



ICCGI 2014

The Ninth International Multi-Conference on Computing in the Global Information
Technology

ISBN: 978-1-61208-346-9

June 22 - 26, 2014

Seville, Spain

ICCGI 2014 Editors

Arno Leist, Massey University, New Zealand

Tadeusz Pankowski, Poznan University of Technology, Poland

ICCGI 2014

Foreword

The Ninth International Multi-Conference on Computing in the Global Information Technology (ICCGI 2014), held between June 22-26, 2014 - Seville, Spain, continued a series of international events covering a large spectrum of topics related to global knowledge concerning computation, technologies, mechanisms, cognitive patterns, thinking, communications, user-centric approaches, nanotechnologies, and advanced networking and systems. The conference topics focus on challenging aspects in the next generation of information technology and communications related to the computing paradigms (mobile computing, database computing, GRID computing, multi-agent computing, autonomic computing, evolutionary computation) and communication and networking and telecommunications technologies (mobility, networking, bio-technologies, autonomous systems, image processing, Internet and web technologies), towards secure, self-defendable, autonomous, privacy-safe, and context-aware scalable systems.

This conference intended to expose the scientists to the latest developments covering a variety of complementary topics, aiming to enhance one's understanding of the overall picture of computing in the global information technology.

The integration and adoption of IPv6, also known as the Next Generation of the Internet Protocol, is happening throughout the World at this very moment. To maintain global competitiveness, governments are mandating, encouraging or actively supporting the adoption of IPv6 to prepare their respective economies for the future communication infrastructures. Business organizations are increasingly mindful of the IPv4 address space depletion and see within IPv6 a way to solve pressing technical problems while IPv6 technology continues to evolve beyond IPv4 capabilities. Communications equipment manufacturers and applications developers are actively integrating IPv6 in their products based on market demands.

IPv6 continues to represent a fertile area of technology innovation and investigation. IPv6 is opening the way to new successful research projects. Leading edge Internet Service Providers are guiding the way to a new kind of Internet where any-to-any reachability is not a vivid dream but a notion of reality in production IPv6 networks that have been commercially deployed. National Research and Educational Networks together with internationally known hardware vendors, Service Providers and commercial enterprises have generated a great amount of expertise in designing, deploying and operating IPv6 networks and services. This knowledge can be leveraged to accelerate the deployment of the protocol worldwide.

We take here the opportunity to warmly thank all the members of the ICCGI 2014 Technical Program Committee, as well as the numerous reviewers. The creation of such a high quality conference program would not have been possible without their involvement. We also kindly thank all the authors who dedicated much of their time and efforts to contribute to ICCGI 2014. We truly believe that, thanks to all these efforts, the final conference program consisted of top quality contributions.

Also, this event could not have been a reality without the support of many individuals, organizations, and sponsors. We are grateful to the members of the ICCGI 2014 organizing

committee for their help in handling the logistics and for their work to make this professional meeting a success.

We hope that ICCGI 2014 was a successful international forum for the exchange of ideas and results between academia and industry and for the promotion of progress in the area of computing in the global information technology.

We are convinced that the participants found the event useful and communications very open. We also hope the attendees enjoyed the charm of Seville, Spain.

ICCGI 2014 Chairs:

ICCGI Advisory Committee

Constantin Paleologu, University Politehnica of Bucharest, Romania

Tibor Gyires, Illinois State University, USA

Luc Vouligny, Institut de Recherche d'Hydro-Québec - Varennes, Canada

John Terzakis, Intel, USA

Yasushi Kambayashi, Nippon Institute of Technology, Japan

Mirela Danubianu, "Stefan cel Mare" University of Suceava, Romania

Teemu Kanstrén, VTT Technical Research Centre of Finland, Finland

Mansour Zand, University of Nebraska, USA

Arno Leist, Massey University, New Zealand

Jean-Denis Mathias, IRSTEA, France

Dominic Girardi, RISC Software GmbH, Austria

ICCGI Special Area Chairs

Knowledge/Cognition

Constandinos Mavromoustakis, University of Nicosia, Cyprus

Tadeusz Pankowski, Poznan University of Technology, Poland

e-Learning/Mobility

José Rouillard, Université Lille Nord, France

Industrial Systems

Beata Czarnacka-Chrobot, Warsaw School of Economics, Poland

Bernhard Freudenthaler, Software Competence Center Hagenberg GmbH, Austria

ICCGI Publicity Chair

Marek Opuszko, Friedrich-Schiller-University of Jena, Germany

ICCGI 2014

Committee

ICCGI Advisory Committee

Constantin Paleologu, University Politehnica of Bucharest, Romania
Tibor Gyires, Illinois State University, USA
Luc Vouligny, Institut de Recherche d'Hydro-Québec - Varennes, Canada
John Terzakis, Intel, USA
Yasushi Kambayashi, Nippon Institute of Technology, Japan
Mirela Danubianu, "Stefan cel Mare" University of Suceava, Romania
Teemu Kanstrén, VTT Technical Research Centre of Finland, Finland
Mansour Zand, University of Nebraska, USA
Arno Leist, Massey University, New Zealand
Jean-Denis Mathias, IRSTEA, France
Dominic Girardi, RISC Software GmbH, Austria

ICCGI Special Area Chairs

Knowledge/Cognition

Constandinos Mavromoustakis, University of Nicosia, Cyprus
Tadeusz Pankowski, Poznan University of Technology, Poland

e-Learning/Mobility

José Rouillard, Université Lille Nord, France

Industrial Systems

Beata Czarnacka-Chrobot, Warsaw School of Economics, Poland
Bernhard Freudenthaler, Software Competence Center Hagenberg GmbH, Austria

ICCGI Publicity Chair

Marek Opuszko, Friedrich-Schiller-University of Jena, Germany

ICCGI 2014 Technical Program Committee

Pablo Adasme, Universidad de Santiago de Chile, Chile
El-Houssaine Aghezzaf, Gent University, Belgium
Werner Aigner, FAW, Austria
Johan Akerberg, ABB Corporate Research, Sweden
Nadine Akkari, King Abdulaziz University, Kingdom of Saudi Arabia

Konstantin Aksyonov, Ural Federal University, Russia
Areej Al-Wabil, King Saud University - Riyadh, Saudi Arabia
Cesar Alberto Collazos, Universidad del Cauca, Colombia
Cristina Alcaraz, University of Malaga, Spain
Jose M. Alcaraz Calero, University of the West of Scotland, UK
Panos Alexopoulos, iSOCO S.A. - Madrid, Spain
Ali Alharbi, The University of Newcastle, Australia
Fernando Almeida, University of Porto, Portugal
Hala Mohammed Alshamlan, King Saud University, Saudi Arabia
José Enrique Armendáriz-Iñigo, Universidad Pública de Navarra, Spain
Stanislaw Ambroszkiewicz, Institute of Computer Science Polish Academy of Sciences, Poland
Christos Anagnostopoulos, Ionian University, Greece
Plamen Angelov, Lancaster University, UK
Josep Arnal Garcia, Universidad de Alicante, Spain
Ezendu Ariwa, London Metropolitan University, UK
Kamran Arshad, University of Greenwich, UK
Mustafa Atay, Winston-Salem State University, USA
Ali Barati, Azad University - Dezful Branch, Iran
Reza Barkhi, Virginia Tech - Blacksburg, USA
Carlos Becker Westphall, Universidade Federal de Santa Catarina, Brazil
Hatem Ben Sta, University of Tunis, Tunisia
Jorge Bernardino, Institute Polytechnic of Coimbra - ISEC, Portugal
Robert Bestak, Czech Technical University in Prague, Czech Republic
Ateet Bhala, Oriental Institute of Science and Technology, Bhopal, India
Fernando Bobillo, University of Zaragoza, Spain
Mihai Boicu, George Mason University - Fairfax, USA
Eugen Borcoci, University 'Politehnica' of Bucharest, Romania
Djamila Boukredera, University of Bejaia, Algeria
José Braga de Vasconcelos, Universidade Atlântica, Portugal
Daniela Briola, University of Genoa, Italy
Francesco Buccafurri, University "Mediterranea" of Reggio Calabria, Italy
Luigi Buglione, Engineering.IT SpA, Italy
Xiaoqiang Cai, The Chinese University of Hong Kong, Hong Kong
Ani Calinescu, Oxford University, UK
George Caridakis, University of the Aegean and National Technical University of Greece, Greece
Laura Carnevali, University of Florence, Italy
Cheng-Yuan Chang, National United University, Taiwan
Maiga Chang, Athabasca University, Canada
Ankit Chaudhary, MUM, USA
Savvas A. Chatzichristofis, Democritus University of Thrace, Greece
Chi-Hua Chen, National Chiao Tung University - Taiwan, R.O.C.
David Chen, University of Bordeaux, France
Shu-Ching Chen, Florida International University - Miami, USA
Tzung-Shi Chen, National University of Tainan, Taiwan

Wen-Shiung Chen (陳文雄), National Chi Nan University, Taiwan
Zhixiong Chen, School of Liberal Arts, Mercy College - Dobbs Ferry, USA
Albert M. K. Cheng, University of Houston, USA
Amar Ramdane Cherif, University of Versailles, France
Dickson Chiu, Dickson Computer Systems, Hong Kong
Gihwan Cho, Chonbuk University, Korea
Michal Choras, University of Technology and Life Sciences Bydgoszcz (UTP), Poland
Franklyn Chukwunonso, Universiti Teknologi Malaysia (UTM), Malaysia
Jose M. Claver, University of Valencia, Spain
Francesco Colace, DIEII - Università degli Studi di Salerno, Italy
Rebeca Cortázar, University of Deusto, Spain
Beata Czarnacka-Chrobot, Warsaw School of Economics, Poland
Kevin Daimi, University of Detroit Mercy, USA
Dimitrios Damopoulos, Stevens Institute of Technology, USA
Mirela Danubianu, "Stefan cel Mare" University of Suceava, Romania
Bernard De Baets, Gent University, Belgium
José de Oliveira Guimarães, Federal University of São Carlos, Brazil
Vincenzo Deufemia, Università di Salerno - Fisciano, Italy
Kamil Dimililer, Near East University - Nicosia, Cyprus
Alexandre Dolgui, Ecole des Mines de Saint-Etienne, France
Ludek Dolihal, Masaryk University - Brno, Czech Republic
Juan Carlos Dueñas, ETSI Telecomunicación, Universidad Politécnica de Madrid, Spain
Christof Ebert, Vector Consulting Services GmbH, Germany
Chanaka Edirisinghe, The University of Tennessee - Knoxville, USA
Hans-Dieter Ehrich, Technische Universität Braunschweig, Germany
Moez Esseghir, Technology University of Troyes, France
Fausto Fasano, University of Molise, Italy
Maria Fazio, University of Messina, Italy
Javier Dario Fernandez Ledesma, Universidad Pontificia Bolivariana - Medellín, Colombia
Kurt B. Ferreira, Sandia National Laboratories, USA
Joerg Fliege, The University of Southampton, UK
Rana Forsati, Shahid Beheshti University - Tehran, Iran
Panagiotis Fotaris, University of West London, UK
Rita Francese, Università di Salerno - Fisciano, Italy
Bernhard Freudenthaler, Software Competence Center Hagenberg GmbH, Austria
Matjaz Gams, Jozef Stefan Institute - Ljubljana, Slovenia
Raghu Ganti, IBM Thomas J. Watson Research Center, U.S.A.
Félix J. García, University of Murcia, Spain
Vanessa Gardellin, Institute for Informatics and Telematics CNR, Pisa, Italy
David Garcia Rosado, University of Castilla-La Mancha, Spain
Joseph Andrew Giampapa, Carnegie Mellon University, USA
Debasis Giri, Haldia Institute of Technology, India
Apostolos Gkamas, University Ecclesiastical Academy of Vella of Ioannina, Greece
Gustavo González, Mediapro Research - Barcelona, Spain

T.R. Gopalakrishnan Nair, Prince Mohammad Bin Fahd University, Saudi Arabia
Feliz Gouveia, University Fernando Pessoa, Portugal
William Grosky, University of Michigan, USA
Francesco Guerra, University of Modena and Reggio Emilia, Italy
Carlos Guerrero, Universitat de les Illes Balears Palma de Mallorca, Spain
Nalan Gulpinar, University of Warwick - Coventry, UK
Chris Guy, University of Reading, UK
Tibor Gyires, Technology Illinois State University, USA
Maki K. Habib, The American University in Cairo, Egypt
Jana Hájková, University of West Bohemia, Pilsen, Czech Republic
Hani Hamdan, École Supérieure d'Électricité (SUPÉLEC), France
Seung-Hoon Han, Chonnam National University, Korea
Petr Hanáček, Brno University of Technology, Czech Republic
Fei Hao, Huazhong University of Science and Technology, China
Sven Hartmann, TU-Clausthal, Germany
Mohammad Mehedi Hassan, King Saud University, Kingdom of Saudi Arabia
Wladyslaw Homenda, Warsaw University of Technology, Poland
Samuelson W. Hong, Zhejiang University of Finance & Economics, China
Jun Hu, Eindhoven University of Technology, The Netherlands
Yo-Ping Huang, National Taipei University of Technology, Taiwan
Tee Sim Hui, Multimedia University, Malaysia
Rosziati Ibrahim, Universiti Tun Hussein Onn Malaysia (UTHM), Malaysia
Larisa Ismailova, National Research Nuclear University "MEPhI" - Moscow, Russia
Kyoko Iwasawa, Takushoku University - Tokyo Japan
Helge Janicke, De Montfort University, UK
Mehrshid Javanbakht, Azad University - Tehran, Iran
Guorui Jiang, Beijing University of Technology, China
Maria João Ferreira, Universidade Portucalense - Porto, Portugal
Paul Johannesson, Royal Institute of Technology - Stockholm, Sweden
Matjaz B. Juric, University of Ljubljana, Slovenia
Imed Kacem, Université de Lorraine, France
Hermann Kaindl, TU-Wien, Austria
Yasushi Kambayashi, Nippon Institute of Technology, Japan
Georgios Kambourakis, University of the Aegean - Samos, Greece
Byoung Uk Kim, US Air Force Research Laboratory, USA
Georgios Kioumourtzis, Center for Security Studies, Greece
Wojciech Kmiecik, Wroclaw University of Technology, Poland
Mehmet Koc, Bilecik Seyh Edebali University, Turkey
Leszek Koszalka, Wroclaw University of Technology, Poland
Janet L. Kourik, Webster University - St Louis, USA
Piotr A. Kowalski, SRI Polish Academy of Sciences and Cracow University of Technology, Poland
Bartosz Krawczyk, Wroclaw University of Technology, Poland
Jakub Kroustek, Brno University of Technology, Czech Republic
Panos Kudumakis, Queen Mary University of London, UK

Robert Law, Glasgow Caledonian University, UK
Gyu Myoung Lee, Telecom SudParis, France
Tracey K.M. Lee, School of EEE Singapore, Singapore
Arno Leist, Massey University, New Zealand
Daniel Lemire, Université du Québec à Montréal (UQAM), Canada
Ricardo Lent, Imperial College - London, UK
Isaac Lera, Universitat de les Illes Balears Palma de Mallorca, Spain
Tiberiu S. Letia, Technical University of Cluj-Napoca, Romania
Bo Li, Beihang University, China
Chendong Li, University of Connecticut - Storrs, USA
Abdel Lisser, University of Paris Sud, LRI - Orsay, France
Angela Locoro, University of Genova, Italy
Francesca Lonetti, CNR-ISTI, Italy
Josip Lorincz, University of Split, Croatia
Pericles Loucopoulos, Harokopio University of Athens, Greece / Loughborough University, UK
Alen Lovrencic, University of Zagreb, Croatia
Szymon Lukasik, Cracow University of Technology, Poland
Stephane Maag, Telecom SudParis, France
Olaf Maennel, Loughborough University, UK
José María Luna, University of Córdoba, Spain
Francesco Maiorana, University of Catania, Italy
Viliam Makis, University of Toronto, Canada
Giuseppe Mangioni, University of Catania, Italy
Gregorio Martinez, University of Murcia, Spain
Maristella Matera, Politecnico di Milano, Italy
Jean-Denis Mathias, National Research Institute of Science and Technology for Environment and Agriculture (IRSTEA), France
Constandinos Mavromoustakis, University of Nicosia, Cyprus
Natarajan Meghanathan, Jackson State University, USA
Angelos Michalas, Technological Educational Institute of Western Macedonia, Greece
Jerzy Michnik, University of Economics in Katowice, Poland
Marius Minea, University Politehnica of Bucharest, Romania
Vladimir Modrak, Technical University of Kosice - Presov, Slovakia
Lars Moench, FernUni-Hagen, Germany
Ghodrat Moghadampour, Vaasa University of Applied Sciences Technology and Communication, Finland
György Molnár, Budapest University of Technology and Economics, Hungary
Valérie Monfort, Université Paris 1 Panthéon Sorbonne, France
Paula Morais, Universidade Portucalense - Porto, Portugal
Mary Luz Mouronte López, Universidad Politénica de Madrid, Spain
Isabel Muench, German Federal Office for Information Security), Germany
Antonio Muñoz, University of Malaga, Spain
Phivos Mylonas, Ionian University, Greece
Pablo Najera, University of Malaga, Spain

Tomoharu Nakashima, Osaka Prefecture University, Japan
Joan Navarro, Ramon Llull University, Spain
Antonio Navarro Martín, Universidad Complutense de Madrid, Spain
Leila Nemmiche Alachaher, ISIMA, France
Marcio Katsumi Oikawa, Universidade Federal do ABC, Brazil
Hichem Omrani, CEPS/INSTEAD, Luxembourg
Chung-Ming Ou, Kainan University, Taiwan
Constantin Paleologu, University Politehnica of Bucharest, Romania
Thanasis G. Papaioannou, École Polytechnique Fédérale de Lausanne (EPFL), Switzerland
Marcin Paprzycki, Systems Research Institute / Polish Academy of Sciences - Warsaw, Poland
Tadeusz Pankowski, Poznan University of Technology, Poland
Mukaddim Pathan, Telstra Corporation Limited, Australia
Al-Sakib Khan Pathan, International Islamic University Malaysia, Malaysia
Kunal Patel, Ingenuity Systems, USA
Jose J. Pazos Arias, Universidad Vigo, Spain
Marek Penhaker, VSB - Technical University of Ostrava, Czech Republic
Yoseba K. Peña, University of Deusto, Basque Country, Spain
Andrea Perego, European Commission - Joint Research Centre, Italy
Fernando Pereñíguez García, Universidad Católica San Antonio Murcia, Spain
Zeeshan Pervez, University of the West of the Scotland, UK
Dana Petcu, Western University of Timisoara, Romania
Willy Picard, Poznan University of Economics, Poland
Selwyn Piramuthu, University of Florida, USA
Kornelije Rabuzin, University of Zagreb - Varazdin, Croatia
Stefan Rass, Universitaet Klagenfurt, Austria
Danda B. Rawat, Georgia Southern University, USA
Marek Reformat, University of Alberta - Edmonton, Canada
Luis Paulo Reis, University of Minho/LIACC, Portugal
Eric Renault, Télécom SudParis, France
Agos Rosa, LaSEEB, Portugal
José Rouillard, Université Lille Nord de France
Pawel Rózycki, University of IT and Management (UITM), Poland
Serap Sahin, Izmir Institute of Technology, Turkey
Ozgur Koray Sahingoz, Turkish Air Force Academy, Turkey
Manuel Filipe Santos, Universidade do Minho - Guimarães, Portugal
Ana Šaša, University of Ljubljana, Slovenia
Peter Schartner, Klagenfurt University, Austria
Isabel Seruca, Universidade Portucalense - Porto, Portugal
Marc Sevaux, Université de Bretagne-Sud, France
Qiao Shaojie, Southwest Jiaotong University, China
Ashok Sharma, TechMahindra, India
Mei-Ling Shyu, University of Miami - Coral Gables, USA
Patrick Siarry, Université Paris 12 (LiSSi) - Creteil, France
Sanjay Singh, Manipal University, India

Spiros Sirmakessis, Technological Educational Institution of Messolongi, Greece
Tomas Skersys, Kaunas University of Technology, Lithuania
Martin Stanton, Manchester Metropolitan University, UK
Kathryn E. Stecke, University of Texas at Dallas - Richardson, USA
Anca-Juliana Stoica, Uppsala University, Sweden
Renate Strazdina, Ernst&Young Baltic SIA, Latvia
Vadim Strijov, Computing Centre of the Russian Academy of Sciences, Russia
Weifeng Sun, Dalian University of Technology, China
Ryszard Tadeusiewicz, AGH University of Science and Technology, Poland
Mohamed Taleb, University of Quebec, Canada
Dan Tamir, Texas State University, Texas
John Terzakis, Intel, USA
Zenonas Theodosiou, Cyprus University of Technology, Cyprus
Ousmane Thiare, Gaston Berger University of Saint-Louis, Senegal
Maria Tortorella, University of Sannio, Italy
Guglielmo Trentin, National Research Council - Genoa & University of Turin, Italy
Chrisa Tsinaraki, European Commission - Joint Research Center (EU JRC) - Ispra, Italy
Meltem Sonmez Turan, NIST, USA
Ion Tutanescu, University of Pitesti, Romania
Theodoros Tzouramanis, University of the Aegean, Greece
Eleni I. Vlahogianni, National Technical University of Athens, Greece
Ante Vilenica, University of Hamburg, Germany
Patravadee Vongsumedh, Bangkok University, Thailand
Luc Vouligny, Institut de Recherche d'Hydro-Québec - Varennes, Canada
Mihaela Vranic, University of Zagreb, Croatia
Chieh-Yih Wan, Intel Corporation, USA
Yufeng Wang, University of South Florida, USA
Gerhard-Wilhelm Weber, METU - Ankara, Turkey
Dietmar Winkler, Vienna University of Technology, Austria
Viacheslav Wolfengagen, Institute "JurInfoR-MSU", Russia
Ouri Wolfson, University of Illinois, USA
Min Wu, Oracle Inc., USA
Mudasser F. Wyne, National University - San Diego, USA
Farouk Yalaoui, University of Technology of Troyes, France
Xia Yan, Hunan University, China
Chao-Tung Yang, Tunghai University, Taiwan
Yugang Yu, The University of Science and Technology of China, China
Fernando Zacarias Flores, Benemerita Universidad Autonoma de Puebla, Mexico
Marcelo Zanchetta do Nascimento, Federal University of Uberlândia (UFU), Brazil
Xuying Zhao, University of Notre Dame, USA
Zuqing Zhu, University of Science and Technology, China
Iveta Zolotová, Technical University of Kosice, Slovakia
Piotr Zwierzykowski, Poznan University of Technology, Poland

Copyright Information

For your reference, this is the text governing the copyright release for material published by IARIA.

The copyright release is a transfer of publication rights, which allows IARIA and its partners to drive the dissemination of the published material. This allows IARIA to give articles increased visibility via distribution, inclusion in libraries, and arrangements for submission to indexes.

I, the undersigned, declare that the article is original, and that I represent the authors of this article in the copyright release matters. If this work has been done as work-for-hire, I have obtained all necessary clearances to execute a copyright release. I hereby irrevocably transfer exclusive copyright for this material to IARIA. I give IARIA permission to reproduce the work in any media format such as, but not limited to, print, digital, or electronic. I give IARIA permission to distribute the materials without restriction to any institutions or individuals. I give IARIA permission to submit the work for inclusion in article repositories as IARIA sees fit.

I, the undersigned, declare that to the best of my knowledge, the article does not contain libelous or otherwise unlawful contents or invading the right of privacy or infringing on a proprietary right.

Following the copyright release, any circulated version of the article must bear the copyright notice and any header and footer information that IARIA applies to the published article.

IARIA grants royalty-free permission to the authors to disseminate the work, under the above provisions, for any academic, commercial, or industrial use. IARIA grants royalty-free permission to any individuals or institutions to make the article available electronically, online, or in print.

IARIA acknowledges that rights to any algorithm, process, procedure, apparatus, or articles of manufacture remain with the authors and their employers.

I, the undersigned, understand that IARIA will not be liable, in contract, tort (including, without limitation, negligence), pre-contract or other representations (other than fraudulent misrepresentations) or otherwise in connection with the publication of my work.

Exception to the above is made for work-for-hire performed while employed by the government. In that case, copyright to the material remains with the said government. The rightful owners (authors and government entity) grant unlimited and unrestricted permission to IARIA, IARIA's contractors, and IARIA's partners to further distribute the work.

Table of Contents

Four Scenarios of Effective Computations on Sum-like Graphs <i>Elena Ravve and Zeev Volkovich</i>	1
Topological Analysis of the Subway Network of Madrid <i>Mary Luz Mouronte</i>	9
QPSOL: Quantum Particle Swarm Optimization with Levy's Flight <i>Ennio Grasso and Claudio Borean</i>	14
Transient State Analysis of the Multichannel EMG Signal Using Hjorth's Parameters for Identification of Hand Movements <i>Michele Pla Mobarak, Juan Manuel Gutierrez Salgado, Roberto Munoz Guerrero, and Valerie Louis-Dorr</i>	24
Realising Duality Principle for Prognostic Models <i>Mohammad Samie, Amir Movahdi Saveh Motlagh, Alireza Alghassi, Suresh Perinpanayagam, and Epaminondas Kapetanios</i>	31
Design and Simulation of Electronic Service Business Process <i>Peteris Stipravietis, Edzus Zeiris, and Maris Zieme</i>	38
Development of Real-time Simulation Models: Integration with Enterprise Information Systems <i>Konstantin Aksyonov, Eugene Bykov, Olga Aksyonova, and Anna Antonova</i>	45
RISK-DET: ICT Security Awareness Aspect Combining Education and Cognitive Sciences <i>Schaff Guillaume, Harpes Carlo, Aubigny Matthieu, Junger Marianne, and Martin Romain</i>	51
ICT in Education: A New Paradigm and Old Obstacle <i>Gyorgy Molnar and Andras Benedek</i>	54
Instructional Approach in Adult Education using Mobile Devices <i>Felix Buendia, Angel Perles, and Juan Vicente Capella</i>	61
WebETL Tool – A Prototype in Action <i>Kornelije Rabuzin and Matija Novak</i>	67
Deductive Data Warehouses and Aggregate (Derived) Tables <i>Kornelije Rabuzin, Mirko Malekovic, and Mirko Cubrilo</i>	72
P2P Integration of Relational Knowledge Bases <i>Tadeusz Pankowski</i>	78

Generating Customized Sparse Eigenvalue Solutions with Lighthouse <i>Ramya Nair, Sa-Lin Bernstein, Elizabeth Jessup, and Boyana Norris</i>	85
Analysis of the Utilization of Web 2.0 Resources in Secondary Education and Advanced Vocational Training Studies <i>Federico Banda Sierra and Antonio J Reinoso</i>	90
Twitter Usage of German Online Retailers <i>Georg Lackermair and Daniel Kailer</i>	95
Encouraging Students to Document Software Development Projects using Blogs <i>Robert Law</i>	101
Ranking Domain Names Using Various Rating Methods <i>Kyriacos Talattinis, Christina Zervopoulou, and George Stephanides</i>	107
Combining Load Balancing with Energy Saving in a Cluster – Based P2P System <i>Minas Tasiou, Konstantinos Antonis, and Theofanis – Aristofanis Michail</i>	115
A Comparative Analysis of Parallel Programming Models for C++ <i>Arno Leist and Andrew Gilman</i>	121
A Study on the Performance Control for Building Facades Applied with Patterns of Hanok Openings <i>Seulki Kim, Kyoung-Hee Kim, and Seung-Hoon Han</i>	128
Performance Improvement in Applying Network Coding to On-demand Scheduling Algorithms for Broadcasts in Wireless Networks <i>G. G. Md. Nawaz Ali, Yuxuan Meng, Victor C. S. Lee, Kai Liu, and Edward Chan</i>	134
Combination of IMS-based IPTV Services with WebRTC <i>Tilmann Bach, Jens Zimmermann, Michael Maruschke, Kay Hansge, and Matthias Baumgart</i>	140
Aspect-Oriented Implementation of Concurrent Processing Design Patterns <i>Shingo Kameyama, Masatoshi Arai, Noriko Matsumoto, and Norihiko Yoshida</i>	146
Analysis of the Development Process of a Mutation Testing Tool for the C++ Language <i>Pedro Delgado-Perez, Inmaculada Medina-Bulo, and Juan Jose Dominguez-Jimenez</i>	151
A Negotiation Model Based on Multi-agent System under Cloud Computing <i>Jian Chen, Xinwei Han, and Guorui Jiang</i>	157
ICT Utilization in Libyan Universities: A Report on Case Study Research <i>Ali Bakeer and Martin Wynn</i>	165

Mutation Testing: Guideline and Mutation Operator Classification 171
Lorena Gutierrez-Madronal, Juan Jose Dominguez-Jimenez, and Inmaculada Medina-Bulo

MeterGoat: A Low Cost Hardware Platform for Teaching Smart Meter Security 180
Jefferson Capovilla, Nelson Uto, Danilo Suiama, and Jose Resende

Four Scenarios of Effective Computations on Sum-like Graphs

Elena V. Ravve and Zeev Volkovich

Department of Software Engineering
Ort Braude,
Karmiel, Israel
Email: {cselena, vlvolkov}@braude.ac.il

Abstract—In this paper, we consider computations on sum-like graphs, which we introduced in our previous works. For such graphs, we proposed a method that allows us to reduce the solution of a Monadic Second Order or First Order definable problem on the graphs to the solution of effectively derivable Monadic Second Order or First Order definable problems on their components, respectively. Now, we describe in great details four particular scenarios, where this method may be applied, and explain how it may improve the complexity of the solution. This lead us to a generalized formulation of Amdahl's style laws for each scenario. Moreover, we consider applications of our method to the fields of repetitive and hierarchical structures, widely used in hardware and software design, as well as parallel and distributed computations.

Keywords— Sum-like graphs; Translation schemes; (Weighted) Monadic Second Order Logic ; First Order Logic; Repetitive and hierarchical structures; Incremental computations; Parallel and distributed computations.

I. INTRODUCTION

Replacing of solution of a problem on a given graph by solution of other problems on derived graphs is widely used in different fields of science. In this paper, we address the case, when the original graph is *sum-like* and we extend the approach of [1], implemented in the framework of model checking. In [2], we describe in great details how the concept of *FSM* may be interpreted as a graph, and how computation of its properties may be expressed as a set of logical formulas. In [3], we investigate the case of computation with weighted automata on *sum-like labeled weighted trees*. In this paper, we are mostly concentrated on the complexity issues of the approach.

Assume we are given a model of an object in the form of graph G and a formalized presentation of a problem to be solved in the form of formula ϕ , possibly with free variables. Assume G is built from components G_i , where $i \in I$ is some index set or structure. Assume that G is a sum-like composition of G_i as defined in [1][2], or it is a sum-like labeled weighted tree, as defined in [3].

The main results of [1][2][3] show how the solution of ϕ on G depends on the components G_i of G and the index structure I . It is an extension of the Feferman–Vaught Theorem, cf. [4], for First Order Logic (*FOL*) to Monadic Second Order Logic (*MSOL*) or Weighted Monadic Second Order Logic

(*WMSOL*). The detailed description of algorithmic use of Feferman–Vaught Theorem may be found in [5].

The Feferman–Vaught theorem covers a very wide class of generalized products and sums of structures and is extremely powerful. We extend these theorems to the case of (Weighted) Monadic Second Order Logic and it works only for a more restricted class of *sum-like* graphs (trees), cf. [1][2][3]. From our main theorems, we derive a method for solving *MSOL* ϕ on sum-like G , which proceeds as follows (the similar treatment of *WMSOL* ϕ on sum-like T see in [3]):

Preprocessing: Given ϕ and Φ , but no G , we construct *at once* a sequence of formulas $\psi_{i,j}$ and a function $F_{\Phi,\phi}$. This construction is polynomial in the size of ϕ and Φ .

Incremental Computation: We compute the values $b_{i,j}$ defined by $b_{i,j} = 1$ iff $G_i \models \psi_{i,j}$.

Final Integration: Our theorems now state that $G \models \phi$ iff $F_{\Phi,\phi}(\vec{b}) = 1$.

In this paper, we investigate how the fact that G is built from several components may be used in order to make the solution more effective in the general case as well as how the general approach is connected to different practical applications. In fact, this analysis gives a generalized formulation of Amdahl's style laws for each described computation model and scenario.

We consider four different scenarios, when the complexity gain may be reached. The first one, described in Section V, represents the case when our graph G is composed by repetition of some basic component(s) $G_i = \tilde{G}$. We may find lots of such constructions, for example, in chemistry: polymers, VLSI design: adders, shifters, memories and other applications, etc.

Very often, after the problem ϕ was solved once on G , it should be solved again on some light modification of G . The situation is investigated in Section VI. More precisely, let j denote the solution on the j th variant of G . We assume that G^j differs from G^{j+1} in one component G_i^{j+1} , otherwise we may consider a line of such variants $G^{j=j_0}, G^{j_1}, \dots, G^{j_i=j+1}$. We consider two situations. The classical one looks at the cost of solving *once* the problem ϕ on G , and uses either the size of G , the size of ϕ or the sum of the two as the relevant input size. In addition, we ask, what our method can gain by repeating this process many times, with small changes at a time. For this purpose we also look at the size of the changed component and the number of iterations.

In Section VII, we consider the case, when computations may be done in parallel rather than sequentially on one computational unit. We show how our general approach leads us to the *BSP* computation model, as introduced in [6].

With the frequent use of the Internet, it becomes customary to have data distributed over many sites. We consider how our approach works in the case of the distributed databases in Section VIII. We show that under some reasonable conditions, our method leads to some variation of *LogP* model, cf. [7].

The paper is structured as follows:

- We start from a list of notations.
- In Section II, we give a motivating example.
- In Section III, we recall general definitions and results taken almost verbatim from [2].
- Section IV provides detailed discussion of the common basis of all cost evaluations.
- In section V, we analyze complexity of single computations on repetitive structures.
- In section VI, we analyze complexity of incremental recomputations on the lightly modified graphs.
- In section VII, we analyze complexity of parallel computations.
- In section VIII, we analyze complexity of computations on distributed databases.
- Section IX summarizes the paper.

LIST OF NOTATIONS

$\pi_A R$	Projection of attributes A from relation R
$\sigma_\theta R$	Selection from R of tuples satisfying θ
$R \bowtie S$	Join of relations R and S
$b_{i,j}$	Boolean values
\vec{b}	Vector of Boolean values
<i>BSP</i>	Bulk synchronous parallel model
E_1, E_2, E_i, \dots	Sets of edges of a graph
$F_{\Phi, \phi}$	Computation, associated with Φ and ϕ
<i>FOL</i>	First Order Logic
<i>FSM</i>	Finite State Machine
$G, G_\Phi, \tilde{G}, G_i, \dots$	Graph structures
I	Index structure
$I(\mathbf{R}) \dots$	Instances of Database schemes
\mathcal{L}	Logic
<i>MSOL</i>	Monadic Second Order Logic
$P_1, P_2, P_i, Q_1, Q_2, Q_i$	One place relations on the set of vertices
$\mathbf{R}, \mathbf{R}_I, \dots$	Database schemes
R_i, \dots	Relation symbols
$\rho(R_i)$	Arity of R_i
<i>SOL</i>	Second Order Logic
Φ	Translation scheme
$\Phi^*, \Phi^\#$	Two mappings of Φ
T	Sum-like tree
v	Vertex of a graph
\vec{v}	Vector of vertices of a graph
$V, V_1, V_2, V_i, V_\Phi, \dots$	Sets of vertices of a graph
<i>VLSI</i>	Very Large Scale Integration
<i>WMSOL</i>	Weighted Monadic Second Order Logic

II. MOTIVATING EXAMPLE

In this section, we consider how verification of a *MSOL*-property over some composition of graphs can be reduced to

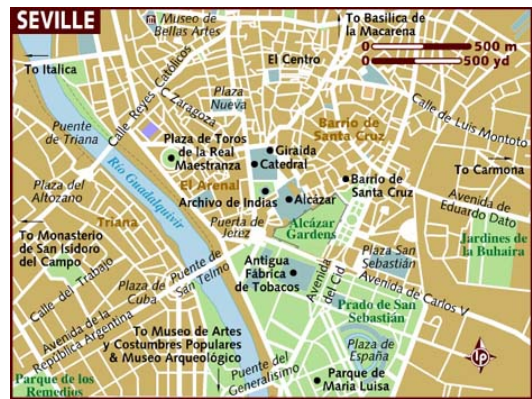


Fig. 1. A city with bridges [8].

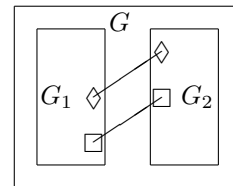


Fig. 2. Composition of two graphs: $G_1 \equiv G_2$.

its components. Assume we consider a city that is divided into two parts by a river, see Figure 1 taken from [8]. These parts are connected by bridges. Without loss of generality, let us assume that there exist only two bridges. Assume we are looking for one-way cycles in the city.

We may formulate the situation in the following way. We are given two finite graphs $G_1 = \langle V_1, E_1, P_1, Q_1 \rangle$ and $G_2 = \langle V_2, E_2, P_2, Q_2 \rangle$, where V_i denotes a set of vertices, E_i denotes a set of edges and P_i, Q_i are one place relations (vertex colorings) respectively. Let G be the disjoint union of G_1 and G_2 with additional edges, corresponding to the bridges, see Figure 2. We define this composition of two colored graphs formally as follows: $G = G_1 \equiv G_2 = \langle V_1 \dot{\cup} V_2, E \rangle$, where $V_1 \dot{\cup} V_2$ denotes disjoint union of sets of vertices, and two vertices v and u of G belongs to E iff

$$\begin{aligned} \Phi : (v, u) \in E_1 \vee (v, u) \in E_2 \\ \vee (v \in Q_1 \wedge u \in Q_2) \vee (v \in Q_2 \wedge u \in Q_1) \\ \vee (v \in P_1 \wedge u \in P_2) \vee (v \in P_2 \wedge u \in P_1) \end{aligned}$$

We want to check whether G has cycles. To do so, we observe that

- (*) G has a cycle iff G_1 has a cycle, or G_2 has a cycle, or there are at least two connected coloured vertices in G_{2-i} and at least one coloured vertex in the same color vertex in G_{i+1} , where $i \in \{0, 1\}$,

and proceed as follows.

- We first write the cyclicity property as a formula ϕ in *MSOL*.
- Then, using (*), which depends only on ϕ and Φ , but not G , we look for formulas $\psi_{1,1}, \dots, \psi_{1,n_1}$ and

$\psi_{2,1}, \dots, \psi_{2,n_2}$ in *MSOL*, which will give us the properties to be checked in G_1 and G_2 respectively.

- Then, again using (*), we look for a boolean function F of $n_1 + n_2$ arguments $b_{1,1}, \dots, b_{2,n_2}$.
- Now we put $b_{i,j} = 1$ iff $G_i \models \psi_{i,j}$ and hope to conclude that $G \models \phi$ iff $F(b_{1,1}, \dots, b_{2,n_2}) = 1$.

Surprisingly, our main theorems from [1][2][3] imply that this method can be mechanized in certain cost, even if (*) is not given in advance.

We have explained our main result by using a very simple example of a disjoint union of two colored graphs with some edges added, defined by Φ . The main theorems of [1][2][3] generalize this approach to combination of more than two structures and more complicated additional relations.

III. GENERAL BACKGROUND

In this section, we recall general definitions and results taken almost verbatim from [2]. The corresponding extension to the case of *WMSOL* may be found in [3].

A. Translation schemes

In this section, we follow [1] and introduce the general framework for syntactically defined translation schemes. A vocabulary is a finite set of relation symbols and constants.

Definition 1: General Translation Schemes.

Let τ and σ be two vocabularies and \mathcal{L} be a logic, such as *FOL* or *MSOL*. Let $\sigma = \{R_1, \dots, R_m\}$ and let $\rho(R_i)$ be the arity of R_i . Let $\Phi = \langle \phi, \psi_1, \dots, \psi_m \rangle$ be formulas of $\mathcal{L}(\tau)$. Φ is *feasible for σ over τ* if ϕ has exactly 1 free first order variable and each ψ_i has $\rho(R_i)$ distinct free first order variables. Such a $\Phi = \langle \phi, \psi_1, \dots, \psi_m \rangle$ is also called a τ - σ -translation scheme or, shortly, a *translation scheme*, if the parameters are clear in the context.

With a translation scheme Φ we can naturally associate a (partial) function Φ^* from τ -structures (graphs) to σ -structures (graphs).

Definition 2: The induced map Φ^* .

Let G be a τ -graph and Φ be feasible for σ over τ . The graph G_Φ is defined as follows:

- 1) The universe V_Φ of G_Φ is the set

$$V_\Phi = \{v \in V : G \models \phi(v)\};$$

- 2) The interpretation of R_i in G_Φ is the set

$$G_\Phi(R_i) = \{\bar{v} \in G_\Phi^{\rho(R_i)} : G \models \psi_i(\bar{v})\};$$

Note that G_Φ is a σ -graph of cardinality at most $|G|$.

- 3) The partial function $\Phi^* : \mathcal{G}(\tau) \rightarrow \mathcal{G}(\sigma)$ is defined by $\Phi^*(G) = G_\Phi$. Note that $\Phi^*(G)$ is defined iff $G \models \exists v \phi$.

With a translation scheme Φ we can also naturally associate a function $\Phi^\#$ from *MSOL*(σ)-formulas to *MSOL*(τ)-formulas.

Definition 3: The induced map $\Phi^\#$.

Let θ be a σ -formula and Φ be feasible for σ over τ . The formula θ_Φ is defined inductively as follows:

- 1) For $R_i \in \sigma$ and $\theta = R_i(x_1, \dots, x_m)$, we put $\theta_\Phi = \psi_i(v_1, \dots, v_m)$.

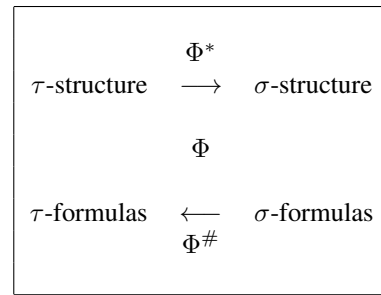


Fig. 3. Translation scheme and its components.

- 2) For the boolean connectives, the translation distributes, i.e., if $\theta = (\theta_1 \vee \theta_2)$ then $\theta_\Phi = (\theta_{1\Phi} \vee \theta_{2\Phi})$ and if $\theta = \neg\theta_1$ then $\theta_\Phi = \neg\theta_{1\Phi}$, and similarly for \wedge .
- 3) For the existential quantifier, we use relativization, i.e., if $\theta = \exists v \theta_1$, we put $\theta_\Phi = \exists v(\phi(v) \wedge (\theta_1)_\Phi)$.
- 4) For second order variables U of arity ℓ and u a vector of length ℓ of first order variables or constants we translate $\theta = \exists U \theta_1$, by treating U like a relation symbol and put $\theta_\Phi = \exists U(\forall u(U(u) \rightarrow (\phi(u_1) \wedge \dots \wedge \phi(u_\ell) \wedge (\theta_1)_\Phi)))$.
- 5) The function $\Phi^\# : MSOL(\sigma) \rightarrow MSOL(\tau)$ is defined by $\Phi^\#(\theta) = \theta_\Phi$.

The following fact holds, see Figure 3:

Proposition 1:

Let $\Phi = \langle \phi, \psi_1, \dots, \psi_m \rangle$ be a τ - σ -translation scheme, G a τ -graph and θ a σ -formula. Then $G \models \Phi^\#(\theta)$ iff $\Phi^*(G) \models \theta$. The proof may be found in [9][10].

B. Sum-like graphs

In this section, we discuss ways of obtaining graphs from components. The *Disjoint Union* of a family of graphs is the simplest example of juxtaposing graphs, where none of the components are linked to each other. For our purpose, we include the index set I in the resulting structure as well.

Definition 4: Disjoint Union of Graphs

Let $\tau_i = \langle R_1^i, \dots, R_{j^i}^i \rangle$ be a vocabulary of graph G_i . In the general case the resulting graph $G = \bigsqcup_{i \in I} G_i$ is $G = \langle \mathcal{V} = I \cup \bigcup_{i \in I} \mathcal{V}_i, R_j^i(1 \leq j \leq j^i), R_{j^i}^i(i \in I, 1 \leq j^i \leq j^i) \rangle$ for all $i \in I$, or rather any graph, isomorphic to it.

We assume existence of the following mappings:

- $h_\nu : \mathcal{V} \rightarrow I, h_\nu(v) = i$ if $v \in \mathcal{V}_i$;
- $h_\nu : PS(\mathcal{V}) \rightarrow PS(\mathcal{V}_i), h_\nu(V) = V_i$ if V_i is a i^{th} component of \mathcal{V} , while PS denotes the power set.

Definition 5: Partitioned Index Structure

Let I be an index structure. I is called *finitely partitioned* into ℓ parts if there are unary predicates $I_\alpha, \alpha < \ell$, in the vocabulary of I such that their interpretation forms a partition of the universe of I .

The following holds:

Theorem 1:

Let I be a finitely partitioned index structure.

Let $G = \bigsqcup_{i \in I} G_i$ be a τ -graph, where each G_i is isomorphic to some B_1, \dots, B_ℓ over the vocabularies τ_1, \dots, τ_ℓ , in accor-

dance to the partition (ℓ is the number of the classes).

For every $\phi \in MSOL(\tau)$ there are:

- a boolean function $F_\phi(b_{1,1}, \dots, b_{1,j_1}, \dots, b_{\ell,1}, \dots, b_{\ell,j_\ell}, b_{I,1}, \dots, b_{I,j_I})$
- $MSOL$ -formulas $\psi_{1,1}, \dots, \psi_{1,j_1}, \dots, \psi_{\ell,1}, \dots, \psi_{\ell,j_\ell}$
- $MSOL$ -formulas $\psi_{I,1}, \dots, \psi_{I,j_I}$

such that for every G, I and B_i as above with

$$B_i \models \psi_{i,j} \text{ iff } b_{i,j} = 1 \text{ and } I \models \psi_{I,j} \text{ iff } b_{I,j} = 1$$

we have

$$G \models \phi \text{ iff}$$

$$F_\phi(b_{1,1}, \dots, b_{1,j_1}, \dots, b_{\ell,1}, \dots, b_{\ell,j_\ell}, b_{I,1}, \dots, b_{I,j_I}) = 1.$$

Moreover, F_ϕ and the $\psi_{i,j}$ are computable from ϕ , ℓ and vocabularies alone, but are exponential in the quantifier depth of ϕ .

The disjoint union as such is not very interesting. However, combining it with translation schemes gives us a rich repertoire of patching techniques. Let τ_0, τ_1, τ be finite vocabularies of graphs. For a τ_0 -model I (serving as index model), τ_1 -graphs are pairwise disjoint for simplicity $G_i (i \in I)$ and a τ -graph G is the disjoint union of $\langle G_i : i \in I \rangle$ with $G = \bigsqcup_{i \in I} G_i$. Now we generalize the disjoint union of graphs to *sum-like* graphs in the following way:

Definition 6: Sum-like Graphs

Let I be a finitely partitioned index structure.

Let $G = \bigsqcup_{i \in I} G_i$ be a τ -graph, where each G_i is isomorphic to some B_1, \dots, B_ℓ over the vocabularies τ_1, \dots, τ_ℓ , in accordance with the partition. Furthermore let Φ be a τ - σ $MSOL$ -translation scheme. The Φ -sum of B_1, \dots, B_ℓ over I is the graph $\Phi^*(G)$, or rather any graph isomorphic to it.

Theorem 2:

Let I be a finitely partitioned index structure and let G be the Φ -sum of B_1, \dots, B_ℓ over I , as above.

For every $\phi \in MSOL(\tau)$ there are:

- a boolean function $F_{\Phi,\phi}(b_{1,1}, \dots, b_{1,j_1}, \dots, b_{\ell,1}, \dots, b_{\ell,j_\ell}, b_{I,1}, \dots, b_{I,j_I})$
- $MSOL$ -formulas $\psi_{1,1}, \dots, \psi_{1,j_1}, \dots, \psi_{\ell,1}, \dots, \psi_{\ell,j_\ell}$
- $MSOL$ -formulas $\psi_{I,1}, \dots, \psi_{I,j_I}$

such that for every G, I and B_i as above with

$$B_i \models \psi_{i,j} \text{ iff } b_{i,j} = 1 \text{ and } I \models \psi_{I,j} \text{ iff } b_{I,j} = 1$$

we have

$$G \models \phi \text{ iff}$$

$$F_{\Phi,\phi}(b_{1,1}, \dots, b_{1,j_1}, \dots, b_{\ell,1}, \dots, b_{\ell,j_\ell}, b_{I,1}, \dots, b_{I,j_I}) = 1.$$

Moreover, $F_{\Phi,\phi}$ and the $\psi_{i,j}$ are computable from $\Phi^\#$ and ϕ , but are exponential in the quantifier depth of ϕ .

Moreover, in [3], we prove that:

Theorem 3:

Let I be a finitely partitioned index structure and a tree is the Φ -sum. For every $\phi \in WMSOL(\tau)$ there are:

- a computation on values $\varpi_{1,1}, \dots, \varpi_{\ell,j_\ell}$

$$F_{\Phi,\phi}(\varpi_{1,1}, \dots, \varpi_{1,j_1}, \dots, \varpi_{\ell,1}, \dots, \varpi_{\ell,j_\ell})$$

and

- $WMSOL$ -formulas $\psi_{1,1}, \dots, \psi_{1,j_1}, \dots, \psi_{\ell,1}, \dots, \psi_{\ell,j_\ell}$ such that $\varpi_{i,j} = \varrho_{i,j}$ iff $[\psi_{i,j}] = \varrho_{i,j}$ we have

$$[\phi] = \varrho \text{ iff } F_{\Phi,\phi}(\varpi_{1,1}, \dots, \varpi_{1,j_1}, \dots, \varpi_{\ell,1}, \dots, \varpi_{\ell,j_\ell}) = \varrho.$$

Moreover, $F_{\Phi,\phi}$ and the $\psi_{i,j}$ are computable from $\Phi^\#$ and ϕ , but are exponential in the quantifier depth of ϕ .

IV. GENERAL COMPLEXITY ANALYSIS

In this section, we discuss under what conditions theorems of [1][2][3] improve the complexity of computations, when measured in the size of the composed graphs (trees) only. Our scenarios are as follows: A (W) $MSOL$ formula (set of formulas) ϕ is given in advance. A sum-like graph (tree) is now submitted to a computation unit and we want to know, how long it takes to check whether ϕ is true on the graph (tree). Now we give the general complexity analysis of the computation on sum-like graphs (trees).

A. Complexity of computation for different logics

Theorems of [1][2] hold for $MSOL$ and, with restrictions, also for FOL . Computation for FOL is polynomial (even in logarithmic space), whereas computation for $MSOL$ is likely to be non-polynomial, as it sits fully in the polynomial hierarchy. Theorems of [3] hold for $WMSOL$. Computation for $WMSOL$ may be done, using Weighted Tree Automata.

More precisely, the complexity of computation (in the size of the graph) of Second Order Logic expressible properties can be described as follows. The class NP of non-deterministic polynomial-time problems is the set of properties, which are expressible by Existential Second Order Logic on finite structures, cf. [11]. Computation for SOL definable properties is in the polynomial hierarchy, cf. [12]. Moreover, for every level of the polynomial hierarchy there is a problem, expressible in SOL , that belongs to this class. The same fact hold for $MSOL$, too, as observed in [13].

Computation for properties, definable in Fixed Point Logic, is polynomial, cf. [14]. CTL^* is a superset of Computational Tree Logic and Linear Temporal Logic. All the problems, which are expressible by CTL^* , can be computed in polynomial time, cf. [15]. Most properties, which appears in real life applications, are stronger than FOL but weaker than $MSOL$, and their computational complexity is polynomial. In [16], it was shown that the similar theorems are valid for Transitive Closure and Monadic Fixed Point Logic. However, it does not hold for all of the languages (logics).

B. General analysis

Assume that G is a sum-like graph (or a sum-like tree). Its components are G_i with index structure I , and we want to check whether ϕ is true in G . Assume that:

- $\mathcal{T}(N)$ or $\mathcal{T}_{old}(N)$ denotes time to solve the problem by the traditional sequential way (N denotes the size of the coding of graph G);
- $\mathcal{E}_{\mathcal{I}}$ denotes time to extract index structure I from G ;
- \mathcal{E}_i denotes time to extract each G_i from G ;

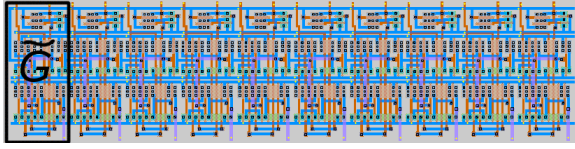


Fig. 4. Layout of a full 10-bit adder [22].

- $\mathcal{C}_{\mathcal{I}}(n_I)$ denotes time to compute all values of $b_{I,j}$, where n_I is the size of I ;
- $\mathcal{C}_i(n_i)$ denotes time to compute all values of $b_{i,j}$, where n_i is the size of G_i ;
- $\mathcal{T}_{F_{\Phi,\phi}}$ denotes time to build $F_{\Phi,\phi}$;
- $\mathcal{T}_{\mathcal{S}}$ denotes time to achieve one result of $F_{\Phi,\phi}$.

According to these symbols, the new computation time is:

$$\mathcal{T}_{new} = \mathcal{E}_{\mathcal{I}} + \sum_{i \in I} \mathcal{E}_i + \mathcal{C}_{\mathcal{I}} + \sum_{i \in I} \mathcal{C}_i + \mathcal{T}_{F_{\Phi,\phi}} + \mathcal{T}_{\mathcal{S}} \quad (1)$$

and the question to answer is: when $\mathcal{T}_{old} > \mathcal{T}_{new}$.

V. SCENARIO A: SINGLE COMPUTATION ON REPETITIVE STRUCTURES

In this section, we consider the underlying structure, the formula and the modularity as well as possible applications of our method for single computations on repetitive structures. We analyze the corresponding complexity gain for the computations of properties, expressible in different logics.

A. The underlying structure, the formula and the modularity

Our underlying structure is a sum-like graph G , where for each v : $G_v = \tilde{G}$. The property, which we want to check on it, is expressible in formula ϕ of *MSOL*, which has an exponential checker. We check whether $G \models \phi$.

B. Applications

Different repetitive combinations of graphs have been under deep investigation for long time, cf. [17][18][19]. Polygraphs were introduced as generalization of polymers in chemistry, cf. [20], and explored in VLSI design, cf. [21].

We restrict ourselves to VLSI design, which widely uses repetition of blocks (Figure 4 taken from [22]) and hierarchical structures. Many basic elements of the design, such that shifters, adders, registers, etc. are build in this manner. Repetition of modules is explored also in control logic and other kinds of hardware design. Memory is another example of a repetitive structure in VLSI design. If we go up in the hierarchy of our design, we found multi-core processors, which are single computing components with two or more independent processors (called "cores").

C. Complexity gain for *MSOL*

Assume that our design is presented as a graph (*FSM*), such that:

- N is a size of G , n is a size of \tilde{G} and l is a size of index structure I .
- The decomposition is given: $\mathcal{E}_{\mathcal{I}} = \mathcal{E}_i = 0$.
- The computation is exponential in the form: $\mathcal{T} = e^{g(x)}$.

In this case $\mathcal{T}_{new} = P^p(\mathcal{T}(n), \mathcal{T}(l))$, where P^p denotes polynomial of degree p , and $\mathcal{T}_{old} = \mathcal{T}(l \cdot n)$.

The question to answer is: when $f(n \cdot l) > P^p(f(n), f(l))$. According to our assumptions, we obtain that the comparison of the computation times in (1) looks like:

$$e^{g(n \cdot l)} > a_p(e^{p \cdot g(n)} + e^{p \cdot g(l)}).$$

Assume that $n = l$. Then $g(n^2) > p \cdot g(n) + ln2 + ln(a_p)$. Assume that $g(x) = ln^2(x)$, then $f(x) = x^{ln(x)}$. In this case we obtain that (1) is transformed to: $ln^2(n^2) > p \cdot ln^2(n) + ln2 + ln(a_p)$ or $ln^2(n) > \frac{ln(2 \cdot a_p)}{2-p}$.

D. Complexity gain for other logics

If we use *FOL* or another logic, where the computational procedure is polynomial in the sizes of G and in I and each G_i too, then we do not obtain any time gain.

VI. SCENARIO B: INCREMENTAL RE-COMPUTATIONS

In this section, we consider the underlying structure, the formula and the modularity as well as possible applications of our method for incremental re-computations. We analyze the corresponding complexity gain for the computations of properties, expressible in different logics.

A. The underlying structure, the formula and the modularity

Our underlying structure is a sum-like graph, and the property, which we want to check on it, is expressible in a formula ϕ that has a polynomial checker. Assume that we change several times (let us denote the number of the times by ς) some fixed component \tilde{G} of G . We check each time whether $G \models \phi$.

B. Applications

We restrict ourselves to the case of hardware verification, where we find the following situation: We are given a mathematical model of a device in form of a finite relational structure G (*FSM* or Kripke model) and a formalized property ϕ . Usually ϕ is given in advance and G is being built with the aim to satisfy ϕ . Checking whether ϕ holds in G is to be atomized. This process is called model checking. The literature is rich in papers addressing this problem, cf. [15].

As a rule, hardware design G is built from components (modules) G_i , where $i \in I$. The modules are the building blocks, installing one in other that gives the design hierarchy. The communication between a module and its environment is executed using ports. All but the top-level modules in a hierarchy have ports. In the process of building G , several candidate structures G^j have to be checked for ϕ , where j denotes the j^{th} attempt of designing G . Often G^j differs from G^{j+1} in one component G_i^{j+1} . It is easy to see (Figure 2) that our motivating example exactly fails in this framework, when we are talking about hardware design rather than about city maps. Combination of graphs in terms of graph grammars may be found also in [23].

C. Complexity gain for FOL

Let \mathcal{T}_{old} be time to solve the given problem by the traditionally applied way. It should be clear that $\mathcal{T}_{old} = \varsigma \cdot \mathcal{T}(N)$. Let \mathcal{T}_{new} be time to solve the same problem, when structure G is viewed as a generalized sum. It is easy to see that

$$\mathcal{T}_{new}(N, n) = \mathcal{T}(N - n) + \varsigma \cdot \mathcal{T}(n) + \mathcal{T}_{F_{\Phi, \phi}} + \varsigma \cdot \mathcal{T}_S.$$

The question to answer is: which value of n provides that $\mathcal{T}_{old} > \mathcal{T}_{new}$. Assume that $\mathcal{T}(x) = x^2$, then (1) becomes to be:

$$\varsigma \cdot N^2 > (N - n)^2 + \varsigma \cdot n^2 + \mathcal{T}_{F_{\Phi, \phi}} + \varsigma \cdot \mathcal{T}_S$$

$$N^2 - 2 \cdot n \cdot N + n^2(\varsigma + 1) + \mathcal{T}_{F_{\Phi, \phi}} + \varsigma \cdot \mathcal{T}_S - \varsigma \cdot N^2 < 0$$

$$n_{1,2} = \frac{N \pm \sqrt{N^2 + (\varsigma + 1)(N^2(\varsigma - 1) - \mathcal{T}_{F_{\Phi, \phi}} - \varsigma \cdot \mathcal{T}_S)}}{\varsigma + 1}.$$

If $n_1 \leq n \leq n_2$ then $\mathcal{T}_{old} > \mathcal{T}_{new}$.

$$n_2 = \frac{N + \sqrt{\varsigma^2(N^2 - \mathcal{T}_S) - \varsigma(\mathcal{T}_S + \mathcal{T}_{F_{\Phi, \phi}}) - \mathcal{T}_{F_{\Phi, \phi}}}}{\varsigma + 1}$$

$$\lim_{\varsigma \rightarrow \infty} n_2 = \sqrt{N^2 - \mathcal{T}_S}.$$

The same consideration can be done for other polynomial dependencies $\mathcal{T}(x)$ for FOL-definable logics.

D. Complexity gain for other logics

Let \mathcal{L} be any proper sub-logic of MSOL stronger than FOL. Our theorems do not hold in the following: if we apply it then $\psi_{i,j}$ are not necessary in \mathcal{L} .

VII. SCENARIO C: PARALLEL COMPUTATIONS

In this section, we consider the underlying structure, the formula and the modularity as well as possible applications of our method for parallel computations. We analyze the corresponding complexity gain for the computations of properties, expressible in different logics.

A. The underlying structure, the formula and the modularity

Our underlying structure is a sum-like graph G , and the property, which we want to check on the structure, is expressible in formula ϕ of a logic. We check whether $G \models \phi$.

B. Applications

Let us consider the following composition of two input graphs H and G . G can be viewed as a display graph, where on each node we want to have a copy of H , such that certain additional edges are added. In practice, this is a model on how a pipeline works. The nodes marked with L^j are the latches.

Let $G = \langle V_G, R \rangle$ and $H = \langle V_H, S, L^j (j \in J) \rangle$ be two relational structures (J is finite), then their composition $C = \langle V_C, L_C^1, \dots, L_C^{|J|}, S_C, R_C^j (j \in J) \rangle$ is defined as follows, see Figure 5:

- $V_C = \dot{\bigcup}_{g \in G} V_H^g$, where each V_H^g is isomorphic to V_H ;
- $L_C^j(w)$ is true if w belongs to L^j ;
- $S_C = \{(w, v) : w \in V_H^g, v \in V_H^g, S(w, v)\}$;
- $R_C^j = \{(w, v) : L^j(w), L^j(v), P(i, w), P(i', v), R(i, i')\}$.

C can be obtained from the disjoint union $\bigsqcup_{g \in G} H$ by a FOL translation scheme. In this example, depending on the choice

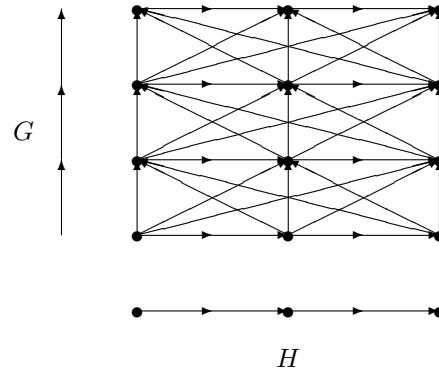


Fig. 5. Uniform graph substitution.

of the interpretation of the L^j 's, more sophisticated parallel computations can be modeled, but not all.

C. Complexity gain

In (1), the new computation time is calculated as:

$$\mathcal{T}_{new} = \mathcal{E}_I + \sum_{i \in I} \mathcal{E}_i + \mathcal{C}_I + \sum_{i \in I} \mathcal{C}_i + \mathcal{T}_{F_{\Phi, \phi}} + \mathcal{T}_S$$

for the case, when all the computations are done sequentially on a single computational unit. In fact, now even personal computers and smartphones have several cores. In this case, the computation may be done in the following way (we assume that there exist enough computational units for total parallelism) :

Extraction Super Step: The extraction of the index structure I from G and each G_i from G may be done in parallel as well as the building of $F_{\Phi, \phi}$. We denote the extraction time by: $\mathcal{E} = \max\{\mathcal{E}_I, \max_{i \in I} \{\mathcal{E}_i\}, \mathcal{T}_{F_{\Phi, \phi}}\}$.

Computational Super Step: The computation of all values of $b_{i,j}$ and $b_{i,j}$ may be done in parallel as well. We denote by $\mathcal{C} = \max\{\mathcal{C}_I(n_I), \max_{i \in I} \{\mathcal{C}_i(n_i)\}\}$. In fact, at this step, even more parallelism may be reached if we compute all $b_{i,j}$ in parallel.

Final Proceeding: \mathcal{T}_S still denotes time to search one result of $F_{\Phi, \phi}$. The new computation time for the case of full parallelism is: $\mathcal{T}_{new}^{BSP} = \mathcal{E} + \mathcal{C} + \mathcal{T}_S$. The computation model fails in the general framework of BSP, cf. [6].

D. Complexity gain for MSOL

The computation is exponential in the form: $\mathcal{T} = e^{g(x)}$. In this case $\mathcal{T}_{old} = \mathcal{T} = f(N) = e^{g(N)}$ and $\mathcal{T}_{new}^{BSP} = \mathcal{E} + PP(e^{g(\frac{N}{k})}) + \mathcal{T}_S$ and the question to answer is: when $f(n \cdot k) > PP(f(n))$. According to our assumptions, we obtain:

$$e^{g(n \cdot k)} > \mathcal{E} + a_p \cdot e^{p \cdot g(n)} + \mathcal{T}_S.$$

Assume that $k = n$, it means that there exist enough computation units for full parallelization. In this case, the condition of the effective computation looks like:

$$e^{g(n^2)} > \mathcal{E} + a_p \cdot e^{p \cdot g(n)} + \mathcal{T}_S.$$

E. Complexity gain for FOL and other logics with polynomial checkers

Assume that again $\mathcal{T}(x) = x^2$, and each G_i are of the same size $\frac{N}{k}$ then:

$$\mathcal{T}_{old} = N^2 \text{ and } \mathcal{T}_{new}^{BSP} = \mathcal{E} + \left(\frac{N}{k}\right)^2 + \mathcal{T}_S$$

$$N^2 > \mathcal{E} + \left(\frac{N}{k}\right)^2 + \mathcal{T}_S ; N^2 - \left(\frac{N}{k}\right)^2 > \mathcal{E} + \mathcal{T}_S$$

Now the condition of the effective computation looks like:

$$N^2 \cdot \frac{(k^2-1)}{k^2} > \mathcal{E} + \mathcal{T}_S$$

Complexity consideration for other logics \mathcal{L} , which are proper sublogics of *MSOL* stronger than First Order Logic, are similar to the given in subsection VI-D.

VIII. SCENARIO D: PARALLEL COMPUTATIONS ON DISTRIBUTED DATA

In this section, we consider the underlying structure, the formula and the modularity as well as possible applications of our method for parallel computations on distributed databases. We analyze the corresponding complexity gain for the computations of queries, expressible in different logics.

A. The underlying structure, the formula and the modularity

Our underlying structure is a sum-like graph G that is stored in the distributed way: each G_i is stored in the i^{th} site. The property, which we want to check on the structure, is expressible in formula ϕ of a logic. We check whether $G \models \phi$.

B. Applications

We restrict ourselves to investigation of distributed databases. In this case, a user sees one (virtual) database instance over a fixed database scheme. Queries and updates are submitted to a central processing site which will compute the required view or transaction by distributing the appropriate sub-tasks among the different sites.

While Datalog is hard to check, cf. [14], in [24], a new variant of Datalog was introduced: Datalog LITE. On the one hand, the deductive query language has linear time model checking. On the other hand, it encompasses modal and temporal logics, such as *CTL* and alternation-free μ -calculus. Moreover, it was shown that Datalog LITE with only unary and binary input predicates is contained in *MSOL*.

Assume we are given a database scheme that contains four relations: $\mathbf{R} = (R_1, R_2, R_3, R_4)$. Assume that we want to define a view that is derived from the database by applying the following query, given in the format of relational algebra: $(\pi_A R_1 \cup R_2) \bowtie (R_3 - \sigma_\xi R_4)$. In this case, the corresponding translation scheme is: $\Phi_{View} = \langle x = x, \phi_{View} \rangle$, where $\phi_{View} = (\pi_A R_1 \cup R_2) \bowtie (R_3 - \sigma_\xi R_4)$.

Let \mathbf{R}_I be an index scheme with finite domain and $|I| = n$ (to simplify the example, let $n=2$) and let $\mathbf{R}_1 = (R(y_1, \dots, y_{r_1}))$, $\mathbf{R}_2 = (R(y_1, \dots, y_{r_2}))$ be database schemes $\mathbf{R}_i (i \in \{1, 2\})$, where $r_j (j \in \{1, 2\})$ be an arity of the corresponding relation. Let $\mathbf{R} = \bigsqcup_{i \in I} \mathbf{R}_i = (P(\iota, x), I(x), R^1, R^2)$ be a \mathbf{R} database scheme, which is the disjoint union of \mathbf{R}_i 's.

We define the following translation scheme Φ_{Join} from the \mathbf{R} -instances to \mathbf{S} -instances, where $\mathbf{S} = (S)$.

$$\Phi_{Join} = \langle y \approx y, \exists y_{r_1}^1 \exists y_{r_1}^2 (R^1(y_1^1, \dots, y_{r_1}^1) \wedge R^1(y_1^2, \dots, y_{r_1}^2) \wedge \phi = (y_{r_1}^1, y_{r_1}^2)) \rangle$$

Assume we are given a sum-like database and we want to compute the view (query), defined by ϕ . Now, given a tuple t over $I(\mathbf{R})$, in order to check whether t belongs to the view, defined by ϕ , we compute the following:

- 1) $View_j^i = \{t \in I_i(\mathbf{R}) : \psi_j\}$, the views at site i , defined by the queries ψ_j .
- 2) $X_j^t = \{i \in I : t \in View_j^i\}$, the set of sites where t belongs to the view, defined by ψ_j ;
- 3) The truth value of $I(\mathbf{R}_I) \models \psi_I(X_1^t, \dots, X_n^t)$.

By our theorems: tuple t belongs to the view, defined by ϕ , iff $I(\mathbf{R}_I) \models \psi_I(X_1^t, \dots, X_n^t)$. In other words

$$\{t : I(\mathbf{R}) \models \phi\} = \{t : I(\mathbf{R}_I) \models \psi_I(X_1^t, \dots, X_n^t)\}.$$

We see that $View_j^i$ is computed at the site i and only the queries ψ_j have to be sent over the net. The $View_j^i$'s can be computed in parallel. Furthermore, when we compute X_j^t , only the tuple t is sent over the net. Finally, evaluating ψ_I can be done in *PSpace* (in the size of I). However, it is likely that for special cases of ψ_I (when ϕ is a pure Datalog query) the complexity (in the size of I) becomes (at least non-deterministically) polynomial.

In distributed databases, where the data are measured in gigabytes (terabytes) and the size of I is a small finite set (say ≤ 100), this method, in spite of its considerable overhead, gives a considerable improvement over any other method, which moves parts of the databases over the net.

Further gains can be achieved by exploiting a hierarchical structure of the way the databases are distributed: we introduce virtual sites (gather sites), which gather the data of its sub-sites and make sure that the number of sub-sites remains bounded (say ≤ 10). At each gathering site, the evaluation of ψ_I remains thus feasible.

As it was shown in [25], our method generalizes the propagation technique from [26] for relational algebra and the incremental re-computation technique from [27] for some kinds of Datalog programs to cases of definable sets of tuples to be deleted or inserted.

In addition, assume that the query language allows us to ask optimization questions. In such cases, our generalized propagation technique may be directly extended to the case of incremental optimization as considered in [28][29][30]. Using our method, the final optimal result is computed from the local (not necessarily optimal) solutions as explained in [3]. Moreover, our approach is directly connected to Parallel Distributed Genetic Programming as introduced in [31].

C. Complexity gain

The full computation process is composed now from the following steps (the above $\mathcal{E} = 0$):

Computational Super Step The computation of all values of $b_{I,j}$ and $b_{i,j}$ is done in parallel in the corresponding sites. We still denote by $\mathcal{C} = \max\{\mathcal{C}_I(n_I), \max_{i \in I}\{\mathcal{C}_i(n_i)\}\}$. Recall that in each site, the $b_{i,j}$ still may be computed in parallel if the corresponding computer has several cores.

Communication Super Step The results $b_{I,j}$ and $b_{i,j}$ must be sent for the final proceeding. We denote by \mathcal{T}_I time to transfer all values of $b_{I,j}$, and by \mathcal{T}_i time to transfer all values of $b_{i,j}$. The communication time now is calculated as $\mathcal{T} = \max\{\mathcal{T}_I, \max_{i \in I}\{\mathcal{T}_i\}\}$.

Final Proceeding \mathcal{T}_S still denotes time to search one result of $F_{\Phi, \phi}$.

The new computation time of (1) for the case of the distributed storage and computation is: $\mathcal{T}_{new}^{distr} = \mathcal{C} + \mathcal{T} + \mathcal{T}_S$. If the computations and the data transfer in each site may be done in parallel then we may combine two first super steps in the above model in one step that, in fact, leads to some variation of LogP model, introduced in [7]. If we denote by $\mathcal{D} = \max\{(\mathcal{C}_I(n_I) + \mathcal{T}_I), \max_{i \in I}\{(\mathcal{C}_i(n_i) + \mathcal{T}_i)\}\}$, then the corresponding computation time is: $\mathcal{T}_{new}^{LogP} = \mathcal{D} + \mathcal{T}_S$.

IX. CONCLUSION

In this work, we consider computations on sum-like graphs. We considered different scenarios, when our method leads to improvement in the complexity of the computations. We have shown several applications of our method in the fields of design and verification of repetitive and hierarchical structures, parallel computations, computations on distributed data.

Each of the considered computational models may be combined with each of scenarios A-D in order to analyze $\mathcal{T}_{old}/\mathcal{T}_{new}$. In fact, this analysis is a generalized formulation of Amdahl's style law for each computation model and scenario.

Acknowledgments

We would like to thank Prof. J. A. Makowsky for valuable discussions as well as for his reading of the contribution and his many suggestions.

Finally we would like to thank the referees for their careful reading and constructive suggestions.

REFERENCES

- [1] J. Makowsky and E. Ravve, "Incremental model checking for decomposable structures," in *Mathematical Foundations of Computer Science (MFCS'95)*, ser. Lecture Notes in Computer Science, vol. 969. Springer Verlag, 1995, pp. 540–551.
- [2] E. Ravve and Z. Volkovich, "A systematic approach to computations on decomposable graphs," 2013, to appear in Proceedings of SYNASC13.
- [3] E. Ravve, Z. Volkovich, and G.-W. Weber, "Effective optimization with weighted automata on decomposable trees," 2013, optimization Journal Special Issue at ECCO 2013.
- [4] S. Feferman and R. Vaught, "The first order properties of products of algebraic systems," *Fundamenta Mathematicae*, vol. 47, pp. 57–103, 1959.
- [5] J. Makowsky, "Algorithmic uses of the Feferman-Vaught theorem," *Annals of Pure and Applied Logic*, vol. 126, pp. 159–213, 2004.
- [6] L. Valiant, "A bridging model for parallel computation," *Communications of the ACM*, vol. 33(B), pp. 103–111, 1990.
- [7] D. Culler, R. Karp, D. Patterson, A. Sahay, K. Schauer, E. Santos, R. Subramonian, and T. von Eicken, "LogP: towards a realistic model of parallel computation," in *POPP '93 Proceedings of the fourth ACM SIGPLAN Symposium on Principles and practice of parallel programming*, vol. 28(7), 1993, pp. 1–12.
- [8] (2013, Jan) Seville map. [Online]. Available: <http://www.lonelyplanet.com/maps/europe/spain/andalucia/seville/>
- [9] H. Ebbinghaus, J. Flum, and W. Thomas, *Mathematical Logic, 2nd edition*, ser. Undergraduate Texts in Mathematics. Springer-Verlag, 1994.
- [10] J. Makowsky, "Translations, interpretations and reductions," 1994, unpublished Manuscript.
- [11] R. Fagin, "Generalized first-order spectra and polynomial time recognizable sets," in *Complexity of Computation*, ser. American Mathematical Society Proc, R. Karp, Ed., vol. 7. Society for Industrial and Applied Mathematics, 1974, pp. 27–41.
- [12] M. Garey and D. Johnson, *Computers and Intractability*, ser. Mathematical Series. W.H. Freeman and Company, 1979.
- [13] J. Makowsky and Y. Pnueli, "Arity vs. alternation in second order definability," in *LFCS'94*, ser. Lecture Notes in Computer Science, vol. 813. Springer, 1994, pp. 240–252, (Extended version to appear in the Annals of Pure and Applied Logic, 1995).
- [14] M. Vardi, "The complexity of relational query languages," in *STOC'82*. ACM, 1982, pp. 137–146.
- [15] E. Emerson, "Temporal and modal logic," in *Handbook of Theoretical Computer Science*, J. van Leeuwen, Ed. Elsevier Science Publishers, 1990, vol. 2, ch. 16.
- [16] J. Makowsky and E. Ravve, "Incremental model checking for fixed point properties of decomposable structures," 1995, technical Report TR844, revised version, April 1995, Department of Computer Science, Technion–Israel Institute of Technology, Haifa, Israel.
- [17] N. Biggs, R. Damerell, and D. Sands, "Recursive families of graphs," *Journal of Combinatorial Theory*, vol. 12, pp. 123–131, 1972.
- [18] B. Courcelle and J. Makowsky, "Fusion in relational structures and the verification of monadic second-order properties," *Mathematical Structures in Comp. Sci.*, vol. 12, pp. 203–235, April 2002. [Online]. Available: <http://dl.acm.org/citation.cfm?id=966880.966886>
- [19] E. Fischer and J. Makowsky, "Linear recurrence relations for graph polynomials," in *Pillars of computer science*, A. Avron, N. Dershowitz, and A. Rabinovich, Eds. Berlin, Heidelberg: Springer-Verlag, 2008, pp. 266–279. [Online]. Available: <http://dl.acm.org/citation.cfm?id=1805839.1805854>
- [20] D. Babić, A. Graovac, B. Mohar, and T. Pisanski, "The matching polynomial of a polygraph," *Discrete Applied Mathematics*, vol. 15, pp. 11–24, 1986.
- [21] E. Cockayne, E. Hare, S. Hedetniemi, and T. Wimer, "Bounds for the domination number of grid graphs," *Congr. Numer.*, vol. 47, pp. 217–228, 1985.
- [22] D. Crandall. (2013, Jan) CSE 477 design project specifications report for the discrete cosine transform decoder. [Online]. Available: <http://www.cs.indiana.edu/djcran/projects/cse477/report3/report3-2.html>
- [23] A. Glikson and J. Makowsky, "NCE graph grammars and clique-width," in *WG'03*, 2003, pp. 237–248.
- [24] G. Gottlob, E. Grädel, and H. Veith, "Datalog LITE: a deductive query language with linear time model checking," *ACM Transactions on Computational Logic*, vol. 3(1), 2002.
- [25] E. Ravve, "Decomposition of databases with translation schemes," Ph.D. dissertation, Department of Computer Science, Technion–Israel Institute of Technology, Haifa, 1998.
- [26] X. Quan and G. Wiederhold, "Incremental recomputation of active relational expressions," *IEEE Transactions on Knowledge and Data Engineering*, vol. 3, no. 3, pp. 337–341, 1991.
- [27] G. Dong and R. Topor, "Incremental evaluation of Datalog queries," in *Database Theory, 4th ICDT'92*, ser. Lecture Notes in Computer Science, J. Biskup and R. Hull, Eds., vol. 646. Springer Verlag, 1992, pp. 282–296.
- [28] R. Neal and G. Hinton, "A view of the EM algorithm that justifies incremental, sparse, and other variants," in *Learning in Graphical Models*. Kluwer Academic Publishers, 1998, pp. 355–368.
- [29] S. Ahn, J. A. Fessler, D. Blatt, and A. Hero, "Convergent incremental optimization transfer algorithms: Application to tomography," *IEEE Trans. Med. Imaging*, vol. 25(3), pp. 283–296, 2006.
- [30] O. Şeref, R. Ahuja, and J. Orlin, "Incremental network optimization: Theory and algorithms," *Operations Research*, vol. 57, pp. 586–594, 2009.
- [31] R. Poli, "Evolution of graph-like programs with parallel distributed genetic programming," in *Proceedings of 7th ICGA*, 1997, pp. 346–353.

Topological Analysis of the Subway Network of Madrid

Mary Luz Mouronte

Departamento de Ingeniería y Arquitecturas Telemáticas.
Escuela Técnica Superior de Ingeniería y Sistemas de Telecomunicación
Universidad Politécnica de Madrid
Email: mouronte.lopez@upm.es

Abstract—In this paper, we study by means of Complexity Science the topological structure of the subway network of Madrid. Different statistic features are analyzed: Degree of a Node i (K_i), Degree Probability Distribution $P(K)$, Nearest Neighbour Degree $K_{nn}(k)$, Clustering Coefficient, Average Path Length $\langle l \rangle$, Mean Service Efficiency ρ , Global Network Efficiency E_{glob} , Correlation Coefficient r_D and robustness of the network. This analysis will allow to obtain a deeper knowledge of this network and it will also help to improve its management: insight about connectivity, most relevant stops, efficiency, vulnerability and way of growth.

Keywords—Network Science; Transport Network; Statistical Analysis

I. INTRODUCTION

This paper calculates different mathematical parameters in the subway network of Madrid. Our research will allow to increase the current knowledge about this network. Structural properties of a subway network are very relevant for an effective transportation management in the urban cities. There are several works that analyze the characteristics of the transport networks:

Chen et al. [1], present an empirical investigation about the urban bus networks of four major cities in China: Hangzhou, Nanjing, Beijing and Shanghai.

Chen et al. [2], investigate the evolution of the dynamic properties in bus networks of Hangzhou, Nanjing, Beijing and Shanghai. New measurements of the average sum of the nearest-neighbors degree-degree correlation $D_{nn}(K)$ and the degree average edges among the nearest-neighbors $L(K)$ are proposed. The obtained results reflect that the considered transport network are organized randomly.

Chang et al. [3], study the subway network of Seoul, Tokyo, Boston and Beijing, by using the global and local efficiency. It is found that the Seoul subway network has a smaller global and local efficiency than the Tokyo network. The authors suggest that the Tokyo subway system is better for an overall distance trip but is weaker regarding incidents of disconnection. It is also shown for the subway networks of Boston, Seoul, Tokyo and Beijing, the global efficiency is inversely proportional to the length of the network. The Boston and Beijing local efficiencies are very low which means that these are somewhat deficient in some routes.

The rest of the paper is organized as follows: Section 2 describes the subway network of Madrid, in Section 3 the method of analysis and the results are presented, and finally

in Section 4 we end with some conclusions.

II. THE SUBWAY NETWORK OF MADRID

Madrid is one of the most populous cities in the world. It has a population of 3,254,950 dwellers on an area of 60,683 hectares, and a high developed public transportation network. The subway network of Madrid is one of the largest subway networks in the world, rivalling other networks such as the Shanghai, Guangzhou, Beijing or Delhi. In 2007, this network became the second largest subway network in Europe after London. The subway network of Madrid has 16 routes and 322 stops. The first route of the network began its operation in 1919.

The subway network of Madrid is operated with more than 2,400 trains and its yearly ridership was 628 millions in 2012. This network has been transformed by means of several improvement plans since 2011. Generally, the subway is open to the public from 6:00 AM to 1:30 AM every day of the year.

III. MATHEMATICAL ANALYSIS

We can map this network in three Topological Spaces: Space P, Space L and Space R. In these Spaces, the network is abstracted in a graph $G = (E;V)$, in which E is the set of nodes and V is the set of links between them. An adjacency matrix of $N \times N$ dimension $A(G)$ can be built as a bidimensional representation of the relationships between nodes, where $A_{ij} = 1$ when a connection between nodes i and j exists and $A_{ij} = 0$ otherwise. N is the number of nodes in E .

In the Space L, one node represents one stop, and one link symbolizes an union between two nodes if one stop is the successor of the other on a subway route. Space L is named Stop Geographical Space. In the Space P, one node represents one stop, and one link joins a pair of stops if at least one route provides direct service. A link means that passengers can take at least one route for a direct travel between two stops. If passengers have to exchange routes then the pair of stops is linked by more than one link. Space P is called Subway-Transferring Space. In the Space R, nodes are defined as routes and common stops determine the links. Space R is called Route Space.

While carrying out a topological study of the network for these three Spaces, some parameters are estimated:

Degree of a Node K_i and Degree Probability Distribution

$P(K)$: The degree of a node i is the number of links connected to it:

$$K_i = \sum_j A_{ij} \quad (1)$$

Not all nodes in the subway network have the same number of links: how the degree is distributed among the nodes is an interesting property which can be analyzed by estimating the Degree Probability Distribution $P(K)$.

Nearest Neighbour Degree $K_{nn}(k)$, which is defined as:

$$K_{nn}(K) = \sum_{K'=0}^{\infty} K' p(K'/K) \quad (2)$$

where $p(K'/K)$ is the conditional probability that a link belonging to a node with degree K links to a node with degree K' . Therefore K_{nn} is the average degree of those nodes that are found by following the links originating from a node of degree K . The evolution of $K_{nn}(K)$ is related to the assortativity of the network, which indicates the tendency of a node of degree K to associate with a node of the same degree K . In an assortative network, $K_{nn}(K)$ increases with increasing K but, in a disassortative network, $K_{nn}(K)$ decreases with increasing K while in a neutral network, $K_{nn}(K)$ does not depend on K .

Clustering Coefficient: given three actors i , j and w with mutual relations between j and i as well as between j and w , Clustering Coefficient is supposed to symbolize the likelihood i and w are also related. This parameter was used by Watts and Strogatz [4] for social networks analysis.

This concept can be explained by defining for $j \in V$,

$$m(j) = |\{i, w\} \in E : \{j, i\} \in E \text{ and } \{j, w\} \in E| \quad (3)$$

and

$$t(j) = \frac{K_j(K_j - 1)}{2} \quad (4)$$

We named $m(j)$ the number of opposite links of j , and $t(j)$ the number of potential opposite links of j .

For a node j with $K_j \geq 2$, the Clustering-Coefficient is defined as:

$$C(j) = \frac{m(j)}{t(j)} \quad (5)$$

and Clustering-Coefficient of a graph $G = (V, E)$ is denoted as:

$$C(G) = \frac{1}{|V|} \sum_{j \in V} C(j) \quad (6)$$

where V' is the set of nodes i with $K_i \geq 2$.

The Clustering Coefficient of a node ranges between 0 and 1.

Average Path Length $\langle l \rangle$, which is the average shortest path between all nodes of the network. We denote by $l(i, j)$ the distance between i and j , i.e., the number of links on a shortest path between them.

$$\langle l(i) \rangle = \frac{1}{n} \sum_j l(i, j) \quad (7)$$

represents the average distance from i to all nodes. The average distance in G is defined as:

$$\langle l \rangle = \frac{1}{n} \sum_i \langle l(i) \rangle = \frac{1}{n^2} \sum_{i, j} l(i, j) \quad (8)$$

We denote by $D = \max_{i, j} l(i, j)$ the diameter (D) of G , i.e., the largest distance between two nodes of the network.

Mean Service Efficiency (ρ) [9], which can be defined for the subway network as:

$$\rho = \frac{Ns}{M\phi} \quad 0 \leq \rho \leq 1 \quad (9)$$

Where Ns , M , and ϕ are the total of stops, the number of routes and the mean number of stops per route respectively. For a specific number of subway stops, a larger magnitude of ρ implies fewer subway routes that the transport company should maintain to satisfy the travel demand.

Global Network Efficiency (E_{glob}) [9], which may be described as:

$$E_{glob} = \frac{\sum_{i \neq j \in G} l_{ij}^{-1}}{Ns(Ns - 1)} \quad 0 \leq E_{glob} \leq 1 \quad (10)$$

Global efficiency is a measure of the performance of the network, under the assumption that the efficiency for sending information between two nodes i and j is proportional to the reciprocal of their distance $l(i, j)$.

Correlation Coefficient, the degree-degree correlation was analyzed as the correlation function between the remaining degrees [5] of the two nodes on each side of a link. Remaining degree means the degree of that nodes minus one. The normalized Correlation Coefficient is defined as:

$$r_D = \frac{1}{\sigma_D(q)^2} \sum_{u, v} uv \{e_D(u, v) - q_D(u)q_D(v)\} \quad (11)$$

where:

- $e_D(u, v)$ is the joint probability that the two vertices on each side of a randomly chosen link have u and v remaining degrees, respectively.
- $q_D(v)$ is the normalized distribution of the remaining degree [6].

$$q_D(v) = \frac{(v+1)P(v+1)}{\sum_u uP(u)} \quad (12)$$

$$\sigma_D(q)^2 = \sum_v v^2 q_D(v) - |\sum_v v q_D(v)|^2 \quad (13)$$

This quantity was named by Newman [6] the Degree Assortative Coefficient. In an assortative network, r_D is positive but, in a disassortative network, r_D is negative while in a neutral network $r_D = 0$.

Some statistical parameters are available in Table I: Total of Stops (Ns); Number of Subway Routes (M); Clustering Coefficient ($\langle C \rangle$), showing the subway routes density near each stop; Network Diameter (D), providing the maximum number of stops (or routes) on the shortest paths between any pair of stops (or routes); and finally, Average Shortest Path Length $\langle l \rangle$, denoting the average number of stops (or routes) on all the shortest paths between any two stops (or routes).

Clustering Coefficient is considered to be a measure of the local connectivity of a graph. High clustering is associated with robustness of a network, that is resilience against random network damage. Considering this parameter the subway network shows moderate resilience.

A node with high K controls the traffic flow, acting as gatekeeper. A node with high k can also act as a link between two distant sectors of the network. The average path length $\langle l \rangle$ can be interpreted as a measure of efficiency in the flow

TABLE I: Empirical data corresponding to the Subway Network of Madrid

Space	Parameter	Value	
Space L	M	16	
	Ns	322	
	$\langle k \rangle$	2.42	
	K_{max}	7	
	$\langle C \rangle$	0.01	
	$\langle w \rangle$	2.03	
	w_{min}	2.00	
	w_{max}	4.00	
	D	30	
	Space P	$\langle l \rangle$	10.19
$\langle k \rangle$		29.39	
K_{max}		99	
$\langle C \rangle$		0.90	
$\langle w \rangle$		2.08	
w_{min}		2.00	
w_{max}		6.00	
D		4	
Space R		$\langle l \rangle$	2.26
		$\langle k \rangle$	6.00
	K_{max}	13	
	$\langle C \rangle$	0.65	
	$\langle w \rangle$	4.88	
	w_{min}	1.00	
	w_{max}	14.00	
	D	4	
	$\langle l \rangle$	1.62	

of the network.

Several researches have estimated the efficiency and vulnerability in networks [7][8]. We analyze the robustness of subway network by calculating the value of the average shortest path length ($\langle l' \rangle$) and the distribution of the number of pairs of nodes Np separated by the shortest distance, in the original network and in the same network but with the highest K degree nodes removed for the tree Spaces. In the Space L, $\langle l' \rangle = 10.93$; in the Space P, $\langle l' \rangle = 2.27$ while in the Space R, $\langle l' \rangle = 2.00$. In Figures 1, 2 and 3 we observe that the distribution of distances changes drastically in the Spaces L and R after the gatekeepers elimination. The network shows low robustness in both spaces (removal of a route or a stop elimination in a route occurs). This feature is due to the current subway network design which, can be improved by means of optimization tasks. We can also notice that the most frequent short path length is 7 in the Space L (passangers should cross 7 stops without changing their route to get a destination in most cases), 2 in the Space P and 1 in the Space R (most routes are linked by a stop).

We also calculate E_{glob} and ρ in the Space P since these parameters lack clear meaning in the other Spaces. E_{glob} represents the total ability of the network to minimize the spatial resistance (or travel impediment), $E_{glob} = 0.298$ and $\rho = 0.847$.

The node degree and its distribution are very important properties for a network. From Figures 4, 5, 6 and Table I several conclusions can be obtained:

- In the Space L, we observe that the number of nodes

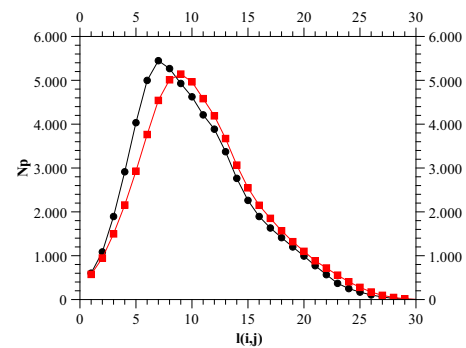


Figure 1: $Np - l(i, j)$ in the original network (black line) and in the same network but with the highest degree nodes removed (red line) in Space L

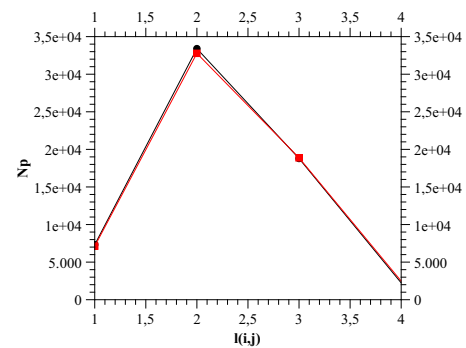


Figure 2: $Np - l(i, j)$ in the original network (black line) and in the same network but with the highest degree nodes removed (red line) in Space P

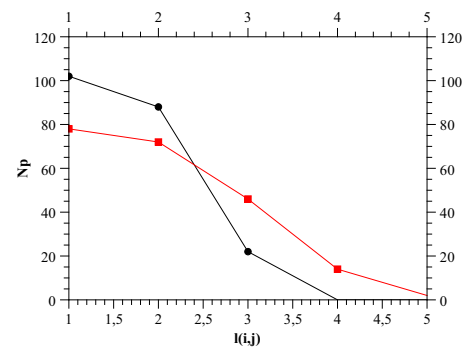


Figure 3: $Np - l(i, j)$ in the original network (black line) and in the same network but with the highest degree nodes removed (red line) in Space R

with degree $K = 2$ are the higher quantity, which means that a typical stop is directly connected to two other stops. In the Spaces P and R there are no nodes with a connectivity degree significantly different from the other nodes; the connectivity distribution is close to an uniform distribution. This happens because the company responsible for the urban transport in the city must ensure the uniform distribution of local equipment so that they are accessible to the entire population. In the Space R, we notice that routes with 9 common stops are the most frequent.

- The subway network constrained in different Geographical Spaces leads to different values of the node

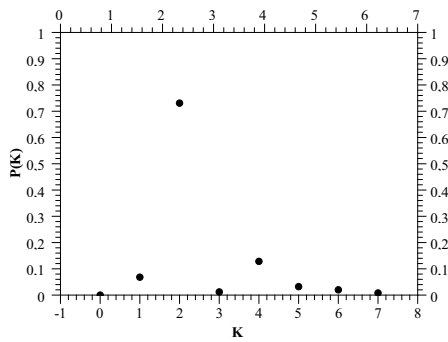


Figure 4: Degree Distribution in Space L

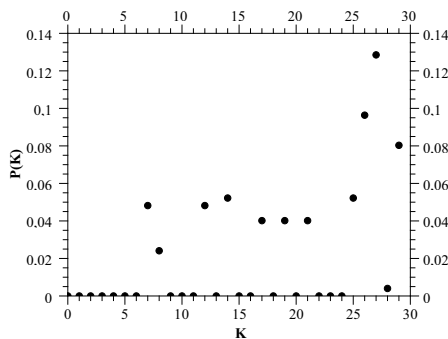


Figure 5: Degree Distribution in Space P

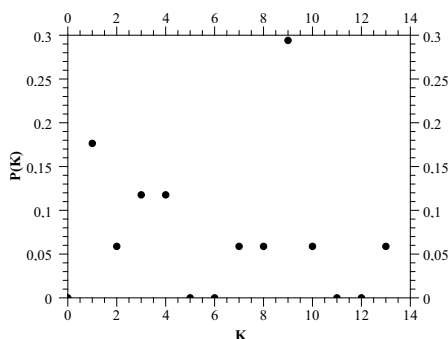


Figure 6: Degree Distribution in Space R

degree: low magnitudes in the Spaces L and R ($\langle K \rangle$ is 2.42 and 6 respectively) but very high value in the Space P ($\langle K \rangle$ is 29.39). This is because there are few common stops to different routes, although, from one origin stop many final stops can be reached. In the Space L, the most connected stops exist: Alonso Martínez and Avenida de América; in the Space P, the most connected node is also Alonso Martínez (this stop can be reached through most of the routes), finally, in the Space R, the most connected route is the route number 10. This happens because these elements are relevant communications centers in the city.

Regarding K_{nn} we can observe in Figures 7 and 8 that in some intervals it is difficult to establish whether the correlations are positive, negative or uncorrelated. The

statistical variations in K_{nn} can be suppressed by estimating its cumulated value. This magnitude decays with increasing K as it is showed in Figures 9 and 10; therefore we conclude that the network is assortative, that is the nodes in the subway network that have many connections tend to be connected to other nodes with many connections. This characteristic is also supported by the positive value of r_D in the spaces L and P (i.e., $r_D = 0.270223$ and $r_D = 0.092046$ respectively). This happens because during the tasks of design and planning to satisfy the traffic needs, was established that the new stops or new routes would be linked to other stops or routes that had similar connectivity in the network.

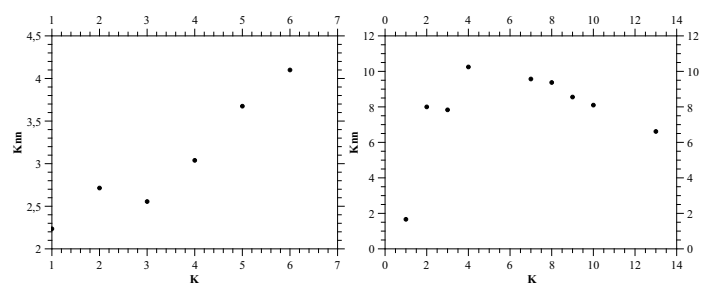


Figure 7: Left side: $K_{nn}-K$ in Space L. Right side: $K_{nn}-K$ in Space P

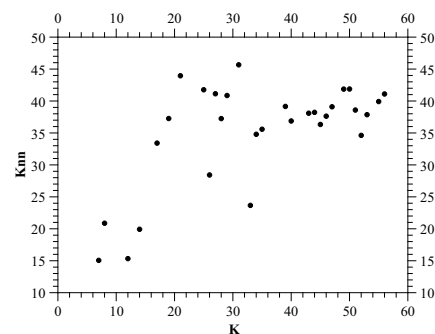


Figure 8: K_{nn} in Space R

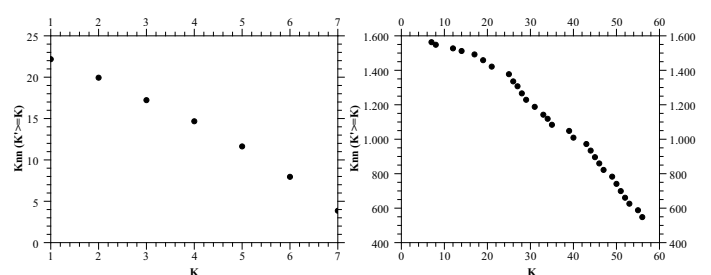
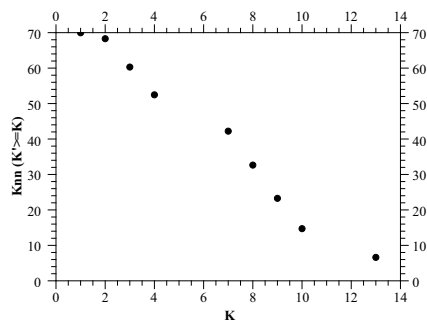


Figure 9: Left side: Cumulated $K_{nn}-K$ in Space L. Right side: Cumulated $K_{nn}-K$ in Space P

Figure 10: Cumulated K_m-K in Space R

The interactions between nodes is higher in the Space R (nodes have a larger average weight). In the Spaces L and P the interaction magnitudes are very similar (2.03 and 2.08 respectively). Link weights in the network represent multiplicity of connections between stops (Spaces L and P) and between routes (Space R). Passengers have one route (one link in each direction) in terms of average to move from one stop to another.

IV. CONCLUSION

This research uses the Network Science as a mathematical method to translate networks into graphs, from which important properties are collected. The underlying structure of a network has relevant consequences for its performance.

The subway network of Madrid is studied in three topological Spaces: Space L (stop geographical space), Space P (subway-transferring space) and Space R (Route space). We can conclude:

The study of Space R allows to know the average number and maximum value of the subway routes that a stop joins. These magnitudes are 6 and 13, respectively. The number of shared stops by two specific routes is defined as the weight of the link joining them. The average weight of a link is 4.88. The Degree Probability Distribution shows that there are no routes with a connectivity degree significantly different from the other routes; the connectivity distribution is close to an uniform distribution. Due to assortativity property routes that have many connections tend to be connected to other routes with many connections.

The analysis of Space P allows to know more precisely the accessibility and convenience of the network. The degree of a node symbolizes the number of stops a passenger can go to directly without any change, while the distance between two nodes represents the shortest path between them. The average degree and maximum degree of a node is 29.39 and 99 respectively; the average distance between nodes is 2.26. The average clustering coefficient of a node is high, 0.9, which means that there is high probability that the neighbours of this node (all other nodes to which it is joined by an link) are also connected to each other. In this Space, the Degree Probability Distribution also shows that there are no stops with a connectivity degree significantly different from the other stops. Nodes with many links tend to join other nodes with many links.

In the Space L, one link between two stops exists if they are consecutive on at least one route. The average degree of

a node is 2.42. In this Space, the network is assortative. The Degree Probability Distribution shows that many nodes have degree equal to 2.

Regarding robustness, we observe that the network is more robust in Space P than in the Spaces L and R. That is, if a failure occurs in one stop, those stops linked to it will be easily reachable by means of routes that provides a direct service between them, although the following stop on that route will be difficult to reach through one direct link; this network is more sensitive to problems in a route than in a stop.

Global efficiency can be an useful parameter for the assessment of the centrality before and after alterations to the network structure; these alterations can be caused by failures or planning changes.

Our future works will investigate deeper the vulnerability of subway network. We will also build a mathematical model that explains the growth of this network.

REFERENCES

- [1] Y. Z. Chen, N. Li, and D. R. He, "A study on some urban bus transport networks", *Physica A*, vol. 376, 2007, pp. 747–754.
- [2] C. Y. Z. Chen, and N. Li, "The randomly organized structure of urban ground bus-transport networks in China", *Physica A*, vol. 386, 2007, pp. 388–396.
- [3] K. H. Chang, K. Kim, H. Oshima, and S. M. Yoon, "Subway networks in cities", *Journal of the Korean Physical Society*, vol. 48, 2006, pp. S143–S145.
- [4] D. J. Watts, and S. H. Strogatz, "Collective dynamics of small-world networks". *Nature*, vol. 393, 1998, pp. 440–442.
- [5] M. E. J. Newman, "Assortative mixing in networks", *Physical Review Letters*, vol. 89, No. 20, 2002, pp. 208701–208705.
- [6] M. E. J. Newman, "Mixing patterns in networks", *Physical Review E*, vol. 67, No. 2, 2003, pp. 026126–026139.
- [7] J. P. Cardenas, Cardenas, R. M. Benito, M. L. Mouronte, and V. Feliu, "The effect of the complex topology on the robustness of spanish SDH network", in *Fifth International Conference on Networking and Services*, IEEE Xplore, 2007, pp. 86–90.
- [8] R. Criado, B. Hernández, and M. Romance, "Efficiency, vulnerability and cost: An overview with applications to subway networks worldwide", *Int. Journal of Bif. And Chaos*, vol. 17, 2007, pp. 2289–2301.
- [9] Z. Zhen-Tao, Z. Jing, L. Ping, and C. Xing-Guang, "An evolutionary model of urban bus transport network based on B-space", *Chinese Physics B*, vol 17, No 8, 2008, pp. 2874–2880.

