Damgov (Bulgarian Academy of Sciences, Bulgaria)
Prof. Menelaos Karanasos (Brunel University, UK)
Prof. P.Borne (Ecole Central de Lille, France)

Additional Reviewers

Jose Flores	The University of South Dakota, SD, USA
Abelha Antonio	Universidade do Minho, Portugal
Lesley Farmer	California State University Long Beach, CA, USA
Takuya Yamano	Kanagawa University, Japan
Miguel Carriegos	Universidad de Leon, Spain
Francesco Zirilli	Sapienza Universita di Roma, Italy
George Barreto	Pontificia Universidad Javeriana, Colombia
Eleazar Jimenez Serrano	Kyushu University, Japan
Tetsuya Yoshida	Hokkaido University, Japan
Philippe Dondon	Institut polytechnique de Bordeaux, France
Genqi Xu	Tianjin University, China
M. Javed Khan	Tuskegee University, AL, USA
Xiang Bai	Huazhong University of Science and Technology, China
Dmitrijs Serdjuks	Riga Technical University, Latvia
Hessam Ghasemnejad	Kingston University London, UK
José Carlos Metrôlho	Instituto Politecnico de Castelo Branco, Portugal
João Bastos	Instituto Superior de Engenharia do Porto, Portugal
Tetsuya Shimamura	Saitama University, Japan
Imre Rudas	Obuda University, Budapest, Hungary
Konstantin Volkov	Kingston University London, UK
Frederic Kuznik	National Institute of Applied Sciences, Lyon, France
James Vance	The University of Virginia's College at Wise, VA, USA
Angel F. Tenorio	Universidad Pablo de Olavide, Spain
Sorinel Oprisan	College of Charleston, CA, USA
Santoso Wibowo	CQ University, Australia
Jon Burley	Michigan State University, MI, USA
Kazuhiko Natori	Toho University, Japan
Shinji Osada	Gifu University School of Medicine, Japan
Francesco Rotondo	Polytechnic of Bari University, Italy
Deolinda Rasteiro	Coimbra Institute of Engineering, Portugal
Alejandro Fuentes-Penna	Universidad Autónoma del Estado de Hidalgo, Mexico
Moran Wang	Tsinghua University, China
Bazil Taha Ahmed	Universidad Autonoma de Madrid, Spain
Andrey Dmitriev	Russian Academy of Sciences, Russia
Masaji Tanaka	Okayama University of Science, Japan
Matthias Buyle	Artesis Hogeschool Antwerpen, Belgium
Kei Eguchi	Fukuoka Institute of Technology, Japan
Zhong-Jie Han	Tianjin University, China
Valeri Mladenov	Technical University of Sofia, Bulgaria
Ole Christian Boe	Norwegian Military Academy, Norway
Yamagishi Hiromitsu	Ehime University, Japan
Stavros Ponis	National Technical University of Athens, Greece
Minhui Yan	Shanghai Maritime University, China

Table of Contents

Prediction of Cancer Behavior based on Artificial Intelligence <i>Shayma M. Al-Ani, Maysma Abbod</i>	15
Modeling and Analysis of Elapsed Time and Energy Consumption of Interactive Applications in Mobile Cloud Computing Environments <i>Young-Chul Shim</i>	20
Lossless, Multiband, on Board, Compression of Hyperspectral Images <i>Bruno Carpentieri, Raffaele Pizzolante</i>	27
Semantic Web Technologies and Model-Driven Approach for the Development and Configuration Management of Intelligent Web-Based Systems <i>Arturs Bartusevics, Andrejs Lesovskis, Leonids Novickis</i>	32
Clock Distribution using a Bi-Dimensional Orthogonal Salphasic Structure <i>Andrei Pasca</i>	40
Synchronous Differential Logic Gate for Low Clock Swing Operation with Standing Wave Clock Distribution Networks <i>Andrei Pasca</i>	48
On-line Monitoring of Yogurt Fermentation using Ultrasonic Characteristics <i>Ahmad Aljaafreh, Ralf Lucklum</i>	56
Automatic Censoring in K-Distribution for Multiple Targets Situations <i>N. Boudemagh, Z. Hammoudi</i>	60
Fuzzy Method for Suppressing of Different Noises in Color Videos <i>Volodymyr Ponomaryov</i>	65
Temporal Data Approach Performance <i>Michal Kvet</i>	75
Alternative Approach to Enable RTSP-based Services with Dynamic Quality of Service over 4G LTE Mobile Networks <i>Andrei Rusan, Radu VasIU</i>	84
Control of Interferograms Image of Deformed Object Samples by Non Destructive Control as Optical Method <i>R. Daira</i>	90
A Modified Adaptive Line Enhancer for Noisy Speech Signals <i>Maha Sharkas, M. Essam Khedr, Amr Nasser</i>	95

A Comprehensive Analysis of XML and JSON Web Technologies <i>Zia Ul Haq, Gul Faraz Khan, Tazar Hussain</i>	102
System for the Detection Earthquake Victims – Construction and Principle of Operation <i>C. Buzduga, A. Graur, C. Ciufudean, V. Vlad</i>	110
Question-Answering Systems in the Specific Domain of E-Government <i>A. Beltrán, S. Ordoñez, S. Monroy, L. Melo, N. Duarte</i>	116
Monitoring Metropolitan City Air-quality using Wireless Sensor Nodes based on ARDUINO and XBEE <i>Ali Al-Dahoud, Mohamed Fezari, Ismail Jannoud, Thamer AL-Rawashdeh</i>	121
Extending the Matrix Vector Transition Net Approach for Modeling Interaction <i>A. Spiteri Staines</i>	126
Location Search by using Phonetic Algorithm with Location-Based Service <i>Kittiya Poonsilp, Attakorn Poonsilp</i>	133
Improved Non-local Algorithm with Reliability of Neighbor Pixel <i>J. Lee, J. Jeong</i>	139
Urban Traffic Management Approach based on Ontology and VANETs <i>H. Touluni, B. Nsiri, M. Boulmalf, T. Sadiki</i>	145
Integrated Visual-Perception Real-Time Monitoring System <i>Jian-Wei Li, Fu-Syuan Yang, Yi-Chun Chang, Yen-Lun Chiu</i>	150
Novel M-ary PPM Time Hopping Scheme for UWB Communications <i>Said Ghendir, Salim Sbaa, Riadh Ajjou, Ali Chemsal, A. Taleb-Ahmed</i>	156
Interoperability for an Obserbatory of Habits and Healthy Life Styles Ralated with Physical Activity <i>Andrea Torres Ruiz, Fernando Prieto B., Jose Arturo Lagos, Nixon Duarte, Rosmary Martinez, Juan Pablo Moreno, Aldo Vilardy, Bryan Toro</i>	161
A Compact Microstrip Lowpass Filter using a Stepped Impedance Hairpin Resonator with Radial Stubs <i>M. Samadbeik, B. F. Ganji, A. Ramezani</i>	167
Modification of the Cryptographic Algorithms, Developed on the Basis of Nonpositional Polynomial Notations <i>Rustem G. Biyashev, Saule E. Nyssanbayeva, Yenlik Ye. Begimbayeva, Miras M. Magzom</i>	170
Experimental Human Machine Interface System based on Vowel and Short Words Recognition <i>Mohamed Fezari, Ali Al-Dahoud</i>	177

On DC/DC Voltage Buck Converter Control Improvement through the QFT Approach	183
<i>Luis Ibarra, Israel Macías, Pedro Ponce, Arturo Molina</i>	
Irregular Segmentation Technique for the Image Compression using Stochastic Models	191
<i>Benabdellah Yagoubi</i>	
Efficient Media Digital Library Design of Summarized Video based on Scalable Video Coding for H.264 (MDLSS)	195
<i>Hesham Farouk, Kamal EIDahshan, Amr Abozeid, Mayada Khairy</i>	
Speech Enhancement using Rao-Blackwellised Particle Filtering of the Real and Imaginary DFT Coefficients Part	200
<i>M. Meddah, A. Amrouche, A. Taleb-Ahmed</i>	
Swarm Intelligence Optimization of Lee Radio-wave Propagation Model for GSM Networks in Irbid	207
<i>M. S. H. Al Salameh, M. M. Al-Zu'bi</i>	
A Recognition and Synthesis Environment for the Arabic Language	213
<i>Tebbi Hanane, Hamadouche Maamar, Azzoune Hamid</i>	
Image Encryption using Development of Chaotic Logistic Map based on Feedback Stream Cipher	220
<i>Hossam Eldin H. Ahmed, Ayman H. Abd El-aziem</i>	
Multi-Element Circuits based on LCLC Resonant Tank - Theory and Application	230
<i>Branislav Dobrucky, Juraj Koscelnik</i>	
Parallel Adaptive Arbiter for Improved CPU Utilization and Fair Bandwidth Allocation	241
<i>M. Nishat Akhtar, Junita Mohamad-Saleh</i>	
Predictive Robots Programming based on Imitation Strategy	253
<i>A. Fratu, M. Fratu</i>	
Unifying Geometric Features and Facial Action Units for Improved Performance of Facial Expression Analysis	259
<i>Mehdi Ghayoumi, Arvind K. Bansal</i>	
Authors Index	267

Prediction of Cancer Behavior Based on Artificial Intelligence

Shayma M. Al-Ani, Maysam Abbod

Abstract— Cancer has been one of the most famous conditions discussed and researched about throughout the human history. Some of the earliest medical records regarding cancer are dated back to around 1600 BC. Cancer is a general condition which is subdivided into a group of conditions that are concerned with an abnormal growth in the cells within an organ or a tissue with the chance of spreading and invading other parts of the body. Nowadays, there is a growing number of cancer patients and with this increase arises the necessity for new techniques to accurately diagnose and predict cancer in its different forms and thus playing a huge part in improving the quality of life. Moreover, techniques that depend on the principle of intelligent systems and artificial neural networks are proven to be very efficient in the field of cancer research.

Keywords—ANN, Cancer prediction, Ensemble model.

I. INTRODUCTION

CANCER has been known all the way through human history. Other names for cancer are malignant tumor or malignant neoplasm [1]. Genetic heritage, tobacco and alcohol intake, obesity, radiation exposure as well as having a poor and inactive lifestyle are some causes for abnormal cell growth and thus providing higher risks of getting cancer. For such reasons, cancer is considered as one of the most dangerous and unpredictable diseases nowadays. Much of research is done on the prevention of cancer and cancer treatment [2]. Innovative and modern technology is being implemented in the goal of providing proper diagnosis, prediction and in some cases treatments for cancer. Artificial intelligence is one of the methods for approaching cancer and understanding its nature [3]. One of the most popular types of cancer being approached by artificial intelligence is breast cancer in females.

An Artificial Neural Network (ANN) is an imitation to the basic human brain operation and it is an interconnected neurons system that is capable of computing values using mathematical functions in which they determine the activation of the neuron[4][5]. To adapt to the environmental changes, a learning system has to change itself. In addition, Multi-layers

ANNs are complex neural networks providing a nonlinear relationship of input-to-output results. Multi-layer ANNs comprise of an input layer, a hidden layer and an output layer as illustrated in Fig. 1. Basically, the input layer provides an input value to the network and each of the input cells has a weighting factor, which identifies the effect of the cell on the network[6]. As for the hidden and output cells, they represent a function where the hidden layer is first computed, and then the results are used in computing the output layer.

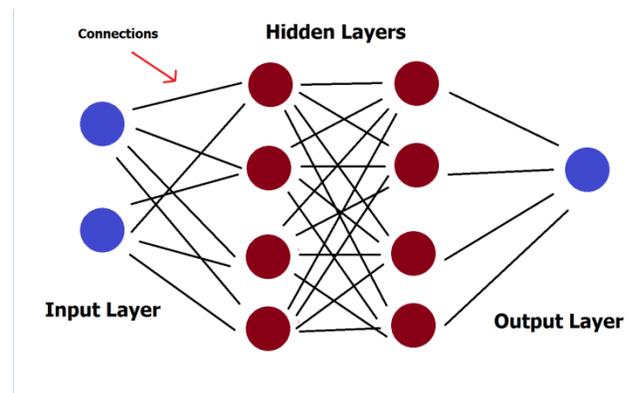


Fig. 1 Artificial neural network model.

The data presented in this paper is patients suffering from bladder cancer . This database is obtained for xx patient, each patient is represented with different information as input data such as the type of tumor, patient details (sex, age, tobacco consumption) in addition to protein expression (p53, msh2, mlh1) and DNA mutations (bat25, bat26, mfd15, apc, d2s123). Furthermore, the output data will be represented with the actual behavior of the tumor such as how long did and how many times the tumor went back to appear. In addition the time it took to advance to other stages, the time the cancer spread and lead to patient's death and whether or not the cancer was the cause of death or there might be other causes of death (e.g. complications).

II. METHODOLOGY

The proposed ANN based prediction algorithm accurately predicts the patient's cancer records output by employing the ensemble method shown in Table I. In In this method, the patient's record is equally divided into 10 window groups. In

S Alani and M Abbod are with the Department of Electronic and Computer Engineering, Brunel University, Uxbridge, UB8 3PH, UK, (email: Shayma.Al-Ani@brunel.ac.uk, maysam.abbod@brunel.ac.uk)

this method, the average of 10 ANN network functions under different combinations of 10 window groups have been used in order to find out the predicted output, in addition to improving the prediction performance of a model with more accurate results.

The proposed model is trained by using three different ANN networks which are cascade-forward back propagation network (NEWCF), feed-forward input time-delay back propagation network (NEWFFTD), and fitting network (NEWFIT), each network is trained by using ensemble methods under different combination of groups by applying two methods which are the averaging and voting methods.

The averaging method is one of the major types of static committee machines. The network design for such method depends upon mean average of the networks. In addition, ensemble averaging depends on the mean average networks results. So in general, the whole idea of averaging method can be summarized by the following; generating N experts each having their initial values which are chosen from a random distribution. After that, each expert will be separately trained separately and finally, they are combined and their values are averaged.

As for the voting method, it does not consider the level of significance by each network. This as a result, allows simple integration of all different sorts of network architectures. Majority voting is a simple voting method in which a group of unlabeled instance are performed depending on the class with the most frequent votes. This technique has been widely used to compare newly proposed methods.

TABLE I. SLIDING WINDOW METHOD.

W1	W2	W3	W4	W5	W6	W7	W8	W9	W10
W2	W3	W4	W5	W6	W7	W8	W9	W10	W1
W3	W4	W5	W6	W7	W8	W9	W10	W1	W2
W4	W5	W6	W7	W8	W9	W10	W1	W2	W3
W5	W6	W7	W8	W9	W10	W1	W2	W3	W4
W6	W7	W8	W9	W10	W1	W2	W3	W4	W5
W7	W8	W9	W10	W1	W2	W3	W4	W5	W6
W8	W9	W10	W1	W2	W3	W4	W5	W6	W7
W9	W10	W1	W2	W3	W4	W5	W6	W7	W8
W10	W1	W2	W3	W4	W5	W6	W7	W8	W9

In the case of Artificial Neural Network model, the neuron behaves as an activation function $f(.)$ producing an output $y = f(net)$, where net is the cumulative input stimuli to the neuron and f is typically a nonlinear function of net, where x_i indicates the inputs and w_i indicate the weighting parameters.

$$net = x_1w_1 + x_2w_2 + x_3w_3 = \sum_{i=1}^3 x_iw_i \quad (1)$$

Output performances of the proposed algorithms are analysed using various parameters such as Sensitivity, Specificity, Accuracy, Receiver Operator Curve (ROC), Area Under the Curve (AUC) and Mean Square Error (MSE) value.

Regression model is a statistical model for estimating the relationships among variables. It includes many techniques for modeling and analyzing several variables and statistical models to compare it with the ANN model, output performances of them are analysed using various parameters like Sensitivity, Specificity, Accuracy, Roc, AUC and MSE value.

III. FINDINGS

The proposed ANN models are trained using three different ANN networks, namely NEWCF, NEWFFTD and NEWFIT[7][8]. First of all, randomly dividing the networks into two groups called training records and testing records[9][10][11]. Training records group contains about 70% of the total records, which are used to train the ANN by using 80% for training and 20% for validation of the ANN networks. The trained ANN networks are used to predict the output parameter of testing records group which contain 30% of total records .

Moreover, three different methods have been used, average, voting, and regression model[12][13][14]. Table II shows the input variables used in the modeling analyses. Tables III and IV show the performance of three methods for various ANN training networks in which it follows the principle of 70% training and 30% testing, while Tables V and VI show the predicted patients records for three methods and three different ANN trained networks using ensemble method.

Sensitivity relates to the test’s ability to identify positive results which measures the proportion of actual positives which are correctly identified as such.

While specificity relates to the test’s ability to identify negative results, which measures the proportion of negatives which are correctly identified.

The accuracy is the proportion of true results (both true positive and true negative) in the population.

$$\text{Sensitivity} = TP / (TP + FN)$$

$$\text{Specificity} = TN / (TN + FP)$$

$$\text{Accuracy} = (TP + TN) / (TP + FP + TN + FN)$$

TABLE II. INPUT VARIABLES USED IN THE MODELING ANALUSES.

Input Variables
1. age
2. sex
3. chest pain type (4 values)
4. resting blood pressure
5. serum cholesterol in mg/dl
6. fasting blood sugar > 120 mg/dl
7. resting electrocardiographic results (values 0,1,2)
8. maximum heart rate achieved
9. exercise induced angina
10. oldpeak = ST depression induced by exercise relative to rest
11. the slope of the peak exercise ST segment
12. number of major vessels (0-3) colored by flourosopy
13. thal: 3 = normal; 6 = fixed defect; 7 = reversible defect

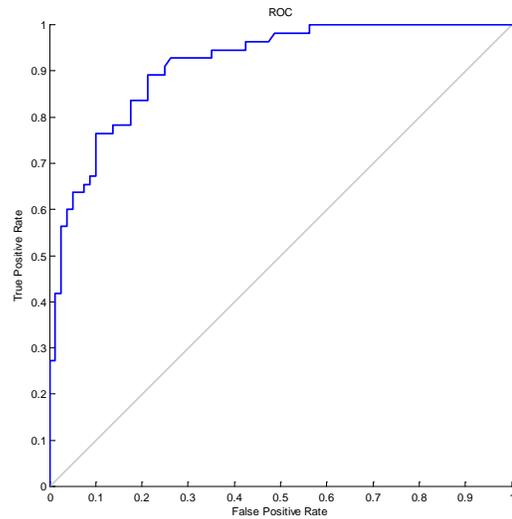


Fig.2 Average method train case.

TABLE III. PERFORMANCE OF ANN NETWORKS TRAIN RECORDS RESULTS ANALYSIS.

Methods	Sensitivity	Specificity	Accuracy
Average	78.1818	86.25	82.963
Voting	74.5455	88.75	82.963
Regression Model	63.6364	73.75	69.630

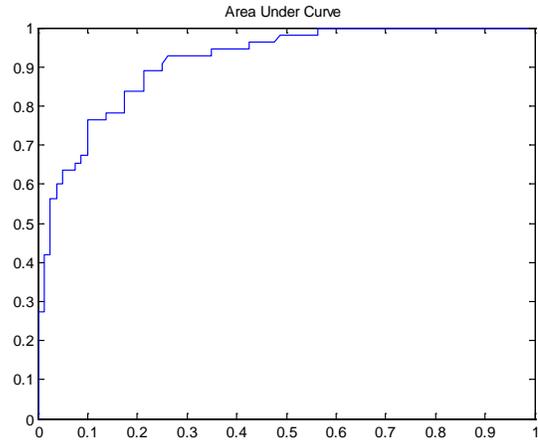


Fig.3. Average method train case.

TABLE IV. PERFORMANCE OF ANN NETWORKS TRAIN RECORDS RESULTS ANALYSIS MSE AND AUC VALUES.

Methods	MSE Value	AUC
Average	0.1378	0.9009
Voting	0.1330	0.9048
Regression Model	0.1995	0.7398

TABLE V. PERFORMANCE OF ANN NETWORKS TEST RECORDS RESULTS ANALYSIS.

Methods	Sensitivity	Specificity	Accuracy
Average	65.2174	74.2857	70.6897
Voting	60.8696	77.1429	70.6897
Regression Model	52.609	50	51.034

TABLE VI. PERFORMANCE OF ANN NETWORKS TEST RECORDS RESULTS ANALYSIS MSE AND AUC VALUES.

Methods	MSE Value	AUC
Average	0.1908	0.7280
Voting	0.1956	0.7193
Regression Model	0.1335	0.5745

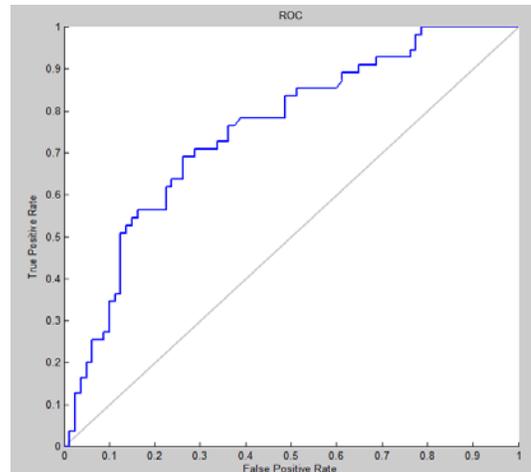


Fig.4 Regression model train case.

Figs. 2-5 show the ROC plot of bladder cancer train records of average method and regression model for NEWCF, NEWFFTD and NEWCF networks, while Figs. 6-9 show the ROC plot of bladder cancer test records of average method and regression model.

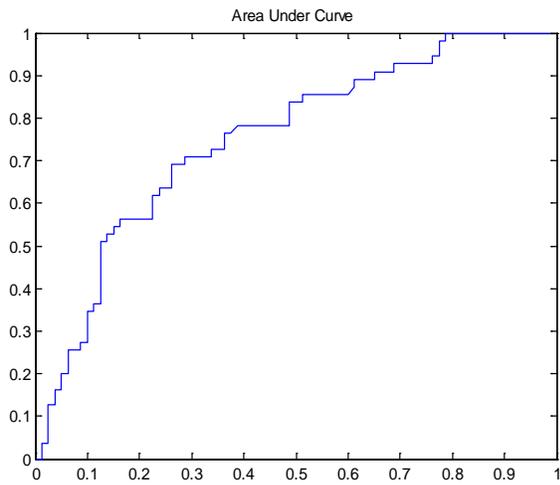


Fig.5 Regression model train case.

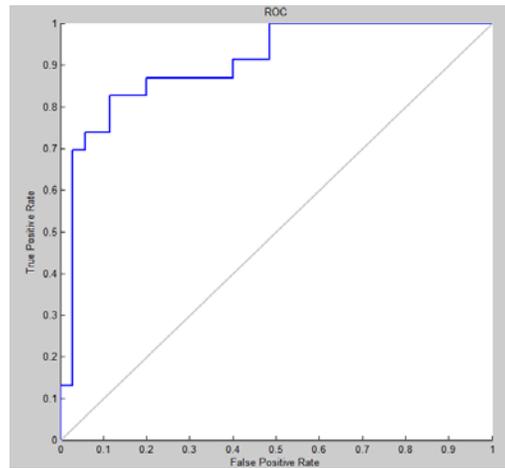


Fig.8 Regression model test case.

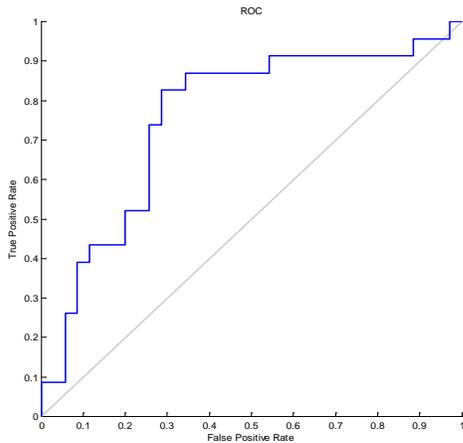


Fig.6. Average method test case.

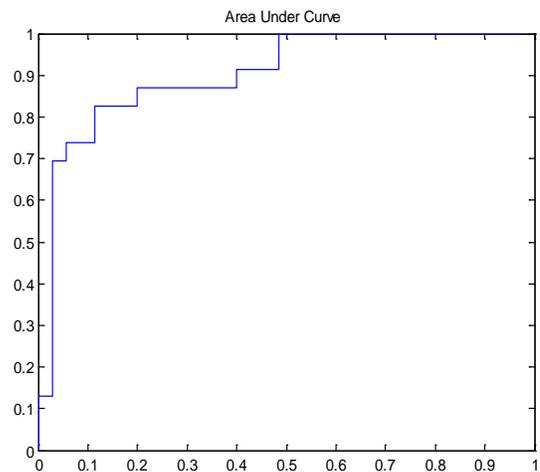


Fig.9 Regression model test case.

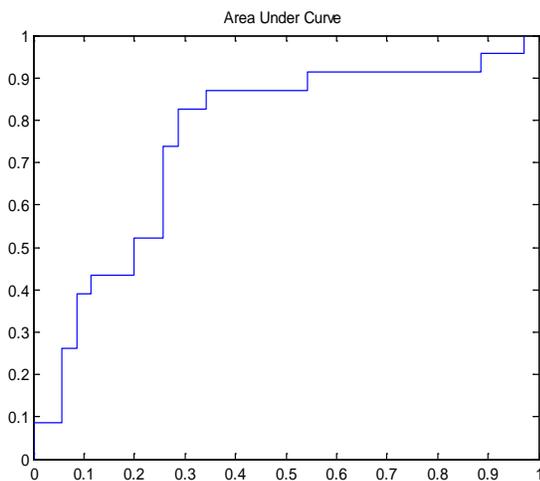


Fig.7 Average method test case.

IV. CONCLUSIONS

The proposed ensemble model, the artificial neural network algorithm using two methods averaging, voting and for various artificial neural network functions such as feed-forward input time-delay back-propagation network, cascade-forward back-propagation network and radial basis network, successfully and accurately predicted patient's records. Output performances of records are analyzed using various parameters such as Sensitivity, Specificity, Accuracy, ROC, AUC and MSE value and the results show that artificial neural network methods obtain better predictive performance than could be obtained from regression models and that was all based on the different validations of the artificial neural networks.

REFERENCES

- [1] (2014). "What is Cancer". *National Cancer Institute*. [Online]. Available: <http://www.cancer.gov/cancertopics/cancerlibrary/what-is-cancer>
- [2] (2014). "Definition of Bladder Cancer". *National Cancer Institute*. [Online]. Available <http://www.cancer.gov/cancertopics/types/bladder>.

- [3] J.A. Cruz and D. S. Wishart. (Feb, 2007). Applications of Machine Learning in Cancer Prediction and Prognosis. *Cancer Informatics*. [Online]. 2(2006), 59-77. Available: <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2675494/>.
- [4] Z. Chi, Z. Lu and F. Chan, "Multi-channel handwritten digit recognition using neural networks", Circuit and Systems ISCAS'97 proceeding of 1997 IEEE International Symposium, Vol.1, 1997, pp. 625-628.
- [5] K. Tadashi; U. Junji; T. Shoichiro, "Hybrid GMDH-type neural network using artificial intelligence and its application to medical image diagnosis of liver cancer", *System Integration (SII), 2011 IEEE/SICE International Symposium on 2011*, pp. 1101-1106.
- [6] P.J.G. Lisboa, T.A. Etchells, I.H. Jarman and M.S.H. Aung, "Time-to-event analysis with artificial neural networks: An integrated analytical and rule-based study for breast cancer", *Neural Networks, 2007. IJCNN 2007. International Joint Conference on 2007*, pp: 2533-2538.
- [7] P. Melville, "Creating Diverse Ensemble Classifiers," PhD proposal, Texas Univ., Texas, United States, 2003.
- [8] Robi Polikar (2009) Ensemble learning. *Scholarpedia*, 4(1): 2776.
- [9] Opitz, D.; Maclin, R. (1999). "Popular ensemble methods: An empirical study". *Journal of Artificial Intelligence Research* **11**: 169–198.
- [10] U. Naftaly, N. Intrator, and D. Horn. "Optimal ensemble averaging of neural networks." *Network: Computation in Neural Systems* 8, no. 3 (1997): 283–296.
- [11] L. Rokach. (Nov, 2009). Ensemble-based classifiers. *Springer Science+Business Media*. [Online]. 33 (2010), 1-39. Available: <http://www.ise.bgu.ac.il/faculty/liort/AI.pdf>.
- [12] (2014). "Regression analysis". Wikipedia. [Online]. Available http://en.wikipedia.org/wiki/Regression_analysis#Regression_models.
- [13] C.Chen, C. Hsu, H.Chiu and H.Rau, "Prediction of survival in patients with livercancer using artificial neural networks and classification and regression trees", *Natural Computation (ICNC), 2011 Seventh International Conference on 2011*, pp. 811-815.
- [14] Janghel, R. R.; Shukla, Anupam; Tiwari, Ritu; Kala, Rahul, 2010. Breast cancer diagnosis using Artificial Neural Network models. *Information Sciences and Interaction Sciences (ICIS), 2010 3rd International Conference on 2010*, pp. 89-94.

Modeling and Analysis of Elapsed Time and Energy Consumption of Interactive Applications in Mobile Cloud Computing Environments

Young-Chul Shim

Abstract—We consider the execution of an interactive application in a mobile cloud computing environment including a local cloudlet and a remote cloud. We introduce models for computing elapsed time and energy consumption of an interactive application when its offloadable portion is executed on a mobile device, a local cloudlet, or a remote cloud. Applying practical numbers to some of the parameters in the models, we derive conditions under which computation offloading to either a local cloudlet or a remote cloud becomes profitable in terms of elapsed time and energy consumption.

Keywords—Computation Offloading, Elapsed Time, Energy Consumption, Mobile Cloud Computing.

I. INTRODUCTION

SINCE the emergence of the concept of cloud computing, it is getting more widely adopted and deployed in the IT industry sector and receiving more attention from computer scientists and engineers. NIST defines cloud computing to be a model for enabling ubiquitous, convenient, on-demand network access to a shared pool of configurable computing resources (e.g., networks, servers, storage, applications, and services) that can be rapidly provisioned and released with minimal management effort or service provider interaction [1]. Examples of well-known cloud computing systems include Amazon AWS, Microsoft Azure, Google AppEngine, and Rackspace CloudServers [2].

Mobile devices such as smartphones and tablet PCs are becoming more and more essential part of our life. People use mobile devices to not only communicate with others but also run application programs and store information gathered from their daily activities. But mobile devices suffer from the critical drawback, lack of resources. They have limited computation and storage capability and, more seriously, are very much restricted in their battery power.

To continue enjoying the convenience of mobile devices while making up for their weaknesses, people introduced the concept of mobile cloud computing [3]-[5]. In mobile cloud computing environments, computation and/or storage requirements on mobile devices can be offloaded to outside cloud computing environments and mobile users can finish their programs faster, store more information, and save the

battery power of their mobile devices.

A large fraction of applications that mobile users run in the mobile cloud computing environment will be interactive. People play a chess game with their mobile devices but finding the optimal next movement requires tremendous amount of computation and, therefore, has been usually calculated on very powerful machines. In an application called content-based image retrieval, people may want to retrieve photos containing a certain image from a large file of photo collections, which will require a large amount of image matching computation. People may design a product such as a house or a machine with their mobile devices and performing this mobile computer-aided design activity usually requires solving many complex partial differential equations and demands extremely fast floating point computation. People can also benefit from the technology of augmented reality with their mobile devices. When they go to a store to buy a certain product such as furniture, they will want to check whether the chosen furniture goes well with their house. They can get the information about the chosen furniture from the store computer, send it to the site which stores information about their house, request to compose these two data, and view the results from many different aspects to make the buying decision.

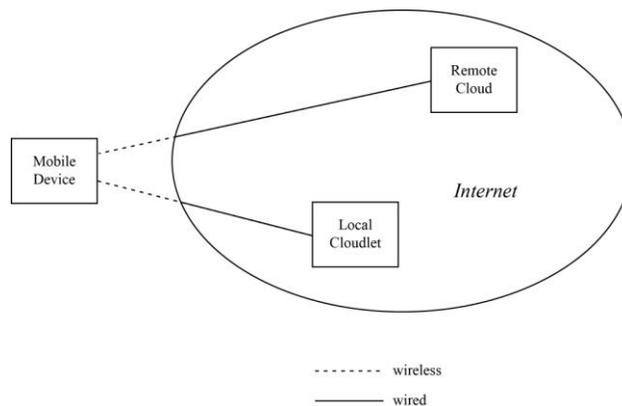


Fig. 1. a mobile cloud computing environment

Cloud computing systems may be located far away from a mobile user and the communication latency between them can become non-negligible. If running an application requires a large number of interactions, the non-negligible latency may produce a negative effect on the response time and battery power saving. Recently the concept of cloudlets is introduced.

This work was supported by 2013 Hongik University Research Fund
Young-Chul Shim is with Hongik University, Seoul, South Korea
(82-2-320-1695; e-mail: yeshim@hongik.ac.kr).

Cloudlets are decentralized and widely-dispersed Internet infrastructure whose compute cycles and storage resources can be leveraged by nearby mobile users [6]. As server machines are becoming cheaper, we can make each node of a cloudlet as powerful as that of a cloud center and deploy a large number of cloudlets. With this we envision a mobile cloud computing environment as Fig. 1 in which a mobile user can choose the location for running his program among a mobile device, a local cloudlet, or a remote cloud.

Some researchers derived conditions which must be satisfied to make computation offloading profitable to a mobile user [7], [8]. But their models about the computation and data transfer requirements and power consumption on mobile devices are too simplistic. And they do not consider propagation delay and, therefore, do not make distinction between local cloudlets and remote clouds. Some researchers developed static or dynamic program partitioning methods which help determine which portion of a program is to be offloaded to reduce response time and energy consumption [9], [10]. But their results apply to only specific programs in specific environments and do not provide general guidelines for when a portion of program should be offloaded to which outside environments, local cloudlets or remote clouds, and how much benefit can be obtained from the computation offloading.

In this paper we consider interactive applications running in a mobile cloud computing environment including both local cloudlets and remote clouds. We provide models for calculating elapsed time and energy consumption of an interactive application. In the model we consider factors including computation and data requirements of an application, processing speed, bandwidth, propagation delay, and energy usage of a mobile device in three states: computation, data transfer, and idle. The application can be run completely on a mobile device or the offloadable portion be executed on either a local cloudlet or a remote cloud. Using the proposed model and applying typical values to the parameters in the model, we derive conditions that must be satisfied if executing the offloadable portion on either a local cloudlet or a remote cloud is to be profitable.

The rest of the paper is organized as follows. Section 2 explains the execution model of an interactive application in a mobile cloud computing environment. The models for computing elapsed time and energy consumption of a program in various environments are provided in Section 3. Section 4 describes the conditions that must be satisfied if computation offloading to be profitable and is followed by the conclusion in Section 5.

II. EXECUTING AN INTERACTIVE APPLICATION IN A MOBILE CLOUD COMPUTING ENVIRONMENT

Interactive applications that will be executed in a mobile cloud computing environment may have different characteristics. An application such as a chess game does not need any initial data representing the initial state to be loaded to invoke the program. But an application such as the content-based image retrieval requires initial data which is usually a very size file of photo collection, on which various

kinds of image matching operations will be executed. A mobile CAD program can lie in between. If a completely new design is to be started, there is no initial design data to be loaded but if an unfinished design is to be resumed, the design data that have been accumulated from the previous design activities should be loaded. Even for the applications for which initial data should be loaded, the location of initial data can vary. They can be on the mobile node with which a user will run the interactive application or they can be located on a remote cloud. Sometimes when data are collected and stored at a mobile node, they can be synchronously copied to a remote cloud server to prevent data loss.

After the initial data loading phase, which is performed only when required, comes the interactive computation phase. In the interactive computation phase, an activity consisting of three steps are repeated until the user terminates the program. The three steps are input, process, and output steps. In the input step, input data are obtained from the input devices such as a microphone, a camera, a keypad, etc and can then be preprocessed. The preprocessing can consist of various kinds of activities and one example is data compression which is performed to reduce the amount of data before being sent to a remote node. Because the input step requires the use of input device of a mobile node, it should be executed on the mobile node. In the second step, the process step, the input data is processed to produce the requested result. Examples of this step include finding the optimal movement in a chess game, retrieving photos matching a certain image from a large file of photo collection, or performing a requested design action which requires a huge amount of complex calculations. This second step tends to be very compute-demanding and is a very good candidate for being migrated to and executed on a remote fast node. The last step, the output step, receives the result from the process step and presents it to the user using the output devices on the mobile node. Therefore, this step should be performed on the mobile node. If the process step is executed on a remote fast node, a single iteration of the three step computation phase will proceed as in Fig. 2.

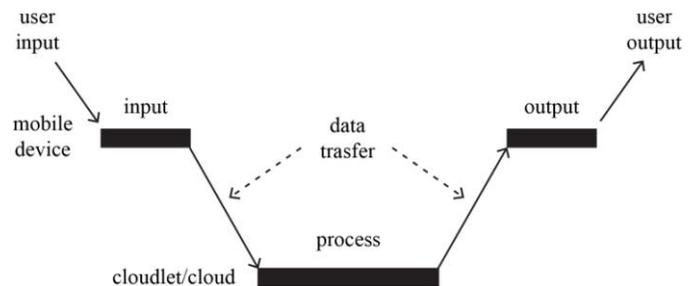


Fig. 2. remote execution of a process step

III. MODELING ELAPSED TIME AND ENERGY CONSUMPTION

If an interactive application requires C instructions for computation, these instructions can be divided into two parts: one that must be executed on a mobile node and the other that can be offloaded to a remote node, possibly a very fast cloud node. With the execution model described in the previous

section, the input and output steps belong to the first part and the process step belongs to the second part. In this paper we consider two candidates for the outside computation node: a local cloudlet and a remote cloud. A local cloudlet may have less computation capacity than a remote cloud but is definitely closer from the mobile node and, therefore, incurs much smaller propagation delay when exchanging data with a mobile node.

We consider two kinds of interactive applications: one which does not require initial data and the other which requires initial data. For those applications which require initial data, such as a large size file of photo collection in the content-based image retrieval application, we consider two places for originally storing the initial data: a mobile node or a remote cloud.

In this section we consider two kinds of interactive applications, two places for originally storing the initial data when required, and three places for executing the instructions that can be performed outside a mobile node. We calculate the total elapsed time and energy consumption requirements for various cases and compare them.

Now we will define symbols that will be used in this paper.

- N : Represents the number how many times the input-process-output stage is repeated during the interactive computation phase.
- C_M : Number of instructions that should be executed on a mobile node
- C_C : Number of instructions that can be offloaded to an outside node such as a local cloudlet or a remote cloud. These instructions are called offloadable instructions in this paper and they correspond to the process step in Fig. 2. If there are N iterations in the computation phase, each process step executes $C_C/N = C_C'$ instructions on average.
- S_M : Speed of a mobile node in terms of instructions/second
- S_C : Speed of a cloud or cloudlet in terms of instructions/second. If we have to distinguish a local cloudlet from a remote cloud, we use S_{CL} for the speed of a local cloudlet and S_{CR} for a remote cloud.
- D_I : The size of the initial data needed for the operation of an interactive application.
- D_C : The amount of data that a mobile node exchanges with an outside node during the interactive computation phase. They are the data moved to and from the process step in Fig. 2. If there are N iterations in the computation phase, each process exchanges data of the amount $D_C/N = D_C'$.
- B_M : The bandwidth between a mobile and an outside node.
- B_C : The bandwidth between a local cloudlet and a remote cloud. This can be much larger than B_M .
- T_P : The propagation delay between a mobile and an outside node. This delay includes not only the delay proportional to the distance between two communicating nodes but also delays at the intermediate communication devices such as switches and routers. If we have to distinguish a local cloudlet from a remote cloud, we use T_{PL} for a local cloudlet

and T_{PR} for a remote cloud.

- P_C : The energy consumed during computation by a mobile node in watts.
- P_I : The energy consumed during idle time by a mobile node in watts.
- P_T : The energy consumed to send and receive data by a mobile node in watts. Usually transmitting power is higher than receiving power, but we assume that they are the same in this paper for simplicity.
- MN : Mobile node
- LC : Local cloudlet
- RC : Remote cloud

Now we will present models for elapsed time and energy consumption for the following 3 cases depending upon the necessity and location of initial data that are required for the operation of an interactive application.

- Case 1: Initial data is not needed
- Case 2: Initial data is located at a mobile node
- Case 3: Initial data is located at a remote cloud

A. No Initial Data

When all the computation is performed on a mobile node, all the (C_M+C_C) instructions are executed with the speed of S_M . There is no data exchange. So, the elapsed time T and the consumed energy E are

$$T(MN) = (C_M+C_C)/S_M \quad (1)$$

$$E(MN) = P_C \times (C_M+C_C)/S_M \quad (2)$$

When all the C_C instructions are executed on either a local cloudlet or a remote cloud, in addition to the time spent to execute instructions, extra time is needed to exchange the data of size D_C . The data exchange time consists of a data transmission time, which is defined to be $(data\ size)/bandwidth$, and the data propagation time between two nodes. Each iteration of a input-process-output step consists of two data transmissions, one for input and the other for output, and we assume that this three step operation is repeated N times. Thus the whole interactive computation phase involves $2N$ data transmissions and requires $2N \times T_P$ seconds for the whole propagation time. Therefore, the elapsed time and the consumed energy at the mobile node with C_C instructions offloaded to either a local cloudlet or a remote cloud become

$$T(LC\ or\ RC) = C_M/S_M + C_C/S_C + 2N \times T_P + D_C/B_M \quad (3)$$

$$E(LC\ or\ RC) = P_C \times (C_M/S_M) + P_I \times (C_C/S_C + 2N \times T_P) + P_T \times (D_C/B_M) \quad (4)$$

Note that in (4) when instructions are executed on a local

cloudlet or a remote cloud and exchanged data are moved to and from the mobile node (this means that data are in the medium), the mobile node stays in the idle mode. For (3) and (4), S_C and T_P become S_{CL} and T_{PL} if a local cloudlet is used and become S_{CR} and T_{PR} if a remote cloud is used.

To compare the elapse time and consumed energy, we compute $T(MN) - T(LC \text{ or } RC)$ and $E(MN) - E(LC \text{ or } RC)$ as follows.

$$T(MN) - T(LC \text{ or } RC) = C_C(1/S_M - 1/S_C) - 2N \times T_P - D_C/B_M \quad (5)$$

$$E(MN) - E(LC \text{ or } RC) = C_C \times (P_C/S_M - P_C/S_C) - P_I \times (2N \times T_P) - P_T \times (D_C/B_M) \quad (6)$$

Now we consider the case in which C_C instructions are offloaded to either a local cloudlet or a remote cloud. In order to compare these two cases we compute $T(RC) - T(LC)$ and $E(RC) - E(LC)$ as follows.

$$T(RC) - T(LC) = C_C \times (1/S_{CR} - 1/S_{CL}) + 2N \times (T_{PR} - T_{PL}) \quad (7)$$

$$E(RC) - E(LC) = P_I \times \{C_C (1/S_{CR} - 1/S_{CL}) + 2N \times (T_{PR} - T_{PL})\} \quad (8)$$

We see that (8) is the same as (7) except that (8) is P_I times of (7).

B. Initial Data at a Mobile Node

In this subsection we consider a case in which initial data is required for the application to operate and the data reside at a mobile node.

If all the computation is performed at a mobile node without any computation offloaded to an outside node, the elapsed time and the consumed energy at the mobile node are the same as (1) and (2), respectively

If C_C instructions are to be executed at an outside node, the initial data D_I should be transferred to that outside node. This incurs more time $2 \times T_P + D_I/B_M$ and more energy $2 \times P_I \times T_P + P_T \times (D_I/B_M)$. For simplicity we assume $N \gg$ and we can ignore term due to $2 \times T_P$. Then, the elapsed time and the consumed energy at the mobile node with C_C instructions offloaded to either a local cloudlet or a remote cloud become

$$T(LC \text{ or } RC) = C_M/S_M + C_C/S_C + 2N \times T_P + (D_I + D_C)/B_M \quad (9)$$

$$E(LC \text{ or } RC) = P_C \times (C_M/S_M) + P_I \times (C_C/S_C + 2N \times T_P) + P_T \times ((D_I + D_C)/B_M) \quad (10)$$

To compare the elapsed time and consumed energy, we compute $T(MN) - T(LC \text{ or } RC)$ and $E(MN) - E(LC \text{ or } RC)$ as follows.

$$T(MN) - T(LC \text{ or } RC) = C_C \times (1/S_M - 1/S_C) - 2N \times T_P - (D_I + D_C)/B_M \quad (11)$$

$$E(MN) - E(LC \text{ or } RC) =$$

$$C_C \times (P_C/S_M - P_C/S_C) - P_I \times (2N \times T_P) - P_T \times ((D_I + D_C)/B_M) \quad (12)$$

If we want to see whether offloading to a local cloudlet is better than offloading to a remote cloud or not, we have to compute $T(RC) - T(LC)$ and $E(RC) - E(LC)$. By replacing S_C and P_T with S_{CR} and P_{TR} for a remote cloud and S_{CL} and P_{TL} for a local cloudlet in (9) and (10), we can easily see that these two equations become the same as (7) and (8), respectively. This is quite an obvious result because irrespective of whether the mobile node transfers initial data to a local cloudlet or a remote cloud, it pays almost the same cost in terms of elapsed time and consumed energy.

C. Initial Data at a Remote Cloud

In this section we assume that the initial data is stored in a remote cloud server which offers both computation and storage service.

We first assume that the initial data is transferred to a mobile node and all the computation is performed at the mobile node. The elapsed time and the consumed energy become

$$T(MN) = (C_M + C_C)/S_M + 2 \times T_{PR} + D_I/B_M \quad (13)$$

$$\approx (C_M + C_C)/S_M + D_I/B_M$$

$$E(MN) = P_C \times (C_M + C_C)/S_M + P_I \times (2 \times T_{PR}) + P_T \times D_I/B_M \quad (14)$$

$$\approx P_C \times (C_M + C_C)/S_M + P_T \times (D_I/B_M)$$

The approximations in (13) and (14) were made assuming that the transmission delay is much larger than the propagation delay when moving a large size initial data.

If the initial data stays as the remote cloud and C_C instructions are offloaded to the remote cloud, the elapsed time and the consumed energy become

$$T(RC) = C_M/S_M + C_C/S_{CR} + 2N \times T_{PR} + D_C/B_M \quad (15)$$

$$E(RC) = P_C \times (C_M/S_M) + P_I \times (C_C/S_{CR} + 2N \times T_{PR}) + P_T \times (D_C/B_M) \quad (16)$$

If the initial data is moved to a local cloudlet and C_C instructions are offloaded to the local cloudlet, the elapsed time and the consumed energy become

$$T(LC) = C_M/S_M + C_C/S_{CL} + 2N \times T_{PL} + D_C/B_M + D_I/B_C \quad (17)$$

$$E(LC) = P_C \times (C_M/S_M) + P_I \times (C_C/S_{CL} + 2N \times T_{PL} + D_I/B_C) + P_T \times (D_C/B_M) \quad (18)$$

In (17) and (18) we ignored the propagation delay required when moving initial data from the remote cloud to the local cloudlet.

To be able to analyze which location is better in terms of elapsed time and consumed energy, we calculate the differences in elapsed time and consumed energy for the possible three combinations as follows.

To compare no offloading and offloading to a remote cloud we have

$$T(MN)-T(RC) = C_C \times (1/S_M - 1/S_{CR}) - 2N \times T_{PR} + (D_I - D_C)/B_M \quad (19)$$

$$E(MN)-E(RC) = C_C \times (P_C/S_M - P_I/S_{CR}) - P_I \times 2N \times T_{PR} + P_T \times (D_I - D_C)/B_M \quad (20)$$

To compare no offloading and offloading to a local cloudlet we have

$$T(MN)-T(LC) = C_C \times (1/S_M - 1/S_{CR}) - 2N \times T_{PL} + (D_I - D_C)/B_M - D_I/B_C \quad (21)$$

$$E(MN)-E(LC) = C_C \times (P_C/S_M - C_C/S_{CR}) - P_I \times (2N \times T_{PR} + D_I/B_C) + P_T \times (D_I - D_C)/B_M \quad (22)$$

Finally to compare offloading to a remote cloud and offloading to a local cloudlet we have

$$T(RC)-T(LC) = C_C \times (1/S_{CR} - 1/S_{CL}) + 2N \times (T_{PR} - T_{PL}) - D_I/B_C \quad (23)$$

$$E(RC)-E(LC) = P_I \times \{C_C \times (1/S_{CR} - 1/S_{CL}) + 2N \times (T_{PR} - T_{PL}) - D_I/B_C\} \quad (24)$$

We see that (23) and (24) are the same equation except a multiplier P_I in (24).

IV. ANALYSIS WITH THE PROPOSED MODEL

In this section we use the model derived in the previous section to analyze which location, among a mobile, a local cloudlet, and a remote cloud, is the best place for executing offloadable C_C instructions for three cases as identified as in the three subsections in the previous section. To make the analysis more amenable, we survey data collected from some real mobile nodes, networks, and clouds and choose realistic numbers for some of the symbols used in the equations.

- S_M : ARM Cortex A7 processor which is used in many mobile nodes including smartphones and tablet PCs has the instruction execution speed of 2.85GIPS at 1.5GHz and we use this number.
- S_{CR} : Intel Xeon processors are popularly used in server machines and the Xeon 5690 processor has the speed of 84GIPS at 3.46GHz[11]. This number is almost 30 times of the speed of ARM Cortex A7. But a cloud server has much faster memory hierarchy and higher performance for floating point calculation and can provide more numbers of cores to an application. Therefore, in a real situation the speedup of a cloud node over a mobile can easily surpass 100~150.
- S_{CL} : Because a cloudlet can be assembled from the same kind of off-the-shelf server machines as in a remote cloud, the speed of each of component server machine can be almost the same. But a cloudlet will have less number of server

machines and the internal physical network connecting them and supporting operating environment software can be slower. Using admission control mechanisms in [12], we can assign almost the same number of cores to an application as in a remote cloud, although the number of simultaneously executable applications will be much lower in a cloudlet. With these observations, we guess that the speedup of a cloud over a cloudlet will not be very high and we can assume around 2~4 and call this speedup factor F in this paper.

- B_M : The bottleneck in the communication path between a mobile node and an Internet server is the wireless link directly connected to the mobile node. WiFi and 3G/4G are popular technologies for the wireless section. But it is well known that 3G/4G has longer delay and consumes more energy than WiFi [13] and, therefore, we assume the use of WiFi. IEEE80211.n has a data rate of 72.2 Mbps with 20 MHz bandwidth. Assuming around 40% throughput, we choose 30Mbps for B_M .
- T_{PR} : In [14], over 90% of users experience not greater than 25msec latency to access the closest Amazon cloud center. We choose T_{PR} to be 25msec.
- T_{PL} : Akamai is a content distribution network consisting of over 20,000 hosts and can be much closer to a user. [14] shows that over 90% of users experience not greater than 5msec latency to access the closest Akamai host. But in this paper, we want to be more aggressive assuming that much more cloudlet sites are deployed than Akamai and choose T_{PL} to be around 1msec.

Table 1 shows energy consumption data in watts for some mobile nodes. Because faster processors are adopted in more recent mobile nodes and consume more energy during computation and transmission mode we choose energy consumption values as follows.

- P_C : 1.0 watt
- P_I : 0.3 watt
- P_T : 2.0 watts

Table 1. energy consumption in mobile nodes

Mobile Devices	P_C	P_I	P_T
HP iPAQ PDA 400MHz [7]	0.9	0.3	1.3
Nokia N810 400MHz [8]	0.8		1.5
Openmoko Neo Freerunner [15]		0.27	
Galaxy S2 1.5GHz [16]		0.36	1.7

A. No Initial Data

For offloading to a local cloudlet to be beneficial, four inequalities $T(MN)-T(LC) = (5) > 0$, $E(MN)-E(LC) = (6) > 0$, $T(RC)-T(LC) = (7) > 0$, and $E(RC)-E(LC) = (8) > 0$ should be satisfied. But we know that last two inequalities are the same in the previous section. We start with the first inequality.

$$C_C(1/S_M - 1/S_{CL}) - 2N \times T_{PL} - D_C/B_M > 0$$

In the above inequality, because $S_{CL} \gg S_M$, $1/S_M - 1/S_C$ can be approximated to $1/S_M$. Then if we divide the inequality by N , it becomes

$$C_C'/S_M > 2 \times T_{PL} + D_C/B_M \quad (25)$$

The second inequality becomes

$$C_C \times (P_C/S_M - P_I/S_C) - P_I \times (2N \times T_P) - P_T \times (D_C/B_M) > 0$$

Using the fact $P_C/S_M \gg P_I/S_C$ and dividing both sides by N , we get

$$P_C \times C_C'/S_M > P_I \times 2T_P + P_T \times D_C/B_M > 0 \quad (26)$$

From the third inequality, we get

$$C_C \times (1/S_{CR} - 1/S_{CL}) + 2N \times (T_{PR} - T_{PL}) > 0$$

By dividing both sides by N we obtain

$$2(T_{PR} - T_{PL}) > C_C' \times (1/S_{CL} - 1/S_{CR}) = C_C' \times (F-1)/S_{CR} \quad (27)$$

After inserting parameter values assumed in the beginning of this section into (25), (26), (27), we obtain the following three inequalities.

$$C_C' > 5.7 \times 10^6 + 95 \times D_C' \quad (28)$$

$$C_C' > 1.7 \times 10^6 + 190 \times D_C' \quad (29)$$

$$C_C' < 4 \times 10^9 / (F-1) \quad (30)$$

We see that C_C' is lower-bounded by (28) and (29) and upper-bounded by (30).

For offloading to a remote cloud to be beneficial, four inequalities $T(MN) - T(RC) = (5) > 0$, $E(MN) - E(RC) = (6) > 0$, $T(RC) - T(LC) = (7) < 0$, and $E(RC) - E(LC) = (8) < 0$ should be satisfied. Following the similar calculations we obtain following three inequalities.

$$C_C' > 1.43 \times 10^8 + 95 \times D_C' \quad (31)$$

$$C_C' > 4.37 \times 10^7 + 190 \times D_C' \quad (32)$$

$$C_C' > 4 \times 10^9 / (F-1) \quad (33)$$

We see that for offloading to a remote cloud to be beneficial, more instructions should be executed during the process step.

From (28) to (30), we see that if the process step has a computation requirement large enough to compensate the cost needed to exchange data between a mobile node and a local

cloudlet, it is better to execute the process step on the local cloudlet. But From (31) to (33), if the computation requirement of a process step gets even larger, then it is better to offload the process step to a remote cloud.

B. Initial Data at a Mobile Node

If the initial data has a modest size, it can be stored at a mobile node. The conditions that must be satisfied to make offloading to a local cloud the best decision are $T(MN) - T(LC) = (11) > 0$, $E(MN) - E(LC) = (12) > 0$, $T(RC) - T(LC) = (7) > 0$, and $E(RC) - E(LC) = (8) > 0$. Following the same approximations and calculations, we obtain the following inequalities.

$$C_C' > 5.7 \times 10^6 + 95 \times D_C' + 95 \times (D_I/N) \quad (34)$$

$$C_C' > 1.7 \times 10^6 + 190 \times D_C' + 190 \times (D_I/N) \quad (35)$$

$$C_C' < 4 \times 10^9 / (F-1) \quad (36)$$

If the conditions $T(MN) - T(RC) = (11) > 0$, $E(MN) - E(RC) = (12) > 0$, $T(RC) - T(LC) = (7) < 0$, and $E(RC) - E(LC) = (8) < 0$, the remote cloud becomes the best place which offload the process step computation to. With the same methods, we get the following inequalities.

$$C_C' > 1.43 \times 10^8 + 95 \times D_C' + 95 \times (D_I/N) \quad (37)$$

$$C_C' > 4.37 \times 10^7 + 190 \times D_C' + 190 \times (D_I/N) \quad (38)$$

$$C_C' > 4 \times 10^9 / (F-1) \quad (39)$$

We see that to make offloading useful, the process step should have more instructions to offset the extra cost to transfer the initial data to a local cloudlet or a remote cloud. But this extra cost becomes smaller as more interactions are made with a user.

C. Initial Data at a Remote Cloud

When the initial data is stored in a remote cloud, we assume that the size of the data is extremely large and, therefore, transferring the initial data to a mobile node will demand too much time, energy, and storage space on a mobile node. So it is not a good idea to transfer the initial data to a mobile node and we only compare offloading to a local cloudlet and offloading to a remote cloud in this subsection.

For offloading to local cloudlet to be beneficial, two inequalities $T(RC) - T(LC) = (13) > 0$ and $E(RC) - E(LC) = (14) > 0$ should be satisfied. From (23) and (24) we see that these two inequalities are the same, and we need to solve only

$$C_C \times (1/S_{CR} - 1/S_{CL}) + 2N \times (T_{PR} - T_{PL}) - D_I/B_C > 0$$

Assuming $B_C = 1\text{Gbps}$ we get

$$C_C' < \{4 \times 10^9 - 84 \times (D_I/N)\} / (F-1) \quad (40)$$

This sets the upper bound on the number of instruction in a

process step so that it may be executed on a local cloudlet. We see that this upper bound grows as more interactions are made with a user and shrinks as the initial data gets larger and/or the local cloud gets slower compared with the remote cloud.

V. CONCLUSION

We consider the execution of an interactive application in mobile cloud computing environment including a local cloudlet and remote cloud. We consider three cases for an interactive application: without initial data, with a medium size initial data stored at a mobile device, and with a large size initial data stored at a remote cloud. We also assume that the program consists of two parts: one that should be executed at a mobile device and the other that can be offloaded to an outside node, either a local cloudlet or a remote cloud. We introduce models for computing elapsed time and energy consumption of an interactive application when its offloadable portion is executed on a mobile device, a local cloudlet, or a remote cloud. Applying practical numbers to some of the model parameters such as processing speed, bandwidth, propagation delay, and energy usage of a mobile device, we derive conditions under which computation offloading to either a local cloudlet or a remote cloud becomes profitable in terms of elapsed time and energy consumption.

REFERENCES

- [1] P. Mell and T. Grance, "The NIST definition of cloud computing," NIST Special Publication 800-15.
- [2] A. Li, X. Yang, S. Kandula, and M. Zhang, "CloudCmp: comparing public cloud providers," *ACM Internet Measurement Conference*, 2010.
- [3] N. Fernando, S. W. Loke, and W. Rahayu, "Mobile cloud computing: a survey," *Future Generation Computer Systems*, vol. 29, pp. 84-106, 2013.
- [4] H. T. Dinh, C. Lee, D. Niyato, and P. Wang, "A survey of mobile cloud computing: architecture, applications, and approaches," *Wireless Communications and Mobile Computing*, 2011.
- [5] L. Guan, X. Ke, M. Song, and J. Song, "A survey of research on mobile cloud computing," 10th *IEEE/ACIS International Conference on Computer and Information Science*, 2011.
- [6] M. Satyanarayana, P. Bahl, R. Caceres, and N. Davies, "The case for vm-based cloudlets in mobile computing," *IEEE Pervasive Computing*, vol. 8, no. 4, pp. 14-23, 2009.
- [7] K. Kumar and Y. -H. Lu, "Cloud computing for mobile users: can offloading computation save energy?" *IEEE Computer*, pp. 51-56, April 2010.
- [8] A. P. Miettinen and J. K. Nurminen, "Energy efficiency of mobile clients in cloud computing," *The 2nd USENIX Conference*, 2010.
- [9] B. -G. Chun, S. Ihm, P. Maniatis, M. Naik and A. Patti, "CloneCloud: elastic execution between mobile device and cloud," *ACM Symposium on Computer Systems (EuroSys)*, pp. 301-314, 2011.
- [10] E. Cuervo, A. Balasubramanian, D. -K. Cho, A. Wolman, S. Saroiu, R. Chandra, and P. Bahl, "Maui: making smartphone last longer with code offload," *The 8th International Conference on Mobile Systems, Applications, and Services (MobisSys'10)*, 2010.
- [11] N. Sorensen, Industry Benchmarks Performance, Available: <http://www.cisco.com>.
- [12] D. T. Hoang, D. Niyato, and P. Wang, "Optimal admission control policy for mobile cloud computing hotspot with cloudlet," *IEEE Wireless Communications and Networking Conference: Services, applications, and Business*, 2012.
- [13] N. Balasubramanian, A. Balasubramanian, and A. Vekaramani, "Energy consumption in mobile phones: a measurement study and implications for network applications," *The 9th ACM SIGCOMM Conference on Internet Measurement (IMC '09)*, 2009.
- [14] J. Sherry, S. Hasan, C. Scott, A. Krishnamurthy, S. Ratnasamy, and V. Sekar, "Makin middleboxes someone else's problem: network processing as a cloud service," *ACM SIGCOMM Conference*, 2012.
- [15] A. Carroll and G. Heiser, "An analysis of power consumption in a smartphone," *USENIX Annual Technical Conference*, 2000.
- [16] M. Y. Malik, "Power consumption Analysis of a Modern Smartphone," *Lecture Notes in Computer Science*, 2012.

Lossless, Multiband, on Board, Compression of Hyperspectral Images

Bruno Carpentieri, Raffaele Pizzolante

Abstract— Hyperspectral remote sensing produces a huge amount of three-dimensional digital data: the hyperspectral images. Hyperspectral images are used to recognize objects and to classify materials on the surface of the earth. They are considered a useful tool in different real-life applications. In this paper we propose a novel approach for the efficient lossless compression of hyperspectral images, which is based on a predictive coding model. Our approach relies on a three-dimensional predictive structure that uses, one or more, previous bands as references to exploit the redundancies among the third dimension. The proposed technique uses limited resources in terms of CPU and memory usage. The achieved results are comparable, and often better, with respect to the other state-of-art lossless compression techniques for hyperspectral images.

Keywords— Hyperspectral images, lossless compression, low complexity, 3-D data.

I. INTRODUCTION

THREE-dimensional data generated by hyperspectral remote sensing are collected from the visible and the near-infrared spectrum of reflected light. The human visual system, can only see visible light: the wavelengths between 360 to 760 nanometers (nm), the hyperspectral data, commonly referred as hyperspectral images, reveal also the frequencies of ultraviolet and infrared rays. Thus, a hyperspectral image is a collection of information derived from the electromagnetic spectrum of an observed area.

Figure 1 shows a graphical representation of an hyperspectral image that highlights its three-dimensional nature. The X-axis indicates the columns, the Y-axis indicates the rows and the Z-axis indicates the spectral channels of the hyperspectral image, often referred as bands.

There are many real-life applications in which hyperspectral data are used: agriculture, mineralogy, physics, surveillance, etc.. In geological applications, for example, the capabilities of hyperspectral remote sensing can be useful to identify various types of minerals, by permitting the search of minerals and oil.

Each hyperspectral sensor generates daily data in the order of many gigabytes, it is therefore necessary to compress these data so to be able to transmit and to store them efficiently. Lossless compression is generally used in order to preserve the original data, because of the high costs involved in the acquisitions and also for the importance of these data in

delicate tasks (as for instance target classification or detection).

In this paper, we propose a novel technique for the lossless compression of hyperspectral images. The proposed algorithm is based on the predictive coding model and the proposed predictive structure uses a multiband three-dimensional structure. Our technique allows to customize the encoding parameters, as for instance the number of the previous bands which will be used as references. We designed our approach to optimize the computational complexity and the memory usage, which depends on the chosen parameters.

The experimental results show that the compression results obtained by this algorithm reach, and often outperform, the performance of the other state of the art approaches, and that the algorithm maintains a good trade-off between computational complexity/memory usage and compression performances.

Our algorithm is suitable for on board implementations: it is highly configurable and it is possible to implement it with limited hardware capabilities, as on airplane or a satellite.

The paper is organized as follow: Section 2 shortly discusses previous work on lossless and lossy compression of hyperspectral images, Section 3 describes the proposed lossless compression approach, Section 4 reports the experimental results and Section 5 highlights our conclusion and future work directions.

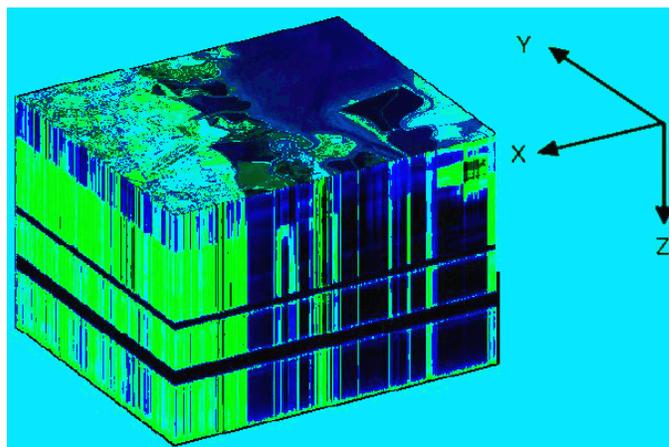


Figure 1. Graphical representation of an hyperspectral image (NASA AVIRIS Moffett Field).

B. Carpentieri and R. Pizzolante are with the Dipartimento di Informatica, Università di Salerno, I-84084 Fisciano (SA) – Italia. (phone: +39 089969500; fax: +39 089969600/1; e-mail: bc@dia.unisa.it, rpizzolante@unisa.it).

II. PREVIOUS WORK

Lossless compression of hyperspectral images is generally based on the predictive coding model. The predictive-based approaches have different advantages: they use limited resources in terms of computational power and memory and achieve good compression performances. Thus, these models are suitable for on board implementations.

Spectral-oriented Least Squares (SLSQ) [20], Linear Predictor (LP) [20], Fast Lossless (FL) [8], CALIC-3D [10], M-CALIC [10] and EMPORDA [21] are among the state-of-art predictive-based techniques.

Other approaches are designed for offline compression, since they use more sophisticated techniques and/or require to have available at once the whole hyperspectral image. These approaches are not suitable for an on board implementation but can achieve better compression performances.

Mielikainen, in [12], proposed an approach for the compression of hyperspectral image through Look-Up Table (LUT). LUT predicts each pixel by using all the pixels in the current and in the previous band, by searching the nearest neighbor, in the previous band, which has the same pixel value as the pixel located in the same spatial coordinates as the current pixel. LUT has high compression performances, but it uses more resources in terms of memory and CPU usage.

Other lossless techniques are based on dimensionality reduction through principal component transform [17].

An error-resilient lossless compression technique is proposed in [1].

For the lossy compression of hyperspectral images, the compression algorithms are, generally, based on 3D frequency transforms: as for examples 3-D Discrete Wavelet Transform (3D-DWT) [9], 3-D Discrete Cosine Transform (3D-DCT) [11], Karhunen–Loève transform (KLT) [16], etc.. These approaches are easily scalable. On the other hand, they require to maintain in memory the entire hyperspectral image at the same time. Locally optimal Partitioned Vector Quantization (LPVQ) [3, 13] applies a Partitioned Vector Quantization (PVQ) scheme independently to each pixel of the hyperspectral image. The variable sizes of the partitions are chosen adaptively and the indices are entropy coded. The codebook is included as part of the coded output.

This technique can be used also in lossless mode, but the high costs required in terms of CPU and memory do not allow an on board implementation

III. MULTIBAND COMPRESSION OF HYPERSPECTRAL IMAGES

Hyperspectral images present a strong correlation among consecutive bands (inter-band) and a high correlation in the spatial context (intra-band).

Figure 2 highlights the correlation among consecutive bands: the X-axis indicates the bands and the Y-axis indicates the Pearson's correlation [15] between the i -th band and the $(i-1)$ -th band. As it is possible to observe, the Pearson's correlation assumes high values in most of the cases.

These characteristics can be exploited by a compression

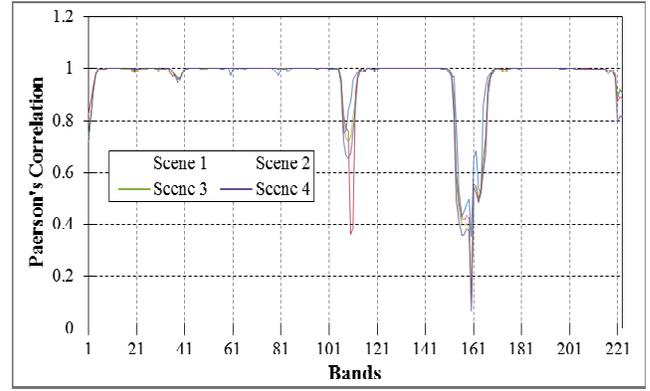


Figure 2. Pearson's correlation among consecutive bands for the fourth scenes of the Moffett Field hyperspectral image.

algorithm that optimizes the redundancy among the third dimension.

The proposed lossless compression technique: named Lossless MultiBand compression for Hyperspectral Images (LMBHI), is based on the predictive coding model.

LMBHI takes as input the hyperspectral image, and, for each pixel X of the hyperspectral image performs the prediction of the current pixel, \hat{X} , by using the appropriate prediction context of X .

Since the pixels of the first band have no reference pixels in the previous bands, they are predicted by using a bi-dimensional predictive structure: the 2-D Linearized Median Predictor (2-D LMP) [19] that uses only the neighboring pixels.

All the other pixels of all the other bands are predicted by using a new three-dimensional predictive approach, which uses for the prediction the neighboring pixels of X and its reference pixels in the previous bands.

After the prediction step, the prediction error:

$$e = [X - \hat{X}]$$

is computed, modeled, and coded.

In the following, we give more details on all the components of the algorithm.

The 2-Dimensional Linearized Median Predictor (2D-LMP) [19] uses as prediction context the three neighboring pixels of X , referred as I_A , I_B and I_C , as shown in Figure 3.

The predictive structure is derived from the well-established 2-D Median Predictor, that is used in JPEG-LS [4].

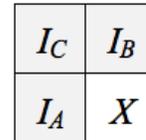


Figure 3. The prediction context of the 2D-LMP predictive structure. The gray part is already coded and the white part is not coded yet.

The 2-D Median Predictor has the following predictive structure:

$$\hat{X} = \begin{cases} \max(I_A, I_B) & \text{if } I_C \geq \min(I_A, I_B) \\ \min(I_A, I_B) & \text{if } I_C \leq \max(I_A, I_B) \\ I_A + I_B - I_C & \text{otherwise} \end{cases}$$

Median Predictor selects one of the above three options, depending on the context.

By combining all the three options, it is possible to obtain the predictive structure of 2D-LMP, defined as:

$$\hat{X} = \frac{2 \cdot (I_A + I_B) - I_C}{3}$$

Our three-dimensional Multiband Linear Predictor (3D-MBLP) uses, instead, N (up to 16) neighboring pixels (of X) for each of the B previous bands, to compute the prediction of X .

In order to define the prediction context, we need to enumerate the neighboring pixels of X in the current and in the previous bands.

For these reasons, we define an enumeration that depends on a distance d , defined as:

$$d((z, u, v), (z, w, z)) = \sqrt{(u-w)^2 + (v-z)^2}$$

When more pixels have the same indices, it is possible to reassign the indices of these pixels in clockwise order with respect to X .

Let $I_{i,j}$ denotes the i -th pixel of the j -th band, according to the above enumeration.

Let $I_{0,j}$ denotes the pixel that has the same spatial coordinates of X , of the j -th band ($j \neq k$), according to the above enumeration.

Figure 4 shows the resulting enumeration of the first $N=16$ pixels for the k -th band.

3D-MBLP is based on least squares optimizations and the prediction is computed as:

$$\hat{X} = \sum_{i=1}^B \alpha_i \cdot I_{0,k-i}$$

The coefficients:

$$\alpha_0 = [\alpha_1 \quad \dots \quad \alpha_B]$$

are chosen to minimize the energy of the prediction error

$$P = \sum_{i=1}^N (I_{i,k} - \hat{I}_{i,k})^2$$

P can be rewritten in matrix notation as: $P = (C\alpha - X)^T \cdot (C\alpha - X)$, where:

$$C = \begin{bmatrix} I_{1,k-1} & \dots & I_{1,k-B} \\ \vdots & \ddots & \vdots \\ I_{N,k-1} & \dots & I_{N,k-B} \end{bmatrix} \text{ and } X = \begin{bmatrix} I_{1,k} \\ \vdots \\ I_{N,k} \end{bmatrix}.$$

By taking the derivate of P and by setting it to zero, we obtain the optimal coefficients:

$$(1) \quad (C^T C) \alpha_0 = (C^T X).$$

		$I_{16,k}$	$I_{14,k}$		
	$I_{11,k}$	$I_{8,k}$	$I_{6,k}$	$I_{9,k}$	$I_{12,k}$
$I_{15,k}$	$I_{7,k}$	$I_{3,k}$	$I_{2,k}$	$I_{4,k}$	$I_{10,k}$
$I_{13,k}$	$I_{5,k}$	$I_{1,k}$	X		

Figure 4. The prediction context ($N=16$) for the current pixel X in the k -th band. The gray pixels have already been coded, the white pixels are not coded yet.

Once the coefficients α_0 , which solve the linear system (1), are obtained, then it is possible to compute the prediction \hat{X} of the current pixel X .

A prediction error can assume positive or negative values. In order to have only non-negative values, similarly to [14], we mapped each prediction error with an invertible mapping function M (which does not alter the redundancy among the errors). The simplified definition of the function M is:

$$M(error) = \begin{cases} 2|error| & \text{if } error \geq 0 \\ 2|error|-1 & \text{otherwise} \end{cases}$$

where $|x|$ means the absolute value of x .

Once mapped, the error is coded through arithmetic coding.

The main computational costs of our approach are due to the resolution of the linear system (1) to generate the optimal coefficients α_0 for the computation of the predicted pixel. By using the normal equation method, the linear system (1) can be solved with $(N + B/3) \cdot B^2$ floating-point operations [6].

Figure 5 shows the trend of the computational complexity of our predictive model, in terms of number of operations (Y-axis) that are required for the solving of the linear system (1), by using configurations with different parameters (X-axis).

If we use only the previous band as a reference ($B = 1$), only about 20 operations are needed to solve the system.

Instead 4 or 9 times more operations are required, if we use two previous bands ($B = 2$) or three previous bands ($B = 3$).

A linear system can have three kinds of solutions: no solutions, one solution and infinity solutions.

In the first and the third scenarios, the proposed predictive structure cannot perform the prediction.

In these cases, it is desirable to use another low-complexity predictive structure and we have used the 3-D

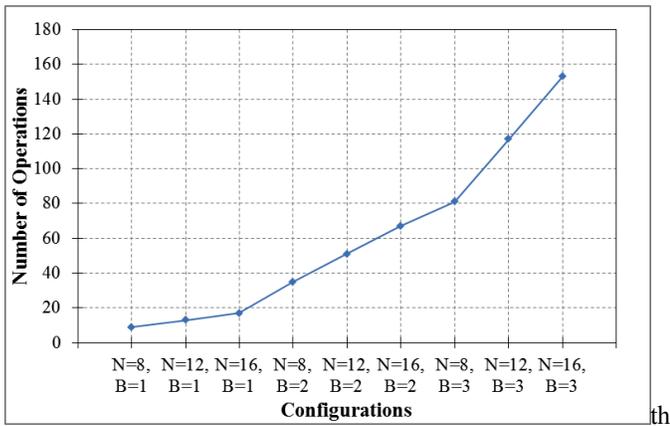


Figure 5. The number of operations (Y -axis) required to solve the linear system (1), by using different parameters (X -axis).

Distances-based Linearized Median Predictor (3D-DLMP) [19].

IV. EXPERIMENTAL RESULTS

We have experimentally tested our approach on five Airborne Visible/Infrared Imaging Spectrometer (AVIRIS) [2] hyperspectral images provided by the NASA Jet Propulsion Lab (JPL) [7], each image is subdivided into scenes.

The test set we have used is composed by the following images: Lunar Lake, Moffett Field, Jasper Ridge, Cuprite and Low Altitude, respectively of 3, 4, 6, 5 and 8 scenes. Except for the last scenes of each image that have a minor number of rows, each scene of the images has 614 columns, 512 lines and 224 spectral bands. Each sample is represented by an integer with 16 bits.

Table 1 reports the results achieved by using our approach with different parameters on all the test images. These results are reported in terms of compression ratio (C.R.) and they are compared with other state of the art lossless compression schemes.

By using two previous bands as references ($B = 2$), LMBHI outperforms, in average, all the state of the art approaches.

By using only the previous band as reference ($B = 1$), LMBHI outperforms all the state of the art techniques, with exception of LPVQ: an algorithm that is not suitable for on board implementation.

In this latter case, LMBHI achieves better results with respect to LPVQ on 3 of the 5 hyperspectral images: Moffett Field, Jasper Ridge and Low Altitude, but LPVQ gains on Cuprite and especially on Lunar Lake.

The high flexibility and adaptability of our approach makes it considerable for on board implementations. In fact, the coding parameters can be customized depending on the hardware available.

Therefore, it is possible to implement the algorithm on different typologies of sensors, by using an appropriate configuration for each one. Moreover, the proposed approach could be easily scaled for future generation sensors, which will have better hardware capabilities.

Methods / Images	Lunar Lake	Moffett Field	Jasper Ridge	Cuprite	Low Altitude	Average
LMBHI ($N=16, B=2$)	3.27	3.23	3.23	3.27	3.07	3.21
LMBHI ($N=8, B=2$)	3.21	3.18	3.18	3.21	3.02	3.16
LMBHI ($N=8, B=1$)	3.18	3.14	3.16	3.19	2.99	3.13
LPVQ	3.31	3.01	3.12	3.27	2.97	3.14
SLSQ	3.15	3.14	3.15	3.15	2.98	3.11
JPEG-2000	2.98	2.99	2.96	2.98	2.82	2.95
LP	3.05	2.88	2.94	3.03	2.76	2.93
JPEG-LS	2.87	2.90	2.87	2.87	2.74	2.85
Diff. JPEG2000	2.94	2.83	2.82	2.92	2.69	2.84
Diff. JPEG-LS	2.93	2.84	2.81	2.91	2.70	2.84
M-CALIC	3.19	3.27	3.06	3.14	N.A.	N.A.
CALIC-3D	3.06	3.08	3.09	3.25	N.A.	N.A.
LUT	3.44	3.23	3.40	3.17	N.A.	N.A.

Table 1. Compression results (C.R.) achieved by LMBHI (by using various parameter configurations), compared to other lossless compression methods.

V. CONCLUSIONS AND FUTURE WORK

In this paper we have proposed a predictive-based scheme to compress hyperspectral images, which uses a multiband three-dimensional predictive structure and that it is suitable for onboard implementations.

The results achieved are comparable and often outperform the other state of the art lossless compression techniques.

Future work will include a more intensive testing of the proposed approach, by taking also into consideration the possibility of pre-processing the hyperspectral image before compression, or by reordering the bands by considering their correlation.

This will possibly improve the compression performance [5, 13, 18].

ACKNOWLEDGMENT

The authors would like to thank our students Dario Di Nucci, Fabio Palomba, Stefano Ricchiuti and Michele Tufano for testing a preliminary version of our algorithm.

REFERENCES

- [1] A. Abrando, M. Barni, E. Magli, F. Nencini, "Error-Resilient and Low-Complexity Onboard Lossless Compression of Hyperspectral Images by Means of Distributed Source Coding", IEEE Trans. on Geosci., vol. 48, no. 4, pp. 1892-1904, April, 2010.
- [2] AVIRIS NASA Page, Available on: <http://aviris.jpl.nasa.gov/>, Accessed on Oct. 2013.
- [3] B. Carpentieri, J.A. Storer, G. Motta, F. Rizzo, "Compression of Hyperspectral Imagery", Proceedings of IEEE Data Compression Conference (DCC '03), Snowbird, UT, USA, pp. 317-324, 25-27 March 2003.
- [4] B. Carpentieri, M. Weinberger, G. Seroussi, "Lossless Compression of Continuous Tone Images", Proceeding of IEEE, vol. 88, no. 11, pp. 1797-1809, November, 2000.
- [5] B. Carpentieri, "Hyperpectral Images: Compression, Visualization and Band Ordering", Proceedings of IPCV 2011; Volume 2, pp. 1023-1029, 2011.

- [6] G.H. Golub, C.F. Van Loan, "Matrix Computations, 3rd ed. Baltimore", MD: The Johns Hopkins Univ. Press, 1996.
- [7] Jet Propulsion Laboratory (JPL) Page, Available on: <http://www.jpl.nasa.gov/>. Accessed on Oct. 2013.
- [8] M. Klimesh, "Low-complexity lossless compression of hyperspectral imagery via adaptive filtering", IPN Progress Report, vol. 42-163, pp. 1–10, 2005.
- [9] S. Lim, K. Sohn, C. Lee, "Compression for hyperspectral images using three dimensional wavelet transform," in Proc. IGARSS, Sydney, Australia, 2001, pp. 109–111.
- [10] E. Magli, G. Olmo, E. Quacchio, "Optimized onboard lossless and near-lossless compression of hyperspectral data using CALIC", Geoscience and Remote Sensing Letters, IEEE, vol.1, no.1, pp.21,25, Jan. 2004.
- [11] D. Markman, D. Malah, "Hyperspectral image coding using 3D transforms," Proc. IEEE ICIP, Thessaloniki, Greece, pp. 114–117, 2001.
- [12] J. Mielikainen, "Lossless compression of hyperspectral images using lookup tables", IEEE Signal Process. Letters, vol. 13, no. 3, pp. 157–160, Mar. 2006.
- [13] G. Motta, F. Rizzo, J.A. Storer, Hyperspectral Data Compression, Springer Science, Berlin, Germany, 2006.
- [14] G. Motta, J.A. Storer, B. Carpentieri, "Lossless Image Coding via Adaptive Linear Prediction and Classification", Proceedings of the IEEE, vol. 88, no. 11, pp. 1790–1796, November, 2000.
- [15] K. Pearson, "Mathematical contributions to the theory of evolution.-III. Regression, heredity and panmixia. Philos.", Trans. R. Soc. Lond., 187, pp. 253–318, 1896.
- [16] B. Penna, T. Tillo, E. Magli, G. Olmo, "Transform Coding Techniques for Lossy Hyperspectral Data Compression", IEEE Trans. on Geosci., vol. 45, no. 5, pp. 1408-1421, May, 2007.
- [17] M. Pickering, M. Ryan, "Efficient spatial-spectral compression of hyperspectral data", IEEE Trans. Geosci. Remote Sens., vol. 39, no. 7, pp. 1536–1539, Jul. 2001.
- [18] R. Pizzolante, B. Carpentieri, "Visualization, Band Ordering and Compression of Hyperspectral Images", Algorithms, 5(1), pp. 76-97, 2012.
- [19] R. Pizzolante, B. Carpentieri, "Lossless, low-complexity, compression of three-dimensional volumetric medical images via linear prediction", International Conference on Digital Signal Processing (DSP) 2013 18th, pp.1,6, 1-3 July 2013.
- [20] F. Rizzo, B. Carpentieri, G. Motta, J.A. Storer, "Low-complexity lossless compression of hyperspectral imagery via linear prediction", Signal Processing Letters, IEEE, vol.12, no.2, pp. 138,141, Feb. 2005.
- [21] J.E. Sánchez, E. Auge, J. Santaló, I. Blanes, J. Serra-Sagristà, A.B. Kiely, "Review and Implementation of the Emerging CCSDS Recommended Standard for Multispectral and Hyperspectral Lossless Image Coding", Proceedings of 2011 First International Conference on Data Compression, Communications and Processing (CCP), Palinuro, Italy, pp. 222–228, 21–24 June 2011.

Semantic Web Technologies and Model-Driven Approach for the Development and Configuration Management of Intelligent Web-Based Systems

Arturs Bartusevics, Andrejs Lesovskis, Leonids Novickis

Abstract- The study provides the model-driven approach for implementation of software configuration management. Provided approach uses technologies of Semantic Web to improve transformations between different models. Software configuration management is a discipline that controls software evolution process and helps to make valid builds of intelligent web-based systems. Nowadays, high level of agility requires to setup process of software configuration management as soon as possible. Provided approach helps to decrease process implementation time by reuse of existing source code for new processes.

Keywords— Software Configuration Management, Model-Driven Approach, Semantic Web.

I. INTRODUCTION

SOFTWARE CONFIGURATION MANAGEMENT (SCM) is the discipline that controls the evolution of large and complex software systems. SCM tools and systems vary and range from small tools such as RCS (Revision Control System) over medium-sized systems such as Subversion to large-scale industrial systems such as Adele ClearCase.

Nowadays software configuration management is not only challenge to choose optimal system for source code management. Complex software development projects with multiple mutually dependent components and high level of agility require two important things: firstly, in context of software configuration management, many tasks should be implemented such as source code manage managements, version control, build and deploy management, accounting of statuses of items, etc.; secondly, the mentioned implementation should be ready as soon as possible because agile methodologies require frequent releases of new versions of product.

Some of the most common problems in the area of software configuration management are the following:

- Use of multiple different configuration management tasks in a single solution. For example, use of a script that performs source code management, build management, and installation

management tasks. Such multifunctionality makes this script specific for one particular project and makes it impossible to reuse it without some modifications.

- Lack of approaches and recommendations on how to design reuse-oriented solutions for configuration management that could be used in the other projects without additional customizations to save up time and resources.

Reusable configuration management solutions should be parameterized and structured by different tasks. It means that, for example, solutions on how to build the product from the source code should be independent from the other tasks like source code management or installation management. The mentioned product build solution should receive a set of parameters and return an executable or an error message. It should not contain any details or hardcoded information like the location of the source code or the address of the server where the executable should be installed.

The paper describes a new model-based approach for implementation of software configuration management. Unlike the other approaches, it is not oriented to any particular tool or script that “should solve any problem” but describes the steps how to increase the reuse of the existing solutions. This approach is independent from the tools being used for tasks like source code management, continuous integration, bug tracking, build management, etc. as it only defines a way to make a solution reusable.

Authors also investigate how the Semantic Web technologies like OWL and SPARQL could be used to improve this approach and to perform transformations between different levels of models.

The paper is structured as follows: the second section describes the work done by other researchers in the fields of Software Configuration Management and Semantic Web technologies. The third section covers the use of the Semantic Web technologies (like RDF, OWL, etc.) in the model-driven Software Configuration Management and potential advantages they could bring. The fourth section contains the detailed description of the proposed model-driven EAF methodology. Fifth section is concerned with how the Semantic Web

technologies could be used to improve EAF methodology. The last section is related to the results of the experimental assessment of the EAF methodology in 5 different software projects.

II. RELATED WORKS

Long term expert in software build management Tracy Ragan in her paper [1] writes that solutions for software configuration management should be model-driven, not static script-oriented as it was normal in 20th century. That paper outlined some of the reasons why model-driven approaches could be better for modern software configuration management processes.

Firstly, a lot of software products are supported by cloud computing technologies where the things like server names, absolute paths, and other information required for static scripts and specific platforms are unknown. In this case, model-driven approach could provide a way from planning to implementation of software configuration management in virtual environments [1]. There are many tools can configure builds and deployments for complex products within a few hours without writing huge and complex static scripts for particular platforms [2].

Secondly, model-driven approach helps to reduce the human factor during selecting solutions for particular task of software configuration management. Usually traditional implementation of software configuration management contains two major steps: process planning and development of static scripts for the planned process. There are risks related to writing source code for scripts, because sometimes executable source code is not in consistence with initial planning. Model-driven approach could reduce these risks by generating source code for configuration management automatically [1][2][4].

Unfortunately, modern tools mentioned in [2] are oriented only for one task of software configuration management: build and deployment management. But it is not possible to prepare a valid build with any tool without correct source code management [3][4][5]. Additionally, tools that mentioned in paper [2] help to improve builds and deployments, but can not fix successful use cases to reuse them in new projects. These tools require the acquisition of additional knowledge and financial investments. However, sometimes enterprises already have tools for software configuration management that they trust and rely on. Because of that these companies are looking for approaches and methodologies that could increase reuse of the existing solutions.

There are some novel approaches related to implementation of software configuration management using models. The study [3] provides semantic integration of different tools. Ontologies are used to create concept of each tool using in software configuration management process.

The paper [4] presents software configuration management model based on ITIL framework. The model is theoretical and no any suggestions about increasing reuse of existing solutions are given. The main advantage of approach [4] is based on

well-known framework that works in real world, but lack of suggestions about tools and approaches that could be used for implementation makes provided approach perspective in theory but not trusted in practice.

Approach described in work [5] provides a method for selecting optimal software configuration management tool by using artificial intelligence methods. The method seems very useful from the theoretical point of view. Empirical evidence provided in [5] shows that method could really improve source code management and version control in particular software development project. As a main disadvantage is a fact, that theory of fuzzy logic is relative difficult for understanding.

In context of software configuration management, very important task is a management of source code and version control. The study [6] describes the approach related to improvement of version control. Nowadays, majority of version control tools supports management of source code. However, there is a lack of approaches that provide version control for model-driven software development. Approach presented in work [6] describes a model of universal version control system that could be used for code and model management.

The current paper provides approach for implementation of software configuration management process. Provided approach is an improved version of the researches described in works [7][8]. Practical assessment of the first version of the approach outlined some disadvantages. In the improved version of approach described in current paper, all useful ideas will be taken from studies [3][4][5][6] and conclusions from the practical experiments of approaches [7][8].

Unlike other approaches, methodology described in this paper:

- Does not impose to use any specific tools for software configuration management process,
- Increase reuse of existing solutions for software configuration management,
- Defines full cycle of models related to step-by-step implementation of software configuration management using existing solutions,
- Transformations between models are improved with semantic web technologies.
- Has an abstract views that allows to design new implementation of provided approach.

This study is not the first attempt to introduce the Semantic Web technologies into software configuration management. In a related study, Falbo [13] proposed an SCM ontology that was used to establish a common conceptualization about the SCM domain in order to support SCM tools integration.

III. SEMANTIC WEB TECHNOLOGIES AND MODEL-DRIVEN SOFTWARE CONFIGURATION MANAGEMENT

The Semantic Web is an extension of the current Web in which information is given well-defined meaning, better enabling computers and people to work in cooperation [10]. It is based on the idea of having data on the Web defined and

linked in a way that it can be used by machines not just for display purposes, but for automation, integration, and reuse of data across various applications.

The Semantic Web is composed of a set of technologies, and it can be defined as a symbiosis of Web technologies and knowledge representation. Semantic Web technologies can be used in a variety of application areas; for example: in data integration, whereby data in various locations and various formats can be integrated in one, seamless application; in cataloging for describing the content and content relationships available at a particular Web site, page, or digital library; by intelligent software agents to facilitate knowledge sharing and exchange; in content rating; in describing collections of pages that represent a single logical “document”; for describing intellectual property rights of Web pages, and in many others [11].

Authors think that the Semantic Web technologies can be efficiently utilized in the software configuration management to ease and improve efficiency of processes like data integration and reuse, transformation, and searching.

Ontologies serve as a key enabling technology for the semantic software configuration management. Ontologies are developed to provide a machine-processable semantics of information sources that can be communicated between different agents (software and humans). Ontology is an explicit formal specification of a shared conceptualization. 'Conceptualization' refers to an abstract model of some phenomenon in the world which identifies the relevant concepts of that phenomenon. 'Explicit' means that the type of concepts used and the constraints on their use are explicitly defined. 'Formal' refers to the fact that the ontology should be machine readable. Hereby different degrees of formality are possible [12].

The type of knowledge used to describe SCM field is very hard to represent formally using traditional approaches. Organizations use different proprietary solutions that make it much harder to reuse encoded software configuration information across different software projects. The lack of a standard for compatible encoding of configuration data is one of the major SCM problems. Web Ontology Language (OWL), a family of knowledge representation languages for ontology authoring and one of the key Semantic Web building blocks, provides developers with features and opportunities that can be used to solve this problem.

OWL is an ontology language designed for use in the Semantic Web and is the language recommended by the W3C for this use. OWL DL and OWL Lite semantics are based on Description Logic (DL). OWL 2 exhibits the desirable features of Description Logics, including useful expressive power, formal syntax and semantics, decidability, and practical reasoning systems, resulting in OWL 2 providing effective ontology representation facilities.

Ontologies provide software developers with a standard and powerful way of representing knowledge not only the configuration management tasks but about software

engineering project in general.

Besides the general benefits of formal specification that can be gained by encoding the software configuration knowledge, the main advantage of using OWL ontologies is an ability to provide a fully automated procedure to detect inconsistencies in the configuration via standard reasoning engines (for example, Pellet or RacerPro). The reasoning service can not only determine, if a certain configuration is valid or not, but it also provides justifications for such decision (i.e. it lists the reasoning steps that lead to this decision).

One of the key benefits of the use of the Semantic Web technologies is that they provide means to reason and query over semantically annotated metadata from the software configuration models. Reasoning provides an opportunity to perform an inference.

Inference is a process to infer a new relationship from the existing resources and some addition information in form of set of rules. Inference base technique is also used to check data inconsistency at time of data integration. The inference engine can be described as a form of finite state machine with a cycle consisting of three action states: match rules, select rules, and execute rules [12]. The use of DL reasoners allows OWL ontology applications to answer complex queries and to provide guarantees about the correctness of the result. This is obviously of crucial importance when ontologies are used in safety critical applications.

Some of the features that can be provided by a standard Description Logic reasoner are the following [9]:

- Consistency checking – ensures that an ontology does not contain any contradictory facts;
- Concept satisfiability - determines whether it is possible for a class to have any instances;
- Classification - computes the subclass relations between every named class to create the complete class hierarchy;
- Realization - computes the direct types for each of the individuals.

The reasoning is especially important when developers are dealing with different versions of software application to support various machines and/or operating systems and in the scenarios where the constantly changing requirements are different for different target groups. For example, it is possible to detect that individual doesn't belong to a certain class or doesn't satisfy constraints and/or restrictions of available configurations.

Another important aspect of using OWL ontologies is that every document or resource can be uniquely identified and referenced using Universal Resource Identifier (URI). This makes tasks like configuration management planning much easier

IV. EAF METHODOLOGY FOR IMPLEMENTATION OF SOFTWARE CONFIGURATION MANAGEMENT PROCESS

Methodology provided in this section is an improved version of works [7][8] and is related to decreasing

implementation time of software configuration management by reuse of existing solutions. The main principles of novel methodology are the following:

- Ready implementation of software configuration management process is a source code that could be executable only from particular software configuration management server. It could be one of well-known continuous integration servers like Bamboo, Jenkins, CruiseControl etc.
- The main result of EAF methodology is generated source code for software configuration management process.
- Generation of the source code for software configuration management is automated process that uses models and existing executable units of source code.

The name of provided methodology (Environment -> Action -> Framework) provides the main steps for generating source code for software configuration management:

- *Defines all environments in particular software development project.* Environments in context of EAF methodologies is a set of servers, applications, virtual machines needed to use software. For example, web application needs database server and application server to make this application ready to use. Usually, the scope of particular environment is particular process in software development project. For example, DEV environment for development, TEST for testing etc. During the first step of EAF methodology, all environment in project should be defined and all

transfers of changes between these environments. The mentioned step is formalized by Environment Model, described in [7].

- *Defines all actions needed to apply Environment Model from the previous step.* In context of this work, actions are tasks of general software configuration management process. For example, to move software changes from DEV environment to TEST environment, the following actions could be defined:
 - PREPARE_BASELINE: merge changes of source code from development to test branch;
 - BUILD: run builds scripts or tools to make executable from prepared source code;
 - INSTALL: deploy ready executables to applications servers.

In the previous versions of provided methodology [7][8], there are two models: Environment Model and Platform independent Action Model. Improved EAF methodology contains merged variant of mentioned models, called PIEM (Platform Independent Environment Model). This model also contains all environments and flows of changes between them, but additionally, configuration manager could define actions needed to apply transfers of changes. Example of PIEM model is given in Fig.1.

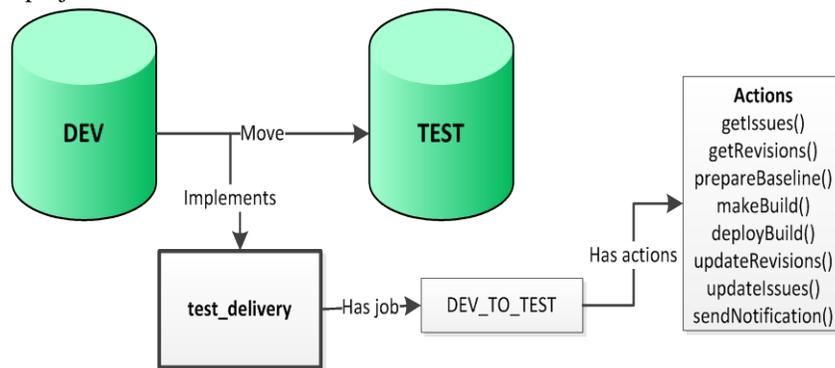


Fig. 1 Example of PIEM model

- Choose the Framework for each action defined in PIEM model. The framework in context of this paper is a set of executable units of source code for implementation of particular action form PIEM model. During this step, configuration manager chooses implementation for each PIEM action

from Solutions Database. Solution Database is a structure there all solutions for configuration management in particular enterprise are stored. The structure of Solution Database is improved version that have been provided in [7][8]. The structure of Solution Database provided in Fig. 2.

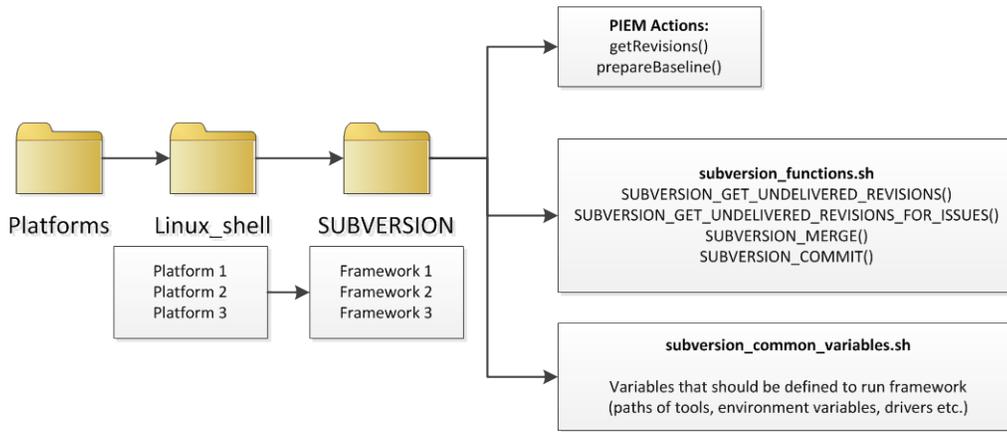


Fig. 2 Solution Database

Solution database provided in Fig. 2 illustrates principles of framework with Subversion example. PIEM model shows actions of configuration management but after actions are defined, configuration manager have to choose a framework for each action. For example, for action “prepareBaseline()” Subversion framework could be selected. Configuration manager should choose platform for implementation and framework. After framework for particular action is defined,

Solution Database provides notes about variables that should be defined to activate framework. Each framework has a set of functions that could be called from other scripts and tools. So, during implementation of software configuration management in particular project, only project specific parts should be developed. Significant part of the source code could be taken from framework. Fig. 3 contains an overview of all steps and principles of provided EAF methodology.

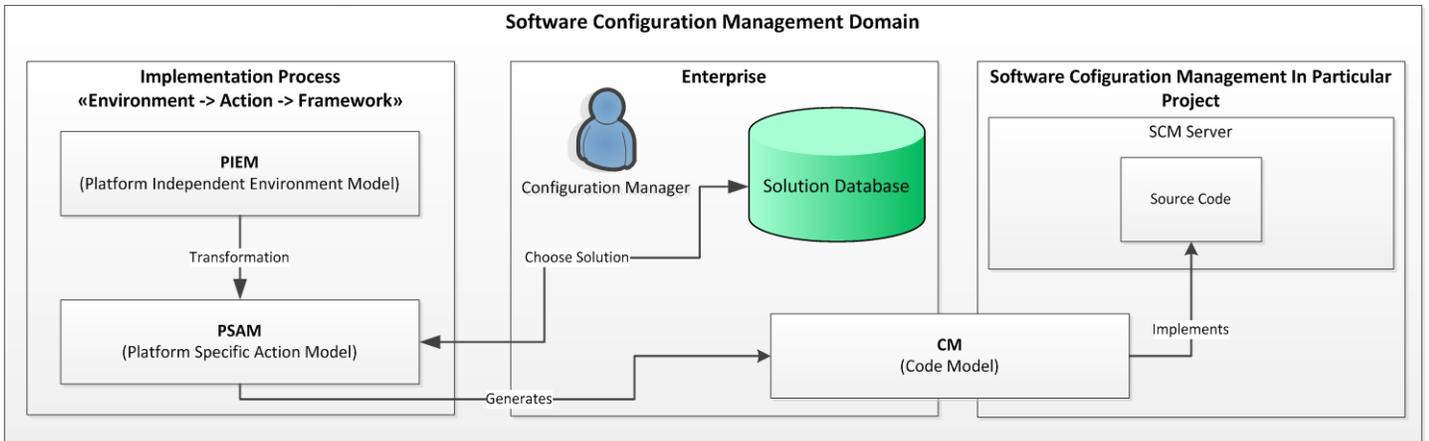


Fig. 3 EAF methodology in enterprise

Fig. 3 contains all main steps of EAF methodology. During the implementation process configuration manager makes environment model (PIEM), defines actions of software configuration management. Ready PIEM model should be fulfilled with information about frameworks. Configuration manager chooses frameworks from Solution Database for each actions. Platform Specific Action Model contains implementation details. This model could be transformed to Code Model and last one could be implemented in software configuration management server to support software configuration management process for particular project.

from the use of the Semantic Web technologies not only in terms of improved efficiency but could also potentially provide additional functionality.

In their research [14], Arantes et al came to the conclusion that ontology from [13] lacked some important concepts mostly related to change and version control. Therefore, they presented what they called an evolution of this ontology that introduced a few new concepts and a taxonomy of change control actions. The new version features concepts like Repository, Branch, Version, Artifact, Change, etc.

Authors believe that Arantes et al’ SCM ontology could be used in the proposed EAF methodology as a base ontology. That would allow the reuse of the valuable expert knowledge encapsulated within this ontology. However, it is necessary to correspondingly modify it according to the needs of the EAF methodology. One of the main reasons for these changes is that the ontology was not designed with a good reasoning support.

V. APPLYING SEMANTIC WEB TECHNOLOGIES TO THE EAF METHODOLOGY

Authors think that EAF methodology could greatly benefit

The Semantic Web Rule Language is a language for the Semantic Web that can be used to express rules as well as logic, combining OWL DL or OWL Lite with a subset of the Rule Markup Language [15]. SWRL complements DL by providing the ability to infer additional information in DL ontologies, but at the expense of decidability. SWRL rules are Horn clause-like rules written in terms of DL concepts, properties, and individuals. SWRL includes a high-level abstract syntax for Horn-like rules in both the OWL DL and OWL Lite sublanguages of OWL.

An SWRL rule is composed of an antecedent (body) part and a consequent (head) part, both of which consist of positive conjunctions of atoms. The SWRL rule syntax is the following:

$$\text{antecedent} \Rightarrow \text{consequent}$$

where both antecedent and consequent are conjunctions of atoms written $a_1 \wedge \dots \wedge a_n$. For example, we can say that if a is a parent of b and b is a parent of c , then a is also a parent of c using the following rule:

$$\text{parent}(?a, ?b) \wedge \text{parent}(?b, ?c) \Rightarrow \text{parent}(?a, ?c)$$

SWRL atom can be either a class, an object property, a data type, a data type property, or a built-in. A rule is satisfied by an interpretation if every binding that satisfies the antecedent also satisfies the consequent.

SWRL has already been successfully used in the quite related field of Network Access Control Configuration Management [16]. This experience suggests that SWRL rules can also be used to implement ontology-based SCM model transformation rules (for example, transformation of PSAM to Code Model).

Given the rich SCM ontology, SWRL provides developers with an opportunity to define the resilient rules for different kinds of model transformation scenarios all while inferring possible new knowledge.

VI. PRACTICAL ASSESSMENT OF THE EAF METHODOLOGY

For the practical assessment, software configuration management process had been implemented at 5 different software development projects by provided EAF methodology. The projects are the following:

- Project 1: Maintenance of Oracle E-Business Suite

system. Development technologies are Oracle Forms, Oracle ADF, PL/SQL, JAVA. Version control system is Subversion, bug tracking system is JIRA and continuous integration server is Bamboo.

- Project 2: Implementation of Oracle Customer Care and Billing Utilities system. Development technologies are Oracle CC&B, PL/SQL, Java, Cobol. In this project also Subversion and JIRA are used, but continuous integration server is Jenkins.
- Project 3: Development of Web-Based ERP (Enterprise Resource Planning) system based on Ruby On Rails platform. In this project Git system have been used for version control. For bug tracking and continuous integration also JIRA and Jenkins have been used.
- Project 4: Development of custom billing system based on Oracle Application Development Framework. This project have been used the same set of tools as Project 1.
- Project 5: Development of Web-services for ERP system using Oracle SOA Suite 11g platform. This project also have been used Subversion as version control system, JIRA for bug tracking and Jenkins for continuous integration.

To underline benefits of EAF methodology, implementation time of software configuration management is fixed and compared by implementation time by old methods (without EAF). The difference between old and new implementation time provided at percent. To save up daily processes of mentioned projects, experimental implementation of configuration management by EAF are completed at new (parallel) software configuration management servers. The Fig. 4 provides overview of difference between implementation time of software configuration management by old methods and by EAF methodology. Difference is provided in percent, for example value '-10%' mean that implementation time by EAF methodology of later for 10%, but value '+15%' means that implementation by EAF methodology is not useful because it needs for 15% more time.

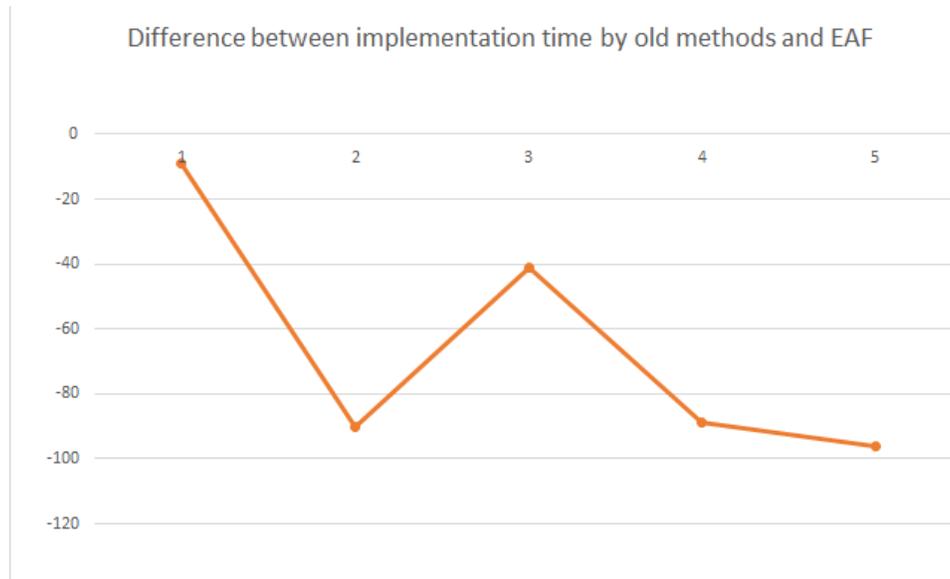


Fig. 4 Difference of implementation time by old methods and EAF

Fig. 4 provides that the first experiment at “Project 1”, while Solution database is empty, benefits of EAF methodology in context of implementation time is only 9%. However, implementation of software configuration management by EAF at Project 2 is for 90% later because Solution Database already has a set of reusable solutions for similar set of tools (Subversion, JIRA). Implementation time at Project 3 is only for 41% later, because Solution Database does not contain framework for Git version control system. The benefits at projects 4 and 5 are great (89% and 96%), because at this moment Solution Database contains all reusable Frameworks for Subversion, Git, JIRA, Jenkins, Hudson. During implementation of software configuration management, mostly existing functions from frameworks are used. In this case, only a small part of source code for particular software configuration management process should be developed.

The main trend provided at Fig. 4 shows that benefits from EAF is relatively small while Solution Database is empty, but if mentioned database contains framework for all most used tools at particular enterprise, implementation time of software configuration management could be for 80% - 95% later.

VII. CONCLUSIONS

The study provides new model-driven approach for implementation of software configuration management. The main scope is to increase reuse of existing solutions and reduce efforts to implement the process in other projects. General picture and principles of new EAF approach are provided, Platform Independent Environment Model and Solution Database with example are introduced. Finally, experiments of implementation of software configuration management by EAF methodology are provided.

Results of this work were used in Latvian Council of Science project "Approach and Generic Methodology for Development of Applied Intelligent Software Based on

Artificial Intelligence, Modelling, and Web Technologies" (project leader prof. J. Grundspenkis) and in European Commission 7th Framework project "eINTERASIA "ICT Transfer Concept for Adaptation, Dissemination and Local Exploitation of European Research Results in Central Asia's Countries"" (project coordinator prof. L. Novickis).

In order to continue research, it is necessary to carry out the following activities:

- Develop additional criteria that evaluate models benefits in software development projects not only from point of implementation time,
- Based on developed criteria, evaluate benefits of designed models,
- Develop criteria to assess whether the developed model-driven approach for configuration management implementation corresponds to guidelines of ISO/IEC 15504, ITIL, CMMI standards.
- Design Code Models and transformation algorithms for other platforms.
- Add and improve tools and frameworks in existing platforms.

The approach provided in this article is abstract and only general stages, kinds of models and basic elements are defined. The authors hope that the new approach will generate new ideas because many useful lessons could be learned from different implementations of this model-driven approach.

ACKNOWLEDGMENT

The research has been partly supported by the project eINTERASIA "ICT Transfer Concept for Adaptation, Dissemination and Local Exploitation of European Research Results in Central Asia's Countries", grant agreement No. 600680 of Seventh Framework Program Theme ICT-9.10.3: International Partnership Building and Support to Dialogues for Specific International Cooperation Actions - CP-SICA-

INFSO.

REFERENCES

- [1] Ragan, T., 21st-Century DevOps--an End to the 20th-Century Practice of Writing Static Build and Deploy Scripts, *Linux Journal*, 230, pp. 116-120, Computers & Applied Sciences Complete, EBSCOhost, viewed 22 October 2014.
- [2] Azoff, R., DevOps: Advances in Release Management and Automation. [ONLINE] Available at: http://electric-cloud.com/wp-content/uploads/2014/06/EC-IAR_Ovum-DevOps.pdf, 2014.
- [3] Calhau R., Falbo R., A Configuration Management Task Ontology for Semantic Integration. Proceedings of the 27th Annual ACM Symposium on Applied Computing Pages 348-353 ACM New York, NY, USA, 2012.
- [4] Giese H., Seibel A., Vogel T., A Model-Driven Configuration Management System for Advanced IT Service Management. Available at: http://www.hpi.unipotsdam.de/giese/gforge/publications/pdf/GSV-MRT09_paper_7.pdf, 2009.
- [5] Yongchang, R., Fuzzy Decision Analysis of the Software Configuration Management Tools Selection. In ISCA 2010. France, 19-23 June, 2010. Information Science and Engineering (ISISE): ACM. 295 - 297., 2010.
- [6] de Almeida Monte-Mor, J., GALO: A Semantic Method for Software Configuration Management. In *Information Technology: New Generations (ITNG)*, 2014. USA, 7-9 April, 2014. ITNG: IOT360. 33 - 39., 2014.
- [7] Novickis, L., Bartusevics, A. Model-Driven Software Configuration Management and Environment Model. In: *Recent Advances in Electrical and Electronic Engineering: Proceedings of the 3rd International Conference on Systems, Communications, Computers and Applications (CSCCA'14)*, Italy, Florence, 22-24 November, 2014. Florence: WSEAS Press, 2014, pp.132-140. ISBN 978-960-474-399-5. ISSN 1790-5117.
- [8] Novickis, L., Bartusevičs, A., Lesovskis, A. Model-Driven Software Configuration Management and Semantic Web in Applied Software Development. Proceedings of the 13th International Conference on Telecommunications and Informatics (TELE-INFO '14), Istanbul, Turkey December 15-17, 2014.
- [9] Clark & Parsia, LLC, Pellet Features. Available at: <http://clarkparsia.com/pellet/features>, 2015.
- [10] Berners-Lee, T., Hendler, J., and Lassila, O. (2001). "The Semantic Web", *Scientific American*, May 2001, p. 29-37. Available at : http://www-sop.inria.fr/acacia/cours/essi2006/Scientific%20American_%20Feature%20Article_%20The%20Semantic%20Web_%20May%202001.pdf
- [11] W3C. Semantic Web Frequently Asked Questions. Available from <http://www.w3.org/2001/sw/SW-FAQ#swgoals>, 2009.
- [12] Fensel D. *Ontologies: Silver Bullet for Knowledge Management and Electronic Commerce*. Springer, 2003, 162 p.
- [13] Falbo, R., A., Calhau, R. F. A Configuration Management Task Ontology for Semantic Integration. In: Proceedings of the 27th Annual ACM Symposium on Applied Computing. ACM, New Yorok, 2012. Pages 348-353.
- [14] Arantes, L., D., Falbo, R. D., Guizzardi G. Evolving a Software Configuration Management Ontology. [ONLINE] Available at: <http://citeseerx.ist.psu.edu/viewdoc/download;jsessionid=C71AC33F802C1644AB292AFD9268ED9F?doi=10.1.1.95.9969&rep=rep1&type=pdf>
- [15] Horrocks, I. et al. SWRL: A Semantic Web Rule Language Combining OWL and RuleML [ONLINE] Available at: <http://www.w3.org/Submission/SWRL/>
- [16] Fitzgerald, William M. and Foley, S. N. and Ó Foghlu, M. (2009) Network Access Control Configuration Management using Semantic Web Techniques. *Journal of Research and Practice in Information Technology*, 41 (2). pp. 99-117.

Arturs Bartusevics currently is a Doctoral Student at Riga Technical University, the Faculty of Computer Science and Information Technology, the Institute of Applied Computer Systems. He obtained BSc (2008) and MSc (2011) degrees in Computer Science and Information Technology,

respectively, from Riga Technical University. His research areas are software configuration management, release building and management process and its optimization. He works at Ltd. Tieto Latvia as a Software Configuration Manager.

E-mail: arturik16@inbox.lv

Andrejs Lesovskis is a Doctoral Student at Riga Technical University, the Faculty of Computer Science and Information Technology. He obtained MSc degree in Computer Science and Information Technology at Riga Technical University in 2009. His research areas are e-Learning and Semantic Web. He works as a researcher at Riga Technical University.

E-mail: andrejl@inbox.lv

Leonids Novickis is a Head of the Division of Software Engineering at Riga Technical University. He obtained Dr.sc.ing. degree in 1980 and Dr.habil.sc.ing. degree in 1990 from the Latvian Academy of Sciences. He is the author of 180 publications. Since 1994, he is regularly involved in different EU-funded projects: AMCAI (INCO COPERNICUS, 1995-1997) – WP leader; DAMAC-HP (INCO2, 1998-2000), BALTPORTS-IT (FP5, 2001-2003), eLOGMAR-M (FP6, 2004-2006) – scientific coordinator; IST4Balt (FP6, 2004-2007), UNITE (FP6, 2006-2008) and BONITA (INTERREG, 2008-2012) – RTU coordinator; LOGIS, LOGIS-Mobile and SocSimNet (Leonardo da Vinci) – partner, eINTERASIA (FP7, 2013-2015)- project coordinator. He was an independent expert of IST and Research for SMEs in FP6 and FP7. He is a corresponding member of the Latvian Academy of Sciences and an elected expert of the Latvian Council of Science. His research fields include Web-based applied software system development, business process modeling, e-learning and e-logistics.

E-mail: lnovickis@gmail.com

Clock Distribution Using a Bi-dimensional Orthogonal Salphasic Structure

Andrei Pasca

Abstract—Salphasic (contraction based on the words phase and saltation) networks assume the creation of a special standing wave pattern along the propagation direction. The specific requirement for the salphasic distribution network is the presence of a voltage anti-node at the generator side, allowing for reduced, virtually no load seen by the driving circuit. An important characteristic of the salphasic pattern (or any other standing wave pattern) is the existence of large same-phase regions where only the signal amplitude varies as a function of the position in the network. Sharp phase transitions (180 degrees steps – from which the attribute saltation is derived) between different regions are created at the nodes of the standing wave pattern. The advantage of standing wave based networks, and in particular of salphasic networks, is the possibility to distribute the clock signal over the entire 2D area of the application with minimal effort. The present article details a theoretical model and then shows a practical implementation of an orthogonal, yet irregularly shaped, bi-dimensional clock distribution network. The model is based on process-specific parameters of selected CMOS technology node.

Keywords—Bi-dimensional clock distribution, salphasic surface, standing wave.

I. INTRODUCTION

ACCORDING to Moore's law, system performance of Integrated Circuits (ICs) or Systems on Chips (SoCs) keeps growing steadily, although clock frequencies remained somewhere around 3.5 – 4GHz. At first sight, from the clock distribution network's perspective, this relatively fixed operating frequency should pose no problem. However, because of the use of massively parallel architectures needed to support the required system performance, the number of elementary loads for the distribution network is constantly growing.

For synchronous devices, this growth comes with more stringent phase alignment requirements, as more and more nodes must have a strictly defined relation between active clock edges and valid signal levels. Standard approaches like those found, for instance, in [1], [2], [3], all need an important fraction from the overall system's power budget. Taking into account the value of the operating clock frequency coupled with the physical sizes of the silicon die, of course, other approaches are also possible, like in [4], [5].

For instance, [1] assumes that skew is an intrinsic parameter of any given clock distribution network and therefore it should be integrated into the design. By introducing additional

programmable skew in critical propagation paths, it will be possible to get an increase in the maximum operating frequency. The programmable skew is calibrated in a post-fabrication step and can be either positive or negative (a technique present also in [3] by the name of cycle stealing).

In [2], the skew is again minimized with a post fabrication step and active programmable delay lines. The delay lines are used to re-align the active clock edges between adjacent clock domains (driven by different clock distribution structures), assuring thus minimal skew between different parts on the same silicon die. The re-alignment structure was so effective, maybe even over-designed, that it concealed a bug in the first silicon. Since the delay lines compensated a supplemental inverter inserted in one of the propagation paths, it can be said that the length of the delay line was more than needed, and hence, more than needed silicon area was dedicated to the clock distribution network.

These two previous approaches used both a post-fabrication adjustment, trying to increase the system performance at the expense of system complexity – i.e. software and hardware layers are added to the SoC. The approach in [3], in fact a classic approach, performs the skew analysis in the final stages of the design process, before the actual tape-out. Symmetric or optimally-partitioned clock trees (with X, H or binary structures) are used to achieve the skew target (set to minimal achievable skew). This comes as a tradeoff between system / design complexity and proper signal alignment / maximum operating frequency.

The solutions depicted in [4], [5] are taking advantage from the range in which the operating frequencies are situated. For instance, considering a moderately loaded transmission line, built on a silicon die, such that the wave speed is about a third of the speed of light, a clock signal of 4.5GHz will have a wavelength of about 22mm. Comparing this with current chip sizes for SoCs in the range of 6 to 10mm, it becomes evident that the clock distribution networks become compatible with techniques that are used in the microwave domain.

Both [4] and [5] are using standing wave based resonant networks. In essence, the clock network acts as a distributed oscillator that uses transmission line segments as LC tanks. The interconnections between the individual oscillator nodes and the sharing of the transmission line segments between these oscillators allow a full synchronous operation of the network over the entire silicon area. In this way, constant-phase clock distribution becomes possible using only minimal resources. However, since the equivalent LC parameters

Andrei Pasca is PhD student in the Electrical and Telecommunications Department at Politehnica University, Timisoara, Romania, (e-mail: and_pasca@yahoo.com).

emulated by the transmission line segments depend on the physical dimensions on the silicon die, the distributed frequency can be placed only in a limited range in which the standing wave pattern can be maintained (given by the electrical lengths of the transmission line segments). A possible extension for the frequency range, as used in [5], consists in using additional inductors (used to artificially increase the length of the transmission lines). In other words, the flexibility comes again with a downside of increased area allocated for the clock distribution network.

At first sight in a similar approach with that of [4] or [5], the approach from [6], [7], [8] and [9] allows the extension of the operating clock frequency range. Basically, the salphasic distribution, as named in the references, assumes the creation of a standing wave pattern in such a way as to present minimal loading for the generator node – i.e. for a voltage driving buffer, this translates in the presence of an anti-node configuration at the input of the distribution network. One important advantage deriving from this configuration is the relaxed range in which the output impedance of the clock drivers can be situated – i.e. no matching between the characteristic impedance of the network and the output impedance is required, the restriction being the length of the transient response.

Unlike the solution from [4] or [5], the salphasic distribution network is not self-oscillating – without an external clock source, no IC operation is possible. Furthermore, as seen in fig. 1, the implementations of [4] and [5] are creating a clock signal only along specific sections. From these points, a classic clock distribution network must be used to spread the signal to all needing logic elements. In other words, the standing wave oscillators are used only to replace the root level and a few branches of a standard H or X distribution network, leaving the bulk of the clock capacitance still driven by standard buffers (for this reason, the present

paper will compare the salphasic network with the last stage of a classic X-tree distribution network).

Bi-dimensional salphasic distribution networks are capable of delivering the clock signal to all nodes without any additional circuitry. The bulk of the clock distribution network capacitance is an integral part of the network and because of this, as it will be shown later, it can be properly charged and discharged with much weaker drivers than normally used in clock distributions.

Regarding the maximum operating frequencies, for the salphasic approach as presented in [6] - [9], the limiting factors are the frequency range in which the needed loads can be adjusted such as to maintain the anti-node at the input or, at least in [8] and [9], where the concept is extended at IC level by compensating the resistive losses of the transmission lines, the range in which the losses can be compensated.

With respect to the first limitations, there are circuit configurations where the transmission line loads are external to the silicon die, so they can be adjusted independently of the post-fabrication parameters of the IC or SoC. Furthermore, since the anti-nodes of a standing wave pattern are rather flat, small load matching errors can be easily tolerated by the clock drivers, resulting in a further extension of the frequency range.

Although [6] and [7] presented practical structures, no full theoretical models were provided for the standing wave configuration so it was not readily apparent what are the disadvantages. Partial solutions were provided in [8], however, still no full analysis was offered. As shown in [9], the salphasic configuration of previous references proved to be only marginally interesting by having Bessel functions of the first and second kind as design equations.

Although interesting and with clear advantages, the salphasic configuration introduced in [9] can be restricted to only some applications – for instance, uniform arrays that require the same clock signal over the entire chip area, like fully systolic, very high frequency, signal processing circuits. For applications that have multiple clock domains or a segmented floorplan with blocks of irregular size, the approach introduced in [9] can be difficult to implement. Additionally, the method used for adjusting the parameters of the salphasic surface increases the system complexity.

The present paper further expands the theoretical model of the salphasic surfaces for the case of orthogonal structures. Like in [9], the model is based on process-specific electrical parameters that can be easily extracted or estimated directly from the technology’s design manual (with good-enough accuracy). Also, field solvers (from inexpensive 2D solvers to full 3D, state of the art solvers) can be used to construct the lumped model of the salphasic structure.

II. ORTHOGONAL SALPHASIC SURFACE MODEL

Assuming two infinite parallel plates as in fig. 2, ideally conductive, separated by a good dielectric material, aligned such that their origins are situated on the same O_z vertical axis and an infinitely long, uni-dimensional signal generator,

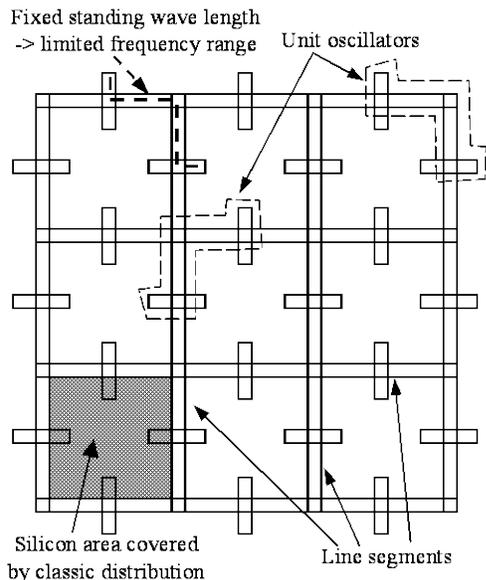


Fig. 1. Standing-wave oscillating clock distribution

connected between the two planes' Oy axes, it can be shown that, if existing, any signal propagation should assume a transversal electromagnetic (TEM) mode. From the orthogonal salphasic surface model's point of view, a TEM propagation mode is essential as it allow a lumped-circuit approach for the system.

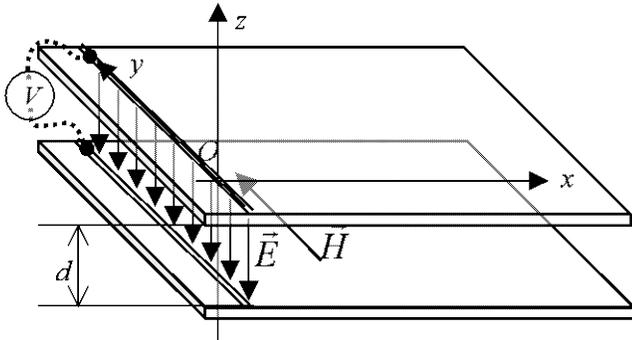


Fig. 2. Infinite parallel plates with extended generator overlapped on Oy axis

If there is a temporal variation for the potential difference created between the two infinite and equipotential generator nodes, the electric field will follow the same variation. This temporal change will generate a displacement current between the two parallel plates. Assuming an arbitrary infinite closed loop, circling the generator nodes and applying the Ampere-Maxwell law for this loop, it can be shown that a variable magnetic field will form around the initial electric field. Because the generator nodes are infinitely long, the electric field that is created between them will be uniform. Considering also that the arbitrary loop on which the magnetic field was calculated is infinite, the magnetic field, at a given distance from the generator, will also have a uniform configuration.

Next, by further selecting arbitrary closed contours contained in planes that are parallel with the xOz Cartesian plane, the previously computed magnetic field lines will produce a variable electric field along this contour (using Maxwell-Faraday law). It can be noted that there will be a variable electric field component that will be parallel with the vertical Oz axis. Continuing this line of reasoning, it can be shown that the initial variable electric field generated along the Oy axes will produce a propagating configuration of electric and magnetic fields along the Ox direction. The direction of the magnetic field will be parallel with the main axis of the generator nodes, while the direction of the electric field will be contained in arbitrary (i.e. all possible) planes that are perpendicular to the extended generator nodes. That is, the electric and magnetic field lines will be orthogonal, suggesting a TEM configuration. The only condition needed for the previous assumption to hold is that the electromagnetic field must be contained in the dielectric material that separates the parallel plates. This can be easily asserted, both for the electrostatic and electrodynamic cases. For the electrostatic case, it's easy to show that it's impossible to have electric field above or below the considered system.

For the electrodynamic case, it can be shown that, if any

magnetic and electric fields external to the dielectric material are present, then the initial electric field established between the extended, infinite generator nodes, must be non-uniform, which is impossible.

With these, it is clear that the propagation would assume a TEM mode, as both the electric and magnetic fields are always perpendicular to the propagation direction and fully contained between the two parallel plates. Furthermore, as the wave travels along a direction parallel to the Ox axis, it will be a plane wave. Even more, as the propagation mode is TEM and the electromagnetic wave is a plane wave, the model for the two infinite conductive parallel plates must be similar with the model of simple transmission lines.

As a final note here, although the system is in fact a three-dimensional one, because the wave propagates only about two of the spatial dimensions, as was done for the simple transmission line, it's possible to reduce the number of Cartesian coordinates for the considered structure to only two spatial dimensions.

A. Lumped-Circuit Model Parameters

The four general model parameters (that is resistance R and conductance G – being the loss inducing terms; capacitance C and inductance L) can be extracted starting from an elementary section of the initial system as seen in fig. 3. The general requirement for the elementary section is to be an orthogonal polygon of arbitrarily small dimensions \hat{dx} and \hat{dy} , either a rectangle or a square. However, since sheet resistance defined in the selected fabrication process' manual is usually given as per-square resistance, the present paper assumes the unit element to be also square.

Additionally, since the structure is symmetric with respect to up / down orientation, the model assumes that all of the conductive material's resistance is given only by the upper plate, the lower plate being an ideal conductor. In other words, the elementary section consists of two parallel plate, one of thickness g , and an ideal one, with no resistive losses and very thin, separated by the dielectric material having a depth of d and a conductance of σ_D .

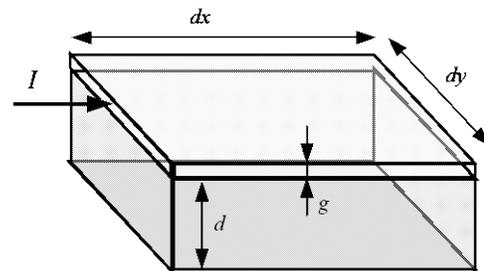


Fig. 3. Elementary model section

The first three stated parameters can be directly extracted based on the area and electrical properties of the conductive and dielectric materials that are used in the silicon fabrication process. Resistance is given directly in the design manual as the sheet resistance of the conductive layer. The conductance

can be derived as the reciprocal of the resistance of the dielectric material enclosed between the two conductive surfaces of the elementary section (this conductance may also be inferred from the leakage currents or from the dielectric losses that are specific to the selected technology node).

To a first approximation, the elementary capacitance of the structure can be computed assuming the parallel plate capacitor model. This approximation is valid for sections that are not close to the edges of the real system. However, as will be detailed further on in the present paper, with careful design for the salphasic structure, it's possible to get same elementary capacitance, independent of the actual position in the clock distribution network. For this reason, the rest of this paper will use the ideal parallel plate capacitor model.

The above parameters can be computed with equations (1) to (3).

$$R = R_S \tag{1}$$

$$G = \frac{1}{d \cdot \sigma_D} \partial x \cdot \partial y = G_S \cdot \partial x \cdot \partial y \tag{2}$$

$$C = \frac{\epsilon_r \cdot \epsilon_0}{d} \cdot \partial x \cdot \partial y = C_S \cdot \partial x \cdot \partial y \tag{3}$$

The *S* subscript in the above equations stands for specific. That is R_S stands for the specific resistance, given in datasheet as the per-square resistance, G_S stands for the specific conductance and, similarly, C_S for specific capacitance. It can be noted that the specific resistance is independent of the size of the arbitrary small elementary section and that both the capacitance and the conductance are showing a similar behavior with respect to the size of the unit section.

The only model parameter that creates some difficulties and for which a field solver would yield better results, is the elementary inductance. For a theoretical analysis, it is possible to assume a further decomposition of the elementary structure into small parallelepipedic loops as depicted in fig. 4. In each loop, a current is assumed having the orientation as depicted in the figure.

It can be noted that with this decomposition, there is no net current circulation between the two parallel plates as at the boundary between any two adjacent loops, one loop will have the current going from the bottom plate to the upper one, while the other loop will have a same magnitude current, but circulating from the upper plate to the bottom one.

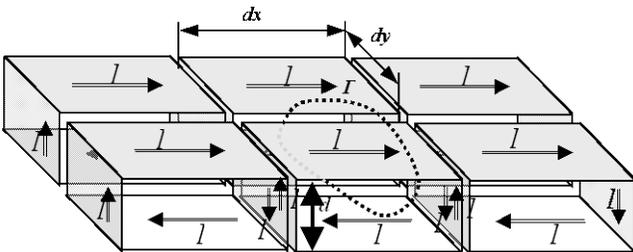


Fig. 4. Elementary section decomposition for inductance calculation

Using Ampere's law for the arbitrary Γ loop of fig. 4, it is possible to get a relation between the current and the magnetic field. After this correlation is made, the inductance can be computed based on the geometry and the magnetic flux. The results are expressed by the series of equations (4), (5) and (6).

$$H \cdot \partial y = I \tag{4}$$

$$\Phi_M = \mu \cdot H \cdot d \cdot \partial x = L \cdot I \tag{5}$$

$$L = \mu \frac{d \cdot \partial x}{\partial y} = L_S \tag{6}$$

In (5), Φ_M stands for the magnetic flux that goes through the elementary loop area. As seen in (6), because the elementary section is considered to be square, also the elementary inductance reduces to the per-square specific inductance, resulting, as expected, to a similar behavior as for the elementary resistance.

It is interesting to note that, if the elementary size ∂y that is not along the propagation direction is absorbed in the specific parameters of the lumped-circuit model, the electrical model parameters will have a similar behavior as those specific to the simple transmission line model – i.e. they could be considered as lineic parameters. This suggests a strong correlation of the surface model to the simple transmission line, that is, the orthogonal surface excited by a generator along one of its edges will be a generalization of the one-dimensional case.

B. The Simple Orthogonal Surface Model

Moving from the elementary section to the elementary model cell, based on (1) to (3) and (6), it is possible to construct the *RLCG* structure of the lumped-circuit model. The next step in constructing the simple orthogonal surface model is to replicate this cell along the Cartesian axes as seen in fig. 5.

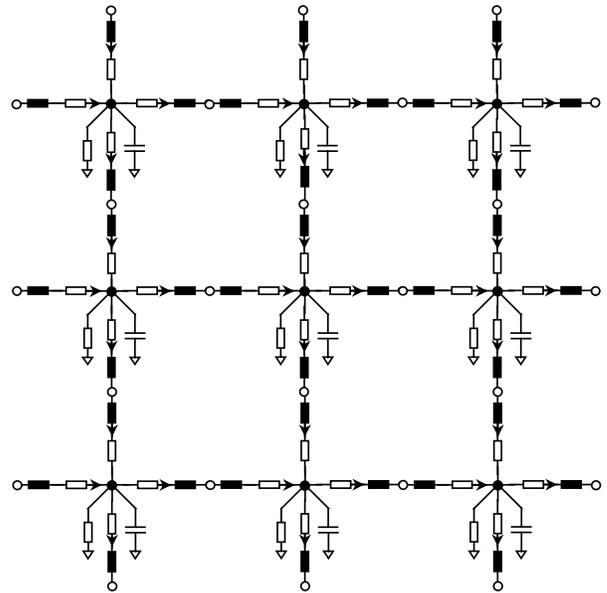


Fig. 5. Sketch of an orthogonal surface model using lumped-circuit elements

Writing the Kirchhoff's laws for current and voltage and taking into consideration the small physical sizes (i.e. some reasonable approximations can be made) it's possible to derive the differential equations for the model.

As seen in (7) and (8), there is a strong resemblance between the equations of the orthogonal surface and those for the transmission line case.

$$u(x, y) = -\frac{I}{(G_S + j\omega \cdot C_S) \cdot \partial y} \cdot \frac{\partial}{\partial x} i(x, y) \quad (7)$$

$$i(x, y) = -\frac{\partial y}{R_S + j\omega \cdot L_S} \cdot \frac{\partial}{\partial x} u(x, y) \quad (8)$$

The equations can be made independent one from the other by further differentiating them by the x spatial variable. In the end, the orthogonal surface model has the following second order differential equations attached to it.

$$\frac{\partial^2}{\partial x^2} u(x, y) - \gamma^2 \cdot u(x, y) = 0 \quad (9)$$

$$\frac{\partial^2}{\partial x^2} i(x, y) - \gamma^2 \cdot i(x, y) = 0 \quad (10)$$

$$\gamma = \sqrt{(G_S + j\omega \cdot C_S) \cdot (R_S + j\omega \cdot L_S)} \quad (11)$$

It's evident from (9) and (10) that the model allows for a variation only along the propagation direction, parallel with the Ox axis. Furthermore, (11) shows a strong similarity to the propagation constant γ of the transmission line case. Even more, given a specific technology (i.e. dielectric material, dielectric thickness, conductive material resistance), the propagation constants for both the one-dimensional (transmission line) and the bi-dimensional case (orthogonal surface) will have the exact same numeric value.

Using complex notation for the voltage and current waves, the solutions for the differential equations are given in (12) for the voltage wave and in (13) for the current wave.

$$u(x, y) = U_0^D \cdot e^{-\gamma \cdot x} + U_0^R \cdot e^{\gamma \cdot x} \quad (12)$$

$$i(x, y) = I_0^D \cdot e^{-\gamma \cdot x} + I_0^R \cdot e^{\gamma \cdot x} \quad (13)$$

It can be seen they are independent of the spatial variable y and are a linear combination of two waves, one propagating from the generator extended node, denoted with the superscript D (from direct), and another one, traveling towards the generator, denoted with superscript R (from reverse). The subscript 0 denotes the amplitudes of the signals from the generator nodes.

It is interesting also to compute the characteristic impedance (14) of the structure and to compare it with the impedance of a simple transmission line.

$$Z_C = \frac{I}{\partial y} \cdot \sqrt{\frac{R_S + j\omega \cdot L_S}{G_S + j\omega \cdot C_S}} \quad (14)$$

In (14), the spatial dimension ∂y can have any arbitrary value – in fact, for the real life implementations, it will be equal with the length of the edge that has the extended generator of the orthogonal surface. From this, it is evident that the surface can be assimilated to a collection of paralleled transmission lines, its resulting characteristic impedance being the equivalent impedance of the parallel grouping.

In the case the surface has very little or no losses at all, the propagation constant becomes purely imaginary. For the IC case, this could be possible by using loss compensation as presented in [8] or [10]. If the load connected opposite to the generator's edge creates a total reflection condition, then the magnitudes of the direct and reflected wave will be equal, resulting in a standing wave pattern established across the orthogonal surface. If the phase of the reflected wave is chosen such that the first voltage anti-node is produced at the generator's terminals, the surface will act as a salphasic distribution network. It should be noted, however, that there is nothing preventing these surfaces to be used also as general waveguides, not only as salphasic structures. The only needed condition is to be able to create extended generator and load nodes. The following section will provide a solution for this.

III. PRACTICAL SALPHASIC ORTHOGONAL SURFACES

Before proceeding to any configuration it should be noted that all depictions are assuming also the presence of an upper and a lower reference (i.e. equivalent ground) layer. These reference layers are implicit and mandatory to the salphasic structures and without them there will be no wave propagation in the structure. However, for ease of presentation, all the drawings of the present paper will not show the reference layers.

Any practical implementation for the salphasic orthogonal structure has to use extended generator and load structures. These can be approximated to a very good accuracy using Huygens' principle, as depicted in fig. 6.

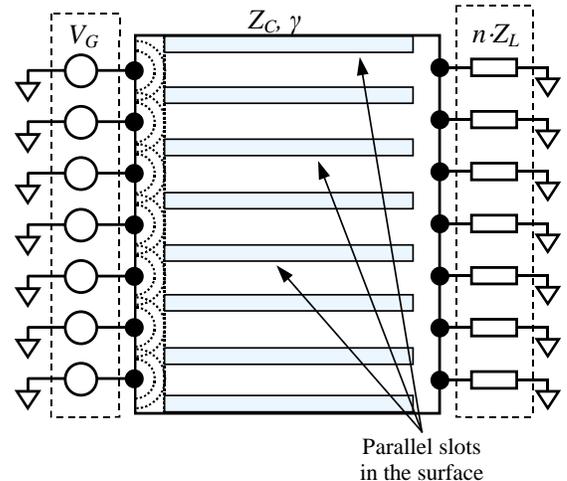


Fig. 6. Practical regular orthogonal salphasic surface

For the generator node, this amounts to several point-like drivers, equally distributed on the input-side edge of the surface, all operating in phase. For the load node, the Huygens approximation translates to equal-value, uniformly distributed loads, connected on the opposite side of the structure.

There are a variety of ways in which same-phase generators can be constructed in a SoC, for instance, a simple, small, binary tree distribution can be used to carry the clock signal from a unique source to all the salphasic surface drivers.

Unlike for standard clock distribution networks where the strength of the drivers is critical (i.e. it must be correlated to the total load capacitance seen by the network) a salphasic network driver can be made quite small. This comes from the fact that the salphasic principle creates a minimal loading condition for the input node. However, an important requirement for the distributed driver is to maintain a very good alignment between individual points, independent of the process, operating conditions and (usually only marginally considered for clock distribution networks but critical here) matching parameters. Random variations in the process parameters can induce offsets in the decision levels of the clock drivers – these offsets, if not minimized by design, will be translated in timing differences between the individual point-sources of the equivalent extended generator. For this reason, the matching parameters are restricting the minimum allowable area for the final stage drivers and hence, they are imposing also the minimum drive strength.