QPSOL: Quantum Particle Swarm Optimization with Levy's Flight

Optimization of appliance scheduling for smart residential energy grids

Ennio Grasso, Claudio Borean

Swarm Joint Open Lab

TELECOM ITALIA

Turin, Italy

e-mail: ennio.grasso@telecomitalia.it, claudio.borean@telecomitalia.it

Abstract— This paper considers the minimum electricity cost scheduling problem of smart home appliances in the context of smart grids. Functional characteristics, such as expected duration and peak power consumption of the smart appliances can be adjusted through a power profile signal. The optimal scheduling of power profile signals minimizes cost, while satisfying technical operation constraints and consumer preferences. Time and power constraints, and optimization cost are modeled in this framework using a metaheuristic algorithm based on a Quantum inspired Particle Swarm with Lévy flights. The algorithm runs on the limited computational power provided by the home gateway device and in almost real-time as of user perception.

Keywords: scheduling, swarm intelligence, metaheuristic smart grids.

I. INTRODUCTION

This paper considers the minimum electricity cost scheduling problem of smart home appliances in the context of the Energy@Home international project [1]. Functional characteristics, such as expected duration and peak power consumption of the smart appliances can be modeled through a power profile signal. The optimal scheduling of power profile signals minimizes cost, while satisfying technical operation constraints and consumer preferences. Time and power constraints, and optimization cost are modeled in this framework using a metaheuristic algorithm based on a Quantum inspired Particle Swarm with Lévy flights. The algorithm runs on the limited computational power provided by the home gateway device and in almost real-time as of user perception.

The innovative Quantum inspired Particle Swarm Optimization (QPSO) with Lévy flights metaheuristic algorithm for scheduling home smart appliances, capturing all relevant appliance operations, is not only described in the paper but also validated, since the results of the implementation of it running on an embedded platform are presented. With appropriately dynamic tariffs and short-term load forecasting, the proposed framework can calculate and propose a schedule for achieving high cost savings and overloads prevention, improving the user experience of energy management services. Good quality approximate solutions can be obtained in a short amount of computation time, in the order of about 2 seconds an almost optimal approximate solution can be obtained, which enables the

usage of this algorithm on very embedded and low cost platforms .

It is also described in the paper how the proposed framework could be extended to incorporate solar power forecasting in the presence of a residential Photovoltaic (PV) system by tuning the objective function and using the solar energy forecaster as further input to the scheduler ([16], [18]).

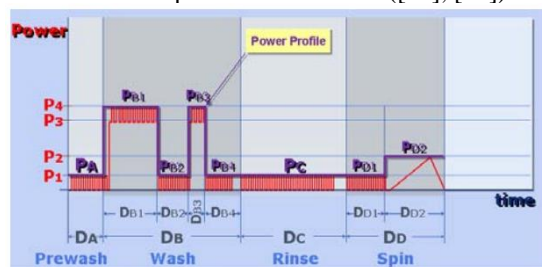


Figure 1. Example Power Profile with its phases generated by a washing machine

The paper is structured as it follows. Section 2 describes and models the problem of scheduling of smart appliance. Section 3 highlights how this problem can be classified as a NP-Hard Combinatorial Optimization Problem. In Section 4 a review of metaheuristic algorithms is performed while in Section 5 the new QPSOL algorithm proposed in the paper is described. Section 6 reports the results of the simulations of the QPSOL algorithm applied to the problem of scheduling of smart appliance while Section 7 presents the conclusions and future work.

II. SCHEDULING PROBLEM OF SMART HOME APPLIANCES

Europe has set the ambitious target of a 20% share of overall energy demand to be supplied from renewable energy by 2020. In order to achieve this target, the share of renewable energy will need to increase to some 35%. Most of the increase will come from wind and solar energy, which are both fluctuating resources by nature.

Electricity consumption varies between different hours of the day, between days of the week, and between seasons of the year. In recent years, the power demand has reached new peak levels and environmental / economic reasons will require more complex power balance scenarios also based on the introduction of residential renewable electricity generation to reduce the carbon footprint and CO2 emission.

One of the major challenges associated with this drastic restructuring of the energy supply with renewables is how electricity networks can cope with the extreme variability of wind and solar energy production. In the past, the ideal load curve was flat in order to allow for the full load operation of conventional power plants. In the future, the ideal demand needs to be variable in order to adapt to the current production from renewable energy sources.

It is expected in the near future that time-varying and dynamic electricity tariffs will increase popularity, especially for the reduction of peak power consumption which are the most detrimental from the grid operators. However, such load balancing is only feasible if consumers are both able and willing to consider tariff information, but it is still unrealistic to expect most consumers to identify the most economical operations of their appliances with dynamic tariffs, or in the presence of a small-scale photovoltaic (PV) power generation system which adds even more complexity in determining the economic convenience between immediate power usage versus selling the power to the main grid operator.

In view of the above considerations, not only is an automatic decision system highly desirable but even necessary in most cases, which either directly takes control of the appliances' operations, or at the very least is capable of providing advice to the home consumers.

A. Smart Appliances in Energy@Home

The Energy@Home (E@H) consortium is based on ZigBee communication between smart appliances in a home mesh wireless network [1]. The "core" element of this home network is the Home Gateway (HG) that coordinates and manages the smart appliances as end devices. Among its functionalities, the HG provides the intelligence for real-time scheduling of residential appliances, typically in the time interval 24 hours ahead, based on the (possibly) varying tariff of the day, the forecasted energy power consumption, and also the forecasted home PV power generation, if available.

The proposed scheduling framework is based on the Power Profile Cluster defined in the E@H specifications [1], which specifies that each appliance operation process is modeled as a Power Profile divided into a set of sequential energy phases, as presented in Figure 1. An energy phase is an uninterruptible logic subtask of the appliance operation, which uses a pre-specified amount of electric energy. The energy phases are sequential since the next phase cannot start until the previous phase is completed, e.g., a washing machine agitator cannot start until the basin is filled with water.

In addition to having a specified energy usage, each energy phase is characterized by peak maximum power, a specific duration, and a possible maximum activation delay after the end of the previous phase. Some phases cannot be delayed and must start soon after the previous phase completes (maximum delay is zero). Other phases may be delayed adding extra flexibility in the scheduling of the Power

Profile, e.g., the washing machine agitator must start within ten minutes of the basin being filled.

In addition, the scheduler needs to take into account user specified time preferences, requiring that certain appliances should be run within some particular time intervals, e.g., the dishwasher must complete washing dishes between 13:00 and 18:00.

The objective of the HG scheduler is to find the least expensive scheduling for a set of smart appliances, each characterized by a Power Profile with its energy phases, while satisfying the necessary operational constraints.

The scheduling execution interval is divided into 1440 1-minute time slots for a 24-hour period. The number of appliances considered for scheduling is denoted N , and the number of energy phases for each appliance is denoted n_i for $i = 1, 2, \dots, N$. The problem dimension, i.e., the number of independent variables that make up the problem, is

$$\sum_{i=1}^N \sum_{j=1}^{n_i} p_{ij} \quad (1)$$

where p_{ij} is the j th phase of power profile i . The objective of the scheduler is to minimize the total electricity cost for operating the appliances based on a given 24-hour ahead electricity tariff while taking into account user comfort criteria (earlier executions are preferable than delayed execution) and respecting time and energy constraints.

B. Modeling Time and Energy Constraints

Even in the presence of the HG scheduler controlling a set of smart appliances, the real number of home appliances and other electric powered devices that consume energy in the house is higher and outside the control of the scheduler. For that reason any sensible scheduling system must be complemented by an appropriate forecasting module that provides good estimation of the overall power consumption based on past statistics.

The energy constraints imposes that for each time slot, the total sum of power required by all phases running in that slot, plus the forecasted power consumption "outside" the control of the scheduler, be less than peak power threshold provided by the grid operator,

$$\sum_{i=1}^N \sum_{j=1}^{n_i} \text{Power}(p_{ij}) + \text{Load}_{\text{forecast}} \leq \text{PeakPower} \quad (2)$$

Time constraints are twofold. On the one hand the user can set up time preference constraints, specifying the time interval that a particular appliance must be scheduled in terms of an earliest start time (EST), e.g., after 13:20, and a latest end time (LET), e.g., before 18:00.

$$\text{EST} \leq \text{PP}_i^{\text{start}} \leq \text{PP}_i^{\text{end}} \leq \text{LET} \quad (3)$$

where $\text{PP}_i^{\text{start}}$ and PP_i^{end} are respectively the scheduling start slot of the 1st phase of Power Profile i and the end slot of the last phase.

The second time constraint is the maximum activation delay of each of the sequential phases that make up the Power Profiles. While the scheduling interval specified in the first

constraint is absolute, the maximum activation delays are relative and therefore the lower and upper bound time limits of each phase need to be adjusted based on the scheduling decisions for the previous phase.

$$p_{ij}^{\text{end}} \leq p_{ij+1}^{\text{start}} \leq p_{ij}^{\text{end}} + p_{ij+1}^{\text{maxdelay}} \quad (4)$$

III. NP-HARD COMBINATORIAL OPTIMIZATION PROBLEMS

Given the problem formulation, the scheduling of Power Profiles, each composed by a set of sequential (and possibly delayable) phases, under energy constraints is classified in the more general family of Resource Constrained Scheduling Problem (RCSP), which is known as being an NP-Hard combinatorial optimization problem [12] [13].

Moreover, the presence of time constraints introduce even another dimension to the complexity of problem, known as RCSP/max, i.e., RCSP with time windows. Combining the inherent complexity of the problem with the fact that the limited computing power of the HG which runs the logic of algorithm, and the almost real-time requirement for finding a solution (typically the user wants a perceived immediate answer), make the formulation of the problem and its solution a challenging endeavor.

From a theoretical perspective, combinatorial optimization problems have a well-structured definition consisting of an objective function that needs to be minimized (e.g., the energy cost) and a series of constraints. These problems are really important for the great amount of real-life applications that can be modeled in this way. For example, areas like routing or scheduling contain plentiful hard challenges that can be expressed as a combinatorial optimization problem.

For easy problems, exact methods can be exploited, such as Branch&Bound and Mixed Integer Linear Programming (MILP), with back-tracking and constraints propagation to prune the search space. However, in most circumstances, the solution space is highly irregular and finding the optimum is in general impossible. An exhaustive method that checks every single point in the solution space would be infeasible in these difficult cases, since it takes exponential time.

As a point of fact, [2] also addresses a similar scheduling problem of smart appliances, and relies on traditional MILP as a problem solver. They provide computation time statistics for their experiments, running on an Intel Core i5 2.53GHz equipped with 4GB of memory and using the commercial application CPLEX and MATLAB. According to their figures, discretizing the time interval in 10-minute discrete slots (for a total of 144 daily slots), takes their algorithm about 15.4 seconds to find a solution. With 5-minute slots the time rises to 83.6 seconds, and with 3-minute slots to 860 seconds. From these figures it is clear that a traditional approach like MILP or B&B is hardly acceptable for scheduling home appliances, and other more efficient methods need to be investigated.

A. Convex Constraints and Smooth Objective Functions

Generally speaking, optimization problems can be categorized, from a high-level perspective, as having either convex or non-convex constraints.

Convex constraints form a series of convex regions where exact methods could be applied (e.g., B&B, linear-programming, etc.). The main idea, in convex optimization problems, is that every constraint restricts the space of solutions to a certain convex region. By taking the intersection of all these regions we obtain the set of feasible solutions, which is also convex. Due to the nice structure of the solution space, every single local optimum is a global one. Most conventional or classic algorithms are deterministic. For example, the simplex method in linear programming is deterministic, and use gradient information in the search space, namely the function values and their derivatives.

Non-convex constraints create a many disjoint regions, and multiple locally optimal points within each of them. As a result, if a traditional search method is applied, there is a high risk of ending in a local optimum that may still be far away from the global optimum. But the main drawback is that it can take exponential time in the size of problem dimension to determine if a feasible solution even exists!

Another definition is that of smooth function, i.e., a function that is differentiable and its derivative is continuous. If the objective function is non-smooth, the solution space typically contains multiple disjoint regions and many locally optimal points within each of them. The lack of a nice structure makes the application of traditional mathematical tools, such as gradient information, very complicated or even impossible in these cases.

Most “real” problems are neither convex nor smooth, so traditional exact methods cannot be applied. Finding a solution not the best one but “acceptable”, or even finding a feasible solution is NP-Hard.

B. An Overview of General Metaheuristic Algorithms

A problem is NP-Hard if there is not an exact algorithm that can solve the problem in polynomial time with respect to the problem’s dimension. In other words, aside from some “toy-problems”, an NP-Hard problem would require exponential time to find a solution by systematically “exploring” the solution space.

A common method to turn an NP-Hard problem into a manageable, feasible approach is to apply heuristics to “guide” the exploration of the search space. These heuristics are based on “common-sense” specific for each problem and are the basis for developing Greedy Algorithms that can build the solution by selecting at each step the most promising path in the solution space based on the suggested heuristics. Obviously this approach is short-sighted since it proceeds with incomplete information at each step. Very rarely do greedy algorithms find the best solution or worse yet they might fail to find a feasible solution even if one does exist.

A better approach for solving complex NP-Hard problems that has shown great success is based on metaheuristic algorithms. The word *meta* means that their heuristics are not problem specific to a particular problem, but general enough to be applied to a broad range of problems. Examples of metaheuristic algorithms are *Genetic and Evolutionary Algorithms*, *Tabu search*, *Simulated Annealing*, *Greedy Randomized Adaptive Search Procedure (GRASP)*, *Particle-Swarm-Optimization*, and many others.

The idea of metaheuristics is to have efficient and practical algorithms that work most the time and are able to produce good quality solutions, some of them will be nearly optimal. Figuratively speaking, searching for the optimal solution is like *treasure-hunting*. Imagine we are trying to find a hidden treasure in a hilly landscape within a time limit. It would be a silly idea to search every single square meter of an extremely large region with limited resources and limited time. A more sensible approach is to go to some place almost randomly and then move to another plausible place using some hints we gather throughout.

Two are the main elements of all metaheuristic algorithms: intensification and diversification. *Diversification* via randomization means to generate diverse solutions so as to explore the search space on the global scale and to avoid being trapped at local optima. *Intensification* means to focus the search in a local region by exploiting the information that a current good solution is found in this region as a basis to guide the next step in the search space. The fine balance between these two elements is very important to the overall efficiency and performance of an algorithm.

IV. CLASSIFICATION OF METAHEURISTIC ALGORITHMS

Metaheuristic algorithms are broadly classified in two large families: *population-based* and *trajectory-based*. Going back to the treasure-hunting metaphor, in a trajectory-based approach we are essentially performing the search alone, moving from one place to the next based on the hints we have gathered so far. On the other hand, in a population-based approach we are asking a group of people to participate in the hunting sharing all information gathered by all members to select the next promising paths for the next moves.

A. Genetic Algorithms

Genetic algorithms (GA) were introduced by John Holland and his collaborators at the University of Michigan in 1975 [3]. A GA is a search method based on the abstraction of Darwinian evolution and natural selection of biological systems, and representing them in the mathematical operators: *crossover* (or recombination), *mutation*, *fitness evaluation* and *selection* of the best. The algorithm starts with a set of candidate solutions, the initial population, and generate new offspring through random mutation and crossover, and then applies a selection step in which the

worst solutions are deleted while the best are passed on to the next generation. The entire process is repeated multiple times and gradually better and better solutions are obtained. GA algorithms represent the inseminating idea of all more recent population-based metaheuristics.

One major drawback of GA algorithms is the “conceptual impedance” that arises when trying to formulate the problem at hand with the genetic concepts of the algorithm. The formulation of the fitness function, population size, the mutation and crossover operators, and the selection criteria of the offspring population are crucially important for the algorithm to converge and find the best, or quasi-best, solution.

B. Simulated Annealing

Simulated Annealing (SA) was introduced by Kirkpatrick et al. in 1983 [5] and is a trajectory-based approach that simulates the evolution of a solid in a heat bath to thermal equilibrium. It was observed that heat causes the atoms to deviate from their original configuration and transition to states of higher energy. Then, if a slow cooling process is applied, there is a relatively high chance for the atoms to form a structure with lower internal energy than the original one. Metaphorically speaking, SA is like dropping a bouncing ball over a hilly landscape, and as the ball bounces and loses its energy it eventually settles down to some local minima. But if the ball loses energy slowly enough keeping its momentum, it might have a chance to overcome some local peaks and fall through a better global minimum.

C. Particle Swarm Optimization

Particle swarm optimization (PSO), introduced in 1995 by American social psychologist James Kennedy, and engineer Russell C. Eberhart [6], represents a major milestone in the development of population-based metaheuristic algorithms. PSO is an optimization algorithm inspired by swarm intelligence of fish and birds or even human behavior. The multiple particles swarm around the search space starting from some initial random guess and communicate their current best found solutions and also share the global best so as to focus on the quality solutions. The greatest advantage of PSO over GA is that it is much simpler to apply in the formulation of the problem. Instead of using crossover and mutation operations it exploits global communication among the swarm particles. Each particle in the swarm modifies its position with a velocity that includes a first component that attracts the particle towards the best position so far achieved by the particle itself. This component represents the personal experience of the particle. The second component attracts the particle towards the best solution so far achieved by the swarm as a whole. This component represents social communication skill of the particles.

Denoting with N the dimensionality of the search space, i.e., the number of independent variables that make up the exploring search space, each individual particle is

characterized by its position and velocity N-vectors. Denoting with x_i^k and v_i^k respectively the position and velocity of particle i at iteration k , the following equations are used to iteratively modify the particles' velocities and positions:

$$v_i^{k+1} = w v_i^k + c_1 r_1 (p_i - x_i^k) + c_2 r_2 (g^* - x_i^k) \quad (5)$$

$$x_i^{k+1} = x_i^k + v_i^{k+1} \quad (6)$$

where w is the *inertia* parameter that weights the previous particle's momentum; c_1 and c_2 are the *cognitive* and *social* parameter of the particles multiplied by two random numbers r_1 and r_2 uniformly distributed in $[0 - 1]$, and are used to weight the velocity respectively towards the particle's personal best, $(p_i - x_i^k)$, and towards the global best solution, $(g^* - x_i^k)$, found so far by the whole swarm. Then the new particle position is determined simply by adding to the particle's current position the new computed velocity, as shown in Figure 2.

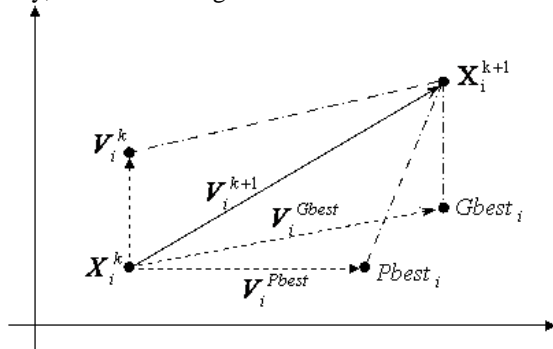


Figure 2. New particle position in PSO

The PSO coefficients that need to be determined are the inertia weight w , the cognitive and social parameters c_1 and c_2 , and the number of particles in the swarm. We can interpret the motion of a particle as the integration of Newton's second law where the component $c_1 r_1 (p_i - x_i^k) + c_2 r_2 (g^* - x_i^k)$ are the attractive forces produced by springs of random stiffness, while w introduces a virtual mass to stabilize the motion of the particles, avoiding the algorithm to diverge, and is typically a number such that $w \approx [0.5 - 0.9]$. It has been shown, without loss of generality, that for most general problems the number of parameters can even be reduced by taking $c_1 = c_2 \approx 2$.

D. Quantum Particle Swarm Optimization

Although much simpler to formulate than GA, classical PSO has still many control parameters and the convergence of the algorithm and its ability to find a near-best global solution is greatly affected by the value of these control parameters. To avoid this problem a variant of PSO, called Quantum PSO (QPSO) was formulated in 2004 by Sun and al. [7], in which the movement of particles is inspired by quantum mechanics.

The rationale behind QPSO stems from the observation that statistical analyses have demonstrated that in classical PSO each particle i converges to its local attractor a_i defined as

$$a_i = (c_1 p_i + c_2 g^*) / (c_1 + c_2) \quad (7)$$

where p_i and g^* are the personal best and global best of the particle. The local attractor of particle i is a stochastic attractor that lies in a hyper-rectangle with p_i and g^* being two ends of its diagonal, and the above formulation can also be rewritten as

$$a_i = r p_i + (1 - r) g^* \quad (8)$$

where r is a uniformly random number in the range $[0 - 1]$.

In classical PSO, particles have a mass and move in the search space by following Newtonian dynamics and updating their velocity and position at each step. In quantum mechanics, the position and velocity of a particle cannot be determined simultaneously according to uncertainty principle. In QPSO, the positions of the particles are determined by the Schrödinger equation where an attractive potential field will eventually pull all particles to the location defined by their local attractors. The probability of particle i appearing at a certain position at step $k+1$ is given by:

$$x_i^{k+1} = a_i + \beta |x_{mbest}^k - x_i^k| \ln(1/u), \text{ if } v \geq 0.5 \quad (9)$$

$$x_i^{k+1} = a_i - \beta |x_{mbest}^k - x_i^k| \ln(1/u), \text{ if } v < 0.5 \quad (10)$$

where u and v are uniformly random numbers in the range $[0 - 1]$, x_{mbest}^k is the mean best of the population at step k defined as the mean of the best positions of all particles

$$x_{mbest} = (1/N) \sum_{i=1}^N p_i \quad (11)$$

β is called *contraction-expansion* coefficient and controls the convergence speed of the algorithm.

The QPSO algorithm has been shown to perform better than classical PSO on several problems due to its ability to better explore the search space and also has the nice feature of requiring one single parameter to be tuned, namely the β coefficient. The exponential distribution of positions in the update formula makes QPSO search in a wide space. Moreover, the use of the mean best position x_{mbest} , each particle cannot converge to the global best position without considering all other particles, making them explore more thoroughly around the global best until all particles are closer. However, this may be both a blessing and a curse; it may be more appropriate in some problems but it may slow the convergence of the algorithm in other problems. Again, there is a very fine balance between exploration / exploitation. How large is the search space, and how much time is given to explore before returning a solution.

E. Dealing with Constraints

Many real world optimization problems have constraints, for example the available amount of certain resources, the boundary domain of certain variables, etc. So an important question is how to incorporate constraints in the problem formulation.

In some cases, it may be simple to incorporate the feasibility of solutions directly in the formulation of a problem. If we

know the boundary domain of a certain dependent variable and the proposed solution violates such domain we can either reject the solution or modify it by constraining the variable within the boundaries. For example, suppose a time variable must satisfy the time interval between 9:00 and 13:00, while the proposed solution would place it at 14:34. One way to deal with the above violation is to constrain the variable to its upper bound (UB) 13:00 and reevaluate the objective function. This will be probably worse than before, but at least it will be feasible and need not be rejected altogether.

A second way is to incorporate the constrains directly in the formulation of the objective function through the addition of a *penalty* element so that a constrained problem becomes unconstrained. If $f(x)$ is the objective function to be minimized, and subject to the constraints x in domain $[x_{lower}, x_{upper}]$, we rewrite the objective function as

$$f(x) = f(x) + \sum_{i=1}^N w_i g_i(x) \tag{12}$$

$g_i(x)$ measure the amount of constraint violation and is zero if x is within the domain boundaries, or it is some function $g(x - x_{lower}), (x_{upper} - x)$ otherwise. w_i are the penalty weights that needs to be large enough to skew the choice of the fittest solutions towards the smallest penalty component, typically in the range $10^9 - 10^{15}$.

Note that the two approaches described above to deal with constraints need not be mutually exclusive and can both be incorporated in the formulation of a problem; some constraints may very well be modeled with the first method, while other are modeled with the penalty method.

F. Nature Inspired Random Walks and Lévy Flights

A random walk is a series of consecutive random steps starting from an original point: $x_n = s_1 + \dots + s_n = x_{n-1} + s_n$, which means that the next position x_n only depends on the current position x_{n-1} and the next step s_n . This is the typical main property of a Markov chain. Very generally we can write the position in random walks at step $k+1$ as

$$x^{k+1} = x^k + s \sigma_k \tag{13}$$

where σ_k is a random number drawn from a certain probability distribution. In mathematical terms, each random variable follows a probability distribution, for example a *Gaussian* (normal) distribution is the most well-known because many physical phenomena obey this distribution and the random walk becomes the *Brownian* motion. But if the step length obeys other non-Gaussian distributions we have to deal with a more generalized form of random walks.

Various studies have shown that the random walk behavior of many animals and insects have the typical characteristics of the *Lévy* probability distribution and the random walk is called Lévy flight [8] [9] [10]. The Lévy distribution has the

nice mathematical feature of being both stable and heavy-tailed. A stable distribution is such that any sum n of random number drawn from the distribution is finite and can be expressed as

$$\sum_{i=1}^n x_i = n^{1/\alpha} x \tag{14}$$

where α is called the index of stability and controls the shape of the Lévy distribution with $0 < \alpha \leq 2$. Notably, two value for α are special cases of two other distribution, the Gaussian distribution for $\alpha = 2$, and the Cauchy distribution for $\alpha = 1$.

The heavy-tail characteristic implies that the Lévy distribution has an infinite variance decaying at large x to $\lambda(x) \sim |x|^{-1-\alpha}$

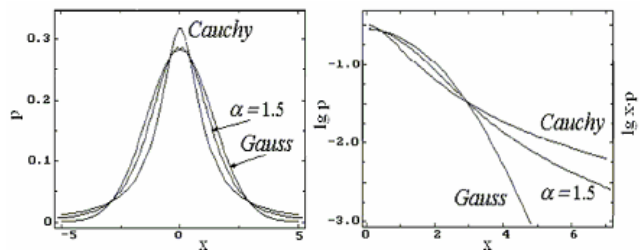


Figure 3. Cauchy

Figure 3. shows the shapes of the Gaussian, Cauchy, and Lévy distribution with $\alpha = 1.5$. The difference becomes more pronounced in the logarithmic scale showing the asymptotic behavior of the Lévy and Cauchy distribution compared with the Gaussian.

Due to the stable property, a random walker following the Lévy distribution will cover a finite distance from its original position after any number of steps. But also due to the heavy-tail (divergence of the variance), extremely long jumps may occur, and typical trajectories are self-similar, on all scales showing clusters of shorter steps interspersed by long excursions, as shown in Figure 4. In fact, the trajectory of a Lévy flight has fractal dimension $d_f = \alpha$.

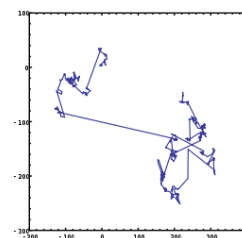


Figure 4. Levy's flight

In that sense, the Gaussian distribution in Figure 5. represents the limiting case of the basin of attraction of the so-called generalized central limit theorem for $\alpha = 2$ and the motion of the walker follows a Brownian path.

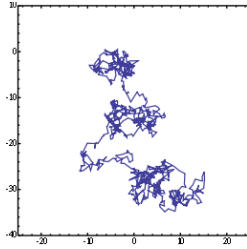


Figure 5. Brownian path

Due to the remarkable properties of stable, heavy-tailed distributions it is now believed that the Lévy statistics provides a framework for the description of many natural phenomena in physical, chemical, biological, economical systems from a general common point of view. For instance, the foraging behavior of bacteria and higher animals relies on the advantages of Lévy distributed excursion lengths, which optimize the search compared to Brownian search giving a better chance to escape from local optima.

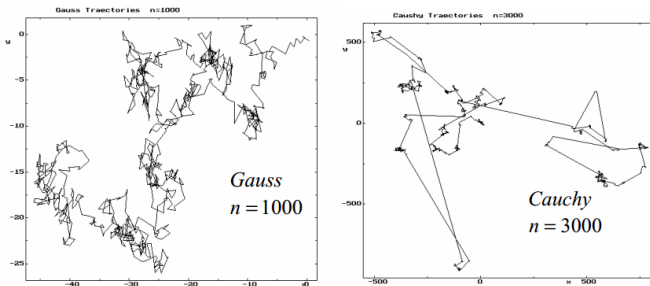


Figure 6. the trajectories of a Gaussian (left) and a Lévy (right) walker

The Figure 6. above shows the trajectories of a Gaussian (left) and a Lévy (right) walker. Both trajectories are statistically self-similar, but the Lévy motion is characterized by island structure of clusters of small steps, connected by long steps.

G. Step Size in Random Walks.

In the general equation of a random walk $x^{k+1} = x^k + s \sigma_k$, a proper step size, which determines how far a random walker can travel after k number of iterations, is very important in the exploration of the search space. The two component that make up the step are the scaling factor s and the length of the random number in the distribution σ_k . A proper step size is very important to balance exploration and exploitation, too small a step and the walker will not have a chance to explore potential better places, on the other hand too large steps will scatter the search from the focal best positions. From the theory of isotropic random walks, the distance traveled after k steps in N dimensional space is

$$D = s \cdot \sqrt{k N} \tag{15}$$

In a length scale L of a dimension of interest, the local search is typically reasonably limited in the region $L / 10$, that is $D = L / 10$, which means that the scaling factor

$$s \approx 0.1 L / \sqrt{k N} \tag{16}$$

In typical metaheuristic optimization problems, we can expect the number of iterations k in the range 100 – 1000. For example, with 100 iterations and $N=1$ (a one dimensional problem) we have $s = 0.01 L$, and to another extreme with 1000 iterations and $N=10$ we have $s = 0.001 L$. Therefore a scaling factor between 0.01 – 0.001 is basically a reasonable choice in most optimization problems. L is still kept independent as each dimension of the problem may very well have a very different length scale.

V. QUANTUM INSPIRED PARTICLE SWARM ALGORITHM WITH LÉVY FLIGHTS

After several experimental and simulated alternative metaheuristic approaches, we have come to the definition of a novel variant of the PSO algorithm that can be described as Quantum inspired PSO with Lévy flights (QPSOL). The algorithm tries to capture and exploit some of the best characteristics of various algorithms described in the previous sections. The result being an algorithm that provides a good balance between exploration and exploitation that gives quasi-optimal solutions within a very short time even with limited computing power. In fact, the Home Gateway (HG) is a low power ARM embedded system running a Java Virtual Machine in the OSGi framework.

The two main assumptions of the QPSOL algorithm are: first, as in Quantum PSO, particles have no mass and move around their attractor within a probability distribution. Secondly, rather than follow the quantum physics that uses the exponential distribution, in QPSOL particles move according to the nature-inspired Lévy distribution. From our experiments and simulations, the quantum inspired PSO, coupled with the Lévy distribution, has proven to outperform the classical PSO and traditional QPSO.

For our purposes, the Lévy distribution coefficient α chosen in QPSOL is actually the Cauchy coefficient $\alpha = 1$. The Cauchy random generator is much simpler than the more general algorithm for Lévy generation and that is a determining factor in runtime execution. Since the random generation needs to be executed for an umpteen number of times (i.e., the dimension of the problem, by the number of particles in the swarm, by the number of iterations of the algorithm), the computing speed of the random generation is of paramount importance. From our experiments, within a given time limit allotted to the algorithm to find a solution, the Cauchy version of the algorithm is able to execute almost twice the number of iterations than the general Lévy version. Therefore, even if there was an optimal coefficient α that provides better results for the same number of iterations, it will be outperformed by the Cauchy variant that with more allowed iterations finds better solutions. Since Cauchy is simply a special case of the general Lévy

distribution, henceforth we will continue to refer to the algorithm as a Quantum PSO with Lévy flights QPSOL.

A. QPSOL for Scheduling Appliances

As any population based metaheuristic algorithm, each particle represent a complete solution to the problem, i.e., a complete schedule for all the Power Profiles of the appliances. Since each Power Profile is itself composed by a sequence of phases, we model each particle (complete solution) as a set of N sub-particles, where N is the number of Power Profiles and where each sub-particle represents the schedule for the energy phases of that Power Profile. Below we report the Java code of the evolution of the sub-particles in the swarm and it represents the core of the QPSOL algorithm.

The Lévy light of a sub-particle is a loop on the sequence of the energy phases. First, the maximum delay for each phase is determined. The maximum delay for the first phase (index = zero) is actually the maximum slack interval of the whole Power Profile as imposed by the user.

The maximum delay of each subsequent phase is the minimum between its maximum delay as per Power Profile specification, and the remaining slack for the whole remaining phases updated at each step after a phase is moved with the Lévy flight.

After calculating the maximum delay for each phase i , the Lévy flight is performed with these equations:

$$a_i = r p_i + (1 - r) g^* \quad (17)$$

$$x_i = a_i + \beta (a_i - x_i) \lambda_i \quad (18)$$

where r is a random number with uniform distribution in $[0 - 1]$, λ_i is a random number with Cauchy (Lévy) distribution, and β is the *constriction coefficient* that controls the step size of the flight. Finally, the maximum delay constraint is enforced on the new position of the sub-particle to keep its feasibility by resetting the delay to zero if the flight exceed the allotted maximum delay.

Borrowing from QPSO, the attractor of the sub-particle a_i can be thought of as a point randomly chosen in the hyper-plane that connects the particle's best position and the global best position. This attractor is the next starting point for the Lévy flight, and the next equation updates the sub-particle position with the value of the attractor modified by the flight, which is itself a random number generated with Cauchy distribution multiplied by the value $(a_i - x_i)$, i.e., the difference between the attractor and the current position.

This value provides the scaling factor of the flight around the attractor and is crucial in the balance between exploration of new solutions and exploitation focusing in the proximity of the current solution.

On the other hand the β parameter need not be modified in the course of the algorithm and is tied to a probability density function "attitude" to generate large numbers, for instance with Cauchy ($\alpha = 1$) we set $\beta = 0.35$, while with a

general Lévy with $\alpha = 1.4$, which we have found as a good Lévy coefficient, we set $\beta = 0.75$.

Finally, note that contrary to QPSO formulation, we always execute the random flight "away from" the current position x_i . In fact $(a_i - x_i)$ is indeed a signed value that provides the direction of the attractor away from the current position. We have found through experiments that this gives better results in the exploration, trying to explore away from current "beaten track".

B. Modeling Constraints and Objective Function

An important aspect of the algorithm is the formulation of the objective function and the time and energy constraints. As described before, feasibility time constraints are enforced directly when updating the particles positions within their lower and upper bound limits.

Energy constraints are instead formulated as penalty components of the objective function to be minimized, which is defined as the sum of three elements with their respective weighting coefficients:

$$\text{minimize: } f(x) = w_1 O(x) + w_2 C(x) + w_3 T(x)$$

where w_1 , w_2 , and w_3 are the weight coefficients assigned respectively to the overload amount $O(x)$, the energy cost $C(x)$, and the tardiness in the execution $T(x)$.

The overload is the penalty component: if there is power overload the constraint is violated and therefore all other components can be ignored as their contribution would me trifle to the whole objective function $f(x)$. As such w_1 is chosen large enough to privilege constraints satisfaction before anything else, $w_1 \approx 10^9$.

w_2 is the weight of the energy cost and is normalized to the value 1. Finally the tardiness component is an added contributing element to the objective function and corresponds to perceived user comfort and tends to privilege schedule solutions that complete sooner rather than later (tardiness of the execution). The relation between w_2 and w_3 is the "sensible" balance between low energy cost on the one hand and low tardiness on the other. We typically set w_3 small enough (e.g., 10^{-3}) so as to attribute much more importance to energy cost, but still prefer earliest completions within a very small cost difference.

VI. SIMULATION AND RESULTS

We ran a number of simulation modeling the same scheduling problem both in the QPSOL algorithm and a pure mathematical model with commercial linear programming (LP) solvers, namely XPress and CPLEX. The scheduling problem was formalized with 4 instances of washing-machine power profiles, each profile being made of 4 phases, and 3 instances of dish-washing-machines each made of 5 phases, for a total of 31 independent variables to optimize in the scheduling problem instance.

Due to the hard problem space for the brute-force exact algorithms, the scheduling horizon was limited to 12 hours and the time slots at multiples of 3 minutes, otherwise, with

one-minute slot time, no feasible solutions were found even in 7 days of uninterrupted run. Running 96 hours, XPress found a solution at a cost of € 2.57358. With the same problem and running 1 hour CPLEX found a solution at € 2.59123. Finally the QPSOL was given a bound time of 15 seconds, and run 10 times to have reliable statistics, finding a best solution at € 2.7877, with an average cost of € 2.9351 for the 10 times.

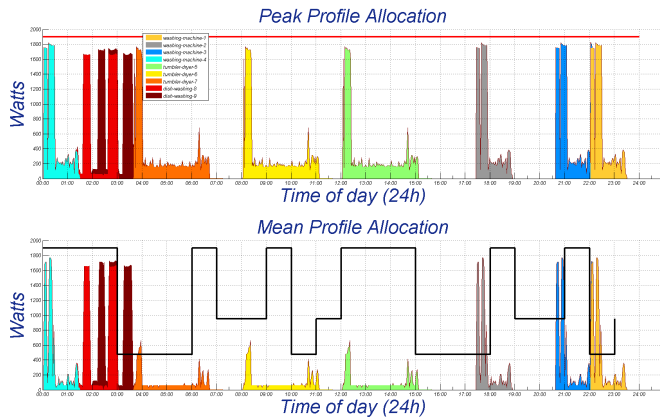


Figure 7. QPSOL simulation results: appliance scheduling with constant overload threshold, variable tariff, no photovoltaic.

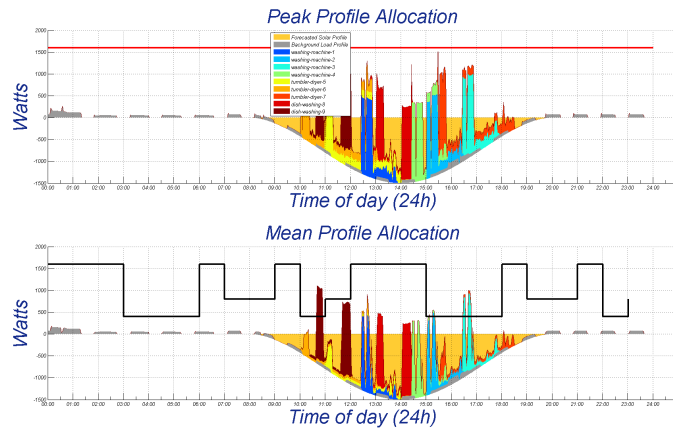


Figure 8. QPSOL simulation results: appliance scheduling with constant overload threshold, variable tariff, photovoltaic.

The results obtained using linear programming and exact solvers are very important as they fix theoretical optima for benchmarking the convergence and performance of the metaheuristic approach of the QPSOL. Results show that although QPSOL finds a worse solution than the theoretical optimum by a 8 – 13 %, the very short allotted time to find a solution is anyway a very promising approach. In Figure 7. and Figure 8. are reported simulation results when considering appliance scheduling with constant overload threshold, variable tariff, and with the absence and presence of photovoltaic generation respectively.

An interesting use case is the scheduling of an entire apartment house where tenants share a common contract with the utility provider in which the energy consumption of the apartment house as a whole must be below a given “virtual” threshold that changes in time. Figure 9. shows such scenario. The curved red line represents the virtual threshold that the apartment house should respect. All energy above such threshold will not cause an overload but its cost grows exponentially with the net effect of encouraging a peak shaving of profile allocation. The case study of Figure 10. is a scheduling of 15 apartments, with 3 appliances each, for a total of 45 appliances. The apartment house is also provided with common PV-panels.

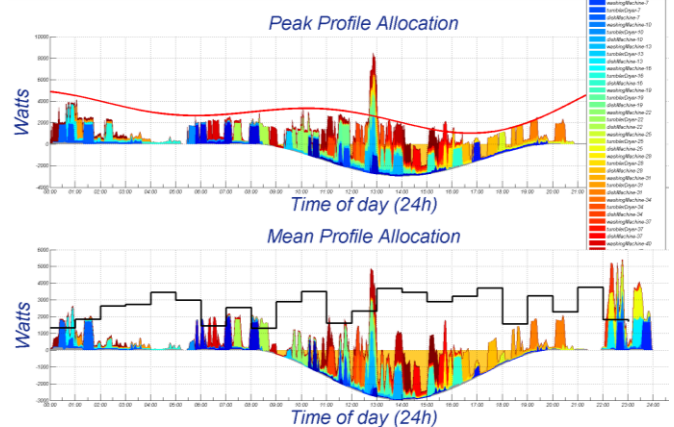


Figure 9. QPSOL simulation results: appliance scheduling for different apartments with variable overload threshold, variable tariff, photovoltaic.

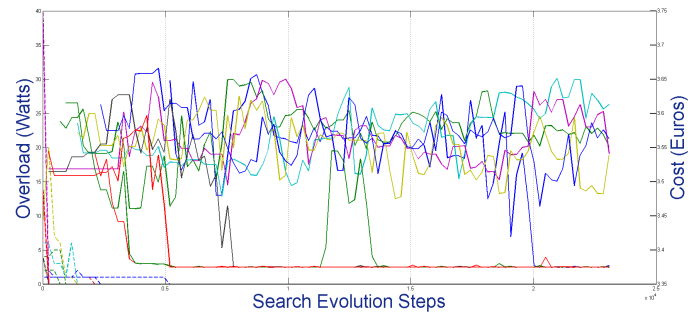


Figure 10. QPSOL simulation results: overload avoidance and optimization of cost

The 3 case studies described here show the remarkable flexibility of the QPSOL algorithm, and many other metaheuristic algorithms for that matter, i.e. the ability to adapt the algorithm to the unique attributes of a given problem and not based on predefined characteristics. In a rapidly changing world, algorithmic paradigms that are the most flexible to new conditions and can contribute to a time-based competitive advantage are more likely to be utilized. In such volatile environments, the utility of an algorithm framework will not be derived from the ability to solve a static problem. Instead it will be the ability to adapt

to changing problem conditions that is likely to define the success or failure in the optimization algorithms of tomorrow.

Exact and formal techniques decompose the optimization problems into mathematically tractable domains involving precise assumptions and well-defined problem classes. However many practical optimization problems are not strictly members of these problem classes, and this becomes especially relevant for problems that are non-stationary during their lifecycle. Mathematical techniques not only place constraints on the current problem definition but also on how that problem definition can change over time. Under these circumstances, long-term algorithm survival / popularity is less likely to reflect the performance of the canonical algorithm and instead more likely reflects success in algorithm design modification across problem contexts [20].

VII. CONCLUSION

This document describes an innovative Quantum inspired PSO with Lévy flights metaheuristic algorithm for scheduling home smart appliances, capturing all relevant appliance operations. With appropriately dynamic tariffs and short-term load forecasting, the proposed framework can propose a schedule for achieving high cost savings and overloads prevention. Good quality approximate solutions can be obtained in a short amount of computation time, in the order of about 2 seconds an almost optimal approximate solution can be obtained.

Finally, the proposed framework can be extended to incorporate solar power forecasting in the presence of a residential PV system by tuning the objective function and using the solar energy forecaster as further input to the scheduler.

ACKNOWLEDGMENT

This work has been partially supported by INTRÉPID, INTelligent systems for Energy Prosumer buildIngs at District level, funded by the European Commission under FP7, Grant Agreement N. 317983.

The authors would like to thank Prof. Della Croce of Operational Research department of the Politecnico di Torino for the valuable insights and contribution on the linear programming solvers.

REFERENCES

- [1] Energy@Home project, "Energy@Home Technical Specification version 0.95," December 22, 2011.
- [2] K. Cheong Sou, J. Weimer, H. Sandberg, and K. Henrik Johansson, "Scheduling Smart Home Appliances Using Mixed Integer Linear Programming," 50th IEEE Conference on Decision and Control and European Control Conference (CDC-ECC), Orlando, FL, USA, December 12-15, 2011.
- [3] J. Holland, "Adaptation in Natural and Artificial systems", University of Michigan Press, Ann Arbor, 1995.
- [4] F. Glover, and M. Laguna, "Tabu Search", Kluwer Academic Publishers, Boston, 1997.
- [5] S. Kirkpatrick, C. D. Gellat, and M.P. Vecchi, "Optimization by Simulated Annealing", Science, 220, pp. 671-680, 1983.
- [6] J. Kennedy, and R. Eberhart, "Particle Swarm Optimization", in: Proc. of the IEEE Int. Conf. on Neural Networks, Piscataway, NJ, pp. 1942-1948, 1995.
- [7] J. Sun, B. Feng, and W. Xu, "Particle swarm optimization with particles having quantum behavior," in IEEE Congress on Evolutionary Computation, pp. 325-31, 2004.
- [8] X. Yang, "Nature-Inspired Metaheuristic Algorithms", Luniver Press, 2008.
- [9] X. Yang "Review of metaheuristics and generalized evolutionary walk algorithm", Int. J. Bio-Inspired Computation, vol. 3, No. 2, pp. 77-84, 2011.
- [10] A. Chechkin, R. Metzler, J. Klafter, V. Gonchar, "Introduction to the theory of lévy flights." In: Klages R, Radons G, Sokolov IM (eds) Anomalous Transport: Foundations and Applications, Wiley-VCH, Berlin, 2008.
- [11] D. Ionescu, A. Juan, J. Faulin, and A. Ferrer, "A Parameter-Free Approach For Solving Combinatorial Optimization Problems Through Biased Randomization Of Efficient Heuristics", in Proceedings of the Conference on Numerical Optimization and Applications in Engineering (NUMOPEN-2010), Barcelona, Spain. October 13-15, 2010.
- [12] R. Kolisch, and S. Hartmann, "Heuristic Algorithms for Solving the Resource-Constrained Project Scheduling Problem: Classification and Computational Analysis", pp. 147-178, Kluwer, Amsterdam, the Netherlands, Kluwer academic publishers, 1999.
- [13] R. Kolisch, and S. Hartmann, "Experimental Investigation of Heuristics for Resource-Constrained Project Scheduling: An Update", European Journal of Operational Research 174, pp. 23-37, Elsevier, 2006.
- [14] J. W. Taylor, "Short-Term Electricity Demand Forecasting Using Double Seasonal Exponential Smoothing" Said Business School University of Oxford - Journal of Operational Research Society, vol. 54, pp. 799-805, 2003.
- [15] J. W. Taylor and P. E. McSharry, "Short-Term Load Forecasting Methods: An Evaluation Based on European Data", IEEE Transactions on Power Systems, vol. 22, pp. 2213-2219, 2008.
- [16] J. W. Taylor "Short-Term Load Forecasting with Exponentially Weighted Methods", IEEE Transactions on Power Systems, vol. 27, pp. 458-464, February 2011.
- [17] Í. Goiri, K. Le, M. E. Haque, R. Beauchea, T. D. Nguyen, J. Guitart, J. Torres, and R. Bianchini "GreenSlot: Scheduling Energy Consumption in Green Datacenters", SC'11, Seattle, Washington, USA, November 2011.
- [18] N. Sharma, J. Gummesson, D. Irwin, and P. Shenoy, "Cloudy Computing: Leveraging Weather Forecasts in Energy Harvesting Sensor Systems", SECON 2010, Boston, MA, June 2010.
- [19] P. Bacher, H. Madsen, and H. A. Nielsen, "Online Short-term Solar Power Forecasting", Sol. Energy, vol. 83, pp. 1772-1783, 2009.
- [20] J. M. Whitacre "Survival of the flexible: explaining the recent dominance of nature-inspired optimization within a rapidly evolving world", Journal Computing, Vol. 93, Issue 2-4 , pp 135-146 2009.

Transient State Analysis of the Multichannel EMG Signal Using Hjorth's Parameters for Identification of Hand Movements

Michele Pla Mobarak, Juan Manuel Gutiérrez Salgado,
Roberto Muñoz Guerrero
Department of Electrical Engineering, Bioelectronics
CINVESTAV
Mexico City, Mexico
e-mail: {mpla, mgutierrez, rmunoz}@cinvestav.mx

Valérie Louis-Dorr
Centre de Recherche en Automatique de Nancy (CRAN)
CNRS : UMR7039 – University of Lorraine
Nancy, France
e-mail: valerie.louis@univ-lorraine.fr

Abstract—Most myoelectric controlled systems are based on the common assumption that there is no information in the instantaneous value of the myoelectric signal and therefore, analysis is made on the steady state of the muscle contraction. However, this control scheme faces two main drawbacks: users need to be trained in order to produce the sustained contractions, and the control signal can only be generated until the steady state is reached. Prosthetic devices with long actuating delays often result in users' frustration and eventually, in the abandonment of the devices. As a proposed solution, analysis of the transient state of the electromyography (EMG) signal would allow classifying movements during the dynamic part of the muscle contractions reducing the time required to generate control commands. This paper proposes a novel method for transient EMG classification based on the use of Hjorth's parameters. Surface multichannel EMG signals were recorded from 10 normally limbed subjects for both the transient and steady EMG states while performing six different hand motions. Comparatively high classification accuracy was obtained from the transient state analysis of the signals suggesting the existence of deterministic information in this part of the muscle contraction and the fact that Hjorth's parameters seem to adapt well enough to the nature of myoelectric signals as to allow extracting highly representative information from them.

Keywords—EMG steady state; EMG transient state; Hjorth's parameters; multichannel EMG; normalized slope descriptors

I. INTRODUCTION

Research on the use of electromyography (EMG) for upper limb prostheses control has been conducted since the 1940's [1]. In the following years, remarkable progress in myoelectric controlled devices has been achieved; however, still numerous challenges remain in signal processing in order to replicate as close as possible the functions of the human limb. Development on this area must be closely related to the understanding of the psychological complexities that the amputee faces [2] due to the fact that the prosthetic device becomes an extension of the patient's body.

High abandonment rates reported for upper limb prosthesis [3] make evident the need to improve myoelectric control. The performance of a myoelectric controlled system is evaluated with regards to three important aspects of

controllability [4]: the accuracy of movement selection, the intuitiveness of the actuating control, and the response time of the control system. A 200 to 300 ms interval is a clinically recognized maximum delay that users find acceptable before they get frustrated with the response time of the prosthesis [4]-[6]. Hence, the motivation to analyze the transient state of the EMG signal arises in order to identify movements while the muscle contraction is being generated and not until it reaches a steady state.

The EMG signal is a non-stationary, non-linear, and stochastic process produced as a result of the summation of several motor unit action potential trains (MUAPTs) [1],[7]-[9]. However, two main states can be recognized during the muscle contraction. The transient state is described as the bursts of myoelectric activity that accompany sudden muscular effort while executing the movement. It is related with the beginning of the activation of the motor units (MUs) that will be involved in the muscle contraction. The steady state corresponds to the part of the contraction when almost every MU that will be involved in the movement is already activated. It will be considered as the muscular effort during a sustained contraction when the movement's final position is reached, and the muscle length is no longer modified, i.e., the myoelectric signal produced by a stable muscle contraction [10],[11].

The aim of this work is to propose a novel method for transient state analysis of the multichannel EMG signal by using normalized slope descriptors (NSDs) as features for classification. The main objectives are proving that Hjorth's parameters are suitable for EMG analysis, and demonstrating the possibility of identifying movements from the beginning of the muscle contraction in order to reduce the delay obtained when waiting until the steady state is reached to generate the control signal.

In the following sections, this paper provides an overview of the state of the art in the use of the transient EMG state for movement identification followed by an introduction to Hjorth's parameters and their use in biomedical applications. Subsequently, it explains the method that was followed, including the data acquisition protocol and the signal processing method, as well as the results obtained from the analysis of both the transient and steady EMG signals recorded from 10 normally limbed

subjects in order to identify a set of hand and wrist movements. Results are discussed and compared to previously reported methods. Finally, the last section presents conclusions and some ideas on the future work that can be done.

II. STATE OF THE ART IN THE USE OF THE TRANSIENT EMG STATE FOR MYOELECTRIC CONTROL

EMG classification has been most often based on the steady state analysis of the muscle contraction. This has greatly simplified commercial myoelectric controlled systems that usually rely on the premise of the accepted myoelectric signal generation models. However, the steady state contains a short temporal structure of the active modification of recruitment and firing patterns involved in the contraction and that can be found within the transient state [1],[5],[6],[10]. In 1993, Hudgins *et al.* [5] were the first to consider the structure in the myoelectric signal (MES) during the onset of the contraction to develop a new control strategy based on the analysis of the transient EMG state. They were able to discriminate between four movements with roughly 90% accuracy. Only a few studies, such as [12]-[15], have reported the use of the transient state to classify EMG signals. Englehart *et al.* [6] introduced the use of Wavelet Transform (WT) and Wavelet Packet Transform (WPT) for classification of transient EMG signals. They classified more accurately the steady state than the transient data. However, Hangrove *et al.* [16] showed that including transient data along with steady state data while training the classifier increases the classification error, but it also increases real-time performance and system usability, which should be considered when evaluating the system.

The use of multichannel EMG provides a better representation of the real muscle activity in the collected signal [17]-[19]. The increase in classification performance while increasing the number of channels was investigated in [20]. Moreover, with multichannel EMG, the positions of the electrodes become less critical [21], making it a promising technique. However, when extrapolating the system to amputees, an excessive number of electrode sites could be hard or even not possible to locate.

Interference and muscle crosstalk introduce non-linearity into the EMG signal. The combination of muscle tissue, adipose tissue, skin, and the skin-electrode interface behaves like a non-linear low pass filter that attenuates and distorts the surface EMG signal; nevertheless, methods for non-linear time series analysis have not been widely applied to EMG [7].

III. HJORTH'S PARAMETERS

Hjorth introduced, in 1970, three parameters based on time domain properties [22]-[25]. They were intended as a clinically useful tool capable of describing quantitatively the graphical characteristics of an electroencephalography (EEG) trace in terms of amplitude, slope, and slope spread, so that they receive the name of normalized slope descriptors (NSDs). These parameters are named "activity", "mobility", and "complexity".

Activity measures the variance of the amplitude of the signal as shown in (1). In the frequency domain, it can be conceived as the envelope of the power spectrum.

Mobility measures the ratio between the standard deviation of the slope and the standard deviation of the amplitude given per time unit; hence, it represents dominant frequency. This ratio depends on the curve shape in such a way that it measures the relative average slope. Its mathematical definition is presented in (2).

Complexity is a dimensionless parameter that quantifies any deviation from the sine shape as an increase from unit. It is calculated as shown in (3). It can be interpreted as a measure of the signal's bandwidth.

$$Activity = m_0 = \sigma_0^2, \quad (1)$$

$$Mobility = \sqrt{m_2/m_0} = \sigma_1/\sigma_0, \quad (2)$$

$$Complexity = \sqrt{(m_4/m_2) - (m_2/m_0)} = \frac{\sigma_2/\sigma_1}{\sigma_1/\sigma_0}, \quad (3)$$

where m_n is the spectral moment at order n , σ_0^2 is the variance from the analyzed segment of the non-linear time series $f(t)$, and σ_1 and σ_2 are the standard deviations of the first and second derivatives of $f(t)$, respectively. It has been shown that the spectral moment of order $2n$ corresponds to the variance σ_n^2 of the derivative of order n [24], so that:

$$m_0 = \sigma_0^2, m_2 = \sigma_1^2, m_4 = \sigma_2^2, m_6 = \sigma_3^2 \dots m_{2n} = \sigma_n^2$$

The spectral moment m_n can also be calculated in terms of its frequency as shown in (4).

$$m_n = \int_{-\infty}^{+\infty} \omega^n \cdot S(\omega) d\omega \quad (4)$$

$S(\omega)$ corresponds to the power density spectrum, and it is obtained from the multiplication of the Fourier Transform, $F(\omega)$, by its conjugate, $F^*(\omega)$, which causes the phase to be excluded. As the frequency description from the Fourier transform is always symmetrical with respect to zero frequency, in a statistical approach to the shape of the frequency distribution, all odd moments will become zero, and the information will be contained in the even moments.

Hjorth's parameters serve as a bridge between a physical time domain interpretation and the conventional frequency domain description [23]. The transformation between both domains is based on the energy equality within the actual epoch and can be calculated by the time-frequency relationship shown in (5)-(7).

$$m_0 = \int_{-\infty}^{+\infty} S(\omega) d\omega = \frac{1}{T} \int_{t-T}^t f^2(t) dt = \sigma_0^2 \quad (5)$$

$$m_2 = \int_{-\infty}^{+\infty} \omega^2 S(\omega) d\omega = \frac{1}{T} \int_{t-T}^t \left(\frac{df(t)}{dt} \right)^2 dt = \sigma_1^2 \quad (6)$$

$$m_4 = \int_{-\infty}^{+\infty} \omega^4 S(\omega) d\omega = \frac{1}{T} \int_{t-T}^t \left(\frac{d^2 f(t)}{dt^2} \right)^2 dt = \sigma_2^2 \quad (7)$$

Hjorth's parameters were originally formulated for EEG analysis and description and have been widely used in sleep EEG processing for data reduction and discrimination of sleep stages [26]-[28]. Other studies related with EEG signal analysis have reported the use of Hjorth's parameters for applications such as psychotropic drug research [26],[27], assessment of postalcoholic diseases [29], temporal lobe seizures lateralization [25], classification of facial movement artifacts in the EEG signal [30], monitoring changes in EEG signals of patients with renal failure before and after hemodialysis [31], creating ink topographic displays for visual monitoring of changes in EEG signals [32], evaluation of performance in channel reduction for EEG classification in emotion assessment [33], among others. Mouz -Amady and Horwat [34] applied NSDs to EMG signals and concluded that they could be used to describe the spectral content of surface EMG during repetitive movements due to their results of high correlation coefficients (ranging from 0.81 to 0.93) between Hjorth's mobility and the FFT mean frequency. Hjorth's parameters have also been applied successfully in non-biomedical fields [22].

IV. METHOD

A. Data Acquisition Protocol

The surface EMG signals used for this study were those recorded in [35] using 8 differential channels (Ag-AgCl surface electrodes model VERMED NeuroPlus A10043 with an inter-electrode distance of 1.5 cm) placed on the dominant forearm of 10 normally limbed subjects, aged between 23 and 50, and with no register of neuromuscular disorders. The electrode disposition is shown in Fig. 1. To ensure the positioning of the electrodes over the muscles of interest, each participant was asked to repeatedly close and open the hand in order to identify the muscles mentioned in table I.

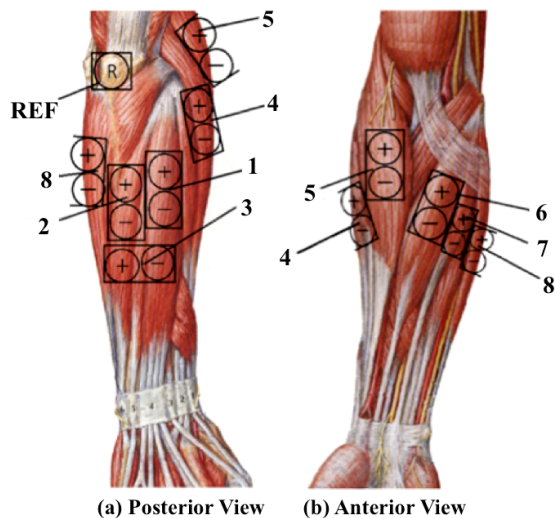


Figure 1. Posterior (a) and anterior (b) views of electrode placement for EMG signal recording. Corresponding muscles are identified in Table I.

TABLE I. FOREARM MUSCLES RECORDED BY EACH EMG DIFFERENTIAL CHANNEL

EMG Channel	Forearm Muscle
1	Extensor digitorum communis
2	Extensor carpi ulnaris
3	Differential measure between extensor digitorum communis and extensor carpi ulnaris
4	Extensor carpi radialis longus
5	Brachioradial
6	Flexor carpi radialis
7	Palmaris longus
8	Flexor carpi ulnaris

The skin was carefully cleaned before electrode placement, and the reference electrode was located on the elbow. Each subject was asked to execute six different hand motions, namely hand opening/closing, wrist pronation/supination, and wrist flexion/extension. The series of movements were repeated five times with a one-minute rest between them in order to avoid muscle fatigue. Each recording was 20 seconds long, starting with the forearm in an inactive position, followed by the dynamic part of the contraction, and finished by sustaining the contraction, once the final position was reached, until the end of the recording.

Written consent was obtained from every subject before starting the study. In a previous session, the protocol was explained to each of the participants and the amplification gain of the eight differential EMG channels was calibrated according to the amplitude of the contraction for each subject.

In order to reduce unwanted variability, every participant was asked to perform the study in a standing position with the dominant arm extended to the front and the hand relaxed.

The acquisition system consisted of 8 differential channels with adjustable amplification gain and a first order analog band-pass filter with a low cut-off frequency of 20 Hz and a high cut-off frequency of 400 Hz. Each analog output was connected to a National Instruments acquisition card (model DAQ-Card 6024WE) for 12-bit A/D conversion. EMG signals were recorded with a sampling rate of 1024 Hz.

B. Data Processing

The multichannel EMG signals were processed and analyzed using MATLAB® (version R2012b). Each recording was divided by supervision in transient and steady states. The transient state was extracted from the part of the EMG recording that corresponded to the dynamic part of the movement. The steady state consisted of the section in the recording where the final position was reached and the muscle contraction was sustained. A Hamming window was applied to segment the signal for feature extraction. Classification performance was tested with two window lengths, 256 and 128 ms, both with 50% sample overlap.

A wavelet shrinkage method at third level of decomposition and based on Stein's unbiased risk estimate

(SURE) was applied to each windowed segment for de-noising purposes and to narrow the signal’s frequency band. The decomposition level was chosen considering that the main concentration of energy in the surface EMG signal is located within the band of 50-150 Hz. The de-noised signals were rescaled using a noise level dependent estimation.

Once the window length was selected, Hjorth’s parameters were calculated per channel for each of the segments using (1)-(3). This was made for every movement. Each of these parameters constituted an independent input; they were arranged in rows in order to build the feature matrix. Four different input matrices were tested to evaluate which one yielded the best performance. The first one contained as inputs the three parameters for each of the channels. The other three consisted of just two independent inputs per channel excluding one of the parameters in each case, i.e., the first one excluded complexity, the second one excluded mobility, and the third one excluded activity.

C. Artificial Neural Network’s Parameters

For classification of the EMG signals, an artificial neural network (ANN) model was trained using the aforementioned feature matrices as inputs, containing 24 independent inputs when using the three parameters and 16 in the cases where two parameters were used. A Bayesian regulation backpropagation algorithm was used to train the model. The final ANN’s architecture depended on the available feature matrix dimension; however, the model consisted of only one hidden layer with 9 neurons in it. This number of neurons was defined based on experimental testing. For classification purposes, the network output was binary codified, i.e., a unique 3-digit combination of zeros and ones was used to identify each of the movements.

In order to evaluate the network’s performance, a k-fold cross-validation process (k=5) was carried out.

V. RESULTS

A. Window Length Selection

Two Hamming windows of different size were tested over the whole transient data set. The first one had a length of 256 ms and the second one of 128 ms. Both windows were applied with a 50% sample overlap. Table II and Fig. 2 show the classification percentage obtained for each window length.

TABLE II. MEAN CLASSIFICATION PERCENTAGE AND STANDARD DEVIATION ACCORDING TO WINDOW LENGTH

Movement Type	Window Length [ms]	
	256 ms	128 ms
Closing	97.71±2.95%	98.71±1.25%
Opening	97.43±2.11%	96.86±1.48%
Pronation	96.86±2.50%	94.86±3.03%
Supination	97.71±3.51%	96.43±3.03%
Flexion	100.00±0.00%	99.71±0.60%
Extension	99.14±1.38%	98.14±1.66%
Mean	98.14±0.99%	97.45±1.10%

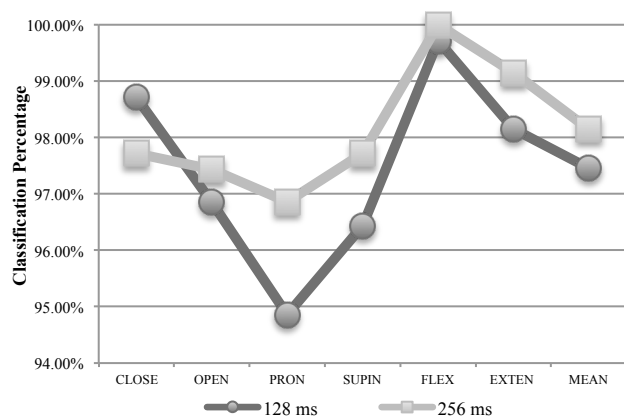


Figure 2. Comparison in classification percentage of the transient EMG state for each movement type using different window lengths.

Classification percentage did not decrease dramatically when using the 128 ms window length as compared to the 256 ms window; therefore, the smaller window size was used for the following tests considering that decreasing processing delays was one of the objectives.

B. Hjorth’s Parameters Selection

Classification performance was evaluated with four different feature matrices in order to select which of Hjorth’s parameters allowed extracting the most representative information from the signals. Each matrix consisted of a different set of Hjorth’s parameters as previously explained in section IV.

The test was applied on the transient state of the EMG signals. The obtained results are shown in Table III. The column labeled ‘A, M, and C’ contains the classification percentages for each movement type using the three Hjorth’s parameters. The three following columns denote the classification percentages obtained when using just two of the parameters.

The mean classification error percentage obtained from using each feature matrix was calculated and is presented in Fig. 3. These percentages consider the mean classification for the whole test population including every movement type.

TABLE III. MEAN CLASSIFICATION PERCENTAGE ACCORDING TO THE COMBINATION OF HJORTH’S PARAMETERS USED IN THE FEATURE MATRIX

Movement Type	Hjorth’s Parameters			
	A, M, and C	A and M	A and C	M and C
Closing	98.14%	96.57%	97.14%	98.71%
Opening	96.86%	95.43%	96.14%	96.86%
Pronation	94.00%	94.43%	93.00%	94.86%
Supination	95.57%	94.29%	93.29%	96.43%
Flexion	99.86%	99.57%	97.87%	99.71%
Extension	97.86%	98.29%	96.71%	98.14%
Mean	97.05%	96.43%	95.69%	97.45%
SD	2.95%	3.57%	4.31%	2.55%

A stands for ‘Activity’, M for ‘Mobility’, and C for ‘Complexity’

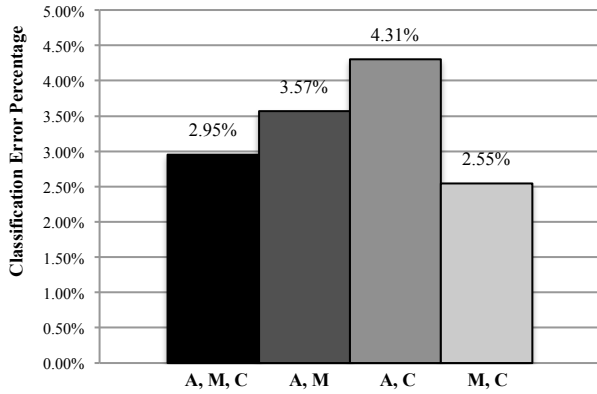


Figure 3. Mean classification error percentage obtained by using different combinations of Hjorth’s parameters as features for classification.

As the lowest classification error was obtained from using mobility and complexity as input features, the following comparison between classification performance using the transient and steady EMG states was made with this selected feature matrix.

The parameter ‘activity’ was no longer considered for the analysis.

C. Comparison in Classification Accuracy of the Transient and Steady EMG States

Previous studies such as [6], have reported higher classification accuracy when using the steady state of the EMG signal as compared to the transient state; therefore, the proposed method was evaluated for both EMG states. The mean classification percentages obtained for each of the subjects are presented in Table IV. The comparison in mean classification accuracy per movement type is illustrated in Fig. 4.

TABLE IV. TABLE TYPE STYLES

Subject	EMG State	
	Transient	Steady
1	97.38%	97.62%
2	99.05%	99.29%
3	97.62%	97.14%
4	95.71%	96.67%
5	97.86%	99.76%
6	96.67%	97.14%
7	98.33%	99.05%
8	96.19%	96.67%
9	96.90%	98.10%
10	98.81%	97.86%
Mean ± SD	97.45±1.10%	97.93±1.11%

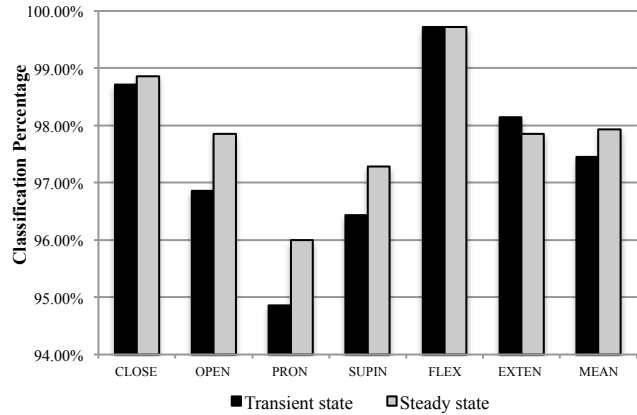


Figure 4. Mean classification percentage per movement type for the transient and steady EMG states.

The black bars in Fig. 4 represent the mean classification percentage of the transient state for each of the movements. The gray bars represent the mean classification of the steady state. These results are discussed in the following section.

VI. DISCUSSION

The proposed method departs from the use of features originally intended for EEG description. The EEG signal is formed by the superposition of characteristic responses; in a similar way, the EMG signal is the superposition of MUAPs, which seems to make Hjorth’s parameters also suitable for their analysis due to the nature of myoelectric sources.

The parameters’ values, except for activity, referring to a single response, are also valid for a superposition of responses [23], which can justify that classification error increases when including this parameter, as shown in Fig. 3. For non-periodic phenomena with a limited complexity, as it is the case of the EMG signal, the basic information is essentially contained in the first few polynomial coefficients; therefore, the number of required (non-redundant) parameters correspond to the complexity of the system under observation [24]. Several studies [25],[26],[28],[31] have reported to find significant information in EEG traces using just two of Hjorth’s parameters; furthermore, three of them reported to find them in mobility and complexity as it was the case in the present study.

Hjorth explained in [24] the way in which Hjorth’s parameters describe first and second order responses of a system. The response of a first order system is an exponentially decaying impulse ($e^{-\alpha t}$), where α is the inverse time constant of the system. When Hjorth’s parameters are computed for this response, mobility is identical to α and hence describes the system. The response of a second order system is modeled as a decaying sinusoid ($e^{-\alpha t} \cdot \sin(\beta \cdot t)$). The relationship between Hjorth’s parameters and the constants of the system are given in (8). By making some algebraic manipulations, the constants can be expressed as a function of Hjorth’s parameters as in (9).

$$\begin{cases} M(mobility) = (\alpha^2 + \beta^2)^{\frac{1}{2}} \\ C(complexity) = 2\alpha \end{cases} \quad (8)$$

$$\begin{cases} \alpha = C/2 \\ \beta = M \cdot (1 - C^2/4M^2)^{\frac{1}{2}} \end{cases} \quad (9)$$

Based on the previous assumptions, the system is completely determined by mobility and complexity. The number of required descriptors to obtain the basic information of the system corresponds to the system's order. The muscle contraction can be modeled as a second order system, which would justify the performance achieved by using these two parameters.

The mean classification percentage obtained from the steady EMG signals was just slightly higher than the one reported for the transient state (refer to Fig. 4). This suggests that similar classification accuracy can be obtained from both EMG states using Hjorth's parameters as compared to previous reported methods [6].

The lowest classification percentages using both the transient and steady states corresponded to pronation and supination movements. This indicates that the muscles recorded by the EMG channels, which are mainly flexors and extensors, had less participation in these two movement types. Classification accuracy did not suffer a significant decrease when using a window length of 128 ms as compared to one of 256 ms, which could allow obtaining a faster response for control purposes and would make it more suitable for on-line applications.

VII. CONCLUSIONS AND FUTURE WORK

It has been widely discussed whether or not the transient state of the muscle contraction contains enough relevant information as to accurately discriminate between different types of motions; however, Hjorth's parameters seem to adapt well enough to the transient MES as to extract highly representative information from it. Most myoelectric controlled devices have been based on the assumption that there is no information in the instantaneous value of the MES; therefore, it is necessary to wait until a sustained stable contraction is reached in order to generate the control signal and start actuating the device, which is not desirable in clinical applications. The present study reaffirms the existence of deterministic components within the onset part of the muscle contraction as initially proposed by Hudgins *et al.* [5], and the fact that the information is relevant enough as to discriminate between the proposed movements. Using the transient EMG state for myoelectric control would allow generating a control signal since the beginning of the muscle contraction, which would also resemble more to the natural movement of the human limb and would allow diminishing the actuating delay of the devices. Even if classification is slightly higher when using the steady state, including transient MES information can lead to more robust usability and performance. If the system is capable of identifying the transient EMG state, the subject is simply prompted to perform a contraction in a natural manner instead of needing long training periods to learn how to make the sustained contractions.

Hjorth's parameters allow characterizing signals in both the time and frequency domains. Although they are based on

spectral moments, they can be calculated using time variances, which implies simpler processing and makes it more suitable for continuous on-line calculations as compared to frequency-domain analysis that normally requires complex transformations.

The analyzed data proved to contain significant information within the first 128 ms of the onset of the contraction. Classification was more accurate when using a 256 ms window for feature extraction; however, this value gets really close to the clinically recognized maximum delay for real-time applications (200 to 300 ms), and the total system's delay has yet to be considered. Thus, a 128 ms window seems to provide a good compromise between system's accuracy and response time.

Using multichannel EMG signals allows recording information from several muscle sites, which allows analyzing the participation of corresponding muscles on a certain movement; however, when using a great number of electrodes, extrapolating the system to an amputee patient becomes very complex and sometimes not feasible.

As future work, we propose to use a greater number of subjects to test the algorithm's performance; to develop a reliable method to automatically identify the onset of the contraction regardless of the noise affecting the system, which would allow automatically initiating feature extraction in order to generate the control signal; to increase the number of movements to identify; and to extrapolate the system to amputee subjects.

ACKNOWLEDGMENT

The authors would like to acknowledge the master scholarship granted to Michele Pla by CONACYT.

REFERENCES

- [1] M. Zecca, S. Micera, M. C. Carrozza, and P. Dario, "Control of multifunctional prosthetic hands by processing the electromyographic signal," *Critical Reviews in Biomedical Engineering*, vol. 30, no. 4-6, 2002, pp. 459-485, doi:10.1615/CritRevBiomedEng.v30.i456.80.
- [2] D. Desmond and M. MacLachlan, "Psychological issues in prosthetic and orthotic practice: a 25 year review of psychology in Prosthetics and Orthotics International," *Prosthetics and Orthotics International*, vol. 26, no. 3, Dec. 2002, pp. 182-188, doi:10.1080/03093640208726646.
- [3] E. A. Biddiss and T. Chau, "Upper limb prosthesis use and abandonment: a survey of the last 25 years," *Prosthetics and Orthotics International*, vol. 31, no. 3, Sept. 2007, pp. 236-257, doi:10.1080/03093640600994581.
- [4] K. Englehart and B. Hudgins, "A robust, real-time control scheme for multifunction myoelectric control," *IEEE Trans. Biomed. Eng.*, vol. 50, no. 7, July 2003, pp. 848-854, doi:10.1109/TBME.2003.813539.
- [5] B. Hudgins, P. Parker, and R. N. Scott, "A new strategy for multifunction myoelectric control," *IEEE Trans. Biomed. Eng.*, vol. 40, no. 1, Jan. 1993, pp. 82-94, doi:10.1109/10.204774.
- [6] K. Englehart, B. Hudgins, and P. A. Parker, "A wavelet based continuous classification scheme for multifunction myoelectric control," *IEEE Trans. Biomed. Eng.*, vol. 48, no. 3, Mar. 2001, pp. 302-311, doi:10.1109/10.914793.

- [7] P. Padmanabhan and S. Puthusserypady, "Nonlinear analysis of EMG signals – a chaotic approach," *Proc. IEEE Eng. Med. Biol. Soc. (IEMBS '04)*, IEEE Press, vol. 1, Sept. 2004, pp. 608-611, doi:10.1109/IEMBS.2004.1403231.
- [8] J. Rafiee, M. A. Rafiee, F. Yavari, and M. P. Schoen, "Feature extraction of forearm EMG signals for prosthetics," *Expert Systems with Applications*, vol. 38, 2011, pp. 4058-4067, doi:10.1016/j.eswa.2010.09.068.
- [9] A. Merlo, D. Farina, and R. Merletti, "A fast and reliable technique for muscle activity detection from surface EMG signals," *IEEE Trans. Biomed. Eng.*, vol. 50, no. 3, March 2003, pp. 316-323, doi:10.1109/TBME.2003.808829.
- [10] H. A. Romo, J. C. Realpe, and P. E. Jojoa, "Surface EMG signal analysis and its applications in hand prosthesis control (Análisis de señales EMG superficiales y su aplicación en el control de prótesis de mano)," *Revista Avances en Sistemas e Informática*, vol. 4, no. 1, Jun. 2007, pp. 127-136.
- [11] T. W. Calvert and A. E. Chapman, "The relationship between the surface EMG and force transients in muscle: simulation and experimental studies," *Proc. IEEE*, vol. 65, no.5, May 1977, pp.682-689, doi:10.1109/PROC.1977.10547.
- [12] U. Kuruganti, B. Hudgins, and R. N. Scott, "Two-channel enhancement of a multifunction control system," *IEEE Trans. Biomed. Eng.*, vol. 42, no. 1, Jan. 1995, pp. 109-111, doi:10.1109/10.362912.
- [13] K. Englehart, B. Hudgins, M. Stevenson, and P. A. Parker, "Classification of myoelectric signal burst patterns using a dynamic neural network," *Proc. IEEE 21st Annual Northeast Bioengineering Conference (NEBC)*, May 1995, pp. 63-64, doi:10.1109/NEBC.1995.513734.
- [14] K. Englehart, B. Hudgins, M. Stevenson, and P. A. Parker, "A dynamic feedforward neural network for subset classification of myoelectric signal patterns," *Proc. IEEE 17th Annual Conference, Engineering in Medicine and Biology Society (IEMBS)*, vol. 1, Sept. 1995, pp. 819-820, doi:10.1109/IEMBS.1995.575359.
- [15] S. Karlsson, J. Yu, and M. Akay, "Time-frequency analysis of myoelectric signals during dynamic contractions: a comparative study," *IEEE Trans. Biomed. Eng.*, vol. 47, no. 2, Feb. 2000, pp. 228-238, doi:10.1109/10.821766.
- [16] L. Hangrove, Y. Losier, B. Lock, K. Englehart, and B. Hudgins, "A real-time pattern recognition based myoelectric usability study implemented in a virtual environment," *Proc. IEEE 29th Annual International Conference of the Engineering in Medicine and Biology Society (EMBS 2007)*, IEEE Press, Aug. 2007, pp. 4842-4845, doi:10.1109/IEMBS.2007.4353424.
- [17] D. Staudenmann, I. Kingma, D. F. Stegeman, and J. H. van Dieën, "Towards optimal multi-channel EMG electrode configurations in muscle force estimation: a high density EMG study," *Journal of Electromyography and Kinesiology*, vol. 5, Feb. 2005, pp. 1-11, doi:10.1016/j.jelekin.2006.08.006.
- [18] M. F. Lucas, A. Gauffriau, S. Pascual, C. Doncarli, and D. Farina, "Multi-channel surface EMG classification using support vector machines and signal-based wavelet optimization," *Biomedical Signal Processing and Control*, vol. 3, no. 2, Apr. 2008, pp. 169-174, doi:10.1016/j.bspc.2007.09.002.
- [19] E. A. Clancy and N. Hogan, "Multiple site electromyograph amplitude estimation," *IEEE Trans. Biomed. Eng.*, vol. 42, no. 2, Feb. 1995, pp. 203-211, doi:10.1109/10.341833.
- [20] K. Davidge, "Multifunction myoelectric control using a linear electrode array," M.Sc. dissertation, Electronics and Electrical Engineering, University of New Brunswick, Canada, 2005.
- [21] B. Karlik, M. O. Tokhi, and M. Alci, "A novel technique for classification of myoelectric signals for prosthesis," CD-ROM Proc. of the 15th Triennial World Congress of the International Federation of Automatic Control (IFAC), vol. 2, Jul. 2002, pp. 978-982, doi:10.3182/20020721-6-ES-1901.00980.
- [22] P. P. Balestrassi, A. P. Paiva, A. C. Zambroni de Souza, J. B. Turroni, and E. Popova, "A multivariate descriptor method for change-point detection in nonlinear time series," *Journal of Applied Statistics*, vol. 38, no. 2, Feb. 2011, pp. 327-342, doi:10.1080/02664760903406496.
- [23] B. Hjorth, "EEG analysis based on time domain properties," *Electroencephalography and Clinical Neurophysiology*, vol. 29, no. 3, Sept. 1970, pp. 306-310, doi:10.1016/0013-4694(70)90143-4.
- [24] B. Hjorth, "The physical significance of time domain descriptors in EEG analysis," *Electroencephalography and Clinical Neurophysiology*, vol. 34, no. 3, Mar. 1973, pp. 321-325, doi:10.1016/0013-4694(73)90260-5.
- [25] T. Cecchin, et al., "Seizure lateralization in scalp EEG using Hjorth parameters," *Clinical Neurophysiology*, vol. 121, no. 3, Mar. 2010, pp. 290-300, doi:10.1016/j.clinph.2009.10.033.
- [26] O. Kanno and P. Clarenbach, "Effect of clonidine and yohimbine on sleep in man: polygraphic study and EEG analysis by normalized slope descriptors," *Electroencephalography and Clinical Neurophysiology*, vol. 60, no. 6, Jun. 1985, pp. 478-484, doi:10.1016/0013-4694(85)91107-1.
- [27] H. Depoortere, D. Francon, P. Granger, and M. G. Terzano, "Evaluation of the stability and quality of sleep using Hjorth's descriptors," *Physiology & Behavior*, vol. 54, no. 4, Oct. 1993, pp. 785-793, doi:10.1016/0031-9384(93)90093-U.
- [28] M. Ziller, et al., "Bivariate global frequency analysis versus chaos theory," *Neuropsychobiology*, vol. 32, no. 1, 1995, pp. 45-51, doi:10.1159/000119211.
- [29] W. Spehr and G. Stemmler, "Postalcoholic diseases: diagnostic relevance of computerized EEG," *Electroencephalography and Clinical Neurophysiology*, vol. 60, no. 2, Feb. 1985, pp. 106-114, doi:10.1016/0013-4694(85)90016-1.
- [30] S. Pourzare, O. Aydemir, and T. Kayikcioglu, "Classification of various facial movement artifacts in EEG signals," *Proc. 35th International Conference on Telecommunications and Signal Processing (TSP)*, IEEE Press, July 2012, pp. 529-533, doi:10.1109/TSP.2012.6256351.
- [31] W. Spehr, et al., "EEG and hemodialysis. A structural survey of EEG spectral analysis, Hjorth's EEG descriptors, blood variables and psychological data," *Electroencephalography and Clinical Neurophysiology*, vol. 43, no. 6, Dec. 1977, pp. 787-797, doi:10.1016/0013-4694(77)90001-3.
- [32] A. Persson and B. Hjorth, "EEG topogram – An aid in describing EEG to the clinician," *Electroencephalography and Clinical Neurophysiology*, vol. 56, no. 5, Nov. 1983, pp. 399-405, doi:10.1016/0013-4694(83)90221-3.
- [33] K. Ansari-Asl, G. Chanel, and T. Pun, "A channel selection method for EEG classification in emotion assessment based on synchronization likelihood," *Proc. 15th European Signal Processing Conference (EUSIPCO 2007)*, EURASIP, Sept. 2007, pp. 1241-1245.
- [34] M. Mouzé-Amady and F. Horwat, "Evaluation of Hjorth parameters in forearm surface EMG analysis during an occupational repetitive task," *Electroencephalography and Clinical Neurophysiology*, vol. 101, no. 2, Apr. 1996, pp. 181-183, doi:10.1016/0924-980X(96)00316-5.
- [35] M. León Ponce, "Classification of myoelectric patterns for the operation of an antropomorphic device (Clasificación de patrones mioeléctricos para la operación de un dispositivo antropomórfico)," Ph.D. dissertation, Bioelectronics Program, Electrical Engineering Department, CINVESTAV, Mexico City, Mexico, Dec. 2012.

Realising Duality Principle for Prognostic Models

Mohammad Samie
School of Applied Sciences
Cranfield University
Bedford, UK
m.samie@cranfield.ac.uk

Amir M. S. Motlagh
Faculty of Science and Technology
University of Westminster
London, UK
w1418752@my.westminster.ac.uk

Alireza Alghassi
School of Applied Sciences
Cranfield University
Bedford, UK
a.alghassi@cranfield.ac.uk

Suresh Perinpanayagam
School of Applied Sciences
Cranfield University
Bedford, UK
suresh.nayagam@cranfield.ac.uk

Epaminondas Kapetanios
Faculty of Science and Technology
University of Westminster
London, UK
E.Kapetanios@westminster.ac.uk

Abstract— The lack of scientific approaches in estimating the remaining useful life (RUL) of various components and devices used in complicated systems, such as airplanes remain to be addressed. Regardless, there has been some progress in demonstrating feasible and viable techniques so far that are relevant to ‘integrated system health management’ (ISHM). ISHM entails a series of techniques and scientific measures that have collaborative self-awareness features to increase the overall reliability of systems. However, these resulting systems were often too expensive and time consuming, as well as requiring a lot of resources to develop. This paper presents a radically novel approach for building prognostic models that compensates and improves on the inconsistencies and problems witnessed in current prognostic models. Essentially, it proposes a state of the art technique that utilizes the physics of a system rather than the physics of a component. An advantage to this approach is; the prognostic model can be generalized such that a new system could be developed on the basis and principles of the prognostic model of another systems. Simple electronic circuits are to be used as an experiment to exemplify the potential success that can be discovered from the development of a novel prognostic model that can efficiently estimate the RUL of one system based on the prognostics of another system.

Keywords-*Prognostic Model, Integrated System Health Management (ISHM), Degradation, Duality, Cuk Converter.*

I. INTRODUCTION

Integrated System Health Management (ISHM) [1] is the next evolutionary step in condition based asset management, endeavoring to build automated prognostic and diagnostic systems to preserve and enhance the safety and readiness obtained from legacy Health and Usage Monitoring Systems. ISHM is to detect, diagnose, predict, and mitigate undesirable events caused by degradation, fatigue and faults in components over a certain period of time. For instance, the presence of such problems may occur during an important

function related to a system’s aircraft, regardless of whether the adverse event was caused by the subsystems. To properly address this problem, it is critical to develop technologies that can integrate large, heterogeneous distributed system [2], asynchronous data streams from multiple subsystems to detect a potential adverse event. The technologies would later be used to diagnose the cause of the event, foresee what consequences the event will have on the remaining useful life of the system (i.e., how it would jeopardize the entire system), and lastly take appropriate precautions to mitigate the event, if necessary [1].

Furthermore, effective estimation of the remaining useful life of devices and systems rely on development of prognostic models. This in turn requires extensive effort being made towards accelerating ageing mechanisms for each component, which ultimately enables us to prepare a sufficient amount of degradation profiles. This therefore makes it necessary to obtain the degradation profiles of every subsystem, including their individual components. This leads to a new degradation profile being devised every time a component is upgraded. The following degradation profile is calculated from either the accumulated damage or the data driven. Consequently, any changes made in the design of the system will both consume time and incur additional costs, considering that the prognostics model will need to be re-upgraded. It is thus apparent that the proposals discussed above are all expensive and time consuming processes that suffer from unreliability, noise, inaccuracies, etc [3].

To effectively overcome these problems, at the highest system level, a System- Level Reasoning (SLR) can be developed to at least provide the system with significant capabilities that can potentially decrease costs by assigning the system prognostics with a System Integrated Prognostic Reasoner (SIPR) [1][4]. A Vehicle Integrated Prognostic Reasoner (VIPR), for instance, is a NASA funded effort for developing the next generation VLRS. A typical functional

module within the SLR is a System Reference Model. This System Reference Model divides the system into partitions; and provides the necessary relationships between subsystems for the inference process. This partitioning enables the inference engine to reuse and link the same prognostic models to multiple subsystems and further minimize certification and qualification costs [1][4].

In summary, various techniques and methods, such as neural network, fuzzy, statistics, semantic computing, graph theory etc. have been utilized for the development of ISHM. However, ISHM still suffers from problems related to inefficient models, uncertainties and inadequate reasoning. In addition, the development of prognostic models still remains to be very costly and time consuming. These problems however still exist, mainly because the prognostics of a system heavily relies on the physics of failure models and degradation profiles that are known to be either inaccurate, inconsistent or very noisy. We believe that the ISHM system will greatly benefit if the prognostic of a component and a system is perceived as a feature rather than a system or component, which allows us to develop the prognostics based on this specific feature of the system instead of having to worry about the physics of the components. An advantage of this approach is that it will enable SLR to develop prognostics for a new subsystem based on a collection of features (encompassing various models/patterns) already known from the previous prognostics of subsystems. In order to fulfill this task, SLR may need to employ various techniques, such as those that involve Soft Computing (SC) including (fuzzy and neural network) in its Inference Engine and System Reference Model units, so that the subsystems properties can be linked to one another. In this proposal, we expect that there may be a duality connection found between the prognostics of dual systems, assuming that the prognostics of the dual systems are also seen as their parameters and features.

The next section shall describe in more detail the prognostics in systems. The principles of duality in electrical systems, along with brainstorming the duality concept of system's prognostics, are covered in Section 3. Section 4 covers the prognostics of Cuk converter and its dual circuit via developed algorithms and simulations with details of test approaches in Section 5. Lastly, the conclusion is covered in section 6.

II. PROGNOSTICS

In condition-based maintenance, prognostics can be defined as a controlled engineering discipline that focuses on the estimation and prediction of the future course of a system or component that attempts to work out at what point it starts to slowly develop irregularities and faults to the point where it eventually malfunctions. As a result of such malfunctions, a system or component can hence no longer meet the desired performance expectation. The predicted lifecycle of a system or component is referred to as the Remaining Useful Life (RUL). RUL is an important concept that is used in decision making for contingency mitigation and maintenance. The prognostics of a system or component are constructed from various scientific techniques including: failure mode

analyses, early detection of aging signs, and damage propagation models. Failure mechanisms are often used in conjunction with system lifecycle management to create prognostics and health management (PHM) disciplines. PHM is also sometimes referred to as system health management (SHM) or within the field of transportation applications; it is either referred to as vehicle health management (VHM) or engine health management (EHM). There are three main technical approaches related to building prognostic models which are broadly categorized into data-driven approaches, model-based approaches, and hybrid approaches [1][4][5].

A. Data-Driven Prognostics

Data-driven prognostics [6] are mainly based on pattern recognition and machine learning approaches in order to identify and detect changes and trends in system state phases. In regards to predicting trends in nonlinear systems, the classical data-driven methods include stochastic models, such as an autoregressive model, the bilinear model, the projection pursuit, etc. Soft computing techniques that involve using various types of neural networks (NNs) and neural fuzzy (NF) systems have also been commonly adopted to deal with data-driven forecasting of a system state [7]. The following prognostic approach concerns applications that have a complicated system; meaning that developing an accurate prognostic model of such a system will be expensive. So by using this particular approach to deal with complex systems will allow the prognostics of a system to be frequently set up much faster and cheaper as compared to other approaches. On the contrary, data driven approaches may have a wider confidence intervals than other approaches which mean it will require a substantial amount of data for training purposes [8].

Various strategies that are used to develop data-driven prognostics involve the analysis of either (1) modeling cumulative damage and then extrapolating out to a damage threshold, or (2) learning directly from the data relating to the remaining useful life.

Since individually failing systems is a lengthy and rather costly process, we thus seek to obtain a run-to-failure data which is the main fundamental setback, especially for new systems. In order to retrieve adequate data-driven prognostics, the accelerated aging data should be carefully extracted from a number of similar products by suitable measuring tools. This means that both quality and quantity aspects of the data driven prognostics will add to expenses; especially since the data sources may have been derived from a wide range of factors including temperature, pressure, oil debris, currents, voltages, power, vibration and acoustic signal, spectrometric data, as well as calibration and calorimetric data. It is therefore important to fully be aware of what parameters and signals are necessary to be measured, and which features must be extracted from noisy, high-dimensional data [6][7][8].

B. Model-Based Prognostics

The attempts made to incorporate a physical model of system which is (either accomplished via micro or macro

levels) into the estimated remaining useful life (RUL) is known as model-based prognostics [5]. The micro level (also known as material level) is often referred to as damage propagation model which is a physical model that is integrated with a series of dynamic equations. These dynamic equations define the very relationships between damage and degradation of a system or component. They further define how the system or component is operated under environmental and operational conditions. As it's almost impossible to measure many critical damage properties, an alternative solution is to use sensed system parameters instead. However, there may be a possibility that the level of uncertainty and inaccuracy are increased. In spite of the uncertainty and inaccuracy added as a result of sensed system parameters, uncertainty management must be considered with the proper assumptions and simplifications, which may overcome the significant limitations caused by that approach [4][5][9].

In contrast to physical expressions used in micro-levels, macro-level models alternatively use mathematical models at a system level in order to define the relationship among system input, system state, and system measure variables. The mathematical model is often a simplified representation of the system. Simplification may help make prototyping faster; but the trade-off to this is that the coverage of the model is increased at the expense of reducing accuracy of a particular degradation mode. In addition, within a complex application, such as a gas turbine engine, there would be a lack of knowledge in attempting to develop the proper mathematics for all subsystems or components. Again, this adds uncertainty and inaccuracy, similar to micro-level models; which means simplifications would need to be accounted for by performing uncertainty management procedures [1][4][9].

C. Hybrid Approaches

In reality, having a purely data-driven or purely model-based approach is almost impossible. However, both models do include some aspects of one another mechanisms. Hybrid approaches intend to bring the strength of both 'data-driven' approaches and 'model-based' approaches into one prognostics strategy. The two well known categories of Hybrid approaches are, 1) Pre-estimate fusion and 2.) Post-estimate fusion. The first technique applied, hardly has any 'ground truth' data or 'run-to-failure' data available. The second technique is more suitable in situations where uncertainty management is required. This means that the second technique helps to narrow the uncertainty intervals of data-driven or model-based approaches while also improving accuracy [10][11].

III. PROGNOSTICS OF DUAL SYSTEMS

Duality is one of the fundamental properties of systems, so that it can be consistently seen in systems that have any kind of physics [12][13]. It has a captivating history in mathematics, engineering and science. Duality relations have been established between geometric objects, algebraic structures, topological constructs and various other scientific constructs. In electrical systems, duality relations have

appeared in the core principles for any theorem in electrical circuit analysis in situations where there is a dual theorem that replaces one of the quantities with dual quantities; examples of dual quantities are current and voltage, impedance and admittance, meshes and nodes (shown in Table 1) [14].

TABLE I. DUALITY PRINCIPLES IN ELECTRICAL SYSTEMS

System	Dual of System
Voltage of nodes or across device	Current of branch or mesh
Current of branch or mesh	Voltage of nodes or across device
Resistor (R)	Conductivity (1/R)
Conductivity (1/R)	Resistor (R)
Capacitor (C)	Inductance (L)
Inductance (L)	Capacitor (C)
Voltage Source (Vs)	Current Source (Is)
Current Source (Is)	Voltage Source (Vs)
Kirchhoff's Current Law	Kirchhoff's Voltage Law
Kirchhoff's Voltage Law	Kirchhoff's Current Law
Mesh/Loop	Node
Node	Mesh/Loop

In regards to duality concepts, there will be a duality relationship between two electrical circuits if the parameters values and topologies of these two circuits are linked to one another based on details in Table 1. From a mathematical point of view, dual circuits have the same mathematical model except for having different parameters. Thereby we want to fully comprehend that if one was to consider that the prognostic of a system or component were to be seen as a parameter, it will thus mean that the prognostics of a system that have different topologies can be assigned to one another, while considering that the systems have the same mathematics model but with dual parameters as shown in Table 1. This provides us with the required facilities to develop the prognostics of a system based on the prognostics of its dual system.

From graph theory [12], it is well established that the behavior and function of a system can be recognized from knowing the topology of a system without the need of knowing the components and devices that are used in the system, considering that the nodes voltages and currents of branches in the circuit are known. Hence, it can be expected that graph theory provides us with the capability to construct the prognostic of a system based on its topology rather than concentrating on the devices and components that are integrated within the system. It is also expected that systems that have the same topology and mathematical models will also share the same prognostics no matter what components are included in the system. Therefore, it is possible to investigate how prognostic models can be designed from the topology of system rather than having to know physics of failure of a system. This makes the process of modeling the prognostics of a system much more ideal and realistic by saving a substantial amount of resources and time, since you wouldn't have to individually test each system to identify its prognostics.

Figure 1 shows an example of dual circuits. Using Kirchhoff's laws, it is evident that both circuits have the same form of mathematical model as shown in (1) for circuit in Figure 1-a; and (2) for circuit in Figure 1-b:

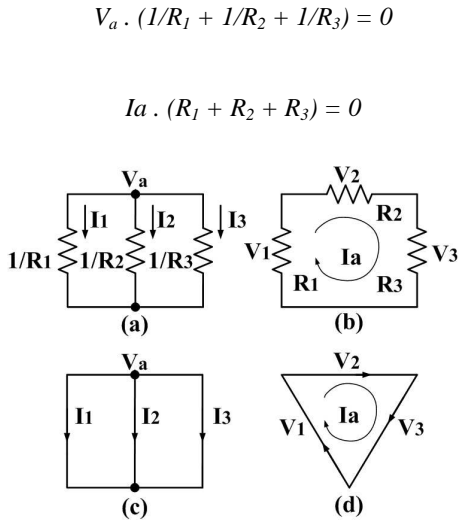


Figure 1. a) Cuk Converter, b) Dual circuit for Cuk converter in 2-a.

If for instance a degradation mechanisms is added, R_2 in circuit of Figure 1-b is aged as short circuit ($R_2 \rightarrow 0$), this is turned as ($1/R_2 \rightarrow \infty$) in circuit of Figure 1-a. This actually represents the duality principles shown in Table 1 in which the resistor is a dual of a conductive; or in regards to this example, it can be known as the short circuit being a dual of an open circuit.

The same rules can be used in more complex circuits where various components including capacitors and inductances are also used. The most critical point that needs to be worried about is the fact that degradation and failure mechanisms of dual components are not truly related to one another. Degradation mechanism of capacitor, for instance, is not related to degradation mechanisms of inductance, at all.

In order to deal with this problem, we rely on the well known physics principles, such as Ohm's and Kirchoff's laws. In reference to these two laws, it's obvious that any electric component can be formulated by using voltage across the component and current through the component. Alternatively, in regards to basic principles in graph theory of circuit and system design, it is well known that the behavior of a system is fully formulated if voltage of all nodes and current through all branches in the circuit are also known. This means that no matter what components are used in the circuit, as long as all the voltages and currents are known, the behavior and function of circuit can be fully formulated. Figure 1-c and 1-d, respectively show the graph of the equivalent circuits in Figure 1-a and 1-b.

From a circuit level point of view, the components details do not necessarily need to be known in order to develop a prognostics model for a circuit. Practically, sensors are used to measure voltages, currents, temperature etc. This allows the experiences of a degraded circuit or system of any form, to be interpreted as a circuit not functioning properly, on the basis of the sensed values meaning. Although this principle can be applied for greater purposes, i.e., to design a device independent prognostic model, this paper will mainly aim to

(1) present a realization of duality principles for the development of prognostics for dual circuits.

In addition, duality concept has already been recommended for diagnosing faults. Reference [15] proposes a fault diagnoser based on the duality principle and the optimal control theory for linear systems. However, this paper will present duality applications in system prognostics.

IV. PROGNOSTICS OF CUK CONVERTER

This section shows how duality concept can be used to develop prognostic models for Cuk converter [16] and its dual circuit. The following simulations were all conducted with Matlab and Orcad. Schematic of Cuk converter and its dual circuit are shown in Figure 2-a and 2-b. We use certain values for Cuk converter devices as well as all the equations depicted in reference [16] for all the simulations in this paper. Cuk is a step-down/step-up converter that shares a similar switching topology with boost-buck. Thus, it presents the voltage ratio of a buck-boost converter:

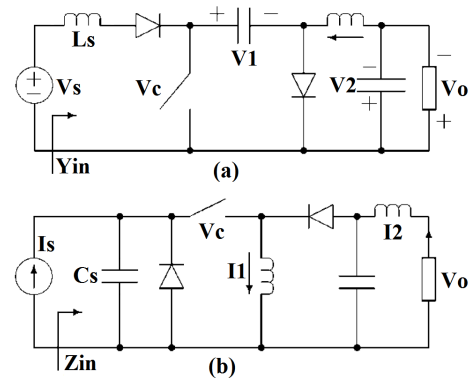


Figure 2. a, b) Resistive circuit with duality relationship, c,d) Graphs for circuits 1-a and 1-b.

$$v_o/v_s = D_s / (1 - D_s). \quad (3)$$

where v_o is output voltage, v_g is the input voltage, D_s is the duty cycle of the switch $t_{on}/(t_{on}+t_{off})$; and t_{on} and t_{off} are durations for when the switch is on and off. Equation (3) is calculated from the principle of conservative energy and the fact that the inductor currents relate to the input and output currents. This equation shows that the output voltage can be controlled by maintaining the duty cycle of the switch. Depending on the switching scheme, output voltage can be higher or lower than the input voltage. The state equations for Cuk converter are:

$$x' = Ax + B_v g + B_c d \quad (4)$$

$$v_o = C_x$$

$$x = [v_2 \ v_1 \ i_2 \ i_1]'$$

The Cuk converter has two inputs, a control input (V_c) and an input from the power supply (v_s) and one output (v_o). Therefore, matrix $[A \ B \ C \ D]$ relates to 'state space matrices' for the open-loop model from the v_s to the v_o . Similarly, $[A \ B_c \ C \ D]$ is the state space matrices from the control input d to the output v_o . Values for A, B, B_c , C, and D are given in [16]. The same equation can be extracted for dual circuit of Cuk converter in Figure 2-b; however, parameters are used in a dual form as shown in Table 1. Switches in Figure 2 are IGBT with a control voltage V_c . Y_{in} and Z_{in} are input admittance and input impedance of Cuk circuit and its dual circuit.

In converters, components that are mainly damaged are IGBTs and capacitors. IGBT experience numbers of failure mechanisms, such as bond wire fatigue, bond wire lift up, corrosion of the wires, static and dynamic latch up, loose gate control voltage, etc. The resulting affects mentioned are too complex, but we assume that these failure mechanisms can cause IGBT to behave as either an open circuit on a collector-emitter or a device encountering malfunction on its gate-emitter control. For instance, IGBTs thermal junction is increased due to solder crack which turns to wire bond lift off that increases the resistor relating to the collector-emitter. On the other hand, hot carrier injection is increased due to electrical stress. This causes short circuit on the IGBTs gate-emitter junction. As a result of this failure, IGBT's gate controllability is missed (loose gate control voltage) that causes IGBT to malfunction. The result of this effect is an increase in current through collector-emitter which means that the resistor of collector emitter is decreased. Therefore, it can be realized that wire bond lift off and loose gate control voltage are failure mechanisms that presents some kind of duality relationship. While one of them increases the resistor, the other one decreases the resistor. Generally, we assume that IGBT's failure and malfunction mechanisms are parameters with duality relationships.

Figure 3 shows IGBT run to failure data for four different IGBTs. This data is too noisy and needs to be filtered, but still there are a number of states that can be seen in the data. These states refer to cracks or wires that were lifted up due to degradation mechanisms. The resulting effects are changes in the IGBT's function; and changes in the channel resistor of that IGBT. We assume that degradation is processed in a form of duality for Cuk and its dual circuit, so that if IGBT of Cuk experiences degradation towards its open circuit, IGBT of dual circuit of Cuk is degraded towards short circuit. By the time that the IGBTs are damaged, C_s and L_s are fully charged as well as the other energy storage components lose energy, so V_o would be 0. It is however impossible to have a real short circuit in IGBT, thus we assume that it may have happened when the current through the collector-emitter exceeds over its limit just before the IGBT is burned out.

Based on the level of accuracy, there are number of models for a real capacitor and an inductance. To simplify simulation, we assume that the capacitor and the inductance can both be modeled like Figure 4 for the purposes of this paper. These models will present duality relationship between capacitance and inductance while also presenting

the energy lost by the resistors. R_1 typically has had very large values, while R_2 has a small value; but due to degradation, these resistors are changed towards either open or short circuits.

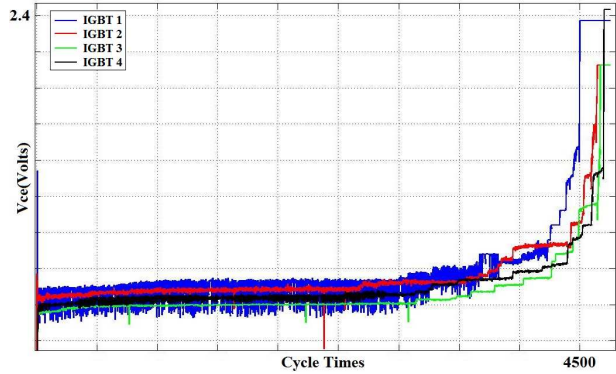


Figure 3. Run to failure data for four different IGBTs.

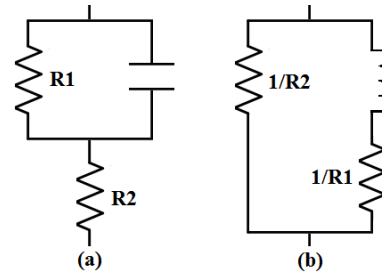


Figure 4. Real model for a) Capacitor, b) Inductance.

Figure 5 illustrates the proposed algorithm devised to develop this prognostics model. The same process that is possessed with different sets of run to failure degradation and malfunction profiles is repeated for both Cuk and its dual circuit. The components of the circuits are initially set to be in a good condition. Then as soon as the time step for the circuit is increased, the values of the components are changed by using a series of values provided in the degradation profile for the new time step. Signals, such as v_1 , v_2 , v_o , i_1 , i_2 , i_o , are measured at each time step phase. These signals are used for calculating systems properties, such as transfer functions, input and output impedances and admittances. Subsequently, the system degradation is turned according to changes encountered in the transfer functions ($Z_c(d,t)$, $Y_c(d,t)$, $Z_{dc}(d,t)$, $Y_{dc}(d,t)$). So where d is an index of a selected degradation profile, c is Cuk and dc is the dual circuit of the Cuk converter. Whenever d is altered, time step (t) is reset to zero which will reset the process of the circuit to a healthy condition for the new degradation scheme. By measuring the mentioned signals and parameters, it would be possible to realize how energy is transferred between capacitances and inductances; and how that transferred energy is lost when the system is also degraded.

We realized that if a degradation profile is used for Cuk, such that it's converted to a malfunction profile for its dual circuit so that the IGBTs in both circuits are always realized in dual forms; then a duality relationship would be seen

between the transferred functions of these two circuits. For instance, $Z_c(t)$ is equal to $Y_{dc}(t)$. This is because as the degradation profile changes the IGBT of Cuk towards an open circuit; its malfunction profile also changes the IGBT of dual circuit towards a short circuit.

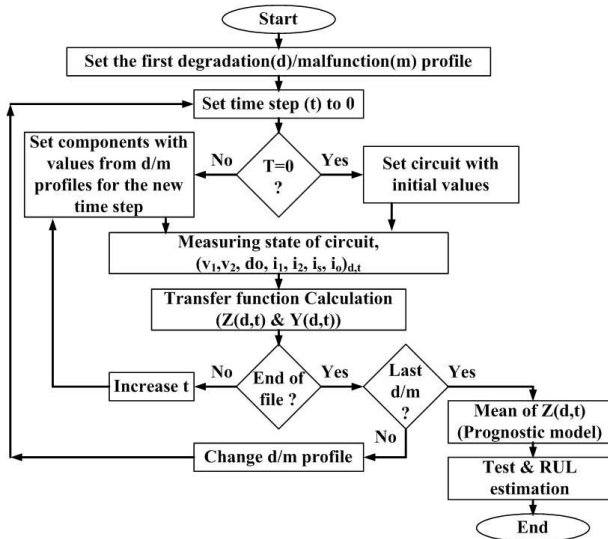


Figure 5. Algorithm used to develop prognostic model.

If the malfunction profile for dual circuit of Cuk is not extracted from the degradation profile of a Cuk circuit, then $Z_c(t)$ is not identical to $Y_{dc}(t)$. However, we come to a conclusion that if the whole process is repeated for number of different degradation and malfunction profiles and that the mean value of $Z_c(t)$ and $Y_{dc}(t)$ are used for comparison; leads to meaningful similarity patterns to be found between $Z_c(t)$ and $Y_{dc}(t)$. $Z_{mc}(t)$ can be used for the mean value of $Z_c(d,t)$ and $Y_{mdc}(t)$ can be used for the mean value of $Y_{dc}(d,t)$, in situations where m refers to the mean value. $Z_{mc}(t)$ and $Y_{mdc}(t)$ can be both used as prognostic models for Cuk and its dual circuit. However, these two transfer function are not exactly identical, but they would be more similar to one another if the process that is required to be executed to obtain the functions is repeated for various numbers of degradation and malfunction profiles for both circuits. By implementing more intelligent algorithms that use stochastic, neural network, fuzzy and other techniques instead of a simple mean value function will increase the accuracy of this prognostic model. Implementing such intelligent algorithms also reflects the future aim and direction of our research. Additionally, we should be aware that prognostics have always been a way to estimate the life time of devices and systems within different confidence levels. Confidence levels provide assurance, so that we can comfortably rely on the performance of an aged system. The point is the accuracy of prognostic models has always been under doubt and remains to be under margins of confidence levels. So in summary, by using the prognostic model of a system for other systems where similarities in their properties (like duality) are found, would give us a more accurate and reliable representation of the state and condition of the system or component. This is

assuming that the prognostics are developed from adequate number of degradation profiles, and that they also have the right minimum and maximum confidence levels.

V. TEST APPROACH

The resulting prognostic model is tested with an additional degradation profile which is used as a test data to estimate the remaining useful life time for the converter. During the testing process, the prognostic model is stimulated via the samples derived from the test data. This causes the parameters of the prognostics model to change, which therefore leads to the degradation of the system. The accuracy of the degradation depends on the number of delayed and differentiated samples that are used to simulate the prognostic model as well its time step t sample.

The tests would be inaccurate, if the model was stimulated one sample at a time, despite there being durations in the test data where the samples remain almost the same. A more accurate testing is achieved, if differentiated samples are also used for stimulation. This thereby allows the prognostics model to follow the test data trend rather than only following one sample at a time.

Therefore to estimate the life time of system at each time interval, the model is stimulated with the sample at t and a set of differences. Once the life time of the system is estimated for that specific sample, it then selects the next sample from the test data provided for simulation, while also updating the differences. In addition to the system degrading at each time step, the next sample test (let's call it S^+) is also calculated from using the model's system. The simulation is then continued by stimulating the prognostic model using a calculated sample (S^+) which in turn degrades the model again and updates the calculated sample (S^+) with a new value. The same process is continued until S^+ reaches a threshold which refers to a specific class in the test data where the device is no longer in good condition for operation. We set the threshold to 7 based on the degeneration profile that we had available, Figure 3. Once S^+ reaches the threshold, the simulation continues with the next sample provided by the test data. This also requires the differences to be updated, accordingly. The life time for each sample of test data depends on the time that it takes for the model to reach the threshold from the time a new sample of test data has been selected for stimulation to the time that the calculated sample test reaches the threshold. This means that the stimulation for each time step starts with a new sample obtained from the test data. As this process is repeated with the sample calculated, the model also eventually degrades. Figure 6 shows a real and estimated RUL with % 10 and % 90 confidence levels.

VI. CONCLUSION

In conclusion, this paper shows that the prognostics of a system can be applied to other systems that share similar properties in the form of duality. A prognostic model is developed in the form of a time dependant transfer function where values are altered over a certain period of time based on the degradation mechanisms of a system's components. By having the prognostics assigned to a system's property

reflects the duality connection of degradation and malfunction of system. This means that if the components of a system are aged, their dual components in the dual circuit will be faced with malfunction. The accuracy of the developed prognostic model is dependent on the number of available degradation profiles; and the method that is used to train the time dependant transfer function. The accuracy of this model is guaranteed and expressed within the minimum and maximum confidence levels. However, we presented our approach just for Cuk converter and its dual circuit, but it seems that the same technique can be used for systems that have slightly similar topologies, degradation mechanisms, and properties. Thereby, further research needs to be conducted for systems that are not in dual forms, especially for the purposes of exploring how the prognostic model of a system could be mapped to the prognostic model of another system.

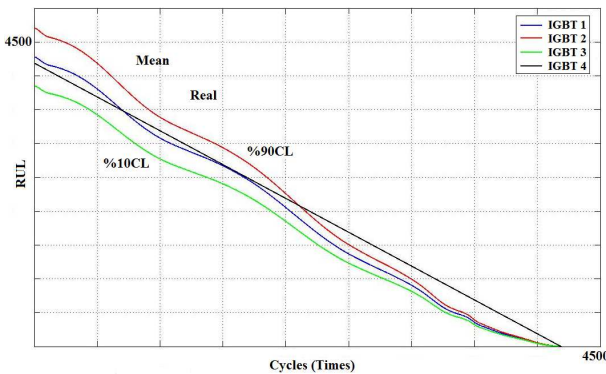


Figure 6. Resulting RUL after testing prognostic model with data test.

The advantage and usage of such a technique is emphasized in the implementation stage of the inference engine for System- Level Reasoning (SLR) and System Integrated Prognostic Reasoner (SIPR). In addition, it provides us with the facility to transfer degradation knowledge and experiences between systems. This means that the development of prognostics for huge systems, such as heterogeneous distributed systems used in applications like aircraft is much faster, while the cost assigned to accelerated aging tests and preparing degradation profile is decreased. We essentially intend on pushing forward with our research, in order to apply this technique to the development of the prognostic inference engine and reasoned for aircraft.

ACKNOWLEDGMENT

The authors would like to sincerely thank Professor C Mark Johnson and Dr Paul Evans from the Power Electronics, Machines and Control Group, University of Nottingham for the contribution of failure data of the IGBTs and the power cycling test rig configuration.

REFERENCES

- [1] I. K. Jennions, "Integrated Vehicle Health Management Perspectives on an emerging field", SAE International, Warrendale, pennsylvania, USA 2011, pp. 100-110.
- [2] A. El-Sayed and M. El-Helw, "Distributed Component-Based Framework for Unmanned", Proceeding of the IEEE International Conference on Information and Automation Shenyang, China, June 2012, pp. 45-50.
- [3] W. Wenbin and M. Carr, "A Stochastic Filtering Based Data Driven Approach for Residual Life prediction and Condition Based Maintenance Decision Making Support" Prognostics & Systems Health Management, IEEE Conference, Macao, China, Jan. 2010, pp. 1-10.
- [4] I. K. Jennions, "Integrated Vehicle Health Management The Technology", SAE International, Warrendale, pennsylvania, USA, 2013, pp. 139-154.
- [5] M. Daigle and K. Goebel, "Model-Based Prognostics under Limited Sensing", IEEE Aerospace Conference, Big Sky Resort, USA, March 2012, pp. 1-12.
- [6] C. Chen and M. Pecht, "Prognostics of Lithium-Ion Batteries Using ModelBased and Data-Driven Methods", 2012 Prognostics & System Health Management Conference, IEEE PHM Conference, Beijing, China, May 2012, pp. 1-6.
- [7] H. Chao, D. Byeng, K. Youn, and K. Taejin, "Semi-Supervised Learning with Co-Training for Data-Driven Prognostics", Conference on Prognostics and Health Management, IEEE PHM Conference, 2012, pp. 1-10.
- [8] S. Sarkar, X. Jin, and A. Ray, "Data-Driven Fault Detection in Aircraft Engines With Noisy Sensor Measurements", Journal of Engineering for Gas Turbines and Power, Vol. 133, ASME, August 2011, pp. 1-10.
- [9] L. Jianhui, M. Madhavi, K. Pattipati, Q. Liu, M. Kawamooto, and S. Chigusa, "Model-based Prognostic Techniques Applied to Suspension System", IEEE Transactions on Systems, Man and Cybernetics, Part A: Systems and Humans, IEEE 2008, Vol. 38, Issue 5, pp. 1156-1168.
- [10] P. Shetty, D. Mylaraswamy, and T. Ekambaram, "A hybrid prognostic model formulation system identification and health estimation of auxiliary power units", Aerospace Conference, IEEE, Big Sky, MT, USA, March 2006, pp. 1-10.
- [11] A. K. Garga, K. T. McClintic, R. L. Campbell, Y. Chih-Chung, M. S. Lebold, T. A. Hay, and et al., "Hybrid reasoning for prognostic learning in CBM systems", Digital Object Identifier: IEEE Proceedings on Aerospace Conference, Big Sky, Montana, USA, vol.6, March 2001, pp 2957-2969.
- [12] L. O. Chua , C. A. Desoer, and E. S. Kuh, "Linear and Nonlinear Circuits", Mcgraw-Hill, 1st Ed. Edition, March 1987.
- [13] D. Y. Gao, "Duality Principles in Nonconvex systems: Theory, Methods and Applications", Kluwer Academic Publishers, USA, ISBN 9781441948250, 2000.
- [14] G. E. Sharpe and, G. P. H. Styan, "Circuit Duality and the General Network Inverse", IEEE Transactions on Circuit Theory, Vol 12, IEEE 1965, pp. 22-27.
- [15] J. Li, "Optimal Fault Diagnosis Based on Duality Principle for Linear Systems", Control and Decision IEEE Conference, July 2008, pp. 573-577.
- [16] F. J. Ryttonen and R. Tymerski, "Modern Control Regulator Design for DC-DC Converters", Electrical and Computer Engineering Department Portland State University. [online]. Available from http://web.cecs.pdx.edu/~tymerski/ece451/Cuk_Control.pdf, 2014.05.14.

Design and Simulation of Electronic Service Business Process

Peteris Stipravietis, Edzus Zeiris, Maris Zieme

ZZ Dats, SIA

41/43 Elizabetes street, Riga, Latvia, LV-1010

[peteris.stipravietis,edzus.zeiris,maris.ziema]@zsdats.lv

Abstract — The paper discusses the identification of common business process design-time problems using Yet Another Workflow Language (YAWL). The approach proposed by the authors is based on the creation of business process in the YAWL environment in order to simulate the process model which could resolve some of the design-time problems, i.e., possible bottlenecks, as well as provide hints on how to correct initial process. The simulation is done using process mining software “ProM Framework” and the Colored Petri nets simulation and analysis framework “CPN Tools”. The process design with YAWL is done with respect to Business Process Execution Language (BPEL) requirements, thus later allowing the transformation from YAWL to BPEL via the intermediate structure. The examples show that it is possible to identify some of the possible faults of the process using the proposed approach.

Keywords – *Electronic service; YAWL; simulation.*

I. INTRODUCTION

E-services are common in information society nowadays, and even though they tend to become more and more accessible and varied, the problems that occur during the design phase of the service remain the same. These problems include, for example, questions on how to facilitate the creation of business process to the user with no specific programming skills, how to define the process in a way that creates the process description abstract but accurate enough at the same time, how to check the created model – to determine the weaknesses, perform the measurements based tuning, and others.

The validation of the process is even more important when the process being changed is already deployed and used in production environment – one must make sure that the changes are implemented correctly, that the new instances of the process can run together with old instances already running. The implementation and validation of changed process also have to be simple and cost-effective enough – hence the conclusion that the solutions of these problems rely heavily on the choice of the language used to describe the process – does it provide the possibilities to validate and simulate the process.

Existing business process modeling languages can be divided in two groups. The languages of the first group are favored by the academic community, but rarely used in real-life solutions. These languages are based on Petri nets, process algebra; they have formal semantics, which allow the validation of the models described by these languages. The languages of the second group are used in real-life projects

much more than in academic researches. Business Process Execution Language (BPEL) [1] and Web Services Flow Language (WSFL) [2] are among these languages. These, so called business languages, often lack proper semantics, which could lead to debate on how to interpret the business models described by these languages. The availability of different implementations of these languages from different vendors does not facilitate the situation either, yet they are used much more, compared to rarely used models described by academic languages. If a situation arises when business process model described by business language needs to be validated using Petri nets, one must either abandon the validation or transform the process model to another model, described in academic language, for example Yet Another Workflow Language (YAWL) [3][4][5]. The authors propose reverse approach – first, a process is created using academic language. The design problems of the process model can then be solved by mathematical means. Second, the verified and updated model is transformed to model described in business language. The advantages of the approach described follows:

- If a model is created using academic language, it is more readable and maintainable than the model, which is a transformation result itself. It is also easier to perform analysis of untransformed model, because the transformation could lose some design information.
- Model, transformed to business language, is already validated and ready to be executed. Of course, the model must be double-checked to make sure if it needs any corrections. The alternatives of the execution environment for the model are much more than the environments for academic languages; in addition to that, they have superior technical support.

The purpose of this paper is to examine the design and simulation stages of the aforementioned approach – can it be used during the design of simple e-service business project; and to check if it helps to identify and resolve most common design-time problems.

The rest of the paper is structured as follows – Section II provides an overview of the proposed design approach. Next, Section III defines the design-time restrictions of the YAWL workflow which must be observed. The simulation phase of the approach is discussed in Section IV, while Section V shows the practical example of the simulation. Finally, Section VI contains conclusion.

II. PROPOSED APPROACH

The approach proposed by authors consists of five consecutive steps – the design, the simulation, the transformation to ‘protostructure’ – simplified, yet fully descriptive notation of the business process control flow –, the optimization based on quality attributes of the service and the transformation to BPEL. Some of these steps may be omitted or repeated as necessary, as shown in Fig. 1 – solid lines show the most common path, dashed lines show alternative paths of execution, lighter boxes represent steps which may be omitted.

A similar approach is proposed by authors of [6], but their solution is based on straightforward conversion of YAWL workflow to BPEL process (straight to step 5 from step 1). While the straightforward transformation is more efficient in terms of development cost and time, the authors’ approach includes simulation and optimization steps which should reduce the costs of maintenance later on.

The first step is the design of the business process using academic language. The initial business process model is created during this phase. Authors use the YAWL as the language of choice – it is based on Extended Workflow Nets (EWF), the workflows described in YAWL can be transformed to colored Petri nets to perform simulations and formal semantic validation. It also supports all workflow patterns [7] – parallel flow, branching, synchronizations and others. Although YAWL supports all the patterns, the business process model must take into account that the process will later be transformed to BPEL – as such, it may not contain patterns which do not translate to BPEL directly or using non-solution specific workarounds.

The second step of the approach is the simulation of the business process model. The analysis of the created business process is very important part of the design – one needs to find bottlenecks, when instance of the process or its part could use up all available resources, thus forcing other instances to wait for these resources; identify dead ends, which could lead to infinite loops and never ending process instances; find deadlocks, when process querying for the same resources effectively block each other; define fault handling and cancellation activities, which cancel all the work done by previous activities.

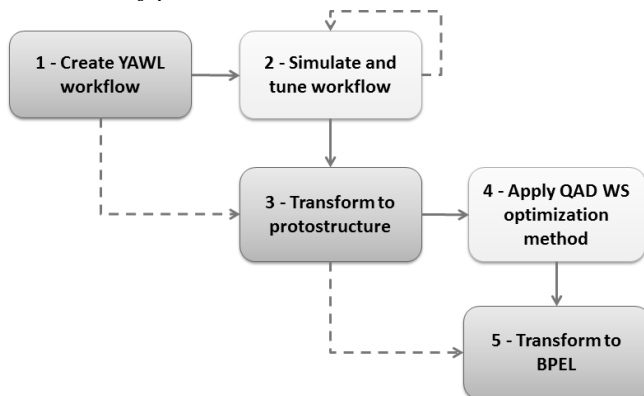


Figure 1. Proposed approach, step by step.

During this step it is also possible to identify reusable structures, for example, audit log activities. Such challenges usually are solved with the simulation.

The third step provides the transformation to primitive structure. Primitive structure is simplified definition of business process control flow, although it can also be used to maintain the data flow. The primitive structure serves as an intermediate between academic and business languages and can be used to create processes described by multiple languages, not only BPEL. The primitive structure may be changed and improved during this phase to facilitate the transition to target language, i.e., restructure its control flow in a way that it becomes well-formed and contains only patterns supported by BPEL.

The fourth step provides the segmentation of primitive structure using the Quality Attributes Driven Web Services Design (QAD WS) method which offers the segmentation of business process, represented as oriented graph. The segmentation result depends on process quality attributes selected by designer and their respective values. The result of this method is Pareto optimality set – the method returns the most suitable segmentations from all possible considering the quality attributes given.

The business language selected by authors and used in their proposed approach is BPEL, and using of QAD WS method on primitive structure would provide the possible structures of BPEL process – which parts of the process would belong to orchestration and which ones would be implemented as web service calls. The work in [8] also discusses the partitioning of Web services into orchestrations based on their Quality of Service (QoS) values, although that approach do not use multicriterial optimization, but rather is based on Petri nets and usage of statistical data.

The QAD WS method perceives the business process as an oriented graph G , whose vertices corresponds to process activities, but edges between them represents the control flow. Using various quality attributes and the structure of graph G , QAD WS method solves multi-criteria optimization task, which results in the segmentation set of initial graph G : $G' = QAD(G)$. Criteria used by the method are:

- Costs of development C ;
- Performance T ;
- Maintenance costs E ;
- Reusability R ;
- Integrity I [9].

The segmentation set G' consists of N most optimal solutions designer can choose from – this method may greatly reduce possible solutions of process architecture, thus aiding the designer. For instance, if the main criterion for segmentation is performance, then the segmented graph should contain the Web service invokes as few as possible, because every invoke adds to total execution time. On the other hand, the reusability will lead to much more segmented graph to allow the components to be reused. Preliminary tests show that if a certain graph G consists of 11 vertices, then, taking into account the segmentation restrictions (mainly to preserve control-flow), one ends up with 1015 possible solutions. Applying the QAD WS method reduces

the count of possible solutions to 8 solutions. The total solution count varies depending on the structure of the graph, and resulting set solution count varies depending on the used criteria and their respective weights.

The result of the 3rd step of authors proposed approach is primitive structure – oriented graph P that corresponds to initial YAWL workflow, which could be used as an input graph G for QAD WS. Authors also note that P is more complicated and restrictive as G – in addition to process activities and links between them it also contains the process control flow.

The last phase of proposed approach is the transformation of primitive structure to business language process which results in the business process defined in business language. This process is not ready to be executed, but its structure corresponds to initial process model described by academic language and maintains its process flow.

III. THE DESIGN OF THE BUSINESS PROCESS

To be able to transform the YAWL workflow successfully, the workflow must conform to some requirements. Firstly, it should not contain patterns, which have no analog constructions in BPEL, for example, the passing of process control to an activity residing outside the synchronized block, i.e. – goto-like construction. BPEL directly supports 13 patterns out of 20 [10], discussed in [11]. The parallel flow with runtime-only knowledge is also not supported.

Secondly, the incoming and outgoing messages are associated with specific process instance using correlation sets. YAWL lacks concept of correlation sets, because each workflow instance (case) is started by its clients (users), thus creating an instance in execution environment. This environment manages the workflows and offers to users corresponding options, based on the state of instance and its specification [4]. The variable which could be used as a correlation set variable must be created in the workflow or during the transformation and finally added to each defined data type used in BPEL messages.

Thirdly, support of human tasks – all BPEL activities related to exchange of information with process partners are perceived as web service operations, i.e., BPEL has no concept of “Human interaction”. To fill this gap several BPEL extensions are proposed, for example, BPEL4People [12][13] – Organization for the Advancement of Structured Information Standards (OASIS) is working on standardizing this extension.

Last but not least, workflow definition must correctly define all the branching conditions to avoid cases when transformed BPEL process’ While, Repeat/Until and If blocks contain incorrect values.

IV. THE SIMULATION

As already mentioned before, the developers and architects may encounter the same kind of problems during the development of the business process of electronic service – bottlenecks, dead ends, and fault handling and cancellation

activities. Part of these problems may be identified using simulation.

The authors of the paper [14] propose simulation which uses process design data, historical data about executed process instances from audit logs and state data of the running process instances from the execution environment. Data from all three sources are combined to create simulation model – design data are used to define the structure of the simulation model, historical data define simulation parameters, state data are used to initialize the simulation model.

Altering the simulation model allows to simulate different situations, for example, to omit certain activities or divert the process flow to other execution channels. Taking into account the state data of running process instances, it is possible to render the state of the system in near future and use the information to make decisions regarding the underlying business process.

The simulation of the workflow is carried out using process data mining framework ProM [15]. To create simulation model, following steps are performed:

- Workflow design, organizational and audit log data are imported from execution environment;
- According to imported data a new YAWL workflow model is created and state data are added;
- The new model is converted to Petri net;
- Resulting Petri net is exported to simulation execution environment CPN Tools [16] as a colored Petri net.

CPN Tools environment provides the process simulation possibilities both in long and short-term, using the state data of chosen process instance. This technique differs from others with its degree of realism. For instance, the work [18] shows the so called mediator approach based on Discrete Event System Specification (DEVS) models, while the work [19] exploits event-based approach on Service-oriented Architecture (SOA) Both of these approaches, however, do not use the statistical data of already finished instances to take into account the availability of the resources; yet this method creates artificial delays based on historical data from audit logs and organizational model.

Returning to design-time problems, mentioned before – the authors now will evaluate the simulation approach, focusing on its capabilities to identify these problems.

A common mistake is to propose that the user of process will provide all data when prompted, for example, fill out all fields in a web form. Some fields could be left blank because user is not interested in sharing particular information (because of privacy concerns, etc.) or other reasons. If a process needs such information to continue, but user does not provide it, it enters in a waiting state and theoretically may never be finished. As an example, a simple YAWL workflow is provided, which expects the input of e-mail address and phone number from the user. The workflow contains AND-flow – it is finished when both branches are finished – see Fig. 2.

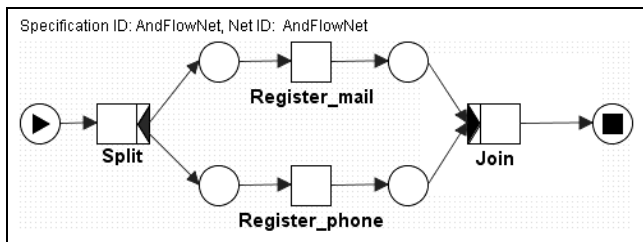


Figure 2. Simple YAWL workflow – parallel flow

The solution is to use OR-flow instead of AND, but the simulation approach discussed is not applicable – Petri nets does not support OR-flows; also the ProM Framework generates syntactically incorrect Standard ML (SML) file, if the net output variables have not been assigned value during the flow – however, it is not possible to assign initial values to them during the design. However, the problem can be identified by analyzing the process execution log, using, for example, ‘Basic Log Analysis’ module from ProM toolset. Fig. 3 clearly shows the difference between activities executed. Naturally, the question arises – why are there so many processes in waiting state and why do the users register their phone numbers far less than their e-mail addresses.

The simulation will not be helpful to identify the problem even after the editing of the SML file – the generated Petri net will contain AND-flow and the short-term simulation will direct the token through both branches. It is not reasonable to fix the net – the problem would be already identified and there would not be need to carry out the simulation once more. As mentioned before, such problem could be fixed introducing the OR-flow, however, both Petri net and BPEL lack the concept of OR-flow – it could be replaced by subsequent XOR-flows (IF-THEN branching).

Another task is to identify possible situations which could lead to infinite looping. The reason behind the looping mostly is incorrectly defined loop exit condition or loop variable does not have correct value assigned. This case can also be identified using ProM tool ‘Basic Log Analysis’ (activities within loop would be far more than others), but the simulation technique discussed can be used too.

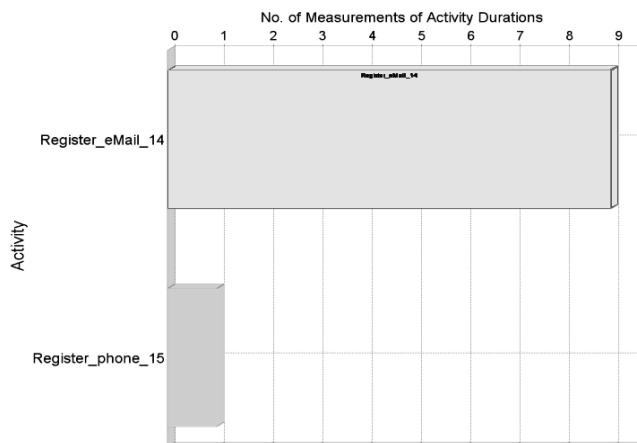


Figure 3. Difference of activities

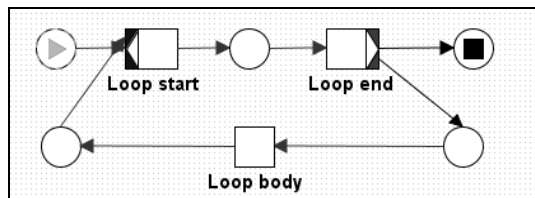


Figure 4. Simple YAWL workflow – possible infinite loop

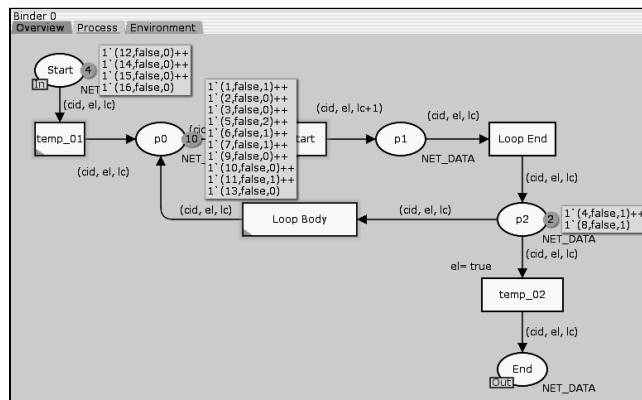


Figure 5. Colored Petri net – infinite loop

Fig. 4 shows the example of loop. The workflow has one local variable – “exitLoop” of Boolean data type, but its value intentionally is not being changed. Fig. 5 shows resulting Petri net – variables “loopCount” and “caseId” were added for demonstrative purposes (token colset NET_DATA is defined as product INT*BOOL*INT). Examining the net simulation, one observes that incoming tokens never leave loop.

The simulation performed to identify the bottlenecks produces similar results. YAWL has mechanism to distribute activities to resources defined in its organizational model – these resources are tokens of different colorset in a colored Petri net. A transition in Petri net may fire if all the places leading to transition have tokens – if a ‘resource’ place has no tokens, transition never fires, and all ‘data’ tokens accumulate in ‘data’ place until resources are freed and transition may fire again.

V. EXAMPLE

The example for simulation is quite simple process – the user is prompted to enter both phone number and email. The request is processed and notification to user is sent about the availability of result. If the user does not retrieve the result within some amount of time, the notification is sent again. “Wait/Check” activity sleeps for some predefined time and then checks the value of the element “exitLoop” of each token to determine if it is possible to exit the loop and end the process instance or does the instance send the reminder once more. The example is not too complex because it contains the possible simulation problems discussed previously and there is no need to make it more complicated.

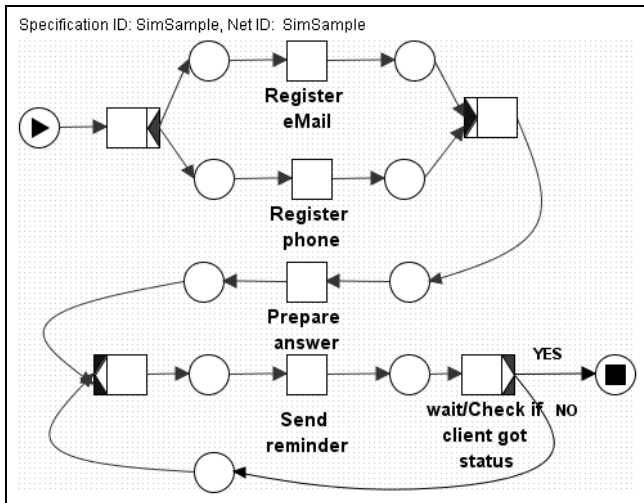


Figure 6. YAWL workflow - the simulation example

The example consists of three blocks – AND-flow, proposed ‘bottleneck’ and possible infinite loop – Fig. 6. The user of e-service is prompted to enter his/hers email address and phone number – AND-flow. Then some worker/service checks the data and prepares answer – possible ‘bottleneck’. Finally the reminder to user is sent until he/she turns up for the answer – possible infinite loop. The activity ‘Prepare answer’ is assigned to YAWL resource ‘User1’.

There were 12 process instances created – in 6 of them both e-mail address and phone number were provided, in other 6 just the e-mail address. As mentioned before, CPN Tools cannot simulate OR-flows – the only way to diagnose waiting processes is using ‘Basic Log Analysis’ from ProM toolset. Fig. 7 illustrates the measurements of the execution count of each activity – the activity ‘Register_eMail_14’ has been executed in all 12 instances, while the activity ‘Register_phone_15’ has been executed only 6 times.

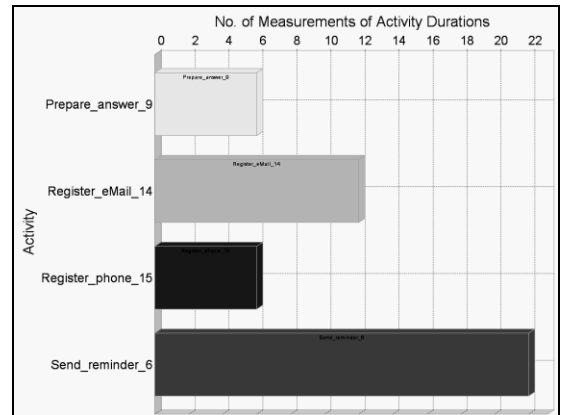


Figure 7. Basic Log Analysis – activity execution total count

After the tweaking the workflow – changing the OR-flow to AND-flow, all instances of the process could complete both activities. Fig. 8 shows the initial Petri net which corresponds to YAWL workflow in example. It has one token in place ‘Resources’, which means that all the concurrent instances will be processed in order by the same user, as seen in Fig. 9. The availability of only one resource token forces a wait in other processes, illustrated by 5 tokens in place ‘p4’, waiting for the resource ‘USER 1’ to become available again. Transition ‘Timeout’ was created to add artificial delay, which simulates the processing of the application.

Of course, it could be implemented using timed data types for all colsets used in net, but for the demonstrative purposes timed net is not used.

Fig. 9 also shows 6 instances moving through loop block – possible candidate for an infinite loop. The loop exit variable ‘exitLoop’ in each of the tokens holds value ‘FALSE’, so until the transition ‘t3’ can accept the token, it is stuck in the loop.

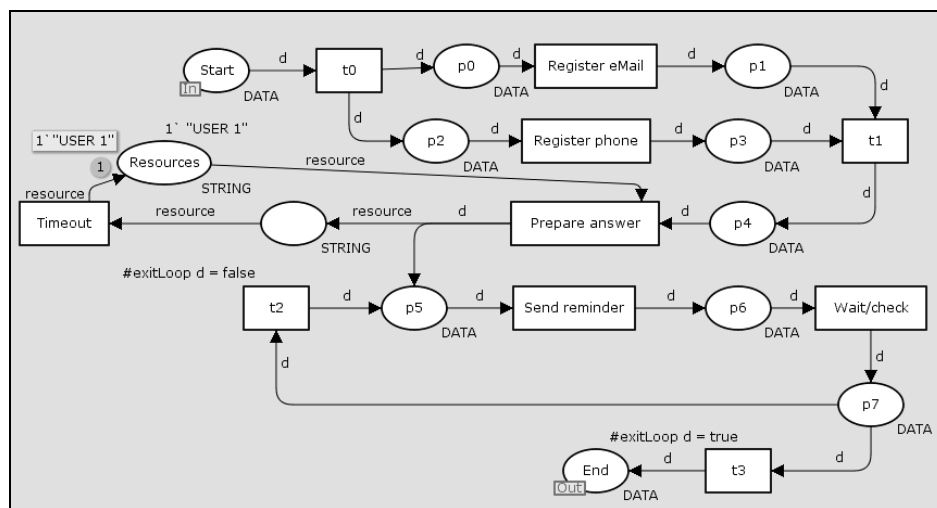


Figure 8. The Petri net corresponding to the example

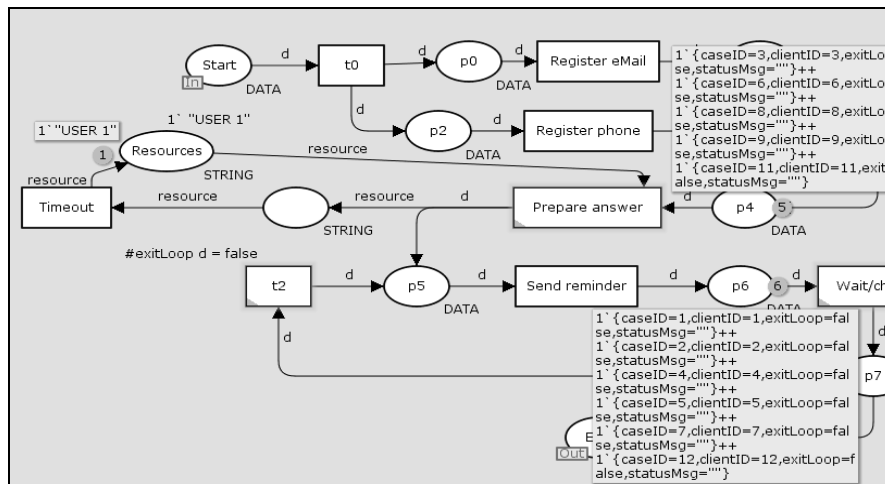


Figure 9. The bottleneck and the infinite loop in the Petri net

Examining the net in Fig. 9, it is clear that the process lacks an activity where the applicant can receive the results, thereby changing the value of variable “exitLoop” in the corresponding token.

VI. CONCLUSION

The proposed approach of business process modeling, when initial business process model is created using academic language and then transformed to the process described in business language, is quite successful. The main benefit of this approach is creation of primitive structure, which would later allow the transformation from YAWL workflow to any other hierarchical language (not only BPEL), both academic and business.

After examining and simulating simple workflow containing three possible problems, it is clear that the approach can identify simpler design problems; for example, the bottleneck and dead end identification can be achieved with this approach. However, more complex problems, such as deadlock identification or the operation of cancellation region could not be resolved. The problem lies with Petri nets, used in the simulation model, because they lack support of multiple simultaneous process instances or cancellation regions. One possible solution that could detect the deadlocks in the process model would be to use XML Process Definition Language (XPDL). Unfortunately, the XPDL supports only 9 out of 20 workflow patterns, while YAWL supports 19 [3][17].

On the other hand, the deadlock identification may be not as important since the deadlock situations would not arise in pure BPEL orchestration. The restrictions imposed upon YAWL workflow would also prohibit the use of OR-flow – another pattern which cannot be simulated. Taking into account the restrictions, authors conclude that the proposed simulation approach may be successfully applied during the development of electronic services.

ACKNOWLEDGMENT

The research leading to these results has received funding from the research project “Information and Communication Technology Competence Center” of EU Structural funds, contract nr. L-KC-11-0003, signed between ICT Competence Centre and Investment and Development Agency of Latvia, Research No. 1.21 ”Research on effective transformation of business processes to the architecture conforming to cloud computing”.

REFERENCES

- [1] Web Services Business Process Execution Language Version 2.0, Online. Available: <http://docs.oasis-open.org/wsbpel/2.0/OS/wsbpel-v2.0-OS.html>, accessed on the 19th of February, 2014
- [2] Web Services Flow Language Version 1.0, Online. Available: <http://cin.ufpe.br/~redis/intranet/bibliography/standards/leymann-wsfl01.pdf>, accessed on the 19th of February, 2014
- [3] W. M. P. van der Aalst and A. H. M. ter Hofstede, „YAWL: Yet Another Workflow Language” Information Systems, vol. 30(4), 2005, pp. 245–275.
- [4] W. M. P. van der Aalst, L. Aldred, M. Dumas, and T. A. H. M. Hofstede, „Design and implementation of the YAWL system”, Proc. of the 16th International Conference on Advanced Information Systems Engineering (CAiSE 04), LNCS, vol. 3084, 2004, pp. 142–159.
- [5] M. Weske, "Business process management", Springer, 2007, p. 169.
- [6] S. Pomudomthap and W. Vatanawood, “Transforming YAWL workflow to BPEL skeleton”, Proc. of the IEEE Software Engineering and Service Science (ICSESS 11), Beijing, China, July 2011, pp. 434-437.
- [7] N. Russell, A. H. M. ter Hofstede, W. M. P. van der Aalst, and N. Mulyar, „Workflow control-flow patterns: A revised view”, BPM Center Report BPM-06-22, BPMcenter.org, 2006.
- [8] S. Rosario, A. Benveniste, S. Haar and C. Jard, “Foundations for web services orchestrations: functional and QoS aspects,

- jointly”, Proc. of the 2nd International Symposium, Leveraging Applications of Formal Methods, Verification and Validation (ISoLA 06), 2006, pp. 309-316
- [9] E. Zeiris and M. Ziema, ”E-Service architecture selection based on multi-criteria optimization”, Proc. of the 8th International Conference PROFES 2007, Riga, Latvia, July 2007, pp. 345-357.
- [10] M. Havey, „Essential business process modeling”, O’Reilly 2005, p. 141.
- [11] P. Wohed, W. M. P. van der Aalst, M. Dumas, and A. H. M. ter Hofstede, “Pattern based analysis of BPEL4WS”, FIT Technical Report, FIT-TR-2002-04, Queensland University of Technology, Brisbane, 2002.
- [12] OASIS WS-BPEL Extension for People, Online. Available: https://www.oasis-open.org/committees/tc_home.php?wg_abbrev=bpel4people, accessed on the 18th of February, 2014
- [13] T. Holmes, M. Vasko, and S. Dustdar, “VieBOP: Extending BPEL engines with BPEL4People”, Proc. of the 16th Euromicro International Conference on Parallel, Distributed and network-based Processing, February 2008, pp. 547-555.
- [14] A. Rozinat, M. T. Wynn, W. M. P. van der Aalst, A. H. M. ter Hofstede, and C. J. Fidge, “Workflow simulation for operational decision support using design, historic and state information”, Proc. of the 6th International Conference on Business Process Management (BPM 08), Milan, Italy, Springer, 2008, LNCS, vol. 5240, pp. 196 – 211.
- [15] W.M.P. van der Aalst et. al., “ProM 4.0: Comprehensive support for real process analysis” In J. Kleijn and A. Yakovlev, editors, Application and Theory of Petri Nets and Other Models of Concurrency (ICATPN 2007), Springer, 2007, LNCS, vol. 4546, pp. 484-494.
- [16] K. Jensen, L. M. Kristensen, and L. Wells, “Coloured Petri nets and CPN Tools for modelling and validation of concurrent systems”, International Journal on Software Tools for Technology Transfer, vol. 9(3-4), 2007, pp. 213-254.
- [17] WFMC, Workflow Management Coalition Workflow Standard, Workflow Process Definition Interface – XML Process Definition Language (XPDL)(WFMC-TC-1025), Technical report, Workflow Management Coalition, Lighthouse Point, Florida, USA, August 2012.
- [18] D. Lee, H. Shin, and B. K. Choi, “Mediator approach to direct workflow simulation”, Simulation Modelling Practice and Theory, vol. 18(5), May 2010, pp. 650-662.
- [19] Y. Zheng, Y. Fan, and W. Tan, “Towards workflow simulation in service-oriented architecture: an event-based approach”, The 1st International Workshop on Workflow Systems in Grid Environments (WSGE ’06), vol. 20(4), March 2008, pp. 315-330.

Development of Real-time Simulation Models

Integration with Enterprise Information Systems

Konstantin Aksyonov, Eugene Bykov, Olga Aksyonova, Anna Antonova

Dept. of Information Technology

Ural Federal University

Ekaterinburg, Russia

wiper99@mail.ru, speedmaster@inbox.ru, K-36398@planet-a.ru, bpsim.dss@gmail.com

Abstract—The paper discusses the integration method of simulation models used within the enterprise information system. The integration problem is presented on one sample class of models – the real-time models that are used in control, diagnostics and decision making processes. The suggested method is based on multi-agent approach with distributed knowledgeable agents. The dynamic model consists in the multi-agent resource conversion process model that supports multi-approach modeling, including discrete-event, agent-based, queuing systems. To substantiate suggested technical decision of the integration module existing message brokers were analyzed. The subject area ontology is presented. It is used for semantic data integration that is required for simulation modeling of technological processes, business processes and logistical processes. In order to achieve a cross-platform system, the Java language is used for development.

Keywords-model integration; resource conversion process.

I. INTRODUCTION

Currently, multi-agent approach is among the perspective directions of enterprise management systems [1][2]. The paper deals with the representation of distributed corporate information systems of a metallurgical enterprise in form of a multi-agent system. Such systems consist of multiple interacting agents that solve the goals, set by analysts [3].

Development of state-of-the-art technologies allows large industrial enterprises to obtain and store vast data volumes that define technological, logistical and commercial problems of an enterprise. These data may be used for simulation modeling of various aspects of its activity. Simulation results may be used for quality control of manufactured products, defect prevention, optimization of logistical and commercial schemes. Integration of these models into the enterprise control circuit by interaction with the corporate information system and development of unified software is a pressing task.

After the introduction in section 2 we start from the state of the art overview for the area of simulation modeling tools with real-time capabilities. We identify the richest systems from the functionality point of view. In part 3 we present the architecture of suggested system, describing role of the agents. We move on to application of the system to metallurgical production. Here we present and discuss the ontology of the metallurgical enterprise. We get to the role of

each agent within the system and discuss their interaction. In part 4 we present our research on how the developed models interact with each other. In part 5 we describe the integrating data model. Part 6 presents the mechanisms for semantic integration of data sources. In part 7 we compare our approach to similar implementations in other simulation systems in terms of performance. Finally in part 8 we present the conclusions from our research.

II. CURRENT STATE OF TOOLS

The development trend of enterprise information systems focuses on wide application of Internet technologies. Currently the commercial simulation systems available on the market, including AnyLogic [4], ARIS [5], G2 [6], are all desktop applications. Additional requirements for simulation modeling tools for team development of comprehensive simulation models include support for multi-user environment, availability of model access on the Internet and running simulation experiments on the Internet.

A comparison of systems [7] showed that the most part of the functionality is included in AnyLogic [4] and BPsim [8] products. Only G2 is developed focusing on the direction of service oriented architecture [7]. At the moment the SaaS (software-as-a-service) technology [9] is the most convenient in use, optimal in performance and client software requirements. The end user in this case is the analyst or decision making person. Thus, the pressing problem is development of the model integration software for simulation modeling servers, using service oriented approach.

III. SYSTEM ARCHITECTURE

A multi-agent system architecture will further be discussed based on a sample integration system of automated system models of metallurgical production. It contains the following software agents:

- Data exchange agent. It is used for actualization of model parameters and data transfer (including experiment results) into corporate information system,
- Modeling agent. It is used to solve process control tasks in real-time on the basis of real-time models,
- Message exchange agents. It provides interaction between data exchange and modeling agents. This

agent decides when to activate real-time model, based on occurring events and activation rules, and also transfers messages into the corporate information system, e.g., into a MES-system or to a corresponding analyst's (specialist, technologist) workstation.

The method of design, development and operation of real-time models is based on the methodology of business process analysis and development of information systems. It includes integration of structural and object-oriented approaches, simulation and multi-agent modeling [10] and consists of the following stages:

- Design of simulation model in the model definition module,
- Running experiments for model verification and adequacy checks in simulation module. BPsim.MAS system is used for this task at the stage of schematic design [10]-[11],
- Design of real-time model for its further use in model integration module and interaction with other sub-systems of the corporate information system. BPsim.SD tool [12] is used for this task at the stage of schematic design. It implements the following design stages:
 - Design of architecture for the model integration modules by dataflow diagrams, use-case diagrams and sequence diagrams of unified modeling language,

- Subject area ontology representation in form of class diagram,
- User interface modeling,
- Testing and debugging the real-time model in corporate information system,
- Operation.

Use of real-time models means that modeling time must be less or equal to a set value, and modeling must be completed before the next portion of data is received from the corporate information system. Thus, the following features need to be considered during the models integration:

- Performance. Architecture of automated system for metallurgical production must be oriented towards maximum use of server resources,
- Scalability. Models must be able to run simultaneously on multiple computers, as well as effectively use multi-core and multi-CPU computers.

To provide these features each model needs to be executed as a standalone process. Special mechanisms, included with the integration module should be used for the interaction of processes.

Integration is suggested to be performed at the data level. Each model performs analysis of data, received from the data storage. Modeling results are transferred either into the data storage or immediately into the corporate information system.

In general case, the following data integration levels may be distinguished [13]: physical, logical and semantic. A

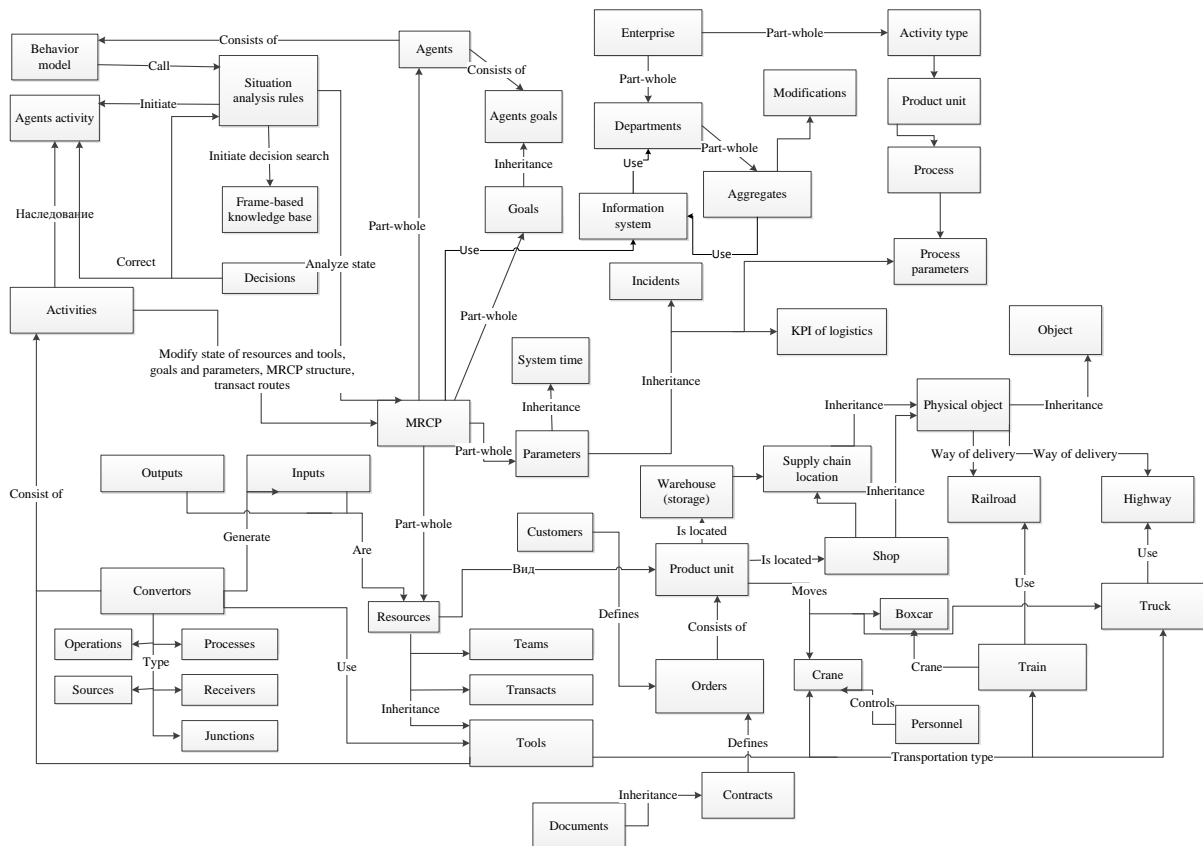


Figure 1. Subject area ontology

single ontology of subject area needs to be developed for consideration of semantic properties.

Ontologies are defined as a result on subject area analysis. In our case, the approach suggested by Girardi et al. [14] has been used. It is based on the Chen's model "entity-relation", since all data is suggested to be stored in relational database. The model has been extended in a way to be able to store other "entity-relation" models and related data.

The method has been extended with such features as availability to process cause-and-effect relations and knowledge of decision making people. Semantic model of the multi-agent resource conversion processes [11] was used for this. It was further extended with the elements of logistical projects ontology, presented by Kowalski et al. [15], and adapted to specific features of logistical problems, related to metallurgical production. Also, the ontology included elements of technological and business process. The designed ontology is presented on Figure 1.

Model integration method focuses on several problems [13]. They are briefly discussed further.

The automated system for monitoring, control, modeling, analysis, and optimization of the full production cycle of metallurgical production, due to specific requirements of the automation object, consists of a large number of various modules, each of which performs a specific task. Together they monitor state of the industrial objects, check parameter validity, model consistency, analyze and prepare recommendations for optimization of the full production cycle of metallurgical production. These recommendations are based on integration of mathematical models of technological, logistical, and business processes of an enterprise. Thus, the automated system for metallurgical production may be considered a distributed multi-agent system. Separate modules are represented with software agents with complex behavior and communicative capabilities.

Automated system for metallurgical production is an open multi-agent system. This consists in bi-directional interaction with multiple information systems of a metallurgical enterprise, related to such classes of information systems as ERP (Enterprise Resource Planning) [16], MES (Manufacturing Execution System) [17], automated technological process control systems, technologists' automated workstations.

An automated system for metallurgical production consists of the following modules (or main agent types):

1. Enterprise automated systems data exchange. Technically this corresponds to the enterprise services bus
2. Data preparation
3. Enterprise processes optimization
4. Models integration – allows use of models in decision making tasks in real-time.

Certain problems make use of data storage, query constructor, and model design modules that are also included in the automated system of metallurgical production.

Architecture of the automated system is implemented in the way that load ration of specific agents may result in

copying of these agents in order to distribute and balance load.

One of the applied directions of multi-agent technologies is planning. The concept of an agent corresponds to hardware or software implemented entity, which is capable of acting for the benefit of goal achievement set by the owner or user. Agents possess certain intelligent capabilities [18][19].

A sample application of the multi-agent system for planning operation of a flexible production system is discussed in [18]. We may name the following advantages of the multi-agent system:

1. Formalization of decision making points in form of the agents. The points include specific situation processing scenarios. Technically this process is a part of knowledge formalization stage.
2. Planner is dynamically embedded by means of interaction of specific element of the multi-agent system and thus is ready to modify the plan in case of delays or unexpected (unintended) situations. The planner works in real-time.
3. Agent network, interconnected with relations, self-coordinates its activity.

An additional benefit of multi-agent planning is the capability of automated information sharing between process individuals about changes of controlled object, which introduces control transparency. Subject area knowledge is being formalized during development and deployment of the planning multi-agent system, the decision making process is automated. Thus we ease activity, related to decision making.

Agents may be separated into three following types: reactive, intelligent and hybrid [18][19]. Reactive agents make decisions on the basis of "Situation-Response" rules. Intelligent agents solve the set tasks according to its goals, using common limited resources and knowledge of external world. Hybrid agents have features of both classes.

Agents of the automated system for metallurgical production, that are immediately operating in control and decision making tasks, may have goals presented on Figure 2. Intersection of their goals may be present. Thus, agents have to co-operate. In order to achieve a common goal, agents use messaging.

Interaction of agents of the metallurgical enterprise corporate system introduces problems, related to identifiers of the very same objects and parameters in different data storages. In addition, there is a dis-synchronization of single object-related processes in time.

In order to fix such problems, the data exchange agent is capable of transforming its internal identifiers into identifiers of other agents and vice-versa. Apart from this, the messages are dispatched, which help other agents to fix the problems, related to process arrangement in time.

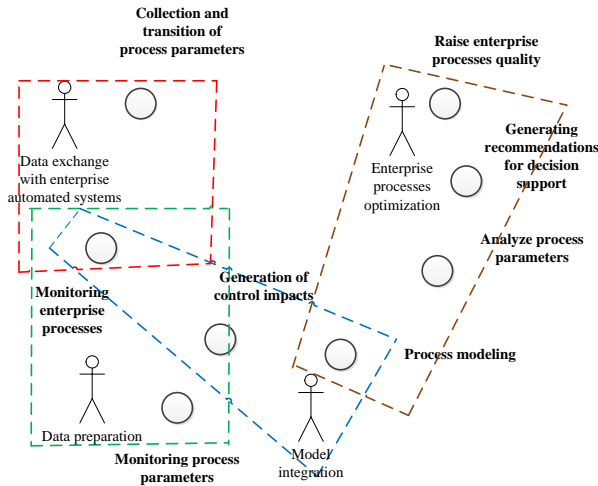


Figure 2. Agent goals

The development of the automated system for metallurgical production is based on decision support method for information system development. This method, in turn, is based on multi-agent approach. The method is supported by the products of BPsim family [11][12], which allows definition of hybrid agents on production and frame-based knowledge bases.

IV. MODEL INTERACTION

The most effective way of interaction of model integration system and automated information system of an enterprise consists in automated obtaining the data required for modeling directly from the automated information system (Figure 3). In order to implement this method, we suggest using the Messages queue system, which itself constitutes the architecture and intermediate level software, which collects, stores and distributes messages between subscribers.

Existing message brokers have been analyzed during research. All of them provided guaranteed message delivery between applications. Analysis results are presented in Table 1.

Since implementation, based on REDIS [20] and Socket.IO message exchange [21], is simpler, they were selected for data exchange between the automated information system and automated system of metallurgical production.

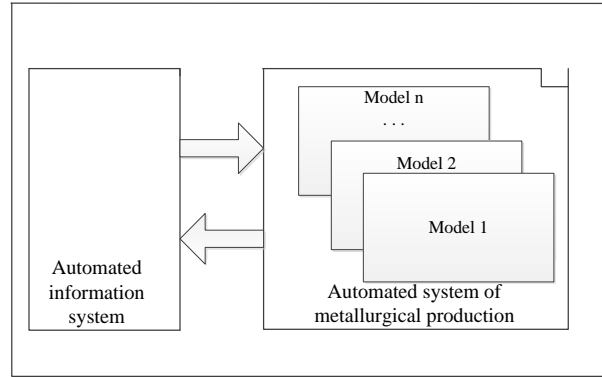


Figure 3. Interaction of integration module with the corporate information system

V. DESIGN OF INTEGRATING DATA MODEL

Integrating data model represents the basis of the common user interface in the integration system. Since the web-interface is suggested for model integration system, a decision, based on JSON (JavaScript Object Notation) [22] and XML (eXtensible Markup Language) standards for the integration model seems reasonable.

The MVC (Model-View-Controller) [23] concept is suggested as the main concept of the model integration system development. The concept utilizes several design templates, which allows the application data model, user interface and user interaction are distributed between three specific components, when modification of just one component has minimum impact on other ones. Model integration system includes common classes that implement typical procedures for data obtaining from the automated information system, as well as presentation of modeling results (parameter values, graphs, etc.).

Since the integration module has the multi-agent structure, the agent elements need to correspond to certain elements of MVC. To make things easier, consider a reactive agent with a single rule: “if a>b, then a=a-b”. Figure 4 shows the dataflow diagram that presents operation of such agent. Data storages correspond to work memory, which is required to store the variable. Operations on the diagram are all If-Then rules. The agent formulae in software implementation that are stored in “If” and “Then” rules of an agent, are transferred into method definition of the corresponding class.

Thus, from the MVC point of view, work memory represents the *Model*, while the logical output machine together with agent rules form the *Controller*. When visualization of modeling results is required by the user, corresponding classes would represent the *View*.

TABLE I. COMPARISON OF QUEUING BROKERS

Criteria	Redis	RabbitMQ	ActiveMQ	Socket.IO
Performance	high	high	high	high
Scalability	high	high	high	high
Clustering	no	yes	yes	no
Java support	yes	yes	yes	yes
Ease of use	high	average	average	high

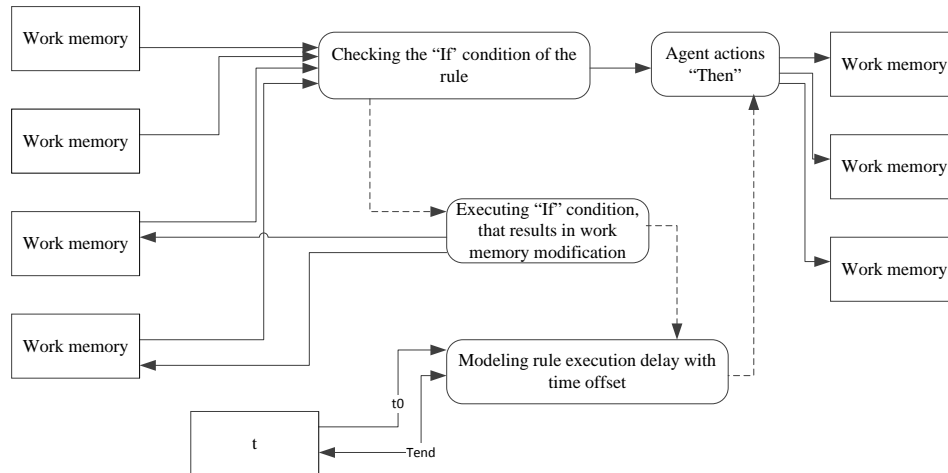


Figure 4. Sample DFD diagram for the reactive agent with one rule

VI. DEVELOPMENT OF MECHANISMS FOR SEMANTIC INTEGRATION OF DATA SOURCES

Object-relational mapping is used as a means of data sources semantic integration. This is a programming technology, which allows conversion of incompatible model types between relational data storage and programming objects. Such technology is implemented in ORM (Object-Relational Mapping) systems [24]. After analysis, two systems have been selected for further development, namely, Morphia [25] and Cayenne [26] due to the following factors:

- Their features fully satisfy the requirements of model integration system,
- Cayenne has a convenient feature of visual development of connection of software classes and entities in the database. This reduces the time required for development and debugging.

The prototype of model integration module for automated information system of metallurgical production has been developed after analysis. Since at this stage, some models of the system are yet to be implemented, testing and running the experiments used the emulated parameter inflow from the technological process. Data required for the simulation model integration module were forwarded to corresponding models for analysis and result output. Model results transfer into the corporate information system for further processing has also been emulated.

VII. PERFORMANCE COMPARISON

An enterprise uses its own quality assurance software, which includes the enterprise process definition module for design of simulation models of processes under research, and process optimization module for experimenting with the models and searching for management decisions. These two modules are based on multi-agent modeling and the concept of big data. Authors compare the metallurgical manufacturing process model definition with the enterprise software and the popular simulation tools Plant Simulation

[27], Simio [28], and AnyLogic [4]. The description of the model itself goes beyond the scope of this work, and only the comparison results are presented.

We assume that models are equivalent and produce averagely the same output. As an effectiveness criterion we use the duration of experiment on the same hardware with animation set to off (Figure 5).

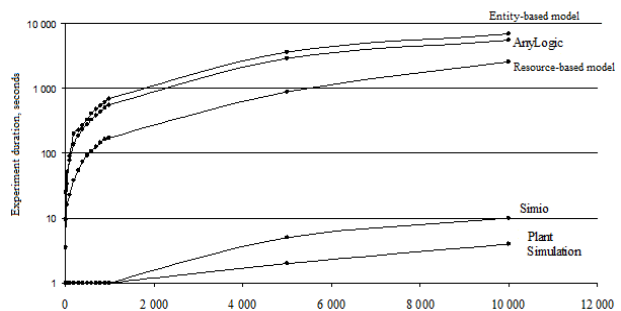


Figure 5. Experiment duration and number of processed product units

The analyzed simulation modeling tools may be separated into fast tools (Simio, PS) and slow tools (AnyLogic and enterprise process optimization module). Optimization module speed is related to detailed journaling of log tables and statistics on model variables and entity instances. No other simulation tool provides these statistics. Resource-based model works faster in the same simulation system. This fact is related to computational resources being spent on queues processing in the process optimization module. The slowest experiment duration was 2 hours and 13 minutes, which may be applicable in case of non-real-time decision making.

After analyzing simulation results we may conclude:

1. All models are adequate to logistical processes of a metallurgical enterprise,
2. Simulation speed is applicable for all simulation systems for various production volumes,

3. Simio and PS have an advantage in simulation speed for simulation of logistical processes of an enterprise,
4. CPU and RAM load are applicable for a short (under 10 minutes) simulation experiment without animation in all systems,
5. CPU and RAM load are applicable for a long non-animation experiment (over 1 hour) for all systems, except PS (due to hang up) and Simio (due to high RAM load),
6. Advantage of the enterprise optimization module from the RAM load point of view for “short” and “long” non-animated experiments.

VIII. CONCLUSION

Use of simulation modeling for analysis of technological, logistical and business problems of an enterprise is a perspective direction. The discussed method of simulation models integration has been implemented in practice and has successfully passed the tests.

The automated system for metallurgical production may assist in the following areas:

1. Collection and storage of information about enterprise products and processes,
2. Analysis of quality of products, diagnosis of production stages with most faulting operations, with full information of production cycle,
3. Application of models in decision making and control tasks. In case a model used in control process diagnoses a significant deviation from quality indicators for a product unit, it generates a signal and forwards in to a MES system, in order to reassign routes for further processing.

ACKNOWLEDGMENT

Research is conducted under the terms of contract № 02.G25.31.0055 (project 2012-218-03-167).

REFERENCES

- [1] L. Cao, V. Gorodetsky, and P. A. Mitkas, Guest Editors' Introduction: Agents and Data Mining, *IEEE Intelligent Systems*, vol. 24, no. 3, pp. 14-15, May/June, 2009
- [2] P. Sridhar, S. Sheikh-Bahaei, S. Xia, and M. Jamshidi, Multi-agent simulation using discrete event and soft-computing methodologies, *Proceeding of: Systems, Man and Cybernetics, 2003. IEEE International Conference on*, Volume: 2
- [3] V. Gorodetsky, O. Karsaev, V. Konushy, W. E. Matzke, E. Jentsch, V. Ermolayev, Multi-agent Software Tool for Management of Design Process in Microelectronics, *Intelligent Agent Technology, IEEE / WIC / ACM International Conference on*, pp. 773-776, 2006 *IEEE/WIC/ACM International Conference on Intelligent Agent Technology (IAT'06)*, 2006
- [4] <http://www.anylogic.com>
- [5] http://www.softwareag.com/corporate/products/new_releases/aris9/overview/default.asp
- [6] <http://www.gensym.com/en/product/G2>
- [7] O. P. Aksyonov a, K. A. Aksyonov, V. D. Kamelsky, and A. L. Nevolina, Analysis of organization of distributed multi-user work in business processes simulation modeling systems // *Modern problems of science and education. – 2012. – № 5*; URL: <http://www.science-education.ru/105-6936> [retrieved: 05, 2014]
- [8] <http://www.bpsim.ru>
- [9] http://en.wikipedia.org/wiki/Software_as_a_service
- [10] K. A. Aksyonov, A. S. Antonova, and I. A. Spitsina, Analysis and synthesis of resource conversion processes based on simulation modeling and intelligent agents. // *Science and technology news, St. Petersburg State Technical University, № 1 (115) 2011. Informatics. Telecommunication. Control. St. Petersburg*, pp.13-20.
- [11] K. A. Aksyonov, Theory and practice of decision support tools. Germany, Saarbrücken: LAP LAMBERT Academic Publishing GmbH & Co. KG, 2011.
- [12] K. Aksyonov, I. Spitsina, E. Bykov, E. Smolij, and O. Aksyonova, Computer-supported software development with BPsim products family – integration of multiple approaches // *Proceedings of the 2009 IEEE International Conference on Information and Automation (ICIA)*. (22-25 June 2009). - Zhuhai/Macau, China, 2009, pp. 1532-1536.
- [13] M. R. Kogalovsky, Methods of data integration in information systems. Moscow, 2010. <http://www.cemi.rssi.ru/mei/articles/kogalov10-05.pdf> [retrieved: 05, 2014]
- [14] D. Girardi, J. Dirnberger, and J. Trenkler, A Meta Model-Based Web Framework for Domain Independent Data Acquisition // *ICCGI 2013: The Eighth International Multi-Conference on Computing in the Global Information Technology*. Nice, France, 2013, pp. 133-138.
- [15] M. Kowalski, S. Zelewski, D. Bergenrodt, and H. Klupfel, Application of new techniques of artificial intelligence in logistics: an ontology-driven case-based reasoning approach // *Proceedings of European Simulation and Modelling Conference 2012 (October 22-24, 2012, FOM University of Applied Sciences)*. — Essen, Germany, 2012, pp. 323-328.
- [16] http://en.wikipedia.org/wiki/Enterprise_resource_planning
- [17] http://en.wikipedia.org/wiki/Manufacturing_Execution_Systems
- [18] N. R. Jennings, On agent-based software engineering // *Artificial Intelligence. — 2000, vol. 117, — P. 277-296.* - URL: <http://www.agentfactory.com/~rem/day4/Papers/AOSE-Jennings.pdf> [retrieved: 05, 2014].
- [19] M. Wooldridge, Agent-based software engineering // *IEEE Proc. Software Engineering*, no. 144 (1), 1997, pp. 26–37.
- [20] <http://redis.io/>
- [21] <http://socket.io/>
- [22] <http://www.json.org/>
- [23] <http://en.wikipedia.org/wiki/Model%E2%80%93view%E2%80%93controller>
- [24] http://en.wikipedia.org/wiki/Object-relational_mapping
- [25] <https://github.com/mongodb/morphia>
- [26] [http://en.wikipedia.org/wiki/Cayenne_\(programming_language\)](http://en.wikipedia.org/wiki/Cayenne_(programming_language))
- [27] <http://www.plm.automation.siemens.com>
- [28] <http://www.simio.com>

RISK-DET: ICT Security Awareness Aspect Combining Education and Cognitive Sciences

Guillaume Schaff, Carlo Harpes, Matthieu Aubigny
itrust consulting (Luxembourg)
{schaff, harpes, aubigny}@itrust.lu

Marianne Junger
University of Twente
m.junger@junger.nl

Romain Martin
University of Luxembourg
romain.martin@uni.lu

Abstract— This paper explains the main innovation of a risk assessment tool, called RISK-DET, which will include an ICT risk awareness aspect supported by a specific application: Voozio 2.0. The design of the RISK-DET tool considers the implementation of the emergent ICT (Information and Communication Technology) Risk Detection Skill (IRDS) concept. Today, the users' inability to detect a risk situation is a real security problem and represents a societal challenge. According to the results of a security experiment based on a malicious smartphone application called Voozio 1.0, the main reason for this problem is the absence of effective ICT security awareness training programs adapted to users' needs. To prove and confirm this hypothesis, we aim to evolve the Voozio application in the 2.0 version. This new version will be able to determine the ability of ICT users to detect a risk situation and improve it by combining cognitive sciences and education technologies. We will describe here the specifications of the new version of Voozio. We also present the Voozio 2.0 implementation framework.

Keywords-E-learning; ICT security awareness; social engineering; cyber-security; cognitive sciences; risk perception; education science; human-computer interaction.

I. INTRODUCTION

The rising use of new Information and Communication Technologies (e.g., smartphones, digital tablets, laptops, etc.) in our daily life has increased our vulnerability to new cyber-attacks [1][2][3]. With the cyber-criminal professionalization, the ICT threats (virus, phishing, scamming) are more sophisticated and their impact can be very significant on our lives [4][5] (personal data theft, ransomware). In parallel, the current security mechanisms are not yet sufficiently adapted to face these new types of ICT attacks. To limit their impacts, several researchers have developed anti-phishing training programs [6][7][8]. However, these programs are not sufficient to limit an ICT attack and that is why the ICT users' ability to detect a risk situation (ICT Risk Detection Skill (IRDS)) should be improved. To improve their ICT Risk Detection Skill (IRDS), users should be able to adopt good security practises when faced with cyber-threats. Here, an ICT risk is the probability that a threat exploits ICT vulnerabilities (e.g., malicious email, phishing link, etc.) which impacts the confidentiality, integrity or availability of information. To limit these impacts, we propose to develop a security awareness aspect to improve the users' ICT Risk Detection Skill level. Research has been done to develop, for instance, security awareness and education programs [9]. A few experiments have been executed. They show that, through training, ICT users develop new skills improving their ability to detect ICT attacks (detect false/malicious email, malicious spam). Still, a lot of work needs to be done in the field of cybercrime prevention, in particular ways to prevent users becoming victims of social engineering [10][11].

Accordingly, the present paper aims to present a new tool which has two aspects: firstly, a measurement aspect to determine the users' ability to detect a risk situation (IRDS); secondly, an ICT security awareness aspect to improve IRDS. This global ICT security awareness solution will not be an isolated solution, but will be integrated into a set of several tools, developed in the framework of the FP7 TREsPASS project [18]. This European project aims to combine technical and social sciences in order to develop methods and tools to analyse and visualise information security risks in dynamic organisations. The expected outcome of this project is an "attack navigator" indicating which attack opportunities are possible in a targeted organisation, which of them are the most critical, and which countermeasures are most effective. To identify potential attack opportunities, the RISK-DET tool contributors (who are also participants of the TREsPASS project) aim to develop an additional risk assessment tool focused on social sciences.

The present paper first explains the role of cognitive sciences in ICT security awareness. Secondly, it presents the security awareness aspect (represented by the Voozio 2.0 application) of the RISK-DET tool. Thirdly, it describes the implementation framework of the Voozio 2.0 application. Finally, we conclude by describing the next steps, the research hypothesis and expected results.

II. COGNITIVE SCIENCES IN ICT SECURITY AWARENESS

With a lack of ICT security awareness, users are strongly susceptible to social engineering and phishing attacks [13]. As part of our previous research [12], we created a malicious application on Google Play and used it to test the ability of a representative sample of informed users to detect a risk situation on their smartphones. Our results show that more than half of the targets submitted their personal information to our malicious smartphone application called Voozio. It seems feasible to conclude that the majority of ICT users are relatively susceptible to IT attacks. We concluded that the cognitive aspect has an impact on the ICT users' reaction when faced with risk situations.

In general, cognitive sciences are based on the study and modelling of users' perception and particularly the risk perception. That is why we have decided to integrate cognitive sciences a.o. in Voozio 2.0, to obtain more in-depth results. In general, scholars have argued that ICT users have a major role to play in enhancing global ICT security [14]. The purpose of the present study is to present an evolved version of Voozio that will integrate cognitive, social and education sciences. As users are one key element to avoid ICT attacks, the cognitive aspect of the ICT users should be strongly developed thanks to efficient security awareness solutions. The development of Voozio 2.0 will be based on the cognitive sciences, in line with our IRDS concept. What do we mean when we use the term "cognitive

sciences”? Cognitive sciences are a set of scientific disciplines dedicated to the description, explanation, and appropriate simulation mechanisms of human thought [13]. Based on the approach of cognitive sciences applied to ICT, the design of Voozio 2.0 should include two following properties: firstly, a function for evaluating the IRDS of the users whilst considering several psychological factors (character, personality, reason ability), and secondly, a function for generating a relevant training program in line with users’ needs. Cognitive science is tied with Educational Technology, and we believe that this discipline should be integrated into the new version of our application by incorporating an education program which can be adapted for individual users. For Voozio 2.0, the cognitive science integration will consist of collecting test subjects’ psychological factors (character, personality, reason ability) with a dedicated questionnaire, validated by an expert panel. The psychological factors will allow the Voozio 2.0 administrator to establish subjects’ psychological profile. This profile will be used in the IRDS measurement phase and refined continuously in the next phase.

III. ICT SECURITY AWARENESS ASPECT PRESENTATION

The Voozio 2.0 application is composed of a Computer Assisted Test (CAT) based on identified ICT risk situations, e.g., malicious email and/or website. The risk situations will be presented to the tested users through short videos, pictures, or games. The users’ aim is to identify what level of risk the ICT situation presents. During this CAT, the researchers will analyse the users’ behaviour based on several factors, such as time taken to answer, uncertainty, etc. In order to obtain an accurate way of scoring, the risk situations will follow a precise graduation depending on the level of danger based on their impacts. The analysis of users’ reactions when faced with an ICT risk scenario will be scored and analysed according to a pre-defined scale. After the IRDS measurement phase, Voozio 2.0 will introduce an educational program composed of e-learning modules. The e-learning modules will be adapted to the users’ risk perception ability and reactions. Here, we aim to improve the IRDS level of the users. As shown in Figure 1, the increase in IRDS level is supported by six disciplines.

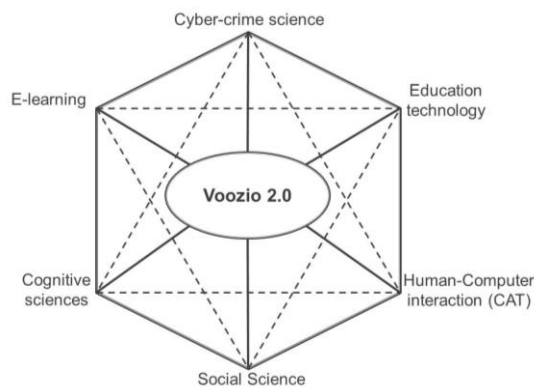


Figure 1. Disciplines of the IRDS increase phase

A sustainable improvement of the IRDS requires as many interactions as possible between these six disciplines. Therefore, Voozio 2.0 will integrate the Cognitive and Social sciences, as well as Human-Computer interaction in

the IRDS level measurement phase. These disciplines will consider test subjects’ psychological profiles in the IRDS measurement. Moreover, the researchers aim to include E-learning, Education and Cyber-crime science [14] in Voozio 2.0 conception to increase the users’ IRDS level. Furthermore, the ICT security awareness aspect will allow Voozio 2.0 administrators to provide training programs adapted to users’ training needs depending on the IRDS test results and subjects’ reactions (user behaviour). The social sciences play an important part in the Voozio 2.0 general process by bringing an additional precision level to the IRDS measurement.

IV. VOOZIO 2.0 IMPLEMENTATION

An adapted IRDS measurement test should be integrated in Voozio 2.0 to determine users’ training needs. This test is intended for a group of users and necessitates a preliminary requirement; the agreement between “the trainer” and the management team of the organisation (e.g., Managing Director, IT Manager, etc.). Here, the trainer corresponds to the Voozio 2.0 administrator who will submit the IRDS measurement and the training program to the targeted staff. The test will consist of sending an email which proposes to install the Voozio 2.0 application to the entire staff of a private/public organisation. After installation, Voozio 2.0 will generate several risk situations (e.g., malicious email and spam). We stress that the risk situations generated will not present any danger for the users’ devices. The test will only measure how many of the users are unaware of the threat and “fall into the trap”. Once the data has been collected and interpreted, the trainer will be able to establish a first IRDS level evaluation for all staff. After the test, Voozio 2.0 will send a training program composed of e-learning modules focusing on good ICT security practices to all the employees. A short time (one or two weeks) after the ICT security awareness program, an additional risk situation will be generated by the Voozio 2.0 application to evaluate the reliability of the provided training programs. A comparison study will be performed by the trainer to note the difference between the IRDS test results obtained before and after training program. Thanks to this test, we will be able to measure the efficiency of the ICT security awareness aspect based on any improvements. To have relevant results, application implementation will be on a large organisation to measure and improve the ICT Risk Detection Skill level of their employees (pool of testers).

If we succeed to reach the critical mass for the test pool, we will be able to collect relevant results needed to establish a statistical study. According to a precise and automatic analysis of the results, the solution will create a user classification depending on their ICT attack vulnerability and will generate dedicated ICT training programs specific to different user groups to enhance ICT Risk Detection Skill. As in our previous research, no personal data will be retained during this test and test subjects’ privacy will be strictly respected.

V. MODEL ANALYSIS

The preliminary work of researchers consisted to identify related works [15][16][17] and establish a formal

state-of-the-art in risk perception domain. Before the experimental phase, we aim to submit Voozio 2.0 to an expert panel which have been worked on the same domain to analyse the proposed model and give their feedback on it. The expert panel feedback will allow us to refine the prototype to be adapted to users' needs. The expert panel will be selected by the Voozio 2.0 designers and will group experts from industrial and academic organisations.

VI. CONCLUSION

In the FP7 TREsPASS project framework, we aim to develop the RISK-DET risk assessment tool which will include an ICT security awareness aspect based on Cognitive and Education Sciences called Voozio 2.0 (based on the results of the previous experience conducted with the smartphone application Voozio 1.0). Analysis of similar works will allow us to define the Voozio 2.0 technical and functional specifications and implement our IRDS concept. After the Proof of Concept development phase (planned 3rd trimester 2014), we suggest that the Voozio 2.0 beta version could be validated by an expert panel and tested on an organisation of 1000 employees during an experimental phase. This phase will allow us to test our innovative IRDS measurement methodology and make our research work (results interpretation and model analysis). In the long term, Voozio 2.0 could provide a commercial ICT risk awareness solution. However, we should keep in mind that the effect of training can decrease over time. Due to this fact, the tool shall be based on a core system using versatile contents which will be set up over time, depending on the context of the ICT threat and on the security maturity of the organisation. Voozio 2.0 will enable us to propose a permanent IRDS measurement and security awareness solution. The first result of our work is the Voozio 2.0 technical and functional specifications document which will be published during the ICCGI 2014 conference.

ACKNOWLEDGEMENTS

The research leading to these results has received funding from the European Union Seventh Framework Programme (FP7/2007-2013) under grant agreement n° 318003 (TREsPASS). This publication reflects only the author's views and the Union is not liable for any use that may be made of the information contained herein.

REFERENCES

- [1] K. Choo, "The cyber threats landscape: Challenges and future research directions", *Computer & Security*, vol. 30, iss. 8, 2011, pp. 719-731
- [2] M. Yar, "The novelty of Cybercrime, an Assessment in Light of Routine Activity Theory", *European Journal of Criminology*, vol. 2, iss. 4, 2005, pp. 407-427
- [3] P. Williams, "Organized Crime and Cybercrime: Synergies, Trends and Responses", *Global Issues*, vol. 6, iss. 2, 2011, pp. 22-26
- [4] T. Jagatic, N. Johnson, M. Jakobsson, and, F. Menczer,, "Social phishing" *Communication of the ACM*, vol. 50, iss. 10, 2007, pp. 94-100
- [5] A. Gazet, "Comparative analysis of various ransomware virii", *Journal in Computer Virology*, vol. 6, iss. 1, 2010, pp. 77-90
- [6] S. Sheng, et al, "Anti-Phishing Phil: the design and evaluation of a game that teaches people not to fall for phish", *Proceedings of the 3rd symposium on Usable privacy and security*, 2007, pp. 88-89, ISBN: 978-1-59593-801-5.
- [7] P. Kumaraguru, et al, "School of phish: a real-world evaluation of anti-phishing training", *Proceeding of the 5th Symposium on Usable Privacy and Security*. Iss. 3, 2009, ISBN: 978-1-60558-736-3.
- [8] X. Luo and, Q. Liao, "Awareness Education as the key to Ransomware Prevention", *Publishing models and article dates explained*, vol. 16, iss. 4, 2007, pp. 195-202.
- [9] M.E. Thomson and, R. von Solms, "Information security awareness: educating your users effectively", *Information Management & Computer Security*, vol. 6 iss. 4, 1998, pp.167 – 173
- [10] B. Claverie, "Cognitive, Science et pratique des relations à la machine à penser", *Revue des sciences de l'éducation*, vol. 34, iss. 1, 2005, pp. 227-228
- [11] Hartel, P.H., Junger, M. and, Wieringa, R.J. "Cyber-crime Science = Crime Science + Information Security", *Technical Report TR-CTIT-10-34*, Centre for Telematics and Information Technology University of Twente, Enschede, 2010
- [12] G. Schaff, C. Harpes, R. Martin and, M. Junger, "An application to estimate the cyber-risk detection skill of mobile device users". *Sixth International Conference on Advances in Human oriented and Personalized Mechanisms, Technologies, and Services, CENTRIC 2013*, November 2013, Venice, Italy.
- [13] E. Albrechtsen. "A qualitative study of users' view on information security". *Computer & Security*, vol. 26, iss. 4, pp. 276-289
- [14] T. Vidas, E. Owusu, S. Wang, C.Zeng, L. Faith and, N. Christin. "QRirshing: The Susceptibility of Smartphone Users to QR Code Phishing Attacks" *Financial Cryptographie and Data Security*, 2013, pp. 52-69,.
- [15] L. Sjoberg, B-E. Moen, and, T. Rundmo, "Explaining risk perception. An evaluation of the psychometric paradigm in risk perception research", *Trondheim*, 2004.
- [16] N. Pidgeon, "Risk assessment, risk values and the social science programme: why we do need risk perception research". *Reliability Engineering & System Safety*, vol. 59, iss. 1, 1998, pp. 5-15.
- [17] B. Fischhoff, et al. "Risk perception and communication". *Oxford textbook of public health*, vol. 2, iss. 5, 2009, pp. 940-953
- [18] P. Hartel and, W. Pieter; "FP7 project TREsPASS press release", January 2013

ICT in Education: A New paradigm and old obstacle

András Benedek

Department of Technical Education
Budapest University of Technology and Economics
Budapest, Hungary
e-mail: benedek.a@eik.bme.hu

György Molnár

Department of Technical Education
Budapest University of Technology and Economics
Budapest, Hungary
e-mail: molnar.gy@eik.bme.hu

Abstract— In this paper, we review the phenomena and processes which, when viewed together, indicate the birth of a new education paradigm for students in higher education. Our approach to learning in particular is changed by digital tools. Our thinking process and learning methods are subject to development. This paper, while making references to this development, analyses the interactions between typical and atypical learning, in other words, formal and informal learning. It also explains the most important characteristics of the new education paradigm, which we symbolically call „Version 2.0” when addressing the new generation. When analysing these changes, we should review and interpret the most important participants and their roles (students and the teachers supporting them), old routines and future expectations and the new ICT solutions and their increasing role in supporting learning.

Keywords-Digital education; Digital learning; E-learning; Competence; Education, Teaching; Learning, Learning styles; Learning network; Learning theory; Knowledge; Web 2.0.

I. INTRODUCTION

The new methods of *e-learning* and „edutainment” could be born as a result of digital development. Edutainment is based on the idea of disseminating professional (or legal) information to students in an entertaining audio-visual environment provided by multimedia tools. The method is based on the realization coming from the latest results of education science that the efficiency of traditional teaching methods (with special regard to oral presentations) and additional electronic devices is multiplied when placed into an entertaining framework. The interest of students may be maintained continuously by impulsive presentations where information transfer is aided by music or visuals added at given points, facilitating emotional identification, active participation and retrieval of information from the memory later on.

Being aware of new ICT trends and their constructive application are essential in modern teaching and learning; this is why a whole chapter is dedicated to these topics as well. Modernizing education and training does not mean that teaching and learning cannot be experiential. Another chapter formulates theories on „experiential education” in a digital environment and its results, also verified by feedback from students.

The section on ICT (Information and Communications Technology) and Web2.0 opportunities in language learning

reflects the reforms in the syllabus. The guide based on the modern methodology of language teaching may be useful for both students and teachers in computer aided language education [1][13].

In our accelerated world, information transfer with special regard to presenting, communicating and receiving useful information items is of paramount importance. One of the most important practical issues in digital education is the development of micro-contents that correspond to the interests of students (culture, sport, specific professional topics) while fitting the screen of ICT tools and complying with content and format requirements. This is why the textbook on the topic written by András Benedek is now completed with a new chapter on developing and presenting micro-contents [2]. The information in this paper, presented without the reference to the books/papers, is our contribution and based on our own research explained in textbook cited.

The study will begin by outlining the new learning theories and contexts by examining new learning environments, forms and specific processes. The following sections will set out the altered time and space of learning as well as the use of modern, digital teaching materials. The paper ends with conclusions.

II. LEARNING THEORIES AND CONTEXT

A. Learning and knowledge – a new context

In our environment, new digital tools appear every day bringing along new trends. This section explains how to recognize these trends and their potential uses.

In the 21st century, people of the modern era found themselves in a new work environment. The social and economic networks surrounding individuals are more complicated than ever. Learning theory analysis typically examines the characteristics of social learning, where new learning methods and techniques meet, as a result of current changes, with special regard to developments in ICT. In such an increasingly rich learning space, learning how to apply new learning methods consciously and efficiently/successfully may prove to be an investment with a good return on the long run.

Physical networks (like the urban environment) and their virtual counterparts have changed our lives in two important aspects.

We are able to connect to and communicate with much more people than previously. This is partly the result of our accelerated lifestyle and partly of the sophisticated hierarchies organized around our multiple roles in life.

A new, virtual dimension is now attached to learning. In this respect we should note that our multiple connections are now much less restricted in time and space. In the developed world in particular, we are now able to contact anybody, anywhere, exchanging information and organizing our lives, once the technical background is available (a smart phone or a broadband Internet connection will suffice; see Figure 1.)

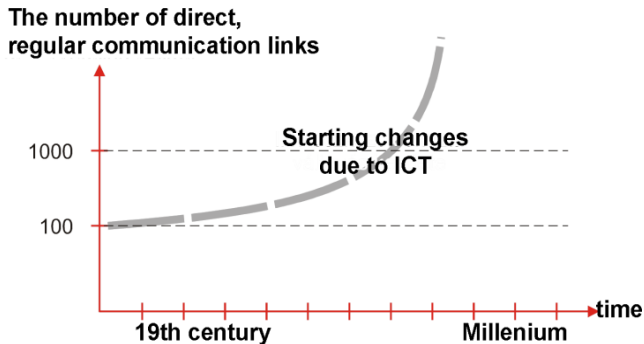


Figure 1. Changes in the system of our communication.

By now, knowledge has become a dynamic concept and acquiring knowledge a process with ever increasing spatial and temporal dimensions. On the one hand, education has diversified, corresponding to increasingly high levels; on the other hand, the time spent in education has almost doubled in the past hundred years, from 6-8 years to 12-16 years [3].

B. New Roles – Old Participants

In education, with special regard to institutionalized education, typically in schools roles (those who teach and those who learn) and participants (teachers and pupils or students) are relatively well defined. According to their age and acquired and acknowledged qualifications, teachers are formally positioned in the organization that has been created to establish the right conditions for teaching and learning. Traditional elements of the organizational structure of education, i.e. schools where the youth is educated, were founded as early as in the Ancient World and their development in the centuries to follow was slow, always paying heed to traditions.

Educational institutions perform several tasks related to preserving, transferring and updating knowledge. In addition to preserving knowledge and archiving its most valuable elements, participation in development is another important role of these institutions. It is only possible if education has a strong link with knowledge development and research, as several worthy traditions illustrate. The mission of higher education in this respect is particularly important; however, the quality of teaching and learning is also of paramount importance in primary and secondary education. In this respect, the direct and indirect, supportive presence of parents in addition to classic participants

(teachers and students) is very important in the first school years, when the role of informal learning is taken over by formal education. In the period when the basics are learnt in particular, pupils, teachers and parents are the most important participants, with their roles and weight varying over time and according to specific conditions.

The progressive approach to education claims the increasing success of student oriented systems that rely on the indirect support of parents and the direct contribution from teachers [4].

Teachers, as representatives of a unique profession model and significant human factors having a strong influence on social performance, stand in the cross section of public education and higher education. This subject remains in the focus of interest, its actuality guaranteed by the interactions of profession, job and profession and the actual practice representing challenges in several fields. Continuous reforms in the teaching profession are necessitated by changes and their management and new communication tasks.

The quality of the applied educational process greatly depends on the qualification, preparedness and actual experience of teachers. Nevertheless, in the case of student centered education, the commitment, motivation and continuous learning activity of students are equally important. These factors are summarized in the table below.

TABLE I. PARTICIPANTS IN THE LEARNING PROCESS.

Most important direct and indirect participants of the learning process – ideal characteristics	
Teacher	prepared, with practical experiences, committed to education
Student	open minded, continuously active in learning, committed to fulfil learning objectives
Parents	support the learning process, also directly but mostly indirectly

C. Learning Environment

In addition to human factors, the direct environment of learning also significantly affects learning success. If we are asked about our learning environment, the answer is relatively simple. Those participating in institutionalized education usually refer to the features of their classrooms and list the tools available for them in the actual technical framework. The definition of „education tools” in Pedagogical Lexicon (Encyclopedia of Education) was phrased some 15 years ago but it still applies today. This definition offers a comprehensive, systemized approach and divides education tools into four major categories according to various criteria (physical nature of tools, their effect on sensory channels, required technical equipment and user).

When defining the learning environment and interpreting the related development directions, the widespread use of the Internet, this unique „public utility” and the Web 2.0 movement may be considered those changes that have had the most significant impact on modern times [5][6].

One distinct feature of the new info-communication environment is this new, *virtual learning* system, which is increasingly capable of systemizing and transferring „learning objects” and organizing communication between teachers and students/pupils.

These initiatives, also developing at institutional level, emphasis the expansion of the „space” for both teachers and students where the elements of teaching and learning that is creating interest, knowledge transfer, receiving knowledge, illustration-experience-experiment-research-practice, conclusion, systemizing may be arranged in some sort of a rigid didactic structure and organized according to pre-set algorithm. This change obviously affects teachers, too, they have to find a way to harmoniously use their various tools (traditionally in the framework of curricula and textbooks but now also the tools available in the new environment) in this expanding space, in institutionalised education in particular. In the new environment, students also occupy a specific position as they are opened to new information yet need orientation and possibilities to develop, see Figure 2.)

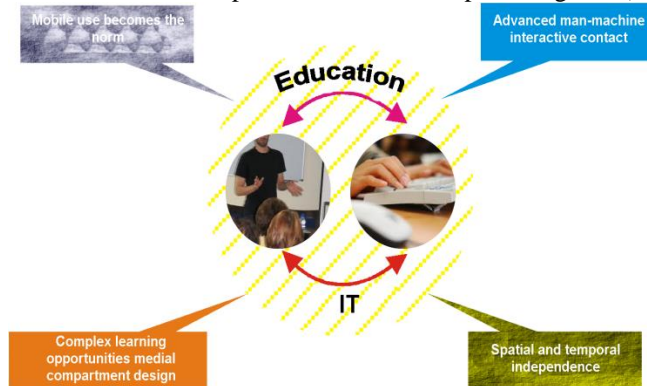


Figure 2. New features of education informatics.

As for future prospects, learning communities are expected to have an increasingly important role. These new communities are characterized by the common interest of the members above all. In these communities, learners interact, they learn together and create a shared collection of information sources.

The learning network is not just educational principle but also an environment where learning efficiently supports acquiring and continuously updating one’s knowledge, using the new educational theory. Major current trends defining the interaction between education and informatics are:

- Developed forms of human-machine interaction.
- Independence in space and time.
- General use of mobile devices.

- Possibility to create complex, medial „learning environments” [7].

D. Learning Theories

There is no simple answer to the question „is there a theory for learning?” The relevant scientific papers, some based on empirical tests, some offering theoretical concepts mostly from psychology experts but also from educational scientists of course could fill volumes. Learning is discussed from several aspects, from cell biology through formal logic to organization theory, to name but a few. Human learning is one of the most complex fields, seconded by the well-known experiments on animal learning. In the present textbook we only discuss general theory though the results of animal tests are often applicable in the case of human learning .

In an educational sense, learning means acquiring or developing some ability. The knowledge thus created is defined as a change to be evaluated (in performance, behaviour or knowledge) and a product of external impact, experience or practicing [8].

Human learning being such a complex discipline, several theories have been conceived about it. Some of these concern age related learning while others to different learning levels as determined by the quality of knowledge. One of the earliest and most comprehensive general theories phrased by psychologists is behaviourism. It is related to Piaget’s theory of cognitive development that emphasizes the importance of intrinsic motivation in learning. The concept states that human beings have an innate urge to see the world as coherent and stable and intellectual processes aim at creating equilibrium.

According to the traditional view, elementary learning in human beings is essentially associative. Our norms and habits have developed for millions of years by successfully or unsuccessfully associating phenomena to activities.

Indeed, the easiest way to learn is when at least two information inputs are active during the time of perception. In this case, the two processes are associated with each other in the brain.

Constructivism is a philosophic approach having a great influence at the same time on institutionalized education and learning in school. The theory states that students build up („construct”) their knowledge relying among others on „learning-by-doing”. This theory also emphasizes the importance of curriculum development and the important role of motivation to inspire students to acquire knowledge. The methodology of constructivism includes collecting and linking sources and the individual and group motivation of students. Constructivism offers a theoretical background for the teaching and learning of natural sciences and other educational programmes where experiments and project work are preferred.

In past decades, institutionalized education was based on the theory of cognitive learning, complying with the development of complex learning systems with diverse

subjects and sophisticated organizational structures. The findings of cognitive psychology yielded numerous data proving the theory that learning is based on the human ability to form a mental picture of perceived phenomena and then manipulate these mental representations. Piaget's work [15] in developmental psychology also indicated the existence of well-defined stages in the cognitive development of children. The location in time and contents of these stages may vary, yet they are a general feature and their order is invariably the same for everybody [9].

Modern educational theory „passed” the slowly spreading, century-old paradigm of continuous learning with one swift step around the Millennium. Influenced by ICT, modern theorists (Downes, Siemens [16]) made a significant, impulsive turn, recognizing the role of networks while still relying on classic psychological and philosophic schools such as behaviorism, cognitivism and constructivism. In the past decade, the taking up of the idea of lifelong learning both in theory and practice contributed to the successful integration of informal and non-formal education into the methodology of early development, socialization and social capital development via educational practice.

Connectivism (network based learning) has become an educational theory with a particularly significant impact providing an option for the educational exploitation of network theories, informatics and Web 2.0. According to the theory of George Siemens [16], published in 2005, connectivism is a learning theory specifically designed for the digital age. By defining the relevant terms, this theory is particularly important for the new educational paradigm as it focuses on network theory and its incorporation into knowledge management.

The modern time uptake of connectivism as a learning theory is also related to the swift in educational methods observed in progressive educational institutions with special regard to higher education institutions, where e-learning has been adapted in its entirety or in parts in the past two decades, complying with student attitudes, habits and new learning forms. Its impact may be detected in the application of electronic environments impacting an increasing number of students and with continuously diversifying functions. (This topic is discussed in detail in the last chapter of our book, Chapter VIII–The Use of Electronic Learning Environments in Education [12]). These online or Web-based learning support systems establish continuous, synchronous or asynchronous communication between the nodes of the network that is participants (teacher-student, teacher-teacher, student-student).

According to connectivism, learning is a process defined by informal information exchange between nodes organized into a network and supported by electronic tools. Acquiring knowledge is a process where specific nodes are linked to sources of information. Participation in the network and access to software that promote the interpretation of information putting it into a context offers a brand new form

of learning based on cooperation and self-organizing; see Figure 3.

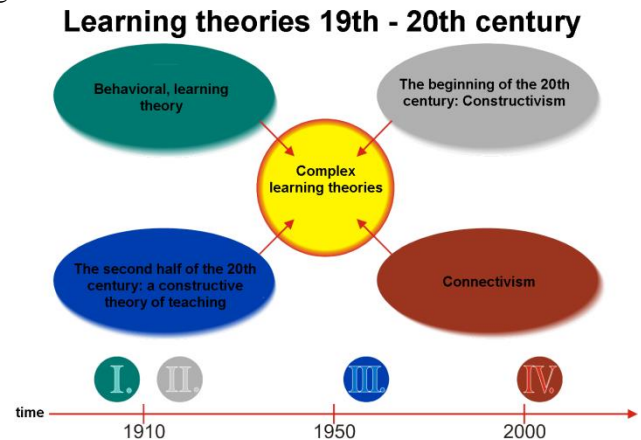


Figure 3. The system and history of complex learning theories.

Though the review in previous sections on educational theories is far from being comprehensive, it illustrates how diverse the topic is. These theories explain actual individual learning activities or the learning processes organized by institutionalized education together, sometimes slightly overlapping. The real processes become even more diversified when developmental stages are interpreted in a broader sense and the institutional framework of formal education is left behind. If we wish to use a simplified systemic approach we should conclude that modern educational theories originate in behaviorism, born more than a century ago, with constructivism catching up and showing an increasing impact with the development of institutionalized teaching and learning until cognitivism became the underlying concept of reforms in the second half of the 20th century. These days, with special regards to the years after the Millennium, the impact of connectivism has become increasingly important, specifically working through Web 2.0 and social activities and essentially defining the development of the new educational paradigm.