A further improvement for the orthogonal salphasic structure, as depicted in fig. 6, is to have slots cut out from the conductive surface along the propagation direction. These slots are allowed as normally, there should be no current circulation in a direction parallel to the extended generator’s terminals. The main role of these slots is to uniform the surface’s parameters in all points. Indeed, because of the edge effects, the parameters for an uncut surface will have different values at the borders than in the center of the surface. Because of this, the wave propagation assumes a quasi-TEM mode at the boundaries and, in the center, a propagation mode that is closer to TEM. By providing these slots, the field lines will show similar fringe effects all over the clock distribution area, resulting in effectively same lumped-circuit model parameters and a uniform, quasi-TEM propagation over the entire surface. Additionally, the slots allow for other signal routing from the upper metallization layers to the active silicon logic area (i.e. the routing of the power supply lines), without interfering with the clock distribution network.

Another possibility to create an extended generator node, as seen in fig. 7, is to use a salphasic transmission line in which equally-spaced taps are placed on the line. From these taps, clock buffers are used to provide same-amplitude, same-phase signals for the orthogonal structure. This is needed as the salphasic transmission line has the specific standing wave amplitude envelope pattern, according to a cosine function.

A critical aspect here is the avoidance of amplitude to phase conversion. For instance, if the buffers used to drive the clock

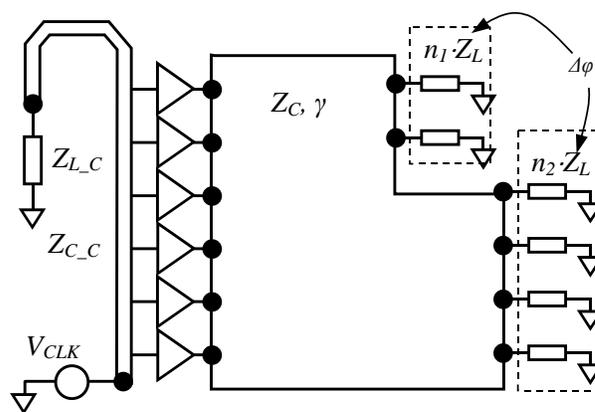


Fig. 7. Practical non-rectangular orthogonal salphasic surface

distribution surface are simple CMOS inverters, in the typical operating situation, the threshold level of all buffers is placed in the middle of the supply domain. Nominally, the salphasic clock distributed by the transmission line is also centered on this same level, so all buffer outputs will exhibit the same transition instants.

Departing from the typical case, if, for instance, the clock signal from the transmission line is centered with a slight offset in the middle of the supply domain, the buffers will convert the cosine amplitude variations of the clock signal in unwanted phase variations for the output signal, destroying the Huygens equivalence – i.e. the extended generator node will not coincide with the generated wavefront.

The easiest way to avoid the conversion from amplitude to phase is to use differential drivers instead of simple CMOS buffers. This greatly reduces the impact of common mode voltage drifts on the equivalent extended generator’s timing, due to process, voltage and operating conditions. Additionally, it also reduces the influences of across-chip variations on the generated wavefront (active components with separations between them in the order of millimeters have different mean values, even if they are on the same silicon die). In this way, the major influence comes only from local random variations of the active devices’ parameters, which can be easily accounted by properly sizing the devices (the larger the area, the smaller the influence)

Fig. 7 depicts a situation that is quite frequent in SoCs, as many times the clock domains have irregular shapes. Here, the structure still behaves as a salphasic one, with the length difference between different paths being corrected by properly selecting the load impedances. In fact in fig. 7, all load impedances must provide a total reflection condition, which in an ideal case, with lossless transmission structures, translates in purely reactive loads. By properly selecting the reactive load value and type (capacitive or inductive) it’s possible to introduce different phase shifts in the reflected wave, making it possible to maintain a proper standing wave configuration with a salphasic behavior. The above explanation holds also for loss-compensated structures as presented in [8] and [10] with the observation that the load impedances will no longer be purely reactive.

The drawback for this non-rectangular configuration is that it can be used only in a salphasic (i.e. total reflection at the load) or matched-impedance (i.e. no reflections at the load) configuration. For the structure of fig. 6 this restriction does not hold, allowing virtually all circuit configurations that are used for the simple transmission line case, with the advantage that signals will be distributed over the entire area of the IC.

The orthogonal bi-dimensional salphasic distribution can also be constructed in a differential configuration as seen in fig. 8. Here, two slotted surfaces are stacked on consecutive metallization layers such that the slots of the upper surface correspond to conductive segments in the lower surface. By controlling the widths of the slots it is possible to adjust the specific capacitance and inductance parameters of the structure. Furthermore, by having a differential salphasic structure, all needed clock buffers (not represented in fig. 8) can all have a differential configuration, with the benefits presented earlier.

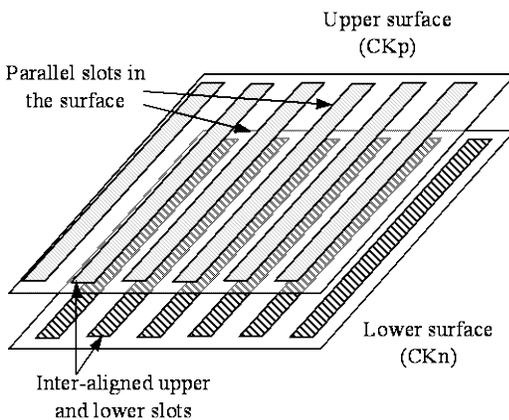


Fig. 8. Differential orthogonal bi-dimensional salphasic structure

IV. SIMULATION RESULTS

The simulated structure was a combination of those of fig. 6 and 7, that is, a slotted surface is used to distribute the clock signal over an irregular shape (but still orthogonal). The active edge of the surface has 10mm while the lengths along the propagation axis are 15mm and 10mm, respectively. As noted before the introduction of the practical salphasic structures, the system also has a lower and upper reference layer. For simulation purposes, no losses were considered in the following simulations.

The chosen technology node was IBM 130nm CMOS process. The distance between two consecutive layers in this processing node is $0.6\mu\text{m}$ (for the upper metallization levels). Considering a relative dielectric constant for the silicon dioxide of $\epsilon_r = 4$, the intrinsic capacitance for the structure results as $118\mu\text{F}/\text{m}^2$. The specific capacitance must also include the clock-driven transistor gates. The total clock load for the simulated structure was 19.5nF , evenly distributed across the entire test surface. This value amounts to about 3.3 millions clock transistors ($4.5\mu\text{m}$ by $0.12\mu\text{m}$). With this, the specific capacitance value for the structure results in

$274.2\mu\text{F}/\text{m}^2$. The specific inductance is determined from (6) to be 0.377pH (this is the per-square value). The result includes the contribution of both upper and lower reference layers.

Given all the above values for the specific parameters and considering also the 10mm width of the salphasic structure, the characteristic impedance results in $3.7\text{m}\Omega$.

The distributed generator consists of 45 inverters, each having typically 22Ω output impedance, evenly distributed along the 10mm active edge. The W/L ratios for the transistors in the unit buffers are $45\mu\text{m}/0.12\mu\text{m}$ for the nMOS and $150\mu\text{m}/0.12\mu\text{m}$ for the pMOS. The equivalent generator impedance is $489\text{m}\Omega$.

The 15mm length of the surface corresponds to the wavelength of a 6.56GHz signal and the shorter, 10mm section, corresponds to the wavelength of a 9.84GHz signal. For simulation purposes, the signal distribution from the clock injection point to the individual clock buffers was not modeled.

The clock distribution network was tested at 3.28GHz, 2.46GHz and 1.23GHz, corresponding to a standing wave pattern of half wavelength, $3/8$ and, respectively, to $3/16$ of the wavelength considered for the longer section. For lower power consumption, the system was operated only at 1V, resulting also in a clock signal amplitude at the driving edge of just 1V. Because of the limited bandwidths of the CMOS buffers, the shape of the signal can be assimilated to a sine wave.

For the first test frequency, since the length of the longer section is equal with half of the wavelength, proper salphasic behavior needs an infinite load. The needed impedance for the shorter section has a capacitive value. The second test frequency needs a capacitive load on the longer surface section and a short-circuit for the 10mm length section. For the last test frequency, all loads had inductive behavior.

For each test frequency, after the structure reached the steady-state condition, .MEASURE directives were used to determine the maximum amplitude of the distributed signal in equally spaced nodes, placed at 1mm intervals. Based on these measurements, it was possible to reconstruct the envelope of the distributed signal across the entire surface. Thus, fig. 9 depicts the envelope of the distributed clock signal on the non-regular shape structure at the 2.46GHz test frequency. As expected, the configuration produced an anti-node at the generator side. Because the test frequency of 2.46GHz is exactly a quarter of the resonant frequency of the shorter section, if an anti-node is formed at the generator side, the load side must have a voltage node in the standing wave configuration. As expected, fig. 9 shows the correct behavior of the simulated structure.

Fig. 10 presents the configuration of the standing wave pattern at all tested frequencies, along the propagation axis of the longer section of the surface. The figure was constructed by choosing measured nodes sharing the same coordinate along the wave propagation direction. Since the length of the structure is longer than a quarter of a wavelength at 2.46GHz and 3.28GHz frequencies, the pattern should also reach 0V.

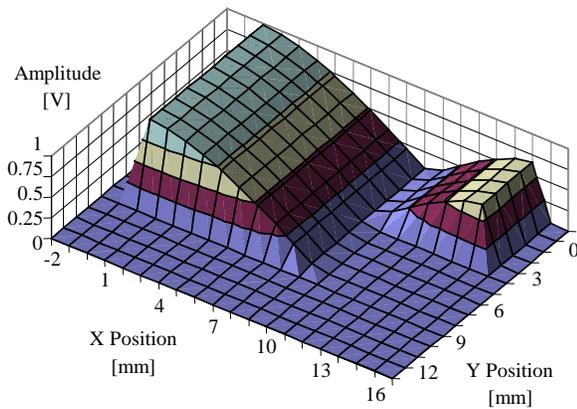


Fig. 9. Envelope of salphasic signal on non-rectangular

Fig. 10 does not show a clear 0V level, however, that is only an artifact, as there were no testprobes placed on the exact required positions of the distribution network, being limited by the finite size of the lumped elements.

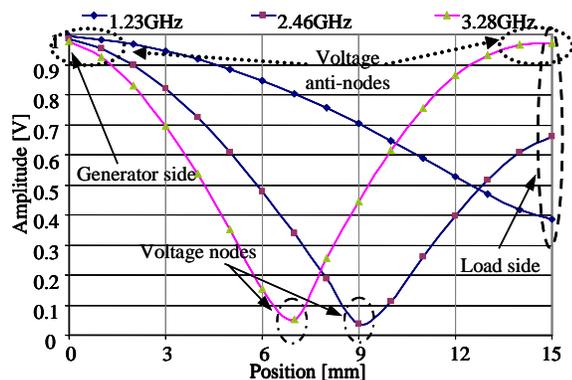


Fig. 10. Amplitude envelope against different operating frequencies

It is interesting to note that for a classic clock-driving network, operating on the same given capacitive load, the needed output impedance at 3.28GHz is $2.48\text{m}\Omega$. For this approach, the total load capacitance was split in 3900 unit capacitors, each having 5pF. Each unit capacitor was driven by a unit clock buffer (W/L ratios of 166/0.12 for nMOS and 555/0.12 for pMOS). No attempt was made to simulate the rest of the distribution network as the impact on the total power is only in the 15 to 30% range.

Given a mean power consumption of 14mW per unit buffer, the total needed power to drive the entire 19.5nF load results in 54.6W. This value can be compared to the theoretical approximation given by the product of the load capacitance, the working frequency and the square of the supply voltage. For the given operating conditions, the theoretical power results in 64W, with the difference coming from the assumptions made about the slopes of the clock signal.

For the simulated salphasic structure, the 45 buffers consumed a total of 79.5mW. However, this figure does not include the power needed for the loss compensation circuit. By evenly distributing compensation cells similar with that of [10] on a grid of 1mm by 1mm (125 cells in total), the power

consumption of the compensation circuit is found to be 1.01W. With these results, the salphasic distribution needs a total of 1.1W.

V. CONCLUSIONS

The present article introduced an expansion of the theoretical model of [9] and a generalization for the simple, uni-dimensional transmission line. By using an extended generator and an extended load, it was shown that it is possible to construct a two-dimensional waveguide, capable to carry signals like any other transmission line, with a planar wave propagation configuration.

Based on the model, a practical circuit implementation was simulated, validating the theoretical results. It was shown that the concept could be extended even for structures that have non-regular contours, as long as proper phase alignment is maintained by careful selection of the load impedance.

Although, the resulting characteristic impedance of the structure seems quite small, the output impedance of the extended generator can be much larger, resulting in each buffer having an output impedance in the range of tens of ohms. With this, the structure lands itself easily to any standard CMOS process implementation.

The total power consumption (including the power of the loss compensation cells) was only 1.1W. This result is explained by the no-load condition seen at the input, specific for the salphasic distribution.

REFERENCES

- [1] E. Takahashi, Y. Kasai, M. Murakawa, T. Higuchi, "Post-Fabrication Clock-Timing Adjustment Using Genetic Algorithms", IEEE JSSC, vol. 39, no. 4, pp. 643-650, April 2004.
- [2] T. Fischer, J. Desai, B. Doyle, S. Naffziger, B. Patella, "A 90-nm Variable Frequency Clock System for a Power-Managed Itanium Architecture Processor", IEEE JSCC, vol. 41, no. 1, pp. 218-228, January 2006.
- [3] E. G. Friedman "Clock Distribution Networks in Synchronous Digital Integrated Circuits", Proceedings of the IEEE, vol. 89, no. 5, pp. 665-692, May 2001.
- [4] A. J. Drake, K. J. Nowka, T. Y. Nguyen, J. L. Burns, R. B. Brown, "Resonant Clocking Using Distributed Parasitic Capacitance", IEEE Journal of Solid-State Circuits, vol. 39, no. 9, pp. 1520-1528, September 2004.
- [5] M. Shiozaki, M. Sasaki, A. Mori, A. Iwata, H. Ikeda, "20GHz uniform-phase uniform-amplitude standing-wave clock distribution", IEICE Electronics Express, vol. 3, no. 2, pp. 11-16, 25th of January 2006.
- [6] V. L. Chi, "Salphasic Distribution of Clock Signals for Synchronous Systems", IEEE Transactions on Computers, vol. 43, no. 5, pp. 597-602, 1994.
- [7] V. L. Chi, "Salphasic Distribution of Timing Signals for the Synchronization of Physically Separated Entities", US5387885, 7 February, 1995.
- [8] A. Pașca, "Probleme specifice ce apar în rețelele de distribuție de clock", MSc thesis, Universitatea "Politehnica", Timișoara, România, 2006.
- [9] A. Pașca, "Bi-dimensional Radially-Salphasic (Standing Wave) Clock Distribution", 2014 IEEE 20th International Symposium for Design and Technology in Electronic Packaging (SIITME), pp. 157-162, Bucharest, 23-26 October, 2014.
- [10] A. Pașca, "Negative Impedance Converter Circuits for Integrated Clock Transmission Lines Loss Compensation", Buletinul științific al Universității "Politehnica" din Timișoara, seria Electronică și Telecomunicații, Tom 54(68), Fascicula 1, pp. 19-24, 2009.

Synchronous Differential Logic Gate for Low Clock Swing Operation with Standing Wave Clock Distribution Networks

Andrei Paşca

Abstract—Standing wave based clock distribution networks – like salphasic (contraction from phase saltation) distribution – are capable of delivering a same-phase clock signal over the entire silicon area. However, the amplitude of the distributed signal follows a standing wave pattern across the clock domain, with logic gates from different parts of the system having different clock amplitudes. The present paper proposes and characterizes a differential logic gate capable to operate with these reduced clock signal amplitudes without any speed penalty or any impact on the synchronous system requirements. Since the standing wave techniques are applicable at high clock frequencies, the focus for the designed logic gate is on speed. However, starting from the model, it is possible to derive also other lower power gates by re-writing the transistor sizing equations for the used circuit topology. Although the proposed logic gate is intended to be used with a salphasic clock distribution network, there is nothing preventing its use also in other low clock swing (system-wise) applications.

Keywords—Low clock swing, standing wave, synchronous differential logic.

I. INTRODUCTION

IN recent years, as the operating clock frequencies reached the GHz range, coupled with a push for low power operation of the integrated circuits (ICs), distributing the clock signal – the highest frequency signal and having the highest capacitive load in the system – through resonant or standing wave techniques became quite a standard approach [1], [2], [3], [4].

Additionally, techniques that were used initially at system level, stemming from the microwave domain – salphasic distribution [5] – also became applicable at the IC level, like seen in [6], with a focus on distributing the signal over the entire silicon bi-dimensional surface.

All these clock distribution networks are capable of easily distributing a same-phase signal over large silicon areas. However, due to the characteristics of the standing wave patterns, the amplitude of the clock signal will vary across the clock domain, with a more pronounced effect for [4] and [6]. However, it is important to note that this amplitude variation happens about a fixed common mode voltage. As it will be seen later, this fixed common mode voltage largely restricts

the possibilities for using typical low clock swing logic gates.

This varying amplitude versus the physical position within the system, if not taken into account in a proper way, can, potentially, completely negate the advantage coming from the same-phase characteristic of the distributed clock.

Although it is possible to use a combination between the standing wave network at a global level and a classic buffered clock distribution tree at a local level, the focus of this paper will be on a logic circuit that is capable of working directly with the reduced swing signal. This approach would allow an efficient use of the bi-dimensional character of the salphasic distribution. Furthermore, the intention in the current paper is to use a differential circuit topology, for high-speed operation with reduced or predictable power supply noise. Although the present paper makes references to different types of logic circuits, they will be assimilated into the generic term of logic gate since its main focus will be on an embedded-logic dynamic latch.

As seen in fig. 1, assuming a fixed value for the logic threshold, the decreasing amplitude results in a dependence between the signal level and the effective sampling moment. It is evident that for a proper skew minimization (or even elimination), the logic threshold should either be very low, either should be adjusted for the decrease in the clock signal's amplitude.

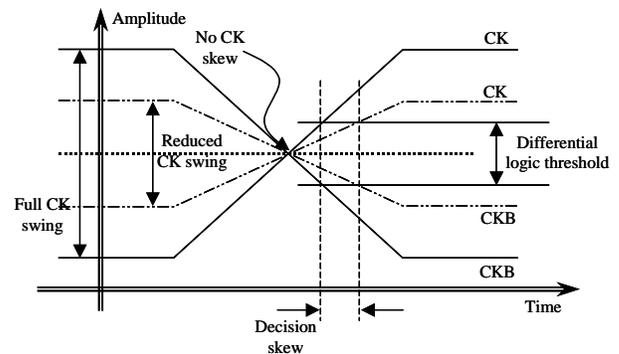


Fig. 1. Clock amplitude to skew conversion

Depending on the clock driver's architecture, the first approach may not be readily applicable - for instance, in case the system uses a global differential clock distribution network, with the signal swing centered around half the supply

Andrei Paşca is PhD student in the Electrical and Telecommunications Department at Politehnica University, Timisoara, Romania, (e-mail: and_pasca@yahoo.com).

voltage like drawn in fig. 1, the resulting clock transistor sizes will be unreasonably large.

The second approach may seem similar with those found in reduced swing clock distribution networks. For this, at least in those areas where the clock signal has a lower level, the logic synchronization elements (latches or flip-flops with or without embedded logic) should be capable of operating with a fraction of the nominal amplitude.

II. EXISTING REDUCED SWING CAPABLE CIRCUITS

In [7], the low clock signal amplitude operation is reached by using only nMOS type transistors for the clock circuit. However the system also needs a separate power supply for the local clock inverter. Additionally, since the output stage is a static CMOS latch stage, there will be limits to the operating frequency. As it was shown in [8], [9], dynamic latches or sense-amplifier based latches can provide a better operating frequency. It is important to stress the focus on the operating frequency, as the standing wave techniques are only feasible in the GHz range.

Reference [10] introduces another type of low clock swing logic gate, however, without a differential topology. Here, the circuit employs a combination of low and high threshold voltage transistors. These optional devices may not be offered by the silicon foundries in the standard fabrication package or may come with cost adders for low volumes, multi-project wafers. Again, the output latching stage is a static CMOS latch, having the same speed penalty. It is interesting to note that since the clock transistors are inserted in the middle of the logic chain, the topology has some potential to work with clock signals centered on a common mode voltage of half the supply voltage, but can be affected by the input logic value.

Another potentially useful reduced clock swing logic gate was introduced in [11]. Here, although referred to the inputs, the proposed circuit is a single-ended one, it can be easily extended to differential operation through the elimination of the locally generated inverted logic inputs. However, the topology still employs a combination of regular, low and high threshold voltage transistors and requires an additional power supply for driving the local clock buffers, adding to system complexity.

The circuit introduced in [12] can also be easily converted for differential operation – in essence, the topology of the circuit is the same with that of a dynamic cascode voltage switch gate, but uses a combination of regular, low and high threshold voltage devices. The static CMOS latch for the output stage comes again with a speed penalty and the differential clock signal is not centered on the middle of the power supply range but rather towards the upper power line. The main advantage for the technique of [12] stems from the use of a single supply voltage.

Besides proposing a reduced clock swing circuit for resonant based distribution networks, [13] also estimates the propagation delay introduced by the reduced amplitude of the clock signal. Basically, the paper quantizes the magnitude of

amplitude to phase-shift conversion as a function of the clock amplitude and clock frequency. In the scope of that paper, the delay introduced by the reduced swing is a global effect as the clock distribution has a uniform amplitude. However, for the standing wave based clock distribution of [4] and [6], the amplitude dependent phase shift translates in skews between different parts of the clock domains. With respect to the logic gate, [13] avoids the usage of another power supply by having a special topology for the local reduced swing clock buffers. The output latch is again a static CMOS latch with the same speed penalty.

One of the highest speed dynamic differential CMOS logic gates is that of [14], in fact, a variation of the sense amplifier or static RAM based gate. As seen in fig. 2, the topological feature, which gives also the speed advantage of the latch, is the elimination of the tail transistors driven by the clock signal. Without the tail transistors, the voltage headrooms in the circuit are increased, with a positive impact on transistor sizes and, hence, transistor capacitances. The penalty of the increased speed is, however, the increased power consumption given by the short circuit current during the equalization phase. Although not specified in [14], the proposed synchronous differential logic (SDL) gate, as it will be shown in the present article, can also be made to work with a reduced swing clock signal. The only needed change is a parametric change, with no impact on the circuit topology or on the type of active devices – i.e. there is no requirement for different threshold voltage transistors. As it will be detailed later on, the direct correlation between the clock signal amplitude and the equalization voltage level allows one to reduce the clock swing by reducing the voltage level during the equalization phase, only through transistor sizes changes.

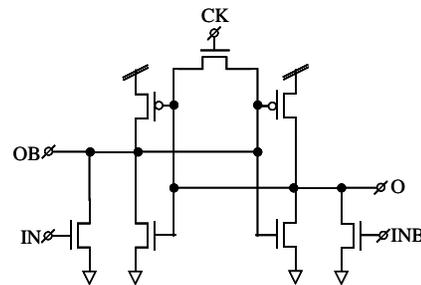


Fig. 2. Synchronous Differential Logic (SDL) gate

As an observation, another advantage, not readily apparent from [14] and only tangential to the scope of the present paper, is the reduced correlation between the data processing performed by the SDL gates and the induced and radiated noise. In cryptographic applications, this advantage can be a major deciding factor, as systems implemented using SDL gates can be made very robust against differential power analysis attacks.

Like it was previously stated, an important aspect of the signal distributed by standing-wave clock distribution networks is the fact that the signal is centered on a fixed voltage and the amplitude only about this fixed common mode changes, according to the coordinates in the network. This

means that both the High and the Low clock signal values are affected. From all the above differential gates, only the SDL gate may be operated with such a signal. For the rest of the gates, using a standing-wave based clock will either keep the clock transistor in ON state for the entire period, either will not allow the clock transistor to fully reach the ON state. This happens like this because the circuits are using as reference for the clock signal one of the supply bars. For the SDL gate, the reference for the clock signal is, in fact, the equalization (built-in) level.

III. PROPOSED REDUCED CLOCK SWING SDL GATE

Before the introduction of the theoretical model it is useful to have an understanding of the SDL gate and its operation.

Topologically, an SDL gate is quite similar with a six transistors SRAM cell. Basically, the gate comprises a CMOS latch circuit, constructed using two cross-coupled inverters, and an input circuitry represented in fig. 2 by the IN and INB transistors. The gate uses both outputs, allowing thus a simple differential utilization. As a difference to the SRAM cell, the input transistors are connected to the ground. By replacing the simple input transistors with other pull-down logic networks it is possible to embed arbitrary logic functions directly in the latch cell.

Another difference to the SRAM cell is the presence of an equalization transistor, connected between the two outputs. This transistor is driven by the clock signal of the SDL gate. Depending on the level of the clock signal, there are two phases of operation for the circuit – when the clock signal is at logic high, the gate is said to be in the equalization phase. In this phase, the clock transistor effectively creates a low impedance path between the complementary outputs. This equalizes the voltages present at the latch outputs, erasing thus any previously stored logic value. In the same time, the clock transistor optimally biases the inverters (in fact gain stages) in the highest gain region. This happens because the voltage level on the two differential outputs becomes centered more or less in the middle of the supply domain, in the region where the slope of the inverter's characteristic is at maximum.

At the end of the equalization phase, the latch is ready to sample the logic value given at its input. The input logic levels are converted by the input transistors or by the pull-down network from voltage levels to input currents for the latch.

In the next phase, when the clock signal is at logic-low, the equalization transistor shuts off, biasing the latch gate in the maximum positive feedback region. In this situation, the latch is in the evaluation phase. Based on the input currents imbalance that was present at the end of the equalization phase (created by the input transistors), the strong positive feedback will regenerate the output logic levels in a very short time. In fact, the current imbalance creates a small voltage difference at the inputs of the cross-coupled inverters. Since the inverters were initially biased at the maximum gain point, this small difference is readily amplified to the proper logic values.

For optimal gate performance, it is necessary to have a

strong enough equalization as to cancel the effect of the positive feedback. This increases the input sensitivity and accelerates the output switching between opposite logic levels. If the equalization transistor cannot cancel the effect of the positive feedback, the cross-coupled inverters will try to re-establish the previously stored logic value. If the stored logic value is opposite to the current one, there will be a race condition between the input transistors and the latch stage. In certain operating conditions (for instance at high temperatures in the slow process corner), this race condition will yield metastable states at the output. With proper equalization, these race conditions are absent as the stored logic value is effectively erased during the equalization phase.

As it will be shown further on, there is a condition past which the equalization transistor effectively cancels the effect of the positive feedback. If that condition is not met, the performance of the gate is impacted at all frequencies – for this reason, the condition will be called the static equalization condition. Even if the condition is met, proper circuit operation is still not achieved at all frequencies – the larger the distance between the static equalization condition and the actual circuit situation, the higher the operating frequency will be, in the limits imposed by the CMOS process node.

Since the SLD gate is in fact an embedded-logic latch, for proper operation, consecutive logic gates must use opposite clock phases – i.e. when one logic stage is in the evaluation phase, the next logic stage is in the equalization phase. It can be noted that during the equalization phase the inputs are properly held at correct logic levels by the previous logic stage.

A. Generic Theoretical Model

Reference [8] showed that it is possible to calculate the transient response of the logic gate based on a small signal model. The focus of that paper was, however, not on the complete theoretical analysis of the latch circuit, but more on the recovery time of the comparator used in the analog to digital converter. Still, it was shown that a reset type latch is capable of reaching higher speeds than a non-reset gate. Although the comparison performed in [8] was between a latch-based amplifier with equalization phase and an amplifier gate with no positive feedback, it can still be shown that the recovery time from one logic level to the opposite is further increased by the presence of the positive feedback. The gate not only has to change its state, but it also has to work against the previous stored logic value. The present paper builds on the referenced model and shifts the focus towards the output latch.

In order to build the small signal model of fig. 3 it is necessary to make some assumptions about the operating conditions of the SDL gate – i.e. there will always be time intervals in which the gate will be biased close to the maximum-gain condition of the cross-coupled inverters (amplifiers). This can happen during an equalization phase, but this can also happen during signal transitions. From this observation, for the small signal model, the cross-coupled

inverters will be considered as amplifiers in the active region (this also covers the topologies where the inverters have a tail transistor, either driven statically to generate a current source, either dynamically, with the clock signal).

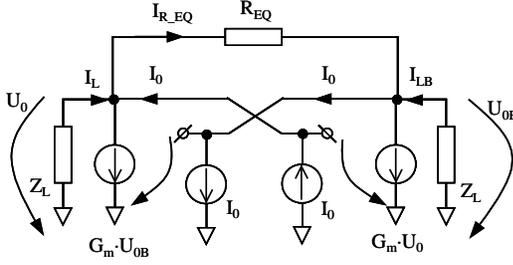


Fig. 3. SDL gate small signal model

For this region, the model of the inverter is assimilated to a voltage-controlled current source G_m working on a load impedance Z_L . This load impedance is basically composed from the equivalent output resistance of the inverter stages (i.e. given by drain diffusions leakage currents, channel length modulation effects) and from capacitive loading (self-loading, parasitic loading and useful loading). The static gain of the stage, given by the controlled current source and the resistive part of the load, can be considered very large, allowing some simplifications in the theoretical model.

The equalization transistor will be modeled as a simple resistor, with its value depending on the operating phase of the logic gate – i.e. quite a small value during the equalization period or a very large value during the evaluation period.

Based on the model of fig. 3, it is possible to derive the system of equations and, in the end, the transistor sizes required at the imposed operating conditions. The model, as drawn, is a generic one as it can represent any sense amplifier based circuit like that of [8], [9], or the circuit introduced in the present paper, or any other precharge/discharge differential cascode voltage switch gate [12].

The small signal model as drawn in fig. 3 represents the latch during the equalization phase. The input current sources are modeling the effect initially stored logic value. Starting from this, it is possible to write the following system of equations:

$$\begin{cases} G_m \cdot U_0 = I_{LB} - I_0 + I_{R_EQ} \\ U_{0B} = -I_{LB} \cdot Z_L \\ I_{R_EQ} = \frac{U_0 - U_{0B}}{R_{EQ}} \\ G_m \cdot U_{0B} = I_L + I_0 - I_{R_EQ} \\ U_0 = -I_L \cdot Z_L \end{cases} \quad (1)$$

In the above system of equations, G_m stands for the total transconductance of the inverter stage. R_{EQ} represents the equivalent resistance of the equalization element (an nMOS transistor for the case of an SDL gate). Z_L corresponds to the loading impedance present at the outputs of the inverter gate

and usually consists of a parallel connection between a load capacitance C_L and load resistance R_L . The rest of the parameters are representing the node voltages and branch currents.

Solving (1) for the difference between the output voltages gives the following result:

$$U_0 - U_{0B} = \frac{2 \cdot I_0}{1/Z_L + 2/R_{EQ} - G_m} \quad (2)$$

Taking into account the parallel composition of the load impedance (the static gain given by the resistance can be considered very large) and solving (2) for the solution of the differential output voltage $U_D = U_0 - U_{0B}$, gives the Laplace transform:

$$U_D = 2 \cdot I_0 \cdot \frac{1}{2 - G_m \cdot R_{EQ}/R_{EQ} \cdot C_L + s} \quad (3)$$

The current pulse I_0 of (3) must provide enough charge as to create the initial differential voltage V_P on the load capacitances. Assuming a step response for the output voltage, the time domain expression for the input current is given by (4), where $\sigma(t)$ stands for the unit step function and $\delta(t)$ stands for the Dirac impulse.

$$i_0 = \frac{d}{dt}(V_P \cdot \sigma(t)) \cdot C_L/2 = V_P \cdot C_L/2 \cdot \delta(t) \quad (4)$$

From (4) the Laplace transform of the input current results directly:

$$I_0 = V_P \cdot C_L/2 \quad (5)$$

Combining (3) and (5) and using the inverse Laplace transform it is possible to write:

$$u_D = V_P \cdot e^{-t \frac{2 - G_m \cdot R_{EQ}}{R_{EQ} \cdot C_L}} \quad (6)$$

As a note, the initial amplitude of the output voltage was assumed to be equal with the peak voltage V_P . At this point, it is important to highlight the static equalization condition. In fact, in (6), if the term $2 - G_m \cdot R_{EQ}$ is zero, the output voltage of the latch becomes undefined. If the same term is negative, the solution (6) gives an ever-growing response, the latch still having a positive feedback. The output voltage decays only if the term from (6) is negative – this translates in a maximum value for the equivalent resistance R_{EQ} . Inequality (7) summarizes the static equalization condition:

$$G_m \cdot R_{EQ} < 2 \quad (7)$$

At the end of the equalization phase, the output voltage must reach a desired value U_{D_EQ} . Regarding the duration of the process, until now, the theoretical analysis used a step response for the clock signal. Assuming a trapezoidal clock signal having transitions equal with a quarter of the clock period $T_{CK}/4$, one approximate way to consider also the effect of non-ideal steps is to use a relation similar with that used for the signal's risetime based on the oscilloscope measurements. The total time is the length of the equalization phase – i.e. $T_{CK}/2$ – and is given by the clock signal risetime and by the exponential decay t_{M_EQ} of (6). With this, the time interval t_{M_EQ} to reach the equalization phase in (6) must not be longer than:

$$t_{M_EQ} = T_{CK} \cdot \sqrt{3}/4 \quad (8)$$

Combining (8) with (6) and solving for the exponential argument gives the final dynamic equalization condition:

$$T_{CK} \cdot \frac{\sqrt{3}}{4} \cdot \frac{1}{\log(V_P/U_{D_EQ})} = \frac{R_{EQ} \cdot C_L}{2 - G_m \cdot R_{EQ}} \quad (9)$$

In (9), $\log()$ represents the natural logarithm function. The above equation sets a correlation between the gain of the inverter stages and the equalization resistor. In order to get also an equation for the required gain at the inverter stages, it is necessary to extend the theoretical analysis also for the evaluation phase of the latch.

Using (6) and letting R_{EQ} be infinite will give the time domain solution for the output transition during the evaluation phase. A difference between the solutions for the equalization and evaluation phases would be the initial and final voltage levels. In fact, the evaluation phase starts from the U_{D_EQ} end-level of the equalization phase and ends when the output voltage reaches again the peak voltage V_P :

$$V_P = U_{D_EQ} \cdot e^{\frac{t_{M_EV} \cdot G_m}{C_L}} \quad (11)$$

The same condition that is used for the decay time estimation is also assumed for the growth time t_{M_EV} of (11) with respect to the clock period and falltime. With this, the final solution for the evaluation phase is given by:

$$T_{CK} \cdot \frac{\sqrt{3}}{4} \cdot \frac{1}{\log(V_P/U_{D_EQ})} = \frac{C_L}{G_m} \quad (12)$$

Since (12) depends only on the current gain G_m of the inverter stages, it can be used to size the transistors from the output latch for proper operation in the required speed conditions. After this step, using (9) and the result given by (12) it is possible also to size the equalization circuit.

B. Transistor sizing for the reduced swing SDL gate

It should be noted that the previous analysis made no assumption regarding the physical construction of the equalization resistor. In fact, the resistor can be made directly as an equalization resistor as seen in [8] or in the SDL gate of [14]. However, by observing that it is possible to split the resistor into two series resistors with a virtual ground node between them, it is evident that the previous demonstration holds also for precharge or discharge logic gates, as seen in [12]. For the rest of the paper, the theoretical analysis assumes the synchronous differential gate of fig. 2. Still, starting from this, the model can be reconstructed for all the other configurations.

As it was said at the end of the theoretical model section, the first step is to size the inverter stages. The current gain G_m of the inverter stage is given by the sum of the gains g_{mn} and g_{mp} of the nMOS and pMOS transistors. Right at the edge between the two phases, additionally to the output transitions, also the logic inputs are switching state. This translates in an intermediate voltage level for the inputs, placed somewhere between ground and the power supply. With this, the input nMOS transistors can also be considered to have the same

operating point as the latch nMOS transistors. However, since they are not included in the positive feedback loop, they are not contributing directly to the gain of the inverter, they only impose a static drain current for the pMOS transistors. This current has a double value compared to that of the nMOS latch transistors. In other words, at equilibrium, the current gain of the pMOS transistors must be equal with the gain of both the input and the latch nMOS transistors – i.e. double the current gain of the latch nMOS transistor. With these, the gains of the nMOS and pMOS transistors may be expressed as a function of the total current gain by the following system:

$$\begin{cases} g_{mn} = 1/3 \cdot G_m \\ g_{mp} = 2/3 \cdot G_m \end{cases} \quad (13)$$

The current gains of the transistors can be written also as a function of transistor sizes, process parameters and operating point values:

$$\begin{cases} g_{mn} = K_n \cdot W_n / L_n \cdot (V_{EQ} - V_{THn}) \\ g_{mp} = K_p \cdot W_p / L_p \cdot (V_{DD} - V_{EQ} - V_{THp}) \end{cases} \quad (14)$$

In (14), K_n , K_p are the transistor process specific current gains for the nMOS and pMOS respectively. V_{THn} , V_{THp} are the threshold voltages of the same transistors – for the circuit configuration of [14] and fig. 2, they are not affected by the body effect. For the latch of [8], since the latch has also a tail transistor, the nMOS transistors will have a voltage difference between the source and the body terminal, and hence exhibit body effect. The operating point parameters are given by the V_{EQ} voltage during the equalization phase and by the power supply V_{DD} . Transistor sizes are given by widths W_n , W_p and lengths L_n , L_p . Combining (14), (13) and (12), the latch nMOS and pMOS transistors sizes are given by:

$$\begin{cases} W_n = L_n \cdot \frac{4 \cdot C_L}{3 \cdot \sqrt{3} \cdot T_{CK}} \cdot \frac{\log(V_P/U_{D_EQ})}{K_n \cdot (V_{EQ} - V_{THn})} \\ W_p = L_p \cdot \frac{8 \cdot C_L}{3 \cdot \sqrt{3} \cdot T_{CK}} \cdot \frac{\log(V_P/U_{D_EQ})}{K_p \cdot (V_{DD} - V_{EQ} - V_{THp})} \end{cases} \quad (15)$$

Note that, for simplicity, absolute values were considered for the pMOS process parameters.

For the equalization transistor it is necessary first to evaluate its equivalent resistance. By inspecting the terminal voltages and the transistor current during the transition from evaluation to equalization, the switch goes from cut-off to saturation and then to the linear region. However, the clock signal rises sooner than the source terminal, resulting in a quite large effective gate to source voltage, allowing for a high initial drain current. When the drain to source voltage becomes equal with the gate overdrive, the transistor enters linear region, but at this point in time, its effective gate to source voltage is about half of that from the beginning. Even more, as the source terminal has risen to the equalization level, the threshold voltage will be strongly affected by the body effect, further reducing the allowable drain current. With all these, the ratio between the drain to source voltage and the drain current remains rather flat during the transition. In the case of the full swing clock signal, the time-averaged value of the resistance is

close to the final value at the end of the equalization interval. When the clock swing is reduced however, the time-averaged resistance increases. This can be explained by noting that at the beginning of the equalization phase, when the clock signal has already reached a value near its maximum, the effective gate overdrive voltage is still small. With this, the available equalization current decreases with decreasing clock swing. The effect can be quantized starting from the drain current in the saturation region and considering the ratio between the achievable currents for full swing and reduced swing.

Starting from the expression of the transistor's drain current operating in the linear region, with a non-negligible drain to source voltage, the equivalent resistance for the full swing situation can be written as:

$$R_{EQ} = \frac{I}{K_n \cdot \frac{W_{n_EQ}}{L_{n_EQ}} \cdot \left(V_{CK} - V_{EQ} - \left(V_{TH_EQ} - \frac{V_{DS}}{2} \right) - V_{DS} \right)} \quad (16)$$

In (16), the operating point parameters are given by the maximum level of the clock signal V_{CK} , the equalization level V_{EQ} and by the drain to source V_{DS} voltage. As process parameters, the relation uses the nMOS current gain K_n and a different threshold voltage V_{TH_EQ} (affected by the body effect). It is interesting to note that the source terminal is lower than the equalization level by half of the drain to source voltage – this is so because the drain to source voltage is the differential output voltage that is centered on the common mode equalization level. Rearranging (16) by separating the V_{DS} term from the rest gives:

$$R_{EQ} = R_{EQ0} \cdot \frac{I}{I - \frac{I}{2} \cdot \frac{V_{DS}}{V_{CK} - V_{EQ} - V_{TH_EQ}}} \quad (17)$$

$$R_{EQ0} = \frac{I}{K_n \cdot W_{n_EQ} / L_{n_EQ} \cdot (V_{CK} - V_{EQ} - V_{TH_EQ})}$$

The step from (17) highlights the influence of the differential voltage on the resulting equivalent equalization resistance. The differential voltage effectively creates a multiplier term for the ideal R_{EQ0} resistance. For practical circuits, the differential voltage V_{DS} is comparable with the gate overdrive voltage $V_{CK} - V_{EQ} - V_{TH_EQ}$ – i.e. the equalization transistor is always close to the saturation region.

As a next step it is necessary to introduce in (17) the correction factor associated with the reduced swing operation for the SDL gate.

In (18), V_{CK_F} stands for the full swing maximum clock value. It is important to note that the correction factor depends on the nMOS threshold voltage with zero body-source bias.

$$R_{EQ} = R_{EQ0} \cdot \left(\frac{V_{CK_F} - V_{THn}}{V_{CK} - V_{THn}} \right)^2 \cdot \frac{I}{I - \frac{I}{2} \cdot \frac{V_{DS}}{V_{CK} - V_{EQ} - V_{TH_EQ}}} \quad (18)$$

$$R_{EQ0} = \frac{I}{K_n \cdot W_{n_EQ} / L_{n_EQ} \cdot (V_{CK} - V_{EQ} - V_{TH_EQ})}$$

This happens because the correction assumes that the output latch had no time to change its output voltages and hence, the source of the equalization transistor is close to the ground

potential. With this and combining (18) with (9), the equalization transistor can be sized using the following relation:

$$W_{n_EQ} = \left(\frac{V_{CK_F} - V_{THn}}{V_{CK} - V_{THn}} \right)^2 \left(3 \cdot W_n \cdot \frac{L_{n_EQ}}{L_n} \cdot \frac{V_{EQ} - V_{THn}}{V_{CK} - V_{EQ} - V_{TH_EQ}} + \frac{4 \cdot C_L}{\sqrt{3} \cdot T_{CK}} \cdot L_{n_EQ} \cdot \frac{\log(V_P / U_{D_EQ})}{K_n \cdot (V_{CK} - V_{EQ} - V_{TH_EQ})} \right) \quad (19)$$

As was said in the second section of the present paper, (19) shows the existing correlation between the clock signal level V_{CK} , the equalization level V_{EQ} and the transistor sizes. In fact, it is evident from (19) that if the clock swing is reduced, changing also the equalization level with a similar amount can compensate this reduction. However, transistor sizes are not kept constant, as the reduction in the equalization level changes also the transistors from the inverter stages.

Regarding V_P , nominally, the amplitude in the system is limited at the V_{DD} power supply level. For the full swing clock, this assumption is true, as one output is rising from half the supply voltage to the V_{DD} , while the other output falls from the middle point to the ground level. However, if the equalization level is decreased, the output switching becomes asymmetric. Since both outputs must reach the steady state value in the allowable time window, the circuit behaves as if the differential amplitude is increased above the power supply. Given an equalization level V_{EQ} , the peak voltage will be calculated by $V_P = 2 \cdot (V_{DD} - V_{EQ})$.

IV. SIMULATION RESULTS

The SDL gate was simulated using an IBM 130nm CMOS process operating at 1.2V supply voltage. All used transistors were of the regular threshold voltage type with the minimum gate length supported by the technology (120nm).

The target clock frequency was set to 3.5GHz. The clock signal was centered around half of the supply voltage – i.e. 600mV and had the amplitudes of 300mV, 450mV and 600mV (going from 50% to 75% and up to 100% swing). Combining the common mode voltage with the signal amplitude gives the following maximum V_{CK} levels: 0.9V, 1.05V and 1.2V. Since the equalization level is correlated to the maximum clock level, the chosen values for the equalization V_{EQ} voltage are 600mV for full clock swing, 450mV for 75% clock swing and 300mV for 50% swing.

In the selected technology, the current gain constants for the transistors are around 465 $\mu\text{A}/\text{V}^2$ for the nMOS and around 65 $\mu\text{A}/\text{V}^2$ for the pMOS. The threshold voltage for the nMOS without body effect is around 360mV and for the pMOS is around 350mV. Considering also the body effect, for the equalization transistor, the threshold voltage is found to be around 420mV at 300mV body to source voltage, 445mV for 450mV and around 470mV for 600mV body to source voltage.

As a first observation, because the nMOS threshold voltage is minimum 360mV, the theoretical model has a problem sizing the latch nMOS transistor for the lowest clock swing. However, simulations showed that the latch can work also at

these operating conditions – it should be stressed that, although it produces good results in most of the cases, the model is only an approximation. Since (19) used as reference the nMOS transistor sizes, the model has again a problem to size the equalization transistor. However, the model can be changed to use the pMOS transistor sizes as reference.

The SDL gate is loaded with a 20fF capacitance in each output node. Additionally, the model also includes self-loading effects (parasitic capacitances, Miller effect), resulting in a total load capacitance of around 50fF per output. The final differential voltage at the end of the equalization phase was assumed to be around 10% of the power supply – i.e. 120mV.

Table I bellow synthesizes the results obtained for the transistor dimensions by applying the theoretical model and through simulation. The simulated transistors are obtained from multi-finger devices, so table I gives the physical size and the equivalent size for the simulated transistors.

Table I. Transistor sizes as function of equalization level

			Equalization level [mV]		
			300	450	600
Width [μm]	nMOS	Model	N/A	0.98	0.33
		Sim	1.8 (4 x 0.45)	0.9 (2 x 0.45)	0.45
	pMOS	Model	2.45	3.14	4.58
		Sim	3 (5 x 0.6)	4.2 (7 x 0.6)	4.8 (8 x 0.6)
	Equalization	Model	11.39 (from pMOS)	5.05	3.68
		Sim	11.05 (17 x 0.65)	7.8 (12 x 0.65)	4.55 (7 x 0.65)

As it can be seen, the theoretical model gives relatively good size estimates for the SDL’s gate transistors. Both pMOS and the equalization transistors are undersized in the model but this may be due to second order effects that were not taken into account in the model.

Fig. 4 shows the simulation results for the three clock signal swings and gate sizes. The simulation was done such as, during the equalization phase, to have the input transistors biased close to the latch nMOS transistors. It can be seen that in all operating conditions, the differential voltage at the end of the equalization phase is close to the desired 120mV. Fig. 5 shows the behavior of the three gates with normal input signals.

Based on the simulated transistor sizes, a universal programmable gate was constructed using CMOS selectors to add or remove additional transistors. For instance, looking at the equalization transistor, for the full swing clock signal, the transistor consists of 7 fingers. For the 75% clock swing, the transistor has another 5 fingers added. For the half swing clock signal, the same transistor has a total of 17 fingers. Using CMOS pass-transistors, the gates of the additional fingers are conditionally added to the clock signal, depending on the configuration bits. No other transistor resizing was done, so the programmable gate used the final values from Table I.

Fig. 6 shows a simplified circuit topology for the universal gate and fig. 7 shows the simulation results for this gate at all clock signal swings. The test set-up for the programmable gate

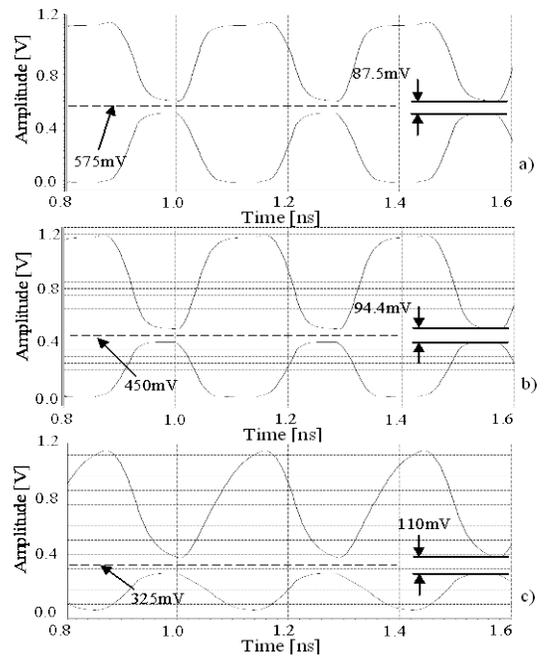


Fig. 4. Simulation results with equal bias inputs: a) full swing clock; b) 75% clock swing; c) half clock swing

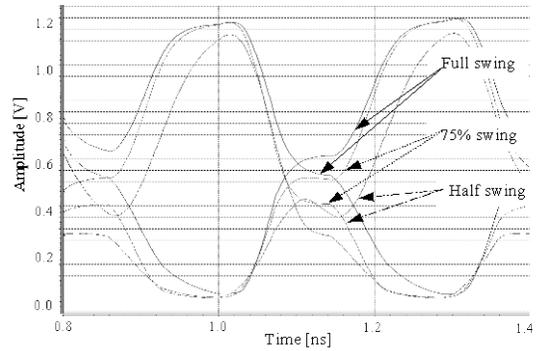


Fig. 5. Time alignment of full, 75% and half swing gates outputs

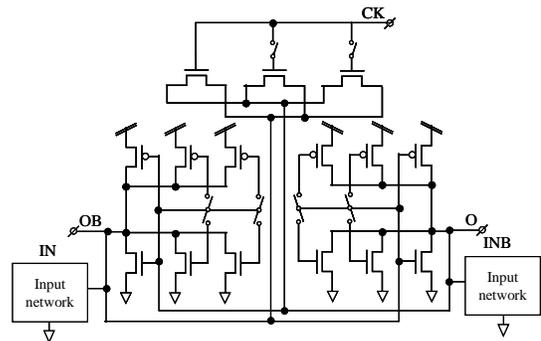


Fig. 6 Block diagram of the programmable universal reduced clock swing SDL gate

is the same as that used to generate fig. 5. As expected, the additional circuitry degrades a bit the performance of the gate. Both the rising and the falling edges are becoming slower and hence, the final voltage levels reached during the evaluation phase are not anymore close to the supply lines. However, the gate still performed the correct decisions. For the scope of the present paper, no attempt was made to compensate for this

degradation.

Although an attempt was made to size also a circuit

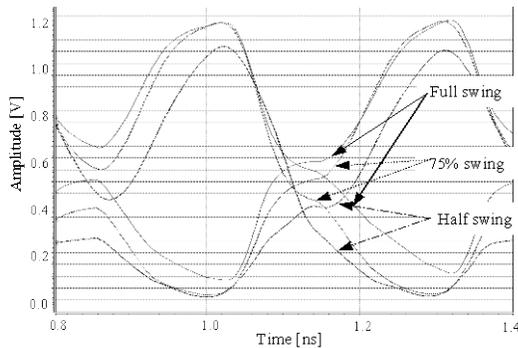


Fig. 7 Universal programmable gate simulation for all clock swings

according to [7], it was found that using the same clock signals as for the SDL gate is impractical – for the low clock swing, the minimum clock level of around 300mV prevents the clock transistor to completely shut-off.

V. ABOUT THE IMPACT OF PROCESS CORNERS

Since the design equations are generic with respect to the process parameters (they make no assumption about the threshold voltages or process gains), it is possible to use them with the worst-case parameters. This will give an indication for the impact on the circuit operation.

Regarding the worst-case combination, at first glance it may seem that slow nMOS coupled with fast pMOS will hinder the operation at low clock swings, as the equalization level will shift upwards. However, if the clock distribution network uses a self-bias technique for the common mode voltage, also the clock common mode voltage will exhibit the same shift.

In fact, the worst-case process corner is slow nMOS / slow pMOS as both devices will have degraded performance while the equalization level will remain at the nominal level. The problem stems from the increased threshold voltage of the nMOS. One way to mitigate it is to switch to lower threshold devices, but this is not always an option. Although uncommon in the digital IC design, many analog ICs are taking profit of the threshold voltage variation versus the channel length. For instance, in the selected technology, by tripling the channel length of the nMOS, the resulting threshold voltage becomes about 270mV. Even if, at first, the impact may seem drastic, the reduction of the threshold voltage can have a big impact on the operating point of the transistor. As an illustration, keeping the same aspect ratio but using a tripled channel length value, the nMOS transistor for the 75% clock swing gate will have the current gain becomes two times larger than that of the initial transistor.

VI. CONCLUSION

The present paper introduced a new low clock swing logic gate capable of operating with a standing wave based clock distribution network. This is an important aspect as, unlike

other clock distribution techniques, the amplitude of the distributed clock signal varies across the IC surface, but the logic gates must still operate in a synchronous fashion. Additionally, the clock common mode voltage remains fixed, independent of the standing wave pattern, posing problem for most of the other low swing logic gates.

For the SDL gate, based on the theoretical model, it was shown that in order to reduce the clock swing, one must also reduce the voltage level during the equalization phase.

Simulations showed a good correlation between the theoretical model gate sizes and the required final sizes. In the end, a programmable universal SDL could be constructed, capable to operate with different clock swing levels, depending on the configuration bits.

REFERENCES

- [1] A. J. Drake, K. J. Nowka, T. Y. Nguyen, J. L. Burns, R. B. Brown, "Resonant Clocking Using Distributed Parasitic Capacitance", *IEEE Journal of Solid-State Circuits*, vol. 39, no. 9, pp. 1520-1528, September 2004.
- [2] S. C. Chan, K. L. Shepard, P. J. Restle, "Uniform-Phase Uniform-Amplitude Resonant-Load Global Clock Distributions", *IEEE Journal of Solid-State Circuits*, vol. 40, no. 1, pp. 102-109, January 2005.
- [3] J. Wood, T. C. Edwards, S. Lipa, "Rotary Traveling-Wave Oscillator Arrays: A New Clock Technology", *IEEE Journal of Solid-State Circuits*, vol. 36, no. 11, pp. 1654-1665, November 2001.
- [4] F. O'Mahony, C. P. Yue, M. A. Horowitz, S. S. Wong, "A 10-GHz Global Clock Distribution Using Coupled Standing-Wave Oscillators", *IEEE Journal of Solid-State Circuits*, vol. 38, no. 11, pp. 1813-1820, November 2003.
- [5] V. L. Chi, "Salphasic Distribution of Clock Signals for Synchronous Systems", *IEEE Transactions on Computers*, vol. 43, no. 5, pp. 597-602, 1994.
- [6] A. Paşca, "Bi-dimensional Radially-Salphasic (Standing Wave) Clock Distribution", 2014 IEEE 20th International Symposium for Design and Technology in Electronic Packaging (SIITME), pp. 157-162, Bucharest, 23-26 October, 2014.
- [7] YS. Kwon, IC. Park and CM. Kyung, "A New Single-Clock Flip-Flop for Half-Swing Clocking", *Proceedings of the ASP-DAC '99 Design Automation Conference, Asia and South Pacific*, vol. 1, pp 117 - 120, Wanchai, 18-21 January 1999.
- [8] M. Choi, A. Abidi "A 6b 1.3 Gs/s A/D Converter in 0.35um CMOS", *IEEE Journal of Solid-State Circuits*, vol.36, pp 1847-1858, December 2001.
- [9] B. Nikolic, V. G. Oklobdzija, V. Stojanovic, W. Jia, J. K.-S. Chiu, M. M. Leung, "Improved Sense-Amplifier-Based Flip-Flop: Design and Measurements", *IEEE Journal of Solid-State Circuits*, vol. 35, no. 6, pp. 876-884, June 2000.
- [10] C. Kim, S. King, "A Low-Swing Clock Double-Edge Triggered Flip-Flop", *IEEE Journal of Solid-State Circuits*, vol. 37, no. 5, pp. 648-652, May 2002.
- [11] D. Levacq, M. Yazid, H. Kawaguchi, M. Takamiya, T. Sakurai, "Half V_{DD} Clock-Swing Flip-Flop with Reduced Contention for up to 60% Power Saving in Clock Distribution", 33rd European Solid State Circuits Conference ESSCIRC 2007, pp. 190 - 193, Munich, 11-13 September 2007.
- [12] K. Mohammad, B. Liu, S. Agaian, "Energy efficient swing signal generation circuits for clock distribution networks", *IEEE International Conference on Systems, Man and Cybernetics SMC 2009*, pp. 3495 - 3498, San Antonio TX, 11-14 October 2009.
- [13] S. Esmaili, A. J. Al-Kahlili, G.E.R. Cowan, "Low-Swing Differential Conditional Capturing Flip-Flop for LC Resonant Clock Distribution Networks", *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, volume 20, issue 8, pp. 1547 - 1551, 2011.
- [14] C. C. Timoc, "Synchronous differential logic system for hyperfrequency operation", US patent, number US6002270 A, issue date 14 Dec. 1999.

On-line Monitoring of Yogurt Fermentation Using Ultrasonic Characteristics

Ahmad Aljaafreh

Communications, Electronics and Computer Dept.
Tafila Technical University
Tafila, Jordan
a.aljaafreh@ttu.edu.jo

Ralf Lucklum

Institute for Micro and Sensor Systems (IMOS)
Otto-von-Guericke-University Magdeburg
Magdeburg, GERMANY
ralf.lucklum@ovgu.de

Abstract— Fermentation is the process where sugars are transformed into lactic acid. pH meters have traditionally been used for fermentation process monitoring based on acidity. Ultrasonic systems can provide a rapid, accurate, inexpensive, simple and non-destructive method to on-line assess and monitor the properties of food during process operations. This paper evaluates the use of ultrasonic measurements to characterize yogurt fermentation process by correlating acoustic properties and fermentation process characteristics. This research shows the correlation between fermentation time and acoustic attenuation as well as acoustic velocity. It also shows the effect of temperature on the received signal attenuation and velocity for yogurt and milk.

Keywords—Ultrasonic; yogurt fermentation; sound attenuation; sound velocity

I. INTRODUCTION

Ultrasonic (US) sensors are proved to be effective in many industrial and medical applications. Most known are nondestructive testing (NDT) and ultrasonic imaging. Ultrasonic sensors are used as flow and level meters in process industry. Recently new fields of ultrasonic sensor technology have emerged, their importance is increasing. This development is partly boosted by improvements in transducer development and signal processing [1], [2], and [3]. Technology roadmap points out that ultrasonic will be among the emerging techniques that solve future problems in process control [4]. Ultrasonic characteristics are used for non-invasive process control by correlating sound parameters and characterizing process parameters. One specific and challenging application is the analysis of liquid multi-phase mixtures like suspensions, emulsions and dispersions [5]. Ultrasonic sensor systems are used instead of industrial chemical sensors because of their fast response, robustness and reliability.

Sound velocity, sound absorption and acoustic impedance can be applied as acoustic properties together with advanced signal processing to measure liquid mixture properties and to continuously control and monitor processes [6]. Ultrasonic measurement is based on the change in the properties of the transmitted acoustic wave, which are influenced by the medium [7] as in shown in Fig. 1. Information that is needed

to characterize the fluid media in time and space is contained in the transmitted or reflected ultrasonic waves [3]. This kind of measurement is considered as indirect method which makes it affected by other unwanted phenomena [3]. Ultrasonic sensors and their advantages, disadvantages and limitations of ultrasonic process are discussed in [8]. Henning describes the state of technology in the field of the computer-assisted ultrasonic transducer development and their limits [9].

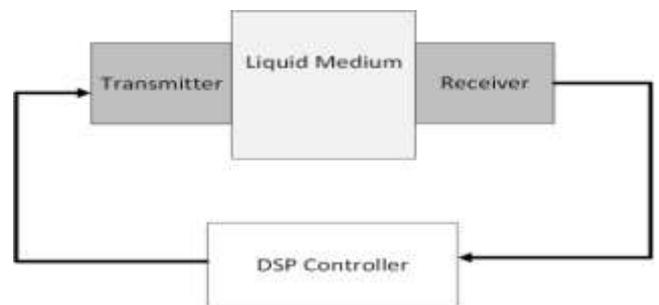


Fig. 1 Generic Ultrasonic System Block Diagram

This paper verifies the use of the two basic sound parameters (velocity, attenuation) for yogurt fermentation process monitoring and control.

II. YOGURT FERMENTATION

The ultimate goal of this research is to automate the butter churning process using ultrasonic sensors which will improve the fermentation and churning process efficiency. Technological innovation is one of the drivers of increased payoffs for the dairy industry [10]. Ultrasonic systems are increasingly being used in the dairy industry. Butter has traditionally been made from yogurt. When a sufficient amount of milk has been collected, it is fermented then churned by shaking until butter granules are formed. Fermentation is the process where sugars are transformed into lactic acid. pH meter has been used for fermentation process monitoring based on acidity [11]. Extensive cleaning and calibration make pH measurement not preferred for on-line monitoring. Ultrasonic systems can provide a rapid, accurate,

inexpensive, simple and non-destructive method to on-line monitor the properties of foods during process operations. Ultrasonic is not an off-the-shelf technology. Thus it needs to be developed and scaled up for each application. Ultrasonic sensor systems can be utilized to continuously monitoring fermentation and churning processes to allow for inline control of the process. For example, ultrasonic sensor have the capability to replace pH sensor in fermentation process provided a correlation can be established to pH, conductivity and/or density values during the churning process. This can be implemented by utilizing artificial intelligence as well as digital signal processing techniques [12].

III. EXPERIMENTS AND MATERIALS

Two ultrasonic transducers have been used as transmitter and receiver as shown in Fig. 2. A 1 MHz 5-period burst has been generated (AWG 520, Sony Tektronics) and amplified to 10 V (AR 75A520, Amplifier Research) before transmitting it through the milk (transmitter: Olympus Panametrics V302, receiver: pico 1.2, Physical Acoustics Corp., Princeton, USA). Alignment has been performed with 3D stages (Newport). An oscilloscope (wavepro 700, LeCroy) has been used to obtain the time of flight which is the time taken by the ultrasonic pulse to travel through the milk. A double-jacket vessel contains the sample. Temperature has been maintained by a temperature controller with an accuracy of 0.1 K. A stirrer has been used for proper shaking and homogenization. Homogenized and fresh milk and plain yogurt were purchased from the local grocery store, Magdeburg, Germany. Raw milk has been provided by Wanke Agrar GmbH, Cobbel, Germany. Plain yogurt was used as the yogurt starter culture. The fermentation has been performed with 400 cc of milk and 40 cc of yogurt starter at 40° C. The ultrasonic velocity was calculated dividing the distance between the transmitter and receiver transducer by the time of flight. The received ultrasonic peak to peak amplitude is measured by the oscilloscope. A pH probe (oMX 3000, WTW Weilheim, Germany) has been used for on-line measurement of yogurt acidity.

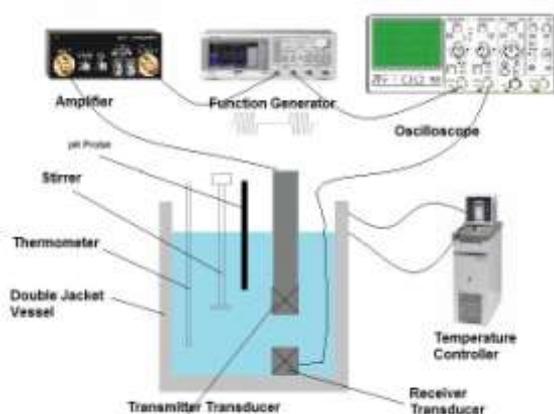


Fig. 2 Experimental setup

IV. RESULTS AND DISCUSSION

A. Acoustic Attenuation

Relative peak to peak amplitude is used to measure the change in the amplitude relative to the peak to peak amplitude at the beginning of fermentation process. Absolute attenuation is not used because it depends on many other factors. Relative peak to peak amplitude provides remarkable information for monitoring the change in yogurt fermentation process. Correlation was found between acoustic attenuation and the fermentation process as shown in Fig. 3. These findings can directly be used to model the yogurt fermentations process. The pH values determined during the fermentation process are shown in Fig. 4. pH decreases with fermentation as attenuation does. The slope of the attenuation is very similar to the slope of the relation of pH with fermentation time. Some distinct differences larger than the confidence range need further analysis. Hence, relative attenuation can be used to monitor the fermentation process as a replacement for pH measurement.

Care should be taken when measuring the amplitude of the received signal. An opposite relation could be measured if stirrer is used in high speed. We assume separation of water from milk although it could not be observed visually. Ultrasonic waves going through water should cause an ultrasound amplitude increase with the fermentation process as shown in Fig. 5 with the stirrer at high speed.

B. Acoustic Velocity

Relative time-of-flight is considered in this research with setting time-of-flight to zero at the beginning of the experiment. In this way uncertainties in alignment and effective distance between the transducers are compensated. Fig. 6 shows a remarkable correlation between the relative time of flight and the fermentation process. The time-of-flight measurements, however, have not been affected when stirrer is used at high speed. The ultrasound velocity increases with the fermentation process.

All the experimental results displayed in Fig. 3-6 are the average of a set of experiments under the same conditions. Moreover, when keeping the milk in a fridge below 10° C, we did not find a systematic tendency within several days.

C. Discussion

Our first steps toward a real-time, on-line, non-contact monitor of the fermentation process with ultrasound have shown an adequate relation between attenuation and speed of sound with fermentation time. Both acoustic characteristics can therefore be used to monitor the fermentation process. However, the difference in the amplitude between the start and the end of fermentation is quite bigger and smoother than the time of flight. The second advantage of the amplitude measurement is that it does not depend on the temperature as shown in Fig. 7, whereas Fig. 8 shows an effect of temperature on sound speed of yogurt and milk.

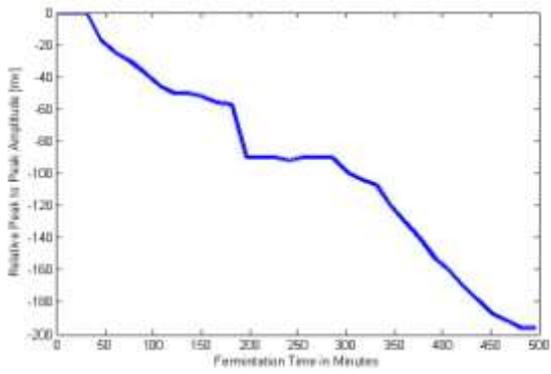


Fig. 3 Relation between amplitude of the received signal and fermentation time. Amplitude decreases with the fermentation time.

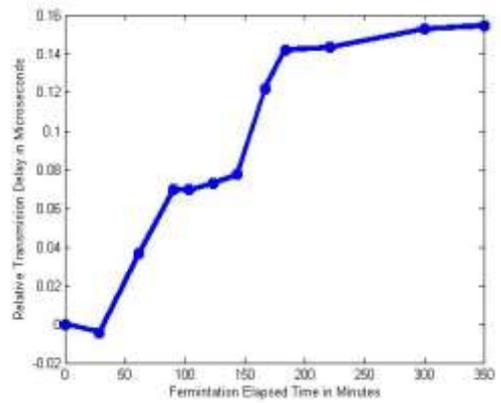


Fig. 6 Relation between relative time of flight and fermentation time.

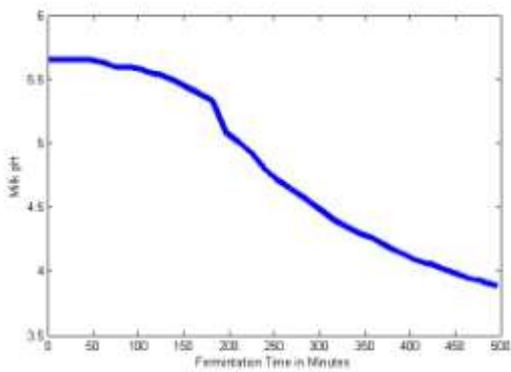


Fig. 4 pH measured using a pH meter. pH decreases with fermentation

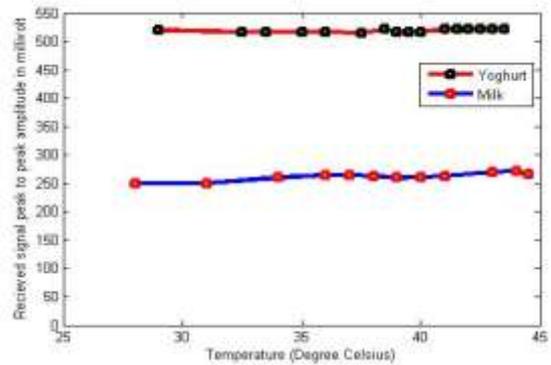


Fig. 7 Temperature effect on the received signal attenuation for yogurt and milk

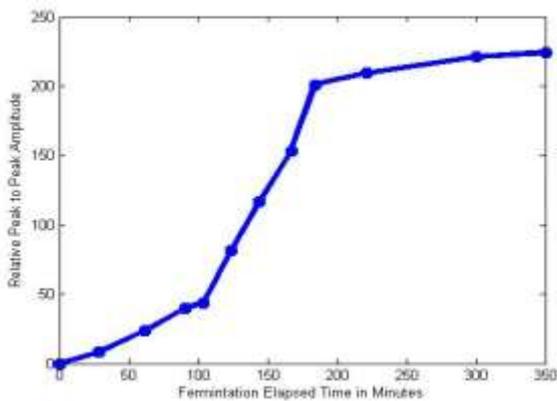


Fig. 5 Relation between amplitude of the received signal and fermentation time when stirrer is used. Amplitude increases with the fermentation time.

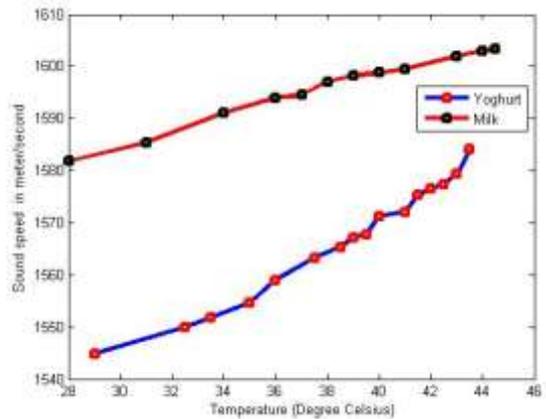


Fig. 8 Temperature effect on the sound speed of yogurt and milk

V. CONCLUSIONS AND FUTURE WORKS

Fermentation is the conversion of lactose to lactic acid by bacteria. When milk is fermented to make yogurt its elasticity increases accordingly. Relative amplitude attenuation and time of ultrasonic flight is measured during fermentation process. Correlation was found between acoustic characteristics and the fermentation process. During the fermentation, acoustic attenuation and velocity change due to the change in the nature of the crossed middle. Results showed that relative ultrasonic measurements (velocity and attenuation) can be used to characterize yogurt fermentation process. The amplitude measurement is less dependent on temperature than ultrasonic velocity measurement. Stirring at high speed can have a big effect on the ultrasonic measurement of fermentation process; we assume separation of water and milk whereas air bubbles are much less likely. This finding has to be considered for future work, where ultrasonic measurement will be evaluated for butter churning process which requires stirring.

References

- [1] Gonzalez Hernandez, J. R., and Chris J. Bleakley. "Low-Cost, Wideband Ultrasonic Transmitter and Receiver for Array Signal Processing Applications." *Sensors Journal*, IEEE 11.5 (2011): 1284-1292.
- [2] Ru, Yan, and Saba Mylvaganam. "Performance monitoring of ultrasonic transducers with laser vibrometers." *Sensors Applications Symposium (SAS)*, 2011 IEEE. IEEE, 2011.
- [3] Woeckel, Sebastian, Ulrike Hempel, and Joerg Auge. "Inline monitoring of liquid multiphase systems." *2012 IEEE International Conference on Instrumentation and Measurement Technology (I2MTC)*, IEEE 2012.
- [4] Püttmer, Alf. "New applications for ultrasonic sensors in process industries." *Ultrasonics* 44 (2006): e1379-e1383.
- [5] Henning, Bernd, and Jens Rautenberg. "Process monitoring using ultrasonic sensor systems." *Ultrasonics* 44 (2006): e1395-e1399.
- [6] Schäfer, Robert, Johan E. Carlson, and Peter Hauptmann. "Ultrasonic concentration measurement of aqueous solutions using PLS regression." *Ultrasonics* 44 (2006): e947-e950.
- [7] Muhamad, I. R., Y. A. Wahab, and S. Saat. "Identification of water/solid flow regime using ultrasonic tomography." *2012 International Conference on System Engineering and Technology (ICSET)*, IEEE 2012.
- [8] Hauptmann, Peter, Niels Hoppe, and Alf Puettmmer. "Ultrasonic sensors for process industry." *2001 IEEE Ultrasonics Symposium Proc.*, Vol. 1. IEEE 2001.
- [9] Henning, Bernd. "Trends in ultrasonic transducer design." *9th International Conference on Electronic Measurement & Instruments, (ICEMI'09)*, IEEE 2009.
- [10] Augustin, M. A., et al. "Towards a more sustainable dairy industry: Integration across the farm–factory interface and the dairy factory of the future." *International Dairy Journal* 31.1 (2013): 2-11.
- [11] R. Meng, J. Zhou, X. Ye, D. Liu, "On-line monitoring of yogurt fermentation using acoustic impedance method.", *Applied Mechanics and Materials* 101-102 (2012): 737-742.
- [12] Henning, Bernd, and Jens Rautenberg. "Process monitoring using ultrasonic sensor systems." *Ultrasonics* 44 (2006): e1395-e1399.

Automatic Censoring in K-Distribution for Multiple Targets Situations

N. Boudemagh, Z. Hammoudi

Abstract—The parameters estimation of the K-distribution is an essential part in radar detection. In fact, presence of interfering targets in reference cells causes a decrease in detection performances. In such situation, the estimate of the shape and the scale parameters are far from the actual values. In the order to avoid interfering targets, we propose an Automatic Censoring (AC) algorithm of radar interfering targets in K-distribution. The censoring technique used in this work offers a good discrimination between homogeneous and non-homogeneous environments. The homogeneous population is then used to estimate the unknown parameters by the classical Method of Moment (MOM). The AC algorithm doesn't need any prior information about the clutter parameters nor does it require both the number and the position of interfering targets. The accuracy of the estimation parameters obtained by this algorithm are validated and compared to various actual values of the shape parameter, using Monte Carlo simulations, this latter show that the probability of censoring in multiple target situations are in good agreement.

Keywords—Parameters Estimation, Method of Moments, Automatic Censoring, K Distribution

I. INTRODUCTION

IN higher resolution radars, the background clutter distribution quite often deviates from Rayleigh and shows long-tail characteristics. The non Gaussian distribution fits will the data acquired from a number of clutter environments, especially in the urban, forest and/or sea environments. The most commonly used models to represent this type of radar clutter are Weibull, log-normal and K, in many practical situations, can be modelled by K-distribution [1], the Probability Density Function (PDF) of the K-distribution has two parameters, the scale a and the shape parameter ν , both of which are usually unknown and need to be estimated the parameters of these distribution from limited quantities of observed data.

In [2], a variety of estimating techniques have been applied to the K-distribution, with most falling under either maximum likelihood (ML) estimation or MOM [3]-[5]. For the K-distribution, ML techniques are typically discarded in favor of a MOM approach as a consequence of their lack of closed-form solution, which typically necessitates a two-dimensional

search or iterative solution. ML techniques for K-distribution can also be computationally demanding owing to the need for evolution of K-Bessel functions for every data sample in the likelihood function. The MOM techniques, however, suffer from a non-zero probability that the moment equations are not invertible [5], with a higher probability when the sample size is small or the shape parameter large.

In practice, the environment is usually heterogeneous due to the presence of multiple targets and/or clutter edges in the reference window cells. In such situations, the main difficulty resides in the parameters estimation of the K-distribution. The automatic censoring techniques have, for their part, many contributed in the improvement of the discrimination between homogeneous and non-homogeneous environments. The well-known approaches proposed in the literature are found in [6]–[8] for a Gaussian clutter and [9]–[11] for a non Gaussian clutter.

In this work, we consider the problem of automatic censoring in presence of an unknown number and position of interfering targets in reference cells for K-distribution. To do this, the censoring algorithm proposed in this work offers a good discrimination between homogeneous and non-homogeneous environments. The first step consists to create the new resamples techniques from the samples of reference cells. Secondly, the proposed algorithm uses the Fine to Coarse (FTC) segmentation algorithm [12] in order to determine the appropriate thresholds to separate the modes in the histogram of the means of each resample. Thirdly, we can make use all resamples corresponding to the means of the first mode (homogeneous population), to estimate the unknown parameters by the classical MOM.

The performances of the proposed algorithm, in terms of the probability of censoring and the estimation accuracy are validated and compared in multiple target situations by means of extensive Monte Carlo simulations.

The remainder of the paper is organized as follows: in section II, the problem formulation and the method of moment for parameter estimation are introduced. The censoring algorithm is presented in Section III. In section IV, we discuss the performance analysis results by means of Monte Carlo simulations. Finally, the conclusions and perspectives are provided in section V.

II. ASSUMPTIONS AND PROBLEM FORMULATION

At the input of the envelope detector, the in-phase and quadrature phase signals (I, Q) are square-law envelope detected and fed into a tapped delay line of length $N+1=2n+1$.

N. Boudemagh is with the Electronic Department, University of Constantine 1, Compus Hamani, Route Ain El Bey, Constantine, Algeria (phone:+213-06-69-93-47-57;fax:+213-38-92-70-02;e-mail:boudemagh.naime@yahoo.com).

Z. Hammoudi, is with the Electronic Department, University of Constantine 1, Compus Hamani, Route Ain El Bey, Constantine, Algeria (e-mail: z_hammoudi@yahoo.fr).

The $(N+1)$ samples correspond to the reference cells $X_i (i = 1, \dots, N)$ surrounding the cell under test X_0 . We consider Rayleigh fading model is assumed for fluctuating targets and correspond to Swerling II (SWII) case in single pulse processing [13]. For a homogeneous noise plus clutter level, we will assume that the clutter signal at output of the envelope detector follows the compounded K-distribution and is independent, identically distributed (IID) in each samples. The statistics of a K-distributed random variable X are described by the probability density function, see [3].

$$f_{X_i}(x) = \frac{2}{a \Gamma(v+1)} \left(\frac{x}{2a}\right)^{v+1} K_v \quad (1)$$

where $x \geq 0$, $v > -1$ and $a > 0$, $\Gamma(\cdot)$ denotes the gamma function, and $K_v(\cdot)$ is the modified Bessel function of the second kind and order v , the K-distribution is completely specified by the shape parameter v and the scale parameter a , which subject estimation.

A. Moment Based Method for Parameters Estimation

In [13], the estimates of the parameters a and v can be obtained by estimating the moments of K-distribution, given by [3]

$$m_k = E[X^k] = \frac{\Gamma(0.5k+1)\Gamma(v+1+0.5k)}{\Gamma(v+1)} (2a)^k \quad (2)$$

By the sample moments

$$\hat{m}_k = \frac{1}{N} \sum_{i=1}^N x_i^k, \quad k \geq 0 \quad (3)$$

where $\{x_i; i = 1, \dots, N\}$ is a set of realizations of N statistically independent random variables $\{X_i; i = 1, \dots, N\}$, it is readily seen that one can estimate the parameters a and v using any two estimate of the moments given in (2).

One of the samples choices for the two moments, as suggested in [15], is the use of the sample mean and sample variance obtained from data. These quantities can be fitted to mean and variance of the K-distribution which are given by

$$E[X] = \frac{\Gamma(1.5)\Gamma(v+1.5)}{\Gamma(v+1)} (2a) \quad (4)$$

$$\text{Var}[X] = 4a^2(v+1) - E[X]^2$$

This approach is numerically inefficient in that the derivation of v or a from (4) involves solving a nonlinear equation.

In radar systems [16], the shape parameter v of the theoretical PDF has been estimated by the classical method of moments, which consists of equating the first- and second-order experimental moments with the corresponding theoretical moments. That is, the shape parameter has been estimated for each range cell by solving for v the equation

$$\frac{m_2}{m_1^2} = \frac{\hat{m}_2}{\hat{m}_1^2} \quad (5)$$

Using (2) and recalling that $\Gamma(v+1) = v \Gamma(v)$ and $\Gamma(1/2) = \sqrt{\pi}$, the estimate \hat{v} has been calculated by solving for v in the equation

$$\frac{4}{\pi} \frac{v \Gamma^2(v)}{\Gamma^2(v+1/2)} - \frac{\hat{m}_2}{\hat{m}_1^2} = 0 \quad (6)$$

Equation (6) does not yield a closed form expression for estimating the shape parameter, but numerical methods can be devised to solve for this equation.

The parameter a can be obtained by simply using one of the moments, for example, the first order moments of X , which yields

$$a = \frac{m_1 \Gamma(v+1)}{\sqrt{\pi} \Gamma(v+1.5)} \quad (7)$$

Estimates of v and a are then obtained from (6) and (7) after replacing the unknown moments by their estimates.

In homogeneous environment, the parameters are best stistical estimates. But when the reference cell contains interfering targets, the estimations of the shape and the scale parameters, calculated by (6) and (7) are far from the actual values. In order to avoid interfering targets in the reference cells, we propose a novel algorithm based on automatic censoring of radar interfering targets.

III. THE PROPOSED CENSORING ALGORITHM

The algorithm discussed here Fig. 1 does not need any prior information about the clutter parameters nor does it require neither the number of interfering targets. This latter is inspired by the technique of "bootstapping" [17]. The censoring procedure starts, first by creating a number B of resamples $\mathcal{X}_1^*, \dots, \mathcal{X}_B^*$. The resample $\mathcal{X}_b^* = [X_1^* X_2^* \dots X_M^*]^T$ is an unordered collection of M sample points drawn randomly from the reference cells $\mathcal{X} = [X_1 X_2 \dots X_N]$ with no replacement, where M is chosen to be inferior to N . Let this new sample be in the matrix $\mathcal{X}^* = [\mathcal{X}_1^*, \mathcal{X}_2^*, \dots, \mathcal{X}_B^*]$ with the size of $(M * B)$.

For each resamples \mathcal{X}_b^* , we calculate the mean value as follows:

$$\hat{\mu}_b^* = \frac{1}{M} \sum_b \mathcal{X}_b^* \quad (8)$$

where $b = 1, 2, \dots, B$.

The mean of all resamples \mathcal{X}_b^* , is then obtained as

$$E[\mathcal{X}^*] = [\hat{\mu}_1^*, \hat{\mu}_2^*, \dots, \hat{\mu}_b^*] \quad (9)$$

The description of the censoring algorithm in this level shows that a large time processing is required. To reduce this computational requirement, we propose a pseudo-random number generator is essential for valid application of the resamples technique. On the other hand, uses the architecture based on a parallel approach.

Note that in the case where M is large and the samples of reference cells are IID the distribution of $\hat{\mu}_b^*$ could be approximated by the Gaussian distribution stated by the Central Limit Theorem [18]. This theorem is not valid in applications where M is too small.

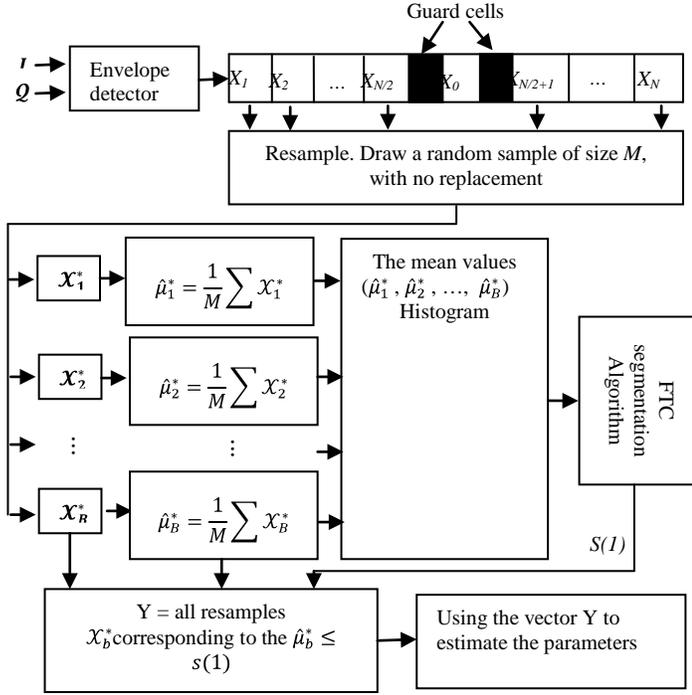


Fig. 1 Block diagram of the proposed algorithm

We remark that when the reference cells contains interfering targets, the histogram of $\hat{\mu}_b^*$ becomes multimodal, where the first mode represents the homogeneous environment and the number of others modes (second mode and upper) are equivalent to the number of interferences estimated (\hat{m}) in the reference cells is given by

$$\hat{m} = NM - 1 \quad (10)$$

where NM is a number of the modes in the histogram.

In this case, the proper segmentation of the histogram can be obtained by computing the appropriate thresholds that separate the modes in the histogram. Therefore, we have the nonparametric approach uses the FTC Segmentation Algorithm [12]. This is based on automatic detection of unimodal intervals in the histogram, which allow us to segment it.

To this effect, the homogenous population corresponding to the mean of the first mode in the histogram as indicated in the following equation

$$\hat{\mu}_b^* \leq s(1) \quad (11)$$

where $s(1)$ denotes the threshold of the first mode in the histogram, given by FTC Segmentation Algorithm.

For example, in Fig. 2 shows that the histogram of $\hat{\mu}_b^*$ for $N=36$, $M=16$, $B=5000$, $k=2$, interfering target-to-clutter ratio ($ICR=15\text{dB}$) and $a = 0.35$, $v = 1$ parameters of K -distribution.

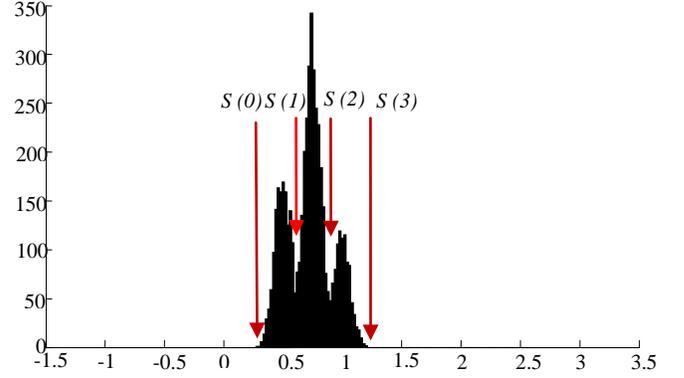


Fig. 2 Histogram of three modes related to $\hat{\mu}_b^*$

Then, we can make use of the mean $\hat{\mu}_b^*$ obtained by (11) to extract all the resamples X_b^* corresponding in (8) noted in the vector $Y=[X_b^*; \dots]^T$. Finally, the vector Y is used to estimate the unknown parameters by the classical MOM (first- and second-moments).

Notice that in a homogenous environment where ($m=0$), the FTC segmentation algorithm detect only the unimodal interval in the histogram of $\hat{\mu}_b^*$. Consequently, the samples obtained Y are equal to $X^*(:)^T$ with the length of ($M * B$).

A. Selection of Censoring Threshold

Remark that when we used the FTC segmentation algorithm in order to determine the appropriate thresholds to separate the modes in the histogram, the algorithm does not detect the small modes in lower ICR , which means that some interfering are masked, also the algorithm can detect small oscillation and considerable interfering targets. In order to solve this problem, the equation of the number of false alarm increasing (NFA_c), see [12] can be used and multiplied by the factor λ then the NFA_c becomes.

$$\frac{L_c(L_c+1)}{2} * \exp(\lambda * r * \left(-\frac{r}{B} \log \frac{r}{Bp} - \left(1 - \frac{r}{B}\right) \log \frac{1-r}{1-p}\right)) \quad (12)$$

In the same way, for the number of false alarm decreasing (NFA_d).

The factor λ is chosen empirically to detect the small modes and not to detect the small oscillation in modes contained in the considered that the interfering targets, after all calculus done we found that the optimum value of λ is equal to 0.36.

IV. RESULTS AND DISCUSSIONS

In this section, we evaluate the performances of the probability of censoring and the estimators of the parameter v . For this task, we normalized the power of the received clutter, so that the second order moment of the clutter process is unity. The normalization procedure leads to the following

relation between the two parameters [14].

$$a = \frac{1}{2\sqrt{v+1}} \quad (13)$$

In order to generate independent and identically distributed K random variants with parameters a and v , we need to use $c = \sqrt{\tau}x$, $c(i) = \sqrt{\tau(i)}x(i)$, where $x \in CN(0,1)$ complex Gaussian models the speckle and τ models the texture. When τ is modeled as a Gamma distributed random variable with shape parameter v and scale parameter a/v , that is $\tau \in \text{Gamma}(v, \frac{a}{v})$, the amplitude of the clutter sample $|c(i)| = \sqrt{\tau}|x(i)|$ belongs to the K-model [16].

We assume in our evaluation that the environments may contain unknown number (m) and position of interfering targets and without any prior knowledge about the clutter parameters. We also assume that the interfering targets are SWII model and have the same ICR . Monte Carlo simulations generate data for different values of the shape parameter in various interfering targets ($m=0, m=1, \text{ and } m=3$) with $ICR=20\text{dB}$ for $M=16$ and $B=5000$.

Fig. 3 shows the probability of censoring for $N=32, M=16$ and 3 interferences ($m=3$) with different values of ICR . λ has been fixed to 0.36. Note that the algorithm has the capability to determine the exact number of interferences $\hat{m} = m$ (\hat{m} is the estimated value of m) with a probability of 88.8% at $ICR=25\text{dB}$, 66.9% at $ICR=20\text{dB}$ and 32.9% at $ICR=15\text{dB}$. It is readily seen that an increase in ICR , yields an increase in probability of censoring.

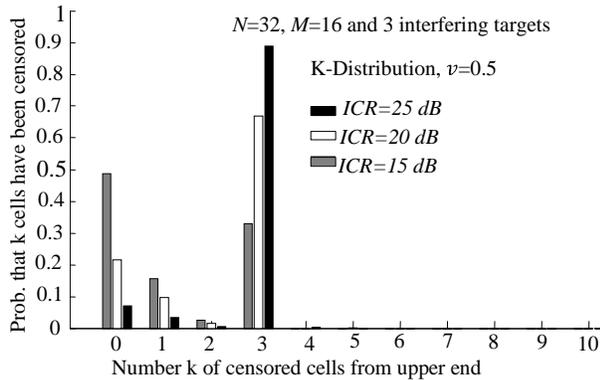


Fig. 3 Probability of censoring in multiple target situations for different values of ICR

Table I gives the corresponding parameter v of K-distribution by the classical MOM, with different shape parameters for multiple interfering targets through Monte-Carlo simulations.

TABLE I

THE ESTIMATION OF THE PARAMETERS OF K-DISTRIBUTION WITH DIFFERENT SHAPE PARAMETERS FOR VARIOUS INTERFERENCES ($M=0, 1$ AND 3) AT $ICR=25\text{dB}, N=32, M=16$ AND $\lambda=0.36$.

shape parameter	$m=0$		$m=1$		$m=3$	
	bias	var[\hat{v}]	bias	var[\hat{v}]	bias	var[\hat{v}]
0.1	0.08	0.0035	0.0805	0.0038	0.0818	0.003
0.5	0.143	0.1448	0.1137	0.1502	0.1179	0.177
1	0.398	0.8077	0.3441	0.8541	0.2936	0.924
1.5	0.2566	1.5654	0.1621	1.4724	0.1412	1.736

Note: the averages are computed over 2000 independent trials in each case.

From Table I, it can be seen that the variance of the estimates is lower for small values of real parameters and the estimates possess lower variance when the interfering targets decrease in different shape parameters, also we can remark that the proposed method has not great bias of estimates.

V. CONCLUSION

In this work, we have considered the problem of automatic censoring for the K-distribution in the presence of an unknown number of interfering targets and unknown clutter parameters. A novel algorithm has been proposed in this work based on automatic censoring techniques of radar interfering targets in K-distribution; this algorithm uses an efficient technique for discriminating between homogeneous and non-homogeneous environments. The homogeneous population is then used to estimate the unknown parameters of the K-distribution by the classical MOM. The effectiveness of the proposed AC algorithm has been assessed by computing the probability of censoring and the accuracy of the estimation parameters for multiple target situations. The simulation results showed that in interfering targets environments, AC can censor interfering targets adaptively and effectively. This later showed that the best statistical estimates where the shape parameter is very small the case such as in radar applications for multiple target situations. The algorithm is also guaranteed a good quality of the estimates of the parameter in multiple target situations. Future work is to search an adaptive threshold of detection for K-distribution.

REFERENCES

- [1] K. D. Ward, "Compound representation of high resolution sea clutter," *Electronics Letters*, vol. 17, pp. 561-563, 1981.
- [2] A.A. Douglas, and P.L. Anthony, "Reliable Methods for Estimating the K-distribution Shape Parameters," *IEEE Journal of Oceanic Engineering*, vol. 35, pp. 288-302, April. 2010.
- [3] R.S. Raghavan, "A method for estimating parameters of K-distributed clutter," *IEEE Trans. Aerosp Electron Syst*, vol. 27, pp. 238-246, Mar. 1991.
- [4] I. R. Joughin, D. B. Percival, and D. P. Winebrenner, "Maximum likelihood estimation of K distribution parameters for SAR data," *IEEE Trans. Geosci Remote Sens*, vol. 31, pp. 989-999, Sep. 1993.
- [5] P. Lombardo, and C. J. Oliver, "Estimation of texture parameters in K-distributed clutter," *Proc Inst Electr Eng.—Radar Sonar Navigat*, vol. 141, pp. 196-204, Aug. 1994.
- [6] A. Farrouki, and M. Barkat, "Automatic Censoring CFAR detector based on ordered data variability for non-homogeneous environments," *IEE Proc. Radar Sonar Navig*, vol. 152, pp. 43-51, Feb. 2005.

- [7] A. Zaimbashi, Y. Norouzi, "Automatic dual censoring cell averaging CFAR detector in nonhomogenous environments," *EURASIP J. Signal Processing*, vol. 88, pp. 2611-2621, Nov. 2008.
- [8] N. Boudemagh, Z. Hammoudi, "Automatic censoring CFAR detector for heterogeneous environments," *Int. J. Electron. Commun. (AEÜ)*, Vol. 68, pp. 1253-1260, 2014.
- [9] M.N. Almarshad, S.A. Ashbeili, and M. Barkat, "A forward automatic censored cell-averaging detector for multiple target situations in log-normal clutter," *PWASET*, vol. 17, pp. 1307-6884, Dec. 2006.
- [10] S. Chabbi, T. Laroussi, M. Barkat, "Performance analysis of dual automatic censoring and detection in heterogeneous Weibull clutter: A comparison through extensive simulations," *Signal Processing*, vol. 93, pp. 2879-2893, 2013.
- [11] F. Soltani and M. Barkat, "CFAR binary integration detection in nonhomogeneous partially correlated clutter", *IEE Proc. Radar Sonar navig.*, Vol.144, No.5, pp.293-300, October 1997.
- [12] J. Delon, A.Desolneux, J.L.Lisani, and A.B. Petro, "A nonparametric Approach for Histogram Segmentation," *IEEE Trans.Image Process*, vol. 16, pp. 253-261, Jan. 2007.
- [13] C.W. Helstrom, *Elements of Signal Detection and Estimation*, prentice hall, 1995.
- [14] D.R. Iskander, and A.M. Zoubir, "Estimation of the parameters of the K-distribution using higher order and fractional moments," *IEEE Trans. Aerosp Electron Syst*, vol. 35, pp. 1453-1456, Oct. 1999.
- [15] D. Blacknell, "Comparison of parameter estimators for K-distribution," *IEE proc. Radar Sonar Navigation*, vol. 141, pp. 45-52, 1994.
- [16] P. Stinco, F. Gini, and M.Rangaswamy, "Impact of Sea Clutter Nonstationarity on Disturbance Covariance Matrix Estimation and CFAR Detector Performance," *IEEE Trans. Aerosp Electron Syst*, vol. 46, pp. 1502-1513, July. 2010.
- [17] A.M Zoubir, and D.R Iskander, *Bootstrap Techniques for Signal Processing*, Cambridge, U.K, 2004.
- [18] E.B. Manoukian, *Modern Concepts and Theorems of Mathematical Statistics*, New York, USA: Springer-Verlag, 1986.

Fuzzy method for suppressing of different noises in color videos

Volodymyr Ponomaryov

Abstract— A novel approach in denoising of color videos corrupted by impulsive and additive noises is proposed. In difference with existing methods, novel method employs designed fuzzy rules selecting high similarity pixels in the vicinity of the central one using the correlation in the RGB channels and in the consecutive frames of a video for better preservation of the features via adjusting the possible local motions, processing separately the areas with different texture behaviors (e.g., smooth regions, edges, and fine details). Numerous simulation experiments have demonstrated the superiority of novel filters presenting better values for objective criteria (PSNR, MAE, NCD, SSIM) as well as in increasing the perceptual vision.

Keywords—Video, Fuzzy logic, Denoising, Similarity.

I. INTRODUCTION

THE presence of noise produces deficiencies during acquisition, broadcast or storage of the color images and videos [1]-[7]. Noise affects not only the performance of an image in a specific problem but also its perceived quality. Therefore, it is a priority task to filter each image or frame of a video prior to other processing in following stages [1]-[4], reducing the amount of noisy pixels. A principal problem here consists of a design of a noise reduction technique while image content (edges, fine features, etc.) should be preserved. Numerous techniques have been proposed that are mainly based on order statistics technique, on fuzzy logic theory, on sparse representation, etc. [8]-[19]. In color video filtering, employing existing interchannel and temporal correlation between the neighboring frames and processing them together it is possible to obtain sufficiently improved performance in comparison with case 2D frame filtering. The principal obstacle encountered when two or more frames are processed together for noise removal is the possible existence of local motions between different frames, which usually introduce motion blur and ghosting artifacts [14]-[17], [21]-[23]. Modern theoretical approaches in denoising of different noises are principally based on a possibility to gather more samples for similar patches in an image. Then, the methods use sophisticated statistical methods, which depend on image/noise model. The principal problem here is how to measure and employ the similarity of group of objects in a color image [20]-[28]. Proposed in this paper approach exploits the similar ideas using fuzzy set type filtering in searching similar patches that permit to gather more samples

This work was supported by the Instituto Politécnico Nacional de Mexico and Consejo Nacional de Ciencia y Tecnología de Mexico.

Volodymyr Ponomaryov is with the Instituto Politécnico Nacional, ESIME-Culhuacan, Av. Santa Ana 1000, San Fco. Culhuacan, 04430, Mexico-city, Mexico (ph: 525556562058; fax: 525556562058 e-mail: vponomar@ipn.mx)

for processing together color channels for a video frame; following, more samples should be found gathering neighboring frames of a video where the local motions in different frames should be adjusted. The proposed approach in difference to other state-of-the-art approaches employs the RGB channels data and fuzzy logic description of semantic properties of image features via designed Fuzzy Rules in all filtering steps, processing several pixel gradients together in neighboring frames.

II. FUZZY APPROACH IN DENOISING OF COLOR VIDEOS

In current paper, we present two novel techniques based on fuzzy logic approach in denoising: for impulsive noise suppression FMINS (*fuzzy multichannel impulse noise suppression*) filter, and for additive noise suppression FMANS (*fuzzy multichannel additive noise suppression*) filter.

A. Impulsive Noise Suppression

The designed method in denoising of impulsive noise is divided in three steps [16], [21]. In the first step, several gradient vector values for a basic gradient and four related gradients are computed. Each pixel is characterized by a level where it can be considered as *noise-free* and a level where it can be considered as *noisy*; the output of this step is denoted as $E(i,j)_1$. In the second step, the noise detection and filtering is based on mutual processing of three *RGB* color channels in a current frame. The output of the second stage is denoted as $E(i,j)_2$. In the final third step, the filtering procedure using spatial and temporal processing in two neighboring frames is performed where the remaining noisy pixels should be removed, guaranteeing edges and fine detail preservation, forming output filtering result $E(i,j)_3$. Details of FMINS framework are presented in the block diagram of Fig.1.

During filtering, a 3x3 sliding window located into a bigger 5x5 window in the novel framework is employed in present approach, applying the gradient values for neighboring pixels in eight different directions $\gamma = (NW, N, NE, E, SE, S, SW, W)$ with respect to a central pixel (see Fig.2).

$$\begin{aligned} \nabla_{(1,1)}^\beta E(i,j) &= \nabla_{SE(B)}^\beta; \nabla_{(0,2)}^\beta E(i-1,j+1) = \nabla_{R1,SE}^\beta, \\ \nabla_{(2,0)}^\beta E(i+1,j-1) &= \nabla_{R2,SE}^\beta, \quad \nabla_{(-1,1)}^\beta E(i-1,j+1) = \nabla_{R3,SE}^\beta, \\ \nabla_{(1,-1)}^\beta E(i+1,j-1) &= \nabla_{R4,SE}^\beta \end{aligned} \quad (1)$$

Two hypotheses are resolved: the central pixel is a *noisy* or it is a *free-noise* pixel. The *LARGE* and *SMALL* fuzzy sets are

introduced with an objective to estimate the noise contamination employing the Gaussian membership functions [6, 8] for membership degrees of gradient values [16]:

$$\rho(\nabla_{\gamma}^{\beta}, LARGE) = \begin{cases} 1, & \nabla_{\gamma}^{\beta} > \nabla_1 \\ \exp\{-[(\nabla_{\gamma}^{\beta} - \nabla_1)^2 / 2\sigma^2]\}, & otherwise \end{cases} \quad (2)$$

$$\rho(\nabla_{\gamma}^{\beta}, SMALL) = \begin{cases} 1, & \nabla_{\gamma}^{\beta} < \nabla_1 \\ \exp\{-[(\nabla_{\gamma}^{\beta} - \nabla_2)^2 / 2\sigma^2]\}, & otherwise \end{cases} \quad (3)$$

The values of the parameters used in (2) and (3) were selected according to optimal values of *PSNR* and *MAE* criteria during several simulation experiments for different video sequences. The found values of these parameters are: $\nabla_1 = 60$, $\nabla_2 = 9$, $\sigma^2 = 1000$; for interchannel RGB filtering $\nabla_{2,inter} = 9$, $\sigma_{inter}^2 = 750$; and in the case of mutual frames filtering: $\nabla_1 = 0.1$, $\nabla_2 = 0.01$, $\sigma^2 = 0.1$ [16].

To resolve the hypothesis: a central pixel is *noisy* or *noise-free* that belongs to image features, several fuzzy rules are proposed. Table 1 exposes the designed fuzzy rules. *Fuzzy Rule 1-1* defines the fuzzy gradient value ∇_{γ}^{nF} that belongs to fuzzy set *LARGE* for γ direction. A color component pixel is considered as *noisy* pixel if its basic gradient value is similar to its related gradients R_3 and R_4 , and differs from related gradients R_1 and R_2 (see Fig. 2). *Fuzzy Rule 1-2* presents the noisy factor r_{β} that gathers eight fuzzy gradient-directional values presented in the *Fuzzy Rule 1-1*. The noisy factor r_{β} is a measure to distinguish between a noisy pixel and a noise-free one. This value determines the level of noise presence in the processed sample in the fuzzy set *LARGE* indicating that this central pixel is corrupted. If a central pixel in a sliding window is considered as noisy one, the special procedure of ranking for all pixel values in ascending order according to its weights ρ_{γ}^{β} is used. Interchannel processing procedures that use the existing correlation between the *R*, *G* and *B* frame components are explained in the *Fuzzy Rules 2-1*, *2-2* and *2-3* (Table 1). *Fuzzy Rule 2-1* defines the condition when the *R* component is *noise-free*. In final spatial-temporal stage, the remaining noisy pixels are now processed gathering data from two neighboring frames. The absolute difference values between (*t*) and (*t-1*) frames $\delta E_{(k,l)}^{\beta}$ are calculated, forming the error frame for time (*t*). The remaining noisy pixels in this step are now processed inside a $5 \times 5 \times 2$ sliding window, gathering two neighboring frames: $E^{t,\beta}(i, j)$, $E^{t-1,\beta}(i, j)$ calculating the difference values between (*t*) and (*t-1*) frames:

$$\delta E_{(k,l)}^{\beta} = \left| E^{t,\beta}(i+k, j+l) - E^{t-1,\beta}(i+k+k_1, j+l+l_1) \right|, \quad (4)$$

$k, l \in (-3, -2, -1, 0, +1, +2, +3)$; $k_1, l_1 \in (-1, 0, -1)$; $\beta = (R, G, B)$

In this step, gradient values for frame difference $\nabla_{\gamma(B \text{ or } Ri)}^{\beta}$ for each of eight directions γ are used.

Additionally, the best match using criterion *MAD* should be found between the central pixel in current frame and pixels at the vicinity of central one in previous frame, increasing size of common sample that consists of pixels from (*t*) and (*t-*

l) frames. Finally, at postprocessing step, the complex (edges, fine details) and plane regions are processed separately (*Fuzzy Rules: 3-1 to 3-4*). *Fuzzy Rule 3-1* that employs the absolute difference gradient values $\nabla \delta_{\gamma}^{\beta}$ determines the first fuzzy gradient difference $(\nabla_{\gamma}^{\beta F})_I$ for a central pixel in respect to its neighbors in a sliding window similar as it has been done in *Fuzzy Rule 1-1*. *Fuzzy Rule 3-2* determines the fuzzy gradient difference $(\nabla_{\gamma}^{\beta F})_{II}$ using the fuzzy gradient differences in the direction γ , distinguishing between homogeneous and non-homogeneous regions. *Fuzzy Rule 3-3* computes the noisy factor r_{β} that gathers the fuzzy gradient-directional values presented in the *Fuzzy Rule 3-1*. *Fuzzy Rule 3-4* introduces the factor η_{β} .

B. Additive Noise Suppression

The designed method in denoising of additive noise uses similar *fuzzy ideology* of FMINS framework as it can be seen in the block diagram in Fig.3. There is used a 5×5 sliding window into bigger 7×7 sliding window to compute the gradient values in eight directions for basic and six related gradients (see Fig.2). The gradient values are introduced for each direction $\gamma = \{N, E, S, W, NW, NE, SE, SW\}$, where (*i, j*) values are $\{-3, -2, -1, 0, 1, 2, 3\}$. For FMANS filter, the *Fuzzy Rules 1-1* and *1-2*, where Gaussian functions are employed for calculating membership degrees of fuzzy gradient values, are used. In following stage of first step, the noise detection and filtering is based on mutual interchannel processing using RGB color representation in a current frame forming output of this stage. *Fuzzy Rule 1-3* defines the condition when the *R* component can be estimated “*noise-free*”. *Fuzzy Rules 1-4* and *1-5* are employed to compute the weights for *noise-free* pixels as well as for *noisy* pixels. In next step, the filtering procedure is applied in spatio-temporal processing for two neighboring frames forming output filtering result (*Fuzzy Rules: 2-1 to 2-3*). The remaining noisy pixels in this step are now processed inside a $7 \times 7 \times 2$ sliding window, gathering two neighboring frames: $E^{t,\beta}(i, j)$, $E^{t-1,\beta}(i, j)$ calculating the difference values between (*t*) and (*t-1*) frames (see eq. (4)). The best match using criterion *MAD* can be found between the central pixel in current frame and pixels at the vicinity of previous frame, increasing size of common sample that consists of pixels from (*t*) and (*t-1*) frames. Finally, at postprocessing step, the complex (edges, fine details) and plane regions are processed separately (*Fuzzy Rules: 3-1 to 3-3*) as presented in Table 2. Two variants of filtering are presented in Block diagram: FMANS_2 and FMANS_H, which uses hybrid processing connecting the “*fuzzy ideology*” of the FMANS technique and multiscale Wiener DCT-based filtering [28] is also performed that increases denoising ability as simulation results show.

III. SIMULATION RESULTS AND PERFORMANCE EVALUATION

The color videos *Flowers*, *Stefan*, *Foreman*, *Tennis* in the CIF format (352x288) and *Carphone*, *Grandma*, *Miss America* and *Saleman* in the QCIF format (176x144 pixels, RGB, 24 bits) [29] were used to evaluate the promising 3D fuzzy algorithms in wide range of impulsive noise intensity (0% to 20%). As shown in recently published articles the FRINR_Seq and 3D FD filtering techniques outperform all other existing state-of-the-art techniques in denoising impulsive noise, so comparing novel filter we justify in correct way the performance of novel FMINS 3D framework. The filtered frames were evaluated according to PSNR (*Peak Signal-to-Noise Ratio*) that describes the noise suppression ability for an algorithm; the MAE (*Mean Absolute Error*) that measures the edge preservation ability [1-3], [5]. Additionally, another metric NCD (*Normalized Color Difference*, in the L^*u^*v color space) is used to measure color preservation properties of filtering results [2], [3]. Recently introduced SSIM (*Similarity Structural Index Measure*) [30] that matches better with human subjectivity is applied to characterize the performance of a chosen algorithm. These objective criteria and subjective perception via human vision are used to characterize the performance of the FMINS 3D filter averaging per 100 frames against mentioned techniques FRINR_Seq [17] and 3D-FD [15], [20], exposing the better values for *Foreman*, *Stefan*, *Grandma*, and *Carphone* color video sequences, guaranteeing its robustness. These results are exposed in Fig. 4a - 4c.

Table 5 presents the average per 100 frames PSNR, MAE, NCD and SSIM values for the proposed FMINS 3D framework against other better techniques FRINR_Seq and 3D FD, exposing the better values for *MA*, *SM*, *F*, and *S* video sequences. The best performance is realized by novel method according to all four objective criteria in wide range of noise intensity. Additionally, all three algorithms applied in filtering of the video sequences with varying the random impulse noise levels in range from 0% to 20%. The results of these experiments in terms of PSNR, MAE in different frames show that novel framework outperforms other better mentioned algorithms (Fig.4). The proposed fuzzy approach combines sufficiently good detail preservation to good noise removal and appear outperforms other compared filters in wide range of noise intensity.

Similar simulations on mentioned color videos have been performed in denoising of additive noise contamination. The proposed algorithm FMANS and other better techniques were evaluated in terms of the PSNR, MAE, NCD and SSIM criteria applied to different videos, presenting averaging values per 50 frames of each a video. As one can see in tables 5 and 6 the proposed (FMANS_2 and FMANS_H) techniques outperform other state-of-the-art techniques according to all criteria. The performance of the NLM is slightly better than for the other comparative filters. The FMANS_2 (window 7x7) and FMANS_H (Table 3) yield better results in comparison with the NLM, where the FMANS_H provides the best performance compared with all other denoising techniques. The same conclusion can be done analyzing behavior of criteria on different frames (see Fig.5).

Figure 6 shows the filtering frames and their error images for different filters in case of color video *Stefan*. In the 50th frame of the *Stefan* video, one can observe the better preservation of details in the field and letters located on the front wall compared with the other competing methods.

Comparing with the related state-of-the-art methods, the principal contributions of the current fuzzy approach are as follow:

- Developed fuzzy rules that permit selection of high similarity pixels in the vicinity of the central pixel employing the correlation in the RGB channels and in the consecutive frames of a video for better preservation of the features via adjusting the possible local motions.
- Separating and processing differently the areas with different texture behaviors (e.g., smooth regions, edges, and fine details).
- Hybrid denoising scheme for additive noise suppression that consists of combining the designed fuzzy framework and the multiscale Wiener filter at the final denoising stage.
- Preservation of the chromaticity properties of the image (such as color balance), avoiding unexpected color combinations after filtering operation.
- Demonstration of superiority in achieved better PSNR, MAE, NCD and SSIM values and in increasing the perceptual quality on the textured, and plain areas of the images.

Finally, while the proposed fuzzy approach is justified in the reduction of additive Gaussian and impulsive noises, nevertheless the fuzzy ideology can be generalized to other kind of noise because the filtering procedures adapt to the characteristics of an image without prior information, and it is also not necessary to have previous information about the type of noise that corrupts the image.

IV. CONCLUSION

The designed filters FMINS and FMANS are based on fuzzy logic approach using the interchannel correlations, matching possible motions in neighboring frames, forming the most similar pixels in different spatial areas of a current and neighboring frames of a video, finally, demonstrating superiority in comparison with better existing fuzzy and non fuzzy techniques in denoising of impulsive (FMINS) and additive (FMANS) noises.

These filters excellently suppress the noises in color videos, preserving edges, fine features, color properties, justifying their efficiency in PSNR, MAE, NCD and SSIM metrics and in subjective perception via human vision system.

Future work will be focused in increasing the ability of current fuzzy proposal in noise suppression for other kind of noises as well as in speed via parallel processing implementation on GPU hardware.

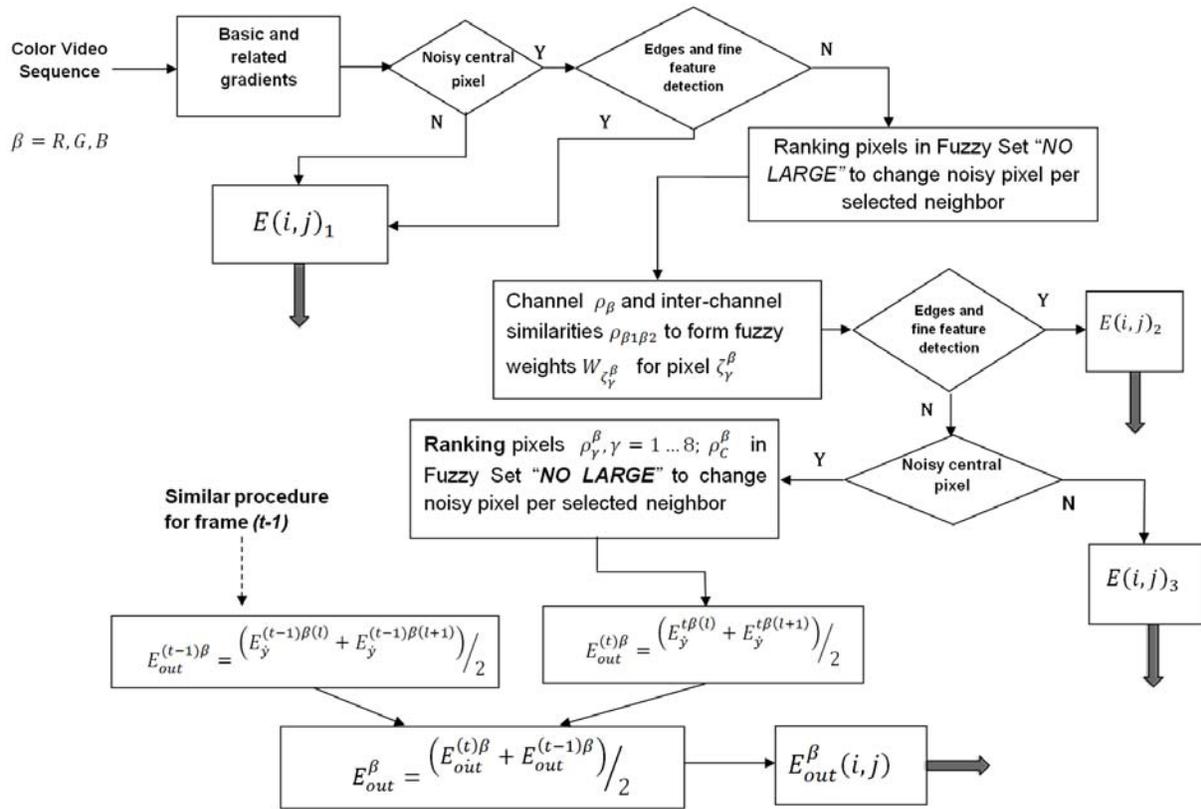


Fig. 1 Block diagram of FMINS denoising filter

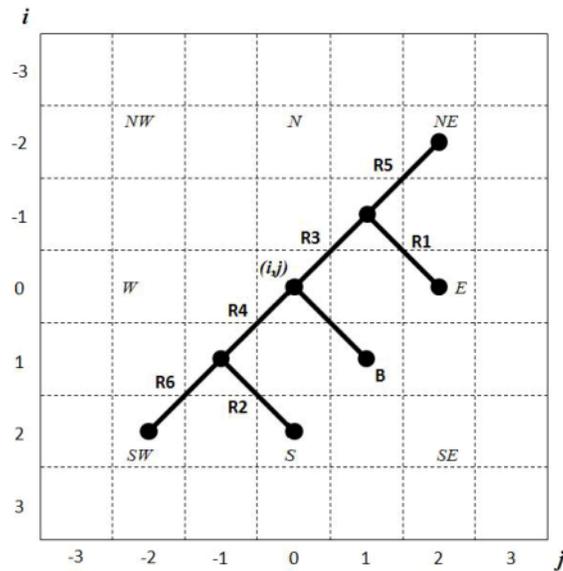


Fig.2 Basic B and several related (R_1 to R_6) gradients applied in sliding window

<p>FR 1-1: Defining fuzzy gradient values $\nabla_{\gamma}^{\beta F}$ into set LARGE</p>	<p>IF ($\nabla_{\gamma B}^{\beta}$ is L AND $\nabla_{\gamma R1}^{\beta}$ is S AND $\nabla_{\gamma R2}^{\beta}$ is S AND $\nabla_{\gamma R3}^{\beta}$ is L AND $\nabla_{\gamma R4}^{\beta}$ is L) THEN Fuzzy Gradient $\nabla_{\gamma}^{\beta F}$ is LARGE, $(A \text{ AND } B) = A \bullet B$ $(A \text{ OR } B) = A + B - A \bullet B$</p>
<p>FR 1-2: Defines fuzzy noisy factor r_{β}</p>	<p>IF MAX (∇_N^{β} is L, MAX (∇_S^{β} is L, MAX (∇_E^{β} is L, MAX (∇_W^{β} is L, MAX (∇_{SW}^{β} is L, MAX (∇_{NE}^{β} is L, MAX (∇_{NW}^{β} is L, ∇_{SE}^{β} is L)))))) THEN r_{β} is LARGE.</p>
<p>FR 2-1: Membership degree for R component ζ_C^R in fuzzy set “noise free” FR 2-2: Defining the weight for R component ζ_C^R</p>	<p>IF (μ^R is L AND μ^{RG} is L AND μ^G is L) OR (μ^R is L AND μ^{RB} is L AND μ^B is L) THEN the noise-free degree of ζ_C^R is LARGE. IF $N(\zeta_C^R)$ is LARGE THEN $W(\zeta_C^R)$ is LARGE</p>
<p>FR 2-3: Defining the weight $W(\zeta_{\gamma}^R)$ for the neighbor of R component ζ_{γ}^R</p>	<p>IF ($N(\zeta_{\gamma}^R)$ is not L AND $W(\zeta_{\gamma}^R)$ is L AND $\mu(\nabla_{\zeta_{\gamma}^G})$ is L AND $W(\zeta_{\gamma}^G)$ is L) OR ($N(\zeta_{\gamma}^R)$ is not L AND $W(\zeta_{\gamma}^R)$ is L AND $\mu(\nabla_{\zeta_{\gamma}^B})$ is L AND $W(\zeta_{\gamma}^B)$ is L) THEN $W(\zeta_{\gamma}^R)$ is LARGE.</p>
<p>FR 3-1: Determines the first fuzzy gradient difference $(\nabla_{\gamma}^{\beta F})_I$ to characterize confidence “movement-noise” FR 3-2: Determines the fuzzy gradient difference $(\nabla_{\gamma}^{\beta F})_{II}$ FR 3-3 computes fuzzy factor fuzzy noisy factor r_{β} for interframe processing</p>	<p>Repeat FR 1-1 changing $\nabla_{\gamma B}^{\beta}$, $\nabla_{\gamma Ri, i=1,2,3,4}^{\beta}$ per $\nabla_{\gamma B}^{\beta}$, $\nabla_{\gamma Ri}^{\beta}$, accordingly. IF ($\nabla_{\gamma B}^{\beta}$ is S AND $\nabla_{\gamma R1}^{\beta}$ is S AND $\nabla_{\gamma R2}^{\beta}$ is S) THEN $(\nabla_{\gamma}^{\beta F})_{II}$ is SMALL Repeat FR 1-2 changing ∇_{γ}^{β} per $(\nabla_{\gamma}^{\beta})_I$</p>
<p>FR 3-4: Determines the fuzzy factor η_{β} defining the confidence “no movement-no noise” in interframe processing</p>	<p>IF MAX ($(\nabla_N^{\beta F})_{II}$ is S, MAX ($(\nabla_S^{\beta F})_{II}$ is S, MAX ($(\nabla_E^{\beta F})_{II}$ is S, MAX ($(\nabla_W^{\beta F})_{II}$ is S, MAX ($(\nabla_{SW}^{\beta F})_{II}$ is S, MAX ($(\nabla_{NE}^{\beta F})_{II}$ is S, MAX ($(\nabla_{NW}^{\beta F})_{II}$ is S, $(\nabla_{SE}^{\beta F})_{II}$ is S)))))) THEN η_{β} is SMALL.</p>

Table.1 Fuzzy rules used in the FMINS filter

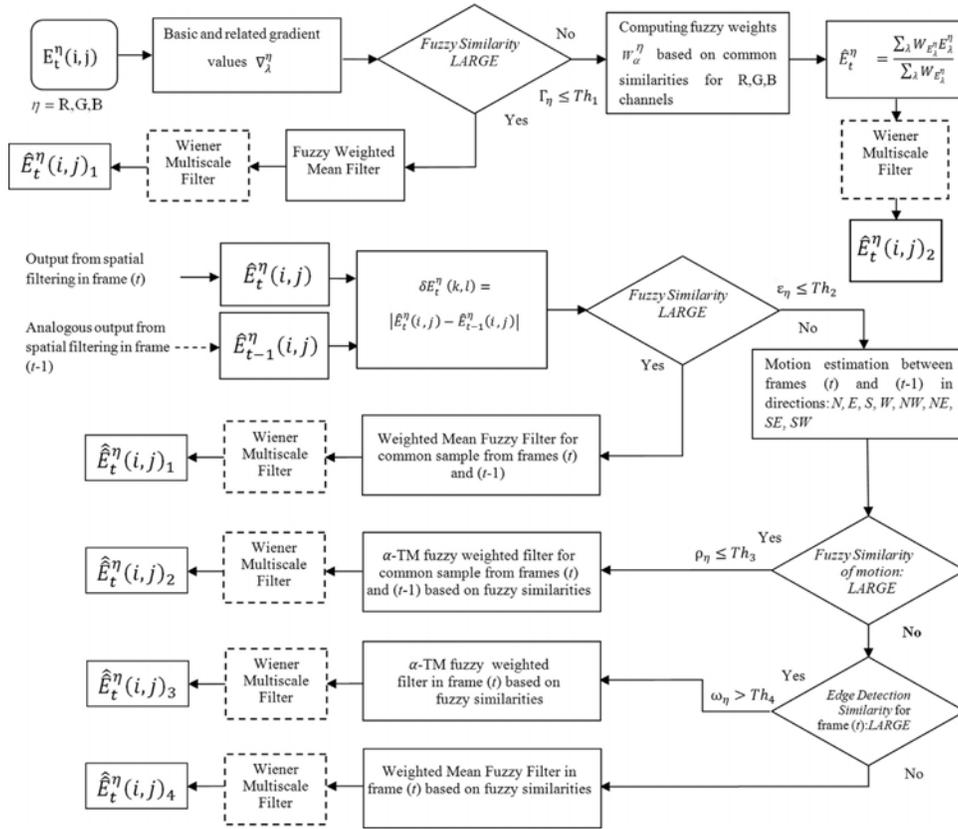


Fig.3 Block diagram of FMANS framework

<p>Fuzzy Rule 1 – 1. Defining the fuzzy gradient value $\nabla_{\lambda}^{\eta F}$ into the fuzzy similarity – set LARGE:</p>	<p>IF $((\nabla_{\lambda}^{\eta B}$ is LARGE AND $\nabla_{\lambda}^{\eta R1}$ is LARGE) OR $(\nabla_{\lambda}^{\eta B}$ is LARGE AND $\nabla_{\lambda}^{\eta R2}$ is LARGE)) AND $((\nabla_{\lambda}^{\eta B}$ is LARGE AND $\nabla_{\lambda}^{\eta R3}$ is LARGE) OR $(\nabla_{\lambda}^{\eta B}$ is LARGE AND $\nabla_{\lambda}^{\eta R4}$ is LARGE)) AND $((\nabla_{\lambda}^{\eta B}$ is LARGE AND $\nabla_{\lambda}^{\eta R5}$ is LARGE) OR $(\nabla_{\lambda}^{\eta B}$ is LARGE AND $\nabla_{\lambda}^{\eta R6}$ is LARGE)), THEN the fuzzy gradient value $\nabla_{\lambda}^{\eta F}$ is LARGE.</p>
<p>Fuzzy Rule 1 – 2. Defining the fuzzy noisy factor Γ_{η}:</p>	<p>IF $\text{MAX}((\nabla_{\lambda}^{\eta N})$ is LARGE, $\text{MAX}((\nabla_{\lambda}^{\eta S})$ is LARGE, $\text{MAX}((\nabla_{\lambda}^{\eta W})$ is LARGE, $\text{MAX}((\nabla_{\lambda}^{\eta SE})$ is LARGE, $\text{MAX}((\nabla_{\lambda}^{\eta SW})$ is LARGE, $\text{MAX}((\nabla_{\lambda}^{\eta NE})$ is LARGE, $\text{MAX}((\nabla_{\lambda}^{\eta NW})$ is LARGE, $\text{MAX}((\nabla_{\lambda}^{\eta SE})$ is LARGE))))))))) , THEN the noisy factor Γ_{η} is LARGE.</p>
<p>Fuzzy Rule 1 – 3. Defining the membership degrees $NE_{E_c^R}$ for the red component E_c^R in the fuzzy set "noise free":</p>	<p>IF $(\tau^R$ is LARGE AND τ^{RG} is LARGE AND τ^G is LARGE) OR $(\tau^R$ is LARGE AND τ^{RB} is LARGE AND τ^B is LARGE), THEN the noise – free degree of E_c^R is LARGE, where the conjunction $(A \text{ AND } B) = A \cdot B$ and the disjunction $(A) \text{ OR } (B) = A + B - A \cdot B$</p>
<p>Fuzzy Rule 1 – 4. Defining the weight WE_c^R for the red component E_c^R:</p>	<p>IF $(NE_{E_c^R}$ is LARGE) THEN WE_c^R is LARGE.</p>
<p>Fuzzy Rule 1 – 5. Defining the weight WE_{λ}^R for the neighbor of the red component E_{λ}^R:</p>	<p>IF $(NE_{E_c^R}$ is NO LARGE AND $NE_{E_{\lambda}^R}$ is LARGE AND $\tau(\Delta E_{\lambda}^R)$ is LARGE AND $NE_{E_{\lambda}^R}$ is LARGE) OR $(NE_{E_c^R}$ is NO LARGE AND $NE_{E_{\lambda}^R}$ is LARGE AND $\tau(\Delta E_{\lambda}^B)$ is LARGE AND $NE_{E_{\lambda}^R}$ is LARGE), THEN WE_{λ}^R is LARGE.</p>
<p>Fuzzy Rule 2 – 1. Defining the vectorial fuzzy gradient value $\nabla \delta E_{\lambda}^{\eta F}$ into the fuzzy –similarity set LARGE:</p>	<p>IF $(\nabla \delta E_{\lambda}^{\eta B}$ is LARGE AND $\nabla \delta E_{\lambda}^{\eta R1}$ is LARGE) OR $(\nabla \delta E_{\lambda}^{\eta B}$ is LARGE AND $\nabla \delta E_{\lambda}^{\eta R2}$ is LARGE) AND $(\nabla \delta E_{\lambda}^{\eta B}$ is LARGE AND $\nabla \delta E_{\lambda}^{\eta R3}$ is LARGE) OR $(\nabla \delta E_{\lambda}^{\eta B}$ is LARGE AND $\nabla \delta E_{\lambda}^{\eta R4}$ is LARGE) AND $(\nabla \delta E_{\lambda}^{\eta B}$ is LARGE AND $\nabla \delta E_{\lambda}^{\eta R5}$ is LARGE) OR $(\nabla \delta E_{\lambda}^{\eta B}$ is LARGE AND $\nabla \delta E_{\lambda}^{\eta R6}$ is LARGE), THEN the fuzzy similarity value $\nabla \delta E_{\lambda}^{\eta F}$ is LARGE.</p>
<p>Fuzzy Rule 2 – 2. Defining the fuzzy noisy factor ϵ_{η}:</p>	<p>IF $(\text{MAX}((\nabla \delta_{\lambda}^{\eta N})$ is LARGE, $\text{MAX}((\nabla \delta_{\lambda}^{\eta S})$ is LARGE, $\text{MAX}((\nabla \delta_{\lambda}^{\eta W})$ is LARGE, $\text{MAX}((\nabla \delta_{\lambda}^{\eta SE})$ is LARGE, $\text{MAX}((\nabla \delta_{\lambda}^{\eta SW})$ is LARGE, $\text{MAX}((\nabla \delta_{\lambda}^{\eta NE})$ is LARGE, $\text{MAX}((\nabla \delta_{\lambda}^{\eta NW})$ is LARGE, $\nabla \delta_{\lambda}^{\eta SE}$ is LARGE))))))))) , THEN the noisy factor ϵ_{η} is LARGE.</p>
<p>Fuzzy Rule 2 – 3. Defining the weights for the Alfa – TM filter WE_{λ}^R in the case of motion for the central pixel located in (t) frame:</p>	<p>IF $(NE_{E_c}$ is LARGE), THEN WE_{λ} is LARGE.</p>
<p>Fuzzy Rule 3 – 1. Defining the value of ∇_{λ}^{NF} for edges detection into the edge detection similarity – fuzzy set LARGE:</p>	<p>IF $(\nabla_{\lambda}^{\eta B}$ is NO LARGE AND $\nabla_{\lambda}^{\eta R1}$ is NO LARGE) OR $(\nabla_{\lambda}^{\eta B}$ is NO LARGE AND $\nabla_{\lambda}^{\eta R2}$ is NO LARGE) AND $(\nabla_{\lambda}^{\eta B}$ is LARGE AND $\nabla_{\lambda}^{\eta R3}$ is LARGE) OR $(\nabla_{\lambda}^{\eta B}$ is LARGE AND $\nabla_{\lambda}^{\eta R4}$ is LARGE) AND $(\nabla_{\lambda}^{\eta B}$ is LARGE AND $\nabla_{\lambda}^{\eta R5}$ is LARGE) OR $(\nabla_{\lambda}^{\eta B}$ is LARGE AND $\nabla_{\lambda}^{\eta R6}$ is LARGE), THEN the fuzzy edge detection similarity gradient value ∇_{λ}^{NF} is LARGE.</p>
<p>Fuzzy Rule 3 – 2. Defining the weights for the Alfa – TM filter WE_{λ}^R in the case of edge detection for the red component E_{λ}^R in the (t_{th}) frame:</p>	<p>IF $(NE_{E_c^R}$ is LARGE), THEN WE_c^R is LARGE.</p>
<p>Fuzzy Rule 3 – 3. Defining the weights for the mean filter WE_c^R in the case of plain areas detection for the red component E_c^R in the (t_{th}) frame:</p>	<p>IF $(NE_{E_c^R}$ is NO LARGE), THEN WE_c^R is NO LARGE.</p>

Table 2 Fuzzy rules used in FMANS filter

REFERENCES

1. Bovik, A., *Handbook of image and video processing (Communications, networking and multimedia)*, Orlando, FL, USA: Acad. Press Inc., 2005.
2. K. Plataniotis and A. Venetsanopoulos, *Color image processing and applications*, New York: Springer-Verlag, 2000.
3. R. Lukac, and K. N. Plataniotis, *Color Image Processing: Methods and Applications*, Boca Raton, FL, USA: CRC Press, 2006.
4. J. W. Woods, *Multidimensional Signal, Image, and Video Processing and Coding*, San Diego, CA, USA: Academic Press, 2011.
5. V. Kravchenko, H. Perez, V. Ponomaryov. *Adaptive Digital Processing of Multidimensional Signals with Applications*, Moscow: FizMatLit Edit., 2009. Available: <http://www.posgrads.esimecu.ipn.mx/> (Libros publicados).
6. V. Ponomaryov. "Real-time 2D-3D filtering using order statistics based algorithms," *J. Real Time Image Proc.*, vol. 1, no.3, pp.173-194, 2007.
7. V.I Ponomaryov., F.J. Gallegos-Funes, A. Rosales-Silva, "Real time color imaging based on RM-filters for impulsive noise reduction," *J. Imaging Sci. Technol.*, vol. 49, no.3, pp.205-219, 2005.
8. S. Morillas, V. Gregori, and A. Hervas, "Fuzzy Peer Groups for Reducing Mixed Gaussian-Impulse Noise From Color Images," *IEEE Trans. on Image Proc.*, vol. 18, no. 7, pp.1452-1465, 2008.
9. C. Lien, C. Huang, P. Chen, Y. Lin, "An Efficient Denoising Architecture for Removal of Impulse Noise in Images," *IEEE Trans. on Comp.*, vol.62, no. 4, pp.631-643, 2013.
10. J. Varghese, M. Ghouse, S. Subash, M. Siddappa, M. S. Khan, O. B.Hussain, "Efficient adaptive fuzzy-based switching weighted average filter for the restoration of impulse corrupted digital images," *IET Image Process.*, vol. 8, no. 4, pp.199-206, 2014.
11. J. Camarena, V. Gregori, S. Morillas and A. Sapena, "A simple fuzzy method to remove mixed Gaussian-impulsive noise from colour images," *IEEE Trans. Fuzzy Systems*, vol. 21, no.5, pp.971-978, 2013.
12. J.G. Camarena, V. Gregori, S. Morillas, A. Sapena, "Two-step fuzzy logic-based method for impulse noise detection in colour images," *Pattern Recogn. Lett.*, vol. 31, no. 13, pp.1842-1849, 2010.
13. B. Smolka, "Peer group switching filter for impulse noise reduction in color images," *Pattern Recogn. Lett.*, vol. 31, no. 6, pp.484-495, 2010.
14. V. Ponomaryov, A. Rosales, F. Gallegos, "3D filtering of colour video sequences using fuzzy logic and vector order statistics," *Lecture Notes in computer Science*, vol. LNCS 5807, pp. 210-221, 2009.
15. V. Kravchenko, V Ponomaryov, V Pustovoi, "Three dimensional filtration of multichannel video sequences on the basis of fuzzy-set theory," *Doklady Phys.* Springer, vol. 55, no.2, pp.58-63, 2010.
16. V. Ponomaryov, H. Montenegro, A. Rosales, G. Duchon, "Fuzzy 3D filter for color video sequences contaminated by impulsive noise," *J. Real Time Image Proc.*, [Online]. Doi:10.1007/s11554-012-0262-9, pp.1-16, 2012.
17. T. Melange, M. Nachtegaele and E. Kerre, "Fuzzy Random Impulse Noise Removal From Color Image Sequences," *IEEE Trans. on Image Processing*, vol. 20, no. 4, pp.959-970, 2011.
18. D. Fevrarev, N. Ponomarenko, V. Lukin, S. Abramov, K.Egiazarian and J. Astola, "Efficiency analysis of color image filtering," *EURASIP J. on Advances in Signal Processing*, doi:10.1186/1687-6180-2011-41, 2011.
19. V. Ponomaryov, F. Gallegos-Funes, A. Rosales-Silva. Fuzzy directional (FD) filter to remove impulse noise from colour images. *IEICE Trans. Fund. Electron. Commun. Comput. Sci.*, vol.E93-A, no.2, pp.570-572, 2010.
20. A. Rosales-Silva, F. Gallegos Funes, V. Ponomaryov, "Fuzzy Directional (FD) Filter for impulse noise reduction in colour video sequences," *J. Visual Commun. & Image Represent.*, vol. 23, no.1, pp.143-149, 2012.
21. L. Jovanov, A. Pizurica, A., A. Schulte, et. al., "Combined wavelet-domain and motion-compensated video denoising based on video codec motion estimation methods," *IEEE Trans. on Circ. Syst. Video Techn.*, vol.9, no.3, pp.417-421, 2009.
22. T. Melange, M. Nachtegaele, E. Kerre, "Video denoising by fuzzy motion and detail adaptive averaging," *J. Elect. Imag.*, vol.17, no.4, 043005, pp.1-19, 2008.
23. H. Yin, X. Fang, Z. Wei, and X Yang, "An Improved Motion-compensated 3-D LMMSE filter with spatio-temporal adaptive filtering support," *IEEE Trans. on Circ. Syst. Video Tech.*, vol.17, no.12, pp.1714-1727, 2007.
24. M. Dong, J. W. Zhang, Y. Ma, "Image denoising via bivariate shrinkage function based on a new structure of dual contour let transform," *Signal Processing*, vol.109, pp.25-37, April, 2015.
25. A. Buades, B. Coll, and J. Morel, "A non-local algorithm for image denoising," in *IEEE Conf. Comput. Vis. Patt. Recogn.*, vol. 2, pp.60-65, 2005.
26. K. Dabov, A. Foi, A., V. Katkovnik, and K. Egiazarian. Image denoising by sparse 3D transform-domain collaborative filtering", *IEEE Trans. Image Proces.*, Vol.16, No.8, 2080-2095, 2007.
27. V. Ponomaryov, H. Montenegro-Monroy, L. Nino de Rivera, H. Castillejos, "Fuzzy Filtering Method for Color Videos Corrupted by Additive Noise," *The Scientific World Journal*, Vol. 2014, pp.1-21, Art. ID 758107.
28. O. Pogrebnayk., V. Lukin, "Wiener discrete cosine transform-based image filtering," *Jf Elect. Imag.*, vol.21, no. 4, paper 043020: pp.1-15, 2012.
29. Video Trace Library, <http://trace.eas.asu.edu/yuv/>, Arizona State University.
30. Z. Wang, A. and Bovik, "Mean squared error: love it or leave it? A new look at signal fidelity measures," *IEEE Sign. Proc. Mag.*, vol.26, no.1, pp.98-117, 2009.

FMANS_H Algorithm

-
- (a) *Spatial Stage*
1. **Input:** RGB video frame $E_\eta^n(i, j)$; where $\eta = R, G, B$.
 2. Calculation of the Basic and related gradient values into a 7x7 sliding window $\nabla_{(k,l)}^n$.
 3. Fuzzy Similarity set:
 - if "LARGE"
 - Weighted Mean Filter
 - Wiener Multiscale Filter
 - Output** $\hat{E}_1^n(i, j)_1$
 - else
 - Fuzzy weights W_n^n are calculated.
 - Weighted Mean Filter
 - Wiener Multiscale Filter
 - Output** $\hat{E}_2^n(i, j)_2$
- (b) *Spatio-temporal denoising stage*
1. **Input** (t and $(t-1)$) video frames
 2. Fuzzy Similarity set is estimated $\delta E_\eta^n(i, j)$ in frames (t) and ($t-1$)
 3. Fuzzy Similarity set:
 - if "LARGE"
 - Weighted Mean Filter for samples from (t) and ($t-1$) frames
 - Wiener Multiscale Filter
 - Output** $\hat{E}_1^n(i, j)_1$
 - else
 - Motion estimation between (t) and ($t-1$) frames in eight directions
 4. Fuzzy Similarity of motion set:
 - if "LARGE"
 - Alfa-TM Filter for common samples from frames (t) and ($t-1$)
 - Wiener Multiscale Filter
 - Output** $\hat{E}_2^n(i, j)_2$
 - else
 - Weighted Mean Filter in frame (t)
 - Wiener Multiscale Filter
 - Output** $\hat{E}_1^n(i, j)_1$
- (c) *Fuzzy Spatial postprocessing Filtering Stage*
1. Edge Detection Similarity set
 - if "LARGE"
 - Alfa-TM filtering is executed in frame (t)
 - Output** $\hat{E}_3^n(i, j)_3$
 - else
 - Weighted Mean Filter in frame (t)
 - Wiener Multiscale Filter
 - Output** $\hat{E}_2^n(i, j)_2$
-

Table 3 Structure of FMANS_H algorithm

Volodymyr Ponomaryov attained his Ph.D. in 1974, D. Sci. in 1981. He is a referee for more than 20 international scientific journals. Dr. Ponomaryov served and currently serves as Associate Editor in *Journal of Real Time Image Processing* (Springer) His current research interests include signal/image and video processing, denoising, pattern recognition, real-time filtering, 3D reconstruction, medical sensors, remote sensing, etc. Over the years, he has also been a promoter of 39 Ph.D. degrees on his scientific areas of interests. He has authored or co-authored more than 150 scientific papers in the international referred journals, more than 300 international referred conference papers, and also 23 patents of ex USSR, Russia and Mexico, and five scientific books in international editorials.