E. Forms and Specific Processes of Learning

It is also a useful approach to discuss learning in a less sophisticated way. How learning is incorporated into our lives? The answer claiming learning is the thing we do in school is oversimplified. It is more realistic and closer to the modern approach if we say learning is a basic human activity, which is present at each and every stage of our lives, even if to a varying extent. It is a fact that learning starts with life. The first and most important learning environment is our family. Basic skills (how to move, how to communicate) and behavioural norms are learnt in the initial stage of our lives.

It is important to note that the key to fast and successful learning is learning in a profoundly informal community. Family in childhood and groups of modern times provide such motivating and evaluating environments.

Institutional education (nursery schools, schools) with its formal approach to learning and organized structure differs from informal learning. Though informal learning related to one's interests is also typical in adulthood, either in the framework of one's family or group of friends, employers and other communities also offer learning options that comply with the needs of adults, known as non-formal learning.

- Forms of learning change with life or complement each other in a way that facilitates both direct and indirect learning.
- A well-defined learning strategy (such as being a student in higher education) manifests as a conscious effort made for career and success when selecting a profession or institution.

Even though technology and its social application change asynchronously, developed educational systems realized around the Millennium that the transfer of knowledge that places the individual into the focus of lifelong learning and its success may be significantly affected by educational framework systems. While traditional educational institutions have been mostly focusing on handing over knowledge, modern learning options and lifelong learning concentrate on individual abilities and the development of learning skills. The concept of lifelong learning focuses on enabling and encouraging people „to learn to learn”.

- Atypical learning is not necessarily related to some educational institution neither does it presume a pre-made schedule timetable or exam schedule. Atypical learning does not rely on classes, lectures or seminars in the classical sense; instead learners decide themselves on how to create their own learning environment to their objectives and goals. Atypical learning is a form of lifelong learning. As adults usually work by learning and also need time for their families they cannot always attend institutional (formal) education.
- In addition to the activities of informal learning performed outside formal education (the school system), atypical learning includes all the methods that influence the intellectual and physical development of an individual (distance education, opened learning).

E-learning—is a learning process that combines traditional education distance education and the new options provided by the Internet. E-learning generally, means Internet applications that aim at communicating contents in new ways, administrating the learning process, monitoring individual learning in a way that facilitates evaluation and support and transparent operation at system level. The next figure, though only schematized, shows the varying impact of the 3 learning forms in time along the classic learning curve. We should note that any of the 3 forms may significantly contribute to acquire and update knowledge in any stage of life. As illustrated by the figure, early life is usually dominated by informal learning, followed by formal learning in the period before minimum school leaving time

and while obtaining the highest school degrees. Adults typically learn in non-formal systems; see Figure 4.

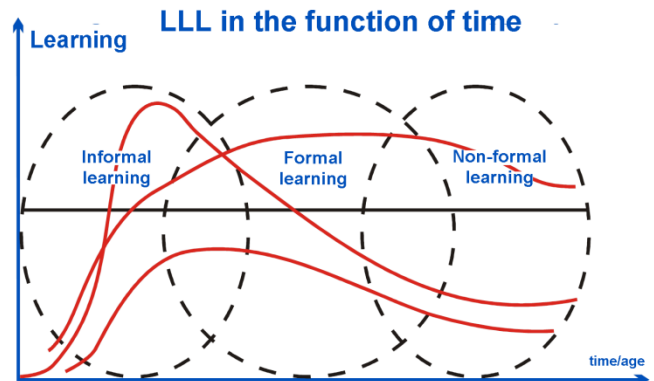


Figure 4. Lifelong learning as a function of time.

In addition to referring the nature of the learning process, we may also use the terms typical learning (learning activities related to formal education) and atypical learning. These do not relate to major teaching or learning forms in the same way. Informal learning typically presumes atypical learning while the terms typical and atypical learning may be used for both formal and non-formal education, depending on the actual activities of students. Terms referring to learning activities have a unique relevance in digital education where ICT use justifies the emphasized role of e-learning and m-learning methods. From the aspect of learning theory or learning forms, these are not more and not less than new, perspective methods that have a special relevance in modernizing education [10].

F. Learning–Time and Space

An important trend in our modern age is rising life expectancy and the increasing ratio of education and training processes within one's lifetime. People of the 21st century are expected to live 7-8 decades spending about one and a half or even two decades in formal education and returning to learning in adulthood increasingly often. In countries developing the learning environment in a broader sense and at the level of the whole society, social stability and coherence have reached a recognized high level while the state of economic development is also above the average (for example, in Northern European, Scandinavian countries). In these countries, the learning trail involves the entire life cycle of the social activities of an individual from the beginning of conscious activities till the end of social activities, all this according to a well-defined social and political strategy. This period may be as long as 60-70 years. Though we think we essentially understand this paradigm, the traditional way institutions operate and are operated is extremely hard to transform into a new educational approach that focuses on the individual. It is the task of the state to create the suitable environment for this, depending on the wide or restricted range of available tools. We should emphasise that in this respect it is not the

magnitude of sources that matter but the approach that considers learning as an important activity in any life stage and is able to support learning by the appropriate tools.

When approaching these issues from social dimensions, we may conclude that a new generation of users committed to online trainings is being born (it may have already very well reached the growing up stage in fact). This generation is much more skilled in navigating in the info-communication space and is becoming increasingly informed and organized. Using this kind of knowledge, people may receive more information and support from each other than from any kind of institution. This is why the role of learning communities has been increasing and it is not too difficult to make projections for the future either. The new communities are essentially characterized by their members sharing the same interests. Learners can interact with each other in these communities, learning together and generating shared sources of knowledge. This newly formulating practice, however, is not incompatible with the learning possibilities offered by higher education structures currently being reformed [11][12].

G. The use of modern, digital teaching materials

One of the main characteristic of the atypical e-learning forms lies in suitable and up-to-date digital teaching materials. Our institute renewed numerous curricula related to a number of university courses to meet needs of e-learning. As a result, in 2012 and 2013 academic years Sharable Content Object Reference Model (SCORM) [1] standards-based e-learning materials were made available to students to support their learning with already familiar Moodle (Modular Object Oriented Electric Learning Environment) [4] e-learning environment. The following screenshot shows its appearance in mentioned learning environments; see Figure 5.



Figure 5. SCORM curriculum in a Moodle System

In autumn 2013, among students using the new curriculum, a satisfaction survey was carried using a digital micro-environment. The survey was carried out after the first semester of curriculum's introduction. The questionnaire was edited by Moodle (Modular Object Oriented Electric

Learning Environment) system tools and the students were asked to send in responses using this system as well [14]; see Figure 6.

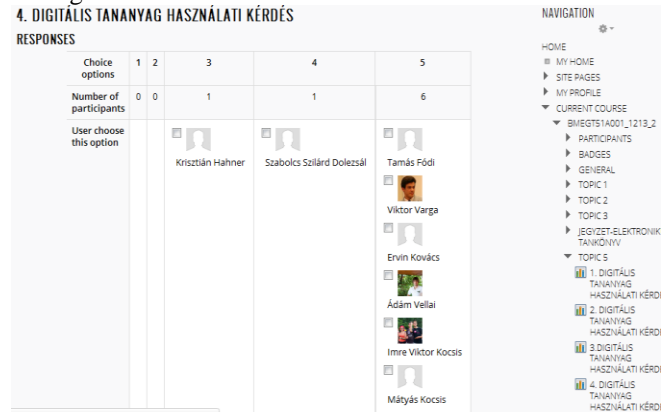


Figure 6. Student replies in Moodle System.

The survey aimed at examining the using efficiency of curriculum presented in new visual form. The survey results are shown below. About 70% of the respondents are fully or greatly satisfied with the use and professional accuracy of new developed curricula; see Figure 7.

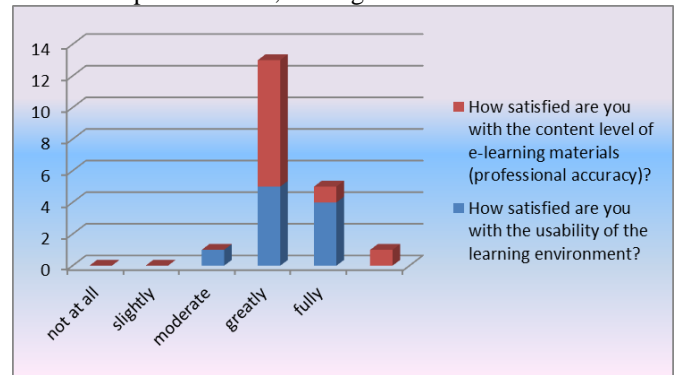


Figure 7. SCORM curriculum in Moodle System

The vast majority of responding students engaged with the curriculum 1-2 hours per day; see Figure 8.

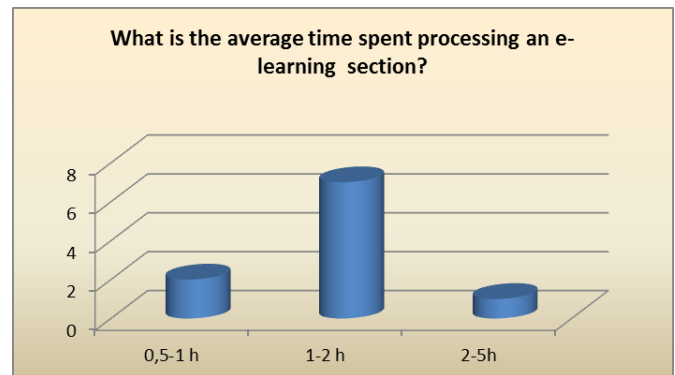


Figure 8. Student responses in Moodle System

Overall, based on the response of the students we found that they are responsive to a new type of digital teaching materials, which the same time motivates them. It is due to the functional design of media objects in teaching materials used.

III. CONCLUSION

In the framework of this paper, we place learning forms into the context of lifelong learning for students reaching the end of a long journey in the world of formal learning. We attempt to investigate why formal learning is losing its dominance and interpret the complementary role of informal and non-formal learning. We also demonstrate the learning forms and educational theories supported by virtual learning frameworks. A number of further key issues have been addressed, such as the possibilities of the rightfully demanded reform in education by means of analysing the impact of visual culture on learning combined with the analysis of the educational effects of demonstration in the context of the history of ideas. The change in approach is closely related to the new ICT application trends that generate changes almost immediately felt in expanding educational space. At the same time, the multimedia environment that can also be efficiently exploited in classroom education has a specific, efficient way to impact our individual worlds. This is what „experiential education” is about, already practiced at BME (Budapest University of Technology and Economics), and the methodology review that attractively describes the world of online language learning.

The opinions above have been of course phrased from the aspect of teachers. Although they address students, a doubt may be raised that we are biased, representing the views of „one side” only. Hence, to promote professional credibility, we close this chapter with the thoughts of a student from not so long ago that may be interesting for both students and teachers.

„... if you like learning, you will be able to get absorbed in almost any topic with the help of Web 2.0, utilizing blogs, Wikipedia and various content sharing sites. You can learn a lot from the experiences and writings of others as they followed the same track back when they were students, they wanted to learn and this way or that they acquired information they wanted. To increase the efficiency of learning, I think providing feedback to the authors is important, like for example sharing our experiences with them. If we are proficient enough in a topic and we feel like it, we can also start a blog or create a site to disseminate relevant items of information we find important so that others can exploit our knowledge. This is indeed the whole point: we should facilitate the flow of information.”

REFERENCES

- [1] P. Kommers, ICT as explicit factor in the evolution of life-long learning. *International journal of continuing*

- engineering education and life-long learning, 20 (1/2010), pp. 127-144., [retrieved: 04, 2014]
- [2] A. Benedek (ed.): *Digital pedagogy 2.0 – Typotex Budapest 2013.*, pp. 18-23
- [3] M. Castells, „The Rise of the Network Society”. *Classics of the Information Society. The Information Age. Economy, Society, Culture. Volume I. Thinking – Infonia, 01.2005*, pp. 489
- [4] <http://index.hu/tech/net/Web1214>, [retrieved: 05, 2014]
- [5] http://ec.europa.eu/education/policies/2010/doc/keyrec_en.pdf, [retrieved: 03, 2014]
- [6] A. Benedek and Gy. Molnár: The empirical analysis of a Web 2.0-based learning platform, In: Constantin Paleologu, Constantinos Mavromoustakis, Marius Minea (ed.): *ICCGI 2011, The Sixth International Multi-Conference on Computing in the Global Information Technology, Luxembourg, June 19-24, 2011.*, ISBN: 978-1-61208-008-6, 06.2011, pp. 56-62
- [7] I. Simonics, “eLearning and Presentation Techniques.” *Óbuda University e-Bulletin, Vol. 1,01.2010*, pp. 211-217.
- [8] Gy. Molnár: Collaborative Technological Applications with Special Focus on ICT based, Networked and Mobile Solutions., *Wseas Transactions on Information Science and Applications 9:09.2012*, pp. 271-281
- [9] Gy. Molnár: Flashes or steady light? Or the potentials of developing networked learning, In: Miguel Baptista Nunes, Maggie McPherson (ed.): *Proceedings of the IADIS International Conference e Learning, IADIS international conference E-learning 2011, Volume II. Rome, Italy, July 20-23, 2011*, ISBN: 978-972-8939-38-0, 07.2011, pp. 405-408
- [10] A. Buda, Attitudes of Teachers Concerning the Use of ICT Equipment In Education, *Journal of social research and policy 1,02.2010*, pp. 131-150.
- [11] Z. Szűts, An Iconic Turn in Art History - The Quest for Realistic and 3D visual Representation on the World Wide Web, In: A. Benedek and K. Nyíri (ed.): *The Iconic Turn in Education (Visual Learning)*, Frankfurt: Peter Lang Internationaler Verlag der Wissenschaften, ISBN:978-3631637715, 2012. pp. 59-66.
- [12] R. Ósz, New technologies mean new methods of learning?, In: Hamido Fujita, Jun Sasaki (ed.): *Proceedings of the 12th International conference on Education and Educational Technology: Recent Advances in Modern Educational Technologies, Iwate, Japan, WSEAS Press, ISBN:978-1-61804-180-07. 2013*, pp. 59-63.
- [13] Z. Szűts, Changes generated and not generated by Printing Press in Europe and Korea, *The Sixth International Conference of KACEEBS, Hungary, Central & Eastern Europe and Korea. Current Issues in Humanities and Social Sciences 17th-19th July, 2006, The Korean Association of Central & Eastern European and Balkan Studies, Korea Literature Translation Institute, Hankuk University of Foreign Studies, Seoul, Korea, 07.2006*, pp. 38-46.
- [14] <http://moodle.appi.bme.hu/> [retrieved: 05, 2014]
- [15] http://hu.wikipedia.org/wiki/Piaget_tanul%C3%A1selm%C3%A9lete [retrieved: 04, 2014]
- [16] Siemens, G., *Connectivism: A learning theory for the digital age. International Journal of Instructional Technology and Distance Learning; 04.2005*, pp. 1-5

Instructional Approach in Adult Education using Mobile Devices

New chances for lifelong learning

Félix Buendía, Angel Perles, Juan-Vte. Capella

Computer Engineering Department

Universitat Politècnica de Valencia,

Valencia, Spain

e-mail: {fbuendia,aperles, jcapella}@disca.upv.es

Abstract—Mobile technologies are currently present in many instructional areas from basic levels to higher education. This work is focused on studying the impact of these technologies in lifelong learning contexts, and more exactly, concerning adult education. To this end, an approach is proposed to organize the different items involved within an instructional process with the purpose of taking advantage from the use of mobile devices. Such approach has been applied in a course addressed to people older than 55 years where the potential of mobile technologies has been assessed to promote a more customized learning adapted to the adult learner profile.

Keywords- Mobile devices; Instructional design; Life-long learning; Adult education

I. INTRODUCTION

Mobile technologies are currently present in many instructional areas from basic levels to higher education. The UNESCO work [1] about Mobile Learning or the Horizon Report [2] recognize the broad extension of tablets, smartphones and similar devices in different learning scenarios. Other initiatives also show the application of these technologies in basic education levels [3], instruction for academic physicians [4] or mobile learning trends for postgraduate students [5]. The current work is focused on studying the impact of these mobile technologies in life-long learning contexts, and more exactly, concerning adult education.

Lifelong learning comprises a wide period of time in which people remain engaged in the process of learning either in formal or informal contexts, by means of training procedures or motivated by personal interests. Therefore, there are many individual situations and circumstances that characterize this kind of learning so instructional and teaching methods have to fit them [6]. Mobile technologies are able to offer special support to these methods and this research work intends to integrate such technologies in a general instructional approach.

In the case of adult education, there are additional aspects to be addressed including those barriers that burden the incorporation of elder people in lifelong learning processes [7]. Such barriers can be removed or lowered by applying technology applications and sound instructional methods.

This paper proposes an instructional approach to guide the application of mobile technologies in adult education contexts. These technologies provide several interface features that allow users to deal with mobile devices in an easier and intuitive way compared to traditional desktop or laptop-based applications. Even, handicapped people or adults with special needs have diverse mechanisms to interact with these devices. This circumstance contributes to a more customized and adaptive learning. Nevertheless, this diversity of interface options provided by mobile technologies makes also difficult to organize or manage their instructional potential. Moreover, the continuous arrival of new mobile devices and those educational contents and tasks based on them is also a factor that complicates the generation of guidelines and procedures to take advantage of such technologies. This situation is particularly critical for adults who are not used to face quick and disruptive changes in their learning process.

The instructional approach presented addresses several issues that range from the technical opportunities offered by mobile devices to the special learning needs that adult users require. All these issues can be framed in a context of *Instructional Design* as the discipline that provides a systematic and rigorous process of managing resources and activities in learning processes [8]. Ozdamli [9] and Park [10] described several pedagogical frameworks in mobile learning contexts to drive technology deployment in an effective way. Sharples et al. [11] proposed a learning theory that fit mobile environments and Elias [12] formulated several instructional principles in the design of such environments. More recently, Dillard [13] collected a set of recommendations to instruct adults by means of mobile technologies. The current approach intends to formalize some of these guidelines and recommendations, and it has been evaluated in the context of a course about *Mobile Technologies* oriented towards adults older than 55. The remainder of the work is organized as follows. Section 2 introduces the instructional approach proposed in this paper and Section 3 describes the case study used to assess it, as well as the evaluation outcomes and their discussion together with some lessons learned. Finally, some conclusions and further works are reported.

II. INSTRUCTIONAL APPROACH

An instructional approach is proposed to organize the different items involved in those learning processes concerning adult education. Such approach deals with technical aspects related to mobile technologies but also takes into account learning requirements for users with specific needs. Therefore, a framework is required that addresses all these technological and instructional issues. In this paper, the ADDIE model [14] provides a sound instructional framework that also has been applied in technological contexts. ADDIE stands for *Analysis, Design, Development, Implementation and Evaluation* and it provides a general model to represent the several stages that are part of instructional processes.

The *Analysis* stage is crucial to gather those learning needs expressed by users. Adult education covers a broad landscape of scenarios ranging from training courses to more informal learning contexts. The current work is focused on such informal scenarios addressed especially to elder people (more than 55 years-old) who are sometimes handicapped users or have special needs. Sometimes, these users also lack knowledge about technological notions though they are usually motivated to learn this kind of topics. From these requirements, the *Design* stage can deal with the formulation of goals that fit such requirements as well as the selection of educational resources and activities adapted for the needs of their users. Figure 1 shows a concept map that displays some of these *Design* components that are part of the proposed approach. For example, procedural goals can be based on achieving psychomotor skills by using keyboard in cellular phones or promoting self-organization abilities by means of scheduling apps. Moreover, resources, such as photo galleries or video recording can be easily accessed in these devices and activities, either individual (e.g., elaborating a personal agenda) or group-based (e.g., exchanging messages in a social network) that can contribute to improve self-esteem and other emotional issues.

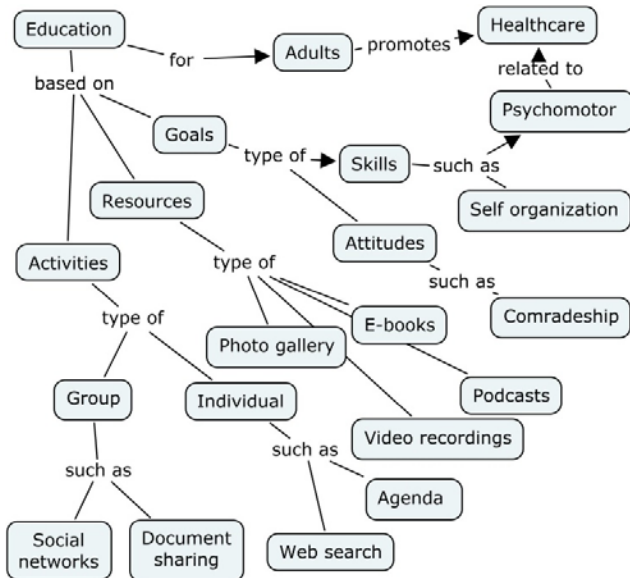


Figure 1.- Instructional process concept map.

The next stage concerns the *Development* of resources and activities, which can be applied in adult education. Several examples have been already mentioned and mobile technologies offer a wide spectrum of tools to develop them. For example, cameras in tablets and smartphones permit to take pictures or record video clips in a simple way. These devices also enable the processing of Quick Response (QR) codes to access different types of contents. Learning activities can be also benefited from the use of mobile technologies through the multiple applications available in distinct platforms and the multimodal interfaces provided by these technologies. Figure 2 shows a concept map that displays some of the technological possibilities offered by mobile devices. During the *Implementation* stage these resources and activities can be delivered by means of several environments and learning systems. In this sense, the most popular Learning Management Systems (LMS) provide adapted versions that fit with mobile features. Moreover, the use of HTML5 and CSS standards are pushing educational Web sites to be adapted to different screen sizes and device configurations applying responsive design principles. This circumstance allows users to easily access their educational contents and services through multiple mobile devices [15].

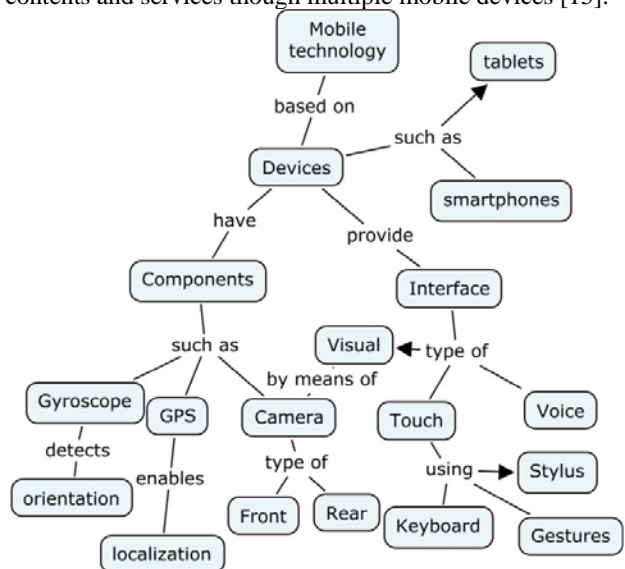


Figure 2.- Mobile technology concept map.

The last stage in the proposed approach deals with the *Evaluation* of the learning processes when mobile technologies are deployed. A first method is used to gather personal opinions by means of questionnaires that are submitted during these processes. This kind of data source is very useful to evaluate the point of view from users who are using mobile devices reporting their qualitative perception about the learning experience. However, a more objective or quantitative way is ever required to get a less “invasive” point of view that provides neutral data from such experience. There are several methods to assess these interactions, such as “eye tracking” techniques, control of user gestures or tracking down the several applications that users run interacting with their mobile devices.

III. CASE STUDY

The current research work has been based on a case of adult education in which the proposed instructional approach has been applied. This case study is integrated in the Senior University program developed at Universitat Politècnica de Valencia (UPV). Next subsections describe the context of this academic program and present some courses about mobile technologies delivered in this context.

A. Academic context

The Universitat Senior at UPV offers more than fifty courses to people older than 55 years. Every year about three thousand students join these courses that address several topics from general cultural issues, such as arts (e.g., music, sculpture or movies) to health and technology. The current case study is part of the courses focused on information and communication technologies and it complements other courses about computer basics, internet or office applications. In general, these courses have a high level of demand and their registration is mostly bounded by technical restraints imposed by those facilities deployed to deliver them. In most of the technological proposals, the number of registered users is about thirty students by course and one special feature is the huge diversity of user profiles that attend these courses. Many of them have completed higher education studies (degrees or masters) though they come from very different disciplines with sparse knowledge in technologies. This aspect sometimes makes teaching these courses difficult but the high motivation and the student’s degree of commitment, generally, compensates that situation.

B. Courses about mobile technologies

There are two types of courses offered to senior students that are addressed to teach mobile technologies:

- *Introductory* courses that covers basic concepts about this kind of technology.
- *Advanced* courses addressed to people with previous experience in these technologies.

The first type of courses allows users to get some starting ideas about the use of mobile phones and tablets and no prior knowledge is required. These users usually own these devices though in most cases they are only used as traditional phones for merely calling people. Table 1 shows a description of the main contents delivered in these courses and their weekly scheduling.

TABLE I. INTRODUCTORY COURSES

Contents	Weeks
Introduction to the world of mobile phones and tablets. Physical handling of the device.	2
Device configuration and customization. Choosing your tablet / smartphone.	2
Communication applications: e-mail, chat, videoconference. Apps instalation. .	3
Multimedia applications: camera use, image gallery, audio and video storage, e-books.	3
Internet: web browsing and search engines.	2
Social networks. Health and travel applications	2

Advanced courses require a previous experience in the management of mobile devices, for example by means of *Introductory* courses or being users of smartphones or tablets. These courses add an extra knowledge about more sophisticated services provided by these technologies, as well as practical activities over them. Table 2 shows a description of course contents and their scheduling by weeks.

TABLE II. ADVANCED COURSES

Contents	Weeks
Advanced settings of mobile phones and tablets. Technical details of mobile fundamentals.	2
Managing several device interfaces. Analyzing multiple tablet / smartphone devices.	2
Mobility applications based on GPS. Agenda and calendar services.	3
Advanced multimedia applications: image edition, audio and video processing.	3
Internet advanced search queries. Apps management.	3
Electronic commerce, shoping using your mobile device. Network security.	3

C. Course resources

The instructional resources and material used in the presented courses are stored on a Web portal implemented by means of a blog tool. Figure 3 shows a screenshot sample of this portal that displays an activity description addressed to teach users how to setup their mobile device. Most of the resources are focused on practical activities that promote the acquisition of different types of skills together with knowledge about mobile technologies.

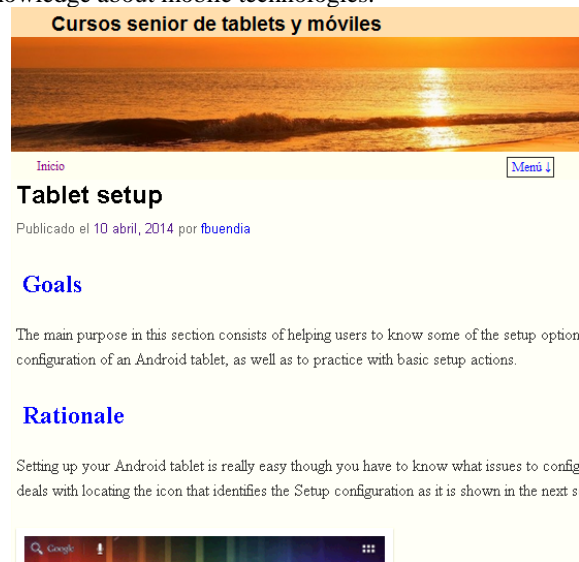


Figure 3.—Web portal for mobile technology courses.

Basic abilities are addressed in *Introductory* courses to enable simple motor skills using the smartphone keyboard or short term memory by means of routine tasks using tablets. Advanced courses encourage other abilities, such as self-organization through calendar activities and systematic search procedures through the Web.

IV. CASE IMPLEMENTATION

The case of study has been implemented during the year 2013 in several courses. In the first term, two *Introductory* courses to mobile technologies were taught to about fifty students. The second term, a new edition of *Introductory* course was delivered, and an *Advanced* course was included to cope with users that joined a previous course and wanted to improve their knowledge about these topics. Both courses were taught by the same lecturers in sessions of two hours by week.

A. Course delivery

The implementation of these courses has been based on a delivery method that mixed teaching in classrooms with learning experiences around the campus. This method allowed instructors to present basic concepts about mobile technologies in a controlled environment, as well as enabling activities in real scenarios outside the classrooms. An example consisted in visiting an outdoor sculpture exhibition that provided information as hot spots accessible from mobile devices. Figure 4 shows a picture that displays a sculpture sample and the QR hot spot used to support this learning activity.



Figure 4.- Sculpture hot spot.

Each course was delivered to about 24 students who attended sessions one day by week. Classroom sessions were taught in a special building prepared to receive multiple Wi-Fi connections. A slide projector was also used to teach basics about mobile technologies and present the different learning activities supported by them. A set of tablets and smartphones was available for course participants though many of them brought their personal devices. This circumstance is further analyzed because the complexity of dealing with a diversity of device configurations and the multiple questions asked by participants about it.

Another critical issue in the course delivery was the format of the documents that reported the technology basics and related activities. These documents were in electronic format (see Figure 3) but the lack of desktop computers in

the classroom led to instructors to provide printed versions of this documentation that helped course participants in their instruction. Activities proposed in the classroom were mostly developed as individual works though some of them, such as sending or receiving messages were performed by pairs or even, in groups through social networks applications. Eventually, these activities were checked in order to detect the personal course progress but avoiding a grading purpose.

B. Course evaluation

The evaluation was mainly focused on analyzing the perception of course participants. Previously, a study of some participants' attributes was performed to achieve a user profile in the several implemented courses. The analysis of user perceptions was organized in several items concerning issues, such as handiness of mobile technologies or their usefulness for course participants. A last aspect consisted in evaluating the user accesses to mobile applications through data logs provided by Android devices used in the courses.

The study of user profiles was carried out using a printed questionnaire fulfilled by participants at the start of the course. This questionnaire was divided in two sections: the first one addressed to get demographic data, such as the user age, gender, or their academic studies. The second one was oriented towards gathering some technical data related to the use of mobile technologies. The next list presents some of the demographic data collected:

- Most of course participant were female (62% in the first term of 2013 and 70% in the second one).
- More than fifty percent of users aged between 63 and 70 years.
- More than sixty percent of users had higher education.
- Three quarters or more were retired people.

The technical study is summarized in Figure 5 that shows some interesting statistics. For example, the percentage of course participants who had internet in their mobile phones was higher than 60% in the first term (2013a) This percentage increased in the next course (2013b) till 85%. In the same way, the percentage of participants who owned smartphone or tablet devices grew until 80% and about 65%, respectively during 2013. Also, the percentage of course participants who daily used their mobile devices in the course 2013b was near the double than the first term.

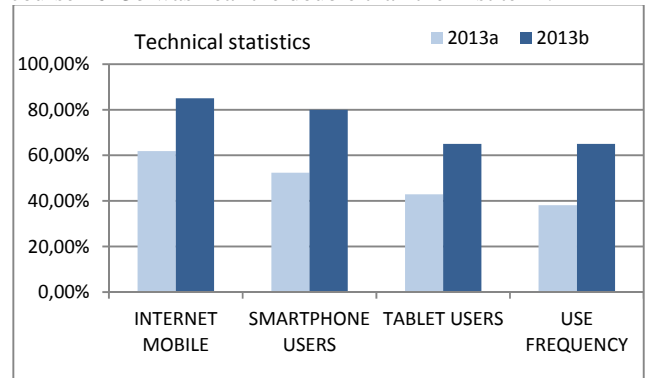


Figure 5.- User technical profile.

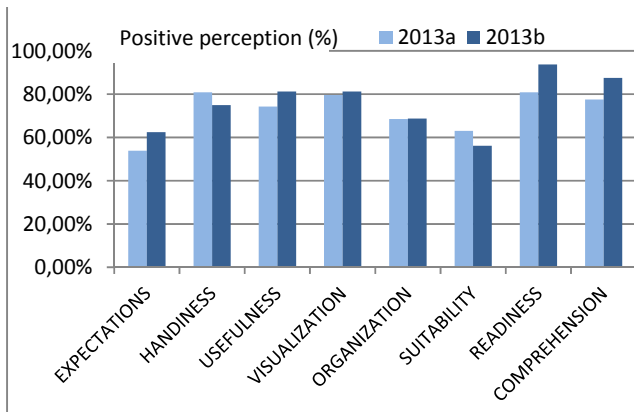


Figure 6.- User technical profile.

Figure 6 shows a bar chart that displays percentage statistics about the user viewpoints of the course. These viewpoints were gathered using an electronic form submitted to users at the end of each term and its answers were based on a Likert scale (from 1 to 5 being 5 the more positive valuation). The graphic view reveals some outcomes, such as the fulfillment of *Expectations* where about half of users stated a positive perception (a value higher than 3 in their answer) in the first term of 2013 and more than 60% in the second term. Higher percentages were drawn in issues, such as the *Handiness* to use mobile devices, their *Usefulness* or the ability to *Visualize* course documents. An aspect that provided a similar percentage (about 70% of positive valuation) was the perception on the use of mobile technologies to improve *Organization* skills.

A second question that was evaluated consisted in the analysis of the log accesses to the mobile devices used during these courses. Such logs were obtained through a *Logcat* tool that provided information about those Android services, which were run in the course activities. Figure 7 shows a sample of timing diagram that displays a sequence of Android events related to an e-mail activity.



Figure 7.- Android service logs.

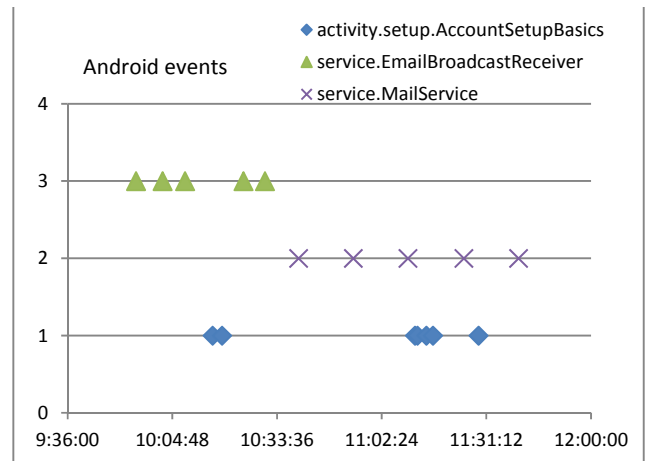


Figure 8.- Android access statistics.

This type of event tracking based on data logs enabled a deeper analysis of user behavior. In the current work, it allowed to measure the activity level of participants in the course when they were using mobile technologies. Such activity could be checked in order to test if the planned scheduling was met or get a performance index for each user. Figure 8 shows a chart that displays this kind of activity registering the timing of the Android events related to sending an e-mail in several devices. In this case, there is sparse distribution of events configuring the e-mail activity (*AccountSetupBasics*) regarding those events of type *MailService*.

C. Lessons learned

The evaluation of the case of study previously introduced has provided significant outcomes and lessons learned concerning the use of mobile technologies in specific scenarios of adult education. A first lesson is the need of a sound framework when learning processes are implemented within adult educational scenarios. This kind of framework is crucial when dealing with users who have different level of expertise and singular educational goals. The proposed framework has been based on a well-known instructional approach that fits quite well with the use of technologies in this type of scenarios. A key issue in such approach is the *Analysis* of the user learning requirements and the conditions that these users request. Mobile devices are relatively new technologies, especially for elder people, and it is important to gather their perspective in the knowledge of such technologies. Afterwards, it would be more suitable to adapt the *Design* of learning resources and teaching activities to these requirements.

This rigorous and systematic approach has also enabled the appropriate *Development* of educational media that fit technical features of mobile devices: The diversity of size and interface features makes difficult to develop instructional materials to be deployed in several devices (e.g., a tablet vs. a smartphone). Therefore, it is important to collect these media and label them to get a catalog of available material adapted to the required learning scenarios. That lesson led to elaborate a database of potential mobile devices that could be

used for adult education and the suitable resources for these devices. An additional lesson is focused on the mechanisms that can be deployed to deliver and supervise the use of mobile applications. There are several mobile platforms in the market and even within a specific platform like Android, multiple versions are available providing a different or customized display of applications. Moreover, users who brought their personal devices asked for tools to ease the access to these mobile applications.

Finally, there is a general agreement about the need to find new *Evaluation* tools that enable a deeper study of the user behavior when they interact with mobile technologies. Questionnaires provide a qualitative perception of the users' point of view but additional measures are required. In the current work, a *Logcat* tool was applied to track Android events in mobile devices though it was complex to manage the huge amount of collected data. Moreover, this type of tool is bounded to older Android versions. Alternative methods, such as eye tracking techniques and touch detection mechanisms could be explored to obtain new data sources in the context of adult education learning scenarios.

V. CONCLUSIONS

The current work has presented an instructional approach that permits to organize learning processes in a context of deploying mobile technologies. These technologies have been used to implement several courses during the year 2013 especially addressed to elder people. The proposed approach has shown its usefulness when implementing such courses since it provides a systematic way to analyze, design and develop learning resources and activities based on the use of mobile devices. A case study has been evaluated concerning two *Introductory* courses about mobile technologies and assessing their impact in the way elder people interact with such technologies. Several questionnaires have been delivered to gather profile information about users who joined these courses. At the end of each course, an additional questionnaire has been submitted to obtain the viewpoints of course participants. Also, data logs extracted from mobile devices deployed during these courses have been analyzed to check timing patterns in the use of mobile apps. The global analysis of the evaluation outcomes reveals a high interest of elder people about learning through this kind of technology although the data logs show the users' problems to follow planned activities.

ACKNOWLEDGMENT

This work is supported by the TEA project (PAIDUPV/2791) and AulaSenior (aulasenor.es).

REFERENCES

- [1] UNESCO, "Mobile Learning Week Report", 2011 [Online] Retrieved March, 2014: <http://www.unesco.org/new/fileadmin/MULTIMEDIA/HQ/ED/ICT/pdf/UNESCO%20MLW%20report%20final%2019jan.pdf>
- [2] NMC "Horizon Report 2012" Higher Education Edition [Online] Retrieved March, 2014: <http://www.nmc.org/publications/horizon-report-2012-higher-ed-edition>
- [3] M. Bjerede and T. Bondi, "Learning is personal: Stories of Android Tablet Use in the 5th Grade" [Online] Retrieved March, 2014: <http://www.learninguntethered.com/>
- [4] J. Sclafani, T. F., Tirrell and O. I. Franko, "Mobile Tablet Use among Academic Physicians and Trainees," *Journal of Medical Systems*, 37 (1), 2013.
- [5] D. M. Ruth, T. M. Mastre and R. Fricker, "A Study of Mobile Learning Trends at the U.S. Naval Academy and the Naval Postgraduate School," *Educase Review*, 2013 [Online] Retrieved March, 2014: <http://www.educause.edu/ero/article/study-mobile-learning-trends-us-naval-academy-and-naval-postgraduate-school>
- [6] J. C. Dunlap and R. S. Grabinger, "Preparing students for lifelong learning: A review of instructional features and teaching methodologies," *Performance Improvement Quarterly*, 16(2), pp. 6-25, 2003.
- [7] A. Perry, C. Shepherd, D. Moore and S. Schmolle, "Scaling up: Achieving a breakthrough in adult learning with technology," Ufi Charitable Trust Report, 2013 [Online] Retrieved March, 2014: http://www.ufi.co.uk/sites/default/files/Scaling%20up_21_5_V3.pdf
- [8] W. Dick and L. Carey, *The systematic design of instruction*. 4th ed. New York, NY: Harper Collin, 1996.
- [9] F. Ozdamli, "Pedagogical framework of m-learning," *Procedia - Social and Behavioral Sciences*, 31(0), pp. 927-931, 2011.
- [10] Y. Park, "A pedagogical framework for mobile learning: Categorizing educational applications of mobile technologies into four types," *The International Review of Research in Open and Distance Learning*, 12(2), pp. 78-102, 2011.
- [11] M. Sharples, J. Taylor and G. Vavoula, "A Theory of Learning for the Mobile Age" In B. Bachmair (ed.) *Medienbildung in neuen Kulturräumen*. Stuttgart, Kohlhammer Verlag, pp. 87-99, 2010.
- [12] T. Elias, "Universal Instructional Design Principles for Mobile Learning," *International Review of Research in Open and Distance Learning*, 12 (2), pp. 143-156, 2011.
- [13] A. Dillard, "Mobile Instructional Design Principles for Adult Learner" Capstone Report. University of Oregon, 2012 [Online] Retrieved March, 2014: <https://scholarsbank.uoregon.edu/xmlui/handle/1794/12253>
- [14] M. Molenda, In search of the elusive ADDIE model. *Performance improvement*, 42(5), 34, 2003.
- [15] K. McCrane, *Content Strategy for Mobile*. Ed. A Book apart, 2012.

WebETL Tool – A Prototype in Action

Kornelije Rabuzin, Matija Novak
 Faculty of Organization and Informatics
 University of Zagreb
 Varazdin, Croatia
 {kornelije.rabuzin, matija.novak}@foi.hr

Abstract – Every grocery store has data about bills, bought items, etc., usually stored in a database (DB). There is often a need to analyze this data. This is not suitable to do on an operational DB, especially when data from two or more stores must be analyzed. Therefore, data from the operational DB is extracted, transformed, and loaded (ETL) into a data warehouse (DW). The ETL process is the most important part of building a data warehouse. It is used to extract data from operational data sources, transform the data (as needed) and to load the data into a data warehouse. Because the ETL process (as such) is very complex and time consuming, a prototype of a web ETL tool was built and tested. The results have shown that it is possible to build a completely web-based ETL tool and that it is faster than manual ETL. Also, because it is completely web based, multiple users can use the tool at the same time, with no installation is needed.

Keywords – ETL; data warehouse; web; ETL tool

I. INTRODUCTION

Data warehouses are used to store data in a so-called star schema (Fig. 1) that consists of dimensional and fact tables. This model is understandable to the end user, thus making it easier to perform data analyses. [1, p.85]

Although many different definitions can be found; according to Kimball & Caserta, data warehouse is defined as follows: “A data warehouse is a system that extracts, cleans, conforms, and delivers source data into a dimensional data store and then supports and implements querying and analysis for the purpose of decision making.” [2, p. 23]

The most important part of a DW system is the ETL process, sometimes called extract, clean, conform, and delivery (ECCD) process, as described by Kimball & Caserta in [2, pp. 18-19]. During the construction of a DW, by Inmon [3, p. 295], 80 percent of the time and resources are consumed by the ETL process.

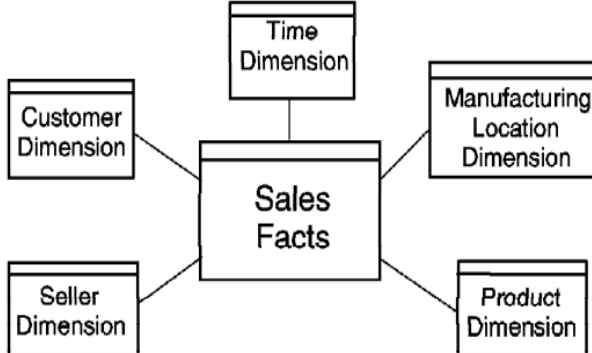


Figure 1. Star model [4]

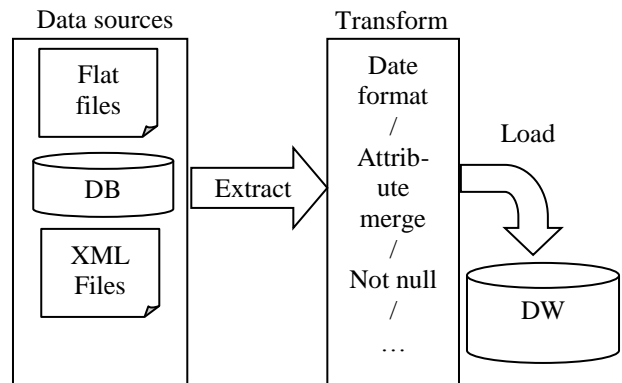


Figure 2. ETL process steps

The ETL process has three important steps (Fig. 2):

- **Extract** - The first step, extracts data from operational sources such as flat files, operational databases, extensible markup language (XML) files, Enterprise Resource Planning (ERP) systems, etc.
- **Transform** (clean and conform in ECCD) – The second step, transforms data (e.g., date format, capitalize names, merge attributes, etc.) and the main purpose of this step is to enhance the quality of the original data. In addition, this step resolves conflicts (e.g., duplicate values) if more than one data source is used.
- **Load** (delivery in ECCD) – The third step, loads transformed data into a data warehouse.

To ease the ETL process, tools are built to support the ETL process. Existing tools are desktop tools, which means they require time for installation and configuration, space to store data created during usage of the tool. In addition, these tools are very complex so user must be very familiar with ETL process and learn how to use the tool. All that requires some time before the tool is efficiently used.

To support the ETL process, and remove the above-mentioned problems, a completely web-based ETL tool was built. “This ETL tool is designed as a web application so that users can save time (and space) required for installation purposes.”[5] To help users during the ETL process, step-by-step guidance throughout the entire process was built into the tool making the whole process much easier, even for less experienced users. The explanation on how to configure and use the tool is described in another article [5] that has been submitted for publication. This article is focused on testing the prototype on a real case scenario.

This paper is structured as follows: In Section 2, the ETL tool basics are described. Section 3 describes the model (a real data warehouse example that we were working on). In Section 4, the test results are presented. Section 5 concludes and presents some open questions, listed for future research.

II. ETL TOOL DESCRIPTION

Since many professional ETL tools exist on the market, it is reasonable to ask why should another tool be used and implemented. The main difference between existing tools and our tool is that our ETL tool is completely web-based.

The main advantages of using a web-based tool are:

- No installation needed.
- More users can use it at the same time.
- Lower maintenance costs, because users do not need hardware like in traditional client ETL application, only a browser is needed.

In the article [6], authors introduce web-based the ETL framework and describe the benefits of using a web-based ETL tool, such as lower maintenance costs.

The main difference between the tool presented in [6], is that our tool implements the ETL process in such a way that users can learn the process as they go (the tool guides a user through the steps (Fig. 3) that have to be carried out in the ETL process). Steps are visible throughout the entire process, so the user at every moment knows where he is. During the process, only the right side of the tool changes, as presented in other snapshots (Figs. 4-6 and Fig. 8 and Fig. 9). Therefore, the main two advantages are as follows:

- User can be a beginner in ETL and still be able to use the tool correctly and learn the ETL process while using it.
- Guidance while using the tool through steps (Fig. 3).

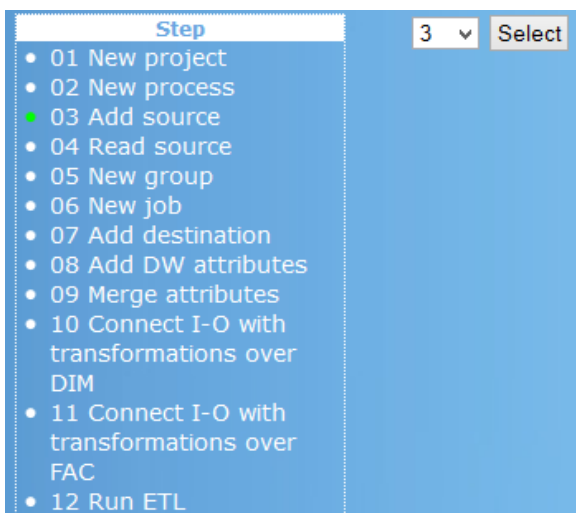


Figure 3. Menu of checkpoint (steps) for the user (left) and form to select number of sources (right)

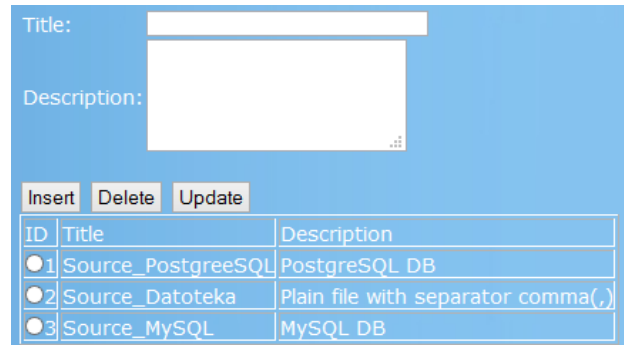


Figure 4. Administration view of source types

The tool is implemented in Java and, for the interface, Hypertext Markup Language (HTML) and JavaServer pages (JSP) sites are used. In order to use the tool, a user has to enter all required metadata. A thread will run and perform extraction, transformation and load of data into a data warehouse. The tool has just one main thread that does the management of three parts: extraction, transform and load. A detailed explanation of the tool is given in [5].

Currently, the tool supports three types of data sources: *Flat source files*, *PostgreSQL* and *MySQL*. However, the tool can be easily expanded to other data sources. The sources that will be used can be located anywhere on the web. The only requirement is an open connection, so that the tool can access the source over Internet.

The tool supports extraction from multiple sources into the data warehouse while removing the duplicates based on the unique key (key can consist of one or more attributes) defined for each dimension.

At this point, four types of transformations are implemented; one can:

- Add default value if the source value is empty.
- Merge two (or more) attributes (e.g., name and surname into "name surname").
- Change date format (change the date and/or time format).
- Change all letters to upper or lower case.

To support the flexible ETL tool, data source and transformation parts were made with dynamic loading. It is possible to add new classes for new data source types as well as new transformations at any time without editing the source code (in order to add new classes, those classes must implement a specific interface). Information about these new classes must be entered on the administration site (Fig. 4). Title in Fig. 4 is the name of the class that implements the functionality for one data source.

As already mentioned, the whole process is guided (step-by-step). Through the guided interface (Fig. 3) the user is asked to enter the information about data sources (Fig. 5, e.g., Internet protocol (IP) address, username, password, DB name). Once all data sources are defined (Fig. 5), the tool reads metadata info about tables and attributes in defined sources (for example, in MySQL, the "information_schema" table). Then, one has to define dimensional and fact tables (with their attributes). Next, the user must match (Fig. 6) attributes from the data source with attributes in the destination, or fact, tables, and define transformations that must be made.

Figure 5. Form for entering new source

This transformation metadata vital and it is stored in the logical data map.

The logical data map is stored as a table that has three main parts [2, pp. 56-71]:

- Destination – described by table name, column name, data type, table type and slow changing dimension (SCD) type
- Source – described by database name, table name, column name and data type
- Transformation – information about what transformations are performed on a specific attribute

Based on all of the metadata, the ETL process can start. One by one, dimensions are processed and for each dimension, data sources are read. When all dimensions are processed, fact tables are processed and connected to corresponding dimensions as defined in metadata.

Figure 6. Attribute matching and transformation define

Because data in the data warehouse is stored in relational tables, this mechanism is called relational online analytical processing (ROLAP). There is another mechanism called multidimensional online analytical processing (MOLAP). More about these two mechanisms can be found in [7].

The next Section will show how was test done and what the results are.

III. TEST AND RESULTS

For testing purposes, we used data from a small data warehouse that we implemented a few years ago (more precisely, data from two grocery stores). Every store had its own operational database and data from those two systems were integrated in order to analyze sales results. This project was implemented in a Microsoft (MS) Access (database management system) and Business Objects XI (business intelligence tool). For testing purposes, we decided to implement the web based ETL tool.

In order to analyze sales results, we had to analyze databases that were used to store data. In each store, the following tables were important for us:

- *Articles* - contains information about products (article name, retail price, purchase price, code, units of measure, etc.),
- *Salesman* – this table contained data on sales people (name, surname, id ...),
- *Bills* – contains data about bills (date and time when the bill was made, id of salesman that made it, bill number, payment method, etc.), and
- *Bought items* - contains data about sold items (bill number, article name (connection to article table, code was not used), quantity, retail price, unit of measure, amount, purchase price, tax ...).

Because the original source was a very old MS disk operating system (DOS) application, MS access was the easiest way to extract data. Although the original data sources had a little less than 500 000 records, for the purpose of this research, we took about 10% of data in this initial testing stage.

```

INSERT INTO sales (bill_number, date_time, article, salesman, item_amount, quantity, profit)

SELECT bills1.bill_number, dates1.id, articles1.code, salesmans1.id, bill_items1.amount, bill_items1.quantity, IIf(cint([bill_items1.tax])=0, [bill_items1].[amount] -[bill_items1].[purchase_price], ([bill_items1].[amount]-1.22*[bill_items1].[purchase_price])/1.22) AS Expr1

FROM salesman INNER JOIN (salesmans1 INNER JOIN (articles1 INNER JOIN (article_helpers1 INNER JOIN ((bills1 INNER JOIN bill_items1 ON bills1.bill_number = bill_items1.bill_nubmer) INNER JOIN dates1 ON bills1.date = dates1.date) ON article_helpers1.article = bill_items1.article) ON articles1.id = article_helpers1.id) ON salesmans1.id = bills1.salesman_id;
    
```

Figure 7. Example 1

Figure 8. Dimension crate

Figure 9. Destination attribute create

Because these two stores were at different locations, they both had information about articles; the problem was that some articles were equally tagged (the code was the same), but the article name was different. Further, some articles were only present in one store and some were present in another store. When we implemented the ETL manually, we merged these two tables and we created a new dimension table with a new id value. Of course, we had to maintain a connection to bills and items sold in order to load the data into a fact table. A query that was used to load data into a fact table is shown below (Fig. 7).

Of course, many other things had to be resolved as well; this was just an example of some transformations that we had to carry out.

In order to test the ETL tool, some data (from two grocery stores) were placed in a MySQL 5.6.11 database and some data was placed in PostgreSQL 9.3 database. For DW, ROLAP model is used; the dimension and fact tables are stored in PostgreSQL 9.3 in star schema format.

As described in previous Section, information of these two sources (Fig. 5) was entered in the ETL tool. We performed the ETL graphically, so no structured query language (SQL) queries needed to be written; the tool generated all queries automatically.

In the first part, we defined data sources (Fig. 5), dimension tables (Fig. 8), dimension table attributes (Fig. 9), fact table, and fact table attributes:

- Dimension Articles (id, article name, retail price, purchase price, code, units of measure, etc.)
- Dimension Salesman (id, name and surname, etc.)
- Dimension Date (id, date, time, date_time, etc.)
- Fact Sales (id, bill number, quantity, dim articles id, dim salesman id, dim date id)

Figure 10. Attribute “name”, matching and transformation define from one source

```
SELECT bill_items1.quantity as quantity,
bill_items1.ammount as amount, bills1.bill_number as
bill_number, bills1.date as date, bills1.time as time,
articles1.code as code, salesmans1.salseman as sales-
man from bill_items1, bills1, articles1, salesmans1
where bills1.bill_nubmer = bill_itmes1.bill_number
and articles1.article = bill_itmes1.article and
bills1.salesman_id = salesmans1.salesman_id
```

Figure 11. Example 2

Then, we defined transformations for attributes, example (Fig. 10) for attribute “name” in dimension table “DIM_articles” read from PostgreSQL source table “Article” it was defined transformation “Fill_not_null” and default value was added “no name”. In this way, all attributes were specified from both sources MySQL and PostgreSQL.

Another example, as shown in Fig. 6, is for one dimension (Salesman) and one source (PostgreSQL). To have the “Name_Surname” field in salesman dimension, the merge transformation was used, and for date dimension attributes the “date_format” transformation.

When all metadata was defined, we started the process. First, based on entered data, extraction was done for every defined source (in our case two sources MySQL and PostgreSQL). Second, data was transformed based on chosen transformations; and finally it was loaded into dimensional tables. When all dimension tables were loaded, fact table must be loaded as well.

TABLE I. SOURCE DATA TABLES AND ATTRIBUTES

Table	Attribute	Value
Articles	Article name	White Bread
	Retail price	8,80
	Unit of measure	Item
	Purchase price	6,92
Salesman	Id	38503177
	Name	John
	Surname	Doe
Bils	Bill Number	33703
	Date	2008-07-07
	Time	19:39
	Amount	8,80
	Salesman Id	01
	Payment method	Cache
Bought items	Article Name	White Bread
	Quantity	1
	Retail price	8,80
	Unit of measure	Item
	Amount	8,80
	Purchase price	6,92
Tax	22	
Bill number	33703	

```

INSERT INTO FAC_sales
(quantiy, amount, bill_number, id_DIM_date,
id_DIM_Articles, id_DIM_salesmans )
VALUES(1.0,8.80,33703.0,4653,473,3);
    
```

Figure 12. Example 3

TABLE II. DIMENSIONAL AND FACT TABLES AND ATTRIBUTES

Table	Attribute	Value
DIM_Date	ID	4653
	Date	2008-07-07
	Time	19:39:00
	Date Time	2008-07-07 19:39:00
DIM_Users	ID	3
	Name	John
	Surname	Doe
	Name Surname	John Doe
	Salesman Id	01
DIM_Articles	ID	473
	Name	White Bread
	Retail price	8,80
	Unit of measure	Item
	Purchase price	6,92
	Article Id	38503177
FAC_Sales	ID	18874
	Id DIM Users	3
	Id DIM Articles	473
	Id DIM Date	4653
	Quantity	1
	Amount	8,8
	Bill number	33703

Example 2 (Fig. 11) shows the generated query that extracted data for fact table. The query presented in Example 2 (Fig. 11), was used to extract the data and to match dimensions (by finding which dimension id should be used for each row) with the fact table. The final query that inserted data into the fact table is shown in Example 3 (Fig. 12).

To give better insight into the data, an example of one row in each source table (Table 1) that is then transformed into one row into dimensional and fact tables (Table 2) is shown.

IV. RESULTS AND FUTURE WORK

The main benefit from using the tool was that it was much faster and easier than before, when we were writing all queries manually. The end-user does not need to write SQL queries manually. Instead, a step-by-step guide ensures that the whole process is easy to follow, so unexperienced users can easily use the tool. With our test we have proven that a completely web based tool is functional and that it has some benefits compared to manual ETL. The main benefit is that the tool is intuitive and easier to use than manual ETL. With these benefits, it is possible to use the tool in education while students learn the ETL process. Because no installation is required and more users can use the tool at the same time, this tool is more convenient than traditional ETL tools

However, although we were able to extract data, there were some limitations:

- No aggregate functions were implemented so the profit value could not be calculated and stored into the data warehouse;
- The source database must use primary keys so that the ETL tool can automatically recognize connec-

tions between tables. Further on, it is expected that the name of the attribute by which two tables are connected is the same in both tables;

- As a result, we can say that PostgreSQL can quickly process queries shown earlier while MySQL is much slower. The solution could be implemented in a form that MySQL extracts and processes a smaller number of records (for example, 10000 records at once);
- Since data contained some Croatian characters, there was no way to change these characters automatically.

We can say that the tool has limitations and there is definitely space for improvement (we could add aggregate functions, we could speed up the process of extraction and load, parallel extraction from multiple sources is also something that one can consider, etc.).

In the future, detailed tests with more concurrent users are needed. In addition, the idea is to expand and upgrade the performances of this tool to be competitive (for example, in speed) with some traditional tools. In addition, it is planned to completely automate the entire process with ontologies.

V. CONCLUSION

In comparison to traditional ETL tools, the presented ETL tool makes it easier to extract data. It also helps during the ETL process by showing hints of what to do next. But still, the knowledge about own data sources is needed to know what is in which table and user must know what he wants to get at the end.

Since the tool is still a prototype, limitations are present; but, at this point in time, it is possible to create a small data warehouse. Flexible implementation ensures that new features can be easily added.

The main benefit is that it is completely web based, allowing multiple users to use it at the same time. Further on, with integrated step-by-step guidance, even users not familiar with the ETL process can use the tool and learn along the way. Finally, no installation is required.

REFERENCES

- [1] K. Rabuzin and M. Novak, "Data warehouses and ETL," Methods and Tools for Information and Business Systems development (Case22), Zagreb, Jun. 2010, pp. 85-89
- [2] R. Kimball and J. Caserta, The Data Warehouse ETL Toolkit: Practical Techniques for Extracting, Cleaning, Conforming, and Delivering Data, Indianapolis: Wiley Publishing Inc., 2004.
- [3] H. W. Inmon, Building the Data Warehouse – Third Edition, New York: John Wiley & Sons Inc., 2002.
- [4] C. Ballard, D. Herreman, D. Schau, R. Bell, E. Kim, and A. Valencic, Data Modeling Techniques for Data Warehousing - IBM redbook, IBM Corporation: International Technical Support Organization, 1998.
- [5] M. Novak and K. Rabuzin, "Prototype of a web ETL Tool," International conference on data warehousing and knowledge discovery, Unpublished
- [6] R. K. Vangipuram, V. Sreekanth, and B. Rangaswamy, "Implementation of web-ETL transformation with pre-configured multi-source system connection and transformation mapping statistics report," International Concurrence on Advanced Computer Theory and Engineering (ICACTE'10), IEEE Press, vol. 2, Aug. 2010, pp. 317-322, doi:10.1109/ICACTE.2010.5579100
- [7] P. Ponniah, Data Warehousing Fundamentals: A Comprehensive Guide for IT Professionals, New York: John Wiley & Sons Inc., 2001

Deductive Data Warehouses and Aggregate (Derived) Tables

Kornelije Rabuzin, Mirko Malekovic, Mirko Cubrilo
 Faculty of Organization and Informatics
 University of Zagreb
 Varazdin, Croatia
 {kornelije.rabuzin, mirko.malekovic, mirko.cubrilo}@foi.hr

Abstract - In one of our previous papers, the idea of deductive data warehouses has been introduced. It was shown how to use Datalog rules to perform Online Analytical Processing (OLAP) analysis on data. In this paper, we show how to use Datalog rules to specify the data warehouse model (data mart) as well as how to add rules that produce aggregate and derived tables that are normally used to speed up the process of retrieving data. Since it is good to have aggregate and derived tables (to speed up queries), the main drawback is that they require extra storage. Consequently implicit definition of such tables may seem interesting.

Keywords: data warehouse; deductive data warehouse; Datalog; aggregate tables; derived tables; data mart.

I. INTRODUCTION

Data warehouses are popular due to the fact they can efficiently store large amounts of data and data can be analyzed by means of front-end business intelligence tools (Business Objects, QlikView, etc.) that support different ways of analysis including drill down, roll up, slice, dice, etc. One can find many different definitions of what a data warehouse is. According to Kimball et al. [2] it is "a system that extracts, cleans, conforms, and delivers source data into a dimensional data store and then supports and implements querying and analysis for the purpose of decision making."

Over the years many companies implemented many (partial) applications and/or information systems that covered only one aspect of business. This was a common scenario in the past any many companies have similar problems today because of such an approach. The main problem is that these companies possess many heterogeneous applications and information systems that were built by means of incompatible technologies (different programming languages were used as well as different ways of storing data) and it becomes hard (expensive) and sometime almost impossible to integrate data from all sources. Now, one may ask why data integration is important? The answer is obvious, especially today, when we know that data must be integrated and compared in order to make good decisions and in order to understand what is really going on.

When building a data warehouse many different steps have to be carried out, but one process that is quite crucial is the so-called Extract Transform Load (ETL) process. The goal of the ETL process is to extract data (from different (non)relational sources), to transform data (data has to be cleaned, business rules have to be obeyed, missing values have to be found, some attributes have to be merged, some values have to be split, different formats have to be unified, etc.) and to load data into the data warehouse. A

few good books [3][4][5][7] have been written on data warehouses and ETL; here we just want to emphasize the importance of the ETL process but we will skip the details. All information technology (IT) people, especially those who were working on data warehousing projects, know what end users are capable of doing and what mistakes and bad data transactional sources do often contain.

Since the ETL process is very important, we tried to develop a web ETL tool that had an educational component as well. The main idea was to build a tool that could help users to build a data warehouse and to guide them during the ETL process. More details about the tool can be found in [6]; this is a paper that is under a review at this point in time.

When we talk about the data warehouse data model, we distinguish two types of tables (Fig. 1):

- Fact tables contain facts, i.e., numbers that are used to quantify business processes (number of sold items, number of items in stock, etc.). They contain much larger number of records and small number of attributes.
- Dimension tables (dimensions) contain much larger number of attributes that are used to analyze data in different ways and much smaller number of records (when compared to fact tables).

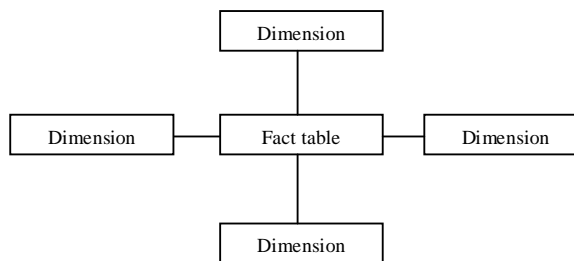


Figure 1. Star schema

When we discuss different data loading techniques, one should keep in mind that data loading is different and it depends on whether a data warehouse already contains some data or not. We have to take care about dimension and fact tables and we have to achieve some level of parallelism, if possible. Because of that we distinguish:

- initial load (a data warehouse is empty and all data have to be inserted),
- incremental load (we change existing records and add new ones) and
- complete reload, i.e., refresh (some or all data are deleted and re-loaded again).

Deductive databases are databases that use deductive rules to produce new information. They usually contain a set of facts (they would correspond to rows in a table), a set of rules (rules produce new piece of information) and certain integrity constraints that have to be satisfied. One of the main advantages of deductive databases (in the past) was the ability to specify recursive rules. Furthermore, it was possible to view in recursive queries as well. In recent years, recursive queries have been implemented in different database management systems as well. Here, we list a few courses and their prerequisites (Fig. 2):

```

prereq(math, databases).
prereq(math, statistics).
prereq(databases, programming).
comes before(X,Y):-prereq(X,Y).
comes before(X,Y):-prereq(X,Z),comes before(Z,Y).
    
```

Figure 2. Deductive database – an example

In [1], the idea of deductive data warehouses has been proposed. Things are similar to deductive databases but some important distinctions exist and they are explained in the paper. The paper also shows how to perform OLAP analysis on data, i.e., how to use Datalog rules to analyze data. In [8], the idea was extended in order to show how some other types of analysis could be implemented by means of Datalog rules as well.

In this paper, we show how Datalog rules can be used in order to reduce the size of a data warehouse. More precisely, we show how to implement implicit aggregate and derived tables. Although they are not physically implemented as such, their existence is sometimes effective to perform reasoning on data. Furthermore, we show that view materialization can significantly improve performances.

There are several different papers that explore the use of Datalog and its role in data warehousing. Boulicaut et al. [10], use rules in a similar way, but they focus on knowledge discovery. Neumayr et al. [11] use Datalog to reason over multidimensional ontologies. Aligon et al. [12] explore how to summarize and query logs of OLAP queries. However the term deductive data warehouses is new and papers on the topic cannot be found.

This paper is structured as follows: first the deductive data warehouse is described. Then we say a few words about aggregate tables, their role and ways how to implement them. In the next section, we say a few words about derived tables. In both sections we show how deductive rules can be used to specify derived and aggregate tables. Then we show some experimental results. Finally, the conclusion is presented.

II. DEDUCTIVE DATA WAREHOUSES

In this section, we say a few words regarding deductive data warehouses. Deductive data warehouses are quite similar to deductive databases, i.e., they both contain facts, rules and integrity constraints. When talking about deductive data warehouses, integrity constraints are not so important as in deductive databases because ETL designers are responsible that rules are obeyed and quality of data is ensured.

In [1], it was shown how Datalog rules can be used to simulate and perform data analysis including slice, drill down, what if, etc. The term used to describe the model was new. The idea seems to be important in the same way that deductive databases are important for data analysis in regular databases. Some rules (from [1]) are given below (Fig. 3). The first rule was used to find users. The second rule was used to perform drill down/roll up analysis on data. The third rule was used to perform what-if analysis on data. In [8], the idea was extended and some more complex rules were added to perform the Recency Frequency Monetary (RFM) analysis.

```

user(X):-users(_,X).

drill_down(X,Y,H,R):- X=year,
group by((dates(B,F,Y,H,I,J), log(A,B,C,D)), [H],
R=sum(D)).

what_if(B,H,R,Z,W):-user_points(B,H,R,I,P),
A is I + Z, W is A * R.

compare(X,R):-addrange(Minv,Q1,Q2,Q3,Q4,Maxv),
((X >= Minv, X < Q1, R is 1);
(X >= Q1, X < Q2, R is 2);
(X >= Q2, X < Q3, R is 3);
(X >= Q3, X < Q4, R is 4);
(X >= Q4, R is 5)).
    
```

Figure 3. Datalog rules

In [1], a small example was used to show how deductive rules can be used to perform OLAP analysis on data. The scenario model was quite simple; users performed some actions on certain dates and we needed to analyze the number of actions (this is a part of a real project). For that purpose several tables were defined (Fig. 4):

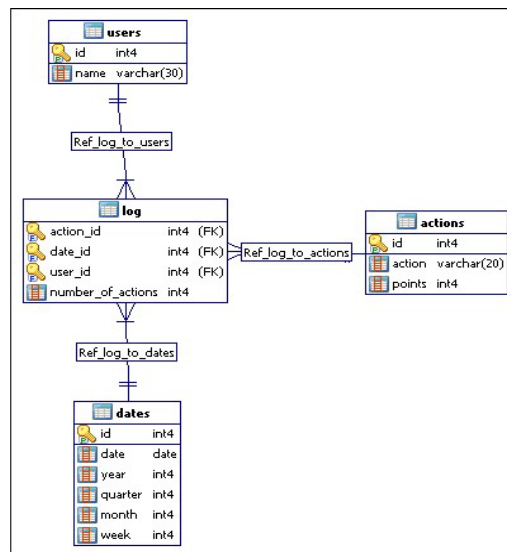


Figure 4. Data model

Some of the rules that were used to analyze data are listed above (Fig. 3). These rules were used to perform simple queries that are not so common in data warehouses, as well as some more complex queries. Some other rules were defined as well, but for more information we refer to [1] and [8].

In this paper we extend the idea and we show that Datalog rules can be used for other purposes as well, primarily for specifying aggregate and derived tables. In the next section, we show how to use rules to specify aggregate tables that are used to improve performances.

III. AGGREGATE TABLES

When we discuss tables that contain aggregate data, their purpose is quite simple. Since fact tables usually contain very large number of records, queries that use fact tables with large number of records can take too long. Although people that use data warehouses know that time needed to get results is much greater than the time needed when a query is posed against a transactional database, this doesn't mean that we are not interested in reducing that time. Since fact tables could contain millions of records, aggregation of a large number of records can last several minutes (or dozens of minutes). That is why we use tables with aggregated data, i.e., data that are pre-calculated and stored in order to improve query performances. By using tables with aggregated data queries can be answered much sooner, but on the other hand the size of a data warehouse grows (in addition to initial load of data and taking into account incremental loads that have to be carried out on a daily or weekly basis). But this is a trade-off that one does in order to speed up the access to relevant information.

Once aggregate tables are added into the system, data in aggregate tables have to be maintained as well. We do not discuss aggregate table maintenance any further (one can extract data from original sources or one can use tables from the data warehouse), but it is important that queries that are executed on aggregate tables return the same results as queries that use fact tables. Just to have in mind, when one uses aggregate and derived tables performances can be improved up to several hundreds or thousands time because the number of records (and I/O operations) becomes significantly reduced [9].

Once aggregate tables are created, the system needs to know how and when to use them. Although aggregate tables are created, that doesn't mean that it makes sense to use them always, but in some occasions (most certainly) they should be used. In order to use aggregate tables (when needed), there should exist a component that is called aggregate navigator [9]. It has a number of tasks, but the most important one is to know which aggregate table to use and when.

One Business Intelligence (BI) tool that is used in this paper is Business Objects XI. This tool has a construct called @Aggregate_Aware that is used to specify that aggregate table exists and that it should be used in certain queries (Fig. 5):

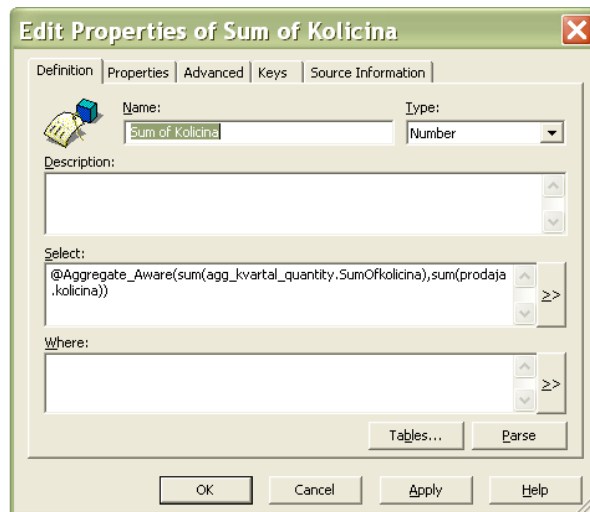


Figure 5. Business Objects XI - @Aggregate_Aware

When talking about aggregate tables, the most common scenario is to create such a table by means of a query that uses GROUP BY clause and SUM() as an aggregate function.

In this example, we show how to create a rule that (in fact) represents an aggregate table (log file analysis scenario), but it is not materialized (Fig. 6). Datalog Educational System (DES) was used to implement the rules and results are presented below the rules (the first row means that user 1, Smith Peter, committed 645 actions in January):

```
user_month_agg(A,B,K,R):-
group by((users(A,B),log(C,D,A,F),dates(D,H,I,J,K,L)),
[A,B,K], R=sum(F)).
```

```
DES> user_month_agg(A,B,C,D).
{
user_month_agg(1,'Smith Peter',1,645),
user_month_agg(1,'Smith Peter',2,637),
user_month_agg(1,'Smith Peter',3,573),
user_month_agg(1,'Smith Peter',4,585),
user_month_agg(1,'Smith Peter',5,642),
user_month_agg(1,'Smith Peter',6,604),
user_month_agg(1,'Smith Peter',7,645),
user_month_agg(1,'Smith Peter',8,581),
user_month_agg(1,'Smith Peter',9,564),
user_month_agg(1,'Smith Peter',10,585),
user_month_agg(1,'Smith Peter',11,626),
user_month_agg(1,'Smith Peter',12,593)
}
Info: 12 tuples computed.
```

Figure 6. Data aggregation

Now, when we want to compute the sum of the number of actions on a quarter (year) level (hierarchy), the idea is that the system uses *user_month_agg* table and not the fact table any more. Namely, the original fact table (called log) could be transactional. This would mean that every action that user did in the past was stored in the fact table on a certain date and time. However, the log table that we have stores the number of actions for each day, so it is a periodic snapshot fact table.

In Business Objects XI it is quite easy to use hierarchies and such reports are easy to produce. In Designer tool one creates a hierarchy (we see Hierarchies Editor in Fig. 7):

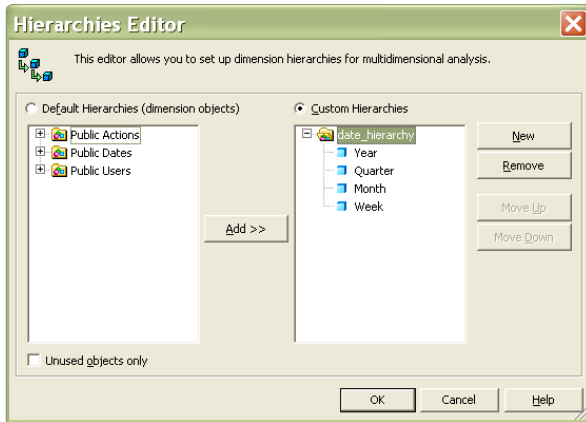


Figure 7. Business Objects XI - Hierarchies Editor

Once hierarchies are specified, one can use them in Desktop Intelligence tool. Mouse over quarter column offers drill down to a lower (Month) level (Fig. 8):

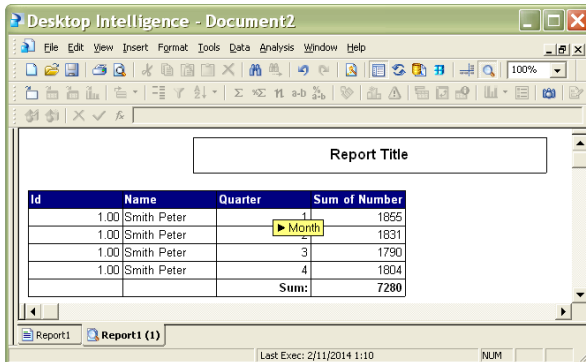


Figure 8. Number of actions (quarter level)

We can see the month level results (Fig. 9) that are only a few mouse clicks away:

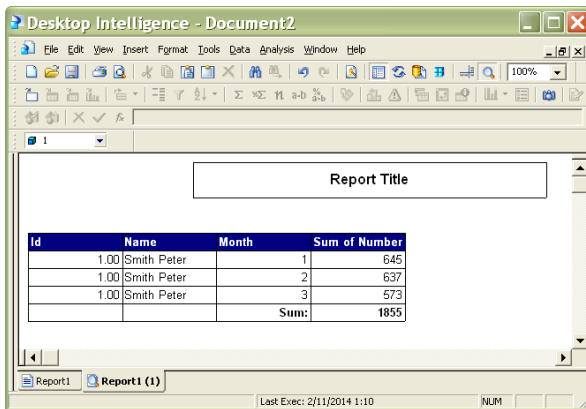


Figure 9. Number of actions (month level)

In order to calculate the number of actions on a quarter level by means of Datalog rules, we could define a rule (*uq*) that looks like this (Fig. 10):

```

uq(A,B,J,R):-
(group_by((users(A,B),log(C,D,A,F),dates(D,H,I,J,K,L)),
[A,B,J], R=sum(F))).

DES> uq(A,B,J,R)
{
uq(1,'Smith Peter',1,1855),
uq(1,'Smith Peter',2,1831),
uq(1,'Smith Peter',3,1790),
uq(1,'Smith Peter',4,1804)
}
Info: 4 tuples computed.
    
```

Figure 10. Number of actions – quarter level

The rule uses the log fact table and groups the records in order to get the result. The first row means that in the first quarter the user 1, Smith Peter, committed 1855 actions, etc. In DES, it took 17 seconds to answer the query. We can see that results in DES and in Business Objects XI are the same, as one could expect (Fig. 8 and Fig. 10).

But since we already have a rule that calculates the sum of actions on a month level, months could be easily aggregated to a higher (quarter) level. For a beginning let us add a rule that merges months and quarters (Fig. 11):

```

user_date(A,B,J,K,R):-
user_month_agg(A,B,K,R),
dates(E,F,G,J,K,I).

DES> user_date(A,B,C,D,E).
{
user_date(1,'Smith Peter',1,1,645),
user_date(1,'Smith Peter',1,2,637),
user_date(1,'Smith Peter',1,3,573),
user_date(1,'Smith Peter',2,4,585),
user_date(1,'Smith Peter',2,5,642),
user_date(1,'Smith Peter',2,6,604),
user_date(1,'Smith Peter',3,7,645),
user_date(1,'Smith Peter',3,8,581),
user_date(1,'Smith Peter',3,9,564),
user_date(1,'Smith Peter',4,10,585),
user_date(1,'Smith Peter',4,11,626),
user_date(1,'Smith Peter',4,12,593)
}
Info: 12 tuples computed.
    
```

Figure 11. Number of actions – month level

The first row means that the first user (Smith Peter) committed 645 actions in the first month of the first quarter, etc. Once months are joined with the date dimension, we can aggregate on other attributes using the date dimension (Fig. 12):

```

user_quarter(A,B,J,X):-
group_by(user_date(A,B,J,K,R), [A,B,J], X=sum(R)).

```

```

DES> user_quarter(A,B,C,D).
{
user_quarter(1,'Smith Peter',1,1855),
user_quarter(1,'Smith Peter',2,1831),
user_quarter(1,'Smith Peter',3,1790),
user_quarter(1,'Smith Peter',4,1804)
}
Info: 4 tuples computed.

```

Figure 12. Number of actions – user quarter

When called, DES needed 17 seconds to calculate the result (we restarted the program in between). So, implicit definition does not seem to be helpful except it just reduces the rule as such. Namely, this rule is used to give the same result (i.e., the number of actions on a quarter level), but it does not use the log fact table any more. It uses the implicit definition aggregate table, i.e., the *user_date* view. In the next section we try to see what could happen if the rule was materialized.

IV. MATERIALIZED VIEWS – EXPERIMENTAL RESULTS

However, if the view was materialized (physically), the time need to calculate the answer should be much smaller. The next SQL statement was used to materialize the view (Fig. 13):

```

SELECT u.id "user", u.name "name",
d.quarter "quarter", d.month "month",
SUM(l.number of actions)
INTO uqa
FROM users u INNER JOIN log l ON(u.id = l.user_id)
INNER JOIN dates d ON(l.date_id = d.id)
GROUP BY 1, 2, 3, 4
ORDER BY 3, 4;

```

Figure 13. View materialization (SELECT INTO statement)

This statement created a table called *uqa*. Once the view was materialized (in the form of a table), the results were calculated much sooner (it took less than a second). The next rule uses *uqa* table to perform the grouping and to calculate the result:

```

user_quarter_m(A,B,C,D,X):-
group_by(uqa(A,B,C,D,E), [A,B,C,D], X=sum(E)).

```

Figure 14. Using materialized view (uqa) in a rule

Of course, one has to have in mind that once data are aggregated to a higher level, some lower level queries cannot be answered any longer because the details are lost. More on materialized views can be found in [14].

Based on the previous discussion these would be some basic prerequisites that the aggregate navigator should possess. It should be capable to recognize that aggregate tables exist and it should be able to use them in situations when it makes sense. Here, we could extend the approach in order to make a Datalog aggregate navigator implemen-

tation but this could be done in our future papers. In the next section, we show how derived tables could be used.

V. DERIVED TABLES

In this section, we present different types of derived tables ([9]) and we show how to implement some of them in Datalog.

When talking about derived tables, several types can be distinguished. Pre-joined table is a table that consists of several tables that are joined together in order to speed up the querying. Further on, one can define derived tables that contain only a portion of data coming from original tables, etc. However, one has to have in mind that they may become quite big and they may require additional space.

The first example is used to demonstrate how Datalog rules can be used to create pre-joined tables. The rule name *dm* stands for *data mart* and it means that the rule would contain data from three different tables (users, log and dates), i.e., it would represent a complete data mart (only several rows are shown in the result):

```

dm(A,B,C,D,F,H,I,J,K,L):
users(A,B),
log(C,D,A,F),
dates(D,H,I,J,K,L).

dm(A,B,C,D,F,H,I,J,K,L).
...
dm(1,'Smith
Peter',4,362,1,date(2012,12,27),2012,4,12,52),
dm(1,'Smith
Peter',4,363,3,date(2012,12,28),2012,4,12,52),
dm(1,'Smith
Peter',4,364,6,date(2012,12,29),2012,4,12,52),
dm(1,'Smith
Peter',4,365,2,date(2012,12,30),2012,4,12,52),
dm(1,'SmithPeter',4,366,5,date(2012,12,31),2012,4,12,1)
}
Info: 1464 tuples computed.

```

Figure 15. Data mart specification

We can see that data from three different tables can be accessed easily, from a single (implicit) data mart. However, the time needed to calculate the data mart was a little less than 20 seconds.

Further on, we can define a rule that contains only a portion of data from the original fact table; one has to add a condition $A=1$ (Fig 16.):

```
portion(A,B,C,D):-log(A,B,C,D), A=1.
```

```

DES> portion(A,B,C,D).
...
portion(1,363,1,1),
portion(1,364,1,10),
portion(1,365,1,1),
portion(1,366,1,2)
}
Info: 366 tuples computed.

```

Figure 16. Partial fact table

Here we select only one portion of the fact table, more precisely only rows that refer to `action_id = 1` (actions in the paper were events such as read, update and insert). There are other possible ways to produce derived tables (one can combine two of the already mentioned approaches) or define other derived tables. For example, one derived table could be used to transform data (certain measures) from original table if there was a need to do so, etc.

VI. CONCLUSION

In this paper, it has been presented how deductive data warehouses could use Datalog rules to specify aggregate and derived tables. On a number of examples it was shown how to use Datalog rules in order to explain how aggregate navigator should behave and to demonstrate how other types of derived tables could be built as well. A few reports were built in DES as well as in Business Objects XI and a small data warehouse was implemented in PostgreSQL database management system. Implicit table definitions may seem to be interesting as they do not require additional space, but only after view materialization we noticed that performances were improved significantly.

Based on the previous results ([1]) it is now clear that deductive data warehouses support OLAP analysis on data and some other (more complex) analysis as well. In this paper, we have shown how to add support for derived and aggregate tables and we showed that view materialization is good for data warehousing purposes.

In future papers one could look at how to implement aggregate navigator in Datalog. Furthermore, one can see that it is not practical to work with dimension tables that have large number of attributes. Because of that one could also explore and see how to create rules more easily and how to pose goals more intuitively.

REFERENCES

- [1] K. Rabuzin, "Deductive data warehouses," IJDWM, in press.
- [2] R. Kimball and J. Caserta, *The Data Warehouse ETL Toolkit: Practical Techniques for Extracting, Cleaning, Conforming, and Delivering Data*. Indianapolis, USA: Wiley Publishing, 2004.
- [3] H. W. Inmon, *Building the Data Warehouse – Third Edition*. New York, USA: John Wiley & Sons, 2002.
- [4] C. Ballard, D. Herreman, D. Schau, R. Bell, E. Kim, and A. Valencic, *Data Modeling Techniques for Data Warehousing*. [Online]. Available from: <http://www.redbooks.ibm.com/redbooks/pdfs/sg242238.pdf>. Retrieved on 16.04.2013.
- [5] F. Silvers, *Building and Maintaining a Data Warehouse*. Boca Raton, USA: CRC Press, 2008.
- [6] M. Novak and K. Rabuzin, "Prototype of a web ETL tool," unpublished.
- [7] P. Ponniah, *Data Warehousing Fundamentals: A Comprehensive Guide for IT Professionals*. New York, USA: John Wiley & Sons, 2001.
- [8] K. Rabuzin, A. Lovrencic, and M. Malekovic, "Using deductive data warehouses to analyze data", *The Business Review*, in press.
- [9] C. Adamson, *Mastering Data Warehouse Aggregates, Solutions for Star Schema Performance*, USA: Wiley Publishing, 2006.
- [10] J. F. Boulicaut, P. Marcel, and C. Rigotti, "Query driven knowledge discovery via OLAP manipulations", [Online]. Available from: <http://liris.cnrs.fr/~jboulica/bda01.pdf>. Retrieved on 30.10.2013.
- [11] B. Neumayr, S. Anderlik, and M. Schrefl, "Towards ontology-based OLAP: datalog-based reasoning over multidimensional ontologies", *DOLAP '12 Proceedings of the fifteenth international workshop on Data warehousing and OLAP*, 2012, pp. 41-48.
- [12] J. Aligon, P. Marcel, and E. Negre, "Summarizing and querying logs of OLAP queries", *Advances in Knowledge Discovery and Management*, 471, pp. 99-124, 2013.
- [13] H. C. Tjioe and D. Taniar, "Mining Association Rules in Data Warehouses", *IJDWM*, vol. 1, pp. 28-62, 2005.
- [14] J. V. Harrison and S. W. Dietrich, "Maintenance of materialized views in a deductive database: an update propagation approach", [Online]. Available from: <http://www.public.asu.edu/~dietrich/publications/MaintenanceOfMaterializedViews.pdf>. Retrieved on 25.04.2013.

P2P Integration of Relational Knowledge Bases

Tadeusz Pankowski

Institute of Control and Information Engineering
Poznań University of Technology

Poznań, Poland

Email: tadeusz.pankowski@put.poznan.pl

Abstract—We discuss some strategies of query answering in a Peer-to-Peer (P2P) knowledge integration system. In such a system, a set of autonomous services (peers) manage knowledge bases, which are connected by means of mappings between signatures of these knowledge bases. A query is issued against an arbitrarily chosen peer (a target peer) and is propagated along semantic paths determined by mappings. Next, partial answers are sent back to the target peer. We discuss a strategy of both query propagation and merging partial answers with possibility of discovering some "missing values". The proposed method guarantees improvement of quality of answers (their completeness) and controlling efficiency of merging partial answers by deciding whether it is useful to involve the whole peer's knowledge base in the process of discovering missing values.

Keywords—knowledge bases; data integration; integrity constraints; data exchange; ontology-based data management.

I. INTRODUCTION

In recent years, we can observe a dynamic development of the Semantic Web technologies and using these technologies to create Web-oriented applications. Semantic Web technologies enable web-wide integration of data coming from various sources. In this way, the Web of Data is created, which can be perceived as a set of interrelated knowledge bases. Extensional layers of these knowledge bases consist of sets of the Resource Description Framework (RDF) triples [20] (or the corresponding Web Ontology Language (OWL) assertions [11]), and intensional layers are sets of axioms (in RDF Schema or OWL). Very often, the knowledge bases expose relational databases – then we can call them *relational knowledge bases* (RKBs). Owners of such RKBs are often interested in making them available to a wide range of users. So, the owners might be interested in including their knowledge resources into knowledge integration or knowledge exchange systems.

Assuming that RKBs independently created resources on the Web of Data, we face the challenging issues of knowledge exchange and knowledge integration across them. In the case of the knowledge exchange, we have to solve problems connected with mappings and restructuring the knowledge (captured by axioms in TBoxes and assertions in ABoxes) [1], and in the case of knowledge integration we have to do with execution and rewriting of queries according to so-called Ontology-Based Data Access (OBDA) paradigm [2].

A. Related work

Some recent results of representing relational databases in the Semantic Web are surveyed in [3] and some solutions were proposed in [4]–[6]. A relationship between relational databases and Description Logics (DLs) knowledge bases

has been studied in [7] and [6]. Similarities and differences between databases and knowledge bases, and combining these technologies in data integration activities, has been an important and attractive field of research since many years [7]–[9]. Usually, a knowledge base is founded on one of DL variants [10] or on one of OWL profiles [11]. Formally, a knowledge base is a pair $\mathcal{K} = (\mathcal{T}, \mathcal{A})$, where \mathcal{T} is a set of axioms modeling the intensional knowledge (the TBox axioms), and \mathcal{A} is a set of assertions forming the extensional knowledge (the ABox assertions). In the case of RKBs, \mathcal{T} is divided into two parts, $\mathcal{T} = (\mathcal{S}, \mathcal{C})$, where: \mathcal{S} is a set of axioms treated as *deductive rules* (so-called *standard TBox axioms*), and \mathcal{C} is a set of axioms treated as *checks* (called *integrity constraint TBox axioms* or *ic-axioms*) [7]. Some standards for representing relational databases by means of RDF and OWL have been under development by a special W3C Working Group [5]. Following these RDB-to-RDF methods of mapping, one can create Web of Data repositories consisting of a set of RKBs representing traditional relational databases (RDBs). It is expected that such a representation obeys some requirements concerning *preservation of information and semantics* of the underlying RDBs [4] [6].

One of the most promising integration architectures is the integration in P2P environment [12] [13]. Much work has been done on data integration systems both with a mediated (global) schema and in P2P architecture, where the schema of any peer can play the role of the mediated schema. There is a number of systems built in P2P data integration paradigm [14] (notably Piazza [15], PeerDB [16]). In these research, the focus was on overcoming syntactic heterogeneity and schema mappings were used to specify how data structured under one schema (the source schema) can be transformed into data structured under another schema (the target schema) [17]. A little work has been paid on how schema constraints influence the query propagation.

B. Contribution

In this paper, we will follow the P2P integration architecture and we will assume that the OBDA paradigm can be applied to any peer. Then, the peer's knowledge base plays the role of a reference ontology and mappings (alignments) between ontologies (knowledge bases) are used in query rewriting. The user issues queries against an arbitrarily chosen peer (the target peer) and expects that the answer will include relevant data stored in all P2P-connected data sources. The data sources are related by means of schema mappings. A query must be propagated to all peers in the system along semantic paths determined by mappings and reformulated accordingly.

The partial answers must be merged and sent back to the calling peer.

In this paper, we will focus on strategies of query propagation and query answering in a P2P knowledge integration system. The aim is to obtain an answer with a maximal information contents. We will show that this contents depends on the way the query is propagated and on the way in which the partial answers are being merged. The important role plays the process of discovering *missing values*, i.e., values denoted by labeled nulls. Missing values can be sometimes discovered – either by merging only partial answers or in result of involving also the whole local RKB. In the latter, the cost of computation can be high, so the decision what way of merging should be applied is significant.

We formulate and prove a necessary condition (Proposition 4.2 in Section IV) stating when it is useful to discover missing values by referring to the whole local RKB, and when it is pointless. We will discuss a strategy of query propagation and the method of merging partial answers with possibility of discovering some missing values. The proposed method guarantees improvement of quality of answers (their completeness) and controlling efficiency of merging partial queries by deciding whether it is useful to involve a whole peer’s knowledge base in the process of discovering missing values. We shortly show how the issues under consideration have been implemented in SixP2P (*Semantic Integration in P2P environment*) system [18] [19].

The structure of the paper is as follows. In Section II, we introduce a running example consisting of relational databases (RDBs), we study some query propagation strategies and merging partial answers. In Section III, we discuss a way of representing RDBs by means of RKBs. The P2P knowledge integration is proposed in Section IV. Section V summarizes and concludes the paper.

II. MOTIVATION SCENARIO

In this section, we motivate our research. We will start with three relational databases (Figure 1) as information resources stored in three peers forming a P2P data integration system (Figure 2). We discuss how a sample query can be propagated and answered in the system, and how these influence the contents of the answer to the query.

Without loss of generality we will assume that names and attributes of tables in databases are pair-wise disjoint. Integrity constraints in DB_1 (analogously in DB_2 and DB_3) are:

- $PKey(Paper1, PapId1)$ – $PapId1$ is the primary key in $Paper1$, i.e., values in $PapId1$ are both unique and not-null;
- $unique(Paper1, Title1)$ – not-null values in $Title1$ uniquely identifies tuples in $Paper1$;
- $FKey(Author1, APapId1, Paper1, PapId1)$ – $APapId1$ in $Author1$ is a foreign key referencing to the column $PapId1$ in $Paper1$.

The databases considered in Figure 1 can be stored in peers constituting a data integration system depicted in Figure 2.

Paper1		
$PapId1$	$Title1$	$Year1$

Author1		
$Name1$	$Univ1$	$APapId1$

Paper2	
$PapId2$	$Title2$
p1	KBs

Author2		
$Name2$	$Univ2$	$APapId2$
John	NY	p1
Ann	NULL	p1

Paper3		
$PapId3$	$Title3$	$Year3$
p1	KBs	2014

Author3	
$Name3$	$APapId3$
Ann	p1

Figure 1. Databases DB_1 (with empty instance), DB_2 and DB_3

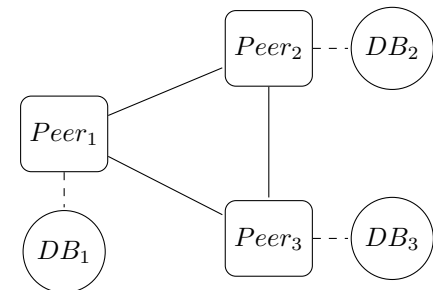


Figure 2. A sample P2P integration architecture with three peers and three local databases

Now, let us consider a query q_1 issued against DB_1 on $Peer1$. Informally, the meaning of the query is:

“Give all information about John.”

The query can be answered using different strategies. We will consider three strategies depicted in Figure 3.

- 1) In strategy (1), see Figure 3(1), q_1 is issued to DB_1 (at $Peer1$). $Peer1$ rewrites q_1 to queries q_{12} and q_{13} against, respectively, DB_2 and DB_3 and sends them to $Peer2$ and $Peer3$, respectively. Next, the answer $q_1(DB_1)$ is obtained and the peer waits for answers $q_{12}(DB_2)$ and $q_{13}(DB_3)$. After receiving all expected answers, the peer merges them producing the final answer. In this case we have:

$$\begin{aligned}
 q_1(DB_1) &= q_{13}(DB_3) = \emptyset, \\
 q_{12}(DB_2) &= \{Paper2(p1, KBs), \\
 &\quad Author2(John, NY, p1)\}.
 \end{aligned}$$

- 2) In strategy (2), see Figure 3(2), q_1 , q_{12} , and q_{13} , are sent to, respectively, $Peer1$, $Peer2$, and $Peer3$, as in the strategy (1). $Peer2$ rewrites q_{12} to q_{123} and sends it to $Peer3$. Now, $Peer3$ has to answer two queries, q_{13} and q_{123} . These queries are not identical because

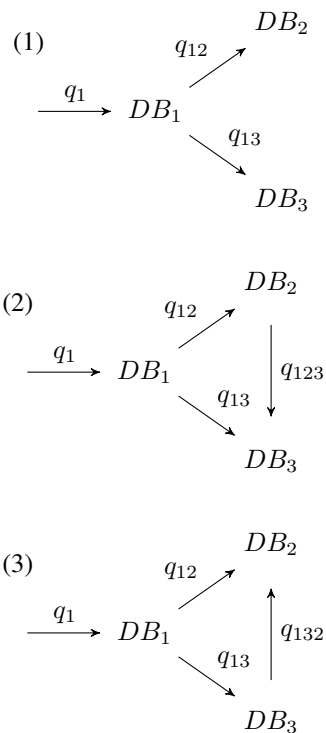


Figure 3. Query propagation strategies

q_{13} expects information about *Year3* whereas q_{123} does not, since the property *Year* is not in its area of interest. Answers $q_{13}(DB_3)$ and $q_{123}(DB_3)$ are expected to be sent to, respectively, *Peer1* and *Peer2*. *Peer2* merges $q_{123}(DB_3)$ with $q_{12}(DB_2)$ and the result returns to *Peer1*. Finally, *Peer1* merges all received data and produces the final result. In this case we have:

$$\begin{aligned} q_1(DB_1) &= q_{13}(DB_3) = q_{123}(DB_3) = \emptyset, \\ q_{12}(DB_2) &= \{Paper2(p1, KBs), \\ &\quad Author2(John, NY, p1)\}. \end{aligned}$$

- 3) Strategy (3), see Figure 3(3), is similar to strategy (2), but now, *Peer3* rewrites q_{13} to q_{132} and sends it to *Peer2*. We have the following answers:

$$\begin{aligned} q_1(DB_1) &= q_{13}(DB_3) = \emptyset, \\ q_{12}(DB_2) &= \{Paper2(p1, KBs), \\ &\quad Author2(John, NY, p1)\}, \\ q_{132}(DB_2) &= \{Paper2(p1, KBs), \\ &\quad Author2(John, , p1)\}. \end{aligned}$$

The difference between $q_{12}(DB_2)$ and $q_{132}(DB_2)$ follows from the fact that *Peer3* does not ask about university of "John".

Now, we will consider some possible ways for merging answers and discovering missing values in the process of data integration. We will use the following two operators:

- \otimes_i – a binary operator that: restructures its operands (if necessary) to the structure of DB_i and merges these operands;

- \oplus_i – a unary operator that discovers missing values in its operand using DB_i .

Merging in strategy (1):

$$q_1(DB_1) \otimes_1 q_{12}(DB_2) \otimes_1 q_{13}(DB_3) = \{Paper1(p1, KBs, NULL), Author1(John, NY, p1)\} \quad (1)$$

Merging in strategy (2). There are two possible merges:

(2a) a partial merge

$$q_1(DB_1) \otimes_1 q_{13}(DB_3) \otimes_1 (q_{12}(DB_2) \otimes_2 q_{123}(DB_3)) = \{Paper1(p1, KBs, NULL), Author1(John, NY, p1)\} \quad (2)$$

(2b) a total merge

$$q_1(DB_1) \otimes_1 q_{13}(DB_3) \otimes_1 (\oplus_2(q_{12}(DB_2) \otimes_2 q_{123}(DB_3))) = \{Paper1(p1, KBs, NULL), Author1(John, NY, p1)\} \quad (3)$$

Merging in strategy (3). Again, two merges are possible:

(3a) a partial merge

$$q_1(DB_1) \otimes_1 q_{12}(DB_2) \otimes_1 (q_{132}(DB_2) \otimes_3 q_{13}(DB_3)) = \{Paper1(p1, KBs, NULL), Author1(John, NY, p1)\} \quad (4)$$

(3b) a total merge

$$q_1(DB_1) \otimes_1 q_{12}(DB_2) \otimes_1 (\oplus_3(q_{132}(DB_2) \otimes_3 q_{13}(DB_3))) = \{Paper1(p1, KBs, 2014), Author1(John, NY, p1)\} \quad (5)$$

We see that merging strategies (1), (2a), (2b) and (3a) produce the same result. However, the strategy (3b) gives more information than the other strategies. It is so, since in (3b) the missing value of *Year1* (the year of publication of the paper KBs) has been discovered (inferred). It is possible, because there is a functional dependency between titles and years of papers.

III. RELATIONAL DB VS RELATIONAL KB

Relationships between RDBs and RKBs can be considered from the following two points of view.

1. *RDB-to-RKB transformation*. RDBs can be naturally represented in the Semantic Web by means of RDF triples or OWL specification [5]. Then, the instance of an RDB is represented by an ABox (\mathcal{A}), and the structure, properties and integrity constraints by a TBox (divided into a set of deductive axioms, \mathcal{S} , and a set of ic-axioms, \mathcal{C}). Then we obtain an RKB $\mathcal{K} = (\mathcal{S}, \mathcal{C}, \mathcal{A})$ [6].

2. *Checking properties of RKB*. For a given set of RDF triples or OWL assertions, treated as an ABox \mathcal{A} , and for a given RKB schema $\mathcal{T} = (\mathcal{S}, \mathcal{C})$, check whether the tuple $(\mathcal{S}, \mathcal{C}, \mathcal{A})$ is a consistent RKB.

Let $\mathbf{C} = \{C_1, \dots, C_k\}$ and $\mathbf{P} = \{P_1, \dots, P_p\}$ be sets of (names of), respectively, *classes* and *properties*. The set $\Sigma = \mathbf{C} \cup \mathbf{P}$ is then referred to as a *signature* of a knowledge base. We say that an RKB $\mathcal{K} = (\mathcal{S}, \mathcal{C}, \mathcal{A})$ has the signature Σ (or \mathcal{K} is over Σ) if all classes and properties occurring in \mathcal{K} are in Σ .