Noise %	3D FD				FRINR_seq				FMINS			
	F		MA		F		MA		F		MA	
	PSNR	MAE	PSNR	MAE	PSNR	MAE	PSNR	MAE	PSNR	MAE	PSNR	MAE
0	30.46	1.68	48.22	0.037	30.99	1.47	49.62	0.022	31.13	1.26	50.14	0.012
5	29.41	2.13	39.36	0.381	30.26	1.95	39.75	0.369	30.47	1.82	40.22	0.349
10	28.46	2.72	35.99	0.752	29.97	2.25	36.52	0.716	30.19	2.11	36.92	0.693
20	26.84	4.16	32.10	1.826	27.21	3.84	32.65	1.610	27.34	3.64	32.74	1.602
	NCD	SSIM	NCD	SSIM	NCD	SSIM	NCD	SSIM	NCD	SSIM	NCD	SSIM
0	0.004	0.882	0.000	0.989	0.003	0.882	0.000	0.989	0.002	0.883	0.000	0.9892
5	0.005	0.847	0.002	0.982	0.005	0.849	0.001	0.982	0.003	0.851	0.000	0.9824
10	0.006	0.816	0.003	0.977	0.005	0.818	0.003	0.976	0.005	0.820	0.001	0.977
20	0.009	0.756	0.009	0.961	0.007	0.757	0.007	0.961	0.007	0.758	0.005	0.962
Noise %	3D FD				FRINR Seq				FMINS			
	S		SM		S		SM		S		SM	
	PSNR	MAE	PSNR	MAE	PSNR	MAE	PSNR	MAE	PSNR	MAE	PSNR	MAE
0	45.23	0.175	47.69	0.106	45.87	0.160	47.71	0.094	46.33	0.152	47.94	0.083
5	41.36	0.184	42.46	0.177	42.11	0.165	42.92	0.168	43.11	0.158	43.67	0.156
10	37.69	0.197	38.19	0.517	38.58	0.174	39.64	0.513	38.82	0.168	40.27	0.505
20	32.01	0.221	36.37	0.738	32.94	0.199	36.73	0.712	33.11	0.188	37.12	0.702
	NCD	SSIM	NCD	SSIM	NCD	SSIM	NCD	SSIM	NCD	SSIM	NCD	SSIM
0	0.009	0.986	0.003	0.962	0.007	0.993	0.002	0.962	0.004	0.994	0.000	0.963
5	0.010	0.959	0.006	0.931	0.009	0.973	0.004	0.930	0.005	0.979	0.002	0.943
10	0.013	0.941	0.009	0.909	0.011	0.952	0.007	0.915	0.008	0.958	0.003	0.921
20	0.019	0.913	0.012	0.872	0.012	0.907	0.009	0.874	0.010	0.915	0.008	0.878

Table.4 Mean per 100 frames values for PSNR (dB), MAE, NCD and SSIM criteria (F, MA, S and SM video sequences)

Filters	FDARTF_G [5]		VBM3D [26]		NLM [25]		FMANS_2		FMANS_H	
Noise variance	PSNR	MAE	PSNR	MAE	PSNR	MAE	PSNR	MAE	PSNR	MAE
0.000	28.13	7.31	28.92	6.54	28.83	6.65	29.38	6.39	29.71	6.21
0.005	26.67	8.80	27.81	7.78	27.69	7.96	28.17	7.54	28.51	7.38
0.010	25.40	10.65	26.44	9.81	26.37	9.94	26.83	9.65	27.09	9.39
0.020	23.42	13.56	24.21	11.58	24.16	11.73	24.58	11.33	24.95	11.17
0.030	22.38	15.28	23.04	13.15	22.93	13.29	23.37	12.96	23.75	12.64
Noise variance	NCD	SSIM	NCD	SSIM	NCD	SSIM	NCD	SSIM	NCD	SSIM
0.000	0.014	0.8532	0.011	0.8869	0.011	0.8861	0.011	0.8882	0.010	0.8896
0.005	0.016	0.7975	0.015	0.8190	0.016	0.8179	0.015	0.8192	0.013	0.8209
0.010	0.023	0.7248	0.019	0.7515	0.019	0.7509	0.019	0.7523	0.018	0.7542
0.020	0.028	0.6395	0.020	0.6665	0.021	0.6659	0.020	0.6672	0.019	0.6694
0.030	0.031	0.6033	0.022	0.6319	0.023	0.6312	0.022	0.6323	0.021	0.6338

Table 5 Mean values per 50 video frames for criteria PSNR (dB), MAE, NCD, SSIM for color video Flowers after filtering

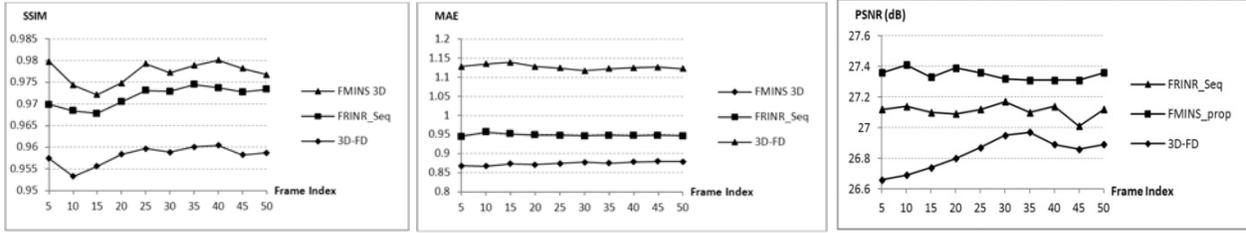


Fig.4 a) SSIM results for the different methods on *Stefan* video ($p_n=5\%$). b) MAE results for the different methods on *Carphone* video ($p_n=20\%$). c) PSNR results for the different methods on *Flowers* video ($p_n=20\%$)

Filters	FDARTF_G [5]		VBM3D [26]		NLM [25]		FMANS_2		FMANS_H	
	PSNR	MAE	PSNR	MAE	PSNR	MAE	PSNR	MAE	PSNR	MAE
Noise variance										
0.000	35.68	5.36	36.46	4.74	36.39	4.82	36.83	4.62	37.04	4.57
0.005	33.25	7.70	34.41	7.07	34.38	7.21	34.80	6.89	35.08	6.58
0.010	31.69	9.70	32.53	8.79	32.41	8.88	32.95	8.65	32.28	8.48
0.020	28.84	11.31	29.65	10.62	29.57	10.69	30.07	10.48	30.32	10.27
0.030	26.26	12.82	27.39	12.03	27.32	12.14	27.83	11.83	28.05	11.54
Noise variance	NCD	SSIM	NCD	SSIM	NCD	SSIM	NCD	SSIM	NCD	SSIM
0.000	0.020	0.9454	0.018	0.9595	0.018	0.9586	0.018	0.9611	0.017	0.9629
0.005	0.022	0.8876	0.020	0.9005	0.021	0.8996	0.019	0.9019	0.018	0.9047
0.010	0.025	0.8279	0.022	0.8486	0.023	0.8472	0.022	0.8498	0.021	0.8524
0.020	0.030	0.7430	0.027	0.7721	0.028	0.7706	0.026	0.7737	0.025	0.7761
0.030	0.033	0.7074	0.029	0.7256	0.031	0.7241	0.028	0.7271	0.026	0.7292

Table 6 Mean values per 50 video frames for criteria PSNR (dB), MAE, NCD, SSIM for color video *Stefan* after filtering

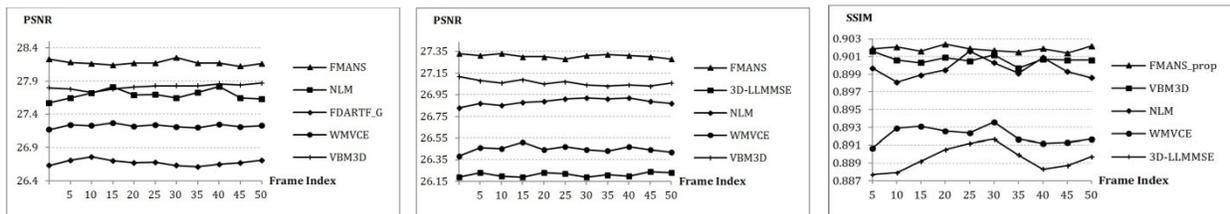


Fig.5a PSNR values for the different methods on *Tennis* video ($\sigma^2 = 0.015$) b) PSNR values for the different methods on *Stefan* video ($\sigma^2 = 0.005$) c) SSIM values for the different methods on *Stefan* video (variance 0.005)

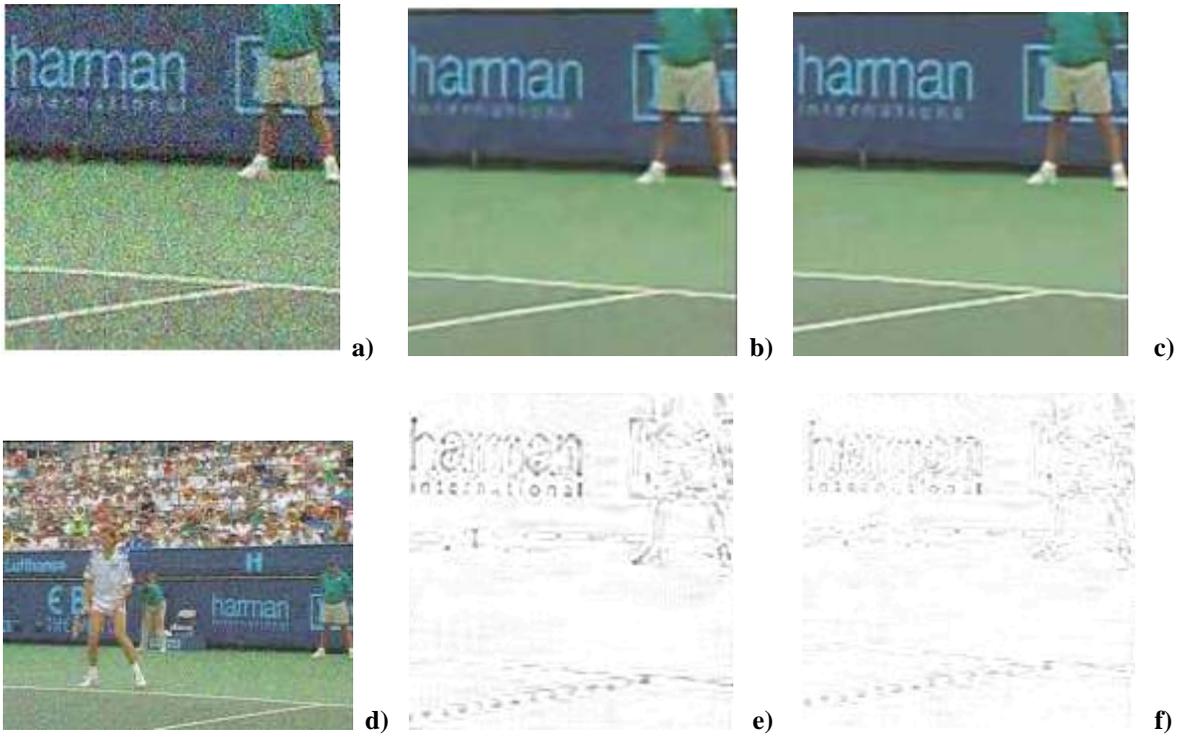


Fig. 6 (d) frame No.50 of color video *Stefan* contaminated by additive Gaussian noise with variance 0.02; a) zoomed part of contaminated frame, and filtered frames by: b) – *NLM*; c) – *FMANS_H*; e) and f) inverted error (amplified in 3 times) after filtering, accordingly for filters b) and c)

Temporal data approach performance

Michal Kvet

Abstract—Temporal database is an extension of the concept of standard databases which process only current valid data. Temporal approach can manage also historical and future valid data and provides management and progress monitoring over the time. Based on granularity and frequency of Update operations, there can be several approaches. This paper deals with object level uni-temporal system and extended column level principle. It describes the principles, required methods, procedures, functions and triggers to provide complex functionality of the system. Defined operations usually process large amount of data, therefore, it is necessary to define indexes for better data access. In the Performance evaluation section, 17 different types of indexes are compared based on speed and disc requirements.

Keywords—temporal database, column level approach, validity, index

I. INTRODUCTION

DATABASE systems are one of the most important parts of the information technology. It is generally known that database system is usually the basic part - the root of any information system. The development of data processing has brought the need for modelling and accessing large structures based on the simplicity, reliability and speed of the system [1] [2] [6]. However, even today, when database technology is widespread, most databases process and represent states of the data valid in this moment. Properties and states of the objects evolve over the time, become invalid and are replaced by new ones. Once the state is changed, the corresponding data are updated in the database and it still contains only the current valid data. However, temporal data processing is very important in dynamically evolving systems, industry, communication systems and also in systems processing sensitive data, which incorrect change would cause a great harm. It can also help us to optimize processes and make future decisions [3] [4] [6].

Historical data management using log files and archives is completely inappropriate. Although the data can be obtained, it is too complicated, it lasts too much time and the output data are obviously not in suitable form (raw material). The main disadvantage is the need of the administrator intervention (operation manager) and possible data loss [4] [6] [8]. Moreover, this system cannot manage future valid data and has very poor support for transaction management.

Although conventional databases and backups are used for actual valid data processing, developed functionality allows historical data management.

Flashback feature available in database systems is primarily intended as a method to restore the database after the incorrect destructive operations. However, it can be used for validity limited data processing. For this purpose, it is necessary to allocate enough disk space to store information about updates. However, the period during which it is possible to obtain object or table status, is directly influenced by the size of disc space and parameter *UNDO_RETENTION* (period during which history is saved, default value is 900 seconds). In principle, it can be said that this is a special type of log file, thus also "unnecessary" information are stored.

Another approach extends the concept of the conventional principle by the definition of the validity as the part of the primary key. Nowadays, spectrum of the data structures for conventional approach can be transformed and used also for temporal characteristics. Rows with the different values of primary key can represent the same object during the validity time frames characterized by the *Update* operation. Transactions for data management must therefore contain not only object itself, but also processed time interval.

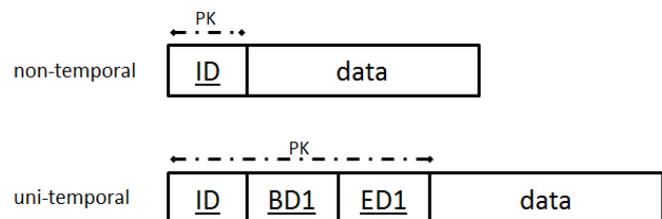


Fig. 1. Object level temporal data management [7]

Time validity interval modelling can be represented by various definitions (mostly closed-closed representation, closed-open representation) [6] [7].

II. EXTENDED COLUMN LEVEL TEMPORAL SYSTEM

The solution described in the previous section manages attributes of the objects, whereas the standard uni-temporal model works with the whole objects. Extended column level

This publication is the result of the project implementation: Centre of excellence for systems and services of intelligent transport, ITMS 26220120028 supported by the Research & Development Operational Programme funded by the ERDF and Centre of excellence for systems and services of intelligent transport II, ITMS 26220120050 supported by the Research & Development Operational Programme funded by the ERDF.

Michal Kvet, Faculty of Management Science and Informatics, University of Zilina, Slovakia. E-mail: Michal.Kvet@fri.uniza.sk

temporal system can be considered as the improvement of the column level temporal system [5] [6] in the term of the performance and the simplicity of the model management for the users. It is extended by the definition of the type of the operation. If the operation is *Update*, there is also the reference to the data type tables with historical values.

Existing applications are connected to the conventional layer with actual values, thus program can continue to operate without any changes. The main part is to manage the table containing information about the changes of temporal columns. Column, which changes need to be monitored, is temporal. If the value is changed, information about the *Update* is stored in the developed *temporal_table* and historical value is inserted into to the table containing historical values. Fig. 2 shows the complete structural model. Application is directly connected to the main tables with current valid values. It means that currently used applications can be used without any change. Historical values are stored in the special section, each temporal data type has one table defined by the identifier (primary key) got using the sequence and trigger and the values themselves. Thus, the principle and system is similar to the column level temporal system, but historical values management and temporal table is different.

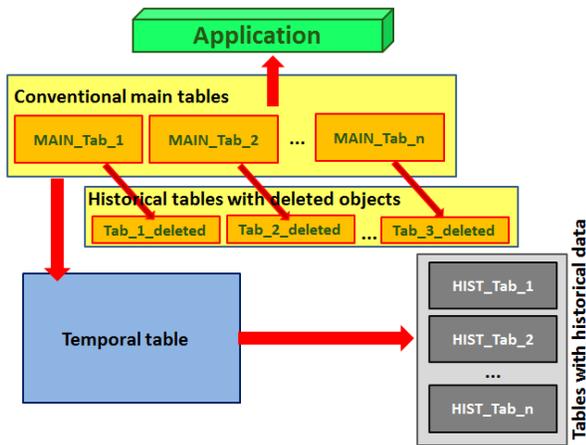


Fig. 2. Extended column level temporal system structure.

Management of the temporal table is completely different. It consists of these attributes (fig. 3) [7]:

- *ID_change* – got using sequence and trigger – primary key of the table.
- *ID_previous_change* – references the last change of an object identified by *ID*. This attribute can also have *NULL* value that means, the data have not been updated yet, so the data were inserted for the first time in past and are still actual.
- *ID_tab* – references the table, record of which has been processed by DML statement (*Insert*, *Delete*, *Update*, *Restore*).
- *ID_orig* - carries the information about the identifier of the row that has been changed.
- *ID_column* – holds the information about the changed attribute (each temporal attribute has defined value for the referencing).

- *Data_type* – defines the data type of the changed attribute:
 - *C* = *char / varchar*
 - *N* = *numeric values (real, integer, ...)*
 - *D* = *date*
 - *T* = *timestamp*

This model can be also extended by the definition of the other data types like binary objects.

- *ID_row* – references to the old value of attribute (if the DML statement was *Update*). Only update statement of temporal column sets not *NULL* value.
- *Operation* – determines the provided operation:
 - *I* = *Insert*
 - *D* = *Delete*
 - *U* = *Update*
 - *R* = *Restore* (renew validity)

The principles and usage of proposed operations are defined in the part of this paper.

- *BD* – the begin date of the new state validity of an object.

Temporal_table		
<i>id_change</i>	Integer	NN (PK)
<i>id_previous_change</i>	Integer	
<i>operation</i>	operation_domain	NN
<i>id_tab</i>	Integer	NN
<i>id_orig</i>	Integer	NN
<i>id_column</i>	Integer	
<i>id_row</i>	Integer	
<i>bd</i>	Date	NN
<i>data_type</i>	data_type_domain	

Fig. 3. Temporal table.

A. Insert

Trigger before insert operation stores information about this operation into *temporal_table*. New value of the primary key is set using the sequence:

```
create sequence seq_temp_tab start with 1
increment by 1;
```

It is new object, therefore attribute *id_previous_change* (based on this object) has *NULL* value. Also attributes referencing temporal attributes changes have *NULL* values (*id_column*, *id_row*, *data_type*). Value 1 references the main table (*I = Tab1*):

```
insert into temporal_table values
(seq_temp_tab.nextval, null, 'I', 1,
:NEW.id, null, null, sysdate, null);
```

B. Update

Update statement is a bit more complicated, because the root of the system granularity is the attribute itself, not the whole object. Thus, the *Update* operation is divided into *Updates* of the temporal columns – each standard update operation causes different number of the temporal updates. However, it is important to mention the main table can also consists of the non-temporal column (conventional column) – it means, that the attribute does not change its value over the time or the changes are not important to be monitored. These data updates do not cause the insert operation into the *temporal_table*.

The next code shows the series of the steps, which must be done to store the history of the temporal attributes complexly:

- Get the value of the last change of the appropriate object.
- Store the old value of the changed attribute in the historical tables based on the data type:
 - $C = \text{char} / \text{varchar} \rightarrow c_tab$
 - $N = \text{numeric values} \rightarrow n_tab$
 - $D = \text{date} \rightarrow d_tab$
 - $T = \text{timestamp} \rightarrow t_tab$
- Insert the information about the update operation and references into the *temporal_table*.

```
--date has been changed
if (:OLD.datum<>:NEW.datum) then

  get maximal value of the change into
  the variable v_max_change

insert into d_tab(val)
values (:OLD.datum)
returning ID into v_id_hist_val;

insert into temporal_table values
(seq_temp_tab.nextval,v_max_change,
'U', 1, :OLD.id, 1, v_id_hist_val,
sysdate, 'D');
end if;
```

This operation is also the limitation of the whole system. It is necessary to know the ratio of the changed attributes at a given time to the total number of attributes (temporal and conventional). If all of the attributes are temporal and have the same granularity for the changes, the uni-temporal system managing the whole object is more convenient. However, a lot of systems used today are characterized by the different granularity. In that case, this solution rapidly decreases the need for the disc storage and process the changes more easily.

C. Delete

Operation for the deleting objects can be divided into two groups.

The proposed scheme handles only the beginning of the validity. If the new valid value for the attribute, respectively for the whole object is not available (or is incorrect), such object must be removed from the table, which contains only current valid objects.

However, this state is often only temporary and short-term, therefore it is convenient to remove the entire object due to time management. Thus, the current time value is set into the attribute validity, which determines the last time point of the validity. If the values are corrected and known later, the object state is rebuilt and the value of the validity has the *NULL* value then. Of course, all of these operations are stored into the *temporal_table*. Thus, we can get the complete image of the object in the system over the time.

Another situation occurs, if the user is sure, that the object should be removed. In that case, the whole object is moved into historical table, which has the same structure

Based on the performance, the system defines the time interval, which determines the maximal time, during the invalid object can be placed in the main table. After the defined time, the object is automatically moved into the historical tables. Database cursor processes the mentioned objects:

```
for cur IN
(select *
from TAB1
where MONTHS_BETWEEN(sysdate, validity)>=1)
LOOP
processing
END LOOP;
```

Processing inside the cycle consists of these steps:

- Insert information about the historical object into the table managing historical (deleted) objects (*tab1_deleted*).
- Remove old value from main table.
- Store information about the operation into the *temporal_table*.

For the simplicity, the previous code does not show the time interval directly, but deletes the objects after 1 months.

D. Restore

The role of the restore operation is to reestablish the validity of the object. If the object is in main table, it deletes the value of the validity. However, if the object is in the table with deleted objects, it must be transformed into the main table (the principle is similar to the operation delete).

```
select count(*) into v_count_tab1
from TAB1
where ID=pID;
if v_count_tab1=1 then
--change validity
--insert into temporal_table
end if;
select count(*) into v_count_tab1_del
from TAB1_deleted
where ID=pID;
if v_count_tab1_del=1 then
--move from historical table with deleted
objects
--insert into temporal_table
end if;
```

Each operation manages *temporal_table*, whereas update operation deals also with the tables for historical values. *Delete* operation updates status in main table (*Tab1*), if necessary. *Restoring* data transforms not actual object into current valid.

III. GETTING STATES – SELECT

The most important part of the system is the optimization for the operation getting states of the objects, which can be considered as the critical part of the whole model. *Select* operation is the mostly used to get the states, to create future prognoses and optimize the whole system. Therefore, the special structures like indexes are developed to improve the performance of this operation.

Select operation can be divided into two groups based on the data management:

- *Get database state.* This method is based on selection all currently valid objects in the database.
- *Get object state.* Monitoring the progress of the changes of the attributes is far more complicated. Using *temporal_table*, the object data are reconstructed and the complete states during the defined time interval or timepoint can be obtained. This procedure has in principle two parameters – identifier of the object and time point determines the validity – the system provides data from the defined time till now. Using cursor, the data about operations from *temporal_table* are processed.

```
select * bulk collect into v_nest_tab1
  from temporal_table
  where id_tab=1
     and id_orig=pID
     and bd>=pBD
  order by BD DESC;
-- add the last operation before the
defined time into collection (if this
operation exists)
```

Then, the changes are processed and written step by step from now to the history.

IV. INDEX STRUCTURES

One of the main features of optimization is based on using index structures. Temporal databases are oriented for state management and monitoring over the time. Getting states and individual changes in the *Select* statement forms the core of a major milestone of efficiency and speed of the system.

Oracle defines an index as an optional structure associated with a table or table cluster that can sometimes speed data access. By creating an index on one or more columns of a table, you gain the ability in some cases to retrieve a small set of randomly distributed rows from the table. Indexes are one of many means of reducing disk I/O. If a heap-organized table has no indexes, then the database must perform a full table scan to find a value.

The absence or presence of an index does not require a change in the wording of any SQL statement. An index is a fast access path to a single row of data. It affects only the speed of execution. Given a data value that has been indexed, the index points directly to the location of the rows containing that value. Database management system automatically maintains the created indexes – changes (*Insert*, *Delete*, *Update*) are automatically reflected into index structures.

However, the presence of many indexes on a table degrades the performance because the database must also update the indexes. Moreover, the optimizer does not have to use them, if the full scan would be faster or if the index is not suitable for the query based on the conditions. Thus, if the user forces the system to use the index, the performance rate can be much worse than without their use. Each index has 2 properties – visibility (*invisible index = maintained, but not used by the optimizer*) and usability (*unusable index = not maintained and not used by the optimizer*). Information about the indexes can be found in *user_indexes* system table [7].

A. Index sequence file

Index sequence data alignment is based on two data sets - the sequential file and the index file, consisting of a key and a pointer to the data row in a sequential file. Finding the record is based on the index scan and access to primary (sequential) file. If the file was too large, it is possible to add another index layer - hierarchical index.

The main disadvantage of index-sequential arrangement is the significant performance decrease in performance (processing time requirements and the response to it) based on amount of data increase. This problem can be partially solved by the data index reorganization.

B. B-tree, B+tree

The index structure of the B+tree is mostly used because it maintains the efficiency despite frequent changes of records (*Insert*, *Delete*, *Update*). B+tree index consists of a balanced tree in which each path from the root to the leaf has the same length.

In this structure, we distinguish three types of nodes - root, internal node and leaf node. Root and internal node contains pointers S_i and values K_i , the pointer S_i refers to nodes with lower values the corresponding value (K_i), pointer S_{i+1} references higher (or equal) values. Leaf nodes are directly connected to the file data (using pointers).

B+tree extends the concept of B-tree by chaining nodes at leaf level, which allows faster data sorting. DBS Oracle uses the model of two-way linked list, which makes it possible to sort ascending and descending, too.

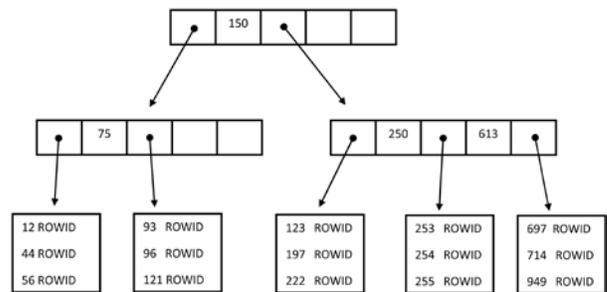


Fig. 4. B-tree

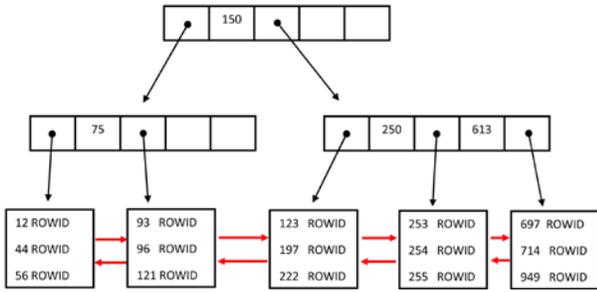


Fig. 5. B+tree

Limitation of this approach is a small number of records (low cardinality). In that case, using index does not produce the desired effect in terms of performance (acceleration). Another disadvantage is the lack of support SQL queries with functions implicitly.

C. Inverted key B+tree

Index B+tree structure with inverted key is used in case of often requirement for tree balancing (column value is obtained using the sequence and the trigger - autoincrement) caused by frequent execution of the *Insert* statement. Indexing will not use original key value, but the inverted variant. For example, for the key 123 is inverted key value 321.

D. Bitmap index

Bitmap indexes are represented by two-dimensional array, the number of rows is identical to the cardinality of the table. The first column contains a reference to a record in the data file, the other columns are called bitmap and represent different values of the indexed column.

The following figure illustrates the general assessment of conditions of the order using bitmap indexes. Let have a table - cars with primary key represented by the registration number (*license_plate*). Non-key attributes include color (*color*), producer (*producer*), construction year (*year*) and price (*price*). The aim is to find red cars produced by the “Škoda” brand in 2012. The solution shows following figure.

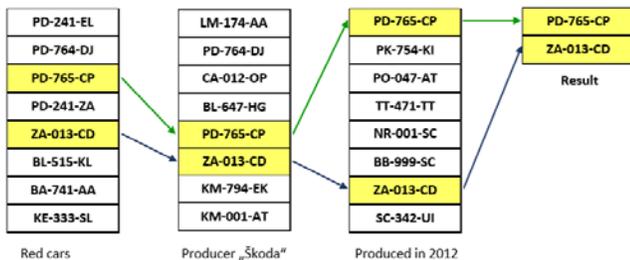


Fig. 6. Bitmap

As can see in fig. 6, this type of index is very effective for the *Select* statement. On the other hand, *Insert*, *Update*, and *Delete* cause problems because implementation may change the size of the bitmap index.

E. Hash index

Another type of indexing is hash approach. Hash index is composed of a plurality of blocks (*buckets*). The basis of the approach is mapping key value taken by a hash function - $h()$ defining the index block that contains a pair - value and pointer. If the block is full, it is necessary to add a overflow block. Using this index is characteristic using the condition $attribute_value = parameter$.

F. Table index

If the table is too small, it is better to keep complete records directly instead of references (pointers) to the leaf level data. This approach reduces the number of *Read* operations and thus accelerates the evaluation of the request. Oracle (version 8 and later) allows you to create a special type of table (*index-only table*), where the entire table is stored in the index structure (B + tree).

G. Cluster index

Cluster index provides the physical layout and location of the data set by attributes constituting the index. Accordingly, each table can have no more than one index of such type.

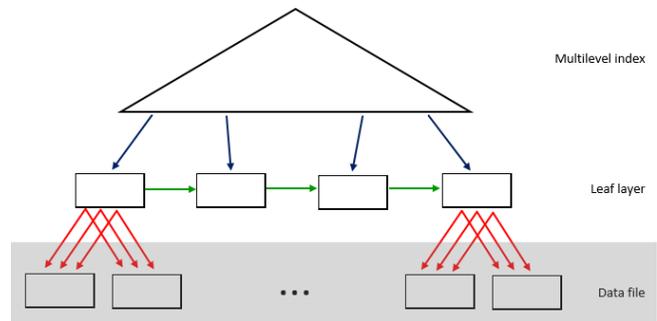


Fig. 7. Cluster index

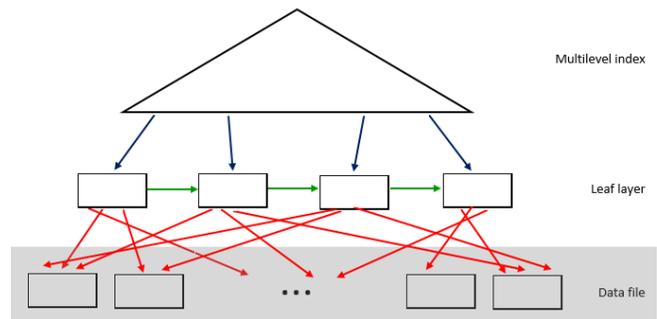


Fig. 8. Non-clustered index

H. Table cluster

Cluster tables is linked to several tables, records with the same key value (*cluster key*) are stored together – these data are usually required together. It means, that the clustered data can be loaded using one *Read* operation. In temporal system, historical and future valid data are clustered together with actual data to provide complete monitoring. Cluster key is the identifier of the data extended by the definition of the table, if the primary key identifier does not provide it automatically.

V. PERFORMANCE EVALUATION

The aim of the developed structure is to address the lack of data management with different granularity and frequency of changes. The model performance is compared based on various developed index structures. *Insert*, *Update*, *Delete* (logical) and *Restore* (renew validity) operations have been monitored to highlight time consumption performance and disc space size requirements. The main criterion of quality is the speed of the processed *Select* statement, which reflects the proposed index structures. Global object changes over the time have been monitored (*Select_obj*). Another fundament is characterized by the duration to obtain actual snapshot of the temporal database (*Select_db*).

This section deals with 17 different indexes, which are compared to declare the quality and usability of the system. Notice, some index can even decrease the performance of the system (optimizer is not forced to used index by hints). The system uses none or one developed index for *temporal_table* (fig. 3):

No index (*ind1*).

- B+tree:
 - ID_orig (*ind2*),
 - ID_orig, id_previous_change (*ind3*),
 - ID_orig, BD (*ind4*),
 - ID_orig, ID_previous_change, BD (*ind5*),
 - Unique – ID_orig, ID_previous_change (*ind6*),
 - BD, ID_orig (*ind7*),
 - BD, ID_orig, ID_previous_change (*ind8*),
- Inverted key B+tree:
 - ID_orig (*ind9*),
 - ID_orig, id_previous_change (*ind10*),
 - ID_orig, BD (*ind11*),
 - ID_orig, ID_previous_change, BD (*ind12*),
 - Unique – ID_orig, ID_previous_change (*ind13*),
 - BD, ID_orig (*ind14*),
 - BD, ID_orig, ID_previous_change (*ind15*),
- Bitmap indexes:
 - ID_orig (*ind16*),
 - ID_orig, ID_previous_change (*ind17*),

Experiment results were provided using Oracle Database 11g Enterprise Edition Release 11.2.0.1.0 - 64bit Production; PL/SQL Release 11.2.0.1.0 – Production. Parameters of used computer are:

- Processor: Intel Xeon E5620; 2,4GHz (8 cores),
- Operation memory: 16GB,
- HDD: 500GB.

Complete number of each operation was 10 000 (*Insert*, *Update*, *Delete*, *Restore*). Number of updated temporal attributes has been generated, total number was 24 965. Minimal number of operations on the object was 3, maximal number was 26, the average value was 5,4965.

The following code characterizes the structure of main table:

```
desc tabl;
Name          Type
-----
ID            NUMBER( 38 )
DATUM        DATE
CISLO        NUMBER
CISLO2       NUMBER( 38 )
RETAZEC      VARCHAR2( 100 )
VALIDITY     DATE
```

Main factor of the system consists of the *Select* statement performance. Therefore, this operation is highlighted. Of course, some operations (mostly *Insert*) can reflect slowdown because of the index construction. Performance measurement is divided into two groups. The first one consists of the states of the objects monitoring over the time. Using this criterion solution prefer index with object identifier as the first category. In comparison with solution without indexes (*reference 100%*), we can get even about 69,01% improvement (*ind6*) and 69,00% (*ind13*). In comparison with other indexes, these results are provided:

- 68,21% (*ind7 – reference 100% vs. ind6*),
- 68,53% (*ind8 – reference 100% vs. ind6*),
- 68,50% (*ind14 – reference 100% vs. ind6*),
- 68,13% (*ind15 – reference 100% vs. ind6*).

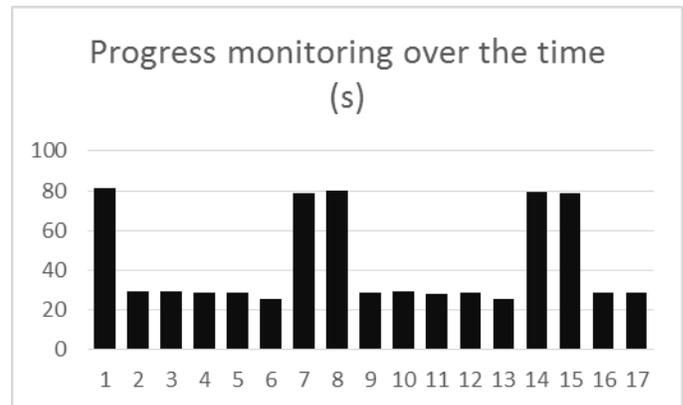


Fig. 9. Progress monitoring over the time

The second group consists of method to get actual state of the database, respectively database table. The results are almost balanced, only *ind6* and *ind13* provides better performance (in comparison with solution without indexes – reference 100%), specifically:

- 19,19% (*ind6*),
- 5,74% (*ind13*).

Now, it is necessary to group these results into other operation. *Insert* operation prefers *ind9 – 10,12%* (reference no index – 100%), *Update* operation prefers *ind13 – 58,89%* (reference no index – 100%), *Delete* operation prefers *ind2 – 57,35%* (reference no index – 100%) and *Restore* operation prefers *ind10 – 62,38%* (reference no index – 100%). At the first sight, there can be the inconsistency and ambiguity in the term of used index. However, there is not the same weight of the operations. In temporal system, mostly used destructive operation is *Update*, therefore the *ind13* is the best, also *ind6* provides very good performance results.

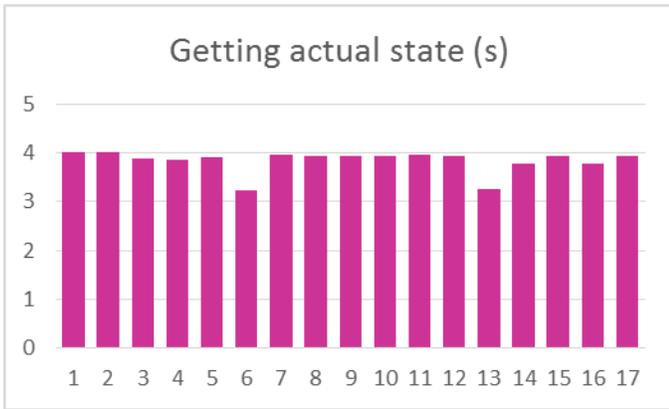


Fig. 10. Getting actual state

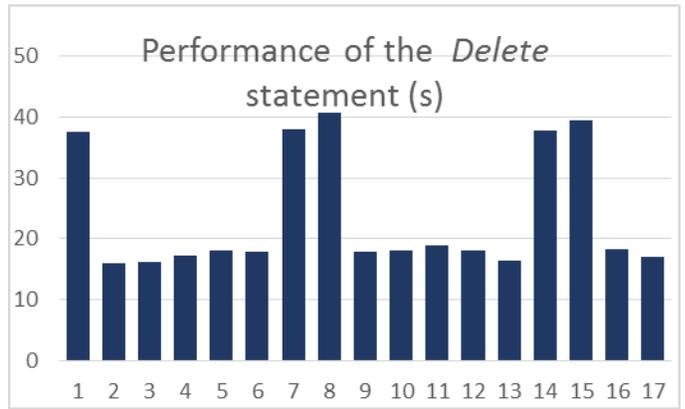


Fig. 13. Delete statement results

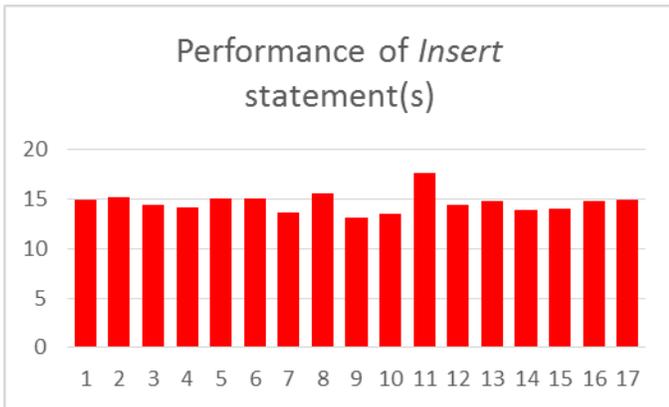


Fig. 11. Insert statement results

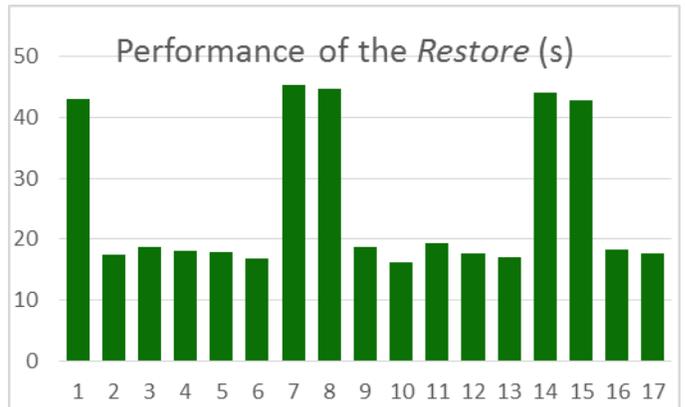


Fig. 14. Restore statement results

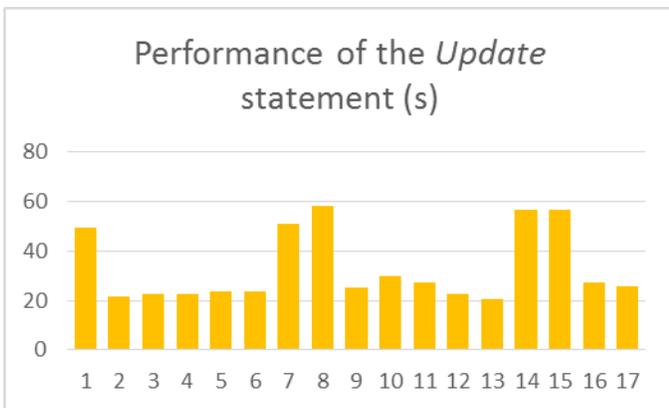


Fig. 12. Update statement results

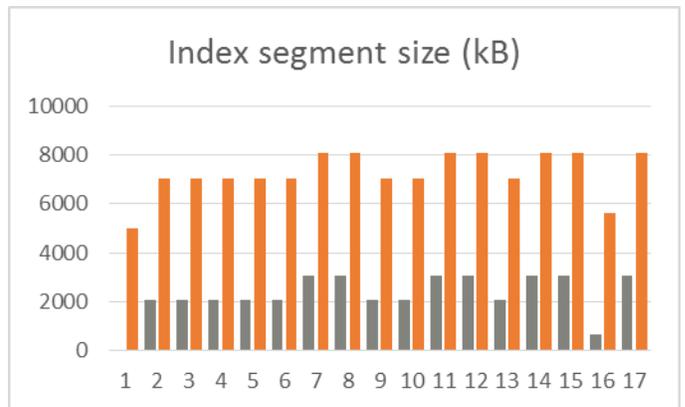


Fig. 15. Index segment size

The last category of tests provides the size of the index structure and global structure.

Using 10 000 rows in main table (*Tab1*), the index structure size varies from 640kB when using bitmap index with object identifier (*ind16*) up to 3072kB, when managing *ind7*, *ind8*, *ind11*, *ind12*, *ind14*, *ind15* and *ind17*. Fig. 15 shows the index size results. Gray color shows the size of the index, orange color summarizes complete structural size of the whole temporal database system.

As we can see, the time performance and size of the structure do not provide the same results, they do not prefer the same index (although the differences are not significant). Therefore, one of these segments must be dominant. In global conclusion, we can say, the best solution for temporal system is *ind13* or *ind6* based on unique index definition containing object identifier (*ID_orig*) and reference to previous change of the state (*ID_previous_change*). It requires 2048kB disc space for index (10 000 objects stored).

Importance of index definition is highlighted by the growth of data rows. Data management requirements based on number of stored data are modelled using exponential function. In *Update* and *Select* statement, the differences are most significant. Fig. 16 and fig. 17 shows the performance of the

system with various number of processed data (100 – blue color, 1000 – orange color, 10 000 – gray color).

Tab. 2 shows complex performance results.

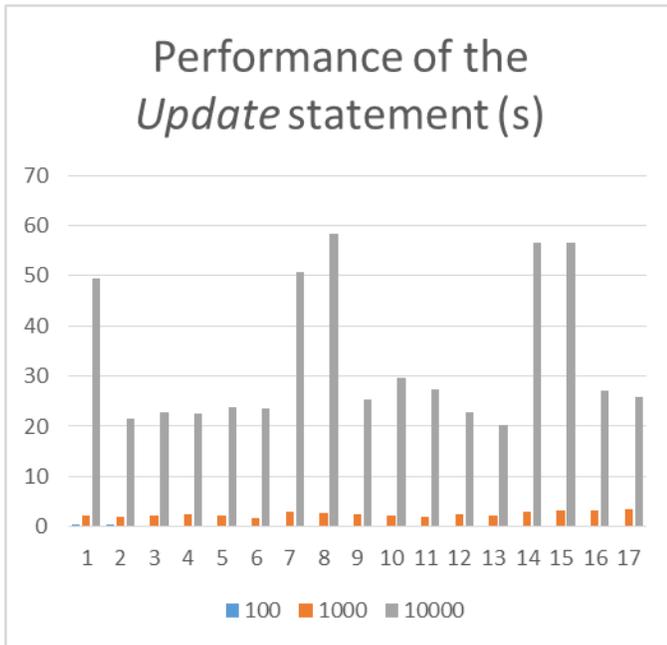


Fig. 16. Performance of the Update statement

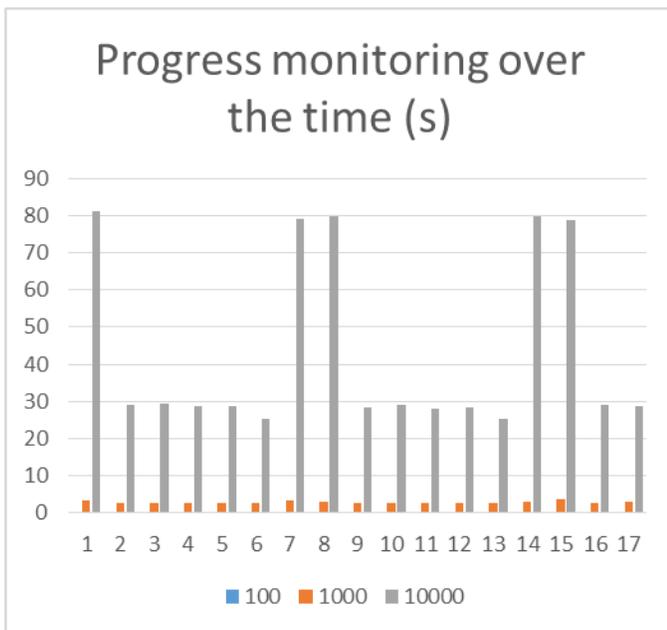


Fig. 17 Object monitoring results

VI. CLASSIFICATION OF TEMPORAL DATABASES

Current issues in this area have not use any classification for temporal model of databases. From this point of view we suggested classification rule in this form:

$$\alpha/\beta/\gamma/$$

where

α – represents kind of DBS:

- N – No database system support (e.g. file system only)
- R – Relational DBS (RDBS)
- X – Object relational DBS
- O – Object oriented DBS
- U – Unspecified DBS

β – represents kind of the temporal model:

- N – Non-temporal (Conventional)
- U – Uni-temporal
- B – Bi-temporal
- M – Multi-temporal
- E – Difference defined value of the attribute (ϵ - Epsilon)

γ – characterizes kind of transaction processing (Online transaction processing - OLTP):

- N – Nontransactional
- L – OLTP only with logs
- O – OLTP with temporal objects
- A – OLTP with temporal attributes of the object

For example, temporal database with type *R/U/O* represents temporal database using RDBS with the uni-temporal model and with an OLTP support.

Temporal querying requires definition of used index structure for each table:

Table 1. Index type

N	<i>No index</i>	Without the use of user-defined index.
SI	<i>Standard index approach</i>	Standard index type.
I_TS	<i>Indexes in multiple table spaces</i>	Indexes in different table spaces.
I_TSMD	<i>Indexes in multiple drives (HDD)</i>	Indexes in different table spaces, which are located on different discs.
I_TSDD	<i>Distributed indexes</i>	Distributed indexes.

Each object can be used either to process:

- object during a time interval,
- all objects within a defined period of time.

Moreover, it should be possible to associate objects into groups based on common characteristics (e.g. geographical area, sensors for one device, ...) [9]. The following figure illustrates the general location of the index.

Table 2. Experiment results

		1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17		
number of data	I	10000	10000	10000	10000	10000	10000	10000	10000	10000	10000	10000	10000	10000	10000	10000	10000	10000	10000	
	U	24965	24965	24965	24965	24965	24965	24965	24965	24965	24965	24965	24965	24965	24965	24965	24965	24965	24965	
	D	10000	10000	10000	10000	10000	10000	10000	10000	10000	10000	10000	10000	10000	10000	10000	10000	10000	10000	
	R	10000	10000	10000	10000	10000	10000	10000	10000	10000	10000	10000	10000	10000	10000	10000	10000	10000	10000	10000
	Min	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3
	Avg	5,4965	5,4965	5,4965	5,4965	5,4965	5,4965	5,4965	5,4965	5,4965	5,4965	5,4965	5,4965	5,4965	5,4965	5,4965	5,4965	5,4965	5,4965	5,4965
time processing	Insert	00:14.989269	00:15.163563	00:14.478481	00:14.146725	00:15.022632	00:15.128661	00:13.639867	00:15.545360	00:13.194747	00:13.571628	00:17.677513	00:14.389833	00:14.837120	00:13.954089	00:14.039283	00:14.829606	00:15.008499	00:14.829606	00:15.008499
	Update	00:49.509872	00:21.520270	00:22.752126	00:22.510033	00:23.728714	00:23.654100	00:50.746861	00:58.252042	00:25.367584	00:29.569732	00:27.426164	00:22.841241	00:20.351274	00:56.559946	00:56.530541	00:27.207213	00:25.776772	00:27.207213	00:25.776772
	Delete	00:37.493929	00:15.992814	00:16.145049	00:17.160087	00:18.126900	00:17.871596	00:37.940560	00:40.734276	00:17.938380	00:18.196805	00:18.977805	00:18.182327	00:16.465309	00:37.722899	00:39.535942	00:18.364344	00:17.055967	00:18.364344	00:17.055967
	Restore	00:43.091336	00:17.521696	00:18.737257	00:18.135006	00:17.960300	00:16.895810	00:45.450542	00:44.713059	00:18.732620	00:16.210235	00:19.417701	00:17.696572	00:17.142672	00:44.063914	00:42.893921	00:18.220083	00:17.669175	00:18.220083	00:17.669175
	Select_obj	01:21.003012	00:28.996935	00:29.4436787	00:28.719686	00:28.745750	00:25.102863	01:18.947076	01:19.767744	00:28.442692	00:29.019777	00:28.180638	00:28.324229	00:25.108355	01:19.683865	01:18.766930	00:28.912447	00:28.703291	00:28.912447	00:28.703291
	Select_db	00:04.010253	00:04.006039	00:03.887190	00:03.851006	00:03.904664	00:03.240508	00:03.954560	00:03.935130	00:03.925098	00:03.944205	00:03.964011	00:03.948297	00:03.245451	00:03.780117	00:03.929227	00:03.780614	00:03.943188	00:03.780614	00:03.943188
segment size (kB)	Temporal_table	3072	3072	3072	3072	3072	3072	3072	3072	3072	3072	3072	3072	3072	3072	3072	3072	3072	3072	
	Tab1	1024	1024	1024	1024	1024	1024	1024	1024	1024	1024	1024	1024	1024	1024	1024	1024	1024	1024	
	N_tab	320	320	320	320	320	320	320	320	320	320	320	320	320	320	320	320	320	320	
	C_tab	448	448	448	448	448	448	448	448	448	448	448	448	448	448	448	448	448	448	448
	D_tab	128	128	128	128	128	128	128	128	128	128	128	128	128	128	128	128	128	128	128
	Index	0	2048	2048	2048	2048	2048	3072	3072	2048	2048	3072	3072	2048	3072	3072	640	3072	3072	

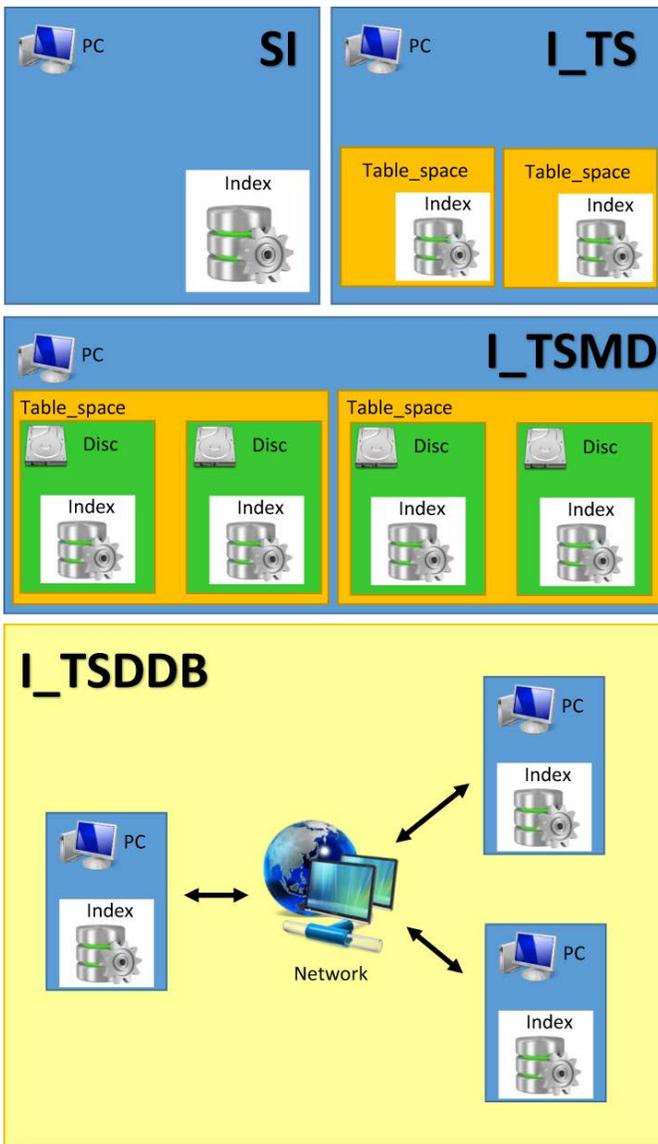


Fig. 18. Index localization

VII. CONCLUSION

Standard temporal system deals with the whole objects as the granularity, which is used, if all temporal columns are updated at the same time. Effective managing temporal data

can be very useful for decision making, analyses, process optimization and can be used in any area – industry, communication systems, medicine, and transport systems and so on. However, a lot of systems manage data with the different character of the granularity. Temporal data management used today is based on object level (one row represents the whole state at defined time point or time interval), which often does not cover the complexity of the data management in time validity. A significant aspect is just performance reflected to processing time and size of the structure. Extended column level temporal system is based on the column (attribute) level developed by us in the past, the whole state is created by the grouping the properties and states of the attributes. Model described in this paper significantly improves the usability thanks to simplifying the structure and manipulation system.

This paper deals with the principles, characteristics and describes implementation methods to provide complex temporal data management. The developed system is compared based on the index performance. In the future research, the system will be extended by the distribution and parallelism management.

REFERENCES

- [1] C. J. Date, "Date on Database". Apress, 2006.
- [2] C. J. Date, H. Darwen, and N. A. Lorentzos, "Temporal data and the relational model", Morgan Kaufmann, 2003.
- [3] Ch. S. Jensen and R. T. Snodgrass, "Temporally Enhanced Database Design", MIT Press, 2000.
- [4] T. Johnston and R. Weis, "Managing Time in Relational Databases", Morgan Kaufmann, 2010.
- [5] M. Kvet and K. Matiaško, "Transaction Management", 2014. CISTI, Barcelona, pp.868-873.
- [6] M. Kvet, K. Matiaško, M. Kvet, "Complex time management in databases", In Central European Journal of Computer Science, Volume 4, Issue 4, 2014, pp. 269-284.
- [7] M. Kvet and M. Vajsová. "Transaction Management in Fully Temporal System", 2014. UkSim, Pisa, pp. 147-152.
- [8] J. Maté, "Transformation of Relational Databases to Transaction-Time Temporal Databases", in ECBS-EERC, 2011, pp. 27-34.
- [9] M. T. Oszu and P.Valduriez, "Principles of Distributed Database Systems", Springer, 2011.

Alternative approach to enable RTSP-based services with dynamic Quality of Service over 4G LTE mobile networks

Andrei Rusan, Radu VasIU

Abstract—In this article, we present an alternative approach to enable Real Time Streaming Protocol (RTSP) based streaming services with Quality of Service (QoS) over 4th generation (4G) Long Term Evolution (LTE) mobile networks. Specifically, we present a proof-of-concept architecture that is more cost effective, easier and more simple to implement than existing proposals using the IP Multimedia Subsystem (IMS), with or without native RTSP support or a Session Initiation Protocol (SIP)-RTSP hybrid implementation. The solution we propose does not require any extension to existing protocols or standards, making use of protocols that are already available and in general use and maintaining compatibility with existing solutions and applications. We discuss the advantage and potential of such an implementation to be used not only in laboratories, proof of concepts or test and demo networks, but also in specialized environments, in commercial (unicast) scenarios and in enabling commercial over the top applications for 3rd party streaming services providers. We exemplify using video streaming scenarios, but this applies to audio streaming as well.

Keywords—QoS streaming, RTSP streaming in 4G LTE, unicast video streaming with QoS.

I. INTRODUCTION

MULTIMEDIA streaming services have evolved into an integral part of our daily lives and technology environment. We take it for granted to be able to access video and audio sources anytime and from anywhere – and our expectations in terms of service quality and perceptual quality are growing with every new generation of devices and networks that we have access to. Wireless networks in particular strive to incorporate new technologies, as wireless represents a more challenging environment for streaming [1]-[4].

It is thus natural that users expect perfect quality of streaming video and audio to be available on their mobile fourth generation (4G) Long Term Evolution (LTE) user equipments (UEs), with the obvious expectation that everything just works, no matter if it is Content on Demand (CoD), IP Television (IPTV) or real time streaming (RTS). Streaming clients are either natively incorporated or easy to install nowadays on all available UEs (be it handsets, laptops

or LTE enabled tablets), driving a fast and broad adoption of streaming technologies and services - and their associated use cases.

Modern mobile communication networks, like 4G LTE, provide enhanced technologies, which enable higher spectral efficiencies and allow for higher data rates and cell capacities. High Definition (HD) streaming services that offer outstanding Quality of Experience (QoE) are now available and affordable also over mobile networks [5]. As these mobile networks get faster and more reliable, they attract more users that make increased use of such services, leading to an exponential increase in traffic. This further drives the need to develop new technologies and to provide enhanced Quality of Service (QoS) for a broader range of streaming services, so that the expected QoE can be achieved while making best use of the available network resources.

In 4G LTE networks (3rd Generation Partnership Project (3GPP) Release 9 and later) [6] this is achieved by enforcing QoS for specific service classes, in accordance with the 3GPP QoS concept and architecture [7] – enabling applications such as Voice over IP (VoIP) with HD voice, broadcast and unicast streaming of HD audio and video (IPTV, real time streaming, etc.). This is of high importance especially over the radio interface, where the main bottleneck is observed – with well-known and also new packet scheduling algorithms [8], [9] driving optimization of resource allocation.

II. CURRENT STATE OF MULTIMEDIA SERVICES OVER 4G LTE

For multimedia services using Session Initiation Protocol (SIP) [10] based session control, 3GPP introduced the IP Multimedia Subsystem (IMS) since Release 5 [11]. It enables a wide range of multimedia services, with Voice over LTE (VoLTE) being the most prominent one [12].

Since 3GPP Release 6, the multimedia broadcast/multicast service (MBMS) was introduced, and it was updated in Release 9 for LTE, where it is called evolved MBMS (eMBMS). It defines a point-to-multipoint service in which data are transmitted from a single source entity to multiple recipients, making it useful for broadcast scenarios where many users access the same data stream(s) in a concentrated area – for example IP Television (IPTV) scenarios or live streaming of events. eMBMS is not attractive for unicast scenarios, with no benefits that would justify eMBMS complexity when using it for such scenarios.

Andrei Rusan is a PhD student, with the Politehnica University of Timisoara, Romania. (e-mail: arusan@gmx.com).

Radu VasIU is professor, with the Politehnica University of Timisoara, Department of Communications. (phone: +40.722.516555; e-mail: radu.vasiu@cm.upt.ro).

Unicast streaming scenarios are specified by 3GPP under Packet Switched Streaming (PSS) services [13] - the detailed architecture to enable PSS with QoS and charging functionality being IMS based and using a PSS Adapter to translate Real Time Streaming Protocol (RTSP) into SIP and vice versa [13], [14]. The main rationale for this is SIP support in IMS (which already supports most multimedia functions and services, interacting with the proper entities for QoS, charging, etc.), coupled with the fact that streaming services are widely supported over the Internet Engineering Task Force's (IETF) suite of protocols for real time streaming. In this paper, we will refer to these services as "RTSP streaming" - with RTSP being the application layer session control protocol common to all these streaming services - and will go into more detail on the rest of the protocols further in this article. It is worth noting that unicast streaming is commonly also delivered as progressive download over the Hyper Text Transfer Protocol (HTTP) [14]. This is most often the case for non real-time scenarios, where higher delay and longer buffering are not bothering users. This brings the disadvantage of making these applications more difficult to differentiate from other best-effort (BE) services - and hence more difficult to monetize [15]. It is out of the scope of this paper to dive into this debate: while we recognize that each solution has its advantages and disadvantages, we will concentrate on RTSP streaming unicast scenarios in this paper.

Unicast streaming will likely stay high in demand even when eMBMS will be deployed - be it for Content on Demand (CoD) (with its most famous flavor of Video on Demand (VoD)) applications or for other specialized ones like security and surveillance. In such scenarios, each user accesses different streaming data at random points in time and the chance of many users being the recipients of the same data is slim, taking away the advantage of using broadcast/multicast infrastructure for such scenarios [16].

It is important to note that RTSP streaming is widely supported by most of the commercially available equipments and applications in a standardized manner, be it personal or industrial ones. This makes over the top (OTT) and 3rd party deployments of such solutions and services attractive and easy to integrate, with a fair chance of monetizing these if they can be differentiated in terms of quality and user experience. With the benefits of increased QoE (building on guaranteed QoS), there is an increased potential of making such applications more and more interesting also for high-end users (as it offers a fair price-quality combination) and service providers (which can now financially benefit from allocating differentiated resources for specific services that they, or 3rd party service providers, offer).

III. CURRENT APPROACHES FOR UNICAST RTSP STREAMING SERVICES ON 4G LTE

A. IETF protocols for streaming: RTSP, RTP, RTCP, SDP

As mentioned earlier, we refer to "RTSP streaming" when talking about streaming services that are supported over the IETF's suite of protocols for real time streaming.

RTSP, specified by IETF Request for Comments (RFC) 2326 [17] (with RTSP 2.0 available as draft [18]) is the application layer session control protocol used for real time streaming. It is widely adopted within the industry, with a wide array of different software and hardware, stand-alone or embedded clients and servers, that support RTSP. RTSP is a bidirectional protocol that can be transported over either UDP or TCP and is used to control the session and playback, from querying for session description (using the DESCRIBE method) and session setup (using the SETUP method), to conveying PLAY, PAUSE, STOP commands or ending a session and freeing the associated connection (TEARDOWN method).

Session description in conjunction with RTSP is commonly provided by the Session Description Protocol (SDP), specified by IETF RFC 4566 [19] (developed out of RFC 2327 [20]). SDP messages are sent in response to RTSP DESCRIBE commands. A SDP message carries details about the multimedia stream that is targeted, its structure (audio and video sub-streams) and properties (bit rate, resolution, frame rate etc.).

Transport is provided by the Real-Time Transport Protocol (RTP) and augmented by the Real-Time Control Protocol (RTCP), both specified by IETF RFC 3550 [21]. RTP provides end-to-end network transport functions suitable for applications transmitting real-time data, such as audio, video or simulation data. RTCP allows monitoring of the data delivery and provides minimal control and identification functionality.

This combination of protocols is supported and implemented in a standardized way for most of the existing video streaming solutions, from software client and server implementations (IPTV, live streaming or VoD) to dedicated hardware for video surveillance (accessible as audio-video streaming sources from either concentration nodes, transcoders or often directly out of the IP cameras themselves). For the sake of simplicity, we will refer to any such source as Streaming Server (SS).

B. Current proposals

While the number of IMS deployments is growing, it is important to note that the main focus for these deployments is voice calling and video calling, as well as other SIP based multimedia services. RTSP streaming services support is not common for any such deployments and also seldom to be demonstrated or tested. This is because RTSP services are not commonly and natively supported by the IMS. Equally, proposed solutions to harmonize RTSP and SIP (and, by extension, the IMS) have also not seen wide industry adoption to this point.

One proposal, as outlined by 3GPP PSS, is to implement a SIP-RTSP translator that is integrated into the IMS. This adds complexity and calls for SIP extensions to support RTSP commands. SIP extensions for media control have been outlined by the IETF SIP working group [22] but are not commonly implemented. Developing such a solution has not been a priority for most equipment providers and is not commonly found in IMSs deployed.

Another proposal was to directly integrate SIP to RTSP (also based on [22]), but due to the major impact on existing streaming clients and servers this also was not a widely adopted solution, with currently very few RTSP implementations offering such functionality.

Yet another proposal was to integrate RTSP support directly into the IMS (along with support for HTTP) and to blend these separate services into one [23], [24]. This proposal has ignited discussion and research, but as of now RTSP support in commercially deployed IMSs is rare, even though the idea existed already for 3G networks [24].

C. The complex IMS

An important detail is the fact that commercial implementations of the IMS are of high complexity and have a high associated cost, not only for acquisition, but also for installation, provisioning and operation. The high complexity of the IMS, paired with the fact that LTE voice services using the IMS (e.g. VoLTE) represent a new technology to be rolled out and optimized, mean that for LTE players (service providers, equipment manufacturers etc.) there is now little focus left to extend IMS capabilities and services offered at this point. The high associated cost, in turn, makes it much more difficult to equip laboratory or demo networks with an IMS Core Network (CN) for integration, testing and optimization of RTSP streaming services. And last, but not least, there is a learning curve for the transition to using IMS and expertise is still below the high demand in this area.

IV. CONSIDERING AN ALTERNATIVE

Considering these aspects that do not favor and promote the adoption of streaming services currently, we started looking for a new approach to ease and accelerate the implementation of RTSP streaming over LTE, keeping in mind that any solution to be considered should provide dynamic QoS and (differentiated) charging functionality at least. This new approach should be kept as simple as possible and should make use of existing network elements with standardized functionality. It should minimize the amount of new developments needed and should not break any existing standards or require complex extensions, which would prevent us from reaching our target in terms of simplicity and (almost) immediate availability of a working solution.

When defining our new approach for streaming with QoS over LTE, we targeted support for RTSP sources like the ones mentioned earlier as SS (last paragraph of *III. A.*). For the test-bed we built in order to demonstrate the approach described in this paper, we used the Alcatel-Lucent 5910 Mobile Streaming Server (MSS), which is an industrial grade 3GPP/3GPP2 standards compliant multimedia server.

We also sought simplicity in the form of existing and widely used software, when choosing the streaming client to be used in our test bed. Since our test UE was a LTE dongle connected to a laptop running Microsoft Windows, we chose to use the QuickTime Player as a streaming client.

We focused our demonstration around video streaming, considering that a scenario with one video stream and a

concurrent audio stream sufficiently proves the feasibility of our approach, as audio streaming is just a simplified use case of the tested scenario.

V. OVERALL STREAMING SCENARIO DESCRIPTION

When streaming video content, a video client requests the audio-video stream from a SS, so that it can play it back on a screen. No matter where the source material originates (live content from a camera, encoded content stored on a VoD server, etc.) and what network or technology it is delivered over (wired, fixed or mobile wireless), the client needs to receive the stream packets in time so that it can decode the stream and present it on the screen for playback.

Video streams are almost never constant bit rate. The first reason is the fact that the original video content that serves as input is not constant. The second reason is related to how encoding is performed with modern codecs and encoding technologies, where only some frames hold full information for a particular scene (key frames or I-frames); other frames only hold delta-data: changes from the content of previous frames (P-frames) or previous and following frames (B-frames). Additionally, the transport network injects delays (often of variable duration, as they are dependent on the network's load) which lead to jitter or out-of-order arrival of RTP packets at the receiving side. In order to smooth these out, the client works with de-jitter and other types of buffers, which are often adaptable in order to best match the conditions of the network and counteract the effects of these less than ideal conditions. These buffers work well for de-jittering and reordering plus preventing early discarding of late packets, but they also come with an obvious drawback: the bigger the buffer (more efficient in correcting any problems), the higher the playback delay (resulting from the amount of data in the buffer). This means that a balance is needed in order to minimize discarded or missing packets and also provide reasonable and acceptable playback delay for the targeted applications [8].

There is not much that can be done about the variable nature of the content and the resulting variable bit rate of the video stream (except for what is already employed in encoders and is out of the scope of this paper). Instead, we are interested in enforcing minimal impairments on the transport network, by using QoS, so that buffers can be kept at a minimum. There are specific QoS mechanisms and associated protocols and technologies that accomplish this, for both wired and wireless networks.

In 4G LTE networks, powerful and flexible QoS mechanisms are available, with various QoS classes that strive to guarantee reliable and efficient transport of data over both the CN and the Radio Access Network (RAN) part of the networks.

The important detail in order to be able to benefit from the advanced QoS capabilities of 4G LTE networks, is to be able to trigger the appropriate QoS for video streaming services when usage of such a service is started – and to free resources when their usage ends.

In accordance with the 3GPP standards and architecture,

LTE networks provide Policy and Charging Control (PCC) as part of the LTE Evolved Packet Core (EPC). The PCC architecture includes several software nodes, like the Policy and Charging Rules Function (PCRF), the Proxy Call Session Control Function (P-CSCF) and the Policy and Charging Enforcement Function (PCEF). Combined, the elements of the PCC provide access, resource, and QoS control.

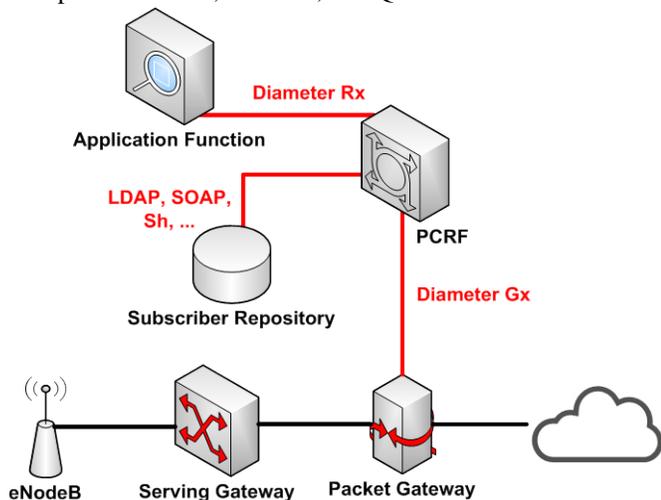


Figure 1: PCRF system architecture

One of the functions of the PCRF (in a QoS context) is to create and destroy dedicated bearer packet data protocol (PDP) contexts (and thus radio bearers) in response to requests coming from an Application Function (AF). The AF provides application session related information, e.g. based on SIP and SDP. The idea is to have one AF per operator-provide service – for example, an AF is found in IMS implementations, where it generates requests in the case of voice sessions.

Considering the above, it is clear to see that an AF is a key functionality that we need to ensure as part of our simplified approach to video streaming with QoS in LTE networks – the difference to the AF in the IMS being that we need an AF that provides session related information based on RTSP and SDP. For the sake of completeness, let’s mention that the AF and PCRF communicate over a Diameter Rx interface (see Fig. 1).

VI. PROPOSED APPROACH FOR STREAMING SCENARIO

Synthesizing the above, for a simplified and working solution to video streaming one needs one or more network elements (or a subsystem) that:

- support RTSP functionality (that extends to SDP, RTP, RTCP)
- interface over Diameter Rx with the PCRF (e.g. provide an AF)
- and can thus trigger the appropriate QoS (including transparent dedicated bearer setup on the RAN part) and charging for video streaming services
- is preferably readily available and with a much lower cost and complexity than an IMS

The network element that we identified as a good fit and that allowed for demonstrating our approach and simplified architecture, is the Session Border Controller (SBC). Our

proposal enables video streaming services with QoS over LTE while avoiding any SIP-to-RTSP translation, any developments on the IMS and any changes to existing streaming clients/servers.

SBCs are common network elements found in most modern networks. As the name suggests, they are usually deployed at the border of the networks, for interconnection and to augment security and control over session based services. Such services are firstly SIP based voice services, but RTSP based streaming services are also commonly supported by SBC functionality – commercially available SBC implementations commonly offer Application Layer Gateways (ALGs) for both SIP (SIP-ALG) and RTSP (RTSP-ALG). In both cases, the SBC acts as a proxy server for the service, with network address translation (NAT) (and usually also corresponding port address translation (PAT)) performed as part of the scenario. SBCs also commonly export a Diameter Rx interface – combined with the session awareness through the RTSP-ALG, this makes them suitable embodiments for a RTSP streaming AF that communicates with the PCRF to provide QoS management and control for RTSP streaming sessions.

It is important to note that SBCs are mature and commercial grade solutions. While implementations vary and SBCs come in all shapes and sizes, they are usually modular (at least from a software point of view, but sometimes also from a hardware point of view). This means that one can find an “integrated SBC” (a single standalone SBC device that is placed in front of existing equipment in the path of all signaling and media traffic on an interface) that can be used as a RTSP proxy and to provide an AF for the RTSP streaming service. Such a device, with only the needed functionality licensed, costs a fraction of an IMS solution. As a bonus, it provides robust functionality and can augment security of the network, all while typically supporting a very high number of users (tens or hundreds of thousands) which also makes it usable in more than just laboratory or demonstrational networks.

For our demonstrational test bed we used an Alcatel-Lucent 5780 Dynamic Services Controller (DSC) as the PCRF, and an Acme Packet 4250 Net-Net Session Director (SD) SBC. While the ALu 5780 DSC is part of the LTE CN anyway, it is worth mentioning that the Acme 4250 SD adds little to the footprint of the core network, as it is a single blade server of 1U height; its cost (depending on licensing and configuration, of course) is also close to the expected range for a single blade network element, making it easier and more likely to be part of a laboratory LTE network or a demonstrational network. Such environments will not have a high number of users, but it is worth noting that the solution can be scaled without any problems also to applications where a high number of users need to be supported, as our approach proposes the use of commercial grade SBCs that offer both robust performance and high capacity when deployed in the field.

A simplified diagram of the proposed architecture is presented in Fig. 2, to illustrate the solution:

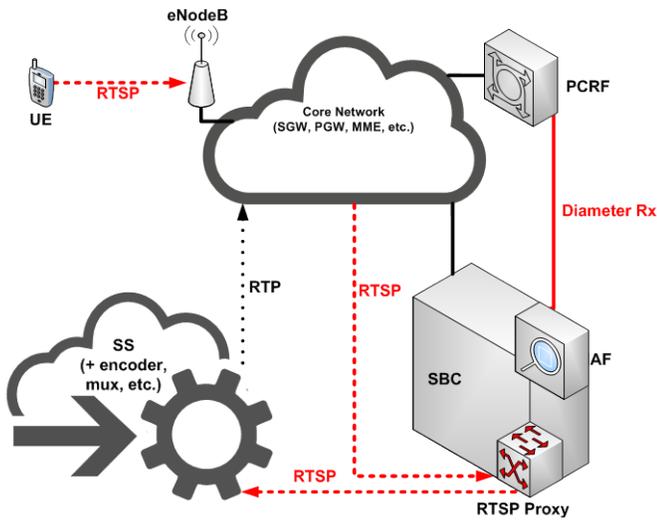


Figure 2: Proposed solution

When considering the approach we propose, unicast RTSP-based streaming is enabled with dynamic QoS in 4G LTE networks, using solely available software and hardware and without any modification or new developments needed to existing solutions and standards.

A user starts a RTSP streaming session on his UE (smartphone, PC, tablet, etc.), using a standards compliant RTSP client. The server can be a CoD server, a live streamer, specialized hardware that has RT SP streaming capabilities, etc. The initial request reaches the RTSP proxy (implemented in the SBC), which provides RTSP-ALG functionality and will forward the request to the RTSP server. The same SBC provides and AF, which will extract session details from the SDP provided by the RTSP server – as the SDP will be forwarded by the RTSP proxy. This AF will then communicate with the PCRF and, if the user has the appropriate subscription level and can perform video streaming with guaranteed QoS, the necessary resources will be reserved and allocated once the session is set up and the streaming begins – this means allocating the appropriate PDP contexts, dedicated bearers on the radio interface, data bearers in the CN, etc.). The AF (functionality provided by the SBC) will identify and interpret each stage of the process, as all messaging flows through the SBC’s RTSP-ALG, and communicate the needed requests to the PCRF over the Rx interface. Charging can also be performed now in a differentiated way for this service, starting with the moment the service usage begins. While the video streaming session is ongoing, all involved network elements will try to observe the QoS constraints of the stream and ensure that priority is correctly considered and resources allocated, to minimize impact and obey restrictions on delay, packet loss and jitter. This means correct prioritization and scaling of queues, buffers and adjusting scheduling especially over the air interface. Once the session is stopped (by either client or server, or ends because of other events), resources are freed and charging is stopped in a similar way to how it was set up.

The call flow in Fig. 3 aims to illustrate the first part of this process, as an RTSP streaming scenario is started:

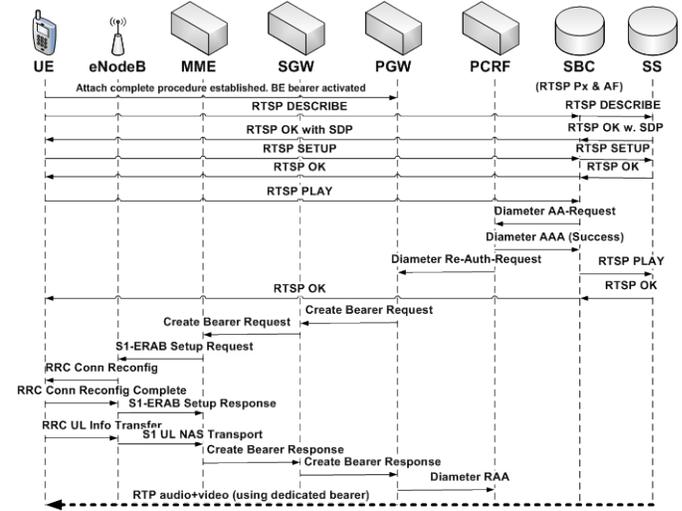


Figure 3: RTSP streaming session start

VII. CONCLUSION

Building on existing elements, this alternative approach proposed for enabling RTSP streaming in 4G LTE networks provides a practical, simplified and economic approach that is immediately available for implementation, testing and even larger scale usage. We performed initial functional testing using such a setup and could confirm its functionality and usefulness in testing and demonstrating RTSP streaming with dynamic QoS over 4G LTE.

While specifications of all involved network elements promise high capacity that would be usable in commercial environments, such testing and validation could represent a next step in validating this solution built with existing commercially available elements.

Last, but not least, the setup allows for further performance testing of the LTE QoS functionality implementation in existing networks - be it laboratory, demo or even live networks. This was the initial need and main driver in developing this approach, representing its true value with the numerous use cases that it now enables.

REFERENCES

- [1] A. Majumdar, D. G. Sachs, I. V. Kozintsev, K. Ramchandran, M. M. Yeung, "Multicast and Unicast Real-Time Video Streaming Over Wireless LANs", in *IEEE Transactions on circuits and systems for video technology*, vol. 12, no. 6, pp. 524-534, June 2002.
- [2] D. G. Sachs, I. Kozintsev, M. Yeung, D. L. Jones, "Hybrid ARQ for robust video streaming over wireless LANs," in *Proc. Inform. Technol.: Coding and Computing, ITCC2001*, 2001, pp. 317–321.
- [3] A. Kassar, A. Schorr, L. Chen, C. Niedermeier, C. Meyer, M. Helbing, M. Talanda: "Multimedia Communication in Policy based Heterogeneous Wireless Networks", in *IEEE Vehicular Technology Conference VTC2004*, Milan, Italy, May 2004.
- [4] Y. Liu, F. Li, L. Guo, B. Shen, S. Chen, Y. Lan, "Measurement and Analysis of an Internet Streaming Service to Mobile Devices," in *IEEE Transactions on Parallel and Distributed Systems*, pp. 2240-2250, 2013.

- [5] K. Piamrat, C. Viho, J. Bonnin, A. Ksentini, "Quality of Experience Measurements for Video Streaming over Wireless Networks", in *IEEE 6th International Conference on Information Technology: New Generations 2009, ITNG '09*, Las Vegas, NV, pp. 1184-1189, April 2009.
- [6] *3GPP Specification Release 9 – LTE*, 3rd Generation Partnership Project [Online]. Available: <http://www.3gpp.org/specifications/releases/71-release-9>.
- [7] ETSI TS 123 107 v12.0.0, "Digital cellular telecommunications system (Phase 2+); Universal Mobile Telecommunications System (UMTS); LTE; Quality of Service (QoS) concept and architecture (3GPP TS 23.107 version 12.0.0 Release 12)", Sept. 2014.
- [8] H. A. M. Ramli, K. Sandrasegaran, R. Basukala, R. Patachaianand, T. S. Afrin, "Video Streaming Performance Under Well-Known Packet Scheduling Algorithms," in *International Journal of Wireless & Mobile Networks (IJWMN)*, Vol. 3, No. 1, February 2011.
- [9] H. A. M. Ramli, K. Sandrasegaran, R. Basukala, R. Patachaianand, X. Minjie, C.-C. Lin, "Resource allocation technique for video streaming applications in the LTE system," in *IEEE 19th Annual Wireless and Optical Communications Conference, WOCC*, Shanghai, pp. 1-5, 2010.
- [10] J. Rosenberg, H. Schulzrinne, G. Camarillo, A. Johnston, J. Peterson, R. Sparks, M. Handley, E. Schooler, "SIP: Session Initiation Protocol," Request for Comments (RFC) 3261, IETF Network Working Group, June 2002, [Online] Available: <https://www.ietf.org/rfc/rfc3261.txt>
- [11] ETSI TS 122 173 v12.8.0, "Digital cellular telecommunications system (Phase 2+); Universal Mobile Telecommunications System (UMTS); LTE; IP Multimedia Core Network Subsystem (IMS) Multimedia Telephony Service and supplementary services; Stage 1 (3GPP TS 22.173 version 12.8.0 Release 12)", 2015-01, [Online] Available: http://www.etsi.org/deliver/etsi_ts/122100_122199/122173/12.08.00_60/ts_122173v120800p.pdf
- [12] GSMA VoLTE Initiative, [Online] Available: <http://www.gsma.com/network2020/volte/>
- [13] TSI TS 126 237 v12.0.0, "Universal Mobile Telecommunications System (UMTS); LTE; IP Multimedia Subsystem (IMS) based Packet Switch Streaming (PSS) and Multimedia Broadcast/Multicast Service (MBMS) User Service; Protocols (3GPP TS 26.237 version 12.1.0 Release 12)", 2014-10.
- [14] F. Gabin, M. Kampmann, T. Lohmar, C. Priddle, "3GPP Mobile Multimedia Streaming Standards [Standards in a Nutshell]", in *IEEE Signal Processing Magazine*, 2010.
- [15] J. Erman, A. Gerber, K. K. Ramakrishnan, S. Sen, O. Spatscheck, "Over the Top Video: The Gorilla in Cellular Networks," in *Proc. ACM SIGCOMM Conf. Internet Measurement Conf. (IMC)*, 2011.
- [16] Motorola Whitepaper, "Opportunity and impact of video on LTE networks", 2009, Available: http://www.lteportal.com/Files/MarketSpace/Download/235_Motorola_LTEVideoImpactWhitePaper5_13_2009.pdf?PHPSESSID=89a7f201ba524b4932578ea159114c79
- [17] H. Schulzrinne, A. Rao, R. Lanphier, "Real Time Streaming Protocol," Request for Comments (RFC) 2326, IETF Network Working Group, April 1998, [Online] Available: <https://tools.ietf.org/html/rfc2326>
- [18] H. Schulzrinne, A. Rao, R. Lanphier, M. Westerlund, M. Stiemerling, "Real Time Streaming Protocol 2.0 (RTSP)," Draft IETF RFC2326-bis-40, IETF MMusic Working Group, February 10, 2014, [Online Draft] Available: <https://tools.ietf.org/html/draft-ietf-mmusic-rfc2326bis-40>
- [19] M. Handley, V. Jacobson, C. Perkins, "SDP: Session Description Protocol," Request for Comments (RFC) 4566, IETF Network Working Group, July 2006, [Online] Available: <https://tools.ietf.org/html/rfc4566>
- [20] M. Handley, V. Jacobson, "SDP: Session Description Protocol," Request for Comments (RFC) 2327, IETF Network Working Group, April 1998, [Online] Available: <https://www.ietf.org/rfc/rfc2327.txt>
- [21] H. Schulzrinne, S. Casner, R. Frederick, V. Jacobson, "RTP: A Transport Protocol for Real-Time Applications", Request for Comments (RFC) 3550, IETF Network Working Group, July 2003, [Online] Available: <https://www.ietf.org/rfc/rfc3550.txt>
- [22] S. Sivasothy, G. M. Lee, N. Crespi, (March 4, 2009) "SIP extensions for media control", IETF SIP Working Group [Online Internet Draft] Available: <http://www.ietf.org/staging/draft-siva-sip-media-00.txt>
- [23] S. Q. Khan, R. Gaglianella, M. Luna, "Experiences with Blending HTTP, RTSP and IMS [IP Multimedia Systems (IMS) Infrastructure and Services]" in *IEEE Communications Magazine*, vol. 45 issue 3, pp.122-128, March 2007.
- [24] H. Montes, G. Gomez, R. Cuny, "Deployment of IP Multimedia Streaming Services in Third-Generation Mobile Networks," in *IEEE Wireless Communications*, 2002.

Control of interferograms image of deformed object samples by non destructive control as optical method

R. DAIRA^{a*}

Abstract—The electronic speckle interferometry technique used to measure the deformations of scatterers process is based on the subtraction of interference patterns. A speckle image is first recorded before deformation of the object in the RAM of a computer, after a second deflection. The square of the difference between two images showing correlation fringes observable in real time directly on monitor. The interpretation these fringes to determine the deformation.

In this paper, we present experimental results of deformation out of the plane of two samples in aluminum, electronic boards and stainless steel.

Keywords— Optical method; CND Control; Deformation.

I. INTRODUCTION

Speckle interferometry technique is usually used to measure the displacements and deformations of rough surfaces. Despite its importance this technique has not seen a great adoption by researchers and industry. Mechanical stability, the need for photographic processing and the difficulty of interpretation fringes are behind this lack of admission [1]. For these reasons, it is clear that the efforts of researchers directs to operate new systems to replace circles holographic recordings. The idea is to use television systems for the detection and treatment of figures speckles. These techniques are generally known as electronic speckle interferometry name (Electronic speckle pattern interferometry [ESPI]). The major feature of the method is that it provides exposure fringes correlations in real time on the monitor. The process consists in recording the interference corresponding to the object before and after deformation camera. A video system is used to generate the correlation fringes that correspond exactly to the movement of the object between the two exposures [1-4]. In this work, a theoretical noticed the ESPI technique is presented with some experimental examples followed by discussion results.

^aDepartment of sciences of mater, University August 20, 1955 of skikda, Road of El Haddeik LP 26, Physico Chemistry of Surfaces and interfaces Research Laboratory of Skikda (LRPCSI), Algeria
daira_radouane@yahoo.fr

II. PRINCIPE OF FORMATION OF FRINGES CORRELATION

All points of the object illuminated by a laser, are consisted and they send out vibrations of the retina capable of interfering.

Have interference from these diffraction patterns that produce this granular appearance called speckle.

The speckle pattern in the image plane resulting from the addition of two fields coherent speckles. Whether: $A_1 = a_1 \cdot \exp(i\phi_1)$ and $A_2 = a_2 \cdot \exp(i\phi_2)$, the complex amplitudes of the two wavefronts, where a_1 , a_2 et ϕ_1 , ϕ_2 are random variables corresponding respectively to the amplitudes and phases of the two speckle fields in the image plane.

The intensity distribution at a point in the image plane is given by:

$$I_1(x, y) = a_1^2 + a_2^2 + 2a_1a_2 \cos \phi \quad (1)$$

Where : $\phi = \phi_1 - \phi_2$

When the object sudden a deformation, the corresponding intensity distribution is given by :

$$I_2(x, y) = a_1^2 + a_2^2 + 2a_1a_2 \cos(\phi + \Delta\phi) \quad (2)$$

And $\Delta\phi$ is the phase difference introduced by the deformation.

The Correlation coefficient $\rho(\Delta\phi)$ between two intensity distributions are defined by [4] :

$$\rho(\Delta\phi) = \frac{\langle I_1 I_2 \rangle - \langle I_1 \rangle \langle I_2 \rangle}{\left(\langle I_1^2 \rangle - \langle I_1 \rangle^2 \right)^{\frac{1}{2}} \left(\langle I_2^2 \rangle - \langle I_2 \rangle^2 \right)^{\frac{1}{2}}} \quad (3)$$

Where : $\langle \rangle$ denotes an average over a set of points in the scattering field. Substituting equations (1) and (2) in equation (3), and assuming that a_1^2 , a_2^2 et ϕ are independent variables and their averages are computed separately and $\langle I_1 \rangle = \langle I_2 \rangle = \langle I \rangle$.

For this must take into account that:

$$\langle \cos \phi \rangle = \langle \cos(\phi + \Delta\phi) \rangle = 0 \text{ et } \langle I^2 \rangle = 2\langle I \rangle^2$$

Then $\rho(\Delta\phi)$ is expressed as follows:

$$\rho(\Delta\phi) = \frac{1 + \cos \Delta\phi}{2} \quad (4)$$

Thus, the correlation between the intensities I_1, I_2 are unitary (i.e. $\rho(\Delta\phi)=1$). So the figures are completely correlated speckle when:

$$\Delta\phi = 2\pi n \quad (5)$$

There's no correlation, (i.e. $\rho(\Delta\phi)=0$), when :

$$\Delta\phi = (2n+1)\pi \quad (6)$$

With :

$$\Delta\phi = 2\pi \delta / \lambda = 2\pi (n \Delta L) / \lambda \equiv \omega t \quad (7)$$

From the equations (5) and (6), variations in the speckle correlation figures appear as a fringe pattern. These fringe patterns are called correlation fringes. For interferometry methods work, it is necessary that the average size of speckle grains that are larger motion components when the object is deformed.

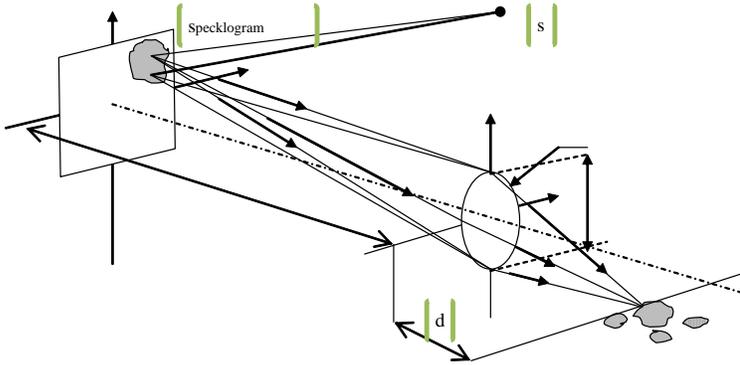


Fig. 1 point by point data analysis by speckles generated from specklogram

III. PRINCIPE OF TECHNIQUE [1-8]

Figure (2) shows the optical interferometer used in this work arrangement. Light emitted by a He-Ne laser source (30 mW, 632.8 nm), the wave front is divided into two beams by a separator cube. The first beam will illuminate a scattering object and second is used as reference. The latter is reflected on the sensor surface of a camera (CCD: 500x582 pixel) by a second beam splitter.

The scattered light from the object is collected by a photographic lens ($f = 75$ mm), and an image is formed in the plane of the sensor surface of the camera where it will interfere with the reference wave.

To measure the ESPI deformation method, a speckle image of the object wave is recorded. After deformation, the recorded image is directly subtracted from the original image, and the square of this difference is displayed on the monitor as fringe correlations.

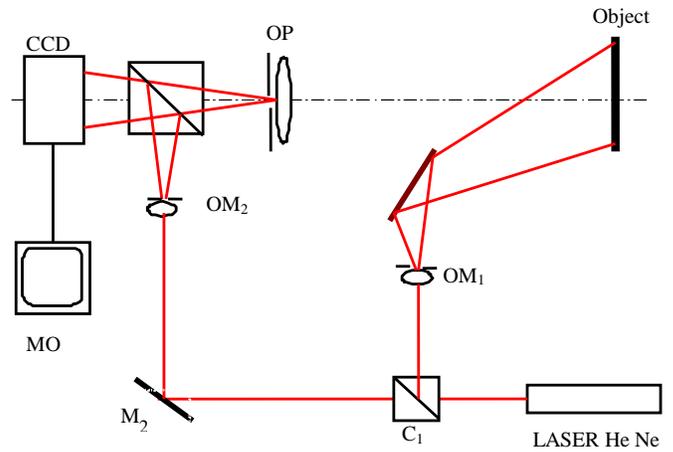


Fig. 2 Experimental setup for the measurement of displacement out of the plane

IV. TEMPORAL PHASE SHIFTING [9]

While the fringe pattern describes the surface deformation of the object, its appearance is not unique. The direction of the deformation, towards the optical setup or away, can only be detected with the determination of the phase ϕ . This is commonly done by application of phase shifting techniques. The optical setup is altered with a phase shifter in one of both beams (Fig.4). This phase-shifter (e.g. piezo element behind mirror) allows adding a known phase shift to the random phase ϕ . Several images (i.e. ≥ 3) are recorded in a temporal manner, using different known phase shifts (i.e. $\pi/2, 3\pi/2, 5\pi/2$). Subsequently, the phase angle ϕ can be calculated utilizing :

$$\phi(x, y) = \tan^{-1} \left[\frac{I_3(x, y) - I_2(x, y)}{I_1(x, y) - I_2(x, y)} \right] \quad (7)$$

Where I_n describes the intensity at x,y in accordance to the additional applied phase shift. The phase difference (after subtraction) is usually displayed in a phase map, which contains still the information of the fringe pattern and in addition the directional information of the deformation. Phase maps are unwrapped using a computer algorithm to display the deformation in e.g. color-coded plots.

ϕ is the phase change between the object and reference beams, and $\Delta\phi$ ($\Delta\phi = (2\pi/\lambda) \cdot d$) is the phase change caused by deformation 'd'.

If the signals V_1 et V_2 (the output of the camera) are proportional to the intensity of the input image, then the resultant of the subtraction signal (V_s) is given by:

$$V_s = (V_1 - V_2) \propto (I_1 - I_2) = 4\sqrt{I_0 I_r} \sin(\phi + \Delta\phi) \cdot \sin \frac{\Delta\phi}{2} \quad (8)$$

This signal has positive and negative values. Negative values are displayed as dark areas on the TV monitor. To avoid this loss of signal and get fringes, the square of the difference V_s runs. Therefore the average intensity in a given B in the monitor image is point:

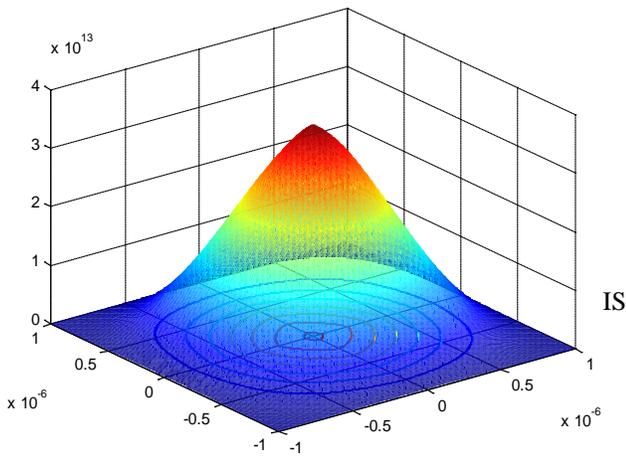
$$B = 8 \langle I_0 I_r \rangle \sin^2 \left(\frac{\Delta\phi}{2} \right) \quad (9)$$

Equation (9) is similar to that obtained in conventional interferometry where the fringes are sinusoidally dependent on the phase difference relative to the deformation.

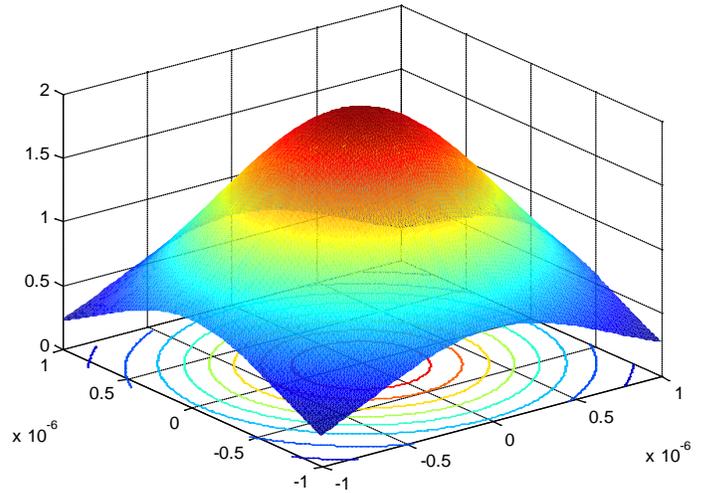
In figure, we can simulate a "speckle pattern". To do that, I propagate first a gaussian beam $b(x,y)$ on a certain distance: the propagation of a beam between planes is given by:

$$g(x, y) = b(x; y) \otimes h(x; y) \quad (10)$$

Where $h(x,y)$ is the Fresnel impulse response.



a

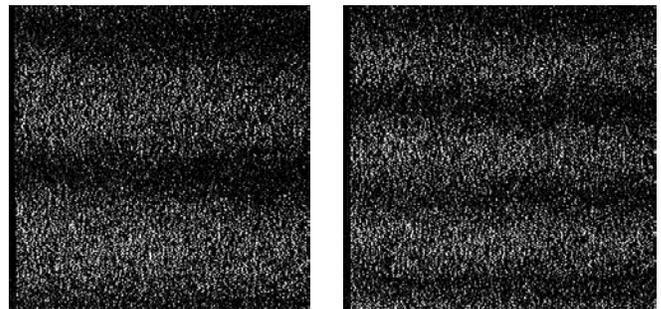


b

Fig. 3 Simulation of speckle pattern: a. Gaussian beam
b. Gaussian beam after a certain distance

V. EXPERIMENTAL RESULTS

To demonstrate the method described previously, deformation out of the plane is formed on a stainless plate of rectangular shape with dimensions $(70 \times 30 \times 5) \text{ mm}^3$. The plate attached to its lower end by a clamp, is bent at its upper end by a displacement of the micrometer screw. The object undergoes a series of strains in the range of 20 microns, and seven images of interference are recorded in succession to each state of deformation. between the interference image on the object before deformation and pictures of each strain state. Specklogramms these are none other than parralleles interference fringes with a constant spacing for each image. We can represent the deformation of the object according to the displacement:



a

b

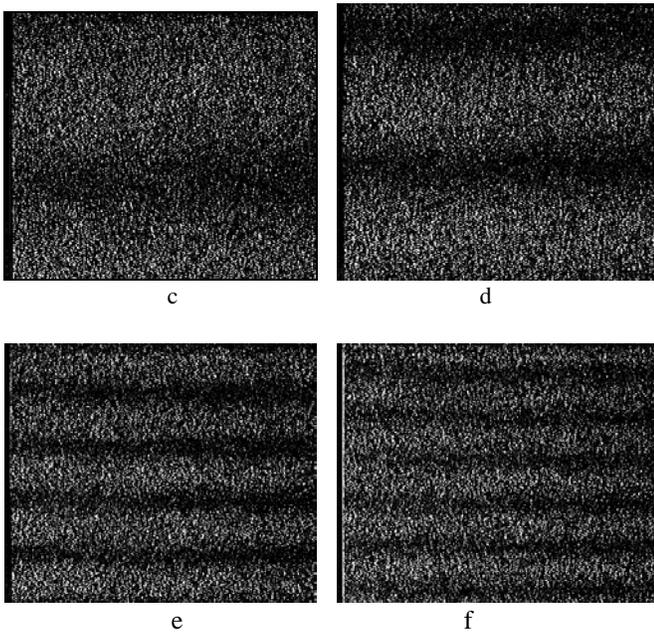


Fig. 4 Specklograms for six states of deformation of a stainless steel plate (70 x 30 x 5) mm³ obtained by ESPI method.

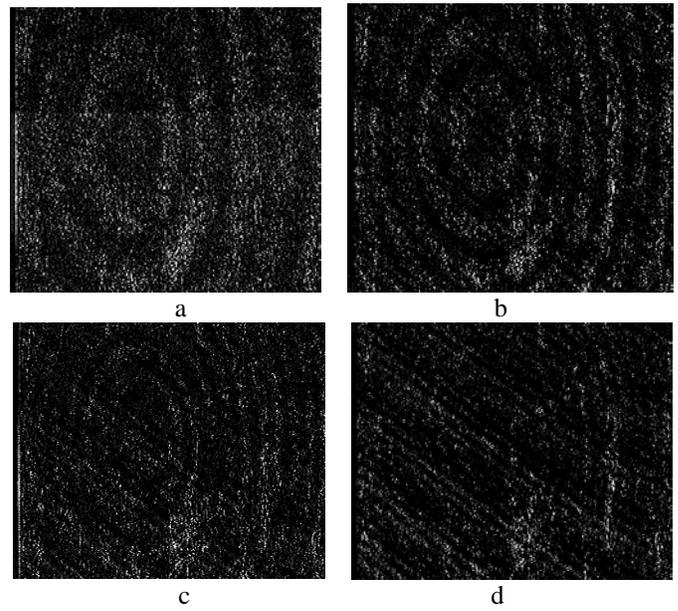


Fig. 6 Figures fringe on four states of deformation of an aluminum plate of dimensions (45 x 45 x 1) mm³.

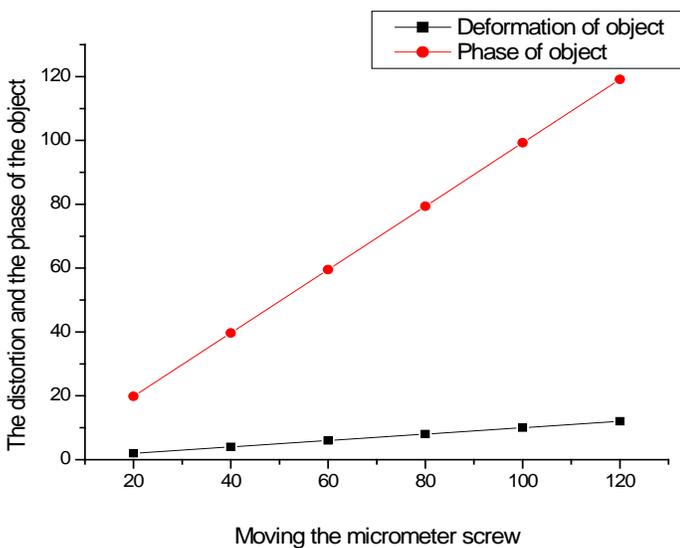


Fig. 5 Representation of deformation and phase for six states of a stainless steel plate (70 x 30 x 5) mm³ obtained by ESPI method.

The same procedure is applied to a sheet of aluminum with dimensions (45 x 45 x 1) mm³ glued to a sample holder with a circular aperture of 35 mm diameter. Is then applied perpendicular to the deformation of the center plate by the displacement of a micrometer screw. At same we can represent the deformation of the object according to the displacement:

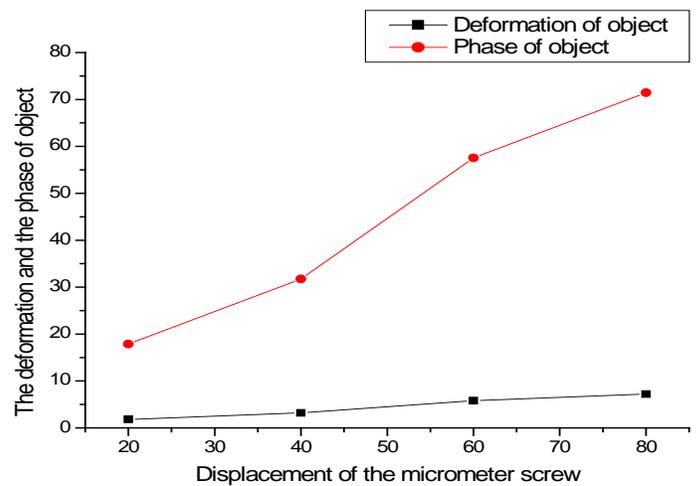


Fig. 7 Representation of deformation and phase on four states of deformation of an aluminum plate of dimensions (45 x 45 x 1) mm³.

It is noted that as the displacement 'd' increases, decreases the fringe

We also see that the visibility of these fringes decreases as the displacement increases and becomes almost zero when the displacement is of the order of magnitude of the diameter of the speckle grains. The numerical value of the deformation can be assessed by the analysis of fringes techniques (eg the technique of shifting phase) [10-12].

The same procedure is applied to a electronic boards with dimensions (40 x 40 x 1) mm³ glued to a sample holder, Is then applied perpendicular to the deformation of the center

plate by the displacement of a micrometer screw. The figure (8) shows the obtained specklogramms.

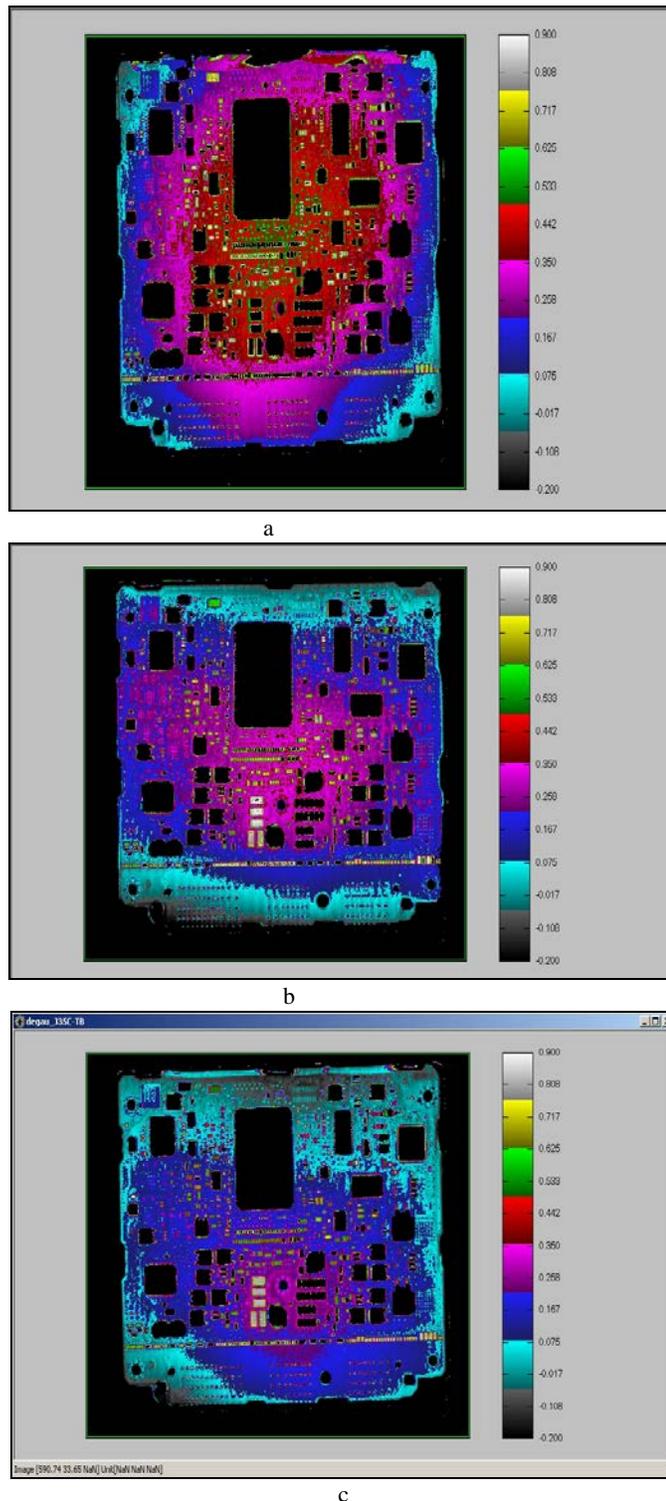


Fig. 8 Measures distorted on electronic boards

The above interferograms show clearly the changes on the object as fringes. The feeding of the board card under deformation changes with increasing time. A membrane under vibration lets appear different resonance frequencies. These

images are only examples of several results obtained with this technique, which means that the technique is useful for any application and any kind of object [13,14].

VI. CONCLUSION

The above interferograms show clearly the changes on the object as fringes. The feeding of the board card under deformation changes with increasing time. A membrane under vibration lets appear different resonance frequencies. These images are only examples of several results obtained with this technique, which means that the technique is useful for any application and any kind of object [13,14].

Electronic speckle interferometry (ESPI) is a fast and accurate non-destructive technique. It allows to make qualitative measurements of static or dynamic deformation at low frequencies, provided that the deformation does not exceed the diameter of the speckle grain. It is also shown that arithmetic operations such as subtraction to generate correlation fringes in real time, similar to those obtained by holography method.

REFERENCES

- [1]. Gary L. Cloud, 'Optical methode of engineering analysis,' Cambridge University Press.(1995).
- [2]. C. Wykes, 'Use of electronic speckle pattern interferometry (ESPI) in the measurement of static and dynamic surface displacements,' Opt. Eng., 21, 400-406, (1982).
- [3]. D. J. Løkbberg, 'Recent developments in video speckle interferometry,' In speckle Metrology. Edited by Rajpal. S. Sirohi. Marcel Dekker Inc., (1993).
- [4]. S. Nakadate, T. Yatagai, and H. Saito, ' Computer- aided speckle pattern interferometry,' Appl. Opt., 22, 237, (1983).
- [5]. Cloud, G., 2003. "Experimental Techniques" 27(4) pp27-30, also 27(5) pp15-17
- [6]. G. L. Cloud, G.L., 1995 "Optical Methods of Engineering Analysis" Cambridge University Press
- [7]. Hariharan, P., 2003 "Optical Interferometry" Second Edition , Elsevier Academic Press
- [8]. J. Tong, D. Zhang, H. Li and L. Li, 'Study on in- plane displacement measurement under impact loading using digital speckle pattern interferometry,' Opt. Eng., 35, 1080-1083, (1996).
- [9]. E.Olivier, Electronic speckle pattern interferometry, temporal vs. Spatial phase shifting, Poland state university
- [10]. K. Creath, 'Phase-shifting speckle interferometry,' Appl. Opt., 24, 3054-3058, (1985).
- [11]. S. Nakadate and H. Saito, 'Fringe scanning speckle pattern interferometry,' Appl. Opt., 24, 2172-2180, (1985).
- [12]. J. I. Kato, I. Yamaguchi, and Q. Ping, 'Automatic deformation analysis by a TV speckle interferometr using a laser diode,' Appl. Opt., 32, 77-83, (1993).
- [13]. R.Daira, B.boudjema, 'Study of shearing and deformation object by optical method', Digest Journal of Nanomaterials and Biostructures Vol. 7, No. 2, April - June 2012, p. 555 – 561.
- [14]. R. Daira, B. Boudjema 'Application of optical method in chemical processus as corrosion', phys.chem.news PCN, volume 68, pp 36-41.2013.

A Modified Adaptive Line Enhancer for Noisy Speech Signals

Maha Sharkas, M. Essam Khedr, Amr Nasser

Abstract— Adaptive filter is a method for enhancing noisy speech signals. In this work, a refinement for the adaptive line enhancer proposed by Widrow and Sambur is introduced. The new algorithm is based on both the exact value of gradient and variable step-size. These values are computed and used for coefficient updating. The algorithm is applied on different cases of noisy speech signals as well as GSM signals.

The new algorithm gave good results even when noise is greater than signal which is considered a novel achievement in the field. It avoids the drawbacks of LMS algorithm introduced by widrow and others. The adaptive line enhancer (ALE) design method updates the FIR filter coefficients in terms of the exact value of gradient and a variable step size.

Keywords— Least mean squares, adaptive filters, adaptive noise canceller.

I. INTRODUCTION

ADAPTIVE filter is a kind of technology which is widely used in the field of the modern signal processing. It can detect and extract useful signal from the strong interference of noise environment. The adaptive noise canceller (ANC) whose objective is noise interference is the typical application of adaptive filters. Also, it restrains and attenuates the noise to improve SNR quality of the transmitting and receiving signal. Adaptive filter includes a digital filter and an adaptive algorithm. It uses least-mean-squares algorithm as its standard, then adjusts the filter coefficients to achieve the best filter characteristics using an adaptive algorithm. The least-mean-squares (LMS) algorithm is the most commonly method as it has a simple structure and a basic architecture but with a disadvantage in stability and convergence speed. Therefore, the adaptive noise canceller is designed adopting an improved LMS algorithm [1].

The ANC Based on Least Mean Square Algorithm has been extensively used in many researches [1-6]. The purpose of ANC is to remove the noise from a signal adaptively to improve the SNR. Because of its simplicity and ease of implementation, the LMS algorithm is the most popular

Maha Sharkas Author is with Arab academy for Science & Technology, Alexandria, Egypt; (e-mail:msharkas@aast.edu).

M. Essam khedr Author, is with Arab academy for Science & Technology, Alexandria, Egypt; (e-mail: khedr@vt.edu).

Amr Nasser. Author is with Arab academy for Science & Technology, Alexandria, Egypt; (e-mail: amrnaser2@yahoo.com).

adaptive algorithm. However, the LMS algorithm suffers from slow and data dependent convergence behavior.

One of the main disadvantages of the LMS algorithm is that it cannot be used when noise is greater than signal [7]. The algorithm can be modified to overcome this problem by using (Overlap and add - Overlap and save) methods as suggested in [7].

The objective of this paper is to implement a novel algorithm to extract signal from noise especially when noise is greater than signal.

Fig.1 gives a simplified block diagram of the system.

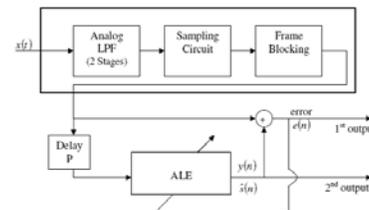


Fig.1 A simplified block diagram of the system

The paper is organized as follows: section 2 reviews the Overlap and add - Overlap and save methods. The algorithm design steps are introduced in section 3. Results are presented in section 4 and the paper is concluded in section 5.

II. THE OVERLAP AND ADD - OVERLAP AND SAVE METHODS

In practice, the length of an input signal may be very large. To process such long sequences, the input signal can be sectioned into blocks of a length small enough to be processed by a given computer, and then add the output resulting from all the input blocks. This procedure can be used because the operation is linear. Sectioning the input allows the output to have a smaller processing delay.

There are two well-known techniques for performing a linear convolution using the FFT algorithm and these are referred to as the overlap-save and overlap-add sectioning methods. By overlapping elements of the data sequences and retaining only a subset of the final DFT product, a linear convolution between a finite-length sequence and an infinite-length sequence is readily obtained [8 and 9].

III. ALGORITHM DESIGN STEPS

The design steps include the following computations [9]:

- (1) The delay time p .
- (2) The maximum number of filter coefficients.
- (3) The exact value of gradient.
- (4) The variable step-size.

A. Computing the delay time (P):

The adaptive line canceller ALE is a special case of the adaptive noise canceller ANC, where there is only one signal $x(n)$ available, the primary input, which is contaminated by noise. In such a case the reference signal $u(n)$ is the delayed replica of $x(n)$ such that $u(n) = x(n - p)$. Suppose the signal $x(n)$ consists of two components: a narrow band component $s(n)$, which has long range correlation, and a broadband component $v(n)$, which will tend to have short range correlation. Also, L_{NB} and L_{BB} are effectively the self-correlation lengths of the Narrowband (NB) and broadband (BB) respectively, as shown in figure 3. Beyond these lags, the respective correlations die out quickly. The delay P is selected so that:

$$L_{BB} \leq P \leq L_{NB}$$

If P is longer than the effective correlation length of BB, then the delayed replica $v(n - p)$, will entirely be uncorrelated with $v(n)$ which is a part of the primary signal. The adaptive filter will not be able to respond to this component. On the other hand, when P is shorter than the correlation length of the NB component, the delayed replica of $s(n - p)$ that appears in the reference input will still correlate with $s(n)$ which is a part of the primary signal. In this case the filter will respond to cancel $v(n)$.

If p is selected to be longer than both correlation lengths, then the reference input will become uncorrelated with the primary input and the adaptive filter will be turned off. In the opposite case, when the delay is selected to be less than both the correlation lengths, then both components of the reference signal will be correlated with the primary signal, and therefore the adaptive filter will respond to cancel the primary signal $x(n)$ completely.

The structure of the proposed ALE is illustrated in figure (3).

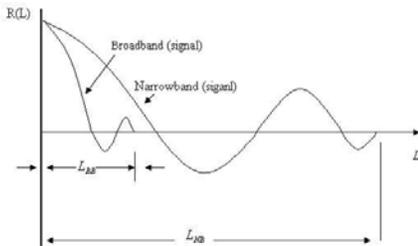


Fig.2 The correlation lengths for the narrow band and the broadband signals.

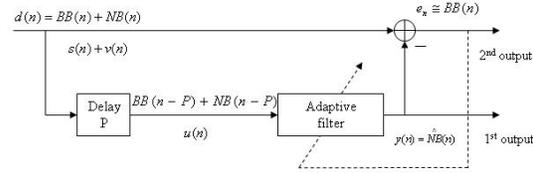


Fig.3 The structure of ALE

B. Computation of Maximum of filter (taps) coefficients

In the following, the maximum number of filter taps, M_{max} is estimated, based on the cross correlation of the primary and the reference input samples. During M intervals of length T , M is the number of samples per observation window. The cross correlation for an assumed delay τ samples will

$$R(\tau) = \sum_{K=1}^M x(k)u(k - \tau)$$

$$= \sum_{K=1}^M x(k)x(k - p - \tau)$$

The estimated NB component, or filter output $\hat{s}(n)$ is obtained from the M forgoing samples $u(n), u(n-1), \dots, u(n-M+1)$. The maximum correlation length τ_{max} is computed as the value of the variable τ at which the correlation between $x(n)$ and $u(n - \tau)$ is negligible with respect to the correlation between $x(n)$ and $u(n)$. Mathematically, $x(n)$ and $u(n - \tau)$ are completely uncorrelated if:

$$\frac{R(\tau)}{R(0)} \leq \varepsilon$$

Where ε a small positive values, typical value for is ε is 0.01. The maximum number of filter taps is given by:

$$L = M_{max} = \tau_{max}$$

C. Computation of the exact value of the gradient

In this subsection, an exact value for the gradient is obtained. Under stationary conditions for the signal $u(n)$, the best estimate of $v(n)$ would be:

$$\frac{\nabla(n)}{2} = -\left[\frac{1}{n} \sum_{i=1}^n e(i)u(i) \right]$$

$$= \frac{-1}{n} [(n-1)\nabla(n-1) + e(n)u(n)]$$

But, in the adaptive situation $u(n)$ is non-stationary, therefore it can be seen that the previous equation would not be a good estimate of $\nabla(n)$ because of its infinite memory. This estimate would become insensitive to changes for large values of n . To provide this effect, the previous equation becomes:

$$\begin{aligned}\frac{\nabla(n)}{2} &= \frac{-1}{n} \left[\sum_{i=1}^n \lambda^{n-i} e(i) U(i) \right] \\ &= \frac{-1}{n} \left[\sum_{i=1}^{n-1} \lambda^{n-i} e(i) u(i) + e(n) u(n) \right] \\ &= \frac{-1}{n} [(n-1)\lambda(n)\nabla(n-1) + e(n)u(n)]\end{aligned}$$

The choice of the forgetting factor λ is often very important:

- Theoretically one must have $\lambda = 1$ (for stationary input) to get convergence.
- For the non-stationary case, the algorithm becomes more sensitive and the parameter estimates change very quickly. For that reason, it is often an advantage to allow the forgetting factor to vary with time. Therefore, substitute λ by $\lambda(n)$. A typical choice is to let λ tend exponentially to 1. This can be written as

$$\lambda(n) = \lambda_0 \lambda(n-1) + (1 - \lambda_0)$$

Typical value for λ_0 and $\lambda(0)$ are

$$\lambda_0 = 0.99 \text{ and } \lambda(0) = 0.95$$

D. Computation of a variable step size

An adaptive step-size must satisfy: 1. The speed of convergence should be fast. 2. When operating in stationary environments [10], the steady state miss-adjustment values should be very small. 3. When operating in the non-stationary environments, the algorithm should be able to sense the rate at which the optimal coefficients are changing. The above goals are achieved when the step-size $\mu(n)$ is adjusted as:

$$\mu'(n+1) = \alpha\mu(n) + \gamma e^2(n)$$

With $0 < \alpha < 1$, $\gamma > 0$

$$\left\{ \begin{array}{l} \text{And} \\ \mu_{\max} > \mu_{\max} \\ \mu(n+1) = \mu_{\min} < \mu_{\min} \\ \mu'(n+1) \text{ Otherwise} \end{array} \right.$$

The initial step-size $\mu(0)$ is usually taken to be μ_{\max} . The

choice of μ_{\max} is based on stability considerations and it must be chosen in the range $\frac{1}{tr[R]} > \mu_{\max} > 0$, where

$$tr[R] = LP_x$$

P_x is the power of the input signal, which is approximated by

$$P_x \cong \frac{1}{L} R(0)$$

On the other hand, μ_{\min} is chosen to provide the desired misadjustment, after the filter has converged. The misadjustment, ratio of the average excess mean-square error

to the minimum mean-square error, is taken to be 5 percent, and μ_{\min} is

$$\mu_{\min} = 0.1\mu_{\max}$$

The value of α appears to be a good choice for all experiments ($\alpha = 0.97$), while the value of γ is chosen arbitrary. Typical values of γ in stationary and non-stationary are $4.8 \cdot 10^{-4}$ and $7.65 \cdot 10^{-14}$ respectively.

IV. EXPERIMENTAL RESULTS

The speech is normally low-pass filtered at frequency of about 1 KHz, which is well above the maximum anticipated frequency range for pitch (500 Hz for female speech). Filtering helps to reduce the effect of higher formants and any high-frequency noise [11].

Two low pass filters are used:

The first one keeps all frequencies below 1000 Hz and eliminates all frequencies higher than 2000 Hz.

The second filter preserves the frequencies less than 500 Hz, and attenuates the frequencies higher than 1000 Hz. The corresponding cut-off digital frequencies in radians for the first filter are as follows.

Pass band frequency:

$$\omega_p = \frac{f_p}{F_s/2} = \frac{1000}{8000/2} = 0.25$$

Stop band frequency:

$$\omega_s = \frac{f_s}{F_s/2} = \frac{2000}{8000/2} = 0.5$$

Similarly, for the second filter we have:

$$\omega_p = 0.125, \omega_s = 0.25$$

The sampling frequency $\frac{f_s}{2}$ is chosen to be fairly close to the upper limit of the ear's sensitivity

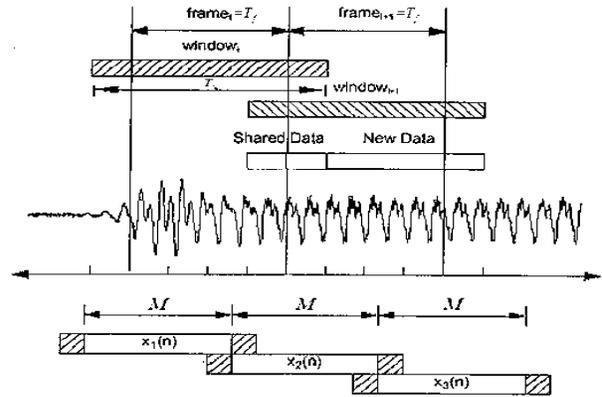


Fig.4 Frame Blocking

Frame blocking is applied as follows:

- If we assume the signal is piecewise stationary, we can analyze the signal using a sliding window. The two used key parameters are:

- Frame duration T_f .

(b) Window Duration T_w .

Typical values are: $T_f = 5 \text{ ms}$, $T_w = 10 \text{ ms}$

(ii) The window duration controls the amount of averaging (or smoothing) used in power calculation.

- Each frame will contain 40 samples.
- Each window will contain 80 samples.

(iii) The amount of the overlap value is given by:

$$(\text{overlap}) 100\% = \left(\frac{T_w - T_f}{T_w}\right) 100 = 50$$

Frame blocking is depicted in figure 4.

The proposed system referred to in figure 1 is simulated, using Matlab. The experiment is carried out in three main steps [10]:

Step 1: Preliminary adjustment.

Step 2: Determination of the optimum filter length.

Step 3: Measurement of ALE performance.

- a) It was necessary to record speech alone and noise alone. The noise speech is obtained by adding the noise to the clean signal according to the required SNR[12]

$$SNR_{in}[dB] = 10 \log_{10} \left[\frac{\text{signalpower}}{\text{noisepower}} \right]$$

- b) The speech signals are uttered by male speaker for the following three cases:

- (i) Case (1): Voiced vowel [a] “Aah”
- (ii) Case (2): For word [Ahmed]
- (iii) Case (3): For complete sentence “How old are you”
- (iv) A GSM signal

- c) The noise signal is Gaussian type.
- d) The sampling frequency used is 10024 Hz.
- e) The duration of vowel to be synthesized is 105 msec sectioned into 7msec duration at which $T_f = 5\text{msec}$ and $T_w = 7\text{msec}$.

The results presented shows for all the above cases a clean signal, signal with added noise then the extracted signal using the modified technique. The SNR between the obtained output and the input is then given. The modified technique succeeded in retaining the original signal with good SNR as depicted in figures 5, 6 and 7.

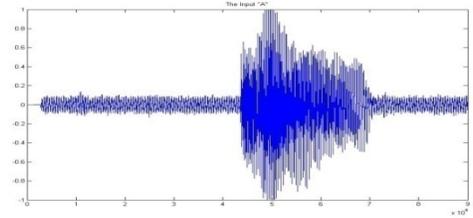


Fig. 5 (a)

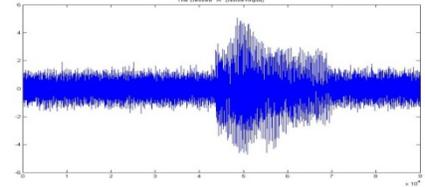


Fig. 5 (b)

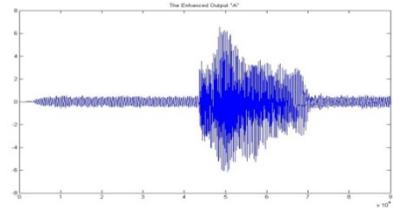


Fig. 5 (c)

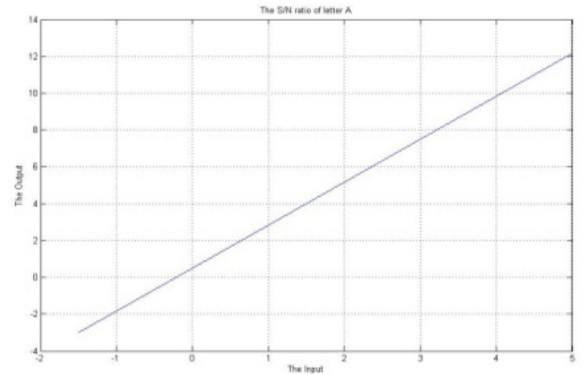


Fig. 5 (d)

Fig. 5: a: the clean signal for the vowel "a", b: the signal with added noise, c: the extracted signal and d: the SNR between the output and the input.

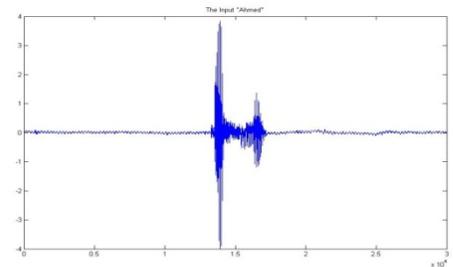


Fig. 6 (a)

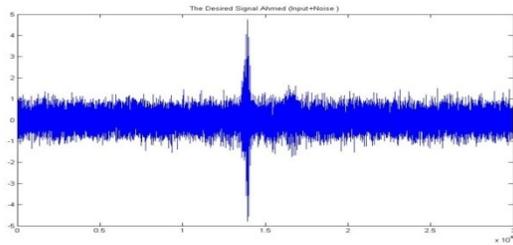


Fig.6 (b)

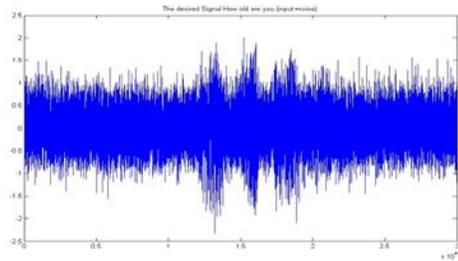


Fig. 7 (b)

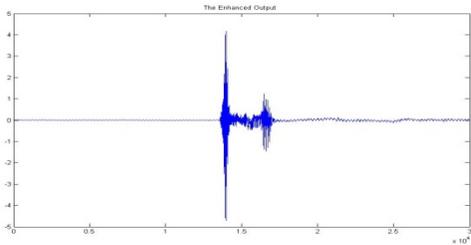


Fig. 6 (c)

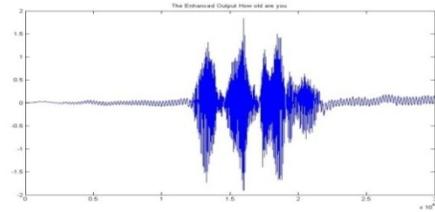


Fig. 7 (c)

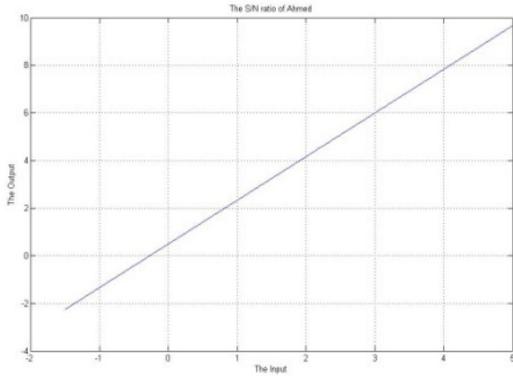


Fig. 6 (d)

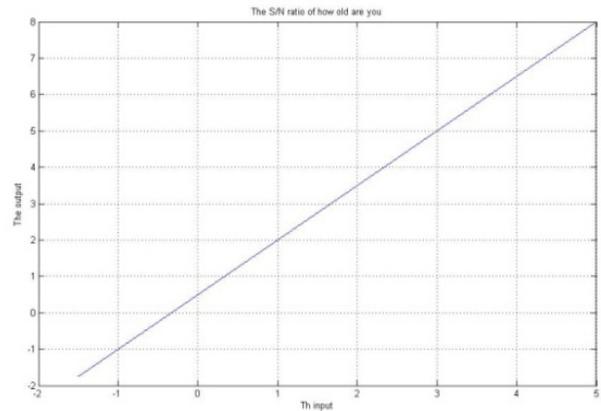


Fig.7 (d)

Fig. 6: a: the clean signal for the word "Ahmed", b: the signal with added noise, c: the extracted signal and d: the SNR between the output and the input.

Fig. 7: a: the clean signal for the sentence "how old are you", b: the signal with added noise, c: the extracted signal and d: the SNR between the output and the input.

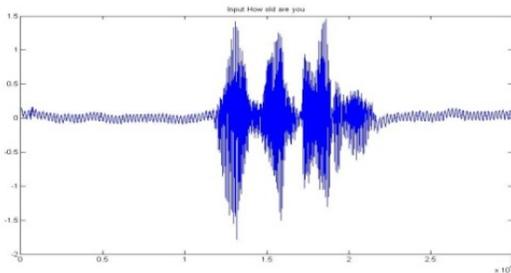


Fig.7 (a)

The technique is also applied on GSM signals namely the GMSK signal as shown in figure 8. The extracted output signal is shown in figures 8 c & d with and without using the modified technique. The SNR without and with using the modified algorithm are compared in figure 8.e indicating the modified performance of the suggested technique.

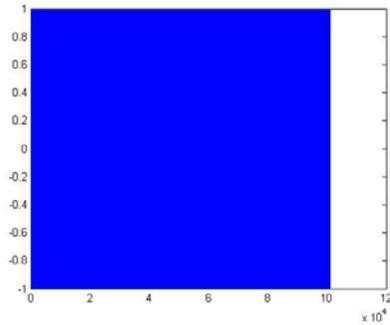


Fig. 8.a The input GMSK signal

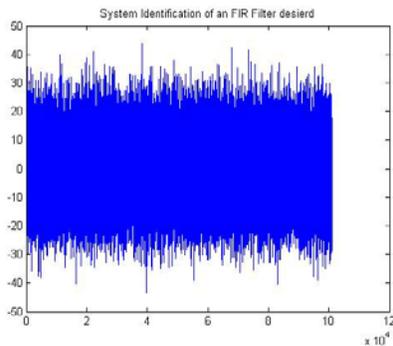


Fig. 8.b The input GMSK signal + noise

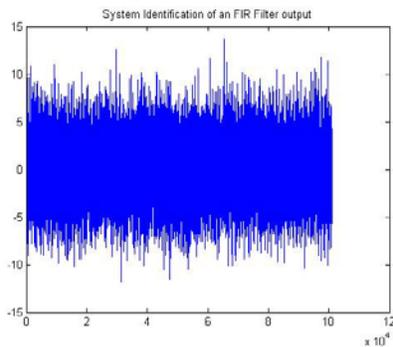


Fig. 8.c The extracted output GMSK signal without the modified algorithm

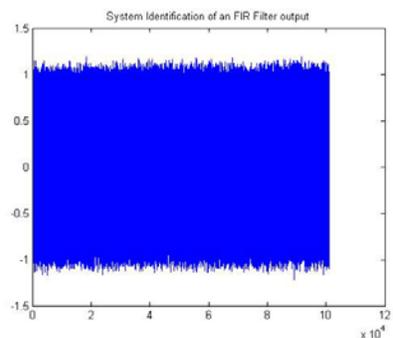


Fig. 8.d The extracted output GMSK signal with the modified Algorithm

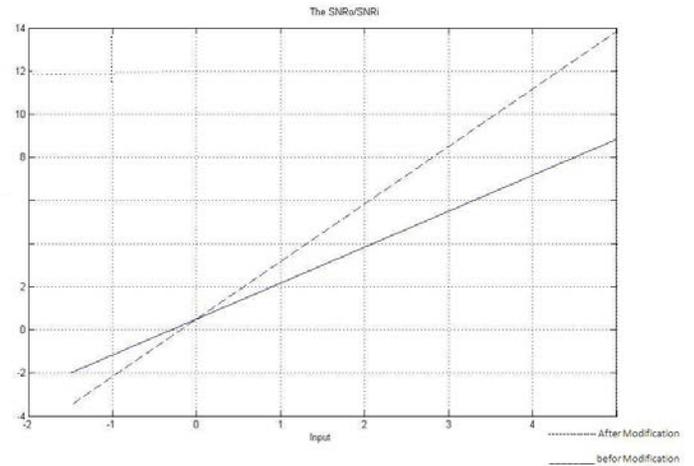


Fig.8.e. SNR in db between input and output signals before and after modification

Fig,8 Results for GSM signal

V.CONCLUSION

An adaptive noise enhancement based on a modified LMS algorithm is implemented and tested. This modified adaptive line enhancer (ALE) can be used to detect and estimate weak signals in noise. Two modifications are introduced to the ALE using the FIR structure which are;

- a) The exact value of the gradient is computed. The exact value emphasized the most recent values of the gradients and ignore the older one according to the value of the forgetting factor λ [λ lies between 0 and 1]
- b) The variable step-size that depends on the square of predication error allows the adaptive filter to track changes in the system as well as produce small steady-state error.

This work demonstrates the ability of ALE to reduce additive periodic or stationary random noise in both periodic and random signals. The suggested modification allowed enhancing weak signals and provided good SNR. It is also applied on GSM signals and succeeded in providing improved results.

REFERENCES

- [1] Xu Yanhong and Zhang Ze, " Design the adaptive noise canceller based on an improved LMS algorithm and realize it by DSP", fifth international conference on intelligent computation technology and automation, 2012.
- [2] Gaurav Saxena, Subramaniam Ganesan, and Manohar Das," Real time implementation of adaptive noise cancellation", IEEE international conference on electro/information technology, 2008.
- [3] Hongbing Li and Hailin Tian, "A New VSS-LMS Adaptive Filtering Algorithm and Its Application in Adaptive Noise Jamming Cancellation System", IEEE circuits and systems international conference on testing and diagnosis, 2009.
- [4] Mohammed Hussein Miry, Ali Hussein Miry and Hussain Kareem Khleaf," Adaptive Noise Cancellation for speech Employing Fuzzy and Neural Network", Ist international conference on energy, power and control, December, 2010.
- [5] Vakulabharanam Ramakrishna and Tipparti Anil Kumar, *SMIEEE*, " Low Power VLSI Implementation of Adaptive Noise Canceller Based

- on Least Mean Square Algorithm", 4th International Conference on Intelligent Systems, Modelling and Simulation, 2013.
- [6] Yuzhong Jiao, Rex Y. P. Cheung, Winnie W. Y. Chow and Mark P. C. Mok, " A Novel Gradient Adaptive Step Size LMS Algorithm with Dual Adaptive Filters", 35th Annual International Conference of the IEEE EMBS Osaka, Japan, 3 - 7 July, 2013.
 - [7] Orfanidis .J."Optimum Signal Processing: An Introduction ". McGraw Hill, 1990.
 - [8] Shynk,J.J., "Frequency-domain and multirate adaptive filtering", IEEE signal processing magazine, vol.9, jan, 1992.
 - [9] Simon Haykin, "Signals and Systems", John Wiley and Sons, Inc, 2003.
 - [10] R.W.Kwong, "A variable step-size LMS Algorithm". IEEE Trans., Signal Proc., Vol.40, No.70, PP. 1663_1642, July 1992.
 - [11] Saeed Vaseghi."Advanced digital signal processing and noise reduction". 4thed, John Wiley, 2008.
 - [12] Yuzhong Jiao, Rex Y. P. Cheung, Winnie W. Y. Chow and Mark P. C. Mok, " A Novel Gradient Adaptive Step Size LMS Algorithm with Dual Adaptive Filters", 35th Annual International Conference of the IEEE EMBS Osaka, Japan, 3 - 7 July, 2013.

A Comprehensive analysis of XML and JSON web technologies

Zia Ul Haq¹, Gul Faraz Khan², Tazar Hussain³

¹Management Information System Department, College of Business Administration,
King Saud University, P.O. Box 71115, Riyadh 11587, Saudi Arabia

^{2,3}Department of Software Engineering College of Computer and Information Sciences,
King Saud University, P.O. Box 51178, Riyadh 11543, Saudi Arabia

¹zparacha@ksu.edu.sa

²gfaraz@ksu.edu.sa

³tazhussain@ksu.edu.sa

Abstract—In this global era, internet plays a vital role to share information all over the world. There are some standard protocols and web technologies to represent the information on internet like HTML, JavaScript, ASP, JSP, XML, JSON etc. All these standards have pros and cons, and also depend on the requirement that which exchange of format more reliable. This research work shows the comparative analysis of XML and JSON using Multi-criteria Decision Support System MCDS, and analyze use of web technology according to the requirement need.

Keywords: Web technology, JSON, XML, MCDS, AJAX

I. INTRODUCTION

Today web technology is on the level that we are talking about synchronous and asynchronous transformation technologies like XML, AJAX, .NET, JSON and so on are in practice. But 15-16 years ago the people was only thinking how to display documents on the internet and for that purpose HTML (Hyper Text Markup Language) had been developed. And then with the time as people get use to with the internet, then the user start demands for new things like to have interactive websites that they can display information as well as can take input from user for that purpose different technologies had been developed that we are using in dynamic websites such that ASP, JSP, PHP, ASP.NET etc. then there was a problem how to transfer data between different platforms so XML has been developed which store plain text base data independent of platform and can easily transfer data between cross platforms.

In the last couple of years AJAX came in to existence which is totally new technology, in the same time JSON has also been introduced, and it is nowadays one of the most debatable topic for developers that “JSON is a fate free alternative to XML?” [10].

In this paper have discussed a little background of the web development starting from HTML, SGML, AJAX, XML and JSON and then XML and JSON is compared in different

aspects and then on the basis of these comparisons conclusion has been derived that which technology is required for whom and for what purpose.

2. Background of XML and JSON

HTML is used and still the people are using it how to display some data or documents, how it will be look like and what will be its layout for that purpose, it uses some different tags like <HTML>, <HEAD>, <TITLE>, and <BODY> etc. to display documents. It was difficult to remember all the tags of HTML but as rapid developments are taking place in this field and now we can have different editors that we can use them how to design a HTML pages but still HTML has some limitations.

When the technology been more developed and the developers started thinking about to store and retrieve data to and from the database so HTML is ok with how to display data and documents but it is not suitable to deal with data structure like some stuff like database and objects hierarchies etc. [11].

2.1. SGML

As we know that HTML as good for displaying documents and their layout but HTML is not suitable for explicit queries. So if we need to go for a bit different things like how to access and manipulate the data objects instead of documents then HTML is not capable of that, So SGML (Standard Generalized Markup Language). As we know that HTML is one of the applications of SGML. SGML is known as the mother tongue of HTML. So to overcome the limitation of HTML, SGML has the capability to overcome the problems of HTML many organizations had defined their own standards using SGML but SGML is very complex as its specifications are about 500 pages. Due to the very complex structure of SGML it turn the researcher to focus on such a tool which should overcome the limitations of HTML and should be

simple where the user can use in their application on the basis of simple standard and XML came into being [11].

2.2. AJAX (Asynchronous JavaScript and XML)

AJAX is tremendous development on the web which enables web application to look like desktop application means like in the normal WebPages when we need a bit of information we need to reload the whole page, but with AJAX it is now possible to transfer information with out reloading or refreshing the page. For example Google Suggest Lab or yahoo search when we are typing the word to search it give us suggestion with out reloading the page. It means that the AJAX technique enable pages to have the capability to transfer a small amount of data and can retrieve the data from the server in the same page, with out reloading the whole page. Now when we transferring this data we need data transfer format either XML or JSON can be used is data transfer formats. In the next section I am giving an overview to XML and then JSON and then I will come to give detail about which one XML or JSON to use with AJAX [7].

2.3. XML (Extensible Markup Language)

XML is a simple standard that can be used to transform and encode both text and data and it can be processed and transformed across different platforms. XML is more simple and manageable as compared to SGML, it is not a full independent markup language but it is a tool that enables users to define their own markup language, user can define their own tags which are more easy and readable by both machine, as well as user [11].

To take care of compatibility the XML is built to in the way to fit the HTML in the new framework and after some time HTML 4.0 which is a very suitable version that is compatible with XML and there it is named as XHTML, and simply we can say that it is nothing more but a special application of XML.

XML is a Meta data and the basic difference between HTML and XML is that, XML allow the user to define their own tags as HTML have a pre defined set of tags but in XML the user can define the tags according to their own use. And the other thing is that, that there are no formal criteria like the HTML have (head, title, body etc), it also take care of contents as well as XML do not defined any contents criteria but as XML is meta language so it allows contents base structuring as well as XLM is functioning different then the level HTML do [11].

So XML will need to be fid with XML based data by the user to get succeed. User can define their own tags according to their own need in their native language. Such that “StartDate”, “TrouserSize”, “EndTime” so XML do not define these tags but actually XML give the way how to define these tags, it means that XML have tags that apply to all documents (XML processing instruction) and also user can define their own criteria that can be stored as a separate documents (DTD) document type definition (Jung 2000). XML provide some grammar rules that apply to all documents, such that

“<identifier>contents</identifier>” then the user can put the contents according to their own use for example <name>khan</khan>. Then everyone can use it according to their own use for instance automobile company can use it <EngineSize>, <ModelOfCar>, <Color> and then can save them in a special schema. The XML capable application can use these definitions directly.

Beside these as XML don’t understand these tags so it means that instead of <StartDate> if say you put <ABC123> XML parser will accept it without any problems it means that it could be anything that is understandable so for this purpose a group of user can be agree on some definitions of tags Document Type Definition (DTDs) can be used then [11].

XML architecture is very flexible documents are written in plain text not in cryptic form so it will be human readable and XML document can be stored, processed and distributed except any special DTD or XML code of the particular author [11].

As XML is flexible so any kind of tags can be defined. But DTD should be use for particular type of data such that real store, Publishing House etc. some of the organization and association already reached on agreement of such type of XML Documents definitions [11].

A simple example of XML data as given below,

```
<teaminfo>
  <players>
    <player>
      <name>Ajmal Khan</name>
      <height>5.8</height>
      <age>24</age>
      <postgrad>true</postgrad>
    </player>
    <player>
      <name>Nadeem shehzad</name>
      <height>5.9</height>
      <age>26</age>
      <postgrad>true</postgrad>
    </player>
    <player>
      <name>Aditya</name>
      <height>6.0</height>
      <age>24</age>
      <postgrad>true</postgrad>
    </player>
  </players>
</teaminfo>
```

The above XML example store information about three players in a team. We can see there are some elements that they are repeating for each piece of information however the actual data is repeated once. And we are interested in the players and their information. The elements <players> and <teaminfo> are not needed but they are just used to define the structure and meaning of the information. Table 1 shows advantages of XML technology over JSON [3].

Table 1. XML strengths over JSON

JSON	XML
There is no grammar support and that's why it is difficult to communicate and enforce interface contracts	While XML have XML schema and Document Type Definition which can be used to define grammar rules
Extensibility is not good as namespaces are not supported	Very strong support for namespaces, schema have more extensibility options
Development tools support is very limited as it is newly introduced	As XML is in the market since long time there for is supported by most of the development tools
JSON is very narrow focused as it is used only for Remote Process Call(RPC), mainly with JavaScript Client	XML is very broad focused, it can be used for Remote Process Call (RPC), Electronic Data Interchange (EDI), Metadata etc.
Very limited support for web services associated stuff (products).	huge hold of web services related products

2.4. JavaScript Object Notation (JSON)

As web services are gaining attractiveness day by day, XML has almost turned into the actual paradigm for data transmission. But still there are different things that people considering that XML is heavy some time it send more bytes through the internet to get the things done which can be done with a much smaller data. To overcome this problem new formats of XML been introduced like binary XML so it means all these solutions are extending XML but still the problem exist when it come to backwards compatibility. Douglas Crockford is software engineer he introduced a new data format which is based on JavaScript called JavaScript Object Notation (JSON) [6].

JSON is simple very lightweight object serialization technique or data format which is based on JavaScript Object initialization syntax, specifically array and object literals. JSON definitions can be incorporated inside JavaScript files and accessed with no further parsing that comes alongside with XML-based languages, because it uses JavaScript syntax. But prior to use JSON, it is essential to know the array and object literals particular JavaScript syntax. JSON the

initialization code is assigned to a string and then is dealt with JavaScript eval() function or JSON parser [6].

The JSON parser is very light weight. JSON is mainly used with different AJAX tools kits and frameworks and provide easy serialization for remote calls. JSON is supported by GWT and DOJO. JSON and web 2.0 technologies must be considered very seriously by Service Oriented Architecture (SOA) (Alexander 2007).

Most of the server languages do not contain JavaScript interpreters there for they wouldn't able to process the JavaScript code and lets suppose that they can evaluate it, but still the developer wouldn't allow the arbitrary code to b run on the server because it can generate a serious security problem. But both of these problems been solved by JSON which is the literal Syntax for JavaScript, as JSON can be run using eval() function of the JavaScript side. JSON allows any JavaScript data types to be transferred and would be faster than the XML-based solutions because the compact encoding allows for much smaller data to be transferred. On server side small parser to be built to serialize the native data types into JSON and also to create native data types from JSON [7]. The following sample describes how to represent JSON data.

```
{ "teaminfo" :
  {
    "players" : [
      {
        "name" : "Ajmal Khan",
        "height" : 5.7,
        "age" : 24,
        "postgrad" : true
      },
      {
        "name" : "nadeem shehzad",
        "height" : 5.8,
        "age" : 26,
        "postgrad" : true
      },
      {
        "name" : "Aditya",
        "height" : 6.0,
        "age" : 23,
        "postgrad" : false
      }
    ]
  }
}
```

We can observe that a lot of redundant information is not present as compare to the XML one; no closing tags are required to match opening tags this will reduce the number of bytes to be send out for same information of to great extent. In the above example excluding spaces the JSON data is 249 bytes while the XML data is 378 bytes so it save more than 120 bytes in this much data. That is the reason why Crockford, JSON inventor said that “JSON is a fat free alternative to XML”.

The draw back of the JSON format is that it is a bit hard to read as compare to XML as XML one is easier to read by a layperson as it is clearer and meaning full but JSON format is reduced by its shorthanded notations and due to that it will be difficult to read it with nicked eye. But no it can be disagreement that why we need to view the data exchange format with naked eye if we can use tools for parsing the data passing back and forth but still the question arising that are such tools available so there are some tools available but still this can be a limitation some where (EICHORN 2006). JSON have some pros over XML as shown in table 2 [3].

Table 2. JSON strengths over XML

JSON	XML
Completely programmed technique for de-serializing and serializing JavaScript objects, with very little coding.	JavaScript code will be written by developer to serialize and de-serialize to and form XML
Most of the browsers have enough support of JSON.	All new browsers have built-in XML parser but it could be a bit tricky when it come to cross-browser XML parsing.
The format is very concise due to having name/value pair-based approach.	Because of tags and namespaces the format is very lengthy.
de-serialization is very speedy in JavaScript	De-serialization is slower in JavaScript
Most of JavaScript libraries and AJAX toolkits have good support of JSON	AJAX toolkits don't have strong support for it.
Having simple API for JS and more other languages	The APIs are very complicated

3. Related work

Data interchange format have significant consequences on data transferring rates and performance, data interchange format generate from mark-up to further support for structural attribute of information using encoding of meta-data. XML and JSON are data interchange format that can be use in different aspects with unique purpose. XML primary uses are

object serialization for transfer of data between application and Remote Procedure Calls RPC [13]. To analyze the impact of management of energy and cost of processing for transferring data in mobile devices of proposed formats of XML, JSON and Protocol Buffer [5]. These two form of data interchangeable using data serialization approach which allows for better communication between applications. Data transmission of web application more secure, powerful in the XML serialization approach and with JSON serialization approach fast and convenient [17].

XML is more complex than JSON in web services of web programming, some time programmer doesn't need to use namespace and mixed content documents. The developers focus on to use simple data structure, compact and exchange format. XML is great in problem domain, namespace, well-formed and mixed content document [15].

The Analytic Hierarchy Process is a method of measurement for formulating and analyzing decisions. AHP is a decision support tool which can be used to solve complex decision problems considering tangible and intangible aspects. Therefore, it supports decision makers to make decisions involving their experience, knowledge and intuition [2],[9].

4. Comparative analysis

In this research work comparison between JSON and XML carried out by using the Make it Rational MCDM tool. MCDM is decision making tool based on comparative analysis and defined steps. For decision making system these are the main steps to be evaluate goal, alternative, criteria, preferences, sub criteria and final result [1].

4.1. Goal

Main goal for this work to be analyzed are JSON and XML technologies based on criteria and sub criteria. This work will help developer to choose proper web technology in certain condition [12].

4.2. Alternatives

There are many web development technologies, this research work focus on alternatives of JSON and XML.

4.3. Main criteria

Hierarchy view show top down approach for alternatives, criteria's and sub criteria's, as shown in Figure 1, in our case main criteria's are Format of Exchange, Validity, Readability, Efficiency, Debugging and Troubleshooting, Ease of data creation, and sub criteria's are Machine and Human Readability.

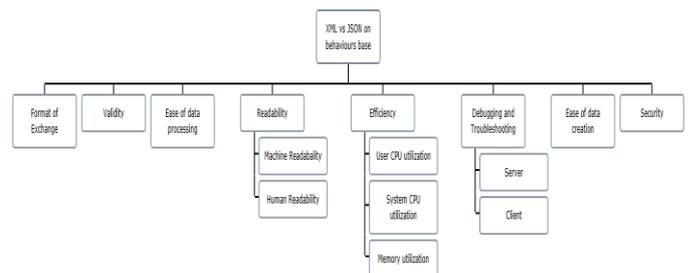


Figure 1: Hierarchy view of XML and JSON

4.3.1. Format of Exchange

JSON format is always smaller than XML, in fact the more tag involves in XML format increase size of XML exchange format. JSON specification excluding unnecessary tag and syntax produce small size of exchange format comparatively to XML. AJAX application can use XML or JSON as transformation format. Now it is an important issue how to select transformation format. XML is more composite structure and could be able to transfer any type of data however JSON is simple data structure that will be all we need to transfer AJAX data, and it is useful enough to use JSON with AJAX such that, JSON is subset of JavaScript, so using eval() method of JavaScript JSON text can be simply changed into JavaScript Object and then we can extract data using JavaScript. So if someone know JavaScript will be easy for them to use.

JSON is subset of JavaScript so it have the JavaScript Data types but it is not same for XML we can define it for XML using XML schema or DTD to define XML structure.

JSON can be parsed as JavaScript, for more security we can use the JSON parser to convert. But if we use to get the as XML it means that we will need to parse it. For that Dom method will be needed which is very comple [16].

4.3.2. Validity

XML content rules conforms the data to be valid documents. These rules describe document organizational structure and accurate data values. Valid XML documents match defined schema, constrains on the structure and content of the document articulated in schema. JSON validator is a program that verifies JSON data with provided schema which contains define validation, documentation, interaction control of JSON data and hyperlink navigation. Research analysis shows that XML validation bit stronger than JSON.

4.3.3. Ease of Data processing

The simple data structure and data standard of XML provide very easy process. But it is same in JSON as it having very simpler structure [10].

4.3.4. Efficiency

XML document includes statements, handling instruction, elements and rich tags. Also composed of root node, the root node contains a root element, nested child element also include chilled element and properties which increase size of the XML documents. Data format of JSON is very simple that can be transmitted with a single array, variable of number or Boolean type or also string type. JSON exchanging data by object while object are tagged by unordered which contains a series of key values and pair key value [14]. Analysis and t-

test observation indicated that XML appears to use less user CPU utilization than JSON, XML use more system CPU utilization than JSON, and memory utilization of the JSON and XML encoded transmissions nearly the same on server [13].

4.3.5. Debugging and Troubleshooting

XML server checks data being sent to client well formed and valid. XML document verify with XML schema, as an alternative to XML, JSON manually involves verification that the response object has the right attributes. On client side it is difficult to spot error in either format; browser would fail to parse XML into the response XML. For small data relatively easy to detect error in JSON but with large data it is difficult to relate the error message to the data.

4.3.6. Ease of Data creation

XML data-binding APIs to create XML in more than a few programming languages, XML APIs have been around for years and may be deal with complex application. JSON APIs are new to create JSON responses but not so far behind than XML. There are so many ways to create XML alternative to JSON [4]. Security

XML language specifications ensure that the form of serialization data of XML has strong security; labels stored all data in the tag closed strictly with the data index. AJAX lightweight application demanding for low security while have high demanding for efficiency, so JSON have good support as an alternative to XML [17].

4.3.7. Extensible

Extensibility reduces the coupling between the producer and the consumer of the data. XML is extensible while JSON is not but there is no need for that because it is not document mark-up language so there is no need for new tags as already store data so there is no need to have tags to store data about data [10].

4.3.8. Reusability of software

As XML claim that the there are plenty of software code is available that developer can use that code and there is no need for recoding but JSON is simpler and there is no need for more programming/ additional software only JSON simple code is enough [10].

4.3.9. Adoptability by the industry

XML is adopted by the wide range of computer industry. Hover JSON is just newly known to the industry and because of it simplicity and it is easy to convert JSON from XML, due to these properties JSON becoming more adoptable [10].

4.4. Preferences

Preference is concerned with the priorities based on importance of criteria/sub-criteria, as assigned priorities are shown in the table 3.

Table 3. Preference based on priorities

Intensity	Importance	Intensity	Importance
1	Equal Importance	6	Strong Importance plus
2	Weak Importance	7	Very Strong Importance
3	Moderate Importance	8	Very Strong Importance plus
4	Moderate Importance plus	9	Extreme Importance
5	Strong Importance		

Each criteria assigned priority value shown in the given table 4.

Table 4. Criteria priority scale based on importance

Criteria	Ratio	Criteria	Ratio
Ease of data processing vs. Format of Exchange	1:1	Efficiency vs. Ease of data creation	3:1
Validity vs. Readability	3:1	Format of Exchange vs. Ease of data creation	2:1
Format of Exchange vs. Security	1:1	Security vs. Efficiency	1:1
Efficiency vs. Readability	2:1	Ease of data processing vs. Ease of data creation	2:1
Ease of data processing vs. Readability	2:1	Ease of data processing vs. Efficiency	1:1
Validity vs. Ease of data creation	3:1	Debugging and Troubleshooting vs. Efficiency	1:1
Security vs. Ease of data creation	2:1	Ease of data creation vs. Readability	1:1
Format of Exchange vs. Debugging and Troubleshooting	1:1	Security vs. Debugging and Troubleshooting	1:1
Debugging and Troubleshooting vs. Ease of data creation	2:1	Validity vs. Security	2:1
Validity vs. Efficiency	2:1	Format of Exchange vs. Efficiency	1:1

Debugging and Troubleshooting vs. Readability	3:1	Ease of data processing vs. Security	1:1
Format of Exchange vs. Readability	2:1	Validity vs. Ease of data processing	2:1
Ease of data processing vs. Debugging and Troubleshooting	1:1	Security vs. Readability	2:1
Validity vs. Format of Exchange	2:1	Validity vs. Debugging and Troubleshooting	2:1

4.5. Result and analysis

4.5.1. Ranking graph

Figure 2 shows rank graph for each criteria of XML and JSON, i.e. format of the exchange rank graph for JSON contribute more than XML, while the validity of an XML value higher than JSON.

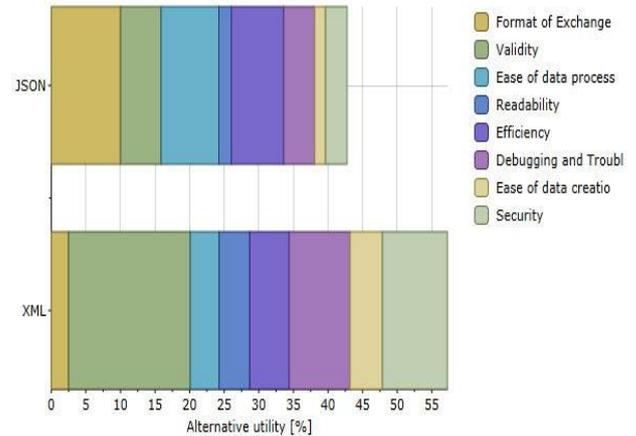


Figure 2: XML and JSON Attributes ranking graph

4.5.2. Values allocation table

Values assign to criteria's in alternatives of XML and JSON as shown in Table 3, i.e. format of exchange 2.5 value for XML and 10.01 value of JSON for same criteria.

Table 3: Values allocation to criteria of XML and JSON

Criteria	XML	JSON
Format of Exchange	2.5	10.01
Validity	17.6	5.87
Ease of data process	4.17	8.34
Readability	4.41	1.81

Efficiency	5.51	7.57
Debugging and Troubleshooting	8.8	4.48
Ease of data creation	4.67	1.56
Security	9.38	3.13
Total	57.23	42.77

4.5.3. Weight Chart

Criteria's Weight distribution shown in Figure 3, Validity has the highest value and lowest value for Ease of data creation.

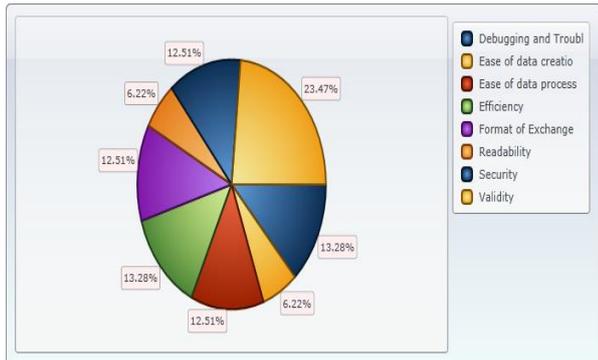


Figure 3: Weight chart of XML and JSON attributes

4.5.4. Alternative Chart of XML and JSON

In readability, security, validity and debugging and troubleshooting XML have edge over JSON, while JSON better in format of exchange, efficiency and ease of data processing than XML as shown in Figure 4.

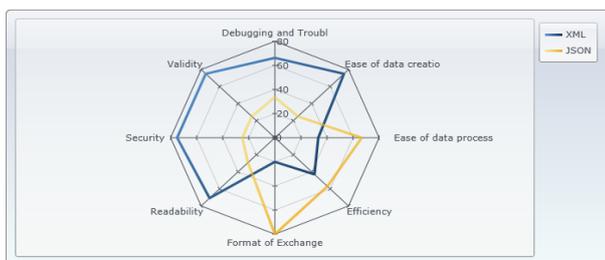


Figure 4: Alternatives chart of XML and JSON

4.5.5. Efficiency ranking graph

For each criterion in alternatives have ranking graph, as figure 5 show graph for efficiency further extended to sub criteria's of user CPU utilization, system CPU utilization and memory utilization. The graph shows that JSON best in user CPU utilization than XML, while in system CPU utilization XML perform well than JSON.

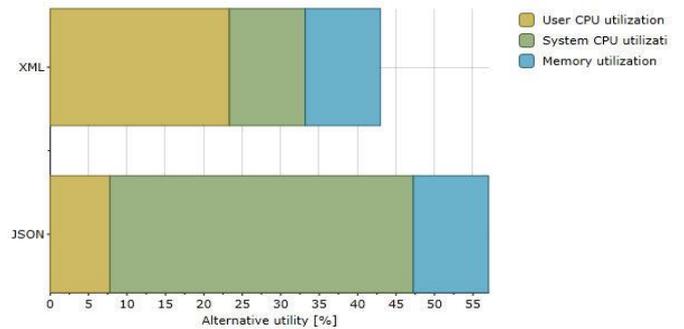


Figure 5: ranking graph for efficiency of XML and JSON

4.5.6. Efficiency Alternation chart

Alternative chart for efficiency in figure 6 shows that system CPU utilized well by alternative JSON, XML efficiently use user CPU and in memory utilization both alternative are same.

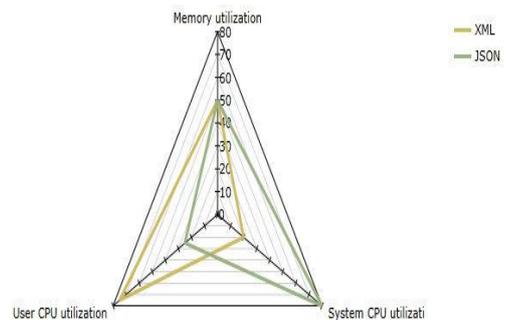


Figure 6: Alternative chart for attribute of XML and JSON Efficiency

5. Conclusion

From the above comparison it is clear that both technologies JSON and XML have their own advantages and drawbacks and it is also clear that both can be use according to the need of the system. The research study, understanding and facts from the above graphs and results conclude that XML have a bit edge over JSON for web technologies. According to my experience and above description I will advice that combination of both should be used depending on the requirement and demand. Both technologies have good properties for different situation as discussed above.

ACKNOWLEDGEMENT

This work was supported by the Research Centre of College of Computer and Information Sciences, King Saud University. The authors are grateful for this support.

References

- [1] Abdullah S. Alghamdi, H. U., Syed Usman Ali (2011). Evaluating Chaos-based vs. Conventional Encryption Techniques for C4I System. International Conference on Computer Communications and Networks (ICCCN 2011). Lahaina, Hawaii, USA.
- [2] Abdullah Sharaf Alghamdi, I. A., Muhammad Nasir (2010). Evaluating ESB for C4I Architecture Framework Using Analytic Hierarchy Process. Software Engineering Research and Practice Las Vegas, Nevada, USA
- [3] Alexander (2007). "JSON Pros and Cons ". Retrieved April 25, 2012, from <http://myarch.com/json-pros-and-cons>.
- [4] Allamaraju, S. (2006). "JSON vs XML." Retrieved May 28, 2012, from <http://www.subbu.org/blog/2006/08/json-vs-xml>.
- [5] Bruno Gil, P. T. (2011). Impacts of data interchange formats on energy consumption and performance in smart phones. OSDOC '11 Proceedings of the 2011 Workshop on Open Source and Design of Communication, NY, USA ACM.
- [6] C. Zakas, N. M., J. and Fawcett, J. (2006). Professional AJAX, University of Huddersfield.
- [7] EICHORN, J. (2006). Understanding AJAX United States: prentice hall, University of Huddersfield.
- [8] Esposito, D. (2007). Introducing Microsoft ASP.NET AJAX Microsoft Press, University of Huddersfield
- [9] Iftikhar Ahmad, A. A., Abdullah Sharaf Alghamdi (2010). "Evaluating Intrusion Detection Approaches Using Multi-criteria Decision Making Technique." International Journal of Information Sciences & Computer Engineering (IJISCE) 1(1): 60-67.
- [10] JSON (2005). "JSON: The Fate-Free Alternative to XML ". Retrieved April 10, 2012, from <http://www.json.org/xml.html>.
- [11] Jung, F. (2000). "Backgrounder technology and application ". Retrieved February 06, 2008, from http://www.softwareag.com/xml/about/e-XML_Backgrounder_WP03E0700.pdf.
- [12] Khalid Alnafjan, T. H., Gul Faraz Khan, Hanif Ullah, Abdullah Sharaf Alghamdi (2012). "Comparative Analysis and evaluation of Software Security Testing Techniques." International Archive of sciences journal.
- [13] Nurzhan Nurseitov, M. P., Randall Reynolds, Clemente Izurieta (2009). Comparison of JSON and XML Data Interchange Formats: A Case Study. 22nd International Conference on Computer Applications in Industry and Engineering, Hilton San Francisco Fisherman's Wharf, San Francisco, California, USA, CAINE.
- [14] Peng Wang, X. W., Huamin Yang (2011). Analysis of the Efficiency of Data Transmission Format Based on Ajax Applications. Information Technology, Computer Engineering and Management Sciences (ICM). Nanjing, Jiangsu, IEEE. **4**: 265 - 268.
- [15] Sporny, M. (2010). "Web Services: JSON vs. XML." Retrieved June 02, 2012, from <http://digitalbazaar.com/2010/11/22/json-vs-xml/>.
- [16] Ullman, C. a. D., L. (2007). Beginning AJAX Wrox press, University of Huddersfield
- [17] Wang, G. (2011). Improving Data Transmission in Web Applications via the Translation between XML and JSON. Third International Conference on Communications and Mobile Computing. Qingdao, IEEE: 182 – 185

System for the detection earthquake victims – construction and principle of operation

C. Buzduga, A. Graur, C. Ciufudean and V. Vlad

Abstract—This paper presents a system for detecting and rescuing victim's natural disasters. The system has three components: device for victim detection, device for detecting the number of persons in a building and PC interface that will provide information about the persons rescued or died and where they were found directly on the Internet or Data base.

Keywords—detection, earthquake, electrostatic sensor, natural disasters, receiver, transmitter, victims.

I. INTRODUCTION

When we say natural disasters we can mention the following events: earthquakes, landslides etc. The term earthquake or seism is a word used for earth movements that make up the vibrations generated in the internal areas of the Earth, which are propagated in the form of waves in the upper lithosphere. These results in the movement of tectonic plates and is often caused by volcanic activity. The earthquakes at larger scale are very strong disasters that can destroy buildings and construction and can generate various natural disasters such as landslides. Also underwater earthquakes that can cause formation of giant waves, so-called tsunami, which reach up to 30 feet tall and reaching speeds of 800 Km/h. Saving the victims is achieved by classical methods which are slow, chances of finding survivors decreased significantly. To this end we thought a simple and efficient detection much easier and the shortest time victims of such events [1], [2]. Other methods used for prevention and/or salvation of population in the earthquake action area are related in literature as using sonic methods, infrared sensors, technology of Geographical Information System (GIS), or dogs especially trained to find human being blocked in damaged buildings [3], [4]. All these rescue methods are applied for emergence decision-making on

This project is co-financed by European Social Fund through Sectorial Operational Programme for Human Resources Development 2007-2013. **Investing in people!**

C. Buzduga is assistant professor Faculty of Electrical Engineering and Computers Science, University "Stefan cel Mare", Suceava, Romania, 720229, e-mail: cbuzduga@eed.usv.ro.

A. Graur is professor Faculty of Electrical Engineering and Computers Science, University "Stefan cel Mare", Suceava, Romania, 720229, e-mail: Adrian.Graur@usv.ro.

C. Ciufudean is associate professor Faculty of Electrical Engineering and Computers Science, University "Stefan cel Mare", Suceava, Romania, 720229, e-mail: calin@eed.usv.ro.

V. Valentin is lecturer Faculty of Electrical Engineering and Computers Science, University "Stefan cel Mare", Suceava, Romania, 720229, e-mail: vladv@eed.usv.ro.

the earthquake prevention and life rescue in cities. We notice that previous methods related in the literature are different from our method as up to our knowledge such method is new.

The first part of this paper presents construction hardware and principle to functionality the system for detecting victim's natural disasters.

The second part of this paper presents a model realized in Matlab & Simulink for the electrostatic sensor.

II. SETUP SYSTEM

The system detection victims of an earthquake are an experimental acting on the basis that if you knows the exact position of a man save to become faster. This system works differently than conventional systems search for victims because it requires implementation in buildings before a disaster so that when this occurs, the system will provide rescuers the location and the exact number of survivors. Block diagram for this system showed in figure 1.

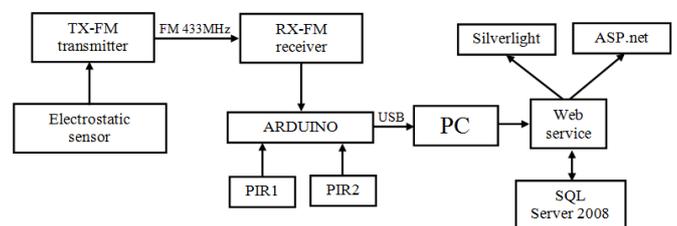


Fig. 1 Block diagram of the system

The survivor's detection is done using a special device designed to survive the collapse of a building to fire immediately after. It consists of a metal cylinder with two compartments isolated from each other. The first compartment is an electrostatic sensor that detects if there are people living in its range. This information is transmitted via radio transmitter TX-FM-TWS-DS, which is located in the second half of the cylinder. The device begins to operate with the production of an earthquake, which causes breakage of conductive liquid ampoules, which in turn is still in the metal cylinder [5], [6].

Conductive fluid will link with a 3 V battery will power the device. The signal detection device will be taken by rescuers equipped with a radio receiver composed of receiver RX-FM-RWS-371 and a development kit Arduino UNO or other model, so they will know exactly where to look for survivors. This device will keep the number of people saved and the number of people found dead information to be stored in a database and displayed on the internet when the device is

connected to a PC. The system will permanently take into account the number of persons inside a building using a device consisting of PIR sensors connected to the same acquisition board containing a microcontroller to ATmega28 [7].

A. The transmitter and receiver

After an earthquake the electrostatic sensor starts to work and when he finds people have to send this information to rescuers. TX-FM-TWS-DS is a RF transmitter module which transmits the band 433 MHz in ASK modulation. The model for this transmitter is presented in figure 2.

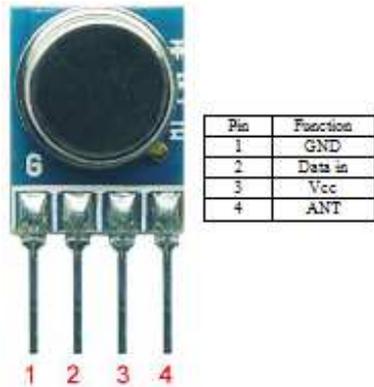


Fig. 2 Model for TX-FM-TWS-DS

Although the transmitter is located in another section of metal casing, thus electrostatic sensor is isolated it is fed to the same battery of 3V and transmits only when the sensor detects persons. In order to transmit the needs of the input terminal voltage may vary from 2,1V to 3.2V. The electrostatic sensor had more than enough of 2.4V forward voltage thus ensuring a strong transmission without interference of any kind. RX-FM-RWS-371 is an ASK modulated digital data receiver with technical characteristics: frequency range 433.92MHz, modulate mode ASK, date rate 4800 bps, selectivity -108 dBm, channel spacing ±500KHz, supply voltage 5V and is directly compatible with the transmitter TX-FM-TWS-DS with technical characteristics: frequency Range 433.92 MHz, modulate mode ASK, date rate: 8Kbps, supply voltage 1.5~12V [8], [9].

The operations of transmission-reception module are simulated in Matlab using mathematical formulas for carrier wave (1) and ASK modulation (2) [10].

$$C(t) = A_c \cdot \cos(2\pi \cdot f_c \cdot t) \tag{1}$$

and

$$S_{ASK}(t) = m(t) \cdot C(t) = m(t) \cdot A_c \cdot \cos(2\pi \cdot f_c \cdot t) \tag{2}$$

where: $C(t)$ – carrier wave
 A_c – amplitude

f_c – frequency carrier
 t – time
 $m(t)$ – signal modulation

For example we obtained the operations of transmission-reception module in ASK modulate and demodulate mode in Matlab. This graphic is presented in figure 3.

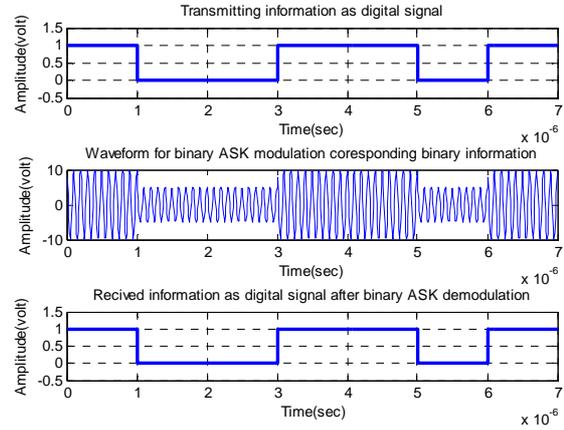


Fig. 3 The operations of transmission-reception module.

This receiver was chosen in part because the work requires a supply voltage of 5V can be supplied easily by development kit Arduino UNO, thus eliminating the need for extra batteries. The model for this receiver RX-FM-RWS-371 showed in figure 4.



Fig. 4 Model for RX-FM-RWS-371

B. Electrostatic sensor

In order to detect people under the ruins of a building our system uses a sensor to be placed inside of an electrostatic metal cylinder to not be affected by the earthquake. Electrostatic sensor is showed in figure 5. The JFET is the sensor of electrostatic field and the other components ensure the gain and adapt the impedance with the wireless transmitter TX-FM-TWS-DS mentioned before.

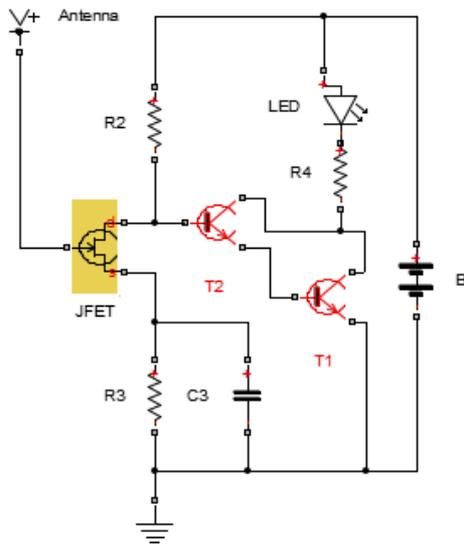


Fig. 5 Electrostatic sensor

It is noted using a field effect transistor (JFET), which presents the rest of the source and drain junction, crossing resistance of very low value of 200 Ω. This transistor has the following output characteristics:

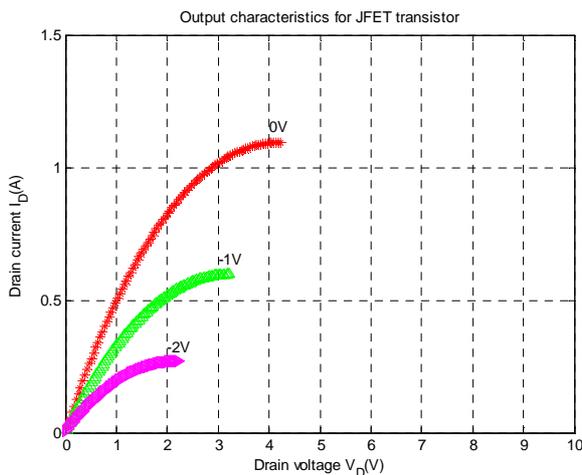


Fig. 6 Output characteristics for JFET transistor

The two NPN transistors with silicon are capable of blocking polarized as they are to drive. Existence 5 kΩ resistor is not sufficient to ensure conduction, because it forms the junction with resistance source - drain field effect transistor, a voltage divider with a voltage to the base of transistor T₂, only 100...200 mV, opening times for transistor T₁ and T₂ respectively coupled Darlington require a voltage of 0.5 V.

This occurs when an electrostatic field is applied gate field effect transistor; when the junction between the source and drain resistance increases considerably, therefore the voltage divider is changed, leading to the opening of transistors T₁ and T₂, resulting in LED lighting, or if our system start transmission transmitter.

The songs used are: JFET can be a BF 245 and BF 256, TIS 34 or any equivalent. The transistor T₁ and T₂ may be of BC

170...173 or equivalent or BC 517 which is a direct Darlington transistor. Electrolytic capacitor must be of good quality without loss [5], [6].

The installation is done on a plastic plate coated with copper foil. For antenna sensor using a piece of copper wire, insulated with polyvinyl of 0,35..1 mm diameter and 10..15 cm long, straight stick with one end terminal of FET's gate. To obtain a higher sensitivity, you can try using a longer antenna, up to half a meter long, extending sensitivity 2...3 meter radius around the antenna.

When making commissioning of the installation, if the pieces are of good quality, light turns on and stays on for a minute, because strong excitation caused by the presence of the user. When he departs, the pilot light goes out; but any approach or production of electrostatic field sensor reacts immediately by switching lights that remain lit as long as the sea was proportionately electrostatic field excitation [10]-[12].

Antenna

C. PIR sensor

To detect the number of people in a room I used a pair of PIR sensors positioned at the entrance so that in the order of their activation system will deduce how many people are in a room. Infrared motion detection may be performed by using:

- a) Infrared barrier
- b) PIR passive sensor

PIR Passive Infrared (PIR), also known as thermal infrared detects the natural radiation emitted by warm objects. Also extremities moving bodies emit infrared radiation more passive than the background that is. People living beings in general and cars with hot engines emit thermal radiation that a PIR detector senses both day and night. The passive infrared radiation should not be confused with the near infrared emitted by remote controls TV; Passive Infrared does not emit any radiation that could be harmful, but we can say that it is a "dark view". Passive infrared radiation is detected by the pyro-electric sensor. This sensor detects changes in temperature of up to one-thousandth of a degree of 10 m caused by the movement of a person. Pyro-electric sensor detectors can be equipped with different types of lenses: volume, Pet, or type curtain corridor, in addition, some detectors and detection using microwave radar principle, be used when there is a risk of false alarms.

A PIR sensor comprises an optical system and a Fresnel lens made of semiconductor crystal that generates electrical charges on their surface when subjected to heat caused by the infrared radiation with a wavelength specific warm-blooded bodies. The tasks collected on the surface of the crystal are applied to a first-stage transistor amplifier with an FET made usually with the sensor encapsulated in a capsule TO5. The capsule is provided with a window covered with a filter which passes radiation in the infrared range (4..8) μm, giving a minimum attenuation for λ=9.4μm wavelength infrared radiation characteristic organisms "Blood warm ". Diagram for PIR sensor is showed in figure 7.

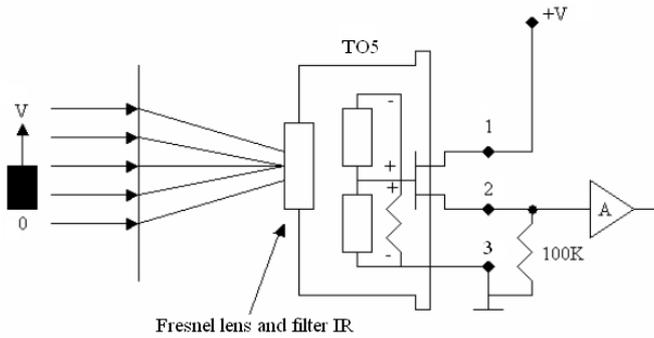


Fig. 7 Diagram for PIR sensor

The drain of the transistor FET is connected to a DC voltage potential and steady and the filter 3 to 15 V and the source is connected to a load resistor of 100 KΩ. The pass band of the amplifier thus produced is limited to 10 Hz, in order to remove high frequency noise. A strong disturbance may generate false alarms, is solar infrared radiation, which is captured as a result of unwanted reflections from the environment. To remove this source of interference, which fortunately has a little speed variation, PIR sensor consists of two crystals that are serially connected differential. The radiation collected by the two crystals in the same point in time, is canceled due to the differential series. In the case of a hot moving body, the two crystals are interlaced in turn, give the information useful here [12], [13].

The optical system for the PIR sensor is a combination of a Fresnel lens made of a plastic film generally translucent white with a thickness of 0.4... 0.5 mm. Specifically Fresnel lens has a series of concentric grooves inclined walls. To understand the role of Fresnel lens optical system made to imagine a source of infrared radiation that travels at a velocity V_0 parallel to the surface S of a PIR sensor. In front of the sensor, suppose there is a C body opaque to infrared radiation. The body C is located a very short distance from the surface of the sensor S is equal with 20...30 mm. In this situation to approximate that both source and the shadow projection on the surface S , moving on circular trajectories with radii R respectively r . How angular velocities of the source are equal that shadow moving in the opposite source speed is:

$$V_v = V_0 \cdot \frac{r}{R} \tag{3}$$

Presence in front of the sensor to a Fresnel lens that focuses the image source S its surface will cause the appearance of a slide "warm" the source, superimposed the background of "cold" ambient, just as many moving shadows generated for your body C is in the form of a grating. If the ambient background is as warm as the source PIR sensor will not cause any electrical signal.

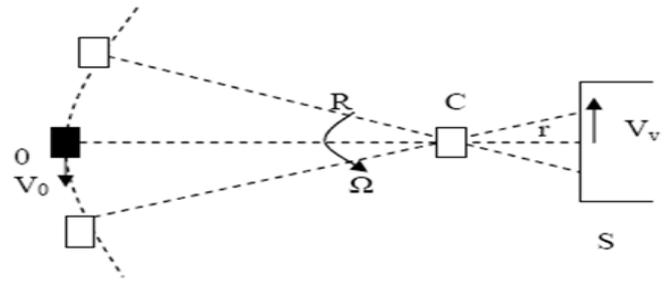


Fig. 8 Trajectory IR rays

The infrared motion detection's principle is precisely scanning performed using Fresnel lens sensor surface contrasts between hot and cold. Channels shape, their number, distribution several Fresnel lenses allow obtaining fascicular-type features, the curtain, etc.

Regardless of the type or characteristic of the detector itself (peripheral alarm system, lighting shutter, door control sensory element to shopping centers, etc.) in most cases using one or more PIR followed floors higher amplification factors of the amplifier and an output validation number of applications depending on the ambient light.



Fig. 9 PIR sensor

D. The Arduino UNO Kit

For information processing system uses the kit Arduino UNO which contains a microcontroller from ATmega328. It takes the information sent by the device to detect people via the receiver RX-FM-RWS-371 such alarm that gives survivors and recover as they are storing this information that is transmitted over the Internet when connected to a PC. Arduino UNO is also responsible for gathering information from the PIR sensor information with which determines how many people were in the building at the time of the earthquake and this information will be transmitted over the Internet.

Arduino UNO is a processing platform open-source software and hardware based on flexible and easy to use. Platform consists of a small (6.8 cm x 5.3 cm - the most common one) built around a signal processor and is able to retrieve data from the environment through a series of sensors and perform actions environment through lights, motors, actuators, and other mechanical devices. The processor is able to run code written in a programming language which is very similar to C++.

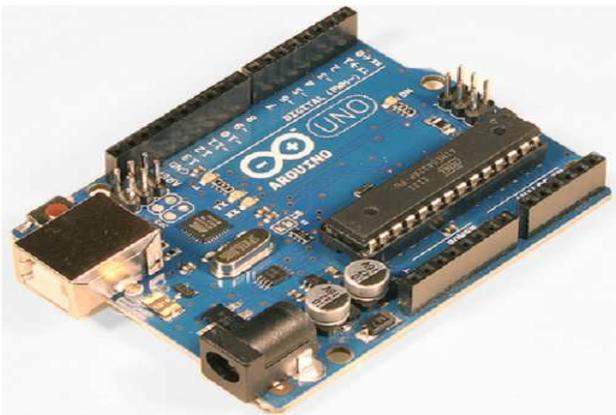


Fig. 10 Arduino UNO

This kit has 14 digital pins of which 6 PWM and 6 analogical pins. Maximum speed program execution allows you to perform hardware and software testing in real time [20].

III. THE SIMULATION ELECTROSTATIC SENSOR

The simulation is realized in Matlab & Simulink using the circuit in figure 11, [16], [17], [21].

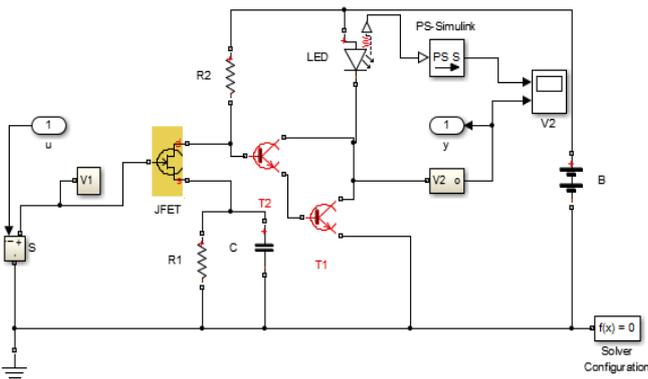


Fig. 11 Circuit for simulation in Matlab & Simulink

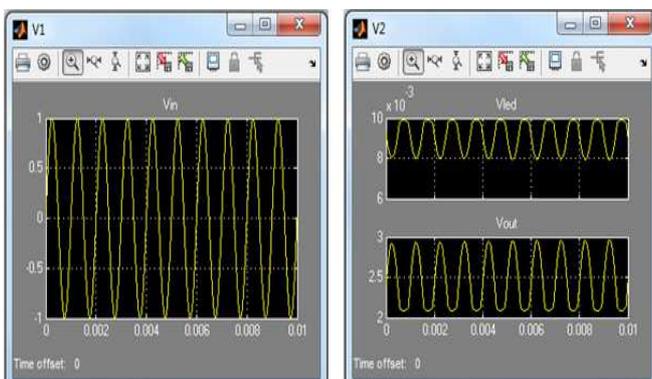


Fig. 12 Signal input, signal for LED and signal output

Simulation results are shown in figure 12. In order to verify the system performance bode diagrams we plotted in figure 13 [21]. From this representation closed loop system is stable

with temperature and differs according to humidity and electrostatic field intensity.

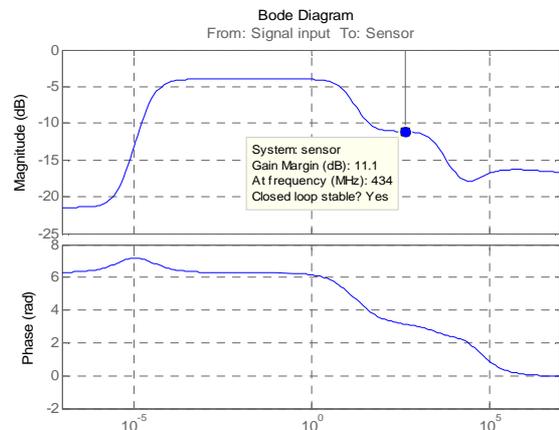


Fig. 13 Bode diagrams for this electrostatic sensor

The device continuously transmits the number of people detected a database that will be displayed on the internet. Web page will be submitted this information graphically displays the relative position of all persons in the room, both those found in life and what could not be saved. In figure 14 is presented a print screen for the internet interface.



Fig. 14 Internet interface.

IV. CONCLUSION

The detection of victims of an earthquake, or land gliding is a revolutionary device's acting differently from existing similar. The principle that was designed argues that when a person knows the exact position of his salvation becomes faster. The difference with classical search systems is to implement his victims in buildings so that when a disaster occurs, the system acts rescuers providing vital information: the location and the exact number of survivors will always be available to anyone. This system is simple and has a low cost with a minimum set of components.

The development of this system will be focused on improvement of placement mode of the electrostatic sensor.

This problem is very important because the position as effectively will give more accurate results.

It should be adapted to the dimensions of the room where the sensor is placed, through the increase or decrease the radius of action. If it is not possible the increase or decrease enough the radius of action, thus four sensors may be mount in each corner of the room and then we have a network of sensors, this method is useful in large buildings. It is also necessary to find a method for the sensor to consume power just when it is necessary.

ACKNOWLEDGMENT

This paper has been financially supported within the project entitled „**SOCERT. Knowledge society, dynamism through research**”, contract number POSDRU/159/1.5/S/132406. This project is co-financed by European Social Fund through Sectorial Operational Programme for Human Resources Development 2007-2013. **Investing in people!**”

REFERENCES

- [1] S. Choi, “The Real-time Monitoring System of Social Big Data for Disaster Management”, *Advances in Automatic Control*, 2014, pp. 110-114.
- [2] R. Villaverde, *Fundamental concepts of Earthquake Engineering*, CRC Press, 2009.
- [3] MA Hao-ran Feng, Qi-min Mo, Shan-jun, Research on the decision support system of the urban seismic emergency, *Journal of World Information On Earthquake Engineering*, 2005 vol. 1, pp.45-53.
- [4] *Huo En-je, Application of GIS in emergence decision-making of urban earthquake prevention and disaster reduction, Journal of Natural Disasters*, 2003, vol.3, pp. 12-124, ISSN: 1004-4574.
- [5] S. Xing, S. Chen. *Unifying Electrical Engineering and Electronics Engineering*, House edition Springer, 2014.
- [6] J. Jaimes-Ponce, I.I. Siller-Alcalá, “Hardware-Software System for laboratory experimentation in electronic circuit”, *Advances in Circuits, Systems, Automation and Mechanics*, 2012, pp. 126-130.
- [7] O. Krini, M. El Bahri, J. Börsök, “Development of Safety Electronic-Components, Devices and Systems-Based on Safety Standard”. *Proceedings of the 12th WSEAS International Conference on Circuits, Systems, Electronics, Control & Signal Processing (CSECS '13)*, Budapest, 2013.
- [8] I. Petrescu, M. C. Surugiu, “Traffic Data Transmission Using Wireless Sensor Networks (WSN) Principles.” *Proceedings of the 16th International Conference on Automatic Control, Modelling & Simulation (ACMOS '14)*, Brasov, 2014.
- [9] L. Merad, F. T. Bendimerad. ”Neural Networks for Synthesis and Optimization of Antenna Arrays”, *Radioengineering*, Vol. 16, No. 1, 2007, pp. 23-30.
- [10] N. Vlajic (2010). *Analog Transmission of Digital Data: ASK, FSK, PSK, QAM, (Garcia 3.7)*, Monograph source online, Available: http://www.eecs.yorku.ca/course_archive/2010-11/F/3213/
- [11] K. D'hoel, A. Van. Nieuwenhuysse, „Influence of Different Types of Metal Plates on a High Frequency RFID Loop Antenna: Study and Design”, *Advances in Electrical and Computer Engineering*, Vol. 9, No. 2, 2009, pp. 3-8.
- [12] M. Sarevska, N. Mastorakis, “Neural Networks and Antenna Arrays”, *Recent Researches in Circuits, Systems, Electronics, Control & Signal Processing*. Athens, Greece December 29-31, 2010, pp. 122-127.
- [13] L. Merad, F. T. Bendimerad, ”Neural Networks for Synthesis and Optimization of Antenna Arrays”, *Radioengineering*, Vol. 16, No. 1, 2007, pp. 23-30.
- [14] T. M. Jamel, “Performance Enhancement of Smart Antennas Algorithms for Mobile Communications System”, *International Journal of Circuits, Systems and Signal Processing*, vol. 8, 2014, pp. 313-320.
- [15] F. Hruska, “Electromagnetic interference and environment”, *International Journal of Circuits, Systems and Signal Processing*, vol. 8, 2014, pp. 22-29.
- [16] S. T. Karris, *Electronic Devices and Amplifier Circuits: With MATLAB/Simulink/SimElectronics Examples*, Orchard Publication, 2012.
- [17] S. T. Karris, *Engineering Electromagnetics with Introductions to S-Parameters, RF Toolbox, and SimRF*, Orchard Publication, 2014.
- [18] Alan McCartney, *Static Electricity & Relative Humidity*, *Fireline* www.asse.org 2012.
- [19] C. Ungureanu, C. Bobric, D. Daniela, “Fuzzy Logic Control of a New Type of Electromagnetic Converter with Rolling Rotor”. *12th International Conference on Applied and Theoretical Electricity, (ICATE)*, 2014 October 23-25, Craiova, pp. 1-4.
- [20] ***www.arduino.cc
- [21] ***www.mathworks.com

Corneliu Buzduga was born in Vicovu de Sus, Suceava, Romania in 1981. In 2012 he became Ph.D. in the field Electrical engineering at University “Stefan cel Mare”, Suceava, Faculty of Electrical Engineering and Computers Science. In present C. Buzduga is assistant professor at the same University.

Question-Answering Systems in the Specific Domain of E-Government

A. Beltrán, S. Ordoñez, S. Monroy, L. Melo and N. Duarte

Abstract— Proposals for answer searching systems or Question-Answering Systems (QAS) are diverse and there is a variety of both the problem at hand, and the variables involved: different user profiles, heterogeneous data sources, implemented technologies, specific application domains and the volume of data. An architecture for a specific domain SQA in e-Government which may be extended to various state agencies is proposed. This proposal is based on an analysis of different frameworks for the development of systems generating answers to questions posed in natural language and adds specific items to the context of e-Government. As an example of evaluation of this architecture an adaptation to a small town of the Colombian state is shown.

Keywords—e-Government, Software Architecture, Question-Answer Systems (QAS).

I. INTRODUCTION

A QAS provides specific answers to questions. Therefore, if a user launches a question on the web, such as "Who is the president of Colombia?" the system should provide as an answer "Juan Manuel Santos," and a set of links to pages where it is possible to find information about the presidents of Colombia and further information about the current president. This is a very different operation from the one performed by current search engines such as Google or Yahoo, where the response is a set of references to websites which the user must filter and inspect for himself to find the information he is looking for (information retrieval task). In addition, these engines also retrieve information that is unrelated to the object of the search and they make more complex the task for the user.

II. CONTEXT

QAS are a type of Information Retrieval System, which process a question in natural language and return or extract an answer from structured (databases) or unstructured (text) sources. Unlike information retrieval systems, queries are performed in natural language and not in keywords, which means that the system must recognize the type of response expected by the user to be able to give the specific answer or a paragraph where the answer might be found. In general, the complexity of these systems lies in establishing the implied relationship between the question and the answer. There are approaches that considered them as systems leading to the

Semantic Web [1].

Following is a description of some frameworks available for QAS, where their typical architectures are reviewed and then analyzed and compared, in order to propose an architecture for the specific domain of interaction with the Colombian state

A. Typical Architectures

A typical first architectural approach for QAS has four main components: the Question Classifier component that is fed by a question in natural language, and is responsible for identifying the kind of question (e.g. what, where), the type of response expected, the focus of the response and its respective semantic context. Next comes the Answer Document component, which is responsible for the search and identification of relevant documents where the answer is expected to be found; this search is performed on unstructured data sources that can be found in various formats. The next one is the Extract Candidate Answer component, in which the possible responses found in the relevant documents or sources selected by the above components are identified. The last component of Answer Selection is perhaps the most important since it is this component the one that generates an answer to the question [2].

A second architectural approach has three main components: a first component of transformation of the question (Question Processing Component) that performs tokenization and tagging tasks, identification of keywords, grammatical analysis of the question, identification of ambiguous words, identification of the type of expected response and expansion of keywords. A second component of search (Search Component) generates queries that will be executed by the search engine on the data sources. And finally the Answer Extraction Component performs tasks of filtering the data found by the search component, entity recognition, response identification and validation [3]

B. Question Answering Frameworks

In this article we will consider the following Question Answering frameworks:

QALL-ME. Provides an architecture consisting of three major modules, reusable and extensible for building QAS on structured data for the specific domain of tourism. This domain is modeled by an ontology used as a primary resource for the use of multilanguage. It also uses spatial and temporal reasoning at the time of processing the question, and reasoning in the classification of responses to determine the most appropriate response to the question analyzed [4].

AQUALOG. Uses a sequential process model, in which the input in natural language is first translated by the linguistic component (Linguistic & Query Classification) into a set of intermediate representations called triple (subject, predicate and object). The linguistic component obtains a set of syntactic annotations associated with the input query (NL Sentence Input) to classify the query. Then the service of similarity ratio (Relation Similarity Service) takes as inputs the triplets to produce queries compatible with the ontologies, called onto-triplet. The inference engine is activated when the onto-triplet is valid, which explores the knowledge base for an answer, otherwise, when the onto-triplet is invalid, a disambiguation is needed by the user, to obtain a valid onto-triplet [5].

QANUS. The QANUS architecture framework adopts a segmented QA approach, dividing the task of quality control in several sub-tasks, including the preparation of the base information, question processing, response retrieval and evaluation [6].

FALCON. It is an open domain search engine, consisting of two main axes; first, the methods of natural language processing (Question Processing and Paragraph Processing) used to identify the semantics of questions, in order to identify candidate responses within the collection of texts. These methods are specially designed with information retrieval techniques intended to retrieve all texts out of relevant paragraphs. Second, in order to extract the correct answer (Answer Processing), the bag of words approach is not enough, in contrast it uses methods of natural language processing that are enriched with pragmatic knowledge to filter out incorrect answers [7].

to help you gauge the size of your paper, for the convenience of the referees, and to make it easy for you to distribute preprints.

III. PROBLEM

The Colombian state seeks to define new channels of communication with citizens to help improve the perception they have of the state and improve the processes of participation, efficiency, effectiveness, and transparency. In this regard, government web portals of state agencies have increasingly supported this process of communication and interaction including information on procedures and services in the web pages. As a result, citizens look for information on procedures, request services and conduct transactions in those websites. Searching for this information is usually performed with traditional navigation menus inside web pages or through options based on keywords searches. Due to the quantity of

information available, it is easy for citizens to get lost in the search for accurate information that meets users' needs. As evidence of this problem can be considered the low use of Internet to interact with government, as shown in the results of the Internet national survey conducted by the national statistical agency (DANE):

The proposed system was implemented on a QA system to improve the interaction between citizens and the state. In this implementation resources from Colombian state agencies and information from the web pages of the small town's local government were used.

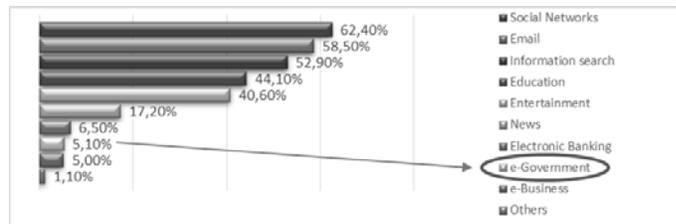


Fig. 1 percentage of Internet Users 2013, Colombia [8]

IV. DEVELOPED SYSTEM

The developed system includes classification techniques for locating both the questions and the answers helped by a defined domain specific ontology for the government. The most important components include question classification, answer classification, and FAQ finder. For text classification two corpus, one of questions and the other of answers, were compiled and organized in 11 categories. The main component consists of text classification based on the categories of the defined taxonomy. For the tests a corpus was constructed from questions posed by actual citizens on different services. Regarding accuracy for the question classification module an average accuracy of 90 % was obtained, for identification of texts of answers related with questions the accuracy obtained was 73 %, and the overall accuracy of the system was 70%. Another component of the system is based in a collection of 8497 FAQs downloaded from government websites. The system used the criterion of cosine similitude to determine if the question posed by the user coincided with the questions in the FAQs.

The FAQ corpus was extracted from all agencies of the state, since from the point of view of citizens the state constitutes a unity. This was evidenced in a sample of questions collected from a section of the population, which showed that most questions are related to general procedures of the state and a lower percentage to local procedures.

A. Definition of EAT Taxonomy (Expected Answer Type)

In an initial phase it seemed appropriate to use a CLEF taxonomy for the system[9]. In compiling the questions it became evident that the taxonomy was not the most appropriate because most user questions were open type that did not match the given taxonomy category. It became necessary to perform a semantic analysis of the texts of the

questions in order to build a taxonomy that would fit most questions posed by users. The proposed taxonomy is showed in Table 1, in which the intentions correspond to the information expected by the user, government sectors to the entities authorized to provide answers, and the EAT taxonomy is the result of combining intentions and government sectors. In general, responses to questions from citizens are not factual, so that the taxonomy proposed by CLEF is not suited to the specific domain.

TABLE I.
EAT TAXONOMY

Intentions	Government Sectors	EAT Taxonomy
AGRARIAN	SUPPORT	AGRARIAN_SUPPORT
EDUCATION	FUNCTION	AGRARIAN_FUNCTION
PROTECTION_SOCIAL	PROCEDURE	EDUCATION_SUPPORT
		EDUCATION_FUNCTION
		EDUCATION_PROCEDURE
		PROTECTIONSOCIAL_SUPPORT
		PROTECTIONSOCIAL_PROCEDURE

B. Question Corpus

The questions in the corpus were collected from the questions that a group of citizens made about government procedures. After a review process 400 questions were selected. There was an effort to consider the intentions of the questions of citizens, which allowed testing the Question-Answer system in a more accurate and complete manner. From the morphosyntactic analysis questions were classified into different categories and intentions according to the EAT taxonomy proposed and questions were divided for training and testing in order to build and test the question classifier.

TABLE II.
QUESTION CORPUS SAMPLE

EAT Category	Question Sample
AGRARIAN_FUNCTION	What is the objective of the agrarian fund from the government.
EDUCATION_SUPPORT	How does the government support students who get a high score on the entrance exam.
PROTECTIONSOCIAL_PROCEDURE	What are the prerequisites to enter the Families in Action program.

C. Answer Corpus

A similar procedure to the construction of the corpus of questions was followed in collecting the texts that provide answers to these questions: they were downloaded from the pages of the entities authorized to provide the answers. Table 3 shows examples of answers corresponding to the EAT categories.

TABLE III.
ANSWER CORPUS SAMPLE

AGRARIAN_ FUNCTION	“The objective of the Agricultural Guarantee Fund is to support the rediscounting loans granted by FINAGRO through special programs of agricultural development aimed at financing agricultural and rural sector projects that are technically, financially and environmentally viable, and are granted to producers who cannot ordinarily provide the guarantees required by the credit granting institutions...” https://www.minagricultura.gov.co/atencion-ciudadano/preguntas-frecuentes/Paginas/Apoyos-Directos.aspx
EDUCATION_ SUPPORT	“The most outstanding graduates in the State Examination Icfes Saber 11 as applied in 2010, have a great incentive in the Ministry of Education Resolution No. 987 of February 15, 2010, which grants the Distinction Andrés Bello to the outstanding students in this test, and the recognition of the country's top 50 scores” http://www.mineducacion.gov.co/cvn/1665/w3-article-265372.html
PROTECTION_ SOCIAL_ PROCEDURE	“The DPS website will also be available to speed up registry in the Families in Action program. Through this service the interested parties can visit www.dps.gov.co and on the right side of the page they will find Appointment for More Families in Action.” http://www.urnadecristal.gov.co/gestion-gobierno/as-puedes-inscribirte-a-m-s-familias-en-acci-n-en-bogot

D. Architecture

The proposed architecture is made of components that facilitate the reuse and extension in the different institutions of the state, according to the philosophy of interoperability that operates in them. It distinguishes between resources (previously described) and components for system operation. Some resources are preprocessed previously to improve performance, the resource name and the contents of the corresponding folders of EAT (Expect Answer Type) can be modified as convenient, in order for it to be easily customizable and reusable. The collection of resources (Resources-Questions, Resources-Answers, Resources-Search, among others) used by the system will be a process parallel to its implementation and will depend on the extent one wishes to impose on the system.

The principal process is called QuestionProcess, which receives the question in natural language and initially applies the NLP component to POST (Part Of Speech), expansion terms with EuroWornet, Ontologies and RDF in specific domain. Then the processed question passes on to the QuestionClassification component in order to perform the search in the collection of documents both in pre-processed and online documents. In the next step, texts filtered in the previous phase are sent to the ExtractionAnswer component which applies the AnswerClassifier to determine the text that answers the user's question in the most accurate manner.

Figure 2 shows the overall proposal from the viewpoint of architecture (components) for the QA system.

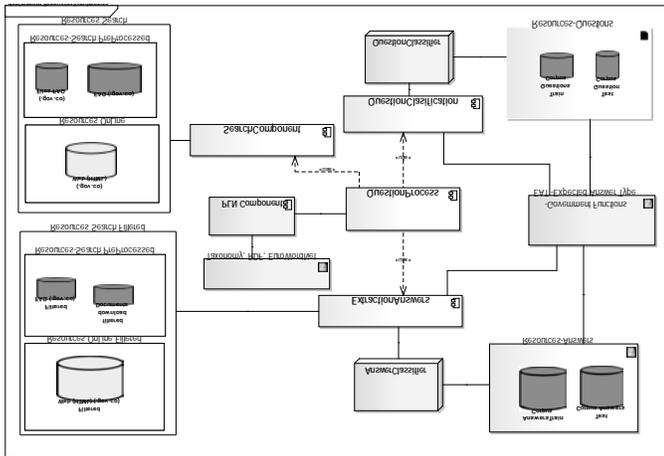


Fig. 2 Proposed architecture. Source: authors

TABLE IV. ARCHITECTURE'S COMPONENTS DESCRIPTION

Resource	Description
EAT (Expect Answer Type)	Defined from the state organization and its services to citizens.
Resources-Questions	Contains all resources related to the corpus of questions according to the categories of the defined EAT. Used to train, test and build the QuestionClassifier
Resources- Answers	Contains the corpus of government texts categorized according to defined EAT. Used to train, test and build the AnswerClassifier.
Resources-Search	Resources where the SearchComponent will perform the search, in this case there are pre-processed elements (previously downloaded documents and FAQ files) and pages from the Web site bearing the domain .gov
Resources –Search Filtered	Items from the Resources-Search that went over the accuracy threshold of the SearchComponent. The ExtractionAnswer component applies the AnswerClassifier to the filtered items to determine which items match the EAT of question.

E. Implementation

It was evident that with the EAT defined taxonomy (Table 1) most answers can be answered correctly and other approaches with more precise taxonomies did not give satisfactory results due in part to the ambiguity of the questions of citizens and in the complexity of related answers that do not allow short answers, but instead retrieve texts where the related questions are given. The tests showed an accuracy of 70% with a corpus collected from citizens of the local government. In Figure 3 the defined

taxonomy for the Colombian context is shown.

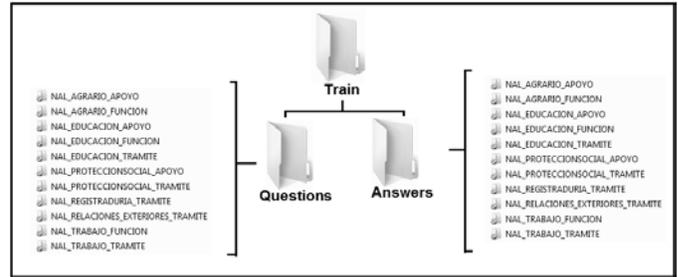


Fig. 3 Taxonomy in Implementation of the architecture in a local government in Colombia.

V. CONCLUSIONS

This architecture intends to contribute in part to the implementation of QA systems in the domain of e-government, considering that the literature review did not retrieve items that discuss the application of these systems to e-Government.

The e-Government domain relates many areas such as education, foreign affairs, health, sports and more. This requires that QA systems deal with texts written in the form of decrees and procedures and generally long, in the Colombian case, since they are subject to many conditions and specific situations, which makes the system unable to answer in short form but in paragraphs related to the questions of citizens.

Another aspect that is taken into account in this architecture is that citizens ask for specific information and procedures in these areas. In addition to the ambiguity of the questions, citizens ask about situations dependent on their actions. In many cases citizens asked for explanations within a specific area of the state, which although cannot be answered in an accurate manner it is possible to extract texts where the citizen might find the answer, thereby improving the interaction of citizens with the state, which in the Colombian case is very low in spite of the high ranking given e-government initiatives by a United Nations survey.

On the other hand, in the different experiments it was evident that the use of ontologies and mapping of terms used by citizens had a positive impact on the accuracy of answers proposed.

ACKNOWLEDGMENT

This research is financed jointly by Colciencias, Ministry of Information and Communication Technologies of Colombia (MinTic), and the consortium of the Universidad Manuela Beltrán, Multimedia Service Ltda and InnoVA&IP Ltda – Project No. 1263-595-37045.

REFERENCES

[1] M. Konopík, o. Rohlík, Question answering for not yet semantic web, in: Brno, 2010, pp. 125-132.
 [2] C. Monz, M. De rijke, Tequesta: the university of amsterdam's textual question answering system, in: Trec, 2001.

- [3] B. Magnini, M. Speranza, V. Kumar, Towards interactive question answering: an ontology-based approach, in, Berkeley, CA, 2009, pp. 612-617.
- [4] O. Ferrández, C. Spurk, M. Kouylekov, I. Dornescu, S. Ferrández, M. Negri, R. Izquierdo, D. Tomás, C. Orasan, G. Neumann, B. Magnini, J.L. Vicedo, The qall-me framework: a specifiable-domain multilingual question answering architecture, *Journal of web semantics*, 9 (2011) 137-145.
- [5] V. Lopez, V. Uren, R. Motta, M. Pasin, Aqualog: an ontology-driven question answering system for organizational semantic intranets, *web semantics: science, services and agents on the world wide web*, 5 (2007) 72-105.
- [6] J.-p. Ng, m.-y. Kan, Qanus: an open-source question-answering platform, in, 2010.
- [7] S. Harabagiu, D. Moldovan, M. Pasca, R. Mihalcea, M. Surdeanu, R. Bunescu, R. Girju, V. Rus, P. Morarescu, Falcon: boosting knowledge for answer engines, in: *Proceedings of trec*, 2000.
- [8] DANE, Encuesta nacional de calidad de vida, 2013, 2013.
- [9] D. Giampiccolo, p. Forner, j. Herrera, a. Peñas, c. Ayache, c. Forascu, v. Jijkoun, p. Osenova, p. Rocha, b. Sacaleanu, r. Sutcliffe, Overview of the clef 2007 multilingual question answering track.

Monitoring Metropolitan City Air-quality Using Wireless Sensor Nodes based on ARDUINO and XBEE

Ali Al-Dahoud, Mohamed Fezari, Ismail Jannoud and Thamer AL-Rawashdeh

Abstract—Wireless sensor networks (WSN) have been experimented in different applications including monitoring many environmental phenomena such as air quality assessment, forest fire monitoring, flood rivers control during last decade. In this paper, we propose an architecture node and a simulation interface for WSN in monitoring air quality in metropolitan cities. Nodes are equipped with gas, temperature and dust sensors, an Arduino-uno as microcontroller have been designed for air quality monitoring in some sensible area at Annaba City East of Algeria. The new design is based on Arduino-uno as microcontroller.

Comparing sensed gas from three different regions in the city to normal gas levels (for the clean air), the obtained results from the several tests and acquired data, indicate that there is a big difference in the gas levels of both gases (LPG, NO₂ and CO). However, the acquired results for the air quality control in some areas in Annaba city show no risky situation to be considered for further actions. In this work we cover the field of Air-quality monitoring electronic Nodes design and wireless transmission of fusion data. Then A GUI has been designed for simulation of the WSN in controlling the environment air Quality. Tests are encouraging; the flexibility, the presence of components and the ease of design facilitates the implementation of this system.

Keywords—AQM, Arduino-Uno, wireless sensors network, air quality in city.

I. INTRODUCTION

AIR pollution can affect many body organs and systems in addition to environment based on report of World Health Organization (WHO), air pollution is significant risk factor for multiple health conditions including: heart disease, lung cancer, pneumonia, difficulty in breathing and coughing due to aggravated asthma [3]. Wireless Sensor Networks (WSNs) technology [4] and [5] is in the front part of the investigation of the computer networks and it could be the next technologic market of with huge sum of money in investment. Sensor nodes can be fixed or mobile, they have limited processing power, storage, bandwidth, limited wireless transmission range and energy powered by battery. This

Ali Al-Dahoud and Thamer AL-Rawashdeh are with Faculty of IT Al-Zaytoonah University Amman JORDAN (e-mails: aldahoud@zuj.edu.jo, Thamer.R@zuj.edu.jo)

Mohamed Fezari is with Badji Mokhtar annaba University Faculty of Engineering Dept. Electronics (e-mail:mohamed.fezari@gmail.com)

Ismail Jannoud is with Al-Zaytoonah University Amman JORDAN, Faculty of Engineering and Technology (email- ismael.jannoud@zuj.edu.jo)

limitation makes provision of the security in sensor networks not an easy task [4]. The availability of cheap, low power, and miniature embedded processors, radios, sensors, and actuators, often integrated on a single chip, is leading to the use of wireless communications and computing for interacting with the physical world in applications such as air quality control.

In Fact, with the increasing number of vehicles on our roads and rapid urbanization air pollution has considerably increased in the last decades in Annaba city (Algeria). For the past twenty years the economic development of this city has been based on industrial activities and the agriculture industry. Hence, there has been the growth of industries and infrastructure works over the island. Industrial combustion processes and stone crushing plants had contributed to the deterioration of the quality of the air.

In this paper, we propose to use a WSN based microcontroller equipped with gas sensors have been actively used for air quality monitoring. The design included several units mainly: Arduino Microcontroller, MQ-2 Gas Sensors, and the current regulator circuit the paper is organized as follow: in second paragraph after introduction we define primary pollutants, in paragraph 3, we present the hardware proposition design with main components. In section 4, format and communication with the special sensor DHT11 is illustrated, then we conclude the paper by presenting results, discussion of simulation in section 5, finally conclusion and perspectives were included in section6.

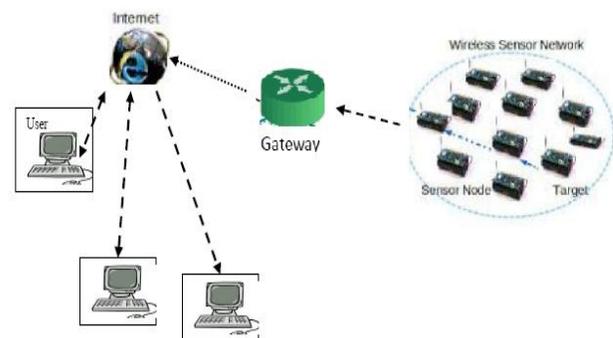


Figure 1: The Hardware Design Schematic Diagram.

II. PRIMARY AIR QUALITY POLLUTANTS

Primary pollutants are those in which the substance emitted is itself hazardous. Some primary pollutants also produce other dangerous substances after undergoing chemical reactions in the atmosphere, and these are known as secondary

pollutants. Primary pollutants include the following substances as mentioned in [15].

Particulates : This includes dust, smoke, aerosols and haze - any finely divided airborne solid material. Particulates are commonly generated by fires, motor vehicles, some industries (particularly road building, quarries and fossil fuel power stations) and various natural sources including volcanoes, plant and animal matter and dirt. Particulates are aesthetically displeasing, can irritate the eyes and cause respiratory problems. In recent years concerns have been raised about the possible health effects of 'fine' particulate matter (less than 10µm diameter). These have been shown to be associated with increases in hospitalization and even deaths from respiratory illnesses and heart disease.

Sulphur dioxide, SO₂ : Sulphur dioxide is often produced by the industrial processes which produce particulates, the primary sources of SO₂ being coal, fuel oil and diesel. Being a corrosive acidic gas, sulphur dioxide damages buildings and other materials, and can cause respiratory problems.

Carbon monoxide, CO : The commonest source of carbon monoxide is motor vehicle emissions, where it results from the combustion of petrol in the presence of insufficient oxygen. It is also a result of some fuel-consuming industries and domestic fires. Carbon monoxide is a color less, odorless, highly toxic gas that displaces oxygen in human blood, causing oxygen deprivation.

The oxides of nitrogen, NO_x : NO_x refers to the mixture of nitric oxide (NO) and nitrogen dioxide (NO₂) formed by the oxidation of nitrogen during the combustion of air. The majority of NO_x is produced in motor vehicle emissions, although other sources can have significant local impact. NO_x is a contributor to several secondary pollutants, and NO₂ is a respiratory irritant that can also corrode metals at high concentrations.

Benzene : Over the last few years leaded petrol have been phased out of use. However this has resulted in higher levels of benzene and other aromatics in the substitute unleaded petrol. Benzene breaks down quickly in the environment and is not stored in the tissues of plants or animals. However, it is still hazardous to humans at high levels as it can cause several diseases of the blood including leukemia (cancer of the white blood cells). Benzene monitoring programs were started in New Zealand in 1994 and are continuing because the levels in some locations were found to be reasonably high.

Hydrogen sulphide, H₂S : Hydrogen sulphide is mainly associated with geothermal activity, where it is responsible for the 'rotten eggs' smell, but it is also formed from the anaerobic decomposition of many organic wastes and is a by-product of paper manufacture and leather tanning (see article). It is highly poisonous (more toxic than hydrogen cyanide), and because it initially anaesthetizes the sensory organs it can build up to high concentrations without warning and cause paralysis and then asphyxiation.

Fluorides: These have two main sources: the Comalco aluminium smelter and fertilizer works . Fluorides can have adverse effects on plants and in some cases concentrate in the leaves so that animals eating the plants ingest significant quantities.

II. PROPOSED AIR MONITORING SYSTEM DESIGN

The complete system design is shown in figure 2, Hardware Design Schematic Diagram. The design

Included the following major hardware components:

A). **Arduino Microcontroller [1]:** this is the core component of the design. Arduino is a flexible programmable hardware platform designed for fast Embedded Systems platform conception. Arduino's little, blue circuit board, mythically taking its name from a local pub in Italy, has in a very short time motivated a new generation of microcontroller users of all ages to make all manner of wild projects found anywhere from the hallowed grounds of our universities to the scorching desert sands of a particularly infamous yearly arts festival and just about everywhere in between.

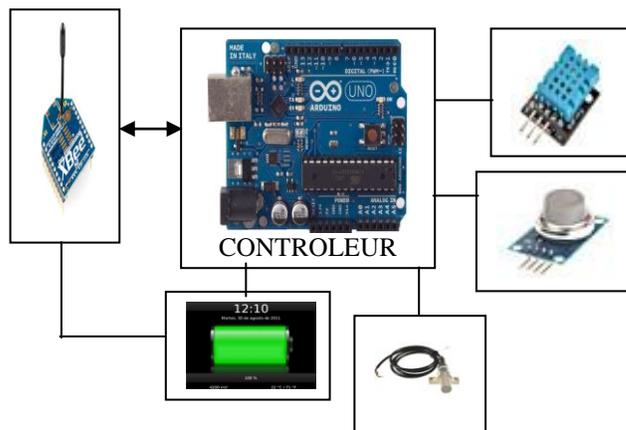


Figure 2: Sensor Node main Components

Usually these Arduino-based projects require little to no programming skills or knowledge of electronics theory, and more often than not, this handiness is simply picked up along the way.

In figure 2.b we can see the main components in the arduino-uno system board.



Figure 2.b: Arduino-uno system board

B). **MQ-2 GAS Sensor [17]:** MQ-2 Sensor is used in gas leakage detecting equipments in family and industry, are suitable for detecting of LPG, i-butane, propane, methane ,alcohol, Hydrogen, smoke.

Resistance value of MQ-2 is difference to various kinds and various concentration gases. So, When using this components, sensitivity adjustment is very necessary. we recommend that you calibrate the detector for 1000ppm liquefied petroleum

gas <LPG>, or 1000ppm iso-butane <i-C₄H₁₀> concentration in air and use value of Load resistance that (RL) about 20 K Ω (5K Ω to 47 K Ω).

This sensor module utilizes an MQ-2 as the sensitive component and has a protection resistor and an adjustable resistor on board. The MQ-2 gas sensor is sensitive to LPG, i-butane, propane, methane, alcohol, Hydrogen and smoke. It could be used in gas leakage detecting equipments in family and industry. The resistance of the sensitive component changes as the concentration of the target gas changes.

C) Temperature and Humidity sensor

The DHT11, DHT21 and DHT22 are relative cheap sensors for measuring temperature and humidity. In reference [6] and [7] there is a description of library for reading both values from these sensors. we contacted the manufacturer to get the details of the differences between the two DHT sensors to build a lib that supports both. The DHT21/22 is quite similar to the DHT11 and has a greater accuracy (one decimal) and range (negative temperatures), however the price of DHT11 is lower. The hardware pins and handshake are identical but they use different data formats.

Communication and format for DHT11: Single-bus data format is used for communication and synchronization between MCU and DHT11 sensor. One communication process is about 4ms.

Data consists of decimal and integral parts. A complete data transmission is **40bit**, and the sensor sends **higher data bit** first.

Data format: 8bit integral RH (Relative Humidity) data + 8bit decimal RH data + 8bit integral T data + 8bit decimal T data + 8bit check sum. If the data transmission is right, the check-sum should be the last 8bit of "8bit integral RH data + 8bit decimal RH data + 8bit integral T data + 8bit decimal T data".

D) Resistance Circuitry: Resistance value of MQ-2 is difference to various kinds and various concentration gases. So, When using this components, sensitivity adjustment is very necessary. we recommend that you calibrate the detector for 1000 ppm liquified petroleum gas <LPG>, or 1000 ppm iso-butane <i-C₄H₁₀> concentration in air and use value of Load resistance that (RL) about 20 K Ω (5K Ω to 47 K Ω).

E) Xbee Transmission module: is based on Zigbee protocol: The Xbee radios can all be used with the minimum four number of connections – power (3.3 V), ground, data in and data out (UART), with other recommended lines being Reset and Sleep [18]. Additionally, most Xbee families have some other flow control, I/O, A/D and indicator lines built in. A version of the XBees called the programmable Xbee has an additional onboard processor for user's code.



Figure 2.c: Transmission Module Xbee

IV. DHT11 SENSOR PROGRAMMING AND PROTOCOL

Source code for DHT11 sensor reading by Arduino uno: in there code lines we illustrate part of the software to be included into the arduino uno memory.

```
#include "dht.h"
int dht::read11(uint8_t pin)
{
  // READ VALUES
  int rv = read(pin, DHTLIB_DHT11_WAKEUP);
  if (rv != DHTLIB_OK)
  {
    humidity = DHTLIB_INVALID_VALUE; // invalid value, or is NaN
    preferred?
    temperature = DHTLIB_INVALID_VALUE; // invalid value
    return rv;
  }

  // CONVERT AND STORE
  humidity = bits[0]; // bits[1] == 0;
  temperature = bits[2]; // bits[3] == 0;

  // TEST CHECKSUM
  // bits[1] && bits[3] both 0
  uint8_t sum = bits[0] + bits[2];
  if (bits[4] != sum) return DHTLIB_ERROR_CHECKSUM;

  return DHTLIB_OK;
}
```

F) Air-quality Index(AQI):

An Air-quality Index (AQI) is used in AQMS. The AQI is an indicator of air quality, based on air pollutants that have adverse effects on human health and the environment. The pollutants are ozone, fine particulate matter, nitrogen dioxide, carbon monoxide, sulphur dioxide and total reduced sulphur compounds. figure 3 illustrate the AQI range.

The Ambient Air-Quality Standards for ANNABA reports that the safe limit for ozone is 100 micrograms per m³ and the safe AQI value set is also 100. Therefore, the AQI itself can, indirectly, be used to measure Ozone concentration in Annaba city.

Air Quality Index Levels of Health Concern	Numerical Value	Meaning
Good	0-50	Air quality is considered satisfactory, and air pollution poses little or no risk.
Moderate	51-100	Air quality is acceptable; however, for some pollutants there may be a moderate health concern for a very small number of people who are unusually sensitive to air pollution.
Unhealthy for Sensitive Groups	101-150	Members of sensitive groups may experience health effects. The general public is not likely to be affected.
Unhealthy	151-200	Everyone may begin to experience health effects; members of sensitive groups may experience more serious health effects.
Very Unhealthy	201-300	Health alert: everyone may experience more serious health effects.
Hazardous	> 300	Health warnings of emergency conditions. The entire population is more likely to be affected.

Figure 3. Description of air quality index categories.

V. RESULTS AND DISCUSSION

The proposed design were used to measure the Air-quality in several places inside the Annaba City and included different gases levels but focused mainly on measuring three main gases: Carbone Monoxide (CO) and Liquid Petroleum Gas (LPG) and NO2. A sample of obtained results from both clean environment close to Seriadi mountains, Annaba city center where there is a crowded circulation and El-Hadjar region a Metal-Steel production firm in Annaba, the results are shown in table 1, 2 and 3 respectively concerning measured area.

Table 1: Situation of air pollution in SERAIDI mountain area in Annaba city

Seraidi Co	Seraidi NO2	Seriadi LPG
0.05	1	2.05
0.8	2.5	3.5
0.75	0.8	2.7
0.48	0.8	1.9
0.87	2.4	2.9
0.79	1.7	3.04
0.61	1.5	2.9

Table 2: Situation of air pollution in Center City area in Annaba city

Center Co	Center NO2	Center LPG
20	16	75
26	13	86.9
24	17	87.4
26.78	15.68	80.6
27.58	19	76
29.15	20	79
30.15	25	78.95

Table 3: Situation of air pollution in Metal Steel factory area in Annaba city

Metal-Steel Co	Metal-Steel NO2	Metal-Steel LPG
35	56	25
34	57	24
36.7	58	38
40.58	55.8	26.25
32.78	50.15	27.8
31.99	52	30
32.58	53	30.5

Simulation results: for simulation of WSN nodes, the area is divided into parts where each part can be controlled by a node, in this case the area is divided into 9 regions, and the transmission circuit is chosen so that it can provide the adjacent nodes with the information with minimum consumption of energy.

Scenario 1: by adjusting the sliders for CO, SO2 and NO2 Gas we obtained the Red color of the region, which illustrates by node 9 the values sensed: Co=206ppm SO2=160 ppm and NO2=200 ppm

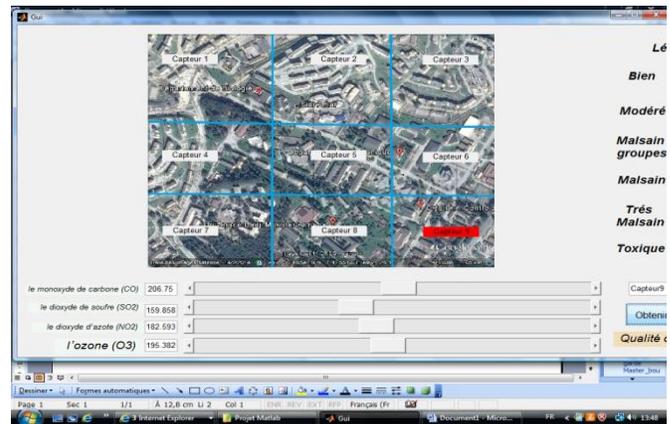


Figure 4: node 9 sensed Co=206ppm SO2=160 ppm and NO2=200 ppm levels the central control unit

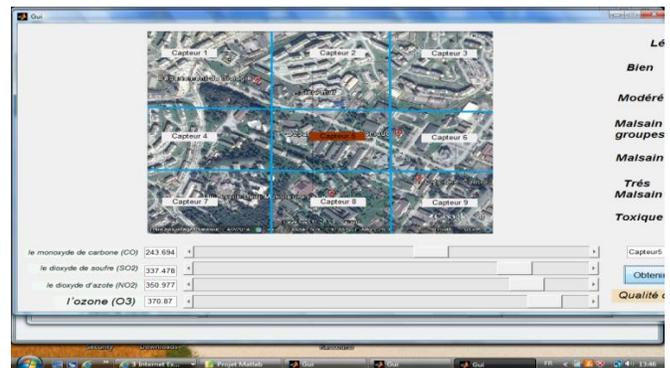


Figure 5: Node 5 sensed Co=243 ppm SO2=337 ppm and NO2=360 ppm levels the central control unit

Based on the normal gas levels of the clean air fig 3., the results indicate that there is a big difference in the gas levels of both gases (LPG and CO) which obtained from the several tests.

The results in figures 4 and 5 show respectively the quality of air in the simulated regions. In figure 4, based on gas concentrations, the AQI indicates that the air is unhealthy for both insensitive and sensitive group. In figure 5, also based on simulation results, the AQI indicates that the quality is hazardous since the quantities of gases exceed 330 ppm.

VI. CONCLUSIONS AND PERSPECTIVES

Air-quality monitoring System Design to assess the pollution of air in some parts of Annaba city using a micro-system, as a node in Wireless Sensor Network (WSN), is proposed in this article. WSN enhanced the process of monitoring many environmental phenomena such as the air pollution monitoring issue in proposed this paper. It provides a real-time information about the level of air pollution in different regions, as well as provides alerts in cases of drastic change in quality of air. Based on collected information, such data can then be used by the authorities to take prompt actions such as evacuating people or sending emergency response team. The proposed design is enhanced by several ways such as: selecting adequacies' sensors, calibrating these sensors for gas detection, integrating them in a WSN system controlled by an Arduino-Uno, and finally transmission to the central unit using Xbee modules. A Graphic user interface (GUI) has been presented in this work to simulate the effect of sensors on selected area . The results are interesting, improvements can be done: in providing a web service page that can provide these data to users, as well as more sophisticated sensors could be used such as MQ-135, MQ-136 and others. We think to improve this work by using mobile sensing system where public transportation infrastructure can be used[19]

7. ACKNOWLEDGEMENTS

Authors appreciate the support of LASA laboratory at Badji Mokhtar Annaba University (BMAU) ,faculty members of Environment department (BMAU) and Dean of Faculty of IT at Al-Zaytoonah University AMMAN, JORDAN.

REFERENCES

- [1] Diego Mendez, Alfredo J. Perez, Miguel A. Labrador, Juan Jose Marron, Mar 2011, P-Sense: A Participatory Sensing System for Air Pollution Monitoring and Control, *IEEE International Conference on PERCOM Workshops*, pp. 344-347.
- [2]. S. Choi, N. Kim, H. Cha, and R. Ha, 2009, Micro Sensor Node for Air Pollutant Monitoring: Hardware and Software Issues ". *Sensors* 2009.
- [3] WHO Library Cataloguing in Publication Data , 1999, Monitoring ambient air quality for health impact assessment , (*WHO regional publications. European series ; No. 85*), **World Health Organization 1999, ISSN 0378-2255**
- [4]. Qasem Abu Al-Haija, 2011, Toward Secure Non-Deterministic Distributed Wireless Sensor Network Using Probabilistic Key Management Approaches, *Journal of Information Assurance and Security* 6 (2011) 010-018.
- [5] I.F. Akyildiz, W. Su, Y. Sankarasubramaniam, E. Cayirci, 2002, Wireless sensor networks: a survey, *Computer Networks* 38 (2002) 393-422, elsevier.
- [6] <http://playground.arduino.cc/Main/DHT11Lib>
- [7] www.micro4you.com/files/sensor/DHT11.pdf
- [8] Raja Vara Prasad Y. et al, June 2011, Real Time Wireless Air Pollution Monitoring System, *ICTACT Journal on Communication Technology*, June 2011.
- [9]. M. Riley, 2012, Programming Your Home: Automate with Arduino, Android, and Your Computer, *The Pragmatic Programmers*, 2012.
- [10] R.A.Roseline, Dr.P.Sumathi, 2012, Local Clustering and Threshold Sensitive routing algorithm for Wireless Sensor Networks, in *the IEEE sponsored International Conference on Devices Circuits and Systems(ICDCS'12)*, March 2012. (Available online at ieeexplore.com).
- [11] Honicky, R.; Brewer, E.A.; Paulos, E.; White, R., 2008, N-smarts: networked suite of mobile atmospheric real-time sensors . In *Proceedings of the 2nd ACM SIGCOMM Workshop on Networked Systems for Developing Regions*, Seattle, WA, USA; ACM: Seattle, WA, USA, 2008.
- [12] Volgyesi, P.; Nadas, A.; Koutsoukos, X.; Ledeczi, A. , 2008, Air Quality Monitoring with SensorMap . In *Proceedings of the 7th International Conference on Information Processing in Sensor Networks*, St. Louis, MO, USA; IEEE Computer Society: St. Louis, MO, USA, 2008.
- [13] Wenhu Wang, Yifeng Yuan, Zhihao Ling, 2011, The Research and Implement of Air Quality Monitoring System Based on ZigBee , 2011 *7th International Conference on Wireless Communications, Networking and Mobile Computing*, pp. 1-4, Sept. 2011.
- [14] Sharma, A.; Golubchik, L.; Govindan, R., 2007, On the Prevalence of Sensor Faults in Real-World Deployments. In *4th Annual IEEE Communications Society Conference on Sensor, Mesh and Ad Hoc Communications and Networks*, San Diego, CA, USA; IEEE Computer Society: San Diego, CA, USA, 2007.
- [15] Saitas, Jeff; 1997, *Ground-Level Ozone*; [Online] Available <http://www.tnrc.state.tx.us/air/monops/ozoneinfo.html>; February 20, 1997.
- [16] Xintaras, Charlie and Perry, Mike; 1997, *Agency for Toxic Substances and Disease Registry*; [Online] Available <http://atsdr.cdc.gov/8080/toxfaq.html>; February 20, 1997.
- [17] Technical Data For MQ-2 Gas Sensor, Website <http://www.seeedstudio.com/depot/datasheet/MQ-2.pdf>
- [18]http://cs.smith.edu/dftwiki/index.php/Tutorial:_Arduino_and_XBee_Communication
- [19] Aruljothi.R , "Air Pollution Measuring System with Mobile Sensor Arrays" *International Journal of Scientific & Engineering Research*, Volume 4, Issue 5, May-2013 ISSN 2229-5518

Ali-Aldahoud: Professor and Dean of Faculty of IT at Zaytoonah University Amman Jordan, IEEE senior member, he has many publication in different journals and is an active TPC on many conferences, he is the chair of ICIT conferences his main research interest field balancing in information technology, and WSN.

Mohamed Fezari: Associate professor at Badji Mokhtar Annaba university faculty of engineering, member of laboratory of Automatic and signals annaba Algeria, he has produced many articles in different journal of computer engineering, he participated at many conferences in technical program committee, his main interesting research domain are speech processing human machine communication and WSN.

Ismail Jannoud: at Zaytoonah University Amman Jordan, Faculty of Engineering and Technology, Has many publications in different journals and conference Proceedings, main research interest: Communication , engineering , Digital Image processing, Multimedia processing, Pattern Recognition and computer vision.

Thamer AL-Rawashdeh: at Zaytoonah University Amman Jordan, Faculty of Information Technology, He has many publications in different journals and conference Proceedings, main research interest: software quality evaluation, embedded systems, and programming languages.

Extending the Matrix Vector Transition Net Approach for Modeling Interaction

A. Spiteri Staines

Abstract—This paper briefly explains the main problems with Petri nets for representing complexity. It introduces the Matrix Vector Transition Net, which is based on Petri net like semantics and explains how this executable structure is more expressive for certain classes of system modeling problems where inputs and outputs can be grouped into matrices or vectors. Some trivial examples are given to explain the salient points of the MVTN and how these can be combined with Petri nets in order to increase their modeling power. Some brief results are findings presented in this paper.

Keywords— Colored Petri Nets, Matrix-Vector Transition Net, Petri Nets, System Modeling and Interaction

I. INTRODUCTION

COMPUTER systems are increasingly popular in various fields in the modern world. The use of such systems varies from complex control in fields like avionics to business organizations and industry. Cloud computing and computer networking are based on grid technologies, complex interacting elements and many different types of time dependent configurations. Formal specification methods and notations have been around for a number of decades. These can be used to prove the reliability and correctness of such structures. Petri net formalisms are types of models that have found extensive use for constructing visual and representational models of various types of systems. Naturally, formal modelling structures are important for various reasons ranging from verification and validation, checking and representation purposes.

Modern systems are not just characterized by computations but also by the correctness of the specification and the interconnectedness of different modules, components or operating parts. These parts transfer information with others. Due to the complex relationships of systems, different structures of networking and connectedness exist. This is evident in modern open systems like computer networking, mobile networking, social networking, grid computing, transport systems, logistics, etc. Sometimes the underlying structures are not obvious, but global events or activities have a much more far reaching effect than can normally be seen. E.g. if a network is considered from a global perspective, each

element in the network can have a state, but all the elements connected together form a composite global state. If one element changes state, the global state is bound to change.

Unfortunately, many formal methods including ordinary Petri nets deal with the axiomatic or lower level parts of global systems. Normal place transition Petri nets and some other classes cannot really model the complex intricacies of communication in complex distributed systems. Still they are useful for symbolic and structural representation in restricted form. When ordinary Petri nets were created they were never intended to model at this level. One disadvantage of Petri nets is that for complex structures it is possible to derive models that have over 50 places and transitions making them very difficult to read and comprehend. Vertices, density and localization of arc connectivities all pose problems to the construction and interpretations of the net.

Previously in [13],[14] another modelling notation based on the MVTN (matrix vector transition net) approach has been suggested. This is based on ordinary Petri net like semantics but ordinary places and transitions are substituted by vectors or matrices are used and the input and output arcs are used with respective functions. This structure is useful for certain problems and systems that have multiple inputs and outputs that can be grouped or even inputs that can have real values. This work builds upon the previous work, showing how the modelling approaches used are suitable to extend for other complex abstract representation purposes.

II. RELATED WORK AND BACKGROUND

Petri nets are expressive formalisms with visual counterparts. These are extensively documented and well supported by over three decades of research. Place transition nets have been used to explain how supervisory control can take place in distributed systems [1]. Communication between elements is like a connection layer based circuit. Petri nets have been used to create Systolic networks in [2]. These are a combination of Petri nets and other notations. They are useful for modeling interconnected processors like grid, circuit or mesh topologies and can be decomposed.

In orthogonal transformations from CPN's (colored Petri nets), matrices with sets are used to represent the networks [3]. Colored Petri nets have extensive use for modeling complex systems like train systems, transportation and logistics. CPN's and higher order nets allow the colored place to contain complex types, possibly even sets and matrices [4]-[10].

Anthony (Tony) Spiteri Staines, is with the Department of Information Systems, Faculty of ICT, University of Malta, (corresponding phone: 00356-21373402,e-mail: toni_staines@yahoo.com)

The actor model is a higher order net structure useful for modeling information systems and distributed systems based on workflow concepts [11],[12]. The concept of a processor element, instead of a simple transition can be used to develop circuits and certain topologies. Petri nets are useful in workflow systems and their ability to properly model the processes at varying levels of detail and formalization [2]-[4], [6]-[8].

In previous work [13],[14] it has been shown how the MVTN notation can successfully model complex communication between different processing elements using a Petri net oriented type of behavior [13]. An example of an abstract switch was used as a case study [14]. The reachability marking graph can be constructed for the firing of each processor/element. Instead of ordinary places the structure presented uses matrices or vectors for inputs and outputs [13]. The inputs and output arcs were assigned respective functions that must match the dimensions of the respective input and output matrices. The operations of the net can be represented using basic matrix algebra [15], [16]. Simultaneously, the structural complexity of the network is kept reduced or simplified. It is possible to apply traditional solutions or modeling methods to many system structures provided that the traditional solutions are used in new ways and preserve certain fundamental properties. Having different views and representation models of a system can offer better insight to what is happening. Thus, multiple views definitely offer better reasoning about complex systems. It is suggested that the MVTN is used in a combined approach for this purpose.

III. MOTIVATION

In today's world there is an ever increasing complexity in the distribution and connectivity of hardware, middleware and software. New approaches to represent these complexities are important if software engineering is to be taken to the next level. This means that new methods of system representation must be developed. One interesting area is that of symbolic modeling structures, that exhibit Petri net like behavior.

The main objective of this work is to provide simple solutions to describe a wide range of applications and systems. As it is impossible to meet optimally all the objectives of formal modeling using diagrammatic notations, the focus is on the reliability and correctness of a system structure. System verification in reduced form can be based on i) proof of correctness, ii) proof or termination of some processing activity and iii) conditional correctness. These could be obtained from the MVTN in various ways [13]. A demand for large scale nets to represent real world systems exists. The large scale nets must exhibit Petri net like behavior. i.e. well foundedness, well formed, some form of execution, conservative behavior, formally verifiable, no loss of information, etc. The matrix vector transition net (MVTN) or the simpler called matrix transition net (MTN) has been precisely defined for these purposes.

The complexity factors between interfacing and linking

components imply a difficulty to represent and model this interaction. Petri nets can be used for this. However ordinary Petri nets were never intended to model at these levels. Colored Petri nets and higher order nets are far better [5]-[8]. If systems evolve it is also necessary that the models for system description evolve. Evolution in systems implies that modeling power must increase. This problem happens because still no proper visual notations are readily available. It is possible to identify a gap of missing information between real world complex models and those that are represented using ordinary Petri nets. This gap can be solved using the MVTN. The MVTN is a higher order net that can be considered to be a type of CPN, but that offers the use of matrices and vectors [13].

Unfortunately, many Petri net models are an oversimplification of the real world system and scenarios. When system representation is done using graph structures and network models, the models do not fully describe all the aspects of the system. Petri nets are appropriate means for describing the partial ordering or sequencing of actions and events. However from the viewpoint of causal ordering Petri nets are representative of the state or states of a system. The graphical and mathematical properties of these nets guarantee that partial event ordering is preserved because of precise rules and restricted actions. The MVTN strives to preserve these properties.

From another aspect large scale Petri nets can be created. However, they are not compact and easily readable [17]-[19]. Because of the limiting factors in ordinary Petri nets, alternative modeling approaches like colored Petri nets, higher order nets, graph structures, structured diagrams, etc. have been suggested. Some classes of higher order nets are just representative and not really executable structures. A Petri net without a processing element is just a graphical description according to modern system architectural views. A Petri net has to be processed or executed for some change in state. Low level behavior cannot be easily replicated by high level structures and vice-versa. Many system structures are based on repeatable patterns and replication of certain parts.

Various reasons can be given for motivating the use of the MVTN and extending it. Matrix structures can be derived from ordinary Petri nets to show their execution. Matrices and vectors have been chosen because they can contain more information than ordinary places. Matrices can represent extensive data sets or results and the global state of a complex system. From an architectural perspective, circuit behavior is easily abstracted. In essence, if a system can be viewed as a set of interconnected, related components passing messages or information, then the inputs and outputs can be grouped into matrices or vectors.

The usage of matrices and vectors, instead of ordinary places allows the possibility of creating new diverse types of models in addition to the standard types. The MVTN is not exclusive and additionally normal Petri net places can be added, increasing the expressiveness and complexity. Matrices

in mathematical viewpoint are considered extraordinary forms of algebra that can express multidimensional views and complex data representations in a compressed form. Combining the MVTN with Petri net structures offers even more modeling possibilities for exploration.

IV. PROBLEM STATEMENT

The main problem, is to represent complex system structures and the connections using more detailed models that can effectively show the detailed communication and its sequence. The usefulness of Petri nets comes from its execution and the possibility to have a change in state. The challenge of the MVTN is that it has to imitate Petri net like behavior and be capable of representing the change in state when an event takes place [13].

Petri nets and the MVTN are based on digraphs having simple edge and node types. The MVTN can include matrices, row and column vectors and other constructs like even elementary Petri net places. Obviously, the MVTN is more complex than an ordinary place transition net and the execution introduces new problems.

The model is similar in principle to higher order net structures and can represent the detail and complexity involved in communication architectures. These structures offer improved and compacted visual representation which is more representative of what is happening in the real world. This is because the MVTN contains the possibility for expansion and expressiveness as regards to complexity.

V. PROPOSED SOLUTION

The proposed solution here is the use of the MVTN in conjunction with other basic Petri net constructs like places, transitions, etc. The complete workings of the MVTN are not shown because it would take up too much space however they are included in [13].

The relationship between Petri net places, vectors and matrices is shown in fig.1. This section outlines the practical implementation of the MVTN. A system can be described as a finite set of interconnected elements (e1, e2, e3,...,en). An element e1 can have an input e1? and output e1! . The system can be defined as [SYSTEM] the set of processing elements or components of the system. The global communication process is defined as a sequence comm. process: Seq SYSTEM, and comm. process =< e1, e2,..> in any given order. # seq SYSTEM > 0. The global process of communication implies that sub processes can be identified. In each process data is transferred from one component to another. The components or elements can be considered to be individual processors. Given that the net is composed of input matrices and output matrices or vectors: The system state change can be given as

(MS,time) (MS',time') where MS implies the matrix set of the net and denotes a single step transition. Formally the MVTN can be expressed as a system that links up a set of

processes or elements. The simple execution of the system can be expressed as a set of (inputs, transitions, outputs) i.e. (Inp1 T1 Out1). The change in state as the result of a transition is reflected by the change in the input and output values.

Some basic properties about the MVTN like i) global transitions, ii) use of matrices and vectors, iii) complete vs partial transition firing, iv) separation of outputs and inputs, v) transition firing by vectors, matrices and places can be taken to be similar to what happens in normal Petri nets. Other issues like i) concurrency , ii) parallelism, iii) choice or conflict, iv) transition enabling and disabling, v) results of firing a

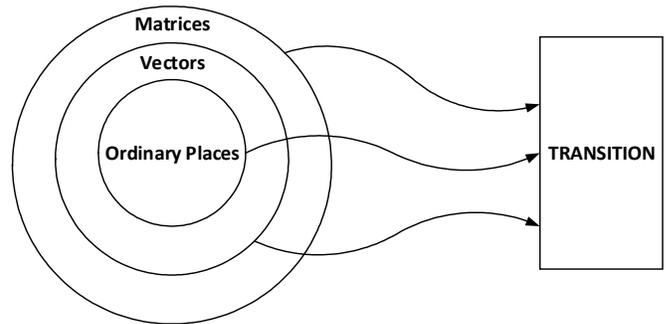


Fig. 1 Relationships between Places, Vectors and Matrices for the MVTN

transition, vi) symbolic reachability or marking graph construction and more can be considered to be similar to what happens in normal Petri nets. These are explained in more detail in [13]. This is because only the input and output types of the net have been changed. The underlying functionality has not been modified at all.

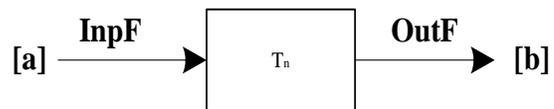


Fig. 2 Generic Form of the MVTN

For proper firing all the input conditions must be satisfied. I.e. the matrices, vectors or places must have sufficient values in them to satisfy the input function. The input and output functions are analogous to Petri net arc weights or arc expressions in colored Petri nets. The input and output functions do not need to be of the same order at all. Different combinations are possible depending on the modeling requirements. When a processing element does not have an even amount of inputs to form a proper matrix, zero values can be added to solve this. When a transition fires the inputs and output values normally change. Thus, if there are some output values and an output function then the output values must change accordingly, when a transition takes place.

Theoretically it is possible to have no output values but just input values as is the case with Petri nets.

For an oversimplified generic or general form of the MVTN represented in fig.2, transition firing can be basically given as $Inp = Inp - Inp_function$ and $Out = Out + Out_function$.

VI. CASE STUDY TOY EXAMPLES

In this section some simple practical examples of the use of the matrix vector transition net combined with Petri net places are used to illustrate the use of this modelling structure for some generic system behavior and structure.

A. Simple MTVN combined with Petri Net Places

To start off a normal Petri net the diagram in fig. 3 is shown and compared with the MVTN counterpart in fig.4. Please note that there is no relationship between fig. 3 and fig. 4. This is just an example. It is obvious that the MVTN structures are more detailed and expressive. The diagram in fig. 4. illustrates how the MVTN can be combined with ordinary Petri net places. The basic idea behind this structure opens many new dimensions and possibilities of modeling. The net in fig. 3 is an executable structure where a token can be removed from P1 and the input matrix has sufficient values for removal also. I.e. the values in the input matrix are equal or greater than the values in the input function. Hence there are sufficient values for firing T1 i.e. transition T1 is enabled.

Fig. 5 shows the new state of the net in fig. 4 after firing transition T1 which is enabled. The net in fig. 5 is non-live in Petri net terminology. I.e. no further activity is possible and markings cannot change.

separate Petri net and MVTN and connected together. The MVTN and the Petri net can be executing in parallel. Hence structures running in parallel are being connected. This example can obviously be extended to other classes of Petri

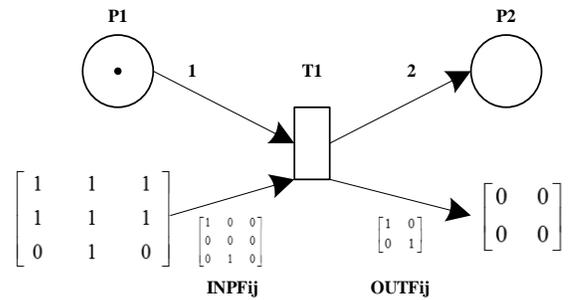


Fig. 4 MVTN with Places Added

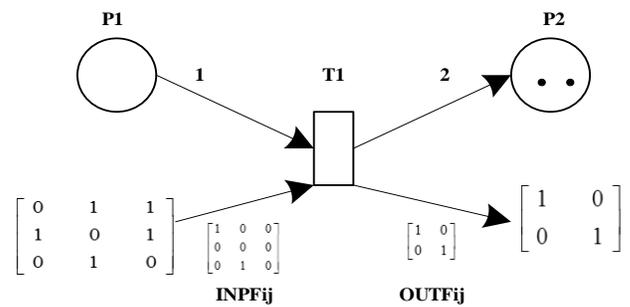


Fig.5 MVTN after Firing T1

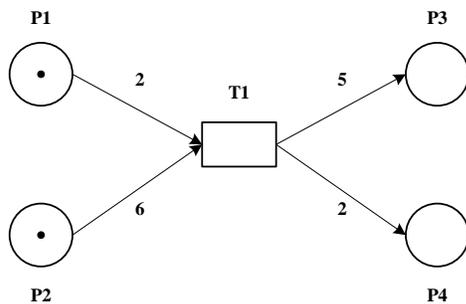


Fig. 3 Normal Place Transition Net

B. Some Simple Examples

The diagram in fig. 6 shows a more complex structure. Even though the inputs to T1 are a matrix and a place P1, the output is a simple place P2 that connects to T2. The net in fig. 2. is executable. However, after T2 fires then it will become a non-live net. The place in P2 is a symbolic representation of a buffer or store or even a possible switch that will activate T2 when the right conditions are present.

The diagram in fig.7 shows the possibility of having a

nets for more detailed modeling as required.

The structures shown are executable. After T1 fires, two tokens are placed in P2, T2 and T3 are activated simultaneously. It is possible that T2 fires twice or T3 fires. If T3 fires no tokens are lost, but whenever T2 fires a token is lost. When transition T3 fires, transition T4 becomes enabled. Other possibilities are present here. The type of behavior in fig. 7 presents the concepts of choice or conflict in Petri nets. This type of structure presents non-determinism, i.e. different types of behavior and outcomes are possible. This structure could be modified or restricted to guarantee more deterministic behavior.

C. Message Communication Decomposition Example

This example is about message communication. Here there is assembly and disassembly of information or composition or decomposition. This example can be compared to packet assembly /disassembly, multiplexing and other forms of networking. The initial network diagram in fig. 8 shows active ports which communicate to a particular source or entity. The same analogy can be applied for a message that has to be split up into different sources and rerouted or packet assembly and

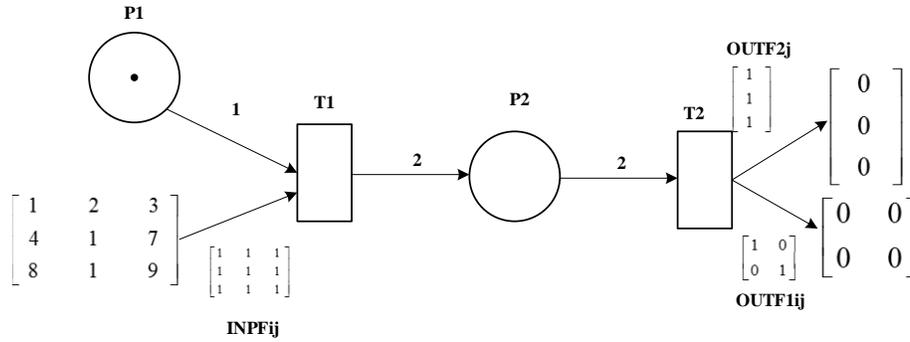


Fig. 6 MVTN structure using a buffer like place connection or communication

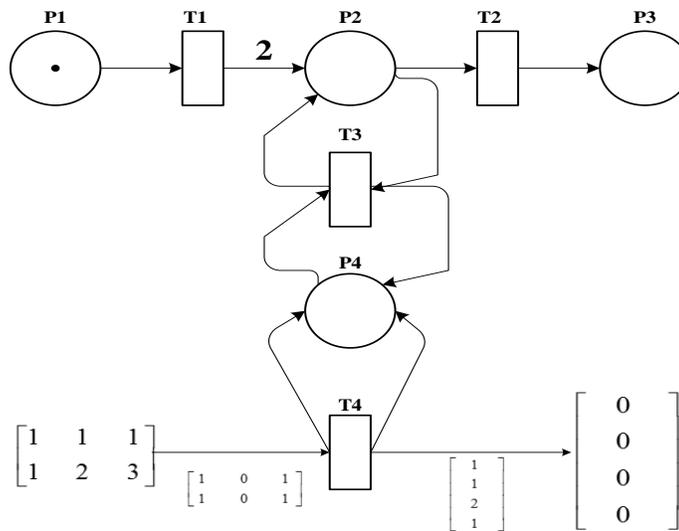


Fig. 7 Separate MVTN and Petri net connecting together

disassembly. The MVTN combined with petri net places for input and output is an executable structure that represents this in fig. 9. Letters are used symbolically to represent the input and output sources respectively. The letters are an abstraction or symbolic representation of possibly real values. At the minimum the MVTN allows for proper representation. The concepts presented in the MTV net can easily be extended to model PTP channels and message splitting algorithms.

D. Message Passing between Three Workstations

A simplified toy example of message passing between three stations is shown in fig. 10. This behavior can be shown using a basic message sequencing chart (MSC). The MSC in fig. 10 describes the sequencing of events. Fig. 11 shows the MVTN for the whole sequence in fig. 10. The MVTN in fig. 11 is an executable model.

VII. RESULTS AND FINDINGS

The example in section 6 A shows how Petri net places have been combined with the MVTN. The structure is an executable one. However after firing the main transition then it becomes a dead structure. I.e. no more states are possible in Petri net terms. Hence the composite marking of this system is rather simple. The example in section 6 B shown in fig.6. is slightly more complex. Here there are more states and the behavior of this net is typical of a buffer, channel or port, but communication is asynchronous, i.e. in one direction only. The MVTN allows complexity to be modeled in this structure. I.e. the actual transfer of information can be depicted. The MVTN can be modified if needed to show the actual data being transferred. This is possible with colored Petri nets but not with ordinary Petri nets.

Fig.7 shows some new challenging behavior where conflict or choice have been included. The outcomes in this net are

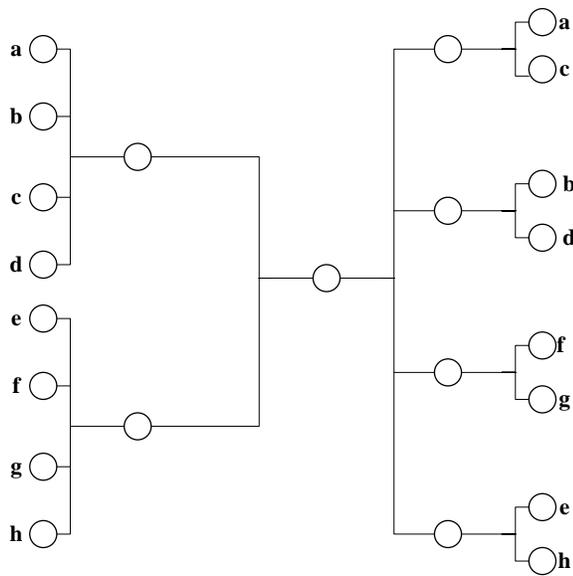


Fig. 8 Symbolic representation of network assembly/disassembly undetermined and many different possibilities exist. Obviously

can be used. This will reduce the complexity of the description. These models are useful for representing different types of architectures and component based systems. The more complex the models the more time consuming is their construction. The models are decomposable into Petri nets which will result in information loss. This approach in fig. 8 can be used for modelling other things like parallel algorithms or processing where specific decomposition is needed. This structure is similar to that of a control flow graph.

In section 6 D fig. 10 a typical message sequencing activity between three work-stations is shown. The message sequence chart in fig. 10 represents this. The model in fig. 10 is not an executable one. The MVTN for fig. 10 is shown in fig. 11. This is an executable model that can be used for modeling and tracing errors. This model is quite complex and detailed. The messages are shown using vectors this time, but obviously the vectors could be replaced by matrices. The input function is omitted and represented symbolically using ?a, similarly for output !a is used. This can also imply that whatever value is inputted is accepted and outputted. This model is also useful to see at which current state the system is in. I.e. station A waits for a reply from station B to continue broadcasting.

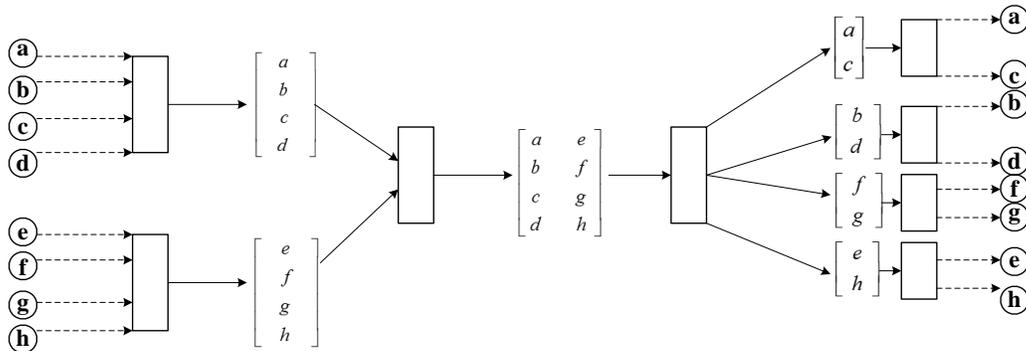


Fig. 9 MVTN for the network assembly/ disassembly structure

the marking graph for this structure will prove to be more difficult to construct. I.e. there is undecidedness in the behavior or the net in fig.7. The result from this is that the MVTN does preserve Petri net like behavior and allows the modeling of different types of conditions associated with normal Petri nets. Both the MVTN and the Petri net can be used to switch or trigger each other or vice-versa. Fig. 7 indicates that communication between the MVTN and ordinary Petri nets is possible at certain synch points as required.

Section 6 C shows an abstract system that represents real world communication is shown. Fig. 8 just shows that this is possible. Fig. 9 is the executable counterpart of fig. 8. For simplicity's sake in fig. 9 and fig. 8 letters which are representative of the actual values have been used. The MVTN structures are fully executable. They have different states and a reachability graph or marking graph can be constructed. If this is too complex a symbolic marking graph

These states can be clearly depicted in this network. It is also possible to identify deadlock and concurrency issues from the net. Many of the major concepts that are applicable in Petri nets can be used for this model. The net in fig. 11 is live but after firing all the transitions it goes into a non-reversible state that depicts exactly what is happening in the message sequencing chart. This modeling approach could be combined with other approaches like those explained in [19].

The MVTN is applicable to real world problems and abstract representation of system structures. This is possible for systems where inputs and outputs can be grouped into matrices. For these systems the MVTN can prove to be quite useful for modeling. The MVTN approach can definitely be combined with other notations from Petri nets and even other Petri net classes if this is required. This can open up a lot of new exploratory modeling. This has been clearly demonstrated in this paper.

If grouping is not possible then this approach might not be

suitable for use. For simple problems it will be better to use Petri nets.

The models can be used for other problems like processor modeling, pipeline architectures [20], hardware modeling [21], concurrent systems etc. and even in other fields like: travel, transportation and logistics modeling, etc. A conclusion section is not required.

VIII. CONCLUSIONS

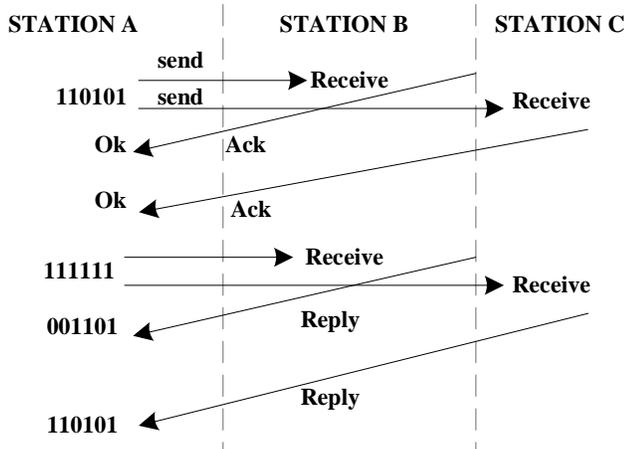


Fig. 10 Typical Message Sequence Chart for Three Workstations

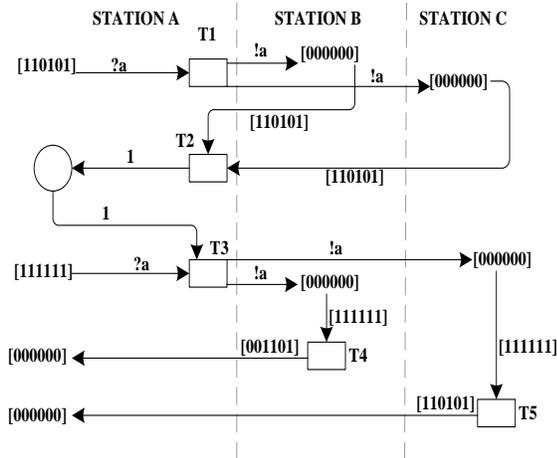


Fig. 11 Typical Message Sequence Chart for Three Workstations

This work has briefly shown the usefulness of the MVTN approach to model interaction and communication in different systems. Obviously the MVTN is based on Petri net like semantics and can be successfully combined with other Petri net components. The best use of this modeling approach seems to be where there are systems that have input and output components or elements that can be easily grouped to form matrices or vectors. The limitations of this type of modeling is that it requires a lot of time to construct more complex and detailed models and the state explosion problem inherent in

traditional Petri nets can occur here also. Thus, when using this approach, it is important to keep the representation simplified. Also, if there is no grouping the approach is not suitable.

New uses for the MVTN can definitely be considered and a lot of exploratory modeling can be done. The models can be simplified or converted into Petri nets.

REFERENCES

- [1] K.C. Wong, J.G. Thistle, R.P. Malhame, H.-H. Hoang, Supervisory control of distributed systems: conflict resolution Decision and Control, 1998. Proceedings of the 37th IEEE Conference on Decision and Control, Volume: 3, 1998, pp. 3275 – 3280.
- [2] A. Abellard, P. Abellard Systolic Petri Nets, Petri Nets Applications, InTech, 2010, ch. 5.
- [3] E. Best, T. Thielke, Orthogonal Transformations for Coloured Petri Nets, ICATPN '97, Springer, 1997, pp. 447 – 466.
- [4] K. Jensen, G. Rozenberg, High-level Petri Nets: Theory and Applications, 1st ed., Springer-Verlag, 1991, ISBN-10: 354054125X, ISBN-13: 978-3540541257.
- [5] A. Spiteri Staines, A Colored Petri Net for the France-Paris Metro, NAUN, International Journal of Computers, Issue 2, Vol. 6., 2012, pp. 111-118.
- [6] CPNTools, CPN Group, Department of Computer Science, University of Aarhus, Denmark <http://cs.au.dk/CPnets/>
- [7] L.M. Kristensen, S. Christensen, K. Jensen, The Practitioner's Guide to Coloured Petri Nets, International Journal On Software Tools for Tech. Transfer (STTT), Vol. 2, Springer-Verlag, 1998, pp. 98-132.
- [8] L.M. Kristensen, J.B. Jorgensen, K. Jensen, "Application of Coloured Petri Nets in System Development", Lecture Notes in Computer Science, Vol. 3098, Springer-Verlag, 2004, pp. 626-685.
- [9] B. Scholz-Reiter, C. Zabel, Integration of Load Carriers in a Decentralized Routing Concept for Transport Logistics Networks, WSEAS Proceedings of the 2nd International Conference on Theoretical and Applied Mechanics (TAM '11) Corfu, Greece, 2011, pp. 259-264.
- [10] J. Cortadella, M. Kishinevsky, A. Kondratyev, L. Lavagno, A. Yakovlev, Hardware and Petri Net Application to Asynchronous Circuit Design, Application and Theory of Petri Nets, Lecture Notes in Computer Science, Springer, Vol. 1825, 2000, pp. 1-15.
- [11] K. van Hee, Information Systems: A Formal Approach, Cambridge Univ. Press, 2009.
- [12] W. Van Der Aalst, K. Max Van Hee, Workflow Management: Models, Methods, and Systems, MIT press, 2014.
- [13] T. Spiteri Staines, Implementing a Matrix Vector Transition Net, British Journal of Mathematics & Computer Science, ISSN: 2231-0851, Vol.: 4, Issue.: 14, 2014, pp. 1921-1940.
- [14] T. Spiteri Staines, F. Neri, A Matrix Transition Oriented Net for Modeling Distributed Complex Computer and Communication Systems, WSEAS Transactions on Systems, Vol 13, 2014, pp. 12-22.
- [15] K.M. Abadir, J.R. Magnus, Matrix Algebra, Cambridge University Press, 2005.
- [16] F. Ayres (jr), Theory and Problems of Matrices, Schaum's Outline Series, Schaum, 1974.
- [17] M.B. Dwyer, L.A. Clarke, A Compact Petri Net Representation and its Implications for Analysis, IEEE Transactions on Software Engineering, vol. 22, issue 11, 1996, pp. 794 – 811.
- [18] R.H. Sloan, U. Buy, Reduction Rules for Time Petri Nets, Acta Informatica, Vol. 33, Issue 5, Springer-Verlag, 1996, pp. 687-706.
- [19] C. Knieke, B. Schindler, U. Goltz, A. Rausch, Defining Domain Specific Operational Semantics for Activity Diagrams, Technical Report Series IfI-12-04, Institut für Informatik, Technische Universität Clausthal, 2012.
- [20] C.V. Ramamoorthy, H.F. Li, Pipeline Architecture, Journal ACM Computing Surveys (CSUR), Volume 9 Issue 1, 1977, pp. 61-102.
- [21] A. Yakovlev, L. Gomes, L. Lavagno, Hardware Design and Petri Nets, Springer, 2000.

Location Search by Using Phonetic Algorithm with Location-Based Service

Kittiya Poonsilp, Attakorn Poonsilp

Abstract—This research was developed from the attempt to support internet searching of location named in Thai language. According to the linguistic complication and ambiguousness of Thai language, for example, the spelling of a word that has tone mark, garun (silenced tone mark—translator) and homonyms. That word can be written into many formats. This complication has an impact on word spelling. It confuses the system, so internet users can't find what they look for. The research team adapted phonetic algorithm to use with location-base service. The proposed procedure would primarily change location names to be in forming of symbol or pronunciation, then convert and compare by using similarity value from Levenshtein Distance calculation method. We prepared 100 words to be used in this research, and found out that this system worked very well with 96% accuracy rate. There were only 4% of all sample words that couldn't be searched. They basically were grammatically complicated. Besides, we tried to develop a location searching program based on the idea mentioned earlier. More useful functions were added. Users can look up maps by inputting latitude and longitude. More specific searching criteria can also be added, such as category, distance and keyword. The search result will be pinned on electronic map and will be seen clearly by users.

Keywords—Location-based Service, Phonetic Algorithm, Soundex.

I. INTRODUCTION

Nowadays, Location-based Service (LBS) [1] becomes more and more popular. There are various LBS devices available for users to use conveniently. For general search on search engine, users normally do it by comparing alphabets, which regularly causes difficulties because of misspelling, especially with Thai words. There are many conditions in Thai language, such as homonyms, tone mark, garun and similar-sound alphabet (ส น ศ, etc.) There are many words which sound the same but are written differently, and it causes misunderstanding and will be spelled incorrectly. For example, a hotel called “บ้านฐานิชา” (Baan Thanicha), can be spelled incorrectly due to the complex spelling. As a result, users can't find what they are looking for.

From the problem above, research team would like to present the new way to perform the location search by introducing phonetic algorithm to solve this problem. The algorithm will define index of each word categorized by pronunciation. Even

though the word is spelled incorrectly by users, they can still find what they look for by comparing keywords input by users to the words in database. Then the result will be shown on the electronic map. We also designed the category search function (universities, hospitals, schools, hotels, etc.) so users can search for nearby facilities. As there could be some places that have similar names, latitude and longitude are used to narrow down the areas of search and make the search more flexible.

II. RELATED RESEARCH

So far, we haven't found any research that share similar theme with ours. We have seen only research about Phonetic Algorithm, Soundex, Location-based Service that was adapted to use with something else. For example, a research [2] presented Location Based Service with simple LBS system which only looks for places (restaurants, gasoline station) and lead users to that location.

Research [3] presented the place of interest (POI) introduction system by taking current position of users to calculate. The research picked up a sample group to study the distance of commuting from current location to desired location. It also studied behaviors of users in the same area. The result was brought to make a POI program.

Most LBS systems basically do this kind of task. For our search system, we increase the ability to LBS so it works better and reach its top notch by enhancing POI name matching using Phonetic Algorithm and smart string comparison. Because this research needed the comparison of names in String variable, when English alphabets of pronunciation are created, we have to compare keyword search from users with our database. They wouldn't be compared alphabet to alphabet. A technique called Approximate String matching will compare string similarity. For example, the word TANISHA and TANICHA are approximately the same word. Many researchers have conducted several studies in this topic. The research [4] focuses on the development of Approximate String Matching. Normally, this matching process is just about looking for words that string needs by comparing similar words. It means that the words don't have to be 100% alike. Two popular techniques are “Edit Distance” and “Soundex”. As Soundex [5] is just sound comparison, it has limitation of not being able to compare string that is not pronunciation, such as string representing DNA. For Edit distance, if it is used with LDAP server, it won't perform

This work was supported in part by the Department of Computer Science, Faculty of Science and Technology, Suan Sunandha Rajabhat University.

Kittiya Poonsilp is with the Department of Computer Science, Faculty of Science and Technology, Suan Sunandha Rajabhat University, Bangkok, Thailand (e-mail: kittiya.po@ssru.ac.th).

Attakorn Poonsilp was with Chulalongkorn University, Bangkok, Thailand. He is now with the Department of Computer Science, Faculty of Science and Technology, Suan Sunandha Rajabhat University, Bangkok, Thailand (e-mail: attakorn@gmail.com).

well enough. So this research was to improve the potentials of the work on LDAP.

Research [6] brought Approximate String Matching to compare people’s names. People’s name can be written in many formats, and they can be ambiguous. For example, case sensitive (JULich and Julich), abbreviation (R.E.Ellis and Randy E. Ellis) or names and last name format (Spike Jones and Jones, Spike). This research transformed all complication into graphs, and analyzed the graphs by Social Network Analysis.

Research [7] focuses on searching for address name from the database. Data on the database were basically wrong and input incorrectly from many sources. Users didn’t know how to spell correctly, specific names of location which might have foreign names, or those names might be homonyms. This research used Pattern matching and Phonetic matching algorithm together to create accuracy. Result showed that the accuracy rate was raised 7% more.

Research [8] concentrated on improving keyword search (locations, people’s name) to be more accurate. They improved by using soundex technique. Before, the system could only look for 658 names from 1,187 names. After adding Code shift and Diagrams into the system, accuracy was raised from 1,140 names from 1,187 names, or around 96%.

III. RESEARCH FRAMEWORK

The research framework is displayed in Fig. 1.

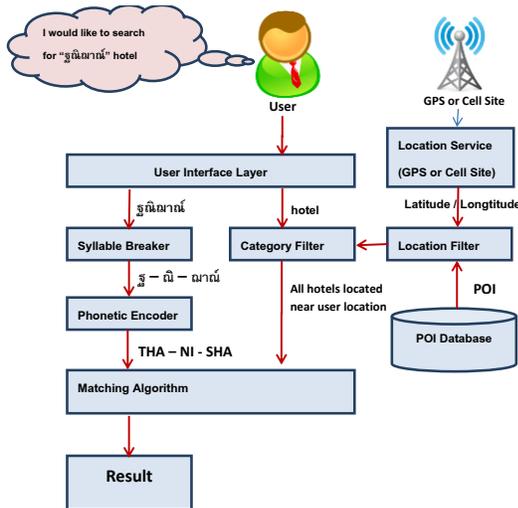


Fig. 1 how system works

From Fig. 1, to get a clearer picture, we would like to explain by an example. Supposing that a user is at Amphawa floating market and would like to look for a place called “ฐนิทาน (Thanisha)” which is in hotel category, the searching system will work in this order.

A. Syllable Breaker

This is the first step after receiving searching request from users. System will break word into syllables. From a single word, it will be cut into 3 syllables ฐ - นิ - ทาน.

B. Phonetic Encoder

Next, the system will encode the syllables into forms of

pronunciation by using English alphabets. Now we have “Tha-Ni-Sha” in our database.

C. Location Filter

When users’ location (from cell-site or GPS) is input, the system will use latitude and longitude to filter the location data in the database. The nearby location will be gathered and prepared for next step.

In this research, we used the technique of locating distance between 2 coordinates (POI and user location), to calculate and find out the most accurate result. Haversine Formula (1) was used here. It was created to find the distance between 2 coordinates, especially on a sphere object.

$$d = 2r \sin^{-1} \left(\sqrt{\sin^2 \left(\frac{\phi_2 - \phi_1}{2} \right) + \cos(\phi_1) \cos(\phi_2) \sin^2 \left(\frac{\lambda_2 - \lambda_1}{2} \right)} \right) \quad (1)$$

Refer to (1), d means distance on earth between 2 coordinates. r means earth radius, which is 6,371 kilometers. ϕ_1, ϕ_2 means Latitude of coordinates 1 and 2. λ_1, λ_2 means Longitude of coordinates 1 and 2 respectively.

After calculating, the result will be distance (in kilometers) between 2 coordinates as desired.

D. Category Filter

After receiving the nearby POIs, the next stage is to filter only kind of POI users want to look for, such as “guest house” or “hotel”. After this stage, we got all nearby hotels.

E. Approximate String Matching Algorithm

This stage is about taking user’s keyword and compare to the POIs we got from step D. The comparison is done syllable by syllable, with the adaptation of Levenstein Distance [9][10]. The search will be more effective. “Cha” and “Sha” will show the same result.

F. Result display with location

After comparison process in step E, the results will be shown to user. Each POI that’s matched will be plotted on Bing map.

The framework presentation is the overall picture of the process. Next is the phonetic algorithm process design, which is considered the main part of this research.

IV. PHONETIC ALGORITHM DESIGN

From Fig. 1, the phonetic algorithm to encode each syllable into pronunciation are considered the most important part of the system because it affects the ability and accuracy of searching process. As a result, this step will be separated into 2 phases.

Phase 1 is the pre-processing phase. The input words will be adjusted to be ready for encoding, such as cutting off garun and tone marks.

In phase 2, State machine will be introduced into this step to operate algorithm encoding to get the complete English pronunciation words.

Fig.2 illustrates the entire process of Phonetic Encoder.

Table 1 alphabet encoding

Thai alphabets	English alphabets		Cluster
	Consonant	Final consonant	
ก	K	K	No
ข	Kh	K	No
ฃ	Kh	K	No
ค	Kh	K	No
ฅ	Kh	K	No
ฆ	Kh	K	No
ง	Ng	Ng	No
จ	J	J	No
ฉ	Ch	T	No
ช	Ch	T	No
ซ	S	S	No
ฌ	Ch	T	No
ญ	Y	N	No
ฎ	D	D	No
ฏ	T	T	No
ฐ	Th	T	No
ฑ	Th	T	No
ฒ	Th	T	No
ณ	N	N	No
ด	D	d	No
ต	T	T	No
ถ	Th	T	No
ท	Th	T	No
ธ	TH	T	No
น	N	N	No
บ	B	B	No
ป	P	P	No
ผ	Ph	P	No
ฝ	F	F	No
พ	Ph	P	No
ฟ	F	F	No
ภ	Ph	b	No
ม	M	M	No
ย	Y	I	No
ร	R	Rn	Yes
ล	L	Rl	Yes
ว	W	W	Yes
ศ	S	S	No
ษ	S	S	No
ส	S	S	No
ห	H	H	No
ฬ	L	Rl	No
อ	A	O	No
ฮ	H	H	No

Thai alphabets	English alphabets		Cluster
	Consonant	Final consonant	
ฤ	Ru	-	No
ຸ	Rue	-	No
ູ	Ru	-	No
ຼ	Rue	-	No

Table 2 simple vowel encoding

Vowel	English alphabet	Position
ะ	a	After consonant
า	a	After consonant
ิ	i	After consonant
ี	ee	After consonant
ุ	ue	After consonant
ู	ue	After consonant
อ	u	After consonant
ู	uu	After consonant
เ	ae	Before consonant
แ	ae	Before consonant
โ	o	Before consonant
อ	or	After consonant
ัวะ	Ua	After consonant
ัว	ua	After consonant
ั	a	After consonant
ำ	am	After consonant
ไ	ai	Before consonant
ใ	ai	Before consonant
อ	o	After consonant

Table 3 compound vowel encoding

Vowels	Position		English alphabets
	Leading	Following	
เ-ะ	เ	ะ	ae
เ-อ	เ	อ	er
เ-า	เ	า	ao
เ-ะ	เ	ะ	er
เ-อาะ	เ	อาะ	o
เ-็	เ	-็	ear
เ-ือ	เ	-ือ	ueu
เ-็	เ	-็	ear
เ-ือ	เ	-ือ	ueu
แ-ะ	แ	ะ	ae
โ-ะ	โ	ะ	o

V. RESULTS

Research team got the collection of places names in each category (hotels, malls, universities, parks, etc.) in total of 100 items. The results are shown as below.

A. Result of locations name search

From the test, we found out that the framework and phonetic algorithm that we designed could find the correct result of location search input in Thai alphabets even though that word wasn't spelled correctly. In Table 4, the name of Natthawaree hotel was used for the test.

Table 4 shows the set of misspelled words that can be matched correctly to “ณัฐวารารี” hotel

Correctly spelled names	Phonetics Code	Examples of misspelled words
ณัฐวารารี	natthawaree	ณัฐวารารี
		ณัฐถวารารี
		ณัฐฐถวารารี
		ณัฐถวารารี
		ณัฐถวารารี
		ณัฐถวาราลี
		ณัฐฐวารารี
		นัตถวารารี
		นัตถวาราลี
		ณัฐวารารี
		ณัฐวาราลี
		ณัฐถวารารี
		ณัฐถวาราลี

According to Table 4, after inputting 20 misspelled keywords, the system could still find that place and locate it correctly as in Fig. 4.

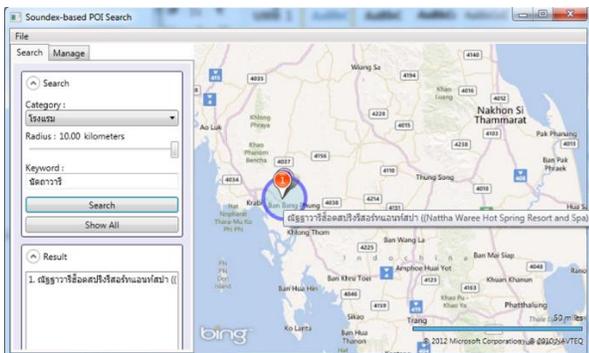


Fig. 4 result for "ณัฐวารารี"

From Fig. 4, after we input the place name and narrow down the search by selecting the category and distance, the system could perfectly show the place and location on the map.