A. Relational knowledge bases

Let Δ_{Cons} be a set of *constants*, and Δ_{Var} be a set of *labeled nulls* or *variables*. We assume that for constants the Unique Name Assumption (UNA) holds while for labeled nulls the UNA does not hold [8]. Constants will be denoted by a, b, c (possibly with subscripts), and labeled nulls by x, y, z (possibly with subscripts), by v (possibly with subscripts) will be denoted elements from $\Delta_{Cons} \cup \Delta_{Var}$. The satisfaction of UNA means that different constants always denote different objects, i.e., the equality $a = b$, for $a, b \in \Delta_{Cons}$, is always false. In contrast, $x = y$ for $x, y \in \Delta_{Var}$ may be true if an interpretation assigns the same object to them. Similarly, $x = a$ can be true when x is interpreted as a . (For the precise definition of interpretation, see [10].)

Further on, we will write ABox assertions in a (simplified) RDF-notation [20]. In particular: $Triple(x, \text{rdf:type}, C)$ or $(x, \text{rdf:type}, C)$ corresponds to OWL class assertion $C(x)$; and $Triple(x, P, v)$ or (x, P, v) corresponds to OWL property assertion $P(x, v)$. If clear from the context, the predicate name $Triple$ will be omitted.

We assume that the following rules must always hold for any RKB. We will call them *general RKB-rules*.

- Classes are subsets of the set of labeled nulls, i.e., for each $C \in \mathcal{C}$, $C \sqsubseteq \Delta_{Var}$.¹
- Domains of properties are subsets of the set of labeled nulls, i.e., for each $P \in \mathcal{P}$, $dom(P) \sqsubseteq \Delta_{Var}$.
- Ranges of properties are subsets of the set of constants or labeled nulls, i.e., for each $P \in \mathcal{P}$, $rng(P) \sqsubseteq \Delta_{Cons} \cup \Delta_{Var}$.
- Any property is a function, i.e., for each $P \in \mathcal{P}$, the specification $(\text{funct } P) \in \mathcal{S}$ holds, which means that

$$(x, P, v_1) \wedge (x, P, v_2) \Rightarrow v_1 = v_2.$$

There are also so-called *specific RKB-rules*, which are used to characterize RKBs. These rules correspond to integrity constraints in relational databases. We restrict ourselves to key and referential constraints defined only on singletons of columns. The set of rules is divided into a set of *deductive rules* (\mathcal{S}), and a set of *integrity constraints* or *check rules* (\mathcal{C}). We assume that $(\text{funct } P) \in \mathcal{S}$. The rest of general RKB-rules belongs to \mathcal{C} . The set of specific RKB-rules includes the following classes of rules (TBox axioms):

- 1) $dom(P) \sqsubseteq C \in \mathcal{S}$ – the domain of a property P is a subset of a class C , i.e.,

$$(x, P, v) \Rightarrow (x, \text{rdf:type}, C).$$
- 2) $(\text{funct } P^-) \in \mathcal{S}$ – the inversion of a property P is a function, i.e., values of P uniquely identify objects in the domain of P

$$(x_1, P, v) \wedge (x_2, P, v) \Rightarrow x_1 = x_2.$$

¹We distinguish between the semantic notion of *inclusion* (\sqsubseteq), and the syntactic notion of *subsumption* (\sqsubset). Then $A \sqsubseteq B$ iff $A^{\mathcal{I}} \subseteq B^{\mathcal{I}}$, where \mathcal{I} is an interpretation function, see [10].

- 3) $C \sqsubseteq dom(P) \in \mathcal{C}$ – a class C is a subset of the domain of P , i.e., P is a total (not-null) on C

$$(x, \text{rdf:type}, C) \Rightarrow \exists v.(x, P, v).$$

- 4) $rng(P) \sqsubseteq C \in \mathcal{C}$ – the range of P is a subset of a class C , i.e., P references objects in C

$$(x, P, x') \Rightarrow (x', \text{rdf:type}, C).$$

B. Transformation of RDBs to RKBs

Each RDB DB_i , depicted in Figure 1, can be translated into an RKB $\mathcal{K} = (\mathcal{S}_i, \mathcal{C}_i, \mathcal{A}_i)$, where \mathcal{S}_i and \mathcal{C}_i , for $i = 1, 2, 3$, can be created from \mathcal{S} and \mathcal{C} in Figure 4 by appropriate extension of the property names and class names with postfixes '1', '2', and '3'. Additionally, in \mathcal{S}_2 do not occur rules involving $Year2$, and in \mathcal{S}_3 are not rules containing $Univ3$.

$$\begin{aligned} \mathcal{S} = \{ & dom(PapId) \sqsubseteq Paper, \\ & dom(Title) \sqsubseteq Paper, \\ & dom(Year) \sqsubseteq Paper, \\ & dom(Name) \sqsubseteq Author, \\ & dom(Univ) \sqsubseteq Author, \\ & dom(APapId) \sqsubseteq Author, \\ & (\text{funct } PapId^-), (\text{funct } Title^-)\}. \\ \mathcal{C} = \{ & Paper \sqsubseteq dom(PapId), Paper \sqsubseteq dom(Year), \\ & Author \sqsubseteq dom(Name), rng(APapId) = Paper\}. \end{aligned}$$

Figure 4. Standard TBox axioms (\mathcal{S}) and TBox ic-axioms (\mathcal{C}) corresponding to schemas of RDBs in Figure 1. Names should be appropriately extended with postfixes '1', '2', and '3'.

Then we obtain the following RKBs, $\mathcal{K}_i = \tau(DB_i)$, $i = 1, 2, 3$, where τ denotes the translation operator from RDBs into RKBs:

- 1) $\tau(DB_1) = \mathcal{K}_1 = (\mathcal{S}_1, \mathcal{C}_1, \mathcal{A}_1)$, where:

$$\mathcal{A}_1 = \emptyset.$$

- 2) $\tau(DB_2) = \mathcal{K}_2 = (\mathcal{S}_2, \mathcal{C}_2, \mathcal{A}_2)$, where:

$$\begin{aligned} \mathcal{A}_2 = \{ & (x_1, PapId2, p1), \\ & (x_1, Title2, K Bs), \\ & (x_2, Name2, John), \\ & (x_2, Univ2, NY), \\ & (x_2, APapId2, x_1), \\ & (x_3, Name2, Ann), \\ & (x_3, APapId2, x_1)\}. \end{aligned}$$

- 3) $\tau(DB_3) = \mathcal{K}_3 = (\mathcal{S}_3, \mathcal{C}_3, \mathcal{A}_3)$, where:

$$\begin{aligned} \mathcal{A}_3 = \{ & (y_1, PapId3, p1), \\ & (y_1, Title3, K Bs), \\ & (y_1, Year3, 2014), \\ & (y_2, Name3, Ann), \\ & (y_2, APapId3, y_1)\}. \end{aligned}$$

A set \mathcal{A} of ABox assertions arises from an instance of relational database in the result of the translation performed by Algorithm 1.

In Algorithm 1, for each relation symbol R and each attribute $A \in att(R)$:

Algorithm 1 Creating ABox assertions

Input: An RDB (R, IC, I) , and an empty ABox \mathcal{A} .
Output: ABox assertions in \mathcal{A} representing I .
for each $R(t) \in I$
 $x :=$ a fresh labeled null in Δ_{Var}
 $U_{R,t} := \{A \in \text{att}(R) \mid t.A \neq \text{NULL}\}$
 if $U_{R,t} = \emptyset$ **then**
 $\mathcal{A} := \mathcal{A} \cup \{(x, \text{rdf:type}, C_R)\}$
 else
 for each $A \in U_{R,t}$
 if $FKey(R, A, R', A') \in IC$ **and**
 exists x' **such that** $(x', P_{A'}, t.A) \in \mathcal{A}$ **then**
 $\mathcal{A} := \mathcal{A} \cup \{(x, P_A, x')\}$
 else
 $\mathcal{A} := \mathcal{A} \cup \{(x, P_A, t.A)\}$
end

- a fresh labeled null x is selected from Δ_{Var} for each tuple t in an instance of R ;
- for each $A \in \text{att}(R)$ if $t.A = a \neq \text{NULL}$ and A is not a foreign key in R , the triple (x, P_A, a) is inserted into \mathcal{A} ; so, fields with NULL values are omitted;
- if A is a foreign key in R , i.e., $FKey(R, A, R', A') \in IC$, and $t.A = a$, then the triple (x, P_A, x') is inserted into \mathcal{A} , where x' is such that $(x', P_{A'}, a) \in \mathcal{A}$;
- if for each attribute $A \in \text{att}(R)$, $t.A = \text{NULL}$ then $(x, \text{rdf:type}, C_R)$ is inserted into \mathcal{A} .

IV. P2P KNOWLEDGE INTEGRATION

Processing a query in a peer-to-peer environment is presented in Algorithm 2. In general, the algorithm is self-explaining. However, some procedures will be discussed deeply later on, namely:

- $\text{mergeAnswers}(\text{curPeer}, \text{query})$ – merge all answers to query gathered in curPeer ;
- $\text{canDiscover}(\text{curPeer}, \text{query})$ – decide whether some missing values in the answer to query can be discovered using the whole RKB stored in curPeer ;
- $\text{discoverUsingRKB}(\text{curPeer}, \text{query})$ – discover missing values in the answer to query , using the whole RKB in curPeer .

We consider conjunctive queries, which are conjunctions of equalities of the form $P = a$, where P is a property and a is a constant.

By $\text{Triple}(x, p, v)$ we denote a triple that either is in \mathcal{A} or can be deduced from \mathcal{A} by means of \mathcal{S} (in the current peer's RKB, $\mathcal{K} = (\mathcal{S}, \mathcal{C}, \mathcal{A})$). An answer to $P = a$ consists of all triples determined by the following recursive datalog program:

$$\begin{aligned} \text{Answer}(x, p, v) &\leftarrow \text{Triple}(x, p, v), p = "P", v = "a" \\ \text{Answer}(x, p, v) &\leftarrow \text{Answer}(x, _, _), \text{Triple}(x, p, v) \\ \text{Answer}(x, p, v) &\leftarrow \text{Answer}(_, _, x), \text{Triple}(x, p, v) \end{aligned}$$

The answer encompasses all triples, which describe an object x with the value a of property P , all other properties of x , and triples connected with x by means of references.

Further on, we will use i instead of Peer_i , and we assume the following denotations:

Algorithm 2 Processing a query in a peer

Input: Set of peers with local RKBs. A peer is connected with its partners by means of mappings between signatures of RKBs.
query - a query against a current peer.
Output: Answer to the *query* containing local answer and answers returned by all peer's partners.
for each srcPeer **in** $\text{partnersOf}(\text{curPeer})$ {
 $\text{query} = \text{rewriteQuery}(\text{curPeer}, \text{srcPeer}, \text{query})$;
 \triangleright query is rewritten using the mapping from curPeer to srcPeer
 $\text{sendQuery}(\text{curPeer}, \text{srcPeer}, \text{query})$;
 \triangleright query is sent from curPeer to srcPeer
 }
for each srcPeer **in** $\text{partnersOf}(\text{curPeer})$ {
 $\text{answer} = \text{answer} \cup$
 $\text{transformAnswer}(\text{curPeer}, \text{srcPeer}, \text{query})$;
 \triangleright answer to query from srcPeer is gathered in curPeer
 }
 $\text{answer} = \text{mergeAnswers}(\text{curPeer}, \text{query})$;
 \triangleright all answers to query are merged in curPeer
if $\text{canDiscover}(\text{curPeer}, \text{query})$ **then** {
 $\text{answer} = \text{discoverUsingRKB}(\text{curPeer}, \text{query})$;
 \triangleright curPeer 's RKB is used to discover missing values
 \triangleright in the answer to query
 }

- 1) $\{i_1, \dots, i_n\}$ – a set of partners of a peer i ;
- 2) $\mathcal{M}_{k,i}$ – a mapping between signatures Σ_k and Σ_i of RKBs stored in peers k and i , respectively;
- 3) $\mathcal{M}_{k,i}(\text{ans}_k(q))$ – the transformation of an answer $\text{ans}_k(q)$ specified by a mapping $\mathcal{M}_{k,i}$.

Then the answer $\text{ans}_i(q)$ is obtained in the following steps:

Step 1. The answer $q(i)$, i.e., the local answer to q w.r.t. the current RKB is returned

$$\text{ans}_i(q) = q(i).$$

Step 2. Answers of all i 's partners are added to $\text{ans}_i(q)$

$$\text{ans}_i(q) = \text{ans}_i(q) \cup \{\mathcal{M}_{k,i}(\text{ans}_k(q)) \mid k \in \{i_1, \dots, i_n\}\}.$$

Step 3. Axioms in $\mathcal{S}_i \cup \mathcal{C}_i$ are used to deduce all possible facts in $\text{ans}_i(q)$

$$\begin{aligned} \text{ans}_i(q) &= \text{ans}_i(q) \cup \mathcal{S}_i \cup \mathcal{C}_i = \\ &\{(x, p, v) \mid \text{ans}_i(q) \cup \mathcal{S}_i \cup \mathcal{C}_i \vdash (x, p, v)\}. \end{aligned}$$

Let us focus on the merge operation in Peer_3 – see the merging strategy (5) in Section II applied to RKBs – we have:

- Step 1 : $\text{ans}_3(q) = \emptyset$;
 Step 2 : $\text{ans}_3(q) = \text{ans}_3(q) \cup \mathcal{M}_{23}(\text{ans}_2(q)) =$
 $\{(x_1, \text{PapId3}, p1), (x_1, \text{Title3}, \text{KBs}),$
 $(x_2, \text{Name3}, \text{John}), (x_2, \text{APapId3}, x_1)\};$
 Step 3 : $\text{ans}_3(q) = \text{ans}_3(q) \cup \mathcal{S}_3 \cup \mathcal{C}_3 =$
 $\{(x_1, \text{PapId3}, p1), (x_1, \text{Title3}, \text{KBs}),$
 $(x_2, \text{Name3}, \text{John}), (x_2, \text{APapId3}, x_1),$
 (a) $(x_1, \text{rdf:type}, \text{Paper3}),$
 (b) $(x_2, \text{rdf:type}, \text{Author3}),$
 (c) $(x_1, \text{Year3}, x_3)\}.$

Triples (a) and (b) were deduced using $\text{dom}(\text{Title3}) \sqsubseteq \text{Paper3} \in \mathcal{S}_3$, and $\text{dom}(\text{Name3}) \sqsubseteq \text{Author3} \in \mathcal{S}_3$. Triple

(c) was deduced by application of the axiom $Paper3 \sqsubseteq dom(Year3) \in \mathcal{C}_3$ to the triple $(x_1, rdf:type, Paper3)$. Note that x_3 is a "fresh" labeled null in Δ_{Var} .

Now, the question is whether the missing value represented by x_3 can be discovered referring to the whole RKB \mathcal{K}_3 . To this order we can proceed as follows:

$$\{(x_1, Title3, KBs), (x_1, Year3, x_3)\} \subseteq ans_3(q),$$

$$\{(y_1, Title3, KBs), (y_1, Year3, 2014)\} \subseteq \mathcal{A}_3.$$

Since $Title3$ is unique in $\mathcal{A}_3 \cup ans_3(q)$, which is stated by the axiom $(funct\ Title3^-) \in \mathcal{S}_3$, we have $x_1 = y_1$. Now we can apply the axiom $(funct\ Year3) \in \mathcal{S}_3$ to the set

$$\{(y_1, Year3, 2014), (y_1, Year3, x_3)\} \subseteq \mathcal{A}_3 \cup ans_3(q).$$

From functionality of $Year3$, we have $x_3 = 2014$. In this way the missing value represented by x_3 has been discovered.

Now, we can describe the process of discovering missing values more precisely. Proposition 4.2 formulates the necessary condition for the possibility of success in the discovery process.

Definition 4.1: Let a triple (x, P, x') be in an set $ans(q)$ of triples constituting an answer to a query q . We say that x' is a *missing value* of P if there is not any triple in $ans(q)$ whose first component (the subject) is x' .

It means that a labeled null x' in a triple (x, P, x') represents a *missing value* if the property P is not interpreted as a reference.

Proposition 4.2: Let P occur in a query q issued against a peer, $\mathcal{K} = (\mathcal{S}, \mathcal{C}, \mathcal{A})$ be an RKB in this peer, and $ans(q)$ be a result of merging all partial answers returned to this peer. Let a triple $(x, P_2, x') \in ans(q)$, where x' is a missing value of P_2 . The necessary condition for discovering the value of x' is:

- (c1) there is a triple $(x, P_1, a) \in ans(q)$, $P_1 \neq P$, and
- (c2) P_1 is unique in its domain, i.e., $(funct\ P_1^-) \in \mathcal{S}$.

Proof: Note that the given condition is a necessary condition.

- 1) Let (c1) and (c2) be true. Assume that the following two triples: (y, P_1, a) and (y, P_2, b) are in \mathcal{A} . Then, from the uniqueness of P_1 it follows that $y = x$, and from functionality of P_2 w obtain $x' = b$. So, b is the discovered value of x' .
- 2) Let $P_1 = P$ and (c2) hold. Because $ans(q)$ contains also the answer to q returned by \mathcal{K} then the triple (y, P_1, a) (mentioned in 1)) occurs either in both \mathcal{A} and $ans(q)$, or in neither of them. Thus, to discover the missing value of x' we can restrict ourselves only to $ans(q)$.
- 3) Now, let (c1) holds and (c2) is not true. Then the equality $y = x$, considered in 1) can not be deduced, and the discovering process fails.

Thus, the thesis of the proposition holds. ■

V. CONCLUSIONS AND FUTURE WORK

We discussed the problem of answering queries issued in a knowledge integration systems in P2P architecture. In such a system, there are many autonomous services (peers), which

collaborates in the process of producing answers to queries. The system is flexible and peers can enter and leave the system dynamically. A peer has a knowledge base and while entering the system it establishes semantic relationships between the signature of its knowledge base and signatures of some knowledge bases already belonging to the system (its partners). A query posed against a peer is propagated to its partners along semantic paths defined by mapping, those partners propagate the query to their partners, etc. Answers to the query flow back in the opposite directions. The target peer merge the answers producing the expected answer. In this merge either only answers are taken into account or also the whole RKB stored in the peer can be involved. This influence the quality (completeness) of the answer as well as the efficiency of the query answering process. We have shown how one can control these factors of the knowledge integration system. In particular, we discussed when the involvement of the whole RKB is useful from the so-called "discovery of missing values" point of view. The proper decision significantly influence both the quality of the answer as well as efficiency of processing. The discussed strategy was inspired by the SixP2P system originally designed for integrating XML data [18] [19]. The future extension of the implementation is intended to capture also semantic-oriented repositories, such as relational knowledge bases organized in a form of RDF triples or OWL specifications.

ACKNOWLEDGMENT

This research has been supported by Polish Ministry of Science and Higher Education under grant 04/45/DSPB/0122.

REFERENCES

- [1] M. Arenas, E. Botoeva, and D. Calvanese, "Knowledge base exchange," in 24th International Workshop on Description Logics (DL 2011), 2011, pp. 1–11.
- [2] D. Calvanese, G. Giacomo, D. Lembo, M. Lenzerini, and R. Rosati, "Tractable Reasoning and Efficient Query Answering in Description Logics: The DL-Lite Family," J. Autom. Reason., vol. 39, no. 3, 2007, pp. 385–429.
- [3] J. Sequeda, S. H. Tirmizi, Ó. Corcho, and D. P. Miranker, "Survey of directly mapping SQL databases to the Semantic Web," Knowledge Eng. Review, vol. 26, no. 4, 2011, pp. 445–486.
- [4] J. Sequeda, M. Arenas, and D. P. Miranker, "On Directly Mapping Relational Databases to RDF and OWL (Extended Version)," CoRR, vol. abs/1202.3667, 2012, pp. 1–17.
- [5] M. Arenas, A. Bertails, E. Prud'hommeaux, and J. Sequeda, "A Direct Mapping of Relational Data to RDF," 2012, <http://www.w3.org/TR/rdb-direct-mapping>, [retrieved: May, 2014].
- [6] T. Pankowski, "Reasoning About Consistency Of Relational Knowledge Bases," in 8th International Multi-Conference on Computing in the Global Information Technology, ICCGI 2013, July 21 - 26, 2013 - Nice, France. IARIA, 2013, pp. 283–288.
- [7] B. Motik, I. Horrocks, and U. Sattler, "Bridging the gap between OWL and relational databases," Journal of Web Semantics, vol. 7, no. 2, 2009, pp. 74–89.
- [8] S. Abiteboul, R. Hull, and V. Vianu, Foundations of Databases. Reading, Massachusetts: Addison-Wesley, 1995.
- [9] R. Reiter, "Towards a logical reconstruction of relational database theory," in On Conceptual Modelling. Perspectives from Artificial Intelligence, Databases, and programming Languages, 1982, pp. 191–233.
- [10] F. Baader, D. Calvanese, D. McGuinness, D. Nardi, and P. Petel-Schneider, Eds., The Description Logic Handbook: Theory, Implementation and Applications. Cambridge University Press, 2003.

- [11] OWL 2 Web Ontology Language Profiles, 2009, www.w3.org/TR/owl2-profiles/, [retrieved: May, 2014].
- [12] D. Calvanese, E. Damaggio, G. D. Giacomo, M. Lenzerini, and R. Rosati, "Semantic data integration in P2P systems," in DBISP2P-2003, LNCS 2944, Springer, 2004, pp. 77–90.
- [13] F. Buccafurri and G. Lax, "Enabling Selective Flooding to Reduce P2P Traffic," in OTM Conferences, LNCS 4803, Springer, 2007, pp. 188–205.
- [14] G. Koloniari and E. Pitoura, "Peer-to-peer management of XML data: issues and research challenges," SIGMOD Record, vol. 34, no. 2, 2005, pp. 6–17.
- [15] I. Tatarinov, et al., "The Piazza peer data management project," SIGMOD Record, vol. 32, no. 3, 2003, pp. 47–52.
- [16] B. C. Ooi, Y. Shu, , and K.-L. Tan, "Relational data sharing in peer-based data management systems," SIGMOD Record, vol. 32, no. 3, 2003, pp. 59–64.
- [17] A. Fuxman, P. G. Kolaitis, R. J. Miller, and W. C. Tan, "Peer data exchange," ACM Trans. Database Syst, vol. 31, no. 4, 2006, pp. 1454–1498.
- [18] G. Brzykcy, J. Bartoszek, and T. Pankowski, "Schema Mappings and Agents' Actions in P2P Data Integration System," Journal of Universal Computer Science, vol. 14, no. 7, 2008, pp. 1048–1060.
- [19] T. Pankowski, "Query propagation in a P2P data integration system in the presence of schema constraints," in Data Management in Grid and P2P Systems DEXA/Globe'08, LNCS 5187, 2008, pp. 46–57.
- [20] Resource Description Framework (RDF) Model and Syntax Specification, 1999, www.w3.org/TR/PR-rdf-syntax/, [retrieved: May, 2014].

Generating Customized Sparse Eigenvalue Solutions with Lighthouse

Ramya Nair*, Sa-Lin Bernstein[†], Elizabeth Jessup* and Boyana Norris[‡]

*Department of Computer Science, University of Colorado Boulder, Boulder, CO, USA

Email: ramya.nair@colorado.edu, elizabeth.jessup@colorado.edu

[†]Computation Institute, University of Chicago and Argonne National Laboratory, Chicago, IL, USA

Email: salin@anl.gov

[‡]Department of Computer and Information Science, University of Oregon, Eugene, OR, USA

Email: norris@cs.uoregon.edu

Abstract—Sparse eigenvalue problems arise in many areas of scientific computing. A variety of high-performance numerical software packages including many different eigensolvers are available to solve such problems. The two main challenges are finding the routines that can correctly solve the problem and implementing the desired solution accurately and efficiently using the appropriate software package. In this paper, we describe an approach that addresses these issues by intelligently identifying the sparse eigensolver that is likely to perform the best for given input characteristics and by generating a code template that uses that solver. The results are delivered to users through Lighthouse, a novel interface and search platform for users seeking high-performance solutions to linear algebra problems. This paper describes the development of the approach with a focus on the analysis of sparse eigensolvers in SLEPc and their integration into Lighthouse.

Keywords—expert systems; sparse eigensolvers; machine learning.

I. INTRODUCTION

Sparse eigenvalue problems are pervasive in scientific computing, engineering, and optimization-based machine learning methods. Such problems appear in various fields ranging from subatomic particle theories in quantum physics [1] to structural engineering [2]. In recent years, they are also increasingly used in search engines [3] and social networks [4]. Scientists and engineers typically use high-performance libraries that were developed over time by teams of numerical linear algebra experts. However, solving sparse eigenvalue problems accurately and efficiently with these packages normally takes significant effort and requires knowledge in applied mathematics, computational methods, and the problem domain. One of the commonly known example of sparse eigenvalue problem is the pagerank algorithm used in Google's search [3]. A modified version of power method was used to solve the search problem.

One of the common challenges faced by the developers is the lack of knowledge of different iterative methods for solving eigenvalue problems. Some methods or algorithms may converge much faster than others for the particular problem at hand. Hence, an inappropriate solver selection or configuration may lead to an unwanted result such as a high residual or no convergence. Another challenge lies in selecting a suitable library and identifying the most appropriate routines for a given problem. These processes, generally, depend on the resources at one's disposal ranging from hardware criteria, such as the number of processors available, to personal convenience, such as the preferred programming language. Optimizing the performance of an implementation is a big challenge for many developers because it requires a solid understanding of the

selected software package framework, numerical computation, compilers, and computer architecture.

Lighthouse [5] is an innovative framework that connects linear algebra software resources with code implementation and optimization. Lighthouse is an ongoing project with ever-increasing content and functionality. This paper focuses on the addition of eigensolvers from the Scalable Library for Eigenvalue Problem Computations (SLEPc) [6] to Lighthouse and is organized as follows. Section II provides the related work. Section III discusses our approach for developing the Lighthouse taxonomy. Section IV presents the details of the integration of SLEPc with Lighthouse to identify and deliver the best sparse eigensolvers to users. Section V summarizes the conclusions and outlines future work.

II. RELATED WORK

A number of existing taxonomies attempt to address the problem of finding and using high-performance numerical software. Perhaps the oldest one (starting in 1985) is the Netlib Mathematical Software Repository [7], which contains freely available software, documents, and databases pertaining to numerical computing including eigensolvers. The information is organized as lists of packages or routines, with or without accompanying documentation. The newer Linear Algebra Software Survey [8] contains over sixty items categorized as support routines, dense direct solvers, sparse direct solvers, preconditioners, sparse iterative solvers, and sparse eigenvalue solvers together with a checklist specifying problem types for each entry. NIST's Guide to Available Mathematical Software (GAMS) [9] includes a wider range of basic linear algebra software along with software for a variety of other numerical applications. While the Linear Algebra Software Survey is a linear list without advanced search capabilities, GAMS allows search by problem solved, package name, module name, or text in the brief module abstract. An earlier Java-based client called HotGAMS [10] allowed an interactive search of the GAMS repository, but is no longer available. Both the Survey and GAMS index into Netlib for software downloads. It is also possible to browse and search Netlib directly.

With existing taxonomies or general-purpose search engines, the user must manually explore the many available packages and learn enough about each of them to be able to make a good choice. For full understanding, the user may also need to read significant portions of the documentation for each candidate software package. After selecting a library that can solve his or her target problem, the user typically

spends considerable time learning how to use it correctly and efficiently before they can encode their solution.

Even after rigorous comparative analysis of the methods and careful selection of software packages, the selected algorithm may or may not work for the problem of interest depending on various factors. When using high-quality software, an unsatisfactory solution may still result for certain inputs. For example, the number of converged eigenvalues may be different from expected or the residual may be greater than expected. Even if no mistakes are made in any of the steps explained above, the chosen solver may fail to produce the desired solution for the given problem, and it may be necessary to repeat some or all of the development steps.

III. PROPOSED APPROACH

Lighthouse is an open source web-based expert system that matches a user’s functional and performance needs with available high-performance linear algebra software. The Lighthouse taxonomy provides a classification of existing linear algebra libraries that currently includes Linear Algebra PACKage (LAPACK) [11], Portable, Extensible Toolkit for Scientific Computation (PETSc) [12], and SLEPc. Built with the Django [13] web application framework, the Lighthouse user interface (UI) is a user-centered design, offering efficient search capabilities that accommodate users with various levels of experience in numerical linear algebra. Figure 1 illustrates the Lighthouse Guided Search for LAPACK linear solver routines. The Guided Search is designed to lead users through increasingly refined subroutine searches until the desired result is attained. The Advanced Search, on the other hand, is recommended for users who are familiar with the library. All users can benefit from the Keyword Search, which allows for subroutine search via an input keyword or phrase.

In addition to the search feature, Lighthouse creates code templates in FORTRAN 90 and C containing working programs that declare and initialize required data structures and call selected subroutines. Users can download and modify the templates to meet their particular project needs. Moreover, Lighthouse provides the ability to automatically generate and tune high-performance implementations of custom linear algebra computations by interfacing with the Build to Order (BTO) [14] compiler. BTO produces highly tuned C implementations based on high-level MATLAB-like input specification of the computation.

The Lighthouse taxonomy is continuously expanding, and SLEPc is one of our current development focuses. SLEPc is based on the PETSc package [15] [12], which provides a comprehensive set of data structures and algorithms for the parallel solution of problems modeled with nonlinear partial differential equations.

To determine what eigensolver is the most appropriate for a given problem, we have analyzed SLEPc solvers by applying machine learning techniques to a large set of different problems to find those that are most likely to yield the best performance for a given set of matrix features. The next section describes the process of integrating SLEPc with Lighthouse.

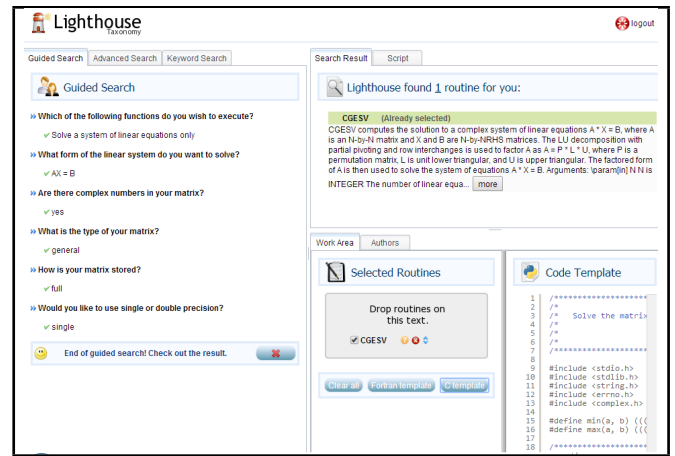


Figure 1. Lighthouse Guided Search for LAPACK linear solver routines.

TABLE I. INPUT FEATURE SET.

Property tested	Range of input parameters tested
Matrix order	112 x 112 to 262,111 x 262,111
Matrix type	Real, Complex
Matrix data	Binary, Double
Matrix characteristics	Hermitian, Non-Hermitian
Number of eigenvalues	1, 2, 5, 10
Portion of spectrum	Largest magnitude, Smallest magnitude, Largest real, Smallest real, Largest imaginary, Smallest imaginary
Tolerance	1.00E-04, 1.00E-08, 1.00E-10
Number of processors	1, 2, 4, 8, 12, 24, 48, 96, 192

TABLE II. SLEPc SOLVERS TESTED.

SLEPc solvers	power, subspace, Arnoldi, Lanczos, KrylovSchur, Generalized Davidson, Jacobi Davidson
---------------	---

IV. SLEPc - LIGHTHOUSE INTEGRATION

We experimented with different eigenvalue problems and available SLEPc solvers to determine the performance of each solver. First, we considered known functional properties of algorithms to narrow down to the solvers that work best for each problem examined (e.g., some methods are only suitable for solving symmetric problems). We then used machine learning techniques to intelligently predict the solvers that will work best for a given set of input features. The experimental results served as training data to the prediction algorithm, which identifies the sparse eigensolvers included in Lighthouse for a given problem.

A. Experiments

To run the experiments, we first obtained matrices from Matrix Market [16] and the Florida Sparse Matrix Collection [17] that cover the problem domains of interest.

Table I shows the range of input parameters tested. For each case of these input parameters, we obtained the performance and result accuracy for different SLEPc solvers. Table II shows the SLEPc solvers used for these experiments. The resultant data set consists of more than 29,000 data points.

B. Training and Prediction

To begin, we experimentally determined the solvers that work best for every case tested, after which, we used machine learning techniques to make intelligent decisions about the best known solvers for any untested case. The predicted results were then verified using standard validation techniques. The following subsections describe these steps in detail.

1) *Finding the best solver - Training data setup:* The analysis to identify the best suited eigensolvers can be split into two steps: elimination and selection. Output characteristics of interest for these two steps are the number and selection of converged eigenpairs, time taken, and residual.

In the elimination step, we first remove the solvers with resultant characteristics completely outside the expected value range. For this, we count the number of eigenvalues that converged with a residual less than the given residual tolerance. If the count obtained is less than the desired number of eigenvalues, the solver is removed from the list.

For the selection step, all of the eigensolvers remaining after the described culling are reasonable choices, but the fastest of them is selected as the best fit for the given problem specifications. Additionally, to avoid losing other efficient solvers, we select from the remaining list those that take at most 10% more time than the best fit eigensolver.

For every input permutation tested, we follow the method detailed above to obtain the resultant table (a subset of the original dataset) which identifies the best solvers for each unique input set.

Non-convergence case: After experimenting with very large upper limits, we set the maximum iteration limit to 1,000 iterations for all eigensolvers except the power method. The time to run an iteration of the power method is much less than the time for other solvers, hence we set the upper limit for it to 5,000. The experiments established that most solutions converge significantly fewer iterations. If the solution does not converge within the specified limits, it is labelled as non-convergent.

In such non-convergent cases, there is still a possibility that a solver converges but not to the exact expected number of eigenvalues or residual tolerance. Solution methods that show some signs of convergence, even though not to the desired values, are retained for the training data. Hence, such input cases have the information of the partially converged solvers along with non-convergent label.

2) *Intelligent solver selection - Prediction:* The previous step produces a reduced set of data which gives the best solvers for each tested case. We apply machine learning techniques to make intelligent predictions for other problems using the data as a training set.

Decision tree induction [18] is a popular prediction model that uses observations with known results, referred to as the training data set, to form a model for predicting the results for any data. It uses the features or attributes of the data set, the matrix properties and expected output parameters to form a pattern that can best fit the training data. The final prediction model is a *classification tree*, where every inner node is an attribute to be selected, every branch from the node

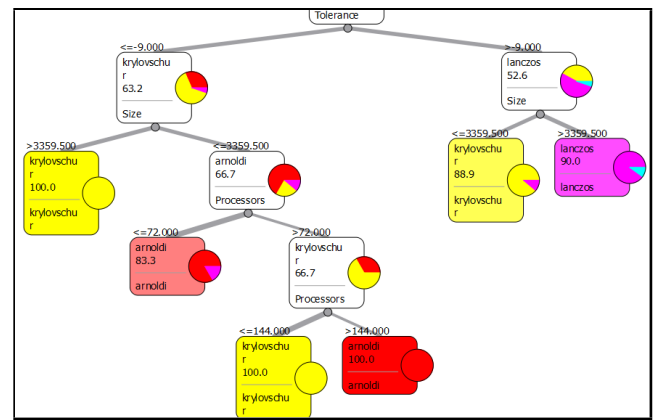


Figure 2. Part of the generated classification tree using ORANGE.

is a selection made, and the leaf nodes lead to the predicted result. For our case, we want the data to be classified into the eigensolvers available in SLEPc.

From the previous step, for every input case tested, we have one or more best solvers. To accommodate multiple solver recommendations, apart from the input parameters, we added a performance index feature to the dataset. For every input case, the performance index is given an integer value of one for the fastest solver and is incremented by one for remaining solvers with increasing time taken. Using the expanded feature set as training data, the obtained decision tree learns the order of the best performing eigenvalue solvers as well. The generated tree can be used to suggest multiple eigensolver classes by simply ignoring the performance index and taking into account all other feature criteria in the query.

We employed two applications to obtain the classification tree from the training data. We first used MATLAB's *ClassificationTree.fit* functionality [19]. We also used the application Orange [20], which is an open-source data mining tool for different learning operations, such as classification, evaluation and prediction. In particular, Orange has good visualization options that make viewing large trees much easier than with other applications like MATLAB or Weka [21].

Using these applications we obtained a classifier in the form of a binary tree, with every node representing a feature, such as matrix size, matrix type and desired eigenvalue spectrum. A portion of the classification tree generated from Orange is shown in Figure 2. The decision to follow the left branch or the right branch of a node is made depending on the value of the feature at that node. The leaf node gives the suggested eigensolver for the path followed. The tree also has the information about the matrix properties and output characteristics which do not converge to a solution for any of the eigensolvers available in SLEPc. Such leaf nodes have the value set to "No Convergence".

3) *Validation:* We employed several methods to validate the results. For the decision tree in MATLAB, cross-validation results were checked using the MATLAB `cvLoss()` function. The decision tree cross-validation classification error (loss) for the training data set was obtained to be 0.1450 which implies a 14.5 percent error in prediction. For the decision tree created using Orange, 10-fold cross-validation was applied,

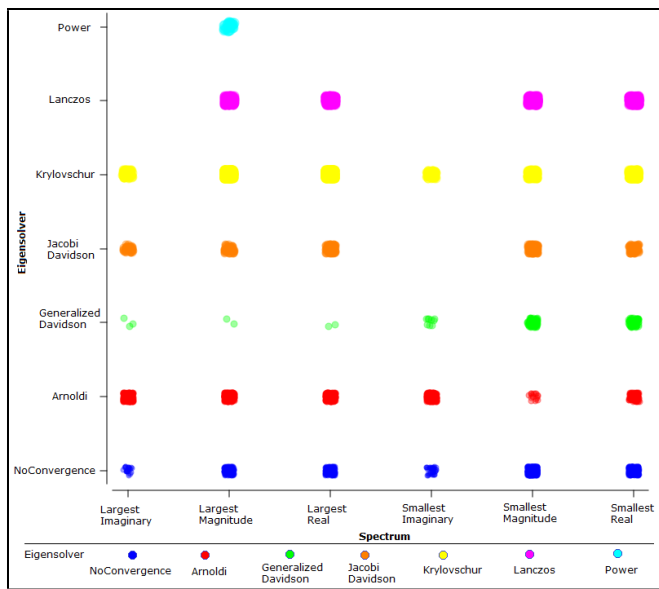


Figure 3. Scatter plot analyzing eigensolver vs portion of spectrum.

and we obtained classification accuracy of 0.8639. Close to the results from MATLAB, the latter result implies a 13.61 percent error in prediction.

The validation results obtained using the above methods take into account the order of the best performing solvers. Lighthouse implementation suggests all the favorable solvers irrespective of the order. Hence, to get more accurate validation results with respect to Lighthouse, we employed a different technique. We used the generated tree to predict the results for each input case tested, varying the performance index. As a result, we formed a favorable solver subset for every input case, each predicted solver in the subset corresponds to a different performance index. If the predicted favorable solver subset is same as the expected subset (in training data), irrespective of the performance index, we considered it to be accurately predicted. This test gives an error of 7.99%, which implies 92.01% input cases were accurately predicted.

4) *Result analysis:* In this section we describe some result analysis with respect to the experiments conducted. Figure 3 shows a scatter plot of the eigensolver versus a portion of spectrum. Note that jacobi davidson (jd) is not a preferred solver for any smallest imaginary case and that generalized davidson (gd) is less likely to be selected as the most efficient solver for the largest portion of the spectrum. Similarly, plotting eigensolvers versus binary and non-binary matrices (where “binary” refers to a matrix with all non-zero values as one), we found that lanczos solver, which predominantly worked for Hermitian non-binary matrices, does not work at all for Hermitian binary matrices.

Plotting these parameters in graph form helps us infer some feature characteristics in two dimensions, whereas the decision tree captures all such results in higher dimensions (for every feature tested).

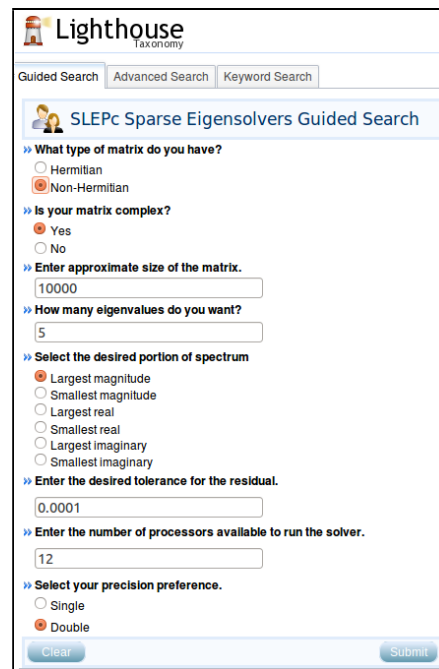


Figure 4. Guided search UI in Lighthouse for sparse eigensolver routines

C. User Interfaces

Once the collection of the data required for integrating SLEPc with Lighthouse is complete, we constructed a MySQL database through Django. The database schema is derived from the information on the decision tree.

We converted the classification tree information into a MySQL datatable using a depth first search algorithm [22]. Every path from the root node to the leaf node forms a single row in the datatable. The columns of the datatable constitute the attributes from the tree including matrix properties, such as matrix type and matrix order, as well as the desired output characteristics, such as the residual tolerance and the number of eigenvalues. The generated decision tree consists of 395 leaf nodes. Thus, the resulting datatable comprises 395 distinct rows with unique features. The user selects the desired features via the Lighthouse-SLEPc interface and the respective best solvers are fetched from the database.

Like the LAPACK UI in Lighthouse, the SLEPc UI also provides Guided Search. The questions correspond to the features in the generated decision tree. Each question is a fundamental query to the database that directs the search toward the most appropriate result. Upon the user’s answering all the questions, Lighthouse suggests the best known eigensolvers, and the user may choose to generate a respective code template in FORTRAN 90 or C. Figure 4 illustrates the Guided Search UI in Lighthouse for sparse eigensolver routines using the SLEPc package. Note that the interface consists of simple questions that lead the user to the desired SLEPc routines.

The Advanced Search implementation contains questions similar to Guided Search, but it delivers all the SLEPc eigensolver routines that are compatible with the user-specified problem (irrespective of the suggestions from the decision tree). The Advanced Search feature is for experienced users who are aware of the routines and would like to make their

own selections. Additionally, Advanced Search also includes some questions requiring extended eigensolver knowledge.

The Keyword Search is also provided so that users can directly search for the routine by name and generate the respective code template.

V. CONCLUSION AND FUTURE WORK

In this paper, we described our methods for easing the most significant task in solving a sparse eigenvalue problem. Our approach includes providing customized SLEPc eigenvalue solutions and generating the respective code through Lighthouse. The intelligent solver suggestions obtained through the machine learning techniques will offer more accurate and efficient solutions and relieve users from having to research the documents of every eigensolver routine available. The automatically generated code will help save time and effort in implementation, thereby improving the productivity of developers and scientists. It will also overcome the barriers caused by unfamiliarity of the programming language and implementation specifics, such as parallel programming.

In future work, we will continue to integrate more SLEPc functionality (e.g., singular value decompositions) into Lighthouse and enable all of the described search features. We believe that our effort will help in reaching out to a wider user base by providing easy access to computationally challenging sparse eigenvalue solutions, thereby further accelerating scientific development and discoveries.

ACKNOWLEDGMENT

This work is supported by the National Science Foundation (NSF) grants CCF-1219089 and CCF-091474. We used the Janus supercomputer, which is supported by the NSF award CNS-0821794 and the University of Colorado Boulder. The Janus supercomputer is a joint effort of the University of Colorado Boulder, the University of Colorado Denver and the National Center for Atmospheric Research. Janus is operated by the University of Colorado Boulder.

REFERENCES

- [1] R. H. Landau, J. Paez, and C. C. Bordeianu, *A survey of computational physics: introductory computational science*. Princeton University Press, 2011.
- [2] K.-J. Bathe and E. L. Wilson, "Solution methods for eigenvalue problems in structural mechanics," *International Journal for Numerical Methods in Engineering*, vol. 6, no. 2, 1973, pp. 213–226.
- [3] K. Bryan and T. Leise, "The \$25,000,000,000 eigenvector: The linear algebra behind Google," *Siam Review*, vol. 48, no. 3, 2006, pp. 569–581.
- [4] X. Ying and X. Wu, "Randomizing social networks: a spectrum preserving approach," in *Proceedings of the 8th SIAM Conference on Data Mining (SDM08)*, vol. 8. SIAM, 2008, pp. 739–750.
- [5] Lighthouse project. <https://code.google.com/p/lighthouse-taxonomy/>. [retrieved: May, 2014]
- [6] V. Hernandez, J. E. Roman, and V. Vidal, "SLEPc: A scalable and flexible toolkit for the solution of eigenvalue problems," *ACM Trans. Math. Software*, vol. 31, no. 3, 2005, pp. 351–362.
- [7] Basic Linear Algebra Subprograms (BLAS). [Online]. Available: <http://www.netlib.org/blas> [retrieved: May, 2014]
- [8] J. Dongarra. Freely available software for linear algebra on the web. <http://www.netlib.org/utk/people/JackDongarra/la-sw.html>. [retrieved: Jan, 2014]
- [9] NIST. Guide to Available Mathematical Software (GAMS). <http://gams.nist.gov>. [retrieved: Jan, 2014]
- [10] ——. HotGAMS: Java interface to Guide to Available Mathematical Software. [retrieved: Jan, 2008]
- [11] LAPACK - Linear Algebra PACKage. [Online]. Available: <http://www.netlib.org/lapack/> [retrieved: May, 2014]
- [12] PETSc - Portable, Extensible Toolkit for Scientific Computation. [Online]. Available: <http://www.mcs.anl.gov/petsc/> [retrieved: May, 2013]
- [13] Django. [Online]. Available: <https://www.djangoproject.com/> [retrieved: May, 2014]
- [14] J. G. Siek, I. Karlin, and E. R. Jessup, "Build to order linear algebra kernels," *Workshop on Performance Optimization of High-level Languages and Libraries*. POHLL, April 2008, pp. 1–8.
- [15] S. Balay, W. D. Gropp, L. C. McInnes, and B. F. Smith, "Efficient management of parallelism in object oriented numerical software libraries," in *Modern Software Tools in Scientific Computing*, E. Arge, A. M. Bruaset, and H. P. Langtangen, Eds. Birkhäuser Press, 1997, pp. 163–202.
- [16] Matrix Market. [Online]. Available: <http://math.nist.gov/MatrixMarket/> [retrieved: June, 2013]
- [17] The University of Florida Sparse Matrix Collection. [Online]. Available: <http://www.cise.ufl.edu/research/sparse/matrices/> [retrieved: Feb, 2014]
- [18] J. Han, M. Kamber, and J. Pei. *Data mining concepts and techniques*, third edition. Waltham, Mass. (2012)
- [19] MATLAB, Release R2013b. Natick, Massachusetts: The MathWorks Inc., 2013.
- [20] Orange. [Online]. Available: <http://orange.biolab.si/> [retrieved: May, 2014]
- [21] M. Hall, E. Frank, G. Holmes, B. Pfahringer, P. Reutemann, and I. H. Witten, "The WEKA data mining software: An update," *SIGKDD Explor. Newsl.*, vol. 11, no. 1, Nov. 2009, pp. 10–18. [Online]. Available: <http://doi.acm.org/10.1145/1656274.1656278>
- [22] S. S. Skiena, *The Algorithm Design Manual*. New York, NY, USA: Springer-Verlag New York, Inc., 1998.
- [23] V. Hernandez, J. E. Roman, A. Tomas, and V. Vidal, "A survey of software for sparse eigenvalue problems," *Universitat Politècnica de València, Tech. Rep. STR-6*, [retrieved: May, 2013]. [Online]. Available: <http://www.grycap.upv.es/slep>
- [24] C. Campos, J. E. Roman, E. Romero, and A. Tomas, "SLEPc users manual," *D. Sistemes Informàtics i Computació, Universitat Politècnica de València, Tech. Rep. DSIC-II/24/02 - Revision 3.3*, 2012.
- [25] V. Hernandez, J. E. Roman, A. Tomas, and V. Vidal, "Single vector iteration methods in SLEPc," *Universitat Politècnica de València, Tech. Rep. STR-2*, [retrieved: May, 2013]. [Online]. Available: <http://www.grycap.upv.es/slep>
- [26] —, "Arnoldi methods in SLEPc," *Universitat Politècnica de València, Tech. Rep. STR-4*, [retrieved: May, 2013]. [Online]. Available: <http://www.grycap.upv.es/slep>
- [27] —, "Lanczos methods in SLEPc," *Universitat Politècnica de València, Tech. Rep. STR-5*, [retrieved: May, 2013]. [Online]. Available: <http://www.grycap.upv.es/slep>
- [28] —, "Krylov-schur methods in SLEPc," *Universitat Politècnica de València, Tech. Rep. STR-7*, [retrieved: May, 2013]. [Online]. Available: <http://www.grycap.upv.es/slep>
- [29] —, "Davidson type subspace expansions for the linear eigenvalue problem," *Universitat Politècnica de València, Tech. Rep. STR-10*, [retrieved: May, 2013]. [Online]. Available: <http://www.grycap.upv.es/slep>
- [30] Y. Saad, *Numerical Methods for Large Eigenvalue Problems*. Halsted Press, NY, 1992.
- [31] L. N. Trefethen and D. Bau, *Numerical Linear Algebra*. SIAM, 1997.
- [32] N. S. Foundation and D. of Energy. BLAS. <http://www.netlib.org/blas/>. [Online]. Available: <http://www.netlib.org/blas/> [retrieved: May, 2013]
- [33] R. Nair, "Customized sparse eigenvalue solutions in Lighthouse," Master's thesis, University of Colorado Boulder, May 2014.

Analysis of the Utilization of Web 2.0 Resources in Secondary Education and Advanced Vocational Training Studies

Federico Banda -Sierra

Department of ICT Engineering University Alfonso X El Sabio. Villanueva de la Cañada. Madrid

fbandsie@myuax.com

Antonio J. Reinoso

Department of ICT Engineering University Alfonso X El Sabio. Villanueva de la Cañada.

areinpei@myuax.com

Abstract— This paper presents a study to address the use of Web 2.0 tools in education. The study focuses on secondary education and is limited to the region of Madrid (Spain). The main objectives of the study are to determine the involvement of both teachers and students in the use of Web 2.0 tools. Moreover, we will also examine the different parameters associated with the acceptance and application of Web 2.0 resources. Finally, we will describe the different patterns of interaction of both students and teachers with respect to these digital platforms. As a result, we will be in position of assessing the convenience and acceptance of artifacts based on the Web 2.0 paradigms to be used as effective resources to improve the students' learning experience.

Keywords-words: Web 2.0, Media Education, School 2.0, Knowledge Society.

I. INTRODUCTION.

The term Web 2.0 Websites comprises those Websites that facilitate information sharing interoperability, user-centered design and collaboration in the World Wide Web. There are several tools that allow to share information in the Web 2.0 such as blogs, wikis or social networks as well as other sites that share resources like Google Drive for documents, Youtube for videos, or Dropbox for online storage. Among these tools, there are also learning platforms like Moodle or Claroline and virtual classrooms such as Edmodo, Virtual Virtual Teacher or Tutor.

Our aims are to analyze the acceptance of Web 2.0 tools in the classroom and the involvement of teachers in their use. Therefore, the main objectives of the research activities reported in this paper are presented below:

1. To perform a study of the use of Web 2.0 tools in the classroom in the region of Madrid and to assess the implication of teachers in their use.
2. To evaluate the acceptance of Web 2.0 tools by teachers.
3. To determine whether the use of Web 2.0 tools is widely accepted by students (according to teachers' opinion).

4. To identify predictive factors for a lower use of Web 2.0 tools by teachers in the classroom.

The rest of the paper is structured as described in the following. First, we will present the state-of-the-art after reviewing the previous research on the use of Web 2.0 tools by teachers. Then, the methodology used to carry out our study will be described. Finally, our main results and conclusions together with some threads to validity will be exposed.

A Web 2.0 site allows users to interact and collaborate to join efforts in the process of creating contents that will remain available for a given virtual community. This approach radically differs from static websites where users are limited to the passive viewing of content that have been created for them.

The main aspects featuring the Web 2.0 technologies are enumerated in the following [9]:

1. The term Web 2.0 groups different interactive approaches that have in mind the social component of the net. Provided with these new mechanisms, Internet has acquire a more participative dimension, enabling the exchange of information and several forms of contact among users through individual blogs, wikis, social network sites (Second Life, Facebook, etc), image or video sharing sites (Flickr, YouTube, etc.)
2. The Web 2.0 is and attitude and not precisely a technology. It's most relevant aspect is that it makes clear the next evolutionary steps of Internet.

Facebook is an example of typical social networking projects that, in general, has been developed using the figure of social networks as the main support for students' learning process. [11]

Free Software is any software respectful to the freedom of its users and to the social solidarity of their communities. It is common to associate Free Software to software with no cost, but this is a mistake, because there are free software that you have to pay for using it. Not every free software is free. And above all not every no cost software is free software. [7]

To sum up, free software is defined by the four freedoms of software users.

Freedom 0: The freedom to run the program for any purpose.

Freedom 1: The freedom to study how the program works, and to change it for any particular purpose.

Freedom 2: The freedom to redistribute copies.

Freedom 3: The freedom to improve the program, and release your improvements (and modified versions in general) to the public, so that the whole community benefits.

Also you should have the freedom to make modifications and use them in your work or free time privately, without even have to announce that this modifications exists [8].

We will take as a reference some figures in the education in Europe and in Spain. That will show us how important can be the use of the web 2.0 tools for the classroom. All of this will give us an idea of the current situation about this topic, in Europe and in Spain [10].

II. STATE OF THE ART

The Society of Information has risen as the result of the implementation of information and communication technologies (ICT) in daily life. The greater use of the ICT paradigms has changed, in many senses, the way of developing many of the activities of the modern society.

The Knowledge Society arises in the context of the Society of Information. The main characteristics of the Knowledge Society are the open access to information, the freedom of expression and the linguistic diversity.

Therefore, the fields of study of our research are Education, Communication and Web 2.0 supporting tools. We considered that it was important to investigate these fields due to the importance that the Regional Administration was giving to the introduction of ICT, particularly the Web 2.0 tools, in education. The incorporation of Web 2.0 tools in teaching is having a great impact in education and the use of these tools is expected to increase according to the principles of the Media Education and School 2.0, which will be explained later in this summary.

In the Society of Information, a change in "what" we teach and "how" we teach is needed. This new Society forces us to rethink the role of teachers, students and taught subjects in education. Therefore, teachers should try to encourage collaborative works in the classroom, which can be supported by mechanisms such as the Web 2.0 tools.

To the best of our knowledge, there are no studies addressing the particular use of Web 2.0 resources in Spanish secondary education. Therefore, our research study covers the lack of information regarding two relevant points: The Media Education and School 2.0.

With respect to the Media Education, it is noteworthy that skills such as the ones related to the use of many languages, to critical thinking, and to the interaction with others in real and/or virtual modes, should be potentiated [1].

Concerning the Education 2.0, it is noteworthy that the advances occurred in education based on the communication model 2.0, i.e., the Web 2.0, is based, inter alia, on the use of social networks. The active interactive joint of teachers and students is also closely related to collaborative learning and authorship for example blogs and wikis.

Education and Education 2.0 are related through the emergence of the communicative model 2.0 in which both teachers and students are active participants in the communicative act, that is, a continuous exchange of roles where both teachers and students can be authors or co-authors of the information and knowledge. According to Aparici et al. [1], the Web 2.0 has changed the rules of the game and allows Internet to contribute collaboratively in the construction of collective knowledge from individual acts of group communication that can occur in cyberspace and real space.

The receiver and the transmitter of the information are transformed into an EMIREC which sends and receives messages [2]. In this model, the role of the students substantially changes. As the learning focuses on the students, they should take an active role building the teaching-learning process. It is expected that the students are able to exercise autonomy, to develop critical thinking, to adopt collaborative attitudes and to use the theoretical knowledge to solve real problems. This model of communication enhances the possibility that the students cast its own messages through different languages, strengthening the educational proposal of personalized education, as well as affective and emotional processes that occur in all educational relationship.

In addition, the digital literacy is a new concept that should be taken into account when considering the training of teachers in the Media School.

Literacy is the education that everyone needs to live in society. A new model of literacy is needed in the digital society: the digital literacy. Furthermore, this new literacy, defined as the basic capacity to understand and to express in different languages and meanings, is constructed collaboratively between teachers and students, as defined in [4]

What we teach and the way we teach it should be changed. Therefore, the role of the teachers, students and contents in learning should be redefined.

Currently, we are part of a networked society, mainly due to the Internet, that has changed the processes, interests, values and social institutions in the way we knew them. Internet has propitiated new ways of relationship, that wouldn't have been possible without the tis technological advance As indicated in [5], Internet is the heart of a new paradigm that currently constitutes the base of our lives and of our ways of relationship, work and communication. Internet processes the virtuality and transforms it into our reality, constituting the network society, which is the society where we live.

We cannot forget the importance of TIC in the school. As specified in [6], The implementation of TIC needs several capacities such as cooperation, capacity of initiative and dynamism in the working places, to be able to work in groups, interactive learning between the members of the group, communication, to be able to work with abstract concepts, to identify and solve problems, aptitude to make decisions, being able to seek and to use the information, predisposition for the permanent training and other.

The appearance of these technologies demands to train teachers, as most of them are not familiar with TIC resources and lack the necessary skills to use them.

The immense majority of the administrators, educational managers and teachers are digital immigrants. (...). These professionals rarely use the digital technologies and are very resistant to modify their conception of the world of the work based on the industrial models of ends of the 19th century [1].

III. METHODOLOGY

In the following, we describe the methodology conducted to perform our study. First of all, this is an observational study in which teachers from secondary schools in the region of Madrid have been electronically surveyed . No exclusion criteria were established.

Data Collection

To achieve the objectives of the study a questionnaire asking about different aspects, such as the professional profile, the opinion on the usefulness of Web 2.0 tools for teaching , the involvement of teachers in the use of Web 2.0 tools and their perception of student interest in the use of these tools, was designed (Figure 1).

Teachers from different schools of the Region of Madrid (both, public and private) were invited to fulfill the questionnaire. The questionnaire will consist of an HTML form heading to a PHP page which will register received

data received into a MySQL database. Later, every data will be exported to the SPSS statistic package.

Statistical analysis

For the descriptive analysis of quantitative variables the mean was calculated. Qualitative variables were expressed as percentages. Comparisons between means were performed using Student's T test for independent samples. Qualitative variables were compared using the chi-square (χ^2) test and the Fisher's exact test. Statistical significance was considered at $p < 0.05$ for the all the comparisons.

A binary logistic regression model was used to estimate the effect of the different variables on the use of Web 2.0 tools. All the variables which reached the statistical significance in the univariate analysis were included in the multivariate analysis. Thus, in the multivariate analysis, the dependent variable was the use of Web 2.0 tools.

The questions asked in the survey were selected considering the objectives of the study and their answers are coded. The data were analyzed with the statistical package SPSS. Questions 1, 2 and 3 aims to determine the most frequently used Web 2.0 tool, while questions 4, 5 and 6 are intended to characterize the use of Web 2.0 tools in the classroom. Finally, questions 7 and 8 will provide different information about the use of the tools. For example, whether they are being used, or not, for individual or teamwork purpose,

Number	Question
1	Are involved in the use of Web 2.0 tools in the classroom ?
2	Which are the Web 2.0 tools that you usually use? Why?
3	Do you think that the use of Web 2.0 tools in the classroom promotes the student learning ? Why?
4	Are students involved in the use of Web 2.0 tools? To work in group or to work individually?
5	Do you recommend the students to use Web 2.0 tools or do you think that they are not necessary for learning?
6	In what subject are you using Web 2.0 tools?
7	When do you use Web 2.0 tools? In group or in individual activities?
8	Do you use Web 2.0 tools in the classroom regularly or sporadically?

Figure 1: The questionnaire fulfilled by the teachers is shown above

IV. RESULTS

Eighty-one teachers have fulfilled the questionnaire so far. The main characteristics of the teachers included are summarized in table 1. Forty-one (51%) were male, 34

(42%) over 50 years and 19 (23,5%) with more than 30 years of experience. The majority of them were officials and 75% worked in public centers. Approximately half of them (54%) were teachers of Computing Science.

Although most of the surveyed teachers (86%) declared that the use of Web 2.0 tools reinforces and benefits the students' learning process, only 63% of them acknowledged using these tools frequently. In addition, only 33% of the teachers being using Web 2.0 tools manifested to use them during whole sessions in the classroom. When focusing on teachers that used the aforementioned tools, only 59% of them did it to work in group.

On the other hand, with respect to the interest shown by the students in the use of Web 2.0 in the classroom, the majority of teachers that used them (79%) considered that students are very enthusiastic about the use of these tools. A small proportion of teachers (6%) answered that the implication in the use of these tools varies widely and it depend on the characteristic of the student. In teachers' opinion, the students that used Web 2.0 tools seemed to do it for both individual and group tasks.

In the multivariate analysis, to be teacher of Computing Science was the only variable independently associated with the use of Web 2.0 tools in the classroom (odds ratio=7.9, confidence interval 95%=2.5-24.5).

Up to 20% of teachers consider that Web 2.0 tools are useful for teaching. However, the majority of them do not use these tools probably because they are not accustomed to these technologies, they have no training or because they do not have enough resources.

Only one third of the teachers that use the Web 2.0 tools do it throughout the whole class.

The majority of teachers that use Web 2.0 tools do it for working in group. This result should be underlined as suggests that teachers know the collaborative character of these tools.

Technological elements always seems attractive to the students (mainly at that age) and that should be use to introduce them in the classroom. Many teachers think that the use of these elements contribute to arouse interest of the students. In addition, these tools offer many facilities to promote and participate in group activities although obviously there are ways of interaction requiring presentality such as debates, exhibitions, etc.

TABLE 1 RESUMES THE MAIN CHARACTERISTICS OF THE TEACHERS THAT FULFILLED THE QUESTIONNAIRE.

Feature	Categories	% of participants
Gender	Male	40 (49%)
	Female	41 (51%)
Age	Over 50 years	34 (42%)
	Under 50 years	47 (58%)
Experience	Over 30 years	19 (23,5%)
	Under 30 years	72 (76,5%)
Type	Officials	55 (68%)
	Interins	26 (32%)
Centers	Public	61 (75%)
	Private	20 (25%)
Area	ICT	44 (54%)
	Non-ICT	37 (46%)

V. CONCLUSIONS

Our results suggested that although most teachers believe that the use of Web 2.0 tools promotes the students learning, only half of them use these tools for teaching and only one third acknowledge to be implementing actively the use of these tools in the classroom. In general, teachers think that the use of Web 2.0 tools in the classroom is well accepted by the students.

To teach Computing Science was the only factor independently associated with a greater use of these tools in the classroom (8-fold higher than among teachers in other specialties). The age of the teacher (older age) is associated with a lower use of Web 2.0 tools.

This study underscores the need for action to generalize the use of Web 2.0 tools in the classroom, with particular emphasis on the older teachers and among teachers from other specialties from Computing Sciences. Future studies should aim to identify the handicaps for the implementation of these tools in order to adopt strategies (procurement, training courses, etc.) to avoid them.

REFERENCES

1. Aparici Marino, R, Campuzano Canales, A, Antonio Ferres, J and García Matilla, A. Madrid. UNED. Julio 2010. *La Educación Mediática en la Escuela 2.0*
2. Clotier Williams, J (1975). *L'ère D'emerec ou la communication audio-ecritovisuelle a' l'heure des self-media* (segunda ed). Montreal: Les Pres de L'Universit e de Montreal.
3. Callejo Gallego, J and Viedma Rojas, A (2006). *Proyectos y Estrategias de Investigacion Social*. UNED. 2005.
4. Guti errez Moraleda, A (2003). *Alfabetizaci n Digital*. Fundaci n AUNA. http://fundacionorange.es/areas/28_observatorio/pdfs/JM_C_8.pdf y <http://www.revistacomunicar.com/pdf/preprint/38/03-PRE-13396.pdf>.
Retrieval date 20:16, June 01, 2014 de.
5. Castells Olivar, M (2009). *Comunicaci n y poder*. Madrid: Alianza999
6. Martin Artilles, A (2005). *Cambios en la organizaci n del trabajo, formaci n, cualificaci n y competencias*. <http://eytsc1011.pbworks.com/f/OTCOMPETENCIAS2004.pdf>. Retrieval date 20:16, June 01, 2014 de.
Feltre Oreja, R. (2006). *La Filosof a del Software Libre*. Vol. I: las Licencias de Software Libre y su desaf o a los modelos vigentes de Propiedad Intelectual. Madrid. UNED. Retrieval date 20:33, June 01, 2014 de.
7. Abella Garcia A and Segovia Romero, M. A. (2003). *Junta de Extremadura. Libro Blanco del Software Libre en Espa a (II)*.
Vicente Hinojosa, A and Le n Mart nez J (2009). *Conceptos sobre el Uso de las Tecnolog as Web 2.0*. (2013).<http://repositoral.cuaed.unam.mx:8080/jspui/bitstream/123456789/813/1/EI%20Uso%20de%20la%20Tecnolog%C3%ADa%20en%20la%20web%202.0.pdf>
Retrieval date 20:45, June 01, 2014 de.
8. Vicente Hinojosa, A and Le n Mart nez J (2009). *Conceptos sobre el Uso de las Tecnolog as Web 2.0*. (2013).<http://repositoral.cuaed.unam.mx:8080/jspui/bitstream/123456789/813/1/EI%20Uso%20de%20la%20Tecnolog%C3%ADa%20en%20la%20web%202.0.pdf>
Retrieval date 21:02, June 01, 2014 de.
9. Instituto Nacional de Tecnolog as Educativas y de Formaci n del Profesorado. *Programa Escuela 2.0* (2013). <http://www.ite.educacion.es/escuela-20>. <http://www.ite.educacion.es>. <http://recursostic.educacion.es/blogs/europa/> Retrieval date 21:10, June 01, 2014 de.
10. Piscitelli, A, Adaime, I, Arribas Urrutia, A, Balestrini, M, Ciuffoli, C, Cobo Romani, I, Gruffat, C, Islas, O, Leal Fonseca, D.E, L pez, G, Primavera, H, Reig Hern andez, D, Schimkus, R, Sued, G, Uman, I, Venesio, M. C.
11. *El Proyecto Facebook y la Posuniversidad. Sistemas operativos Sociales y Entornos Abiertos de aprendizaje* (2010).http://www.fundacion.telefonica.com/es/que_hacemos/convocatorias/sociedad_informacion/2010/mayo/pdf/EVEN_DYC_ARG_EI%20proyecto%20Facebook%20y%20la%20posuniversidad_07_05_10.pdf
Retrieval date 21:15, June 01, 2014 de.

Twitter Usage of German Online Retailers

Georg Lackermair, Daniel Kailer
 Munich University of Applied Sciences
 Department of Computer Science and Mathematics
 Munich, Germany
 Email: {georg.lackermair, dkailer}@hm.edu

Abstract—In this paper, we investigate the role of Twitter for German online retailers based on an empirical study. We analyze Tweets of the best selling German online retailers with qualitative and quantitative methods. Based on our results, we derive a conceptual model that can be used to classify different interaction strategies for online retailers. We identify four different interaction strategies. Three strategies are based on social interactions either on Twitter or by redirecting users to other social networks. The fourth strategy is applied to promote products of the retailers' online store. Another result is that accounts which apply a user interaction focused strategy inside of Twitter through User Mentions, have a significantly higher number of followers compared to the accounts that use other strategies.

Keywords—E-commerce; Twitter; Study; Social Web.

I. INTRODUCTION

E-commerce plays an ever increasing role in the German economy. With the rise of the social web, the importance of social interactions with customers has also increased for online retailers. Traditional tools to reach customers include for example discussion boards, weblogs or newsfeeds. But in recent years, more and more online retailers have turned to the microblogging platform Twitter in order to reach out to new or existing customers.

A lot of research is conducted to understand the usage of Twitter. For example, the study of Java et al. investigated the user intentions and community structures of Twitter [1]. But to our best knowledge, there is no study that explicitly investigated the Twitter usage of online retailers. In this paper, we contribute knowledge to this area of research by presenting an empirical study that investigates the Twitter usage of German B2C online retailers.

In our study, we applied qualitative and quantitative methods to answer the following research questions:

- *RQ1*: How many online retailers are actually using Twitter and how active are they?
- *RQ2*: How can the Twitter interaction strategy of online retailers be classified?

The results to research question RQ1 will show the acceptance of Twitter as a tool to reach new and existing customers. Secondly, we answer research question RQ2 by analyzing Tweets in regard to different Twitter entities, e.g., URLs or User Mentions. We will present a conceptual model that can be applied to categorize the used Twitter interaction strategies of the conducted online stores.

The remainder of this paper is organized as follows. In Section II, we present the theoretical background. We then

explain the design of our empirical study in Section III. After that, we show the results of the study in Section IV and discuss them in Section V. Finally, the article concludes with an outlook for future research.

II. THEORETICAL BACKGROUND

The main research areas addressed in this paper comprise the social web and e-commerce. The combination of social media and e-commerce are often denoted as *Social Commerce* [2], [3]. Most of the studies in social commerce investigate the customers' perspective to the platform Twitter. In this study, we will focus on the retailers' perspective instead.

Twitter is the most popular microblogging service in the web. Every day, about 500 million Tweets are published on this platform. Reasons for this success can be found in its simplicity, scalability, ubiquity, and interactivity. Due to its publish/subscribe capabilities, more and more users shift from traditional newsfeeds based on RSS and Atom to Twitter.

This publish/subscribe capability is Twitter's fundamental pattern: Users subscribe either to other users or to *Hashtags* (*HTs*). The use of HTs is a communication convention, that enables authors to post a message either to a community or to add a content information [4]. Other conventions are *User Mentions* (*UMs*) (denoted by an @-sign) and Retweets [5]. From the Twitter API's view, Hashtags, User Mentions and URLs are treated as special entities. The use of those entities are examined by [6] in a large scale study.

Previous works have investigated that Twitter is used for various communication purposes [7], [4]. Twitter entries can be classified in five different genres [8]: (a) Personal Updates, (b) Directed Dialog, (c) Real-time Sharing, (d) Business Broadcasting and (e) Information Seeking. These communication genres are discussed below.

The genre of personal updates contains mostly "daily chatter" [1], which means information about, what a user is thinking or doing at the moment. According to [1] this is the most common use of Twitter. According to [8] those posts are mainly issued by sparsely connected users.

A directed dialog is a message that is directed to a certain user by the use of an @-sign (User Mention). Measures to detect UMs are not exact, since there are other conventions that use the @-sign as well, but in Twitter traffic the vast majority of occurrences are UMs [9]. According to [10, p.21-25], UMs are an indicator about the interactivity of online relationships in Twitter.

In regard to business broadcasting, it is interesting to classify and measure the influence of a user. Mainly based

on the ratio between followers and friends, Twitter users can be categorized as broadcasters, acquaintances, miscreants and evangelists [11]. The number of posts as a function of the number of friends is increasing without any saturation [12]. Nonetheless, the classification of miscreants/spammers is controversial, as this group would also include users, that use Twitter primarily as newsfeed. The measurement of influence was refined by [13] and adopted to take Twitter's characteristic features into consideration. The indegree influence is the number of followers of a user. This is the size of audience a user can reach, without directing a post via HTs to further channels. A Retweet influence indicates the ability to generate content, that is redistributed by the followers. The User Mention influence indicates the number of User Mentions containing the name of the examined user, which indicates the ability to engage the audience in a conversation.

The information seeking behavior is examined in [14]. They gathered and analyzed data for three different question types. They found out that the most popular type of questions are rhetorical questions with an overall very low response rate.

Another interesting aspect about the communication on Twitter are URLs that are embedded in Tweets. URLs are by default shortened by the platform's own shortening service <http://t.co> [15]. The targeted URLs can be categorized as self-links, social media links and other external links. For our study, the former two contain interesting information. A *self-link* points to the own website of an online retailer and indicates the promotion of a product. *Social media* URLs direct users to discussions on other social networks, e.g., Facebook. This indicates a more community-centric activity than links to product pages.

III. STUDY DESIGN

In this section, we will present the underlying design of our empirical study.

A. Data Collection

To acquire a sample of e-commerce related communications, we used a list of the 115 best-selling online-retailers in Germany [16]. We matched their Twitter accounts by querying search engines. If we were not able to find an account, we tried to collect the information manually by examining the online retailer's website. We then evaluated manually, whether their Twitter account is targeting the German market. For this purpose we checked the profiles' descriptions and timelines.

Based on the found Twitter accounts, we collected two different datasets. First, we retrieved common account information for the Twitter account, e.g., the lifetime of the account, number of followers, number of friends, number of favorites and the number of status updates issued by the account. Second, for each account we collected the timeline consisting of last 100 status updates. We used the Twitter REST API to access this data.

For the analysis of Tweets, we queried the REST API for the last 100 status updates issued by the acquired accounts (see Section IV-A). From this sample we extracted a smaller subset of 200 Tweets from 10 different retailers, which was used to identify genres by a manually performed content analysis. The data collection was carried out on 14th of February, 2014.

B. Analysis

First, we analyzed the data set containing the common information of the accounts. The lifetime of an account in days is defined as L . To calculate the Tweet rate R_T per lifetime as an indicator for broadcasting activity, we use the total number of Tweets since the creation of the account (T) and relate it to L (see (1)).

$$R_T = \frac{T}{L} \quad (1)$$

In order to analyze the links to other users, the number of followers f_{in} , the number of friends f_{out} , and the listed count l_{in} can be used. As stated in Section II, those values are considered as in- and outdegree measures. To reflect the lifetime of an account, we relate those values to L and, thus, define indegree rate R_{in} and R_{out} in (2) and (3).

$$R_{in} = \frac{f_{in} + l_{in}}{L} \quad (2)$$

$$R_{out} = \frac{f_{out}}{L} \quad (3)$$

We used the timeline data set to analyze the data set quantitatively to derive information about the communication strategy implemented by a retailer. We were particularly interested in distinguishing between the purposes to engage users in a dialog or to promote products directly. User Mentions and URLs are indicators for these purposes. A UM indicates the attempt to start a conversation with a customer or to reply to a post issued by a user. Depending on the target, a URL serves as an indicator for both categories: A link to a site of the retailers' online shop follows the purpose of promoting a product, whereas a link to the companies blog or a social network is clearly a community management activity. To enrich the coarse grained strategic categories, we analyzed the sample of Tweets qualitatively to identify e-commerce related Tweet genres, that represent common interaction activities between retailers and consumers.

IV. RESULTS

In this section, we will present the results of our study. First, we present a common analysis of the investigated accounts. Second, we describe the found interaction strategies of the investigated accounts and illustrate those strategies with Tweet genres.

A. Analysis of accounts

First, we have identified the accounts that are managed by the online retailers. Figure 1 shows the results of the identification and retrieval process. As bar (a) shows, we were able to identify 88 accounts for our list of 115 retailers. The removal of accounts not targeting German users is represented by bar (b). We excluded 11 accounts from further retrieval. Bar (c) visualizes the retrieval of account information, which depend on the accounts settings. We successfully retrieved the information for 76 shops. As bar (d) shows, 59 of those accounts were classified as being active in the last month.

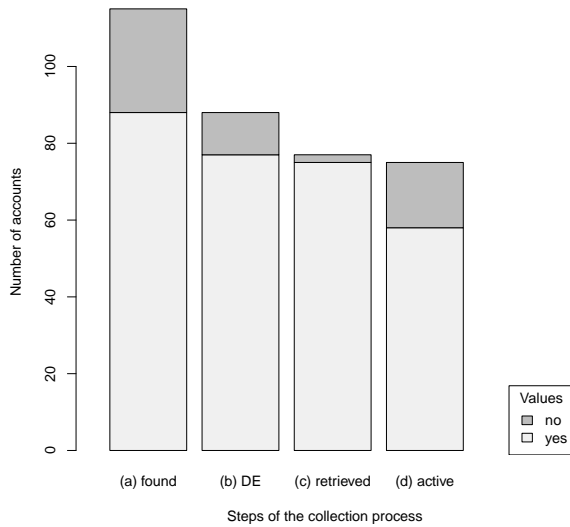


Figure 1. Identification and retrieval of Online-Retailers' Twitter accounts.

Most of the account names reflect the shop's name, which can consist of the company's name or a brand name. Besides that, many ids also include regional attribute, which can be a country name, a region name or a country code. The latter could also reflect the shop's URL.

The variables lifetime, status updates per day, indegree and outdegree derived from the profile information are summarized in Table I. We calculated the .25, .50 and .75 quantiles and added Geary's skewness indicator. The mean account lifetime is slightly above 4 years, while the values are slightly left-skewed. Since its creation, an account issued on average about seven Tweets per day, whereas the data is strongly right-skewed. The majority of accounts issued less than two Tweets per day. The mean increase of indegree is about 1.6 followers per day, while the data is right-skewed. For 75% of the retailers, this value is at about 1.5 or less. The outdegree value for 75% of the accounts is at about 0.6 or less.

B. Twitter interaction strategies

To determine the Twitter interaction strategies of the conducted accounts, we first define the following sets:

- A : All Twitter accounts whereas each element a represents an online retailer from the sample.
- T : All observed Tweets.
- T_M : All Tweets that address other Twitter users, i.e., Tweets that contain at least one User Mention.
- T_S : All Tweets that contain at least one URL that is referring to another social network, e.g., Facebook.
- T_P : All Tweets that contain at least one URL that is referring to the online store, e.g., URLs that link to certain products or special offers.

Based on the Tweets that contain User Mentions (T_M), URLs to social networks (T_S) and URLs to the online store

of the account owner (T_P), we derived our conceptual model as depicted in Figure 2. As shown in the figure, we identified four different strategies, which will be explained below.

The first strategy S_1 is characterized by a frequent communication with other Twitter users. Accounts that apply this strategy (A_{S_1}) make use of User Mentions in at least two-thirds of their Tweets (see (4)).

$$A_{S_1} = \left\{ a \mid a \in A; \frac{T_{M_a}}{T_a} \geq 0.66 \right\} \quad (4)$$

Strategy S_2 is applied by accounts that intend to direct Twitter users either to the weblog of the company or to a website of another social network (e.g., Facebook) to continue communication there. We assigned accounts to this strategy when at least two-thirds of their Tweets contain URLs to other social networks or to a company weblog (see (5)).

$$A_{S_2} = \left\{ a \mid a \in A; \frac{T_{S_a}}{T_a} \geq 0.66 \right\} \quad (5)$$

Strategy S_3 is categorized by accounts that use strategy S_1 and S_2 moderately, i.e., accounts that make moderate use of User Mentions and moderate use of URLs to other social networks or weblogs. Moderate use means that User Mentions and social network URLs are present in at least one third, but no more than two-thirds of the Tweets (see (6)).

$$A_{S_3} = \left\{ a \mid a \in A; 0.33 \leq \frac{T_{M_a}}{T_a} < 0.66 \right\} \cap \left\{ a \mid a \in A; 0.33 \leq \frac{T_{S_a}}{T_a} < 0.66 \right\} \quad (6)$$

The last strategy S_4 is based on the URLs in Tweets that refer to the website of the online retailer. An account is using this strategy, when at least two-thirds of the Tweets contain URLs to the retailers online store, i.e., URLs for promotional purposes (see (7)).

$$A_{S_4} = \left\{ a \mid a \in A; \frac{T_{P_a}}{T_a} \geq 0.66 \right\} \quad (7)$$

Finally, all accounts that did not fit into the four above strategies were classified as having no clear strategy (see (8)).

$$A_{S_x} = A \setminus (A_{S_1} \cup A_{S_2} \cup A_{S_3} \cup A_{S_4}) \quad (8)$$

A short summary of the above mentioned strategies, their relative occurrences and their average number of followers are displayed in Table II. As shown in the table, strategies S_1 to S_3 can be classified as an interactive strategy, while strategy S_4 can be best described as a promotional strategy.

Based on these interaction strategies, we have qualitatively identified Tweet genres, that represent various communication purposes in the e-commerce. To perform that task, we have analyzed the timeline containing the last 20 Tweets of 10 accounts ($n=200$). We illustrate each genre with synthetic

TABLE I. VARIABLES CHARACTERIZING THE RETAILERS' TWITTER PROFILES.

	<i>min</i>	<i>Q</i> _{.25}	<i>Q</i> _{.50}	<i>mean</i>	<i>Q</i> _{.75}	<i>max</i>	<i>skewness</i>
Lifetime in years (<i>L</i>)	0.893	3.748	4.563	4.139	4.822	5.777	-1.350
Tweet rate (<i>R_T</i>)	0.057	0.418	1.208	7.711	1.894	367.300	7.345
Indegree rate (<i>R_{in}</i>)	0.055	0.482	0.987	1.575	1.471	8.258	2.209
Outdegree rate (<i>R_{out}</i>)	0.004	0.068	0.198	0.620	0.629	8.008	4.314

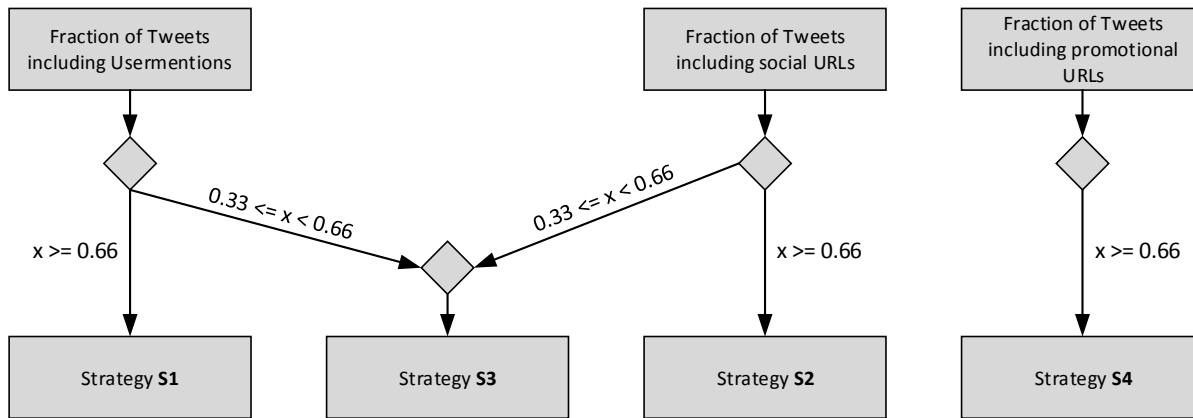


Figure 2. Conceptual model derived from our exploratory research.

TABLE II. SUMMARY OF THE IDENTIFIED TWITTER INTERACTION STRATEGIES.

Strategy	Description	Occurrences	Ø followers
<i>S</i> ₁	High interaction with other Twitter users through User Mentions	19%	6048
<i>S</i> ₂	Redirecting Twitter users to other social network sites (e.g., Facebook) through URLs	24%	1175
<i>S</i> ₃	Moderate interaction through User Mentions and moderate redirection to social networks	5%	3097
<i>S</i> ₄	Posting URLs to promote products of the own online store	21%	2398
<i>S</i> _x	No clear strategy could be determined	31%	1994

TABLE III. GENRES AND SYNTHETIC EXAMPLES FROM OUR QUALITATIVE ANALYSIS.

Genre	Example(s)
News	1 "Today would be Mozart's 258th birthday"
	2 "Report: Namm Show News 2014 http://shorten.er/URL"
Review	3 "RT @user1: very fast delivery by @retailerA"
	4 "RT @user2: testing guitar manufactured by #ESP http://shortener/URL @retailerA"
	5 "Test: Steinberg Cubase 7.5, Digitale Audio Workstation http://shorten.er/URL"
Promotion	6 "under #offer: #TOUCHLET X7Gs Tablet-PC for only EUR 89,90! http://shorten.er/URL"
Service	7 "@user3 Our customer service help@retailerC sure has an idea. They will contact you."
	8 "@user4 We've ordered 20 pieces in November, but we still don't have a delivery date"
Dialog	9 "@user5 What exactly are you planning? #renovate ^cm"
	10 "@user5 Wow, that's quite a project! How many squaremeters? ^cm"
	11 "@user5 You'll certainly cope with it! ^cm"
Recruiting	12 "#retailerD Do you need a job? Apply now! http://shorten.er/URL"
Request for collaboration	13 "To Celebrate #NationalPuzzleDay we are giving away £100 of puzzles for chance to win just RT&Follow ends 6pm 31/01/14"
Agent	14 "@user5 You'll certainly cope with it! ^cm"
	15 "@user6 Please send an email with your customer number to twitter@retailerE.de. We are looking forward to help you soon as possible.~Tom"
	16 "#3: coupon book: http://shorten.er/URL"
	17 "@user7 I need to ask a colleague. I'll get back to you (ar)"

examples in Table III. Because Twitter's terms of service prevents republication of Tweets, we needed to create synthetic examples of our sample for publication. Those Tweet genres are explained below.

News: A Tweet that contains news that are related to the business domain of the retailer. Both examples for this category are published by a retailer for musical instruments. Both contain news that might be interesting for their customers. Since the Tweet contains neither HTs nor UMs, we assume that the Tweet is intended to provide an additional incentive to follow the account.

Review: A Tweet transporting a review for a product, service or overall quality of service. The issuer of the Tweet could be a customer or the retailer itself. The Tweet could be self-contained or could link to media located on the retailers' site or on external media. Examples 3 and 4 are retweeted by the company and contain a review generated by a user.

Promotion: The intention of these Tweets are to promote products or special offers. The Tweets might contain an URL linking to product page and the respective price of that product. Both attributes are contained in the corresponding examples.

Service: Tweets that answer questions posted by customers either before or after purchasing a product. Example 7 indicates that the retailer answers to a message of a customer after purchasing a technical device. Example 8 is also an answer to a directed post by a user, but in contrast to the former the user was asking about the availability of an article that was sold-out at this time.

Dialog: An retailer's agent engages in a dialog in requesting more information, e.g., about a project. The corresponding example shows a dialog that consists of 5 entries and which was initiated by customer with a directed Tweet.

Recruiting: Companies try to promote job offers online. In the corresponding example, a URL linking to a more detailed job description is contained. To increase the audience reached by such an advertisement, usually HTs are added. In example 12, the author misleadingly tags the company's name, which is not very likely to increase the number of direct followers.

Request for collaboration: In order to enhance the attention and visibility, companies try to motivate users to collaborate in retweeting posts issued by a retailer. In example 13 the company uses a competition to motivate users to retweet.

Tagging agent: Adding an attribute, which identifies the person issuing a Tweet on behalf of a company. The corresponding examples illustrate three different variations that were found in the sample. Examples 14 and 15 use the circumflex sign to mark this attribute, whereas the former uses an abbreviation. In example 16, the issuer is identified by a numeric id, which is annotated with a hash notation. The example shows an abbreviated name in brackets.

V. DISCUSSION

In Section IV-A, we analyzed the existence and activity of online retailers on Twitter. In our sample, about 77% of the retailers maintain a Twitter account and about 78% of these accounts were active at the time of retrieval. According to our data, the majority of the accounts were created since 2010, which means that the adoption in German e-commerce started roughly two years after the dramatic growth in popularity in the U.S. in 2008 [9]. This reflects the slower adoption of the German population [6], [17], [18]. Besides that, we observed that the indegrees of those accounts are higher than the outdegrees. Which means that online retailers do not tend to follow back each of their followers, to improve the reputation of their customers. This characteristic is consistent with the category of *broadcasters* defined in [11]. Surprisingly, there are also outliers that do follow back all of their followers.

For the identification of different communication strategies, we have analyzed the use of URLs and UMs in Section IV-B. A limitation of this part of our work is, that we analyzed only the locations of the URLs and not the content behind it. We can't rule out the possibility that, for example, a Facebook URL points to a post, that contains a link to a product promotion. This was not further examined in our study.

We found out that 31% of the accounts did not follow a clear communication strategy. These accounts are not regularly involved in user interactions and are not using Twitter often to promote products. The other 69% of the accounts were

following a specific strategy according to our conceptual model (see Figure 2). Although our model allows the occurrence of multiple strategies per account, we did not find accounts that actually applied more than one strategy. This shows that our defined strategies are disjunctive and clearly separated from each other.

We found that Twitter accounts with a frequent user interaction through User Mentions, have a significantly higher number of followers. On the other hand, Twitter accounts that only post URLs to redirect users to other social networks (e.g., Facebook), have the lowest number of followers. These two results leave room for interpretation. One possible explanation is that a direct contact to customers through User Mentions is a motivator for users to follow the online retailers' account. An interpretation for the low follower count of accounts that only redirect users could be that the targeted users are already connected to the online retailer in another social network (e.g., they are friends in Facebook). Because the interaction happens in another social network, it could be of less value for a customer to follow an online retailer on Twitter.

Besides that, we derived genres of Tweets from a subset of our sample, to describe the coarse grained strategies in detail. Those genres build a tool for classification and analysis of direct interaction between companies in the e-commerce and users on Twitter. Due to the exploratory character of this study and the relatively small data set, we did not quantify the frequency of those genres. Our results could build the conceptual base for such a study.

VI. CONCLUSION AND FUTURE WORK

In this paper, we made three contributions. First, we have determined the share of retailers, that maintain an active account on Twitter. Second, we have identified genres and communication patterns used in the e-commerce for direct communication with users. Those categories can be used as an analytic framework for studying microblogging in the Social Commerce. Third, we have derived basic interaction strategies from our data set and created a conceptual model based on these strategies. We also showed the usage of User Mentions and URLs and their purpose inside of Tweets.

In our future work, we primarily aim at working on the limitation of our results, as stated in Section V. An important aspect will be the validation and refinement of our approach to classify the communication strategies as mentioned in Section IV-B. We aim to perform an in-depth content analysis on a subset of our data to evaluate our result. Besides that, we plan to collect and label a larger data set and quantify the use of the Tweet categories defined in this article. Another limitation of our approach lies within the composition of our sample with a focus on the German market. Thus, we plan to perform a similar study with samples focusing on different markets and compare the results to the results of this work.

REFERENCES

- [1] A. Java, T. Finin, X. Song, and B. Tseng, "Why we twitter: Understanding microblogging usage and communities," 2007, URL: <http://ais1.umbc.edu/get/softcopy/id/1073/1073.pdf> [retrieved: April, 2014].

- [2] M. Bächle, "Economic Perspectives on the Web 2.0 - Open Innovation, Social Commerce and Enterprise 2.0 (Ökonomische Perspektiven des Web 2.0 - Open Innovation, Social Commerce und Enterprise 2.0)," *WIRTSCHAFTSINFORMATIK*, vol. 50, no. 2, 2008, pp. 129–132, URL: <http://dx.doi.org/10.1365/s11576-008-0024-2> [retrieved: April, 2014].
- [3] A. Richter, M. Koch, and J. Krisch, *Social commerce: An analysis of the change in the E-commerce (Social commerce: eine Analyse des Wandels im E-commerce)*. Fak. für Informatik, Univ. der Bundeswehr München, 2007.
- [4] L. Yang, T. Sun, M. Zhang, and Q. Mei, "We know what @you #tag: Does the dual role affect hashtag adoption?" in *Proceedings of the 21st International Conference on World Wide Web*, ser. WWW '12. New York, NY, USA: ACM, 2012, pp. 261–270, URL: <http://doi.acm.org/10.1145/2187836.2187872> [retrieved: April, 2014].
- [5] A. Bifet, G. Holmes, B. Pfahringer, and R. Gavaldà, "Detecting sentiment change in twitter streaming data," in *Workshop on Applications of Pattern Analysis (WAPA) 2011 Proceedings*, 2011, pp. 1–15.
- [6] L. Hong, G. Convertino, and E. H. Chi, "Language matters in twitter: A large scale study," in *ICWSM*, 2011, pp. 518–521.
- [7] H. Kwak, C. Lee, H. Park, and S. Moon, "What is twitter, a social network or a news media?" in *Proceedings of the 19th International Conference on World Wide Web*, ser. WWW '10. New York, NY, USA: ACM, 2010, pp. 591–600, URL: <http://doi.acm.org/10.1145/1772690.1772751> [retrieved: April, 2014].
- [8] S. Westman and L. Freund, "Information interaction in 140 characters or less: Genres on twitter," in *Proceedings of the Third Symposium on Information Interaction in Context*, ser. IIX '10. New York, NY, USA: ACM, 2010, pp. 323–328, URL: <http://doi.acm.org/10.1145/1840784.1840833> [retrieved: April, 2014].
- [9] C. Honeycutt and S. C. Herring, "Beyond microblogging: Conversation and collaboration via twitter," 2009, URL: <http://ella.slis.indiana.edu/~herring/honeycutt.herring.2009.pdf> [retrieved: April, 2014].
- [10] H. Edman, "Twittering to the top: A content analysis of corporate tweets to measure organization-public relationships," 2010, URL: <http://etd.lsu.edu/docs/available/etd-04292010-162453/unrestricted/edmanthesis.pdf> [retrieved: April, 2014].
- [11] B. Krishnamurthy, P. Gill, and M. Arlitt, "A few chirps about twitter," 2008, URL: <http://www2.research.att.com/~bala/papers/twit.pdf> [retrieved: April, 2014].
- [12] B. A. Huberman, D. M. Romero, and F. Wu, "Social networks that matter: Twitter under the microscope," 2008, URL: <http://www.uic.edu/htbin/cgiwrap/bin/ojs/index.php/fm/article/view/2317/2063> [retrieved: April, 2014].
- [13] M. Cha, H. Haddadi, F. Bevevenuto, and K. P. Gummadi, "Measuring user influence in twitter: The million follower fallacy," 2010, URL: <http://snap.stanford.edu/class/cs224w-readings/cha10influence.pdf> [retrieved: April, 2014].
- [14] S. A. Paul, L. Hong, and E. H. Chi, "Is twitter a good place for asking questions? a characterization study," in *ICWSM*, 2011, pp. 578–581.
- [15] "About twitter's link service (<http://t.co>)," URL: <https://support.twitter.com/entries/109623#> [retrieved: April, 2014].
- [16] D. Kailer, P. Mandl, and A. Schill, "An empirical study on the usage of social media in german b2c-online stores," *International journal of advanced Information technology*, vol. 3, no. 5, 2013, pp. 1–14.
- [17] "Twitter reaches half a billion accounts," 2012, URL: http://semiocast.com/en/publications/2012_07_30_Twitter_reaches_half_a_billion_accounts_140m_in_the_US [retrieved: April, 2014].
- [18] "Why do germans shun twitter?" 2013, URL: <http://www.economist.com/blogs/babbage/2013/12/social-media> [retrieved: April, 2014].

Encouraging Students to Document Software Development Projects using Blogs

Robert Law

Computer, Communications and Interactive Systems
Glasgow Caledonian University
Glasgow, Scotland
e-mail: robert.law@gcu.ac.uk

Abstract— This paper will present an ongoing project to encourage students to document Software Development Projects using Blogs. The use of blogs in the documentation of software development projects is gaining acceptance within the industry itself. As such the use of blogging within the industry will be examined to identify the potential areas of use. The pedagogical issues associated with such a project will be investigated and an attempt made to incorporate these into the student experience. The aim is to increase the student's engagement with the documentation process by using a medium that is familiar, media rich and gaining acceptance within the industry.

Keywords- *Blogs; Documentation; Software Development.*

I. INTRODUCTION

Quality documentation is an integral part of any software development project. This is true for both students' software development assignments and industry projects. Within the author's institution a trend towards a steady decline in the volume and quality of documentation submitted as part of these assignments has been detected. One hypothesis for this is the modularization of subject areas within the curriculum paints a disjointed view of the syllabus as a whole, leading students to miss the connection between analysis, design and programming. In the author's institution teaching of these subject areas tends to be as separate modules. Students perception of programming modules is such that their focus is mainly on the programming leading to assignment submission with no or extremely sparse documentary evidence of the analysis, design and development of the software product.

Since the inception of blogging, circa 1994, the past two decades have seen an almost meteoric rise in their popularity leading to academics across many disciplines incorporating blogs into their teaching [1]. Blogging, in this context, can also be used to enhance the students' capacity for reflection, opening the door to active learning [3].

Thus, the concept of using a blog was employed to engage the students in the process of recording and reflecting on their progress during the production of a software product. The hope was that the informal nature of blogging would forge a relaxed environment, in which the student could convey their thought processes. Blogs also allow for

the Lecturer to monitor the student's posts offering constructive feedback in a timely manner.

The remainder of this paper is organized as follows: Section II will give an overview of the author's use of blogs within teaching; indicating the nature of the cohort and the subject are studied. Section III gives information about pedagogical issues related to blogging, Section IV introduces the uses of blogging within the industry of software development. Section V presents the results from the use of blogging with the student cohort. Section V discusses the test results, while Section VI attempts to draw conclusions from the data and the ideas presented in Sections III and IV. Section VII discusses issues encountered during the use of blogs. Section VIII offers ideas for future work based on the synthesis of Sections III, IV and VI.

II. STUDENTS DEVELOPMENT BLOG

The author predominately teaches Games Programming using C++, DirectX and OpenGL to Year 2 and Year 4 Degree students. Thus, the participating students for this research are Year 2 and Year 4 students on the BSc Computer Games Software Development degree. The students were undertaking the modules Games Programming 1 and Games Programming 3 respectively. Since the initial test group, in academic year 2009-2010, the number of participants has reached 218.

The inherent nature of game software development is the production of a multi-media rich artifact. As such the multi-media rich environment of blogs is an ideal partner for students to express their analysis, design and software development progress in a vibrant and colourful manner. Thus, the belief was that the students would engage with the documentation process more avidly as the ability to produce a media rich online portfolio of their software artifact.

Both cohorts of students were issued with course work that involved the development of a game of their choosing. Such an open brief makes it crucial that the students embark on the proper analysis, design and planning before attempting to write any code.

The coursework was assessed in three main parts: coding (40 marks), documentation (40 marks) and extension work (20 marks). The study has been set-up to integrate into

the documentation section of the course work. The documentation section is subdivided into the following components: traditional paper based Code Explanation (10 marks), Class Diagram (5 marks), Storyboards (10 marks), References & Documentation formatting (5 marks) and Development Blog (10 marks). The expectation on the student was to blog regularly building a development log. This development log would demonstrate the student's ability to analyse a problem, define a solution to the problem and implement the defined solution. Thus, blog entries would detail the processes the student undertook to develop their game. A typical set of blog entries would include: a user specification, a proposed solution, class diagrams, storyboards, code snippets, screen captures and the final post should include a video demonstrating their game being played.

III. PEDAGOGICAL ISSUES

A number of pedagogical issues have been identified as a justification for incorporating blogging into the curriculum, irrespective of the subject area.

Zinger and Sinclair identify the property of blogging tools to have the ability to be "enjoyable yet educational" [4]. In the setting outlined above this is a key factor in the attainment of engagement with the software development process.

Zinger and Sinclair observe that blogging engages the student in a cognitive process transforming education from its fact based approach [4].

An important observation (by Zinger and Sinclair) notes that students' use of blogs exhibited the students "abilities to apply the knowledge and skills learned in one setting to another" [4]. This helps address the hypothesis postulate in the introduction about the disjointed nature the student possess of the curriculum.

A vital soft skill for any student is the ability to communicate clearly and effectively. This is particularly important in the production of technical documents such as that produced in the software development process. Using blogs has been identified as means to enhance and improve the student's communication skills [4][8].

Griffith et al. observe the blogs potential to offer an environment for collaboration, participation, interaction and writing skills as such implementing blogging within their own institution to enhance student engagement [8]. Again a list of soft skills that is essential for any graduate. The outcome of their study has resulted in their idea of the "5C's: Conversations, Creativity, Community, Collaboration and Connections", which can be mapped by the use of blogs both educationally and by the software development industry [8].

Game programming in itself is a very creative process so the fact that blogs lend themselves well to this is a definite boon for the student. This can be utilized by the student to produce feature rich and multi-media rich documentation.

Collaboration and community are interlinked and blogging plays a major part in this cohesion become engaged within the software development community through collaboration [8].

Griffith et al. also implement the idea of blogs as an e-portfolio [8] confirming the author's identification of this possibility. Chong identifies the usefulness of blogs as a just in time diary with a peer and lecturer feedback mechanism [2][5][9][13].

Two issues identified by Chong as positives of using blogging software educationally are "RSS delivery" and the ability to archive posts [5] these are also key components in the use of blogs in the software development industry.

Feedback in the form of comments are seen by Chong as vital to the "interactive nature of blogs" [5]. Von Krogh reiterates the tenants of blogging as the ability to share efficiently and effectively knowledge [11].

Chu, Kwan and Warning agree that blogs can support and facilitate communication, self-reflection, idea sharing and information organization [15]. Hsu and Lin suggest "blogs can be considered as one of the major ways of knowledge sharing" [14].

IV. BLOG USE IN INDUSTRY

Storey et al. notes that software developers make use, in general, of social media and in particular blogs and as such "the use of these mechanisms influences software development practices" [7]. This reinforces the author's decision to implement blogging as part of the student's documentation process.

Storey et al. identifies collaboration as a key component of software development; the synergistic nature of software development and the crucial role communication plays "when designing large scale modern software systems" [7]. Pagano and Maalej also indicate that the software development community has seen the potential of blogging and social media to improve collaboration and communications within software development projects [12].

Developers appear to use blogs for a number of different documentary processes: new release features, "how to's" and requirements engineering [7].

Begel, DeLine, and Zimmermann identify the fact that "social media has changed the way that people collaborate and share information" stressing how this can aid software development teams to find new ways to work together [9].

The idea that social media can play its part in the conception, design, development and successful deployment of software products is stressed by Begel, DeLine, and Zimmermann [9].

Noted by Begel, DeLine, and Zimmermann is the use of blogs by Microsoft to share "technical information and opinions with their employees" providing a synchronous communication channel and the ability to react to new information quickly [9].

Parnin and Treude indicate that software development blogs "are changing the way software is documented" allowing developers the ability to create and communicate their knowledge and experiences [10].

A crucial observation made by Parnin and Treude is "all too often documentation is absent or incomplete" this is borne out by the student population too [10].

Pagano and Maalej analysis of the content of developers' blog posts identified "functional requirements and domain

concepts” as the most favoured posts [12]. From their study, they have drawn the conclusion that the blog posts made by developers tend to be based around high level concepts [12].