B. Result after the test of ability to display similar information sorting by information from the least different to the most different

Since this research conducted the comparison of name similarity by using Approximate String Matching, the system can display related places that have similar names with the keyword and order by similarity as shown in Fig.5.

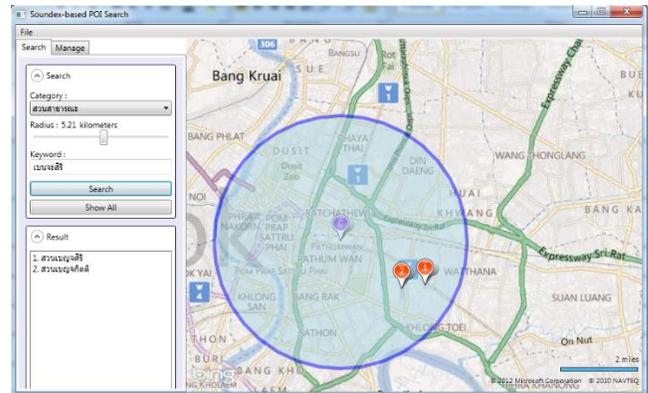


Fig. 5 the display of the search showing nearby places from the least difference to the most difference

C. Results of false keywords which were too different from the original keywords

From Table 4, we can see that no matter how wrong we misspelled the keyword, the system will still work and find the place. It is because the system already encoded the word and replaced alphabets with English, as well as getting help from Approximate String Matching which locates the similar pronunciation.

However, it doesn't mean that users can input "any words" into the search box without any format. The system will not work well if the input word is too much different from the original. Table 5 shows some examples.

Table 5 example of misspelled word that is too much different from the original

Correct names	Phonetics Code	Misspelled word that are too different from original
ณัฐวารารี	natthawaree	นัตถวารารีนาง
		Natthawareenang
		ณัฐวารารีนางง
		Natthawareenanggg
		นัตถวารารีมันนา
Natthawareemanna		
นัตถวารารี		นัตถวารารี
		Nattamanee
		ณัฐวารารี
		Nattha

From Table 5, we can see that all keywords were too much different from the original word, and it is too much for the system to search for. For example, the word *นัทธาวารีนาจ* (Natthawareenanggg), even though it has the similar pronunciation with the original, the word contains more syllables which are different from the original “Natthawaree” It is not accepted by the system.

D. Result of system test in terms of searching for places by categories

The research team tested the system to prove that the designed framework can work well with misspelled Thai keyword. It is hoped to solve that problem and help users to find correct results on search engine.

In the research, we used 100 location names and 100 misspelled items to test the system.

Table 6 result of 100 misspelled items

Units used in the research	Search ability		Total
	Find the result	Can't find the result	
Items	96	4	100

From Table 6, the system could find 96 items out of 100.

Table 7 some misspelled words and its result

Correct names	Phonetics code	Misspelled words	Matched	Not matched
ธรรมรินทร์	thornmarinth	ธรรมาริน	✓	
		ธรรมมารินทร์	✓	
		ท่ามาริน		✓
นครพรหม	nokroprohmly	นะคอนพรม		✓
		นาคอนพม		✓
		นครพม		✓
ธรรมภีร์กัษ	thornmaphirak	ธรรมะพีร์กั	✓	
		ท่าพีร์กั		✓
		ธรรมมาภีร์กัษ	✓	

From Table 7, we could see that the item containing “จ” are good example of unsearchable items. It is because the algorithm could not separate the “ร” words which could be cluster, vowel (รร) and consonant.

When words that have “ร” as final consonant, such as “ธรรมรินทร์” (Thammarin), is encoded, it will be pronounced “Thornmarinth”, which is not the correct name. On the other

hand, if we write “ท่ามาริน”, we will get “Tammarin”, which is still different from the original meaning and can't be searched for. Thus, if we use “ธรรมาริน” or “ธรรมารินทร์”, it can still be searched for after the approximate string matching process because it has the similarity with original word form.

VI. CONCLUSION

Location search engine using phonetic algorithm and LBS aims to improve the ability of search engine to solve problems in searching for locations in Thai language which can be homonyms, grammatically complex and ambiguousness. Thai language has garun (silence mark), tone marks and clusters which can cause difficulties in writing. This issue causes problem when users input misspelled words and they can't get the answer. The research team adapted phonetic algorithm to use with LBS by conducting the conversion of word into phonetic syllable, and compare them by Levenshtein Distance calculation to support the search that is not correctly spelled.

From the research, the framework showed 96% of accuracy, while only 4% could not be searched for. It was because the complication of Thai grammar. Although the technique used in this research can solve most misspelled search, some words, especially with “ร” and “รร” are hard to decode because they can be vowel, final consonant and cluster.

However, the framework that we designed can be adapted to use with other types of search engine in Thai language. It can also be used and developed in mobile device, such as car navigator, booking system and products location device on shopping carts in big shopping malls.

REFERENCES

- [1] Shu Wang, Jungwon Min and Byung K. Yi. “Location based services for mobiles: Technologies and standards.” *IEEE International Conference on Communication (ICC)* Beijing, China, 2008.
- [2] Weifeng Lv, Fei Wang, Yuan Zhang and Tongyu Zhu. “A distributed location based service framework of ubiquitous computing.” *IEEE Computer Society*. 201, page: 43-47.
- [3] Horozov T., Narasimhan N. and Vasudevan, V. “Using location for personalized POI recommendation in mobile environment.” *IEEE Computer Society*. 2006, page : 129 – 134.
- [4] Chi-Chien Pan, Kai-Hsiang Yang and Tzao-Lin Lee. “Approximate string matching in LDAP based on edit distance.” *IEEE Computer Society*. 2002.
- [5] Holmes D., McCabe, M.C. “Information technology: Coding and computing”, 2002. *Proceedings. International Conference on ITCC.2002*, Page : 22- 26.
- [6] Byung-Won On. “Social network analysis on name disambiguation and more.” *IEEE Computer Society*. 2008, page : 1081-1088.
- [7] Cihan Varol, John R. Talburt “Pattern and phonetic based street name misspelling correction.” *IEEE Computer Society*. 2011, page : 553-558.
- [8] Holmes D., McCabe M.C. “Improving precision and recall for soundex retrieval.” *International Conference on ITCC*, 2002, page : 22-26.
- [9] Li Yujian, Liu Bo, “A normalized Levenshtein distance metric”, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 29, no. 6, June 2007, page. 1091-1095.
- [10] Rane S., Wei Sun, “Privacy preserving string comparisons based on Levenshtein distance”. *Information Forensics and Security (WIFS), IEEE International Workshop*.2010, page : 1-6.

Improved Non-local Algorithm with Reliability of Neighbor Pixel

J. Lee, and J. Jeong

Abstract— The Non-local mean algorithm, which is one of the famous algorithms used in area of image denoising, uses the weighted sum of Euclidean distances between pixels. Euclidean distance is attained by calculating the difference between values of pixels or patches. If current pixel or patch has a similar value to the reference pixel or patch, it will have a large weight value. But if the patterns of current and reference are different, it will have small weight value. Let's assume there is a clear image where the patterns of current and reference pixels are the same. Then, after a pixel or patch(current or reference)attains noise through certain channel(for example, image acquisition or transmission), the weight value will be small and it will have a weak effect on the denoising calculation. Here, we need to pay attention to this "weak effect". Even if there is a noise or a big difference between pixels or patches, it will not have a significant effect on the denoising algorithm. Therefore, if we control this effect from noise and difference of patch values, we will be able to improve the result from denoising algorithm.

Keywords—image denoising, non-local mean algorithm, neighbor pixel, reliability of neighbor.

I. INTRODUCTION

MANY denoising algorithms are proposed for removing the noise in image. Previous denoising algorithms use neighbor-pixels in local for image denoising. Buades[1][2] developed a non-local mean algorithm(called NL-mean algorithm) that explores every patches on the image. Also the method uses weighted average, like previous algorithms. The weighted average is calculated by estimating similarity with other pixel or patch, and the similarities are calculated between a current pixel and all of the other reference pixels.

Darbon[3] proposed a fast NL-mean algorithm that makes the discrete integration of the squared difference of the image. When compute integral image of squared difference, it used shifting the image with respect to x and y coordinates. Then it shows good performances in lower computations and short

This research was supported by the MSIP(Ministry of Science, ICT & Future Planning), Korea, under the "Establishing IT Research Infrastructure Projects" supervised by the NIPA(National IT Industry Promotion Agency) (I2221-14-1005).

J. Lee is with the Department of Electronics and Computer Engineering, Hanyang University, 222, Wangsimni-ro, Seongdong-gu, Seoul, Korea (e-mail: mizaling@gmail.com).

J. Jeong is with the Department of Electronics and Computer Engineering, Hanyang University, 222, Wangsimni-ro, Seongdong-gu, Seoul, Korea (corresponding author to provide phone: +82-2-2220-4369; e-mail: jjeong@hanyang.ac.kr).

timing for denoising.

In this paper, we proposed "The Reliability of Neighbor - pixel" and the method to improve performance using the reliability. The other parts of this paper are organized as follow. Introduce some points of the NL-mean algorithm and fast NL-mean algorithm for approaching to proposed algorithm in Section II. The proposed algorithm is specified in Section III. Show the experimental results and analysis in Section IV, and Section V concludes the paper.

II. THE NON-LOCAL MEAN ALGORITHM

A. Non-local mean algorithm

The NL-mean algorithm uses the weighted sum of Euclidean distance between the pixels with neighbor pixels were called patch unit. And There are two images, one that is named image u and consists of pixels u_k is the noisy image and another is named image \hat{u} and consist of pixels \hat{u}_k is the result image of denoising. Follow NL-mean algorithm, the denoised pixel \hat{u}_i is computed as :

$$\hat{u}_i = \sum_{j \in I} w_{ij} u_j \quad (1)$$

where the given current pixel u_i which is the object of the noise reduction and the weight w_{ij} is the similarity that is calculated by comparison of patch i and j . The calculation of weight progress all of pixel in image I, as non-locally. But actually we use the enough widen window region that has fixed window size instead of all of pixel in image I. The weight w_{ij} satisfies two conditions $0 \leq w_{ij} \leq 1$ and $\sum_{j \in I} w_{ij} = 1$. And the weight w_{ij} is defined as,

$$w_{ij} = \frac{1}{z_i} \exp\left(-\frac{1}{h^2} \|P_i - P_j\|^2\right) \quad (2)$$

where P_k denotes a patch of fixed size and centered at pixel u_k , h is a smoothing constant, and z_i is the normalizing constant

$$z_i = \sum_j \exp\left(-\frac{1}{h^2} \|P_i - P_j\|^2\right). \quad (3)$$

It can express the flowchart of figure 1 without the proposed part.

B. Fast non-local mean algorithm

Previous NL-mean algorithm, the most time taking part is the weight computation. (2) represents as patch unit, if the weight equation represents as pixel unit in patch, then

$$w_{ij} = \frac{1}{z_i} \exp\left(-\frac{1}{h^2} \sum_{\delta \in \Delta} (v(i + \delta) - v(j + \delta))^2\right) \quad (4)$$

where Δ means the patch, and δ means the sequence of pixel.

Let a new image S_{d_x} is the discrete integration of the squared difference of the image v and it has translation by given translation vector d_x , then S_{d_x} is defined as,

$$S_{d_x}(p) = \sum_{k=0}^p (v(k) - v(k + d_x))^2 \quad (5)$$

where $p \in [[0, n - 1]]$ and n is the number of pixels of the image. Recall the patch $\Delta = [[-P, P]]$ as 1D form, and let the vector $d_x = j - i$, $\hat{p} = i + \delta_x$, and substitute the parameters on the weight (4). Then we get the equation as,

$$w_{ij} = \frac{1}{z_i} \exp\left(-\frac{1}{h^2} \sum_{p=s-P}^{s+P} (v(\hat{p}) - v(\hat{p} + d_x))^2\right). \quad (6)$$

Use (5), and split the sum then the weight equation represents as,

$$w_{ij} = \frac{1}{z_i} \exp\left(-\frac{1}{h^2} (S_{d_x}(s + P) - S_{d_x}(s - P))\right). \quad (7)$$

Finally, the (7) makes us to calculate the weight for a pair of pixels in constant time and we can save time.

III. PROPOSED ALGORITHM

A. Reliability of neighbor pixel

In previous NL-mean algorithm, every neighbor pixels of patch is used for compute weight. If the patch P_i has the similarity with P_j , then the weight w_{ij} is a large value. But if two patches do not have similar pattern, then w_{ij} will be smaller.

If two patches are too different, then zero weight is better than small weight value. So let the weight of two differently patches is “unreliability”. The other way, if two patches have similarity, then it is good “Reliability”.

Assume the condition that two patches are same, if a patch is polluted, then unfortunately the reliability of weight will be bad. And it has a bad influence on denoising process.

B. Elimination of bad neighbors in NL-mean algorithm

The bad reliability is worse than any effect. So if the reference pixel has bad reliability, because it is took noise or is different in the beginning, than it is better than bad reliability to be zero effect. So in this paper, we propose to change the

bad effect weight to zero effect by elimination the bad reliability of effect.

In previous NL-mean algorithm, change zero the last k-th value of weight, where k is user defined variable, then decrease the effect of small weight. It means that smallest k-th value is excluded in the weighted sum z_i and NL-mean calculation as Figure 1.

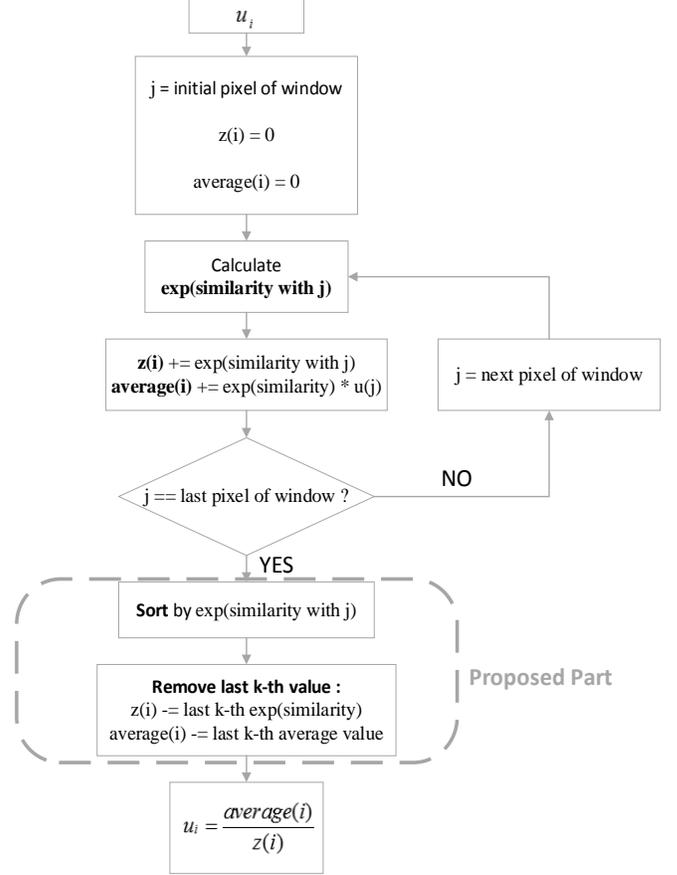


Fig. 1 Flowchart for proposed algorithm

C. Elimination of bad neighbors in fast NL-mean algorithm

The fast NL-mean algorithm uses a system of integration of the squared difference of the image. It has a difference that it uses the image unit, and previous algorithm uses the pixel unit in calculation of weight. So it can't use the flowchart in figure 1. Specifically, “calculate **exp(similarity with j)**” part is the repeat section by pixel unit in previous NL-mean algorithm, but these part is the calculation by 1D patch unit with pair of pixels in fast NL-mean algorithm. At this time, fast NL-mean algorithm makes the new image(1D form) that contains integral image of the input image. Although it can not remove a few pixels that are unreliability, use the thresholding method and we can eliminate some of pixels which has the lower reliability. Along this assumption, the hard thresholding is suitable method. But by experimental results the soft thresholding makes better results than the hard thresholding.

Table 1. Experimental results on previous NL-mean algorithm in PSNR(in dB)

	σ_{noise}	10	20	30	40	50	60	70	80	90	100
Lena	NLM	31.7503	31.5685	29.5833	27.5571	26.1599	25.1718	24.3678	23.6387	22.9512	22.3080
	Proposed $k=2$	31.7502	31.5599	29.5758	27.5833	26.2181	25.2466	24.4477	23.7164	23.0256	22.3775
	Proposed $k=3$	31.7502	31.5495	29.5674	27.6093	26.2755	25.3162	24.5177	23.7825	23.0868	22.4315
	Proposed $k=4$	31.7502	31.5378	29.5593	27.6333	26.3243	25.3765	24.5775	23.8381	23.1349	22.4748
	Proposed $k=5$	31.7501	31.5247	29.5491	27.6555	26.3700	25.4282	24.6270	23.8829	23.1747	22.5064
	Proposed $k=6$	31.7501	31.5117	29.5378	27.6742	26.4089	25.4724	24.6694	23.9177	23.2035	22.5292
	Proposed $k=7$	31.7500	31.4972	29.5273	27.6899	26.4445	25.5101	24.7011	23.9443	23.2234	22.5436
		σ_{noise}	10	20	30	40	50	60	70	80	90
Barbara	NLM	30.9347	30.0076	27.1362	24.3918	23.0086	22.2417	21.6759	21.1688	20.6808	20.2102
	Proposed $k=2$	30.9347	30.0042	27.1451	24.4432	23.0882	22.3305	21.7636	21.2519	20.7580	20.2812
	Proposed $k=3$	30.9347	30.0005	27.1560	24.4968	23.1680	22.4170	21.8461	21.3266	20.8257	20.3415
	Proposed $k=4$	30.9347	29.9961	27.1671	24.5489	23.2430	22.4958	21.9200	21.3935	20.8852	20.3939
	Proposed $k=5$	30.9346	29.9917	27.1793	24.6024	23.3154	22.5695	21.9878	21.4537	20.9363	20.4393
	Proposed $k=6$	30.9346	29.9866	27.1910	24.6535	23.3857	22.6390	22.0498	21.5073	20.9810	20.4775
	Proposed $k=7$	30.9346	29.9810	27.2037	24.7033	23.4515	22.7030	22.1070	21.5552	21.0199	20.5107
		σ_{noise}	10	20	30	40	50	60	70	80	90
Boat	NLM	27.2676	25.0202	23.7105	22.9366	22.3749	21.8869	21.4292	20.9932	27.2676	25.0202
	Proposed $k=2$	27.2729	25.0547	23.7713	23.0099	22.4482	21.9578	21.4954	21.0549	27.2729	25.0547
	Proposed $k=3$	27.2796	25.0911	23.8320	23.0788	22.5163	22.0202	21.5524	21.1061	27.2796	25.0911
	Proposed $k=4$	27.2851	25.1258	23.8894	23.1412	22.5774	22.0748	21.5997	21.1466	27.2851	25.1258
	Proposed $k=5$	27.2912	25.1611	23.9435	23.1999	22.6319	22.1229	21.6417	21.1810	27.2912	25.1611
	Proposed $k=6$	27.2968	25.1942	23.9947	23.2532	22.6796	22.1636	21.6739	21.2074	27.2968	25.1942
	Proposed $k=7$	27.3034	25.2253	24.0424	23.3023	22.7233	22.1989	21.7026	21.2287	27.3034	25.2253

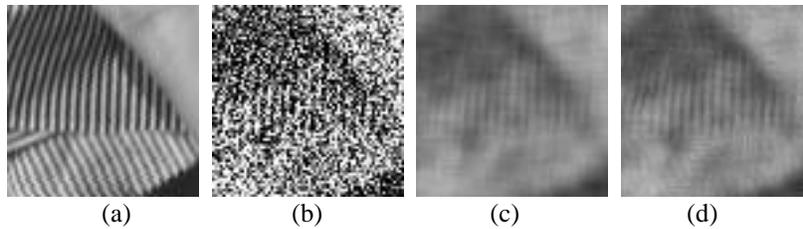


Fig. 2 The part of image - barbara with $\sigma = 100$

(a) original image (b) noisy image (c) result of NL-means (d) proposed on previous algorithm($k=7$)

IV. EXPERIMENTAL RESULTS

We implement the proposed algorithm in MATLAB and use 3 images(Lena, Barbara, and Boat) that have 512x512 size. Experiments use Gaussian noise that has zero mean and the 10 variances value which are from 100 down to 10 by unit 10. When we experiment on previous NL-mean algorithm, the number of elimination is defined by variable k which is 2 to 7. In experiment on fast NL-mean algorithm, both soft thresholding and hard thresholding are used to integration image and the threshold value is approached by empirical case study.

A. Experiment on previous NL-mean algorithm

There are the results on Table 1. If the noise variance increases, then the number of bad neighbors increases, too. Along Lena and Barbara, when the noise variance is small, the

elimination of least weight is an opposite effect for denoising, because usually neighbors don't have noise or it means the original image has many repeat patterns.

The other way, when noise variance is large, especially the proposed algorithm is effective in repeat patterns like Fig. 2. The elimination of bad neighbors brings the effect of reliable neighbors into relief.

When the noise variance is large, there are better performances in PSNR from 0.0262 dB to 0.3383 dB in image Lena, from 0.0089 dB to 0.4613 dB in image Barbara, and from 0.0053 dB to 0.3657 dB in image Boat.

When noise variable $\sigma=50$ and the proposed algorithm execute with $k=7$, unfortunately the result doesn't have remarkably point in Fig. 3. But if you get this result 512 x 512 size image on display, you can find the specific result exactly.

B. Experiment on fast NL-mean algorithm

In this experiment, the threshold value is a variable which is

Table 2. Experimental results on fast NL-mean algorithm in PSNR(in [dB])
(Th 1 = 5e-5, Th 2 = 1e-4, Th 3 = 5e-4, Th 4 = 1e-3, and Th 5 = average of integration difference)

	σ_{noise}	10	20	30	40	50	60	70	80	90	100	
	Lena	NLM	31.7503	31.5685	29.5833	27.5571	26.1599	25.1718	24.3678	23.6387	22.9512	22.3080
Fast NLM		30.5100	29.5547	28.4873	27.0489	25.9776	25.1240	24.3626	23.6430	22.9580	22.3127	
Hard Th.		Th 1	30.5094	29.5545	28.4873	27.0489	25.9776	25.1240	24.3626	23.6430	22.9580	22.3127
		Th 2	30.5079	29.5545	28.4873	27.0489	25.9776	25.1240	24.3626	23.6430	22.9580	22.3127
		Th 3	30.4890	29.5506	28.4871	27.0488	25.9777	25.1240	24.3425	23.6430	22.9580	22.3127
		Th 4	30.4634	29.5430	28.4867	27.0489	25.9777	25.1241	24.3425	23.6429	22.9580	22.3127
		Th 5	31.9822	29.2888	28.3873	27.0362	25.9818	25.1300	24.3654	23.6445	22.9583	22.3132
Soft Th.		Th 1	30.5117	29.5551	28.4873	27.0489	25.9776	25.1240	24.3625	23.6430	22.9580	22.3127
		Th 2	30.5146	29.5557	28.4874	27.0489	25.9776	25.1240	24.3625	23.6430	22.9580	22.3127
		Th 3	30.5595	29.5640	28.4878	27.0489	25.9776	25.1239	24.3625	23.6429	22.9580	22.3127
		Th 4	30.6419	29.5788	28.4883	27.0488	25.9774	25.1239	24.3624	23.6429	22.9579	22.3127
		Th 5	32.4275	30.1172	28.3648	26.9079	25.8763	25.0529	24.3144	23.6111	22.9355	22.2977
Barbara		σ_{noise}	10	20	30	40	50	60	70	80	90	100
		NLM	30.9347	30.0076	27.1362	24.3918	23.0086	22.2417	21.6759	21.1688	20.6808	20.2102
		Fast NLM	29.8891	28.2418	26.5242	24.3667	23.0797	22.3047	21.7217	21.2011	20.7033	20.2263
	Hard Th.	Th 1	29.8888	28.2418	26.5242	24.3667	23.0797	22.3047	21.7217	21.2011	20.7033	20.2263
		Th 2	29.8883	28.2415	26.5243	24.3667	23.0797	22.3047	21.7217	21.2011	20.7030	20.2263
		Th 3	29.8762	28.2388	26.5243	24.3668	23.0797	22.3047	21.7218	21.2011	20.7033	20.2263
		Th 4	29.8597	28.2338	26.5248	24.3670	23.0797	22.3047	21.7218	21.2011	20.7033	20.2262
		Th 5	30.8882	28.0792	26.4852	24.3989	23.1036	22.3182	21.7283	21.2042	20.7050	20.2270
	Soft Th.	Th 1	29.8907	28.2420	26.5242	24.3667	23.0797	22.3047	21.7217	21.2011	20.7033	20.2263
		Th 2	29.8928	28.2424	26.5241	24.3666	23.0797	22.3046	21.7217	21.2011	20.7033	20.2262
		Th 3	29.9203	28.2488	26.5239	24.3663	23.0795	22.3046	21.7217	21.2011	20.7033	20.2262
		Th 4	29.9719	28.2599	26.5232	24.3659	23.0794	22.3045	21.7216	21.2010	20.7032	20.2262
		Th 5	31.1186	28.5785	26.1038	24.0607	22.9116	22.2105	21.6649	21.1644	20.6792	20.2091
	Boat	σ_{noise}	10	20	30	40	50	60	70	80	90	100
		NLM	30.8755	29.8155	27.2676	25.0202	23.7105	22.9366	22.3749	21.8869	21.4292	20.9932
Fast NLM		29.8907	28.2420	26.5242	24.3667	23.0797	22.3047	21.7217	21.2011	20.7033	20.2263	
Hard Th.		Th 1	29.8608	28.2591	26.6070	24.8786	23.7313	22.9784	22.4089	21.9114	21.4466	21.0053
		Th 2	29.8593	28.2590	26.6069	24.8786	23.7313	22.9784	22.4089	21.9114	21.4466	21.0053
		Th 3	29.8467	28.2567	26.6070	24.8787	23.7313	22.9783	22.4089	21.9114	21.4466	21.0053
		Th 4	29.8319	28.2522	26.6071	24.8788	23.7313	22.9784	22.4089	21.9115	21.4466	21.0053
		Th 5	30.3549	27.9575	26.5100	24.8761	23.7422	22.9868	22.4133	21.9141	21.4475	21.0061
Soft Th.		Th 1	29.8623	28.2594	26.6070	24.8786	23.7313	22.9784	22.4089	21.9114	21.4466	21.0053
		Th 2	29.8644	28.2595	26.6069	24.8785	23.7312	22.9783	22.4088	21.9114	21.4466	21.0053
		Th 3	29.8919	28.2650	26.6070	24.8784	23.7311	22.9783	22.4088	21.9114	21.4465	21.0053
		Th 4	29.9440	28.2745	26.6067	24.8780	23.7310	22.9782	22.4087	21.9114	21.4465	21.0053
		Th 5	30.3997	28.2491	26.2595	24.6363	23.5818	22.8888	22.3537	21.8765	21.4232	20.9896

is influential with the proposed algorithm. Darbon [3]’s fast NL-mean algorithm makes worse result than previous NL-mean algorithm from the viewpoint of PSNR, in small noise variance. But in the viewpoint of subjective image quality evaluation, it can’t be looked differently from previous algorithm to Darbon’s, as Fig. 4. For demonstration of the proposed algorithm on fast NL-means, we compared Darbon’s algorithm with ours.

Belong assumption of the proposed algorithm, the hard thresholding is appropriate, however sometimes the soft

thresholding has better results than hard thresholding when the noise variance is small(10 or 20) in Table 2. However if the noise variance is larger than almost $\sigma = 30$, then the hard thresholding is effective for image denoising.

By hard thresholding on large noise variance, the PSNR performance is better than fast NL-means from 0.0003 dB to 0.0060dB in image Lena, from 0.0007 dB to 0.0322 dB in image Barbara, and from 0.0008 dB to 0.0109 dB in image Boat. At that time, threshold value is used average of integration. It means that almost half of the neighbors have

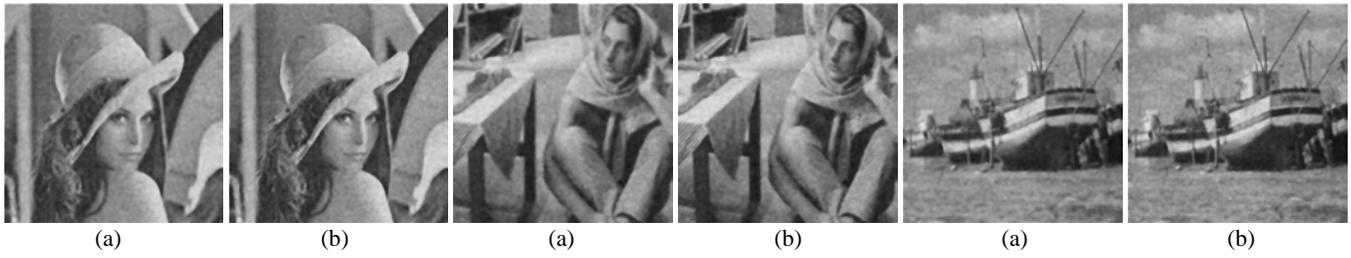


Fig. 3 Comparing experimental result with previous NLM : (a) the previous NLM (b) the proposed algorithm with $k=7$



(a) previous NLM (b) Fast NLM (c) proposed algorithm (soft thresholding)

Fig. 4 Comparing previous NLM, Darbon's fast NLM, and proposed algorithm(soft thresholding) on $\sigma=10$

unreliability and denoising is disturbed by them.

By soft thresholding on small noise variance like $\sigma = 10$, the PSNR performance is better than fast NL-means 1.9175 dB in image Lena, 1.2295 dB in image Barbara, and 0.5383 dB in image boat. When the noise variance is small, the soft thresholding removes bad neighbors and reconstructs the weight to be increased relative importance.

In Fig. 4, the proposed algorithm makes smoother in flat region. Because it eliminates speckled shapes, it has higher PSNR than Darbon's fast NLM.

V. CONCLUSIONS

Many denoising methods, especially NL-mean algorithm, use the weighted sum of neighbors. It uses similarity, and we confirmed the reliability that means all of the weights are helpful for image denoising by the experiments.

The elimination of bad neighbors which is unreliable improves the denoising effects. In this paper, we removed the bad neighbors by threshold techniques easily. From this experiments, if we develop other methods that distinguish bad neighbors exactly, remove them, and calculate upon the effect of reliable neighbors like the soft thresholding in Fig. 4, then we can expect more denoising effects.

By this experiments, if we choose adaptively hard thresholding or soft thresholding by the noise variance, then we can get more clear image by image denoising.

REFERENCES

- [1] A. Buades, B. Coll, and J. Morel, "A Non-Local Algorithm for Image Denoising," *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on*, vol. 2, 2005.
- [2] a. Buades, B. Coll, and J. M. Morel, "Image Denoising Methods. A New Nonlocal Principle," *SIAM Rev.*, vol. 52, no. 1, pp. 113–147, Jan. 2010.
- [3] J. Darbon, A. Cunha, T. Chan, S. Osher, and G. Jensen, "Fast nonlocal filtering applied to electron cryomicroscopy," *Biomedical Imaging: From Nano to Macro, 2008. ISBI 2008. 5th IEEE International Symposium on*, pp. 1331–1334, 2008.

Junghyun Lee received a B.S degree the Department of Electronics and Computer Engineering from Hanyang University, Korea, in 2014. He is currently puersuing a M.S in Electronic and Computer Engineering at Hanyang University. His research interests include video processing, and high efficiency video codec.

Je-Chang Jeong received a BS degree in electronic engineering from Seoul National University, Korea, in 1980, an MS degree in electrical engineering from the Korea Advanced Institute of Science and Technology in 1982, and a PhD degree in electrical engineering from the University of Michigan, Ann Arbor, in 1990. From 1982 to 1986, he was with the Korean Broadcasting System, where he helped develop teletext systems. From 1990 to 1991, he worked as a postdoctoral research associate at the University of Michigan, Ann Arbor, where he helped to develop various signal-processing algorithms. From 1991 through 1995, he was with the Samsung Electronics Company, Korea, where he was involved in the development of HDTV, digital broadcasting receivers, and other multimedia systems. Since 1995, he has conducted research at Hanyang University, Seoul, Korea. His research interests include digital signal processing, digital communication, and image and audio compression for HDTV and multimedia applications. He has published numerous technical papers.

Dr. Jeong received the Scientist of the Month award in 1998, from the Ministry of Science and Technology of Korea, and was the recipient of the 2007 IEEE Chester Sall Award and 2008 ETRI Journal Paper Award. He was also honored with a government commendation in 1998 from the Ministry of Information and Communication of Korea.

Urban Traffic Management Approach Based on Ontology and VANETs

H. TOULNI, B. NSIRI, M. BOULMALF and T. SADIKI

Abstract— Everyone knows the important role of transport in the economic and social development, at the level global this sector is considered among the principal criteria for country rankings. But control and management of traffic in urban areas remains a major challenge to rise, especially accuses rapid development of the transport sector. On the other hand, many alternative solutions have been proposed within the intelligent transportation systems (ITS) to optimize traffic safely. To address this problem, among them Vehicular Ad-hoc NETWORK (VANET), that will implement new innovative applications and services. In this paper, we propose an approach using ontologies and VANET to enable more efficient and optimal use of road infrastructure.

Keywords—Ontology, Traffic Management, Vehicular Ad-hoc Networks, Safety.

I. INTRODUCTION

No one can ignore today the crucial role of transport in developing countries, whether it's on the economic or social plane. At the global level this sector is considered among the principal criteria of development, also this sector is growing rapidly and continuously. However, the transport sector is currently confronted with significant challenges, particularly in urban areas.

These transport problems reduce the economic development opportunities, and quality of life of citizens who are affected psychologically and physically. Among these problems that may be mentioned traffic congestion, increased energy consumption, waste of time, limited mobility and degradation of air quality. Moreover, transport accidents are the most serious problems because of their socio-economic damages, including property damage and human losses.

To remedy this problem, intelligent transportation systems (ITS) become an alternative to optimize traffic safely.

ITS have emerged as an effective way to improve circulation, it will help to use less energy in the travel, less distance to reach the desired position with a time and money,

Hamza TOULNI is with LIAD Laboratory. Faculty of Sciences Ain Chock Casablanca, University Hassan II, Morocco. (E-mail: hamza.toulni07@etude.univcasa.ma).

Benayad NSIRI is with LIAD Laboratory. Faculty of Sciences Ain Chock Casablanca, University Hassan II, Morocco. (E-mail: b.nsir@fsac.ac.ma).

Mohammed BOULMALF is with International University of Rabat, Sala El Jadida, Morocco. (E-mail: mohammed.boulmalf@uir.ac.ma).

Tayeb SADIKI is with International University of Rabat, Sala El Jadida, Morocco. (E-mail: tayeb.sadiki@uir.ac.ma).

while respecting nature. And one of the effective ways possible to have these benefits at lower cost, are Vehicular Ad-hoc NETWORK (VANET) [1].

VANET are based on mobile ad hoc networks (MANET), where each vehicle is equipped with wireless communication devices, which are designed to ensure communication between vehicles and the road infrastructure, and therefore play a crucial role in providing innovative applications and services [1, 2, 3] in the road transport sector. These applications and services are designed not only to improve road safety but also for comfort, support and entertainment.

VANETs uses specialized short-range protocol communications (DSRC) [4, 5] to broadcast messages at high speed in several directions [6, 7] because its latency is low, but the coverage of this solution is very limited. To overcome this problem, researchers proposed V2V communication [8, 9], so that vehicles can further communicate with Road Side Unit (RSU), and VANETs does not require a significant investment for implementation. In addition to high-speed connectivity at lower cost, vehicles equipped with VANET devices can take advantage of multiple location technologies with high accuracy [10], either with a relative location [12,11, 13] or even a global location [14, 15,16].

In this paper, we propose an approach using ontologies and VANET to enable more efficient and optimal use of road infrastructure. The rest of this article is organized as follows. First, we give a literature review of traffic management systems. Second, we give an overview of the ontologies we show the main solutions based ontologies. In section 3 we present our approach. Finally, the conclusions and future research are shown in section 4.

II. LITERATURE REVIEW

In recent years, several articles have been published about the management of urban traffic, these articles fall into two broad categories: Estimated circulation and the optimization of traffic.

The traffic estimate is mainly based on analytical modeling of data collected by sensors installed all along the roads, or even vehicles.

The information collected in real time are also used in the optimization of traffic, but the techniques and methods change.

For example, the works presented in [17,18] the road is

defined congested state when the vehicle travel time exceeds the normal travel time of this road, the normal travel time of each segment of a road must be calculated by the vehicle for a day and then stored in a centralized entity.

In [19,20] the authors introduced mechanisms to detect traffic jams, which are mainly based on messages regularly broadcast by vehicles. The estimate of the traffic and the status of various routes are evaluated by analysis of the information broadcast messages.

However, the problem that arises in these mechanisms is overload of the communication channel, because they requires the exchange of a large number of packets.

In[21] authors propos an Adaptive Traffic Control Systems (ATCSs) utilize real time traffic data in an attempt to optimize the timing and length of the traffic light signals. As a result, effective ATCSs aim to minimize stop times and delays in a bid to reduce traffic congestion in major urban areas.

Another strategy in urban traffic management is to optimize traffic signals [22, 23, 24] deployed at intersections by analyzing the data collected in real-time traffic. The goal of this optimization is to minimize waiting times in an intersection and increase the number of vehicles crossing the intersection. Then you have to synchronize the lights different intersections to improve traffic in all directions

However, a local synchronization for an intersection influence on all other intersections of the road network, thus the optimizing desired goals include the minimizing of the waiting time and the length of the queue will not be achieved in other intersections, which could cause more congestion. For this, researchers have proposed to favor roads with high demand but for special events or temporary changes such as road closures due to construction or other, which results inadequacy of this strategy if the traffic is huge.

III. ONTOLOGIES IN INTELLIGENT TRANSPORTATION SYSTEM

Nowadays, ontologies are highly valued in almost all areas, for this reason we find many definitions of an ontology, the simplest and most popular since 1993 until now is the definition of Tom Gruber [25] who has said: "An ontology is a formal, explicit specification of a shared conceptualization."

In other words, an ontology is a structured set of terms and concepts representing the information and the relations between them, in a specific domain, these relations can be semantic relations, or relations of composition and inheritance. The power and usefulness of ontology is the reuse of information and the definition of a common vocabulary, and in addition any domain can be modeled using ontologies.

The main element required for the construction of ontology is language, it is designed to describe the information and allow their reuses. In the last few years, many languages have been developed to the implementation, these languages are classified into four levels: informal, semi-informal, semi-formal and formal, this is why the ontologies are not all built by the same way, but the choice of language is a challenge for construction.

Otherwise, several articles have been published on ontologies as solutions to the problems and the challenges of ITS. In [26] authors present the VEhicular ACCident ONtology (VEACON) designed to improve traffic safety, and for enabling interoperability between vehicles, RSUs, authorities and emergency vehicles. This ontology combines the information collected when an accident occurs, and the data available in the General Estimates System (GES) accidents database.

In [27] authors present an ontology for a reliable Traffic Information System. This ontology had been developed in OWL, and it is based on road traffic, and on possible scenarios of vehicles traveling in a highway. It is composed by classes, properties, attributes and relations between classes. The ontology is included in each agent executing the Traffic Information System. Each agent may ask for traffic information based on the ontology, and also based on its knowledge base.

In [28] authors propose a method to increase situation awareness during emergency transportation of patients. Their approach combines semantic reasoning with the emerging Car-2-X technology. The developed system continuously matches data retrieved from inter-vehicular communication with structured knowledge from vehicular ontologies and OpenStreetMap.

In [29] authors present the Car Accident lightweight Ontology for VANETs (CAOVA). The instances of our ontology are filled with: (i) the information collected when an accident occurs, and (ii) the data available in the General Estimates System (GES) accidents database. We assess the reliability of our proposal in two different ways: one via realistic crash tests, and the other one using a network simulation framework.

In [30, 31] authors propose ontology-based approaches for adding reasoning capabilities to autonomous vehicles. The main use case is at self-assessment of the perception system to monitor co-driving. The module designed for situation assessment formalizes knowledge such as: environment conditions, moving obstacles, driver state, navigable space, which are also relevant concepts for VANET.

IV. OUR APPROACH

Improving road safety requires constant supply of traffic information to the driver, this information should also improve the driving quality and keep traffic moving. But it is difficult to acquire all information and interpreted by the pilot in this context our approach is involved to enable more effective use of road infrastructure safely.

VANET can provide information faster and more pertinently in real time, but the interpretation of this information by the driver and his reflexes are not always precise. For this we propose an ontology that will ensure the best presentation of collected information.

A. Overview

We propose an ontology using VANET in order to facilitate

driving and interpretation of messages to the driver. This ontology is integrated directly into each vehicle, and it also communicates with the infrastructure to obtain traffic information in real time.

So our proposed approach consists of three main phases:

A learning phase is to collect information on the infrastructure, in order to reconstitute the map and the connections between roads.

A phase of knowledge acquisition of acquiring the information necessary for the driver.

And finally a Knowledge Representation phase.

This ontology “Fig. 1” consists of four subclasses: vehicle, Infrastructure, Traffic Control and Message. These concepts that relate to each other.

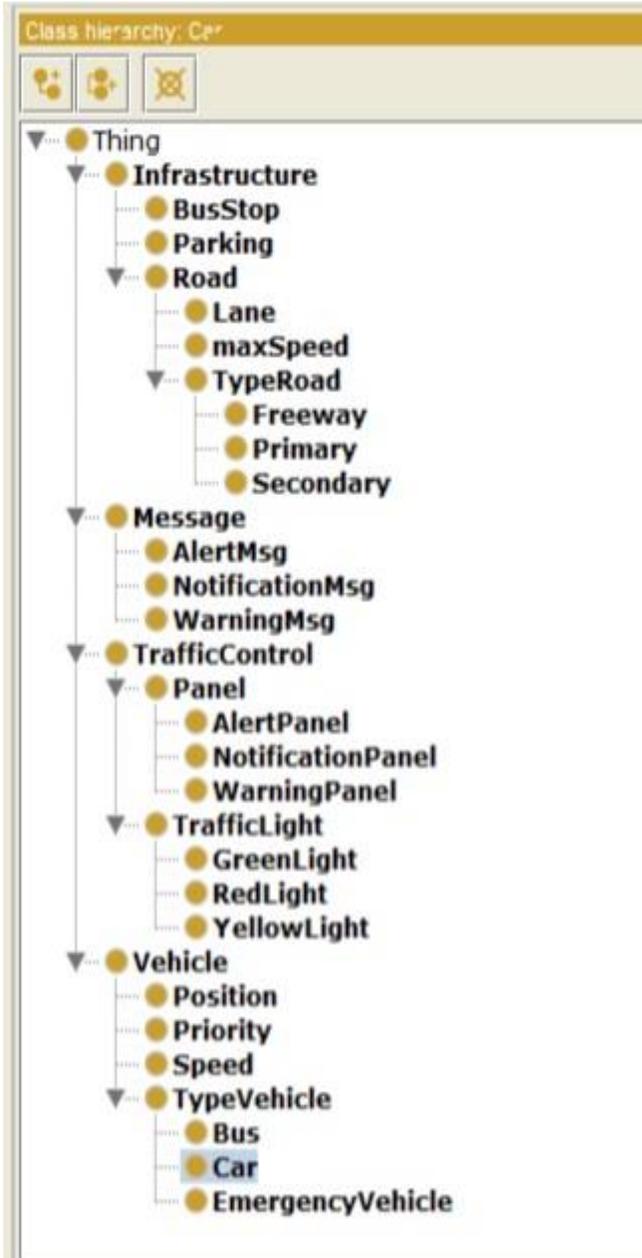


Fig. 1 Component of the Ontology

B. Language

We choose the Web Ontology Language (OWL) as a language to describe and organize knowledge for our ontology, it is developed and recommended by the World Wide Web Consortium (W3C). OWL is designed for use by applications that need to process the content of information instead of just presenting information to humans. OWL facilitates greater machine interpretability of the content that supported by XML, RDF, and SRDF.

C. Design

Our ontology was designed using Protégé [32], it begins with a super class named Thing “Fig. 2”, which all other classes are subclasses. This brings us directly to the concept of inheritance, therefore inherited classes are: Vehicle, Infrastructure, Message and TrafficControl.

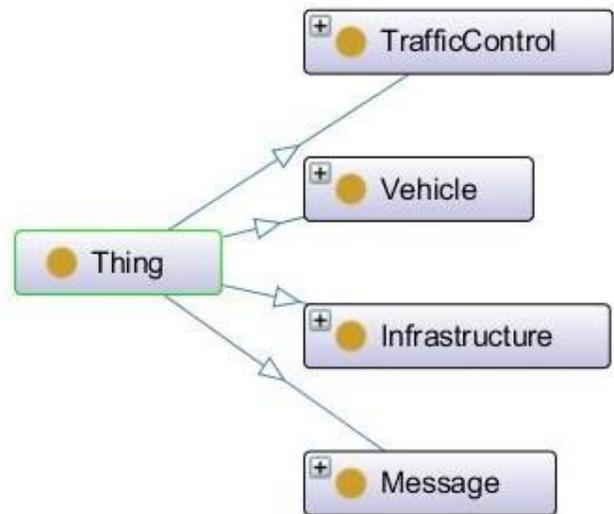


Fig. 2 High Level Ontology

The first class is Vehicle “Fig. 3”, which include the properties of Vehicle, such as Priority, Position and Speed. Vehicle comes in three types: simple car, bus and emergency vehicles. TypeVehicle describes the vehicle's physical properties and their priority.

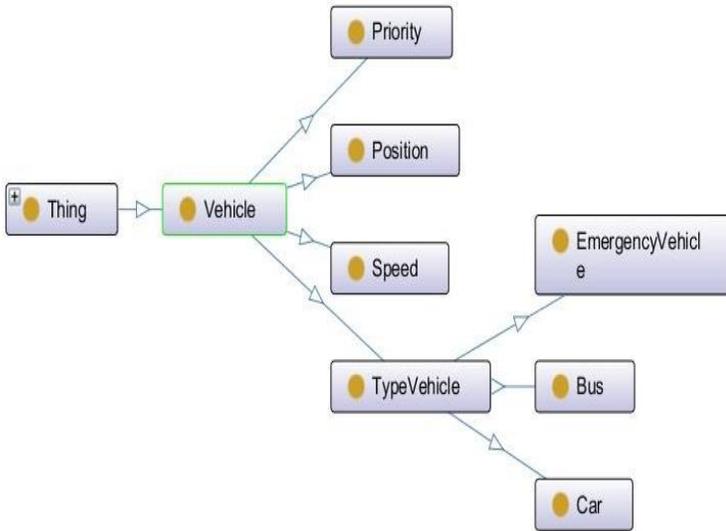


Fig. 3 Description of vehicle

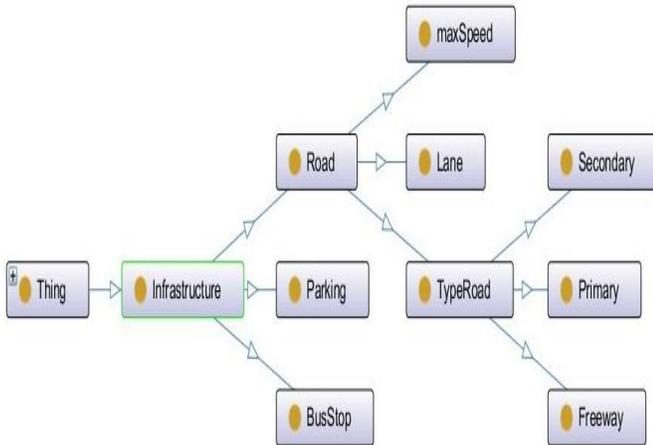


Fig. 4 Description of Infrastructure

The class Infrastructure “Fig. 4”, which is composed of Road, Parking and BusStop for bus station, and each Road has a number of lanes for the rolling of the vehicle, and a maximum speed not to exceed by vehicles, and her type which include the properties of road.

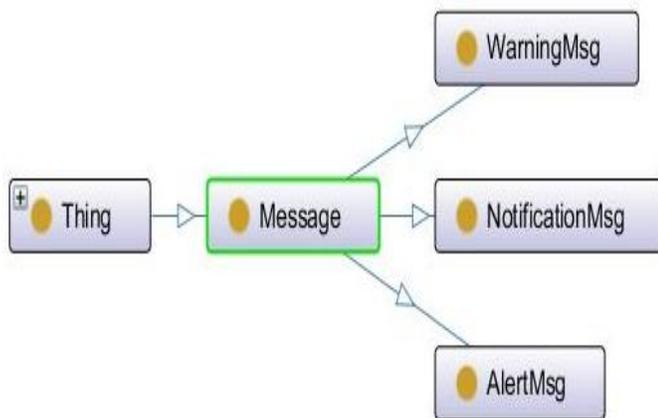


Fig. 5 Description of Message

The class Message “Fig. 5”, which include the type of Message, it can be an AlertMsg for emergency situations, Warning for unpredictable situations or NotificationMsg for the information. These messages are sent by the other driver in the event of a request or change of situation.

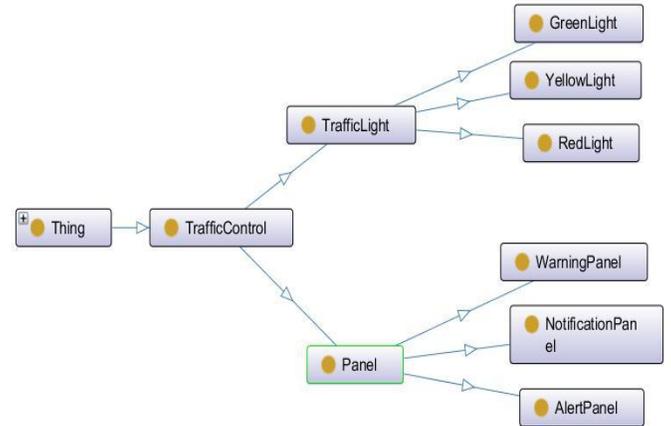


Fig. 6 Description of Traffic Control.

The class Traffic Control “Fig. 6”, which include the Panel of Traffic Control and Traffic Light to provide important information as a message to help drivers to respect traffic law.

V. CONCLUSION AND PERSPECTIVES

Traffic management is the most critical problems in urban areas. Advanced techniques and methods as VANET have the potential to solve this problem, but the interpretation by the driver does not reach the desired goal. Our solution is to integrate ontology in vehicles, to facilitate the interpretation of the information collected by the driver. Our approach will also allow the implementation of traffic management solutions more efficient and reliable.

Currently our ontology does not contain sufficient concepts for complex scenarios, our ontology can be extended so that it includes several concepts. Thereafter we propose to build a platform for validation of our approach, which could show the effectiveness of our approach, which provides a way to improve safety and traffic fluidity

REFERENCES

- [1] H. Hartenstein, K. P. Laberteaux, VANET Vehicular Applications and Inter-Networking Technologies, John Wiley & Sons, Ltd, 2009
- [2] J Kargl F., “Vehicular Communications and VANETs,” Talks 23rd Chaos Communication Congress, 2006.
- [3] Mahlknecht S. and Madani S., “On Architecture of Low Power Wireless Sensor Networks for Container Tracking and Monitoring Applications,” in Proceedings of 5th IEEE International Conference Industrial Informatics, 2007,pp. 353-358.
- [4] Yu, Fan F. and Biswas, S. A Self-Organizing MAC Protocol for DSRC based Vehicular Ad Hoc Networks, ICDCS Workshops 2007.

- [5] Cseh C., "Architecture of the Dedicated Short-Range Communications (DSRC) Protocol," in proceedings of IEEE Vehicular Technology Conference, 1998, pp. 45-49.
- [6] W. Whyte, "Safe at Any Speed: Dedicated Short Range Communications (DSRC) and On-road Safety and Security," RSA Conference 2005
- [7] Vehicle Safety Communications Consortium, "Vehicle Safety Communications Project Task 3 Final Report, Identify Intelligent Vehicle Safety Applications Enabled by DSRC," March 2005.
- [8] D. W. Franz, "Car-to-car Communication –Anwendungen und aktuelle Forschungsprogramme in Europa, USA und Japan", in Kon-gressband zum VDE-Kongress 2004 – Innovationen für Menschen, VDE, October 2004.
- [9] S. Eichler et al., "Strategies for context-adaptive message dissemination in vehicular ad hoc networks", in Proceedings of the 2nd International Workshop on Vehicle-to-Vehicle Communications (V2VCOM), July 2006.
- [10] R.Parker and S.Valaee, «Vehicle Localization in Vehicular Networks », Vehicular Technology Conference, 2006. VTC-2006 fall. IEEE 64th, pp. 1–5.
- [11] A.Benslimane, « Localization in vehicular Ad-hoc networks », Proceedings of the 2005 Systems Communications (ICW'05).
- [12] V. Kukshya, H. Krishnan, and C. Kellum, « Design of a system solution for relative positioning of vehicles using vehicle-to-vehicle radio communications during GPS outages », in Proceedings of IEEE Vehicular Technology Conference (VTC), September 2005.
- [13] A.Boukerche, H.A.B.F.Oliveira and E.F.Nakamura, «Vehicular Ad Hoc Networks: A New Challenge for Localization-Based Systems». Computer Communications, Vol. 31, No. 12. (30 July 2008), pp. 2838-2849.
- [14] B. Hofmann-Wellenho, H. Lichtenegger, J. Collins, Global Positioning System: Theory and Practice, 4th ed., Springer-Verlag, 1997.
- [15] E.D. Kaplan, Understanding GPS: Principles and Applications, Artech House, 1996.
- [16] M. Chen, D. Haehnel, J. Hightower, T. Sohn, A. LaMarca, I. Smith, D. Chmelev, J. Hughes, F. Potter, Practical metropolitan-scale positioning for GSM phones, in: Proceedings of 8th Ubicomp, Orange County, California, 2006, pp. 225–242.
- [17] Marfia, G.; Rocchetti, M., "Vehicular Congestion Detection and Short-Term Forecasting: A New Model With Results," Vehicular Technology, IEEE Transactions on, vol.60, no.7, pp.2936,2948, Sept. 2011.
- [18] R. Bauza, J. Gozalvez, and J. Sanchez-Soriano, "Road traffic congestion detection through cooperative vehicle-to-vehicle communications," IEEE 35th Conference on Local Computer Networks (LCN), 2010.
- [19] I. Leontiadis, G. Marfia, D. Mack, G. Pau, C. Mascolo, and M. Gerla, "On the effectiveness of an opportunistic traffic management system for vehicular networks", IEEE Transactions on Intelligent Transportation Systems, Vol. 12, Issue 4, 2011.
- [20] A. Lakas and M. Cheqfah, "Detection and dissipation of road traffic congestion using vehicular communication," Mediterranean Microwave Symposium (MMS), 2009.
- [21] Zhao, Y. & Tian, Z., An Overview of the Usage of Adaptive Signal Control System. Applied Mechanics and Materials Volumes 178-181, pp. 2591-2598. 2012
- [22] Maslekar, N., Boussedjra, M., Mouzna, J., Labiod, H., "VANET based Adaptive Traffic Signal Control", IEEE 73rd Vehicular Technology Conference (VTC Spring), pp. 1-5, 2011.
- [23] Gradinescu, V., Gorgorin, C., Diaconescu, R., Cristea, V., Iftode, L., "Adaptive Traffic Light Using Car-to-Car communications", IEEE 65th Vehicular Technology Conference (VTC Spring), pp. 21-25, 2007.
- [24] Barba, C.T.; Mateos, M.A.; Soto, P.R.; Mezher, A.M.; Igartua, M.A., "Smart city for VANETs using warning messages, traffic statistics and intelligent traffic lights," Intelligent Vehicles Symposium (IV), 2012 IEEE, vol., no., pp.902,907, 3-7 June 2012
- [25] Gruber T. "A Translation Approach to Portable Ontology Specifications", Knowledge Acquisition, pp. 199-220, 1993.
- [26] J. Barrachina, P. Garrido, M. Fogue, F. Martinez, J.-C. Cano, C. T. Calafate, and P. Manzoni, "Veacon: A vehicular accident ontology designed to improve safety on the roads," Journal of Network and Computer Applications, vol. 35, no. 6, pp. 1891–1900, 2012.
- [27] Sérgio Gorender, Ícaro Silva; "AN ONTOLOGY FOR A FAULT TOLERANT TRAFFIC INFORMATION SYSTEM", 22nd International Congress of Mechanical Engineering (COBEM 2013), November 3-7, 2013, Ribeirão Preto, SP, Brazil.
- [28] A. Groza, A. Marginean, B. Iancu, Towards improving situation awareness during emergency transportation through Ambulance-2-X communication and semantic stream reasoning, MEDITECH2014, Cluj-Napoca, Romania, 5-7 June 2014.
- [29] J. Barrachina, P. Garrido, M. Fogue, F. J. Martinez, J.-C. Cano, C. T. Calafate, and P. Manzoni, "Caova: A car accident ontology for vanets," in Wireless Communications and Networking Conference (WCNC), 2012 IEEE. IEEE, 2012, pp. 1864–1869.
- [30] Morignot, P.; Nashashibi, F., "An ontology-based approach to relax traffic regulation for autonomous vehicle assistance," CoRR, vol. abs/1212.0768, 2012.
- [31] Pollard, E.; Morignot, P.; Nashashibi, F., "An ontology-based model to determine the automation level of an automated vehicle for co-driving," In 16th International Conference on Information Fusion, Istanbul, Turquie, July 2013.
- [32] Protégé <http://protege.stanford.edu/>

Integrated Visual-perception Real-time Monitoring System

Jian-Wei Li, Yang,Fu-Syuan, Yi-Chun Chang*, and Yen-Lun Chiu

Abstract—People may access information through mobile devices and Internet anytime and anywhere because of development of modern technologies and networks; more and more multimedia applications as well as services characteristic of low latency and high bandwidth are expectably available to people in the future via mobile devices based on the 4G network matured gradually. In addition, the business operators in prosperous manufacturing and logistics industries developed with improved living standards have to control their vehicle fleets for logistics but still face some problems deriving from the existing fleet management system. Against this background, we installed the user end on the mobile device to effectively manage vehicles under surveillance through traffic cameras and process latest information collected from vehicles for post tracking conveniently.

Keywords—3G/4G network; fleet management system; traffic camera.

I. INTRODUCTION

THE Internet-based activities such as purchasing household articles and books or placing orders with manufacturers have been available to people on the strength of the progress and applications of science and technology and prompted service suppliers in manufacturing and logistics industries which serve buyers merchandise fast through freight transport or express delivery in recent years. In this regard, the numerous vehicles owned by a logistics service provider should be controlled with a fleet management system for goods delivered to buyers punctually and accurately. Moreover, a fleet management system is critical to a business operator responsible for transportation of valuables.

Based on GPS (Global Positioning System), a fleet management system should be competent to tracking positions and conditions of a vehicle under surveillance, recording relevant information such as route and dwell time in a back-end system [1][2], and allowing the driver or the car occupant to inform the surveillant of any accident en route and further the surveillant to deploy another vehicle for transportation of cargos punctually and effectively. Currently, the fleet management system is common in the mass transit

system and taxi fleets in addition to the logistics industry and facilitates current positions of a vehicle and waiting time learned by a consumer.

Despite convenience, the existing vehicle management system still derives some problems from time to time, for example, a surveillant depends on a report from the driver or the car occupant in a monitored vehicle to learn an accident or fails to hear about an accident timely in the case of no report from the driver or the car occupant forthwith or the monitored vehicle's position replied by GPS different to its actual position. It can be seen from above descriptions the surveillant plays a passive role during an accident ongoing because of the existing vehicle management system which is based on GPS only for surveillance and makes no contribution to service quality and safety of a transportation service provider.

In this study, we create and install a visual-perception vehicle management system with traffic cameras and GPS integrated wherein the traffic cameras as eyes of a surveillant are able to transmit live images by which a surveillant is informed of a monitored vehicle exactly traveling on roads corresponding to the GPS system. In addition, a surveillant may decisively arrange a rescue action by observing real-time conditions of a monitored vehicle from video cameras adjacent to an accident. The visual-perception vehicle management system developed in the Android system facilitates applications at distinct positions and networks by a surveillant through mobile devices and is integrated with Google Map and RTSP (Real Time Streaming Protocol) [3][7] for visual-perception real-time monitoring, as shown in Figure 1 for its system.

The content of the paper is further explained in the following sections: Section II is background knowledge related to development of a visual-perception real-time monitoring system; Section III presents architecture of a visual-perception real-time monitoring system; Section IV presents the method to effectuate a visual-perception real-time monitoring system and flowcharts; Section V is practical implement; Section 6 is conclusions.

II. BACKGROUND KNOWLEDGE

In this study, a user depends on the Android-based system installed in a mobile device to immediately monitor vehicles, a back-end server to record monitoring information, Google Map to display traffic conditions, and RTSP to transmit real-time information from video cameras and familiarizes

The research is supported by the Ministry of Science and Technology of the Republic of China under the grant number MOST 103-2221-E-324 -012 and MOST 103-2221-E-241 -008.

Jian-Wei Li, Yang,Fu-Syuan, and Yen-Lun Chiu are with the Dept. of Information and Communication Engineering, Chaoyang University of Technology, Taiwan (R.O.C).

Yi-Chun Chang is with the Dept. of Computer Science and Information Engineering, Hungkuang University, Taiwan (R.O.C) (corresponding author to provide e-mail: changyc@livemail.tw).

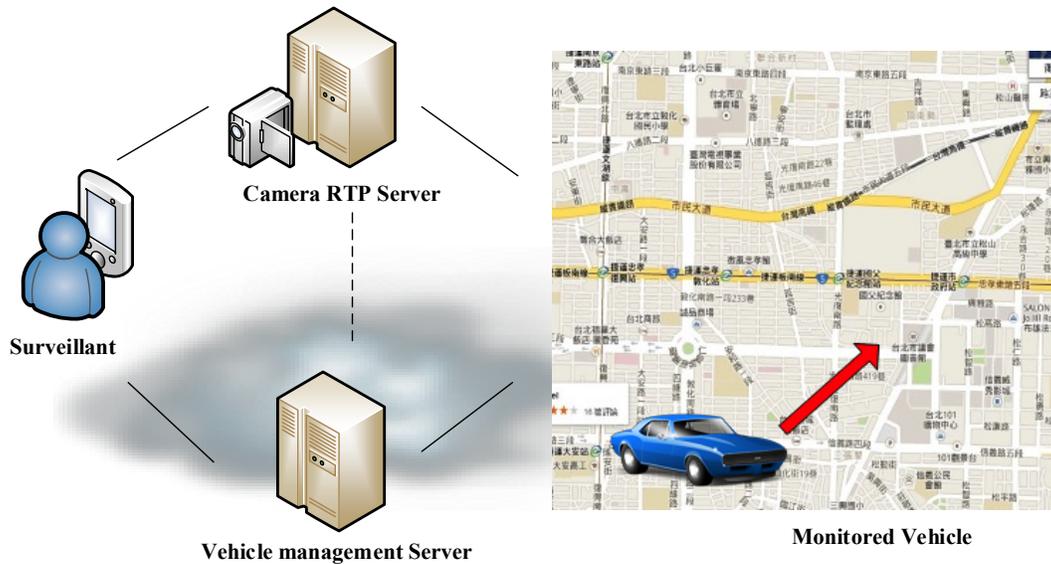


Figure 1. Schematic view of a visual-perception real-time monitoring system

him/her with knowledge related to the existing vehicle management system. For this matter, Android, Google Map, RTSP and the vehicle management system are introduced hereinafter.

A. Android

Android is architected in the form of a software stack for mobile devices, which is developed based on Linux kernel. The main layers in android software stack depicted in Figure 2 as follows [13]:

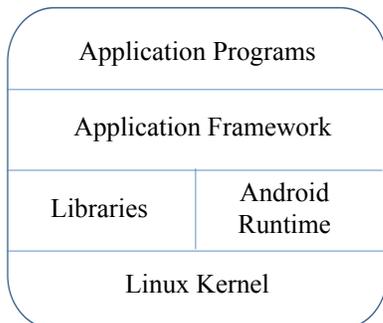


Figure 2. Android software stack.

- (1) Linux Kernel: It provides an abstraction layer between hardware and the rest of the stack, which is responsible for memory management, process management, network stack, device driver, and so on.
- (2) Libraries: Android includes a set of C/C++ libraries, such as libc (System C library), SQL lite, Graphics libraries OpenGL, used by various components.
- (3) Android Runtime: This layer includes core libraries and the Dalvik Virtual Machine.
- (4) Application Framework: This layer provides the Application programming interfaces (API). Developers make use of these APIs to develop their own applications.

- (5) Application Program: All Android applications are built on the application framework using the same API.

B. Google map

Google maps is a cross platform web mapping service provided by Google, which provides online maps, topographic maps and satellite imagery, and other services such as global positioning search, traffic information query, traffic route, and street view. Based on JavaScript technology, Google Map launched API (Application Programming Interface) in 2005 to allow developers to integrate Google maps into their website. After that, Google released Google maps for mobile for applications run on any Java-based phone or mobile devices. In addition to these applications, other categorizes, e.g., map, marker, longitude/latitude, poly line and coordinate, are included in Google Map and used in presenting corresponding information on the online map and developing other applications by programmers as required [4][5][6][14].

C. RTSP

RTSP, which is defined in RFC 2326, is a text-based protocol for multimedia streaming used to remotely control audio/video media, create a link between a user and a media server through RTSP information, describe resources, and own rights to control media [6][7]. As a protocol of multimedia transmission, RTSP cannot be employed independently and should be coordinated with RTP (Real Time Transport Protocol) [9].

For the process of signal transmission, a RTSP client issues request instructions including control parameters such as SETUP, PLAY, STOP, TEARDOWN [8] to a RTSP server. As shown in Figure 3, a RTSP session typically is initialized with SETUP message, and starts the stream with PLAY message. Finally, the session is closed with TEARDOWN

message.

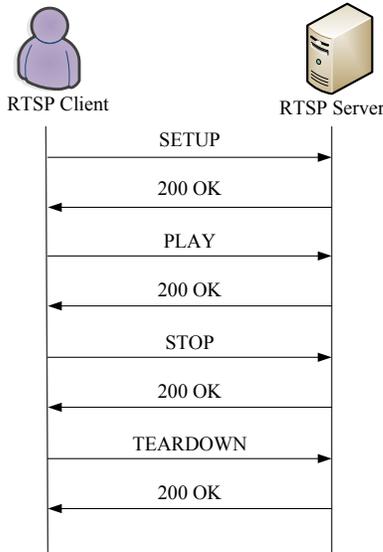


Figure 3. Process of transmitting RTSP signals [7]

In a RTSP SETUP example message shown in Figure 4, the SETUP request sent by a client for a URI (Uniform Resource Identifier) specifies the transport mechanism for the stream media. Then, the server generates and responds session identifiers for the SETUP request [7].

```

C->S: SETUP rtsp://example.com/foo/bar/baz.rm RTSP/1.0
      CSeq: 302
      Transport: RTP/AVP;unicast;client_port=4588-4589

S->C: RTSP/1.0 200 OK
      CSeq: 302
      Date: 23 Jan 1997 15:35:06 GMT
      Session: 47112344
      Transport: RTP/AVP;unicast;
              client_port=4588-4589;server_port=6256-6257
    
```

Figure 4. RTSP SETUP example from RFC 2326.

D. Vehicle management system

The vehicle management system is able to identify

information from single service devices and automatically update positions of moving vehicles into a corresponding database for management of vehicle conditions [10] such as vehicle maintenance, remote processing of vehicles (diagnosis and maintenance) and safety management [11]. Moreover, the vehicle management system emphasizing accurate temporal synchronization relies on PLL (Phase-locked loop) [12] for timekeeping in order to realize precise positioning based on GPS and process information effectively.

III. SYSTEM ARCHITECTURE

Figure 5 presents the system architecture in this study including user end, monitoring client, RTSP streaming server, and visual-perception administrator server. The user end represents a surveillant who monitors objects under surveillance via a mobile device on which Google Map indicates the monitored objects and clicks on the icon of camera on the map to watch live images of the monitored vehicles; the monitoring client includes any object under surveillance and functions simply for transmitting current GPS information via an application.

The visual-perception administrator server runs three processes as follows: (1) Monitoring Client Data presents information of objects under surveillance such as returned GPS coordinates, traffic routes of a vehicle previously recorded and time of a vehicle passing through video cameras; (2) User Data presents information of a surveillant such as user ID and vehicles under surveillance; (3) Camera List presents a list for traffic cameras which corresponds to all video cameras used to monitor passing vehicles.

As traffic cameras, the video cameras provides a surveillant with real-time monitoring images or saves recorded video images in the RTSP streaming server in which two steps are existed: (1) Saving Function to save video images from traffic cameras; (2) Stream Function to transmit saved multimedia files to a surveillant if necessary.

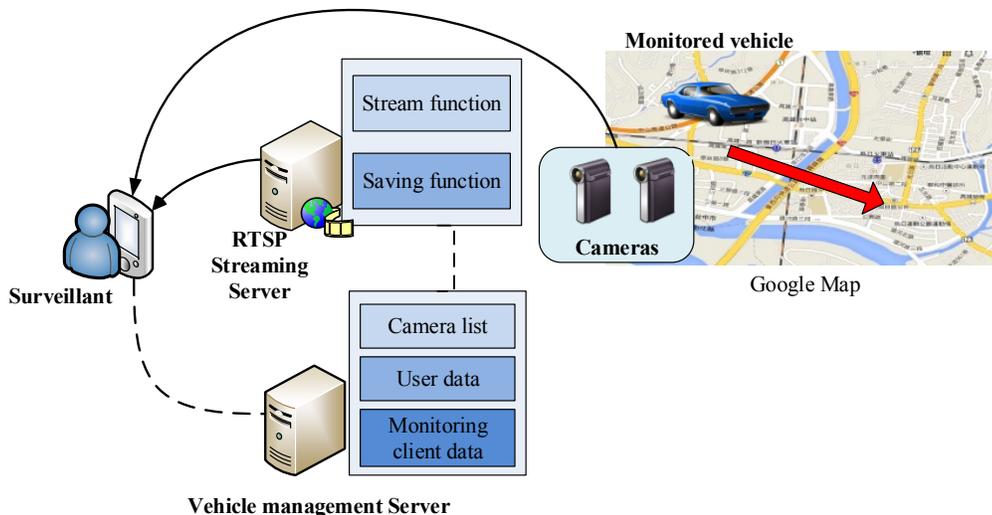


Figure 5. Architecture for the visual-perception monitoring system

IV. FUNCTIONS TO EFFECTUATE VISUAL-PERCEPTION

The Android client environment is divided into two parts, user end and monitoring client. Both clients work with the visual-perception administrator server, which are described hereinafter:

A. Monitoring client

The monitoring client should upload and save GPS information to the visual-perception administrator server from which a surveillant can access information. The two functions for the monitoring client are described here:

- Transmission of GPS coordinates: GPS coordinates and vehicle speeds are transmitted to and saved in the back-end administrator server from a mobile phone.
- S.O.S. signal: The function to proactively call the user end or indicate markers contributes to fast identifying an accident confronted by the monitoring client.

B. User end

A surveillant accesses necessary information displayed on Google Map in a personal mobile device via the link to a server. The user end contains applications as follows:

- Select a vehicle to be monitored: A surveillant accesses information from the back-end administrator server to ascertain conditions displayed on the online map with a vehicle selected.
- Display basic function: The basic function contributes to decisions of a surveillant for a vehicle's conditions displayed.
- Enable status processing – S.O.S. long: Alarms along with video images from corresponding cameras will be enabled for the real-time conditions of a vehicle which stops traveling or issues an S.O.S. signal too long.
- Access a vehicle's traffic routes over a period of time: Information such as a vehicle's traffic routes over a period of time can be displayed on the online map from the back-end server according to time entered by a surveillant.
- Access video cameras passed by a vehicle over a period of time: Based on a video camera to draw a circle, any video image recorded by the camera over a period of time can be checked from the back-end server.
- Enable/disable monitoring conditions: A surveillant decides to monitor a vehicle at his/her disposal.
- Change another vehicle to be monitored: A surveillant can change another vehicle to be monitored at his/her disposal.

V. PRACTICAL IMPLEMENTATION

The practical implementation is based on the Android system and divided into the user end and the monitoring client. We employ JAVA language with eclipse IDE (Integrated Development Environment) for JAVA developers for program development under Android environment. Involved devices

contain (i) android tablets used as mobile devices of the end user and monitoring client, and (ii) IP cameras used as traffic cameras. Besides, we create the visual-perception administrator server based on Apache, MySQL, and RTSP server. Functional operations of the user end and the monitoring client are demonstrated.

A. APP activated at the monitoring client

The configurations at the monitoring client are simple. In this regard, the plate number of a vehicle to be monitored should be entered at the monitoring client for activating the GPS function corresponding to the user end. Then, the GPS information will be uploaded to and recorded in the administrator server from the vehicle-borne device at the monitoring client, as shown in Figure 13 and Figure 7.

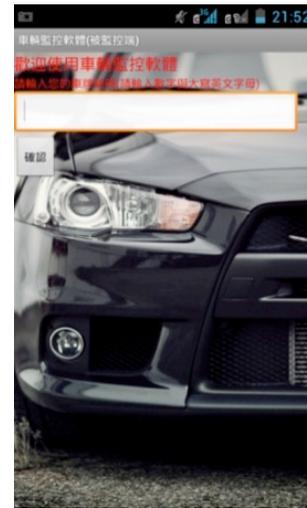


Figure 6. Image for a plate number entered at the monitoring client



Figure 7. GPS information automatically uploaded from the monitoring client

B. APP activated at the user end

When more one vehicles are monitored by the visual-perception administrator server, a surveillant logging in

the system at the user end can decide a vehicle to be monitored by which the information related to the vehicle is accessed from the administrator server for ascertaining conditions and displaying them on Google Map, as shown in Figure 8. With a vehicle decided, the vehicle's current conditions/positions and traffic cameras nearby will be displayed on the upper left corner of the monitoring interface via a mobile device of the surveillant, as shown in Figure 9.



Figure 8. The vehicle to be monitored is selected by a surveillant.

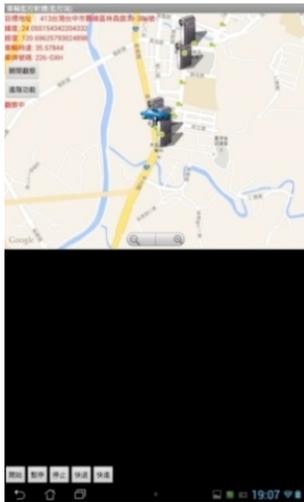


Figure 9. The user end interface is activated with the vehicle decided.

The conditions of a vehicle may be one of two possibilities: the vehicle stopped or issuing an S.O.S. signal. As shown in Figure 10, a monitored vehicle stopped is indicated as the speed equal to 0 on the user end interface and marked as the vehicle stopped at the bottom of the interface. In case of an accident, the S.O.S. signal can be issued by the driver or the car occupant through APP at the monitoring client and indicated on the user end interface, as shown in Figure 11. Besides, the user end interface displays the positions of traffic cameras nearby, as shown in Figure 12. With the traffic camera at the traffic accident identified, a surveillant is able to click on the camera for transmitting and displaying live images at the bottom of the user end interface, as shown in

Figure 13.



Figure 10. The monitored vehicle which stops traveling is displayed on the user end interface.



Figure 11. S.O.S. signal from the monitoring client displayed on the user end interface



Figure 12. Positions of video cameras around an accident vehicle



Figure 13. Live images at the scene displayed on the user end interface

In addition to the function of real-time monitoring, the other function of post tracking is designed in this study. As shown in Figure 14, when a surveillant specifies a vehicle and enters the time interval to be checked, the traffic routes of a monitored vehicle saved in the administrator server along with video cameras passed by the vehicle can be accessed and displayed on the user end interface. Each video camera icon presents the specified video segment that recorded the specified vehicle. Furthermore, the video camera icons clicked by a surveillant are able to present video segment of the specified vehicle traveling then.



Figure 14. Traffic routes of a monitored vehicle displayed in post tracking

VI. CONCLUSION

The transmission of real-time images on the 4G network with the wider bandwidth is progressively robust in the burgeoning Internet era. Thus, a surveillant is able to ascertain regular traffic routes of monitored vehicles and is informed of actual conditions at the scene for rapid rescue with mobile devices and traffic cameras integrated in a real-time

monitoring system.

In this study, a visual-perception real-time monitoring system developed in the Android system and featuring GPS, IP Camera, Google Map, RTSP and MySQL integrated allows GPS tracks of monitored vehicles and traffic cameras to link each other, contributing to real-time surveillance as well as post tracking for restoring initial conditions based on saved images and corresponding GPS positions and time.

REFERENCES

- [1] R. Zantout, M. Jrab, L. Hamandi, & F.N. Sibai, "Fleet Management Automation Using the Global Positioning System", IIT '09. International Conference on Innovations in Information Technology, pp.30-34, 15-17 Dec. 2009.
- [2] S. T. S. Thong, Tien Han Chua & T. A. Rahman, "Intelligent Fleet Management System with Concurrent GPS & GSM", The 7th International Conference on Telecommunications, pp.1-6, 6-8 June 2007.
- [3] Yan Liu, Guo-Hui Zhong, Yu Liu, Hua-Qiang He & Fu-Rong Wang, "The research of streaming media mutual digest authentication model based on RTSP protocol", International Conference on Wavelet Analysis and Pattern Recognition, vol.2, pp.838-842, 30-31 Aug. 2008.
- [4] Hao Zhang, Manchun Li, Zhenjie Chen, Zhiliang Bao, Qiuhao Huang & Dong Cai, "Land Use Information Release System Based on Google Maps API and XML", 2010 18th International Conference on Geoinformatics, pp.1-4, 18-20 June 2010.
- [5] He Li & Zhijian Lai, "The study and implementation of mobile GPS navigation system based on Google Maps", 2010 International Conference on Computer and Information Application (ICCIA), pp.87-90, 3-5 Dec. 2010.
- [6] Juan Li, Xie Yong & Zhang-Yi Lai, "Real-time monitoring system of water resources based on Google Maps", 2012 International Conference on Information Management, Innovation Management and Industrial Engineering (ICIIE), vol.2, pp.239-242, 20-21 Oct. 2012.
- [7] H. Schulzrinne, A. Rao, R. Lanphier, "Real Time Streaming Protocol (RTSP)", RFC 2326, IETF, April 1998.
- [8] Jee Changwoo, K.G. Shin, "A DAVIC Video-on-Demand System Based on the RTSP", 2001. Proceedings. 2001 Symposium on Applications and the Internet, pp.231-238, 2001.
- [9] H. Schulzrinne, S. Casner, R. Frederick & V. Jacobson, "RTP: A Transport Protocol for Real-Time Applications", RFC 3550, IETF, July 2003.
- [10] A. Aljaafreh, M. Khaleel, I. Al-Fraheed, K. Almarahleh, R. Al-Shwaabkeh, S. Al-Etawi & W. Shaqareen, "Vehicular data acquisition system for fleet management automation", 2011 IEEE International Conference on Vehicular Electronics and Safety (ICVES), pp.130-133, 10-12 July 2011.
- [11] M. Cao, M. Alvarez-Campana, R.P. Leal, J. Navarro, E. Polo, "IMS-Based Testbed for Fleet Management", 2010 Fourth International Conference on Next Generation Mobile Applications, Services and Technologies (NGMAST), pp.67-72, 27-29 July 2010.
- [12] G.S. Singh, D. Singh, S. Moorthi, "Low power low jitter phase locked loop for high speed clock generation", 2012 Asia Pacific Conference on Postgraduate Research in Microelectronics and Electronics (PrimeAsia), pp.192, 196, 5-7 Dec. 2012.
- [13] Android Open Source Project, <http://developer.android.com/reference/android/os/PowerManager.html>.
- [14] Google Maps API <https://developers.google.com>

Novel M-ary PPM time hopping scheme for UWB communications

Said GHENDIR, Salim SBAA, Riadh AJGOU, Ali CHEMSA and A. TALEB-AHMED

Abstract—In this paper, we evaluate the performance of Time Hopping-Ultra-Wide Band (TH-UWB) communications, employing M-ary Pulse Position Modulation (PPM) in terms of Bit Error Rate (BER) over a multipath channel model, CM1 through CM4 adopted by IEEE802.15.3a standard. In this work, we propose a novel M-ary Pulse Position Modulation Time Hopping (M-ary PPM TH) scheme, robust to multipath fading, a good means to increase the data rate and improve the BER performances. This scheme based on keeping the chip/frame duration and increasing the order of M-ary PPM modulation. The proposed scheme, leads to increase the bit rate using positively the modulation order. The results have been verified by means of computer simulations proving the efficacy of this scheme.

Keywords—M-ary PPM, S-Rake, TH, UWB.

I. INTRODUCTION

ULTRA-Wide Band is a radio technology that can be used for short-range communication systems with low power consumption at short ranges [1]. The fundamental characteristic of UWB signal uses pulses baseband a very short period of time of the order of a few hundred picoseconds [2]. One of the most widely studied schemes for UWB communications employs pulse position modulation (PPM) combined with time hopping (TH). The FCC (Federal Communications Commission) allows a maximum bandwidth of 7.5 GHz, from 3.1 – 10.6 GHz, for UWB communication. The generated UWB signals are able to communicate by baseband short pulses. Thereafter, the technique was immediately applied to radar, communications, automobile collision avoidance, positioning systems, and liquid-level sensing.

The conventional works for M-ary PPM transmission already use a scheme where the frame duration is variable and increases when M increases. This scheme has some known advantages, but it has a major disadvantages such as in higher modulation order the frame duration increases while bit duration remains constant which means constant data rate for

2PPM or 4PPM.

Widely adopted in UWB systems today and became challenging to implement at high rates with large M.

This paper, an extension of our work presented in [3]; where, we have done the channel estimation by considering the maximum likelihood approach via a data aided method, and using a Selective Rake (S-Rake) receiver. However, this work differentiates itself from other previous works by introducing a novel M-ary PPM TH scheme. This scheme is a new idea that has not been implemented before, where the development of this work will be given as an implementation by means of computer simulations, and we will take this scheme in detail in extended works.

The purpose of this paper is to propose a novel scheme in order to increase the bit rate using M-ary PPM in a fixed chip/frame duration. The UWB pulses are time hopped within a fixed time slot (chip/frame) and each transmitted symbol is spread over several pulses in order to facilitate multiple users. Therefore, this scheme can be considered as a means of transmitting multiple bits per symbol in an intensity modulated system to reach the Gbits/s.

Motivated by the above, we will describe Time Hopping UWB System Model with our proposed scheme finishing by conclusion and some remarks.

II. PROPOSED TIME HOPPING UWB SYSTEM SCHEME

The Time Hopping UWB system modulated with M-ary pulse position modulation (PPM) for the n^{th} user is expressed by

$$s^{(n)}(t) = \sum_{k=-\infty}^{+\infty} b(t - kN_f T_f - \delta_{(M)} a_k^{(n)}) \quad (1)$$

Where, one bit is

$$b(t) = \sum_{j=0}^{N_f-1} \sqrt{E_b^{(n)}} w(t - jT_f - c_j^{(n)} T_c) \quad (2)$$

And E_b is the received energy per pulse from user n ($1 < n < N_u$), assuming that users transmit asynchronously, $w(t)$ is the pulse shape that nominally begins at time zero on the transmitter clock; and T_f is the frame duration, i.e., N_f is a number of frames in one symbol. The sequence $c_j^{(n)}$ is the time hopping code associated at desired user in j^{th} frame and it is between $\{0, N_c - 1\}$; N_c is a number of chips in one frame; the parameter T_c is the duration of the chip.

Said GHENDIR, Salim SBAA, Riadh AJGOU and Ali CHEMSA Authors are with, LIBCUB Laboratory, Electric engineering department, University of Biskra. B.P 145 R.P, 07000 Biskra ALGERIA. And are also with Department of Sciences and Technology, El-Oued University, PO Box 789 39000 El-Oued, ALGERIA. (Email: said-ghendir@univ-eloued.dz, s.sbaa@univ-biskra.dz, riadh-ajgou@univ-eloued.dz, chemsadoct@yahoo.fr)

A. TALEB-AHMED Author is with LAMIH Laboratory University UVHC Mont Houy - 59313 Valenciennes Cedex 9 FRANCE. (Email: abdelmalik.taleb-ahmed@univ-valenciennes.fr)

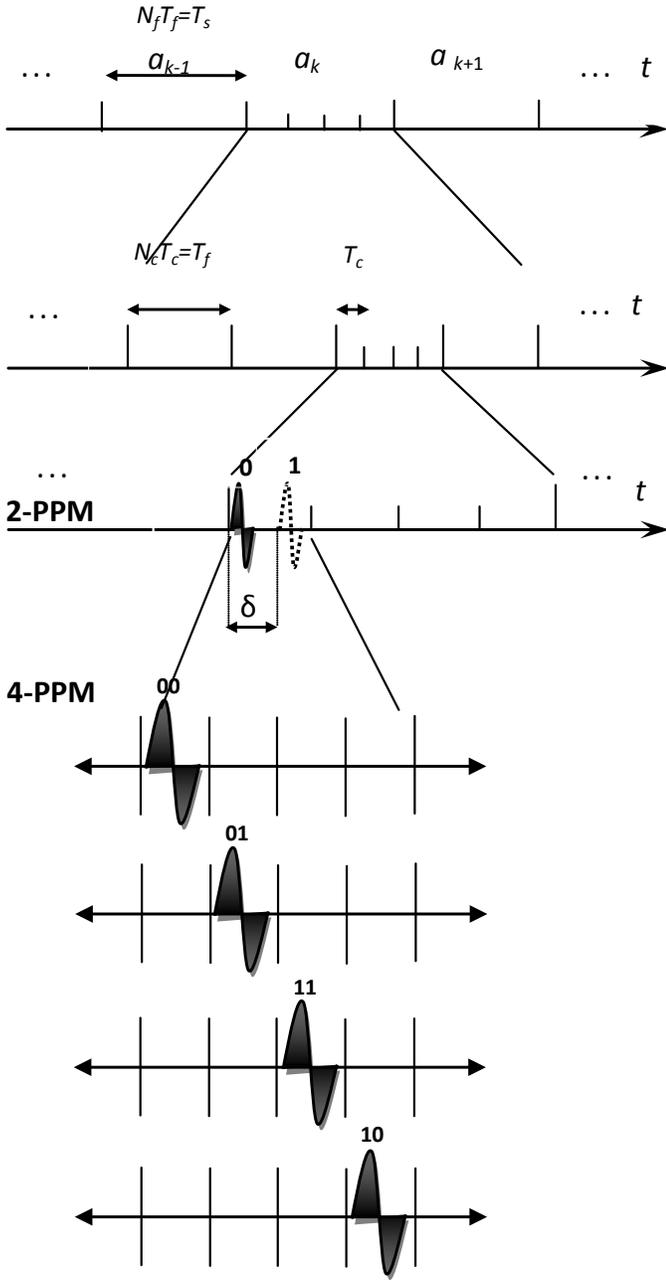


Fig.1 Structure of the proposed M-ary PPM TH-UWB scheme

From equation (1), $a_k^{(n)}$ is the transmitted symbol that has a duration T_s . An evenly spaced M-PPM scheme is considered, where δ is the time shift difference between two subsequent PPM signals. Correspondingly, $\delta_{(M)}$ is a proposed modulation factor for M-ary PPM in kept chip/frame, and finally the parameter T_w is the monocycle duration.

We give in Fig. 1 our proposed scheme for different M-ary PPM modulations ($M = 2^b$), knowing that b is the number of bits and M can take the following values (2, 4, 8, 16 ... etc.).

With reference to equations (1) and (2) the following assumptions have been made. We have chosen the monocycle shape as the first derivative of a gaussian impulse as in Fig. 2.

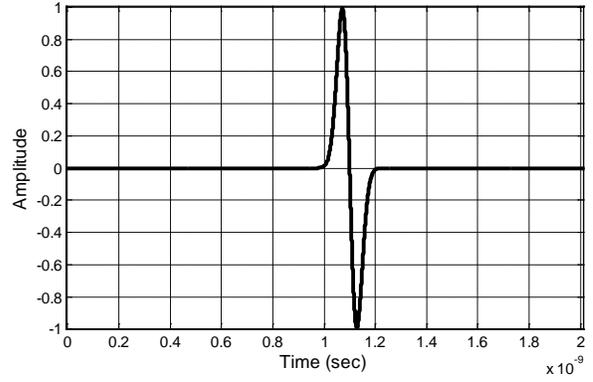


Fig.2 Shape of the monocycle

And designed by the following relation:

$$w(t) = \left(\frac{t - T_w/2}{T_w} \right) \exp \left[- \left(\frac{t - T_w/2}{T_w} \right)^2 \right] \quad (3)$$

On the other hand, four different channel models were defined by IEEE802.15.3a standard, namely CM1, CM2, CM3, and CM4. CM1 describes a LOS (line-of-sight) scenario with a separation between transmitter and receiver of less than 4m. CM2 describes the same range, but for a No-LOS situation. CM3 describes a No-LOS scenario for distances between transmitter and receiver 4-10m. Scenario 4 finally describes an environment with strong delay dispersion.

The UWB multipath channel can be expressed as

$$h(t) = X \sum_{l=0}^{L_c} \alpha_l \delta(t - \tau_l) \quad (4)$$

- L_c : number of paths
- X : represents the log-normal shadowing
- α_l : are the multipath gain coefficients
- τ_l : is the delay of the l^{th} multipath component

Several time-hopping signals are simultaneously transmitted over a channel with L_c paths, the composite waveform at the output of the receiver antenna may be written as:

$$r(t) = \sum_{l=0}^{L_c} \alpha_l s(t - \tau_l) + v(t) + u(t) \quad (5)$$

In this equation, $r(t)$ is the desired user's signal, α_l and τ_l are the attenuation, and the delay affecting its replica traveling through the l^{th} path, $v(t)$ is thermal noise and $u(t)$ represents the MAI caused by the other users.

The optimum detection strategy for this multiple access system leads to a multiuser receiver which, however, is too complex to implement. More feasible schemes are of interest. The simplest suboptimal receiver is obtained making two

approximations. First, the MAI is thought of as a white Gaussian process [4] and as such, it can be lumped into the thermal noise in (5). The Gaussian approximation is justified by the central limit theorem if the users are many and have comparable powers. Second, a dominant path exists that conveys the major part of the desired user's energy. Under these assumptions (5) becomes:

$$r(t) = \sum_{l=0}^{L_c} \alpha_l s(t - \tau_l) + n(t) \quad (6)$$

Where $n(t) = v(t) + u(t)$ is still a Gaussian and white process.

The receiver is considered perfectly synchronized to the desired user, so the output signal of the receiver can be given by:

$$r^{(n)}(t) = X \sqrt{E_b^{(n)}} \sum_{k=-\infty}^{+\infty} \sum_{j=-\infty}^{N_f-1} \sum_{l=0}^{L_c} \tilde{\alpha}_l \tilde{w}(t - kN_f T_f - jT_f - c_j^{(n)}(j)T_c - \delta_{(M)} a_k^{(n)}(i) - \tilde{\tau}_l) + n(t) \quad (7)$$

Where $\tilde{w}(t)$ is the received pulse (distorted from that emitted $w(t)$) and is affected by the UWB multipath channel.

The parameters $\tilde{\alpha}_l$ and $\tilde{\tau}_l$ (the amplitude and relative delay of the l^{th} component of the received signal) are not known a priori and must be estimated as in [3].

The block diagram of a Rake receiver with L_c fingers (Rake- L_c) that exploits L_c signal echoes, is depicted in Fig 3.

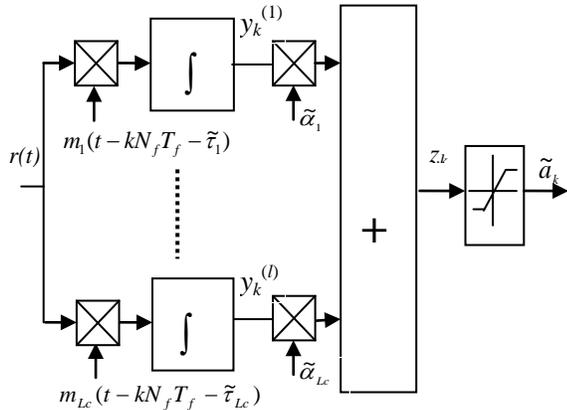


Fig. 3. Block diagram of a Rake- L_c receiver.

In addition, several options are also possible in the number of branches or fingers of the rake receiver, i.e the number of components included, that define the type of rake receiver.

The template signal $m(t)$ depends not only on the user's time hopping code, but also on $w(t)$. For purposes of analysis, we assume that the shape of the monocycle is known [5].

In order to detect TH-UWB signals through multipath diversity existing in the received signal, the S-Rake receiver with L_c fingers is used. Where, in our case the paths that are within 10 dB of the peak amplitude will be captured by fingers.

We used in the receiver, M matched filters with m_i template

functions to have an optimum demodulator based on a pulse correlation technique, for that reason, a proposed formula for template signals is used:

$$m_i(t) = b(t - i \delta_{(M)}) \quad (i=0,1,\dots,M-1) \quad (8)$$

Then, followed by an energy detector which determines the maximum signal output at the end of the symbol duration, using a comparison test of the M possible slots, where a decision is made based on which slot contains the most energy.

The received signal is integrated over a period corresponding to the duration T_s of an information symbol. The weighted sum of the correlators is then applied to a detector which determines the transmitted symbol \tilde{a}_k according to M-ary PPM modulation.

The decision statistic is the maximal-ratio combination of the outputs of the correlators and is given by the following expression [6]:

$$z_k = \sum_{l=1}^{L_c} \alpha_l \int_{kN_f T_f}^{(k+1)N_f T_f} r(t) m(t - kN_f T_f - \tau_l) dt \quad (9)$$

The Fig. 4 shows as an example, how a 2PPM modulated signal is demodulated for a single finger rake receiver.

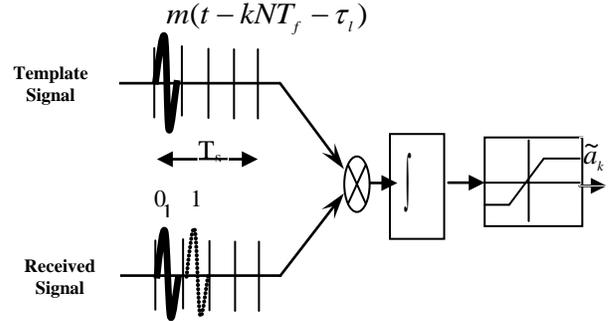


Fig. 4. Demodulation for 2PPM modulated signal

In this work, we established all functional blocks of the UWB communication link cited above, where we have included our proposed scheme in the transmission and reception in order to show the performance evaluation of our novel M-ary PPM time hopping scheme.

The block of the receiver shown in Fig. 5 is given for one finger to facilitate the view. After generating the UWB signal according to M-ary PPM TH scheme, a shaping filter is chained to designate the pulse shape. By passing the signal through the UWB multipath channel, it is important to note that because any delay of the order of nanoseconds in the time domain may cause an erroneous information. The UWB signals affected by the multipath are cumulated constructively or destructively causing different attenuations and distortions due to a different phase shifts. Furthermore, we get to the step of equalization and correction signals via the rake receiver, where this latter, uses a matched filter to estimate the channel. Then

performance performs badly.

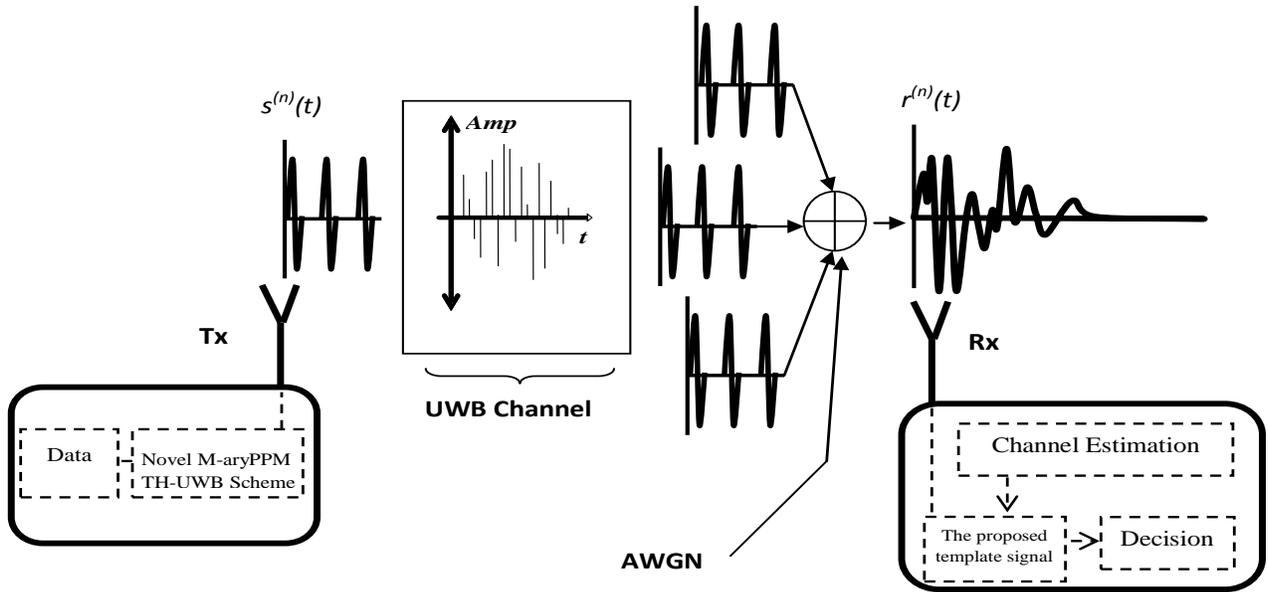


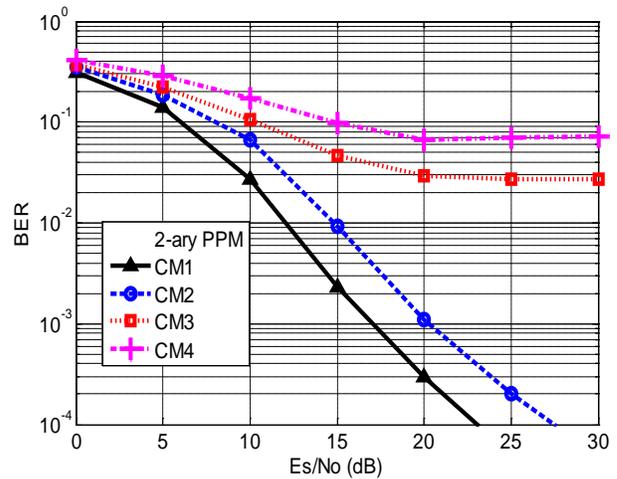
Fig. 5 Block diagram of the UWB communication link

uses this estimate for correlating the received signal with a predefined template signal according to M-ary PPM modulation used or the desired user. Also, the template signal uses the same THC (Time Hopping Code) used by the transceiver to detect the temporal position of the impulses transmitted and then integrate over the duration of one symbol T_s . Finally, by applying the steps above on all the fingers of the rake receiver and using a threshold detector to distinguish the symbols received we obtained the simulation results below.

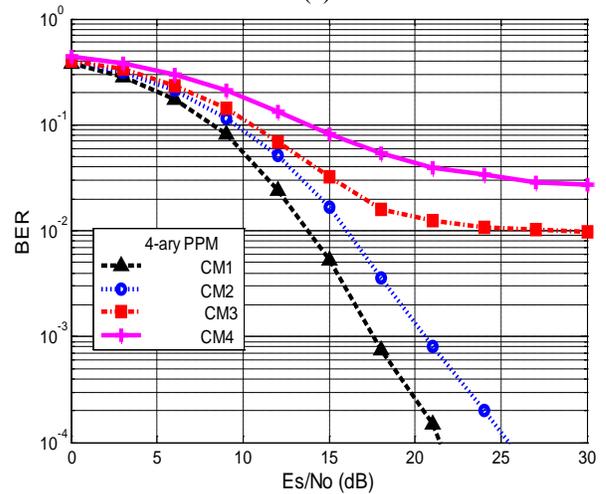
III. SIMULATION RESULTS

All results have been assessed using MATLAB simulations to verify our novel scheme with own suppositions and other parameters assumed in the literature. In fact, the sampling frequency is chosen $f_c=25\text{GHz}$, that means 25 samples per monocycle pulse according to [5] [7]. These results are given by considering one user. Where, the symbol duration $T_s=20\text{ns}$; the fixed chip/frame duration is respectively $T_f=10\text{ns}$, $T_c=2\text{ns}$; the number of frames $N_f=2$ per symbol, and each impulse occupies one frame; the number of chips in one frame is $N_c=5$, $c_j = \{0, 0\}$; the energy of bit E_b is normalized to "1"; the pulse duration $T_w=0.2\text{ns}$, and the modulation factor $\delta_{(M)}=0.2\text{ns}$. Finally, the results have been done over 100 impulse response realizations for each channel model (CM1, CM2, CM3, CM4) to reach the transparency results.

As can be seen from Fig. 6 (a), (b), we evaluated the BER versus E_s/N_0 (signal-to-noise ratio) for channel models (CM1, CM2, CM3, CM4) using 2-ary PPM, 4-ary PPM. It is clear that via the proposed scheme we get two major benefits which are the BER is going lower and the data rate is going higher as we go for higher order modulation. Furthermore, it is shown that when the order of the channel model increases the BER



(a)



(b)

Fig.6 BER evaluation versus E_s/N_0 , a) 2-ary PPM. b) 4-ary PPM

IV. CONCLUSION

The proposed scheme is a new idea that has not been implemented before, which leads unlike the previous literature to increase the data rate and improve at the same time the BER performance for M-ary PPM (e.g. M=2,4), where

2-ary PPM \rightarrow 50Mbps \rightarrow 1bit/Sym

4-ary PPM \rightarrow 100Mbps \rightarrow 2bit/Sym

Finally, we note that the results make in evidence that we can use M-ary PPM positively over our proposed scheme in order to enhance the performance of (TH-UWB) communications in high data rate

REFERENCES

- [1] Jawad, M.S., Ismail, W., Hajjawi, A., Rani, O.A., Hussain, A.-S.T. and Saleh, A. (2014), Review of the State of Art of Tunable Impulse Ultra-Wideband Technology as Integrator for Wireless Sensing and Identifications Short-Range Networks. *Wireless Sensor Network*, vol 6, pp. 137-156
- [2] A. Darif, R. Saadane, and D. Aboutajdine, Short-Range Networks. *Wireless Sensor Network*, 6, 137-156. IR-UWB: An Ultra Low Power Consumption Wireless communication Technologie for WSN, *TELKOMNIKA Indonesian Journal of Electrical Engineering*, Vol. 12, No. 8, August 2014, pp. 5699 ~ 5708
- [3] S. Ghendir, S. Sbaa, A. Ajjou, et al, High Bit Rate UWB Communication in Dense Multipath Channels, in *Proc. 18th International Conference on Communications (part of CSCC '14)*, Santorini Island, Greece, July 2014, pp. 74-79.
- [4] M. Z. Win and R. A. Scholtz, "Ultra-wide bandwidth time-hopping spread-spectrum impulse radio for wireless multiple-access communications," *IEEE Trans. Commun.*, vol. 48, pp. 679-691, Apr. 2000.
- [5] V. Lottici, A. D'Andrea, and U. Mengali, Channel Estimation for Ultra-wideband Communications, *IEEE J. Select. Area. Commun*, Vol.20, No.9, Dec. 2002, pp. 1638-1645.
- [6] G. L. Turin, Introduction to Spread-Spectrum anti Multipath Techniques and Their Application to Urban Digital Radio, *Proc. IEEE*, Vol.68, Mar. 1980, pp. 328-353.
- [7] A. Deleuze, Contributions à l'étude des systèmes ultra large bande par impulsions, *PhD thesis, the Superior National School of Télécommunications*, Paris 2006.

Interoperability for an observatory of habits and healthy life styles related with physical activity

Andrea Torres Ruiz, Fernando Prieto B., Jose Arturo Lagos, Nixon Duarte, Rosmary Martinez, Juan Pablo Moreno, Aldo Vilardy, Bryan Toro

Andrea.torres@umb.edu.co, fernando.prieto@docentes.umb.edu.co, jose.lagos@docentes.umb.edu.co, nixon.duarte@docentes.umb.edu.co, rosmary.martinez@docentes.umb.edu.co, juan.moreno@umb.edu.co, aldo.vilardy@umb.edu.co, bryan.toro@umb.edu.co

Abstract --- It is proposed a design of an interoperable platform for the record and analysis of data related to habits and healthy lifestyles, aiming physical activity for adult and elder populations in rural regions in the country. This tool seeks to encourage the development of strategies and public policies, working as a framework for institutions that make decisions in the field of public health and sports.

Key words: *E-health, HL7, HEVS, interoperability*

I. INTRODUCTION

Nationally, the effort to promote healthy habits and lifestyles in the population is led by COLDEPORTES, this entity is the head of sports national system at the public level and has presence in a great part of the national territory. Organizations of this system of sports that run HEVs programs represent a primary source of information that can be gained through a web platform, the analysis of this information constitutes an input for decision-making and allows direct resources and policies for vulnerable populations. [1][2].

Among the main risk factors for the development of chronic non-communicable diseases (NCDs) are, arterial hypertension, high cholesterol index, poor diet, overweight and obesity, physical inactivity, consumption of snuff and other environmental factors which may be modified; so that interventions aimed at primary risk factors can potentially reduce the risk to about 80% in cardiovascular disease and type 2 diabetes, so as to 40% of the events of cancer [4].

Physical activity is widely recognized as a health protector factor with great influence on NCD prevention factor. Yet everyday people are less active, partly by current lifestyles, estimated that at least 60% of the world population does not perform the minimum recommended physical activity. The recognition of physical inactivity as an independent risk factor, related with morbimortality [4], [5], has generated strategies all around the globe aiming to increase levels of physical activity of populations, and the healthy habits and lifestyles (HEVs), as despite being recognized and documented the benefits of it, every day the proportions of physical inactivity are higher.

The need for interoperability between health information systems has become visible, constituting this, the way to save resources and use all available sources of information, thus avoiding having to repeat the process of data collection that may have already run another information system. Internationally, organizations such as ANSI, CEN, ISO and Health Level Seven International are working on interoperability standards; one of the most important and widespread in the world is the HL7 standard.[3].

This standard in version 3 allows interoperability from the transfer of messages and documents in XML format; HL7 has taken into account the value of information to public health and therefore has a domain of quality measures in health, in which a document to request reports, the Health Quality Report Format (HQMF), is defined. This type of document eases the punctual information amongst systems since it defines the message components, in order to balance functionality and complexity of the message, the semantic interoperability is added as the HQMF comes entirely from information model defined by the HL7 RIM organization. [1]

The response to requests made by the document HQMF generate a message response based on the structure of the Quality Reporting Data Architecture (QRDA), this standard reports use Clinical Data Architecture (CDA) as the foundational standard for the specification of the report (1).[1]

The equipment of interoperability in the web platform that acquires and analyzes information, opens the doors of this information system to the integration with databases of other systems, such as health. Because the HL7 standard has been adopted internationally this information may be used by other health observatories such as the OMS.

As other information systems within the health system include interoperability through the HL7 standard, the management of information may be easier for monitoring the implementation and impact of public policies on population.

II. CONCEPTUAL FRAME

A. Estandard HL7 (Health Level 7)

The HL7 was defined by the foundation with the same name, which has dedicated since its foundation to create standards for the health sector.

In the beginning, the HL7 foundation was dedicated to elaborate specifications of messages for the sending of information between institutions of health. Through time, and observing that the institutions that wanted to establish communication had to spend large quantities of money to add units or to redesign their systems (due to their lack of an standard to define the events and elements related with patients), HL7 began the task to generate standards not only for communication but also for the medical information structure. That's how HL7 V2 was born. Following, due to the lack of a clear model of implementation, the foundation define the HL7 V3 standard, which helps the designer to don't get sidetrack with the wide quantity of possibilities that V2 allowed before, by this way, it was accomplished in the year 2004 to generate a clear standard, easy to implement and manage [21]. In the figure 2 it is shown the time line related to the standard evolution HL7.

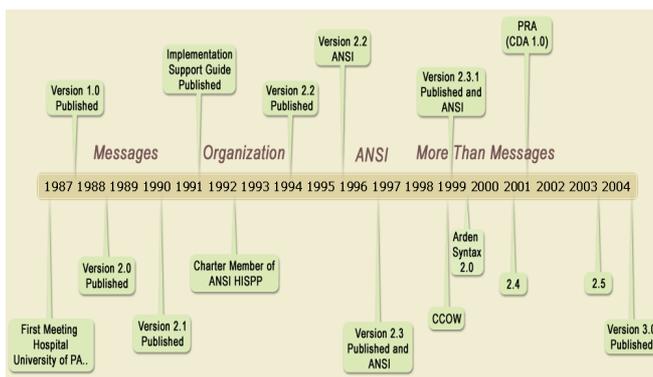


Figure 1. standard evolution HL7¹.

Like is shown in the figure 2, it was in the year 2001, when the foundation HL7 launched its first version of the CDA HL7 standard (Clinical Data Architecture), which define the structure of the electronic medical documents. For their 3 version of the HL7 standard, same as the EN1606, it was seen the necessity to add an unit of references denominated (RIM-Reference Information Model), which corresponds to the semantic representation of the elements and information recorded in clinical messages.

It is important to take in count that the Reference Model HL7 does not define the structure of the electronic medical documents, this function belongs to the CDA HL7; it is

¹ Image taken from Introduction to Health Level Seven (HL7). History of HL7. 2007. Disponible en la URL: http://www.hl7.org/documentcenter/public_temp_CCB1EE66-1C23-BA17-0C50E3F5C517F93F/training/IntroToHL7/player.html

important to understand the relation that exists between the RIM and the CDA, the tags used in the documents XML of the CDA are defined in the RIM, each version of the CDA is related directly with its corresponding RIM, which takes the classes that are needed for each specification of the document parameter to define and construct what is known as the redefine reference model or RMIM [22].

In the model of the HL7 reference there are 3 superior classes for the definition of clinical domain: event (act), number of participants in the event (role) and the subject involve in the event (entity) [23].

A. Interoperability between health observatories

The interconnection of the different departments of the health institutions and between them is an actual necessity. Each time in a greater quantity is implemented in Latin America systems of file images (PACS), the information in radiology (RIS) and in laboratory (LIS), electronic clinical histories (HCE), connected to administrative systems (HIS) and the patient's administration (ADT). This interconnection in various levels demands the use of informatic standards; for example, DIDCOM, HL7, etc. The institutions most demand, when they acquire an informatic solution, compatibility with the different standards, and its need it to know what they are about.[27][28].

The standards are protocols used by the software industry (regular form), to facilitate the interoperability or integration. The exist in the various layers of communication: - the transportation of data, that allows to transport messages with any semantic: XML, XML-HTTP, etc.; - messaging: HL7; images: DIDCOM;- vocabulary, that allows to define the controlled specific vocabulary for each domain: laboratory, pathologic anatomy, diagnose by images, nursery, procedures, etc.; - label of documents, to differentiate the different kinds of documents that can interchange and their possible content; of communications: for example, the wired, TCP/IP routers, etc. [25][26][29].

III. METHODS AND MATERIALS

The project centralizes in the use of the TIC to generate a system of interoperability that allows the consolidation of the data related with habits and styles of a healthy life, and the practice of physical activity, with the objective to promote a source of information for the surveillance of risk factors that determine the most common chronic disease. And generate an strategy to promote the practice of physical activity in regions where the access and presence, of specialized personal is limited.[5]

This is accomplished through the design and implementation of a system to manage standard data. That allows to study in detail the necessities of the population related to the acquisition of habits and healthy life styles and practices

related with the physical activity by the population. This platform includes a developed form based in the STEPwise method, defined by the OMS, for the acquisition of information; this information is stored in a data base.

The objective of the platform is to recollect, store and process the information, the result of this analysis comprehends a battery of indicators that will be the posted information in the observatory of habits and healthy life styles (OHEVS).[7].

A. Definition of the range of interoperability

System of Interoperability is established for the messages that are generated from the web platform that are directed to other systems of external information, which allows the observatory OHEVS to be a source of information, that can be consulted by a web service by other systems of information. In overall is contemplated the generation of 32 measurements. [11][16].

B. Standard selection.

Beginning from the information to manage inside the web platform of the observatory OHEVS, is define that the standard for interoperability of the system should achieve with the following requisites:

- Specialized standard in the transference of health information.
- The standard most allowed the prosecution of messages and documents b machines.
- The standard most be accepted in the international community.
- Most exist an evidence of the use of standards in a national level.
- It has to contain information about interchanging information and reports of public health.

From this requirements it was made a revision of available standards and it came to a conclusion that the most adequate standard for the system of interoperability is the version 3 of HL7.

C. Domain identification HL7 V3.

Once is identify the standard that is going to be used, it was necessary to identify the domain closest to the necessity of the system, thanks to the time of development the HL7 have a wide range of domains that cover the principal functions of the health information systems.

The messages or reports that are generated by the observatory contain information of entire populations, by which is necessary to evaluate the information of more than one person in the consult of the date base, for the cases the HL7 counts with the domain UVQM, that standardizes the format for the representations of measurements of qualities in health HQMF. Next is shown the structure of the document HQMF.

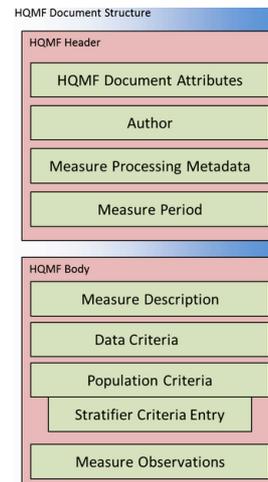


Figure 2. document structure HQMF (2).

The document is divided into two principal sections, the headline records the principal information of the document and its measure; the body of the document carries specific information that describes the measurement, and complementary information that improves the posted information. This type of messages can be utilized for the application of reports generated form the observatory. [8].

The HQMF contemplates in its domain the UVQM the generation for answers through the QRDA, the process of the delivery of messages is illustrated in the next figure.



Figura 3. surroundings of messages of quality of health.

In the figure that is being observed in the superior part the development of the application of the measurement, this is generated by a system of external information besides the observatory, the specification of the measurement is defined in the document HQMF that arrives to the providers of information, in this case the observatory OHEVS corresponds to the platforms that consults in the data base, to finally generate a message of respond or Quality Report with a message QRDA.[18].

The utility of this model of data is shown in the next illustration.

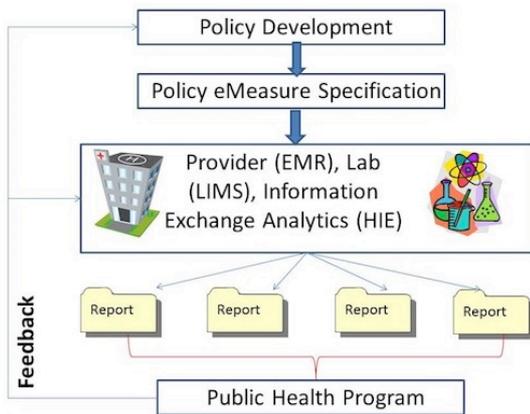


Figure 4 surrounding health public messages.

In HL7 is contemplated the use of the document HQMF and the QRDA for surveillance of the implementation of public politics, it also represents how tools or supplies for the choosing of decisions and the generation and the adjustment of politics directed to improve the population health

It's important to know the existence of three types or categories of messages QRDA:

- QRDA Category I – Single Patient Report: this is a report at an individual level of measurements of health quality; the elements of data content in the report are particular reported measures.
- QRDA Category II – Patient List Report: reports that allowed to include measures of quality in health for multiple patients, this types of reports refers to one or many measures for each patient of the list.
- QRDA Category III – Calculated Report: generate reports of quality measures in health with the estimate of each measurement for one population inside a period of specific measurement. Allows the extraction of information without referring to personal data with each patient, which prevents the possibility that someone access to patients personal data without permission.

Taking in count that the observatory generates the quality reports in health starting from the available information of patients in a population, the category that is closer to the requirements of the system is the category III that shows the estimate measurements of the population and protects the identity of the patients that make up the population measure.

The information flow and the interaction between messages HQMF are shown in the next figure.

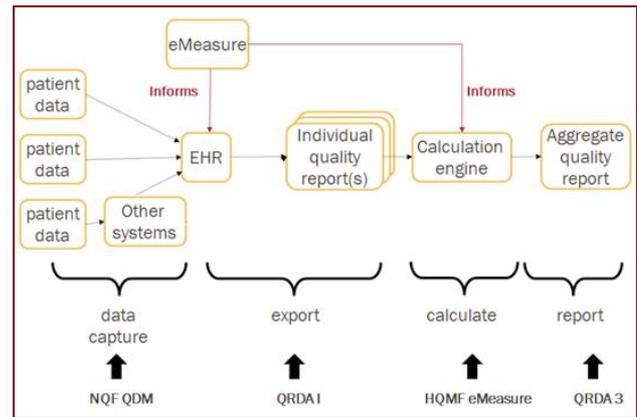


Figure 5. Report of measures using HQMF y QRDA.

D. Development methodology HL7 V3.

The version 3 of the HL7 standard base its programming from Unified Modeling Language (UML), this allows you to model real-life situations to generate messages according to the situation or event that generates them. Interoperability with HL7 is thanks to the definition of a vocabulary for all applications that may derive from the standard, all classes and concepts that can be used in HL7 messages within the Reference Information Model (RIM).

Thanks to the development time that the standard has, a revised documentation for each domain is featured, it is also possible to download the legislative package from the official website of HL7, in this are tools and examples for the implementation of the standard in specific applications. To speed the development process it was started from the refined model reference information coded message or RMIM coded "POQM_RM00001UV", this model uses the entry point named eMeasure to reference the main class QualityMeasureDocument.[17].

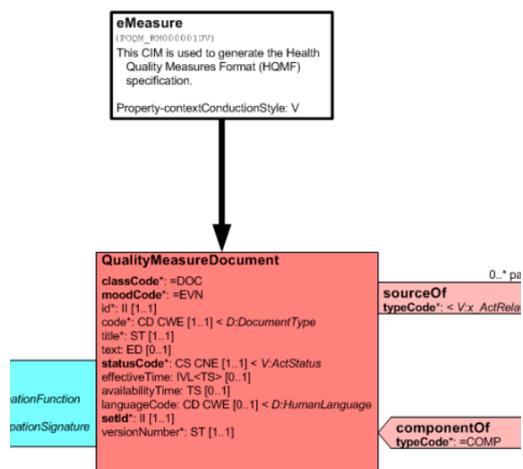


Figure 6. principal class eMeasure.

Within the methodology of implementation of HL7 messages begins with RIM, to this will be applied restrictions and DMIM is obtained which include messages for a specific domain, from the latter and with more restrictions you get the refined model or RMIM containing the classes to be included in a specific message type, with the tools available to HL7 is possible to pass the class diagram RMIM to the hierarchical description of the message or HDM, this is a document in which RMIM classes are obtained with their attributes in tabular form, from this hierarchical model the application can generate messages, assigning the required values for each attribute in a tag file in XML format.

Not always the model classes of the refined reference are used, or it can happen that is necessary to include other classes or attributes within the model, the standard has the flexibility to allow such changes as long as you do not go RIM.

The response message QRDA is based on CDA standard, hence the reference information model RMIM is the same as for clinical documents based on CDA, these contemplate a header with the document information as identification, author, custodian. Then there is the measurement section that has a subsection of reporting parameters with relevant information QRDA measures.

IV. RESULTS

After formulating this proposal the results that this would bring would be linked to the development of the QRDA document adapted to the requirements of the reports generated by the observation.

Likewise the implementation of the web service for automatic generation of the response message from the queried information in the database is achieved.

The development of this project will allow each of the entities that implement it improve the incursion of new technologies in view of improving response times.

V. CONCLUSIONS

This project will allow to understand the problems that often are not taken into account as it is in health computer science. Compression performed with the standard HL7 is currently not easy but what is sought is to be applied to public health and diagnoses that apply to it. The issue of health has many topics that are very difficult to address in a single proposal, however it would be interesting to take this proposal as the beginning of a full investigation that may contribute to the knowledge of this important subject for Colombia.

HL7 standard is beginning to become strong in Colombia, but for being a private entity, the different institutions health-providers are beginning to use it in their systems gradually, and in this moment, there are few that have adopted it.

Within the development of this proposal, the capacity is given to the government to monitor risk-factors of ECNT to guide the resources in preventive measurements that promote physical activity.

This project will allow to take as source information databases of other information systems such as the National Health System in an automated way so that the observatory has a greater capacity to be integrated with other systems and its response could be more effective and closer to the reality of the study population.

VI. BIBLIOGRAPHY

- [1] Principles of health interoperability HL7 and SNOMED. Tim Benson
- [2] International Committee of the Red Cross. Humanitarian Situation Report of Activities 2011 <http://www.icrc.org/spa/assets/files/2012/informe-Colombia-2011.pdf>
- [3] Veléz, Alba L. Public Policy Considerations of interest to health professionals. He documented pdf available at: http://promocionsalud.ucaldas.edu.co/downloads/Revista%2010_4.pdf
- [4] World Health Organization. Implementation of the Global Strategy on Diet, Physical Activity and Health. A guide to approaches to increase population levels of physical activity 2008
- [5] World Health Organization, Global Recommendations on physical activity for health. Geneva, 2010
- [6] Colombian Institute of Family Welfare. National Health and Nutrition Examination Survey ENSIN 2010.
- [7] World Health Organization. First International Conference on Health Promotion. Ottawa 1986.
- [8] Republic of Colombia, Ministry of Health and Social Protection. Methodological Guide for records, observatories, monitoring systems and national health situation rooms. 2013.
- [9] Health Crisis in Colombia. Social Emergency Economic Statement. November 2009. Site: Health Consultant. Carlos Felipe Muñoz Paredes. <http://www.consultorsalud.com/biblioteca/documentos/2009/Crisis%20de%20la%20Salud%20en%20Colombia%20-%20emergencia%20social%202009%20Consultorsalud.pdf>
- [10] Law 1419 2010 Guidelines for the development of telehealth in Colombia. <http://wsp.presidencia.gov.co/Normativa/Leyes/Documents/ley141913122010.pdf>
- [11] Law Reform 2011. 1438 general social security health and other provisions. <http://guajiros.udea.edu.co/fnsp/cvsp/ley1438.pdf>
- [12] Law 100 of 1993. Creation System Integral Social Security. <http://www.colombia.com/actualidad/images/2008/leyes/ley100.pdf>
- [13] Decree 1281 of 2002. Rules governing the cash flows and the timely and efficient use of health resources and their use in the provision. http://www.saludcolombia.com/actual/htmlnormas/Dec1281_02.htm
- [14] Resolution 3374 of 2000. essentials that must report health care providers and managers benefit plans on health services provided entities. <http://www.fosyga.gov.co/LinkClick.aspx?link=MarcoNormativo%20FECAAT%20Resolucion+3374+of+2000.pdf&tabid=310&mid=598>

[15] Glossary. FOSYGA.
<http://www.fosyga.gov.co/AcercaDelFOSYGA/GlosariodeT%C3%A9minos/tabid/324/Default.aspx>

[16] Benefit Plan provision of health services. Annex 20. Protocol reference and counter.
http://www.contratos.gov.co/archivospucl/DA/124004000/07-2-77212/DA_PROCESO_07-277212_124004000_242552.pdf

[17] EN 13606 Association. The CEN / ISO standard.
<http://www.en13606.org/>

[18] Usefulness of ISO 13606 Archetypes to represent detailed clinical models. Pablo Serrano, David Moner, Thomas Sebastian, Jose Maldonado, Rafael Navalon, Montserrat Robles, Angel Gómez.
 INFOLAC2008-AAIM.
<http://revistaesalud.com/index.php/revistaesalud/article/view/308/641>

[19] EN 13606 Association. CEN / ISO EN 13606 Archetype Model.
<http://www.en13606.org/>

[20] Using the European standard EN13606 an electronic medical record system federated. José Alberto Maldonado Segura, Montserrat Robles Viejo, Pere Crespo Molina, Carlos Fernández Angle Andres Sanchis Estruch, Saura Alfonso Herranz.
<http://www.ibime.upv.es/bie/docs/I+S2005.pdf>

[21] The HL7 Evolution. Comparing HL7 Version 2 to Version 3, including a History of Version 2. Corepoint Health. 2010. Available at URL:
<http://www.corepointhealth.com/sites/default/files/whitepapers/hl7-v2-v3-evolution.pdf>

[22] HL7 FAQs. Health Level Seven International. Available at URL:
<http://www.hl7.org/about/FAQs/index.cfm>

[23] http://www.shopcreator.com/mall/Abiescouk/Downloads/chapter_7_the_hl7_v3_rim.pdf

[24] Semantic Integration and standardization of clinical data based on archetypes. Moner D., J. A. Maldonado, D. Boscá, JT Fernández, C. Angulo, PJ Vivancos, M. Robles. XXIV Annual Congress of the Spanish Society of Biomedical Engineering. November 2006.

[25] LinkEHR-Ed: A tool for standardization of electronic medical records. J. A. Maldonado, D. Moner, D. Boscá, C. Angulo, I. Abad, D. Pérez, P. Serrano, E. Reig, M. Robles.

[26] LinkEHR Normalization Platform. Available at URL:
<http://www.linkehr.com/>

[27] Methodological aspects of the process of adopting the standard HL7v3 in Colombia: the experience of the Technical Committee on Use Cases Clinical Laboratory. Gabriel Tamura, Nhora Villegas, Fernando Portilla.

[28] Heon - Health On Line. Available at URL: <http://www.heon.com.co/>

[29] List of institutions providing health services registered for direct shift from nation, ministry of health and social protection, June 2012 accounts.

[30] Law 1430 of 2011, amends the general social security health, ministry of social protection.

A compact microstrip lowpass filter using a stepped impedance hairpin resonator with radial stubs

M. Samadbeik, B. F. Ganji, A. Ramezani

Abstract— A compact microstrip lowpass filter with a stepped impedance hairpin resonator is proposed. New filter has been designed, simulated and analysed to achieve wide stop band, low return loss, minimal transition band and small size. The cut off frequency is 2 GHz, return loss is better than -20 dB at 0.8 GHz. This compact lowpass filter with sharp cut off frequency response and broad stopband better than -15 dB should be useful in many wireless communication systems.

Keywords— Lowpass filter; step impedance hairpin resonator; stopband range

INTRODUCTION

MICROWAVE lowpass filters (LPFs) with the demand for compact size, low insertion loss, sharp transition, and wide stopband are highly desirable in wireless communication systems to suppress harmonic and spurious signals. Since resonators are the basic components of planar filters, it is key to select a proper resonator type. In order to decrease the size of resonators, planar resonators have been studied for many years, such as stepped impedance resonators [1,5] To obtain a sharp cutoff frequency and wide stopband response, most conventional filters need more sections, but increasing the number of sections also increases the size of the filters and insertion loss [4].

The conventional low-impedance stubs are replaced by radial stubs to realize a wide stopband rejection. A filter prototype has been designed, fabricated and measured to verify the validity of the proposed technique [1]. In [1] the stepped impedance hairpin resonator (SIHR) was first introduced to design lowpass filters (LPFs), but it has problem about wide stopband. In [3,5] both single hairpin resonators with additional line and cascades of stepped impedance hairpin resonators were used but sizes isn't suitable for a small filter because of use more resonators. In [4], a compact structure with two transmission zeros is used by tap-connecting the

stepped-impedance hairpin unit with an interdigital structure to introduce internal coupling between the two low-impedance sections to achieve wide stopband and sharp cut off frequency. The radial stub has intrinsic wide stopband parameter, as demonstrated in [5,6].

In [7] A compact lowpass filter using a stepped impedance hairpin resonator with radial stubs is designed and its result is good in size of filter but the wide stopband is limited, the return loss is about -12 dB, the sharp cut off frequency is about 0.7.

In this paper, the conventional step-impedance hairpin resonator by radial stubs are changed to realize a wide stopband rejection, sharp transition, low insertion loss, lower return loss and keep its size small.

I. FILTER DESIGN

A. study proposed filter in [7]

The stepped impedance hairpin resonator with radial stubs proposed in [7] is shown in Fig. 1.

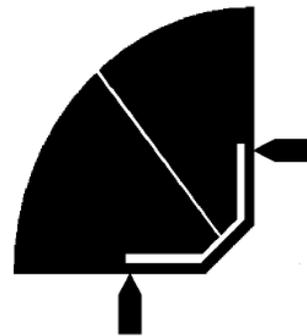


Fig. 1 Layout of stepped impedance hairpin resonator structure in [7]

Fig. 2 shows the simulated results of resonator proposed in [7].

Corresponding author : M. Samadbeik (samadbeik.m@gmail.com).

¹ Mahya Samadbeik, Electrical Engineering Department, Saveh Science and Research Branch Islamic Azad University, Saveh, Iran (samadbeik.m@gmail.com).

² Behrooz Fathi Ganji, Electrical Engineering Department, Lorestan University, Kamalvand, Khorramabad, Lorestan, Iran (ganji.behrooz@yahoo.com).

³ Abbas Ramezani, Electrical Engineering Department, Lorestan University, Kamalvand, Khorramabad, Lorestan, Iran (ramezani.ab@lu.ac.ir).

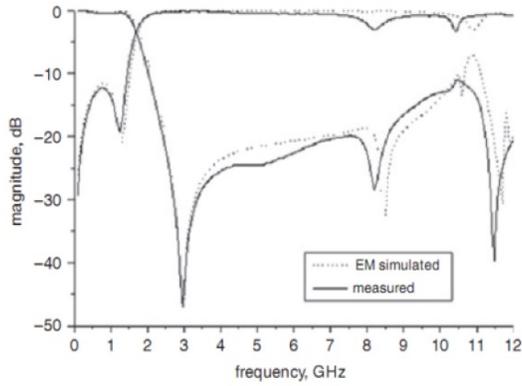


Fig. 2 Simulation results of conventional LPF

This filter has a better than -15 dB stopband rejection from 2.3 GHz to 10 GHz, that show the bandwidth is low. Cut off frequency is equal to 1.67 GHz so it has gradual cut off frequency but size of filter is 100 mm² and acceptable.

B. Design a new low-pass filter

In this article we are going to design a microstrip low-pass filter to meet the needs of today's modern technology for this purpose, we need to simulate a designed filter, as a result we should use of software.

The ADS software is one of the microwave analysis software. In this software we can use layout part directly and design microstrip filter and observe result of simulation according to different parameters.

Designed resonator is based on a hairpin resonator was first introduced in 1972 and having a simple structure and small size is the most characteristics.

Fig 6 shows the structure of the hairpin resonator. The shape of the LC circuit of resonator is also provided.

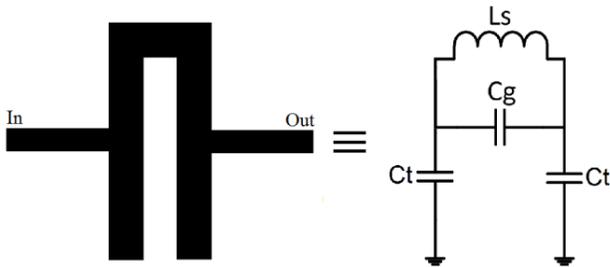


Fig. 3 The basic structure of the hairpin resonator

In the new layout our purpose is designed resonator with radial prongs to improve the gradual transition bands like [7] article and with a change in the structure of the resonator, provide new filter with the appropriate parameters. For this purpose, resize angle resonator from 90 ° to 45 °. Fig 4 shows this resonator.

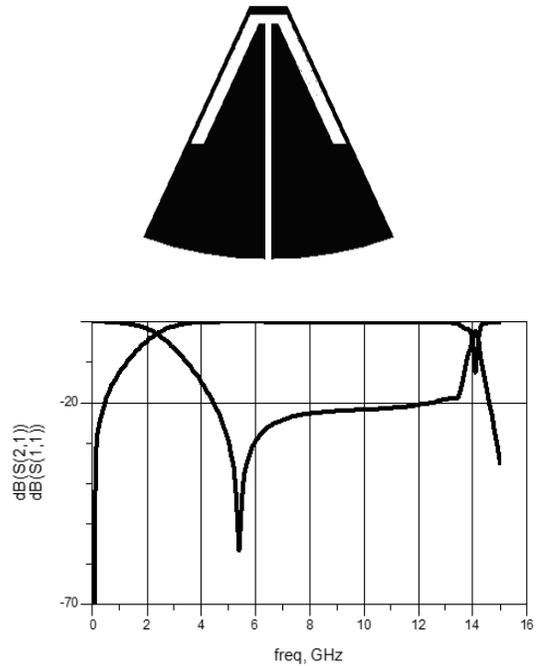


Fig. 4 The new resonator and The frequency response curve

In this resonator cut off frequency is 2.4 GHz, it has a better than -20 dB stopband rejection from 4.4 GHz to 13.4 GHz, filter size is 60 mm² but there is main problem in this resonator with sharp response which is about 2 and it show that it has gradual cut off so introduced new filter with two resonator in Fig 5.

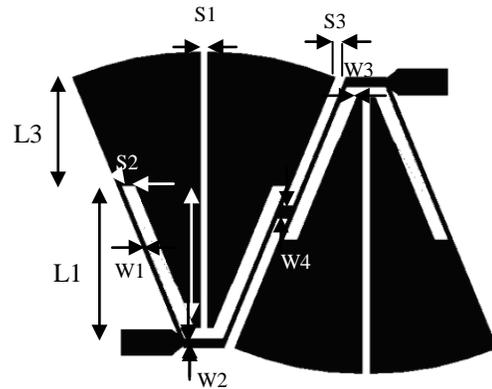


Fig. 5 Layout of new LPF with parameter sizes

The design has been optimized using an EM-simulator of ADS. All the parameters of dimensions are determined as follows: W1= 0.13 mm, W2 = 0.2 mm, W3 = 0.1 mm, W4 = 0.4 mm, L1 = 4.5 mm, L2 = 4.18 mm, L3 = 3.36 mm, S1 = 0.3 mm, S2 = 0.5 mm, S3 = 0.4 mm.

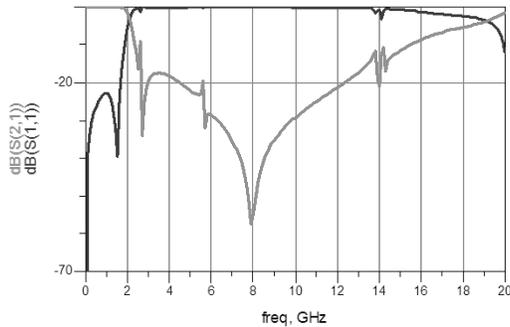


Fig. 6 Simulation results of proposed LPF

This filter compared with conventional filter in [7] and by observing the figures, we can see that the proposed LPF has a sharp transition and it is 36.5% better than conventional LPF, and the 3 dB cutoff frequency is 2.1 GHz. The proposed LPF has a wider stopband and its frequency range is 2.5 – 13 GHz referred to a criterion of better than -18 dB and it is seven times of cut-off frequency. But frequency range in the conventional LPF is 2.3-10 GHz referred to a criterion of better than -15 dB with the cut-off frequency 1.67 GHz. The results show that stopband in proposed LPF improved to 44.6%. The return loss in the proposed LPF is -23 dB that is better than -20 dB but in the conventional LPS, return loss is about -12 dB. After viewing the results, we found that the return loss in the proposed LPF decreases 48%.

The important point is the filter size that is so considerable. Filter size is 800 mm² and shows that its size is smaller with other better parameter.

II. CONCLUSION

We proposed a compact and small LPF designed with a stepped impedance hairpin resonator with radial stubs that used of 45° angle and different structure.

The results indicate that the demonstrated filter achieves very good performance in terms of compact size and low insertion loss in the passband. The proposed filter also exhibits a very wide stopband, and is able to suppress the harmonic response. With this good performance, the proposed structure has potential applications in modern communication systems.

REFERENCES

- [1] H. L. Hsieh, and K. Chang, "Compact lowpass filter using stepped impedance hairpin resonator," *Electron. Lett.*, pp. 899-900, 2001.
- [2] H. L. Hsieh, and K. Chang, "Compact elliptic-function low-pass filters using microstrip stepped-impedance hairpin resonators," *IEEE Trans. Microw. Theory Tech.*, pp. 193-199, 2003.
- [3] H. L. Hsieh, and K. Chang, "Compact, broad stop-band elliptic lowpass filters using microstrip stepped impedance hairpin resonator," *IEEE MTT-S int. Microw. Symp. Dig.*, pp. 1775-1778, 2003.
- [4] M. H. Yang, J. Xu, "Design of compact, broad-stopband lowpass filter using modified stepped impedance hairpin resonator," *Electron. Lett.*, pp. 1198-1200, 2008.

- [5] F. Giannini, R. Sorrentino, and I. Vrba, "Planner circuit analysis of microstrip radial stub," *IEEE Trans Microw. Theory Tech.*, pp. 1652-1655, 1984.
- [6] R. Sorrentino, and L. Roselli, "A new simple and accurate formula for microstrip radial stub," *IEEE Microw. Guide. Wave Lett.*, pp. 480-482, 1992
- [7] X. B. Wei, P. Wang, and M. Q. Liu, "Compact wide-stopband lowpass filter using stepped impedance hairpin resonator with radial stubs," *Electron Lett.*, vol 47, July 2011.

Modification of the cryptographic algorithms, developed on the basis of nonpositional polynomial notations

Rustem G. Biyashev, Saule E. Nyssanbayeva, Yenlik Ye. Begimbayeva, Miras M. Magzom

Abstract — Cryptographic systems, developed on the basis of nonpositional polynomial notations (NPNs), are called nonconventional, nonpositional or modular. In this paper, models of modified nonconventional encryption systems and digital signature are described. The creation of the model of block encryption system includes the development of a modified nonpositional block encryption algorithm using the analogue of Feistel system and application mode of this modified algorithm. The model of a digital signature based on the scheme of the Digital Signature Algorithm (DSA) and NPNs. Application of NPNs allows creating effective cryptographic systems of high reliability, which enables the confidentiality, authentication and integrity of stored and transmitted information. Synonyms of NPNs – classical notations in residue number system (RNS), polynomial notations systems in RNS, modular arithmetic.

Keywords — cipher mode, digital signature, encryption, nonpositional polynomial notations.

I. INTRODUCTION

THE basis for the creation of the proposed models of cryptosystems are nonconventional systems of encryption and digital signature. These systems are developed on the algebraic approach base, using nonpositional polynomial notations (NPNs) or polynomial notations in residue classes (polynomial RNS). Classical RNS (modular arithmetic) is based on the Chinese remainder theorem, which states that any number can be represented by their remainders (residues) from its division by the base numbers systems, which are formed pairwise coprime numbers [1]-[2]. Then in RNS a positive integer is represented by a sequence of remainders or residues

$$A = \alpha_1, \alpha_2, \dots, \alpha_n \quad (1)$$

The conducting research is funded by the Ministry of Education and Science of the Republic of Kazakhstan (MES RK).

R. G. Biyashev is with the Institute of Information and Computational Technologies of MES RK, 125 Pushkin str., Almaty, 050010, Republic of Kazakhstan. (e-mail: brg@ipic.kz).

S. E. Nyssanbayeva is with the Institute of Information and Computational Technologies of MES RK, 125 Pushkin str., Almaty, 050010, Republic of Kazakhstan (phone: +77017743730, e-mail: sultasha1@mail.ru, snyssanbayeva@gmail.com).

Ye. Ye. Begimbayeva is with the Institute of Information and Computational Technologies of MES RK, 125 Pushkin str., Almaty, 050010, Republic of Kazakhstan (e-mail: enlikb89@gmail.com).

M. M. Magzom is with the Institute of Information and Computational Technologies of MES RK, 125 Pushkin str., Almaty, 050010, Republic of Kazakhstan (e-mail: magzomzn@gmail.com).

from dividing this number by the given positive integer numbers p_1, p_2, \dots, p_n , which are called bases of the system.

Numbers α_i are formed in the following way:

$$\alpha_i = A - [A / p_i] p_i, \quad i = \overline{1, n}, \quad (2)$$

where $[A / p_i]$ denotes the integer part of the division A by p_i . From (2) follows, that the number α_i of i -th digit of A is the smallest positive remainder of division A by p_i , and $\alpha_i < p_i$. In this case, the formation of each digit number performed independently. According to the Chinese remainder theorem, representation of A in the form of (1) is unique, in case numbers p_i are pairwise coprime. The range of representable numbers in this case is $P = p_1 p_2 \dots p_n$. Here, similar to a positional number system, the range of representable numbers growing as the product of base numbers, and the digit capacity of the number is growing as the sum of the digit capacity of the same base numbers.

In NPNs (polynomial RNS) bases are used as irreducible polynomials over field $GF(2)$ [3]-[4]. Using NPNs allows reducing the length of the key, to improve durability and efficiency of nonpositional cryptographic algorithms [4]-[5]. Improving the efficiency is provided by the rules of NPNs in which all arithmetic operations can be performed in parallel to the base module NPNs. In developed nonconventional cryptographic algorithms the encryption and formation of digital signature is carried out for an electronic message of the given length. In nonpositional cryptosystems as a criterion of cryptostrength is used cryptostrength of algorithms of encryption and formation of digital signature, which is characterized by a complete secret key. Cryptostrength in this case depends not only on the length of a key sequence, but also on choice of a system of polynomial bases. With the growth of the order of irreducible polynomials with binary coefficients, their number also grows rapidly. Therefore, a wide choice of polynomial bases is possible. Cryptostrength of proposed encryption algorithm which using NPNs significantly increases with the length of the electronic message.

In [3] the arithmetic of nonpositional number systems with polynomial bases and its application to problems of improving reliability are developed. As it is shown, the algebra of polynomials over a field in modulus of the irreducible polynomial over this field is a field and the representation of

the polynomial in the nonpositional form is the only (analogous to the Chinese remainder theorem for polynomials). The rules of performing arithmetic operations in NPNs and restoring the polynomial by its residues are defined. According to the Chinese remainder theorem, all working base numbers should be different.

II. NONPOSITIONAL POLYNOMIAL NOTATIONS

A. Constructing of NPNs

The process of forming of NPNs for an electronic message M of the given length N bits is as follows. Polynomial bases with binary coefficients are selected

$$p_1(x), p_2(x), \dots, p_s(x), \tag{3}$$

where $p_i(x)$ - irreducible polynomial with binary coefficients of degree m_i respectively, $i = \overline{1, S}$. These bases are called working base numbers. The main working range in NPNs is a polynomial $P(x) = p_1(x)p_2(x) \dots p_s(x)$ of the degree $m = m_1 + m_2 + \dots + m_s$. According to the Chinese remainder theorem, all the base numbers must be different even if their degrees are equal.

In NPNs any polynomial $F(x)$, which degree is less than m , has a unique nonpositional representation in a form of sequence of residues of its division by the working base numbers $p_1(x), p_2(x), \dots, p_s(x)$:

$$F(x) = (\alpha_1(x), \alpha_2(x), \dots, \alpha_s(x)), \tag{4}$$

where $F(x) = \alpha_i(x) \pmod{p_i(x)}$, $i = \overline{1, S}$.

In NPNs a message (or its block) of the given length N bits is represented as follows. It is interpreted as a sequence of remainders of division of some polynomial (let us denote it as $F(x)$) by working base numbers $p_1(x), p_2(x), \dots, p_s(x)$ of degree not greater than N , that is, in the form of (4). Each working base number should have a degree not exceeding value of N . These base numbers are selected from all irreducible polynomials with degrees varying from m_1 to m_s , providing that the following equation is satisfied [6]:

$$k_1 m_1 + k_2 m_2 + \dots + k_s m_s = N. \tag{5}$$

Here $0 \leq k_i \leq n_i$ are unknown coefficients and the number of selected irreducible polynomials of degree m_i . One certain set of these coefficients is one of the solutions of (5) and specifies one system of working base numbers, n_i is the number of all irreducible polynomials of degree m_i , $1 \leq m_i \leq N$, $S = k_1 + k_2 + \dots + k_s$ is a number of selected working base numbers. In the system of working bases the order of these

bases is also taken into account.

Equation (5) defines the number S of working bases, which produce residues that covers the length N of the given message. Complete residue systems modulo polynomials of degree m_i include all polynomials with the degree not exceeding $m_i - 1$. The representation of polynomials of degree $m_i - 1$ requires m_i bits.

As shown in Table I, with growth of degrees of irreducible polynomials, their amount rapidly increases, and, as a result, the number of solutions of (5) also considerably increases.

Calculations for finding irreducible polynomials were conducted in two ways: by dividing a particular polynomial to other polynomials and using analog of the sieve method for finding prime numbers. The results of these calculations matched by both quantitative and qualitative composition.

The properly checked table of irreducible polynomials over field $GF(2)$ for the degrees from 1 to 15 was published in [7].

TABLE I. DEPENDENCE OF NUMBER OF IRREDUCIBLE POLYNOMIALS ON THEIR DEGREE

Degree of Irreducible Polynomials	Number of Irreducible Polynomials
1	1
2	1
3	2
4	3
5	6
6	9
7	18
8	30
9	56
10	99
11	186
12	335
13	630
14	1161
15	2182
16	4080
17	7710
18	14532
19	27594
20	52377

Remainders $\alpha_1(x), \alpha_2(x), \dots, \alpha_s(x)$ are selected in the way where binary coefficients of remainder $\alpha_1(x)$ correspond to the first l_1 bits of the message, the next binary coefficients of remainder $\alpha_2(x)$ correspond to the next l_2 bits, etc., and binary coefficients of remainder $\alpha_s(x)$ correspond to the last l_s binary bits.

The positional representation of $F(x)$ is reconstructed from its form (4) [3]-[4]:

$$F(x) = \sum_{i=1}^s \alpha_i(x) B_i(x), B_i(x) = \frac{P_s(x)}{p_i(x)} M_i(x), i = \overline{1, S}. \quad (6)$$

Polynomials $M_i(x)$ are chosen so as to satisfy the congruence in (6).

B. Hashing an electronic message in NPNs

In NPNs it is possible to hash (compress) an electronic message of the given length N to the length of N_k bits [3]-[4]. This is done by introducing redundancy, that is, the message in NPNs is expanded by redundant bases $p_{s+1}(x), p_{s+2}(x), \dots, p_{s+U}(x)$. The system of redundant bases is formed independently of the choice of working base numbers $p_1(x), p_2(x), \dots, p_s(x)$. Note that some bases among the U redundant bases may coincide with some of the working base numbers.

Selection of redundant bases is carried out by analogy with a choice of working bases. These bases are chosen randomly from all irreducible polynomials of degree not exceeding the value of N_k . Denote the degree and the number of irreducible polynomials used in their selection as a_1, a_2, \dots, a_U and d_1, d_2, \dots, d_U respectively. The number of selected redundant bases in this case is determined from the equation (the analogue of (5)):

$$t_1 a_1 + t_2 a_2 + \dots + t_U a_U = N_k, \quad (7)$$

where $0 \leq t_j \leq d_j$, $0 \leq a_j \leq N_k$, $j = \overline{1, U}$, t_j - the number of selected redundant bases of degree a_j , $U = t_1 + t_2 + \dots + t_U$ - the number of selected redundant bases, which produce residues that covers the hash value of length N_k . Solution of the (7) defines a single system of redundant bases.

Further redundant residues (remainders) $\alpha_{s+1}(x), \alpha_{s+2}(x), \dots, \alpha_{s+U}(x)$ are calculated by dividing reconstructed polynomial $F(x)$ by redundant bases $p_{s+1}(x), p_{s+2}(x), \dots, p_{s+U}(x)$. Then the hash value $h(F(x))$ of length N_k bits can be interpreted as a sequence of these residues:

$$h(F(x)) = (\alpha_{s+1}(x), \alpha_{s+2}(x), \dots, \alpha_{s+U}(x)), \quad (8)$$

where $h(F(x)) \equiv \alpha_{s+j}(x) \pmod{p_{s+j}(x)}$, $j = \overline{1, U}$. The sum of the lengths of redundant residues is the length of hash value.

III. NONCONVENTIONAL SYMMETRIC ENCRYPTION ALGORITHM

The encryption algorithm of an electronic message of the given length N bits based on NPNs includes the following steps. Initially nonpositional polynomial number system is formed (this procedure is described in Subsection A). Then a key (pseudo-random) sequence is generated, and the plaintext is encrypted.

Suppose that for encryption from the set of all irreducible polynomials of degree not exceeding N a system of working base numbers (3) is selected. The message of length N bits is represented as a sequence of residues (4) from the division of a polynomial on the working bases.

Then the encryption key length of N bits is also interpreted as a system of residues $\beta_1(x), \beta_2(x), \dots, \beta_s(x)$, but from division of other polynomial $G(x)$ by the same working base numbers:

$$G(x) = (\beta_1(x), \beta_2(x), \dots, \beta_s(x)), \quad (9)$$

where $G(x) \equiv \beta_i(x) \pmod{p_i(x)}$, $i = \overline{1, S}$.

After encrypting the message $F(x)$ using the key $G(x)$ a cryptogram is obtained. This cryptogram is considered as a function $H(x)$:

$$H(x) = (\omega_1(x), \omega_2(x), \dots, \omega_s(x)), \quad (10)$$

where $H(x) \equiv \omega_i(x) \pmod{p_i(x)}$, $i = \overline{1, S}$. In (10) the first l_1 bits of cryptogram are assigned to binary coefficients of remainder $\omega_1(x)$, l_2 bits of cryptogram are assigned to binary coefficients of remainder $\omega_2(x)$, etc. The last l_s bits of cryptogram are assigned to binary coefficients of the last remainder $\omega_s(x)$.

In software implementation of this nonconventional algorithm of encryption of the message the nonconventional method will be used [8]. The usage of different methods allows obtaining different encryption models.

Key length is one of the system strength indicators. In nonconventional encryption the strength of cryptographic algorithm characterized by complete (private) key is used as a cryptostrength criterion. In this algorithm a complete key is the polynomial $G(x)$ and the certain set of working base numbers chosen from the set of irreducible polynomials whose degree does not exceed N [9].

Statement 1. The cryptostrength of an encryption algorithm developed on the basis of NPNs is determined by total number of possible and distinct from each other variants of choice of key sequences and systems of working base numbers.

To prove the above fact, the combination number of choice of base numbers for each base number degree determined by the (5) is calculated. Then the number of combinations of

system forming from S base numbers with the degrees m_1, m_2, \dots, m_S with allowance for their arrangement is determined by expression

$$(k_1 + k_2 + \dots + k_S)! C_{n_1}^{k_1} C_{n_2}^{k_2} \dots C_{n_S}^{k_S}.$$

The encryption is performed by imposing on the message of the generated key sequence of the same length N bits. Therefore for encryption the choice of one system from S base numbers is defined by the formula:

$$(k_1 + k_2 + \dots + k_S)! C_{n_1}^{k_1} C_{n_2}^{k_2} \dots C_{n_S}^{k_S}. \quad (11)$$

Then the encryption cryptostrength of the message of length N bits is determined as the inverse value for (11):

$$p_{kr} = 1 / (2^N \sum_{k_1, k_2, \dots, k_S} (k_1 + k_2 + \dots + k_S)! C_{n_1}^{k_1} C_{n_2}^{k_2} \dots C_{n_S}^{k_S}). \quad (12)$$

In expression (12) the summation is performed over all possible combinations of coefficients k_1, k_2, \dots, k_S satisfying the (5).

Nontraditional method in which the elements of the sequence of residues $\omega_1(x), \omega_2(x), \dots, \omega_S(x)$ in the cryptogram are the smallest remnants of division of products $\alpha_i(x)\beta_i(x)$ by respective bases $p_i(x)$ is used for encryption if multiplication operation is used as the function $H(F(x), G(x))$ [8]:

$$\alpha_i(x)\beta_i(x) \equiv \omega_i(x) \pmod{p_i(x)}, \quad i = \overline{1, S}. \quad (13)$$

For deciphering cryptogram by the known key $G(x)$ for each value $\beta_i(x)$ the calculation of the reverse (inverse) polynomial $\beta_i^{-1}(x)$ is made as follows from (13) provided that the following equation is satisfied:

$$\beta_i(x)\beta_i^{-1}(x) \equiv 1 \pmod{p_i(x)}, \quad i = \overline{1, S}. \quad (14)$$

The result is the polynomial $G_i^{-1}(x) = (\beta_1^{-1}(x), \beta_2^{-1}(x), \dots, \beta_S^{-1}(x))$ inverse to the polynomial $G(x)$. Then the elements of the sequence of residues (4) in accordance with (13) and (14) are restored as compared with:

$$\alpha_i(x) \equiv \beta_i^{-1}(x)\omega_i(x) \pmod{p_i(x)}, \quad i = \overline{1, S}.$$

Thus, in the present model of the encryption algorithm of electronic message of the specified length N bits in NPNs, the complete key is:

- the chosen system of polynomial working bases $p_1(x), p_2(x), \dots, p_S(x)$;
- the key $G(x) = (\beta_1(x), \beta_2(x), \dots, \beta_S(x))$;
- the key $G_i^{-1}(x) = (\beta_1^{-1}(x), \beta_2^{-1}(x), \dots, \beta_S^{-1}(x))$ needed for deciphering and inverse to $G(x)$.

Examples of determination of cryptostrength by the formula (12).

1. Key length equals 100 bits: system of base numbers includes 6 irreducible polynomials of degree 16 and 1 irreducible polynomial of degree 4. $S=7$. For this system of base numbers we obtain $p_{kr} \approx 10^{-53}$.

2. Key length equals 200 bits: system of base numbers includes 12 irreducible polynomials of degree 16 and 1 irreducible polynomial of degree 8. $S=13$. $p_{kr} \approx 10^{-106}$.

3. Key length equals 128 bits: system of base numbers includes 8 polynomials of degree 16. $S=8$. $p_{kr} \approx 10^{-69}$.

4. Key length equals 256 bits: system of base numbers includes 16 polynomials of degree 16. $S=16$. $p_{kr} \approx 10^{-135}$.

Cryptostrength of AES standard for the keys of length 128 and 256 bits is $2^{-128} \approx 10^{-38}$ and $2^{-256} \approx 10^{-77}$, respectively. Cryptostrength of encryption algorithms is also by tens of orders greater (examples 3 and 4).

The State Standard of the Republic of Kazakhstan ST RK 1073-2007 specifies the 1st, 2nd, 3rd and 4th security levels for the means of cryptographic protection of information. Key length of symmetric algorithms for these levels should be at least 60, 100, 150 and 200 bits respectively [10]. Minimum cryptostrength values for the keys of 100 and 200 bits equal to $2^{-100} \approx 10^{-29}$ and $2^{-200} \approx 10^{-60}$, respectively. As is seen from examples 1 and 2, the cryptostrength of nonconventional encryption is by tens of orders greater.

Thus, use of NPNs in creation of symmetric encryption algorithms help to achieve the required levels of reliability specified by the Standard ST RK 1073-2007 with significantly shorter secret key lengths. Nonpositional nature of notations also helps to provide high performance and prevent propagation of errors.

IV. MODELING OF THE SYSTEM OF NONCONVENTIONAL BLOCK ENCRYPTION

Developed nonconventional encryption algorithm is the basis for solving the problems of cryptography.

For the purpose of its practical application, the scientific research is carried out on the development of:

- modified algorithm based on a Feistel network to improve the statistical characteristics of the nonpositional cryptogram (10) - (13);
- models of operation modes of the modified nonconventional encryption algorithm are performed.

A. Modification using the Feistel encryption scheme

If the length of the full key encryption in NPNs larger, then there are more choices of systems of working base numbers. In this regard, one can use several models of Feistel scheme.

Models could differ by the number of sub-blocks as well as by the number of rounds (or iterations). The functions of cryptographic transformation of sub-blocks in scheme models may also be different.

The input data block may be divided into different even number of sub-blocks according to its length. On each step of the iteration the possible variants of key sequences using and systems of working base numbers will be researched.

In computer modeling, the cryptostrength of the developed modified algorithms will be analyzed.

B. Operating modes of the modified block cipher

There is a potential possibility of information leaks about recurring parts of data which encrypted on the one and the same key, in view of the fact that the block ciphers encrypt data by fixed-size blocks [2]-[6], [11]. Therefore, for using block cipher algorithms various modes are developed [12]. Encryption modes in the process of cryptographic transformations are used to provide the required conditions for encrypted messages. The main condition is that the encryption result of each block must be unique regardless of the encrypted data.

It is supposed to consider one of the cipher operation modes models - the Cipher Block Chaining (CBC) mode. In CBC mode, each block of plaintext is XORed with the previous block of the cryptogram and then the result is being encrypted. This way, each ciphertext block depends on all plaintext blocks processed up to that point. Moreover, for computing the first ciphertext block is using a random initialization vector (IV). It is necessary to guarantee the uniqueness of IVs for each encryption in order to identical two messages is not encrypted identically. The message that encrypted in CBC mode can only sequentially decrypt, starting with the first block. IVs are also needed to decrypt the data.

The modification of the CBC mode is Propagating cipher-block chaining (PCBC). The main difference of this mode from the CBC is that changes in the ciphertext propagate to all blocks when decrypting, as well as when encrypting. Changing one bit of plaintext affects all subsequent blocks of ciphertext. Distortion of one bit in the cryptogram leads to damage of all subsequent plaintext blocks.

Software implementation of the proposed models of the CBC mode allows choosing the required operating mode of the modified algorithm based on NPNs.

V. ASYMMETRIC SYSTEM OF DIGITAL SIGNATURE BASED ON NPNs

The ElGamal digital signature (DS) scheme is based on the complexity of the problem of computing discrete logarithms in the finite field [13]-[14]. On the basis of this scheme the standards of digital signature DSS (Digital Signature Standard, USA, 1994) and GOST R 34.10-94 (Russian, 1994) are constructed [15]-[16]. Standard DSS based on the hashing algorithm SHA and formation algorithm of the digital signatures DSA (Digital Signature Algorithm). This algorithm

has been accepted in 1994 as the USA standard of digital signature and is the variation of a digital signature of the ElGamal scheme and K. Schnorr. The length of the signature in DSA system is 320 bits.

DSA algorithm is a "classic" example of DS scheme based on the using of hash functions and asymmetric encryption algorithm. The strength of the system in general depends on complexity of finding discrete logarithms in the finite field.

The essence of DSA electronic signature scheme is the following.

Let sender and recipient of the electronic document in computation of digital signature use large prime integers p and q : $2^{L-1} < p < 2^L$, $512 \leq L \leq 1024$, L multiple of 64, $2^{159} < q < 2^{160}$, q - prime divisor of $(p-1)$ and $g = h^{(p-1)/q} \bmod p$, where h - arbitrary integer, $1 < h < p-1$ such that $h^{\frac{p-1}{q}} \bmod p > 1$.

Key b is randomly selected from the range $1 \leq b \leq q$ and keeping in secret. Calculated value $\beta = g^b \bmod p$. The algorithm parameters p, q, g are the public key and published for all users of the information exchange system with DS.

Consider the formation of the DS for the message M .

Determine hash value h from the signed message M : $h = h(M)$.

Choose integer r by some random method, where $1 \leq r \leq q$. This number stored in secret and varies for each signature.

Calculate: $\gamma = (g^r \bmod p) \bmod q$.

By using the private key of the sender $\delta = (r^{-1}(h + b\gamma)) \bmod q$ is calculated, where r^{-1} satisfies the condition $(r^{-1}r) \bmod q = 1$.

Digital signature for the message M is the pair of numbers (γ, δ) , which passed along with the message by open communication channels.

Verification of DS. Let denote M', δ', γ' obtained by the addressee version of M, δ, γ .

Checking the conditions $0 < \delta < q$ and $0 < \gamma < q$. Reject the signature if any one of the conditions of the digital signature is not satisfied.

Calculate hash value $h_1 = h(M')$ from the received message M' .

Calculate value $v = (\delta')^{-1} \bmod q$.

Calculate the expressions: $z_1 = (h_1 v) \bmod q$ and $z_2 = (\gamma' v) \bmod q$.

Calculate value: $u = ((g^{z_1} \beta^{z_2}) \bmod p) \bmod q$.

The DS is valid if $\gamma' = u$, i.e. in the transfer process the integrity of the message was not compromised: $M' = M$. At

default of equality DS is invalid.

Cryptostrength of DSA scheme against "brute force" attacks is primarily dependent on the size of the parameters p and q . Accordingly, cryptostrength against "brute force" attacks on the parameter p in case of 512 and 160 bits is equal 2^{160} . A successful attack on the parameter q is only possible, if the attacker can calculate discrete logarithms in Galois field $GF(2^{512})$.

One of the theoretically possible attacks on DSA scheme is a compromise of the parameter r . For each signature is required a new value of r , which should be chosen randomly. If the attacker finds the value of r , then the secret key b may be disclosed. Another possible embodiment - two signatures were generated on the same value of r . In this case, the attacker is also able to recover b . Consequently, one of the factors that increase the safety of using DS schemes is the existence of a reliable random number generator.

In DSA length conversion module is approximately 1024 bits. To the same length increased key lengths. In this regard, increasing the computational complexity of cryptographic transformations, but decreases the computational speed. Reducing the key length and increasing computing speed, possible in the development of the modifying of this DS scheme on the basis of NPNs.

The modular system of DS with the public key, in creation that will be used a modified algorithm of DSA based on NPNs are be developed. Initially DSA algorithm written as, in which no number q and all calculations are performed only in one modulo p . Then developed a modification of the scheme on the basis of NPNs.

The formation process of NPNs for electronic message M of the given length N bits and calculating the hash value for this message given in Section II.

The modification of DSA digital signature scheme based on NPNs is carried out as follows.

The digital signature computation. Let formed NPNs with working base numbers $p_1(x), p_2(x), \dots, p_s(x)$. For each of the working base numbers the corresponding generating elements (polynomials) $g_1(x), g_2(x), \dots, g_s(x)$ are selected. Generating polynomials are analogous to primitive elements in finite field modulo prime number.

The sender's secret key b in the range $[1, 2^m]$ is chosen.

Calculates the value of the public key $\beta(x)$:

$$\beta(x) = (\beta_1(x), \beta_2(x), \dots, \beta_s(x)).$$

In the modified DS algorithm based on NPNs, the procedure for calculating the hash value will be used in the NPNs.

The random integer r from a range of $[1, 2^m]$ is selected.

In NPNs polynomials $\gamma(x)$ and $\delta(x)$ has nonpositional representation in the form of sequence of residues from its division by the base numbers of:

$$\gamma(x) = (\gamma_1(x), \gamma_2(x), \dots, \gamma_s(x)),$$

$$\delta(x) = (\delta_1(x), \delta_2(x), \dots, \delta_s(x)).$$

Digital signature for the message M is a pair of polynomials $(\gamma(x), \delta(x))$.

Verification of the digital signature is carried out by analogy of the given DSA verification.

Using algebraic approach based on NPNs will reduce the key length for digital signature without significantly lowering its cryptostrength.

VI. CONCLUSION

Cryptostrength of the developed modified encryption systems and digital signature based on NPNs is characterized by the full secret key. This key is dependent not only on key length (pseudorandom sequence), but also on the chosen system of polynomial bases of NPNs, and also on the number of all possible permutations of bases in the system.

Research and application of encryption modes is aimed at eliminating potential vulnerabilities in the processing of large blocks of messages. In connection with this, models, that applicate CBC cipher mode, will be considered. This mode allows to eliminate the disadvantages of using a single key for encryption all plaintext blocks without significantly reducing the speed of its capacity, as the delay in executing of XOR operation is small. The developed modified system of digital signature, based on DSA algorithm and NPNs, is characterized by improvement of the basic characteristics of the digital signature. Computer modelling of the modified cryptosystems based on NPNs will allow developing recommendations for their secure usage and generation of full secret keys.

REFERENCES

- [1] I. Ya. Akushskii, D. I. Juditskii, "Machine Arithmetic in Residue Classes [in Russian]," Moscow: Sov. Radio, 1968.
- [2] W. Stallings, "Cryptography and Network Security (4th Edition)," Prentice Hall, 2005.
- [3] R. G. Biyashev, "Development and investigation of methods of the overall increase in reliability in data exchange systems of distributed ACSs," Doctoral Dissertation in Technical Sciences, Moscow, 1985.
- [4] R. G. Biyashev, S. E. Nyssanbayeva, "Algorithm for Creation a Digital Signature with Error Detection and Correction," *Cybernetics and Systems Analysis*, 4, 489-497, 2012.
- [5] R. Biyashev, S. Nyssanbayeva, N. Kapalova, "The Key Exchange Algorithm on Basis of Modular Arithmetic," *International Conference on Electrical, Control and Automation Engineering (ECAE2013)*, Hong Kong – Monami, S. 2014. – P.501-505, December 1-2, 2013.
- [6] Gr. C. Moisil, "Algebraic Theory of Discrete Automatic Devices," [Russian translation]. Inostr. Lit., Moscow, 1963.
- [7] N. A. Kapalova, S. E. Nyssanbayeva, R. A. Khakimov, "Irreducible polynomials over the field $GF(2n)$," *Proceedings of Scientific and Technical Society "KAKHAK"*, Almaty, Kazakhstan, № 1. P. 17-28, 2013.
- [8] R. K. Nyssanbayev, "Cryptographical method on the basis of polynomial bases," *Herald of the Ministry of Science and Higher Education and National Academy of Science of the Republic of Kazakhstan*, 5, 63-65, 1999.
- [9] R. Biyashev, M. Kalimoldayev, N. Kapalova, R. Khakimov, S. Nyssanbayeva, "Program Modeling of the Cryptography Algorithms on Basis of Polynomial Modular Arithmetic," *Proceedings. The 5th International Multi-Conference on Complexity, Informatics, and*

Cybernetics. The 5th International Conference on Society and Information Technologies (IMCIC'14 - ICSIT 2014). – Orlando, Florida, U.S.A. 2014. – P. 49-54.

- [10] ST RK 1073-2007 "Means of cryptographic protection of information. General technical requirements", Astana: 2009.
- [11] N. Ferguson, B. Schneier, T. Kohno, "Cryptography Engineering: Design Principles and Practical Applications," Wiley Publishing Inc, 2010.
- [12] Recommendation for Block Cipher Modes of Operation. NIST Special Publication 800-38A. Technology Administration U.S. Department of Commerce. 2001 Edition.
- [13] W. Diffie, M. Hellman, "Privacy and Authentication: An Introduction to Cryptography," Proc. of the IEEE [Russian Translation]. 3, 71–109, 1979.
- [14] T. ElGamal, "A Public-Key Cryptosystem and a Signature Scheme Based on Discrete Logarithms," *IEEE Transactions on Information Theory*, v. IT-31, n. 4, 1985. P. 469-472.
- [15] FIPS PUB 186. Digital Signature Standard (DSS).
- [16] Information technology. Cryptographic protection of information. Hash function GOST 4.11-94, State Standard of the Russian Federation, Moscow, 1994. Available: <ftp://ftp.wtc-ural.ru/pub/ru.crypt/> GOCT 34.11/: 10.01.2015.

Experimental Human Machine Interface system Based on vowel and short words Recognition

Mohamed FEZARI¹, Ali Al-Dahoud²

¹Badji Mokhtar Annaba University, Faculty of Engineering, BP: 12, Annaba, 23000, Algeria

²Faculty of IT University of Al-Zaytoonah Amman, Jordan

e-mail : mohamed.fezari@univ-annaba.dz, aldahoud@zuj.edu.jo

Abstract— HCI (Human Computer Interfaces) applications are generally based on using mouse , keyboard or joystick; in this paper we experimented the use of vowels to activate some input device such as to control the movement of mouse pointer on the screen. The control of the windows icon mouse pointer (WIMP) by voice command is currently based on using vowel utterances, this category of letters is easy to recognize and to be pronounced, especially for individuals who are physically disabled or have a partial voice disorder. So this type of MCI might be used by a category of disabled person. In addition, vowels are quite easy to model by automatic speech recognition (ASR) systems. In this work we represent the design of a system for the control of mouse cursor based on voice command, using the pronunciation of certain vowels and short words. The Mel Frequency Cepstral Coefficients (MFCCs), fundamental frequency (F_0) and Formants (F_1, F_2) are selected as features. The TDW with Euclidian Distance and Hidden Markov Models (HMMs) have been tested as classifiers for matching components (vowels and short words). Comparison between different features and classifiers were tested and results are presented on tables, finally a GUI has been designed for user applications.

Keywords- Human machine communication; vowels; windows icons mouse pointer; MFCC; DTW; HMMs.

I. INTRODUCTION

Existing human-computer interfaces are not suited to individuals with upper limb motor impairments. Recently, a lot of interest is put on improving all aspects of the interaction between human and computer especially for this category of persons, however these devices are generally more expensive example sip-and-switches [1] eye-gas and eye tracking devices[4] , head mice [2,3] chin joystick[5] and tongue switches [6]. Here is some related works on human computer interaction, based on voice activation or control, which can be invested for individuals with motor impairments. Most of concepts of vocal commands are built on the pronunciation of vowels [5, 6, and 7], where the particularity of vowels used is the simple and the regular pronunciation of these phonemes. Many vocal characteristics are exploited in several works, but the most used are: energy [1, 2, 3 and 5], pitch and vowel quality [9,10] speech rate (number of syllables per second) and volume level [7]. However, Mel Frequency Cepstral Coefficients (MFCCs)

[11, 12 and 13] are used significantly of speech processing as bio-inspired feature for automatic speech recognition of isolated words [15-16].

The paper is organized as follows: in section 2, presentation of an overview on related works of mouse cursor control based on voice control and commands. In section 3, we showed LPC and MFCC computation and use as features extraction techniques. Then we describe used classifiers: DTW then HMM in section 4. In section 5, we present tests and results. And finally, we provide graphic user interface as an application.

II. RELATED WORKS:

We describe some related works with vocal command system in the literature review. Voice recognition allows you to provide input to an application with your voice. In the basic protocol, each vowel is associated to one direction for pointer motion [1]. This technique is useful in situations where the user cannot use his or her hands for controlling applications because of permanent physical disability or temporal task-induced disability. The limitation of this technique is that it requires an unnatural way of using the voice [5] [6]. Control by Continuous Voice: In this interface, the user's voice works as an on/off button. When the user is continuously producing vocal sound, the system responds as if the button is being pressed. When the user stops the sound, the system recognizes that the button is released. For example, one can say "Volume up, ahhhhhh", and the volume of a TV set continues to increase while the "ahhh" continues. The advantage of this technique compared with traditional approach of saying "Volume up twenty" or something is that the user can continuously observe the immediate feedback during the interaction. One can also use voiceless, breathed sound [6].

Alex Olwal et al. [7] have been experimenting with non verbal features in a prototype system in which the cursor speed and direction are controlled by speech commands. In one approach, speech commands provide the direction (right, left, up and down) and speech rate controls the cursor speed. Mapping speech rate to cursor speed is easy to understand and allows the user to execute slow. The cursor's speed can be changed while it is moving, by reissuing the command at a different pace. One limitation of using speech features is

that they are normally used to convey emotion, rather than for interaction control.

The detection of gestures is based on discrete pre-designated symbol sets, which are manually labeled during the training phase. The gesture-speech correlation is modeled by examining the co-occurring speech and gesture patterns. This correlation can be used to fuse gesture and speech modalities for edutainment applications (i.e. video games, 3-D animations) where natural gestures of talking avatars are animated from speech [7] [8].

J. Bilmes et al. [9] have been developed a portable modular library (the Vocal Joystick"VJ" engine) that can be incorporated into a variety of applications such as mouse and menu control, or robotic arm manipulation. Our design goal is to be modular, low-latency, and as computationally efficient as possible. The first of those, localized acoustic energy is used for voice activity detection, and it is normalized relatively to the current detected vowel, and is used by our mouse application to control the velocity of cursor movement. The second parameter, "pitch", is not used currently but it is left for the future use. The third parameter: "vowel quality", where the vowels are characterized by high energetic level. The classification of vowels is realized by extraction of two first formants frequencies, tongue height and tongue advancement [9, 10]. Thus, the VJ research has focused on real time extraction of continuous parameters since that is less like standard ASR technology [9]. The main advantage of VJ is the reaction of the system in real time.

In [14], Thiang et al., described the implementation of speech recognition system on a mobile robot for controlling movement of the robot. The methods used for speech recognition system are Linear Predictive Coding (LPC) and Artificial Neural Network (ANN). LPC method is used for extracting feature of a voice signal and ANN is used as the recognition method. Backpropagation method is used to train the ANN. Experimental results show that the highest recognition rate that can be achieved by this system is 91.4%. This result is obtained by using 25 samples per word, 1 hidden layer, 5 neurons for each hidden layer, and learning rate 0.1.

III. FEATURE EXTRACTION

In order to implement the HMI application on embedded system in future, and to get good results in automatic speech recognition is to select better and easy to compute features, so the features would be robust and fast to compute. The LPC, MFCC with energy and derivatives were selected based on literature reviews [15, 16] and [17].

A. MFCC Feature extraction[11]

The extraction of the best parametric representation of acoustic signals is an important task to produce a better recognition performance. The efficiency of this phase is important for the next phase since it affects its behavior. MFCC is based on human hearing perceptions which cannot perceive frequencies over 1Khz. In other words, in MFCC is based on known variation of the human ear's critical

bandwidth with frequency. MFCC has two types of filter which are spaced linearly at low frequency below 1000 Hz and logarithmic spacing above 1000Hz. A subjective pitch is present on Mel Frequency Scale to capture important characteristic of phonetic in speech. The overall process of the MFCC can be presented in the following steps:

1. After the pre-emphasis filter, the speech signal is first divided into fixed-size windows distributed uniformly along the signal.
2. The FFT (Fast Fourier Transform) of the frame is calculated. Then the energy is calculated by squaring the value of the FFT. The energy is then passed through each filter Mel. S_k : is the energy of the signal at the output of the filter K, we have now m_p (number of filters) S_k parameters.
3. The logarithm of S_k is calculated.
4. Finally, the coefficients are calculated using the DCT (Discrete Cosine Transform).

$$c_i = \sqrt{\frac{2}{m_p}} \left\{ \sum_{k=1}^{m_p} \log(S_k) \cos \left[i \left(k - \frac{1}{2} \right) \frac{\pi}{m_p} \right] \right\} \quad (1)$$

pour $i = 1 \dots \dots N$

N: is the number of MFCC coefficients.

B. Fundamental Frequency and formants extraction

Linear predictive analysis of speech has become the predominant technique for estimating the basic parameters of speech. Linear predictive analysis provides both an accurate estimate of the speech parameters and also an efficient computational model of speech.

The basic idea behind linear predictive analysis is that a specific speech sample at the current time can be approximated as a linear combination of past speech samples. Through minimizing the sum of squared differences (over a finite interval) between the actual speech samples and linear predicted values a unique set of parameters or predictor coefficients can be determined.

LPC computation basic steps can be presented as follow [14]:

a) *Pre-emphasis*: The digitized speech signal, $s(n)$, is put through a low order digital system, to spectrally flatten the signal and to make it less susceptible to finite precision effects later in the signal processing.

b) *Frame Blocking*: The output of pre-emphasis step $\tilde{s}(n)$ is blocked into frames of N samples, with adjacent frames being separated by M samples. If $x_l(n)$ is the l^{th} frame of speech, and there are L frames within entire speech signal.

c) *Windowing*: After frame blocking, the next step is to window each individual frame so as to minimize the signal discontinuities at the beginning and end of each frame. If we define the window as $w(n)$, $0 \leq n \leq N-1$, then the result of windowing is the signal:

$$\tilde{x}_l(n) = x_l(n)w(n) \quad (2)$$

d) *Autocorrelation Analysis*: The next step is to auto correlate each frame of windowed signal in order to give:

$$r_l(m) = \sum_{n=0}^{N-1-m} \tilde{x}_l(n)\tilde{x}_l(n+m) \quad (3)$$

$$m = 0, 1, \dots, p$$

e) *LPC Analysis*: which converts each frame of $p + 1$ autocorrelations into LPC parameter set by using Durbin's method.

f) *LPC Parameter Conversion to Cepstral Coefficients*: LPC cepstral coefficients, is a very important LPC parameter set, which can be derived directly from the LPC coefficient set. The recursion used is:

$$c_m = a_m + \sum_{k=1}^{m-1} \left(\frac{k}{m}\right) \cdot c_k \cdot a_{m-k} \quad 1 \leq m \leq p$$

And:

$$c_m = \sum_{k=m-p}^{m-1} \left(\frac{k}{m}\right) \cdot c_k \cdot a_{m-k} \quad (5)$$

$$m > p$$

The LPC cepstral coefficients are the features that are extracted from voice signal and these coefficients are used as the input data for the classifier (Euclidian Distance or DTW). In this system, voice signal is sampled using sampling frequency of 8 kHz and the signal is sampled within 1.5 seconds, therefore, the sampling process results 1200 data. Because we choose LPC parameter $N = 200$, $m = 100$, and LPC order = 10 then there are 119 vector data of LPC cepstral coefficients.

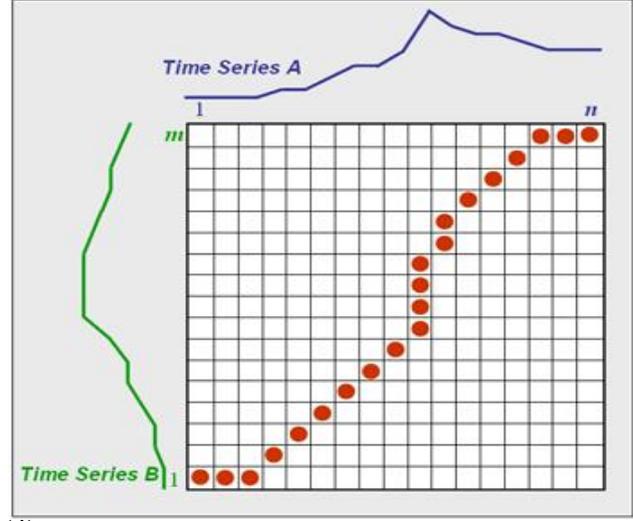
IV. CLASSIFIERS

In pattern recognition in general, automatic speech recognition, speaker Identification, image or shape recognition we need some how an algorithm to classify.

A. DTW(Dynamic Time Warping)

DTW algorithm is based on Dynamic Programming techniques .This algorithm is for measuring similarity between two time series which may vary in time or speed. This technique also used to find the optimal alignment between two times series if one time series may be "warped" non-linearly by stretching or shrinking it along its time axis.

This warping between two time series can then be used to find corresponding regions between the two time series or to determine the similarity between the two time series.



(4)

Fig. 1. The optimal warping path from [22]

B. Euclidian distance formulat:

The **Euclidean distance** between points \mathbf{p} and \mathbf{q} is the length of the line segment connecting them ($\overline{\mathbf{pq}}$).

In Cartesian coordinates, if $\mathbf{p} = (p_1, p_2, \dots, p_n)$ and $\mathbf{q} = (q_1, q_2, \dots, q_n)$ are two points in Euclidean n -space, then the distance (d) from \mathbf{p} to \mathbf{q} , or from \mathbf{q} to \mathbf{p} is given by the Pythagorean formula:

$$\text{dist}((x, y), (a, b)) = \sqrt{(x - a)^2 + (y - b)^2} \quad (6)$$

$$= \sqrt{\sum_{i=1}^n (q_i - p_i)^2}.$$

(7)

The position of a point in a Euclidean n -space is a Euclidean vector. So, \mathbf{p} and \mathbf{q} are Euclidean vectors, starting from the origin of the space, and their tips indicate two points. The **Euclidean norm**, or **Euclidean length**, or **magnitude** of a vector measures the length of the vector:

$$\|\mathbf{p}\| = \sqrt{p_1^2 + p_2^2 + \dots + p_n^2} = \sqrt{\mathbf{p} \cdot \mathbf{p}} \quad (8)$$

where the last equation involves the dot product.

A vector can be described as a directed line segment from the origin of the Euclidean space (vector tail), to a point in that space (vector tip). If we consider that its length is actually the distance from its tail to its tip, it becomes clear that the Euclidean norm of a vector is just a special case of Euclidean distance: the Euclidean distance between its tail and its tip.

The distance between points \mathbf{p} and \mathbf{q} may have a direction (e.g. from \mathbf{p} to \mathbf{q}), so it may be represented by another vector, given by

$$\mathbf{q} - \mathbf{p} = (q_1 - p_1, q_2 - p_2, \dots, q_n - p_n) \quad (9)$$

If $D(x,y)$ is the Euclidean distance between frame x of the speech sample and frame y of the reference template, and if $C(x,y)$ is the cumulative score along an optimal alignment path that leads to (x,y) , then:

$$C(x,y)=\text{MIN}(C(x-1,y),C(x-1,y-1),C(x,y-1))+D(x,y) \quad (10)$$

C. HMMs Basics [13]

Over the past years, Hidden Markov Models have been widely applied in several models like pattern, or speech recognition. To use a HMM, we need a training phase and a test phase. For the training stage, we usually work with the Baum-Welch algorithm to estimate the parameters (π,A,B) for the HMM. This method is based on the maximum likelihood criterion. To compute the most probable state sequence, the Viterbi algorithm is the most suitable.

An HMM model is basically a stochastic finite state automaton, which generates an observation string, that is, the sequence of observation vectors, $O = O_1, \dots, O_t, \dots, O_T$. Thus, a HMM model consists of a number of N states $S=\{S_i\}$ and of the observation string produced as a result of emitting a vector 'Ot' for each successive transitions from one state S_i to a state S_j . 'Ot' is d dimension and in the discrete case takes its values in a library of M symbols.

The state transition probability distribution between state S_i to S_j is $A=\{a_{ij}\}$, and the observation probability distribution of emitting any vector 'Ot' at state S_j is given by $B=\{b_j(O_t)\}$. The probability distribution of initial state is $\Pi=\{\pi_i\}$.

$$a_{ij} = P(q_{t+1} = \frac{S_j}{q_t} = S_j) \quad (11)$$

$$B = \{b_j(O_t)\} \quad (12)$$

$$\pi_i = P(q_0 = S_i) \quad (13)$$

Given an observation O and a HMM model $\lambda=(A,B,\Pi)$, the probability of the observed sequence by the forward-backward procedure $P(O/\lambda)$ can be computed. Consequently, the forward variable is defined as the probability of the partial observation sequence O_1, O_2, \dots, O_t (until time t) and the state S at time t , with the model λ as $\alpha(i)$. and the backward variable is defined as the probability of the partial observation sequence from $t+1$ to the end, given state S at time t and the model λ as $\beta(i)$. The probability of the observation sequence is computed as follow:

$$p(o/\lambda) = \sum_{i=1}^N \alpha_t(i) * \beta_t(i) = \sum_{i=1}^N \alpha_T(i) \quad (14)$$

And the probability of being in state I at time t , given the observation sequence O and the model λ is computed as follow:

$$\pi_i = P(q_0 = S_i) \quad (15)$$

V. DESCRIPTION OF APPLICATION

The application is designed to control the mouse cursor by using the pronunciation of certain phonemes and words, which we chose as vocabulary: "aaa", "ooh", "iii", "eeu", "ou", "uu", "Clic" and "stop".

The choice of these vowels and short words is based on the following criteria:

- Easy to learn.
- Easy to pronounce.
- can be pronounced persons with voice disorder.
- Easy to recognize by automatic speech recognition system.

A. DataBase Description

The database consists of 10 women (age 20 to 50 years), 10 men (age 20 to 60 years), and 5 children (age from 5 to 14 years) and category of persons with voice disorder from German database of the PTSD Putzer's voice in [18], each speaker had: 5 trials for each phoneme or word. Collection of the database is performed in a quiet room without noise.

B. The parametrization

According to the tests, we found that the parameters more robust to noise than other parameters are the LPC coefficients and Mel Frequency Cepstral Coefficients (MFCCs).

The input signal is segmented by a window of 25 ms overlapping 10ms, from each segment parameters were extracted by both methods LPC (the order of the prediction: 10) then MFCC (42 coefficients: Energy and derivative and second derivatives).

C. Classification

For this moment, we have tested two classifier, first one has been used for simplicity in order to be implemented in future on DSP circuit of microcontroller: Dynamic Time Warping (DTW) with Euclidian distance and Hidden Markov chains (HMM) for classification phase.

For Hidden Markov models, in our system, we utilize left-to-right HMM structures with 3 states and 3 mixtures are used to model MFCCs coefficients.

D. Application

Our application is used to control the mouse cursor by voice, pronouncing a vowel or short words above. The vowels are mapped to directions of movement cursor and push buttons on mouse as follow and presented in figure 2:

- Up:" ooh"
- Down:" aah"
- To the right:" iii"
- Left:" eeu"
- To double-click (open):" click" or "eke"

- To exit the application by voice command:" stop" or"abe".
- Left-Click : "ou"
- Right-Click:" "uu"

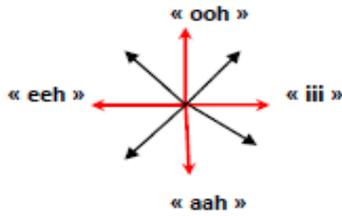


Fig. 2. Directions of cursor mouse mapping from vowels

VI. RESULTS AND DISCUSSIONS

For the testing phase, 20% of recorded sounds are selected for each vowel or short word from the vocabulary.

In order to see the effect of training and making the system speaker independent, different scenarios for the tests were done, where we choose the results of recognition of three users out of database.

Some vowels and short words were correctly classified with some confusion, where a phoneme (or word) test classified as another phoneme (or word), the misclassification is presented in the tables below (I, II). And it is clear that the confusion is higher in LPC features with DTW classifier while it is reduced using MFCC with HMM classifier.

TABLE I. CONFUSION TABLE USING (MFCC/HMM)

Pronounced Vowel	Classified as:						
	aaa	ooh	eeu	iii	clic	stop	ou
aaa	o	x	-	-	-	x	-
ooh	x	o	-	-	-	x	-
eeu	-	x	o	x	-	-	x
iii	-	-	-	o	x	-	-
ou	-	x	x	-	-	-	o
uu	-	-	x	-	-	-	x
Clic or "eke"	-	-	-	-	o	x	-
Stop or "ebe"	-	x	-	-	-	o	-

x: means that pronounced phoneme classified as an other

TABLE II. CONFUSION TABLE USING (LPC/DTW)

Pronounced Vowel	Classified as:						
	aaa	ooh	eeu	iii	clic	stop	ou
aaa	o	x	-	-	x	x	-
ooh	x	o	x	-	-	-	x
eeu	x	x	o	x	-	-	x
iii	x	-	x	o	x	-	-
ou	-	x	-	-	-	x	o
uu	-	x	x	-	-	-	x
Clic or "eke"	-	-	x	x	o	-	-
Stop or "ebe"	-	x	x	-	x	o	x

x: means that pronounced phoneme classified as an other

TABLE III. CLASSIFICATION USING LPCs, MFCC AND DTW AS CLASSIFIER

Vowel	LPC (%)	MFCC (%)
aaa	76	81
ooh	58.33	62
eeu	57	59
iii	61	73
Clic or "eke"	54.55	79
Stop or "ebe"	55.56	81

TABLE IV. CLASSIFICATION USING LPC, MFCCS AND HMM AS CLASSIFIER

Vowel	LPC (%)	MFCC (%)
aaa	85	92
ooh	78	83
eeu	74	84
iii	79	87
clic	87	94
stop	90	95

According to the results presented above (Tables: III, IV), the recognition rates using MFCCs parameterization classification with DTW or HMM classification is better than: LPCs and MFCCs with DTW. So we can say that the MFCCs / HMM system is partially independent of the speaker.

Results using MFCC and HMM, on German database vowels (sounds) for persons with chronic inflammation of the larynx and vocal fold nodules[19], are presented in Table V.

TABLE V. CLASSIFICATION USING LPC AND MFCC USING HMM FOR VOWEL FORM GERMAN DB [18]

Vowel	LPC (%)	MFCC (%)
aaa	55	67
ooh	42	53
eeu	43	49
iii	53	72
Clic or "eke"	57	64
Stop or "ebe"	54	70

We can see that the recognition rate is little bit lower than for healthy persons, we conclude that in this case other special features might be necessary to include on the application.

In addition, we must consider the preprocessing for noise in future work, as well as the database training models need from the category of children.

CONCLUSION

According to the results, we note that the classification using HMM is better than the DTW, and the decision based on MFCC coefficients is more certain than the coefficients LPCs.

From experimental results, it can be concluded that MFCC features and HMM as classifier can recognize the speech signal well. Where the highest recognition rate that can be achieved in the last scenario. This result is achieved by using MFCCs and HMM. Moreover, we need to get better features to improve classification of vowel and short words pronounced from voice disabled persons; in fact this can be resolved by inserting Jitter and Shimmer as features.

We notified that the variety of signals, collected for database from different age and gender, the recording conditions and the environment, have a considerable impact in classification results.

REFERENCES

[1] "Pride mobility products group sip-n-puff system/head array control", 2005, <http://pridemobility.com>

[2] "origine instruments sip/puff switch and head mouse", 2005, orin.com/access/headmouse/index.com

[3] "Headmaster head mouse", 2003, <http://wati.com/headmaster.htm>

[4] "assistive technologies's eye gaze system for computer access", 2003, <http://www.assistivetechologies.com/proddetails/EG001B.htm>

[5] C. de Mauro, M. Gori, M. Maggini, and E. Martinelli, "A voice device with an application-adapted protocol for Microsoft windows," In Proc. IEEE Int. Conf. on Multimedia Comp. and Systems, vol. 2, pp. 1015-1016, Firenze, Italy, 1999.

[6] T. Igarashi and J. F. Hughes, "Voice as sound: Using non-verbal voice input for interactive control," In ACM UIST 2001, November.

[7] Alex Olwal and Steven Feiner, "Interaction techniques using prosodic features of speech and audio localization," In IUI '05: Proc. 10th Int. Conf. on Intelligent User Interfaces. New York: NY, USA, 2005. ACM Press, pp. 284- 286.

[8] M.E. Sargin,O. Aran,A. Karpov,F. Ofli1,Y. Yasinnik,S. Wilson,E. Erzin,Y. Yemez and A.M. Tekalp, "Combined Gesture-Speech Analysis and Speech Driven Gesture Synthesis," ICME 2006 : IEEE International Conference on Multimedia and Expo, July 2006, pp: 893-896.

[9] J. Bilmes, X. Li, J. Malkin, K. Kilanski, R. Wright, K. Kirchhoff, A. Subramanya, S. Harada, J. Landay, P. Dowden, and H. Chizeck, "The vocal joystick: A voice-based human-computer interface for individuals with motor impairments," in Human Language Technology Conf. and Conf. on Empirical Methods in Natural Language Processing, Vancouver, October 2005.

[10] S. Harada, J. Landay, J. Malkin, X. Li, J. Bilmes, "The Vocal Joystick: Evaluation of Voice-based Cursor Control Techniques", *ASSETS'06*, October 2006.

[11] Lindsalwa Muda, Mumtaj Begam and I. Elamvazuthi, "Voice Recognition Algorithms using Mel Frequency Cepstral Coefficient (MFCC) and Dynamic Time Warping (DTW) Techniques", Journal of Computing, Volume 2, Issue 3, March 2010, pp : 138-143.

[12] Mahdi Shaneh and Azizollah Taheri, "Voice Command Recognition System Based on MFCC and VQ Algorithms" ,World Academy of Science, Engineering and Technology 57 2009, pp: 534-538.

[13] A Bala, A Kumar, N Birla - Anjali Bala et al., "Voice command recognition system based on MFCC and DTW International Journal of Engineering Science and Technology, Vol. 2 (12), 2010, pp :7335-7342.

[14] Thiang, S. Wijoyo, "Speech recognition using linear predictive coding and artificial neural network for controlling movement of mobile robot", International Conference on Information and Electronics Engineering IPCSIT vol.6, 2011, 179-183.

[15] C. Snani, "conception d'un system dereconnaissance de mots isolés à base de l'approchestochastique en temps réel : Application commmande vocale d'une calculatrice ", Mémoire de magister ,Institut d'electronique univ. Badji mokhtar Annaba,2004.

[16] C. HAdri, M boughazi and M fezari, "improvement of Arabic digits recognition rate based in the parameters choice", in proceedings of international conf. CISA Annaba, june 2008.

[17] M. Fezari and A. Al-dahoud, "An Approach For: Improving Voice Command processor Based On Better Features and Classifiers Selection," pp. 1-5. The 13th International Arab Conference on Information Technology ACIT'2012 Dec.10-13 ,2012.

[18] Manfred Putzer & Jacques Koreman " A german databse for a pattern for vacal fold vibration " Phonus 3, Institute of Phonetics, University of the Saarland, 1997, 143-153.

[19] I.M. M. El Emary,M. Fezari, F. Amara," Towards Developing a Voice Pathologies Detection System", in Jouranal of Electronics and Communication, 2014 Elsevier.

[20] Eamonn J. Keogh, Michael J. Pazzani, "Derivative Dynamic Time Warping" In Proc. Of the 1st SIAM Int.Conf. on Data Mining (SDM-2001).

[21] H. Sakoe, S. Chiba, "Dynamic programming algorithm optimization for spoken word recognition". IEEE Transaction on Acoustics, Speech and Signal Processing, Vol 26, NO1, pp. 43-49. February 1978.

[22] <http://cst.tu-plovdiv.bg/bi/DTWimpute/DTWalgorithm.html>

On DC/DC voltage buck converter control improvement through the QFT approach

Luis Ibarra, Israel Macías, Pedro Ponce, and Arturo Molina

Abstract—The use of DC/DC voltage converters is widespread and has been studied and improved for a long time. Its nonlinearities and uncertainties have increased the complexity of associated controllers so desired performance is achieved. However, cascade-PI controllers are still used to address this problem due to its relative implementation easiness, relegating the complexity to the tuning strategy. This paper presents the parallel design of a cascade PI controller through LQR tuning method, voltage mode QFT, and current mode QFT, offering comparative conclusions and showing that QFT approach surpasses PI performance.

Index Terms—Robust control, QFT, LQR, Voltage converter, Buck

I. INTRODUCTION

A DC/DC converter is a power electronics circuit used to modify the output characteristics of a DC source (voltage, impedance) at high efficiency and stable operation. The first registered patent is from 1978 [1]. A voltage converter is a time variant system as its dynamical behavior depends on a switch controlled through PWM; moreover, the relation between the PWM duty cycle and the output voltage is not linear.

Besides DC/DC converters have been successfully controlled in the past, it was until 90s when its non-linear characteristics were formally discussed and some advanced control techniques were used to improve their performance. The control objectives have been met before the system was thoroughly understood as stated in [2]. However, the convenience of modeling them in a simplified manner has made researchers also to follow this path despite of the need of two different models dependent on current conditions: continuous and discontinuous conduction modes (CCM and DCM). One of the most used methods to achieve linear representation of voltage converters is the Small-signal state-space averaging [2].

Although simpler linear models allow the designer to consider well-known frequency-domain constraints and design techniques, its validity is restricted to a determined bandwidth and can not attain non-linear behavior; as the linear model is desired to be kept simple, the control loop complexity must be increased through a more dependable controller [3]. This has lead to an increasing number of works related to

control implementation under parametric variations, uncertain environments, and ambiguous measurements which commonly adopts one single control technique and a determined set of tests to validate converter's performance.

The most commonly used control schemes are voltage mode control and current mode control [4]. The former takes the output voltage as its only feedback signal; however, its performance degrades on DCM. The current mode control effectively alleviates the sensitivity of the converter dynamics and could offer near uniform loop gain characteristics for both CCM and DCM operation. The key feature of current-mode control is that the inner loop changes the inductor into a voltage-dependent current source at frequencies lower than crossover frequency of the current loop.

A commonly used way to implement a current mode control is using two Proportional Integral (PI) controllers; one for the inner current loop and one more for the outer voltage one. In this paper, the LQR approach is employed to tune it. The algorithm proposed for PI/PID controller tuning via LQR approach and selection criteria of the Q and R matrices were taken from [5], [6]; this approach aims to control PWM-type switching DC-DC converters independently from their circuit topologies and open-loop pole-zero locations [5].

QFT (Quantitative feedback theory) approach, allows the designer to quantify how demanding a set of plants are in terms of further control design, to deal with uncertainties and disturbances, and to set the problem in commonly used frequency-domain equations [7]. These characteristics seem to suit perfectly to a linear model as the one aforementioned. QFT technique has been effectively used for voltage converters control as in [8]–[12]; however, its use is not popular and little literature can be found about specific problems. Robustness is commonly addressed through fuzzy and sliding-mode approaches [3], [13], [14].

A similar comparison has been previously presented in [15]; however, the design process is not completely presented or explained, and the PI tuning is made mostly arbitrarily. In addition, controllers are designed only for voltage mode so no conclusions about dynamical tracking, output ripple, and rising times are offered beyond overshoot comparison. This paper tries to cover those issues on larger extension and providing enough arguments to effectively expose QFT as a valid and better controller for voltage converter systems.

II. SMALL-SIGNAL STATE-SPACE AVERAGING

This method was developed by [16] and its aim is to describe the dynamics of the converter as a group of time-invariant equations which are valid for the whole commutation

This work was partially supported by a scholarship award from Tecnológico de Monterrey, Campus Ciudad de México and a scholarship for living expenses from CONACYT.

Luis Ibarra, Pedro Ponce, and Arturo Molina are with the School of Engineering and Sciences of the National School of Postgraduate Studies at Tecnológico de Monterrey, Campus Ciudad de México, Mexico City 14380 Mexico (corresponding author e-mail: ibarra.luis@itesm.mx)

Israel Macías is with the SEPI of ESIME at Instituto Politécnico Nacional, Mexico City 1000 Mexico

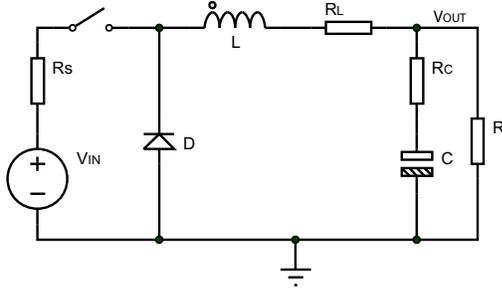


Fig. 1. Buck converter circuit

cycle. Final results are obtained based on the small-signal transfer function so their validity is restrained to relatively small voltage or load perturbations [17].

This modeling is very adequate for analysis in both stable and transient states; it is now considered an essential design tool for component selection and control objectives achievement for a particular group of specifications [17]. In order to obtain a mathematical model, a state-space representation of the circuit is achieved through (1), given that at each stage of commutation (open/closed switch) the circuit is linear and time-invariant. This means that during each time sub-interval, the system can be described by a group of differential ordinary equations which comply with the energy conservation laws [18].

The exact description of the system is obtained by averaging the state variables acquired at each state. While the averaging eliminates the variation in time for the whole commutation cycle, it does not linearize the model, making the small-signal assumptions to be necessary. This considers the circuit feedback to be disabled while a perturbation (with DC and AC components) is added at the voltage input so a frequency analysis is driven. Resulting non-linear function is approximated by Taylor series to its second term. Finally, superposition allows only the DC components to be considered equation solutions; a detailed description of this process can be found in [19].

If the preceding method is applied to the buck converter shown in Figure 1 under CCM, (1) and (2) can be derived as current and voltage transfer functions respectively, according to the duty cycle. The same technique under DCM delivers the transfer function shown in (3) and (4), again, for current and voltage.

$$G_{iD}(s) = \frac{(s(CR+Cr_C)+1)V_{IN}}{CLRs^2+As+R+r_L+r_s} \quad (1)$$

$$A = (L + CRr_L + C(Rr_s + r_C[R + Ls + r_L + r_s]))$$

$$G_{VD}(s) = \frac{V_{IN} * Rr_C(s(CR+Cr_C)+1)}{As^2+Bs+(R+r_C)(R+r_L+r_s)} \quad (2)$$

$$A = (R + r_C)(CLR + CLr_C)$$

$$B = (R + r_C)(L + CRr_C + Cr_Cr_L + Cr_Cr_s + CR(r_L + r_s))$$

 TABLE I
 PHYSICAL PARAMETERS USED FOR MODEL CALCULATION

Symbol	Description	Value
V_{IN}	Input Voltage	24V
L	Inductor	300 μ H
C	Capacitor	220 μ F
R	Load	12 Ω
T_s	PWM Period	10 μ s
D	PWM Duty cycle	1
r_s	t_{ON} switch resistance	0.01 Ω
r_L	Inductor resistance	16.3m Ω
r_C	Capacitor resistance	0.305 Ω

$$G_{iD}(s) = \frac{2V_{IN}(1/CR+s)}{CL(s^2+As+\frac{4(2-2/B)}{CDRT_sB(1-2/B)^2})} \quad (3)$$

$$G_{VD}(s) = \frac{2V_{IN}}{CL(s^2+As+\frac{4(2-2/B)}{CDRT_sB(1-2/B)^2})} \quad (4)$$

$$A = \left(\frac{1}{CR} + \frac{4}{DT_sB(1-2/B)} \right)$$

$$B = 1 + \sqrt{1 + \frac{8L}{D^2RT_s}}$$

III. PID TUNING THROUGH LQR APPROACH

LQR tuning technique was selected as it is a reliable and widely used control approach so further evaluation and comparison is possible towards a different method and conclusions can be elaborated from solid foundations. PI tuning through trial-error or parameter adjustment dependent on fractional order is completely avoided.

Using the Lyapunov's method, the LQR design problem reduces to the Algebraic Riccati Equation (ARE) which is solved to calculate the state feedback gain for a chosen set of weighing matrices that regulate the penalties due to state variables and control signal trajectories deviations. The method used to obtain Q and R matrices was genetic algorithms (GA) as used by [20] while the performance index was taken from [5], [6].

$$Q = \begin{pmatrix} 1e6 & 0 & 0 \\ 0 & 0.025e-2 & 0 \\ 0 & 0 & 1e9 \end{pmatrix} \quad (5)$$

$$R = (0.001)$$

Consider a linear process described by standard state-space representation

$$\dot{x}(t) = Ax(t) + Bu(t) \quad (6)$$

$$y(t) = Cx(t)$$

and

$$x_1 = \int e(t)dt, \quad x_2 = e(t), \quad x_3 = \frac{de(t)}{dt}. \quad (7)$$

From the block diagram of Figure 2,

$$\frac{-E(s)}{U(s)} = \frac{K}{s^2 + as + b}; \quad (8)$$

thus, equations turn into

$$[s^2 + as + b]E(s) = -Ku, \quad (9)$$

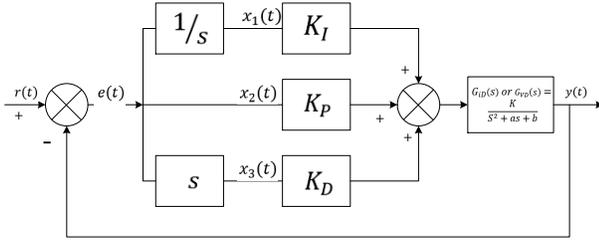


Fig. 2. PID loop controller

which can be written in the time domain as

$$\ddot{e} + a\dot{e} + be = -Ku. \quad (10)$$

Substituting (7), (10) can be rewritten as

$$\dot{x}_3 + ax_3 + bx_2 = -Ku, \quad (11)$$

so the state space formulation becomes

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \dot{x}_3 \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & -b & -a \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ -K \end{bmatrix} u. \quad (12)$$

In order to have a LQR formulation of (6) the cost function (13) is minimized.

$$J = \int_0^{\infty} [x^T(t)Qx(t) + u^T(t)Ru(t)]dt \quad (13)$$

Its result provides the state feedback control law as stated in [21]

$$u(t) = -R^{-1}B^T Px(t) = -Fx(t), \quad (14)$$

where P is the symmetric positive definite solution of the Continuous ARE:

$$A^T P + PA - PBR^{-1}B^T P + Q = 0. \quad (15)$$

The weighing matrix Q is symmetric positive semi-definite and the factor R is a positive number. If the plant's transfer functions are used on (12), A and B matrices can be derived; in addition to matrix Q from (5) and coefficient R from (6), preceding values can solve matrix P form (15) through an optimization algorithm.

If the solution for P is considered to be unique, the state feedback gain matrix becomes (17), corresponding to the optimal control signal.

$$\begin{aligned} F &= R^{-1}B^T P = R^{-1} \begin{bmatrix} 0 & 0 & -K \end{bmatrix} \begin{bmatrix} P_{11} & P_{12} & P_{13} \\ P_{21} & P_{22} & P_{23} \\ P_{31} & P_{32} & P_{33} \end{bmatrix} \\ &= -R^{-1}K \begin{bmatrix} P_{13} & P_{23} & P_{33} \end{bmatrix} \\ &= - \begin{bmatrix} -K_I & -K_P & K_D \end{bmatrix} \begin{bmatrix} x_1(t) \\ x_2(t) \\ x_3(t) \end{bmatrix} \\ &= K_I \int e(t)dt + K_P e(t) + K_D \frac{de(t)}{dt} \end{aligned} \quad (16)$$

therefore,

$$\begin{aligned} K_I &= kP_{13}/R \\ K_P &= kP_{23}/R \\ K_D &= KP_{33}/R. \end{aligned} \quad (17)$$

 TABLE II
 RESULTING PI COEFFICIENTS DUE TO LQR TUNING

Parameter	Value
Current mode	
K_P	20.8593
K_I	63244.6
Voltage mode	
K_P	2.90115
K_I	7071.06

A. Specific design for Buck converter

If the values in Table I are substituted, CCM equations are obtained as (18) and (19), while DCM ones as (20) and (21).

$$G_{iD}(s) = \frac{s + 2.95e7}{s^2 + 1415.19s + 1.47e7} \quad (18)$$

$$G_{vD}(s) = \frac{s + 3.54e8}{s^2 + 1415.19s + 1.48e7} \quad (19)$$

$$G_{iD}(s) = \frac{s + 7.57e4}{s^2 + 4.01e7s + 3.09e9} \quad (20)$$

$$G_{vD}(s) = \frac{-1.54e11}{211.77s^2 + 2.37e7s + 2.29e10} \quad (21)$$

Coefficients from preceding equations are substituted in (12) so A and B matrices of the ARE are obtained as follows:

$$A = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & -1.47e7 & -1415.19 \end{bmatrix} \quad (22)$$

$$B = \begin{bmatrix} 0 \\ 0 \\ -3.5462e8 \end{bmatrix}$$

By substituting the values in equations (22),(5), and (6) and solving (15) the following matrix P is obtained:

$$P = \begin{bmatrix} 415.81 & 0.0363 & 3.98e-7 \\ 0.036 & 14.74e-6 & 1.63e-10 \\ 3.98e-7 & 1.63e-10 & 1.43e-14 \end{bmatrix} \quad (23)$$

By using (23) on (17), the PI tuning coefficients are derived. Results are shown on Table II.

IV. QFT CONTROLLER DESIGN

This technique was first proposed by Horowitz [22] as a frequency-domain technique to analyze a given plant with uncertainties in terms of a desired frequency behavior; moreover, once the conditions to be met were established, a controller could be derived to satisfy those restrictions. It was called QFT after few years and a survey was published on 1982 [23] with referenced works and commentaries about its use.

This methodology examines the close loop (CL) effects of open loop (OL) variations; in this way, a given plant $P(s)$ assumed to be in a unitary feedback CL can be graphically forced to attain certain conditions by adding a controller $G(s)$ over the Nichols chart. This implies $L(s) = G(s)P(s)$ can be fitted so CL $T(s) = L(s)/(1 + L(s))$ fulfills magnitude low and high boundaries $B = \{b_l(s), b_h(s)\}$ obtained from the desired tracking conditions, magnitude and phase margins, crossover frequencies, and sensitivity limits. The topology considered so far is shown in Figure 3.

Systems uncertainty can be expressed as a set of n plants $P = \{P_1(s) \dots P_n(s)\}$; simultaneously controlling all P

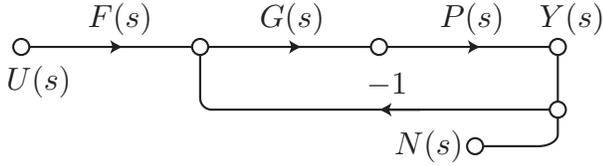


Fig. 3. Canonical form of a closed loop single-input/single-output system

implies to guarantee that all $P_i(s)$ meet the required conditions at some frequencies of interest for which the boundaries are calculated. This implies that instead of a single point over the Nichols chart, the evaluation of P at a certain frequency will provide a closed area named *plant template*. It is evident that the template can be translated but not rotated through the effect of $G(s)$ so the QFT controller design implies moving these templates so the desired conditions are attained.

In order to provide a reference for calculations a single plant within the template is taken as the *nominal plant* regardless its behavior is somehow nominal or not; the nominal category is given as it is used as pivot point to move the templates and to calculate the controller after the desired conditions are met. This nominal plant $P_A(s)$ is commonly selected to match some corner of the templates but it could be indistinctly selected among the template. For every selected frequency through all phases of interest $[-180^\circ, 0^\circ]$ an OL magnitude $|L_A|$ must be found so the whole P magnitude variations are inside CL B so a margin can be drawn on the Nichols chart with respect to $P_A(s)$.

As these margins represent the minimum point where the template can still fulfill B , they are called tracking boundaries as B was obtained precisely from step time-response wanted characteristics. Whenever the template can always remain within B for a given phase, there is no need of a margin on the Nichols Chart. However, there are other design conditions which must be faced like noise rejection, having this particular disturbance represented in Figure 3 as $N(s)$, it is clear that its gain is solely dependent on the CL behavior $T(s)$. A noise gain limit can be set as a constant so $|T_A|$ must always be set below it. The CL representation of this margin is a circle or an open concave curve on the Nichols chart within which the CL response of the template surpasses the noise gain limit, so it must be placed outside.

Depending on the specific system to be controlled an additional margin must be added; whenever the behavior of the system at very high frequencies is needed to be controlled, a template calculated at a ω_h much larger than the last selected frequency can also provide a margin called the *high-frequency boundary*. This margin is not supposed to ever be trespassed by the resulting $L(s)$; however, it is possible to find such a $L(s)$ that touches it or even passes it at $\omega \gg \omega_h$ so guaranteeing desired performance as the system will never attain such ω .

There is not an exclusive way to face the adjustment of $L_A(s)$; it depends greatly on designer's experience and in frequency-domain representation familiarity. This process is commonly called *loop shaping* and can be solved from many different points of view (Graphically, trial and error, genetic algorithms, etc.). Once the resulting $L_A(s)$ complies with

mentioned margins, the controller can be easily derived by remembering $L_A(s) = G(s)P_A(s)$. The last step in this design is to add a pre-filter $F(s)$ which adjusts the resulting $G(s)P$ into the actual boundaries; notice that the boundaries were always taken as magnitudes $\Delta|B(\omega)|$ so resulting $|T(s)|$ is surely contained (The maximum CL gain difference between plants within the template is always below B), but the actual magnitudes of their placement may differ from $b_i(s)$ and $b_s(s)$.

Pre-filter selection can be driven arbitrarily; nevertheless, a certain ω can be found so the resulting CL systems' magnitudes are found 3dB above $b_i(s)$, thus finding the cut-off frequency. Taking this as an initial value, a posterior tuning can deliver exact results.

A. Specific design for Buck converter

The uncertain model to use as input to QFT controller design is obtained through the technique described in Section II by varying the voltage input parameter from 20V to 30V and the output load from 1.2Ω to 60Ω . As many other techniques, QFT design must be aware of the two different modes the Buck converter can operate: CCM and DCM [8]. For sake of convenience, a controller for CCM operating mode will be designed based on the following uncertain transfer function (24), obtained through the variations discussed above.

$$\frac{V(s)}{D(s)} = \frac{[1.62e4, 3.03e4]s + [2.41e8, 4.52e8]}{s^2 + [1.14e3, 3.88e3]s + [1.22e7, 1.50e7]} \quad (24)$$

Equation (24) relates the input duty cycle to the output voltage; in this case the current models are ignored and a single voltage loop is assumed instead of a cascade control approach. Current control is known to alleviate sensitivity to system dynamics [24] and to provide an effective way to limit the output current [10]. Current control mode will be covered later; however, voltage mode controlled systems have also been reported like in [8], presenting good results.

According to the dynamic response of the systems contained in (24), a settling time of 1ms can be achieved so the tracking boundaries are defined to fulfill this specification as (25). As QFT approach will ensure that the magnitudes difference along the whole bandwidth will be constrained to those imposed by the boundaries, it is important to detect those frequencies at which the set of plants perform different. Figure 4 shows that the main variation occurs at 3.9krad/s, so the set of test frequencies can be chosen to be [390, 3.9k, 39k]rad/s.

$$\begin{aligned} b_h(s) &= \frac{2.95e09}{s^2 + 5.40e5s + 2.95e09} \\ b_l(s) &= \frac{1.48e12}{s^3 + 63.60e3s^2 + 5.89e08s + 1.48e12} \end{aligned} \quad (25)$$

Under the aforementioned conditions, the QFT design can proceed to find the tracking boundaries and the noise rejection boundaries by considering a threshold of 1dB. The resulting margins and the nominal plant open loop are shown in Figure 5.

In order to make the resulting behavior of $L_A(s)$ to be consistent to the margins different actions must be taken. A

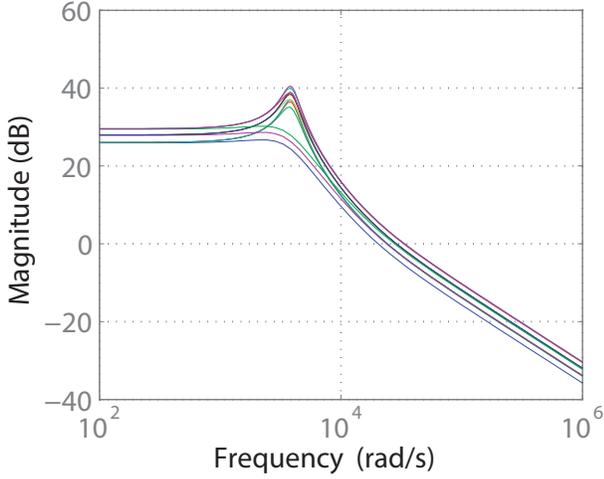


Fig. 4. Open loop Bode plot for some plants contained in (24)

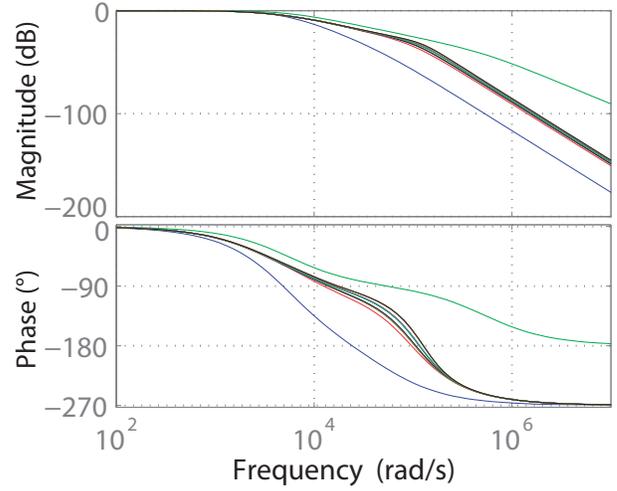


Fig. 6. Resulting Frequency response after applying QFT controller

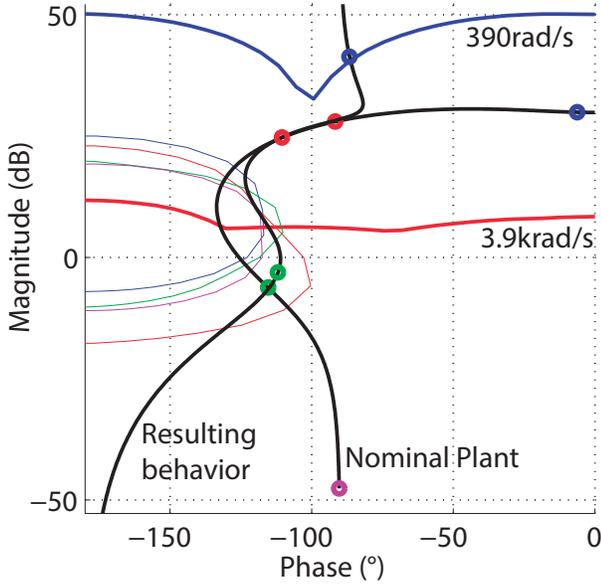


Fig. 5. Open loop and controlled nominal plant Nichols chart for CCM

pole at zero must be added to reach the lowest frequency margin along with two pole-zero pairs to border the sensitivity boundaries. The loop shaping, in this case, delivered the result shown at Figure 5 which fully meets the desired performance by using (26). Resulting plant's trajectory over Nichols chart presents a phase margin of 54.6°, and infinite magnitude margin.

$$G(s) = \frac{5.44e12s^2 + 5.44e16s + 1.142e20}{2.1e07s^3 + 3.675e12s^2 + 7.875e16s} \quad (26)$$

The last part of the design implies the incorporation of a pre-filter (27) to fit the set of plants within the desired boundaries; moreover, the correct operation of the controlled system has been evaluated towards a discrete set of frequencies, so a complete spectral view is needed to confirm a reliable design. Hence, a bode plot is obtained from the resulting loop; results are shown in Figure 6.

$$F(s) = \frac{3450}{s + 3450} \quad (27)$$

The current mode control must be derived from a set of plants which relates the control input as a duty cycle to the output current as (28). This set of plants were obtained under the same uncertain conditions than the voltage mode case based on Section II; however, the implications of these plants are totally different as they can naturally rise to its setting point in about 0.3ms. In this way, the boundaries can not be the same and must be redefined.

$$\frac{I(s)}{D(s)} = \frac{[0.66e5, 1e5]s + [0.05e8, 3.02e8]}{s^2 + [1.14e3, 3.88e3]s + [1.22e7, 1.5e7]} \quad (28)$$

Expected boundaries are constructed so a step response reaches a settling time around 0.4ms as shown in (29). The dynamical implications of this restriction change confirms the expected improvement in dynamical tracking if compared towards voltage mode controller; nevertheless, an additional controller will be needed to complete the cascade outer voltage loop. For this particular case and based on Figure 7, the selected frequencies are [100,1e3,3.9e3,10e3,30e3]rad/s.

$$b_h(s) = \frac{1.155e10}{s^2 + 1.11e06s + 1.15e10} \quad (29)$$

$$b_l(s) = \frac{9.07e12}{s^3 + 87e3s^2 + 1.65e09s + 9.07e12}$$

Considering both, the tracking boundaries imposed by (29) and assuming a noise rejection threshold of 1dB, the Nichols chart can be plotted so the original OL development of the nominal plant is analyzed (Figure 8). It can be seen that a controller similar to the one used on voltage mode control is required as the plant needs to fulfill a very similar trajectory.

After loop shaping considering a pole at zero and two pole-zero pairs, the resulting controller can be derived as (30) and its related pre-filter as (31). The resulting phase margin is of 74.3° and again, an infinite magnitude margin is achieved.

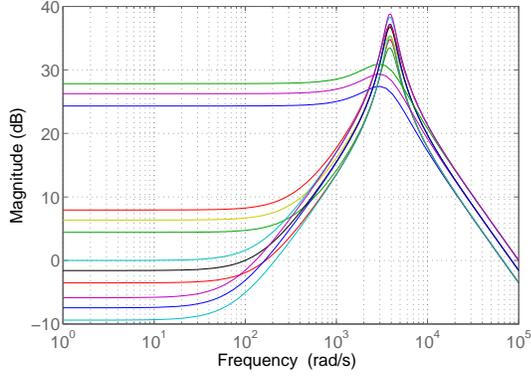


Fig. 7. Open loop Bode plot for some plants contained in (28)

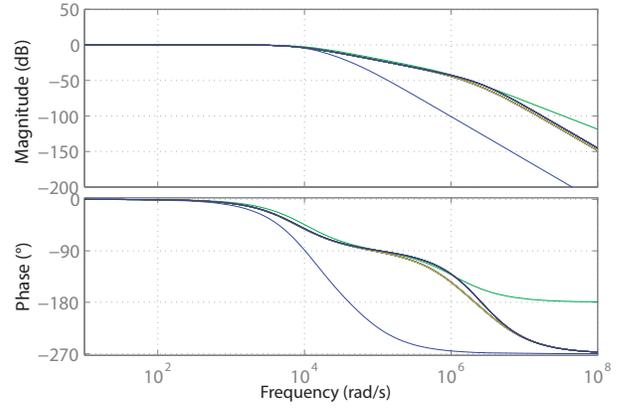


Fig. 9. Resulting frequency response after applying QFT controller

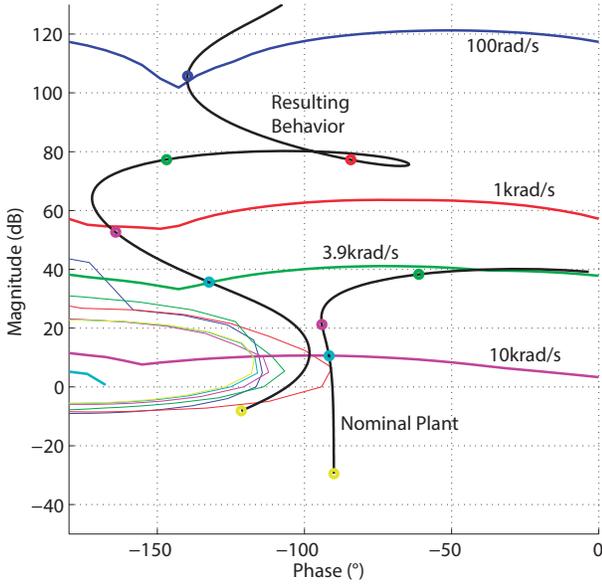


Fig. 8. Open loop and controlled nominal plant Nichols chart for CCM

$$G(s) = \frac{1.92e14s^2 + 1.82e18s + 3.36e21}{1.75e07s^3 + 1.75e13s^2 + 1.4e15s} \quad (30)$$

$$F(s) = \frac{7800}{s + 7800} \quad (31)$$

Besides more frequencies were taken this time to analyze tracking and sensitivity boundaries, a continuous frequency sweep is needed to see if the controlled loop together to the pre-filter are able to fulfill the boundary requirements. The resulting Bode plot is shown in Figure 9 where it is clear that the boundaries are respected along the whole frequency span.

V. SIMULATION RESULTS

In order to evaluate the controllers designed through QFT hitherto, a test is applied to the circuit structure shown in Figure 1. Desired output voltage is set to 10V while the input voltage and the output impedance are varied; the input voltage will be varied from 22V to 28V @ 60Ω while the load will be changed in a factor of 30 @ 22V; this is, from 60Ω to 2Ω. All

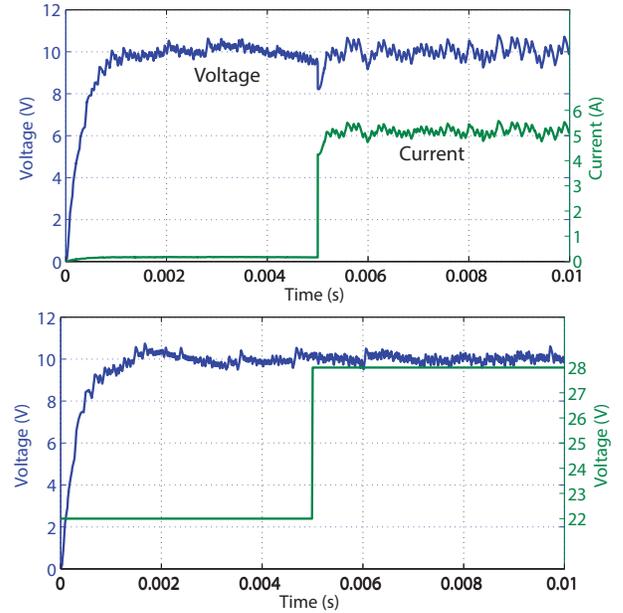


Fig. 10. Voltage mode controlled Buck converter towards load and input voltage variations

tests were programmed through SimPowerSystem library in Simulink® with simulation step-time of 1μs.

Regarding to the QFT controller process described on Section IV for CCM voltage mode, results are shown in Figure 10, where the output voltage and output current are shown. Notice that variability due to uncertainties is greatly minimized but the voltage ripple is high. However, this response shows that the effects of plant variations are properly faced and that the solution is reliable.

The current mode control was tested for the same variations but in different moments as a complete cascade voltage would be needed to do it dynamically; however, results can be effectively compared for the whole load span as shown in Figure 11 where the output voltage is shown together with inductor's current. Notice that ripple is effectively reduced as system dynamics are faced in an improved manner; having direct control of current allows the system to react faster as discussed in Section IV. The model was built considering

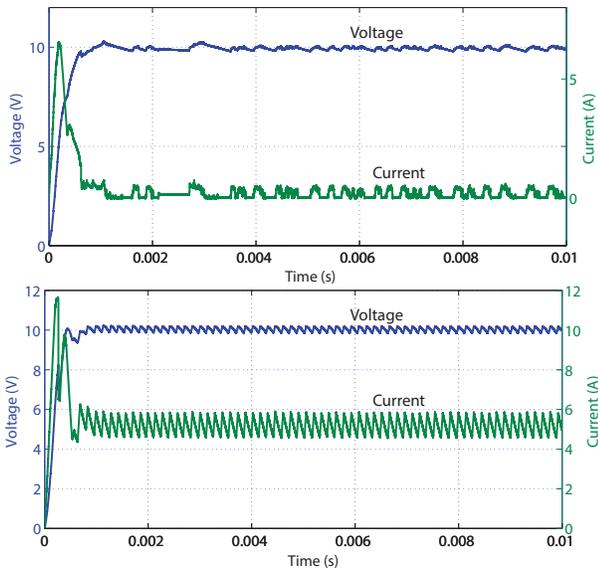


Fig. 11. Current mode controlled Buck converter at different loads: 60Ω and 2Ω

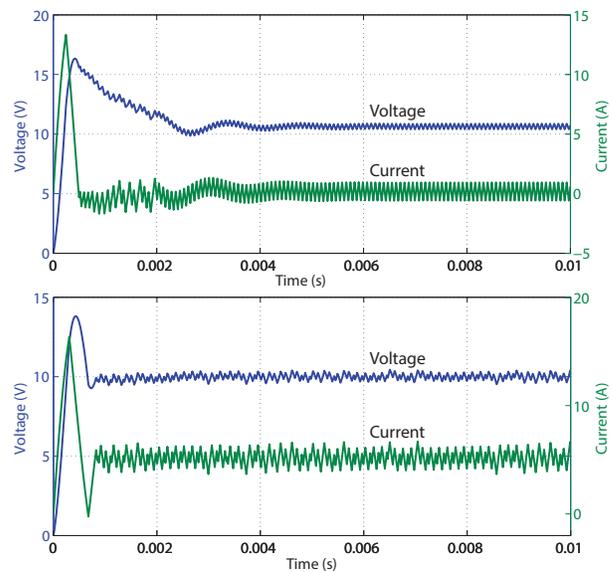


Fig. 12. LQR controller under different loads: 60Ω and 1.2Ω

CCM operation; however, it can be seen that DCM appears for the highest resistance tested. This experiment shows that QFT approach is able to partially deal with DCM even when configured for CCM operation; nonetheless, this is not a conclusive test and results must vary among converters.

LQR controller is tested slightly different as it can not guarantee the same operating points as QFT; additionally, it was designed to work properly under a specific regime so variability is reduced for sudden load changes and is kept the same on input voltage change and step response tests. The step responses for marginal loads are shown in Figure 12 where it is evident that the proper operating point of this controller is set to achieve high load conditions. Considering this fact, the test about input voltage variability is done taken this maximum load characteristic as set-point. Results are shown in Figure 13.

As told before, LQR controller is not capable of working effectively along the whole load span, so the sudden load change test is done for 12Ω and 1.2Ω and it is shown in Figure 13. Notice that all LQR tests delivered a low ripple operation and fast rising time; however, the overshoot and low settling time make this controller to be operating below the specification.

Synthetic information about results on the worst case scenario is briefly presented in Table III.

VI. DISCUSSION

QFT approach overpasses cascade PI performance by offering a robust guarantee of operation under known or expected variable conditions; however, its design can be seen as complicated for the loop shaping process which needs an additional effort to correctly place poles and zeros. This process was first automated in 1998 [25] through GA so can be neglected for complexity considerations, even more towards LQR tuning as it actually needs GA to find Q matrix.

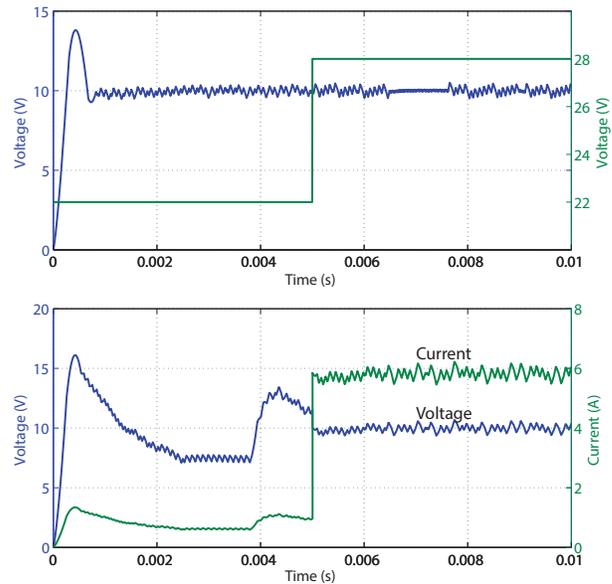


Fig. 13. LQR controller under sudden changes on input voltages and loads: 12Ω and 1.2Ω

Besides QFT voltage mode control shows a poor ripple treatment it still is able to attain variability issues and offers a very simple control topology. It has been confirmed that current mode controllers improves this subject and also permits faster responses. Although a cascade controller with inner current QFT controlled loop was not presented in this work it is a feasible solution which could achieve a better operation. The voltage outer loop could actually be built in a similar manner as the one shown here by considering the whole current closed loop as an uncertain plant and applying QFT methodology based on identified plant characteristics.

PI tuning process could embrace some robust considerations made for the QFT approach so some degree of robustness is achieved while preserving PID perspective easiness like

TABLE III
WORST SCENARIO RESULTS COMPARATIVE TABLE

Test	Magnitude	Test	Magnitude
QFT Voltage		QFT Current	
Overshoot	4%	Overshoot	0%
Ripple [max]	0.79V	Ripple [max]	0.17V
Sensitivity	$\frac{1.77V}{\Delta 30 \Omega / \Omega}$	Sensitivity	N/A
Settling time	0.8ms	Settling time	0.4ms
PI-LQR			
Overshoot	61%		
Ripple [max]	0.23V		
Sensitivity	$\frac{3V}{\Delta 10 \Omega / \Omega}$		
Settling time	2.9ms		

early proposed by [26]. Another possible solution is to embed QFT current mode controller into a cascade controller with an optimal approach so both characteristics are simultaneously achieved.

VII. CONCLUSIONS

The design process and simulated performance results of a DC/DC Buck converter are presented with detail so both, QFT and LQR-PI could be parallel compared. In addition, observations about the effectiveness of cascade control for this type of voltage converters is also provided in terms of easiness and output ripple.

QFT approach greatly overpasses LQR-PI performance in terms of overshoot, voltage ripple, sensitivity to input voltage and load, and settling time. This simulated results need further implementation to be completely validated.

If parametric or operational variability is expected, QFT controlled converters can be a better solution despite of their relative complexity; moreover, for non-varying conditions, QFT can still be considered as its overshoot and settling time characteristics are better than LQR-PI cascade controller and unexpected variations are covered.

REFERENCES

- [1] C. Lindmark, "Switched mode power supply," U.S. Patent US4097773 A, Jun., 1978, united States: US4097773 A. [Online]. Available: <http://www.google.com/patents/US4097773>
- [2] C. Tse and M. di Bernardo, "Complex behavior in switching power converters," *Proceedings of the IEEE*, vol. 90, no. 5, pp. 768–781, May 2002.
- [3] T. Gupta, R. Boudreaux, R. Nelms, and J. Hung, "Implementation of a fuzzy controller for DC-DC converters using an inexpensive 8-b microcontroller," *IEEE Transactions on Industrial Electronics*, vol. 44, no. 5, pp. 661–669, Oct. 1997.
- [4] L. Dixon, "Current-Mode Control of Switching Power Supplies," vol. SM400. United States: Unitrode, 1985, pp. 1–9. [Online]. Available: <http://www.smeps.us/Unitrode2.html>
- [5] F. Leung, P. Tam, and C. Li, "The control of switching DC-DC converters-a general LWR problem," *IEEE Transactions on Industrial Electronics*, vol. 38, no. 1, pp. 65–71, Feb. 1991.
- [6] F. Leung, P. Tam, and C. Li, "An improved LQR-based controller for switching DC-DC converters," *IEEE Transactions on Industrial Electronics*, vol. 40, no. 5, pp. 521–528, Oct. 1993.
- [7] C. Olalla, R. Leyva, and A. El-Aroudi, "control for DC-DC buck converters," in *International Symposium on Circuits and Systems*, May 2006, pp. 4 pp.–5642.
- [8] A. Altowati, K. Zenger, and T. Suntio, "based robust controller design for a DC-DC switching power converter," in *European Conference on Power Electronics and Applications*, Sep. 2007, pp. 1–11.
- [9] A. Basim, P. Kiran, and R. Abraham, "based robust controller for DC-DC Boost Converter," in *International conference on Circuits, Controls and Communications*, Dec. 2013, pp. 1–6.
- [10] C. Olalla, R. Leyva, and A. El-Aroudi, "design for current-mode PWM buck converters operating in continuous and discontinuous conduction modes," in *32nd Annual Conference on IEEE Industrial Electronics*, Nov. 2006, pp. 1828–1833.
- [11] C. Olalla, C. Carrejo, R. Leyva, C. Alonso, and B. Estivals, "Digital QFT robust control of DC-DC current-mode converters," *Electrical Engineering*, vol. 95, no. 1, pp. 21–31, Mar. 2013. [Online]. Available: <http://search.ebscohost.com/login.aspx?direct=true&db=a9h&AN=85386272&site=ehost-live>
- [12] A. Saxena and M. Veerachary, "based robust controller design for fourth-order boost dc-dc switching power converter," in *Joint International Conference on Power Electronics, Drives and Energy Systems*, Dec. 2010, pp. 1–6.
- [13] P. Mattavelli, L. Rossetto, G. Spiazzi, and P. Tenti, "General-purpose fuzzy controller for DC-DC converters," *IEEE Transactions on Power Electronics*, vol. 12, no. 1, pp. 79–86, Jan. 1997.
- [14] S. Tan, Y. Lai, and C. Tse, "General Design Issues of Sliding-Mode Controllers in DC-DC Converters," *IEEE Transactions on Industrial Electronics*, vol. 55, no. 3, pp. 1160–1174, Mar. 2008.
- [15] B. Jayakrishna and V. Agarwal, "implementation of QFT based controller for a buck type DC-DC power converter and comparison with fractional and integral order PID controllers," in *11th Workshop on Control and Modeling for Power Electronics*, Aug. 2008, pp. 1–6.
- [16] R. Middlebrook and S. Ćuk, "A general unified approach to modelling switching-converter power stages," *International Journal of Electronics*, vol. 42, no. 6, pp. 521–550, Jun. 1977. [Online]. Available: <http://dx.doi.org/10.1080/00207217708900678>
- [17] D. Maksimovic, A. Stanković, V. Thottuvelil, and G. Verghese, "Modeling and simulation of power electronic converters," *Proceedings of the IEEE*, vol. 89, no. 6, pp. 898–912, Jun. 2001.
- [18] S. Sanders, "On limit cycles and the describing function method in periodically switched circuits," *IEEE Transactions on Circuits and Systems I: Fundamental Theory and Applications*, vol. 40, no. 9, pp. 564–572, Sep. 1993.
- [19] B. Choi, "Step load response of a current-mode-controlled DC-to-DC converter," *IEEE Transactions on Aerospace and Electronic Systems*, vol. 33, no. 4, pp. 1115–1121, Oct. 1997.
- [20] M. Poodeh, S. Eshtehardiha, A. Kiyoumars, and M. Ataei, "Optimizing LQR and pole placement to control buck converter by genetic algorithm," in *International Conference on Control, Automation and Systems*, Oct. 2007, pp. 2195–2200.
- [21] D. Naidu, *Optimal control systems*, ser. Electrical engineering textbook series. Boca Raton, Fla: CRC Press, 2003.
- [22] I. Horowitz and M. Sidi, "Synthesis of feedback systems with large plant ignorance for prescribed time-domain tolerances," *International Journal of Control*, vol. 16, no. 2, pp. 287–309, Aug. 1972. [Online]. Available: <http://dx.doi.org/10.1080/00207177208932261>
- [23] I. Horowitz, "Quantitative feedback theory," *Control Theory and Applications, IEE Proceedings D*, vol. 129, no. 6, pp. 215–226, Nov. 1982.
- [24] D. Kim, B. Choi, D. Lee, and J. Sun, "Dynamics of Current-Mode-Controlled DC-to-DC Converters with Input Filter Stage," in *Power Electronics Specialists Conference*, Jun. 2005, pp. 2648–2656.
- [25] W. Chen and D. Ballance, "Automatic loop-shaping in qft using genetic algorithms," Tech. Rep., 1998.
- [26] A. Zolotas and G. Halikias, "Optimal design of PID controllers using the QFT method," *Control Theory and Applications, IEE Proceedings*, vol. 146, no. 6, pp. 585–589, Nov. 1999.

Irregular segmentation technique for the image compression using stochastic models

BENABDELLAH YAGOUBI
 Laboratory of signals and systems
 Department of Electrical Engineering
 Faculty of Sciences and technology
 Mostaganem University, ALGERIA.
 E-mail: yagoubibenabdellah@yahoo.com

Abstract—the goal of the current paper is to develop a simple technique for the image irregular compression. This technique consists in segmenting each image matrix line and representing the resulting segments by adequate stochastic models. The irregularly compressed image, in this case, is represented by the parameters matrix of the corresponding mono-dimensional stochastic model. The latter could be either the normal distribution model, the uniform distribution model or else depending on the image reconstruction quality. We have, rather, used the two former models due to their simplicity in implementation as well as the good irregular compression results obtained in this work.

Keywords- irregular compression; image; segmentation; stochastic model.

1. Introduction

It is well known that image compression is the procedure of reducing the size in terms of bytes of a graphics file without degradation of the image quality to a certain reasonable level. This size file minimization allows more images to be saved or stored in a memory space. It also reduces the time required for transmitting images.

Many successful methods and algorithms in image modeling and compression are performed using statistical models, and it is therefore of interest to improve models in order to obtain a higher compression rate as well as to accurately reconstruct an image that is as close as possible to the corresponding original image. Performing such models is usually a difficult task due mainly to the image data to be processed. To overcome this difficulty, two important assumptions are usually pointed out to simplify model analysis; a) the probability of a pixel is conditioned only on very nearest neighborhood and deemed independent from the remaining pixels of the image. This assumption is called Markovianity. b) The local density is thought of being independent of its absolute position in the image, in other words the density is homogenous. Any model that is characterized by these two assumptions is called homogenous Markov random field (hMRF). The non-Gaussian statistics of images, in addition, led some authors to develop non-Gaussian MRF models [1, 2]. But so far the most successful, may be, is the fields of experts model

which has been recently developed by S. Roth and M. J. Black [3] and has shown a quite good performance. The local statistical properties of images have, also, been modeled using Gaussian scale mixtures (GSMs). Despite their global consistency and the interesting results provided by such field models, difficulties in their implementation and processing may hamper their performances. We suggest, therefore, in this paper an alternative and simple method based on matrix line segmentation instead of a Gaussian field realization for the image modeling and particularly for the image compression. This mode of representation corresponds to the regular compression and has been applied for the line by line processing of the images, in particular for the coding, the filtering and the storage [4]. In our case, however, we aimed to obtain a higher compression rate using an approach called irregular compression which will be explained in more detail in the following.

2. Preliminary Results

Many authors [5,...,10] have demonstrated that image statistics are not Gaussian realization, and hence they do not follow, particularly, Gaussian distribution. For example, decomposition of images using wavelet transforms provides coefficients that are non-Gaussian as indicated by their histograms. We believe, however, that this case may be slightly different when considering the image matrix line by line treatment and performing segmentation such that any segment may

be approximated by a corresponding Gaussian model [11, 12] as formulated in the following:

We started first by selecting a type of model with reasonable results and can be applied to a possible large number of images. After many tests on several types of models, we have decided to choose either the normal distribution model or the uniform distribution model for every resulting segment, mainly because of their simple computing and the good results they provide.

A)-The proposed irregular compression method

The principle of our proposed method of compression consists mainly of two steps: first, segment each matrix line into even length stationary intervals and seek the optimal model parameters (variance and mean in the case of a Gaussian model) to represent each of these intervals by an adequate corresponding model. If the reconstructed image is reasonably close to the original one, using this particular model, then we go to the second step which is grouping as many adjacent stationary segments of the same matrix line with very close variances and means as possible in order to obtain stationary segments with longer lengths, and each of them represented by two parameters only. When the image reconstruction is reasonably perfect by modeling every segment using the selected model, then it means that all the segments belong indeed to this model. And by gathering the adjacent segments that have approximately even variances and means, and hence they come from the same distribution, we determine the different lengths of the stationary segments at the same time. We have represented in figure.1 below both the 300th line (column 201 to column 320) of the original image and its reconstruction version as well as the stationary segments which are delimited by the symbols 'o'.

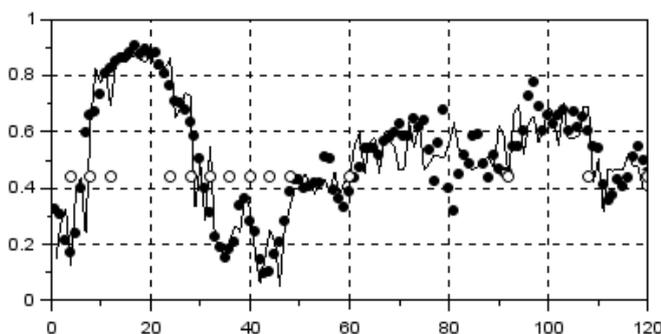


Fig.1. Original line (dotted line) and its reconstructed (continuous line), the symbols 'o' delimit the stationary segments.

So, this procedure is not only an alternative test to Kolmogorov-Smirnov test, but it represents also a

simple way of determining the stationary and ergodic segments [13] for any matrix line. Notice that this second step is called an irregular compression since the stationary segments, obtained using our algorithm are not of even lengths, whereas the former is known as the regular compression. The obtained irregular segments are delimited by the breakpoints separating the adjacent segments. We then gather all these optimal parameters and their corresponding segments breakpoints indices in a matrix form with a smaller size than that of the original image matrix. The goal of our method for compression is, therefore, to reduce the size of the parameters matrix as possible as we could without degrading too much the information in the original image.

B)-The suggested algorithm for the irregular compression

The basic idea of the irregular (compression) segmentation is to seek for longer stationary segments by joining as many adjacent smaller regular stationary segments which are determined using statistical models as possible. This method is applied to each image matrix line. The parameters of the resulting adequate statistical model corresponding to the irregular compression are then gathered in the parameters matrix which is supposed to represent the irregular compressed image.

- 1- Divide each image matrix line into segments with an even given length L .
- 2- Compute the parameters; the variance and the mean of every segment.
- 3- Use these parameters to reconstruct segments [14] using the selected model, and hence to reconstruct the image.
- 4- If the reconstructed image is reasonably perfect, then go to step 5, else reduce the segment length L and go back to step 1.
- 5- Gather as many adjacent segments with approximately even variance and mean as possible to obtain longer stationary segments.
- 6- Groupe the different breakpoints indices (lengths) of the obtained stationary segments as well as their corresponding parameters; variances and means in a matrix form whose size should be not less than that of the original image matrix only but also less than that of the parameters matrix of the regular compression as well. Notice that the latter corresponds to the first step up to the fourth, whereas the irregular compression starts from the fifth to sixth step.

The following scheme in figure.2 describes briefly the six steps of our suggested algorithm for the irregular compression by the matrix line segmentation.

-The original image matrix size is $I*J=348*620=215760$, where I and J are the rows and the columns numbers respectively.

-The regularly compressed image parameters matrix size obtained is $2*Ni*I= 2*155*348= 107880$, where Ni is the number of regular segments (same length) and finally the irregularly compressed image size obtained using our algorithm is $3*d= 3*30605= 91815$, where d is the number of all the irregular segments. As results we have obtained the following compression rates: the regular compression rate $Tr=215760/107880=2$ and the irregular compression rate $Tirr=215760/91815=2.35$. These results which are summarized in the following table 1 show clearly the improvement of the compression rate corresponding to the irregular compression.

Table.1:

Image nature	Image size	Compression Rate
Original	215760	1
Regular compression	107880	2
Irregular compression	91815	2.35

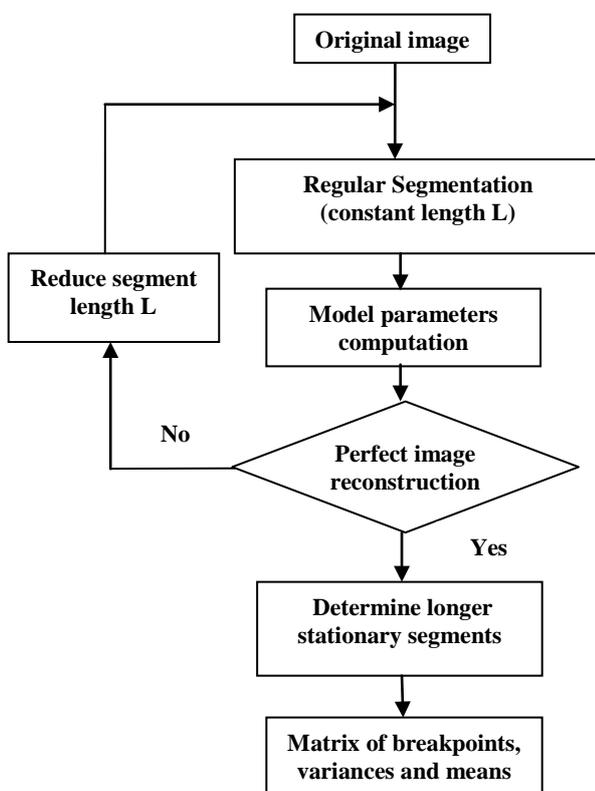
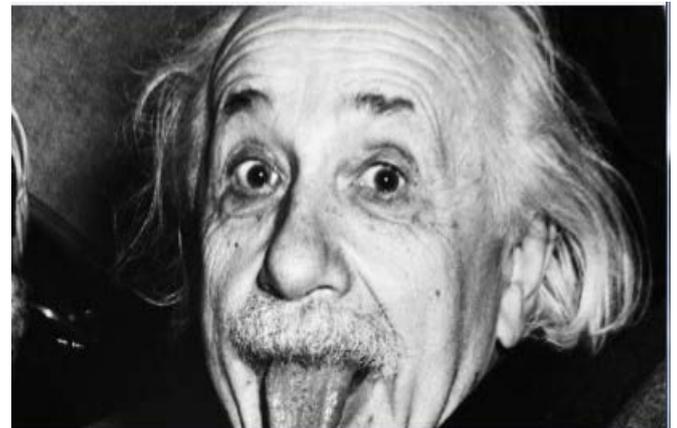
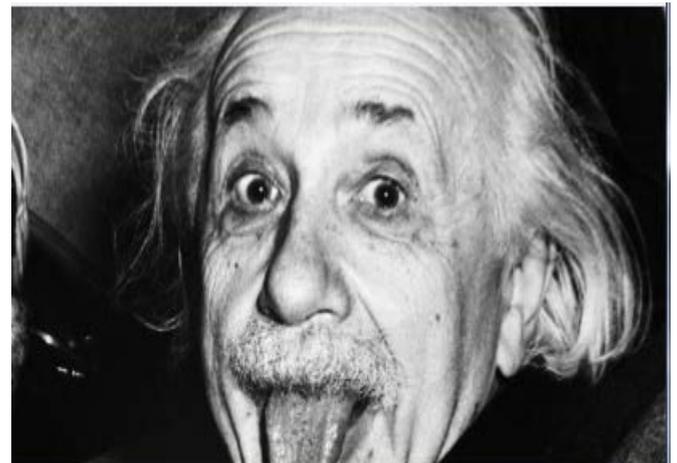


Fig. 2. The image irregular compression algorithm.

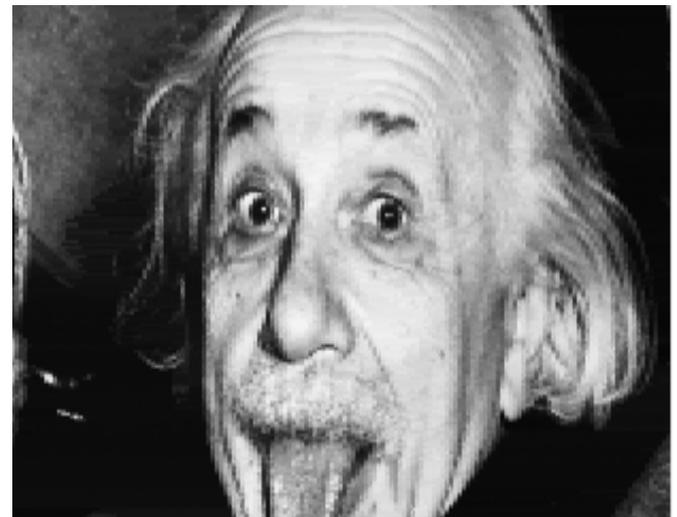
The results of our algorithm are shown in fig.3 using the uniform distribution and the normal distribution. We can see that the results obtained with the latter distribution are better than those obtained using the former distribution.



a



b



c

Fig.3. a) the original image, b) the reconstructed with normal distribution and c) the reconstructed with uniform distribution

The parameters matrix for the irregular compression, obtained using our algorithm, is 3 rows by a number of columns which is equal to the number of the breakpoint indices. Each breakpoint index represents a stationary interval length. The first, the second and the third row of the parameters matrix represent, respectively, the breakpoint indices, the variances and the means of the corresponding segments.

CONCLUSION

The advantage of our technique for the irregular compression over most non-Gaussian fields based methods lies in the ease of computing the one-dimensional Gaussian representation of the stationary segments of each matrix line. It should be noted, however, that in our line by line analysis of the image matrix, we have assumed that the pixels are very weakly dependent in order to reconstruct the image using the independent random variables joint probability. Our good results of the irregular compression and the quality of the reconstructed image show indeed that this assumption is quite reasonable.

References

- [1] S. C. Zhu, Y. Wu, and D. Mumford. Filters, random fields and maximum entropy (FRAME): Towards a unified theory for texture modeling. *Int'l. J. Comp. Vis.*, 27(2):107–126, 1998.
- [2] P. Gehler and M. Welling. Products of "edge-perts". In *Adv. in Neural Info. Proc. Systems (NIPS*05)*. MIT Press, 2006.
- [3] J.S. Roth and M. J. Black. Fields of experts: a framework for learning image priors. In *IEEE Conf. on Comp. Vis. and Pat. Rec.*, volume 2, pages 860–867, 2005.
- [4] Zhong , J. and Sclaro, S.... Segmenting foreground objects from a dynamic textured background via a robust Kalman filter. In *Proceedings of the International Conference on Computer Vision (ICCV)*, pp. 44-50, 2003.
- [5] D. J. Field. Relations between the statistics of natural images and the response properties of cortical cells, *J. Opt. Soc. Am.* , 4(12):2379–2394, 1987.
- [6] D. Ruderman. The statistics of natural images. *Network : Comp. in Neural Sys.* , 5:598–605, 1994.
- [7] M. J. Wainwright and E. P. Simoncelli. Scale mixtures of Gaussians and the statistics of natural images. In *Adv. Neural Info. Proc. Sys. (NIPS*99)* , volume 12, pages 855–861, May 2000.
- [8] E.P. Simoncelli and E.H. Adelson. Noiseremoval via Bayesian wavelet coring. In *Int'l Conf on Image Proc.*, pp379–383, Lausanne, Sep. 1996.
- [9] S. G. Mallat, A theory for multi-resolution signal decomposition: The wavelet representation, "IEEE Pat. Anal. Mach. Intel l.11, pp. 674-693, July 1989.
- [10] X. Descombes, R. D. Morris, J. Zerubia, and M. Berthod. Estimation of Markov random field prior parameters using Markov chain Monte Carlo maximum likelihood. *IEEE Trans. Image Process.*, 8(7):954 - 963, July 1999.
- [11] Lafarge, F, Descombes, X, Zerubia, J, Mathieu, S., (2007). Forest fire detection by statistical analysis of rare events from thermal infrared images, *Traitement du signal*, 24(1),1-12
- [12] Lafarge F, Descombes, X, Zerubia, J, Mathieu, S., (2007). Forest fire detection based on Gaussian field analysis, *European Signal Processing Conference (EUSIPCO)*.
- [13] Kay, S.M., 1998. *Fundamentals of Statistical Signal Processing, Detection Theory*, vol. 2, Prentice–Hall.
- [14] Mumford, D. and B. Gidas: 2001, 'Stochastic Models for Generic Images'. *Quarterly of Applied Mathematics* 59(1), 85-111.

Efficient Media Digital Library Design of Summarized Video based on Scalable Video Coding for H.264 (MDLSS)

Hesham Farouk *, Kamal ElDahshan**, Amr Abozeid **, Mayada Khairy*

* *Computers and Systems Dept., Electronics Research Institute, Cairo, Egypt.*

Hesham@eri.sci.eg, Mayada@eri.sci.eg

** *Dept. of Mathematics, Computer Science Division, Faculty of Science, Al-Azhar University, Cairo, Egypt.*

Dahshan@gmail.com, Amrapozaid@gmail.com

Abstract— With the fast advancement of wireless networks bandwidth and mobile devices, large scale digital video library systems are growing rapidly. However, the huge increasing of content and the data intensive nature of video make the management and browsing of video collections, as well as their search and retrieval, increasingly difficult. The need of having a media digital library is essential these days with intelligent tools for indexing the video with allocating the suitable metadata that describe the content of such videos and at the mean while tools for retrieving the archived video with fast techniques.

These will be achieved across 3 steps, working on the stream coding with multi bit rates and methods of handling, representing the video with summarized stream carrying the same information of the full stream and deriving a media digital library for indexing and retrieval process.

The first step, stream handling, will be across implementing scalable video techniques which set the bit rate according to the required application and the delivery devices. Scalable video coding offers a solution for meeting such heterogeneous requirements. The second step, video summarization, which plays an important role in this context; it makes navigation easier and provides the user with a quick idea about the content. Another issue is that the same video content can be accessed from a wide variety of terminal devices which differs with respect to bandwidth limitation, decoding complexity, power constraints and screen size. The third step is implementing a media digital library for storing the code and/ or the summarized video based on Media Asset management system.

The main innovation of this project is to explore the use of scalable video coding and video summarization techniques to enhancing a digital video library and the integration between these 3 modules.

Keywords— Video Processing, Scalable video coding, Video summarization, Key Frame Extraction, Video skimming, Home video, Mobile computing. Video indexing and retrieval

I. INTRODUCTION

Nowadays, multimedia communications has significantly facilitated and enriched people's daily life. People have witnessed the fast development of various wireless multimedia applications, such as video content distribution (e.g., YouTube) and live video communications (e.g., Skype, MSN, etc.). As a result, the volume of video data is rapidly increasing, over 6 billion hours of video are watched each

month on YouTube and more than 100 hours of video are uploaded to YouTube every minute [1]. Moreover, the increased popularity of mobile devices and wireless networks and their ubiquitous use for video recording and streaming leads to dramatically increases traffic of videos on such devices. Cisco stated that "Mobile Video will generate over 69 Percent of Mobile Data Traffic by 2018". Mobile makes up almost 40% of YouTube's global watch time [2].

A. Problem identification

Video delivery especially via mobile wireless networks faces diverse challenges, including limited bandwidth, dynamic network conditions with low stability, variety of relay equipment, different terminal decoding speeds, various display screen resolution, and limited battery capacity, etc. [3]. Therefore, the video coding system must encode the video sequence in different frame sizes, frame rates, and bit rates to supply such heterogeneous demands [4]. Another problem is that, the increasing amount of content and the intensive nature of video data make the management and browsing of stored video collections, as well as their search and retrieval, this increases the system handling difficulties [5].

B. Need for the system

Video is increasingly becoming one of the most pervasive technologies in terms of everyday usage, both for entertainment and in the enterprise environments. Mobile video is responsible for a majority of the growth seen in mobile broadband data volume. The demand for better video services (streaming, storing, retrieving, browsing and etc.) for mobile devices is a key challenge. The proposed system aims to solve some of these challenges.

This paper is organized as follows: Section II introduces the research approach and methodology. Section III presents the MDLSS design architecture and discusses its modules. Finally, in section IV we conclude the paper and suggest a future work.

II. RESEARCH APPROACH AND METHODOLOGY

A. The first motivation of this system

Today there is a wide range of different devices available for viewing video content, including smartphones, tablets, laptops and televisions. Every client's requirement differs with respect to bandwidth limitation, decoding complexity, power constraints and screen size. Scalable video coding offers a solution for meeting such heterogeneous requirements [6].

A video bit stream is called scalable if a part of the stream can be removed in such a way that the resulting bit stream is still decodable. The three types of scalabilities are [7, 8]:

1. Temporal (frame rate) scalability: the motion compensation dependencies are structured so that complete pictures (i.e. their associated packets) can be dropped from the bit stream. Temporal scalability is already enabled by H.264/MPEG-4 AVC. SVC has only provided supplemental enhancement information to improve its usage.
2. Spatial (picture size) scalability: video is coded at multiple spatial resolutions. The data and decoded samples of lower resolutions can be used to predict data or samples of higher resolutions in order to reduce the bit rate to code the higher resolutions.
3. SNR/Quality/Fidelity scalability: video is coded at a single spatial resolution but at different qualities. The data and decoded samples of lower qualities can be used to predict data or samples of higher qualities in order to reduce the bit rate to code the higher qualities.

This work is represented as module 1 of the proposed work given in Fig. 1.

B. The second motivation of this system

The content of video may be huge and crowded with much redundant information so that it often takes a long time to browse the content from the beginning to the end. Also, the user may not have sufficient time to watch the entire video or the video content, as a whole, may not be of interest to the user. In such cases, the user may just want to view the summary of the video instead of watching the whole video [9].

Video summarization is a mechanism for generating compact representation of a video sequence, which includes only the

important parts in the original video [10]. Video summarization is useful when a system is operating under tight constraints (e.g. Limited bandwidth, watching time or memory size). For example, in surveillance applications the video may be recorded nearly for 24 hours per day, a summary version of the original video may be useful to watch the important events only in such case. Also, video summarization is useful when we need to transmit an important video segment to another device in real time [11]. Video summarization techniques target different domains of video data, such as sports, news, movies, documentaries, e-learning, surveillance, home videos, etc., And discuss various assumptions and viewpoints to produce an optimal or good video summary [9].

This work is represented as module 2 of the proposed work given in Fig. 1.

There are two fundamental types of video summaries [12]: static video summary (also called representative frames, still-image abstracts or static storyboard) and dynamic video skimming (also called video skim, moving image abstract or moving storyboard). The static video summary is a collection of video frames extracted from the original video. The dynamic video summary is a set of short video clips, joined in a sequence, and played as a short video clip. Usually, from the user's viewpoint, a dynamic video summary may provide a more good choice since it contains both audio and motion information that makes the summarization more interesting and natural, while static video summary may provide a glance of video contents in a more concise way. In addition, once video frames are extracted, there are further possibilities of organizing them for browsing and retrieving purposes [13].

C. The Third motivation of this system

The Digital library for media file under Media Asset Management (MAM) system which will be the main storage system for the processed video by module 1 and module 2 and will be based on the MAM purchased by ERI through the EQUIPME initiative issued by scientific research academy 2 years ago.

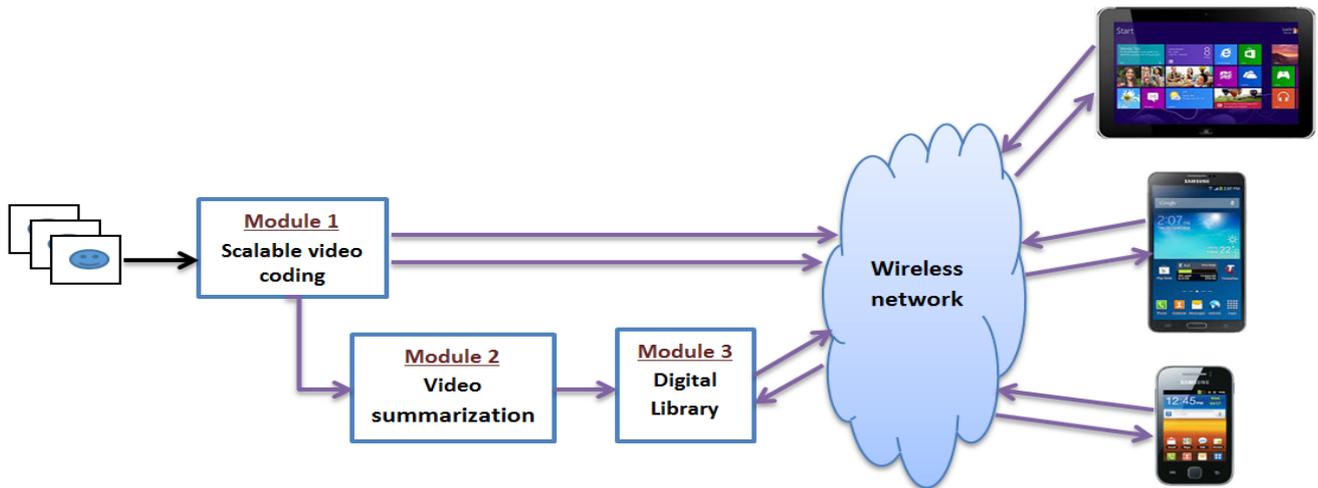


Fig.1 MDLSS architecture

III. THE PROPOSED SYSTEM

A. The proposed System architecture

The architecture of MDLSS consists of three modules; scalable video coding, video summarization and digital library as shown in Fig.1. MDLSS aims to provide better video services (streaming, storing, retrieving and browsing) for mobile devices which increase the interactions and activities between users and digital libraries. In other words, the main goal of MDLSS is to explore the use of scalable video coding and video summarization techniques to enhancing digital video library. This goal can be further specified in the following:-

- Design and develop Scalable Video Coding (SVC) algorithm to meet the requirements of applications and devices heterogeneities.
- Design and develop an automatic video summarization algorithm, which engage in providing concise and informative video summaries to help in browsing and managing video files efficiently.

B. The system design Methods and Procedures

For module 1:

General adapted methods for SVC

1. Determine number of layers for scalable video.
2. Determine number of bitrates available.
3. Analysis the video stream.
4. Select the type of scalability according to 3 steps before.
5. Implement the scalable video types according to previous steps.

For module 2:

A general adapted method for video summarization module is shown in Fig. 2. Each step is described as follows:

1- Frames sampling

The first step towards automatic video summarization is splitting the video stream into a set of meaningful and manageable basic elements (e.g., shots, frames) that are used as basic elements for summarization. Most of existing methods for automatic video summarization have focused on split the video stream into frames. The video sequence is decoded and each frame is extracted and treated separately [14].

2- Feature extraction

Digital video contains many features like color, motion, and voice etc. Color feature is considered an important aspect of video. That's why it has been used quite often for video summarization. Color based summarization techniques are very simple and easy to use. However, their accuracy is not reliable, as color based techniques may consider noise as part of the summary [15].

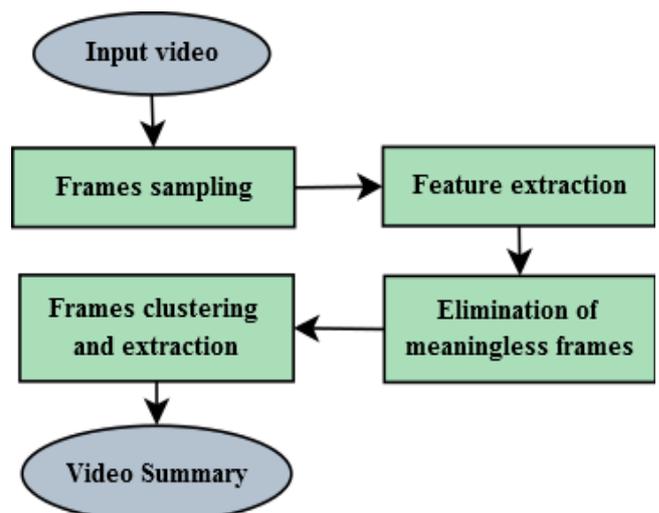


Fig. 2 Flowchart of video summarization module

3- Elimination of meaningless frames

The goal of this step is to avoid possible meaningless frames in a video summary. It has been generally observed that a video usually has some meaningless frames such as totally black frames, totally white frames (a monochromatic frame) and faded frames [13].

4- Frames clustering and extraction

The goal of this step is to group similar video frames together and to select a representative frame per each group, to produce the video summary. The effectiveness of grouping similar frames depends on the suitable choice of a similarity metric used for comparing two frames [16].

Video summarization is hot research filed in recent years due to its important role in many video services (e.g. browsing, indexing and streaming). The reader can find a comprehensive review of video summarization techniques in [9, 17-19]. Moreover, the authors in [15] introduce an analysis and comparative study between various techniques proposed in literature for the summarization of video content, which can be useful for mobile applications.

For module 3:

In this step we will manage video storage after editing and applying Meta data through the Media Asset Management we already have in our Digital Signal Processing Lab in Electronics Research Institute.

IV. CONCLUSIONS AND FUTURE WORK

The scalable video coding as well as video summarization plays an important role in many video services. So, in this paper we present an efficient media digital library design of summarized video based on scalable video coding for H.264 (MDLSS). The proposed design will utilize the conjunction between scalable video coding and video summarization techniques to enhance the digital video library.

In the future we arrange to develop this proposed project (MDLSS). Our implementation activities will be organized as follows:

- Literature survey of related works and highlight our goals. Which we did a lot of this by publishing comparative paper [15]
- Analysing the system requirement for each module and for integration
- Develop a system prototype.
- Testing the system and updates.

REFERENCES

- [1] YouTube Statistics, last access date is 10-12-2014; <http://www.youtube.com/yt/press/statistics.html>.
- [2] V. N. I. (VNI). Cisco Visual Networking Index: Global Mobile Data Traffic Forecast Update, 2013–2018, last access date is 10-12-2014, 2014; http://www.cisco.com/c/en/us/solutions/collateral/service-provider/visual-networking-index-vni/white_paper_c11-520862.html.
- [3] N. V. Uti, and R. Fox, The Challenges of Compressing and Streaming Real Time Video Originating from Mobile Devices, Multimedia Services and Streaming for Mobile Devices: Challenges and Innovation, pages. 1, 2011.
- [4] P.-C. Wang, G.-L. Li, S.-F. Huang, M.-J. Chen, and S.-C. Lin, Efficient mode decision algorithm based on spatial, temporal, and inter-layer rate-

distortion correlation coefficients for scalable video coding, ETRI journal, volume(32), issue(4), pages. 577-587, 2010.

- [5] L. Herranz, and J. M. Martínez, Combining MPEG Tools to Generate Video Summaries Adapted to the Terminal and Network, The Computer Journal, volume(56), issue(5), pages. 529-553, 2013.
- [6] M. Ransburg, E. M. Graciá, T. Sutinen, J. O. Murillo, M. Sablatschan, and H. Hellwagner, Scalable video coding impact on networks, Mobile Multimedia Communications, pages. 571-581, 2012.
- [7] S. Ibrahim, A. H. Zahran, and M. H. Ismail, SVC-DASH-M: Scalable video coding dynamic adaptive streaming over HTTP using multiple connections, Telecommunications (ICT), 21st International Conference on, pages. 400-404, 2014.
- [8] D. Zhang, H. Li, and C. Chen, Robust Transmission of Scalable Video Coding Bitstream over Heterogeneous Networks, Circuits and Systems for Video Technology, IEEE Transactions on, 2014.
- [9] B. T. Truong, and S. Venkatesh, Video abstraction: A systematic review and classification, ACM Transactions on Multimedia Computing, Communications, and Applications (TOMCCAP), volume(3), issue(1), pages. 37, 2007.
- [10] G. Guan, Z. Wang, S. Mei, M. Ott, M. He, and D. D. Feng, A top-down approach for video summarization, ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM), volume(11), issue(1), pages. 4, 2014.
- [11] L. Zhu, Z. Fan, and K. Aggelos K, Joint video summarization and transmission adaptation for energy-efficient wireless video streaming, EURASIP Journal on Advances in Signal Processing, 2008.
- [12] S. E. F. de Avila, and A. P. B. Lopes, VSUMM: A mechanism designed to produce static video summaries and a novel evaluation method, Pattern Recognition Letters, volume(32), issue(1), pages. 56-68, 2011.
- [13] J. Almeida, N. J. Leite, and R. d. S. Torres, Online video summarization on compressed domain, Journal of Visual Communication and Image Representation, volume(24), issue(6), pages. 729-738, 2013.
- [14] J. Niu, D. Huo, K. Wang, and C. Tong, Real-time generation of personalized home video summaries on mobile devices, Neurocomputing, ScienceDirect, volume(120), pages. 404-414, 2013.
- [15] H. Farouk, K. ElDahshan, and A. Abozeid, The State of the Art of Video Summarization for Mobile Devices: Review Article, Graphics, Vision and Image Processing GVIP, volume(14), issue(2), pages. 37-50, 2014.
- [16] S.-H. Ou, C.-H. Lee, V. S. Somayazulu, Y.-K. Chen, and S.-Y. Chien, Low complexity on-line video summarization with Gaussian mixture model based clustering, pages. 1260-1264.
- [17] A. G. Money, and H. Agius, Video summarisation: A conceptual framework and survey of the state of the art, Journal of Visual Communication and Image Representation, volume(19), issue(2), pages. 121-143, 2008.
- [18] [18] R. Pal, A. Ghosh, and S. K. Pal, Video Summarization and Significance of Content: A Review, Handbook on Soft Computing for Video Surveillance, pages. 79, 2012.
- [19] [19] M. Ajmal, M. H. Ashraf, M. Shakir, Y. Abbas, and F. A. Shah, Video summarization: techniques and classification, Computer Vision and Graphics, Springer, pages. 1-13, 2012

ACKNOWLEDGMENT

The work give a great acknowledgement to Dr Alaa Hefnawy, Dr. Alaa Hamdy and Dr. Amr Mohamed for their great support and honest consultation in the core technical specilaization for each of them

BIOGRAPHIES

Assoc. Prof. Farouk is an associate Prof. since 2012. He joined the Electronics Research Institute, Egypt, in 1993. His fields of research are signal processing, mobile systems, Neural Networks, image compression, video processing, video compression, video indexing and retrieval, video on demand, pattern recognition and machine vision. Dr. Farouk received his Ph.D. at 2001 from Electronics & Communications Dept., Faculty of



Engineering, Cairo Univ. and his M.Sc. at 1996 from Electronics & Communications Dept., Faculty of Engineering, Cairo Univ. Dr. Hesham participated in many national projects in MCIT developed based on portals and digital libraries. He also participated in some strategic studies as mobile for development.. Then he is a Research and innovation dept acting manager in ITI. Meanwhile, In ERI he is managing Technology Transfer office since June 2013 and the ERI technical office.



Prof. Kamal Abdelraouf EIDahshan is a professor of Computer Science and Information Systems at Al-Azhar University in Cairo, Egypt.

An Egyptian national and graduate of Cairo University, he obtained his doctoral degree from the Université de Technologie de Compiègne in France, where he also taught for several years.

During his extended stay in France, he also worked at the prestigious Institute National de Télécommunications in Paris.

Professor EIDahshan's extensive international research, teaching, and consulting experiences have spanned four continents and include academic institutions as well as government and private organizations. He taught at Virginia Tech as a visiting professor; he was a Consultant to the Egyptian Cabinet Information and Decision Support Centre (IDSC); and he was a senior advisor to the Ministry of Education and Deputy Director of the National Technology Development Centre. Prof. EIDahshan has taught graduate and undergraduate courses in information resources and centers, information systems, systems analysis and design, and expert systems.

Professor EIDahshan is a professional Fellow on Open Educational Resources as recognized by the United States Department of State.

Prof. Eldahshan wants to work in collaboration with the Ministry of Education to develop educational material for K-12 levels. Prof. Eldahshan is interested in training instructors to be able to use OER in their teaching and hopes to make his university a center of excellence in OER and offer services to other universities in the country.



Amr Abozeid received B.Sc. degree in computer science and mathematics from Al-Azhar University, Cairo, Egypt, in 2005. He received the M.Sc. in computer science from department of mathematics and computer science,

faculty of science, Ain Shams University, in 2012. He is now a Ph.D. student in department of mathematics and computer science, faculty of science, Al-Azhar University. His main research topic is focused on video processing and computer vision. He is working as lecturer assistant of computer science and has eight years' experience in the field of teaching computer science topics, tools and technologies, supervising students' projects and leading software development

projects.



Mayada Khairy Shehata graduated from the Telecommunications and Electronics Department, Faculty of Engineering, Helwan University, Cairo, Egypt in 2004. She received the M.Sc. in video processing from Faculty of Engineering, Helwan University, in 2010 ,Registered PhD in video

processing, Faculty of Engineering, Helwan University in 2014 .Mayada was working on the field of broadcast (TV& Radio) as Projects Manger in Systems Design company from 2004 to 2013 dealing with Egyptian Radio Television Union(ERTU) , Egyptian Media Predication City (EMPC) ,Nielsat, and Private satellite TV Channels. She is dealing with international companies in the field such as; Miranda , Grass Vally Harris, Avid, DMTsys, etc... .Attending international Radio and Television exhibitions, Technical seminars and training for the mentioned companies in England, Dubai, Lebanon, ,Oman ,Jordan. Also working as Assistant Lecturer in Faculty of Engineering Thebes Academy (Part Time).Now working as Assistant Researcher in Computers and Systems department Electronics Research Institute (ERI) from July 2013 up till now.

Speech enhancement using rao-blackwellised particle filtering of the real and imaginary DFT coefficients part

M. Meddah, A. Amrouche, A. Taleb-Ahmed

Abstract— In this work the speech enhancement is achieved by filtering the spectral coefficients of the noisy signal, where both real and imaginary parts are filtered separately, using the Rao-Blackwellized sequential Monte-Carlo (SMC) algorithm. A low order time-varying autoregressive (TVAR) model is adopted for each channel part, the performance of this method to enhance a noisy speech signal from an additive white Gaussian noise are compared to the use of the same algorithm in the time domain under the same conditions. The objective measures evaluation, show that the proposed concept, present a lower log-likelihood ratio (LLR) measure, and high perceptual assessment of the speech quality (PESQ) score, with keeping the same overalls SNR improvement.

Keywords— Speech enhancement, sequential Monte-Carlo, Rao-Blackwellized particle filter, the real and imaginary DFT coefficients part.

I. INTRODUCTION

Processing of the speech signal that has been degraded by additive background noise is of great interest in variety of context. Single channel speech enhancement algorithms attempt to recover a clean speech signal from a degraded signal containing additive noise. The objectives of the enhancement task is to improve speech intelligibility and to reduce the perceptual impact of the noise, improving perceived speech quality, and reducing listener fatigue. This algorithms can be classified either according to the processing domain (sub-space, spectral, temporal, discrete cosine transform,...), or according to processing approach (statistics, Thresholding,...), or given the adopted speech model (parametric, non-parametric).

According to the statistical approach for speech enhancement in [1] they used the maximum likelihood ML (Maximum Likelihood) approach, for estimating the amplitudes of the spectral components of the speech signal in the Fourier transform domain, where the likelihood probability density is defined with a Gaussian distribution according to the adopted linear Gaussian model. Considering the same model [2] have defined a MMSE estimate (minimum mean square error) of the amplitude of the spectral components of the

speech signal, assumed to be independent, despite a priori SNR estimation that exploit the dependence between the components of the frequency trajectory.

To consider the dependency of samples, the Auto-Regressive (AR) model is used, this late exploits the correlation between the components, by performing a prediction of the current sample, with a linear combination of the immediately preceding samples. In the time domain, and to keep the stationary for the AR model the speech signal is processed in segments duration less than 30ms.

In [3], the equations of Wiener-Kalman were defined based on the Bayesian approach, however, a more rigorous proof is provided in [4]. In [5], the AR model of the speech signal is invested in the development of equations of the space state, and the Kalman filter is used to estimate the states in a linear Gaussian model, the results are presented with a pre-knowledge of the parameters of AR model.

However, segmentation does not take into account the variation of the speech signal, in relation to the information flows. In addition, the variation of the vocal tract is not constant, hence time-varying Auto-Regressive (TVAR) model is more appropriate for modelling the speech signal. Thus a recursive parameter estimation of the AR model is imposed [6].

In [7] iterative EM (expectation maximization) method, is associated with the Kalman filter to estimate the vector of parameters, the noise is modelled by an AR model, which takes into account the case of colored noise. In [8] a variety of recent methods derived from the latter are presented.

In [6], the implementation of the sequential Monte-Carlo (SMC) called also particle method, for speech enhancement in the time domain is presented. Where Rao-Blackwellization approach is associated to the sequential Monte Carlo method to filtering noisy speech (RBPF). So, the parameter vector, and hidden states are estimated. Thus a hybrid filter is obtained, such as a part of the calculation is performed by SMC approaches, and the other part is made analytically [9].

In [10] the performance analysis of the particle filter (PF) and the RBPF-speech enhancement in time domain are presented. Where they find that the residual noise level is modulated by speech power.

Several effort was investigating to reduce this residual noise, thus in [11] hybrids approaches are proposed. In [12]

M. Meddah, A. Amrouche are with LCPTS, FEI, USTHB, B.P. 32 El Alia, Bab-Ezzouar, 16111, Algeria (mmeddah@usthb.dz, namrouche@usthb.dz).

A. Taleb-Ahmed is with LAMIH, UMR, CNRS 8530, UVHC Valenciennes, France (Abdelmalik.Taleb-Ahmed@univ-valenciennes.fr).

the RBPF approach to enhance the speech signal in the cosine domain is presented, in [13] the PF approach is exploited to estimate the amplitude and phase of the speech coefficients in the DFT domain.

The SMC, characterized by the weak restrictions on the considered model, is based on the implementation of the Monte-Carlo approximation of the a posterior probability density of interest, where the latter is approximated by independent particles (samples) identically distributed according to it. Unknown sampling according to the distribution of interest is replaced by sampling, from another density that is most similar to the density of interest and having the same support, which sets the sampling importance approach. The Sequential character is obtained by the propagation of the particles and their weight. However, this approach as it was presented, suffers from the growth of the variance of weight with the time. After few propagation step the weight will be concentrated on a single particle and the rest of the particles will have negligible weight. To remedy this, a resampling step is introduced before the propagation of particles, where the latter are multiplied according to their weight, and in order to not lose the exploration diversity of the state space, resampling is applied depending on the effective sample size of particles[9].

Also, based on the fact that, in general, the noise will not affect the speech signal uniformly over the whole spectrum [14]. In this work our contribution is the resulting improvement from the use of the RBPF method to enhance DFT noisy speech channel, where the real and imaginary part are processed separately.

At the best known of the others this approach has not before been applied. In Section II, the considered state space model is presented in section III an overview of SMC and the RB-SMC approach is presented. Subsequently, in section V the proposed approach based on the algorithm [6] is exposed. The performance of the proposed method are compared to the standard Time-RBPF in section VI.

II. TIME SPACE MODEL

After windowing, in which the correlation is maintained for the considered AR model order [15], the Fourier transform of the observation, $Y_{m,k}$ at the channel k , derived from the frame m , is modelled as a linear sum of the spectral component of the clean signal $X_{m,k}$ and the noise $V_{m,k}$:

$$Y_{m,k} = X_{m,k} + V_{m,k} \quad (1)$$

And taking into account the same independence assumption of real and imaginary part for the observation made by [15] and [16]. (1) becomes:

$$Y_{m,k}(\text{real}) = X_{m,k}(\text{real}) + V_{m,k}(\text{real})$$

$$Y_{m,k}(\text{imag}) = X_{m,k}(\text{imag}) + V_{m,k}(\text{imag}) \quad (2)$$

Thereafter the development, will be exposed equally to both parties. So, each spectral component part is assumed following a TVAR model, such as:

$$X_{m,k} = \sum_{i=1}^p a_{i,m,k} X_{m-i,k} + U_{m,k} \quad (3)$$

With; $U_{m,k} \sim \mathcal{N}(0, \sigma_{U_{m,k}}^2)$ is a white Gaussian process, with zero mean and variance $\sigma_{U_{m,k}}^2$ uncorrelated with all previous values of $X_{m,k}$, this distribution is justified by the Gaussian assumption adopted in this work, in order to keep the optimality of the estimate. In the literature, we find that the Laplace and Gamma distributions are closer to the spectral coefficient distribution of speech signal than the Gaussian distribution.