Black, Harrison and Baldwin study singles out communication as the prime benefactor in the software development process of blogging and social media [15].

Reinicke and Cummings show that blue chip companies such as “IBM, Cisco and Sap” are leading the way with the use of blogs and social media as a communication and collaboration tool in software development projects [2]. An interesting observation they make is the connection between software development projects using agile methodologies and the importance of blogs and social media to the success of such projects [2].

V. RESULTS OF STUDENT BLOGGING

The results were determined by recording the data in a two stage process. Firstly the frequency of each student’s blog posts is recorded in a spreadsheet on a weekly basis. The content of the student’s weekly blog entries is reviewed and feedback is provided to the student based on the following subjective criteria: length of entry, quality of English, relevance, and ability to use technical terms in the correct context.

Secondly to determine an individual student’s final mark, the previously collected data on frequency and content is reviewed and combined to arrive at a subjective mark.

The first cohort to undertake the integration of blogging within the Year 2 module Game Programming 1 was in academic year 2009-2010. The results from this initial test showed promise and a decision was made to integrate the blog concept into the Year 4 module Game Programming 3. It has since been used in other modules that the author teaches thus, covering 218 students since its initial inception.

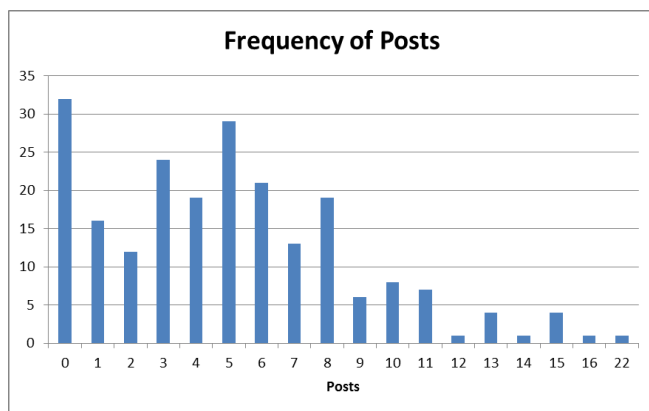


Figure 1. Frequency of Posts.

Figure 1 gives an indication of the frequency of posts made. Unfortunately 32 students from the 218 failed to engage with the process and made no post. Anecdotally, this translated to a very lightweight submission of the remaining documentation required and a low overall mark for the documentation section of the assignment.

The majority of the posts were in the range 3 to 8 with around 125 students in this range. Only around 33 students engaged to a level that could be deemed enthusiastic.

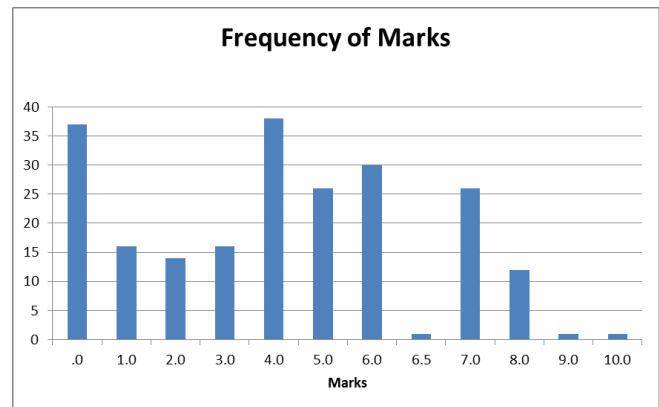


Figure 2. Frequency of Marks.

Figure 2 shows the frequency of marks gained by the students. Marks were awarded based on the number of posts and the level content of the posts. It was not sufficient to make a number of single sentence posts therefore it is conceivable that the number of posts does not match the amount of marks awarded.

Drilling down into the data, reveals the ability to compare a cohort’s performance in Year 2 against their performance in Year 4. This in itself can be revealing. The first cohort to undertake the blogging assessment can be compared again in academic year 2011-2012.

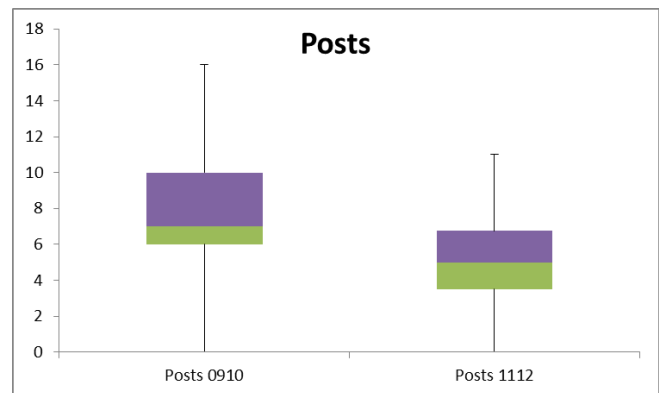


Figure 3. Comparison of posts for same cohort in Years 2 & 4.

Figure 3 indicates a comparison of posts for the same cohort in Year 2 and Year 4. Year 4 is the final year of their Degree. It is quite evident the stark drop in the number of posts being made.

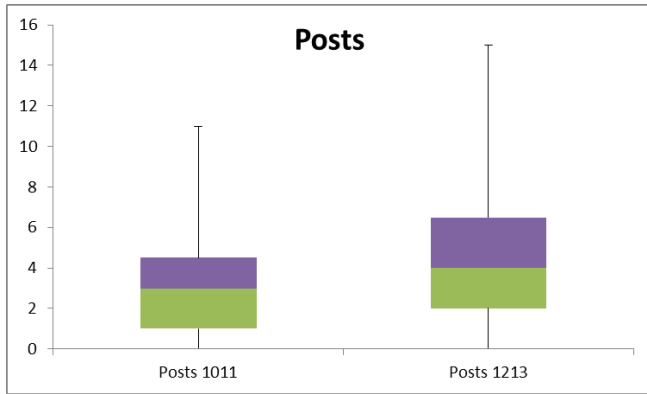


Figure 4. Comparison of posts for same cohort in Years 2 & 4.

Figure 4 indicates how the cohort for academic years 2010-2011 fared when compared again in academic year 2012-2013. On this occasion, it can be seen that there has been a positive increase in the number of posts.

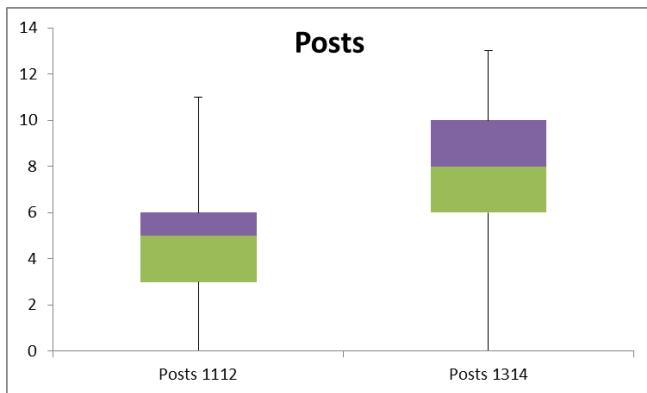


Figure 5. Comparison of posts for same cohort in Years 2 & 4.

Figure 5 indicates how the cohort for academic years 2011-2012 fared when compared again in academic year 2013-2014. On this occasion it can be seen that there has been a positive increase in the number of posts.

Comparing the marks for each cohort in the same fashion as the posts should yield a similar pattern.

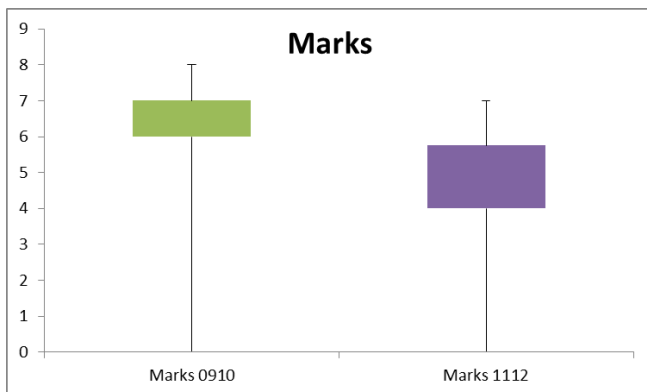


Figure 6. Comparison of marks for same cohort in Years 2 & 4.

Figure 6 indicates how the cohort for academic years 2009-2010 fared when compared again in academic year

2011-2012. On this occasion it can be seen that the marks have decrease in line with the decrease in posts shown in figure 3.

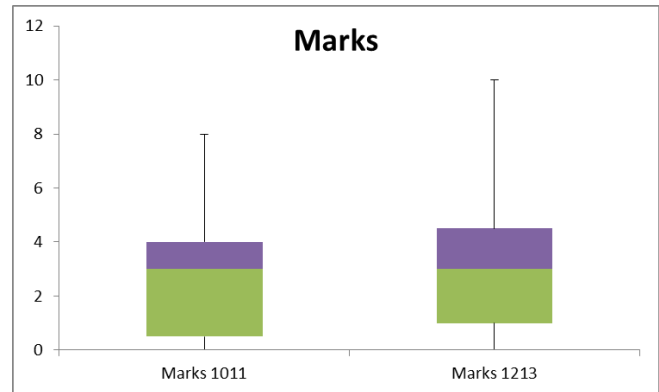


Figure 7. Comparison of marks for same cohort in Years 2 & 4.

Figure 7 indicates how the cohort for academic years 2010-2011 fared when compared again in academic year 2012-2013. On this occasion it can be seen that the marks have increased marginally following the increase in posts shown in figure 4.

Figure 8 indicates how the cohort for academic years 2011-2012 fared when compared again in academic year 2013-2014. On this occasion it can be seen that the marks have increased marginally following the increase in posts shown in figure 5.

Table 1 shows the correlation of marks and posts for each of the pairs of years. All pairs show a positive correlation between the number of posts made and the marks gained.

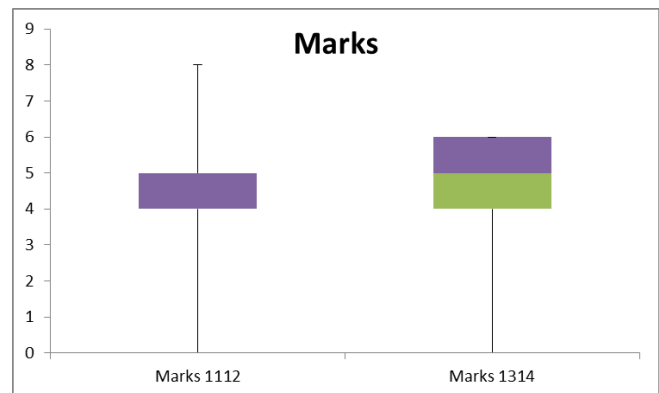


Figure 8. Comparison of marks for same cohort in Years 2 & 4.

TABLE I CORRELATION BETWEEN MARKS AND POSTS

Module		Module	
GP1 09/10	0.790539	GP3 11/12	0.927967
GP1 10/11	0.954243	GP3 12/13	0.967233
GP1 11/12	0.88226	GP3 13/14	0.745314

In the following section, an attempt will be made to draw conclusion between the results obtained and the literature review in Sections III and IV.

VI. CONCLUSIONS

The pedagogical issues revealed in Section III would indicate that the use of blogs should engage the student. However, in the data sample present around 15% of students did not engage at all, registering no posts. This was higher than expected and very disappointing.

Although the quality of the posts could be deemed to be subjective, it was surprising at the number of students that conspired to do the minimum content that they could. This complements the evidence suggested by Chu, Kwan and Warning [13].

From the data presented, it was difficult to quantify if there had been any improvement from Year 2 to Year 4 in the quality of written English. The blog posts in Year 4 did indicate a wider range of content than the previous posts in Year 2.

The drop in posts within the same cohort shown in figure 3 is interesting and unexpected. Our expectation was for the posts to increase in Year 4, as the maturity level of the students has increased and they have experienced the benefits of previous exposure to blogging.

Figures 4 and 5 offer the expected result of an increase in posts from Year 2 to Year 4.

By the end of the module, the blog that the student had created should have been in a fit state to be considered a reasonable example of an ePortfolio of their work. Unfortunately, only a small percentage of the blogs created by the student sample discussed here could conceivably be regarded as in a fit state to be a presentable ePortfolio. This was disappointing, as game programming is a very creative area awash with rich multi-media artifacts that could have been used to create rich content blog posts.

Figure 9 gives an example of the content from one of the better blog posts.

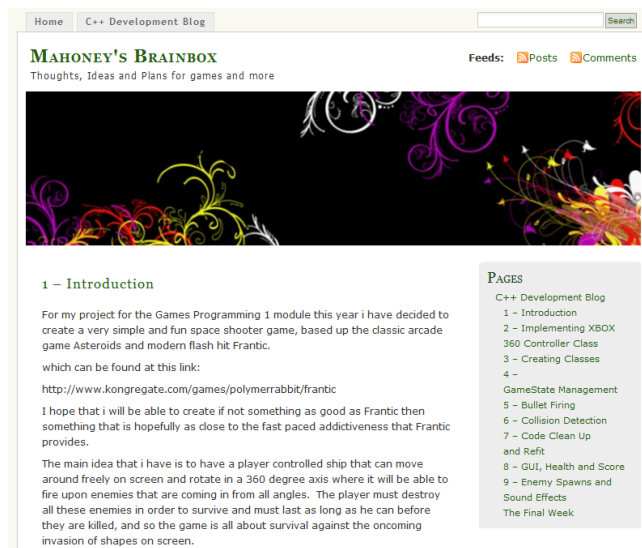


Figure 9 Example Blog Post

Table 1 shows the correlation between the mark awarded and the number of posts. All are highly positive. Thus, indicating that the students that put in the effort to post

regularly and their posts contain a level of content that can be deemed acceptable they will in turn receive a mark commensurate with the number of post they have made.

The one aspect of the blog that worked very well was the chronological aspect as the author was able to monitor and note the periods of effort applied to the assignment. The typical student tended to have top heavy posts clustered around the final few weeks of the assignment. A good student tended to have an even spread of posts illustrating good time management.

For us, the data is an indicator that there are positives in continuing to ask students to blog during software development projects. This is evidence in Section IV which highlights that Blue chip companies and independent developers are turning to the use of blogs and other social media for collaboration, communication and sharing of technical knowledge.

VII. ISSUES

Incorporating blogging within the curriculum leads to a number of issues including feeding back to students in a timely manner [8]. The author set aside approximately two hours a week to read and feed back to the students posting comments on each blog post. This creates an additional load [8] for the member of staff involved but the feeling of the author was that it was worth the extra effort.

Picking the correct environment for the blog is crucial as there is some indication that this has a bearing on how often the student posts. If the environment is perceived to be too difficult to use then the student will procrastinate on its use.

At our institution, the Virtual Learning Environment created by Blackboard tm is used. For future presentations of the module, a change of blogging environment will be sought.

Our expectation was for students to engage with the multi-media rich environment making use of the ability to post images, code fragments and video clips instead most students produced rather bland vanilla heavily text based posts.

VIII. FUTURE WORK

We will investigate expanding from blogs to include other forms of social media to encourage a range of skills. From the literature reviewed the software developers in industry use a number of diverse social media including wiki's, micro-blogging and social networking sites for creating, maintaining, communicating and collaborating.

This will lead to the student creating and maintaining the documentation required for a software development project in a purely electronic format.

The literature reviewed indicated the rise of tools that can be integrated or plugged in to popular Integrated Development Environments (IDE) to enable developers to use social media direct from their IDE.

REFERENCES

- [1] R. McDermott, G. Brindley, and G. Eccleston, "Developing tools to encourage reflection in first year students blogs," in Proceedings of

- the fifteenth annual conference on Innovation and technology in computer science education, 2010, pp. 147–151.
- [2] B. Reinicke and J. Cummings, “Can Social Media Aid Software Development?” in Proceedings of the Conference for Information Systems Applied Research ISSN, vol. 2167, 2013, p. 1508.
- [3] C. Safran, “Blogging in higher education programming lectures: an empirical study.” in Proceedings of the 12th international conference on Entertainment and media in the ubiquitous era, 2008, pp. 131–135.
- [4] L. Zinger and A. Sinclair, “Using Blogs To Enhance Student Engagement And Learning In The Health Sciences,” Contemporary Issues in Education Research (CIER), vol. 6, no. 3, pp. 349–352, 2013.
- [5] E. K. Chong, “Using blogging to enhance the initiation of students into academic research,” Computers & Education, vol. 55, no. 2, pp. 798–807, 2010.
- [6] H. N. Kim, “The phenomenon of blogs and theoretical model of blog use in educational contexts,” Computers & Education, vol. 51, no. 3, pp. 1342–1352, 2008.
- [7] M.-A. Storey, C. Treude, A. van Deursen, and L.-T. Cheng, “The impact of social media on software engineering practices and tools,” in Proceedings of the FSE/SDP workshop on Future of software engineering research, 2010, pp. 359–364.
- [8] M. Griffith, D. Simmons, W.-L. Wong, and S. Smith, “The 5 C’s of Literacy and Literary Skills Development: Conversations, Community, Collaboration, Creativity, and Connection,” in ASCILITE-Australian Society for Computers in Learning in Tertiary Education Annual Conference, vol. 2012, no. 1, 2012.
- [9] A. Begel, R. DeLine, and T. Zimmermann, “Social media for software engineering,” in Proceedings of the FSE/SDP workshop on Future of software engineering research, 2010, pp. 33–38.
- [10] C. Parmin and C. Treude, “Measuring api documentation on the web,” in Proceedings of the 2nd international workshop on Web 2.0 for software engineering, 2011, pp. 25–30.
- [11] G. Von Krogh, “How does social software change knowledge management? Toward a strategic research agenda,” The Journal of Strategic Information Systems, vol. 21, no. 2, pp. 154–164, 2012.
- [12] D. Pagano and W. Maalej, “How do developers blog?: an exploratory study,” in Proceedings of the 8th working conference on Mining software repositories, 2011, pp. 123–132.
- [13] S. K. Chu, A. Kwan, and P. Warning, “Blogging for Information Management, Learning, and Social Support during Internship,” Journal of Educational Technology & Society, vol. 15, no. 2, 2012.
- [14] C.-L. Hsu and J. C.-C. Lin, “Acceptance of blog usage: The roles of technology acceptance, social influence and knowledge sharing motivation,” Information & Management, vol. 45, no. 1, pp. 65–74, 2008.
- [15] S. Black, R. Harrison, and M. Baldwin, “A survey of social media use in software systems development,” in Proceedings of the 1st Workshop on Web 2.0 for Software Engineering, 2010, pp. 1–5.

Ranking Domain Names Using Various Rating Methods

Talattinis Kyriacos, Zervopoulou Christina, Stephanides George

Department of Applied Informatics

University of Macedonia

Thessaloniki, Greece

e-mail: ktalattinis@uom.edu.gr, zervopoulou@uom.edu.gr, steph@uom.gr

Abstract—This paper deals with the ranking of domain names, which is considered important because it is associated with their selling price. For this purpose, four well-known rating methods were used, the Massey method, the Colley method, the Keener method and finally, a method based on finite Markov chains and therefore, called Markov method. Although Massey's, Colley's and Keener's methods have their origin in sports teams rating, they can be modified and successfully applied for domain names ranking. Our effort to correctly rank domain names is based on search volume of the keyword of each domain and therefore, we used Google trends. We have also considered other factors, such as Alexa rank and keyword popularity. Information was collected via Internet and implementation of the models took place using computing tools. Our study is directly related to the global online information and for this reason allows us to do a more sophisticated rating model.

Keywords-domain names; domain name ranking; rating; Massey method; Colley method; Keener method; Markov method; Kendall's tau; Google trends.

I. INTRODUCTION

This paper presents modifications of some sports teams rating methods, in order to use them in ranking of domain names. This kind of ranking is considered important because it is associated with the formation of the price at which domain names can be sold. Specifically, these two figures are proportional amounts, i.e., the higher the rank of a domain name, the higher its selling price will be.

It is a fact that with the growth of Internet, multiple sources of profitability have appeared by its use [2][9]. Thus, the concept of domain name emerged, which mainly refers to names of websites and their extensions. It has also been proved that the ownership of domain names can be particularly lucrative for their owners.

Specifically, each domain name can create value for its owner, through revenues from an active website or even without the existence of it [9]. Nowadays, due to the rapid development of e-commerce on a global level, the domain name market has already grown into a robust and profitable industry, where millions of customers search time after time for high quality domain names in order to promote their businesses. Currently, the actual value of a domain name is

difficult to be accurately determined. However, there is a number of objective factors involved in determining the final selling price. The ownership of a domain name grants its owner two types of rights:

- 1) Managerial flexibility
- 2) Legal protection of trademarks.

Therefore, value can be created from a domain name in two ways: either by the expected profits or by options for action, such as the creation of an active website.

Creating an active site is not so easy, because the development of its content requires hard work, thoroughness and imagination, contrary to domain name fortification, which is achieved by a few "clicks" at the website of the pertinent regulatory authority. Acquiring a domain name has always been speculative. Aspiring investors taking advantage of new profit opportunities offered by the Internet, register a domain name and place a simple graphic like "page under construction". Then, they only have to wait for someone who has an exploitation plan for the domain name, but has however not acquired the appropriate website.

Another profitable and efficient strategy is to use a synonym for a domain name in conjunction with an intensive advertising campaign. The phenomenon of "cybersquatting" has also come of which concerns the creation of a website with a name, closely related to the name of an already popular website in order to exploit its reputation.

Domain names purchase is usually made via an auction. Investors often need to know which domain trading is the most profitable. A domain name can be considered as an investment [3] similarly to the real estate market. However, it is not clear how to estimate domain names value, because this market is relatively new. Consequently, some domain name sellers set selling prices arbitrarily without taking into account the actual value of the domain name. Domain name ranking can help investors to choose which domain to negotiate. Ranking refers only to the domain name and not to the active website.

In this paper, after Section I, which is a short introduction to the subject, follows Section II, where we discuss about some factors that determine domain names' rate. From these factors we chose the most frequently used by the majority of

people involved in domain names market. Section III is an overview of the rating methods presented in the paper, while Section IV consists of some illustrative examples of how these methods can be applied in order to rank certain domain names. Finally, Section V presents the experimental results generated by these methods that refer to the ranking process of domain names. Indicatively, we present top 25 domain names as a partial list of all the domains we tested.

II. DETERMINANT FACTORS

In order to rank a group of domain names, we must first clarify which are the factors that affect their importance, their value and consequently, their rank. It is worth mentioning that this market is in embryonic stage, i.e., there is no enough literature referring to the selection criteria of these factors and no other approaches for domain name ranking have yet been proposed.

Though there are many factors that determine domain names' rank, we indicatively mention these that are usually used by the majority of people (domain traders) involved in domain names market. These factors can be easily computed and are:

1. keyword popularity: the number of search results on Google for a key-word is a good indicator of how efficient is the keyword.
2. search volume of the keyword: the comparison of keyword popularity over a period of time. Google trends is the most popular and free tool used to accomplish this task. In Google trends up to five keywords can be queried simultaneously.
3. traffic: classification of domain name in Alexa. Alexa.com is currently the most reliable counter of Internet traffic and the most popular service which publishes information on the popularity of a website. It calculates the global ranking of a domain name from the traffic it has. This calculation can be done per day, per week, per month etc. The higher the ranking is, the greater the value of the name.
4. domain name extension: the extension of a domain name, in other words, the top level domain name can affect the value and the rank of the domain name. The most dominant extension is .com. Below .com, come .net, .org and domestic extensions.
5. the size of the domain name word: Names with many characters are usually hard to memorize so those with the least possible characters are more preferred.

Some other factors that also affect domain name rank but are difficult enough to be expressed quantitatively are industry popularity and brandability. Industry popularity relates to the market volume to which a specific domain name can be applied, while brandability refers to the case that someone comes up with such an interesting new word that can become a trademark [5].

For this first approach of the subject and inspired by [12], we thought that keyword popularity, search volume of the keyword and traffic measures will have the greatest importance among five determinant factors mentioned above.

III. OVERVIEW OF METHODS

First of all, we should define what the terms rating and ranking exactly imply and realize the difference between them. Rating refers to the evaluation or assessment of an item in terms of quality, quantity or some combination of both and thus, assigns a numerical value to it. Ranking is a relationship between a set of items, i.e., for any two items, the first is either 'ranked higher than', 'ranked lower than' or 'ranked equal to' the second. Therefore, a ranking vector is a permutation of the integers 1 through n or, in other words, a sorted rating vector [1]. The methods presented in this paper are due to K. Massey, W. Colley, J.P. Keener and finally, to finite Markov chains. They are being used many years ago and had initially been invented for very different purposes. Nevertheless, they can all be used for domain names ranking or webpages ranking and more generally, for the ranking of any set of objects.

A. Massey's Method

This method was proposed by Kenneth Massey in 1997 for ranking college football teams. Apart from numbers of wins and losses of a team, it also considers game scores in the ratings, i.e., spread of points, via a system of linear equations [4]. Massey's method is based on the mathematical theory of least squares, which can be represented by the following equation:

$$r_i - r_j = y_k, \quad (1)$$

where r_i and r_j are the ratings of teams i and j , respectively and y_k is the margin of victory for a game k between these teams. Each game k can be given by an equation of this form, so a system of m linear equations and n unknowns is created, where m is the number of the games that have already been played and n is the number of teams [1]. This system can be written as:

$$X \cdot r = y \quad (2)$$

and is overdetermined, because $m \gg n$, i.e., there are more equations than unknowns. To deal with this problem, Massey proposed the use of matrix $M = X^T \cdot X$ instead of X , therefore, a least squares solution is obtained [7]:

$$X^T \cdot X \cdot r = X^T \cdot y \quad (3)$$

Massey matrix M can be easily filled considering that every diagonal element M_{ii} is the total number of games played by team i and every off-diagonal element M_{ij} , for $i \neq j$, is the negation of the number of games played by team i against team j . Consequently, the Massey least squares system now becomes:

$$M \cdot r = d \quad (4)$$

where $M_{n \times n}$ is the Massey matrix described above,

$r_{n \times 1}$ is the vector of unknown ratings and

$d_{n \times 1}$ is the total difference in scores for each team.

Apart from its simple formation and its much smaller size than X , the columns of matrix M are linearly dependent, which leads to $\text{rank}(M) < n$ and so, the linear system $M \cdot r = d$ does not have a unique solution [1]. Massey solved this problem by replacing any row in M with a row of all ones and the corresponding value of d with a zero. The row in M chosen by Massey is the last one.

Summarizing the Massey Rating Method, firstly we have to form the Massey matrix M and the vector d , which represents the total difference in scores for team i , then we have to force matrix M to have full rank by making some replacements and finally, we have to solve the linear system generated by these replacements in order to take ratings vector r [4]. More specifically, we can form the Massey matrix $M = X^T \cdot X$ using $M_{ij} = -n_{j,i}$, if $i \neq j$ and $M_{ij} = n_i$, if $i = j$, where n_i is the number of games played by team i and $n_{j,i}$ is the number of games played by team i against team j . The vector d of the total difference in scores for team i is given by equation $d = X^T y$. We can make the rank of matrix M full either by replacing it with $M + e^T e$, where e is a vector of all ones or by replacing one of the rows of M with e and the corresponding entry in d with c [4]. Finally, we compute the Massey rating vector r by solving the linear system generated by the previous replacement.

B. Colley's Method

This method was proposed by astrophysicist Dr. Wesley Colley in 2001 for ranking sports teams. Colley's method is based on very simple statistical principles. In fact, it is a modified form of one of the oldest rating systems, which uses the percentage of victories of each team. This percentage is given by:

$$r_i = \frac{n_w}{n_{tot}} \tag{5}$$

where n_w are the victories of group i and n_{tot} is the total number of games played for team i [14].

Colley's method makes use of an idea from probability theory, known as Laplace's 'rule of succession', which transforms the standard winning percentage as below [14]:

$$r_i = \frac{1+n_w}{2+n_{tot}} \tag{6}$$

As follows from the above, the only information used by this model are wins, losses and number of games each team played, assuming no ties. Thus, the generated ratings are bias free, which implies that certain points gained by each team in a game are not included [4]. In other words, a win is a win regardless of the score [13]. Due to the use of Laplace's 'rule of succession', Colley's method has several advantages over the traditional rating formula:

1. At the beginning of the season, each team has a rating of $\frac{1}{2}$, instead of the preseason rating $\frac{0}{0}$ of the traditional system, which does not make any sense.
2. Colley's method takes into consideration the strength of schedule, which is the strength of a

team's opponents. This implies that, if a team beats a strong opponent it ought to receive a greater reward than if it has beaten a weaker one [1].

Then follows a summary of the Colley Rating Method: At first, we can form the Colley matrix C using $C_{ij} = -n_{j,i}$, if $i \neq j$ and $C_{ij} = 2 + n_i$, if $i = j$, where n_i is the number of games played by team i and $n_{j,i}$ is the number of games played by team i against team j . Then, we compute vector b given by:

$$b_i = \frac{1+(w_i-l_i)}{2} \tag{7}$$

where w_i is the number of wins by team i and l_i is the number of loses by team i . Finally, we solve the linear system:

$$C \cdot r = b \tag{8}$$

where the r is the rating vector for the teams [4].

C. Keener's Method

This method has been proposed by James P. Keener in 1993 for football teams ranking in uneven paired competition [6]. Keener's method is based on the theory of nonnegative matrices and forms a smoothed matrix of scores [4] generated by Laplace's rule of succession:

$$\frac{S_{ij}+1}{S_{ij}+S_{ji}+2} \tag{9}$$

Laplace's rule of succession refers to computing the entry i of the Keener matrix, where S_{ij} is the points that team i scored and S_{ji} is the points scored by team j . The reason that Keener uses Laplace's rule of succession ratio is to ensure that if a team scores 0 points, the other team does not get the entirety of the points [4].

In contrast to Colley's method, Keener's method is biased, implying that a team can boost its ranking by running up its score in a game. In other words, score points do matter.

Summarizing this method, we can form Keener matrix K using:

$$K_{ij} = h \left(\frac{S_{ij}+1}{S_{ij}+S_{ji}+2} \right) \tag{10}$$

if team i played against team j , otherwise 0, where S_{ij} is number of points scored by team i against team j and

$$h(x) = \frac{1}{2} + \frac{1}{2} \text{sgn} \left(x - \frac{1}{2} \right) \sqrt{|2x - 1|} \tag{11}$$

Finally, we can solve $K \cdot r = \lambda \cdot r$ to get Perron vector of matrix K , i.e., rating vector r . In the linear system given above, λ is the spectral radius (dominant eigenvalue) of K [4].

D. Markov's Method

This method utilizes finite Markov chains theory and therefore, it is called Markov Method. It was first used by graduate students, Angela Govan and Luke Ingram to successfully rank NFL football and NCAA basketball teams

respectively [1], where NFL is the National Football League and NCAA is the National Collegiate Athletic Association of the United States.

Markov’s method is known as Generalized Markov (GeM) ranking model and is, indeed, an adjustment of the famous PageRank algorithm that Google uses for webpage ranking. Similarly to PageRank, GeM uses parts of finite Markov chains and graph theory in order to generate ratings of n objects in a finite set. Not only sports but also any problem that can be represented as a weighted directed graph can be solved using GeM model [4].

The main idea behind the Markov Method is voting. In every game between two teams the weaker team casts a vote for the stronger team. There are many ways for a team to vote another. The simplest method uses wins and losses, implying that a winning team gains a vote by each team that has beaten. A better model would take into account game scores, namely, a winning team gets as much votes by a weaker opponent as is the margin of victory in the game between them. To make the voting method even more advanced both teams should be allowed to cast votes equal to the number of points given up in the game [1].

The main advantage of Markov’s method towards the other rating methods is the combination of more than one statistics to generate rating vector r. In order to get the GeM rating vector r, we first form G using voting matrices for the p game statistics of interest [4]. This can be done by:

$$G = a_0S_0 + \dots + a_pS_p \tag{12}$$

where $0 \leq a_i \leq 1$ and $\sum a_i = 1$.

Each stochastic matrix S_i is called a *feature* matrix and will be formed using another statistic. Finally, we compute rating vector r, the stationary vector or dominant eigenvector of G. If G is reducible, we use the irreducible

$$\bar{G} = \beta G + (1 - \beta)/nE, 0 < \beta < 1 \tag{13}$$

where E is the matrix of all ones.

IV. ILLUSTRATIVE EXAMPLES

Methods referred at the previous section have a wide variety of applications except of sports. Our thought was to apply these methods for ranking domain names. Therefore, we used one of the most significant determinant factors shown at Section II, which is Google trends.

Google Trends provides relative numbers. In fact, it analyzes a portion of searches done in Google in order to compute how many of them have been done for the terms entered, compared to the total number of searches done on Google over time. Google does not reveal absolute numbers for competitive reasons, but also because those numbers would not be exact. The fact that Google trends are relative numbers implies that there may have been done more searches for object A than for object B, but these searches may be less than those of another object C. For example, assuming that object A is Gauss and object B is Markov, the winner is Gauss with 54 Google trends points average against Markov’s 27 points average. However, if object C is

Shannon, Gauss becomes the underdog with 9 points average, while Shannon is given 54. This much similarity between Google trends and points in a game is exactly the reason we decided to use Google trends as determinant factor for the ranking methods presented here.

In the example described below, there are five domain names that have been sold in early 2014, which are jean.com, desirous.com, authorization.com, true.com and finally, peaked.com. We will attempt to rank these domains based on search volume average they get by Google trends during 2013. The question is how can search volume average be related to the points that a team succeeded against another?

There are many ways to define the notion of a game for domain names. For example, if statistics on domain names are given by Google trends, then we can say that domain i beats domain j if $d_i > d_j$, where d_i and d_j are the Google trends measures for these domains. Therefore, $d_i - d_j$ represents the difference in trends’ value between domains i and j. Table I shows Google trends data for the five domains of our example.

TABLE I. GOOGLE TRENDS DATA

Domain i	Domain j	Trends i, j
jean.com	desirous.com	88, 0
jean.com	authorization.com	88, 4
jean.com	true.com	76, 73
jean.com	peaked.com	88, 0
desirous.com	authorization.com	1, 93
desirous.com	true.com	0, 73
desirous.com	peaked.com	20, 80
authorization.com	true.com	3, 73
authorization.com	peaked.com	93, 5
true.com	peaked.com	73, 0

Adjusting the Massey rating method for domain names, we start with the same idealized function (1). Then, the Massey domain ranking method proceeds as usual, according to (4). Below this linear system is showed:

$$\begin{pmatrix} 4 & -1 & -1 & -1 & -1 \\ -1 & 4 & -1 & -1 & -1 \\ -1 & -1 & 4 & -1 & -1 \\ -1 & -1 & -1 & 4 & -1 \\ 1 & 1 & 1 & 1 & 1 \end{pmatrix} \cdot \begin{pmatrix} r_1 \\ r_2 \\ r_3 \\ r_4 \\ r_5 \end{pmatrix} = \begin{pmatrix} 263 \\ -313 \\ 26 \\ 213 \\ 0 \end{pmatrix} \tag{14}$$

Table II gives rating and ranking data generated by Massey method:

TABLE II. RATING AND RANKING BY MASSEY METHOD

Ranking	Domain	Rating
1	jean.com	52.6
2	true.com	42.6
3	authorization.com	5.2
4	peaked.com	-37.8
5	desirous.com	-62.6

As we conclude from the above matrix, jean.com has beaten all the other four domain names and, thus, it terminates at first position of ranking. Contrary to jean.com, domain desirous.com has been defeated by all others,

therefore, it has the lowest rating of all and so, it takes the last position of ranking.

Using the Colley rating method (8) for domain names, we get the results below:

$$\begin{pmatrix} 6 & -1 & -1 & -1 & -1 \\ -1 & 6 & -1 & -1 & -1 \\ -1 & -1 & 6 & -1 & -1 \\ -1 & -1 & -1 & 6 & -1 \\ -1 & -1 & -1 & -1 & 6 \end{pmatrix} \cdot \begin{pmatrix} r_1 \\ r_2 \\ r_3 \\ r_4 \\ r_5 \end{pmatrix} = \begin{pmatrix} 3 \\ -1 \\ 1 \\ 2 \\ 0 \end{pmatrix} \quad (15)$$

Table III gives rating and ranking data generated by Colley method. Rating values have been rounded at four decimal digits.

TABLE III. RATING AND RANKING BY COLLEY METHOD

Ranking	Domain	Rating
1	jean.com	0.7857
2	true.com	0.6429
3	authorization.com	0.5
4	peaked.com	0.3571
5	desirous.com	0.2143

As we may see in the above table, jean.com terminates again first, while desirous.com gets again the last rank. Positions of the three other domains also remain the same.

Then, we continue with the Keener rating method for domain names. Below, the values of Keener matrix K are shown:

$$\begin{pmatrix} 0 & 0.9944 & 0.9727 & 0.5705 & 0.9944 \\ 0.0056 & 0 & 0.0105 & 0.0067 & 0.1165 \\ 0.0273 & 0.9895 & 0 & 0.0263 & 0.969 \\ 0.4295 & 0.9933 & 0.9737 & 0 & 0.9933 \\ 0.0056 & 0.8835 & 0.031 & 0.0067 & 0 \end{pmatrix} \quad (16)$$

Table IV gives rating and ranking data generated by Keener method:

TABLE IV. RATING AND RANKING BY KEENER METHOD

Ranking	Domain	Rating
1	jean.com	0.7391
2	true.com	0.6604
3	authorization.com	0.1253
4	peaked.com	0.039
5	desirous.com	0.0192

Table IV indicates more significant difference between true.com and authorization.com than before. This can be explained by the fact that Keener’s method is not bias-free, which means that all the points succeeded in a duel are taken into account for ranking.

The last method we used in order to rank the five domain names of our example is Markov method. As we have mentioned in Section II, the main idea behind this method is voting. In Table V, follows the trends voting matrix:

TABLE V. TRENDS VOTING MATRIX

	jean	desirous	authorization	true	peaked
jean	0	0	4	73	0
desirous	88	0	93	73	80
authorization	88	1	0	73	5
true	76	0	3	0	0
peaked	88	20	93	73	0

Below, stochastic matrix G is shown, which is generated by normalizing the rows of the above voting matrix:

$$\begin{pmatrix} 0 & 0 & 0.0519 & 0.9481 & 0 \\ 0.2635 & 0 & 0.2784 & 0.2186 & 0.2395 \\ 0.5269 & 0.006 & 0 & 0.4371 & 0.03 \\ 0.962 & 0 & 0.038 & 0 & 0 \\ 0.3212 & 0.073 & 0.3394 & 0.2664 & 0 \end{pmatrix} \quad (17)$$

Table VI gives rating and ranking data generated by Markov method:

TABLE VI. RATING AND RANKING BY MARKOV METHOD

Ranking	Domain	Rating
1	jean.com	0.4801
2	true.com	0.4746
3	authorization.com	0.0435
4	peaked.com	0.0014
5	desirous.com	0.0004

In the small example described in this paper, all four methods generate same ranking results. However, Markov’s method has a vital difference of the other three. This difference comes from the fact that, as we have mentioned in Section III, Markov method allows the use of more than one statistics. Therefore, Markov method can be characterized as more representative than the others. The weights we have set for Google trends, Google results and Alexa rank were 0.4, 0.3 and 0.3 respectively, due to our intention to rely mostly on Google trends, though these weights may vary in all the possible ways.

At this point, we will see the results generated by some other determinant factors, such as Google Result and Alexa Rank.

TABLE VII. OTHER DETERMINANT FACTORS

Domain	Google Result	Alexa Rank (ar)	Adjusted Alexa Score (ads)
jean.com	358,000,000	1,782,928	560.8751
desirous.com	1,660,000	-	1
authorization.com	72,700,000	3,566,076	280.4203
true.com	676,000,000	220,338	4538.4818
peaked.com	5,780,000	-	1

As we have mentioned in Section II, Alexa classifies domain names counting the Internet traffic. Alexa ranks only domain names with an active website, thus dashes shown in Table VII represent domain names with a non-active website. At this point, we should deal with a major issue, which is that though Alexa is a ranking system, we want to turn it into one that uses points.

This issue implies that, though Alexa assigns to rank “1” the most visited website, supposedly A, this value as a number is smaller than ranks of much less visited websites, supposedly B, C, etc. Therefore, we should modify Alexa rankings given in Table VII, so that the ordering of the ranking values is reversed. One solution to this problem is, for n items ranked by Alexa, to compute each item’s rank by:

$$ads_i = \frac{\max\{ar_i, \dots, ar_n\} + 1}{ar_i} \quad (18)$$

where ads_i is the adjusted Alexa rank for website i and ar_i is the rank given by Alexa for website i . However, this solution is not so fair for some items of the set. For instance, in case we have to rank domains google.com, facebook.com and desirous.com, rank given by Alexa is “1”, “2” and “no enough data to rank this website”, respectively. Thus, when dividing the total number of items ranked, namely “3”, with each item’s rank we get the results below:

- Three points are assigned to google.com
- Two points is assigned to facebook.com and
- One point is also assigned to desirous.com

The unfairness of this solution lies in that points of facebook.com, which is an active website, are very close to points of a non-active website, as is desirous.com. Thus, we should use a better solution, which is described below.

This better solution might be obtained using the following equation:

$$ads_i = \frac{tna}{ar_i} \tag{19}$$

where tna is the total number of websites ranked by Alexa. Despite of our thorough research, we have not found any official source referring exactly how many websites are currently ranked by Alexa. Therefore, we have chosen a typically large number, namely, $1 \cdot 10^9$ for tna variable mentioned before, i.e., we divide $1 \cdot 10^9$ with each item’s rank to turn the ordering of the ranking values into descending. In Table VII, we have written the adjusted Alexa score value of each domain.

As we can see in Table VII, true.com has the higher traffic amount of all, while before adjustment it was at the third position of Alexa rank. Then comes jean.com and authorization.com, which is the domain showed up to be at the first position before adjustment. Domains desirous.com and peaked.com take both the value “1” for their traffic, which is the value assigned to non-active websites.

Finally, in Table VIII, we tested the five domains of our example applying Markov method with three determinant factors, which are Google trends, Google results and Alexa rank. The results generated are shown in the table below.

TABLE VIII. MARKOV RATING - RANKING WITH 3 STATISTICS

Ranking	Domain	Rating
1	true.com	0.4649
2	jean.com	0.4242
3	authorization.com	0.1057
4	peaked.com	0.0042
5	desirous.com	0.0011

Table IX is a summary table, which shows the domain rankings generated by all four methods described in Section III and their selling prices:

TABLE IX. SUMMARY TABLE

Domain	Massey	Colley	Keener	Markov 1	Markov 3	Prices
jean	1	1	1	1	2	50000
desirous	5	5	5	5	5	2600
author ization	3	3	3	3	3	35100
true	2	2	2	2	1	350000
peaked	4	4	4	4	4	4000

As we see, due to the use of more than one determinant factors in this ranking, rank positions between the first two domains have interchanged and this agrees with the selling prices.

V. EXPERIMENTAL RESULTS

At first, we should describe our database, which contains information on transaction prices collected from publicly available tenders and values from databases of domain names coming from closed auctions data. Currently the database consists of 75,000 transaction prices that occurred during the period between 1999 and 2013. The database is updated regularly and it is worth mentioning that the collection of data required some effort since selling prices are not always available in digital form, even when they are published. In order to gather our data from Internet resources, we have implemented a web crawler in Java. Our crawler is typically programmed to visit sites, which contain domain names’ selling prices, Google trends, keyword popularity and traffic measures. The crawling process was held by taking into account the reliability of information. The data gathering process also involved parsing data files.

In any case, data gathering and parsing had to be automated. We have conducted thorough research and have already implemented some techniques for parallelization of collecting data, in order to keep our database updated in time. For more details about crawling, its parallelization and parsing processes we refer the reader to [8].

The rest of this paper presents the empirical results generated by the four methods described at Section III when we applied them to our database data. The numerical computations of the ratings were done using Matlab. Tables XI to XIV have been constructed in the following format: the first column represents the ranking of domains, the second column is the domain name itself, the third column represents the rating of each domain, the fourth column is the price at which the domain name was sold and finally, the last column consists of the date on which each domain was sold. Though we refer to selling price, it cannot be a reliable measure of comparison, due to its dependance of the time that happened.

In Tables XI to XIV, we present top 25 domain names as a partial list of all the domains we tested, using Google trends determinant factor for the year 2013. Each table shows the top 25 domains, according to one of the methods. We tested domain names with same Top Level Domain (TLD), i.e., .com. Table XIV shows the results of Markov method using three determinant factors, Google trends, Google results and Alexa rank. Similarly to Section IV, the weights we used for these factors were 0.4, 0.3 and 0.3 respectively. Due to the reliability of method Markov that takes into account three determinant factors, we have also included in its table the importance degree given to the domain by Google PageRank (PR). Briefly, PageRank algorithm states that a website is important if it is shown by other important websites. This degree gets values

from 1 to 10 (PR1 - PR10). The higher the PageRank obtained from a website, the higher its ranking position in search results [11]. The comparison among Markov method and Google PageRank shows that there are many results in common between them.

In order to compare the generated ranking lists, we make use of Kendall 's correlation measure τ , which gives the degree to which one list agrees (or disagrees) with another [1] and is computed by:

$$\text{Kendall's } \tau = \frac{n_c - n_d}{n(n-1)/2} \quad (19),$$

where n_c is the number of concordant pairs and n_d is the number of discordant pairs.

Kendall 's tau value varies between -1 and 1, i.e., $-1 \leq \tau \leq 1$. If $\tau = 1$, then the two lists are in perfect agreement, while if $\tau = -1$, the two lists are totally opposite to each other [10]. Comparing the methods described in this paper, according to Kendall 's tau, we get the results below:

TABLE X. KENDALL'S TAU TABLE

Pair of Methods	Kendall's Tau Value
Massey - Colley	0.942
Massey - Keener	0.9451
Massey - Markov	-0.451
Colley - Keener	0.9969
Colley - Markov	-0.4431
Keener - Markov	-0.44

From Table X, we conclude that methods Massey, Colley and Keener are very alike, while Markov is differentiated due to the use of more than one determinant factors that provides.

Then, follow the tables that show the top 25 domain names according to each of the four methods we described.

TABLE XI. MASSEY RANKING

Ranking	Domain name	Rating	Price	Selling Date
1	fb.com	78.6078	8,500,000	1/1/2010
2	phone.com	69.8627	1,200,000	1/2/2003
3	shop.com	61.2745	3,500,000	1/11/2003
4	photo.com	58.7451	1,250,000	6/5/2010
5	men.com	54.3333	1,320,000	1/2/2000
6	software.com	52.5882	3,200,000	1/12/2005
7	find.com	51.2549	1,200,000	1/3/2004
8	pizza.com	51.1176	2,605,000	3/4/2008
9	express.com	41.8235	2,000,000	1/3/2000
10	call.com	38.6275	1,100,000	2/9/2009
11	tom.com	35.9608	2,500,000	1/12/1999
12	zip.com	34.5294	1,058,830	28/10/2010
13	candy.com	28.6078	3,000,000	10/6/2009
14	vista.com	21.0196	1,250,000	14/11/2007
15	ticket.com	20.5686	1,525,000	16/10/2009
16	coupons.com	20.2746	2,200,000	1/1/2000
17	fly.com	19.9804	1,500,000	1/11/1999
18	wine.com	16.5882	3,300,000	1/9/2003
19	webcam.com	15.3725	1,020,000	10/6/2009
20	beer.com	13.9216	7,000,000	1/1/2004
21	england.com	12.9412	2,000,000	1/12/1999
22	casino.com	12.1961	5,500,000	1/11/2003
23	telephone.com	9.6275	2,000,000	1/1/2000
24	VIP.com	8.5294	1,400,000	1/12/2003
25	autos.com	8.2746	2,200,000	1/12/1999

TABLE XII. COLLEY RANKING

Ranking	Domain name	Rating	Price	Selling Date
1	fb.com	0.9717	8,500,000	1/1/2010
2	phone.com	0.9528	1,200,000	1/2/2003
3	shop.com	0.934	3,500,000	1/11/2003
4	photo.com	0.9151	1,250,000	6/5/2010
5	pizza.com	0.8962	2,605,000	3/4/2008
6	men.com	0.8774	1,320,000	1/2/2000
7	express.com	0.8585	2,000,000	1/3/2000
8	call.com	0.8302	1,100,000	2/9/2009
9	software.com	0.8302	3,200,000	1/12/2005
10	find.com	0.8019	1,200,000	1/3/2004
11	tom.com	0.783	2,500,000	1/12/1999
12	candy.com	0.7641	3,000,000	10/6/2009
13	zip.com	0.7453	1,058,830	28/10/2010
14	coupons.com	0.7264	2,200,000	1/1/2000
15	wine.com	0.7075	3,300,000	1/9/2003
16	ticket.com	0.6887	1,525,000	16/10/2009
17	webcam.com	0.6698	1,020,000	10/6/2009
18	england.com	0.6509	2,000,000	1/12/1999
19	vista.com	0.6321	1,250,000	14/11/2007
20	fly.com	0.6132	1,500,000	1/11/1999
21	korea.com	0.5849	5,000,000	1/1/2000
22	beer.com	0.5849	7,000,000	1/1/2004
23	casino.com	0.5566	5,500,000	1/11/2003
24	telephone.com	0.5377	2,000,000	1/1/2000
25	VIP.com	0.5189	1,400,000	1/12/2003

TABLE XIII. KEENER RANKING

Ranking	Domain name	Rating	Price	Selling Date
1	fb.com	0.3224	8,500,000	1/1/2010
2	phone.com	0.2813	1,200,000	1/2/2003
3	shop.com	0.2685	3,500,000	1/11/2003
4	photo.com	0.2547	1,250,000	6/5/2010
5	pizza.com	0.246	2,605,000	3/4/2008
6	men.com	0.2422	1,320,000	1/2/2000
7	express.com	0.2378	2,000,000	1/3/2000
8	software.com	0.2196	3,200,000	1/12/2005
9	call.com	0.2194	1,100,000	2/9/2009
10	find.com	0.2137	1,200,000	1/3/2004
11	tom.com	0.1964	2,500,000	1/12/1999
12	candy.com	0.1771	3,000,000	10/6/2009
13	zip.com	0.1716	1,058,830	28/10/2010
14	coupons.com	0.152	2,200,000	1/1/2000
15	wine.com	0.1429	3,300,000	1/9/2003
16	ticket.com	0.1391	1,525,000	16/10/2009
17	webcam.com	0.1374	1,020,000	10/6/2009
18	england.com	0.133	2,000,000	1/12/1999
19	vista.com	0.1313	1,250,000	14/11/2007
20	fly.com	0.128	1,500,000	1/11/1999
21	korea.com	0.1169	5,000,000	1/1/2000
22	beer.com	0.1158	7,000,000	1/1/2004
23	casino.com	0.1127	5,500,000	1/11/2003
24	telephone.com	0.1014	2,000,000	1/1/2000
25	VIP.com	0.1001	1,400,000	1/12/2003

TABLE XIV. MARKOV RANKING WITH 3 STATISTICS

#	Domain name	Rating	Google PageRank	Price / Selling Date	
1	coupons.com	0.1055	6	2,200,000	1/1/2000
2	photo.com	0.1053	4	1,250,000	6/5/2010
3	shop.com	0.0816	5	3,500,000	1/11/2003
4	VIP.com	0.0539	3	1,400,000	1/12/2003
5	find.com	0.0472	4	1,200,000	1/3/2004

6	casino.com	0.046	5	5,500,000	1/11/2003
7	phone.com	0.0416	5	1,200,000	1/2/2003
8	express.com	0.0383	5	2,000,000	1/3/2000
9	fb.com	0.0346	0	8,500,000	1/1/2010
10	men.com	0.034	3	1,320,000	1/2/2000
11	tom.com	0.0306	7	2,500,000	1/12/1999
12	software.com	0.0288	5	3,200,000	1/12/2005
13	call.com	0.0268	4	1,100,000	2/9/2009
14	pizza.com	0.0217	4	2,605,000	3/4/2008
15	feedback.com	0.0199	7	1,230,000	1/2/2003
16	zip.com	0.018	0	1,058,830	28/10/2010
17	savings.com	0.0172	5	1,900,000	1/2/2003
18	wine.com	0.0166	6	3,300,000	1/9/2003
19	fly.com	0.0163	5	1,500,000	1/11/1999
20	candy.com	0.0149	5	3,000,000	10/6/2009
21	vista.com	0.0146	0	1,250,000	14/11/2007
22	webcam.com	0.0132	0	1,020,000	10/6/2009
23	england.com	0.0127	2	2,000,000	1/12/1999
24	auction.com	0.0122	5	1,700,000	27/3/2009
25	ticket.com	0.0122	0	1,525,000	16/10/2009

VI. CONCLUSIONS – FUTURE PLANS

In this paper, we saw how we can rank domain names with four different methods. These methods are mainly used in sports industry. Compared to the others, Markov method allows ranking based on more than just one factor and as we saw in Section IV, this is crucial in the formation of ranking values.

For generating our empirical results using Massey, Colley and Keener method, we used Google trends as determinant factor, while in Markov method, we used Google trends, Google results and Alexa rank. Determinant factors were chosen according to what are people involved in domain name market looking for. In our empirical results, we cannot use the selling price as criterion to check if rankings lists generated via different ranking methods match. This is due to the fact that selling prices were formed based on past factors or data, while our ranking is based on current factors or data. This is also confirmed by PageRank indicator.

In conclusion, rating methods presented in the paper may be used by many groups of people, such as domain traders, portfolio managers and investors. Concerning to decision making process, i.e., if someone decides to purchase a domain name according to its rank, the methods presented in this paper can be a utility tool, but not the only one.

Finally, when we compared the ranking lists according to Kendall's tau correlation method, we conclude that Massey, Colley and Keener have much in common, while Markov is different enough due to the factors it takes into account.

Our goal for the future is to use more determinant factors for domain names ranking, such as brandability and internet popularity, but also to test even more rating methods, such as Elo's system or the Park-Newman method.

REFERENCES

- [1] A. N. Langville and C. D. Meyer, *Who's #1?: The Science of Rating and Ranking*, Princeton University Press, Princeton, NJ, USA, 2012.
- [2] A. Tajirian, "Thoughts on Domain Name Investing for Newbies", 2008 [Online]. Available from: http://domainmart.com/news/Thoughts_on_Domain_Name_Investing_for_Newbies.pdf, [retrieved: May, 2014].
- [3] A. Tajirian, "Valuing Domain Names: Methodology", 2005 [Online]. Available from: <http://domainmart.com/news/methodology.htm> [retrieved: May, 2014].
- [4] A.Y. Govan, "Ranking Theory with Application to Popular Sports", unpublished. PhD thesis, North Carolina State University, 2008.
- [5] <http://ezinearticles.com/?Domain-Appraisal-Guide---20-Factors-That-Decide-the-Selling-Price&id=1436181> [retrieved: May, 2014].
- [6] J. P. Keener, "The Perron-Frobenius Theorem and the Ranking of Football Teams". *SIAM Review*, vol. 35, No 1, March. 1993, pp. 80-93.
- [7] K. Massey, "Statistical Models Applied to the Rating of Sports Teams", unpublished. Bachelor's thesis, Bluefield College, 1997.
- [8] K.Talattinis, A Sidiropoulou, K.Chalkias, and G.Stephanides, "Parallel Collection of Live Data Using Hadoop", *IEEE 14th PanHellenic Conference on Informatics (PCI)*, ISBN: 978-1-4244-7838-5, Sept. 2010, pp. 66-71.
- [9] M. Jindra, "The market for Internet domain names", in *Proc. 16th ITS Regional Conf.*, Porto, Portugal, 2005 [Online]. Available from: <http://userpage.fu-berlin.de/~jmueller/its/conf/porto05/papers/Jindra.pdf> [retrieved: May, 2014].
- [10] M. Kendall, "A new measure of rank correlation", *Biometrika*, 1938.
- [11] S. Brin, L. Page, R. Motwami, and T. Winograd, *The PageRank citation ranking: Bringing order to the Web*. Technical Report 1999 – 0120, Computer Science Department, Stanford University, 1999.
- [12] T. Preis, H. S. Moat, and H. E. Stanley, "Quantifying trading behavior in financial markets using Google Trends". *Sci. Rep.* 3, 1684, 2013.
- [13] T.P. Chartier, E. Kreutzer, A.N. Langville, and K.E. Pedings, "Bracketology: How can math help?". In: Gallian, Joseph (Ed.), *Mathematics and Sports Dolciani Mathematical Expositions*, vol. 43. Mathematical Association of America. Jan. 2010, pp. 55-70.
- [14] W.N.Colley, "Colley's Bias Free College Football Ranking Method: The Colley Matrix Explained", 2002 [Online]. Available from: <http://www.colleyrankings.com> [retrieved: May, 2014].

Combining Load Balancing with Energy Saving in a Cluster – Based P2P System

Minas Tasiou, Konstantinos Antonis

Dept. of Computer Engineering
Technological Institute of Central Greece
Lamia, Greece
mtasios@teilam.gr, k_antonis@teiste.gr

Theofanis – Aristofanis Michail

Computer Technology Institute
Rion, Patras, Greece
michail@westgate.gr

Abstract – Load balancing is an interesting domain for research in Peer-to-Peer (P2P) systems, because of the large scale, heterogeneity and dynamic nature of the peers. In this work, we properly change the load balancing algorithm presented in [2] for an unstructured P2P system, that is based on a partially centralized architecture, to lower the power consumption of the whole system, without performance unacceptable loss in load balancing measures. In this algorithm, we consider some peers of the system as “green” nodes, we locate the peer to serve a request with the use of heuristic methods combined with a simple mathematical model, and we use simulation to compare our results with the results of [2]. As the reader will see, we can save much energy with a very low percentage loss in load balancing measurements.

Keywords – Energy saving, load balancing, cluster-based P2P systems, simulation

I. INTRODUCTION

A lot of research has been done during the last years on the problem of energy saving in large distributed systems like clouds, computational grids, peer-to-peer (P2P) systems, etc. In this work, we study the problem of energy saving in P2P systems, where users contribute their resources (storage space, computation time) and content (files, etc.) to the community. The users act both as a client and server. Although there are various potential domains of P2P systems, the file sharing systems (Napster [5], Kazaa, Gnutella [6], BitTorrent [7], etc.), received the greatest popularity among internet users [4].

The operation of any peer-to-peer content distribution system relies on a network of peer computers (nodes), and connections (edges) between them. This network is formed on top of—and independently from—the underlying physical computer (typically IP) network, and is thus referred to as an “overlay” network. Overlay networks can be distinguished in terms of their centralization and structure [4].

Considering network centralization the following three categories are identified:

- Purely Decentralized Architectures. All nodes in the network perform exactly the same tasks, acting both as servers and clients, and there is no central coordination of their activities.
- Partially Centralized Architectures. The basis is the same as with purely decentralized systems. Some of the nodes, however, assume a more important role, acting as local central indexes for files shared by local peers. The way in which these superNodes are assigned their role by the

network varies between different systems. It is important, however, to note that these superNodes do not constitute single points of failure for a peer-to-peer network, since they are dynamically assigned and, if they fail, the network will automatically take action to replace them with others.

- Hybrid Decentralized Architectures. In these systems, there is a central server facilitating the interaction between peers by maintaining directories of metadata, describing the shared files stored by the peer nodes.

By structure, overlay networks are classified as following:

- Unstructured. The placement of content (files) is completely unrelated to the overlay topology. Unstructured systems are generally more appropriate for accommodating highly-transient node populations. Some representative examples of unstructured systems are Napster, Gnutella, Kazaa.
- Structured. In structured networks, the overlay topology is tightly controlled and files (or pointers to them) are placed at precisely specified locations. These systems essentially provide a mapping between content (e.g. file identifier) and location (e.g. node address), in the form of a distributed routing table, so that queries can be efficiently routed to the node with the desired content. Typical examples of structured systems include Chord [11], Content Addressable Network (CAN) [12], Tapestry [13] among others.

P2P systems offer a lot of interesting research domains, like searching, infrastructure, security, load balancing, etc. Load balancing is an interesting domain for research in P2P systems, because of the large scale, heterogeneity and dynamic nature of the peers. The research in this domain falls in two categories. In the first, techniques for better item distribution in the name space so as improvements in routing and searching can be accomplished (e.g. [11, 12, 13]). In the second, techniques for items’ replicas placement to the network nodes, for improving the throughput and Quality of Service (QoS) provided to the users, can be included (e.g. [14]).

Another research domain that we consider here, is energy saving. Recently the energy consumption of Information Technology (IT) and networking infrastructure has attracted more and more attention. Energy consumption is now considered as an important factor of IT and communication system design [8, 9]. In a 2007 report ([10]) towards the Congress, the Environmental Protection Agency (EPA) stated that: “During the past five years, increasing demand for computer resources has led to significant growth in the

number of data center servers, along with an estimated doubling in the energy used by these servers and the power and cooling infrastructure that supports them”.

In this work, we properly change the load balancing algorithm presented in [2] for an unstructured P2P system, that is based on a partially centralized architecture, to lower the power consumption of the whole system, without performance unacceptable loss in load balancing measures. The original algorithm uses replication to achieve better performance in load balance and to prevent the system from users acting as free-riders. We maintain the replication feature in our algorithm, too. But, the original algorithm, as well as the other algorithms for unstructured systems, concerns more to optimize the fairness index in order to achieve better balancing of the existing load, without taking into account the energy efficiency, which is very crucial nowadays. To solve the problem, we consider some peers of the system as “green” nodes and we use some heuristic methods combined with a simple mathematical model to locate the peer able to serve a suitable request. To assess the performance of our algorithm we use simulation and we compare our results with the results of [2]. As the reader will see, we can save much energy with a very low percentage loss in load balancing measurements. Our algorithm can solve the problem of energy saving, if the nodes’ energy efficiency is known. The reader can find how a node’s energy efficiency can be computed in section IV.

The article is organized as following: section 2 presents related works in load balancing of various P2P systems under the energy saving assumption. Section 3 describes briefly the load balancing algorithm appeared in [2] for an unstructured P2P system. Section 4 presents our algorithm which combines the algorithm of [2] with the energy saving feature. Section 5 describes the experiment and the derived simulation results. Finally, section 6 presents our conclusions and presents our plans for future work.

II. RELATED WORK

In [15], a novel scheme called energy-efficient peer-to-peer caching with optimal radius for 4G hybrid networks is introduced, that reduces energy consumption and distributes load equitably among mobile peer nodes. In this scheme, a P2P overlay network is built among mobile nodes to facilitate cooperative data sharing in order to relieve the traffic bottleneck at the base station.

In [21], a P2P Minimum Boundary Rectangle (PMBR) is proposed, which is a spatial index specifically designed for mobile P2P environments. A node that contains desirable data item (s) can be easily identified by reading the PMBR index. Then, a selective tuning algorithm is proposed, called Distributed exponential Sequence Scheme (DSS, for short), that provides clients with the ability of selective tuning of data items, thus preserving the scarce power resource. The proposed algorithm is simple but efficient in supporting linear transmission of spatial data and processing of location-aware queries. The results from theoretical analysis and experiments show that the proposed algorithm with the PMBR index is scalable and energy efficient in both range queries and nearest neighbor queries.

In [16], the Round Robin method for reducing power consumption is discussed and then a load balancer method with queue system for reducing the total power consumption of a server peer in a peer-to-peer system is presented.

In [17], the legacy BitTorrent protocol is compared to EE-BitTorrent, a proxy based version recently proposed for energy efficiency, in a residential scenario. It is shown that the performance achieved by users is strongly influenced by the uplink throughput allowed by the access network. When the available uplink rate is low, the legacy BitTorrent protocol performs poorly and EE-BitTorrent outperforms it, in terms of average download time and energy consumption at the user's PC. The opposite occurs when the uplink rate is good. Motivated by these results, the researchers designed and implemented AdaBT, an adaptive algorithm that dynamically selects the most efficient BitTorrent option (i.e., legacy or proxy-based), depending on the operating conditions experienced by the user. The experimental results show that AdaBT is able to reduce significantly the download time provided by either the legacy BitTorrent or EE-BitTorrent, reducing this way the power consumption.

In [18], a new architecture for sharing resources amongst home environments is proposed. This approach relies on complete decentralization in a peer-to-peer like manner, and above all, aims at energy efficiency. Energy metrics are defined, which have to be optimized by the system. The system itself uses virtualization to transparently move tasks from one home to another in order to optimally utilize the existing computing power. An overview of the proposed architecture is presented as well as an analytical evaluation of the possible energy savings in a distributed example scenario where computers share downloads.

In [19], the power saving potential of P2P file sharing in two cases is revealed; popular and unpopular files. For popular files, it is derived, with regard to BitTorrent, an expression for the optimal time seeders should support leechers. For unpopular files, an existing model is extended by taking into account leechers’ power consumption dependent on the load. Leechers are assumed to build a temporary cluster within the P2P-overlay. The required number of active leechers is determined to cope with a given load and results from an analytical model to simulation are compared. It is demonstrated that it is possible to reach almost optimal energy efficiency for the download scenario by comparing the local case without cooperation with the distributed case where leechers cooperate.

In [20], a mobile P2P video streaming and benchmarking platform is presented, which enables to assess and compare the energy consumption of different approaches in a precise manner through live assessments at runtime. The demonstrated platform includes a simple, yet high-performance tree-based mobile P2P streaming overlay which can be utilized to easily implement and assess further streaming overlay approaches.

In [23], a novel P2P overlay for Energy Level discovery in a sensor network is designed, the so-called Energy Level Distributed Tree (ELDT). Sensor nodes are mapped to peers based on their energy level. As the energy levels change, the sensor nodes would have to move from one peer to another

and this operation is the most crucial for the efficient scalability of the proposed system. Similarly, as the energy level of a sensor node becomes extremely low, that node may want to forward it's task to another node with the desired energy level. It is experimentally verified that it achieves the best-known query performance of the above operation via an appropriate simulator designed for this purpose.

III. LOAD BALANCING IN A CLUSTER – BASED P2P SYSTEM

In this section, we shortly describe the model assumed and the algorithm we plan to reconsider as presented in [2], where a reader must refer for more details.

We consider a system with N Peers ($P_1 \dots P_n$), organized in M Clusters ($C_1 \dots C_m$), with each peer belonging to only one cluster [2] (see Figure 1). After joining the system a peer advertises its characteristics such as connectivity bandwidth, processing capabilities, mean online time of past connections, which are used for peer classification as leaf, SuperNode or Candidate SuperNode. Candidate SuperNodes are chosen to maintain copies of the above information, in order to replace the SuperNodes they are connected to, if they leave the system. Inside a cluster, the SuperNodes form a highly interconnected network, exchanging important information such as file index entries, peer's load and decision messages. Moreover, the SuperNodes maintain connection with SuperNodes of neighboring clusters, in order to exchange information and decide on peer migration or even clusters merging. Each SuperNode maintains a maximum number of connections with leaf nodes, while each cluster holds a maximum number of SuperNodes [2].

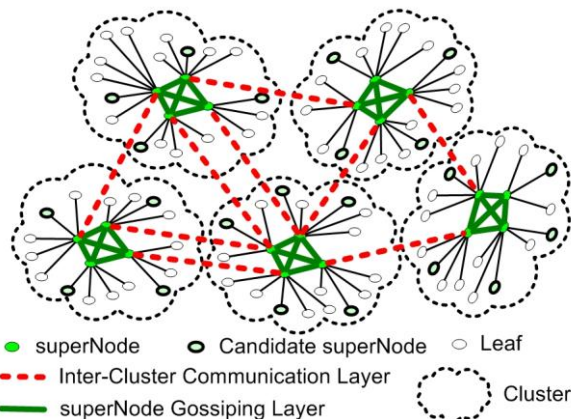


Figure 1. The architecture of the assumed P2P system.

When a peer joins the system, he publishes to the connected SuperNode information of documents he is willing to share and any other information needed to help a SuperNode to build a local metadata repository with information about all leaf nodes connected to him. A gossiping protocol like PlanetP [3] can be used to diffuse information to other SuperNode in the cluster [2] (see Figure 2). So, each SuperNode updates its repository with the received data.

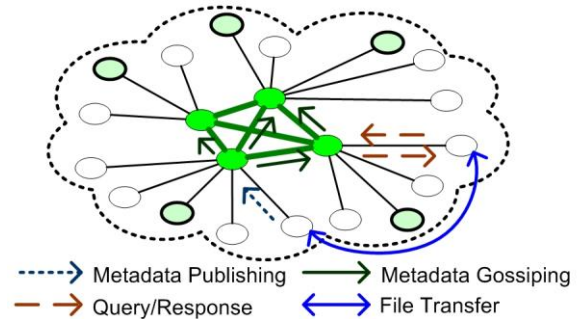


Figure 2. Communication in the P2P system.

To locate the node which will serve the response to a query (download) inside a cluster, a SuperNode makes a decision based on information about peers connection bandwidth or received load, leading to two different response strategies, as following:

MaxBw response strategy. When using this strategy, superNodes try to locate among all peers hosting the requested document, the one with the maximum connection bandwidth. This is not a load balancing response strategy but an approximation to the users' greedy behavior. We implement this strategy to use it as a basis to compare the performance of the rest strategies.

MinL response strategy. When using this strategy, superNodes try to locate among all peers hosting the requested document, the one with the minimum received load. This strategy may be unfair as far as the used bandwidth and the transfer rate is concerned but it distributes the load to more peers. The load is defined based on the consumed bandwidth for data transfers [2].

IV. COMBINING LOAD BALANCING WITH ENERGY SAVING

First of all, we need to introduce the concept of energy efficiency, which is based on [1]. The research presented in [1] is focused in energy efficiency in a data center environment. However, it considers every single computing system, which allows us to adjust it and take useful results out of it.

The final goal is to compute Peer Performance Per Energy (PPPE), which means work production per carbon energy of each individual peer. It is possible to calculate PPPE from 3 sub-metrics, Equipment Utilization (EU), Energy Efficiency (EE) and Green Energy Coefficient (GEC).

EU is a sub-metric to promote reduction in energy consumption by improving utilization rate of equipment and reduction of surplus equipment investment, computed as following:

$$EU = \frac{\text{Total measured energy (Wh)}}{\text{Total specification energy (Wh) (nameplate power rating)}} \quad (1)$$

EE represents an average of rated energy efficiency specific to the equipment. It is a metric showing efforts to procure energy saving equipment, given by the following equation:

$$EE = \frac{\text{Equipment rated Work capacity}}{\text{Total rated energy (Wh)}} \quad (2)$$

where

$$\begin{aligned} \text{Equipment rated Work capacity} = \\ \alpha * (\text{CPU Watts/CPU MTOPS}) + \\ \beta * (\text{Watt/Gbyte of memory capacity}) + \\ \gamma * (\text{Watt/Gbps of Network traffic}) \end{aligned} \quad (3)$$

where MTOPS in (3) means Million Theoretical Operations per second.

Factors α , β and γ are determined so that EE should become 1, if all the equipments have standard energy saving performance based on 2005 standard [1]. If all equipments have the performance that doubles the standard performance based on 2005 standard as of 2009 [22], the EE becomes 2 [1].

GEC provides ratio of renewable energy generated on-site to total energy consumed. Green energy purchased from external organization is not included in this metric. The maximum value of GEC should be limited to 0.8 [1]. GEC is computed as following:

$$GEC = \frac{\text{Total Measured Green Energy (Wh)}}{\text{Measured Energy Consumption (Wh)}} \quad (4)$$

PPPE, mentioned earlier can be calculated using (1), (2) and (4) sub-metrics, where higher value means better energy efficiency. Taking into consideration the standards of 2009 mentioned in [1], results that the range of PPPE is from 0 to 10.

$$PPPE = EU * EE * \frac{1}{1 - GEC} \quad (5)$$

We try to evaluate the energy consumption of each individual peer by using PPPE, as shown in (5) and try to form the connections between them based on that particular number. This will hopefully lead to a peer-to-peer network whose total energy consumption will be lower than a usual one.

In the original load balancing algorithm, two strategies were developed to response to search queries made by nodes [2]. A node that requests a specific document contacts the superNode which responds with the node which will serve the transfer. This decision was based on information about peers' connection bandwidth or received load.

We introduce three additional response strategies, except of those presented in [2], that involve the factor of energy efficiency.

- Energy response strategy. When using this strategy, superNodes try to locate among all peers hosting the

requested document, the one with the maximum PPPE value. This response strategy performs poorly regarding load balancing metrics, as it will be shown later, because it is solely based on PPPE. In this case, peers with high energy efficiency degree collect the majority of requests. However, it offers the maximum possible energy efficiency and it is used as a basis to compare the performance of the rest strategies.

- MaxBw and Energy response strategy. When using this strategy, superNodes try to locate among all peers hosting the requested document, the one which combines the maximum bandwidth and maximum PPPE, according to the following equation:

$$\text{MaxBw_and_Energy} = \alpha * \text{maxBw} + (1 - \alpha) * \text{PPPE} \quad (6)$$

Coefficient α in (6) defines the weight of each response strategy.

- MinL and Energy response strategy. When using this strategy, superNodes try to locate among all peers hosting the requested document, the one which combines the minimum load and maximum PPPE, according to the following equation:

$$\text{MinL_and_Energy} = \alpha * \text{minL} + (1 - \alpha) * \text{PPPE} \quad (7)$$

Coefficient α in (7) defines the weight of each response strategy. This method tries to balance the load between the peers, but also keeps an acceptable value in energy efficiency.

In order to tackle the free-riders or users that share documents with low popularity, the idea of replication has been adopted, where a document is copied to another peer in order to increase it's availability in the cluster [2]. The replication strategies used are the same with the response strategies mentioned above. That means for each simulation in our experiment, we have to decide which response and replication strategy will be used. Besides replication being an important factor of the simulation, it's search method has less effect than the response strategy on the overall results.

V. THE EXPERIMENT

A. Simulation Program Details

We adopted the same simulation algorithm that was used in [2]. It is a time-stepped-event flow-based simulation developed in MATLAB. The simulation process proceeds in rounds. In every round each node may query for documents. Then, nodes respond to these search queries triggering the execution of the replication algorithm. At the end of each round, there is a second phase during which the simulation of the flows among the network's nodes takes place.

The network parameters that are used for the experiments are based on the measurements on real systems. We simulate and evaluate the performance of a single cluster

consisting of 500 nodes in total, where 5 of them are considered as SuperNodes and 25 of them as energy efficient nodes whose PPPE value was randomly selected. The servicing node's energy value is set on querying node's queue, each time the last requests a document.

Factor α variable was set to 0.5. That provides equal participation of both strategies and avoids situations where the one technique conquers the other.

The rest parameters are exactly the same as the original algorithm described in [2].

Metrics

Fairness Index (FI): This is a system related metric which shows the distribution of load among peers. This index is measured based on fairness index as defined on by the formula [2]:

$$f(x) = \frac{\left[\sum_{i=1}^n x_i \right]^2}{n \sum_{i=1}^n x_i^2} \quad (8)$$

In our system n stands for the number of peers in cluster and x_i the received load on peer i . The value of the fairness index is always between 0 and 1. The closer the value is to 1 the fairer the load distribution becomes.

Quality of Availability: QoA is defined as the fraction of accepted requests for data transfer on the total requests for data transfer.

Throughput: This metric is defined as the total number of finished document transfers. The document transfers initiated by the replication process are not counted here.

Energy efficiency: This metric is defined as the average energy value, which is PPPE multiplied by 10, of the servicing nodes. The higher the value, the higher energy efficiency is observed on the system.

B. Comparative Results

Keeping in mind the original results, we run the three following scenarios. The replication method was set as minimum Load with a $P_p = 0.2$.

Sc-1 Energy response strategy: This simulation is based solely on Energy factor, it is trying to achieve the maximum Energy efficiency level on the cluster and thus FI is left unattained.

Sc-2 MaxBw and Energy response strategy: We simulate a greedy nodes' behavior combining also Energy factor.

Sc-3 MinL and Energy response strategy: There we merge minimum Load and Energy.

Sc-4 MinL (Original) response strategy: This comes from the original simulation based purely on minimum Load. We compare with the results of this strategy, because it was derived the best results among other heuristic strategies, as it was shown in [2].

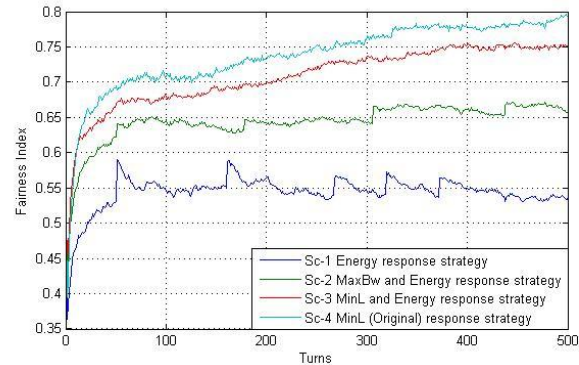


Figure 3. Fairness Index vs Simulation Turns

In Figure 3, we show the FI during the simulation. As we can see the original MinL approach performed better than the others. But as we can see, the combined MinL and Energy response strategy achieves competitive results, with a little loss in fairness.

TABLE I. SIMULATION METRICS

	#DocReq	#AccReq	QoA	En. Eff.	Thr/put
Sc-1	48977	17289	35.3%	4.4251	15025
Sc-2	38116	27684	72.6%	2.8237	22060
Sc-3	34504	23265	67.4%	3.5885	20106
Sc-4	33951	24713	72.8%	2.199	20804

As far the QoA and Throughput are concerned, **TABLE I** summarizes the results. **TABLE I** illustrates also the number of Documents Requested (#DocReq), the number of Accepted Requests (#AccReq) and the energy efficiency. As it was expected, pure Energy response strategy behaves poorly in comparison with the other 3. QoA of Sc-3 loses 5.4% from the original results (Sc-4) but there is a gain of 63% in energy efficiency. The reason that Sc-4 shows a high degree of energy efficiency is due to the fact that some nodes with high PPPE are chosen for their other attributes, such as low load in this case. Throughput's minor differences have not a significant impact on performance. So, it is clearly shown from the results of **Σφάλμα! Το αρχείο προέλευσης της αναφοράς δεν βρέθηκε.** that we can save much of energy with a little loss in load balancing metrics.

VI. CONCLUSIONS AND FUTURE WORK

In this work, we consider a load balancing algorithm [2] for an unstructured P2P system that is based on a partially centralized architecture. We change it properly to lower the power consumption of the whole system, without performance unacceptable loss in load balancing measures. To calculate the power consumption of a peer, we use the energy efficiency metric, appeared in [1] and introduce three new heuristics to balance the load between the peers of the system. Simulation results show that we can save much of energy with a little loss in load balancing metrics.

The plans for our future work involve the creation/alteration of load balancing algorithms, to address the power consumption problem for cloud systems. We also plan to measure the amount of saved energy in real life and not only via simulation experiments.

ACKNOWLEDGEMENT

This research is implemented through the Operational Program "Education and Lifelong Learning" and is co-financed by the European Union (European Social Fund) and Greek national funds.

REFERENCES

- [1] Web Reference: Takao Shiino. "Green IT by all parties" Nomura Research Institute (NRI), March 4, 2010 <http://www.oecd.org/sti/ind/45009540.pdf>, [retrieved: 4,2014].
- [2] J. Garofalakis and T. A. Michail, "Load Balancing in a Cluster-Based P2P System", Fourth Balkan Conference in Informatics, IEEE, Thessaloniki, 17-19 Sept 2009, pp. 133-138.
- [3] F. Cuenca-Acuna, F. Peery, R. P. Martin, and T. D. Nguyen, "PlanetP: Ysing Gossiping to Build Content Addressable Peer-to-Peer Information Sharing Communities", Technical Report DCS-TR-487, Rutgers University, Sept. 2002.
- [4] S. Androutsellis-Theotokis and D. Spinellis, "A Survey of Peer-to-Peer Content Distribution Technologies", ACM Computing Surveys, Vol.36, No.4, December 2004, pp.335-371.
- [5] S. Saroiu, P. Krishna Gummadi and S. Gribble, "Measuring and analyzing the characteristics of Napster and Gnutella hosts", ACM Multimedia Systems Journal, Vol. 9, Issue 2, August 2003, pp. 170-184.
- [6] Web Reference: Gnutella Protocol Development, http://rfc-gnutella.sourceforge.net/src/rfc-0_6-draft.html, [retrieved: 4/2014].
- [7] Web Reference: BitTorrent Protocol Specification, http://www.bittorrent.org/beps/bep_0003.html, [retrieved: 4/2014].
- [8] V. Valancius, N. Laoutaris, L. Massoulie, C. Diot, and P. Rodriguez, "Greening the Internet with Nano Data Centers," Proc. of ACM CoNEXT'09, 2009.
- [9] A. Qureshi, R. Weber, H. Balakrishnan, J. Guttag, and B. Maggs, "Cutting the electric bill for Internet-Scale systems," ACM SIGCOMM 2009, Barcelona (Spain), August 17-21 2009, pp. 123-134,.
- [10] EPA Report to Congress on Server and Data Center Energy Efficiency, US Environmental Protection Agency, ENERGY STAR Program, 2007.
- [11] I. Stoica, R. Morris, D. Karger, M. Kaahoeck, and H. Balakrishnan, "Chord: A scalable peer-to-peer lookup service for internet applications", In Proceedings of SIGCOMM 2001, San Diego (California), August 27-31, pp 149-160.
- [12] S. Ratnasamy, P. Francis, M. Handley, and R. Karp, "A scalable content-addressable network", In Proceedings of SIGCOMM 2001, San Diego (California), August 27-31, pp. 161-172.
- [13] B. Zhao, J. Kubiawicz, and A. Joseph, "Tapestry: An infrastructure for fault-tolerant wide-area location and routing", Tech. Rep. UCB/CSD-01-1141, Computer Science Division, University of California, Berkeley, 94720. (April).
- [14] G. On, J. Schmitt, and R. Steinmetz, "The effectiveness of realistic replication strategies on quality of availability for Peer-to-Peer systems", In. Proc. Of the 3rd International Conference on Peer-To-Peer Computing, Linkoping (Sweden), 1-3 September 2003, pp. 57-64.
- [15] M. Azami and B. Bhargava, "Energy-Efficient Peer-To-Peer Caching and Mobility Management in 4G Hybrid Networks", TR08-030, Purdue University, 2008.
- [16] A. Bhole and B. Nandwalkar, "Reducing Power Consumption in Peer-to-Peer System", International Journal on Computer Science and Engineering (IJCSSE), Vol.4, No. 06, June 2012, pp. 1199-1203.
- [17] I. Giannetti, G. Anastasi, and M. Conti, "Energy-Efficient P2P File Sharing for Residential Bit Torrent Users", IEEE Symposium on Computers and Communications (ISCC), Cappadocia (Turkey), 1-4 July 2012, pp. 524-529.
- [18] H. Hlavacs, K. Hummel, R. Weidlich, A. Houyou, A. Berl, and H. de Meer, "Energy Efficiency in Future Home Environments: a Distributed Approach", Home Networking Conference, Paris, 10-12 December 2007, pp. 69-84.
- [19] H. Hlavacs, R. Weidlich, and T. Treutner, "Energy Efficient Peer-to-Peer File Sharing", The Journal of Supercomputing, Vol. 62, issue 3, December 2012, pp. 1167-1188.
- [20] M. Wichtlhuber, J. Ruckert, D. Stingl, M. Schulz, and D. Hausheer, "Energy- Efficient Mobile P2P Video Streaming", IEEE 12th International Conference on Peer-to-Peer Computing, 3-5 Sept. 2012, Tarragona (Spain), pp. 63-64.
- [21] K. Park and P. Valduriez, "Energy Efficient Data Access in mobile P2P Networks", IEEE Transactions on Knowledge and Data Engineering, Vol. 23, No. 11, November 2011, pp. 1619-1634.
- [22] ENERGY STAR, "Program requirements for Computers", Version 5.0
- [23] S. Sioutas, K. Oikonomou, G. Papaloukopoulos, M. Xenos, and Y. Manolopoulos, "Building an efficient P2P overlay for energy – level queries in sensor networks", Proceedings of the ACM International Conference on Management of Emergent Digital EcoSystems (MEDES), 2009.

A Comparative Analysis of Parallel Programming Models for C++

Arno Leist Andrew Gilman
 Institute of Natural and Mathematical Sciences
 Massey University
 Auckland, New Zealand
 {a.leist, a.gilman}@massey.ac.nz

Abstract—The parallel programming model used in a software development project can significantly affect the way concurrency is expressed in the code. It also comes with certain trade-offs with regard to ease of development, code readability, functionality, runtime overheads, and scalability. Especially the performance aspects can vary with the type of parallelism suited to the problem at hand. We investigate how well three popular multi-tasking frameworks for C++ — Threading Building Blocks, Cilk Plus, and OpenMP 4.0 — cope with three of the most common parallel scenarios: recursive divide-and-conquer algorithms; embarrassingly parallel loops; and loops that update shared variables. We implement merge sort, matrix multiplication, and dot product as test cases for the respective scenario in each of the programming models. We then go one step further and also apply the vectorisation support offered by Cilk Plus and OpenMP 4.0 to the data-parallel aspects of the loop-based algorithms. Our results demonstrate that certain configurations expose significant differences in the task creation and scheduling overheads among the tested frameworks. We also highlight the importance of testing how well an algorithm scales with the number of hardware threads available to the application.

Keywords—parallel programming models; performance; TBB; Cilk Plus; OpenMP 4.0

I. INTRODUCTION

With the widespread availability of parallel processors, and the focus on even more parallel execution units in upcoming chips from all major manufacturers, the need for simple but efficient programming models specifically designed for parallel computing is becoming increasingly apparent.

For software developers, this means a rethinking of the way code is written. We can no longer expect that our applications will automatically run faster with every new processor generation, unless they are able to dynamically scale to large numbers of threads that can be assigned to processing units by a runtime scheduler in such a way that a good load balance is achieved. And all of this should be done with as little overhead as possible, both from the perspective of the programmer and from a performance perspective. While some overheads are algorithm specific, such as a need for thread local storage, others are inherent to the programming model and parallel runtime system. It is important to choose a parallel programming model that not only performs well on the current generation of widely deployed quad to hexa-core processors, but also on new generations of many-core processors.

We investigate how well three of the most popular task-based parallel programming models for C++ — Intel Threading Building Blocks (TBB) [1], Intel Cilk Plus [2] and OpenMP 4.0 [3] — perform when faced with three different but equally common parallel scenarios: divide-and-conquer style algorithms (merge sort), for-loops with no data contention (matrix multiplication) and loops performing a reduction operation (dot product). We want to determine whether the way concurrent tasks are

spawned has a significant effect on the runtime schedulers. Secondly, we are also interested in the support for parallel reduction operations offered by the frameworks, as these are difficult to do in parallel while at the same time retaining good scalability to the large numbers of execution units found in many-core devices. Finally, we also look at the support for vectorisation offered by Cilk Plus and OpenMP 4.0, and how it affects the results of the loop-based algorithms.

While there are already a number of existing publications comparing some subset or superset of TBB, Cilk Plus and OpenMP, they tend to focus on a particular parallelisation strategy, such as subdividing the range of a loop or a recursive algorithm. Most of these articles also do not yet utilise the constructs for expressing opportunities for vectorisation to the compiler, which have been added in Cilk Plus and OpenMP 4.0. These extensions are becoming increasingly important, given the trend towards wider vector units in general purpose CPU cores, and the tendency towards heterogeneous architectures that tightly integrate data-parallel graphics processing units with the CPU to use them as compute accelerators [4]. The ability to easily and efficiently combine task and data-parallelism will be essential for good performance scaling on upcoming hardware.

The following section puts our work into context with existing publications. Section I-B then provides a brief introduction to each of the chosen programming models, which is followed by the implementation of the selected algorithms in Section II. Section III gives the results of the tests, comparing the performance of each of the parallel implementations with respect to a serial reference implementation. In Section IV, we discuss the results in more depth, before we draw some conclusions from the findings in Section V.

A. Related Work

Articles that focus on only one type of parallel problem are able to go into a significant level of detail with regard to implementation decisions, but do not demonstrate how well the programming models deal with a variety of use cases, all of which can easily occur in a single application. For example, matrix operations, in particular matrix multiplication, are a common example for the parallelisation of loop-constructs, as demonstrated in [5] for TBB, Cilk++ (the predecessor to Cilk Plus, developed by Cilk Arts before their acquisition by Intel), OpenMP, the Open Message Passing Interface (Open MPI) [6] and POSIX (Portable Operating System Interface) threads (Pthreads) [7]; or in [8] with TBB, Cilk++, OpenMP, Pthreads, Intel Array Building Blocks (ArBB) [9] and a number of less well known programming models.

The authors of [10], on the other hand, chose the well known Mandelbrot algorithm to compare TBB, OpenMP and Pthreads.

This is similar to matrix multiplication, in that the algorithm's outermost loop is straight forward to parallelise. While the authors briefly discuss an implementation that uses SSE vector intrinsics to further speed-up the computation, none of the chosen programming models offered the necessary support on its own at the time of publication.

Jarp et al. [11] compare the performance of an algorithm for high energy physics data analysis when implemented in TBB, OpenMP and Cilk Plus. Where possible, the code is vectorised using either Cilk Plus array notation or the auto-vectorisation capabilities of the Intel C++ compiler. While more complex than in the articles mentioned so far, the parallelisation is also based on executing loop iterations concurrently, with no interdependence between the iterations of the first few loops and the need for a reproducible reduction operation to sum the results of the final loop.

In [12], the authors compare the performance of irregular graph algorithms implemented using TBB, Cilk Plus, and OpenMP 3.0 on Intel's MIC architecture based "Knights Ferry" prototype, the predecessor of the first officially released generation of Intel Xeon Phi coprocessors. The algorithms implemented are the memory bandwidth intensive graph colouring and breadth-first search algorithms, as well as a synthetic algorithm that artificially increases computation over the same data. Again, loops over vertex and adjacency sets are used to expose parallelism to the programming models.

In [13], Krieger et al. provide a good discussion of asynchronous parallelism represented using task graphs. They compare implementations of a sparse matrix solver and a molecular dynamics benchmark using TBB, which offers explicit support for task graphs, to OpenMP and Cilk Plus. Neither of the latter two models supports the concept of task graphs out of the box. Instead, the authors map the concepts of loop-based and task-spawning based parallelism to task graphs.

Podobas et al. [14] focus on recursive algorithms and specifically the overhead of task creation, spawn and join operations when using Cilk++, several different implementations of OpenMP 3.0, as well as a multi-tasking library developed by one of the authors of the paper.

Although this is not an exhaustive overview of the literature, it illustrates how our article approaches the problem from a more general perspective than what has been done before.

B. The Programming Models

This section gives a brief introduction to the programming models. Specific features that are relevant to the article are discussed in more detail in the implementation section.

1) *Threading Building Blocks (TBB)*: is a portable C++ template library for the development of concurrent software using task parallelism. A runtime scheduler is responsible for maintaining a thread pool and efficiently mapping tasks to worker threads for execution. In doing so, it abstracts the programmer from the platform-specific threading libraries, and also from the details of the hardware, such as the number of physical cores available to the application.

TBB offers a variety of constructs to express concurrency, ranging from simple parallel loops to flexible flow graphs. It also provides a scalable memory allocator and a number of concurrent data structures, such as a thread-safe vector, queue, and hash map implementation.

One of the advantages of TBB compared to Cilk Plus and OpenMP is that it does not require any special compiler

support. Any reasonably modern C++ compiler is able to build a program that uses TBB.

2) *Cilk Plus*: introduces three new keywords to C/C++: `_Cilk_for`, `_Cilk_spawn`, and `_Cilk_sync`. In usage and throughout much of the official documentation, it is common to include the header `cilk/cilk.h`, which defines the alternative and more streamlined keywords `cilk_for`, `cilk_spawn`, and `cilk_sync` as synonyms for the language keywords. As the names suggest, the keywords are used to express concurrent code sections that can be executed in parallel by the runtime scheduler. Code sections that are bounded by these parallel control statements, or the start and end of the application, but that do not contain any such constructs, are called "strands" in Cilk terminology.

Cilk Plus also adds array notations, which provide a natural way to express data-parallelism for array sections. For example: `c[0 : n] = a[0 : n] + b[0 : n]` performs an element by element addition of arrays `a` and `b` and writes the results to `c`. Furthermore, `#pragma simd` can be used to enforce loop vectorisation.

Cilk has to be explicitly supported by the compiler. At the time of writing, this includes stable support in the Intel `icc/icpc` compiler, and a `cilkplus` development branch for GCC 4.9. We would have liked to test the performance of the Cilk extensions in GCC, but the use of any SIMD functionality—array notations or `#pragma simd`—caused internal compiler errors when using the current development version of GCC 4.9 at the time of writing this article (January 2014).

3) *OpenMP*: consists of a number of compiler directives, introduced into the code using the `#pragma` keyword, runtime library routines, and environment variables. Concurrency is generally expressed with parallel regions that are denoted by a `#pragma omp parallel` directive. The master thread is forked at the beginning of each such region into a number of worker threads. At the end of each parallel region, the worker threads are joined back in after implicit synchronisation. Inside the regions, OpenMP provides an option to use a selection of work-sharing constructs that distribute the work amongst the worker threads.

II. IMPLEMENTING THE ALGORITHMS

We choose three algorithms with very different parallel characteristics to compare how each of the programming models and associated runtime libraries handles the given challenges.

A. Merge Sort

The first is merge sort, representing the widely applicable category of recursive divide-and-conquer algorithms. The implementation can spawn recursive tasks during both the sorting and merge phases to increase the amount of work that can be done in parallel. This ensures that even in the last steps of the algorithm, at the top of the merge sort hierarchy, where the longest, sorted subsets of the input data are merged, plenty of concurrent work is available to keep all worker threads busy.

Algorithm 1 gives the pseudo-code for function `sort`. The implementation recursively subdivides the input range until the threshold of length ≤ 32 is reached. At this point, it switches to the insertion sort algorithm, which has a higher worst case computational complexity of $\mathcal{O}(n^2)$, but only a small constant overhead, which makes it more efficient for very short arrays.

As can be seen, the recursive calls on lines 8 and 9 operate on non-overlapping sections of the data array, [low, mid]

Algorithm 1 The implementation of the merge sort algorithm subdivides the input range until the threshold of length ≤ 32 is reached. It then switches to the insertion sort algorithm, which offers better performance for small inputs.

```

function sort(data, aux, low, high)
1: if low < high then
2:   length  $\leftarrow$  high + low + 1
3:   if length > 32 then
4:     mid  $\leftarrow$  (low + high)/2
5:     //Note: This is where the parallel implementations
6:     //check if (mid-lo) > cutoffsort before they spawn
7:     //the following calls to sort() as concurrent tasks.
8:     call sort(data, aux, low, mid)
9:     call sort(data, aux, mid + 1, hi)
10:    call merge(data, aux, low, mid, high)
11:   else
12:     call insertion_sort(data, low, high)
13:   end if
14: end if

```

and [mid+1, hi] respectively. Therefore, they can be safely executed in parallel. Depending on the programming model and implementation, either both calls can be spawned as child tasks while the parent task is suspended until they have completed their work, or only the first call gets spawned and the parent processes the second recursive call itself. The second version intuitively introduces a smaller overhead, as less tasks are spawned. In either case, the spawned tasks and their parent must be synchronised before `merge` is called, which combines the two independently sorted sub-ranges.

The serial implementation uses a straight forward merge function that expects the two ranges to be consecutive in memory. It copies the first half into an auxiliary data array before it combines the two ranges in the original data array. This is quite efficient, because memory is accessed sequentially, but it is not obvious how to parallelise this operation. Therefore, the recursive merge algorithm from [15] is implemented and used for all parallel implementations. The pseudo-code in Algorithm 2 illustrates the procedure.

The merge operation picks the median index of the first range, idx_1 , and performs a binary search for the value found at this index in the second range. Assuming an ascending sort order, the search returns an index, idx_2 , according to these rules:

- If the search range is empty, then it returns low_2
- If $in[idx_1] \leq in[low_2]$, then it returns low_2
- If $in[idx_1] > in[low_2]$, then it returns the *largest* index in the range $[low_2, high_2+1]$ such that $in[idx_2-1] < in[idx_1]$

Note that our implementation uses C++ templates to pass a function pointer or functor to the procedures. The operation defined by this function is used for all comparisons, and, therefore, defines how data of a given type is sorted.

On line 10, the algorithm copies the value from the median index to the output array, and then recursively calls itself twice. The first call passes the lower parts of the two input ranges to `merge_recursive`, while the second call passes the upper parts. Just like before, the recursive calls operate on non-overlapping data regions and can be safely performed in parallel. However, since this algorithm does not traverse memory sequentially and introduces computational overhead, it runs slower than the simple serial merge implementation. Therefore, once sufficient tasks have been spawned, it is beneficial to switch to a more efficient serial merge operation for the shorter data ranges on the lower levels of the hierarchy. A cutoff is defined (line 6) that causes the implementation to do exactly

Algorithm 2 The implementation of the recursive merge function is based on the algorithm described in [15]. It merges two sub-ranges of array `in`, $[low_1, high_1]$ and $[low_2, high_2]$, and writes the result to array `out` beginning at index low_{out} .

```

function merge_recursive(in, low1, high1, low2, high2,
                        out, lowout)
1: length1  $\leftarrow$  high1-low1 + 1
2: length2  $\leftarrow$  high2-low2 + 1
3: if length1  $\leq$  length2 then
4:   swap low1  $\leftrightarrow$  low2, high1  $\leftrightarrow$  high2, length1  $\leftrightarrow$  length2
5: end if
6: if length1 + length2 > cutoffmerge then
7:   idx1  $\leftarrow$  (low1 + high1)/2 //pick the median index
8:   idx2  $\leftarrow$  call binary_search(in[idx1], in, low2, high2)
9:   idx3  $\leftarrow$  lowout + (idx1 - low1) + (idx2 - low2)
10:  out[idx3]  $\leftarrow$  in[idx1]
11:  call merge_recursive(in, low1, idx1 - 1, low2, idx2 - 1,
                        out, lowout)
12:  call merge_recursive(in, idx1 + 1, high1, idx2, high2,
                        out, idx3 + 1)
13: else if length1 > 0 then
14:  call merge(in, low1, high1, low2, high2, out, lowout)
15: end if

```

that. This third merge function works like the first one, except that the two input ranges do not have to be consecutive, as this is the case when switching from the recursive merge operation part way down the hierarchy.

Now that the merge sort algorithms have been described, we can have a look at the different ways used to spawn the recursive calls in the programming models.

TBB offers a number of ways to spawn concurrent tasks that call functions, two of which are well suited for this algorithm: manual task management and the template function `tbb::parallel_invoke`. The former, more flexible but also more complex option, is used in function `sort`. It defines two types derived from `tbb::task`, one to perform the recursive subdivision of the sort routine, the other to initiate the recursive merge operation. This code is quite verbose and is therefore not included here, but examples of the concept can be found in the reference documentation of the task scheduler [1] under continuation-passing style.

The implementation of `recursive_merge`, on the other hand, uses `tbb::parallel_invoke` with the C++11 lambda function syntax to spawn the recursive calls and block until they complete. The resulting code to implement lines 11 and 12 of Algorithm 2 is concise, as the following listing shows (with template arguments omitted for readability):

```

tbb::parallel_invoke (
  [&]{ merge_recursive(in, low1, idx1-1,
                      low2, idx2-1, out, lowOut); },
  [&]{ merge_recursive(in, idx1+1, high1,
                      idx2, high2, out, idx3+1); } );

```

The **Cilk Plus** code for spawning functions is even simpler, as the following listing demonstrates (with both template and function arguments omitted for brevity):

```

cilk_spawn merge_recursive(...);
merge_recursive(...);
cilk_sync;

```

The code essentially behaves like one would expect: the first call to `merge_recursive` is spawned as a separate task, whereas the second call is executed by the current thread, until the `cilk_sync` statement synchronises both execution paths. The actual implementation of Cilk Plus works slightly

differently, in that the worker thread that performs the spawn immediately executes the spawned task, while the so called *continuation* of the function, that is the statements following the spawn, can be stolen by an idle worker. The assumption is that, most of the time, other threads are busy with their own work, and the continuation is not stolen but simply executed by the initial worker once it has completed the spawned function call. This is more efficient, as it takes advantage of data locality in the processor caches.

In the **OpenMP** implementation, we utilise the tasking facility introduced in specification version 3.0. An explicit task can be created using the `task` directive, followed by a structured block. When this directive is encountered by a thread, the task may be executed right away, or placed into a pool from which all worker threads in the current team can take tasks to execute. The following extract from the `merge_recursive()` subroutine demonstrates the syntax simplicity:

```
#pragma omp task
    merge_recursive (...);
#pragma omp task
    merge_recursive (...);
#pragma omp taskwait
```

Two tasks at the next level of recursion will be created and the execution of the current task will be halted at the `taskwait` directive, until all tasks created at the current level of recursion have been completed.

B. Matrix Multiplication

The second algorithm in our comparison is matrix multiplication. Given an $n \times m$ matrix A and an $m \times p$ matrix B , the algorithm computes matrix $C = AB$ of size $n \times p$. The most straight forward implementation uses three nested loops, iterating over the n rows of matrix A , the p columns of matrix B , and, on the innermost level, over the m columns of A while multiplying the values with the corresponding cells from B .

A simple but significant optimisation is to swap the two inner loops, which allows the data for matrix B to be read in row-major order and, thus, from sequential memory addresses. This has a big impact on performance. The pseudo-code for this implementation is given in Algorithm 3. While there are even more efficient serial algorithms for matrix multiplication, this implementation is straight forward to parallelise, because it does not have any data races across the rows of matrix A , given the precondition that the memory pointed to by C does not overlap with A or B . As such, it is used to represent the category of algorithms that can be parallelised by executing independent loop iterations.

The concurrent implementations simply subdivide the iteration space of the outermost loop into chunks that can be executed in parallel. The concurrency is limited to at most n chunks of one iteration each. If we do not swap the inner two loops, then the outer two loops can be collapsed into one loop of length np , but the improved memory access pattern of the optimised code easily makes up for the more limited concurrency, as even the serial version of this implementation runs faster than the multi-threaded code for the basic approach.