The additive observation noise process are assumed, white Gaussian with zero-mean and independent of the speech signal:

$$V_{m,k} \sim \mathcal{N}(0, \sigma_{V_{m,k}}^2) \quad (4)$$

Where the matrix form:

$$\begin{bmatrix} X_{m-p+1,k} \\ \vdots \\ X_{m,k} \end{bmatrix} = \begin{bmatrix} \mathbf{0}_{p-1 \times 1} & & & \\ & I_{p-1 \times p-1} & & \\ & & \dots & \\ & & & a_{1,m,k} \end{bmatrix} \begin{bmatrix} X_{m-p,k} \\ \vdots \\ X_{m-1,k} \end{bmatrix} + \begin{bmatrix} \mathbf{0}_{p-1 \times 1} \\ \sigma_{U_{m,k}} \end{bmatrix} u_{m,k}$$

With: $u_{m,k} \sim \mathcal{N}(0,1)$, Thus we write:

$$\alpha_{m,k} = A_{m,k} \alpha_{m-1,k} + B_{m,k} u_{m,k} \quad (5)$$

As the same way the observation, is written:

$$U_{m,k} \sim \mathcal{N}(0, \sigma_{U_{m,k}}^2)$$

$$Y_{m,k} = \begin{bmatrix} \mathbf{0}_{p-1 \times 1} & 1 \end{bmatrix} \begin{bmatrix} X_{m-p+1,k} \\ \vdots \\ X_{m,k} \end{bmatrix} + \begin{bmatrix} \sigma_{V_{m,k}} \end{bmatrix} v_{m,k}$$

With: $v_{m,k} \sim \mathcal{N}(0,1)$, Thus we write:

$$Y_{m,k} = C_{m,k} \alpha_{m,k} + D_{m,k} v_{m,k} \quad (6)$$

If the TVAR model parameters are known at all pairs (m,k) and the hidden state is assumed to follow a first order Markov model, with independent observations conditionally to the states. A sequential Bayesian filtering of the hidden state for this linear Gaussian model is achieved in two steps. (i) In the first time a posteriori probability density is predicted, (ii) then, the latter is corrected by taking into account the observation at (m,k) , thus the state estimate is obtained either with MMSE or MAP (maximum a posteriori) estimator. An equivalent procedure is performed by the adoption of the Kalman filter, where the estimate is obtained directly in terms of MMSE. Note that the assumption of Gaussian noise is not restrictive to the functioning of the Kalman filter, but if optimality is the ability of the filter to deduce the expression of

a posterior density in this case the Kalman filter is not certainly optimal [17].

III. SEQUENTIAL MONTE-CARLO METHOD

A rigorous and extensive introduction of the SMC method is presented in [17], however, in what follows we will take a few developments reported from [18] and [19].

The a posteriori probability density without channel index can be written as follows:

$$p(\alpha_{m':m''}/y_{1:m}) = \frac{p(\alpha_{m':m''}, y_{1:m})}{p(y_{1:m})} \quad (7)$$

Ones (7) defined, the estimate $\hat{x}_{m':m''}$ of the stat $x_{m':m''}$ conditionally to the observations $y_{1:m}$, is derived according to optimization criterion [20]. For $m' = m'' = m$ we have filtering. 2- for $m' = m'' > m$ prediction. 3- for $m' = m'' = m - l$ fixed lag smoothing. 4- For $m' = 0, m'' = m$ fixed interval smoothing.

A. Importance sampling (IS)

If the expression (7) cannot be derived analytically, the approximation by Monte Carlo method is used to represent this latter, such as for $m' = 0, m'' = m$, we have :

$$p(\alpha_{0:m}/y_{1:m}) = \frac{1}{N} \sum_{i=1}^N \delta(\alpha_{0:m} - \alpha_{0:m}^i) \quad (8)$$

Where $\delta(\cdot)$ denote Dirac function, N is total number of the used particles (samples), $\alpha_{0:m}^i$ are independent identically distributed samples, according to $p(\alpha_{0:m}/y_{1:m})$, we note $\alpha_{0:m}^i \sim p(\alpha_{0:m}/y_{1:m})$, So, taking (8) into account, The a posterior expectation of the function $f(\alpha_{0:m})$, becomes:

$$E_{\hat{p}_N(\alpha_{0:m}/y_{0:m})}[f(\alpha_{0:m})] = \frac{1}{N} \sum_{i=1}^N f(\alpha_{0:m}^i) \quad (9)$$

However, the posterior density is not known. So, let consider, $q(\alpha_{0:m}/y_{1:m})$, the importance density which is most similar to the probability density of interest, with $p(\alpha_{0:m}/y_{1:m}) > 0 \Rightarrow q(\alpha_{0:m}/y_{1:m}) > 0$ thus (7) can be written as following:

$$\begin{aligned} p(\alpha_{0:m}/y_{1:m}) &= \frac{\frac{p(\alpha_{0:m}/y_{1:m})}{q(\alpha_{0:m}/y_{1:m})} q(\alpha_{0:m}/y_{1:m})}{\int \frac{p(\alpha_{0:m}/y_{1:m})}{q(\alpha_{0:m}/y_{1:m})} q(\alpha_{0:m}/y_{1:m}) d\alpha_{0:m}} \\ &= \frac{w(\alpha_{0:m}, y_{1:m}) q(\alpha_{0:m}/y_{1:m})}{\int w(\alpha_{0:m}, y_{1:m}) q(\alpha_{0:m}/y_{1:m}) d\alpha_{0:m}} \end{aligned} \quad (10)$$

And taking into account the MC approximation of the importance density (10) becomes:

$$\hat{p}_N(\alpha_{0:m}/y_{1:m}) = \frac{\frac{1}{N} \sum_{i=1}^N w(\alpha_{0:m}^i, y_{1:m}) \delta(\alpha_{0:m} - \alpha_{0:m}^i)}{\frac{1}{N} \sum_{i=1}^N w(\alpha_{0:m}^i, y_{1:m})} \quad (11)$$

$$= \frac{1}{N} \sum_{i=1}^N W_{0:m}^i \delta(\alpha_{0:m} - \alpha_{0:m}^i) \quad (12)$$

With: $\alpha_{0:m}^i \sim q(\alpha_{0:m}/y_{1:m})$

$$W_{0:m}^i = \frac{w_{0:m}^i}{\sum_{j=1}^N w_{0:m}^j} \quad (13)$$

$$w_{0:m}^i = w(\alpha_{0:m}^i, y_{1:m}) = \frac{p(\alpha_{0:m}^i, y_{1:m})}{q(\alpha_{0:m}^i/y_{1:m})} \quad (14)$$

$W_{0:m}^i$ is the normalize importance weight, which defines a measure of the importance of sample i at the instant m . Therefore, (9) can be written as follows:

$$E_{\hat{p}_N(\alpha_{0:m}/y_{1:m})}[f(\alpha_{0:m})] = \sum_{i=1}^N W_{0:m}^i f(\alpha_{0:m}^i) \quad (15)$$

A. Sequential importance sampling (SIS)

Let choose the importance density in the class of importance distributions of the following form:

$$\begin{aligned} q(\alpha_{0:m}/y_{1:m}) &= q(\alpha_{0:m-1}/y_{1:m-1}) q(\alpha_m/\alpha_{0:m-1}, y_{0:m}) \\ &= q(\alpha_0) \prod_{j=1}^m q(\alpha_j/\alpha_{0:j-1}, y_{0:j}) \end{aligned} \quad (16)$$

Thus:

$$\begin{aligned} w(\alpha_{0:m}, y_{1:m}) &= \frac{p(\alpha_{0:m-1}, y_{1:m-1}) p(\alpha_m/\alpha_{m-1}) p(y_m/\alpha_m)}{q(\alpha_{0:m-1}/y_{1:m-1}) q(\alpha_m/\alpha_{0:m-1}, y_{0:m})} \\ &= w(\alpha_{0:m-1}, y_{1:m-1}) \frac{p(\alpha_m/\alpha_{m-1}) p(y_m/\alpha_m)}{q(\alpha_m/\alpha_{0:m-1}, y_{0,m})} \end{aligned} \quad (17)$$

Thus; (14) becomes:

$$w_{0:m}^i = w_{0:m-1}^i \frac{p(\alpha_m^i/\alpha_{m-1}^i) p(y_m/\alpha_m^i)}{q(\alpha_m^i/\alpha_{0:m-1}^i, y_{0,m})} \quad (18)$$

To minimize the variance of the importance weight, a resampling step is applied in which; large weight particles are duplicated by against the others particle are not propagated again. Then, all the normalized weight are set to $(1/N)$, witch define the sequential version of importance sampling with resampling (SISR). However, the resampling is applied only on indication. By measuring the effective sample size N_{eff} , in order to not lose the exploration diversity of the state space.

$$N_{eff} \equiv \frac{1}{\sum_{i=1}^N (w_{0:m}^i)^2} \quad (19)$$

B. Rao-blackwellisation for sequential importance simplig:

To reduce the variance of the MC estimator, the Rao-blackwellization approach is used. Where hybrid filter is obtained, such that a part of the calculation is obtained analytically and the other part is formed by MC.

Assuming that we can split the state $\alpha_{0:m}$ as $(\alpha_{0:m}, \theta_{0:m})$, such as:

$$E_{p(\alpha_{0:m}, \theta_{0:m}/y_{1:m})} [f(\alpha_{0:m}, \theta_{0:m})] = \frac{\iint f(\alpha_{0:m}, \theta_{0:m}) p(\alpha_{0:m}, \theta_{0:m} / y_{1:m}) d\alpha_{0:m} d\theta_{0:m}}{\iint p(\alpha_{0:m}, \theta_{0:m} / y_{1:m}) d\alpha_{0:m} d\theta_{0:m}} \quad (20)$$

Hence, two estimates of the posterior expectation are possible, either by the method of SMC. Or by RB-SMC approach. Where (20) can be written as following:

$$\frac{\iint f(\alpha_{0:m}, \theta_{0:m}) p(\alpha_{0:m} / \theta_{0:m}, y_{1:m}) p(\theta_{0:m}, y_{1:m}) d\alpha_{0:m} d\theta_{0:m}}{\iint p(\alpha_{0:m} / \theta_{0:m}, y_{1:m}) p(\theta_{0:m}, y_{1:m}) d\alpha_{0:m} d\theta_{0:m}} \quad (21)$$

Thus, for the importance density $q(\theta_{0:m} / y_{1:m})$, (21) becomes:

$$\frac{\iint f(\alpha_{0:m}, \theta_{0:m}) p(\alpha_{0:m} / \theta_{0:m}, y_{1:m}) \frac{p(\theta_{0:m}, y_{1:m})}{q(\theta_{0:m} / y_{1:m})} q(\theta_{0:m} / y_{1:m}) d\alpha_{0:m} d\theta_{0:m}}{\iint p(\alpha_{0:m} / \theta_{0:m}, y_{1:m}) \frac{p(\theta_{0:m}, y_{1:m})}{q(\theta_{0:m} / y_{1:m})} q(\theta_{0:m} / y_{1:m}) d\alpha_{0:m} d\theta_{0:m}} \quad (22)$$

Taking into account the MC approximation of the considered important function, (22), becomes:

$$\frac{\sum_{i=1}^N w(\theta_{0:m}^i) \int f(\alpha_{0:m}, \theta_{0:m}^i) p(\alpha_{0:m} / \theta_{0:m}^i, y_{1:m}) d\alpha_{0:m}}{\sum_{i=1}^N w(\theta_{0:m}^i)} \quad (23)$$

Under condition that $\int f(\alpha_{0:m}, \theta_{0:m}^i) p(\alpha_{0:m} / \theta_{0:m}^i, y_{1:m}) d\alpha_{0:m}$, can be calculated analytically.

IV. RBPF FOR ENHANCING NOISY SPEECH DFT CHANNELS

The proposed method table I, Consist in to filtering each frequency trajectory, were the real and imaginary part are processed separately, under the assumption of their independence. The filtering is done using the Rao-Blackwellization sequential importance sampling filtering algorithm, presented in [6] to enhance the speech in the time domain.

Depending on the adopted model, the estimation of the hidden state requires the definition of model parameters. Let: $\theta_{m,k} = [a_{m,k} \quad \sigma_{v_{m,k}}^2]$ be the vector of parameters, with $a_{m,k} = [a_{p,m,k} \quad \dots \quad a_{1,m,k}]$ is the vector of prediction coefficients. Thus for each channel part we have:

$$p(\alpha_m, \theta_{0:m} / y_{1:m}) = p(\alpha_m / \theta_{0:m}, y_{1:m}) p(\theta_{0:m} / y_{1:m}) \quad (24)$$

So, from section III., if $p(\theta_{0:m}, y_{1:m})$ is defined using SMC, the resulting model is linear and Gaussian conditionally

to each parameter vector and for $f(\alpha_{0:m}, \theta_{0:m}^i) = \alpha_{0:m}(\theta_{0:m}^i)$, the equation (23), becomes:

$$\frac{\sum_{i=1}^N w(\theta_{0:m}^i) \int \alpha_{0:m}(\theta_{0:m}^i) p(\alpha_{0:m} / \theta_{0:m}^i, y_{1:m}) d\alpha_{0:m}}{\sum_{i=1}^N w(\theta_{0:m}^i)} \quad (20)$$

Which is equivalent to the modulation of the Kalman filter output, by the weight corresponding to each vector parameter.

Thus:

$$\hat{\alpha}_{m,k}^{RBPF} = \sum_{i=1}^N W_m^i \int \alpha_m(\theta_{0:m}^i) p(\alpha_m / \theta_{0:m}^i, y_{1:m}) d\alpha_m \quad (26)$$

A posteriori marginal density $p(\theta_{0:m}, y_{1:m})$, [21], can be written as follows:

$$p(\theta_{0:m}, y_{1:m}) = p(\theta_{0:m-1}, y_{1:m-1}) p(\theta_m, y_m / \theta_{0:m-1}, y_{1:m-1}) \quad (27)$$

And for a conditionally linear Gaussian state space model, the equation (27), becomes:

$$p(\theta_{0:m}, y_{1:m}) = p(\theta_{0:m-1}, y_{1:m-1}) p(y_m / y_{1:m-1}, \theta_{0:m}) p(\theta_m / \theta_{m-1}) \quad (28)$$

Thus, if the a priori probability density function is adopted as importance density, i.e.

$$q(\theta_{0:m} / y_{1:m}) = p(\theta_{0:m}) \quad (29)$$

The weight of the particles in (25) is written using (18) as following:

$$w_m^i = w_{m-1}^i p(y_m / y_{1:m}, \theta_{0:m}^i) \quad (30)$$

So, depending on the adopted model, the likelihood density is Gaussian, such as:

$$p(y_{1:m} / y_{1:m-1}, \theta_{1:m}^i) = \mathcal{N}(E[C_m \alpha_m / y_{1:m-1}, \theta_{1:m}^i] + \dots E[D_m v_m / y_{1:m-1}, \theta_{1:m}^i] \text{Var}[C_m \alpha_m / y_{1:m-1}, \theta_{1:m}^i] + \dots \text{Var}[D_m v_m / y_{1:m-1}, \theta_{1:m}^i]) \quad (31)$$

With; $\theta_{0:m}^i \sim p(\theta_{0:m})$, and according to (28) and taking into account the independence of the components of the parameter vector, we have:

$$p(\theta_{0:m}) = p(a_0) p(\phi_{U_0}) \\ p(\theta_m / \theta_{m-1}) = p(a_m / a_{m-1}) p(\sigma_{U_m}^2 / \sigma_{U_{m-1}}^2) \quad (32)$$

Thus, giving the resulting model, the probability density in (31) is derived from the prediction step on the Kalman filter. And once the parameters are resampled before the propagation step, (30) becomes:

$$w_m^i = p(y_m / y_{1:m-1}, \theta_{0:m}^i) \quad (33)$$

Furthermore, the klaman filter output (a posteriori covariance, a posteriori mean) are towers duplicated relatively to the weight that they generate.

The parameters are assumed to evolve randomly according to Gaussian random walk with first order Markov model. The Stability of the TVAR model, is maintained by keeping the instantaneous poles of the model lie strictly within the unit circle. Similarly the variance of the excitation, is assumed to evolve according to a Gaussian random walk, for keeping it positive, propagation is done on its logarithm.

i.e. $\phi_{U_m} = \log \sigma_{U_m}^2$, $p(a_0) = \mathcal{N}(0, \Delta_{a_0} I_{p \times p})$ under condition to be inside the region of convergence, $p(\phi_{U_0}) = \mathcal{N}(0, \delta_{U_0}^2)$, $p(a_m/a_{m-1}) = \mathcal{N}(a_{m-1}, \Delta_a I_{p \times p})$ under condition to be inside the region of convergence, and $p(\phi_{U_m}/\phi_{U_{m-1}}) = \mathcal{N}(\phi_{m-1}, \delta_U^2)$.

The values $\{\Delta_{a_0}, \Delta_a, \delta_{U_0}^2, \delta_U^2\}$, are prefixed and are the same for all the channels parts.

Table. I RBPF Algorithm for enhancing real and imaginary DTF coefficient of the noisy speech channel.

<p>Determine the order of prediction for all channels.</p> <ul style="list-style-type: none"> - Fix the number of particles - Set the parameters $\{\Delta_{a_0}, \Delta_a, \delta_{U_0}^2, \delta_U^2\}$. - Segment with overlapping noisy signal - separate the real part from the imaginary part (later the treatment is the same for both sides). - for $m = 0$ (frame index) - for $k = 1 : L$ (channel index) - for $i = 1 : N$ (N the total number of particles) $a_{0,k}^i \sim p(a_0)$ (With stability condition) $\phi_{U_{0,k}}^i \sim p(\phi_{U_0})$ - end for i - initiate the matrix for equations (5) - end for k - for $m = 1 : T$ (T number of the frames) - achieve the DFT for the courant segment. - for $k = 1 : L$ - for $i = 1 : N$ - Run the prediction step of the Kalman filter - Calculate the weight according to (31) $\{w_m^i = p(y_m / y_{1:m-1}, \theta_{0:m}^i)\}_{i=1}^N$ - Run the correction step of the Kalman filter - end for i - calculates the sum of the weights - Normalize weight - Calculate the sum of the weighted outputs <li style="padding-left: 20px;">such as: $\hat{\alpha}_{m,k}^{RBPF} = \sum_{i=1}^N W_{m,k}^i \hat{\alpha}_{m,k}^{mmse,i}$ <li style="padding-left: 20px;">according to (05) we can just take : $\hat{X}_{m,k}^{RBPF} = \sum_{i=1}^N W_{m,k}^i \hat{X}_{m,k}^{mmse,i}$ - Resample according to the normalized weight: the
--

parameters and the covariance matrix and a posteriori mean for the next iteration

<ul style="list-style-type: none"> - for $i = 1 : N$ - $a_{m+1,k}^i \sim p(a_{m+1,k} / a_{m,k})$ (With stability condition) - $\phi_{U_{m+1,k}}^i \sim p(\phi_{U_{m+1,k}} / \phi_{U_{m,k}})$ - end for i - end for k - add the imaginary part to the real part - perform the IDFT of the current frame. - add with overlapping to the previous frame - end for m

V. SIMULATION AND RESULTS

The following speech enhancement algorithms have been considered. All the simulations are implemented in Matlab:

1- RBPF [6], with a TVAR (10) model and 1000 particles for each sample and $\{\Delta_{a_0} = 5 \times 10^{-3}, \Delta_a = 5 \times 10^{-3}, \delta_{U_0}^2 = 10^{-4}, \delta_U^2 = 5 \times 10^{-3}\}$ (windowing is not required), supposing that the noise variance are known at each instants

2- The proposed concept with a TVAR (2) for each DFT coefficient part, with 100 particles and $\{\Delta_{a_0} = 10^{-3}, \Delta_a = 10^{-3}, \delta_{U_0}^2 = 10^{-1}, \delta_U^2 = 10^{-1}\}$, using a Hanning window with a length of 16 ms, and a 50% overlapping between successive windows. The noise variance for each channel coefficient part are supposed to be known.

The clean signals sampled at 8 Khz are drawn from NOIZEUS database [14] (also available online). A white Gaussian noise (WGN), is added to clean signal, with different SNR.

The evaluation of the algorithms are presented in the table II. Where the expressed results for this algorithm; are deducted from an average of 10 trials.

The used objective measures are: CSII (coherence speech intelligibility index), LLR (log-likelihood ratio), PESQ (perceptual assessment of the speech quality) and the overalls SNR in dB, Segmental SNR in dB, the implementation are reported from [14].

From the table II we can see that the proposed concept outperform the time domain- RBPF speech enhancement, by exhibiting a large efficiency relatively to the resulting low log-likelihood ratio measure and high PESQ score.

Unlike the Time-RBPF, the DFT-RBPF don't present a residual noise relatively to the speech period activity at high frequency as depicted in fig. 02. Nevertheless, the proposed algorithm introduce signal distortion at spectral complements with low energy especially at very low peaked formant localization, But no musical noise are detected.

Table. II evaluation of the enhanced signal using the considered algorithm at different SNR.

		CSII	LLR	PESQ	O-SNR	SEG-SNR	
Input SNR dB	-5	Noisy signal	0.13	1.68	1.31	-4.95	-6.95
		Time-RBPF	0.68	1.44	1.94	6.77	2.58
		DFT-RBPF	0.51	0.85	2.25	6.44	-0.35
	0	Noisy signal	0.41	1.66	1.45	0.06	-4.67
		Time-RBPF	0.88	1.34	2.15	9.45	3.99
		DFT-RBPF	0.81	0.70	2.50	9.32	1.57
	+5	Noisy signal	0.73	1.55	1.77	4.95	-2.15
		Time-RBPF	0.96	1.19	2.49	12.27	5.58
		DFT-RBPF	0.93	0.62	2.66	12.28	3.54
+10	Noisy signal	0.92	1.39	2.06	10.07	0.78	
	Time-RBPF	0.98	1.04	2.68	15.64	7.65	
	DFT-RBPF	0.98	0.47	2.81	15.94	6.09	

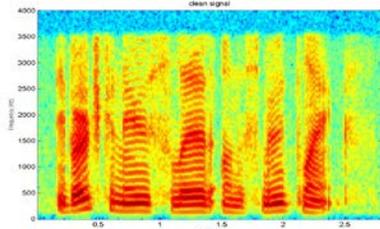


Fig. 1 Clean speech spectrogram.

Also, from the presented spectrogram. The proposed method present a noise at the beginning of the enhancement process. This observation can be explained by the adopted architecture itself, were each RBPF bloc attempted to find the alright region of high weight. To overcome that, this part can only be overlook were the most of the speech enhancement algorithm consider this period as silence period. Or at least, process a segment of the noisy signal before starting to process the all of the noisy signal, as shown in fig. 3.

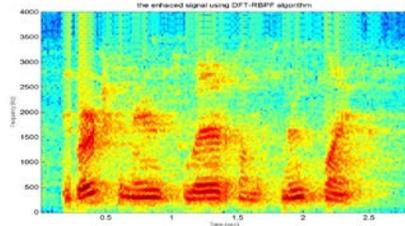


Fig. 1 Spectrogram of the DFT-RBPF enhanced 0dB noisy signal after processing an extracted segment of 0.5s

VI. CONCLUSION

The proposed concept performs better in terms of LLR measure and offer a better speech quality, as compared to the corresponding time-domain algorithm. The DFT-RBPF, introduce signal distortion at very low SNR, but no musical noise are observed. However this late penalize the intelligibility compared to Time-DFT, however the residual noise at this late is considerable at that level of comparison.

The proposed method, based on the TVAR model for each DFT, coefficient part, perform the basic Time-RBPF, where the proposed concept lead to more appearance of the formant peak at frequency domain witch, interpret the resulting performance.

Future works will be focalize in convenient estimation of the observation noise variance for the adopted concept, in goal to evaluate it performance for others noise model.

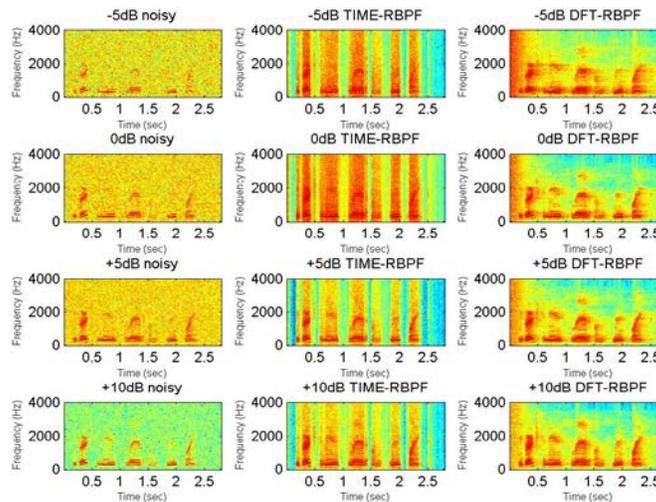


Fig. 2 Spectrogram of the enhanced signal

REFERENCES

- [1] R. J. McAulay and N. L. Malpass, "Speech enhancement using a soft-decision noise suppression filter", *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-28, pp.137 -145 1980
- [2] Y. Ephraim and D. Malah, "Speech enhancement using a minimum mean-square error short-time spectral amplitude estimator", *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-32, pp.1109 - 1121 1984
- [3] Y. C. Ho and R. C. K. Lee "A Bayesian approach to problems in stochastic estimation and control" *IEEE Transactions on Automatic Control*, vol. AC-9, pp.333 -339 1964
- [4] Karim Dahia. " Nouvelles méthodes en filtrage particulière. application au recalage de navigation par mesures altimétriques", Thèse de doctorat de l'Université Joseph Fourier, Janvier 2005
- [5] K. K. Paliwal, A. Basu, "A Speech Enhancement Method Based on Kalman Filtering", *Proc. ICASSP'87*, Dallas, Texas.
- [6] J. Vermaak, C. Andrieu, Doucet A., and Godsill S.J., "Particle methods for Bayesian modeling and enhancement of speech signals," *IEEE Trans. Speech, Audio Proc.*, March 2002.
- [7] B. Koo , J. D. Gibson and S. D. Gray "Filtering of colored noise for speech enhancement and coding", *IEEE Trans. Signal Processing*, vol. 39, pp.1732 -1742 1991
- [8] S. Gannot, "Speech Enhancement," ch. "Application of the Kalman Filter in the Estimate-Maximize (EM) Framework," pp. 161-198, *Springer*, 2005.
- [9] Arnaud Doucet, Simon Godsill And Christophe Andrieu, "On Sequential Monte Carlo Sampling Methods for Bayesian Filtering," *Statistics and Computing*, vol. 10, no. 3, pp. 197-208, 2000.
- [10] F. Mustière, M. Bouchard, and M. Bolić, "Quality assessment of speech enhanced using particle filters," in *Proc. IEEE ICASSP, Honolulu, HI*, Apr. 2007, pp. 1197–1200.
- [11] F. Mustière, M. Bouchard, M. Bolić, "Low-cost modifications of Rao–Blackwellized particle filters for improved speech denoising" *Signal Processing*, 88 (11) (2008), pp. 2678–2692
- [12] B. Laska, M. Bolić, and R. A. Goubran, "Discrete cosine transform particle filter speech enhancement," *Speech Communication*, Volume 52 Issue 9, September, 2010, Pages 762-775
- [13] B. Laska, M. Bolić, and R. A. Goubran, "Particle Filter Enhancement of Speech Spectral Amplitudes", *IEEE Transactions on Audio, Speech and Language Processing*, Volume 18, Issue 8, November 2010, pp. 2155 – 2167.
- [14] Phillips C. Loizou, "Speech enhancement theory and practice"1st ed. Boca Raton, FL.: CRC, 2007. *Releases Taylor & Francis*, pp. 120-121.
- [15] E. Zavarehei, S. Vaseghi, Q. Yan, "Speech enhancement using Kalman filters for restoration of short-time DFT trajectories". in: *Automatic speech Recognition and Understanding (ASRU)*, IEEE Workshop, pp. 219–224. 2005
- [16] R. Martin, "Speech enhancement using MMSE short time spectral estimation with Gamma distributed speech priors", *Proc. IEEE Int. Conf. Acoustics, Speech, Signal Processing*, vol. I, pp.253 -256 2002
- [17] M. S. Arulampalam , S. Maskell , N. Gordon and T. Clapp, "A tutorial on particle filters for online nonlinear/non-Gaussian Bayesian tracking", *IEEE Transactions on Signal Processing*, vol. 50, no. 2, pp.173 -188 2002.
- [18] Arnaud doucet, "algorithmes monte Carlo pour l'estimation Bayésienne de modèle markovienne, application au traitement de signaux rayonnement" Thèse de doctorat de l'Université de Paris XI Orsay, France, Décembre 1997
- [19] Arnaud doucet, "Sequential Monte Carlo methods" http://videolectures.net/mlss07_doucet_smcm/#c5980
- [20] H. Van Trees, "Detection, Estimation, and Modulation Theory", vol. I, 1968 :Wiley, pp
- [21] M. Johansen, Nick Whiteley, Arnaud Doucet, "Exact approximation of Rao-Blackwellised particle filters". *System Identification* 16 488-493. 2012

Swarm Intelligence Optimization of Lee Radio-wave Propagation Model for GSM Networks in Irbid

M. S. H. Al Salameh¹
 Department of Electrical Engineering
 American University of Madaba
 King's Highway, Madaba, Jordan
 m.salameh@aum.edu.jo

M. M. Al-Zu'bi
 Research Assistant
 Jordan University of Science and Technology
 Irbid 22110, Jordan

Abstract— Measurements for GSM cellular phone networks (the 1800 MHz as well as the 900 MHz bands) in different areas of Irbid city, Jordan, were carried out, for over a year and under varying weather conditions, by the authors. To find a radiofrequency propagation model that can correctly predict the propagation in this environment, various path loss models are compared with the measurements. These models include: Lee, COST-231 Hata, COST-231 WI, and Egli models. The results show that the COST-231 Hata model is the most accurate model whereas Lee model is the least accurate model for Irbid city. To enhance the accuracy of the least accurate model for this area, i.e., Lee model, this paper suggests a path loss model based on optimizing Lee model using the swarm intelligence optimization (PSO) technique. It is worth noting that optimization of Lee model for Irbid measurements using the least squares method produced the same results as the PSO optimization; therefore, only the PSO results will be presented here. The accuracy of the optimized Lee model is verified by comparison with measurements in other locations in Irbid. Furthermore, measurements made in Amman city in Jordan confirm the usefulness and validity of the measurements and predictions of Irbid city. The root mean square error (RMSE) between the measured and predicted values for the proposed model is significantly improved by up to 37 dB compared with Lee model.

Keywords— Lee model, Wireless; Measurements; Path loss; Radiowave; Swarm intelligence.

I. INTRODUCTION

Cellular phone networks are the most widely used systems all over the world. In order to efficiently plan new communication networks and improve the existing networks, accurate radio frequency (RF) path loss models are required. Although various models are available, only some of these models will match the area considered because every model was derived from measurements performed for certain area conditions. Therefore, we need to determine the appropriate model that can accurately predict the path losses in Jordan.

Different path loss models were investigated and then modified to match different environmental conditions in the world. For example, Lee model was optimized for global

system for cellular communications at 850 MHz frequency band using the least squares (LS) method [1], where the radio frequency measurements were collected using commercial measuring equipments in suburban and urban areas with flat terrains in Florida, USA. In [2], Lee model was calibrated by the least squares method for Jiza town in Jordan based on data supplied by mobile operators in Jordan. Another study was conducted on cellular communications at 900 MHz for different areas in Istanbul in Turkey where the Bertoni-Walfisch model was optimized using the mean square error (MSE) method [3]. Hata model was optimized for mobile communications using the least squares method for suburban areas within Cyberjaya and Putrajaya areas in Selangor state in Malaysia [4]. Different path loss models, including Lee model, were compared for mobile communications in Kuala Lumpur, Malaysia [5]. Okumura Model was optimized by the use of the regression fitting method and measured data for communication network in urban area in Kuala Lumpur in Malaysia [6]. Hata model was tuned and fitted to measurements using the mean square error method for cellular communications in Brno area in Czech Republic [7]. In [8], the COST-231 Walfisch-Ikegami (WI) path loss model was tuned by the particle swarm optimization (PSO) method for communication networks in the south-western part of Amman, Jordan. The COST-231 Hata model was optimized for communication networks in Banciao city, Taiwan, by fitting this model with measured data using the dual least-squares approach [9]. Also, Okumura model was found to be suitable for cellular communication systems in the suburban area Pathum Thani of Thailand [10].

This paper is intended to find a radiofrequency propagation model that can correctly predict the path loss in the environment of Irbid city and extend this model to other areas in Jordan. To that end, various path loss models are compared with extensive measurements, at both bands of GSM cellular communications (the 1800 MHz as well as the 900 MHz bands), performed by the authors in different areas of Irbid. These models include: Lee, COST-231 Hata, COST-231 WI, and Egli models. The results show that the COST-231 Hata model is the most accurate model whereas Lee model is the least accurate model for Irbid city. To improve the accuracy of the least accurate model, i.e., Lee model, for Jordan environment, this paper suggests a path loss model based on optimizing Lee model using the swarm intelligence

¹ On leave from Jordan University of Science and Technology, Irbid 22110, Jordan, salameh@just.edu.jo

optimization technique. The accuracy of the optimized Lee model is verified by comparison with measurements in other locations in Irbid. Furthermore, measurements made in Amman city in Jordan confirm the usefulness and validity of the measurements and predictions of Irbid city. The root mean square error (RMSE) between the measured and predicted values for the proposed model, is significantly improved by up to 37 dB compared with the Lee model.

II. PATH LOSS MODELS AND THE DRIVE TEST

The path loss models investigated here are COST-231 WI [11, 12], Lee [13], [14], [15], COST-231 Hata [7, 9], and Egli models [5]. More detailed information about these models can be found in [7], [9], [11],[12], [5], [13], [14], [15].

In this paper, the received power from the base station is measured for different distances. After that, the path loss in dB is calculated from the measured received power using the equation [11]:

$$L = P_t + G_t + G_r - P_r - L_t - L_r \quad (1)$$

Where, P_t is the transmitter power, P_r is the received power, G_t is the gain of the transmitting antenna, G_r is the gain of the receiving antenna, L_t is the feeder losses of the transmitter (e.g., connector losses and cables), and L_r is the feeder losses of the receiver (e.g., cables and body losses of the car used in the drive test). Here, these parameters have the values: $P_t= 42$ dBm, $G_t= 18$ dB, $G_r= 2.15$ dB, $L_t= 2$ dB, and $L_r=8$ dB. The loss L_r is due to the penetration loss through the car body with an average value of 8 dB; based on experiments. This value of L_r is similar to the value provided by [16], [17]. The values of L_t , P_t , and G_t were obtained from the operators of the cellular communications in Jordan.

The parameters used in this paper have the following meanings: h_b is the transmitter antenna height in meters, f_c is the operating frequency in MHz, d is the distance between the transmitter and the receiver in km, h_m is the height of the receiving antenna in meters.

The measurements have been performed, for over a year, by using RF measuring tools while driving a car (drive test) on many paths in Irbid city around the GSM (Global System for Mobile Communications) cellular phone base stations at the 1800 MHz and the 900 MHz frequency bands. The drive test is performed to measure the received signal strengths from the base stations, from which the path losses are calculated by equation (1). The measuring tools consist of TEMS (Test Mobile System) RF measuring software [18], GPS receiver, Laptop, and mobile phone. The mobile phone is equipped with RF measuring firmware [19] in order to extract the received RF signal strengths and send these readings to the laptop. The measured data involve the received signal strength levels of the serving base stations for each ARFCN (Absolute Radio Frequency Channel Number) scanned channel, cell-ID, and mobile station (MS) location coordinates.

III. OPTIMIZING THE PROPAGATION MODEL

A. General Form of the Path Loss Model

The purpose of the optimization process is to improve the accuracy of the path loss model in order to fit the area considered. To that end, the general form of the path loss according to Lee model is expressed as follows [14, 15]:

$$L = A + B \log(d) + 10 n \log \frac{f_c}{900} - 10 \log(\alpha) \quad (2)$$

Where f_c is the operating frequency in MHz, d is the distance between the transmitter and receiver in km, A is the path loss at reference distance, B is the slope of path loss curve in dB/decade, the parameter n is the frequency path loss exponent, and the factor α involves correction factors for the heights and gains of the transmitter and receiver antennas. More details on equation (2) can be found in [14] and [15].

Here, the main focus is on how to utilize the general form of the Lee path loss, equation (2), in order to obtain an optimized equation that matches our measurements. In other words, we need to find the optimum values of A , B , and n such that the predictions of equation (2) become as close as possible to the average of the measurements made in Irbid. The particle swarm optimization (PSO) technique is used to find the optimum values of the three parameters: A , B , and n .

B. Swarm Intelligence Optimization

The particle swarm optimization (PSO) algorithm is a global optimization method for the solution of non-linear problems [20]. The idea is related to the swarm intelligence of organisms such as swarms of bees, flocks of birds, schools of fish, and colonies of ants. Each particle in the swarm of the PSO method involves position vector (\mathbf{x}), velocity vector (\mathbf{v}) and personal best vector (i.e., best previous position encountered by each particle) and its fitness value. In the PSO algorithm, the initial velocity and initial position vectors of each particle are randomly assigned in n -dimensional search space. Each particle moves (i.e., modifies its velocity and position) according to its best previous experience in order to reach the best possible solution. The best solution reached by each particle is called the personal or local best (X_{pbest}). The best solution obtained among all the swarm of particles is called the global best (X_{gbest}). Each particle updates its velocity and position according to the following equations [21]:

$$v^{n+1}_{id} = C [\omega v^n_{id} + c_1 r^n_{1d} (X^n_{pbest\ id} - X^n_{id}) + c_2 r^n_{2d} (X^n_{gbest\ d} - X^n_{id})] \quad (3)$$

$$X^{n+1}_{id} = X^n_{id} + v^{n+1}_{id} \Delta t \quad (4)$$

Where,

- v^{n+1}_{id} , v^n_{id} : Velocity component along d^{th} coordinate of i^{th} particle at the $(n+1)^{th}$ and n^{th} iterations, respectively.
- X^{n+1}_{id} , X^n_{id} : d^{th} coordinate of i^{th} particle at the $(n+1)^{th}$ and n^{th} iterations, respectively.
- $X^n_{pbest\ id}$: Personal best position along d^{th} coordinate of i^{th} particle at n^{th} iteration.
- $X^n_{gbest\ d}$: Global best position along d^{th} coordinate at n^{th} iteration.
- Δt : Time step value; usually chosen to be 1 s.

- $i = 1, \dots, N_p$, and N_p is the number of particles in the swarm
- $d = 1, \dots, N_d$, where N_d is the search space dimension

The inertia weight ω and the convergence factor C in equation (3) are as follows [21]:

$$\omega = \omega_{\max} - \frac{(\omega_{\max} - \omega_{\min})}{N_j} j \quad (5)$$

$$C = \frac{2}{\left| 2 - \alpha - \sqrt{\alpha^2 - 4\alpha} \right|} \quad (6)$$

Where the maximum number of iterations is N_j , and the current iteration number is j . The parameter values in this paper are, $\omega_{\min} = 0$, $\omega_{\max} = 1$; these values are chosen in order to obtain largest convergence speed of PSO, where the experiments showed that the best value of ω_{\max} falls between 0.8 and 1.2.

The acceleration constants c_1 and c_2 are used to increase the new velocity towards the best solutions. The parameter $\alpha = c_1 + c_2$, and c_1, c_2 are chosen to be: $c_1 = 2, c_2 = 2$, accordingly $\alpha = 4, C = 1$. The experiments showed that, for best performance, $c_1 = c_2 = 2$. The random numbers, r_{1d}^n and r_{2d}^n , are uniformly distributed between 0 and 1.

Each iteration of the PSO algorithm includes evaluating the fitness value for each particle, where this value replaces the personal best value if the currently obtained fitness value is better than the personal best value. Similarly, the best fitness value obtained by the swarm replaces the global best value if the currently obtained fitness value among all particles is better than the global best value.

IV. RESULTS

The parameters of the optimized Lee model are determined by means of the particle swarm optimization (PSO) method. The fitness function of the PSO method is considered to be the root mean square error RMSE given by equation (7) below [5]. In other words, this root mean square error RMSE is used to evaluate the accuracy of the path loss models as compared with the measurements.

$$RMSE = \sqrt{\frac{\sum_{i=1}^N [L_{mi} - L_i]^2}{N-1}} \quad (7)$$

Where, N is the number of path loss data points, L_{mi} is the measured value of the path loss at position i in dB, and L_i is the predicted path loss at position i in dB. Due to the very large number of measured data samples for the site of each base station and in order to eliminate the effects of fast fading, the measured data were averaged over every 1 m of the path between the transmitting base station and the mobile receiver.

The swarm size, of the PSO method, is chosen in this paper to be 10; based on many trials to obtain accurate results. The optimal solution is reached, on average, after 25 iterations, for all the base stations. As an example, Fig. 1 shows the convergence of the fitness function (RMSE value) for base station 1, sector 1, where the RMSE decreases rapidly from 120 dB to 7.9 dB after about 25 iterations. The locations of the

cellular base stations in Irbid, for which the measurements were performed, are shown in Fig. 2.

Fig. 3 compares the existing path loss models, the measured path loss data, and the proposed optimized Lee model for base station 6, sector 1. The results of the other base stations 1, 2, 3, 6, 8, 10 are summarized in Table 1 which also shows the values of the optimized Lee model parameters, i.e., A, B , and n of equation (2). Table 1 shows significant improvement of the optimized Lee model as compared with Lee model. The improvement ranges between 17.6 dB for base station 10 sector 3 to 35 dB for base station 6 sector 1. It is to be noted here that the optimized Lee model was extracted from the measurements performed in these base stations, i.e., base stations 1, 2, 3, 6, 8, and 10. The measurements in the remaining base stations, i.e., base stations 4, 5, 7, and 9 are used to verify the optimized model in Irbid city, as shown in Table 2 and Fig. 4. Fig. 4 compares the existing path loss models, the measured path loss data, and the proposed optimized Lee model for base station 9, sector 2 which shows significant improvement in the accuracy of the optimized Lee model in comparison with the Lee model. The accuracy enhancement, as can be seen in Table 2, ranges between 14.5 dB for base station 9 sector 4 to 37.2 dB for base station 7 sector 2. Thus, the overall improvement for Irbid city associated with the proposed model compared with the Lee model ranges between 14.5 dB for base station 9 sector 4 to 37.2 dB for base station 7 sector 2.

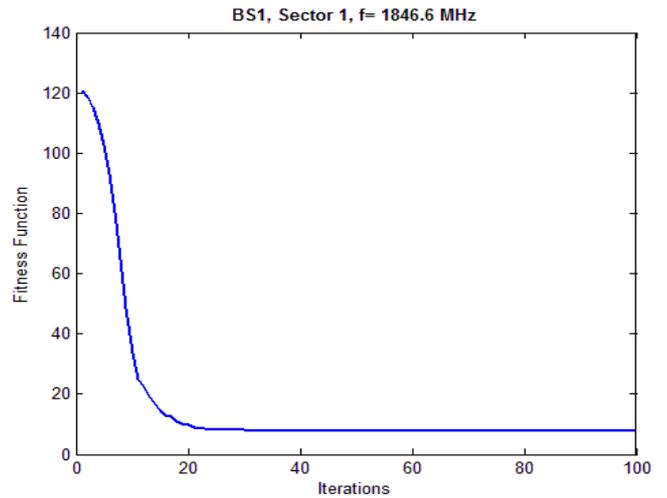


Fig. 1: Solution convergence of the PSO algorithm for base station 1. Other base stations have almost similar convergence.



Fig. 2: GSM Cellular Base Stations under study in Irbid City.

From Table 1 and Table 2, it is clear that the proposed optimized Lee model has the best average RMSE values compared with the other examined models. Substituting the average values of the optimized parameters (A, B, and n) from

the last row of Table 1 into equation (2), the proposed optimized Lee model can be written as follows:

$$L = 124.64 + 23.2 \log(d) + 26.94 \log \frac{f_c}{900} - 10 \log(\alpha) \quad (8)$$

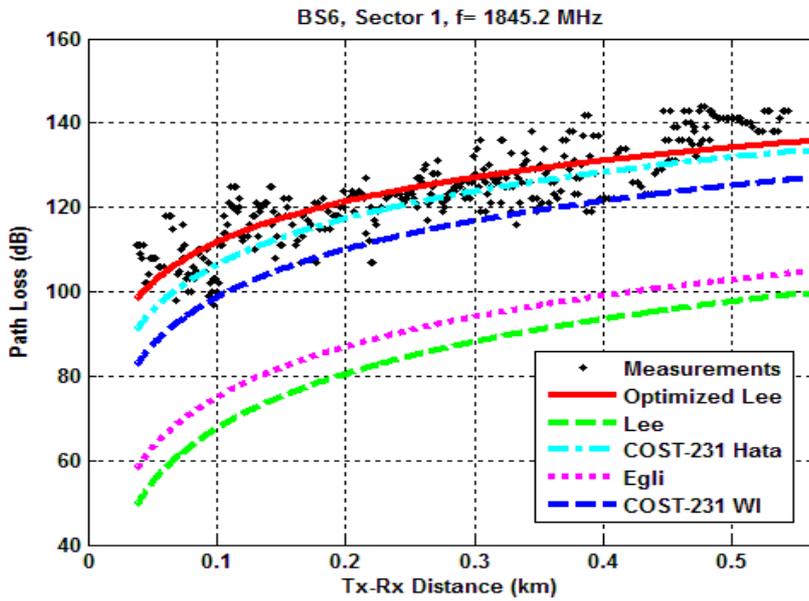


Fig. 3: Path loss vs. transmitter to receiver (Tx-Rx) distance in kilometres for base station 6 (BS6).

Table 1: The root mean square errors (RMSE) of the models compared with measurements, in addition to the optimized Lee model parameters for the base stations used to build the optimized Lee model.

Info		Root Mean Square Error (RMSE) dB					Optimized Lee model Parameters		
BS	Sector No.	Optimized Lee	Lee	Hata	Egli	WI	A	B	n
1	Sec. 2	8.2861	27.2091	13.8918	22.2750	10.0340	122.4484	24.6108	2.6
	Sec. 1	7.9004	32.9196	11.3244	27.7726	13.2643	125.8005	19.3576	2.6634
10	Sec. 1	10.0668	27.9662	14.1491	22.8772	10.7640	122.3311	27.8632	2.6
	Sec. 2	5.4516	31.3264	8.7931	25.4180	8.7172	122.9291	25.5341	2.7512
	Sec. 3	13.2322	30.8646	16.3376	25.7568	14.3335	119.7959	23.2512	2.8075
2	Sec. 3	7.0682	30.9001	9.1672	25.3527	7.79778	124.9121	28.4576	2.8592
	Sec. 2	7.4808	36.8599	7.8158	31.2849	10.6795	130.6544	26.0445	2.5480
	Sec. 1	6.2116	38.5488	7.9642	32.7864	10.7139	131.5981	25.3403	2.6
3	Sec. 2	7.4931	35.6255	9.6318	29.7714	9.7865	124.0177	20.7344	2.7671
	Sec. 1	8.8160	35.6336	12.0170	29.8440	12.0280	121.1776	15.7935	2.5
8	Sec. 1	8.3021	32.4521	10.5352	26.7466	9.8155	122.9000	23.8026	2.7602
	Sec. 2	8.0540	36.9021	8.5430	30.8139	10.8794	128.7991	27.9144	2.6
	Sec. 3	6.8632	39.0120	9.9891	32.9072	13.1312	121.4220	13.918	2.7943
	Sec. 4	5.7996	32.1626	11.7094	26.6179	10.9171	118.9979	12.8889	2.6
6	Sec. 1	5.8158	40.7997	7.0539	34.5448	12.5532	135.3671	32.1887	2.7016
	Sec. 2	10.4525	33.0037	12.7395	27.3376	11.9023	120.0212	19.9881	2.7
	Sec. 3	6.5744	39.5954	7.9736	33.2684	11.2226	128.3335	24.6807	2.7349
	Sec. 4	10.4085	33.4385	11.8711	27.5894	11.2466	121.9733	25.1409	2.9
Average		8.01538	34.1789	10.6393	28.4980	11.0992	124.6377	23.1950	2.6938

Table 2: The RMSE values of the models compared with the measurements for the base station sites which were used to validate the optimized Lee model.

Info		Root Mean Square Error (RMSE) dB				
BS	Sector No.	Optimized Lee	Lee	Hata	Egli	WI
7	Sec. 1	11.4694	42.4688	12.3352	36.0470	12.9086
	Sec. 2	8.1812	45.4074	12.4173	38.3363	11.5059
	Sec. 3	9.3195	44.5295	10.6993	37.9189	11.0074
4	Sec. 1	10.4311	44.8020	9.9025	38.4007	18.3551
	Sec. 2	11.4253	41.6590	10.6715	35.7684	16.8768
	Sec. 3	8.7615	39.9167	6.5816	33.8914	13.7637
	Sec. 4	8.7154	42.0617	10.1256	35.4723	16.1118
5	Sec. 1	9.7138	35.6519	10.6986	30.0986	14.4432
	Sec. 2	9.5448	30.1527	13.2996	24.9481	11.7719
	Sec. 3	12.6058	40.8144	10.5086	35.0020	18.0704
	Sec. 4	7.5550	37.5223	9.2033	31.7698	15.1924
9	Sec. 1	9.5604	28.5422	14.1265	23.1041	11.4169
	Sec. 2	8.9551	44.1671	11.5784	37.5854	16.2953
	Sec. 3	7.2975	39.6647	9.4433	33.5609	13.0553
	Sec. 4	13.3604	27.8148	17.0125	22.7166	14.0077
Average		9.7931	39.0117	11.2403	32.9747	14.3411

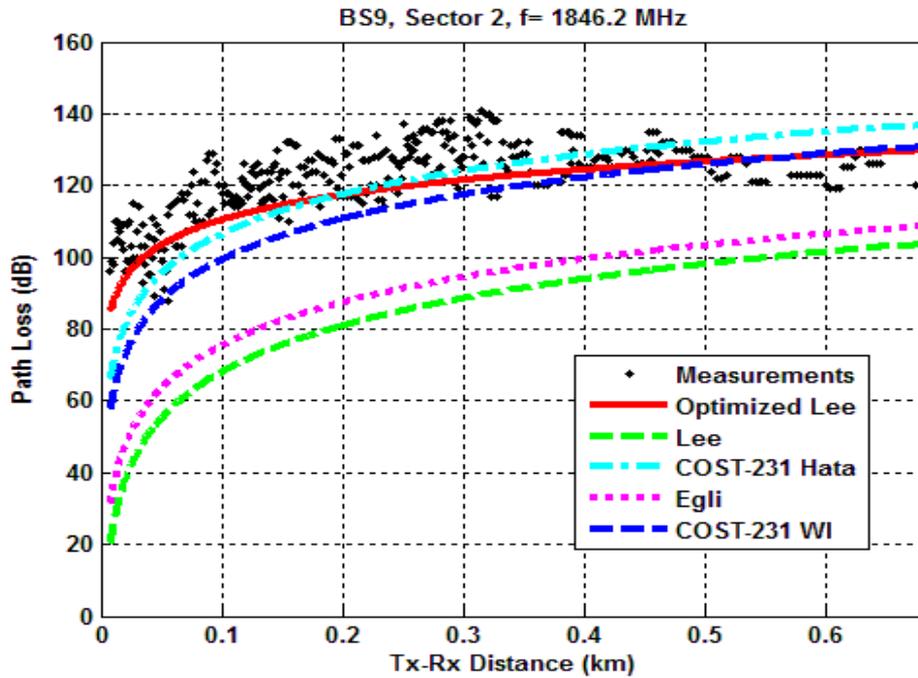


Fig. 4: Path loss models compared with measurements vs. transmitter to receiver (Tx-Rx) distance in kilometres for base station 9 (BS9) sector 2, used to verify the optimized Lee model.

V. CONCLUSIONS

This study proposes an optimized Lee path loss model for Irbid city in Jordan. The Lee model is optimized using the swarm intelligence technique and the measurements at several GSM base stations performed by the authors for over a year under varying weather conditions. It is worth mentioning that optimization of Lee model for Irbid measurements using the least squares method produced the same results as the PSO optimization; therefore, only the PSO results were presented in this paper. The validity of the optimized Lee model is demonstrated by comparison with measurements at other base stations. The proposed model showed high degree of agreement with the measurements where the accuracy of the proposed Lee model is improved by up to 37 dB compared with the Lee model.

Also, measurements performed in Amman city supports the results presented in this paper, and thus this optimized model is expected to suit other areas in Jordan. The proposed model may help in the design and expansion of the cellular communication networks.

REFERENCES

1. N. Mijatovic, I. Kostanic, G. Evans, "Use of scanning receivers for RF coverage analysis and propagation model optimization in GSM networks," 14th European Wireless Conference, Prague, Czech Republic, June 2008, pp. 1-6.
2. L.A. Nissirat, M. Ismail, M. Nisirat, M. Singh, "Lee's path loss model calibration and prediction for Jiza town, south of Amman city, Jordan at 900 MHz," IEEE International RF and Microwave Conference (RFM), Seremban, Negeri Sembilan, Dec. 2011, pp. 412-415.
3. B.Y. Hanci, I.H. Cavdar, "Mobile radio propagation measurements and tuning the path loss model in urban areas at GSM-900 band in Istanbul-turkey," IEEE 60th Vehicular Technology Conference, Los Angeles, USA, Sept. 2004, vol. 1, pp. 139-143.
4. R. Mardeni, K.F. Kwan, "Optimization of Hata propagation prediction model in suburban area in Malaysia," Progress In Electromagnetic Research, vol. 13, pp. 91-106, 2010.
5. J. Chebil, A.K. Lwas, M.R. Islam, A. Zyoud, "Comparison of empirical propagation path loss models for mobile communications in the suburban area of Kuala Lumpur," 4th International Conference On Mechatronics (ICOM), Kuala Lumpur, Malaysia, May 2011, pp.1-5.
6. R. Mardeni, Lee Yih Pey., "Path loss model development for urban outdoor coverage of code division multiple access (CDMA) system in Malaysia," International Conference on Microwave and Millimeter Wave Technology (ICMMT), Chengdu, China, May 2010, pp. 441-444.
7. L. Klozar, J. Prokopec, "Propagation path loss models for mobile communication," In Proceedings of 21st International Conference Radioelektronika, Brno, Czech Republic, April 2011, vol. 48, pp. 1-4.
8. A. Tahat, M. Taha, "Statistical tuning of Walfisch-Ikegami propagation model using particle swarm optimization," IEEE 19th Symposium on Communications and Vehicular Technology in the Benelux (SCVT), Eindhoven, Netherlands, Nov. 2012, pp. 1-6.
9. Y.H. Chen, K.L. Hsieh, "A dual least-square approach of tuning optimal propagation model for existing 3G radio network," IEEE 63rd Vehicular Technology Conference (VTC), Melbourne, Australia, May 2006, vol. 6, pp. 2942-2946.
10. W. Bhupuok, K. Dejhan, "A new method for prediction 3G path loss propagation in suburban of Thailand," The International Conference on Electrical Engineering/ Electronics, Computer, Telecommunications and Information Technology (ECTI-CON), Phetchaburi, Thailand, May 2012, pp. 1-4.
11. S. R. Saunders, Alejandro Aragón-Zavala., Antennas and propagation for wireless communication system, 2nd ed. Wiley, 2007.
12. M.V.S.N. Prasad, S. Gupta, M.M. Gupta, "Comparison of 1.8 GHz cellular outdoor measurements with AWAS electromagnetic code and conventional models over urban and suburban regions of northern India," IEEE Antennas and Propagation Magazine, vol. 53, pp. 76-85, 2011.
13. Lee WCY. Mobile communications engineering: theory and applications. 2nd ed. New York: McGraw-Hill; 1997: 689.
14. Gordon LS. Principles of mobile communication. 3rd ed. Germany: Springer; 2012: 819.
15. Alshami M, Arslan T, Thompson J, Erdogan AT. Frequency analysis of path loss models on WIMAX. 3rd Computer Science and Electronic Engineering Conference (CEEC); 2011 July 13-14; Colchester, UK. 1-6.
16. H. Holma, A. Toskala, "WCDMA for UMTS: HSPA evolution and LTE," 5th ed. John Wiley & Sons, 2010.
17. H. Holma, A. Toskala, "LTE for UMTS: OFDMA and SC-FDMA based radio access," UK: John Wiley & Sons, 2009.
18. TEMS Investigation [Online]. Available: <http://www.ascom.com/nt/en/index-nt/tems-products-3/tems-investigation-5.htm#overview>
19. TEMS Pocket [Online]. Available: <http://www.ascom.com/nt/en/index-nt/tems-products-3/tems-pocket-5.htm#overview>
20. J. Kennedy, R. Eberhart, "Particle swarm optimization," IEEE International Conference on Neural Networks, Perth, WA, Nov./Dec. 1995, vol. 4, pp. 1942-1948.
21. Q. Bai., "Analysis of particle swarm optimization algorithm," Computer and Information Science, vol. 3, pp. 180-184, 2010.

A Recognition and Synthesis Environment for The Arabic Language

Tebbi Hanane, Hamadouche Maamar and Azzoune Hamid

Abstract: In this work we present a Human-Machine Interface in which we investigate both of the automatic recognition process and the automatic synthesis process applied on the Arabic digits in one system that we named ARSSAD (Automatic Recognition and Synthesis System of Arabic Digits). The system is therefore divided into two sub-systems; a recognizer and a synthesizer. The main task of the recognizer is the automatic recognition of the pronounced digit, so it transforms the input sound wave into a text representing the appropriate digit, the second sub-system perform the opposite process of the first sub-system; indeed, it transforms the text (digit) produced by the recognizer to a speech generated by the synthesizer. The methodology used for the system design is based on three essential stages: the creation of the acoustic database (corpus development), the recognition of the read signal and the generation of the synthetic speech. We explain the basics modules that compose it, starting from the signal acquisition and finishing by the decision taken. For the recognition sub-system we make the choice to use the Dynamic Time Warping (DTW) method for the comparison task. ARSSAD contains a front-end and a back-end module, the front-end module converts the input sound into feature vectors that are based on Mel Frequency Cepstral Coefficients (MFCCs), to be used in the DTW method. The back-end module uses the Concatenative method to perform the synthesis of the recognized digit, for this end we create a sound database that contained diphones of the Arabic alphabets. The obtained results show that the system presents a success rate of 94.85% on the three corpuses which we recorded in a noised environment.

Keywords: analysis techniques, speech recognition, speech synthesis, synthesis by diphones, synthesis by phonemes, PRAAT, MFCC, DTW, Standard Arabic.

I. INTRODUCTION

THE automatic speech processing (ASP) is an area of research for which a significant effort has been undertaken over the past three decades. The challenges are considerable and have fundamental nature. They are also multidisciplinary: signal processing, pattern recognition, artificial intelligence, computer science, phonetics, linguistics, ergonomics and neurosciences; which behave at varying degrees in the solutions found.

H.Tebbi, Author is with the USTHB University, Algiers, Algeria, corresponding author, phone: +21369990892; e-mail: tebbi_hanane@yahoo.fr; htebbi@usthb.dz.

M.hamadouche, Author, e-mail : hamadouchemaar@yahoo.fr.

H.Azzoune, is now with the Department of Electronics and Informatics, USTHB University, Algiers, Algeria, e-mail: Hazzoune@usthb.dz

These long-standing works nevertheless give birth at the present time to intermediate products which find their place in practical applications in the context of the Man-machine communication, as shown in [1], [2] and [3].

However, Automatic Speech Recognition/Synthesis (ASRS) systems dedicated to the Arabic language are at the moment still very modest. In this article, we will introduce our ASRS system of the first ten digits of the Standards Arabic language (SA) in mono mode speaker. We are interested exclusively to the step of analysis of the speech signal which allows us to extract the acoustic vectors characterizing it. This step is very important and primordial in the process of automatic recognition, since it produces in output a set of parameters considered pertinent and efficient for the high-quality operation of the speech signal, on this same set we will apply the algorithms of recognition and comparison.

In speech recognition, the step of feature extraction, commonly known as the step of analysis, can be achieved in several ways. Indeed, the acoustic vectors are usually extracted using methods such as temporal encoding predictive linear (Linear Predictive Coding LPC) or Cepstral methods as the MFCC encoding (Mel Frequency Cepstral Coding), as well as the encoding PLP (Perceptual Linear Predictive coding) which is an example of the application of knowledge of the auditory system in human speech recognition. The extraction of characteristics is a key element for the development of an ASR system.

The second part of our system represents a Text To Speech (TTS) system, in which the main, commonly known, techniques used in it design are: Articulator synthesis, Formant synthesis, and Concatenative synthesis [4]. Articulatory synthesis attempts to model the human speech production system directly.

Formant synthesis, which models the pole frequencies of speech signal or transfer function of vocal tract based on source-filter-model. Concatenative synthesis, which uses different length pre-recorded samples derived from natural speech.

In our case, we have used the concatenation method for the synthesis implementation which represents, in our opinion, the method that produces a synthetic voice the most natural and intelligible compared to the others. This result came from the fact of using a set of recorded units pronounced by a real speaker, priory collected and embedded within our sound database.

So, for the recognizer, we have to deal with two essential problems, the first one is the choice of the technique of analysis used, and the second one is the

choice of parameters and their number to extract the relevant parameters of the voice signal. The purpose is to determine which gives the best recognition rate. Whereas, for the synthesizer, we have to face two other problems; the choice of the transcription method (rule-based method or lexicon-based method) in one hand, and the co-articulation problem to improve the quality of the generated speech, in the other hand.

II. SYSTEM DESIGN

When designing a system, two broad ways could be taken into account, the first one is to design the whole system using the known theories, and use it as it is designed, in the real conditions. An alternative way would be to subdivide the system into modules that can be independently created and tested, to eventually be used in other systems to perform several functionalities.

To facilitate the implementation and improvement of our system, we have used the modular approach; this concept makes the program understandable on one part and decreases the cost of development of each module in another part. We have also used the concept of the object-oriented programming which is particularly suitable with the modular technique. We must therefore make out different modules which structure the system as shown in the following diagram (Fig.1):

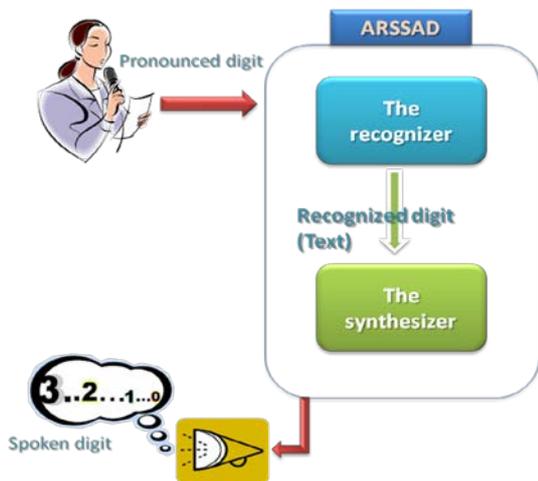


Fig. 1. General architecture of our system ARSSAD

The objective here is to describe the role of each module, explaining in the same time the interest of links which provide the cooperation between them.

III. THE RECOGNIZER

This module represent the front-end of the whole system, it is also composed of a set of sub-modules that can be shown in the following diagram (Fig.2)

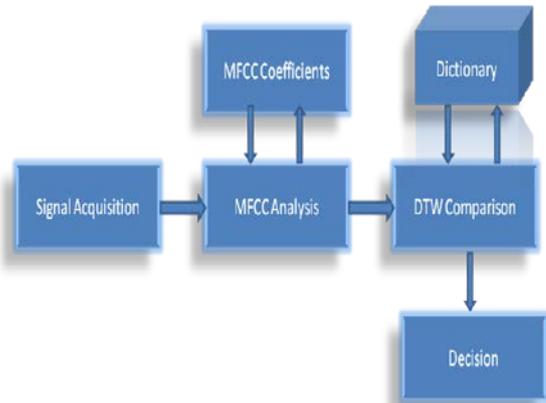


Fig. 2. General schema of the recognizer

We give now the principal functionalities of each sub-module one after another.

A. Signal Acquisition

This module carries out the acquisition of the acoustic signal recorded by a microphone and converts it into a digital form that can be used directly by a machine. There are many types of microphones but all of them provide the same function: transform the pressure fluctuations caused by the acoustic wave of speech into an electrical signal.

This signal will be converted from the analogical form to the digital one, i.e. it will be discrete both in time (sampling) and value (quantification) [5]. As a result we obtain a digital signal in the form of a sequence of samples which measure the amplitude of the microphone's signal in regular spaced moments, and the amplitude of each sample is represented in its digital form. The choice of the sampling frequency is usually determined by the application and referred to the platform used [5].

B. Sampling Frequency

Some thoughts on the frequency of sampling are required in first. According to the theorem of Shannon [6]: « a bandlimited function can be perfectly reconstructed from a countable sequence of samples if the bandlimit, B , is no greater than half the sampling rate (samples per second)".

Sounds that are made by the human voice normally contain relatively insignificant frequency components at or above 10 kHz [6]. Sampling such an audio signal at 20k sam/sec, or more, provides an excellent approximation to ensure that the Shannon criterion is met. But often the sample-rate is pre-determined by other considerations, such as an industry standard format (e.g. 8k sam/sec).

In this situation, the human voice should be filtered, to remove frequency components above 4 kHz, before being sampled. So we consider in our work that the acoustic signal is located mainly in the bandwidth (50 Hz -8 kHz), the frequency of sampling should therefore be at least equal to 16 kHz, according to the theorem of

Shannon. For the case of our application, we have used a sampling frequency of the order of 22050Hz, the default value taken by the software used in this operation PRAAT [7].

C. The corpus preparation

Most of the works carried out in the field of Man-Machine communication often require the registration, and the manipulation of corpus of continuous speech, and this to carry out the studies on the contextual effects, on the phonetic indices, and on the variability intra and inter speakers. There were three recorded corpuses each one containing ten sounds of ten prime numbers of Standards Arabic (Wahid (one), Ithnane (two), Thalatha (three), Arbaa (four), Khamsa (five), Sita (six), Sabaa (seven), Thamania (eight), Tisaa (nine), Aachara (ten)) in a noisy environment and we have changed the speed of elocution from a corpus to another without changing the speaker. The step of analysis may therefore begin.

D. The Cepstral Analysis (MFCC)

The aim of the analysis of the voice signal is to extract the acoustic vectors which will be used in the stage of recognition follows. In this step the voice signal is transformed into a sequence of acoustic vectors on the way to decrease redundancy and the amount of data to be processed. Then a spectral analysis by the discrete Fourier transform is performed on a signal frame (typically of size 20 or 30 ms) [6]. In this frame, the voice signal is considered to be sufficiently stable and we extract a vector of parameters considered to be enough for the good operation of the voice signal, in our work we choose to use of MFCCs coefficients resulting from a Cepstral analysis of the read signal.

The method of extraction of the MFCCs coefficients is a famous method used for the acoustic vectors extraction, in the field of automatic speech recognition. We have also decided to use it in our context of application and we chose a set of 12 coefficients. We expose the different steps leading to Cepstral analysis using the tool of speech analysis, PRAAT, we show the different parameters required for the analysis that we have chosen, and in the end the exploitation of the MFCCs coefficients resulting.

1) Step1: reading the file to analyze and the choice of MFCC method

- i) Start PRAAT
- ii) Open the sound file:
- iii) Read > Read from file (open a sound file)
- iv) Edit (for the view)
- v) File > Extract Selection (for "cutting" the sound)
- vi) Write > Write to .wav file (to save a sound file)
- vii) Select the file to analyses
- viii) Choose the Cepstral method:
- ix) Formants&LPC > To MFCC

2) Step 2: determination of the parameters

required for the analysis

- i) Number of coefficients: 12
- ii) Duration of windows: 30 ms
- iii) Duration between the windows: 10 ms

3) Step 3: analysis Results

It remains now to save the results in a text file format with the extension .MFCC (Write > Write to txt file), to be used in the following stage.

E. The use of DTW method

Our speech recognition sub-system is based on the algorithm of DTW (Dynamic Time Warping), this method is based on an evaluation of the distance between an observation and a list of references (dictionary). As well the reference for which this distance is minimal allows us to decide what word is it. The evaluation of the distance between two signals is not performed with the signals themselves. This would lead to lot of calculations. It is therefore in a prime time to find a better representation of the signals. Here MFCC analysis shines.

So we have programmed the DTW method using, for the comparison, the MFCC coefficients. The training part concerns the recording of the sounds corpuses in order to design our dictionary which will be used as reference in the comparison of the signals tested. Problems of recognition may appear depending on the conditions in which the signal to test is recorded. If the word is pronounced more or less close to the microphone, recognition rates can vary greatly. However if the user says the word always at the same distance and with the same intensity, the rate of recognition is very acceptable. We judge, though, that the representation using the MFCC coefficients provides better results, and it supports better the limitations related to the problem of the capture of the signal. The common skeleton of the DTW algorithm has three steps illustrated as follow:

- 1) Acquisition of the sound file to test
- 2) Extraction of the MFCC coefficients
- 3) Comparison with the dictionary of references

F. The decision

This last module of our recognizer plays two essential roles; it represents the interface in which the user interacts with the system. After the user has entered his voice signal, he starts the search and awaits the results. The system displays the recognized digit written in both Arabic and French language.

The second role is that this same decision (the displayed digit) represents the input (text) of the second module of our system, which is the synthesizer.

IV. THE SYNTHESIZER

It is based essentially on two principal parts; a front-end and a back-end. The front-end is composed of two modules, the first is for the sound database creation

and the second is for the conversion text-to-phoneme or grapheme-to-phoneme. The back-end part represents the speech generation module or in other words the synthesizer itself. So the different modules that compose the system are as follow:

- 1) The sound database creation (segmentation): we have recorded a set of pieces of speech and store it in our database, this set is composed of phonemes and diphones which are the basic units utilized within the back-end module in order to generate voice using the concatenation method.
- 2) The grapheme/phoneme conversion: before achieving this process, a text normalization or preprocessing operation has to be done. After that the module assigns to each word in entry its phonetic transcription, and then divides and marks the text into prosodic units like syllables. This process of assigning phonetic transcription to words is called text-to-phoneme or grapheme-to-phoneme conversion. The output of the front-end module is a symbolic linguistic representation resulting from the phonetic transcription and prosody information together, which represents the input of the back-end module.
- 3) The synthesizer: the back-end module uses information provided by the front-end to convert the symbolic linguistic representation to speech using a specific method. In literature, there are two kinds of synthesis methods; rule-based method and Concatenative corpus-based method.

Like we have mentioned before, we have used the Concatenative method of phonemes and graphemes previously stored in our sound database. The general architecture of this module could be shown in Fig.3 as follow:

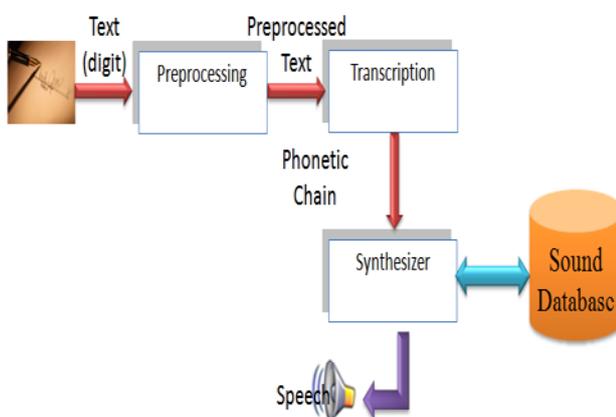


Fig. 3. General architecture of the synthesizer

A. The corpus description

We have created two corpora; the first contains phonemes: It is composed of a set of basic sounds (which consists of the phonemes corresponding to the 28 consonants and 6 vowels, and other additional characters (corresponding to the three sounds of tanwiin ([an], [a], [in]), and the silence).

To improve the quality of the words synthesized by the method of concatenation of phonemes, and to reduce the effects of co-articulation, the solution is to record the transition that exists between phonemes instead of recording the phonemes themselves; we talk about diphones which are an adjacent pair of phones.

Indeed, the transition (diphones) is the bearer of a significant quantity of acoustic information in relation to the phoneme itself. Each transition or dihone also varies from the stable part of a phoneme up to the stable part of phoneme that follows.

B. The phonetic and orthographical transcription "POT":

Transcription provides a phonetic text from the alphabetic text. To accomplish this, it must apply to many pronunciation rules. French language has a few thousands of basic rules; English language has tens of thousands of rules. Therefore, during the passage from the written form to the spoken form two approaches can be used which are: the lexicon-based approach and the rule-based approach [8] [9];

- 1) The use of rules

In this approach each grapheme is converted to phoneme depending on the context and this is thanks to the use of a set of rewriting rules [10]. The main advantage of this approach is the ability to model the linguistic knowledge of human beings by a set of rules that can be incorporated into expert systems. Each of these rules has the following form:

[Phoneme] = {LC (Left Context)} + {C (Character)} + {RC (right context)}

Our transcription module grapheme-phoneme is based on a set of rules;

The rule of tanwin, al madd, etc... Prioritized, and organized in the form of a tree list. Each rule is written in the graphics context in which it is applied.

Here is a concrete example of transcription rule "The rule of Tanwin":

```

If (grapheme [char] == 'Tanwin')
{If (API [position] [0] == ' ') ◌
Phoneme = phoneme + "an";
Else
{If (API [position] [0]
[= = ' ') ◌
Phoneme = phoneme + "in";
Else
Phoneme=phoneme+ "a";}
}
    
```

- 2) The use of the lexicon

In this case we must assign to each word in entry the pronunciation which corresponds to it without taking into account its context. The speed, flexibility and simplicity are the main advantages of this approach.

C. The acoustic generation sub-module

This is the heartbeat of the synthesizer module, in fact the user after that he had pronounced the digit he will see the recognized digit displayed on the screen, and will hear the system spelling back the recognized digit. This is the task of the acoustic generation sub-module.

To accomplish this task, we have implemented a reading function which is exposed below:

```

Position= seek (grapheme [ig], API);
If ((grapheme [ig] == '1') && (grapheme [ig+1] ==
'ل'))
{
MP2- >FileName= "C: \\son_hanane\\alif.wav";
MP2- >Open ();
MP2- >Wait=true;
MP2- >Play ();
IG=ig+2;
Position=seek (grapheme [ig], PLC);
If (API [position] [1] == ')
{
MP2- >FileName= "C: \\son_hanane\\l.wav";
MP2- >Open ();
MP2- >Wait=true;
MP2- >Play ();
MP2- >FileName=API [position] [2]);
MP2- >Open ();
MP2- >Wait=true;
MP2- >Play ();
IG++;
}
Else
{
If (API [position] [1] == 'S'
{
MP2- >FileName=API [position] [2]);
MP2- >Open ();
MP2- >Wait=true;
MP2- >Play ();
IG++;
}
}
}
}

```

V. TESTS AND RESULTS

The main interface of our system, with an example of the recognition of the digit ten (AACHARA) is shown in the following figure (Fig.4):

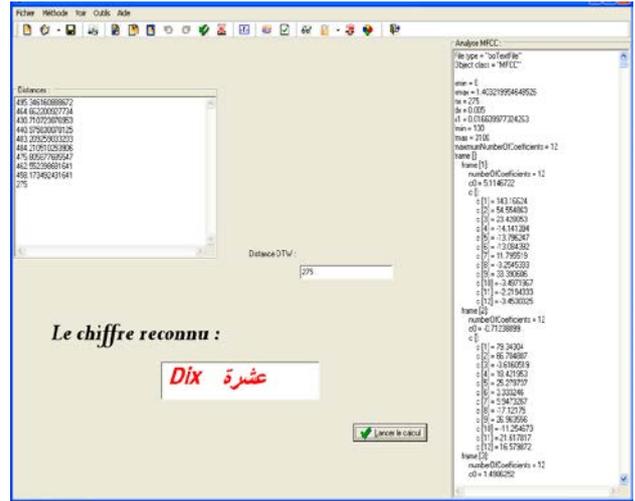


Fig. 4. Main interface of our system

We have applied the recognition on a corpus containing Arabic digits from one to ten pronounced by a male sex speaker in the Standard Arabic language.

To evaluate the performance of our system, we have illustrate two formula for each module; a recognition rate (RR) for the recognizer, and a success rate (SR) for the synthesizer.

We have fixed the number of tests performed to recognize a digit to twelve times. The recognition rate (RR) for each digit is calculated by the following formula:

$$RR = \frac{Nb_recognized_digit}{Nb_tested_digit(12)} * 100\%$$

In the other hand, to calculate the success rate (SR) associated with each digit tested; we got the following formula:

$$SR = \frac{Nb_well_pronounced_digit}{Nb_tested_digits} * 100\%$$

The results obtained for each digit are summarized in the following table:

Table 1 : Recognition/ Success rate for the ten digits

The word in Arabic	Transcript	The word in English	RR	SR
واحد	WAHID	ONE	85.7 %	100%
اثنان	ITHNAN	TWO	100%	100%
ثلاثة	THALATHA	THREE	100%	85%
أربعة	ARBAA	Oven	100%	100%
خمسة	KHAMSA	FIVE	100%	100%
سنة	SITTA	SIX	85.7 %	100%
سبعة	SABAA	SEVEN	100%	86%
ثمانية	THAMANIYA	EIGHT	100%	80%
تسعة	TISAA	NINE	100%	86%
عشرة	AACHARA	TEN	88.6 %	100%

When investigating across the natural language processing field, we haven't found a lot of works dealing with the automatic recognition and speech synthesis in a same work, especially for the Arabic language. Therefore, in the comparison with previous works, we take into account just the success accuracy of the automatic recognition.

The comparison results obtained are summarized in the following table:

Table 2 : Comparison with previous work

ASR using CMUSphinx [11]	85.55 %
DTW-Based ArSR [12]	86%
DHMM-Based ArSR [12]	92%
Heuristic Method [13]	86.45 %
Heuristic Method with RNN [13]	95.82 %
Monophone-Based ArSR [14]	90.75 %
Triphone-Based ArSR [14]	92.24 %
Syllable-Based ArSR [14]	93.43 %
Word-Based ArSR [14]	91.64 %
VQ and HMM Rrna [15]	91%
MCCF-based Rrna [15]	61 %AP -92%
Wavelet-based Rrna [15]	76 %AP -92%
LBC-based FPGA Rrna [16]	91 % -96%
MCCF-based FPGA Rrna [16]	95 % -98%
ARSSAD	96%

The recognition sub-system achieved 96% correct digit recognition in the case of mono-speaker mode. On the other hand, the speech synthesis sub-system achieved 93.7% correct well synthesized digit. So the system present in general 94, 85 % of success rate.

VI. CONCLUSIONS

We set several objectives for this research: that of discover the definitional character of the human voice, to describe the various stages and components used in the production of the voice and to dissect an ASRS system in its main floors. To that end, we have detailed our system of recognition and synthesis of Arabic digit as well as the results obtained. The system presents, using isolated words and in the absence of noise, a success rate quite honorable and acceptable. The acoustic variability of the voice signal, and in particular that due to the effects of coarticulation, is better apprehended by the modeling of its production.

In fact, the voice signal is not an ordinary acoustic signal and the Anatomical constraints may explain the

effects of coarticulation, for example, in the framework of the articulatory phonology.

At the end of this rapid assessment on the voice recognition and synthesis, it has been noted that this area is particularly broad and that there is no miracle product capable of responding to all applications. The noise, for example, remains a brake to the generalization of recognition systems. The voice recognition is still a compromise between the size of the vocabulary, its possibilities multi-speaker, its rapidity, training time, etc... The power of the current calculating tools and the integration capabilities of systems have caused a resurgence of interest in the recent years among the industrials. In fact, they see in the voice recognition or synthesis, "the more commercial", allowing making the difference with the competition.

A quick tour of horizon on the very numerous publications allows us to set the ideas on the nature of the work in progress. Apart from the products dedicated to the voice recognition, the systems with analytical approach (HMM and ANN) give today the best results [11], and currently have the wind in their sails.

As regards the future prospects, the optimism is more measured than in the past. Without risk, we can say that the general problem of the automatic processing of the voice signal will probably not rule before the middle of the next century. We can as even quote a few perspectives to our work in the following points:

- 1) Enlargement of the vocabulary for all digits;
- 2) Recognition of continuous speech;
- 3) Recognition in speaker independent mode;
- 4) Use of the HMM, neural networks and hybrid methods.

REFERENCES

- [1] Halima Bahi, a NeuroExpert system for the recognition of the voice; NESSR: Neural Expert System for Speech Recognition, LRI Laboratory, Department of Computer Science, University of Annaba, Algeria, bahi@lri-annaba.net
- [2] Ali Sadiqui, Nouredine Chenfour, realization of a system of automatic speech recognition Arabic based on CMU Sphinx, Anal. Computer Science Series. 8TH Volume 1st Parts, Faculty of Sciences Dhar El Mehraz, University Sidi Mhamed Ben Abdellah of Fez BP.1796- Fes- Morocco, 2010.
- [3] H. Satori, N. Chenfour, MR. Harti, Introduction To Arabic speech recognition using CMU Sphinx System, International Journal Of Computer Science, 2007.
- [4] Othman. O. Khalifa, M.Z. Obaid, A.W. Naji and Jamal I. Daoud, "A Rule-Based Arabic Text-To-Speech System Based On Hybrid Synthesis Technique", Electrical and Computer Engineering Department, International Islamic, University Malaysia Gombak, P.O Box 10, 50728 Kuala Lumpur, Malaysia, Australian Journal of Basic and Applied Sciences, 5(6): 342-354, 2011.
- [5] S. Dektelaere, "automatic speech recognition", MULTITEL - Department automatic speech recognition, Park Initialis-Avenue Copernic, Mons Belgium, site: www.multitel.be
- [6] L. V. Tray, "automatic recognition of digits in English in noisy conditions", University Joseph Fourier, U .F .R Informatics & Applied Mathematics, June 20, 2002
- [7] Paul Boersma and David Weenink, "PRAAT: doing phonetics by computer" Phonetic Sciences, University of Amsterdam Spuistraat 210, 1012 VT Amsterdam The Netherlands.
- [8] Pierre DRAGICEVIC "a model of interaction in input for interactive systems multi-devices highly configurable ", Ph.d.

thesis from the University of Nantes, the National College of Industrial Technology and Mines of Nantes, France, March 09, 2004

- [9] <http://www.crisco.unicaen.fr/description-des-differentes.html>, last access time : April 24th, 2014
- [10] P. Boula of Mareuil, "Synthesis of the floor from couriers and evaluation of conversion grapheme-phoneme". LIMSI-CNRS [http://www.limsi.fr/ Individu/ mareuil/](http://www.limsi.fr/Individu/mareuil/)
- [11] H. Satori, MR. Harti, N. Chenfour, "Introduction to Arabic Speech Recognition Using CMUSphinx System," Proceedings of Information and Communication Technologies International Symposium (ICTIS' 07), Fes, Morocco, pp. 139-115, July 2007.
- [12] Z. Hachkar, A. Farchi, B. Mounir, J. El Abbadi, "A Comparison of DHMM and DTW for Isolated Digit Recognition System of Arabic Language," International Journal on Computer Science and Engineering, vol. 3, no. 3, pp. 1002-1008, March 2011.
- [13] Khalid Saeed and Mohammad Nammous, Heuristic Method of Arabic Speech Recognition, Bialystok University of Technology, Poland, <http://aragorn.pb.bialystok.pl/~zspinfo/>
- [14] Mohamed Mostafa Azmi, Hesham Tolba, Sherif Mahdy, Mervat Fashal, " Syllable-Based Automatic Arabic Speech Recognition", Proceedings of WSEAS International conference of Signal Processing, Robotics and Automation (ISPRA' 08), University of Cambridge, UK, pp. 246-250, February 2008.
- [15] H. Bahi and Mr. Sellami, "Combination of Vector Quantization (and hidden Markov Models for Arabic Speech Recognition," Proceedings of the ACS/IEEE International Conference on Computer Systems and Applications (AICCSA 2001), Beirut, Lebanon, pp: 96-100, June 2001.
- [16] F. A. Elmisery, A.H. Khalil, A. E. Salama, H. F. M'hammed, "A FPGA Based HMM for a discreet Arabic Speech Recognition System," Proceedings of the 15th International Conference on Microelectronics (ICM 2003), Cairo, Egypt, December 9-11, 2003.

Tebbi Hanane, was born in Algiers, Algeria, 1981. She received her B.Sc (engineer) degree from University of Blida, Algeria in 2004, and her M.Sc. degrees from University of Saad Dahleb de Blida, Algiers, Algeria, in 2007. She is currently the Ph.D. student in department of Computer Science, University of Science and Technology of Houari Boumediene de Bab Ezzouar, Algiers, Algeria. Her research interests include Expert Systems, Natural Language Processing and Systems Engineering.

Hamadouche Maamar, Was born in Chlef, Algeria, 1981. He received his B.Sc. (Engineer) from University Hassiba Benbouali of Chlef, Algeria, in 2004, and M.Sc. degrees from University of Saad Dahleb de Blida, Algeria, in 2008. His research interests include Pattern Recognition, Natural Language Processing, Systems engineering and Data Base.

Dr. Azzoune Hamid, was born in Algiers, Algeria, 1959. He received his B.Sc (engineer) degree in Computer Science from University of Science and Technology of Houari Boumediene de Bab Ezzouar, Algiers, Algeria, in 1984, and his DEA from ENSIMA, Grenoble, France in 1985, and his Ph.D from INPGrenoble, France in 1989. Presently working as Researcher Professor in Department of Computer Science at University of Science and Technology of Houari Boumediene de Bab Ezzouar, Algiers, Algeria, since 1990. His research interests include: AI, DB, logic, CLP and web service.

Image Encryption Using Development of Chaotic Logistic Map Based on Feedback Stream Cipher

Hossam Eldin H. Ahmed

Dept. Of Electronics Comm. Eng., P. Dean of the Faculty of Electronic Eng. Menouf-32952, Menufiya University, Egypt.

Ayman H. Abd El-aziem

Ph. D Student, Dept. Of Electrical Engineering, Shubra Faculty of Engineering, Benha University, Egypt. e

Abstract—Recently due to the development of computer technology many multimedia content as digital image need to be transmitted over network, digital image are used in several application as medical image, confidential video conferences and military image data base of this digital images need to be protect from unauthorized. We need to encrypt these contents of images when transmitted over unsecured network.

This paper focuses on protecting digital images through using developed chaos-based encryption/decryption algorithms. We propose an image encryption based on development of chaotic logistic map and on feedback stream cipher. The proposed algorithm uses a developed chaotic logistic map and an external secret key of 256-bit. Further more, the proposed algorithm obtain solution by iteration, data dependent inputs, inclusion of three feedback mechanisms are verified to provide high security level. Our proposed algorithm has advantages,

1-Extend the range of the variable r by developed the chaotic logistic map. 2-New features for our proposed algorithm such as inputs(key, image). 3-The combine of feedback property with external secret key make cipherimage not depend on key only but depend on key and previous cipherimage pixel, this give the algorithm robustness against any cryptanalysis. 4-The experimental result of our proposed algorithm proof that it is an efficient method and secures way for real time image encryption. Furthermore a simple implementation of our algorithm achieves high encryption rates on general-purpose computers.

Keywords—Image encryption, stream cipher, development of logistic map, information security.

I. INTRODUCTION

In recent years, owing to frequent flow of digital images across the world over the transmission media, it has become essential to secure them from leakage that require reliable, fast, and robust security system to store and transmit digital images. The requirements to fulfil the security needs of digital images have led to the development of good encryption techniques. During the last decade, numerous encryption algorithms [1-2-10-11] have been proposed in the literature based on different principles. Among them, chaos based encryption techniques are considered good for practical use as these techniques provide a good combination of speed, high security, complexity, reasonable computational overheads and computational power.

In this paper we propose a new approach for image encryption based on feedback steam cipher by using a

developed chaotic logistic map it consist of two modules, first the encryption module which encrypt the image pixel-by-pixel, by using external secret key 256-bit are consider in each iteration, the values of the previously encrypted pixels, the second module is decryption module which decrypt the cipherimage pixel by pixel to retrieve the original image using the same key.

The feedback property, combined with the external secret key of 256-bit are verified to provide high security level, also, makes our proposed stream cipher robust against cryptanalytic attacks. The results of security analysis show that the proposed model provides an efficient and secure way for real-time image encryption and transmission. Our proposed chaotic logistic map with generated session key is compared with some different map as Bernoulli map, Genhous map, tent map and logistic map 1. Compare between them in several experimental, statistical analysis and key sensitivity tests. Also our proposed chaotic algorithm is compared with the RC5 , RC6 algorithms.

The rest of this paper contains different chaotic map and its analysis in Section 2 ,the proposed algorithm which is consist of two modules are mention in section 3, Section 4 test and verify of algorithm by applied algorithm using different map and verification for encryption and decryption. Section 5 the security analysis, the Section 6 comparing between proposed algorithm and RC5, RC6 and summery of paper is in Section 7.

II. DIFFERENT CHAOTIC MAP AND ITS ANALYSIS

We introduce some different chaotic maps and analysis the results simulation of these maps.

A. Bernoulli Map

Bernoulli map is chaos function and we use it in cryptography application. Its function is expressed as:

$$X_{n+1} = r \times X_n \text{ mod } 1 \quad (1)$$

Where

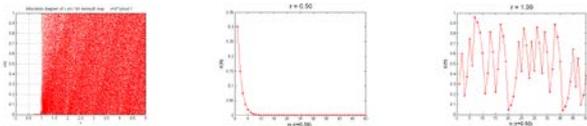
- X_n takes values from interval 0, 1, $r \in [0, 1]$.
- r is a variable and takes values from 0 to ∞ , $r \in [0, \infty]$.
- Initial value $X_n = 0.1$.
- Loop iteration = 8000 for r incremented by 0.001.

The simulation result is shown in figure 1 the parameter r can be divided into two segments which can be experiments on the following condition,

1. When $r \in [0, 1]$ as shown in figure 1.b the calculation result come to the same result after several iteration without chaotic behavior.
2. When $r \in [1, \infty]$ it become a chaotic system without periodicity as shown in figure 1.c.

From the previous discussion we can conclude that:

- When $r \in [0, 1]$ the point concentrate on several values could not use in image cryptosystem.
- When $r \in [1, 4.99]$ Bernoulli map have small change in the range of r to exhibit chaos behavior and hence the property of sensitive dependence so it can use for image cryptosystem in this small range.
- We obliged that not use integer value of r when use Bernoulli map in image cryptosystem, it must use a fraction value for r and this is one of its disadvantage.



a) Bifurcation for $r \in [0,4.99]$, $X_0 = 0.1$ b) Iteration property when $r = 0.50$ c) Iteration property when $r = 1.99$
 Fig 1: Analysis of Bernoulli Map

B. Genhous map

Genhous map is chaos function and we use it in our cryptography applications. Its function is chaos generator with a recursive structure expressed as,

$$X_n = f(r_1 \times X_{n-1} + r_2 \times X_{n-2} + r_3 \times X_{n-3}) \quad (2)$$

And

$$f(X_n) = X - 2 \text{ floor}((X+1)/2) \quad (3)$$

Where

- r_1 arbitrary, $r_2 = r_3 = 1$.
- X_n take values from interval 0, 1, $X_n \in [0,1]$.
- r_1 is variable and takes value from 0 to ∞ $r_1 \in [0, \infty]$.
- Initial value $X_1 = 0.1$.
- Loop iteration=6000 for r incremented by 0.001.
- Floor is function takes the integer part of number as floor(5.5) = 5.

The simulation are shown in figure 2, for different values of parameter r it become a chaotic system without periodicity except in integer value of r as shown in figure 2.b,c. Also we use MATLAB software to graph the bifurcation diagram of Genhous map as show in Figure 2.a, for $r \in [0,4.99]$ from the simulation result the chaotic behaviour occurs as shown in From the previous discussion we can conclude that:

- when $r \in [0,499]$ Genhous map exhibit chaos behavior for integer r and hence the property of sensitive dependence .It can be used for image cryptosystem with disadvantage of r must be fraction and have small range , so we recommend that not use integer value of r when use Genhous map in image cryptosystem.

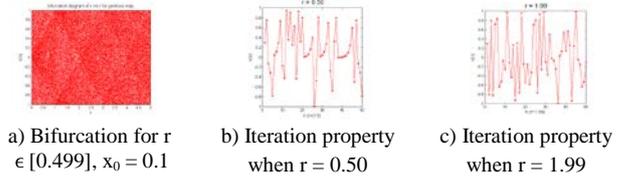


Fig 2 Analysis of Genhous Map

C. Tent Map

Tent map is a chaos function can be used in cryptography application. It expressed as:

$$X = r \times (1 - |1 - 2 \times X|) \quad (4)$$

Applying under the following conditions:

- X_n take value from interval 0, 1 $X_n \in [0,1]$
- r is variable, $r \in [0,4]$.
- Initial value $X_n = 0.3$.
- Loop iteration = 6500 for r incremented by 0.001.

The simulation are shown in figure 3, the parameter r can be divided into three segment which can be experiments on the following condition.

1. When $r \in [0, 0.5]$ as shown in figure 3.b the calculation result come to the same result after several iteration without any chaotic behavior.
2. When $r \in [0.5, 0.7]$ as shown in figure 3.c the phase space concludes several point the Systems appear as periodic behavior.
3. When $r \in [0.7, 1]$ it become a chaotic system without periodicity as shown in figure 3.d.

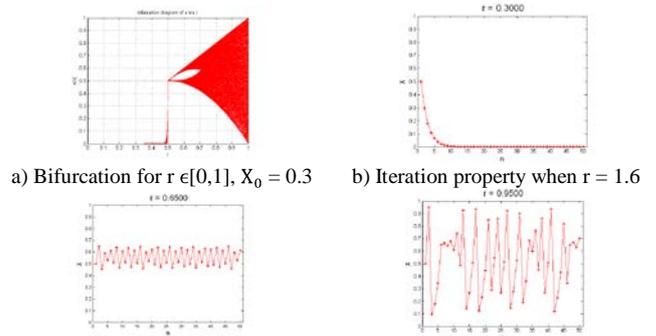


Fig 3. Analysis of tent map

Also we use MATLAB software to graph the bifurcation diagram of tent map as show in figure 3.a

From the previous discussion we can conclude that,

- When $r \in [0, 0.7]$ values could not use in image cryptosystem.

- When $r \in [0.7, 1]$ Tent map exhibit chaos behavior and hence the property of sensitive dependence so it can be used for image cryptosystem for $r \in [0.7, 1]$, but this range of applied r is very small and not enough for a secure cryptosystem

D. Logistic Chaotic Map 1

Logistic map is developed to give a chaos function can be used it in cryptography application. This logistic map 1 is expressed as:

$$X_{n+1} = r^2 \times X_n \times (1-X_n) \times (1-2 \times X_n) \quad (5)$$

Applying under the following conditions:

- X_n take value from interval 0, 1, $X_n \in [0, 1]$.
- r is variable, $r \in [0, 4]$
- Initial value $X_0 = 0.3$.
- Loop iteration = 9000 for r incremented by 0.001.

The simulation are shown in figure 4, the parameter r can be divided into three segment which can be experiments on the following condition.

- When $r \in [0.2, 1]$ as shown in figure 4.b the calculation result come to the same result after several iteration without any chaotic behavior.
- When $r \in [2.1, 2.47]$ as shown in figure 4.c the phase space concludes several points the Systems appear as periodic behavior.
- When $r \in [2.47, 4]$ it become a chaotic system without periodicity as shown in figure 4.d.

Also we use MATLAB software to graph the bifurcation diagram of logistic chaotic map1 as show in figure 3.a From the previous discussion we can conclude that:

- When $r \in [2.47, 4]$ values could not use in image cryptosystem.
- When $r \in [2.47, 4]$ logistic map1 exhibit chaos behavior and hence the property of sensitive dependence so it can be used for image cryptosystem for $r \in [2.47, 4]$ but this range of applied r is small and not enough for a secure cryptosystem .

also but we can be used it in cryptography applications under certain conditions. This logistic map function is express as:

$$X_n = r \times X_{n-1} \times (1-X_{n-1}) \quad (6)$$

And

$$X_{n+1} = 4 \times X_n \times (1-X_n) \quad (7)$$

We substitute from equation (6) into equation (7) and Applying under the following conditions,

- X_n take value from interval 0, 1, $X_n \in [0, 1]$
- r is variable, $r \in [0, 4]$.
- Initial value $X_1 = 0.3$.
- Loop iteration = 6000 for r incremented by 0.001.

We use Matlab software to simulate our logistic map applied to images. We start our program by the initial value , $X_1 = 0.3$ the simulation are shown in figure 4, Where the parameter r can be divided into three segment which can be resume to these three following condition:

1. When $r \in [0, 1.1]$ as shown in figure 5.b, the calculations results come to the same results after several iteration without any chaotic behavior.
2. When $r \in [0, 1.5]$ as shown in figure 5.c, the phase space concludes several points the system appear as periodic behavior.
3. When $r \in [1.5, 4]$ it become a chaotic system without periodicity as shown in figure 5.d.

Also we use MATLAB software to graph the bifurcation diagram of logistic map 2 as show in figure 4.a and we can conclude that:

1. The case of $r \in [0, 1.]$ from the simulation result the trajectory of equation convergence to fixed point as illustrate in figure 5.a.
2. The case of $r \in [1.1, 1.5]$ from the simulation result the phenomena of period double bifurcation as shown in figure 5.a.
3. The case of $r \in [1.5, 4]$ from the simulation result a chaotic behavior as shown in figure 5.a.

We refer to these three regions as convergences, bifurcations, and have chaos behavior respectively from the previous discussion one can conclude that,

1. For $r \in [0, 1.5]$ the point concentrate on several values could not use in image cryptosystem.
2. For $r \in [1.5, 4]$ the logistic chaotic map2 exhibit chaos behaviour and hence the property of sensitive dependence so it can be use for image cryptosystem.

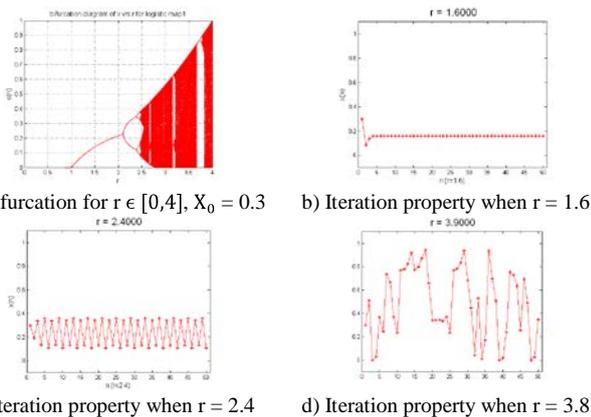
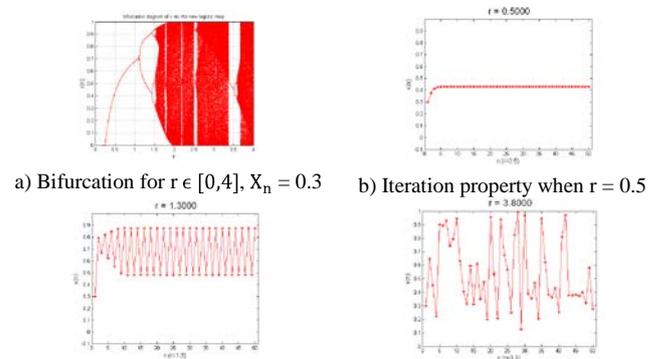


Fig 4: Analysis of logistic chaotic map1



E. Our Proposed Logistic Chaotic Map 2

A proposed logistic map 2 which is a chaos function

c) Iteration property when $r = 1.3$ d) Iteration property when $r = 3.8$

Fig 5. Analysis of our proposed logistic chaotic map 2

F. Another Proposed Logistic Chaotic Map 3

We introduce another developed logistic chaotic map 3 where the range of r increased chaotic area for all applications of as will as shown below.

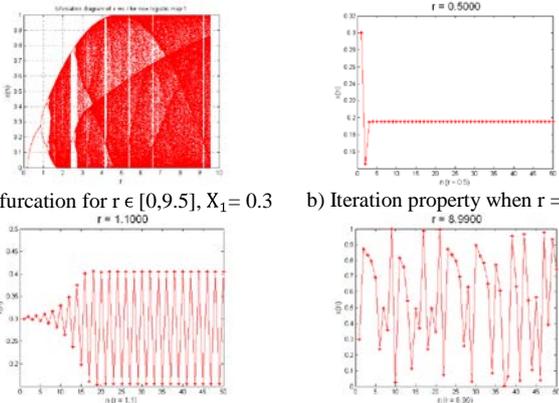
$$X_{n+1} = r \times X_n \times (1 - X_n) \times (1 - 1.2 \times X_n)^2 \quad (8)$$

And substitute from equation (8) into equation (7) Applying these under the following conditions,

- X_n take value from interval $0, 1$ $r \in [0, 1]$.
- The variable r takes value from 0 , to 9.5 $r \in [0, 9.5]$.
- Initial value $X_1 = 0.3$.
- Loop iteration = 8000 for r incremented by 0.001 .

By using Matlab software to simulate bifurcation and iteration results where the initial value $X_1 = 0.3$. These simulations are given in figure 6 where the parameter r divided into three segments given as follow:

- When $r \in [0, 1.1]$ as shown in figure 6.b, the calculation result come to the same result after several iteration without any chaotic behavior.
- When $r \in [1.1, 1.5]$ as shown in figure 6.c, the phase space concludes several point where the system appear as periodic behavior.
- When $r \in [1.1, 9.5]$ it become a chaotic system without periodicity as shown in figure 6.d.



a) Bifurcation for $r \in [0, 9.5]$, $X_1 = 0.3$ b) Iteration property when $r = 0.5$
 c) Iteration property when $r = 1.1$ d) Iteration property when $r = 8.99$
 Fig 6. Analysis of our proposed logistic chaotic map 3

- The case of $r \in [0, 1.1]$ from the simulation result the trajectory of equation convergent to fixed point as illustrate in figure 6.a.
- The case of $r \in [0, 1.3]$ from the simulation result, the phenomena of this period is double bifurcation, as shown in figure 6.a.
- The case of $r \in [1.3, 9.5]$ from the simulation result the chaotic behavior occurs as shown in figure 6.a.

We can refer to these three regions as convergences, bifurcations and exhibit chaos behavior respectively. From the above discussion we can conclude that.

- When $r \in [0, 1.3]$ the point concentrate on several values could not use in image cryptosystem.
- For $r \in [1.3, 9.5]$ which is a long range for r , the proposed new logistic map1 exhibit chaos behaviour and hence the property of sensitive dependence, so as an advantage to use it for a wide range r image cryptosystem analyses and applications.

From the above discussion we can conclude that, these three proposed logistic chaotic maps can be used for different forms of logistic maps and for wide range of r that give a good chaotic behaviour and can be applied in our proposed algorithm.

III THE PROPOSED ALGORITHM

We propose a new approach for image encryption based on developed of chaotic logistic maps in order to meet the requirements of the secure image transfer, we propose algorithm based on feedback steam cipher by using a developed chaotic logistic map it consist of two modules, first the encryption module, the second module is decryption module.

A. The Encryption Module

The proposed is a simple block cipher with block size of 8-bit and 256-bit secret key. The key is used to generate a pad that is then merged with the plaintext a byte at a time, as shown in figure 6.

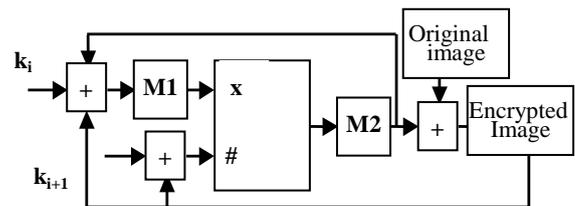


Fig 6: Diagram of Encryption Module

We can obtain the cipherimage from the following steps:

1. For the encryption/decryption, we divide plainimage/cipherimage into blocks of 8-bits (pixel). Plainimage and cipherimage of i blocks can be presented as:

$$P = P_1 P_2 P_3 P_4 P_5 \dots P_i \quad (9)$$

$$C = C_1 C_2 C_3 C_4 C_5 \dots P_i \quad (10)$$

2. The proposed image encryption process utilizes an external secret key of 256-bit long. Further, the secret key is divided into blocks of 8-bit each block referred as session keys the secret key can be represented in ASCII mode as,

$$K = K_1 K_2 K_3 K_4 K_5 \dots K_i \quad (10)$$

Where, each K_i represents one 8-bit block of the secret key i.e. session key.

- The initial condition (X_0) for the chaotic map and the initial code C_0 are generated from the session keys as:

$$R = \sum_{i=1}^{32} M1[K_i] \quad (11)$$

$$X_0 = R - [R] \quad (12)$$

$$C_0 = [\sum_{i=1}^{32} [K_i]] \text{ mod } 256 \quad (13)$$

Where K_i , $\lfloor \cdot \rfloor$, and $M1$ are, respectively, the decimal equivalent of the i th session key, the floor function (result the integer part of number), and mapping from the session, key space, all integers between 0 and 255, into the domain of the logistic map, all real numbers in the interval $[0,1]$.

- Read a byte from the image files (that represent a block of 8-bits) and load it as plainimage pixel P_i .
- Encryption of each plainimage pixel P_i to produce its corresponding cipherimage pixel C_i can be expressed mathematically as:

$$C_i = \left(P_i + M2 \left[\sum_{i=1}^{\#_i} \text{chaoticmap} \right] \right) \text{ mod } 256 \quad (14)$$

Where chaotic map is one of map which we mention in Section 2 the output of each cipherimage pixel is feedback to input for chaotic map to calculate the input of chaotic map and the iteration of chaotic map. Where represents the current input for logistic map and computed as:

$$X_i = M1[k_i + C_{i-1} + X_{i-1}] \quad (14)$$

$$\#_i = k_{i+1} + C_{i-1} \quad (15)$$

Where $\#_i$ is the number of iteration of chaotic map for its current input X_i and calculated as:

And $M2$ maps the domain of the chaotic map $[0,1]$ back into the interval $[0,255]$.

- Repeat steps 4-5 until the entire image file is exhausted.

We have three feedback in our module first the output of cipherimage pixel to the input of chaotic map input second to the number of iteration third the output of chaotic map is input to input of chaotic map this three feedback make the cipherimage pixels not depend on key only but on the previous cipherimage pixel and the number of iteration depend on next session key and previous cipherimage pixel, we applied this module by using the different chaotic map which is mention in section 2 and using the optimum value which in listed in table 1 to generate the pad which is merged with plain image pixel to produce cipherimage.

We introduce the development of chaotic logistic map to extend the range of variable r which gives chaotic properties for logistic map it is from 3.57 to 4 [7, 8], after we develop it becomes

- for logistic map 2 it become from 1.5 to 4
- And for logistic map 3 it become from 1.3 to 9.55.

We use one of chaotic map which mention in section 2 to generate session key which is a pad which is merged to plainimage to generate cipherimage we add the decimal value of K_i to value of the output of pervious chaotic map X_{i-1} to the value of pervious Ciphertext C_{i-1} and mapping this value to chaotic map domain to produce X_i which is input to new logistic map and use number of iteration equal the value of next session key to the previous cipherimage pixel And map the output of this chaotic map to interval $[0,255]$.to generate a pad which is merged with plaintext byte to produce Ciphertext byte we repeat that until we finish the image file.

TABLE I THE RANGE OF R FOR DIFFERENT CHAOTIC MAP

Chaotic map	Range for chaotic	Optimum Value of r
Bernoulli map	1 : ∞	1.99
Genhous map	0 : ∞	1.99
Logistic map 1	2.5 : 4	3.6
Logistic map 2	1.5 : 4	3.8
Logistic map 3	1.3 : 9.5	8.99

B. The Decryption Module

Decryption is very simple it is similar to encryption module except in this case the same pad is generated but this time un-merged with the cipherimage to retrieve the plainimage, the decryption module receives an encrypted image (cipherimage) and the 256-bit secret key and returns the original image (plainimage).

$$P_i = \left(C_i - M2 \left[\sum_{i=1}^{\#_i} \text{chaotic map} \right] \right) \text{ mod } 256 \quad (16)$$

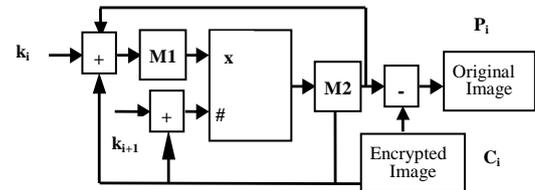


Fig 7 Diagram of decryption Module

The proposed algorithm by using development logistic map it appear to sensitive to any change because each cipherimage depend on session key and previous cipherimage pixel so that any change in plainimage makes great change in cipherimage And it is data dependent iteration because the input of chaotic map map are computed as function of session key and previous computed cipher pixel and previous new logistic output.

IV TEST, VERIFICATION AND EFFICIENCY OF OUR PROPOSED ALGORITHM

We apply our proposed algorithm by using different chaotic map to image Lena, which size 256 x 256 Gray-scale

(0-255) as original image (plainimage) and use secret key (k1="1234578901234567890123456789012") (in ASCII) is used for encryption whose long is 256-bit. As shown in Fig 9 Application of proposed algorithm using different chaotic map and RC6 algorithm to original image Lena , it is clear that the results encrypted images (cipherimages) regions are totally invisible by applied our algorithm using different chaotic map and by applied RC6 algorithm, the decryption method takes (cipherimage) as input together with the same secret key (k1="1234578901234567890123456789012") (in ASCII) and return the plainimage, One of the important examining of encrypted image is the visual inspection, where the more hidden features of the image are, the better the encryption algorithm. But it is not suffusion to determine the quality of our proposed algorithm by visual inspection so that we evaluate different testes to determine the efficiency of our proposed algorithm.

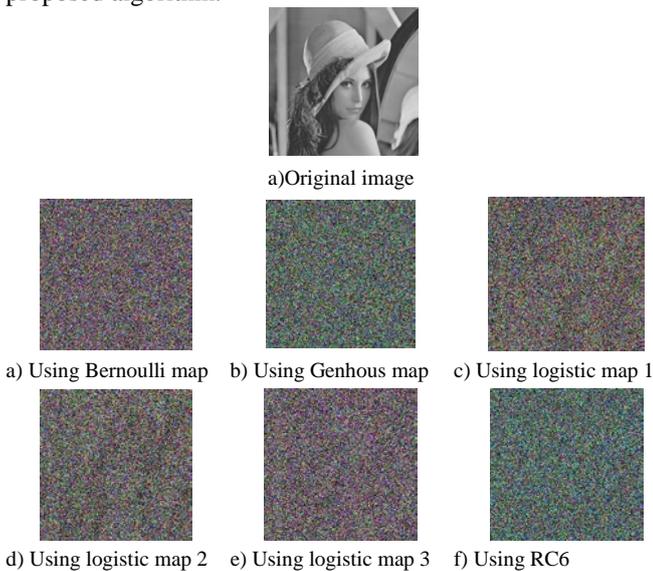


Fig 9. Application of proposed algorithm using different chaotic map and RC6 algorithm to image Lena

V SECURITY ANALYSES

A good encryption procedure should be robust against all kinds of cryptanalytic, statistical and brute-force attacks. In this section, we discuss the security analysis of the proposed image encryption scheme such as statistical analysis, sensitivity analysis with respect to the key and plaintext and key space analysis. To prove that this proposed cryptosystem is secure against the most common attacks [3, 4].

A. Statistical Analysis

It is well known that many ciphers have been successfully analyzed with the help of statistical analysis and several statistical attacks have been devised on them. Therefore, an ideal cipher should be robust against any statistical attack. To prove the robustness of the proposed image encryption procedure by using different chaotic map, we have performed statistical analysis by calculating the histograms, the correlations between two adjacent pixels in the original images and its corresponding encrypted images.

1) *Histogram Analysis:* An image-histogram illustrates how pixels in an image are distributed by graphing the number of pixels at each color intensity level. We have calculated and analyzed the histograms of the several encrypted image that the histograms of the encrypted image are fairly uniform and significantly different from the histogram of original image we analyze and calculate the histogram of original images and its cipherimages, by using different chaotic map. We find that as shown in fig 10 the histogram of plainimage and cipherimage by applied our algorithm using different chaotic map are fairly uniform and significantly different from the histogram of original and using it bears no statistical resemblance to the plainimage, so that it is appear no statistical attack against our proposed by using different chaotic map. Also histogram of applied RC6 algorithm give good result its histogram is fairly uniform and different from the histogram of original image.

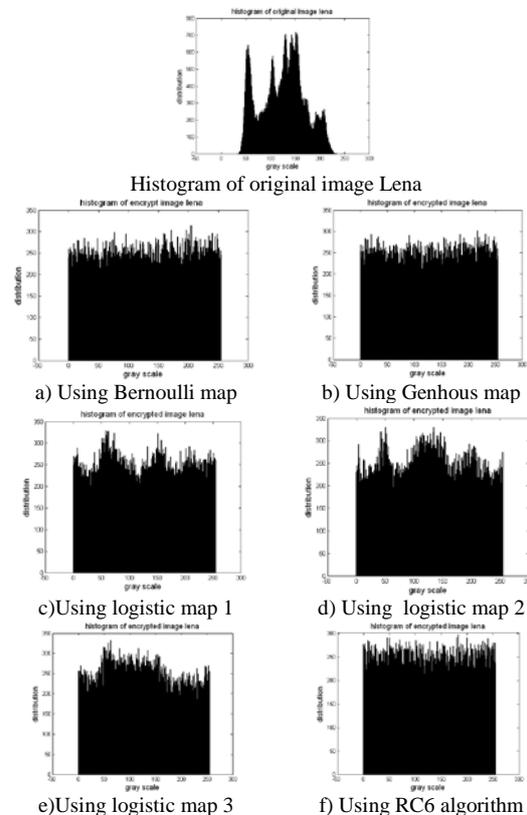


Fig 10. Histogram of cipherimage by using different chaotic map in our proposed algorithm and RC6 algorithm

2) *Correlation Coefficient Analysis:* We analyzed correlation between two vertically adjacent pixels, two horizontally adjacent pixels and two diagonally adjacent pixels in plainimage/cipherimage respectively, by select 1000 pairs randomly of two adjacent pixels from an image. Then, calculate their correlation coefficient using the following two formulas

$$\text{cov}(x, y) = E(x - E(x))(y - E(y)) \quad (17)$$

$$r_{xy} = \frac{\text{cov}(x, y)}{\sqrt{D(x)}\sqrt{D(y)}} \quad (18)$$

Where x and y are the values of two adjacent pixels in the image. In numerical computations, the following discrete formulas were used, Fig 11 shows the correlation distribution of two horizontally adjacent pixels in plainimage/cipherimage for our proposed using logistic map 3 chaotic map and. The correlation coefficients are and respectively for both plainimage/cipherimage.

It is clear that there is no correlation between two adjacent pixels in cipherimage and there are high correlated between adjacent in plainimage.

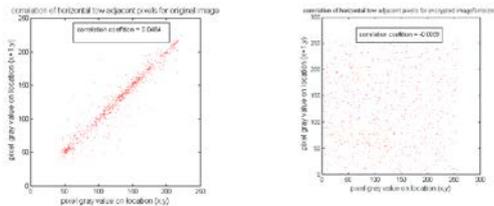


Fig 11. Two horizontally adjacent pixels Correlation in plainimage/cipherimage Using chaotic logistic map 3

In the next table we list the correlation coefficients of two horizontally, vertical and diagonal adjacent pixels in plainimage/cipherimage for our proposed in different chaotic map and RC6 algorithm in plainimage and cipherimage.

TABLE II CORRELATION COEFFICIENTS IN PLAINIMAGE/CIPHERIMAGE FOR OUR PROPOSED ALGORITHM USING DIFFERENT CHAOTIC MAP AND RC6 ALGORITHM

Chaotic map	Direction of Adjacent pixels	Plainimage	Cipherimage
Bernoulli map	Horizontal	0.9575	0.0177
	Vertical	0.9353	0.0
	Diagonal	0.9155	0.0
Genhous	Horizontal	0.9686	0.0586
	Vertical	0.9363	0.0330
	Diagonal	0.9191	0.0
Logistic map 1	Horizontal	0.0169	0.0160
	Vertical	0.9216	0.0036
	Diagonal	0.9129	0.0148
logistic map2	Horizontal	0.0721	0.0
	Vertical	0.9294	0.0
	Diagonal	0.8894	0.0307
logistic map 3	Horizontal	0.0484	0.0
	Vertical	0.9307	0.0
	Diagonal	0.9166	0.0153
RC6	Horizontal	0.0091	0.0571
	Vertical	0.9296	0.0161
	Diagonal	0.9071	0.0405

in previous table we calculate C.C between to adjacent pixel in horizontal , vertical and diagonal. For our proposed algorithm using different chaotic map and RC6 we find that:

- Our proposed algorithm using logistic map 3 have the smallest correlation coefficient compared to other algorithm , also our proposed algorithm using

Bernoulli map has small C.C , the rest algorithm has small C.C but his values greater than the proposed algorithm using logistic map 3 and Bernoulli map.

- There is negligible correlation between the two adjacent pixels in the encrypted image in all directions. However, the two adjacent pixels in the original image are highly correlated that for our proposed algorithm and RC6.

B Sensitivity Analysis

An ideal image encryption procedure should be sensitive to any small change in plainimage and secret key.

1)Key Sensitivity Analysis: An ideal image encryption procedure should be sensitive with respect to the secret key i.e. the change of a single bit in the secret key should produce a completely different encrypted image. For testing the key sensitivity of the proposed image encryption by using chaotic logistic map 3, we have performed the following steps:

- As shown in fig 12.a is original image, b) is in encrypted image by using the secret Key(k1='12345678901234567890123456789012' (in ASCII).
- the same original image is encrypted by making the slight modification in the secret key to become (k2=12345678901234567890123456789013) the least significant bit is changed in the secret key) as shown in fig 13.c.
- Again, the same original image is encrypted by making the slight modification in the secret key to become (k3=11345678901234567890123456789012) the most significant bit is changed in the secret key) and the resultant image referred as encrypted image in Fig. 13.d.
- Finally, the three encrypted images A, B and C are compared.

We have shown the original image as well as the three encrypted images produced it is not easy to compare the encrypted images by simply observing these images. So that to compare between three images, we have calculated the correlation between the corresponding pixels of the three encrypted images. For this calculation, we have used the same Formula as given in Equation (17). Except that in this case x and y are the values of corresponding pixels in the two encrypted images to be compared. In table III, we have given the results of the correlation coefficients between the corresponding pixels of the three encrypted images A, B and C. It is clear from the table 3 that no correlation exists among three encrypted images even though these have been produced by using slightly different secret keys.

TABLE III CORRELATION COEFFICIENTS BETWEEN THE CORRESPONDING PIXELS OF THE THREE DIFFERENT ENCRYPTED IMAGES OBTAINED BY USING SLIGHTLY DIFFERENT SECRET KEY

Image 1	Image 2	Correlation coefficient
Encrypted image A	Encrypted image B	0.0
Encrypted image B	Encrypted image C	0.0186
Encrypted image C	Encrypted image A	0.0211

We have also measured the number of pixels change rate (NPCR) to see the influence of changing a single pixel in the original image on the encrypted image by the proposed algorithm. The NPCR measure the percentage of different pixel numbers between the two images. We take two encrypted images, C1 and C2, whose corresponding original images have only one-pixel difference. We define a two-dimensional array D, having the same size as the image C1/C2. The D(i,j) is determined from C1(i,j) and C2(i,j). If C1(i,j) = C2(i,j) then D(i,j) =1 otherwise D(i,j) = 0. The NPCR is defined by the following equation.

$$NPCR = \frac{\sum_{i,j} D(i, j)}{W \times H} \times 100\% \quad (18)$$

Where w and h are the width and height of encrypted image. We obtained NPCR for a large number of images by using our encryption scheme and found it to be over 99% as listed in table .4 so that the encryption scheme is very sensitive with respect to small changes in key secret.

TABLE IV. NPCR FOR THREE DIFFERENT ENCRYPTED IMAGES IN FIG. 13

Image 1	Image 2	NPCR
Encrypted image A	Encrypted image B	% 99.4812
Encrypted image B	Encrypted image C	% 99.5041
Encrypted image C	Encrypted image A	% 99.4522

Moreover, in Fig. 13, we have shown the results of some attempts to decrypt an encrypted image with slightly different secret keys than the one used for the encryption of the original image, in fig 13.a is the original image and b) is the encrypted image produced using the secret key ‘12345678901234567890123456789012’ in (ASCII), (c) the decrypted image with the secret keys ‘12345678901234567890123456789012’ (in ASCII) and decrypted image with the secret keys ‘12345678901234567890123456789011’ respectively, the images after the decryption of It is clear that the decryption with a slightly different key fails completely and hence the proposed image encryption procedure is highly key sensitive. High key sensitivity is required by secure image cryptosystems, which means that the cipherimage cannot be decrypted correctly although there is only a slight difference between encryption and decryption keys. It is clear that the decryption with a slightly different key fails completely and

hence the proposed image encryption procedure is highly key sensitive.

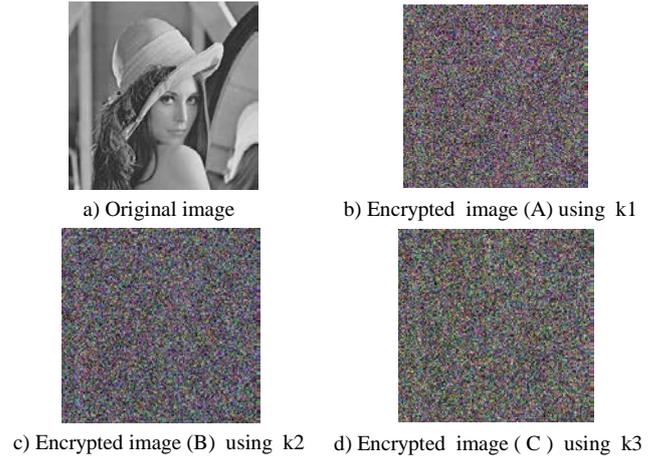


Fig 12: key sensitive result with proposed using chaotic logistic map 3

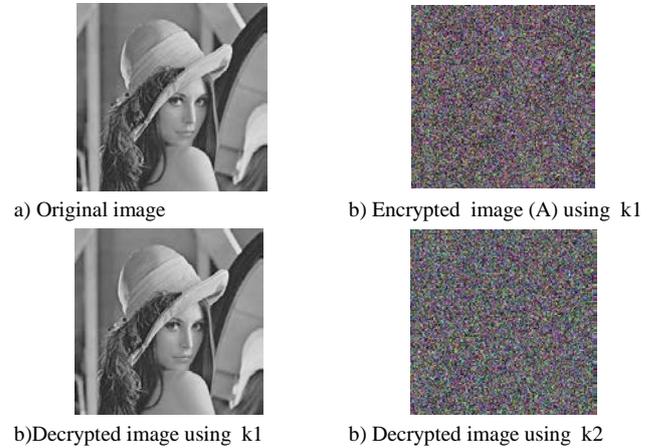


Fig 13:key sensitive result with by using chaotic logistic map 3

2) *Plainimage Sensitivity Analysis:* We have also measured the number of pixels change rate (NPCR) to see the influence of changing a single pixel in the original image on the encrypted image. The NPCR measure the percentage of different pixel numbers between the two images. We take two encrypted images, C1 and C2, whose corresponding original images have only one-pixel difference. We define a two-dimensional array D, having the same size as the image C1/C2. The D (i,j) is determined from C1(i,j) and C2(i,j). If C1(i,j)=C2(i,j) then D(i,j)=1 otherwise D(i,j)=0.

$$UACI = \frac{1}{W \times H} \left[\sum_{i,j} \frac{C_1(i, j) - C_2(i, j)}{255} \right] \times 100\% \quad (18)$$

The UACI is defined as the measures of the average intensity of differences between the two images. One performed test is on the one-pixel change influence on a 256 grey-level Lena image of size 256 x 256.

TABLE V. NPCR AND UACI FOR PROPOSED ALGORITHM USING DIFFERENT CHAOTIC MAP AND RC6 ALGORITHM.

Chaotic map / RC6	NPCR		UACI	
	Lena	Eivel	Lena	Eivel
Bernoulli map	% 99.21	% 99.59	84.90	85.65
Genhous map	% 99.23	% 99.62	85.42	86.14
Logistic map 1	% 98.87	% 99.27	83.90	84.32
logistic map 2	% 99.33	% 99.39	82.00	83.28
logistic map 3	% 99.08	% 99.05	82.05	84.16
RC6	% 99.24	%99.60	84.82	85.06

From table 5 we find that the proposed algorithm is sensitive to small change in plainimage show that the proposed algorithm is very sensitive with respect plainimage (plainimage have only one pixel difference).

C Key space analysis

Key space should be long enough to make brute force infeasible we use key length 256 -bit so that we have 2²⁵⁶ different key. Additionally the number of iterations supported by the logistic map module is between 0 and 767, as cipher pixels take values in the interval [0,512] and the session keys take values in the interval [0,255].

C. speed analysis

Apart from the security consideration is measure running speed for real-time image encryption/decryption. By measure the time required to encrypted/decrypted image we applied our proposed algorithm using development chaotic logistic map 3 to image in dimension 256 ×256 by using the simulator compiler Borland C++ Development Suite 5.0. Performance was measured on a 2.16 GHz Core 2 Duo with 1 GB of RAM running Windows XP. to improve the accuracy of our timing measurements, was executed 10 times, and we take the average, the time for encryption = 0.335 sec, the time for decryption = 0.4106 sec so that our proposed algorithm success in real time application.

VI COMPARISION BETWEEN OUR PROPOSED ALGORITHM AND RC5, RC6 ALGORITHMS

We compare between our proposed algorithm and RC5 and RC6 in optimum parameter[7,9] by measuring two encryption evaluation metrics to determine the quality of encryption, first we calculate the histogram deviation (DIV-1) Measures the deviation between the original and the

encrypted image[5], the higher value of D_H is the better quality of the encrypted image second we measures the irregular deviation it Measures how much the deviation caused by encryption on the encrypted image (is irregular)[5].The lower value of D_I (DIV-2) is the better encryption quality.

BY Analyzing the results of the images in the table 6 which measure the (DIV-1) and (DIV-2) we can conclude the following:

- ❖ For lena.bmp image, the greater value of
 - (DIV-1) are at our proposed algorithm which mean that the proposed algorithm has high maximum deviation rather than the other algorithms and the smallest value of (DIV -2) at RC6 in CFB mode which mean the best irregular deviation at this algorithm.
 - For evel.bmp image the greater value of (DIV-1) is at RC5 in ECB mode which has high maximum deviation rather than the other algorithms, (DIV -2) have small value at RC6 in CBC mode which mean the best irregular deviation at this algorithm.
 - For camera man.bmp image the greater value of (DIV-1) are at proposed algorithm which mean that the proposed algorithm has high maximum deviation rather than the other algorithms and (DIV -2) small value at RC6 in CFB mode and RC5 in OFB mode which has equal value which mean the best irregular deviation at this two algorithms.

By analyzing the results of the images given by the table V which measure the encryption quality we have the following:

- The encryption quality may be expressed in terms of the total changes in pixels values between the original image and the encrypted one [6].
- For lena.bmp image encryption quality has a good result with our proposed algorithm comparing with other ciphers and has a little good result in RC5 in OFB mode.
- For evel.bmp image encryption quality has a good result in RC5 for ECB mode than the other ciphers.
- For cameraman.bmp image encryption quality has a good result for our proposed algorithm than the other ciphers as RC6 CBC mode.

TABLE VI. THE MAXIMUM DEVIATION (DV-1) AND THE IRREGULAR DEVIATION (DV-2) FOR RC5, RC6 AND OUR PROPOSED ALGORITHM

		RC5				RC6				Proposed
		CBC	CFB	ECB	OFB	CBC	CFB	ECB	OFB	
Lena	DEV-1	46245	46561	46457	46006	46888	46746	46355	46138	49005
	DEV-2	46508	46356	46282	46230	46204	46114	46540	46264	50730
Eivel	DEV-1	71879	71611	74262	71637	71772	71529	71868	71696	71313
	DEV-2	36056	35834	39042	36146	35794	36080	37708	36124	41144
C_man	DEV-1	64322	64195	64144	64456	63931	63930	64414	64212	72401
	DEV-2	26788	26776	26778	26584	26868	26584	27048	26726	38122

TABLE VII ENCRYPTION QUALITY MEASURES FOR RC6, RC5, AND PROPOSED ALGORITHM USING CHAOTIC LOGISTIC MAP 3

	RC5				RC6				Proposed
	CBC	CFB	ECB	OFB	CBC	CFB	ECB	OFB	
Lena	182.757	182.757	183.546	181.750	185.164	184.468	183.109	182.281	193.578
Eivel	283.554	282.406	292.726	282.677	283.101	282.054	283.359	282.796	281.460
C_man	252.421	251.750	251.542	252.843	250.664	250.726	252.710	251.929	283.648

From previous two comparisons table we can conclude that: Our proposed algorithm is better than RC5 and RC6 when we use it at these cases

- Encrypted Lena .bmp and cameraman.bmp that from measuring encryption quality.
- Has very low C.C in cameraman encrypted.
- Has higher Maximum Deviation when use it in encrypt Lena and cameraman.

RC5 and RC6 are better than our proposed algorithm in these cases

Encrypted Eivel.bmp that from measuring encryption quality for RC5 in ECB mode

- Has very low C.C in Lena and cameraman encrypted.
- Has lower Irregular Deviation when use it in encrypt Lena, eivel and cameraman.

VI CONCOLUTION

In this paper, a new way of image encryption scheme using development of chaotic logistic map based on feedback stream cipher using an external secret key of 256-bit.

We use a developed logistic map to increase the range of the variable r which change from 3.57 to 4 in logistic map [8] to a wide chaotic range from 1.3 to 9.55 for logistic map 3.

Several test images are used for inspecting the validity of the proposed algorithm. The robustness of the proposed algorithm based on a feedback mechanism, which leads the cipher to a cyclic behaviour so that the encryption of each plain pixel depend on the output of the used chaotic map and the previous cipher pixel.

We have carried out key space analysis, statistical analysis, and key sensitivity analysis to demonstrate the security of the new image encryption procedure. According to the results of our security analysis, we conclude that the proposed algorithm is expected to be useful for real-time image encryption and transmission application.

REFERANCES

- [1] G. Chen, Y. Mao, C.K. Chui, A symmetric image encryption based on 3D chaotic maps, *Chaos Solitons Fractals* 21 (2004) 749–761.
- [2] N. Bourbakis, C. Alexopoulos, Picture data encryption using SCAN pattern, *Pattern Recogn.* 25 (1992) 567–581.
- [3] S. Lian, J. Sun and Z. Wang. "Security analysis of a chaos-based image encryption algorithm," *Physica A: Statistical and Theoretical Physics*, vol. 351, Issues 2-4, 15 June 2005, pp. 645-661.
- [4] T. Paraskevi, N. Klimis, K. Stefanos. "Security of Human Video Objects by Incorporating a Chaos-Based Feedback Cryptographic Scheme," *ACM Multimedia '04*, October, 10-16, 2004, New York, NY USA.
- [5] O. S. Faragallah, "Utilization of Security Techniques for Multimedia Applications", Ph. D. Thesis, Department of Computer Science and Engineering, Faculty of Electronic Engineering, Menofia University, 2007.
- [6] I. F. Elashry, "Image Encryption", Ms. D. Thesis, Department of Computer Science and Engineering, Faculty of Electronic Engineering, Menofia University, 2010.
- [7] Hossam El-din H. Ahmed, Hamdy M. Kalash, and Osama S. Farag Allah "Encryption Efficiency Analysis and Security Evaluation of RC6 Block Cipher for Digital Images" *International Journal of Computer, Information, and Systems Science, and Engineering* 1:1 PP.33-39, 2007, ISSN1307-2331.
- [8] Hossam El-din H. Ahmed, Hamdy M. Kalash, and Osama S. Farag Allah, "An Efficient Chaos-Based Feedback Stream cipher (ECBFSC) for Image Encryption and Decryption". *International Journal of Computing and Informatics*, VOL. 31, No. 1 PP. 121-129, 2007, ISSN 0350-5596242007.
- [9] Hossam El-din H. Ahmed, Hamdy M. Kalash, and Osama S. Farag Allah, "Encryption Quality Analysis of RC5 Block Cipher Algorithm for Digital Images." *Journal of Optical Engineering*, vol. 45(10), 107003(1-7), 2006.
- [10] Fridrich. "Symmetric ciphers based on two-dimensional chaotic maps," *Int. J. Bifurcation and Chaos*, 8(6):1259–1284, 1998.
- [11] . K. Pareek, V. Patidar and K. K. Sud. "Cryptography using multiple one-dimensional chaotic maps," *Communications in Nonlinear Science and Numerical Simulation*, vol. 10, Issue 7, October 2005, pp. 715-723.

Multi-element Circuits Based on LCLC Resonant Tank - Theory and Application

BRANISLAV DOBRUCKY, JURAJ KOSCELNIK

Department of Mechatronics and Electronics

University of Zilina

Univerzitna 1, 010 26

SLOVAKIA

branislav.dobrucky@fel.uniza.sk; juraj.koscelnik@fel.uniza.sk

Abstract: - The paper deals with novel of multi-element resonant circuits and its modification and application. Main circuits consist of series LC resonant branches and of parallel LC sinusoidal output filters. Review of multi-element circuits which contains more accumulation tanks is presented. Its mathematical description; design of accumulation elements and simulation experiments are given in main chapters of paper. Furthermore, mainly is focus given for application field and on possibility of its variable use. Multi-element circuits meet the requirements of the current market as are: small value of THD, high power density and very high efficiency. Also, such special circuits manifest inners self-regulation which provides resistance to short circuit. The paper also shows analysis of transient properties too. Base on the selected cirucits are suggested control methods. All simulation results are verified by experimental measurement created on physical sample. On the end the paper is discussed the application of this circuits and its possible variable use in industry.

Key-Words: - multi-element circuits; LCLC; NF modulation; power electronic systems; transient analysis; non-linear; short-circuit proof system.

1 Introduction

Concept of the resonant converters had greatly expanded into various industrial and consumer applications, such as power supplies for distributed systems, electrical drives, laptops, LCD televisions as well as for aerospace, automotive and energy systems sector. Focusing on the development of power electronic semiconductors, the involvement of resonant converters into various applications is a must due to continual improvements in technology of power semiconductors manufacture. In order to improve electrical properties of power converters suited for power energy systems, the paper deals with investigation novel of multi-element resonant circuits and its modification and application. First, standard multi-resonant circuits are analyzed, as well as mathematical model is developed. After it, modifications of LCLC circuit for better dynamic behavior during start-up and short-circuit operation are described, whereby both operation modes have been verified and compared with standard solution.

The main criterion was achievement of low value of THD of output variables (voltage and current). Based on analysis, physical model was constructed, whereby key results are experimentally supported at the end of the paper. Each design is confirmed with precise measurements [1].

Investigation of a short circuit in multi-resonant network circuit is also included in the analysis. Proposed topologies are based on LCLC resonant

circuits. The main focus is on its ability to withstand short circuit. During short circuit the output current is limited by the properties of resonant network which creates internal self-regulation. Is necessary, to determine appropriate control method in case of linear behaviour of the system. Each modification of standard LCLC converter is mathematically supported in order to understand basic design of multi-tank resonant converters for modern industrial applications [2], [3].

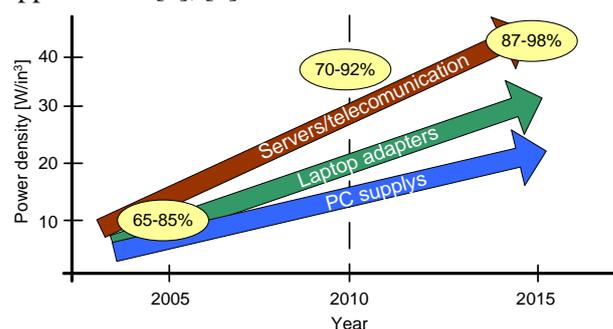


Fig.1 Trends course of development of power density [3]

A growing demand for saving energy and reducing the size of power systems have stimulated substantial research and development efforts towards high-efficiency and high-power density power supplies (Fig.1). The most effective way to achieve high power density in converters is to

increase the switching frequency so that the size of the passive components, such as the capacitor and inductor, as well as the transformer can be reduced, as they occupy a large portion of the overall size. Main design property of proposed converter is possibility to achieve low value of THD (below 5%) of output variables (voltage and current), at very high efficiency (over 97%). High system efficiency is one of the main quality indicators of power supplies. Therefore this parameter was investigated in wide region of switching frequency, in order to meet future demands on second quality factor – power density [3].

2 Multi-element Resonant Circuits

The group of resonant and quasi-resonant topologies consists of serial, parallel and serial-parallel resonant circuit. By combining the basic resonant circuits rise modified multi-elements resonant circuits. Resonant converters use two kinds of the switching technique: Zero voltage switching (ZVS) and Zero current switching (ZCS) [4]. Those techniques are now as soft switching. The converter can operate in ZVS and/or ZCS. The basic scheme of the resonant converter is given in the fig. 2 [12].

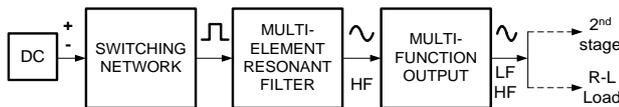


Fig.2 Block scheme of the resonant circuit

The scheme describes basic connection of the resonant converter composed by the DC source; switching network; resonant filter and multifunction output connected to load. Also, is possible to connect second stage converter (e.g. matrix converter). It's important to choose proper control method base on connected 2nd stage. Essence of the MRC's concept is combining the positive properties of conventional topologies in one device. Multi-resonant circuit can absorb all of the parasitic elements. Therefore, enable operation with low switching losses at high switching frequencies [3], [13].

2.1 Comparison of selected multi-element topologies

One of the novel types of converters are LCLC converters based on LLC resonant scheme, and LCTLC inverter consisting of DC/DC buck converter LCLC resonant filter and HF transformer. The HF transformer can also be connected after the LCLC filter, if necessary, and can also be used to

boost converter types. The inverter (LCTLC) is usually used as power supply for either HV rectifiers or HF cyclo-converters or matrix converters for 2-phase motor applications respectively [3], [5].

2.1.1 LCTLC circuit

The circuit is based on LCLC circuit where between serial (L_1, C_1) and parallel (L_2, C_2) accumulation tank is inputted HF transistor. It is fed by DC source and the shape of input voltage is switched in half-bridge connections [6]. In this case is output of the filter the HF harmonic waveform of the voltage (and current) is direct output mode with THD no more than 5% [5].

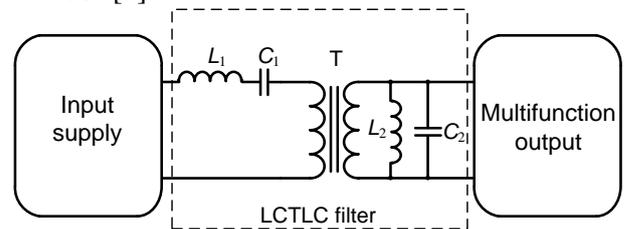


Fig.3 Block scheme of the LCTLC circuit

Equivalent scheme has been created to analysis circuit of proposed circuit. The equivalent parameters of the HF transformer ($L_\sigma, R_\sigma, L_m, R_{Fe}$ and inter-winding capacitance C_{iw} and inter-turn capacitance C_{it} are included into resulting component parameters. More about equivalent scheme is given in the paper [6]. Base on it the state-space equations for equivalent circuit with R-L load will be [4]:

$$\frac{di_{L1}}{dt} = \frac{1}{L_1}u(t) - \frac{R_1}{L_1}i_{L1} - \frac{1}{L_1}u_{C1} - \frac{1}{L_1}u_{C2} \quad (1)$$

$$\frac{di_{L2}}{dt} = \frac{1}{L_2}u_{C2} \quad (2)$$

$$\frac{du_{C1}}{dt} = \frac{1}{C_1}i_{L1} \quad (3)$$

$$\frac{du_{C2}}{dt} = \frac{1}{C_1}i_{L1} - \frac{1}{C_2}i_{L2} - \frac{1}{C_2 \cdot R_2}u_{C2} - \frac{1}{C_2}i_{LL} \quad (4)$$

$$\frac{di_{LL}}{dt} = \frac{1}{L_{load}}u_{C2} - \frac{R_{load}}{L_{load}}i_{LL} \quad (5)$$

where i_{L1}, i_{L2} are currents through the inductors L_1 and L_2 , respectively; i_{LL} is current through the load R_2, L_2 ; u_{C1}, u_{C2} are capacitors voltages of C_1 and C_2 , respectively, $u(t)$ is output voltage of the converter (filter input voltage). Using suitable numerical method or directly MATLAB functions the time waveforms of the quantities of LCTLC inverter can be obtained, Fig. 4.