To further optimise memory access, each row of the matrix arrays is padded such that the memory address of the first element is aligned to a 64-byte boundary. This means that the data in a row is *memory aligned* for vector instructions of up to 512-bit width. It also coincides with the common cache line size on x86 processors. This is relevant, because the algorithm

Algorithm 3 The algorithm for matrix multiplication.

```
1: for rowa ← 0 to n - 1 do
2:   for colb ← 0 to p - 1 do
3:     C[rowa][colb] ← 0 //initialise rowa in C
4:   end for
5:   for cola ← 0 to m - 1 do
6:     for colb ← 0 to p - 1 do
7:       C[rowa][colb] ← C[rowa][colb] +
                        A[rowa][cola] × B[cola][colb]
8:     end for
9:   end for
10: end for
```

lends itself to vectorisation of the innermost loop. The padding, along with the compiler specific alignment guarantee given by `__assume_aligned`, in theory enables the compiler to use the more efficient vector instructions for aligned memory access. In theory, because Intel's compiler refuses to do so for unknown reasons during our tests when targeting AVX, where it uses the unaligned instructions instead, even though it complies and uses aligned instructions when targeting SSE.

The programming model specific changes to the code that allow the iteration range of the loop on line 1 of Algorithm 3 to be processed in parallel, as well as the changes to vectorise the loop on line 6 where applicable, are explained below.

TBB offers the convenient `tbb::parallel_for` function, which uses a partitioner to recursively split the specified range into smaller chunks until a certain condition is fulfilled. While different partitioning strategies can be implemented, the default `tbb::auto_partitioner`, which attempts to perform sufficient splitting to balance load, generally performs quite well. This partitioner is used here. The following code listing shows the TBB implementation.

```
tbb::parallel_for(0,n,[&](const size_t rowA) {
    // loop body
});
```

The vectorisation of the innermost loop uses the Cilk Plus constructs as explained below.

The **Cilk Plus** keyword `cilk_for` is the obvious choice to parallelise the for-loop when using this programming model. The compiler converts the loop body into a function that is called recursively using a divide-and-conquer strategy, thus turning it into a directed acyclic graph of strands that each execute a chunk of up to *grain size* consecutive iterations. The grain size can be defined with the following pragma placed just before the loop: `#pragma cilk grainsize = 10`, but just like discussed for TBB, there are advantages to leaving it up to the runtime to decide what the grain size should be, especially since the pragma is a compile time construct, where the number of worker threads is usually not yet known. The signature of the `cilk_for` loop used in our example is:

```
cilk_for (size_t rowA = 0; rowA < n; ++rowA)
```

Once again, the change to the serial implementation is minimal. In fact, it is sufficient to change the definition of `cilk_for` from `_Cilk_for` to `for` to turn it into a regular for-loop, which can be very useful when debugging a concurrent program.

The second aspect we want to parallelise is the innermost loop on lines 6 to 8 using the Cilk Plus array notation:

```
cRowPtr[0:p] += aRowPtr[aCol] * bRowPtr[0:p];
```

A scalar value from A is multiplied with every element from an entire row of B , and the resulting vector is added to the

matching elements in the current row of C . It should be noted that the second argument in the array notation does not specify the last index in the range, but rather the length of the range.

The **OpenMP** loop work-sharing construct allows for simple parallelisation of the serial matrix multiplication code. Adding the `#pragma omp for` directive immediately before the outer loop indicates to the compiler that each loop iteration can be executed concurrently. The iteration range will be subdivided into equal chunks and executed concurrently by the threads in the current team.

Vectorisation of the inner loop can be achieved with the `simd` directive, introduced in specification version 4.0, by inserting `#pragma omp simd aligned(aRowPtr, bRowPtr:64)` directly before the inner loop. The aligned clause indicates to the compiler that addresses pointed to by `aRowPtr` and `bRowPtr` are aligned to 64-byte boundaries.

C. Dot Product

The dot product of two vectors a and b is defined as $a \cdot b = \sum_{i=1}^n a_i b_i$. The difficulty in parallelising this simple algorithm lies in the sum operation, which updates a shared variable across concurrent loop iterations, introducing data races unless specific precautions are taken. This is a very common problem in concurrent programming.

A common approach is to use some form of thread local storage to accumulate the values of the shared variable across all iterations of the loop executed by a given worker thread, before a reduction strategy is used to safely merge—or reduce—these intermediate results into the final value. The three frameworks used here all offer support for reduction operations, which are discussed in the following paragraphs.

TBB provides the function `parallel_reduce`, which works similar to `parallel_for`. The form we use here expects the following arguments: (range, identity, func, reduction). This shows that, instead of passing the first and last indices of the iteration range directly to the loop construct, as we did in the matrix multiplication code, `parallel_reduce` expects an object that models TBB's range concept as its first argument. This is more flexible and also an alternative option for `parallel_for`. Most commonly, the predefined `tbb::blocked_range` is used, which takes the begin and end values for the range as its arguments, as well as an optional grain size.

The second argument is the left identity value used for the reduction. This is followed by a functor (function object) that implements the body of the loop, taking a reference to the range type and an initialisation value for the reduction variable as its arguments. The last argument is another functor that is called whenever two intermediate values from different tasks need to be joined. This functor accepts two references of the type of the reduction variable, and it returns the merged value, which, in this example, is the sum of the arguments. Conveniently, the C++ standard template library defines the binary functor type `std::plus` in header `functional`, which does exactly what is needed for the join operation.

The following code listing shows the implementation of the reducer in TBB, where T is the type of the reduction variable:

```
const T total = tbb::parallel_reduce(
    tbb::blocked_range<size_t>(0, length),
    T(), /* the left identity value */
    [=](const tbb::blocked_range<size_t> &range,
```

```
    T init)
{
    for (size_t i = range.begin();
        i < range.end(); ++i)
    {
        init += vector1[i] * vector2[i];
    }
    return init;
}, std::plus<T>() /* the join operation */
);
```

Cilk Plus reducers are based on the concept of so called hyperobjects [16], which return a view of the given type when they are dereferenced. If the strand that is doing the dereferencing was stolen (i.e., it is being executed by a different worker thread than its parent), then the hyperobject returns a new instance of the view, otherwise it can safely return the same instance as before. This minimises the overhead of creating view instances. When a spawned strand finishes and merges back into its parent, the reduction operation is invoked to reduce the values of the two view instances, leaving the result in the surviving strand's view. At the end of the concurrent section, all strands have been merged back into the leftmost view, which now contains the final value. The following listing illustrates how this works in practice.

```
cilk::reducer<cilk::op_add<T>> sum{ T() };
cilk_for (size_t i = 0; i < length; ++i) {
    *sum += vector1[i] * vector2[i];
}
const T total = sum.get_value();
```

The type of the reduction operation is specified as a C++ template argument (the C syntax for reducers depends on macros and typedefs instead) to the `cilk::reducer` type, which, when dereferenced, returns a reference to a value of type T that is guaranteed not to be in use by another worker. The reduction operation must be associative, but the order of the operands is the same as in the serial computation.

OpenMP provides a facility for the reduction operation through an additional `reduction` clause attached to the loop work-sharing construct. The following listing demonstrates the ease of expressing the concurrent loop for the dot product operation using OpenMP constructs, including reduction on variable `sum`:

```
#pragma omp parallel for reduction(+:sum)
for (size_t i = 0; i < length; ++i) {
    sum += vector1[i] * vector2[i];
}
```

III. PERFORMANCE RESULTS

The test system runs Ubuntu 13.10 on a 3.4 GHz quad-core Intel Core i7-3770 and 16 GBytes of PC-1600 system memory. The processor supports 2-way simultaneous multithreading (SMT), and is, therefore, seen as having eight logical cores by the system and applications. We use the Intel C++ compiler version 14.0.1, as this is the only compiler that offers stable support for all three programming models at the time of writing.

All tests are repeated 30 times and the fastest result is selected for each, as this most accurately reflects the actual time taken by the test case, with the least interruptions from system processes that are running in the background. Non-essential processes, such as the window manager, were stopped for the test runs.

We first test the merge sort algorithm, which uses two cutoff values to decide when to stop spawning concurrent tasks as

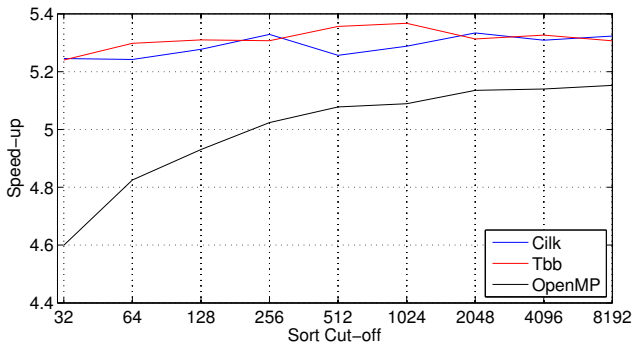


Fig. 1. Scaling the parallel sort cutoff value from 32 to 8192.

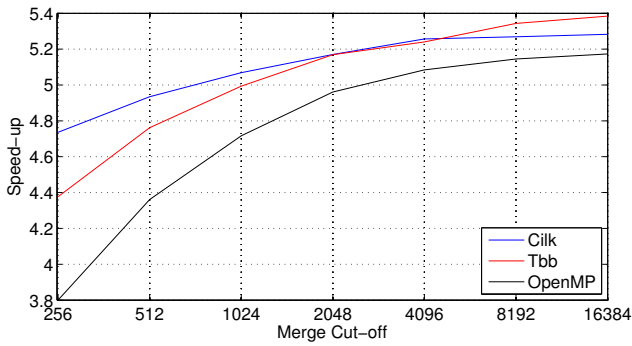


Fig. 2. Scaling the parallel merge cutoff value from 256 to 16384.

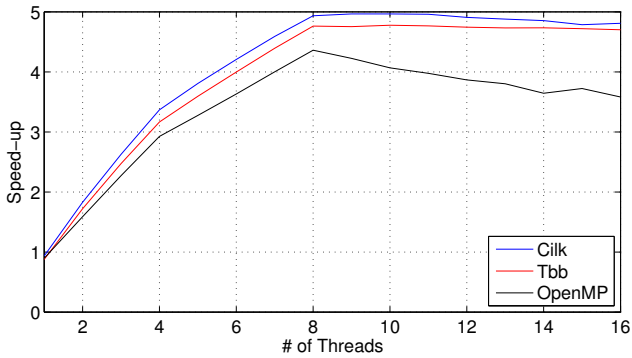


Fig. 3. The merge sort algorithm running with 1-16 worker threads. The parallel cutoffs are set to 512 and 8192 for sorting and merging respectively.

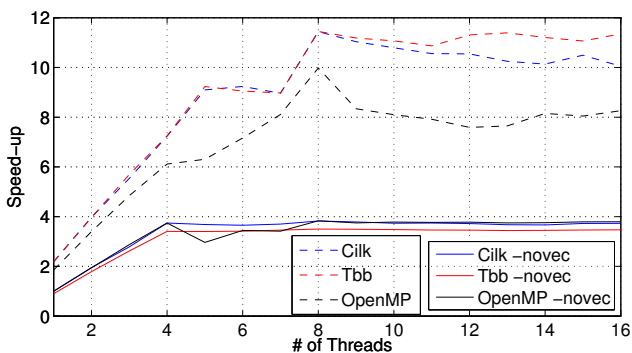


Fig. 4. Matrix multiplication with 1-16 worker threads.

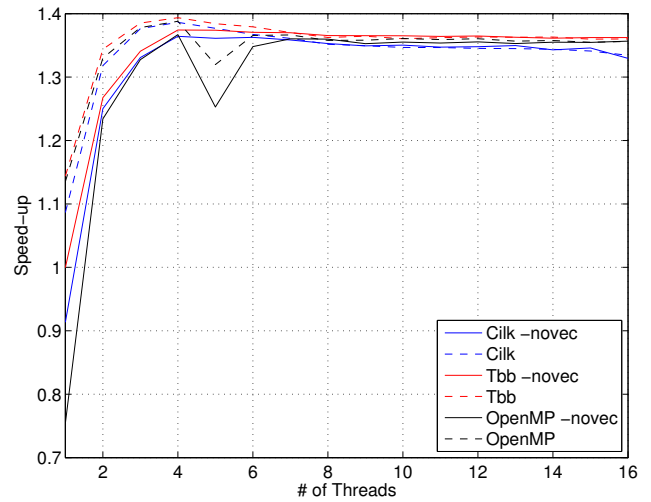


Fig. 5. Dot product with 1-16 worker threads.

indicated by $cutoff_{sort}$ on line 5 of Algorithm 1 and by $cutoff_{merge}$ on line 6 of Algorithm 2. Figure 1 gives the results for $cutoff_{sort}$ values from 32 to 8192, with the merge cutoff set to a constant 8192; and Figure 2 gives the results for $cutoff_{merge}$ values from 256 to 16384, this time with the sort cutoff set to a constant 64 for Cilk Plus, 512 for TBB, and 4096 for OpenMP. These numbers for the sort cutoff were chosen as they are the smallest values that, on average, approach the best performance achieved with the respective framework. It is important to set the cutoffs as low as possible without causing too much overhead, as they limit the amount of concurrency available at runtime.

Now that the effects of the cutoffs are established, we look at how the number of worker threads available to the process affects the performance. Figure 3 shows the results for merge sort. The data array to be sorted contains 50×10^6 integers. This test is also performed for the matrix multiplication and dot product, with the results given in Figures 4 and 5, respectively. The tests for matrix multiplication use two integer matrices of size 2500×2500 each, and the dot product implementations are tested with two vectors of 10^9 integers as input. These latter two plots also differentiate between code compiled with and without the `-novect` flag, which disables vectorisation when set.

IV. DISCUSSION

The results show that all three programming models discussed here can perform well in a wide range of common parallel scenarios, but that there are also some caveats to be considered. The OpenMP implementation in the Intel compiler tends to be a little slower than TBB and Cilk Plus. It is also more susceptible to a loss of performance than the other frameworks in four different ways.

Firstly, the results for the scaling of parallel cutoff values in the merge sort implementation (Figures 1 and 2) show that the task creation and scheduling overheads introduced by a very large number of small tasks affect OpenMP more severely than they affect Cilk Plus or TBB.

Secondly, Figure 3 clearly illustrates that the OpenMP performance suffers significantly from oversubscription of the processor when more than eight threads (i.e., the number of

logical cores) are used in the divide-and-conquer algorithm. Interestingly, the same can be seen for matrix multiplication with vectorisation enabled, but not without vectorisation. Cilk Plus also takes a slight dip in these situations, whereas TBB remains unperturbed. Only for the dot product, where performance is clearly limited by the memory bandwidth, does TBB also take a small performance hit once the thread count exceeds four (i.e., the number of physical cores).

Thirdly, the static division of the problem space for parallel for-loops only works well when all cores of the processor are fully dedicated to the current task, the number of worker threads is equal to the number of logical cores, and each loop-iteration consists of the exact same set of instructions, that is, the amount of work to be done in every iteration is the same. If the operating system, or indeed a concurrent part of the application itself, schedules some other work on one of the cores, then the static work division can unnecessarily delay the completion of the loop, potentially letting some cores go idle, as the remaining threads are not able to take some of the load from the busy core and distribute it among the idle cores. The effects of load balancing issues are evident in the significant dip of the OpenMP performance in Figures 4 and 5 with five threads, and to a lesser extent six and seven threads, only that in our tests it is caused by a number of worker threads that is not a multiple of the core count. Static work load balancing with its minimal scheduling overhead certainly has a place in the repertoire of a parallel framework, but a dynamic scheduling algorithm would appear to be a more robust default choice. It should be noted that OpenMP does offer dynamic scheduling policies as options.

The final aspect we would like to highlight with regard to OpenMP is the lower SIMD performance when using the respective pragma for matrix multiplication (Figure 4). However, it is unclear why it is slower than the other implementations or auto-vectorisation, and we intend to investigate this further by comparing the machine instructions generated by the compiler for the different implementations.

The performance of TBB is impressive given that it does not have the advantage of being integrated into the language, which means that certain optimisations that can only be done when the compiler is aware of the parallel constructs are not open to it. The drawback of TBB is that it is somewhat more verbose than the other approaches, which makes the code more difficult to read. The introduction of lambda functions into the C++11 standard has helped tremendously in this regard, allowing for much more streamlined constructs than before, but TBB still does not reach the level of integration of Cilk Plus or OpenMP. Especially the former integrates beautifully into the language, to the point where it almost becomes difficult to spot the parallel sections. However, this minimalistic feel in Cilk Plus is also partly due to the smaller number of configuration options, which may be seen as a drawback when one needs more control. OpenMP pragmas are plentiful and concise, allowing for many optional parameters. OpenMP 4.0 even adds support for a number of additional modes of operation, such as off-loading of parallel sections to an accelerator, that are not covered by either of the other frameworks.

While TBB does not include support for vectorisation, it does not hinder it being added to the algorithm by other means either, as demonstrated by the use of Cilk Plus array notation in the vectorised TBB results. These array notations are another example of Cilk's seamless integration into the language, as the resulting code is both concise and easy to understand.

V. CONCLUSION

To conclude, we would like to emphasise several important characteristics of parallel programming models: ease of development, code readability, functionality, runtime overheads, and scalability. The last point is significant because of the trend towards integrating ever higher numbers of parallel execution units into new processor architectures.

In the future, we would like to compare the performance of Cilk Plus and OpenMP with their respective implementations in GCC, and potentially other compilers as well, to determine how consistent the results given here are across different implementations. The performance of TBB is expected to offer less surprises between compilers, as it is a template library rather than a language extension.

We would also like to run the same tests on the Intel Xeon Phi, Intel's many-core architecture, to find out how well the frameworks scale to the significantly larger number of cores. This would give a good indication of how future-proof the parallel programming models are.

REFERENCES

- [1] Intel® Corporation, "Threading Building Blocks 4.2 update 2," <https://www.threadingbuildingblocks.org/>, December 2013, retrieved: April, 2014.
- [2] —, "Cilk™ Plus," <https://www.cilkplus.org/>, retrieved: April, 2014.
- [3] OpenMP Architecture Review Board, "OpenMP 4.0," <http://openmp.org/>, July 2013, retrieved: April, 2014.
- [4] A. Leist, D. P. Playne, and K. A. Hawick, "Exploiting Graphical Processing Units for Data-Parallel Scientific Applications," *Concurrency and Computation: Practice and Experience*, vol. 21, no. 18, pp. 2400–2437, December 2009.
- [5] E. Ajkunic, H. Fatkic, E. Omerovic, K. Talic, and N. Nosovic, "A Comparison of Five Parallel Programming Models for C++," in *MIPRO, 2012: Proceedings of the 35th International Convention*, P. Biljanovic, Z. Butkovic, K. Skala, S. Golubic, N. Bogunovic, S. Ribaric, M. Cicin-Sain, D. Ciscic, Z. Hutinski, M. Baranovic, M. Mauher, and J. Ulemek, Eds., Opatija, Croatia, May 2012, pp. 1780–1784.
- [6] The Open MPI Project, "A High Performance Message Passing Library," <http://www.open-mpi.org/>, retrieved: April, 2014.
- [7] *IEEE Standard 1003.1c-1995: Threads Extension*, IEEE, 1995.
- [8] P. Michailidis and K. Margaritis, "Computational Comparison of Some Multi-core Programming Tools for Basic Matrix Computations," in *Proceedings of the 14th IEEE International Conference on High Performance Computing and Communications*, G. Min, L. Lefevre, J. Hu, L. Liu, L. T. Yang, and S. Seelam, Eds., Liverpool, UK, June 25–27 2012, pp. 143–150.
- [9] Intel® Corporation, "Array Building Blocks," <http://software.intel.com/en-us/articles/intel-array-building-blocks>, retrieved: April, 2014.
- [10] W. Tristram and K. Bradshaw, "Investigating the Performance and Code Characteristics of Three Parallel Programming Models for C++," in *Southern Africa Telecommunication Networks and Applications Conference (SATNAC)*, Stellenbosch, South Africa, September 2010.
- [11] S. Jarp, A. Lazzaro, A. Nowak, and L. Valsan, "Comparison of Software Technologies for Vectorization and Parallelization," CERN openlab, Tech. Rep., September 2012.
- [12] E. Saule and U. Catalyurek, "An Early Evaluation of the Scalability of Graph Algorithms on the Intel MIC Architecture," in *Parallel and Distributed Processing Symposium Workshops PhD Forum (IPDPSW), 2012 IEEE 26th International*, Shanghai, May 2012, pp. 1629–1639.
- [13] C. Krieger, M. Strout, J. Roelofs, and A. Bajwa, "Executing Optimized Irregular Applications Using Task Graphs within Existing Parallel Models," in *High Performance Computing, Networking, Storage and Analysis (SCC), 2012 SC Companion*, Salt Lake City, Utah, US, November 2012, pp. 261–268.
- [14] A. Podobas, M. Brorsson, and K.-F. Faxon, "A Comparison of some recent Task-based Parallel Programming Models," in *Proceedings of the 3rd Workshop on Programmability Issues for Multi-Core Computers*, Pisa, Italy, January 2010.
- [15] T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein, *Introduction To Algorithms*, 3rd ed. MIT Press, 2009.
- [16] M. Frigo, P. Halpern, C. E. Leiserson, and S. Lewin-Berlin, "Reducers and Other Cilk++ Hyperobjects," in *Proceedings of the Twenty-first Annual Symposium on Parallelism in Algorithms and Architectures*. Calgary, AB, Canada: ACM, 2009, pp. 79–90.

A Study on the Performance Control for Building Façades Applied with Patterns of Hanok Openings

Seulki Kim
School of Architecture
Chonnam National University
Gwangju, Korea
selky5@naver.com

Kyoung-Hee Kim
School of Architecture
University of North Carolina Charlotte
Charlotte, NC, United States
kkim33@unc.edu

Seung-Hoon Han
School of Architecture
Chonnam National University
Gwangju, Korea
hshoon@jnu.ac.kr
(Corresponding Author)

Abstract — This study focuses on window patterns of Hanok, Korean traditional residence, and their effects on building performance and adaptabilities as building skins. The purpose of this study is to investigate the possibility of Hanok patterns as a new type of building façade that satisfy both aesthetic and the efficiency of building space, and to suggest a new method based on Information and Communication Technology (ICT) for optimizing the building performance. For this research, a variety of Korean traditional window patterns were reviewed first, and a typical façade was applied for a comparative analysis in aspect of the effect on the building indoor environment. In addition, this study examines selected façade patterns that could minimize heat loss and make it easy to adjust indoor light environment by changing the intervals. ICT has played an important role as both a controller for the spatial environment and a simulator for evaluating various experimental factors; Wireless Smart Sensor Network (WSSN) and Wireless Fidelity (Wi-Fi) have been utilized for adjusting configuration of the opening patterns, and popular simulation tools like Project Vasari and Green Building Studio (GBS) have been used for verification of the feasibility, investigating relationship between daylighting capability of proposed skins and their energy performance. Finally, building skins satisfying both Korean traditional identity and efficiency of energy performance are suggested and examined in this study.

Keywords - Hanok; Korean Traditional Window; Opening Patterns; Energy Efficiency; Energy Simulation; Sustainability

I. INTRODUCTION

Modern architecture offers not only aesthetic and user experience but also provides indoor environments that satisfy various requirements, such as visual and thermal comforts. In other words, both design and function should be taken into consideration because contemporary man’s demands are going to be greatly diversified. Nowadays, building façades world-wide look very similar without explicitly presenting their identities and tend to focus on simply available form and design. Consistent skins mean that they do not offer any specific characters responding to country, area, and/or building location. This situation causes building users to be uncomfortable in the indoor environment.

Recently, Hanok, which is a Korean traditional residence, has received attention by changing people’s awareness

because Hanok has its own beauty and is environmentally friendly [1]. Hanok consists of unique structures and components that can accommodate the four seasons of Korean climates. In particular, Hanok window is a more important component than other elements, and it occupies most of the external wall. The Hanok window in this research has been developed as control equipment of external environment that Korean season is hot and humid in the summer and cold and dry in the winter [2]. Figure 1 shows the scheme of this study that Hanok window can be used for new type of façade and increase Energy Efficiency. Hanok window also has a completely distinguished design and tectonic, consisting of a series of wooden ribs and creating unique patterns.

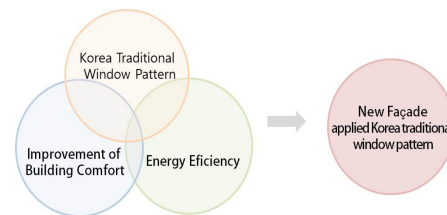


Figure 1. Scheme of the Study

The target of the experiment for this study is a general office building. A comparative study was carried out between an office building with a typical window system and one with Hanok window pattern in the study area of indoor environment and energy efficiency. Indoor environment analysis includes illuminance and insolation and energy efficiency includes fuel consumption, carbon emission quantity, and heating and cooling load [3]. There are details of analysis tools and target factors in Table I.

TABLE I. ANALYSIS TOOLS AND TARGET FACTORS

	<i>Tools</i>	<i>Factors</i>
Analysis for Building Performance	Project Vasari, GBS	Fuel Consumption, Carbon Emission Quantity, Heating and Cooling Load
	Autodesk Ecotect, Radiance	Illuminance, Insolation

A simulation was carried out by changing the interval and the scale of Hanok window pattern ribs. Project Vasari

and GBS was used to confirm energy efficiency, and indoor environment analysis was simulated by Ecotect and Radiance. An office modeling unit is a basic form provided from the Project Vasari program with code complied construction, located in Seoul, Korea where it has various climatic conditions depending on the four seasons.

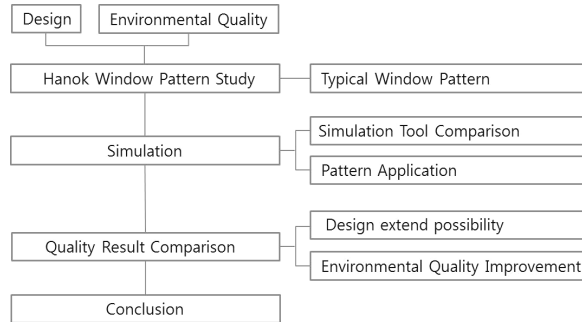


Figure 2. Overall Flow of the Study

Accordingly, this study suggests the architectural skin that is applicable to various climates and also satisfies identical design. Simulation tools were used to prove that the typical Korea window pattern presented in Section II has a positive effect on indoor environment and improves the energy efficiency. ICT such as WSSN and Wi-Fi were applied with window to do a important role as both a controller for the spatial environment and a simulator for evaluating various experimental factors as described in Section III. All data and contents were then integrated to conclude the paper in Section IV. Figure 2 is the flow chart of this study.

II. HANOK-STYLED FAÇADE AND OPENING PATTERNS

Hanok is Korean traditional residence style and is different from western-style housing in structure, space organization, and design [4]. However, previous studies on Hanok façades are limited in that they focused on improving efficiency only or applying design to modern buildings and Table II shows those examples. These studies were mostly related to replacement materials on finishing materials and windows. These existing studies are important, but integrated improvement and eco-friendly design of Hanok and Korean traditional design elements are also important.

TABLE II. EXAMPLE FOR A STANDARD HANOK DESIGN AND EXPERIMENT FOR BUILDING PERFORMANCE WITH HANOK WINDOWS

Design Alternatives for Hanok Window Patterns		
Experimental Devices for Performance Monitoring		

Therefore, the study of Hanok window is needed with a holistic integrated design approach because it acts as not only a mediator between indoor and outdoor climates, but also serves an alternative to suggest the architectural skin considering Korea identity and energy efficiency. Hanok window physically connects the indoor and outdoor of a building and can be a major influence on the indoor environment such as lighting, heat gain and ventilation. In the design, various window styles can be selected based on client's needs, and this shows a possibility for developing new façade typologies with Hanok patterns.

Hanok window patterns can be configured with different typologies according to climate condition, window location, function, and size. For example, the grid pattern is interlaced by both vertical and horizontal ribs, and formed as intersected rectangular shapes. A window graced in the form of a flower is called Flower-rib window (Kkot-sal-moon in Korean) and many other types also exist such as Wan-ja, Ah-ja, Beaum-sal, Se-sal, Kyo-sal, Ti-sal [5].

TABLE III. KOREAN TRADITIONAL WINDOW PATTERNS

Grid	Jeong-ja	Young-ja	Wan-ja	Ah-ja
Kkot-sal	Gu-gab	Sut-dea	Gu-ja	Bit-sal

These names are determined by the number of ribs, size, distance between grids, and angle, but the most affective element is a form composed by crossed ribs. These compositions of window grids play an important role as a structural element to support the window and show the beauty of window design. Figure 3 is that grid patterns as a representative typology of Hanok window is not entirely typical in the aspects of vertical and horizontal proportion, but they have a 3.6~6.0cm gap in average between grids and thickness, width, and angle can be changed by grid types.

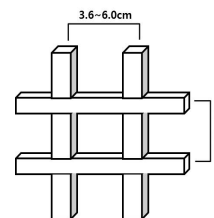


Figure 3. Typical Dimension of the Korean Window Pattern

The grid pattern, a typical layout of Hanok window pattern was used as the basic form for this study. The number of ribs, size, distance between grids, and angle of grid pattern were adjusted to find the most ideal pattern for energy efficiency and design. The adjusted grid pattern was applied to a general office building and simulated to compare its performance with a typical office building.

III. ANALYSIS FOR BUILDING PERFORMANCE APPLIED WITH HANOK-STYLED FAÇADE

A. Methods and Tools

For this study, Project Vasari, GBS, Ecotect, and Radiance were selected for simulation to compare indoor environment and energy efficiency of the target buildings mentioned on the previous section. Project Vasari uses DOE-2 as an analysis engine based on American Society of Heating, Refrigerating and Air-Conditioning Engineers (ASHRAE) standards to input parameters, so it can make relatively accurate results from the simulation. In addition, a Vasari model can produce data in gbXML/IFD format and compute various experimental settings easily by interlocking with GBS; Creating building models and analyzing given problems are executed by Project Vasari, and then GBS performs building energy analysis in this study [6].

For the indoor illuminance, Radiance and Ecotect were used to analyze indoor environment. A lighting research team at Lawrence Berkeley National Laboratory has developed Radiance, which is a program to evaluate lighting performance. This tool uses a ray-tracing logic which starts calculation on the observer's view and investigates complicated phenomenon such as reflection, curve of lighting and diffusion reflection between object surfaces. This program is free for use and can simply be plugged into Ecotect.

Ecotect is developed by Square One Research, Australia, and taken over by Autodesk in 2008. This program can predict the simulation result simply and make the data available for different building performance. This program also can estimate various elements such as sun's radiation, natural lighting, air current, heating, shade and shadow, building heating and cooling load, and so on, allowing designers to use it from the pre-design phase to the design development process. The weather tool and the solar tool are also provided by Ecotect, and make it possible to interpret weather situations and solar movement for loaded building models. Analysis data generated by this tool can easily be exported to various formats for interlocking with more applications [7].

B. Scope and Detail Levels

Since there are no specified guidelines about detail levels of building models for performance evaluation, Level of Development (LOD) 100 suggested by American Institute of Architects (AIA) has been used for modeling in this study. LOD 100 is known as an equivalent level of the conceptual design in Korea that volume of the building mass and the

building type are defined and overall modeling process is established with fundamental building parameters such as area, height, volume, place, axis and so on [8]. In addition, this detail level determines not only a type of the project ordering as suggested by Integrated Project Delivery (IPD) and an initial step for design process with Building Information Modeling (BIM) as well; The AIA documentation, E202TM-2008, includes all the steps of the IPD design process and their detail levels towards BIM [9].

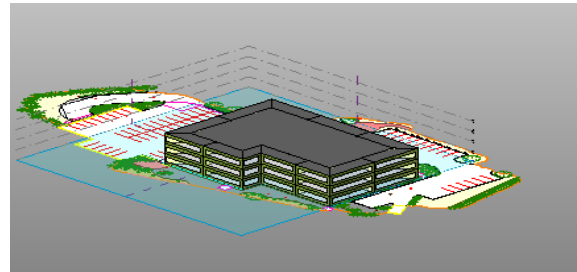


Figure 4. Building Model for Simulation Based on LOD 100

This study has utilized LOD 100 as shown on Figure 4, because the target for the performance evaluation is not for building facility and utilities, but for building façade as an architectural component. Then, building performance focusing on environmental and energy factors has been analyzed with a few alternatives of building skins chosen from both a typical curtain wall and Hanok-styled window patterns for comparison.

C. Settings and Factors

First of all, it is important to set up the climate condition to perform the building performance in aspect of the building skins facing surroundings. The greatest percentage of people on the Earth live in the temperate climate zone along the westerly belt, because this warm area has four clear seasons caused by ideal solar positions normally distributed between 20 to 60 degrees in altitude. Korea is also located in a temperate climate zone in the middle latitude belt, and Seoul, the capital city of Korea, has been selected for analysis, because about 40 percent of the office buildings in Korea are distributed in this city as shown in Table IV.

Most settings for simulation, including weather and climate factors, were extracted from the database supported by official statistics archives such as ASHRAE, and then Ecotect and Project Vasari were utilized to import this information, run simulations and visualize the results on Google Maps directly. Users are able to input standard information about any selected area provided by ASHRAE anytime in this way.

TABLE IV. DISTRIBUTION OF OFFICE BUILDINGS BY METROPOLITAN CITIES IN KOREA (2012)

	Seoul (40%)	Busan (18%)	Daegu (14%)	Incheon (13%)	Gwangju (8%)	Daejeon (7%)
No. of Office	140,987	66,494	49,462	47,654	28,906	28,999

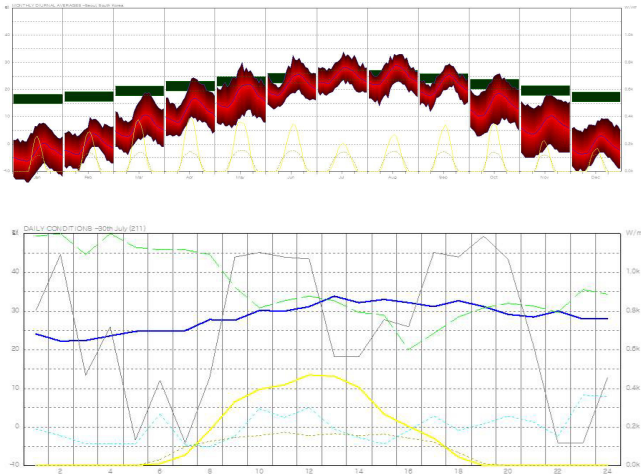


Figure 5. Monthly Climate Data for Seoul

Figure 5 shows detailed information about the Seoul climate illustrated by the weather tool in Ecotect as an example. It is easily noticed that temperature and isolation are the highest on July 30th and this data can be accepted as Summer data. The building modeling for the standard office is performed by the standard of LOD 100 and then, various types of the building façade including Hanok-styled openings have been applied one by one for the same analysis process repeatedly. Standard office settings used for simulation are shown on Table V, and they include operation hours, lighting equipment schedule, user activity level, and heat gain.

TABLE V. DATA TYPES FOR OFFICE BUILDINGS

Parameters	Default Values
Occupancy Schedule	Common Office 8am - 5pm
Lighting/Equipment Schedule	Office lighting 6am - 11pm
People/100 sq. M	3.5
People Activity Level	Standing, Light work, Walking
People Sensible Heat Gain(W/Person)	73
People Latent Heat Gain(W/person)	59
People Sensible Heat Gain(Btu/Person)	250
People Latent Heat Gain(Btu/Person)	200
Lighting Load Density(W/sq. ft.)	1.00
Equipment Load Density(W/sp. Ft.)	1.30
Electrical Equipment Radiant Percentage	0.3
Condition Type	Heated and Cooled
OA L/S Person	10
OA Flow per Area(Cu. M./hr/sq. M)	3.7
Unoccupied Cooling Set point	82

TABLE VI. BUILDING PERFORMANCE FACTORS

Location	Seoul, Korea
Weather Station	555181
Outdoor Temperature	Max:35°C/Min:-14°C
Floor Area	3,721m ²
Exterior Wall Area	1,488m ²
Average Lighting Power	10.87W/m ²
People	130
Exterior Window Ratio	0.24
Electrical Cost	\$0.06/kWh
Fuel Cost	\$1.21/Therm

D. Data Analysis

This study has suggested four types of grid patterns from the most typical Hanok-styled façade for simulation analyses in Figure 6. The first type has a general scale with a standard rate of proportion between vertical and horizontal ribs [10]. This means that a rate of thickness, width and angle of the crossed components are exactly the same as the traditional form. The second type has prominent vertical ribs instead of the original feature which is adjusted to fit in the typical curtain wall mostly used for current office buildings.

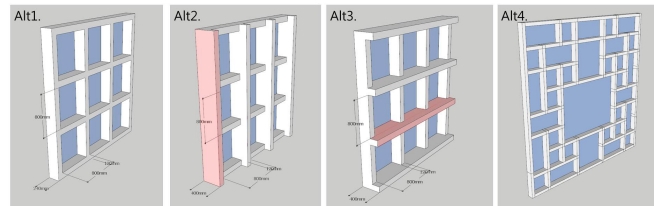


Figure 6. Operational Alternatives for Simulation

The next type has increased in depth for horizontal ribs that can play a role as building louvers. Various rib rates can be selected for both design values and energy efficiencies. The last type shows a combined configuration with both horizontal and vertical ribs that provides more possibilities for design variation. This type looks very similar to Korean Wan-ja or Sut-dea window patterns as shown in Table III, and it was assumed that the most ideal composition can be found with the analysis procedure suggested in this study.

Table VII summarizes the simulation results generated by GBS and Radiance applications for the above four types. It is found that all cases applied with Hanok opening patterns on the building façade have relatively better scores than the controlled group with a façade type of the general office building. Especially, the annual carbon emissions showed the biggest difference among these, and Hanok-styled façade types seemed relatively helpful for increasing energy efficiency in the aspect of cost. This result emphasizes that the use of Hanok-styled façade may possibly have

advantages for both the design creation and the building performance.

E. Proposed Performance Control with ICT

With the above investigation, it turns out that applying Hanok window patterns can increase the effectiveness for indoor environmental quality and energy savings. Furthermore, proposed patterns will possibly be more adaptable to the commercial building façade, if building users can control parameters for the grid patterns conveniently depending on their design preference or any given specific environmental situation.

ICT can play a important role in this sense, and WSSN can be proposed and utilized for adjusting grid intervals and/or component scales of the façade. In addition, WSSN is applicable with mobile devices like smart phone used by most people in daily life. For example, sensors installed onto the building receive signals from smart phones, and then kinetic façade components are operable for controlling their

intervals and scales by the preset scenarios with various contexts for optimizing building performance.

WSSN system makes people possible to activate their smart phone applications or computer programs to adjust daylight, and then Wi-Fi signal via TCP/IP arrives to router stack on each window panel as shown on Figure 7. Its transparency and color can be changed through a series of functions by Window Property Controller. Eventually, users can adjust their indoor environment by using WSSN.

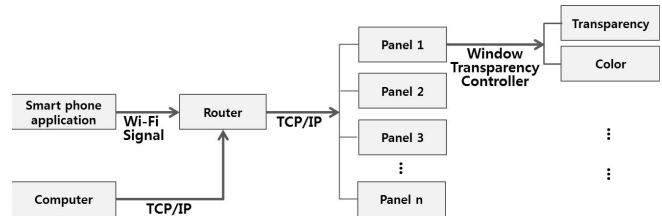


Figure 7. Distribution Diagram of Wireless Smart Sensor Network

TABLE VII. SIMULATION RESULTS

		Control Group	Alt1	Alt2	Alt3	Alt4
Annual Carbon Emissions		129	112	112	109	116
Annual Energy	Use (Electricity: kWh /Fuel: MJ)	555.615 / 597.501	505.753 / 551.923	502.605 / 552.556	498.255 / 541.729	517.350 / 657.167
	Cost (\$)	40,103	36,699	36,638	36,149	37,429
Energy use	Fuel (MJ)	697.501	551.822	662.555	641.728	657.167
	Electricity (kWh)	551.162	501.300	498.152	493.813	512.898
Monthly Heating / Cooling Load						
Radiation	Direct Radiation Average	675.503	34.5808	34.892	1.063	174.382
	Diffuse Radiation Average	656.262	361.973	322.499	304.592	413.698
	Total Radiation Average	1331.766	396.553	357.391	305.656	588.08
	Simulation Figure					

This research also suggests Polymer Dispersed Liquid Crystal (PDLC) for increasing the indoor environmental performance of Hanok-styled façade. This concept was originally invented by Dr. Edwin Land and has been commercialized since 2003 [11]. The proposed façade system in this paper is composed of the Hanok-styled frames and the responsive window glasses that are usually on the opaque state with only 5% daylight penetration and also convertible to the transparent status passing 70% lights from the outside. These variant situations are made possible by controlling quantity and intensity of the electricity flown on the glasses and signals sent from WSSN based on Wi-Fi and/or Zigbee communication module determine the situation. We are quiet certain that usability of the Hanok-styled façade would be doubled, if ICT such as WSSN and PDLC examined in this paper is successfully integrated.

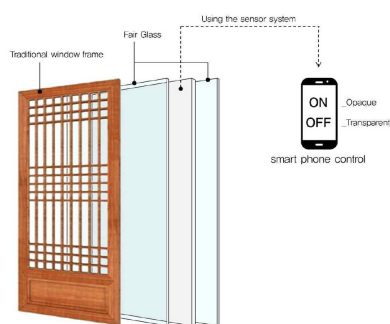


Figure 8. Proposed Performance Control for Hanok-styled Façade

IV. CONCLUSION

The purpose of this study is to investigate characteristics of Hanok opening patterns as the building façade and to prove their energy efficiencies and indoor environmental qualities as important parts of the building performance. With utilizing scientific simulations for various alternatives of the building skins mentioned previously, it was found that Hanok opening patterns may have great advantages to be a type of practical building façade presented worldwide.

For this study, analyses for building performance have been done for selected patterns that may minimize heat loss and make it easy to adjust indoor light environment by changing the intervals. ICT also played an important role as both a controller for the spatial environment and a simulator for evaluating various experimental factors; WSSN has been utilized for adjusting configuration of the opening patterns, and popular simulation tools like Project Vasari and Green Building Studio have been used for verification of the feasibility. The results from the simulations show that Hanok opening patterns can reduce the annual carbon emission from the energy consumption approximately up to 15% in comparison to the generic office buildings, and has the possibility to be used as an alternative for resolving problems related to air pollution. Finally, the building enclosures satisfying both Korean traditional identity and efficiency of the building performance have been suggested as the Hanok-styled façade.

This advanced concept of Hanok-styled façade can be expanded to residential and commercial facilities including office buildings, because it has been proved by simulations that overall energy use and its lifecycle cost are significantly reduced. Future study will include more building components using Hanok-styled elements that make architectural design and fabrication possible to diversify and environmental qualities much higher.

ACKNOWLEDGMENT

This research was supported by a grant (13AUDP-B070244-01) from Urban Architecture Research Program (Development of Hanok Technology, Phase II) funded by Ministry of Land and Transport Affairs of Korean government.

REFERENCES

- [1] Han, S., Im, O., Lee, M., and Cheon, D., A Study on the Establishment of an Evaluation System for Integrative Comfort Performance of Hanok Residence, *Journal of the Korean Housing Association*, vol. 24, no.3, June, 2013, pp. 27-35.
- [2] Lee, T., Kim, H., Song, G., and Kim, S., An Experimental Study on the Sound Insulation and Sound Absorption Characteristics of Korean Traditional Windows, *Journal of the Korea Institute of Ecological Architecture and Environment*, vol. 8, no. 5, October, 2008, pp. 3-10
- [3] Ahn, E. and Kim J., Computational Analysis of Air Flows Inside Korean Traditional House, *Journal of Korea Multimedia Institution*, vol. 15, no. 3, March, 2012, pp. 380-397.
- [4] Lee, K., Kim, I., and Choo, S., A Study on Improvement of Energy Performance Index in Green Building Certification System using BIM, *Journal of the Architectural Institute of Korea*, vol. 27, no. 9, September, 2011, pp. 13-21.
- [5] Joo, N., *Korea Architecture Huge Plan, Window*, Boseonggak, 2007.
- [6] Im, O., A Study on the Kinetic Façade System and the Energy Performance Evaluation of KLSU System, Master's Thesis, Chonnam National University, 2014.
- [7] Lee, J., The Daylight Simulation According to the Arrangement of Apartment Housing in Cheong-ju City, Master's Thesis, Chungbuk National University, 2011.
- [8] Choo, S., Lee, K., and Park, S., A Study on LOD (Level of Development) for Development of Green BIM Guidelines - Focused on Energy Performance Estimation, *Journal of the Architectural Institute of Korea*, vol. 28, no. 6, June, 2012, pp. 37-47.
- [9] Im, O., Kim, K., and Han, S., A Kinetic Light Shelf Unit as an Integrated Intelligent Control Device for Optimizing Interior Illumination, *Proceeding of the Eighth International Multi-Conference on Computing in the Global Information Technology (ICCGI 2013)*, Nice, Paris, August, 2013, pp. 295-298.
- [10] Joo, N., *Korean Gate, Door and Windows*, Daewonsa, 2001.
- [11] Kim, S., Chang, Y., Kang, J., and Han, S., Implementation of the Wireless Smart Sensor Network for Spatial Comfort Performance of Hanok Residence, *Proceeding of the 7th International Conference on Anti-counterfeiting, Security, and Identification (ASID 2013)*, Shanghai, China, October, 2013, pp. 44-47.

Performance Improvement in Applying Network Coding to On-demand Scheduling Algorithms for Broadcasts in Wireless Networks

G. G. Md. Nawaz Ali ^{*}, Yuxuan Meng ^{*}, Victor C. S. Lee ^{*}, Kai Liu [†] and Edward Chan ^{*}

^{*} Department of Computer Science

City University of Hong Kong, HKSAR, China

Email: gnawazali2-c@my.cityu.edu.hk, ymeng7@student.cityu.edu.hk, csvlee@cityu.edu.hk, cshedchan@cityu.edu.hk

[†] College of Computer Science, Chongqing University, Chongqing, China

Email: liukai0807@cqu.edu.cn

Abstract—Due to its ability to satisfy multiple data items through a single broadcast, on-demand broadcasting has enjoyed wide usage in wireless data dissemination. Recently, applying network coding in broadcast has received much interest, because using this technique a number of items can be served through a single broadcast, which further maximizes the channel bandwidth utilization and improves the overall system performance. A generalized network coding based encoding model has been proposed which helps to migrate a number of existing scheduling algorithms to the network coding assisted version while preserving their original scheduling criteria. However, the proposed system only studies the homogeneous system environment. In this work we have done an extensive simulation based on the proposed generalized encoding model, both in the homogeneous and heterogeneous environment, to analyze the efficiency and adaptability of network coding assisted scheduling algorithms against a number of performance metrics in the real-time environment. Simulation studies reveal some interesting results.

Keywords—Network coding; on-demand scheduling; wireless broadcast

I. INTRODUCTION

In wireless and mobile networks, data broadcasting is an efficient means for data dissemination. The server disseminates data items on the shared broadcast channel, and the clients can be served simultaneously by listening to this channel. In general, there are three broadcasting categories: push-based, pull-based, and hybrid broadcast [1][2]. Pull-based broadcast is also referred to as on-demand broadcast. In on-demand broadcast, which is the focus of this research, the server broadcasts according to the requests sent by clients, by compiling requests in queue and broadcasting requested data items based on the various attributes of pending data items at the server. Clients listen to the broadcast channel, and receive the data items that they need.

Some on-demand data scheduling algorithms have been proposed by researchers [1][3]-[6], which have different application-specific performance objectives. Ahlswede et al. [7] proposed network coding that further improves performance. It utilizes previously received data items cached at clients. Clients requesting different data items can be satisfied simultaneously. The impact of network coding is determined by the encoding decision of the server at each broadcast tick. This paper is based on a generalized encoding framework

proposed and described in [8] which incorporates a flexible and adaptive network coding into data scheduling algorithms for on-demand broadcast. Our contribution is in providing extensive simulation experiments to study the performance of scheduling algorithms with and without network coding.

This paper is structured as follows. Section II covers related work in this area. Section III covers the system model used in the paper, followed by coverage of extensive simulation experiments in Section IV. The paper concludes in Section V.

II. RELATED WORKS

A number of data scheduling algorithms have been proposed for on-demand broadcast, such as FCFS by Wong [1], Most Requests First (MRF) and Longest Wait First (LWF) in [2], and $R \times W$ by Aksoy et al. [3]. They focused on different metrics, such as reducing access time, deadline miss ratio, or stretch. The well-known EDF algorithm [9] is a scheduling algorithm in real time systems. Xuan et al. [5] proved its good performance in on-demand broadcast. Xu et al. [6] proposed “Slack time Inverse Number of pending requests” (SIN), which can achieve good performance in the real-time environment.

Network coding encodes different data items, broadcasts the encoded data in a broadcast tick, and improves performance. The simplest coding schemes, linear coding is studied by Li et al. [10] which examined the network capacity of multi-cast networks. Park et al. [11] showed that network coding can achieve even 65% higher throughput than conventional multi-cast in a typical ad hoc network scenario. Fragouli et al. [12] formulated an analytical model for energy-efficient broadcast in wireless ad hoc networks with coding .

Most works [1][3][4][6] assumed that through one broadcast, only clients requesting the same data item can be satisfied, where the broadcast bandwidth cannot be fully utilized. Chu et al. [13] proposed a multi-data delivery algorithm using the XOR operator to encode the data. However, in their encoded data item at most two items can be encoded. Recently, Zhan et al. [8] proposed a generalized coding framework which can be used to determine the set of request items which can be broadcast to the maximum number of clients, where an encoded items can use maximum coding opportunity.

In this paper, based on the generalized framework in [8], we conduct extensive simulation experiments of a number of

existing scheduling algorithms with coding and without coding to explore the performance in both same and heterogeneous item size to analyze the coding adaptivity efficiency.

III. SYSTEM MODEL

In this section we describe the system model as well as the encoding framework used in this research.

A. System Architecture

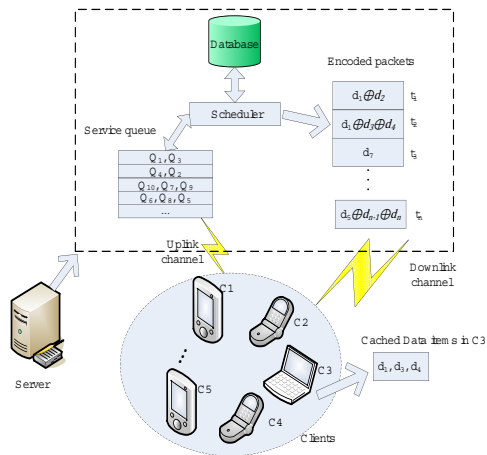


Fig. 1. System Architecture

Our system architecture is the typical on-demand system model for data broadcasting in wireless environment [3] as shown in Figure 1. The system consists of one server and many clients. When a client requires a item and cannot locate it in its local cache, it will send a request through the uplink channel to the server. After sending the request, a client listens to the broadcast channel. A client can decode the requested data item from an encoded packet when it has all the other encoded data items in the packet in its cache.

The server has a received queue where it stores the generated requests received from clients. A request is feasible if it has enough slack time to be served. An infeasible request is one whose deadline has missed and will be removed from the received queue. After invoking a certain scheduling algorithm the server retrieves the data items from the local database and according to the clients' cache, it forms the encoded packet for broadcasting through the broadcast channel in the next serving cycle. We use simple XOR operations for encoding and decoding due to its lower overhead functionality [8][14]. The primary goal of real-time scheduling is to serve as many requests as possible before their deadlines and to maximize the broadcast channel bandwidth utilization.

B. Graph Model

A graph model by Zhan et al. [8] can be constructed as follows. The system has a data server s and n clients, $R = \{c_1, c_2, \dots, c_n\}$. Let $W(c_i)$ be the set of data items requested by client c_i , and $H(c_i)$ be the set of data items cached at client c_i . A database containing data items is in the server, where d_j

TABLE I
THE UNIFIED ENCODING MODEL

Scheduler	Input
FCFS	$\phi(t_i, l_i, T_i) = 2^{N-j+1} t_i$, t_i is the j -th largest in $\{t_1, \dots, t_N\}$, $\chi(\mathbf{t}, \mathbf{l}, \mathbf{T}, \mathbf{x}) = 1$
MRF	$\phi(t_i, l_i, T_i) = 1$, $\chi(\mathbf{t}, \mathbf{l}, \mathbf{T}, \mathbf{x}) = 1$
LWF	$\phi(t_i, l_i, T_i) = t_i$, $\chi(\mathbf{t}, \mathbf{l}, \mathbf{T}, \mathbf{x}) = 1$
$R \times W$	$\phi(t_i, l_i, T_i) = 1$, $\chi(\mathbf{t}, \mathbf{l}, \mathbf{T}, \mathbf{x}) = \max_{1 \leq i \leq N} \{x_i \times t_i\}$
EDF	$\phi(t_i, l_i, T_i) = 2^{N-j+1} \frac{1}{T_i}$, $\frac{1}{T_i}$ is the j -th largest in $\{\frac{1}{T_1}, \dots, \frac{1}{T_N}\}$, $\chi(\mathbf{t}, \mathbf{l}, \mathbf{T}, \mathbf{x}) = 1$
SIN	$\phi(t_i, l_i, T_i) = 1$, $\chi(\mathbf{t}, \mathbf{l}, \mathbf{T}, \mathbf{x}) = \max_{1 \leq i \leq N} \{x_i \times \frac{1}{T_i}\}$

is the j -th data item, and m is the total number of data items in the database.

Definition 3.1: Given $R = \{c_1, c_2, \dots, c_n\}$, $D = \{d_1, d_2, \dots, d_m\}$, $W(c_i) \subseteq D$, $H(c_i) \subseteq D$, $W(c_i) \cap H(c_i) = \emptyset$, we construct a graph $G(V, E)$ as:

$V = \{v_{ij} | \text{client } c_i \text{ requests for item } d_j, 1 \leq i \leq n, 1 \leq j \leq m\}$
 $E = \{(v_{i_1 j_1}, v_{i_2 j_2}) | j_1 = j_2; \text{ or } j_1 \neq j_2, d_{j_2} \in H(c_{i_1}), d_{j_1} \in H(c_{i_2})\}$

Accordingly, there are two rules to construct the graph $G(V, E)$, by connecting an edge between two vertices in the graph:

the first rule is, $(v_{i_1 j_1}, v_{i_2 j_2})$ with $j_1 = j_2$, which means that if client c_{i_1} and client c_{i_2} require the same data item, there is a link between two vertices $v_{i_1 j_1}$ and $v_{i_2 j_2}$.

The second rule is, $(v_{i_1 j_1}, v_{i_2 j_2})$ with $j_1 \neq j_2$, $d_{j_2} \in H(c_{i_1})$, and $d_{j_1} \in H(c_{i_2})$, which means that if client c_{i_1} 's cache contains the data item being requested by client c_{i_2} and vice versa, there is a link between vertices $v_{i_1 j_1}$ and $v_{i_2 j_2}$.

A clique is a subset of the vertices in an undirected graph, such that every two vertices in the subset are connected by an edge in graph theory. To apply the network coding to broadcast, needs to find the maximum clique δ_{max} of the graph $G(V, E)$. By broadcasting the set of requested data items in δ_{max} with coding, the maximum number of clients is served in each broadcast tick.

C. Encoding Framework

Besides satisfying a number of requests by each broadcast unit, scheduling also should pay attention to some other applications specific requirements such as the longest waiting time or current stretch, minimal slack time, etc. A generalized encoding framework is proposed in [8] as follows:

The vertices in the graph are noted as v_1, v_2, \dots, v_N , where N is number of vertices and v_i corresponds to a request q_i . Three weights, t_i, l_i and T_i , are associated with vertex v_i , where they are respectively the waiting time, the size of data item and the slack time of the request q_i corresponding to v_i .

$x_i \in \{0, 1\}$ is to denote whether vertex v_i is selected or not in the optimal clique C , whereas $x_i = 1$ means v_i is selected. Here $\mathbf{T} = \{T_1, T_2, \dots, T_N\}$, $\mathbf{x} = \{x_1, x_2, \dots, x_N\}$, $\mathbf{t} = \{t_1, t_2, \dots, t_N\}$, $\mathbf{l} = \{l_1, l_2, \dots, l_N\}$. $\phi(t_i, l_i, T_i)$ is the weight function of attribute(s) associated with v_i and $\chi(\mathbf{t}, \mathbf{l}, \mathbf{T}, \mathbf{x})$ is the optimized function in different applications.

TABLE II
 SIMULATION PARAMETERS

Parameter	Default	Range	Description
<i>ClientNum</i>	300	—	Number of clients
λ	20	—	Mean request generation rate
<i>DBSIZE</i>	1000	—	Size of the database
<i>SIZEMIN, SIZEMAX</i>	1, 10	—	Min. and Max. data item size
<i>CacheSize</i>	60	30-180	Client cache size
θ	0.6	0.0-1.0	Zipf distribution parameter
μ^-, μ^+	120, 200	—	Min. and Max. laxity

maximize

$$\psi(C) = \chi(\mathbf{t}, \mathbf{l}, \mathbf{T}, \mathbf{x}) \sum_{i=1}^N \phi(t_i, l_i, T_i) x_i \quad (1)$$

subject to

$$(x_i + x_j) \leq 1, (v_i, v_j) \notin E(G), 1 \leq i \neq j \leq N, x_i = \{0, 1\}$$

The weight function $\phi(t_i, l_i, T_i)$ and the optimized function $\chi(\mathbf{t}, \mathbf{l}, \mathbf{T}, \mathbf{x})$ are used with different settings to meet various application requirements.

The strategies for various scheduling algorithms with network coding in a unified model described in [8] are summarized in Table I, where each of the scheduling disciplines tries to maximize the functions $\psi(C)$ in (1).

IV. PERFORMANCE EVALUATION

A. Simulation Setup

Our simulation model is implemented in CSIM19 using the default parameters shown in Table II which are based on [8]. A closed system model is used. The item access pattern follows the zipf distribution, where skewness is controlled by the parameter θ . For a real-time environment, we set the relative deadline (RD) of a request R_i as follows:

$$RD_i = (1 + \text{uniform}(\mu^-, \mu^+)) * T_i^{\text{serv}} \quad (2)$$

where μ^- and μ^+ are the minimum and maximum laxity for calculating the relative deadline, respectively, and T_i^{serv} is the service time of R_i . The deadline of a request R_i is Dl_i , which is computed by:

$$Dl_i = AT_i + RD_i \quad (3)$$

where AT_i is the arrival time of R_i .

The channel bandwidth is 1.0. For the homogeneous environment each item size is 1.0, but for the heterogeneous environment different item sizes are generated using the random item size distribution (RAND) [15].

B. Performance Metrics

- 1) Deadline Miss Ratio (DMR): The ratio of the number of requests which missed their deadline over all submitted requests..
- 2) Average Encode Length (AEL): The average number of data items encoded in each encoded packet. A high AEL means more clients can decode their expected data items from an encoded broadcast packet.

3) Saved Bandwidth Ratio (SBR):

$$SBR = \frac{\#Satisfied\ Item - \#Broadcast\ Item}{\#Satisfied\ Item} \quad (4)$$

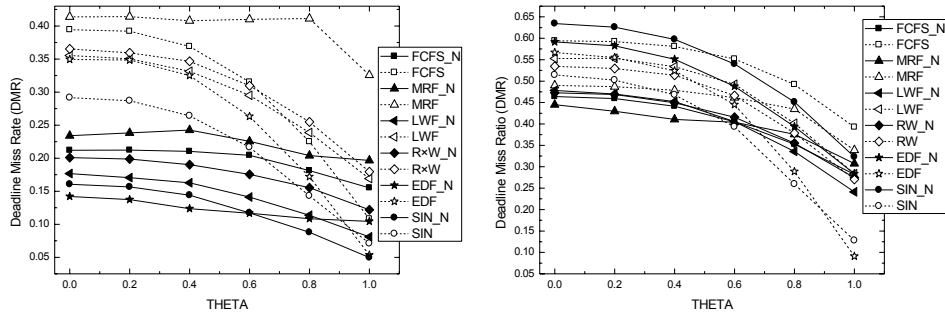
where *Satisfied Item Number* means the number of data items satisfied and *Broadcast Item Number* the number of encoded items being broadcast.

C. Performance Analysis

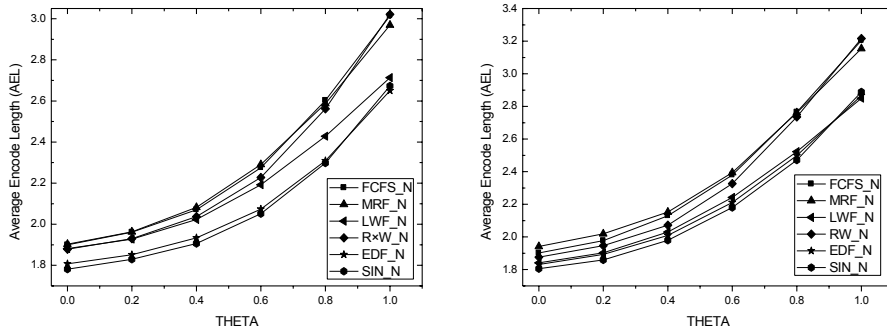
In this section, we discuss the performance analysis of coding and non-coding versions of different existing scheduling algorithms under both the same and different item size distribution. Simulation experiments continue until 98% confidence interval has been achieved. The different scheduling algorithms can be divided into two types: algorithms such as SIN and EDF which target time constraints or urgency (TC) and those that do not focus on time constraints (NTC).

1) *Impact of Skewness Parameter θ* : Figure 2 shows the impact of data item access pattern θ . When $\theta = 0.0$, item access pattern is uniformly distributed. But with increasing θ , the probability for accessing popular data items increases. In serving such a popular item, a number of requests can be served concurrently, which explains why the performance of all algorithms improves with increasing θ . From Figure 2(a)(I), for the same item size distribution, network coding assisted TC algorithm SIN_N and EDF_N show better performance than NTC network coding assisted algorithms MRF_N, FCFS_N, $R \times W_N$ and LWF_N. To see the impact of considering urgency as a request selection, consider Figure 3(I). This figure shows the percentage of served requests over total submitted requests against the slack time (ST) of requests in the system under default settings. The percentage value in the y-axis against the value '0-1/4 of RD' in the x-axis means percentage number of requests served when these requests have remaining deadlines i.e., slack times (ST) less than 1/4 of the original assigned relative deadlines (RD) of the respective requests. Similarly, there are total four categories. To understand this metric, let us have an example. Suppose a request R_i 's relative deadline RD_i is 12. So, initially R_i 's slack time ST_i is also 12. As time passes ST_i also decreases. Now, if R_i is served when $ST_i = 10$, it will fall in the 4th category (that is 3/4 - RD = 9 - 12). Similarly, if R_i is served when $ST = 7$ or $ST = 1$ it will fall in the 3rd and 1st category respectively. Accordingly, if more number requests served in the 1st and subsequent categories, the algorithm is more aware of the scheduling of deadline urgent requests. From Figure 3(I) we can see that, in SIN_N more requests are served in the first category when compared to MRF_N, in which more requests are served in the last category. It proves that MRF_N does consider the slack time of the requests, which is why fewer requests are served in the first category, hence urgent requests may not meet their deadlines resulting in higher DMR in MRF_N.

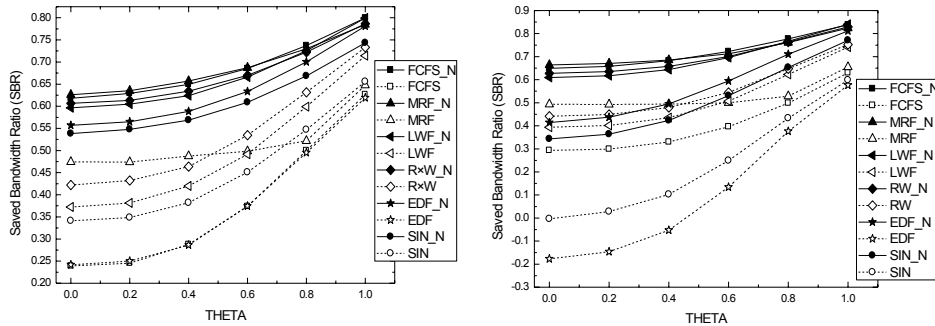
On the contrary, from Figure 2(a)(II), surprisingly TC algorithms SIN_N and EDF_N can not retain their supremacy in terms of DMR for RAND item size distribution. Recalling



(a) DMR for (I) Same, and (II) RAND item size distribution



(b) AEL for (I) Same, and (II) RAND item size distribution



(c) SBR for (I) Same, and (II) RAND item size distribution

Fig. 2. Impact of skewness parameter (θ).

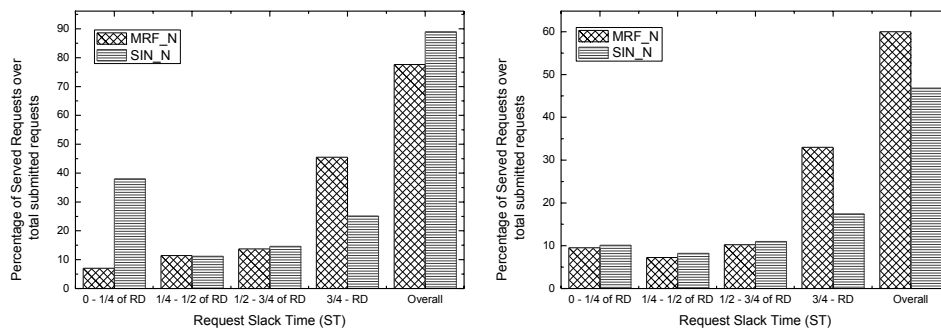
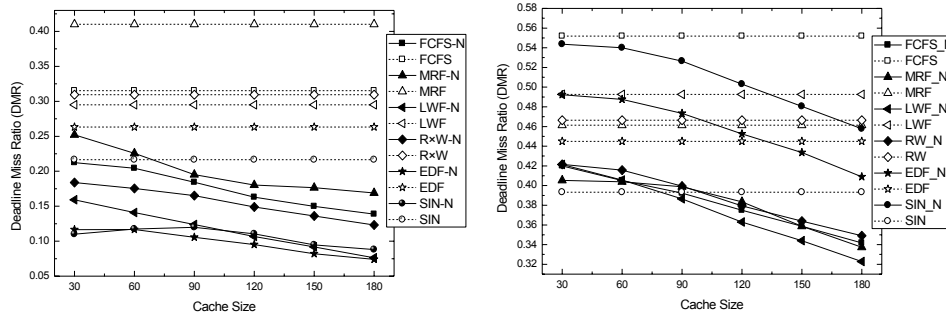
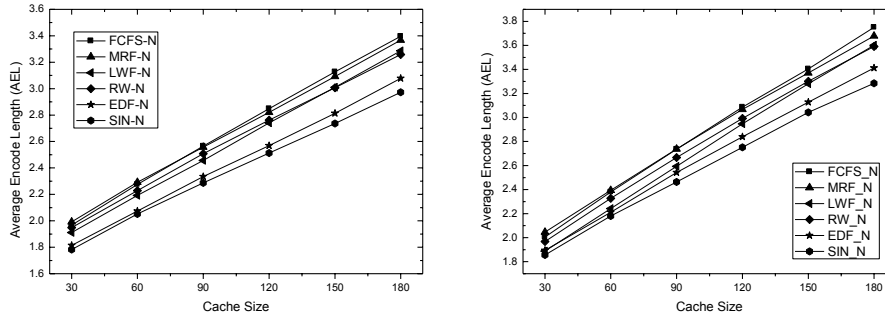


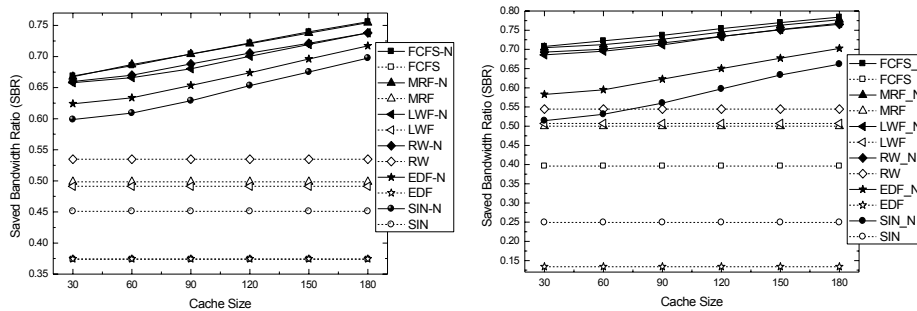
Fig. 3. Distribution of serving requests over total submitted requests for (I) Same, and (II) RAND item size distribution.



(a) DMR for (I) Same, and (II) RAND item size distribution



(b) AEL for (I) Same, and (II) RAND item size distribution



(c) SBR for (I) Same, and (II) RAND item size distribution

Fig. 4. Impact of Cache Size.

from the network coding view, for different item sizes, the size of the encoded packet will be the maximum item size in the selected clique. Now, a TC algorithm uses slack time for selecting the optimal clique. Unlike productivity based algorithm, it may select a clique having smaller clique size with urgent request as a clique vertex. Regardless of the selected clique size, the service time of the encoded packet will be as large as the service time of the largest data item in the selected clique, hence the service time is increased, which makes it more difficult to meet the deadline. This statement is supported by Figure 3(II), which shows that unlike same item size distribution, SIN_N does not outperform MRF_N in the first category (0 - 1/4 of RD), and overall requests serving percentage also lower than MRF_N. Nevertheless, the non network coding version of the algorithms do not have this clique size problem, and hence show superior performance

for skewed access pattern (Figure 2(a)(II)). On the other hand, the coding version of a NTC based algorithm typically tries to serve a clique having maximum clique size, hence although the service time is increased, many requests can be served concurrently. Therefore, although such an algorithm in heterogeneous item size environment has worse performance than the same item size environment, it still has a better performance than its corresponding non-coding version as well as the network coding version.

Except for FCFS_N, all the NTC algorithms use item productivity as their request selection criterion, i.e. use size of a clique ((in other words, number of vertices in the clique)), but ignore the urgency of a request. However, FCFS_N selects the optimal clique, including the weight of a corresponding vertex, waiting time of all the vertices in a clique is counted, which indirectly supports the clique size. For example, it is

more likely that a clique having larger size, might have more total weight than that of a smaller size clique. So, the network coding version of FCFS, i.e. FCSF_N, has similar performance to other WTC algorithms. It is evident from Figure 2(a)(I), network coding version of an algorithm has better performance than its non-network coding version. An exception is EDF, when $\theta = 1.0$; EDF_N shows better DMR than EDF. This is because in the skewed item access pattern the non-network coding version of an algorithm can gain more than network coding environment.

From Figures 2(b)(I) and (II), with increasing θ , average encode length (AEL) of all the algorithms increases, because with increased access pattern more requests ask for the popular data items, hence being a requested item of a client could be a cached item of another client with a higher probability, which helps to form a bigger clique and increase the coding opportunity as well. NTC algorithms have better AEL and saved bandwidth ratio (SBR) than TC algorithms irrespective of item size distribution as shown in Figures 2(b)(I), and 2(b)(II), 2(c)(I) and 2(c)(II). Regarding SBR, regardless of item size distribution, coding version of an algorithm has better SBR than its non-coding version. This is because, as the client gradually accumulates more items in the cache, the client can exploit its cache in the coding version to satisfy more requests. This is clearly not possible in the non-coding version of the algorithms.

2) *Impact of Cache Size:* Figure 4 shows the impact of the cache size of a client for both same and different item size distribution. Since the non-coding version of an algorithm does not consider client's cache, changes in cache size has no impact on performance. Increasing cache of a client provides more room to store more items, which increases the probability of one client's cached item to be another client's requested item. This is the key to increase the coding flexibility and AEL which provides more opportunity for performance gain. For this reason DMR declines with increasing cache size for both same and rand item size distributions (Figures 4(a)(I), 4(a)(II)), AEL increases (Figures 4(b)(I), 4(b)(II)), and SBR increases (4(c)(I), 4(c)(II)).

V. CONCLUSION AND DISCUSSION

In this paper, based on the generalized encoding model proposed in [8], we conducted extensive simulations for same and different item size distributions in the real-time environment. We analyze the simulation performance of a number of existing scheduling algorithms both in their coding and non-coding versions against a number of performance metrics to examine the efficiency and adaptability of coding assisted scheduling algorithms. Simulation results show that network coding assisted algorithms improve their performance with increased cache size and skewed access pattern. Generally speaking, NTC algorithms have superior performance in AEL and SBR than TC algorithms irrespective of data item size distribution. In addition, regardless of item size distribution, all the coding assisted algorithms outperform their corresponding non-coding version in SBR. On the other hand TC based

algorithms perform better in DMR in the real-time setting for the same item size distribution. However, surprisingly, in the different item size distribution their performance decline dramatically. A plausible reason is that in the coding version, the scheduling algorithm needs to serve the maximum sized item of the selected clique irrespective of the clique size in each broadcast, hence a clique may have urgent vertex but also may have smaller clique size. On the contrary, the non-coding version simply selects the best item based on the underlying scheduling principle.

We are exploring the use of other metrics in deepening our understanding of the behavior of different algorithms when network coding is used, and will hopefully present the results in a future paper.

REFERENCES

- [1] J. W. Wong and M. H. Ammar, "Analysis of broadcast delivery in videotex system," *Journal of IEEE Transactions on Computers*, vol. 34, no. 9, September 1985, pp. 863–866.
- [2] J. W. Wong, "Broadcast delivery," *Proceedings of the IEEE*, vol. 76, no. 12, 1988, pp. 1566–1577.
- [3] D. Aksoy and M. Franklin, " $R \times W$: A scheduling approach for large-scale on-demand data broadcast," *Journal of IEEE/ACM Transactions on Networking (TON)*, vol. 7, no. 6, December 1999, pp. 846–860.
- [4] S. Acharya and S. Muthukrishnan, "Scheduling on-demand broadcasts: new metrics and algorithms," in *Proceedings of the 4th annual ACM/IEEE international conference on Mobile computing and networking (MobiCom'98)*, 1998, pp. 43–54.
- [5] P. Xuan, S. Sen, O. Gonzalez, J. Fernandez, and K. Ramamritham, "Efficient and timely dissemination of data in mobile environments," in *Proceedings of the 3rd IEEE Real Time Technology and Applications Symposium (RTAS'97)*, Montreal, Canada, 1997.
- [6] J. Xu, X. Tang, and W. Lee, "Time-critical on-demand data broadcast algorithms, analysis and performance evaluation," *IEEE Transactions on Parallel and Distributed Systems*, vol. 17, no. 1, 2006, pp. 3–14.
- [7] R. Ahlswede, N. Cai, S.-Y. Li, and R. Yeung, "Network information flow," *IEEE Transactions on Information Theory*, vol. 46, no. 4, July 2000, pp. 1204–1216.
- [8] C. Zhan, V. Lee, J. Wang, and Y. Xu, "Coding-based data broadcast scheduling in on-demand broadcast," *IEEE Transactions on Wireless Communications*, vol. 10, no. 11, November 2011, pp. 3774–3783.
- [9] C. Liu and J. Layland, "Scheduling algorithms for multiprogramming in hard real-time traffic environments," *Journal of the Association for Computing Machinery (ACM)*, vol. 20, no. 1, 1973, pp. 46–61.
- [10] S.-Y. Li, R. Yeung, and C. Ning, "Linear network coding," *IEEE Transactions on Information Theory*, vol. 49, no. 2, February 2003, pp. 371–381.
- [11] J.-S. Park, D. Lun, F. Soldo, M. Gerla, and M. Médard, "Performance of network coding in ad hoc networks," in *Proceedings of the IEEE Military Communications Conference*, Washington, DC, October 2006, pp. 1–6.
- [12] C. Fragouli, J. Widmer, and J.-Y. Le Boudec, "A network coding approach to energy efficient broadcasting: From theory to practice," in *Proceedings of the 25th IEEE International Conference on Computer Communications (INFOCOM'06)*, Barcelona, Spain, April 2006, pp. 1–11.
- [13] C.-H. Chu, D.-N. Yang, and M.-S. Chen, "Multi-data delivery based on network coding in on-demand broadcast," in *Proceedings of the 9th International Conference on Mobile Data Management (MDM'08)*, Beijing, April 2008, pp. 181–188.
- [14] J. Chen, V. C. S. Lee, and C. Zhan, "Efficient processing of real-time multi-item requests with network coding in on-demand broadcast environments," in *Proceedings of the 15th IEEE International Conference on Embedded and Real-Time Computing Systems and Applications (RTCSA '09)*, Beijing, August 2009, pp. 119 – 128.
- [15] J. Xu, Q. Hu, W. Lee, and D. L. Lee, "Performance evaluation of an optimal cache replacement policy for wireless data dissemination," *IEEE Transactions on Knowledge and Data Engineering*, vol. 16, no. 1, January 2004, pp. 125–139.

Combination of IMS-based IPTV Services with WebRTC

Tilman Bach
Hochschule fuer Telekommunikation
Leipzig (HfTL)
University of Applied Sciences,
Leipzig, Germany
Email: tilmann.bach@hftl.de

Michael Maruschke
Hochschule fuer Telekommunikation
Leipzig (HfTL)
University of Applied Sciences,
Leipzig, Germany
Email: maruschke@hftl.de

Jens Zimmermann
Hochschule fuer Telekommunikation
Leipzig (HfTL)
University of Applied Sciences,
Leipzig, Germany
Email: jens.zimmermann@hftl.de

Kay Hänsge
Telekom Innovation Laboratories
Deutsche Telekom AG
Berlin, Germany
Email: kay.haensge@telekom.de

Matthias Baumgart
Telekom Innovation Laboratories
Deutsche Telekom AG
Berlin, Germany
Email: matthias.baumgart@telekom.de

Abstract—This paper describes the potential to combine IP Multimedia Subsystem (IMS)-based IP Television (IPTV) services with a future-oriented web browser technology, the Web Real-Time Communication (WebRTC). To enrich the quality of experience for residential customers using this upcoming technology, the article focuses on the merging of the technical capabilities arising from both the IMS-based IPTV services and the WebRTC clients. Advantages of these different technologies are introduced and involved in the authors' concept. The presented proposal reuses the standardized IMS-based IPTV architecture as well as a state of the art WebRTC browser client. The ongoing WebRTC standardization process is considered. A first prototype has been developed successfully involving the Open IMS Core testbed, several IPTV typical components and the WebRTC client.

Keywords—WebRTC; IMS-based IPTV; CoD; Telco

I. INTRODUCTION

Nowadays, Telecommunication Carriers (Telcos) provide IP-based voice and video telephony services to their customers. In order to deploy an IP-based core network, an IMS-based infrastructure is increasingly used. With decreasing sales in the field of conventional legacy voice businesses, the Telcos are forced to develop new business areas. In order to offer new multimedia services like IPTV, it is reasonable to reuse the well established IMS-based core network infrastructure with Telco specific assets like guaranteed Quality of Service (QoS), service interoperability or mobility.

In the current state of the Internet, web technologies are steadily gaining popularity and browser-based real-time communication is an essential feature for future web applications. Using new web technologies like Hypertext Markup Language Version 5 (HTML5), WebSockets, etc., an easy and timely development of new web applications is possible. Furthermore, the new upcoming technology named WebRTC enriches ordinary web browsers with real-time communication functionalities [1]. With these new opportunities, it is possible to implement real-time communication applications within a web

browser. Thus, from the end-user's point of view, the installation of separate communication software or browser plug-ins is not necessary anymore. Instead, the desired communication features can be used in the browser immediately.

A browser with WebRTC features is also capable to deal with real-time streaming data which is used in Content on Demand (CoD) services like video or audio on demand. The combination of browser-based streaming capabilities and IMS-based IPTV services generates benefits for both the end-user and the Telcos. End-users have the ability to enjoy the advantages of IMS-based services like session mobility, QoS or Single Sign On mechanisms and Telcos can deliver their own applications and new features to the customers easily and directly by using web applications instead of legacy clients.

This paper discusses the combination of IMS-based IPTV services like CoD with WebRTC clients. We propose a concept reusing the IMS-based IPTV architecture to offer the CoD service to WebRTC end-users while only requiring a standard web browser. To verify this proposal, a first proof of concept has been implemented.

The present paper is structured as follows: Section II offers an overview to the current status of the standardization of the considered technologies, IPTV based on IMS core networks and WebRTC. Section III describes the authors' concept to combine WebRTC technologies with the IPTV service. The architecture and their specifics are considered and the proof of concept is presented. In Section IV, a conclusion and next steps in evolving this idea are pointed out.

II. STATUS QUO

A. IMS-based IPTV

The European Telecommunications Standards Institute (ETSI) Telecommunications and Internet converged Services and Protocols for Advanced Networking (TISPAN) working group has standardized a comprehensive IMS-based IPTV architecture [2][3]. Using the core IMS subsystem, various IPTV functions are supported.

The most common services and features described are Broadcast TV, Time Shifted TV, Content on Demand (CoD), Network-Personal Video Recorder (N-PVR), Pay-Per-View (PPV), Electronic Program Guide (EPG), parental control and advertising. Having finished this standardization process in 2011, ETSI TISPAN took note of the Telcos carrier grade network capabilities, such as high availability, QoS, mobility and even more. Most of the standardized IPTV architecture components make use of the common protocols utilized in an IMS-based network such as the Session Initiation Protocol (SIP), the Session Description Protocol (SDP), the Hypertext Transfer Protocol (HTTP) and the Real-Time Transport Protocol (RTP). Thus, it has big potential to combine the IMS-based IPTV services with established IMS real-time communication like voice and video telephony [4], presence [5] and other IMS services in order to create additional and more personalized value [6]. The standardized architecture of an IMS-based IPTV service is depicted in Figure 1 and is described as follows:

- Core IMS, core network components as specified in [7]
- Service Discovery Function (SDF), provides Service Attachment Information (SAI) with information about available services and related SSF
- Service Selection Function (SSF), provides Service Selection Information (SSI) containing the metadata of the available content.
- Service Control Function (SCF), is a SIP Application Server (AS) and the reference point for IMS UEs to start and control the IPTV sessions, moreover the SCF assigns the corresponding MCF and forwards the session information to it
- Media Control Function (MCF), controls media transport of MDF and receives instructions of SCF and UE
- Media Delivery Function (MDF), contains media data and transmits them to the UE
- User Equipment (UE), interacts as IMS-based IPTV end-user.

To realize IMS-based IPTV, ETSI [2] defines Generic IPTV Capabilities, a set of typically general signaling functions like service discovery and service control. All corresponding interfaces and used protocols are described more detailed in Section 4 of ETSI [3].

B. Web real-time communication with WebRTC

WebRTC [8] is an open project initiated by Google Incorporated. The purpose of the project is to integrate voice and video real-time communication into standard web browsers. From the first implementations, both web developers and telecommunication companies saw the potential for the future and prototyped new and promising applications. Telcos perceive the project as a risk for their classical voice business and thus they are also interested in expanding their own telecommunication service portfolio with the new WebRTC technology.

WebRTC provides an Application Programming Interface (API) definition which enables real-time communication in web browsers without the need of any additional browser plugin or additional software [1]. It empowers the browser to capture video and voice inputs of the client’s device. WebRTC is still in the standardization process. The World Wide Web Consortium

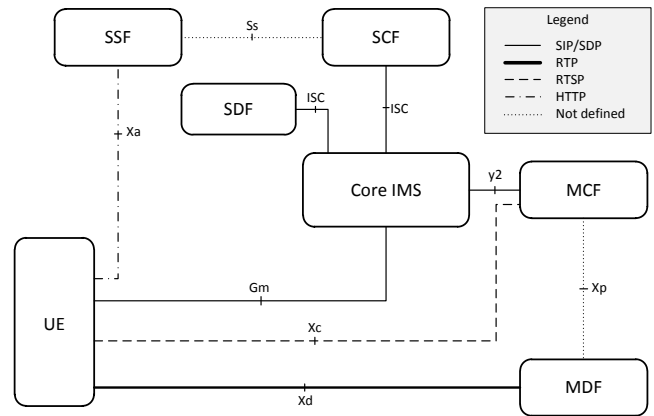


Figure 1: Simplified IMS-based IPTV functional architecture [3]

(W3C) is responsible for the web developer API and the Internet Engineering Task Force (IETF) for all corresponding protocols in an active working group named “Real-Time Communication in WEB-browsers - RTCweb” [9].

This browser extension enables developers to easily implement voice and video call web applications [10]. It also features components for file sharing. The browser implements the video, the voice and the transport engines. While there is still a discussion in the standardization process regarding the adequate audio and video codecs to be used current WebRTC implementations utilize the VP8 video compression format for video and the Opus codec for voice [11][12].

WebRTC does not define any particular signaling protocol. That is why developers can choose the most appropriate protocol for their special use case. So it is possible to implement new communication features, faster.

WebRTC requires secure transport of the RTP packets with the Secure Real-Time Transport Protocol (SRTP) [13] based on the mandatory to implement Datagram Transport Layer Security (DTLS) encryption protocol [14] used for key negotiation [15]. For solving Network Address Translation (NAT) problems, WebRTC also provides Session Traversal Utilities for NAT (STUN) [16], Traversal Using Relays around NAT (TURN) and Interactive Connectivity Establishment (ICE) [17] capabilities. WebRTC requires SDP for the negotiation of the session properties and uses the whole SDP’s Offer/Answer Model. Furthermore, SDP is also used for exchanging

- The fingerprint of the certificate used in the DTLS Certificate exchange procedure
- And ICE specific parameters like the ICE Candidate objects.

The generic architecture of a WebRTC client is described by Alvestrand [18] and illustrated in Figure 2. The components can be described as follows:

- Web server, provides the web application to load and includes a server for the client to connect to for handling the whole signaling flow
- Browser, a generic web browser
- Web application, application source code executed by the web browser

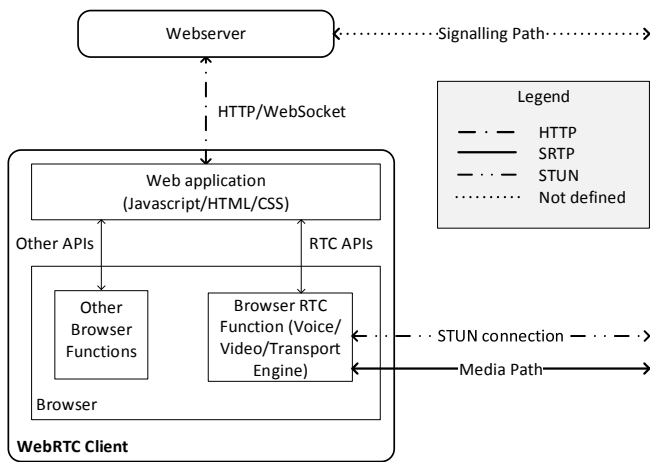


Figure 2: WebRTC client based on [18]

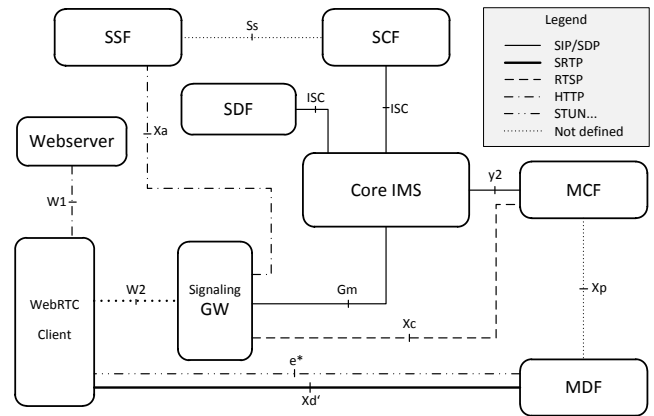


Figure 3: Architecture of proposed concept

- Browser RTC Function, WebRTC component in the web browser with voice, video and transport engines
- Signaling Path, is not specified but is needed to transfer the SDP information
- Media Path, transports the payload
- STUN Connection, is a mandatory component to bypass NAT restrictions.

For running a WebRTC client successfully, it is necessary to use a capable web browser. That means the browser has to implement the Browser RTC Function. Currently, web browsers like Google Chrome, Mozilla Firefox and Opera provide this component by default. Therefore, all devices which are able to run one of these browsers have the ability to use web based real-time communication. This includes all, desktop and tablet computers, laptops and smartphones. At the moment, there are restrictions in some operation systems like Apples iOS or Microsofts Windows Phone.

C. WebRTC access to IMS network based architecture

A 3rd Generation Partnership Project (3GPP) study suggests several solutions for accessing an IMS-based network architecture with WebRTC clients [19]. This clients could be connected either via wireless or wired access network technologies like Long Term Evolution (LTE), Wireless Local Area Network (WLAN) or any Digital Subscriber Line (xDSL). Overall, this study focusses on conversational real-time communication services like audio or video telephony. The technical report scopes the accessibility of typical IMS network characteristics or features like:

- Identity Management
- Accounting and Billing
- Interoperability with legacy networks like PSTN/ISDN
- Enabling of an application-oriented QoS (e.g., for voice telephony)

for Over The Top (OTT) applications such as plain WebRTC browser-to-browser applications if interconnected with a Telco network.

In contrast to [19], our paper is focused on IMS-based IPTV services, which should be made available to WebRTC Clients.

III.CONCEPT

A. Consolidate IMS-based IPTV with WebRTC

With the help of the Browser RTC Function (shown in Figure 2) the web browser is able to handle RTP packets without the need of any separate software or plug-in which was exposed in Section II. However, due to differences in the used media codecs and payload transport protocols the technical parameters of WebRTC and IMS-based IPTV do not match out-of-the-box. Assuming that the browser implementation of the WebRTC cannot be influenced by the Telcos, a modification of the architecture has to be enforced on the carrier side’s network. This section handles this idea and introduces a potential architecture of such a consolidation.

B. Architecture

The proposed architecture is based on the simplified IMS-based IPTV functional architecture. Instead of an IMS IPTV UE the endpoint of this service is a WebRTC client. For the combination of both, a translation for the different signaling and user data is necessary in various network components to gain compatibility. The consolidated architecture of the proposed concept is depicted in Figure 3. Components and interfaces, which are new or modified are listed as follows:

- Components:
 - Web server (new)
 - WebRTC client (new)
 - Signaling Gateway (SGW) (new)
 - MCF (modified)
 - MDF (modified)
- Interfaces:
 - W1 (new)
 - W2 (new)
 - e* (new)
 - Xp (modified)
 - Xd (modified to Xd’).

These changes and modifications made are described below. The web server is only needed for providing the WebRTC application sources which are fetched anew every time the end-users’ browser accesses the web application. The WebRTC application is executed in a WebRTC capable browser. The

application provides signaling functions for the communication with the core network via the inserted Signaling Gateway (SGW). Therefore, to make IMS-based IPTV services accessible to WebRTC clients, Generic IPTV Capabilities described in [2] are supported. The SGW implements the following generic capabilities:

- Service discovery and selection,
- Service control,
- Service interaction and
- Media control.

This gateway function converts session control messages coming from the WebRTC client side into SIP messages for the IMS core network side and vice versa. The SGW generates and forwards SIP messages towards the IMS core network and acts in place of the WebRTC client as a SIP capable signaling endpoint. As appears in the Figure 3 the SGW also converts the session control messages from the WebRTC client into HTTP and Real-Time Streaming Protocol (RTSP).

WebRTC strictly defines the media delivery, which requires a modification of the components of the standardized IMS-based IPTV architecture. These modifications specifically apply to the MCF and the MDF. According to [2], the MCF receives information about sessions from the SCF, finds and chooses the right MDF for media delivery and sends a response back to the user. The choice of the right MDF is based on codec information or geographical location. Information about the used protocol for media delivery is also part of the selection. Afterwards, the MCF transmits session information to the selected MDF. The important parts of information which need to be transmitted are the resource identifier of the media file to be streamed, the DTLS certificate fingerprint, the generated ICE candidates of the WebRTC client and specific connection information for establishing a transport channel.

Hence, the MDF can support the WebRTC's requirements, it shall support several new features. The concept of the MDF is depicted in Figure 4. One of the features is the MDF control engine. This engine dispatches the signaling from the MCF towards the internal MDF functions and vice versa. Another feature includes the audio and video codecs, which are also supported by the WebRTC client. Current WebRTC implementations prefer the VP8 video codec and the Opus audio codec. The codec handling is integrated in the streaming server. Further, WebRTC only allows a secure transport channel. That is why the MDF is required to support SRTP. To establish such a secure channel, the DTLS certificate fingerprint and a successful DTLS key exchange is needed [20]. SRTP keys are obtained through the DTLS key exchange [21]. The streaming server and the security capabilities of the media path are aggregated in the streaming engine. Another additional functionality of the MDF is an ICE agent. This agent supports the ICE methods regarding the SDP Offer/Answer negotiation and procedures for doing connectivity checks, which are similar to the functions documented in [19]. It also has to implement a STUN server functionality to support the STUN keep-alive usage as defined in [16]. This is used by the WebRTC client to preserve the NAT bindings.

With the help of these modifications, usage of a special gateway for media transcoding, which is described in [19], is not necessary.

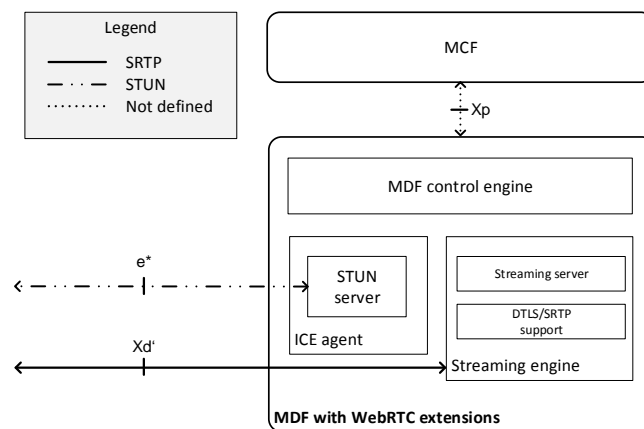


Figure 4: Detailed concept of the MDF

C. Interfaces

The interfaces shown in Figure 3 differ from the standardized architecture. The interfaces for the service interconnection are relocated from the user side to the SGW, which resides in the Telco's network infrastructure. The functionality and the used protocols of the Xa, Xc and Gm interfaces between the SGW and the service or the core functions are still conform to their specification [3]. The Ss, ISC and y2 interfaces remain unaffected despite the consolidation. In addition to these interfaces, some are modified or added and differ from the specification. These are described in more detail below.

The formerly undefined Xp interface between the MCF and the MDF is extended in the range of functions, respectively the extension of the MDF.

The W1 interface is a reference point between the WebRTC client and the web server. It is used to download HTML5, JavaScript, Cascading Style Sheets (CSS) and image files using HTTP. Via this interface, the user receives the latest WebRTC web application.

The W2 interface is located between the WebRTC client and the SGW. The used protocol for this reference point has been deliberately left open because the WebRTC does not define a signaling protocol. Therefore, the developer can choose one out of several state-of-the-art client-server protocols. Protocols like the WebSocket Protocol or the HTTP 2.0 specification could apply [22]. The meaning of protocol messages regarding this interface must cover the sense of the transferred protocol messages from the interfaces Xa, Xc and Gm.

The Xd' interface between the WebRTC client and the MDF is responsible for media delivery using SRTP [23]. The original Xd interface only supports RTP/RTCP or HTTP for media delivery, so the modified interface for WebRTC interconnection is called Xd'. This modification results from the mandatory use of a secure connection in WebRTC [13].

The e* interface is a second reference point between the WebRTC client and the MDF. This interface is added in the proposed concept. It is used for STUN connectivity checks between the both components to preserve the NAT bindings of the client. STUN is a mandatory to implement feature of WebRTC because the most WebRTC clients are behind NAT firewalls.

D. Proof of concept

To verify the functionality and the usability of the proposed concept, a testbed is prepared. With this implemented testbed, the content on demand use case (audio and video) is realized and tested. This includes the following procedures:

- A successful registration of the WebRTC client with the IMS core network is implemented.
- IMS-based IPTV generic capabilities like service discovery and selection, service control and media control are realized.
- The media delivery procedure with the WebRTC specifics like secure RTP transport.

For testing the concept the Google Chrome browser in version 31, which supports WebRTC, is used. The basis of this testbed is formed by an open-source IMS core network implementation originating from Fraunhofer FOKUS institute [24]. An Apache HTTP Web Server provides the web application. The WebRTC client is implemented by using HTML5 and JavaScript. Based on this the Graphical User Interface (GUI) of the client is a responsive web site design using the jQuery mobile framework. This framework makes web sites accessible on all smartphone, tablet and desktop devices [25]. The client's source code, based on JavaScript, utilizes the WebRTC API.

For the W2 interface, located between the WebRTC client and the SGW, a proprietary signaling protocol is defined, which is formatted in JavaScript Object Notation (JSON) and is transmitted through a WebSocket connection.

The SGW is written in C# and designed to handle several WebRTC Client sessions simultaneously. The prototyped SGW provides the main functionalities for the interaction with the Gm and the Xa interface. Based on the sipsorcery project, an enhanced SIP protocol stack supporting IMS specific extensions is implemented [26].

The implementation of the Xc interface for session controlling is not considered yet, because it is necessary only for advanced media streaming control functions, such as 'pause' or 'fast forward'.

Also the IMS-based IPTV components are prototyped, which are results from a related student's bachelor thesis [27]. All prototyped IPTV components, written in Java, are based on the technical specification [2].

The successful implementation of the MDF considers the modifications presented in Figure 4. The control engine parses the session information, passed by the MCF, and operates the ICE Agent and the streaming engine. Open-source frameworks are used for

- The ICE agent with the STUN functionalities (icedjava) [28],
- The DTLS key exchange (BouncyCastle) [29],
- The SRTP implementation (srtplight) [30],
- And the streaming server (FFmpeg) [31].

IV. CONCLUSION AND FUTURE WORK

The presented concept of the combination of IMS-based IPTV and WebRTC has a huge sustainable potential. The authors believe that the web based real-time communication is inevitable for future telecommunication. For Telcos it is possible, with these mentioned modifications of the IMS-based IPTV architecture, to deliver the IPTV service to a wider range of end-user's devices. Future developments of new combined services will be realized, using the advantages

of both technologies. The possibility and practicability of the concept is verified by an implemented prototype. The main advantages of the proposed architecture are:

- No media gateway function for live transcoding of the content is needed.
- The modified MCF and the MDF can coexist to other original media functions.
- The scalability of the IMS-based IPTV architecture remains unaffected.

The presented use case, audio and video content on demand, has a big potential for multimedia users: an online video streaming platform without using any proprietary software on the users' devices could be realized. Thus, the user can access this service having the same experience everywhere, with any device.

Currently, the implementation of session mobility for content on demand services is in progress. This includes a dialog state awareness service. Thereby users are able to obtain information about active sessions of their own several devices. The next step, the authors are focusing on, is to implement QoS characteristics known from IMS-based core networks to work with web application based WebRTC clients. For future work, combined services like IMS-based IPTV with presence services are conceivable.

ACKNOWLEDGMENT

The Telekom Innovation Laboratories and the Hochschule fuer Telekommunikation Leipzig are actively cooperating since 2011. Both are working together on common topics in the area of IMS-based services participating in ongoing projects relating WebRTC with Telco Assets. This paper also presents some of the results and acquired experience arising from various final students' theses.

The authors would like to thank Mr. Sven Walter to provide IMS-based IPTV typical components like SCF, SDF and MCF for the testbed.