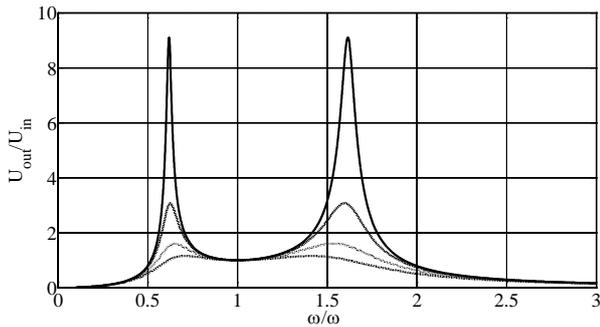


Fig.4 Voltage transfer of the LLCLC circuit

Voltage gain characteristic for proposed LCLC converter is shown in fig. 4. Principle, the shape is similar to standard LLC converter (left side of characteristic), whereby whole characteristic is clear combination of LLC and LCC converter (right side of characteristic). The circuit can be operated at higher or lower frequencies in ZVS region, achieving wide range of voltage gain [7].

2.1.2 LLCLC circuit

The circuit is well known and well described in the scientific literature. The resonant circuit is composed by serial-parallel LLC circuit and one parallel circuit. With one additional resonant element, a second band pass filter is created. A novel LLCLC resonant tank is proposed as an example. The structure is similar to the previously proposed four element resonant tank, but an extra resonant inductor is inserted.

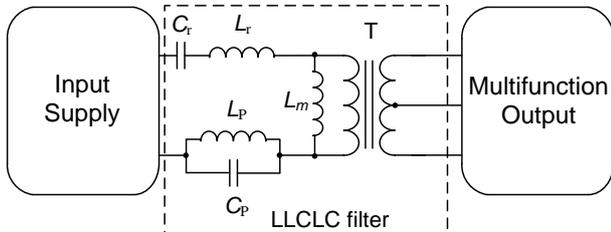


Fig.5 Block scheme of the LLCLC circuit

The half-bridge circuit is adopted as the primary-side structure. It is easy to extend to other types of input structures, including fullbridge, stacked half-bridge, and three-level structures. Similarly, it is easy to use other types of output structures, such as full-bridge, voltage-doubler and current-doubler structures.

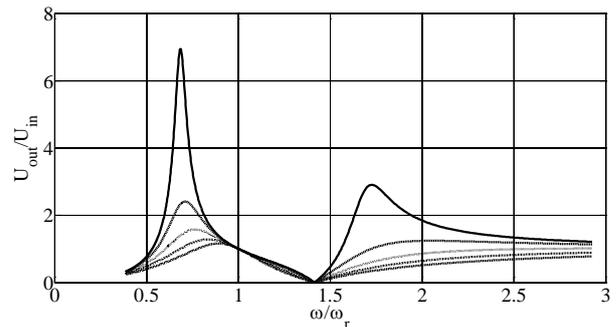


Fig. 6 Voltage transfer of the LLCLC circuit

The voltage gain of the proposed LLCLC resonant tank is illustrated in fig.6 (in range of load 10-100%). Conceptually, L_r , C_r and L_p contribute to the first band pass filter at low frequencies. The second band pass filter consists of L_r , C_r and C_p , which dominate at high frequencies. The first band pass filter can help to deliver the fundamental component to the load. It functions as the traditional resonant converters [6]. The second band pass filter enhances the power delivery with utilization of higher harmonics. Consequently, with the injection of higher-order harmonics, the reactive power of the resonant tank can be reduced and lower RMS current and lower conduction loss can be achieved. The output signal contains higher harmonics what is increasing RMS value of rectified voltage on output. However, THD of output voltage (before rectifier) is higher due to injected higher harmonics [3].

2.1.3 LCL2C2 circuit

Based on previously analyzed multi-element topologies was created new resonant circuit LCL2C2. Circuit is composed by one serial resonant tank and two parallel resonant tanks. This circuit is proposed as non-isolated circuit with brought out zero leg. Block configuration of components in LCL2C2 is given in fig.5.

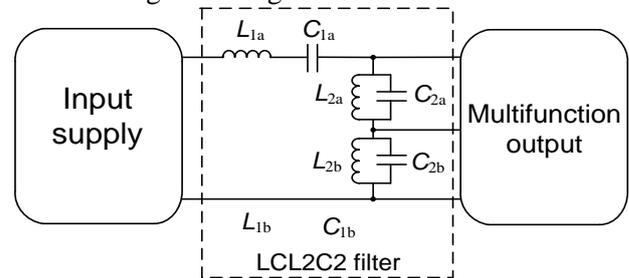


Fig.7 Block scheme of the LCL2C2 circuit

The "LCL2C2" circuit is one of the possible hybrid connections of resonant circuits. The main difference between the LCL2C2 and LCTLC converter is that the second one uses transformer to change the value of output voltage. In this kind of

circuit connection the control of output can be difficult. If the DC-AC inverter is considered, may be used the frequency or the asymmetric input control.

The multifunction output brings the possibilities for AC and DC output as well. Basically, output of the converter can be considered in three ways, i.e.:

- a) direct AC output
- b) diode rectifier
- c) AC output with variable or constant frequency

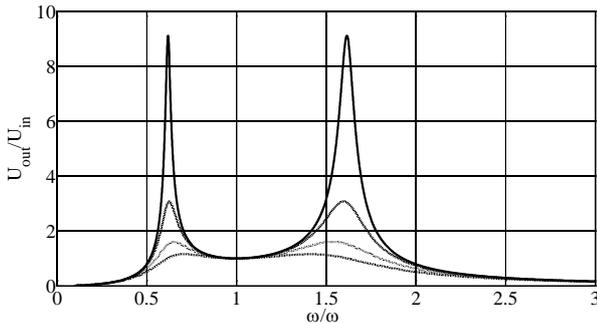


Fig.8 Voltage transfer of the LCL2C2 circuit

Voltage transfer is identical to the LCLC circuit. Fig. 8 shows gain curves in dependency on the load change. Based on this characteristic it is able to determine the proper operation regions of LCL2C2 converter. Above resonant frequency, which is point of f_{rel} is equal 1, the region with ZVS conditions for switching transistor are achieved, whereby boarder between ZVS and ZCS region is limited by the peak gain values of each gain curve. Similar relation is valid below resonant frequency, whereby ZVS and ZCS region are mirrored compared to region above f_{res} [12].

2.2 General Design of Accumulation Components

The resonant frequency of LC components should be the same as basic fundamental frequency of the converter and is governed by load requirements. Thus, based on the Thomson relation

$$\omega_{rez} = \frac{1}{\sqrt{L \cdot C}} \quad (1)$$

or, respectively

$$L \omega_{res} = \frac{1}{\omega_{res} C} \quad (2)$$

where ω_{res} is equal $2\pi \times$ fundamental frequency of the converter. Values of storage LC components and their parameters are important for properties of LCLC filter. Theoretically, $\omega_{res} L_1$ and other values of the converter can be chosen from a wide range [2], [6]. For our first design approximation we suppose a simple resonant circuit with a resonant

frequency equal to the switching input frequency ($\omega_{res} = \omega_{sw}$).

The LC design process can be considered from 3 different points of view or criteria:

1st: nominal voltage and current stresses at steady-states,

2nd: minimum voltage and current stresses during transients,

3rd: required value of total harmonic distortion of the output voltage.

In order to not exceed nominal voltages of the storage elements has been used value of internal impedance of the storage element equal to the nominal load $|Z_N|$.

Let's define the nominal design factor q_N for LC components as [4], [7]

$$q_N = \frac{L \omega_{res}}{|Z_N|} = \frac{1}{\omega_{res} C |Z_N|} \quad (3)$$

The above equation is similar to quality factor defined by $q = L_{load} \omega_{res} / R_{load}$, however q_N does not depend on the load R_{load} .

The design formulas for LC accumulation elements can obtain:

$$L = \frac{U_1^2}{\omega_1 P_1} q_N \quad C = \frac{P_1}{\omega_1 U_1^2} \frac{1}{q_N} \quad (4a,b)$$

The voltage on storage elements at nominal steady-state is defined as

$$U_C = \frac{1}{\omega_{res} C} I_N q_N = \frac{1}{\omega_{res} C} \frac{P_1}{U_1} q_N \quad (5)$$

That means that for q_N equal to one, the voltages on storage elements will be nominal values, and are proportionally depend on q_N factor.

Going back to LCLC filter, then

$$L_1 = \frac{U_1^2}{\omega_1 P_1} q_N \quad C_1 = \frac{P_1}{\omega_1 U_1^2} \frac{1}{q_N} \quad (6a,b)$$

$$L_2 = \frac{U_1^2}{\omega_1 P_1} \frac{1}{q_N} \quad C_2 = \frac{P_1}{\omega_1 U_1^2} q_N \quad (7a,b)$$

where U_1, P_1, ω_1 are nominal output voltage, power and frequency, respectively (fundamental harmonic)[3], [5].

2.2.1 Experimental Simulation of Multi-element Circuit

Simulation model was crated according to the designed parameters of multi-element resonant circuits. MATLAB environment has been use to provide all the simulations experiments using suitable numerical method or directly preprogrammed

functions. Time waveforms are given in following figures.

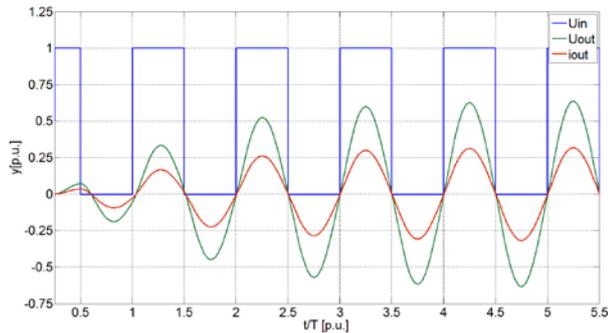


Fig.9 Simulated waveforms of input and output voltage and current (per unit)

Based on theoretical assumptions, the system is operating in ZV/ZC mode. Waveforms of current and voltage on the switching transistor during operating process are given on fig. 4. The converter switches in zero voltage (ZVS) what is preferred operating area for the MOSFET transistors. ZVS conditions have been achieved moving the switching frequency above the resonant frequency.

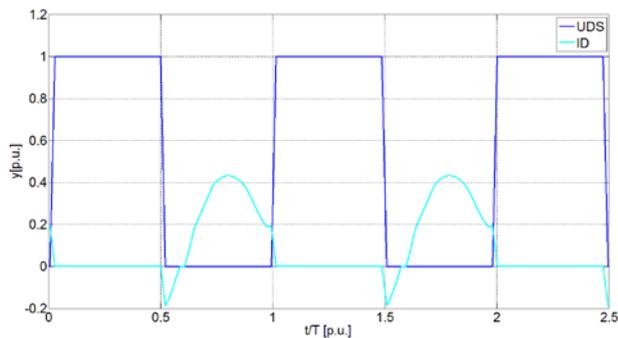


Fig.10 simulated waveforms of current and voltage on the switching transistor

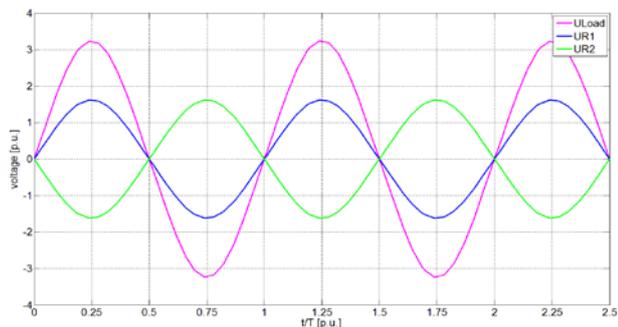


Fig. 11 symmetrical output of LCL2C2 circuit

The waveforms in the fig. 11 are showing the output voltage and voltage on the both of branches with symmetric output of LCL2C2 circuit. The load voltage (also voltage on the branches) has harmonic shape with low THD value. The simulated waveforms are matching with the theoretical

assumptions. The determination of the THD value is given below.

Using Fast Fourier Transformation (FFS) was possible to calculate THD of output voltage. The resonant components of the filter are toned on basic harmonic; therefore the higher harmonic contents are suppressed.

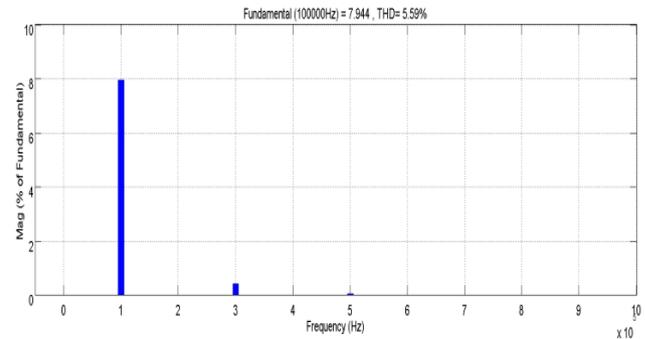


Fig.12 The harmonic content of LCLC converter

The total harmonic distortion (THD) was 5.59%. Because the simulation model considered with parasitic elements the value of THD raised over 5%.

3 Control methods of multi-resonant topologies

It is necessary, to determine appropriate control method in case of linear behavior of the system. Under the condition of non-linearity and taking parasitic into account, occurs the change of f_{res}/f_{sw} ratio.

It is possible to control output voltage by the classical frequency control method connecting corresponding kind of converter on the output side of system. Corresponding converter are rectifier or cyclo-converter. This regulation is suitable for circuits with transformer. However, transformer brings additional losses to system. LCL2C2 circuit is transformer less circuit so; it may be considered other kind of control. By considering the different ways the inverter output can be choose the control method [8], [10].

3.1 Frequency ratio change control

One of the simplest ways how to control resonant circuit is change of ration between switching and resonant frequency. Maximal gain of output voltage is when f_{res}/f_{sw} is equal 1. Changing the switching frequency is possible to change magnitude of output signal. However, increasing the ration will grow number of harmonics contained in output signal.

3.2 Nonsymmetrical control method

The real output voltage of inverter waveform has a wide spectrum of harmonic components. Using nonsymmetrical control the output voltage of inverter (Fig. 13) comprises all harmonic components.

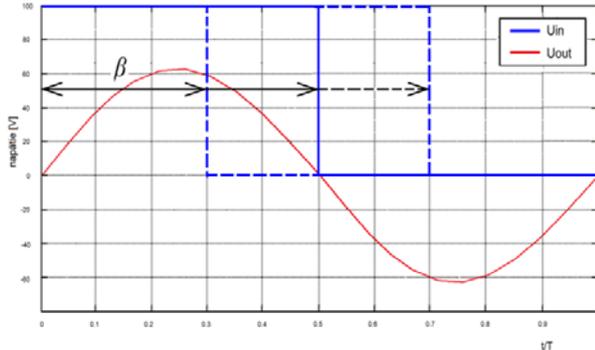


Fig.13 principal method of nonsymmetrical control

Used resonant circuit is tuned on fundamental harmonic, but it should be tuned on switching frequency of converter as well. By the nonsymmetrical change of angle is possible to control magnitude of output voltage.

$$U_{1M}(\beta) = \frac{4}{\pi} U \cdot \sin(\beta/2) \quad (8)$$

where $U_{1M}(\beta)$ is magnitude of fundamental harmonic depend on pulse width and U is maximal value of output voltage. With growing β angle is raising number of even harmonics included in input signal.

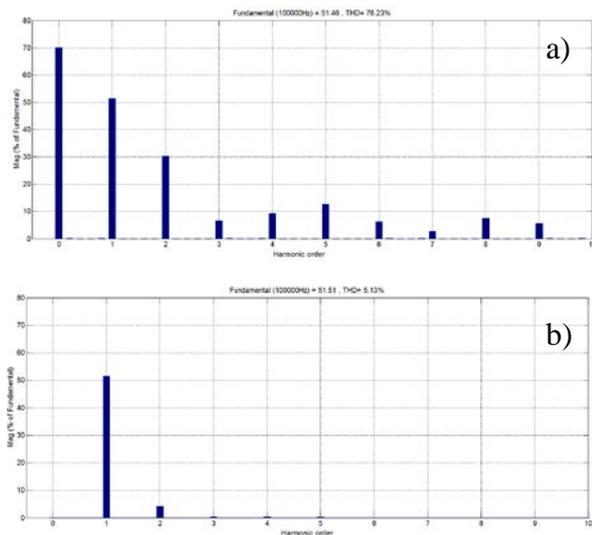


Fig.14 Harmonic order included in input a) and output b) voltage for $\beta=30/70$ [%]

3.3 LF modulation of input voltage

One of possible way is to control output voltage on input side of converter. In this case would be voltage regulated prior to entering to resonant

circuit. This is possible provide by bipolar PWM with LF modulation.

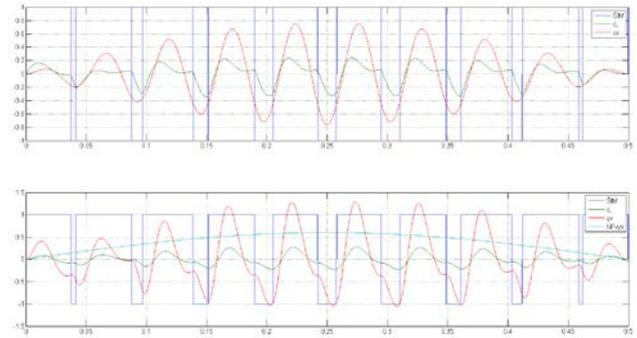


Fig.15 NF modulation of output signal

Basic harmonic HF signal with NF modulation passing resonant circuit (tuned to the HF signal-switching frequency) with unit gain and the output voltage is rectified (HF Schottky diode rectifier with respectively. MOSFET transistors in inverse mode) or modified (cycloconverter, matrix converter) and HF frequency signal component is removed by simple passive LC filter.

3.4 Multi-Resonant converter's inner self-feedback

Let's consider LCLC circuit as common model for all multi-resonant circuits presented in the paper.

Determine appropriate control method is necessary in case of linear behavior of the system. Under the non-linearity condition and taking parasitic into account, occurs the change of f_{res}/f_{sw} ratio. This change creates circuit's inner self-feedback and it provides limitation of short circuit current. It causes saturation of magnetic elements and the change of the inductor inductances values. Therefore, the ratio between switching and resonant frequency is changing. Frequency ratio change impacts on the point of maximal gain. These phenomena can be considered the method of self-regulation due to own internal feedback. This phenomenon is caused by non-linearities of these circuits [11].

3.4.1 Multi-resonant non-linear circuitry

Modelling deals with Euler - and Taylor expansion methods for consequent numerical solution in Matlab environment. As an example, electrical circuit with serial rectifier diode [10], [11] includes thenonlinearities as:

- input voltage
- non-linear capacitance of diode (important for PV applications)

- magnetic circuit of serial inductance and HF transformer magnetic circuit of transformer and parallel inductance

Created model of multi-resonant circuit was updated with nonlinear inductor. Mathematical method used in model nonlinear inductor was fictitious exciting functions method. More about this method is given in [11].

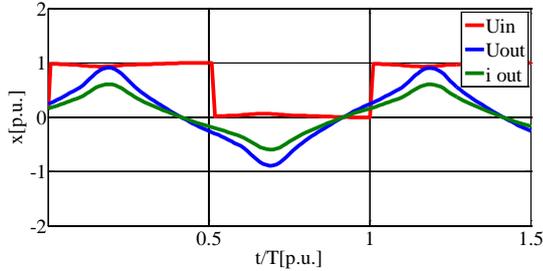


Fig. 16 Model of LCLC with nonlinear inductor (in state of saturation) [10]

The disadvantage of this regulation is that it causes the change in shape of output current and its THD increases by 3-5%. Also, is possible to consider over-dimensioned the accumulation elements where will saturation not occur.

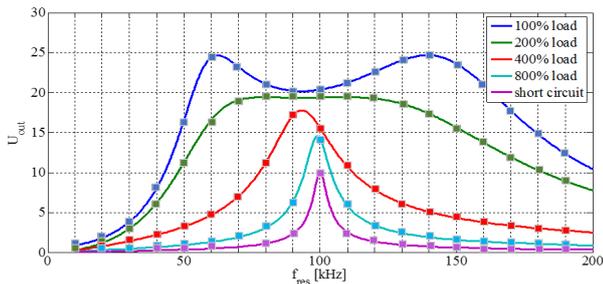


Fig. 17 Voltage transfer of LCLC (parasitic and non-linear elements included in model)

Short circuit causes that output current increase and saturation of magnetic elements. Result is change of the inductor inductances values. Therefore, the ratio between switching and resonant frequency is changing. What is based on Thomson relationship (1). Also, change of frequency ration impacts on the point of maximal gain and moving resonant frequency (fig.17).

4 Transient Properties

Simulation model of multi-element circuit is built by applying knowledge of basic resonant circuits. Non-linear electronic elements as semiconductor devices and ferromagnetic inductors are included in model. Voltage transfer functions and impedance - frequency dependencies are theoretically derived, calculated, computationally simulated and analysed.

Besides, the output voltage value does not depend on the load value.

Simulation model is based on the equations (1-7) for the design of accumulation components from previous chapter [8].

4.1 Analysis in frequency domain

Let's define nominal impedance for series *resistive-inductive* load

$$|Z_N| = \sqrt{R_{sload}^2 + (\omega L_{sload})^2} = \frac{U_{outN}^2}{P_{outN}} \quad (9)$$

and nominal admittance for parallel *resistive-inductive* load

$$|Y_N| = \sqrt{\left(\frac{1}{R_{pload}}\right)^2 + \left(\frac{1}{\omega L_{pload}}\right)^2} = \frac{P_{outN}}{U_{outN}^2} \quad (10)$$

On the beginning will be defined simple resistive load. Impedance of series and parallel part of the LCLC filter is defined by the following equations

$$\begin{aligned} Z_1(\omega) &= R_1 + j\left(\omega L_1 - \frac{1}{\omega C_1}\right) = \\ &= \frac{R_1}{|Z_N|} |Z_N| + j|Z_N| q_{N1} \left(f_{rel} - \frac{1}{f_{rel}}\right) \end{aligned} \quad (11)$$

Where in R_1 is substitute the sum of resistance of series part of the filter (e.g. resistance of series filter coil; of filter capacitor; ...).

Thus

$$\frac{|Z_1(\omega)|}{|Z_N(\omega)|} = \sqrt{r_1^2 + \left[q_{N1} \left(f_{rel} - \frac{1}{f_{rel}}\right)\right]^2} \quad (13)$$

Edited mathematical model of input impedance of LCLC looks:

$$\frac{|Z_{in}(\omega)|}{|Z_N(\omega)|} = \sqrt{\left[r_1 + \frac{\left(\frac{1}{r_2} + \frac{1}{r}\right)^2}{DEN}\right]^2 + \left[\left(f_{rel} - \frac{1}{f_{rel}}\right) \left(q_1 - \frac{q_{N2}}{DEN}\right)\right]^2} \quad (14)$$

where denominator marked DEN is defined as:

$$DEN = \left(\frac{1}{r_2} + \frac{1}{r}\right)^2 + \left[q_{N2} \left(f_{rel} - \frac{1}{f_{rel}}\right)\right]^2 \quad (15)$$

Impedance presentation in frequency domain in different states of load will be:

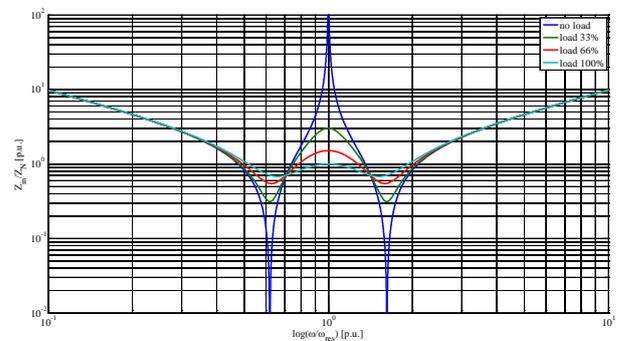


Fig.18 Filter input impedance vs. frequency in range 0-100% of load

Fig. 18 shows impedance dependence on frequency ratio. The ratio is composed by resonant frequency f_{res} and switching frequency f_{sw} what creates relative frequency. In point, where is $f_{rel}=1$ and load=0 impedance value grows to infinity. Impedance transfer is equal 1 where is f_{res}/f_{sw} equal 1, what ensures that circuit operates in resonance. Also, there is possible to choose proper operation area.

As well as impedance transfer is possible to create voltage transfer model.

$$F(\omega) = \frac{|Z_2(\omega)|}{|Z_{in}(\omega)|} = \frac{\sqrt{\frac{1}{DEN}}}{\sqrt{\left[r_1 + \left(\frac{1}{r_2} + \frac{1}{r}\right) \frac{1}{DEN}\right]^2 + \left[\left(f_{rel} - \frac{1}{f_{rel}}\right) \left(q_1 - \frac{qN_2}{DEN}\right)\right]^2}} \quad (16)$$

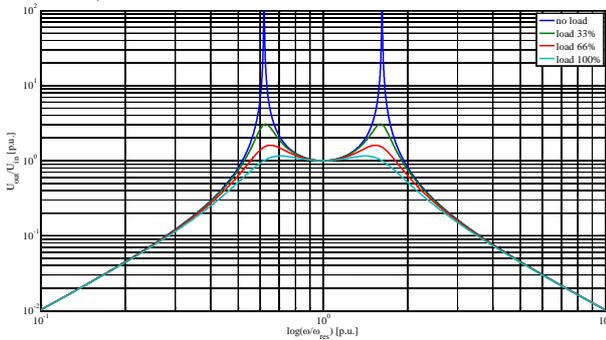


Fig.19 Voltage transfer function U_2/U_1 of LCLC filter in range 0-100% of load

Voltage transfer function of LCLC resonant circuit is given in fig. 7. The transfers curves are changing depend on the load (0-100%). However, in the resonance point ($f_{res}/f_{sw}=1$) is voltage transfer equal 1, what means that the system is no depend on load.

Adding inductance into the final model, voltage transfer and impedance transfer for complex RL load will look:

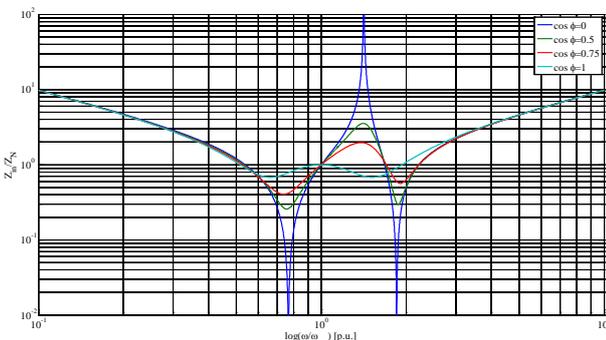


Fig.21 Impedance transfer function U_2/U_1 of LCLC filter in range 0-100% of RL load;

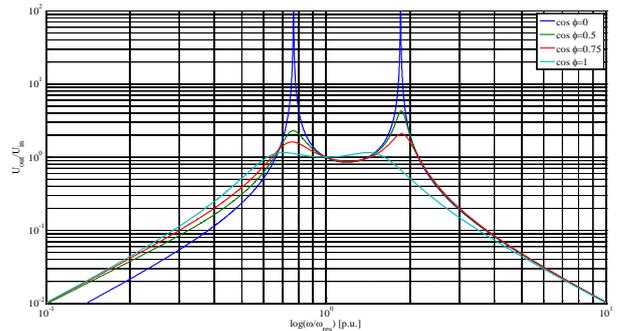


Fig.21 Voltage transfer function U_2/U_1 of LCLC filter in range 0-100% of RL load

Fig. 20 and 21 show impact of impedance included in RL load. Figures show change of $\cos \phi$ value in its impact on transient properties.

4.1.1 Choosing proper operation area

Based on input previous analysis is possible to choose two izo-impedance (invariant impedance) operational points for switching frequency. In this case input impedance is not depending on the load of the inverter. Two mirror trajectories with minimal input impedance of the resonant circuit depending on the load. First point is when impedance is proportional depended on the load (fig.2).

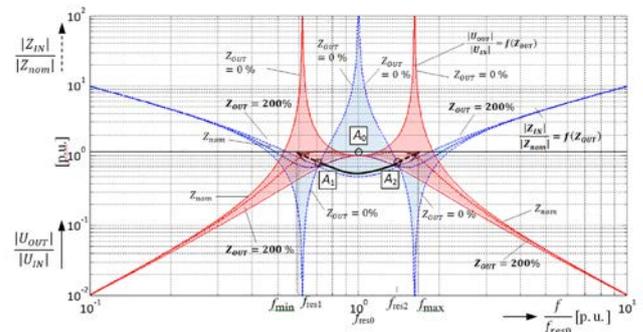


Fig.22 Input impedance and voltage transfer frequency log-characteristics [8]

Similarly, voltage transfer frequency characteristic of the LCL2C2 resonant circuit offers two mirror trajectories with maximal output voltage of the circuit depending on the load, and also one point (A0) when the output voltage of the inverter does not depend on the inverter's load. Also, is possible to determine the optimal operation frequencies for other value of overloading and functional relation is

$$\begin{aligned} |f_{min}|_{overload} &= f(Z_{overload}) \text{ or} \\ |f_{max}|_{overload} &= f(Z_{overload}), \end{aligned} \quad (17)$$

to input current was be the same as nominal one. Carried-out results are original ones, and in spite of non-linear circuitry the output voltage THD is staying rather small, about 7-11 %.

In special operation states, multi-element circuits may present by inner self-feedback. This self-feedback provides limiting of short circuit current. It helps prevent destruction of the device [10].

5 Experimental verification

The paper deals with novel of multi-element resonant circuits. Higher discussed topologies and its properties were verified no physical samples. Experimental measurements fully respond theoretical assumptions and simulation experiments.

5.1 Voltage and current quantities of multi-resonant circuits

For MOSFET device is preferred operation area ZVS, what has been achieved during experimental measurement.

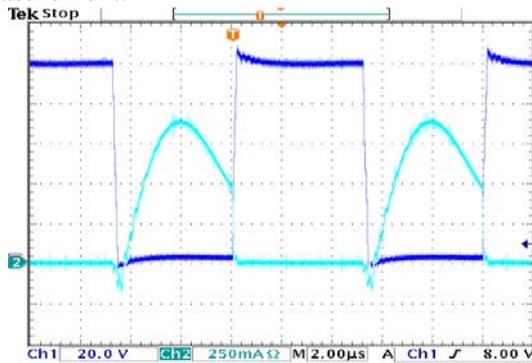


Fig.22. the waveforms of current and voltage on the switching transistor

Fig. 22 is showing current and voltage of the switching transistor during operation in recommended region (slightly above resonant frequency). Nevertheless ZVS conditions are achieved, thus transistor is operated with very low switching losses.

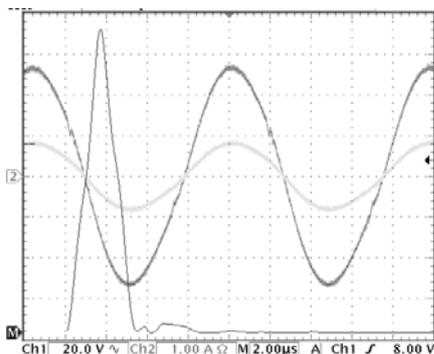


Fig.23 Switching waveforms at the output (current - grey, voltage - black) with FFT analysis of output voltage waveform

Based on the FFT analysis of output voltage we proceeded to calculate the real THD value. For this purpose we used next equation:

$$THD (\%) = \frac{\sqrt{U_2^2 + U_3^2 + U_4^2 + U_5^2 + \dots + U_n^2}}{U_{rms}} \cdot 100\%$$

, where U_2, U_3, U_4, U_5, U_n are parts of higher harmonic order, and U_{rms} is root mean square value of output voltage. Based on this equation the computed THD value of output voltage of proposed converter is 4.02 [%].

5.2 Transient analysis experiments

Transient analysis was prepared on special physical sample –frequency tester.

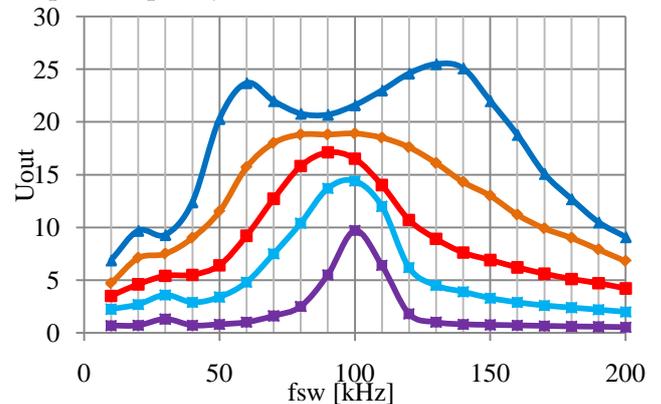


Fig.6. Voltage transfer of LCL2C2 converter (experimental verification; from 100% of load until short circuit)

The shape of transfer waveforms is similar to the simulation results (fig. 17). Output voltage values compared to simulation are similar too. The perceptual difference is from 3 till 15%. The biggest difference is in case of short circuit where the values in 80-90 kHz and between 110-120 kHz. Best match is observed at nominal load. Frequency ratio is changing as it was in simulation experiments where considered nonlinearity and parasitic elements was [9], [10].

These phenomena can be considered the method of self-regulation due to own internal feedback. I case of short circuit is output current limited by the converter self-regulation. Also, the regulation causes the current shape distortion and its THD increases about 3-5%. THD values about 6-9%.

5.2.1 Multi-resonant non-linear experiments

Shape of voltage and current on the load during short circuit is shown on following picture.

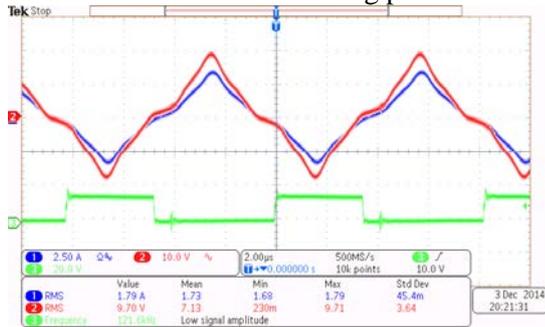


Fig. 5 Output voltage and current quantities in state of short circuit (experimental)

Experiments show that the real device is output current limited during short circuit. The shape has changed and THD value of output current and voltage increased. This impacts to the transfer properties and creates inner self-feedback.

Experimental measurements show that all mathematical models and theoretical assumptions analyzed in article was successfully verified.

6 Conclusion

In the paper was discussed about multi-resonant circuits' their theory and application. Selected circuits are described and analyzed in second chapter. All tree topologies are based on LCLC circuit. LCTL and LLCLC are general know circuits. However, LCL2C2 can be considered as new multi-resonant circuit whit many advantages. Specific control methods are briefly analyzed as possible way to regulate output voltage of these circuits. In particular, LF modulation of input voltage appears to be promising way to control output voltage, especially for LCL2C2 circuit. Base on mathematical models of those multi-resonant circuits was carried out transient analysis. Circuits have been investigated in different stages of load (in wild range). Under the obtained results was possible to choose proper operation areas of multi-resonant circuits. Everything is depends on design and final application. Simulation results shown, that LCL2C2 circuit has inner self-feedback. This feedback provides prevention against short circuit. To understand better this method of self-protection was necessary to create system with non-linear elements. Using fictitious exciting functions method was possible to simulate this system. Under the non-linear condition occurs the change of f_{res}/f_{sw} ratio. It provides limitation of short circuit current. It causes

saturation of magnetic elements and the change of the inductor inductances values.

Theory, models and simulation result was verified by experimental measurements provide on physical samples prepared in our laboratories. All simulation, including non-linearity and inner self-feedback were confirmed by real experiments. The article can serve as a guide for the analysis of multi-resonant circuit.

References:

- [1] M. M. Jovanovic; Technology drivers and trends in power supplies for computer/telecom, *APEC 2006*, Plenary session presentation.
- [2] I. Batarseh, Resonant Converter Topologies with Three and Four Storage Elements. *Power Electronics, IEEE Transaction on*, Vol. 9, No.1, Jan 1994, pp. 64-73.
- [3] J. Koscelnik, M. Frivaldsky, M. Prazenica, R. Mazgut, A review of multi-elements resonant converters topologies, *ELEKTRO, 2014 Publication Year: 2014*, Page(s): 312 - 317
- [4] A.K.S. Bhat, Analysis and design of LCL-type series resonant converter, in *Proc. IEEE INTELEC*, 1990, pp: 172 - 178.
- [5] P. Imbertson, N. Mohan, Asymmetrical Duty Cycle Permits Zero Switching Loss in PWM Circuits with No Conduction Loss Penalty. *Industry Applications, IEEE Transaction on*, Vol. 29, No. 1, Jan/Feb 1993.
- [6] Y.A. Ang, M.P. Foster, C.M Bingham, D.A. Stone, H.I. Sewell and D. Howe, "Analysis of 4th-order LCLC Resonant power converters", *IEE Proc. vol. 131*, no. 2, pp. 169-181, 2004.
- [7] Dianbo Fu: Novel Multi-Element Resonant Converters for Front-end DC/DC Converters, *Power Electronics Specialists Conference, 2008. PESC 2008. IEEE*, 15-19 June 2008, ISBN 978-1-4244-1668-4.
- [8] B. Dobrucky, M. Firvaldsky, J. Koscelnik, Choosing operational switching frequency of LCTL resonant inverter, *In-Tech 2014*, Proceedings, 10-12 september, Leiria, Portugal, pp. 187-190.
- [9] Y. Ang, C.M. Bingham, M.P. Foster, D.A. Stone, Analysis and Control of Dual-Output LCLC Resonant Converters, and the Impact of Leakage Inductance, *Power Electronics and Drive Systems, 2007. PEDS '07. 7th International Conference*, 27-30 Nov. 2007, pp.145-150.
- [10] M. Condon, R. Ivanov, Nonlinear systems – algebraic gramians and model reduction, *COMPEL: The International Journal for Computation and Mathematics in Electrical*

and Electronic Engineering, Vol. 24, Iss. 1, pp.202 – 219, 2005

- [11] J. Koscelnik, J. Sedo, B. Dobrucky, Modeling of Resonant Converter with Nonlinear Inductance, *Applied Electronics 2014*, 19th International Conference, Sept. 9-10, Pilsen(CZ), 2014
- [12] H.M. Suryawanshi, S.G. Tarnekar, Modified LCLC-Type Series Resonant Converter with Improved Performance. *IEE Proc. on Electrical Power Applications*, 1996, 143, (5), pp. 354–360.
- [13] H. Daocheng,: A Novel Integrated Multi-Elements Resonant Converter, *Energy Conversion Congress and Exposition (ECCE)*, 2011 IEEE, 17-22 Sept. 2011, ISBN 978-1-4577-0542-7.

Parallel Adaptive Arbiter for Improved CPU Utilization and Fair Bandwidth Allocation

M. Nishat Akhtar¹ and Junita Mohamad-Saleh²
^{1,2}School of Electrical and Electronics Engineering,
Universiti Sains Malaysia, 14300 Nibong Tebal,
Penang, Malaysia
¹nishat_akhtar2000@yahoo.com, ²jms@usm.my

Abstract: Nowadays, task parallelism is recognized to be a huge challenge for future extreme scale computing system. Advancement in parallel computing system necessitates solving the bus contention in a most efficient manner along with high computation rate. An arbiter receives bus requests from master components to grant the bus access. Therefore, an arbiter plays an important role in solving the bus contention in any System-on-Chip (SoC). This article presents a new technique for arbitration called Parallel Adaptive Arbitration (PAA) to maximize the usage of CPU cores along with fair and moderate bus bandwidth allocation. This arbitration algorithm is developed for heterogeneous masters designed according to the different traffic behavior to enable high degree of task parallelization. The results reveal that PAA is more advantageous than the other conventional arbitration algorithms for several reasons including utilization of CPU cores up to its maximum extent using synchronization of heterogeneous masters. The proposed arbitration technique could be a promising approach for designing SoC for future applications.

Keywords: Parallel adaptive arbiter, multi-core, CPU utilization, parallel arbitration, System-on-Chip, bandwidth allocation

I. INTRODUCTION

Task parallelism could be an elegant programming paradigm for symmetric multiprocessing to reveal uneven parallelism in an efficient manner. In a multiprocessor system, task parallelism could be achieved if each processor is designed to execute a unique thread on a similar or different data. Symmetric multiprocessing has become an operable option for computing in the world of embedded systems where technology is blended with complex chips that incorporate multiple processors dedicated for specific computational needs. In order to realize this complex multiple SoC in the environment of

intellectual-property based methodologies, communication architecture plays a major role. In any SoC, arbitration algorithm plays a major role in order to solve bus contention. Fairness is a property that plays a very crucial role among the various criteria of arbitration algorithms in solving bus contention. The performance of multiprocessors systems depends more on efficient communication among processors and on the balanced distribution of computation among them, rather than on pure speed of the processor. Since arbiters are invoked for every transfer on the bus, they are considered to be in the critical path of bus-based communication architecture and must be designed with great care [1]. An efficient contention resolution scheme is required to provide fine-grained control of the communication bandwidth allocated to individual processor and avoid starvation of low priority transactions [2].

Tuning task parallel application still remains a challenge in the world of parallel computing, however languages and tools like OpenStream, OpenMP, Java, OpenMPI, hadoop etc. can only make the development of application more simple. In symmetric multiprocessing the major architectural bottleneck is the internal bus which connects the processors and peripherals to the memory using an arbitrary network of shared channels. In most cases, the bus bandwidth becomes a dominant barrier because of improper bandwidth allocation. To maintain the bus bandwidth in an efficient manner, the process of memory arbitration cannot be neglected as it is one of an essential factor for concurrent-computing. In an enhanced arbitration environment of SoC, the communication architecture should be fair enough to offer high performance to a wide range of masters according to their traffic

behavior, because masters on a SoC bus may issue multiple requests at the same time. Thus, an arbiter plays an important role to decide on which master should be granted bus access first. Hence, this paper presents a PAA technique, designed in such a way that it suits a system by maintaining high throughput, low starvation among the different masters and attaining high degree of task parallelization.

The rest of the paper is organized as follows; Section II gives a brief overview of related works done by the researchers. Section III gives the concept of masters designed according to traffic behavior. Section IV elaborates on the proposed PAA technique. Finally, section V and section VI present the results and conclusion, respectively.

II. RELATED PREVIOUS WORKS

In recent years, many researchers focused on developing multi-level arbitration scheme in order to reduce system latency and to achieve fair bandwidth allocation. Yi et al. [3] proposed an arbiter called an adaptive dynamic arbiter. They proposed a lottery bus algorithm approach where an arbiter can adjust the bandwidth proportion and automatically assign it to its associated processor. Compared with conventional architectures, their architecture reduces the system latency but it does not allocate fair bandwidth to the processors. Moreover, their architecture is unable to maximize the CPU cores utilization as it is not implemented using parallel programming approach. Aravind [4] presented an algorithm which is a fully distributed software solution to the arbitration problem in multi-port memory systems. His algorithm is purely based on first in first out and least recently used fairness criteria. However, the algorithm does not deal with fair bandwidth allotment to the different masters which may become a barrier to obtain better performance. Moreover, their arbitration technique follows a sequential programming approach and therefore does not make use of the multiple CPU cores leading towards high task execution time.

Massimo and Poncino [5] proposed a novel method of automatic synthesis of easily scalable bus arbiters with dynamic priority assignment strategies.

They emphasized more on those arbitration mechanisms which can be implemented on silicon as a digital circuit, rather than being concerned about how the selected arbitration policies can affect the performance of a multiprocessor system. Their arbitration technique was fair in terms of bandwidth and latency. Nonetheless, it puts the least concern regarding CPU utilization as it did not take parallel approach into account. The major disadvantage of common-bus multiprocessor system is the reduction of throughput caused by conflict among processors requiring access to the shared memory. Ideally, throughput should increase directly with the number of processors but the bus contention diminishes this increasing trend [6]. There is a critical number above which the processors show no improvement and this critical number depends naturally, on the extent of bus used by the processors.

Chen et al. [7] designed a real-time and bandwidth guaranteed arbitration algorithm for SoC bus communication in which RT_Lottery algorithm has been used to meet both hard real-time and bandwidth requirements. However, in terms of fair bandwidth allocation it cannot compete with adaptive arbiter (which will be discussed further) as its bandwidth allocation is quite diverse. Their work demonstrated a two-level arbitration scheme which comprised of time division multiple access algorithm and lottery-based algorithm. They developed master cores according to the traffic behavior of the data flow which consists of both heavy traffic masters and light traffic masters. On the other hand, their masters did not show synchronization among them and were implemented using sequential programming method. Therefore the masters were unable to maximize the CPU utilization. However, in terms of diverse bandwidth allocation, their arbitration technique was superior and was able to handle hard real-time bandwidth requirement.

A unique algorithm was proposed by Li et al. [8], called adaptive arbitration algorithm in which an arbiter can automatically adjust priority to provide the best bandwidth for different master according to their real time bus bandwidth needs. They showed that, it is possible to allocate fair bandwidth to a given set of processors with a very high degree of fairness. In their case, an arbiter records the number

of time each master has requested for the bus and the total time that all master have requested for the bus access. Using these two values, the arbiter can calculate the bus access probability of the corresponding master by the division operation method. The priority weight of the master is decided by its probability of getting the bus access. A master with bigger weight owns higher priority. This way, it is unnecessary for an arbiter to recalculate all the probabilities and weights, and to reorder the priority of masters when a new bus access request appears. The solution to this problem is to reduce the frequency of weight calculation and priority reordering [8]. Their arbiter worked well in terms of fair bus bandwidth allocation. On the other hand, it did not exploit multi-core parallelism.

Akhtar and Sidek [9] proposed a novel arbitration technique called Intelligent Adaptive Arbitration (IAA) which works on the principle of parallel computation and offers moderate bandwidth allocation to its masters along with low latency. IAA is promising arbiter in a sense that it solves the problem of task parallelization by treating each master with moderate bandwidth. However, there is a limitation at the synchronization level of its masters which still remains unexplored.

III. DESIGN OF MASTERS RELATED TO TRAFFIC BEHAVIOR

In terms of traffic behavior, three types of masters has been designed to implement real-time and bandwidth-guaranteed arbitration using non-preemptive RT_Lottery algorithm [7,9]. More emphasis was given on variable bandwidth allotment to handle real-time requirements rather than parallel implementation of the masters. Moreover, in terms of critical bandwidth requirements, RT_Lottery algorithm holds good. For the RT_Lottery algorithm, the weight function (W) was used to decide the bus access priority by each master. The weight function was calculated using the following equation:

$$W = \frac{\text{Required Bandwidth}}{\text{Max Bandwidth (set by user)}} \quad (1)$$

This weight function was further analyzed or tuned to check the extra bandwidth requirement for each master. To implement RT_Lottery algorithm, 6 masters were used to process high beat burst data [7]. The bandwidth requirement for each master in the case of RT_Lottery algorithm was set manually. However, manual bandwidth allotment is unfavorable for parallel implementation of masters, as conflicts may arise among the masters in terms of bandwidth sharing. Table 1 shows allotted bandwidth values using RT_Lottery algorithm [7]. In the table 1, M1, M2, M3, M4, M5 and M6 represents the respective masters.

The three types of masters are appropriate to be utilized for real-time bandwidth-guaranteed algorithms. Recently, similar kind of masters were used by Akhtar and Sidek [9] for their proposed arbitration technique called IAA. The results were appropriate and satisfying in terms of bandwidth fluctuation. The average bandwidth fluctuation in case of IAA was $\pm 1.35\%$ which was better than other conventional arbitration techniques [9].

The first type of master is D_Type master where 'D' stands for dependent. This type of master has no real-time requirements. Its upcoming request is totally dependent on the completion time of a current request. A time interval is the duration between the times when the request is issued till the time it is finished [7,9]. Figure 1 shows an example of interval time. A burst is generated of 4 beat at cycle 13. However, the acceptance of 4 beat burst comes between cycle 16 to cycle 20. If suppose the interval time is 10 cycles, then the next request can be issued only at cycle 30. The uni-directional line represents that no request can be issued between cycle 20 and cycle 30. Chen et al. [7] used this master to implement RT_Lottery algorithm. This master worked well in terms of processing high beat burst data. However, there was a huge difference between the required bandwidth of master and the maximum allotted bandwidth as shown in Table 1. D_Type master was also used to implement IAA developed by Akhtar and Sidek [9]. For IAA, this master acquired fair and moderate bandwidth to enhance the overall bandwidth optimization [9].

Table 1. Allotted bandwidth values using RT_Lottery algorithm

	(D_Type) M1	(D_Type) M2	(DR_Type) M3	(DR_Type) M4	(NDR_Type) M5	(NDR_Type) M6
Maximum Bandwidth (%)	63	18	63	19	17	2
Required Bandwidth (%)	20	5	40	10	17	2

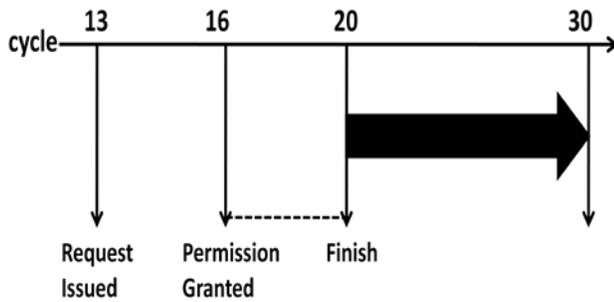


Fig 1. D_Type master

The second type master, DR_Type master is the same as D_Type master, except for its assigned real-time parameter (R) [7,9]. Let us suppose that R is set to 10 cycles for the master and its associated request is issued at cycle 5. This means that the request issued at cycle 5 has to be completed by cycle 15 or else a real-time violation occurs as shown in Figure 2. The bi-directional line represents first request which has to finish before cycle 15. The uni-directional line represents a second request which was continuously being issued. However, the permission is granted only after the completion of the first request at cycle 12. For RT_Lottery algorithm implemented by Chen et al. [7], the difference between the required bandwidth of master and the maximum allotted bandwidth was less if compared to D_Type master as shown in Table 1. Moreover, this master was also used by Akhtar and Sidek [9] to implement IAA. For

IAA, it was observed that DR_Type master acquired moderate and fair bandwidth and it showed high degree of synchronization with its associated masters [9].

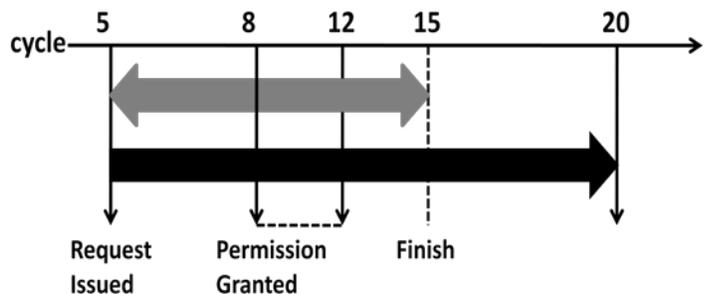


Fig 2. DR_Type master

Third type of master is NDR_Type. Request time for NDR_Type of master does not depend on the finish time of its previous request. The interval time is the clock cycles between two successive requests [7,9]. In Figure 3, the time interval assumed is 15 cycles. At cycle 21 second request is granted permission and is executed which is represented by uni-directional line. This request directly depends on cycle 7 of the first request but not its completion time at cycle 14 which is represented by bi-directional line. In this case the real-time parameter R is supposed to be smaller than minimum possible interval time because the current request must be

finished before the issue of next request. NDR_Type master worked well for RT_Lottery algorithm developed by Chen et al. [7], as the required bandwidth of the master was equal to the maximum allotted bandwidth. Moreover, for IAA developed by Akhtar and Sidek [9], this master acquired optimum bandwidth. However for IAA, NDR_Type master showed tight time constraints due to the synchronization limitation [9]. Hence, an improvement over IAA is deemed necessary.

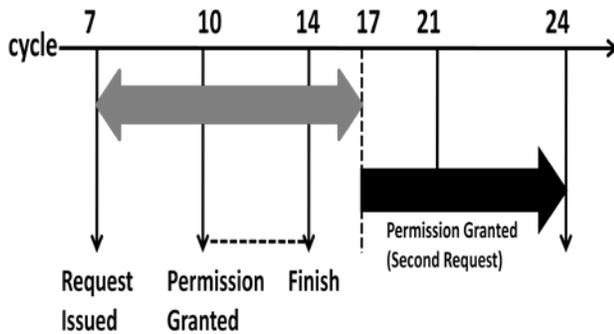


Fig 3. NDR_Type master

IV. PROPOSED PARALLEL ADAPTIVE ARBITRATION (PAA)

The proposed PAA has been designed in a manner that it can synchronize the implementation of the aforesaid masters to maximize the CPU cores utilization along with moderate bandwidth allocation for different masters according to their real time bandwidth requirements. Whenever a thread worker grows as $O(n^2)$, the communication cost only grows as $O(n)$, where O stands for order function [9]. If C is the communication time, then the total time (T) can be expressed as [9]:

$$T = \frac{P}{n} + C \quad (2)$$

Where P stands for parallel component and n stands for number of threads.

For the above equation, runtime complexity of P approaches $O(n^2)$ while C grows as $O(n)$ [10].

Therefore, the new equation for the total time (T) could be expressed as [10]:

$$T = \frac{t_p n^2}{n_t} + t_c n \quad (3)$$

Where t_p is constant processing time and t_c is constant communication time.

Figure 4 shows the connectivity of the designed heterogeneous masters to the arbiter. The benchmark application which is used to test these masters is known as "STREAM". STREAM (sustainable memory bandwidth in high performance computers) is a simple synthetic benchmark which is designed to measure the sustainable memory bandwidth and their corresponding computation rate [11].

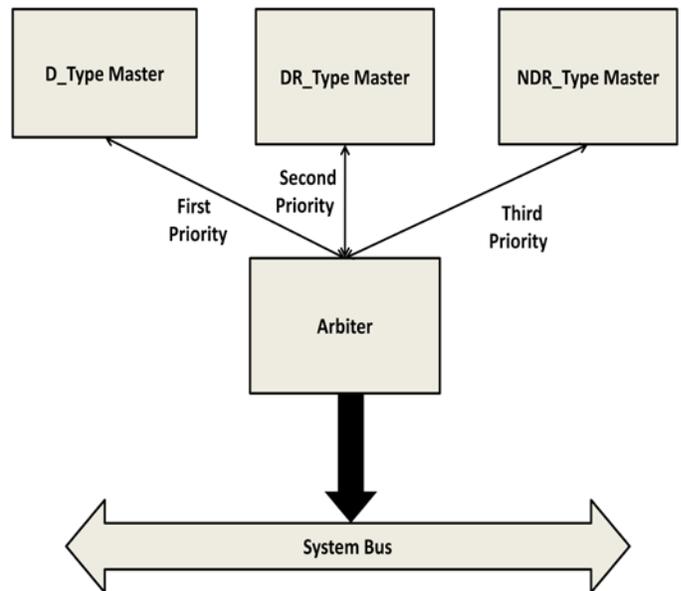


Fig 4. Proposed connectivity of masters to the arbiter

The following are the four STREAM modules which add independent information to the results:

- (i) "stream_copy" measures transfer rates without arithmetic operations.
- (ii) "stream_scale" includes a simple arithmetic operation. (A term scalar is defined in STREAM which is assigned a specific value)

(iii) "stream_sum" includes another operand to allow multiple load/store ports on vector machines to be tested and also.

(iv) "stream_triad" includes chained/overlapped/fused multiply/add operations .

In order to analyze the impact of PAA in terms of performance on masters designed according to their traffic behavior, an experiment have been conducted using OpenMP and SystemC whose libraries were ported on visual studio integrated development environment (VSIIDE) 2010 (Express Edition) platform where AMD Athlon™ II X2 260 processor with 1 MB dedicated L2 cache was used to measure the CPU cores utilization and bandwidth. Moreover, comparison of PAA is done with IAA which will be discussed in the results section. Another experiment was conducted using OpenMP and SystemC whose libraries were ported on Amazon EC2 cloud setup composed of suitable profiler tools to measure the CPU usage and bandwidth allotment. . In order to exploit the multiple cores of the processor, these masters have been synchronized to run parallel implementation of STREAM tasks [11].

Figure 5 shows the pseudo-code for the arbiter which is used in the proposed PAA technique. The arbiter gives first priority to the D_Type master as it is considered to be the blocking master where as DR_Type and NDR_Type has been given second and third priority by keeping the synchronization factor into consideration.

```

Arbiter (){
set priority:
1st priority: D_Type master
2nd priority: DR_Type master
3rd priority: NDR_Type master
if request == D_Type master {arbiter grnt permission with high
priority && D_Type master synchronizes with DR_Type master}
else
if request == DR_Type master {arbiter grnt permission with 2nd
priority && DR_Type master synchronizes with NDR_Type
master} else
if request == NDR_Type master {arbiter grnt permission with 3rd
priority && NDR_Type master synchronizes with DR_Type
master}
}

```

Fig 5. Arbiter pseudo-code

The implementation of the aforesaid masters on proposed PAA technique is discussed as follows:

The D_Type master is considered to be the highest priority master responsible for initiating the task application. For D_Type master, the initiation of next request depends on the finish time of the current request as this master does not have extra real-time requirement. This master is designed according to the typical time constrained thread process, as any method is invoked, it executes completely until it returns a value. As soon as D_Type master finishes its execution, the arbiter excites DR_Type master to start its execution.

As mentioned, the DR_Type master is same as D_Type master except that they have an extra real time requirement. Data sharing is enabled within the parallel region of the code so as to implement DR_Type master in the most appropriate manner. This master deals with two real time parameter to satisfy the extra real time requirement; dependent and independent. The independent real time parameter has to finish its execution within a specific count value where as the dependent real time value has to synchronize with NDR_Type master and has to finish its execution until the next request is issued by NDR_Type master.

The initiation of the NDR_Type master is independent of the finish time of the DR_Type master. However, request to initiate NDR_Type master is issued by DR_Type master to implement its dependent real-time parameter. The function of NDR_Type master is to implement stream functions with DR_Type master by synchronizing with each other. This implementation is done using the data sharing region of the OpenMP to attain high degree of parallelization using `omp_parallel` function.

The level of synchronization between D_Type master and DR_Type master is sequential as D_Type master initiates the implementation of application with high priority. Therefore there is a data dependency constraint between D_Type master and NDR_Type master. Figure 6 shows the synchronization between masters.

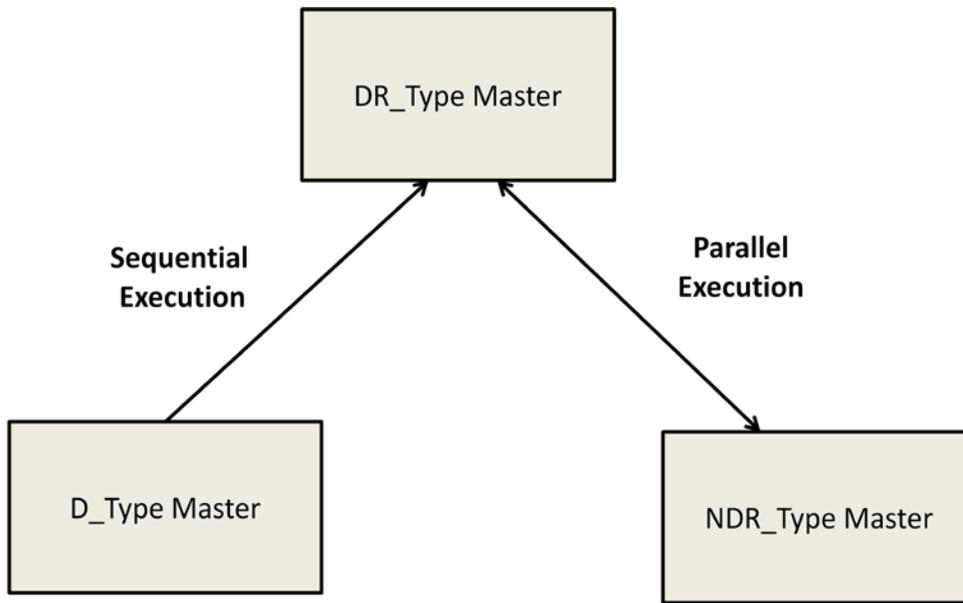


Fig 6. Synchronization between proposed arbitration masters

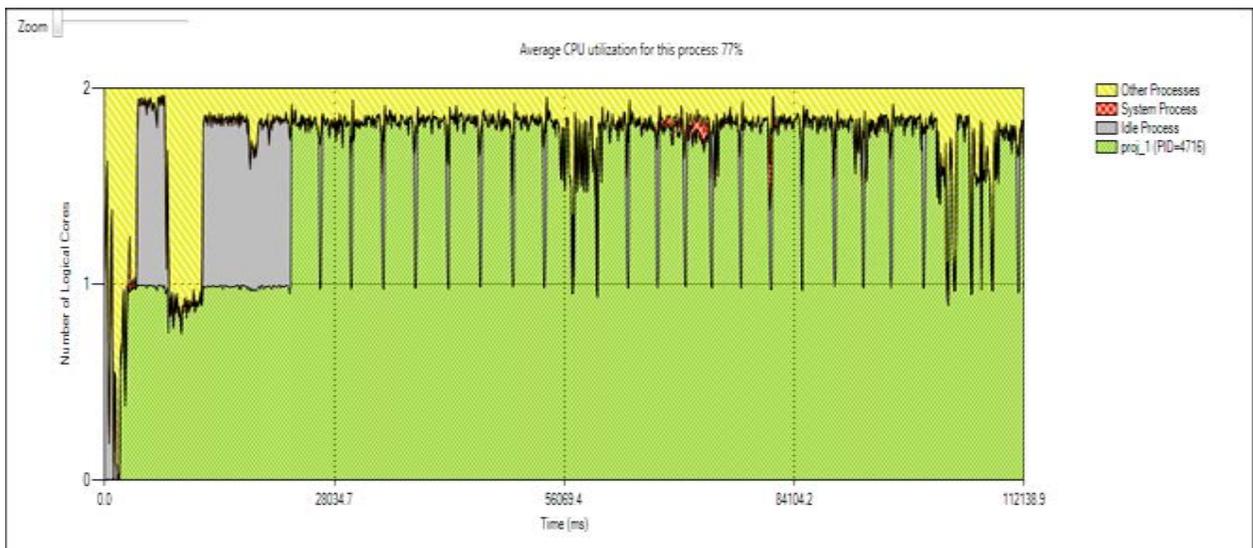


Fig 7. CPU cores utilization rate for PAA on VSIDE

V. RESULTS AND DISCUSSION

The essence of high performance computing lies within the concept of parallel programming. Figure 7 shows the rate of CPU cores utilization for the threads running multiple modules of STREAM using PAA in VSIDE using OpenMP and SystemC

libraries. Thread safety is one of the major criteria of multi-threading support for these masters. This means that communication in a multi-threaded application can be performed in multiple threads. Appropriate techniques should be used to utilize the multiple cores in order to make non-blocking communication primitives to progress in the background. Various thread-scheduling policies try to achieve optimal utilization of the CPU as well as the bus bandwidth during each quantum. The average CPU core utilization observed for PAA implemented on VSIDE is 77%. This shows high degree of multiple thread synchronization.

The graph in Figure 8 gives the comparison of CPU cores utilization for PAA implemented on Amazon EC2 instance using OpenMP default thread mode and IAA implemented using VSIDE. The graph shows the rate of CPU cores utilization for the threads running multiple modules of STREAM using IAA and PAA. It can be observed that PAA outplays IAA. The average CPU core utilization for PAA is 84.26% whereas for IAA is 74% which is also shown in Figure 9 [9]. Moreover in terms of latency, PAA implemented on Amazon EC2 is superior to IAA implemented on VSIDE. In addition, due to the constraints in the synchronization of masters, IAA is incompatible to be implemented on Amazon EC2 platform [9].

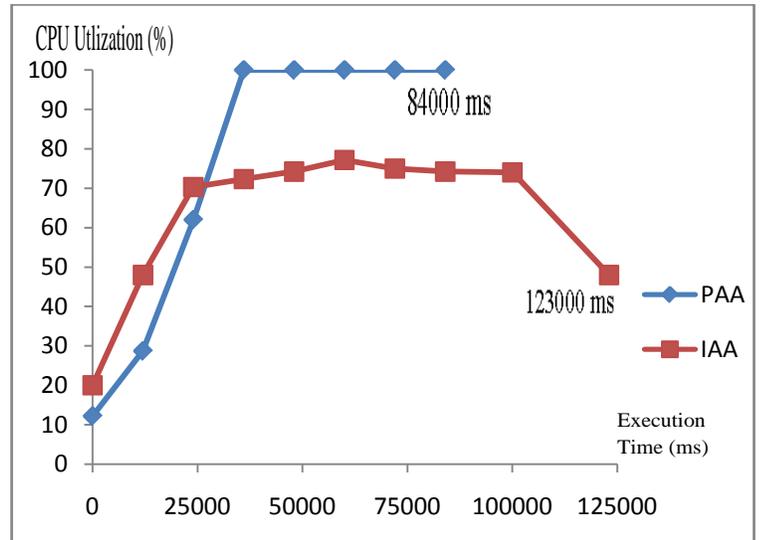


Fig 8. CPU cores utilization rate for PAA and IAA

Figure 9 shows an average CPU utilization for the threads running multiple modules of STREAM using the arbitration technique proposed by Akhtar and Sidek [9]. In this case all three masters i.e. D_Type, DR_Type and NDR_Type are synchronized with each other so as to enable parallel implementation. These masters implement the four functions of stream (stream_copy, stream_scale, stream_sum and stream_triad), where stream_copy is implemented by D_Type master, stream_scale is implemented by DR_Type master and stream_sum along with stream_triad is implemented together by NDR_Type master. If Figure 7 is compared with Figure 9, then it can be observed that PAA implemented on VSIDE is slightly better than IAA as the average CPU core utilization for PAA is 77% whereas for IAA is 74%. Moreover, in terms of latency, PAA is superior to IAA as the execution time is less for PAA as compared to IAA implemented on VSIDE platform.

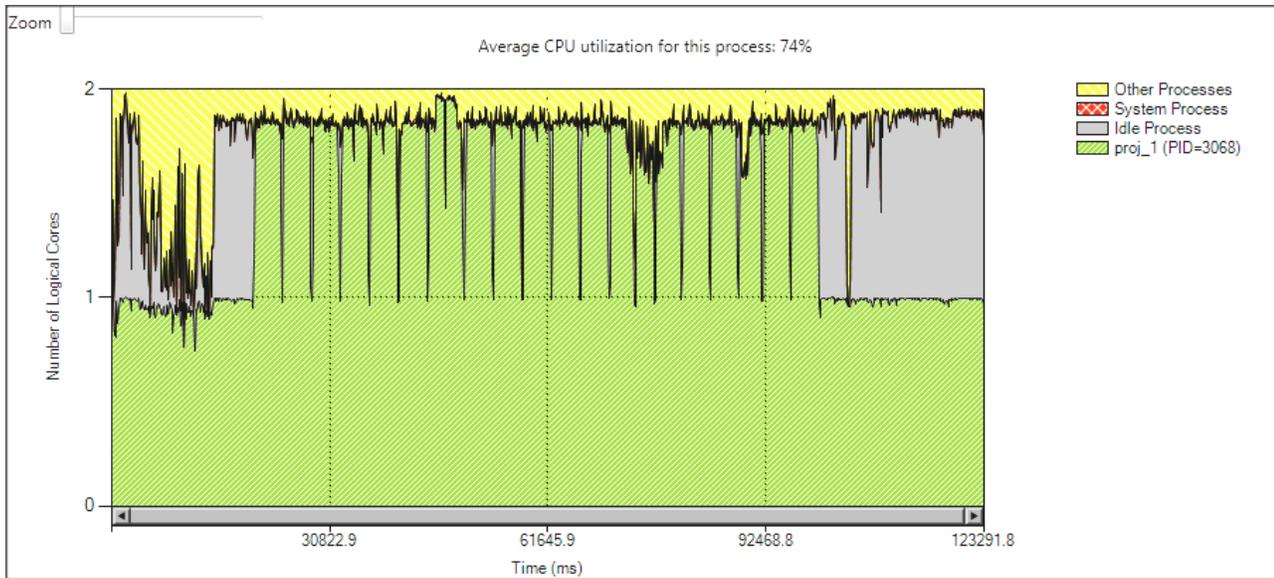


Fig 9. CPU cores utilization rate for IAA

In order to analyze the proposed PAA technique, three major benchmarks have to be dealt simultaneously; CPU utilization, bus bandwidth consumption and system latency. Processor architects have to trade-off the speed versus a lot of the features a processor may offer. In order to maximize the instruction per clock cycle, the instruction, operands and destination must be accessible at the same time parallelly rather than sequentially. In order to analyze the bandwidth fluctuation for heterogeneous masters, Li et al. [8] presented an experimental study for four masters to analyze moderate bandwidth allocation using adaptive arbitration. The set of masters taken were M1, M2, M3 and M4 and all four masters required different bandwidth in the proportion of 40%, 30%, 20% and 10% respectively. In Figure 10, the graph shows the distribution in the bus bandwidth for the different arbitration algorithms when all the four masters requests for variable bus bandwidth. It is observed from the graph that in the case of adaptive arbitration algorithm the difference between the requested bandwidth and the allotted bandwidth is $\pm 3.5\%$ whereas in the case of round robin algorithm its $\pm 4.75\%$, in case of IAA its $\pm 3.4\%$, whereas for lottery bus algorithm its $\pm 3.7\%$ and the minimum is in the case of static fixed priority algorithm that is $\pm 2.15\%$.

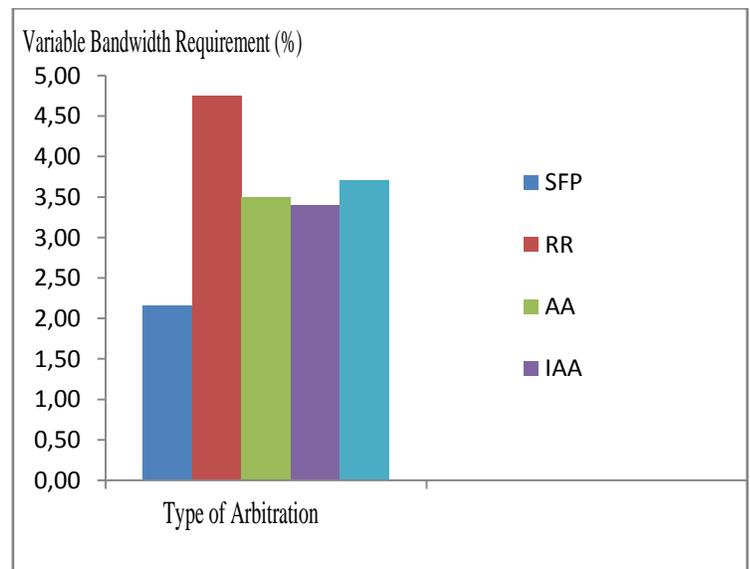


Fig 10. Bandwidth distribution for different arbitration algorithm

For the proposed PAA technique, the masters are designed in such a way that they require bandwidth according to the data traffic for the four functions of stream. Stream benchmark computes bandwidth for each master with its inbuilt functions using the following equation [11]:

$$\text{Allotted bandwidth} = \frac{\text{Required busbandwidth}}{\text{Total bus bandwidth}} \quad (4)$$

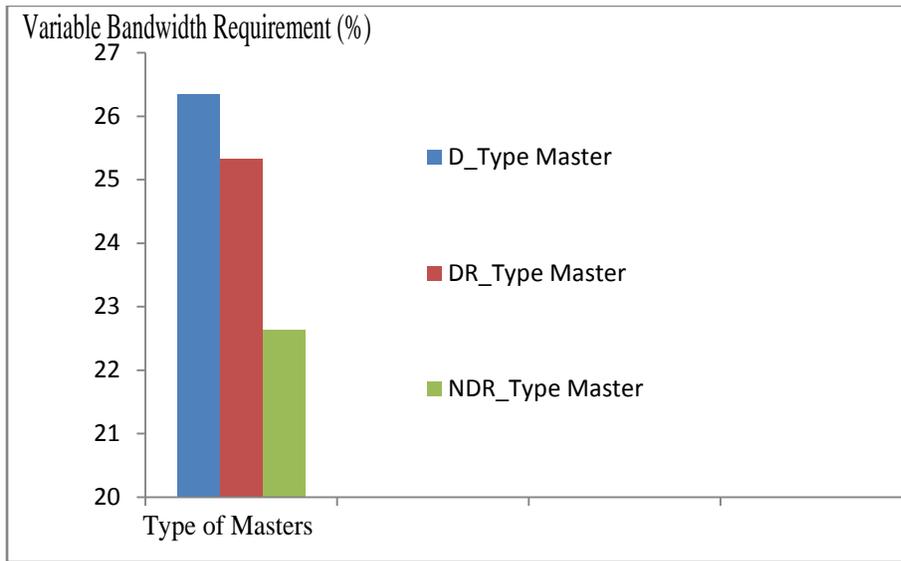


Fig 11. Bandwidth distribution using PAA on VSIDE

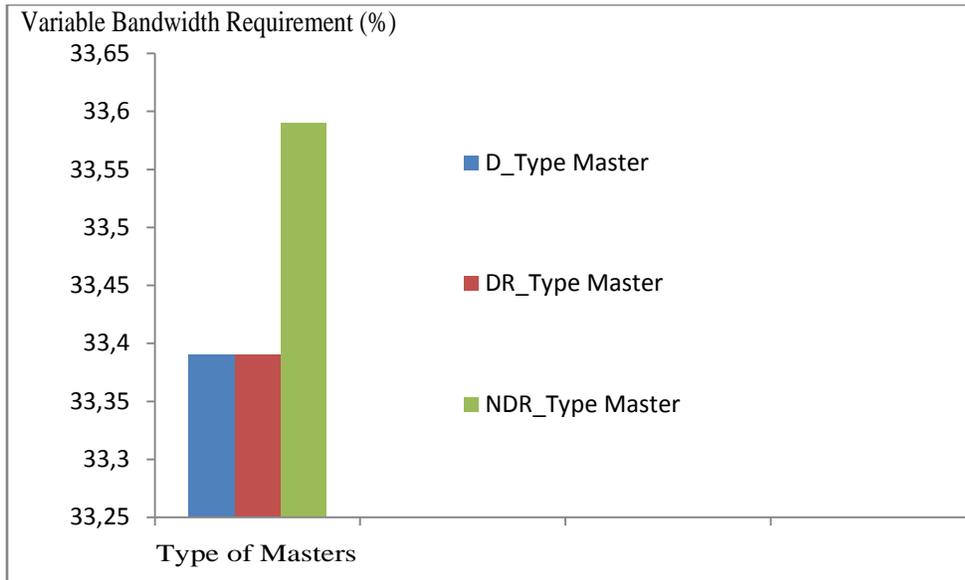


Fig 12. Bandwidth distribution using PAA on Amazon EC2 instance

The graph in Figure 11 shows the allotted bandwidth values for the masters implemented using PAA technique on the platform comprising of OpenMP and SystemC whose libraries were ported in VSIDE.

The graph in Figure 12 shows the allotted bandwidth values for the masters implemented using PAA technique on the platform comprising of OpenMP and SystemC whose libraries were ported in Amazon EC2 instance.

The fluctuation in the bandwidth recorded for PAA implemented on VSIDE platform was $\pm 1.35\%$. However, fluctuation in the bandwidth recorded for PAA implement on Amazon ec2 instance was just $\pm 0.20\%$. This clearly implies involvement of multiple numbers of threads to create high degree of synchronization between these three heterogeneous masters to implement the task in parallel where the bandwidth requirement may be variable.

Compared with other arbitration techniques in terms of bandwidth allocation, PAA is relatively better. After analyzing the above graphs it can be said that the PAA is good for task parallelization as it maximizes the CPU cores utilization with fair and moderate bandwidth allocation.

VI. CONCLUSION

An attempt has been made on maximizing the CPU utilization cores to achieve high degree of task parallelization by keeping moderate bandwidth allocation factor using the proposed PAA. Compared with the IAA, PAA has higher average CPU utilization rate of 84.26% compared to IAA with just 74%. In terms of bandwidth fluctuation, PAA shows minimalistic fluctuation of just $\pm 1.35\%$ on VSIDE and $\pm 0.20\%$ on Amazon ec2 instance. PAA could be used in high performance applications where ever task parallelization is essential. In a nutshell, PAA technique helps in achieving a fair bandwidth allocation for the critical requirements by utilizing the CPU cores to its maximum to achieve high degree of task parallelization. It also reduces the system latency up to an adequate margin through synchronization of masters. Nonetheless, PAA technique has to be tested on various intense image

processing application for enhanced performance evaluation.

ACKNOWLEDGEMENT

This research is supported by the School of Electrical and Electronic Engineering of Universiti Sains Malaysia. The authors would like to acknowledge the Institute of Postgraduate Studies (IPS), Universiti Sains Malaysia for the Global Fellowship [USM.IPS/USMGF(06/14)] financial support to carry out this research.

REFERENCES

- [1] S. Pasricha, and N. Dutt, *Communication Architectures-System On-Chip Interconnect*. USA: Morgann Kaufmann, 2008.
- [2] F. Poletti, D. Bertozzi, L. Benini, and A. Bogliolo, "Performance Analysis of Arbitration Policies for SoC Communication Architectures," *Design Automation for Embedded Systems*, vol. 8, pp. 189-210, 2003/06/01 2003.
- [3] X. Yi, L. Li, G. Ming-Lun, Z. Bing, J. Zhao-Yu, D. Gao-Ming, "An Adaptive Dynamic Arbiter for Multi-Processor SoC," in *Solid-State and Integrated Circuit Technology, 2006. ICSICT '06. 8th International Conference on*, 2006, pp. 1993-1996.
- [4] A. A. Aravind, "An arbitration algorithm for multiport memory systems," *IEICE Electronics Express*, vol. 2, pp. 488-494, 2005.
- [5] E. Massimo. and M. Poncino, "The design of easily scalable bus arbiters with different dynamic priority assignment schemes," presented at the IEEE-Signals, Systems and Computers, Italy, 1995.
- [6] J. R. López-Blanco, R. Reyes, J. I. Aliaga, R. M. Badia, P. Chacón, and E. S. Quintana-Ortí, "Exploring large macromolecular functional motions on clusters of multicore processors," *Journal of Computational Physics*, vol. 246, pp. 275-288, 8/1/ 2013.
- [7] C. H. Chen, L. Geeng-Wei, H. Juinn-Dar, and J. Jing-Yang, "A real-time and bandwidth guaranteed arbitration algorithm for SoC bus communication," in *Design Automation, 2006. Asia and South Pacific Conference on*, 2006, p. 6 pp.
- [8] H. Li, Z. Ming, Z. Wei, and L. Dongxiao, "An Adaptive Arbitration Algorithm for SoC Bus," in *Networking, Architecture, and Storage, 2007. NAS 2007. International Conference on*, 2007, pp. 245-246.
- [9] M. N. Akhtar, O. Sidek, "An Intelligent Adaptive Arbiter for Maximum CPU Utilization, Fair Bandwidth Allocation and Low Latency," *IETE Journal of Research*, vol. 59, pp. 48-52, 2013.
- [10] S. J. Wong. and A. P. Rendell, "The SCore cluster enabled OpenMP environment: performance prospects for computational science," presented at the Proceedings of the 5th international conference on Computational Science - Volume Part I, Atlanta, GA, 2005.

- [11] J. D. McCalpin, "Memory bandwidth and machine balance in current high performance computers," *IEEE Computer Society Technical Committee on Computer Architecture, (TCCA)* pp. 19-25, 1995.

About the authors

M. Nishat Akhtar received his B.E in Computer Science from VTU, Karnataka during the year 2010 and MS in Electrical and Electronics from Collaborative Microelectronic Design and Excellence Centre, Universiti Sains Malaysia during the year 2013. Currently he is a PhD student in School of Electrical and Electronics at Universiti Sains Malaysia. His research interests include High Performance Computation, Parallel Computing, embedded systems and System-on-Chip.

Junita Mohamad-Saleh received her B.Sc (in Computer Engineering) degree from the Case Western Reserve University, USA in 1994, the M.Sc. degree from the University of Sheffield, UK in 1996 and the Ph.D. degree from the University of Leeds, UK in 2002. She is currently an Associate Professor in the School of Electrical & Electronic Engineering, Universiti Sains Malaysia. Her research interests include computational intelligence, tomographic imaging and parallel processing.

Predictive robots programming based on imitation strategy

Prof. Eng. A. Fratu, PhD¹, Lecturer Eng. Mariana Fratu, PhD¹

¹University „Transilvania” of Braşov, Brasov, Romania

Abstract: In this paper, based on original idea, the authors propose a new strategy for physical robot programming using predictive control strategy. To program the desired motion sequence for the physical robot, one captures the virtual reference paths from the virtual robot model and maps these to the joint settings of the physical robot. The control system of the physical robot reproduces the virtual reference paths of the virtual prototype. This requires transfer of a dynamical signature of a movement of the virtual robot to the physical robot, i.e. the robots should be able to imitate a particular path as one with a specific velocity and/or an acceleration profile. Moreover, the virtual robot must cover all possible contexts in which the physical robot will need to generate similar motions in unseen context. The physical robot acts automatically, communicating with corresponding virtual prototype and imitating its behavior.

Keywords: Virtual robots, Learning by imitation, Predictive programming, Behavior control

I. INTRODUCTION

An imitation task can be decomposed into the serial implementation of two processes: an *observation process* and an *imitation process*. The observation process consists of extracting relevant (i.e. in our case, virtual robot prototype) features from a *demonstrated dataset*. The imitation process consists of generating an *imitated dataset* that minimizes the discrepancy between the demonstrated and imitated datasets.

In our prior works, we have addressed *what the imitation* question by developing a general architecture to extract the relevant features of a given task. The method relied on computing the virtual trajectories of each joint of the manipulator.

In this paper we present an extension of this works to address the *how to imitate* problem for the control of manipulator motion, using a virtual manipulator prototype. This leads us to tackling more generic issues of servomotor control, namely that of optimizing the arm controller given specific constraints.

Specifically, we extend our original method for to determine the optimal imitation strategy, i.e. the strategy that satisfies best all the constraints of a given task. The issue of how to imitate also referred to as the transfer of the virtual trajectories to the corresponding physical

joints. The solution to the corresponding transfer problem in the latter works was, however, constrained to a particular arm and did not provide a general solution for robotic arms with an arbitrary number of degrees of freedom

The concept of robot motion prediction was introduced to clearly understand what the robot must to do when trying to localize visual objects. His suggestion was that the physical robot can predict the situation (position and orientation) using her virtual prototype rather than physical sensory signals. For a system with a demand of reacting as precisely as possible, its past information is not suitable for control planning any more.

We should predict the future behavior, at the time when the control command arrives at the physical robot and is executed.

The ability of predicting of the behavior of robots is important in design; the designers want to know whether the robot will be able to perform a typical task in a given time frame into a space with constraints [2].

The control engineer cannot risk a valuable piece of equipment by exposing it to untested control strategies. Therefore, a facile strategy for contact detection and collision avoidance, capable of predicting the behavior of robotic manipulators, becomes imperative.

When the robots need to interact with their surrounding, it is important that the computer can simulate the interactions of the participants, with the passive or active changing environment in the graphics field, using virtual prototyping. The actions for the each task are computed for virtual robot and are transferred, with a central coordination to corresponding physical robot which must imitate her virtual homonym.

In this paper we develop a formally analyze of the concept of robot motion prediction using the behavior of the virtual robot, to control the physical robot in the real world.

II. PHYSICAL ROBOT MOTION CONTROL BASED ON SIMULATED VIRTUAL MODEL

We present a description of the theoretical aspects of the physical robot motion control using a virtual motion prediction model. The advantages of such approach as an alternative to the classical methods (e.g. vision guided trajectory imitation [12]) are on-line adaptation to the motion of the virtual prototype.

A solution to the above problem is to construct a virtual

prototype model and transfer the virtual trajectory by interacting with the physical model.

Designing a virtual model would be an option; however, the behavior of the robots is very difficult to model. Moreover, the use of system knowledge is contrary to our research aim. Therefore we focus on creating a virtual prototype model from experimental data obtained from the physical robot model [6].

Users interact with the simulation environment through the visualization. This includes, but not limited to, computer screen. Optimization of the real robots behavior is performed in the low dimensional virtual space using the virtual robot prototypes.

In the virtual space one simulate even the intersecting of the virtual robot and her environment. The intersecting of two virtual objects is possible in the virtual world, where the virtual objects can be even intersected and there no exist the risk to be destroyed [4]. The visualization provides an interface to develop interactive implementations based on simulated behavior of the virtual prototype.

In our work, we assume that learning of the deterministic part for description motion dynamics should be sufficient to design the corresponding robot control.

We particularly refer to the ability of the system to react to changes in the environment that are reflected by motion parameters, such as a desired target position and motion duration. Therefore, the system is able to manage with uncertainties in the position of a manipulated object, duration of motion, and structure limitation (e.g., joint velocity and torque limits) [3].

The proposed method aims at adapting to spatial and temporal perturbations which are externally-generated. This aspect will be investigated in our future works.

It is easy to recuperate kinematic information from virtual robot motion, using for example motion capture [1]. Imitating the motion with stable robot dynamics is a challenging research problem [9].

In this paper, we propose a predictive control structure for physical robots that uses capture data from their virtual prototypes and imitate them to track the motion in the real space.

We will demonstrate the tracking ability of the proposed controller with dynamics simulation that takes into account joint velocity and torque limits. We apply the controller to tracking motion capture clip to preserve the original behavior of virtual robot.

First, a motion capture system transforms Cartesian position of virtual robot structure to virtual joint angles based on kinematic model [7]. Then, the joint angles are converted in binary words and transferred to real robot. We employ the control loops structure to establish relationships between the virtual and real robot control systems.

We present results demonstrating that the proposed approach allows a real robot to learn move based exclusively on virtual robot motion capture, viewed as predictive control strategy.

III. ONLINE BEHAVIOR IMITATION FOR PREDICTIVE CONTROL

In robotics, one of the most frequent methods to represent movement strategy is by means of the learning from imitation. Imitation learning is simply an application of supervised learning [8]. One goal of imitation of the dynamical systems is to use the ability of coupling phenomena to description for complex behavior [10]. Anything a robot does is called a behavior. Moving, turning, stopping, picking things up, putting them down, delivering a message are all behaviors that a robot can perform. In this paper, we propose a generic modeling approach to generate virtual robot prototype behavior in experimental scenery. The actions for the each task are computed for virtual robot prototype and are transferred online, with a central coordination, to corresponding physical robot, which must imitate her virtual "homonymous". Notice the similarity between moves of the virtual robot prototype in the virtual work space and the "homonymous" moves in the real work space of the physical robot. We assume to use the virtual robot prototypes and the motion capture systems to obtain the reference motion data, which typically consist of a set of trajectories in the Cartesian space.

The paper relates to a method for robot programming by combining off-line and on-line programming techniques. The method consists in using a programming platform on which there is carried out the virtual prototype of the physical robotic arm to be programmed and the real working space wherein it is intended to work.

In the robot program there is written a source code intended to generate the motion paths of the virtual robotic arm prototype [16]. The numerical values of the prototype joints variables are sent to the data register of a port of the information system which, via a numerical interface, are on-line transferred into the data registers of the controllers of the servo system of the physical robotic arm. Finally, there are obtained tracking structures due to which the moving paths of the virtual robotic arm joints are reproduced by the physical robotic arm joints, thereby generating motion within the real working space [13].

IV. PREDICTIVE CONTROL MODEL

A virtual trajectory over a virtual environment is able to predict the evolution of the physical robot in the physical environment under any selected constraints. Once an virtual trajectory has been computed; this is used as predictive model. In our experiments various synthetic behavior models, close to real robot behavior, have been created. However, the use of realistic models for interaction between the virtual world and physical world requires significant computation time and large amounts of data transfer, and most currently existing virtual prototypes have a limited number of predefined expressions.

In this paper we propose a structure that could contribute to answering these three questions and enhance cooperative interaction within virtual worlds; the Virtual world Maker to create realistic scene-referenced 3D environments, a 3D animation system to display tasks' evolution, and the Analyzer to automatically initiates the motion of physical robot.

The question about how to create believable environments is answered by the automatic creation of identifiable scenes.

The question about the creation of realistic virtual robot prototypes is answered by the development of an automated 3D modeling system for to display of the virtual prototype behavior.

The question of how to set in motion the physical robot of appropriate behaviors is answered by a system for automating transfer the virtual robot' poses to physical robot servo system.

A. Platform structure

Initially, a set of virtual postures is created for the virtual robot and the pictures' positions are recorded for each posture, during motion. These recorded pictures' positions provide a set of Cartesian points in the 3D capture volume for each posture.

To obtain the physical robot postures, the virtual pictures' positions are assigned as positional constraints on the physical robot. To obtain the physical joint angles one use standard inverse kinematics (IK) routines. The IK routine then directly generates the physical joint angles on the physical robot for each posture.

We start with a 3 degree-of-freedom (DOF) discrete movement system that models point-to-point attained in a 3D Cartesian space.

Figure 1 shows our experiment involving the imitation learning for a physical robotic arm with 3 degrees-of-freedom (DOFs) for performing the manipulate tasks [17]. We demonstrated the imitation of elbow, shoulder and wrist movements. Importantly, these tasks required the coordination of 3 DOFs, which was easily accomplished in our approach.

The imitated movement was represented in joint angles of the robot. Indeed, only kinematic variables are observable in imitation learning. The physical robot was equipped with a controller (a PD controller) that could accurately follow the kinematic strategy (i.e., generate the torques necessary to follow a particular joint angle trajectory, given in terms of desired positions, velocities, and accelerations) [7].

Figure 1 also displays (left image) the user interface of a virtual robotic manipulator arm, which has been created which a dynamical simulator.

Referring the Figure 1 we comment the following: on programming platform, a robot program is carried out off-line, and one sends into the data registers of a port of the hardware structure, the numerical values of the joint variables of the virtual prototype of the robotic arm (BRV) and displays on a graphical user interface, the

evolution of the virtual prototype during the carrying out of the robotic task. Via numerical interface (IN) the virtual joint dataset, from the data registers of the port of the hardware structure of the programming platform are transferred into the data registers of the numerical comparators of the controllers [9]. These datasets are reference inputs of the pursue loops, resulting a control system (SC).

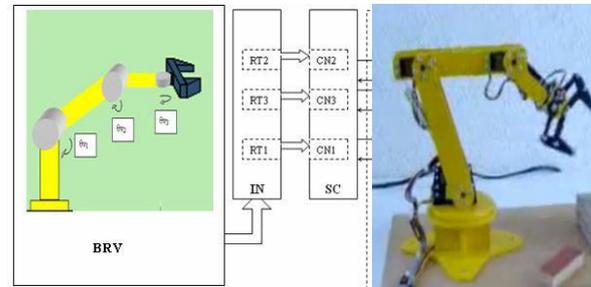


Fig. 1 Imitation software platform structure

The reference datasets are obtained using a motion capture channel taking into account the joints motion range.

The easiest way to generate the spatial relations explicitly is the interactively programming of the behavior of the virtual prototype in his virtual environment, in order to specify suitable positions θ_{v1} , θ_{v2} , θ_{v3} .

This kind of specification provides an easy to use interactive graphical tool to define any kind of robot path; the user has to deal only with a limited and manageable amount of spatial information in a very comfortable manner.

The applicable robot tasks are designed and the desired pathways are programmed off-line and stored in the buffer modules RT1, RT2, RT3.

The comparative modules CN1, CN2, CN3 furnish, to the pursuit controllers, the datasets involving the expected state of the virtual robot prototype and the measured state of the physical robot.

Our system requires an essential step in that one converts the position errors into motor commands by means of the PD controller.

We assume to use the virtual robot prototypes and the motion capture systems to obtain the reference motion data, which typically consist of a set of trajectories in the Cartesian space.

The data is obtained using a motion capture channel taking into account the joint motion range. Due to the joint limits and the difference between the kinematics of the virtual robot and real physical robot, the joint angle data are pre-processed.

In our pre-processing, we assume that both virtual and physical robots are on the scene at the same time and estimate the correct arms position and orientation. We then compute the inverse kinematics for new posture to obtain the cleaned joint angles and retain the difference from original joint angles.

At each frame during control, we add the difference to the original data to obtain the cleaned reference joint

angles. This correction is extremely simple and our controller does not require supplementary cleanup.

B. Basic servo system

We aim at developing controllers for learning by imitation with a virtual robot demonstrator. For this purpose, we assume a simple control system where the position and velocity of the 1 DOF discrete dynamical system, drives the time evolution of joint variable, which can be interpreted as the position controlled by a proportional-derivative controller.

The resolved acceleration controller [5] is picked up in a servo system. The computed torque control method is used for nonlinear control of robotic manipulator, which is composed of a model base portion and a servo portion. The servo portion is a close loop with respect to the position and velocity [11].

Closed-loop servo system are created out of one relatively simple set of equations; based only on the capture of the trajectory of a virtual robot prototype.

1) *Computed torque control*

The dynamic model of a robotic manipulator without friction term is generally given by

$$M(q)\ddot{q} + H(q, \dot{q}) + G(q) = \tau \tag{1}$$

where $M(q) \{6 \times 6\}$ is the inertia term in joint space. $H(q, \dot{q}) \{6 \times 1\}$ and $G(q) \{6 \times 1\}$ are the Coriolis/centrifugal term and gravity term in joint space, respectively. $q \{6 \times 1\}$, $\dot{q} \{6 \times 1\}$ and $\ddot{q} \{6 \times 1\}$ are the position, velocity and acceleration $\{6 \times 1\}$ dimensional vectors in joint coordinate system, respectively. The vector $\tau_i \{6 \times 1\}$ is the joint driving torque $\{6 \times 1\}$ dimensional vector. In the case that the resolved acceleration control law is employed in the servo system of a manipulator, desired position, velocity and acceleration vectors in Cartesian coordinate system are respectively given to the references of the servo system.

2) *Desired Trajectory*

For instance, in order to apply the computed torque control method to the manipulator, the desired trajectory composed of x_r , \dot{x}_r and \ddot{x}_r in Cartesian coordinate virtual system must be prepared. First of all, the desired trajectory in Cartesian coordinate system is designed, in which the manipulator moves from the initial pose to the final pose.

The desired trajectory x_r , \dot{x}_r and \ddot{x}_r are calculated from virtual joint angle θ_v , virtual joint velocity $\dot{\theta}_v$ and virtual joint acceleration $\ddot{\theta}_v$ respectively, using the robot analytical models.

The trajectory in virtual joint space makes the physical robot follows of the homonym virtual prototype as shown

in Fig. 1.

The joint driving torque is calculated from

$$\tau = \hat{M}(q)J^{-1}(q) \times [\ddot{x}_r + k_v(\dot{x}_r - \dot{x}) + k_p(x_r - x)] - \hat{J}(q)\dot{q} + \hat{H}(q, \dot{q}) + \hat{G}(q) \tag{2}$$

where, the symbol $\hat{}$ denotes the modeled term. The $x \{6 \times 1\}$, $\dot{x} \{6 \times 1\}$ and $\ddot{x} \{6 \times 1\}$ are the position / orientation, velocity and acceleration $\{6 \times 1\}$ dimensional vectors in Cartesian coordinate system, respectively. The $x_r \{6 \times 1\}$, $\dot{x}_r \{6 \times 1\}$ and $\ddot{x}_r \{6 \times 1\}$ are the desired position / orientation, velocity and acceleration $\{6 \times 1\}$ dimensional vectors, respectively. The diagonal matrix $K_v = \text{diag}(k_{v1}, \dots, k_{v6})$ and $K_p = \text{diag}(k_{p1}, \dots, k_{p6})$ are the feedback gains of velocity and position, respectively.

The matrix $J(q)$ is the Jacobian matrix which gives the relation $\dot{x} = J\dot{q}$. Note that q , \dot{q} , x and \dot{x} in (2) are actual values, i.e., controlled variables.

Figure 2 shows the block diagram of the computed torque control method, in which $F_{kine}(\cdot)$ is the function to obtain the forward kinematics.

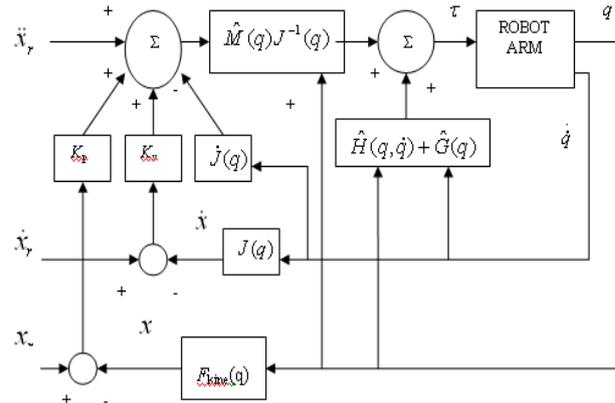


Fig. 2 Block diagram of the resolved acceleration control method, with desired position, velocity and acceleration vectors in Cartesian coordinate system.

The nonlinear compensation terms $\hat{H}(q, \dot{q})$ and $\hat{G}(q)$ are calculated to cancel the nonlinearity and are effective to achieve a stable trajectory control [15].

The tracking controller is responsible for making every joint track the desired trajectory. It solves an optimization problem that respects both joint tracking and desired inputs to the simplified model and obtains the joint torques to be commanded to the real robot.

To obtain satisfactory and safe control performance without falling a singularity, K_v and K_p are approximately tuned in advance with trial and error, considering the combination around critically damped condition. We call this the initial manual tuning process.

Two search ranges for K_v and K_p are obtained after the manual tuning process, so that K_v and K_p must be further tuned finely within the searched spaces to achieve a desirable motion without large overshoots and

oscillations.

In the next section, we propose impedance following force control model which can be applied after the manual tuning process.

V. IMPEDANCE MODEL TO FORCE CONTROL OF THE PHYSICAL ROBOT

When a computer is used to control a robotic manipulator, the control law is generally represented by a discrete-time control system. In this section, it is described on how to evaluate the velocity-based discrete-time control system which is implemented as basic strategy for predictive motion control.

Let's consider the impedance model following force control as an example of velocity-based discrete-time control systems. In order to conduct a simulation with a robotic dynamic model, manipulated variables written by velocity commands in discrete-time domain must be transformed into joint driving torques.

Impedance control is one of the effective control strategies for a robotic manipulator to desirably reduce or absorb the external force from an environment [8]. It is characterized by ability which controls the mechanical impedance such as mass, damping and stiffness acting at joints.

Impedance control does not have a force control mode or a position control mode but it is a combination of force and velocity.

In order to control the contact force acting between the arm tip and environment, we have proposed the impedance model following force control methods that can be easily implemented in robotic manipulators with an open architecture controller [9].

The desired impedance equation in Cartesian space for a robotic manipulator is designed by:

$$M_d(\ddot{x} - \ddot{x}_d) + B_d(\dot{x} - \dot{x}_d) + SK_d(x - x_d) = SF + (I - S)K_f(F - F_d) \quad (3)$$

where $x \{6 \times 1\}$, $\dot{x} \{6 \times 1\}$ and $\ddot{x} \{6 \times 1\}$ are the position, velocity and acceleration vectors, respectively. M_d , $B_d \{3 \times 3\}$ and $K_d \{3 \times 3\}$ are the coefficient matrices of desired mass, damping and stiffness, respectively. F is the force vector. $K_f \{3 \times 3\}$ is the force feedback gain matrix. $x_d \{6 \times 1\}$, $\dot{x}_d \{6 \times 1\}$, $\ddot{x}_d \{6 \times 1\}$ and F_d are the desired position, velocity, acceleration and force vectors, respectively.

The S is the switch matrix to select force control mode or compliance control mode. If $S = 0$, (3) becomes force control mode in all directions; whereas if $S = I$ it becomes compliance control mode in all directions. Here, matrix I is the identity matrix. M_d

B_d , K_d and K_f are set to positive-definite diagonal matrices.

When force control mode is selected in all directions,

i.e., $S = 0$, defining $X = (\dot{x} - \dot{x}_d)$ gives:

$$\dot{X} = -M_d^{-1}B_d X + M_d^{-1}K_f(F - F_d) \quad (4)$$

In general, (4) can be resolved as:

$$X = e^{-M_d^{-1}B_d t} X(0) + \int_0^t e^{-M_d^{-1}B_d(t-\tau)} M_d^{-1}K_f(F - F_d) d\tau \quad (5)$$

In the following, we consider the form in the discrete time k using a sampling time Δt . If it is assumed that M_d , B_d , K_f , F and F_d are constant at $\Delta t(k - 1) \leq t < \Delta t k$, then defining $X(k) = X(t)/t = \Delta t k$ leads to the recursive equation.

Remembering $X = (\dot{x} - \dot{x}_d)$ giving $\dot{x}_d = 0$ in the direction of force control, and adding an integral action, the equation of velocity command in terms of Cartesian space is derived by:

$$\begin{aligned} \dot{x}(k) = & e^{-M_d^{-1}B_d \Delta t} \dot{x}(k - 1) - \\ & (e^{-M_d^{-1}B_d \Delta t} - I)M_d^{-1}K_f \{F(k) - F_d\} + \\ & K_i \sum_{n=1}^k \{F(n) - F_d\} \end{aligned} \quad (6)$$

where $K_i \{3 \times 3\}$ is the integral gain matrix and is also set to a positive-definite diagonal matrix. The impedance model following force control method is used to control the force which an robotic manipulator gives an environment. As can be seen, the force is regulated by a feedback control loop.

CONCLUSION

This paper explores the robot control based on predictive control model. Prediction in robot behavior control has become an increasingly important subject. There are two types of commands, one is the commands generated by the trajectory planner, and the other is generated by the actually virtual trajectory, which are used for prediction behavior of the physical robot.

The commands sent to our robot are logged for each cycle. So if the robots always executed just as what the command tell them to do, the consequences could be predicted easily. When the physical robot executes commands, its actual position and orientation should be the consequences of the information from a virtual path, transferred from virtual homonym prototype.

The problem of the physical robot behavior control is better analyzed on the virtual prototypes in the virtual environment where one may predict there behavior.

Learning approach such as learning by imitation is more flexible and can adapt to environmental change. This method is typically directly applicable due the

possibility to transfer the virtual joint trajectories from virtual space to the real space of the physical robots.

As we have known the relationship between the expectation and the actual consequence, we can modify the actual command sent to the robots, making the robots behaves just as what we expect. This new innovative method for behavior predictive control is attractive for implementation. Not similar to most other methods, our method not only makes a good prediction, but also improves the precision of motion control.

Unfortunately, because of the inherent mechanical constraints, there are always some variations between what we tell our robot to do and the result gotten from the execution. Using the homonym virtual prototype, it was also shown that predictive control model's computational requirement could be reduced. It is desired in future works to address disturbances and physical obstacles with this method.

ACKNOWLEDGMENT

The author wishes to thank, for cooperation and engagement in research activity the entire team of Services and Products for Intelligent Environment Laboratory, within the Research & Development Institute ICDT-PRO-DD of the Transylvania University of Brasov. We hereby acknowledge the structural funds project PRO-DD (POS-CCE, 0.2.2.1., ID 123, SMIS 2637, ctr. No 11/2009) for providing the infrastructure used in this work.

REFERENCES

- [1] V. Zordan and J. Hodgins, *Motion Capture-Driven Simulations that Hit and React*, In Proceedings of ACM SIGGRAPH Symposium on Computer Animation, San Antonio, TX, July 2002, pp. 89–96.
- [2] D. Silver, *Cooperative path finding*, In Proceedings of the 1st Conference on Artificial Intelligence and Interactive Digital Entertainment (AIIDE'05), pp. 23–28, 2005.
- [3] J. Pettre, J. Ondrej, A.-H. Olivier, A. Cretual and S. Donikian, *Experiment based modeling, simulation and validation of interactions between virtual walkers*, In Proceedings of Symposium on Computer Animation, ACM, 2009.
- [4] F. Nagata, Y. Kusumoto, Y. Fujimoto and K. Watanabe, *Robotic sanding system for new designed furniture with free formed surface*, In Robotics and Computer-Integrated Manufacturing Journal, vol. 23, no. 4, pp. 371–379, 2007.
- [5] F. Nagata, K. Kuribayashi, K. Kiguchi and K. Watanabe, *Simulation of fine gain tuning using genetic algorithms for model based robotic servo controllers*, In Proceedings of the IEEE Int. Symposium on Computational Intelligence in Robotics and Automation, pp.196–201, 2007.
- [6] K. Gold, *An information pipeline model of human-robot interaction*, In Proceedings of the 4th ACM/IEEE international conference on Human robot interaction, pp. 85- 92, New York, NY, USA, 2009. ACM. ISBN 978-1-60558-404-1. doi:http://doi.acm.org/10.1145/1514095.
- [7] A. Powers, S. Kiesler, S. Fussell, and C. Torrey, *Comparing a computer agent with a humanoid robot*, In Proceedings of the ACM/IEEE international conference on Human-robot interaction (HRI '07). ACM, New York, USA, pp. 145-152, 2007.
- [8] B. Price and C. Boutilier, *Accelerating reinforcement learning through implicit imitation*, Journal of Artificial Intelligence Research, vol. 19, 2003, pp. 569–629.
- [9] D.Grimes, R. Chalodhorn and R. Rao, *Dynamic imitation in a humanoid robot through non-parametric probabilistic inference*, In Proceedings Robotics: Science and Systems (RSS), Philadelphia, Pennsylvania, USA, August 2006.
- [10] F. Nagata, I. Okabayashi, M. Matsuno, et al., *Fine gain tuning for model-based robotic servo controllers using genetic algorithms*, In Proceedings of the 13th Int. Conf. on Advanced Robotics, 2007, pp. 987–992.
- [11] M. Nicolescu, M. Mataric, *Natural methods for robot task learning- Instructive demonstrations, generalization and practice*, In Proc. Intl Joint Conf. on Autonomous Agents and Multi-agent Systems (AAMAS), Melbourne, Australia, 2003, pp. 241-248.
- [12] S. Calinon, F. Guenter and A. Billard, *On learning, representing and generalizing a task in a humanoid robot*, In IEEE Trans. on Systems, Man and Cybernetics, Part B, Vol. 37, No. 2, 2007, pp. 286–298.
- [13] M. Pardowitz, R. Zoellner, S. Knoop, and R. Dillmann, *Incremental learning of tasks from user demonstrations, past experiences and vocal comments*, In IEEE Trans. on Systems, Man and Cybernetics, Part B, Vol. 37, No. 2, 2007, pp. 322–332.
- [14] J. Ijspeert, J. Nakanishi and S. Schaal, *Movement imitation with nonlinear dynamical systems in humanoid robots*, In Proceedings of the IEEE International Conference on Robotics and Automation, 2002, pp. 1398–1403.
- [15] J. Kober, J. Peters, *Learning motor primitives for robotics*, In Proceedings of the IEEE International Conference on Robotics and Automation, pp. 2112–2118, 2009, Piscataway, NJ: IEEE.
- [16] A. Fratu, *Method and installation for joints trajectory planning of a physical robot arm*, Proposal patent: A 00482 / 28. 06. 2013, available at: http://worldwide.espacenet.com/?locale=en_E

Unifying Geometric Features and Facial Action Units for Improved Performance of Facial Expression Analysis

Mehdi Ghayoumi¹, Arvind K Bansal¹

¹Computer Science Department, Kent State University,
{mghayoum,akbansal}@kent.edu

Keywords: Facial Action Unit, Facial Expression, Geometric features.

Abstract: Previous approaches to model and analyze facial expression analysis use three different techniques: facial action units, geometric features and graph based modelling. However, previous approaches have treated these technique separately. There is an interrelationship between these techniques. The facial expression analysis is significantly improved by utilizing these mappings between major geometric features involved in facial expressions and the subset of facial action units whose presence or absence are unique to a facial expression. This paper combines dimension reduction techniques and image classification with search space pruning achieved by this unique subset of facial action units to significantly prune the search space. The performance results on the publicly facial expression database shows an improvement in performance by 70% over time while maintaining the emotion recognition correctness.

1 INTRODUCTION

Your Emotion represents an internal state of human mind [28], and affects their interaction with the world. Emotion recognition has become an important research area in: 1) the entertainment industry to assess consumer response; 2) health care industry to interact with patients and elderly persons; and 3) the social robotics for effective human-robot interaction. Online facial emotion recognition or detection of emotion states from video has applications in video games, medicine, and affective computing [26]. It will also be useful in future in auto-industry and smart homes to provide right ambience and interaction with the occupying humans. Emotions are expressed by: (1) behavior [28]; (2) spoken dialogs [22]; (3) verbal actions such as variations in speech and its intensity including silence; (3) non-verbally using gestures; (4) facial expressions [11] and tears; and (5) their combinations. In addition to the analysis of these signals, one has to be able to analyze and understand the preceding events and/or predicted future events, individual expectations, personality, intentions, cultural expectations, and the intensity of an action. There are many studies to classify primary and derived emotions [4, 6, 16, and 28].

During conversation, people scan facial expressions of other persons to get a visual cue to their emotions. In social robotics, it is essential for robots to analyze facial expressions and express a subset of human emotions to have a meaningful human-robot interaction.

There are many schemes for the classification of human emotions. One popular theory for social robotics is due to Ekman [10, 11] that classifies human emotions into six basic categories: surprise, fear, disgust, anger, happiness, and sadness. In addition, there are many composite emotions derived by the combination of these basic emotions. The transitions between emotions that require continuous facial-expression analysis.

Three major techniques have been used to simulate and study human facial expressions: FACS (Facial Action Control System) [5], GF (Geometric Features) and GBMT (Graph Based Modeling Techniques) [7]. FACS simulates facial muscle movement using a combination of facial action units (FAUS or AUs). Different combinations of AUs model different muscle movements and specific facial expressions. FACS has found a major use in realistic visualization of facial-expressions through animation [1, 13]. Facial expression analysis techniques are based upon geometrical feature

extraction [8], modeling extracted features as graphs, and analyzing the variations in the graph for deviation.

Previous techniques [9] start afresh every time they analyze the emotion, and the accounting for expectations of emotions is not important. They also do not take into account the fact that a subset of features are unique to the presence or the absence of specific facial-expression. Identification of these subsets of unique features during facial expression analysis can prune the search space.

In this paper, we identify subsets of action units (AUs) that uniquely characterize the presence or the absence of a subset of emotions and map these AUs to the geometric feature-points to prune the search space. The technique extends previous facial expression identification techniques based upon LSH (Locality Sensitive Hashing) [17] that employ LSH for efficient pruning of search space.

The technique has been demonstrated using a publicly available image database [21, 23] that has been used by previous approaches. Results show significant improvement in performance over time (70% improvement in execution time) compared to previous techniques while retaining the accuracy in the similar range. The proposed technique is also suitable for fast emotion recognition in videos and real-time robot-human interaction, as the scheme recognizes emotion transitions.

The major contributions of this paper are:

Applying the subset of action units for pruning the search space for different emotions during interactive communication.

Combining these subsets of action units with geometric modeling to reduce the number of feature points and transformation, thus improving execution efficiency.

The rest of the paper is organized as follows. Section 2 describes the background of FACS, geometric features for emotion, Principal Component Analysis (PCA) and Support Vector Machine (SVM). Section 3 describes the proposed approach. Section 4 refers to the algorithm and the implementation. Section 5 demonstrates the dataset and the results. Section 6 explains the related works and the last section concludes the paper and describes the future works.

2. BACKGROUND

2.1. FACS - Facial Action Control System

Contractions of a subset of facial muscles generate human facial expressions. A set of 66 AUs (Action Units) [11] have been used to simulate the

contractions of facial muscles. An action unit simulates the activities of one or several muscles in the face. Tables 1 and 2 describe a relevant subset of action units needed for the simulation of the facial expressions for the six basic emotions.

Table 2. Set of action units needed for basic emotions

Basic expressions	Involved Action Units
Surprise	AU 1, 2, 5,15,16, 20, 26
Fear	AU 1, 2, 4, 5,15,20, 26
Disgust	AU 2, 4, 9, 15, 17
Anger	AU 2, 4, 7, 9,10, 20, 26
Happiness	AU 1, 6,12,14
Sadness	AU 1, 4,15, 23

2.2. PCA - Principal Component Analysis

PCA is a dimension reduction technique that transforms the data-points to a new coordinate system using orthogonal linear transformation, such that the transformed data-points lie with greatest variance on the first coordinate. Only the dimensions with major variations are chosen for further analysis, reducing the feature space. It is based upon finding out the eigenvectors and eigenvalues [14].

2.3. SVM - Support Vector Machine

SVM creates a set of hyper-planes in a high-dimensional space, which can be used for classification, regression, or other tasks. A good separation is achieved by a hyperplane that has the largest distance to the nearest training data point of any class (functional margin). The operation of the SVM algorithm is based on finding the hyperplane that gives the largest minimum distance to the training examples, and it is called *margin*. Support Vectors (SV) are the elements of the training set that would change the position of the dividing hyperplane that are critical elements of the training set and are closest to the hyperplane. In general, the larger margin makes lower generalization error of the classifier [3, 25].

2.4. Geometric features in facial expression

Face-features can be modeled as a graph [15, 18]. Face movement is a combination of all facial feature points, but some points have a main role in facial expression. There are three types of nodes (feature-points): *stable*, *passive* and *active*. Stable feature-points are fixed. Passive feature-points do not have

significant muscle movement. Active feature-points are most affected by muscle movements and the change in position of active-points makes the change in facial-expressions.

The number of feature points have been reduced from 62 points to 24 major points without loss of information as these 24 major points are present in all basic emotions. In the modified model, there are 6 points in eyebrows

{bl₁...bl₃, br₁...br₃}, 8 points in eyes {el₁...el₄, er₁...er₄} and 10 points on mouth {ml₁...ml₃, mm₁...mm₄, mr₁...mr₃}. The subset {er₁, el₁} represents stable points, the subset {er₄, el₄, mr₂, ml₂} represents passive points, and the subset {br₁, br₂, br₃, bl₁, bl₂, bl₃, er₂, er₃, el₂, el₃, mr₁, mr₂, mr₃, ml₁, ml₂, ml₃, mm₁, mm₂, mm₃, mm₄} represents the active-points. Figure 1 shows the geometric model with the feature points.

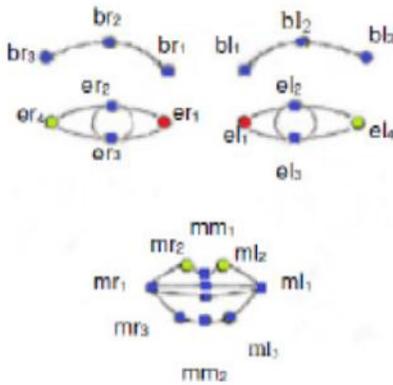


Figure 1. Feature-points in a geometric model of a face

The transformation matrix maps muscles and face movement to a formula based upon the movement of the feature points. Each movement is a combination of translation, rotation, and scaling. This transformation is caused due to head-movements, and the change in coordinates of the fixed inner eye corners er₁ with coordinate (x_r, y_r) and el₁ with coordinate (x_l, y_l). The transformations are given in equations (1) thru (5). The abbreviations *norm*, *Trans*, *rot*, and *sc* denote *normalize*, *transform*, *rotate* and *scale* respectively.

$$norm(x, y) = sc(x, y) \times rot(x, y) \times trans(x, y) \quad (1)$$

$$trans(x, y) = \begin{bmatrix} -\frac{x_l + x_r}{2} \\ \frac{y_l + y_r}{2} \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} \quad (2)$$

Where (x_l, y_l) and (x_r, y_r) are the coordinates of left and right inner eye corners el₁ and er₁ respectively:

$$rot(x, y) \begin{bmatrix} \cos(-\theta) & -\sin(-\theta) \\ \sin(-\theta) & \cos(-\theta) \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} \quad (3)$$

Where θ is the angle between the intervals joining the inner eye corners and the horizontal x-axis.

$$sc(x, y) = \frac{1}{2x_r} \begin{bmatrix} x \\ y \end{bmatrix} \quad (4)$$

$$sc(x, y) \frac{1}{2x_l} \begin{bmatrix} x \\ y \end{bmatrix} \quad (5)$$

Where x_r and x_l are the x-coordinates of right and left eyes respectively.

Table 3 describes the deviations of various facial feature-points that are needed to simulate facial expressions. Various movements of facial feature-points are *left*, *right*, *up*, *down*, *stretch* and *tighten*.

Table 3. Actions of feature-points in Figure 1

Feature Points	Deviation
Brow points (br ₁ , br ₂ , br ₃ , bl ₁ , bl ₂ , bl ₃)	up, down
Mid points of eyes (er ₂ , er ₃ , el ₂ , el ₃)	up
Outer lip points (mr ₁ , ml ₁)	stretch, tighten
Midpoint of upper lips (mm ₁ , mm ₂)	up, down

In order to separate the intensities of feature points for different emotions, the equation for cumulative difference is defined as follows:

$$diff = \sum_{i=1}^{n-1} (E_{i+1} - E_i) - \sum_{i=1}^{n-1} (N_{i+1} - N_i) \quad (6)$$

Where E_i (0 ≤ i ≤ n - 1) represents the feature point of an expressive face-state and N_i (0 ≤ i ≤ n - 1) represents feature point of a neutral face-state respectively. In the equation 6, the outcome *diff* > 0 means muscle-elongation and *diff* < 0 means muscle-contraction.

2.5 New definitions

Facial expressions are modelled using a single action-unit or a composite action-unit made of more than one action-units. For example, the facial expression for "happiness" is characterized by any of the three single AUs: 6, 12, and 14 (see Table 4), while the facial expression for "fear" is defined by a composite action-unit consisting of AUs 4 and 5. Composite action units are modeled as tuples. Thus the facial expression for "fear" is characterized by an AU-tuple (4, 5) (see Table IV).

3. APPROACH - UNIFIED METHOD

The integrated method is based upon:

- 1) Identifying the subset of AUs that are unique to basic emotions as shown in Table IV;
- 2) The subset of AUs that are absent in basic emotions as shown in Table V;
- 3) The subset of the AUs that will clearly shows transition from one basic emotion to another emotion as shown in Table VI.

Mapping these subsets to the change in geometric features (see Figure 1) and restricting the runtime check for the changes in the subset of geometric features significantly reduces the execution time of the facial expression analysis. A minimal subset of at least seven AUs are needed to check for the presence of any emotion uniquely. For example, Table 1 shows AU 1, 2, 5, 15, 16, 20 and 26 and for recognition of surprise. However, AU 16 is not needed in other emotions. It means AU 16 is sufficient to recognize the state “surprise”. Table IV describes a reduced subset to identify six basic facial expressions. The confidence factor can be improved further by:

- 1) Checking for additional AUs that characterize facial-expressions as in the case of happiness;
- 2) Checking for the absence of facial-expressions showing by the presence of AUs (see Table V).

Table 4. Subsets of unique AUs in basic emotions

State	AUs
Surprise	{16}
Fear	{(4, 5)}
Disgust	{17}
Anger	{10}
Happiness	{6, 12, 14}
Sadness	{23}

Table 5 lists sets of major AUs for each state. Using these sets, unique subset of action units present/absent in the specific facial expressions can be predicted. For example, about NOT surprise (absence of surprise) is given by the subset {4, 6, and 23}. A reduced subset T = {1, 2, 4, 5, 6, 7, 9} has been used to check for the absence of any uniquely emotion using a simple decision tree as shown in Figure 2.

Table 5. Subsets of AUs absent in basic emotions.

State	AUs
Not surprise (NSur)	{4, 6, 23}
Not fear (NF)	{6, 9, 16, 23}
Not disgust (ND)	{1, 7}
Not anger (NA)	{1, 5, 23}
Not happiness (NH)	{2, 4, 5, 9, 10, 16, 17, 20}
Not sadness (NSad)	{2, 5, 6, 9, 10, 16, 20}

To handle the transition of emotions from an existing emotion to another emotion, the tables of the difference between emotions is utilized. Table 6 gives the subsets of actions units that are present or absent when emotion transitions from the state surprise to other basic facial expressions. For example, to derive the transition from the facial expression *fear* to *surprise*, AU 4 should be present, and the subset {7, 9, 10, 17, and 23} should be absent. Similarly, to see transition from the state *surprise* to *happy* any of the three AUs 6 or 12 or 14 are sufficient. The symbol “P” denotes *presence*, and the symbol “A” denotes *absence*.

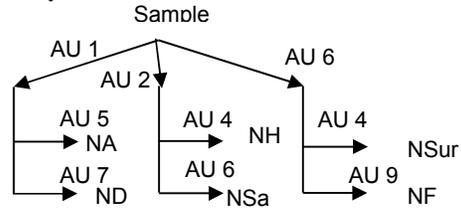


Figure 2. Classification-tree for absence of the emotions.

Table 6. Differences in emotion pairs involving surprise.

Emotions Pairs	AUs
Surprise → Fear	P: {4}; A: {7, 9, 10, 17, 23}
Surprise → Disgust	P: {4, 9, 17}; A: {9, 10, 23}
Surprise → Anger	P: {4, 7, 9, 10}; A: {17, 23}
Surprise → Happiness	P: {6 / 12 / 14}; A: {4}
Surprise → Sadness	P: {4, 23}; A: {7, 9, 10, 17}

As shown in the classification tree in Figure 3 by using the minimal subset of {4, 6, 7, 10, 17, and 23} the separation between emotion transitions can be done.

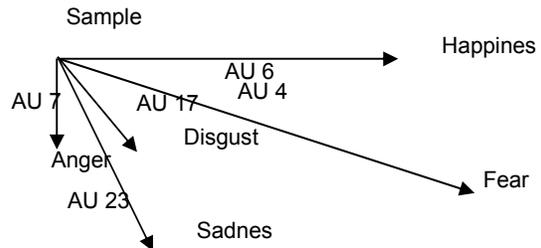


Figure 3. Decision-tree for emotional state after surprise.

3.1 Mapping action units to geometric features

Since the image analysis system only sees the changes in the geometric features of the face, the effect of AUs have to be mapped to the observable changes in the geometric features. Table 7 describes a mapping between AUs and the movement of geometric feature points. The mapping shows that many muscle movements map to the same geometric features. For example, AU # 6, 12, and 14 all are involved in stretching mr_1 and ml_1 ; AU # 4 and #9 pull down the inner brow points br_1 and bl_1 down; and AU # 16 and #26 pull mm_3 down. . While, multiple emotions may map to the movement of the same feature points, the magnitude of movement is different, and is derived by the SVM training using *diff* equation as explained in section 2.

TABLE 7. MAPPING OF ACTION UNITS TO GEOMETRIC FEATURES

Action units ↔ Features				Action units ↔ Features			
AUs		Features		AUs		Features	
AU	Action	Id	Action	AU	Action	Id	Action
1	up	br_1, bl_1	up	12	pull	mr_1, ml_1	stretch
2	up	br_3, bl_3	up	14	dimple	mr_1, ml_1	stretch
4	down	br_1, bl_1	down	16	down	mm_3, mm_4	down
5	up	mm_1, mm_2	up	17	up	mm_3, mm_4	up
6	up	mr_1, ml_1	stretch	20	stretch	mr_1, ml_1	stretch
7	tight	mr_1, ml_1	tight	23	tight	mr_1, ml_1	tight
9	wrinkle	br_1, bl_1	down	26	down	mm_3, ml_3	down

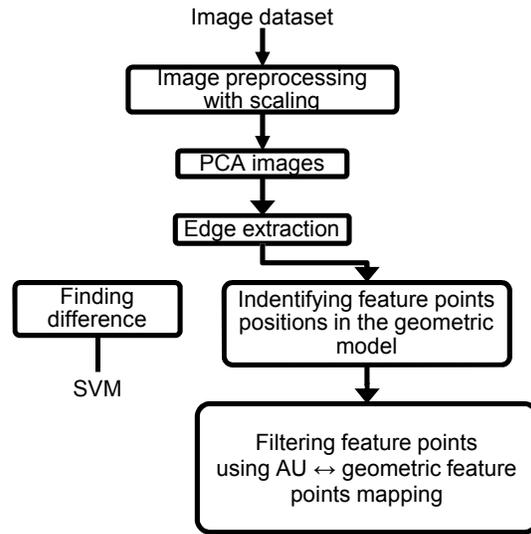


Figure 4. Flow of the algorithm

4. IMPLEMENTATION

Facial behaviour of image sequences has been chosen from CK+ database [23]. We compared emotion recognition results on the CK+ database over different dimensions with those produced by SVM. For each emotion category, one-third of the 653 images in the database were selected for training and the remaining images were used for testing. Images of size 490*400 were transformed into 196000*1 dimensional column vectors. Input features sizes are high and a PCA was used to reduce image dimension. All experiments have been implemented using MATLAB. Figure 4 shows the process. For the first step, a part of the dataset will be selected and after that, some image pre-processing such as scaling is applied on images. Then PCA is used to find component images and then canny filter is applied for edge extraction from PCA components. Now the pixels around the geometric feature points should be extracted in each image for finding measures and difference for SVM training. The number of pixels are filtered by using the AU and geometric feature point mapping as explained in Section 3.

The inputs of SVM are some vectors that are related to each image in the database and extracted using AUs and the geometric feature model.

5. Experimental results

In the figure 4, the flow of algorithm has presented. At first some processing on the images is done that include some resizing and finding face in

each image. Then the PCA is applied on the image that its outputs are some transformed images by eigenvectors. Edge detection of PCA images is the next step. Then, feature points and their position are extracted. In follow some points that are related to the geometric model are selected and finally the differences of extracted pint are calculated for training and testing SVM.

Figure 5 shows a sample of randomly selected emotion-state images after principle component analysis.

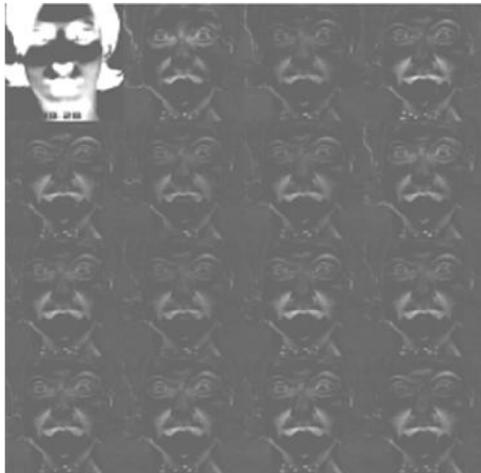


Figure 5. PCA components



Figure 6 shows the output after the corresponding images after edge extraction.

Table 8 presents the percentage of correctness, number of SVs and margin for emotion recognition by the general approach using just SVM. The abbreviation *NSur* denotes "not surprise", *NA* denotes "not angry", *NSad* denotes "not sad", *NH* denotes

"not happy", *ND* denotes "not disgust", and *NF* denotes "not fear".

Table 8. Emotion recognition using just SVM

Emotions	SVnumber	Margin	Correctness
NSur	8.1	1.95	75 %
NA	7.3	1.83	73%
NSa	8.3	1.88	76%
NH	6.5	1.95	84%
ND	8.1	1.79	79%
NF	7.3	1.83	74%

Table 9 shows the correctness using the proposed unified approach. It is clear in Table 9 that correctness is significantly better for *NSur* and *NA*, and it is comparable for *NSad*, *NH*, *ND* and *NF*.

Table 9. Emotion recognition using proposed method.

Emotions	SV number	Margin	Correctness
NSur	7.1	1.67	83
NA	7.9	1.60	81
NSa	7.5	1.88	74
NH	7.3	1.56	82
ND	7.6	1.70	75
NF	8.8	1.83	78

Based on the results shown in the Tables 8 and 9, the average SV number is 7.6 in previous method and 7.7 in the proposed method. The average margin is 1.83 in previous method and 1.71 in the method described in this paper. Clearly, the new method is more time-efficient. Table 10 compares the execution efficiency of different approaches. It is clear that with this strategy processing time improves by 70% due to the reduction of the number of AUs and the corresponding feature-points in a face.

Table 10. Execution efficiency of the proposed method.

Emotions	Execution time (old method)	Execution Time (new method)
NSur	7.5	1.3
NA	6.7	1.3
NSa	7.3	1.3
NH	6.7	1.3
ND	7.4	1.3
NF	7.5	1.8

6. Related works

FACS system has been used for emotion generation by many researchers using AU based simulation [1, 13]. Many researchers use a geometric model [18, 29,

and 31] and try to improve geometric models. Some articles present a framework for recognition of facial action unit (AU) combinations by viewing the classification as a special representation problem [29, 31], and others present heuristic methods for achieving better performance [23, 27, 30, 32]. Here a modified geometric model has been used that reduces the facial points and uses another metric for finding distances. In some researches, a modified coding system has been presented [19], and some others' work researches integrate coding system, improved strategies and methods [20]. There are many situations that need to recognize an emotion that does not exist in the coding tables. In the real world and implementing algorithms we compare or combine some other coding or emotion state to find that special state. This article presents a coding system based on the basic emotions that can apply directly in such these situations. Many researchers have developed the version of a computer vision system that are sensitive to subtle changes in the face [21, 24]. In this paper, we use some statistical method to reduce the dimension of training space. Our proposed scheme significantly improves the performance while retaining accuracy and is suitable for real-time analysis of facial expressions and for real-time human-robot interaction.

7. Conclusion and future work

In this paper a unified method for facial expression detection has been presented. The technique maps the AUs for specific emotions to geometric feature point movements, and uses the characterizing feature points based upon AUs mapping to prune the number of pixels being processed, improving the execution time. Actually, here the correctness improves or are same in three of the emotions, and are within the range of traditional techniques for the remaining three emotions. The major gain is the 70% performance improvement over time due to pruning of the number of pixels being processed. The improved performance makes it suitable for real-time robot-human interaction. This performance in the database with the large number of image or images with the high dimension makes it more applicable. The scheme can be further improved by analyzing the duration of various emotions when the emotion does not change, and needs minimal analysis. We are extending the current scheme to incorporate the duration of various emotions. We are also extending the proposed scheme to handle secondary emotions.

References

- [1] A. K. Bansal and Y. Hijazi. Low Bandwidth Video Conversation Using Anatomical Reconstruction of Facial Expressions Over the Internet. IADIS International Conference on WWW/Internet, Murcia, Spain, 154-161, 2006.
- [2] M. S. Bartlett, G. Littlewort, M. G. Frank, C. Lainscsek, I. Fasel, and J. R. Movellan. Recognizing facial expression: Machine learning and application to spontaneous behavior. Proceedings of the IEEE Conf. Comput. Vis. Pattern Recognition, 2:568- 573, 2005.
- [3] C. J. C. Burges. A Tutorial on Support Vector Machines for Pattern Recognition. Bell Laboratories, 1998.
- [4] E. Cambria, A. Livingstone, and A. Hussain. The Hourglass of Emotions. Cognitive Behavioural Systems, Springer Verlag, 144 – 157, 2012.
- [5] S. W. Chew, R. Rana, P. Lucey, S. Lucey, and S. Sridharan. Sparse Temporal Representations for Facial Expression Recognition. in Advances in Image and Video Technology, Springer-Verlag, 7088:311 – 322, 2012.
- [6] G. Colombetti. From affect programs to dynamical discrete emotions. Philosophical Psychology, 22: 407- 425, 2009.
- [7] T. F. Cootes, G. J. Edwards, and C. Taylor. Active appearance models. IEEE Trans. Pattern Anal. Mach. Intell., 23(6):681 – 685, 2001.
- [8] G. Donato, M. S. Bartlett, J. C. Hager, P. Ekman, and T. J. Sejnowski. Classifying Facial actions. IEEE Trans. Pattern Anal. Mach. Intell., 21(10): 974 – 989, 1999.
- [9] F. Dornaika and F. Davoine. Simultaneous facial action tracking and expression recognition in the presence of head motion. *Int. Journal Comput. Visualization.*, 76(3): 257-281, 2008.
- [10] P. Ekman and W. Friesen. Facial Action Coding System: A Technique for the Measurement of Facial Movement. Consulting Psychologists Press, Palo Alto, 1978.
- [11] P. Ekman, Facial expression and emotion. *American Psychologist*, 48(4):384 – 392, 1993.
- [12] J. M. Fellous and M. A. Arbib. Who needs emotions? The brain meets the robots, Oxford press, 2005.
- [13] M. Fratarcangeli. Computational Models for Animating 3D Virtual Faces. Linköping Studies in Science and Technology Thesis, No. 1610, Division of Image Coding, Department of Electrical Engineering, Linköping University, SE-581 83 Linköping, Sweden, 2013, available at <http://www.diva-portal.org/smash/get/diva2:646028/FULLTEXT02.PDF>
- [14] K. Fukunaga. Introduction to Statistical Pattern Recognition. Academic Press, San Diego, USA, 1990.

- [15] L. Gang, L. Xiao-hua, Z. Ji-liu, and G. Xiao-gang. Geometric feature based facial expression recognition using multiclass support vector machines. *IEEE International Conference on Granular Computing*, 318 - 321, 2009.
- [16] M. Gendron and F. Lisa. Reconstructing the Past: A Century of Ideas about Emotion in Psychology. *Emotion Review*, 1(4): 316-339, 2009.
- [17] M. Ghayoumi and A. Bansal. Exploiting Locality Sensitive Hashing for Improved Emotion Recognition. *International Conference on Signal Processing and Multimedia Applications (SIGMAP 14)*, Vienna, Austria, 211-219, 2014.
- [18] K. Hong et al, A Component Based Approach for Classifying the Seven Universal Facial Expressions of Emotion. *Proceedings of the IEEE Symposium on Computational Intelligence for Creativity and Affective Computing*, 1 - 8, 2013.
- [19] K.E. KO and K. B. Sim. Development of a Facial Emotion Recognition Method based on combining AAM with DBN. *International Conference on Cyber worlds*, 87 - 91, 2010.
- [20] I. Kotsia and I. Pitas. Facial Expression Recognition in Image Sequences Using Geometric Deformation Features and Support Vector Machines. *IEEE Transaction on Image Processing*, 16(1): 172 - 187, 2007.
- [21] J. J. Lien, T. Kanade, J. F. Cohn, and C. Li. Detection, tracking, and classification of action units in facial expression, *J. Robot. Auto. Syst.*, 31(3):131 - 146, 2000.
- [22] C. M. Lee and S. S. Narayanan. Towards Detecting Emotions in Spoken Dialog. *IEEE Trans. on Speech and Audio Processing*, 13(2):293-303, 2005.
- [23] P. Lucey, J. F. Cohn, T. Kanade, J. Saragih, Z. Ambadar, and I. Matthews. The Extended Cohn-Kanade Dataset (CK+): A complete expression dataset for action unit and emotion-specified expression. *Proceedings of the Third International Workshop on CVPR for Human Communicative Behavior Analysis (CVPR4HB)*, San Francisco, USA, 94-101, 2010.
- [24] M. H. Mahoor, M. Zhou, K. L. Veon, S. M. Mavadati, and J. F. Cohn. Facial action unit recognition with sparse representation. *Proceedings of the IEEE Int. Conf. Autom. Face Gesture Recognition*, pp. 336 - 342, 2011.
- [25] H. A. Moghaddam and M. Ghayoumi. Facial Image Feature Extraction using Support Vector Machine. *International Conference of Vision Theory and Applications*, Lisbon, Portugal, 480-485, 2006.
- [26] R. R. Pagariya, M. M. Bartere, and ?. Facial Emotion Recognition in Videos Using HMM. *International Journal of Computational Engineering Research*, 3(4-3):111-118, 2013.
- [27] M. Pantic and I. Patras. Dynamics of facial expressions: Recognition of facial actions and their temporal segments from face profile image sequences. *IEEE Trans. Syst., Man, Cybern. B, Cybern.* 36(2):433 - 449, 2006.
- [28] R. Plutchik. The Nature of Emotions. *American Scientist*, 89:344 - 350, 2001.
- [29] M. Rogers and J. Graham, Robust active shape model search, *Proceedings of the Eur. Conf. Comput. Vis.*, 517 - 530, 2002.
- [30] C. Shan, S. Gong, and P. W. McOwan. Facial expression recognition based on local binary patterns: A comprehensive study. *Image Vis. Comput.*, 27(6): 803 - 816, 2009.
- [31] Y. Tian, T. Kanade, and J. F. Cohn. Recognizing action units for facial expression analysis. *IEEE Trans. Pattern Anal. Mach. Intell.*, 23(2):97 - 115, 2001.
- [32] Y. Tong, W. Liao, and Q. Ji. Facial action unit recognition by exploiting their dynamic and semantic relationships. *IEEE Trans. Pattern Anal. Mach. Intell.*, 29(10):1683 -1699, 2007.
- [33] Z. Wang, Y. Li, S. Wang, and Q. Ji. Capturing Global Semantic Relationships for Facial Action Unit Recognition. *IEEE International Conference on Computer Vision*, 2013.

Authors Index

Abbod, M.	15	Fratu, A.	253	Molina, A.	183
Abozeid, A.	195	Fratu, M.	253	Monroy, S.	116
Ahmed, H. E. H.	220	Ganji, B. F.	167	Moreno, J. P.	161
Ajgou, R.	156	Ghayoumi, M.	259	Nasser, A.	95
Akhtar, M. N.	241	Ghendir, S.	156	Novickis, L.	32
Al Salameh, M. S. H.	207	Graur, A.	110	Nsiri, B.	145
Al-Ani, S. M.	15	H. Toulni	145	Nyssonbayeva, S. E.	170
Al-Dahoud, A.	121, 177	Hamid, A.	213	Ordoñez, S.	116
Aljaafreh, A.	56	Hammoudi, Z.	60	Pasca, A.	40, 48
Al-Rawashdeh, T.	121	Hanane, T.	213	Pizzolante, R.	27
Al-Zu'bi, M. M.	207	Haq, Z. U.	102	Ponce, P.	183
Amrouche, A.	200	Hussain, T.	102	Ponomaryov, V.	65
Bansal, A. K.	259	Ibarra, L.	183	Poonsilp, A.	133
Bartusevics, A.	32	Jannoud, I.	121	Poonsilp, K.	133
Begimbayeva, Y. Y.	170	Jeong, J.	139	Prieto, F. B.	161
Beltrán, A.	116	Khairy, M.	195	Ramezani, A.	167
Biyashev, R. G.	170	Khan, G. F.	102	Ruiz, A. T.	161
Boudemagh, N.	60	Khedr, M. E.	95	Rusan, A.	84
Boulmalf, M.	145	Koscelnik, J.	230	Sadiki, T.	145
Buzduga, C.	110	Kvet, M.	75	Samadbeik, M.	167
Carpentieri, B.	27	Lagos, J. A.	161	Sbaa, S.	156
Chang, Y.-C.	150	Lee, J.	139	Sharkas, M.	95
Chemsa, A.	156	Lesovskis, A.	32	Shim, Y.-C.	20
Chiu, Y.-L.	150	Li, J.-W.	150	Staines, A. S.	126
Ciufudean, C.	110	Lucklum, R.	56	Taleb-Ahmed, A.	156, 200
Daira, R.	90	Maamar, H.	213	Toro, B.	161
Dobrucky, B.	230	Macías, I.	183	Vasiu, R.	84
Duarte, N.	116, 161	Magzom, M. M.	170	Vilardy, A.	161
El-Aziem, A. H. A.	220	Martinez, R.	161	Vlad, V.	110
ElDahshan, K.	195	Meddah, M.	200	Yagoubi, B.	191
Farouk, H.	195	Melo, L.	116	Yang, F.-S.	150
Fezari, M.	121, 177	Mohamad-Saleh, J.	241		