REFERENCES

- [1] C. Jennings, A. Narayanan, D. Burnett, and A. Bergkvist, "WebRTC 1.0: Real-time communication between browsers," W3C, W3C Working Draft, Sep. 2013, <http://www.w3.org/TR/2013/WD-webrtc-20130910/> [retrieved: Dec., 2013].
- [2] ETSI, "Telecommunications and Internet converged Services and Protocols for Advanced Networking (TISPAN); IPTV Architecture; IPTV functions supported by the IMS subsystem," European Telecommunications Standards Institute (ETSI), TS 182027 v3.5.1, Mar. 2011, Available: http://www.etsi.org/deliver/etsi_ts/182000_182099/182027/03.05.01_60/ts_182027v030501p.pdf [retrieved: Dec., 2013].
- [3] ETSI, "Telecommunications and Internet converged Services and Protocols for Advanced Networking (TISPAN); IMS-based IPTV stage 3 specification," European Telecommunications Standards Institute (ETSI), TS 183063 v3.5.2, Mar. 2011, Available: http://www.etsi.org/deliver/etsi_ts/183000_183099/183063/03.05.02_60/ts_183063v030502p.pdf [retrieved: Dec., 2013].
- [4] 3GPP, "IP Multimedia Subsystem (IMS); Stage 2," 3rd Generation Partnership Project (3GPP), TS 23.228 v8.12.0, Mar. 2010, Available: <http://ftp.3gpp.org/specs/html-info/23228.htm> [retrieved: Dec., 2013].
- [5] 3GPP, "3rd Generation Partnership Project; Technical Specification Group Services and System Aspects; Presence Service; Architecture and functional description," 3rd Generation Partnership Project (3GPP), TS 23.141 v10.0.0, Mar. 2011, Available: <http://ftp.3gpp.org/specs/html-info/23141.htm> [retrieved: Dec., 2013].
- [6] E. Mikoczy, S. Schumann, and P. Podhradsky, "Personalization of internet protocol television (iptv) services in next-generation networks (ngn) architectures," Proceedings of the 8th International Conference on

- Advances in Mobile Computing and Multimedia, November 2010, pp. 366 – 369.
- [7] 3GPP, “Network architecture (Release 8),” 3rd Generation Partnership Project (3GPP), TS 23.002 V8.7.0, Dec. 2010, Available: <http://3gpp.org/ftp/Specs/html-info/23002.htm> [retrieved: Jan., 2014].
- [8] Google Inc. Webrtc. [Online]. Available: <http://www.webrtc.org/> [retrieved: Jan., 2014]
- [9] IETF. Rtcweb status pages. [Online]. Available: <http://tools.ietf.org/wg/rtcweb/> [retrieved: Jan., 2013]
- [10] C. Jennings, T. Hardie, and M. Westerlund, “Real-time communications for the web,” *Communications Magazine*, vol. 51, no. 4, April 2013, pp. 20 – 26.
- [11] J. Valin and C. Bran, “WebRTC Audio Codec and Processing Requirements draft-ietf-rtcweb-audio-04,” Internet-Draft, Internet Engineering Task Force, Jan. 2014, Available: <http://tools.ietf.org/html/draft-ietf-rtcweb-audio-04> [retrieved: Jan., 2014].
- [12] IETF. Ietf discussion forum for video in webrtc. [Online]. Available: <http://www.ietf.org/mail-archive/web/rtcweb/current/maillist.html> [retrieved: Jan., 2014]
- [13] C. Perkins, M. Westerlund, and J. Ott, “Web Real-Time Communication (WebRTC): Media Transport and Use of RTP draft-ietf-rtcweb-rtp-usage-11,” Internet-Draft, Internet Engineering Task Force, Dec. 2013, Available: <http://tools.ietf.org/id/draft-ietf-rtcweb-rtp-usage-11.txt> [retrieved: Jan., 2014].
- [14] E. Rescorla and N. Modadugu, “Datagram Transport Layer Security Version 1.2,” RFC 6347 (Proposed Standard), Internet Engineering Task Force, Jan. 2012, Available: <http://www.ietf.org/rfc/rfc6347.txt> [retrieved: Jan., 2014].
- [15] E. Rescorla, “WebRTC Security Architecture draft-ietf-rtcweb-security-arch-07,” Internet-Draft, Internet Engineering Task Force, Jul. 2013, Available: <http://tools.ietf.org/id/draft-ietf-rtcweb-security-arch-07.txt> [retrieved: Jan., 2014].
- [16] J. Rosenberg, R. Mahy, P. Matthews, and D. Wing, “Session Traversal Utilities for NAT (STUN),” RFC 5389 (Proposed Standard), Internet Engineering Task Force, Oct. 2008, Available: <http://www.ietf.org/rfc/rfc5389.txt> [retrieved: Jan., 2014].
- [17] J. Rosenberg, “Interactive Connectivity Establishment (ICE): A Protocol for Network Address Translator (NAT) Traversal for Offer/Answer Protocols,” RFC 5245 (Proposed Standard), Internet Engineering Task Force, Apr. 2010, updated by RFC 6336. Available: <http://www.ietf.org/rfc/rfc5245.txt> [retrieved: Jan., 2014].
- [18] H. Alvestrand, “Overview: Real Time Protocols for Browser-based Applications draft-ietf-rtcweb-overview-08,” Internet-Draft, Internet Engineering Task Force, Sep. 2013, Available: <http://tools.ietf.org/id/draft-ietf-rtcweb-overview-08.txt> [retrieved: Jan., 2014].
- [19] 3GPP, “Study on Web Real Time Communication (WebRTC) access to IP Multimedia Subsystem (IMS); Stage 2,” 3rd Generation Partnership Project (3GPP), TS 23.701 v12.0.0, Dec. 2013, Available: <http://www.3gpp.org/DynaReport/WiVsSpec--580062.htm> [retrieved: Jan., 2014].
- [20] J. Fischl, H. Tschofenig, and E. Rescorla, “Framework for Establishing a Secure Real-time Transport Protocol (SRTP) Security Context Using Datagram Transport Layer Security (DTLS),” RFC 5763 (Proposed Standard), Internet Engineering Task Force, May 2010, Available: <http://www.ietf.org/rfc/rfc5763.txt> [retrieved: Jan., 2014].
- [21] D. McGrew and E. Rescorla, “Datagram Transport Layer Security (DTLS) Extension to Establish Keys for the Secure Real-time Transport Protocol (SRTP),” RFC 5764 (Proposed Standard), Internet Engineering Task Force, May 2010, Available: <http://www.ietf.org/rfc/rfc5764.txt> [retrieved: Jan., 2014].
- [22] I. Fette and A. Melnikov, “The WebSocket Protocol,” RFC 6455 (Proposed Standard), Internet Engineering Task Force, Dec. 2011, Available: <http://www.ietf.org/rfc/rfc6455.txt> [retrieved: Jan., 2014].
- [23] M. Baugher, D. McGrew, M. Naslund, E. Carrara, and K. Norrman, “The Secure Real-time Transport Protocol (SRTP),” RFC 3711 (Proposed Standard), Internet Engineering Task Force, Mar. 2004, updated by RFC 5506. Available: <http://www.ietf.org/rfc/rfc3711.txt> [retrieved: Jan., 2014].
- [24] Fraunhofer FOKUS. Fokus open ims playground. [Online]. Available: <http://www.fokus.fraunhofer.de/en/fokus/index.html> [retrieved: Jan., 2014]
- [25] jQuery Mobile: a HTML5-based user interface system. [Online]. Available: <http://jquerymobile.com/> [retrieved: Feb., 2014]
- [26] T. Bach. sipsorcery-fork. [Online]. Available: <https://github.com/hftl-ims-research/sipsorcery-fork> [retrieved: Apr., 2014]
- [27] S. Walter, “Design and Implementation of an IMS-based IPTV Content on Demand Service,” Bachelorthesis, Hochschule fuer Telekommunikation Leipzig, HfTL, Aug. 2013.
- [28] inspired social. Open source ICE implementation. [Online]. Available: http://code.google.com/p/inspired-social/source/browse/trunk/StunServer/net/mc_cubed/icedjava?spec=svn20&r=20#icedjava [retrieved: Apr., 2014]
- [29] Legion of the Bouncy Castle Inc. Bouncy Castle Crypto API. [Online]. Available: <http://www.bouncycastle.org> [retrieved: Apr., 2014]
- [30] steely gint. set of classes implementing a simple (S)RTP stack. [Online]. Available: <https://github.com/steely-gint/srtplight> [retrieved: Apr., 2014]
- [31] FFmpeg: a open source cross-platform solution to record, convert and stream audio and video. [Online]. Available: <http://www.ffmpeg.org/index.html> [retrieved: Feb., 2014]

Aspect-Oriented Implementation of Concurrent Processing Design Patterns

Shingo Kameyama, Masatoshi Arai, Noriko Matsumoto, Norihiko Yoshida
 Graduate School of Science and Engineering
 Saitama University
 Saitama, Japan
 {shingo, arai, noriko, yoshida}@ss.ics.saitama-u.ac.jp

Abstract—A variety of design patterns are now widely used in software development as their catalog is a collection of knowledge on design and programming techniques and namely elaborated patterns. However, as each design pattern is described in the forms of texts, charts, and simple code examples, it has some limitations in applicability and formal treatment. One of its reasons is that the design patterns include some crosscutting concerns. To solve this problem, aspect-oriented implementation of the so-called “Gang of Four” (GoF) design patterns, which are cataloged for component reuse has been proposed. In this paper, we propose aspect-oriented implementation of design patterns for concurrent processing, so as to improve and accelerate design and development processes of, for example, network systems, embedded systems, and transaction systems. Our aspect-oriented implementation tailors hierarchical or inclusive relationships among design patterns well which are not found in the patterns for component reuse, but found in the patterns for concurrent processing.

Keywords—Design patterns; aspects; concurrency

I. INTRODUCTION

A variety of design patterns (patterns in short, sometimes, hereafter) [1] are now widely used in software development as their catalog is a collection of knowledge on design and programming techniques and namely elaborated patterns. Their catalog enables novices to refer to experts’ knowledge and experiences in software design and development. It accelerates software productivity and improves qualities as well as it helps communication within a development team.

Many researchers have proposed various patterns. However, patterns have sometimes difficulties in formal or systematic treatment. It is because patterns are described in the forms of texts in a natural language, charts, figures, and code fragments and samples, not in any formal description.

To solve this problem, aspect-oriented [2] implementation of patterns for component reuse has been proposed [3][4][5][6][7]. Patterns for component reuse, or sometimes called GoF (Gang of Four) design patterns, were patterns promoting component reuse in object-oriented software [1]. Aspects are a technique to modularize codes which are scattered to several modules [2].

In this paper, we propose aspect-oriented implementation of concurrent processing design patterns [8][9][10]. Different from GoF design patterns, most concurrent processing patterns use other concurrent processing patterns related to them, namely, they are not independent to each other. We implement them by a combination of the implementation of related concurrent processing pattern. This approach can be applied also to aspect-oriented implementation of other design patterns where there are some mutual dependency among them.

This paper is organized as follows: Section 2 gives an overview of design patterns. Section 3 presents an overview and some functions of aspect-orientated implementation. Section 4 introduces an example of aspect-oriented implementation of GoF design patterns. Section 5 and Section 6 explain aspect-oriented implementation of concurrent processing design patterns. Section 7 gives some considerations, and Section 8 contains some concluding remarks.

II. DESIGN PATTERNS

Design patterns are a catalog of typical solutions for some typical problems in designing and programming in software development. Erich Gamma, Richard Helm, Ralph Johnson, and John Vlissides, so-called the Gang of Four (GoF), introduced 23 design patterns for reuse (so-called the GoF design patterns) [1]. Design patterns have been getting widely accepted as a technique to improve software development processes. Many researchers have proposed various design patterns for concurrent processing, real-time processing, and web applications, etc.

A. Benefits

Design patterns accelerate software productivity and improves qualities in the following regards.

- A workload for programming and verification can be reduced reusing codes in design patterns.
- A development beginner can use it as a guide in which development experts’ know-how is accumulated.
- A developer can tell a software design to another developer concisely and precisely by describing the name of a design pattern.

B. Problems

A design pattern has the following problems because its codes are scattered to more than one module.

- A developer must understand the structure of a design pattern to apply, extract some necessary codes from its sample codes, and apply to a program.
- It is difficult to maintain or extend the program with design patterns because codes relating the patterns are scattered across several modules, and a developer must keep all of them in mind correctly.

C. Concurrent Processing Design Patterns

Besides GoF design patterns, concurrent processing design patterns [8][9][10] have been proposed. A concurrent processing design pattern offers a typical solution to the following problems in concurrent processing.

- There are problems, such as race conditions and deadlocks, that do not happen in single threaded processing. A race condition happens when more than one thread read and write a shared resource without mutual exclusion control. A deadlock happens when more than one thread lock more than one object and wait for unlocking with each other.
- Verification is difficult because any problem may or may not occur, but possibly not always.

Different from the GoF design patterns for component reuse, most design patterns for concurrent processing have relations with each other, i.e., they are not independent to each other, and most patterns use other patterns related to them. In other words, there is an hierarchy among them. We explain concurrent processing design patterns below. Figure 1 shows the hierarchical relation among them. An arrow denotes that a design pattern on the head side uses a design pattern on the tail side.

- Single Threaded Execution (or Critical Section)
 - It ensures safety by exclusive control. Only a single thread can access a thread-unsafe object at a time.
- Immutable
 - It improves throughput by eliminating changing a state of an object and exclusive control.
- Guarded Suspension (or Guarded Waits or Spin Lock)
 - It ensures safety by blocking a thread until a state of an object changes if a precondition is not met.
- Balking
 - It ensures safety and improves responsibility by not executing a processing if a precondition is not met.
- Producer-Consumer
 - It ensures safety and improves throughput by passing an object indirectly.
- Read-Write Lock (or Readers and Writers)
 - It ensures safety and improves throughput by allowing concurrent access for read and requiring exclusive access for write.

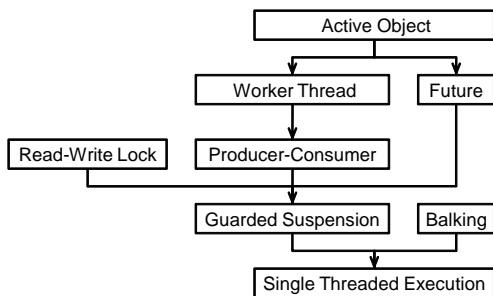


Figure 1. Inclusive relation of concurrent processing design patterns.

Thread-per-Message (or Thread-per-Method)

It improves responsibility by creating a thread every processing request, and leaving the processing to the thread.

Worker Thread (or Thread Pool)

It improves responsibility, throughput, and capacity by leaving a requested processing to another thread, and reusing the thread.

Future

It improves responsibility by enabling a thread to receive a result when it is needed, instead of waiting for a result.

Two-Phase Termination

It ensures safety by terminating another thread indirectly.

Thread-Specific Storage (or Thread-Specific Data)

It enables a thread to execute the thread-specific processing by providing the thread-specific storage.

Active Object (or Actor)

It improves responsibility and enables more than one thread to request a processing to a thread-unsafe object at a time by leaving the processing to a single thread.

III. ASPECTS

Object-oriented programming is a technique that modularizes a concept and a concern as an object, and improves maintainability and extensibility. However, modularization by objects has a limitation because codes of a crosscutting concern related to more than one module are scattered on them. A crosscutting concern is a functionality of the system whose definition appear in several classes. Examples of crosscutting concerns include logging and caching.

Aspect-oriented programming [2] is a technique that compensates the above-mentioned limit in object-oriented programming and modularizes crosscutting concerns.

A. Overview

In object-oriented programming, a module calls a method in another module to execute. When a module is added or removed, it is necessary to adjust all the method calls related to it in all the other modules. Figure 2 shows a class diagram depicting this scheme. An arrow between methods expresses a processing flow.

In aspect-oriented programming, a method or a code fragment specifies positions in other method definitions in any module where it must be called. This method or code fragment runs when execution reaches the position. The specification of the position is called a pointcut. The code fragment is called an advice. This module composed of pointcuts and advices is called an aspect.

When a module is added or removed, it is not necessary to adjust other modules that call it because there are no explicit method calls related to it. Figure 3 shows a class diagram depicting this scheme. A pointcut and an advice are described at the bottom part of boxes which represent aspects. A dashed line expresses a specification by a pointcut.

B. AspectJ

AspectJ is an aspect-oriented language extension adding the aspect-oriented features to the object-oriented language Java. Java and AspectJ are used in a related study [3] and this study. We explain some functions of AspectJ here.

Pointcut

A position that can be specified by a pointcut is limited to a position where a method is called and executed, and where a field value is retrieved or assigned.

Advice

An advice includes a before advice, an after advice, and an around advice. A before advice is executed just before processing of the position specified by a pointcut. An after advice is executed just after it. An around advice is executed instead of it. In an around advice, *proceed* which means execution of the original method can be specified.

Aspect

Like a class, an aspect can include fields and methods, and can be defined as an abstract aspect or a concrete aspect. An aspect can also inherit another aspect.

IV. RELATED WORKS

A design pattern includes several classes in general, therefore a pattern itself is a crosscutting concern. Aspects are expected to enable modularization of the design pattern as a crosscutting concern. Aspect-oriented implementation of the GoF design patterns has been proposed in related studies [3].

In this section, we explain aspect-oriented implementation of the *observer* pattern which is one of the GoF design patterns as an example of the related study [3]. The *Observer* is to execute a method whenever a state of another object changes.

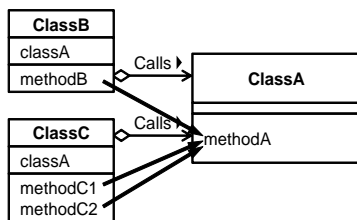


Figure 2. Modularization using objects.

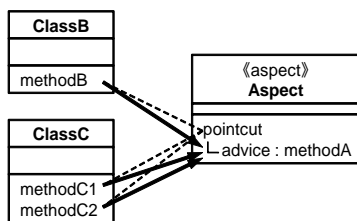


Figure 3. Modularization using aspects.

A. Aspect-Oriented Implementation of Observer

In object-oriented implementation, the *observer* is implemented as follows: a method *notifyObservers* is called just after any change of a state, and calls linked methods. The following codes of *observer* are defined in more than one module. Figure 4 shows a class diagram of this scheme.

- *NotifyObservers* and method calls to execute it.
- A field and methods to manage linked objects.
- A super-class of linked objects.

In aspect-oriented implementation, *observer* is implemented as follows: the code to execute is defined as an after advice instead of *notifyObservers*, and the codes of *observer* is modularized in a single aspect. Figure 5 shows a class diagram of this scheme. *Observer* is actually implemented as an abstract aspect and a concrete aspect. The abstract aspect defines an advice which does not depend on a target program where a design pattern is applied. The concrete aspect defines a pointcut which depends on it. This improves reusability of the abstract aspect.

B. Benefits

There are the following benefits because codes of a design pattern are modularized.

- A developer can apply a design pattern to the position in any target system specified by a pointcut even if he or she does not understand its structure.
- A developer need to update only one module regarding the design pattern when the target system to which the pattern is applied is updated.
- The design pattern is defined separately from the target system, therefore it is easier to maintain and reuse the system.

V. ASPECT-ORIENTED IMPLEMENTATION OF READ-WRITE LOCK

We implemented all the concurrent processing design pattern mentioned earlier as an abstract aspect and a concrete aspect like *observer*, and confirmed that this implementation works correctly.

In this section, we explain aspect-oriented implementation of *read-write lock* as an example of our study. A set of locking

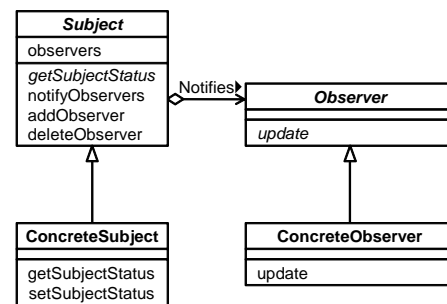


Figure 4. Observer implemented using objects.

and unlocking is a typical example of crosscutting concern, and they are to be done before and after read and write in *read-write lock*.

A. Read-Write Lock Defined Using Objects

In object-oriented implementation, *read-write lock* is implemented by testing a precondition and locking just before read and write, and unlocking and notifying a change just after reading and writing under exclusive control. Figure 6 shows a class diagram of this scheme.

Locking and unlocking are done by counting reading and writing threads, and the counters are used for testing a precondition. The precondition for read is that there is no writing thread. The precondition for write is that there is no reading and writing thread. A thread waits until the state of the counter changes if the precondition is not met. Exclusive control is done by an instance lock that is a feature of Java. The instance lock is also necessary to execute *wait* and *notifyAll*.

B. Read-Write Lock Defined Using Aspects

Guarded suspension is used for testing a precondition and notifying a change of the state in *read-write lock*, and *single threaded execution* is used for exclusive control in *guarded suspension*. This inclusive, or hierarchical relation is shown in Figure 1.

In aspect-oriented implementation, we implement the *read-write lock* including aspect implementations of *single threaded*

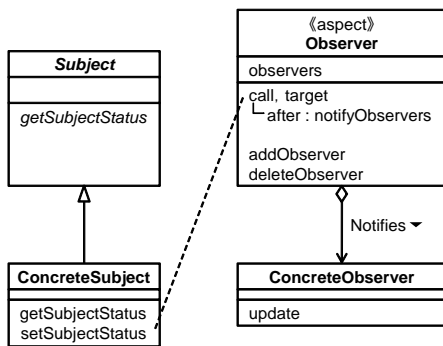


Figure 5. Observer implemented using aspect.

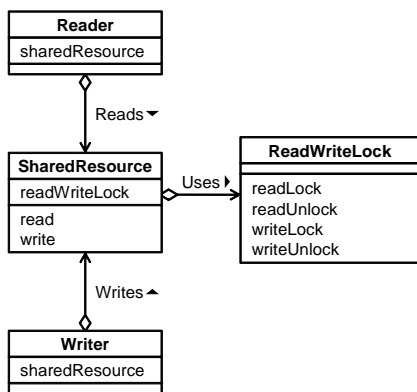


Figure 6. Read-Write Lock using objects.

execution and *guarded suspension*. Figure 7 shows a class diagram of *read-write lock*.

It is necessary to apply the aspect of *single threaded execution* before the aspect of *guarded suspension*. The first reason is that testing a precondition and notifying a change in state are done under exclusive control. The second reason is that *wait* and *notifyAll* are used. An abstract aspect of *guarded suspension* defines a precedence of concrete aspects by adding + after names of the abstract aspects because if a name of a concrete aspect is used, the abstract aspect need to be updated when adding, deleting, or changing it.

In object-oriented implementation, because a class manages a lock and counters, an instance of the class is created for every object to be read and written. In aspect-oriented implementation, because the aspect manages the lock and counters. A test of preconditions for reading and writing is defined in two concrete aspects, and the aspects are defined as privileged to access to the counter which is a private field in another aspect.

VI. ASPECT-ORIENTED IMPLEMENTATION OF OTHER CONCURRENT PROCESSING DESIGN PATTERNS

In this section, we present brief summaries of aspect-oriented implementation of the other concurrent processing design patterns. If a design pattern uses other concurrent processing design patterns, we implemented the design pattern by including the aspects for them as shown in the *read-write lock*.

We are not successful yet in implementing *thread-specific storage* in aspects. The *Thread-specific storage* is a design pattern to re-implement an object as an object with the same API, and such a major structural change is difficult to implement using aspects.

Immutable

Immutable cannot be implemented by an advice because the pattern does not execute any code. Instead, We implemented a checking feature whether a field of a target object is assigned from outside the object or not.

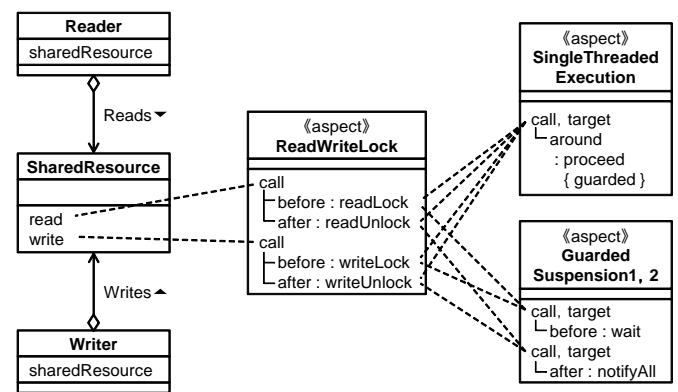


Figure 7. Read-Write Lock using aspects.

Balking

The pattern needs a feature to terminate a method execution, however an advice has no feature to do it. Therefore, we implemented *balking* using an around advice. *Proceed* (i.e. the original code in the target object) is executed if a precondition is met. The around advice is terminated if a precondition is not met.

Producer-Consumer

Codes for a producer and a consumer must specify codes for send and receive. However, it is difficult for a pointcut to specify an advice because an advice has no name. In our implementation of *producer-consumer*, an advice executes a method to send or receive an object. An pointcut specifies the name of the methods.

Thread-per-Message

When applying *thread-per-message* to a program afterward, in object-oriented implementation, a developer needs to modify the program so that another thread may execute a requested processing. In aspect-oriented implementation, a developer need not modify the program. An around advice handles the processing, and another thread executes the original code.

Worker Thread

In aspect-oriented implementation, an around advice is executed instead of a requested processing, and an abstract aspect defines an inner class that executes the original code using the *proceed* feature.

Future

In aspects implemented *thread-per-message* and *worker thread*, an advice returns the same instance as a return value of a requested processing. The instance and the processing result are related using a map feature.

Two-Phase Termination

In aspect-oriented implementation, a thread terminates itself executing *stop* at a safe position. It is because an advice cannot terminate the execution.

Active Object

When adding a method to a thread-unsafe object, in object-oriented implementation, a developer must implement a new task to execute the method, and add a method to create the task. In aspect-oriented implementation, a developer need only to implement a new concrete aspect of *worker thread*.

In any implementation summarized above, an abstract aspect only defines the abstract structure, and a concrete aspect defines the number and the type of an object to be created. The types of an argument and a return value of the advice are not specific but general *Objects*. These disciplines are to improve reusability of the abstract aspects.

VII. CONSIDERATION

We express the inclusive relations among concurrent processing design patterns as a combination of aspects. Regarding this, we had to do some refactoring on aspect implementations of the patterns.

For example, when including an aspect for *single threaded execution* in an aspect for *guarded suspension*, first implementation did exclusion control using an instance lock of the aspect for *single threaded execution*, and the aspect for *guarded suspension* could not execute *wait* and *notify*. Therefore, we modified the aspect for *single threaded execution* so that an instance lock of an object can be used.

However, when including the aspects in an aspect of *read-write lock*, a change in state could not be notified between concrete aspects because *wait* and *notify* are instance methods of the aspect of *guarded suspension*. Therefore, we modified the aspect of *guarded suspension* moreover so that an instance method of a target object can be used.

This approach can be applied also to other aspect-oriented implementations where a design pattern includes other design patterns, as well as concurrent processing design patterns.

VIII. CONCLUSION AND FUTURE WORK

In this paper, we propose aspect-oriented implementation of concurrent processing design patterns, and if a design pattern uses other concurrent processing design patterns, we implemented the design pattern by including aspects implementing them.

We are still at the starting point in this research, and there is still much to do. Below are some future research directions:

- Refinement of aspect-oriented implementation of concurrent processing design patterns.
- Methodology of aspect-oriented implementation and refactoring in inclusion of aspects.
- Categorization of concurrent processing design patterns from the point of aspect-oriented implementation.

REFERENCES

- [1] E. Gamma, R. Helm, R. Johnson, and J. Vlissides, Design Patterns, Addison-Wesley, 1994.
- [2] R. Laddad, AspectJ in Action, Manning, 2003.
- [3] J. Hannemann and G. Kiczales, "Design Pattern Implementation in Java and AspectJ", Proc. ACM OOPSLA, 2002, pp. 161–173.
- [4] M.L. Bernardi and G.A. Di Lucca, "Improving Design Pattern Quality Using Aspect Orientation", Proc. IEEE STEP, 2005, pp. 206–218.
- [5] A. Garcia, C. Sant'Anna, E. Figueiredo, U. Kulesza, C. Lucena, and A. von Staa, "Modularizing Design Patterns with Aspects: a Quantitative Study", Proc. ACM AOSD, 2005, pp. 3–14.
- [6] M. Bynens and W. Joosen, "Towards a Pattern Language for Aspect-Based Design", Proc. ACM PLATE '09, 2009, pp. 13–15.
- [7] Z. Vaira and A. Caplinskas, "Case Study Towards Implementation of Pure Aspect-Oriented Factory Method Design Pattern", Proc. IARIA PATTERNS '11, 2011, pp. 102–107.
- [8] D. Lea, Concurrent Programming in Java, Addison-Wesley, 1999.
- [9] A. Holub, Taming Java Threads, Apress, 2000.
- [10] M. Grand, Patterns in Java, Volume 1, Wiley, 2002.

Analysis of the Development Process of a Mutation Testing Tool for the C++ Language

Pedro Delgado-Pérez, Inmaculada Medina-Bulo, Juan José Domínguez-Jiménez
 UCASE Software Engineering Group, Department of Computer Science and Engineering
 University of Cádiz, Cádiz, Spain
 Email: {pedro.delgado, inmaculada.medina, juanjose.dominguez}@uca.es

Abstract—Mutation testing is a fault-based software testing technique to measure the quality of a test suite depending on its ability to detect faults in the code. This technique has been applied to an assortment of languages of very diverse nature since its inception in the late 1970s. However, the researchers have postponed its development around C++ in favor of other mainstream languages. This paper aims to survey the mutation testing research regarding C++, studying the existing tools and approaches. To the same extent, we discuss the different aspects that should be taken into account in the construction of a comprehensive mutation tool for this language, from the analysis of the code to the execution of the mutants. In addition, we expound how the technique can be assessed so that it can contribute effectively in the composition of a complete test suite. The findings in this paper pose that the construction of a mutation tool for this language is complex, but still realizable.

Keywords—Mutation testing; Mutation tool; C++.

I. INTRODUCTION

Mutation testing is a fault injection technique used to determine the ability of a test suite locating errors in the code [1]. The effectiveness of error detection of a test suite is defined as the percentage of faults that can be detected by their test cases. The technique involves the creation of *mutants*, i.e., versions of the original program with a simple syntactic change. These faults are injected in the code through the mutation operators, which are based on common mistakes made by programmers in a certain programming language.

Mutation testing aims to ensure that a test suite is able to detect all those typical mistakes when comparing the output of the original program and the mutated version for the same test cases; a mutant is *killed* when the output is different for at least one test case, but remains *alive* if the output keeps unaltered. In this latter case, either a new test case is needed to detect the fault or the mutant is completely *equivalent* to the original program.

Mutation testing is a white-box testing technique, so it has to be studied around a particular language. Thus, the technique has been successfully applied to several languages, finding tools that automate the generation of mutants for a wide range of them [2]. Unlike other languages, C++ is clearly behind in the research and in practice. In literature, most research on mutation testing has been focused on procedural programming paradigm. For object-oriented (OO) languages, the research has put the focus on Java or C# [3]. In contrast, the few existing works about mutation testing and C++ [4], [5] are unfinished and various matters are pending.

The overall goal of this paper is to analyze the state of the art in order to devise the construction of a mutation tool for C++. The matters that should be handled for that purpose have been surveyed, from the composition of a catalog of mutation operators to the execution of tests and the manner to assess the results. Consequently, the information introduced throughout the sections lays the foundations for the application of mutation testing to C++, exposing an approach to follow in the future based on the abstract syntax tree for the insertion of errors.

The current state regarding the mutation operators, the mutation tools and the language itself are explored in Section II. Section III deals with the operators at different levels, the steps to accomplish in the creation of a mutation tool and the techniques to mutate the code. In Section IV, the way to gear the evaluation of the results in order to judge the operator behavior is commented, giving special emphasis to the issues that can arise in the analysis of results. Finally, the last section presents the conclusions and also the future work.

II. RELATED WORK

This section looks in depth the existing background around C++ and mutation testing.

A. Mutation Operators

To illustrate the underlying idea in mutation testing, we can consider the C++ code fragment below:

```
if(a > 100){...}
```

If we have defined a mutation operator replacing the relational operator '>', the fragment above can be modified creating a mutant like the following:

```
if(a < 100){...}
```

A set of mutation operators can be defined for each of these levels [2], [6]:

- **Unit level:** Standard operators applied indistinctly to a function or method, checking its correctness. These operators are usually known as *traditional operators*.
- **Class level:** This level deals with the mutation of OO features.
- **Integration level:** Intermediate level between the unit and the system levels, checking the function invocations.

- **Multi-class level:** Operators at this level are intended to test a complete program: interactions among functions, classes, etc.

As advanced in the introduction, we can state that the development of mutation testing with respect to C++ is underrepresented. Regarding a particular set of mutation operators for this language, the research accomplished is really scarce as we cannot find a comprehensive catalog of operators.

However, two attempts have been performed up to now. The first work in [4] composed a set of traditional operators; it is named in [7], a paper aiming to supply the equivalence of operators among different languages. These operators are based on the operators defined for Ada and the Fortran operators used by the tool *Mothra* [8]. This approach cataloged the operators in four blocks: *operand replacement*, *operator insertion*, *arithmetic operator replacement* and *relational operator replacement*. These groups and their operators are shown in Table I. On the other hand, several mistakes regarding OO features have been enumerated around C++ in [5]. This paper poses five categories of possible faults: *Inherit*, *Associate*, *Access*, *Object* and *Member*. Nevertheless, the three first blocks are applied to the Unified Modeling Language (UML) specification and only the errors belonging to the *Object* and *Member* groups refer to the C++ code. The faults exposed in that paper are summarized in Table II. The research regarding both approaches seems given up as no new progress has been published since then.

B. Mutation Tools

At present there are a variety of tools for several languages implementing mutation testing through different techniques [2], as *Mothra* for Fortran [8], *MuJava* for Java [9]

TABLE I. STANDARD OPERATORS PROPOSED BY ZHANG [4].

Block	Operator	Description
Operand Replacement Ops.	OVV	Variable replaced by a variable
	OVC	Variable replaced by a constant
	OVA	Variable replaced by an array reference
	OVP	Variable replaced by a pointer reference
	OVC	Constant replaced by a variable
	OCC	Constant replaced by a constant
	OCA	Constant replaced by an array reference
	OCP	Constant replaced by a pointer reference
	OAV	Array reference replaced by a variable
	OAC	Array reference replaced by a constant
	OAA	Array reference replaced by an array reference
	OAP	Array reference replaced by a pointer reference
	OAN	Array name replaced by an array name
	OPV	Pointer reference replaced by a variable
	OPC	Pointer reference replaced by a constant
	OPA	Pointer reference replaced by an array reference
OPP	Pointer reference replaced by a pointer reference	
OPN	Pointer name replaced by a pointer name	
Operator Insertion Ops.	IBO	Binary Operators Insertion
	IOU	Unary operator insertion
Arithmetic Operator Replac. Ops.	AOR	Arithmetic operator replacement
Relational Operator Replac. Ops.	ROP	Relational operator replacement

TABLE II. FAULTS IDENTIFIED BY DEREZIŃSKA [5] FOR THE *Object* AND *Member* CATEGORIES.

Category	Description
Object	Calls a same function member from a different object of the same class.
	Calls a function from an object of a different class, but both classes have the common base class.
	Calls a function from the derived class instead of the base class.
Member	Calls a different (complementary) function member.
	Calls a function inherited from the base class.
	Swaps calling of function members in a class.
	Swaps calling of functions inherited from one class.
	Accesses the different data in the same object.

or *MILU* for C [10]. Concerning C++, only commercial tools can be found tackling the mutation analysis for the C++ code: *Insure++* [11] from Parasoft, *PlexTest* [12] from ItRegister and *Certitude* [13] from SpringSoft. These three tools appear in this known survey around mutation testing [2]. In that paper, it is also stated that the *ESTP* tool can be applied to C++, but it is only devoted to the C language actually.

Table III displays a summary with the main characteristics of the tools discussed (shown in [2]). All of them are applicable to C and C++, as both languages share much of their syntax. A description of these commercial products is given below:

- *Insure++*: This tool uses mutation testing as one more technique to enhance the software quality. Its approach is somewhat different from classical mutation testing because it only creates *functionally equivalent mutants*. These equivalent mutants are expected to pass the tests (are still alive after their execution against the test suite) instead of failing. *Insure++* only performs some standard mutations as mentioned in [14].
- *PlexTest*: This product implements a *highly selective* mutation testing. Thus, it only performs the mutation of deletion, i.e., the removal of an instruction. This approach tries to avoid the generation of equivalent mutants present in non-selective and full-selective mutation testing.
- *Certitude Functional Qualification System*: This tool combines the mutation testing technique with static analysis, qualifying a program functionally and discovering faults that might not be detected otherwise. Although this product has also been used for the analysis of software systems, it is now addressing the microelectronics industry in order to simulate the performance of a digital circuit before the design is finally accomplished.

As a conclusion of the information shown, we can state that these tools are different in terms of the operators supported and the techniques to accelerate the testing process, but they are not applicable to the whole language. Moreover, the products presented are not absolutely centered on mutation testing (as *Insure++*) or they are not only used for C++ (as *Certitude*).

Java and C# are the OO languages that have been mostly tackled by researchers. The construction of frameworks for these OO languages are mainly based on the insertion of faults directly in the bytecode, like *MuJava* [9], and also on the reflection mechanism to analyze the original program and

TABLE III. EXISTING MUTATION TESTING TOOLS FOR C++.

	Application	Year	Character	Available
Insure++	C/C++	1998	Source Code Instrumentation	Commercially
Plextest	C/C++	2005	General	Commercially
Certitude	C/C++	2006	General	Commercially

determine where the operators can be used. Derezińska et al. [14] has proposed and built two different tools: *CREAM*, using a parser-based approach, and *ILMutator*, which manipulates both the meta-data and the intermediate code generated from C# to insert the faults modeled by the operators. The former approach is better in identifying where the operators can be applied and complying with the correctness conditions, but the latter is more efficient as no recompilation is needed.

C. The C++ Language

C++ is a complex programming language when compared to other similar general purpose languages, such as Java or C#. This language includes a great variety of alternatives. Besides, many of the features are confusing for the programmers, being usual that they neglect using some of them. Thus, we found surprising that this popular language has been omitted in the field of mutation testing up to now. Because of being so popular, we consider that it is worthwhile addressing the application of this technique to C++ in order to harness the contributions of mutation testing for this language, which presents various complicated characteristics for the programmers.

In this regard, we can mention several particular features, such as the destruction and construction of objects, the existence of pointers and references, the use of exception handling or the inclusion of templates. Multiple inheritance is another characteristic which can lead to the creation of new specific mutation operators as well as should be taking into account in general when tackling class-level operators relating inheritance, as more than a single class needs to be considered.

C++ is a language in continuous change and some standards have been approved since 1998, when the first standard appeared. A new standard, *C++11* [15], was ratified in 2011 to replace the previous standard, *C++03*. It represents the first substantial change since 1998. Besides some minor modifications, the *C++11* standard presents important changes. This standard also provides new features and extends the *C++ Standard Library*. Currently, the changes introduced by this standard are taking place gradually and even compilers are not fully adapted to the new proposals. Nonetheless, the *C++14* standard is already being prepared to provide the language with new functionalities.

III. DEVELOPMENT PROCESS OF A C++ MUTATION TOOL

Mutation testing confronts two main challenges when constructing a mutation tool. The technique to insert mutations in the code can be useful to handle certain features of the language, but might not cover all its elements. In the case of C++, as a mainstream language which provides a great range of alternatives (unlike specification languages, for instance), it is necessary to tackle the structures of the language in an uniform manner so that the mutations are the expected ones in every operator addressing the same elements. Secondly, the

generation and execution of mutants to obtain the results suppose a high computational cost, especially when considering the size of today's programs. These and other concerns are commented in the following subsections.

A. Catalog of Operators

C++ is a multi-paradigm programming language which is considered as an enhancement of C with new features, such as the manipulation of objects [15]. Mutation testing can be applied at different levels of the language, as mentioned in Section II-A. Taking into account the dimension of C++, it seems challenging to apply this technique to every level of the language all at once. In addition, the more operators are included in the process, the more mutants will be produced, with a consequent increase in the computational expenses. Thus, a new catalog of operators can be composed for each different level so that they can be used conveniently.

The research around C is the starting point concerning traditional operators. The similarity between these two programming languages is well-known, as C++ is derived from C. Most of the works and tools created for C are based on mutation operators described by Agrawal et al. in [16]. In the aforementioned paper, the operators are separated into four main categories: *statement*, *operator*, *variable* and *constant mutations*. Because of the compatibility with C, all operators can be adapted to the language (see Table IV). However, some of them should be slightly modified in their implementation with respect to some characteristics in C++, such as the possibility to declare a variable at any point in a block.

According to the class level, we should survey the existing mutation operators regarding the main OO features in other languages, undertaking the two following steps:

- 1) To observe how the particular characteristics of C++ alter the operators, determining if they can be adapted to the language or they have to be rejected instead.
- 2) Whenever possible, to add specific operators for this language and its distinguishing features (see Section II-C).

Many papers centered on this level can be found and sets of mutation operators have been defined around the OO programming [17], [18]. A similar survey to this paper around OO languages studies the class mutation operators [3]. In Table V, we expose the catalog of operators defined by Offutt et al. for Java [17], which is one of the most prominent and complete set of class mutation operators. Regarding this list and the aforementioned step (1), most of these operators could be adapted with some changes but the same purpose (for example, the *super* keyword is not used in C++), but others are likely to be excluded, such as the last four operators: they refer to methods for which there is no a convention in C++.

To conclude, we have to note that the number of levels of the language is not fixed at all. For instance, another level of operators could be created only for the *C++ Standard Library*, which provides a great range of extra functionalities commonly used. Thus, the level chosen by the user depends on the program characteristics and the type of testing intended to perform. In other words, the subset of operators to be applied should be selected after a preprocessing of the program under

TABLE IV. MUTATION OPERATORS FOR THE C LANGUAGE PROPOSED BY AGRAWAL ET AL. [16].

Block	Operator	Description
Statement mutations	STRP	Trap on Statement Execution
	STRI	Trap on if Condition
	SSDL	Statement Deletion
	SRSR	return Statement Replacement
	SGLR	goto Label Replacement
	SCRb	continue Replacement by break
	SBRC	break Replacement by continue
	SBRn	Break Out to nth Enclosing Level
	SCRn	Continue Out to nth Enclosing Level
	SWDD	while Replacement by do-while
	SDWD	do-while Replacement by while
	SMTT	Multiple Trip Trap
	SMTC	Multiple Trip Continue
	SSOM	Sequence Operator Mutation
	SMVB	Move Brace Up or Down
	SSWM	Switch Statement Mutation
Operator mutations	<i>Obom (Binary Operator Mutations)</i>	
	Ocor	Comparable operator replacement
	Oior	Incomparable operator replacement
	<i>Ouor (Unary Operator Mutations)</i>	
	Oido	Increment/Decrement
	OLNG	Logical Negation
Variable mutations	OCNG	Logical context negation
	OBNG	Bitwise Negation
	OIPM	Indirection Operator Precedence Mutation
	OCOR	Cast operator replacement
	Vsrr	Scalar Variable Reference Replacement
	Varr	Array Reference Replacement
	Vtrr	Structure Reference Replacement
	Vprr	Pointer Reference Replacement
	VSCR	Structure Component Replacement
	VASM	Array reference subscript mutation
Constant mutations	VDTR	Domain Traps
	VTWD	Twiddle Mutations
	CRCR	Required Constant Replacement
	Cccr	Constant for Constant Replacement
	Ccsr	Constant for Scalar Replacement

test (PUT), so that the operators can be adjusted as much as possible to the concrete application.

B. Applying Mutations

Mutation testing process can be subdivided into three principal stages: the analysis of the code, the generation of the mutants and the execution of the test suite. For the development of a mutation tool, the implementation of these three steps is necessary, so we describe the purpose of these parts below:

- **Analyzer:** This module first acts getting the program under test and the set of operators defined. Its role is to determine the mutation operators that can be applied, that is, the mutation locations for each of the operators.
- **Generator:** Taking into account the analysis accomplished in the previous phase, the function of the generator is to create the mutant for its later execution. In the production of the mutant, the generator should take care of the correctness conditions so that the mutant generated is not an invalid mutant, i.e., it can be compiled because the program has no errors.
- **Test suite runner:** This module executes the mutants against a test suite defined for the PUT, finally classifying them into alive (the mutant has not been detected by any test case), dead (the mutant has been killed by one or more test cases) or invalid.

The two first phases are quite related as they require an implementation technique that actually allows us to detect and then

TABLE V. MUTATION OPERATORS AT THE CLASS LEVEL PROPOSED BY OFFUTT ET AL. [17].

Block	Operator	Description	
Encapsulation	AMC	Access modifier change	
	Inheritance	IHI	Hiding variable insertion
		IHD	Hiding variable deletion
		IOD	Overriding method deletion
		IOP	Overriding method calling position change
		IOR	Overriding method rename
		ISI	<i>super</i> keyword insertion
		ISD	<i>super</i> keyword deletion
	IPC	Explicit call of a parent's constructor deletion	
	Polymorphism	PNC	<i>new</i> method call with child class type
PMD		Member variable declaration with parent class type	
PPD		Parameter variable declaration with child class type	
PCI		Type cast operator insertion	
PCD		Type cast operator deletion	
PCC		Cast type change	
PRV		Reference assignment with other comparable variable	
OMR		Overloading method contents replace	
OMD		Overloading method deletion	
OAC		Arguments of overloading method call change	
Java-specific features	JTI	<i>this</i> keyword insertion	
	JTD	<i>this</i> keyword deletion	
	JSI	<i>static</i> keyword deletion	
	JSD	<i>static</i> keyword deletion	
	JID	Member variable initialization deletion	
	JDC	Java-supported default constructor creation	
	EOC	Reference assignment and content assignment replacement	
	EOA	Reference comparison and content comparison replacement	
	EAM	Accessor method change	
	EMM	Modifier method change	

create a faulty version of the program intended to test (the third step will be addressed in the Section III-C). The options shown in Section II-B to mutate the code are not available in C++, that is, there is nothing similar to the bytecode and a parser-based approach cannot be used resorting to the reflection mechanism in order to check the state of objects at runtime. In the case of the approach used in *ILMutator* [14], neither that kind of intermediate code exists in C++. However, a similar technique to insert faults into the code can be followed as the compilers use an internal representation in the form of abstract syntax tree (AST). Bearing in mind the complexity of C++, we consider that reusing this representation can be fruitful instead of parsing the high-level source code; the AST seems to be much more comprehensible than the intermediate code used in C#, which limits the application of the operators. Notwithstanding, recompilation will be needed, so we cannot leverage this approach for a more efficient mutation system as in C#.

Regarding the generator, we discern two basic options for the mutant creation. The first is the generation of the mutants so that each one can be run as a standalone program, making a duplicate of the original version but containing the mutated files. The second option supposes creating new files only for those where the mutation was included, using the same build system of the original program. In the former, a copy of every file involved in the project will be needed for each mutant, but we can resort to different techniques to save space in disk. In the latter, the mutation tool has to be aware of the build system to perform the appropriate changes so that the mutated files are compiled, which could be a cumbersome task.

C. Execution of Tests

The goal of mutation testing is to determine how good is a test suite defined for a program. Thus, we need to automate the execution of the mutants against the test suite. The process starts assuring that the test suite is successfully passed by the original program. Then, the mutants are run and a comparison can be performed in order to observe whether the changes introduced in the program have been detected or not.

Nevertheless, the connection between the test suite defined by the tester and the mutation tool can be a concern when it comes to automate the test suite execution: the test suite should be able to provide the results so that the tool can retrieve and process them. On the one hand, there is no a prevailing testing framework for C++ currently, unlike Java for instance, where *JUnit* is broadly used. Hence, we can decide to give support for one or more frameworks, but we advocate the design of a library which can be used whatever the way the test suite was developed. This latter approach allows for the execution of tests consistently as the same methods will be used to report the output. On the other hand, when testing mutants generated from class mutation operators, posing different test scenarios of the usage of objects is needed. Besides, various kinds of test cases are required within the scenarios depending on the need of testing that concrete situation. This aspect is different from other unit tests or languages where some values are provided to the program and simply an output is expected.

IV. HOW TO EVALUATE THE RESULTS

In this section, we discuss several matters to be considered when assessing the results obtained with mutation testing.

A. Quality of Operators

The evaluation of the catalog of mutation operators is a significant step in the application of mutation testing to check if the operators are really effective for the purpose of the technique. The most basic and used calculation for the adequacy level of the test suite is the mutation score in each particular mutation operator. The mutation score indicates the percentage of dead mutants versus the number of non-equivalent mutants. In other researches, further dimensions have been studied. Smith et al. [19] proposed to determine the quality of mutation operators by relating the possible states of the mutants after the test suite execution: killed by the initial test suite, killed by a new test case, killed by a new test case specifically defined to kill another mutant or not killed.

Estero-Botaro et al. [20] defined some terms to evaluate the operators for Web Services Business Process Execution Language (WS-BPEL) 2.0, like “weak mutant”, which is killed by every test case in the test suite, or “resistant mutant”, which is killed by a single test case. Thus, the quality of a mutation operator can be determined analyzing some conditions. In summary, the operator should generate a low number of invalid and equivalent mutants (as they do not help in the mutation analysis), but produce few weak mutants as well as the more resistant mutants the better.

Derezińska [18] exposed an idea of the effectiveness for the class mutation operators (specifically for C#), posing various questions (gathered in Table VI) that should be answered to

TABLE VI. QUESTIONS POSED BY DEREZIŃSKA [18] TO ASSESS THE EFFECTIVENESS OF MUTATION OPERATORS.

Questions	
1	Does an operator can be applied in real programs to simulate faults of programmers?
2	Are any invalid mutants generated by an operator?
3	Does an operator generate many equivalent mutants?
4	Is an operator effective in assessing the quality of given test cases? If a mutant is not killed by a given test suite, is it easy to create test cases which kill it?

deem the usefulness of an operator. In addition, this author qualifies a test suite calculating the quotient of the number of test which killed mutants generated by an operator over all test runs performed on these mutants (see Equation 1).

$$Effectiv. = \frac{Killed\ test\ runs}{(Total\ mut. - Equiv.\ mut.) * Total\ tests} * 100 \quad (1)$$

These different kinds of measuring the quality of the testing process should be taking into account when applying this technique to C++, depending on the evaluation intended to perform as well. The operators producing a great amount of mutants should be also studied in depth because they can entail much computation.

B. Issues in the Assessment of Results

When evaluating the results obtained from the application of mutation testing, several issues should be considered. Firstly, the different behavior of the programs can alter in a great degree the results produced with an operator. For example, an operator can involve the creation of many mutants in some cases whereas there are no mutants in others. Thus, we have to take care of this matter in order to generalize results.

Secondly, the initial test suite deserves special attention. Especially in the mutants at the class level, the test suite may not cover every class or member that has been injected a fault. Hence, the mutant will not be killed by any test case, but they cannot be considered as equivalent. This fact occurs very often in C++, where a program is formed by different source files and classes, some of them providing secondary functionalities. This issue has been tackled by Segura et al. in [21], defining the term “uncovered mutant” as that mutant whose fault is not exercised by a test suite. These uncovered mutants should not be computed when assessing the results. Likewise, the authors of that paper also notice the possibility of executing some duplicated mutants when the fault is inserted in a class which is reused in more than a single file. Therefore, either the creation of duplicated mutants is handled in order to avoid them in the generation stage, or they are manually omitted in the execution of the tests.

Finally, the identification of equivalent mutants is still an undecided problem, so the alive mutants are visually analyzed to determine which of them are actually equivalent. This is a harsh and tedious task and, in several situations, it is not easy to assure whether a mutant is definitely equivalent or not. This matter can be even more pronounced in the case of C++ because of its features. According to the great size of current

programs leading to a high number of mutants, the analysis of each one of them can be rather costly in time. Because of this, Segura et al. [21] creates a new classification of mutants: “undecided”. Within this group are cataloged those mutants whose study exceeds an established threshold of time without reaching a final conclusion. Anyway, the implementation of a technique to reduce the time execution or the number of mutants should be considered.

V. CONCLUSION

This paper has supplied a comprehensive survey of the state of mutation testing with regard to the C++ programming language, which has not been almost tackled in the literature and neither in the practice so far. The paper provides useful information about the mutation operators and various research fields where further investigation should be accomplished before conceiving the creation of a mutation tool for this language: the implementation technique to inject the syntactic faults in the program, the execution of tests and the evaluation of the empirical results.

With respect to all these matters, we can summarize that a set of operators can be composed at different levels of the language, being possible to adapt the operators from other similar languages. The class level regarding OO features is probably the area most interesting; this paradigm is widely used in C++ and further research is needed at this level to obtain more concluding results around class mutation operators. Likewise, we consider promising the approach of reusing the abstract syntax tree generated by a compiler because it ensures a complete coverage of the features of the language.

In the future, we intend to start out the development of a mutation tool, first composing a set of class mutation operators. In addition, we aim to study the possibilities of the abstract syntax tree to found the mutation locations as well as appropriately mutate the code. Once the mutants can be created, a complete and automated tool can be constructed to perform some experiments to evaluate the operator quality and the usefulness of the technique in C++. Another salient issue is the high number of mutants that can be produced, so the usage of a cost reduction technique may direct the future research.

VI. ACKNOWLEDGMENTS

This paper was partially funded by the research scholarship PU-EPIF-FPI-PPI-BC 2012-037 of the University of Cádiz and by the MoDSOA research project (TIN2011-27242) under the National Program for Research, Development and Innovation of the Ministry of Science and Innovation (Spain).

REFERENCES

- [1] M. R. Woodward, “Mutation testing - its origin and evolution,” *Information and Software Technology*, vol. 35, no. 3, Mar. 1993, pp. 163–169. [Online]. Available: [http://dx.doi.org/10.1016/0950-5849\(93\)90053-6](http://dx.doi.org/10.1016/0950-5849(93)90053-6)
- [2] Y. Jia and M. Harman, “An analysis and survey of the development of mutation testing,” *Software Engineering, IEEE Transactions on*, vol. 37, no. 5, Oct. 2011, pp. 649–678.
- [3] Z. Ahmed, M. Zahoor, and I. Younas, “Mutation operators for object-oriented systems: A survey,” in *Computer and Automation Engineering (ICCAE)*, 2010 The 2nd International Conference on, vol. 2, feb. 2010, pp. 614–618.

- [4] H. Zhang, “Mutation operators for C++,” retrieved: April, 2014. [Online]. Available: http://people.cis.ksu.edu/~hzh8888/mse_project/
- [5] A. Derezińska, “Object-oriented mutation to assess the quality of tests,” in *Proceedings of the 29th Conference on EUROMICRO*. Belek, Turkey: IEEE Computer Society, 2003, pp. 417–420.
- [6] P. R. Mateo, M. P. Usaola, and J. Offutt, “Mutation at the multi-class and system levels,” *Science of Computer Programming*, vol. 78, no. 4, 2013, pp. 364–387, special section on Mutation Testing and Analysis (Mutation 2010) and Special section on the Programming Languages track at the 25th ACM Symposium on Applied Computing.
- [7] J. Boubeta-Puig, A. García-Domínguez, and I. Medina-Bulo, “Analogies and differences between mutation operators for WS-BPEL 2.0 and other languages,” in *Proceedings of the 2011 IEEE Fourth International Conference on Software Testing, Verification and Validation Workshops (ICSTW)*, IEEE, Berlin, Germany: IEEE, 2011, p. 398–407, print ISBN: 978-1-4577-0019-4. [Online]. Available: <http://dx.doi.org/10.1109/ICSTW.2011.52>
- [8] K. N. King and A. J. Offutt, “A FORTRAN language system for mutation-based software testing,” *Software - Practice and Experience*, vol. 21, no. 7, 1991, pp. 685–718.
- [9] Y. S. Ma, J. Offutt, and Y. Kwon, “MuJava: an automated class mutation system,” *Software Testing, Verification and Reliability*, vol. 15, no. 2, 2005, pp. 97–133.
- [10] Y. Jia and M. Harman, “MILU: a customizable, Runtime-Optimized higher order mutation testing tool for the full c language,” in *Practice and Research Techniques*, 2008. TAIC PART '08. Testing: Academic Industrial Conference, Aug. 2008, pp. 94–98.
- [11] “Insure++: C/C++ testing tool, detect elusive runtime memory errors - Parasoft,” retrieved: April, 2014. [Online]. Available: <http://www.parasoft.com/insure>
- [12] “PlexTest ITRegister,” retrieved: April, 2014. [Online]. Available: <http://www.itregister.com.au/products/plextest>
- [13] M. Hampton and S. Petithomme, “Leveraging a commercial mutation analysis tool for research,” in *Testing: Academic and Industrial Conference Practice and Research Techniques - MUTATION*, 2007. TAICPART-MUTATION 2007, Sep. 2007, pp. 203–209.
- [14] A. Derezińska and K. Kowalski, “Object-oriented mutation applied in common intermediate language programs originated from C#,” in *Software Testing, Verification and Validation Workshops (ICSTW)*, 2011 IEEE Fourth International Conference on, 2011, pp. 342–350.
- [15] S. B. Lippman, J. LaJoie, and B. E. Moo, *C++ Primer, Fifth Edition*, 5th ed. Addison-Wesley, 2013.
- [16] H. Agrawal and et al., “Design of mutant operators for the C programming language,” *Technical Report SERC-TR-41-P*, Software Engineering Research Center, Purdue University, West Lafayette, Indiana, Tech. Rep., Mar. 1989.
- [17] J. Offutt, Y. S. Ma, and Y. R. Kwon, “The class-level mutants of MuJava,” in *Proceedings of the 2006 International Workshop on Automation of Software Test*, K. Anderson, Ed. Shanghai (China): ACM, May 2006, pp. 78–84.
- [18] A. Derezińska, “Quality assessment of mutation operators dedicated for C# programs,” in *Proceedings of VI International Conference on Quality Software*, P. Kellenberger, Ed. Beijing (China): IEEE Computer Society, Oct. 2006, pp. 227–234, ISSN 1550-6002.
- [19] B. Smith and L. Williams, “On guiding the augmentation of an automated test suite via mutation analysis,” *Empirical Software Engineering*, vol. 14, no. 3, 2009, pp. 341–369. [Online]. Available: <http://dx.doi.org/10.1007/s10664-008-9083-7>
- [20] A. Estero-Botaro, F. Palomo-Lozano, and I. Medina-Bulo, “Quantitative evaluation of mutation operators for ws-bpel compositions,” in *Software Testing, Verification, and Validation Workshops (ICSTW)*, 2010 Third International Conference on, 2010, pp. 142–150.
- [21] S. Segura, R. M. Hierons, D. Benavides, and A. Ruiz-Cortés, “Mutation testing on an object-oriented framework: An experience report,” *Information and Software Technology*, vol. 53, no. 10, 2011, pp. 1124–1136, special Section on Mutation Testing. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0950584911000826>

A Negotiation Model Based on Multi-agent System under Cloud Computing

Jian Chen, Xinwei Han, Guorui Jiang

The Economics and Management School,
Beijing University of Technology
Beijing, China

emails: {cjmpw@163.com, jianggr@bjut.edu.cn, hanxinwei@bjut.edu.cn}

Abstract—With the development of cloud computing, the amount of data processing and the ability of information sharing in e-commerce are increasing. Negotiation based on multi-agent is an essential approach to accomplish e-commerce. How to make the negotiation based on multi-agent to adapt to the change brought by cloud computing is an important problem. By considering the degree of market competition pressure, negotiation time, opponent's negotiation historical information under cloud computing, the paper constructs a negotiation model. Finally, the negotiation model's effectiveness is verified by simulation experiment on CloudSim.

Keywords—Negotiation; Multi-agent System; Cloud Computing; Negotiation Framework; Negotiation Model.

I. INTRODUCTION

The emergence of cloud computing represents the arrival of new era of the Internet. Under the circumstance of cloud computing, the methods of requiring information and communication and so on have changed. In cloud computing, all kinds of resources on the Internet could be packaged into service. The packaged resources could supply limitless resource services for requesters [1]. The advantage of cloud computing is that the platform combines enormous resources and could supply variable resources based on actual requirements of users [2]. For suppliers, the process of supplying cloud resource to users is a process of service trade in nature. Negotiation holds an important position in service trade. With the rapid development of Distributed Artificial Intelligence, Multi-agent systems and Autonomy-Oriented Computing, lots of researchers devoted into the research of multi-agent based negotiation [3]. According to the theoretical basis, the multi-agent based negotiation includes negotiations based on game theory [4][5], negotiations based on heuristic [6][7], and negotiations based on argumentation [8][9].

The multi-agent based negotiation has good abilities of distribution and autonomy, it is suitable for the trade under the circumstance of cloud computing especially for cloud resource trade. There have been some researchers who investigated the multi-agent based negotiation under the circumstance of cloud computing. Multi-agent based negotiation under the circumstance of cloud computing has been concerned by researchers [10][11]. There are two markets in cloud computing, i.e., cloud service market between users and service intermediaries, and resource market between service intermediaries and cloud suppliers, and proposed a negotiation mechanism to accomplish dynamic SLA (Service Level Agreement) negotiation in cloud computing [10]. The supply-demand relationship under cloud resource allocation was modeled by game theory [12]. Distributed negotiation mechanism was

proposed for leasing contracts between cloud suppliers and users [13]. Under the market of cloud computing, negotiation would proceed successfully by take good advantage of resource-level information. Then, the negotiation would promote the accomplishment of business targets. Consequently, a negotiation model based on non-addition utility function is proposed to promote the business trade under cloud computing [14]. Service level agreement should be established to resolve the conflicts of participates' different preferences for cloud service. A multi-issue negotiation mechanism is established to resolve the multi-issue negotiation for price, time and service quality under cloud computing. Moreover, corresponding negotiation agreement is established [15]. During the process of cloud resource allocation, the resource suppliers and users are all self-interest agents. The amount of suppliers' resource and the requirement of customers are changing consistently. Facing with these problems, a distribute negotiation mechanism is proposed. While using this mechanism, supplier agent and customer agent could negotiate according to contract price and penal sum. The agents could adapt to the changing environment. Then, the negotiation's accomplishment will promote the cloud resource allocation [16]. The cloud resource suppliers provide large amount of cloud resource to customers according to customers' requirements on IaaS layer. A negotiation mechanism of decision making for cloud resource allocation is proposed by extending the current appointment arithmetic [17].

However, current multi-agent based negotiation under the circumstance of cloud computing mainly used the existing multi-agent based negotiation theory, and aiming at maximizing the economic benefits for users and suppliers. They ignored the influencing factors during the negotiation process, such as degree of competition, time of negotiation, historical information of trade and so on. Moreover, current research mainly used static negotiation process, which may cause the waste of resource and may be lack of interaction between cloud resource suppliers and users.

Based on the above analysis, the article designs a multi-agent based negotiation model under the circumstance of cloud computing. Firstly, the negotiation framework under the circumstances of cloud computing is constructed. Intermediary agent is added to the framework to filter the resource. Secondly, considering the degree of market competition pressure, negotiation time, opponent's negotiation historical information during the negotiation, the multi-agent based negotiation model under cloud computing is established. Thirdly, the negotiation model's effectiveness is verified by simulation experiment. Finally, we summarize our work and propose our future work.

II. MULTI-AGENT BASED NEGOTIATION FRAMEWORK AND WORKFLOW UNDER THE CIRCUMSTANCE OF CLOUD COMPUTING

In this section, we will design a negotiation framework under cloud computing and construct a negotiation workflow correspondingly.

A. Negotiation Framework

The framework of multi-agent based negotiation under the circumstance of cloud computing is designed as showed in Figure 1:

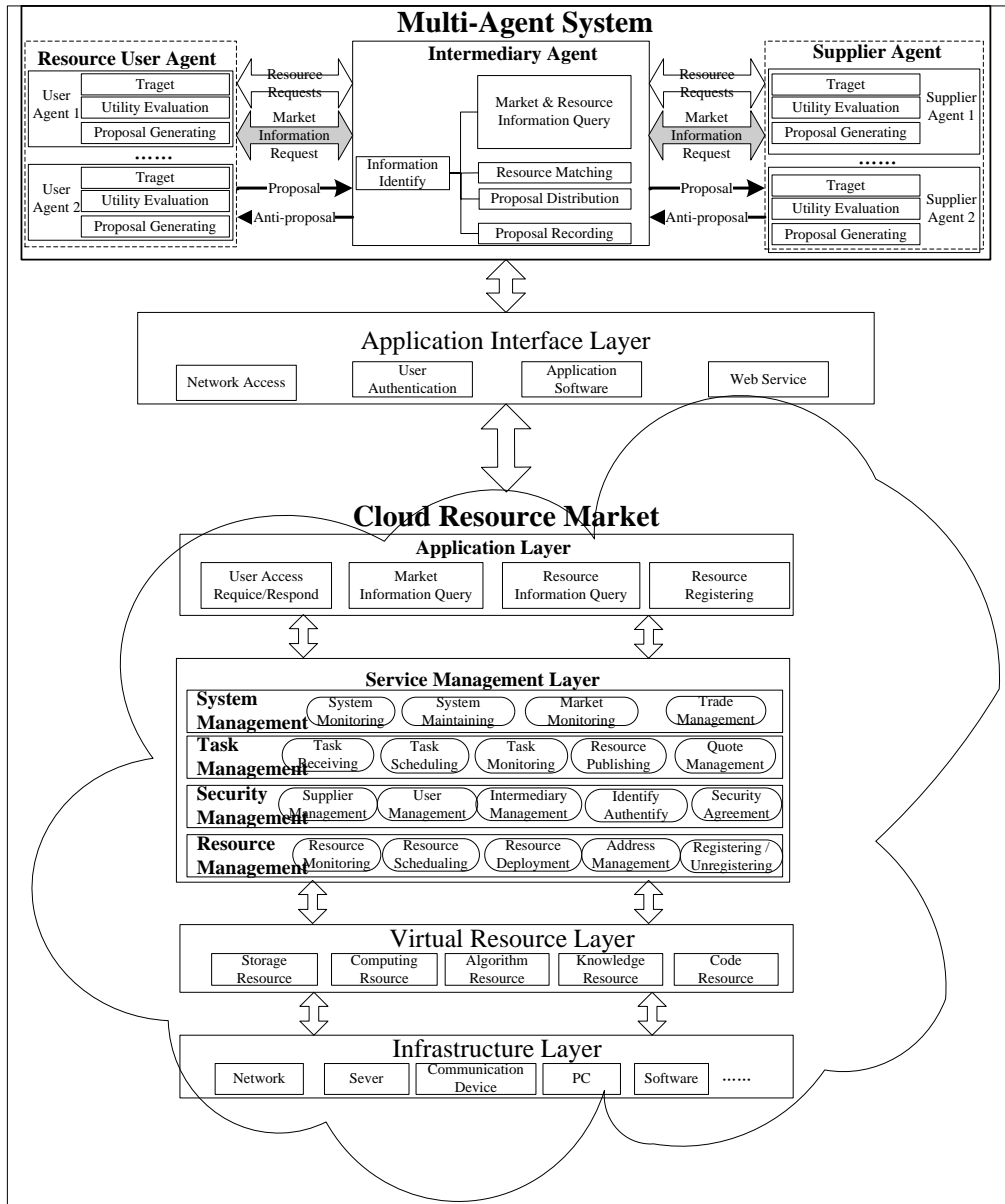


Figure 1. Multi-agent based negotiation framework under the circumstances of cloud computing.

The negotiation framework contains 3 components: multi-agent system, application interface layer and cloud resource market. Multi-agent system is the platform for the service trade, agents represent resource user, resource supplier and intermediary. Application interface connect the multi-agent system to cloud resource market. Cloud resource market contains all the resources used for service trade. The detailed description of the 3 components are introduced as bellow:

1) Multi-agent System

a) Resource Use Agent and Resource Supply Agent

Resource use agents and resource supply agents are the main participants of Multi-agent negotiation, who possess different targets, get information through intermediary agent, and negotiate with opponent agents. Resource supply agents possess the resource in the cloud computing market.

b) Intermediary Agent

In order to improve the efficiency of matching users' requirements to resources, we add intermediary agent to the negotiation framework. The intermediary agent is a third party in cloud computing market which is trusted by participants of cloud resource trade. It interacts with cloud resource market through application interface layer to get market information.

2) Application Interface Layer

The intermediary agent in multi-agent system connects with cloud resource market through application interface layer. The application interface layer supplies web service, user authentication, application software and so on to intermediary agent.

3) Cloud Resource Market

Cloud resources are stored in cloud resource market which mainly contains infrastructure layer, virtual resource layer, service management layer and application layer. The relationship of the layers and the layers' components are shown in Figure 1.

B. Workflow of Multi-agent Based Negotiation

We construct workflow of automated negotiation (shown in Figure 2) based on Figure 1. As intermediary agent is an important part for connection, we also introduce it in this section.

1) Workflow

The workflow of Multi-agent negotiation is shown in Figure 2.

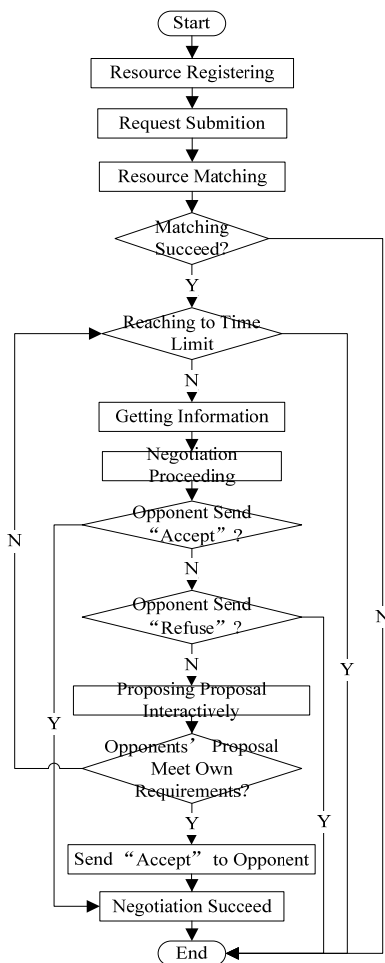


Figure 2. Workflow of multi-agent negotiation under the circumstance of cloud computing.

①Requests submitting:under the circumstance of cloud computing, resource use agents submit their requests to intermediary agent, and intermediary agent gets resource information from cloud resource market through application interface layer.

②Resource Matching:the intermediary agent matches the users' requests to the resource information and sends results to negotiators.

③Getting Information:the matched agents get market information and opponent's proposal history through intermediary agent.

④Negotiation proceeding:participants of negotiation send proposal to opponents. When agent gets the proposal that meets its requirement, negotiation succeed. Otherwise, negotiation proceeds until reaching to the time limit.

If negotiator could not get expected result within the negotiation time limit, the negotiation fails.

2) Function of Intermediary Agent

Under circumstance of cloud computing, in order to improve the efficiency of matching the requests and resource information, we use intermediary agent to match the requests to resource information rapidly. The matching of users' requests to cloud resource information includes three steps:selecting, evaluating and recommendation.

a) To select. The resources use agents submit resource requests to intermediary agent. Intermediary agent acquires service information from cloud resource market and compares the requests with resource information, then selects the resource that match to users' requests.

b) To evaluate. Because there are lots of elements can be evaluated, we only analyze price for the convenience of research. Let F_k represents the resource that matched successfully, k is the serial number of F_k and U represents the utility of F_k 's price.

$$U(F_k) = (P_{c_{max}} - P_{p_{min}}) / P_{c_{max}} \quad (1)$$

$P_{c_{max}}$ is the user's maximum price to accept; $P_{p_{min}}$ is the minimum price that supplier could accept. If $P_{c_{max}} \ll P_{p_{min}}$, it means $P_{c_{max}}$ is far below $P_{p_{min}}$, there is no space to negotiation for user and supplier. If $P_{c_{max}} \rightarrow P_{p_{min}}$, then $U(F_k) \rightarrow 0$, which means $P_{c_{max}}$ is nearly to $P_{p_{min}}$, the space for negotiation is small. If $P_{p_{min}} \rightarrow 0$, then $U(F_k) \rightarrow 1$, which means the gap between $P_{c_{max}}$ and $P_{p_{min}}$ is very large, the negotiation space is large. Intermediary agent should evaluate the value of F_k .

c) To recommend. Resource use agents and supply agents should set γ ($0 < \gamma < 1$), which is the minimum $U(F_k)$. The intermediary agent uses the minimum $U(F_k)$ to select an appropriate F_k : ①If $U(F_k) < \gamma$, relieve the match; ②If all $U(F_k)$ is bigger than γ , relieve the match of minimum $U(F_k)$. Then, the intermediary agent sends the results to the negotiators.

III. NEGOTIATION MODEL BASED ON COMPETITION-TIME-HISTORY

During negotiation under the circumstance of cloud computing, resource use agents and resource supply agents are affected by some influencing factors. The pressure of market competition, negotiation time, opponent's historical information may be the most main factors. We combine the influencing factors under cloud computing circumstance and construct the multi-agent based negotiation model.

A. Formal Description of Negotiation Model

The negotiation model can be described as following:

$$M = \langle A, T, U, m, n, H, S \rangle$$

In the model,

A : Set of Agents, $A = \{\text{Resource Use Agents, Resource Supply Agents, Intermediary Agent}\}$.

T : Negotiation time limit, $T = \langle T_r, T_p \rangle$.

U : Set of Agents' price utility function, $U = \langle U_r, U_p \rangle$. Users' price utility function is

$$U_r = \frac{P_{r_{\max}} - P_{r_t}}{P_{r_{\max}} - P_{r_{\min}}}, \text{ resource suppliers' price utility function is}$$

$$U_p = \frac{P_{p_{\max}} - P_{p_t}}{P_{p_{\max}} - P_{p_{\min}}}. \text{ The negotiators decide whether to accept}$$

opponent's proposal by the utility function.

m_t : No. of competitors in t , gained through intermediary agent. t is the negotiation round.

n_t : No. of opponents in t , gained through intermediary agent. t is the negotiation round.

H : Opponent's historical information of proposals in negotiation gained from intermediary agent. $H = \langle H_{r_j}, H_{p_j} \rangle$. $j \geq 1$, j is the length of proposal history

S : Negotiation strategy agent uses during negotiation.

During each round of negotiations, the elements in the above model would be updated. During the negotiation model, the pressure of market competition is decided by m_t and n_t , the negotiation time is decided by T , the opponent's historical information is decided by H . The specific connotation of three elements is introduced as below.

B. The Three Influencing Factors of Negotiation

1) Pressure of Market Competition

The pressure of market competition should be evaluated in time during negotiation. The evaluation function of competition pressure is defined as:

$$C(m_t, n_t) = 1 - \left(\frac{m_t - 1}{m_t}\right)^{n_t} \quad (2)$$

The agent could get the market information through intermediary agent. Through analysis, we know that if competition pressure $C(m_t, n_t)$ is bigger, the probability of agent being considered as the best opponent is bigger, then the probability of reaching good results is bigger and the agent's competitiveness is bigger. So the environment is advantage for agent and the agent should make smaller concession. Otherwise, the agent should make bigger concession.

When only consider the pressure of market competition, the cloud resource use agent's proposal in the next round is:

$$\begin{aligned} P_{r_{t+1}} &= P_{r_t} + f^{c\text{-user}}(m_t, n_t) \\ &= P_{r_t} + (1 - C(m_t, n_t))(P_{r_{\max}} - P_{r_t}) \end{aligned} \quad (3)$$

P_{r_t} and $P_{r_{t+1}}$ is the agent's proposal at t and $t+1$, $f^{c\text{-user}}$ is the agent's function based on the pressure of market competition. $P_{r_{\max}}$ is the maximum price that resource use agent could accept.

The cloud resource supply agent's proposal in the next round is:

$$\begin{aligned} P_{p_{t+1}} &= P_{p_t} - f^{c\text{-supplier}}(m_t, n_t) \\ &= P_{p_t} - (1 - C(m_t, n_t))(P_{p_t} - P_{p_{\min}}) \end{aligned} \quad (4)$$

P_{p_t} and $P_{p_{t+1}}$ is the agent's proposal at t and $t+1$, $f^{c\text{-supplier}}$ is the cloud resource supply agent's function based on the pressure of market competition. $P_{p_{\min}}$ is the minimum price that resource supply agent could accept.

2) Time

Time limit is usually set during negotiation, negotiators usually take different concession as the time goes. The different concession based on time is summarized and the time constraint equation is proposed [18]:

$$k_t = [1 - (t/T)]^\lambda k_0 \quad (5)$$

k_0 denotes the gap between cloud resource use agent and cloud resource supply agent at initial time. k_t the gap between cloud resource use agent and cloud resource supply agent at time t , $t \leq T$. λ is the nonnegative time factor, it affects degree of concession, and is preset by negotiators and not changes during negotiation.

From (5), k_{t+1} at time $t+1$ is:

$$k_{t+1} = \frac{1 - [(t+1)/T]^\lambda}{1 - (t/T)^\lambda} k_t \quad (6)$$

Suppose the function of time as

$$T(t, t+1, T, \lambda) = \frac{k_{t+1}}{k_t} = \frac{1 - [(t+1)/T]^\lambda}{1 - (t/T)^\lambda} \quad (7)$$

where $T(t, t+1, T, \lambda) < 1$. The bigger $T(t, t+1, T, \lambda)$ would be, the gap between the user's proposal and supplier's proposal would be bigger, the probability of reaching negotiation success would be smaller, in that case, the agent should make bigger concession. When only consider time during negotiation, the cloud resource use agent's proposal in the next round is:

$$\begin{aligned} P_{r_{t+1}} &= P_{r_t} + f^{t\text{-user}}(t) = P_{r_t} + T(t, t+1, T, \lambda)(P_{r_{\max}} - P_{r_t}) \end{aligned} \quad (8)$$

$f^{t\text{-user}}(t)$ is the resource use agent's function based on time.

The cloud resource supply agent's proposal in the next round is:

$$\begin{aligned} P_{p_{t+1}} &= P_{p_t} - f^{t\text{-supplier}}(t) = P_{p_t} - T(t, t+1, T, \lambda)(P_{p_t} - P_{p_{\min}}) \end{aligned} \quad (9)$$

$f^{t\text{-supplier}}(t)$ is the resource supply agent's function based on time.

3) Opponent's historical information

Negotiation opponent's proposals during negotiation have some characteristics. If negotiators could take advantage of these characteristics, the negotiation would be advantage for them.

Suppose H_{a_j} is agent a 's negotiation history that contains the proposal in previous negotiation $\langle P_{a_1}, P_{a_2}, \dots, P_{a_j} \rangle$. P_{a_j} is the proposal that agent a proposed in round j . According to opponent's historical information, we could divide opponent's concession into absolute minimum concession, absolute average concession, absolute maximum concession and relative average concession. All of them can be described as below.

a) Absolute Average Concession

$$\Delta_1^a = (\text{last}(H_{a_j}) - \text{first}(H_{a_j})) / \text{len}(H_{a_j}) \quad (10)$$

where $\text{last}(H_{a_j})$ is the last proposal in H_{a_j} , $\text{first}(H_{a_j})$ is the first proposal in H_{a_j} . $\text{len}(H_{a_j})$ is the length of H_{a_j} .

b) Absolute Minimum Concession

$$\Delta_2^a = \min |P_{a_j} - P_{a_{j-2}}|, 2 < j \leq \text{len}(H_{a_j}) \quad (11)$$

c) Relative Average Concession

$$\Delta_3^a = |P_{a_j} - P_{a_{j-2k}}| / k, 1 \leq j - 2k \leq j \leq \text{len}(H_{a_j}) \quad (12)$$

d) Absolute Maximum Concession

$$\Delta_4^a = \max |P_{a_j} - P_{a_{j-2}}|, 2 < j \leq \text{len}(H_{a_j}) \quad (13)$$

Agent could get opponent's negotiation historical information from intermediary agent. When considering opponent's behaviors, agent's proposal in next round is:

$$P_{a_{j+1}} = P_{a_j} + f^b(H_{b_j}) \quad (14)$$

$f^b(H_{b_j})$ is the function agent a whose opponent is agent b . H_{b_j} is the negotiation history of agent b . The function based on opponent's historical information of could resource use agents and supply agents are $f^{b\text{-user}}(H_{p_j})$ and $f^{b\text{-supplier}}(H_{r_j})$. H_{p_j} is the history of supplier and H_{r_j} is the history of user

$$f^{b\text{-user}}(H_{p_j}) \in \{\Delta_1^p, \Delta_2^p, \Delta_3^p, \Delta_4^p\} \quad (15)$$

$$f^{b\text{-supplier}}(H_{r_j}) \in \{\Delta_1^r, \Delta_2^r, \Delta_3^r, \Delta_4^r\} \quad (16)$$

C. Proposal Generating Based on COMPETITION-TIME-HISTORY

We combine the weighted influencing factors to generate new proposals by summation. Agents in the negotiation could decide each factor's weight by themselves. The weight reflects agent's preferences to each factor. When using the above negotiation model to negotiate, agent could generate new proposal as the following:

a) Colud resource use agent's proposal generating

$$P_{r_{t+1}} = P_{r_t} + \omega_1 f^{c\text{-user}}(m_t, n_t) + \omega_2 f^{t\text{-user}}(t) + \omega_3 f^{b\text{-user}}(H_{r_j}) \quad (17)$$

where P_{r_t} and $P_{r_{t+1}}$ represent the use agent's

proposal at t and $t+1$. ω_i is resource use agent's

preference to the i th influencing factor, and $\sum_{i=1}^3 \omega_i = 1$.

b) Colud resource supply agent's proposal generating

$$P_{p_{t+1}} = P_{p_t} - [\omega_1 f^{c\text{-supplier}}(m_t, n_t) + \omega_2 f^{t\text{-supplier}}(t) + \omega_3 f^{b\text{-supplier}}(H_{p_j})] \quad (18)$$

where P_{p_t} and $P_{p_{t+1}}$ represent the supply agent's

proposal at t and $t+1$. ω_i denotes resource supply

agent's preference to the i th influencing factor, and

$$\sum_{i=1}^3 \omega_i = 1.$$

IV. EXPERIMENT AND RESULTS

CloudSim is a cloud computing simulator developed by research group in the University of Melbourne. The simulator aims at simulating constructing the infrastructure of cloud computing and comparing difference service scheduling and allocation strategies. By this way, CloudSim could control the resources in cloud computing.

A. Targets of Experiment

In cloud computing market, the resource users want to solve their problems by lower cost, while service suppliers want to get more profit by supplying resources. Consequently, the process of cloud resource allocation is a process of service trade in nature. As negotiation holds an important position in service trade, cloud resource allocation is a good field to apply negotiation. Moreover, negotiation could improve the cloud resource allocation's flexibility, interaction and autonomy.

During simulation experiment on CloudSim, we apply the negotiation model proposed in this paper by modifying the class of VmAllocationPolicy in CloudSim. By comparing with the default resource allocation method and negotiation strategy based on time in CloudSim, we could verify the effectiveness of negotiation and the effectiveness of proposed negotiation model.

B. Experimental Parameters Setting

Hardware environment setting: Intel Core 1.86GHz CPU, 2GB RAM, 160G Hard Disk. Software environment setting: operating system is windows XP, development tools are Java 1.7.0, Eclipse 3.2 and CloudSim 3.0.

Environment settings of CloudSim: the number of virtual machine's CPU pick up 1 or 2 randomly. The CPU's capability of processing is 200MIPS-400MIPS. 1G RAM. Network bandwidth is 2M/s-4M/s. Hard disk is 2G-4G.

During experiment, we assume that cloud resource users only request storage resource and virtual nodes only supply storage resource. The experiment will simulate how virtual nodes which are on a same data center deal

with 20 tasks. Each task represents a user's request (that means there are 20 cloud resource users during experiment). Each virtual node represents a cloud resource supplier and there are 100 virtual nodes on a data center during experiment. The data center will use default method and the proposed negotiation model to allocate the cloud resource. We will verify the advantage of proposed negotiation model through comparison.

During experiment, the cloud resource users' expected price will choose from [10,60] randomly and reserved price will choose from [200,250] randomly. Virtual nodes' expected price will choose from [200,250] randomly and

reserve price will choose from [10,60] randomly. The price utility of resource users and suppliers for selecting targets is 0.1. Time strategy is chosen from 1/3, 1.0 and 3.0 randomly. The maximum negotiation round is 20. Other attributes is the default value of CloudSim.

C. Results

The experiment results are shown in Figure 3-5. We will analyze the experiment results from the angle of users and suppliers.

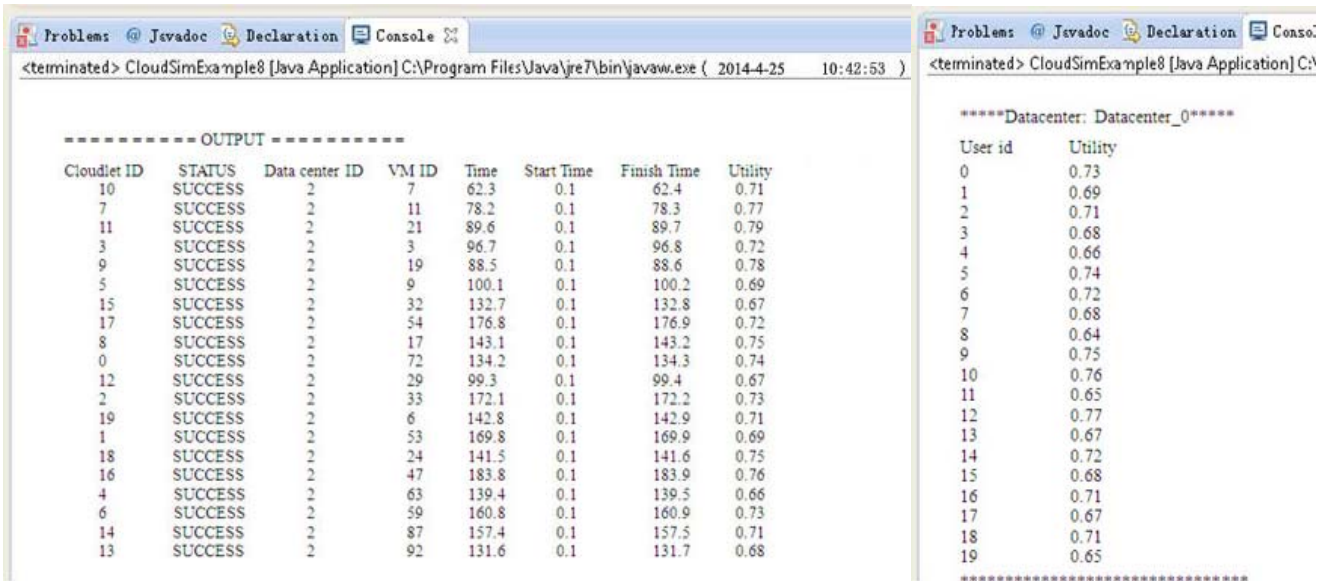


Figure 3. Results of cloud resource allocation using CloudSim default method

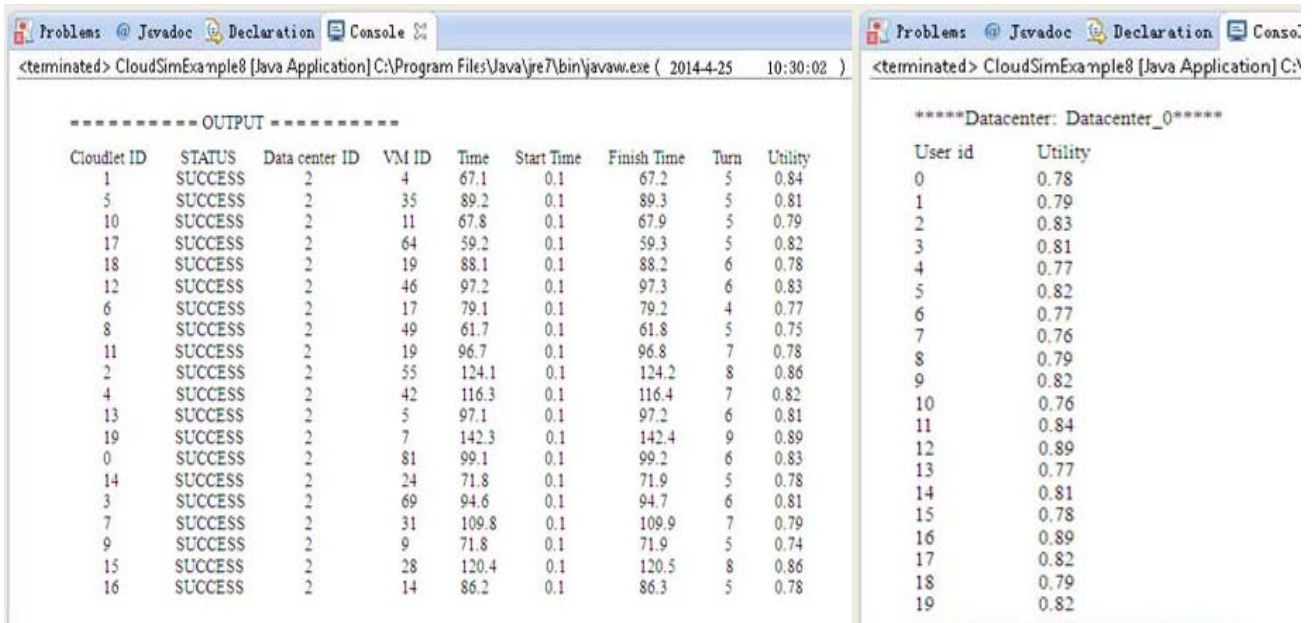


Figure 4. Results of cloud resource allocation using negotiation strategy based on time

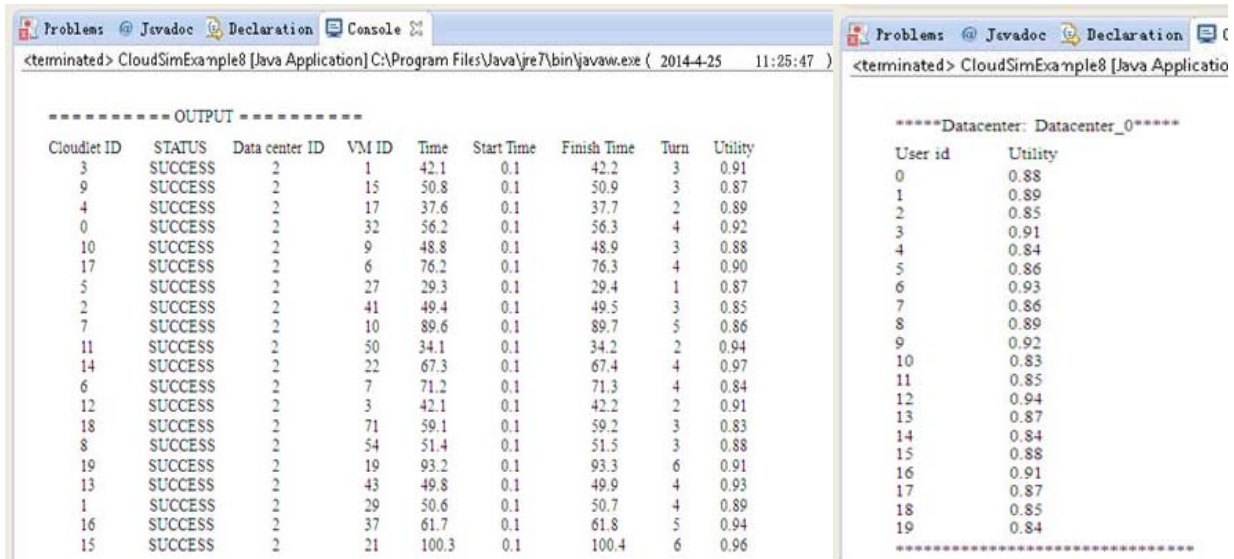


Figure 5. Results of cloud resource allocation using negotiation model based on competition-time-history

When using the CloudSim’s default method to allocate resources, the average negotiation time is 130.035ms, the suppliers’ average utility is 0.7215 and the users’ average utility is 0.6995. When using the negotiation strategy based on time, the average negotiation time is 91.98ms, the suppliers’ average utility is 0.807, the users’ average

utility is 0.8055 and the average negotiation rounds is 6. While CloudSim using the proposed negotiation model to allocate resources, the average allocating time is 58.04ms, the suppliers’ average utility is 0.8975, the users’ average utility is 0.8755 and the average negotiation rounds is 3.55.

TABLE1 RESULTS OF EXPERIMENT

	Default Method of CloudSim	Negotiation Strategy based on Time	Negotiation Model based on COMPETITION-TIME- HISTORY
Average Time/ms	130.035	91.98	58.04
Suppliers’ Average Utility	0.7215	0.807	0.8975
Users’ Average Utility	0.6995	0.8055	0.8755
Average Negotiation Round /turn	/	6	3.55

Through Table 1, we could see that under cloud computing, while CloudSim uses negotiation strategy based on time, cloud resource allocation would use less time, resource users and suppliers will get higher utility. This means negotiation is more effective than CloudSim default method in cloud resource allocation.

While comparing with the negotiation strategy based on time during cloud resource allocation, the proposed negotiation model based on competition-time-history proposed could shorten the resource allocation time and improve the final effectiveness. Consequently, we could hold the view that negotiation is suitable for cloud resource allocation, the proposed negotiation model is better than traditional negotiation methods.

V. CONCLUSION AND FUTURE WORK

The paper designed the multi-agent based negotiation framework under the circumstance of cloud computing. The intermediary agent could shorten the negotiation time and enhance the success rate of negotiation. The multi-agent based negotiation model based on competition-time-history proposed in the paper considers multiple influencing factors during negotiation, and generates reasonable proposal according to current market by combining all the factors. Finally, the negotiation model was applied to the ResourceAllocation of CloudSim and accomplish cloud resource allocation in the simulation experiments. Simulation experiment proved that the

proposed negotiation model could applied to cloud resource well and could get higher effectiveness. With the development of cloud computing and the rapid increase of information technology, negotiation under cloud computing will face more problems such as credit problems, scheduling problems. We will consider how to resolve these problems in the future.

ACKNOWLEDGMENTS

The work in this paper is supported by the National Natural Science Foundation of China (71371018) and the Social Science Fund of Beijing-13JDJGB037.

Corresponding Author: Guorui Jiang, Email: jianggr@bjut.edu.cn

REFERENCES

- [1] X. Zheng, P. Martin, and K. Brohman, Cloud service negotiation:concession vs. tradeoff approaches[C]. Proceedings of the 2012 12th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (ccgrid 2012). IEEE Computer Society, 2012, pp. 515-522.
- [2] K. Chen and W. M. Zheng, Cloud computing:system instances and current research[J]. Journal of Software, 20(5). 2009, pp. 1337-1348.
- [3] M. He, N. R. Jennings, and H. F. Leung, On agent-mediated electronic commerce[J]. Knowledge and Data Engineering, IEEE Transactions on, 15(4). 2003, pp. 985-1003.
- [4] J. S. Rosenschein and Zlotkin, Rules of encounter:design

conventions for automated negotiation among computers [M]. Cambridge MA:MIT Press,1994.

- [5] N. R. Jennings, K. Sycara, and M. Wooldrige, A roadmap of agent research and development [J]. *Autonomous Agents and Multi- Agent Systems*, 1.1998, pp. 275- 306.
- [6] I. Rahwan, S. Liz, and N. R. Jennings, A methodology for designing heuristic agent negotiation strategies [J]. *Applied Artificial Intelligence*, 21(6). 2007, pp. 489-527.
- [7] T. D. Nguyen and N. R. Jennings, A heuristic model of concurrent bilateral negotiations in incomplete information settings [C]. *Proc 18th Int Joint Conf on AI, Acapulco, Mexico, 2003*, pp. 1467- 1469.
- [8] S. D. Ramchurn, C. Sierra, and L. Godo, Negotiating using rewards[J]. *Artificial Intelligence*, 171(10). 2007, pp. 805-837.
- [9] G. Zhang, G. R. Jiang, and T. Y. Huang, Cognition model of argumentation-based multi-Agent business negotiation[J]. *Computer Engineering*, 37(1). 2011, pp. 28-31, 33.
- [10] K. M. Sim, Towards complex negotiation for cloud economy[J]. *Advances in Grid and Pervasive Computing*, 2010, pp. 395-406.
- [11] S. Son and K. M. Sim, A price-and-time-slot-negotiation mechanism for cloud service reservations[J]. *Transactions on Systems, Man, and Cybernetics*, 42(3). 2012, pp. 713-728.
- [12] G. Wei, A. V. Vasilakos, Y. Zheng, and N. Xiong, A game-theoretic method of fair resource allocation for cloud computing services[J]. *The Journal of Supercomputing*, 54(2). 2009, pp. 1-18.
- [13] D. Sun, G. Chang, C. Wang, Y. Xiong, and X. Wang, Efficient nash equilibrium based cloud resource allocation by using a continuous double auction[C]//*Computer Design and Applications (ICDDA), 2010 International Conference on*. IEEE, 1. 2010, pp. 94-99.
- [14] M. Macias and J. Guitart, Using resource-level information into nonadditive negotiation models for cloud market environments[C]//*Network Operations and Management Symposium (NOMS), 2010*, pp. 325-332.
- [15] S. Son, K. M. Sim, A negotiation mechanism that facilitates the price-timeslot-QoS negotiation for establishing SLAs of cloud service reservation[M]//*Networked Digital Technologies*. Springer Berlin Heidelberg, 2011, pp. 432-446.
- [16] B. An, V. Lesser, and D. Irwin, Automated negotiation with decommitment for dynamic resource allocation in cloud computing[C]//*Proceedings of the 9th International Conference on Autonomous Agents and Multiagent Systems: International Foundation for Autonomous Agents and Multiagent Systems*, 1. 2010, pp. 981-988.
- [17] J. Akhiani, S. Chuadhary, and G. Somani, Negotiation for resource allocation in iaas cloud[C]//*Proceedings of the Fourth Annual ACM Bangalore Conference*, 2011.
- [18] P. Faratin, C. Sierra, and N. R. Jennings, Negotiation decision functions for autonomous agents [J]. *Robotics and Autonomous Systems*, 24(3-4). 1998, pp. 159-182.

ICT Utilization in Libyan Universities: A Report on Case Study Research

Ali Bakeer

School of Information Technology
University of Misurata
Misurata, Libya
Ali.m.bakeer@gmail.com

Martin Wynn

School of Computing and Technology
University of Gloucestershire
Cheltenham, UK
MWynn@glos.ac.uk

Abstract - There is a dearth of literature on the use of Information and Communication Technologies (ICTs) in Libyan universities and this paper aims to help address this imbalance by exploring and analyzing how ICTs are used in Libyan universities. Process maps and systems profiling are employed to examine the current and potential uses of ICTs and a new model for assessing ICT utilization in Libyan universities is put forward and applied at the University of Misurata in northern Libya. This innovative approach to assessing ICT deployment in Libya has emerged from initial case study research that will be further developed as other universities in Libya are investigated.

Keywords – *information and communication technologies; ICTs; Libyan universities; process mapping; systems profiling; SCALE model*

I. INTRODUCTION

With the development and application of Information and Communications Technologies (ICTs), organisations have started to accept that it is important for them to embrace new technologies in order to achieve competitive advantage [1]. The term ICT emerged in the new millennium and is generally considered to encompass the recent developments in communications and e-business infrastructure as well as the more traditional core IT functions of an organisation. It thus can also be considered to include an organisation's Information Systems (IS). "The term IT is used interchangeably with IS, or it may even be used as a broader concept that describes a collection of several information systems, users, and management for an entire organisation" [2]. However, "information and communications technology is often used as an extended synonym for Information Technology (IT), but is usually a more general term that stresses the role of unified communications and the integration of telecommunications (telephone lines and wireless signals), computers, middleware as well as necessary software, storage and audio-visual systems, which enable users to create, access, store, transmit, and manipulate information" [3]. It is this broad definition of ICT that is adopted for this study.

The study investigates how ICT utilization can contribute to universities' overall performance. It is one of the few studies focusing on ICT utilization in the Higher Education Institutions sector in the developing World, in spite of the

growing importance of this sector [4]. This research attempts to support the universities in developing more efficient ways of managing core organizational processes and associated information flows, and will also explore the cultural and operational implications of using ICTs more widely and more effectively.

Following this introduction, an overview of some of the existing models relevant to this study is provided, leading to a statement of the three research questions that are subsequently addressed. The qualitative case study approach employed in this research is then briefly discussed, followed by a section which summarises the findings to date, in the main derived from analysis of processes and systems at Misurata University. From these findings, a new conceptual model for analyzing ICT deployment in Libyan universities is presented, and the concluding section discusses how this exploratory research will be advanced through further studies of ICT deployment in the university sector.

II. ORGANISATIONAL AND THEORETICAL FRAMEWORK

In Libya, ICT has had an undoubted impact on university operations [5]. In recent years, a national strategic plan for developing ICT infrastructure in universities has provided new impetus for change in university administration, making its future exploitation a key component of university development plans [6]. Libyan universities still face significant challenges that affect their management and services, which in turn affect their reputation locally and internationally. These barriers include management capability, poor processes and procedures, lack of accountability, lack of technology strategy, poor technology skills, and budgetary constraints [7].

There are eleven public universities in Libya (see Table I), and, in addition, there is one fledgling private university that has yet to receive approval from the Libyan Ministry of Higher Education (LMOHE). For this reason, it is excluded from this study. The LMOHE fund and manage the eleven public universities in Libya, with a common management, financing and regulation system that aims to improve the universities' management and services. In terms of using ICT applications, managers in Libyan universities generally face a number of barriers such as lack of skills to deal with educational problems and they tend to struggle with the management of institutional inefficiencies [8]. It has been suggested [9] that, in order for the Libyan universities to

meet international standards, it will be necessary to upgrade their ICT infrastructure. Barriers to the increased take-up of ICT applications in Libyan institutions include a lack of ICT infrastructure, a lack of qualified personnel and an institutional resistance to change [7]. The root cause of the current situation can be traced back to financial and cultural issues that have severely impeded technology adoption and associated process change.

Top level process mapping and systems profiling have been used in a number of studies in the UK to assess the ICT status of organizations [10][11]. This approach has been applied in a Libyan context in studies of information systems deployment in Libyan banks [13] and Libyan oil companies [14].

III. RESEARCH METHODOLOGY

Research strategy can be defined as the general plan of how the researcher will go about answering the research questions [18]. The philosophical foundation of this study is based on the ontology of subjectivism, while the epistemological position will be interpretivism. The researcher is centrally involved in the phenomena being studied, and is a key player in the process of data collection and analysis to answer the research questions. In terms of the methodological approach, case studies are applicable when research addresses either a descriptive question such as “what is happening or has happened?” or an explanatory question such as “how or why did something happen?” [19].

TABLE I. PUBLIC LIBYAN UNIVERSITIES IN ACADEMIC YEAR 2012/13 [12]

N	The University	Establishment date	Region	Number of Students	University academic staff		
					Full Time		Part Time
					Libyan	Foreign	Libyan
1	Benghazi	1955	Benghazi	84026	1639	240	808
2	Tripoli	1957	Tripoli	83855	2595	120	909
3	Omar Al-Mukhtar	1974	Albayda	33035	822	691	203
4	Sebha	1976	Sebha	15945	685	146	214
5	Misurata	1983	Misurata	16206	578	75	349
6	Al-Zawia	1983	Al-Zawia	35500	658	59	341
7	Al-Mergab	1986	Al-Khums	31030	749	128	661
8	Al-Zaitona	1986	Soq Al-Ahed	10626	768	28	373
9	Sirte	1989	Sirte	10811	264	152	145
10	Al-JabalAl-gharbi	1991	Gherian	17649	600	53	825
11	Al-Asmarya	1999	Zliten	4112	167	35	366
The Total				310845	9525	1727	5194

In addition, since the turn of the millennium, and influenced by the growth of the internet, a number of models have been developed to gauge organizations' use of ICT and/or e-business [15][16]. Most of these e-business models have been designed in Western countries, where the technological and organizational environment is still significantly different to that in Libya. Heeks [17] is one of several authors to point out that systems designed and implemented in the West are often not appropriate for use in the developing world. This ‘design-actuality gap’ arguably applies equally to some of these models, and this research attempts to build a new conceptual framework better geared to developing world organizations.

The following Research Questions (RQs) will be addressed:

RQ1: What is the current level of ICT utilization in Libyan universities?

RQ2: How appropriate are existing models for assessing the current and potential use of ICTs in Libyan universities?

RQ3: Can a new conceptual and operational model be developed for ICT utilization in universities in Libya?

It has also been asserted [18] that if a case study strategy incorporates multiple cases, then the resulting data can provide greater confidence in the research findings. The study population in this research is Libyan public universities, and the cases are likely to either produce similar results (literal replication), or produce contrasting results but for predictable reasons (theoretical replication).

The research approach will be qualitative, which is in accordance with others similar studies [13][20]. Six universities are being considered, selected on the basis of size and geographical location. At each university, a range of activities is being undertaken to gather and analyze the data and information. These activities include assessment of overall university strategy, assessment of information availability, process mapping and systems profiling.

In terms of research philosophy, this study is based on an inductive approach. The time horizon is cross-sectional as data is collected only once. For data collection, the study uses multiple-sources of evidence these include structured questionnaire, open qualitative semi-structured interviews with many different university employees and students,

document analysis, observation and personal experience. The rationale for using multiple-sources is the triangulation of evidence. Triangulation increases the reliability of the data and the process of gathering it [19].

IV. FINDINGS

In this section, the first two RQs are addressed: What is the current level of ICT utilization in Libyan universities? And: How appropriate are existing models for assessing the current and potential use of ICTs in Libyan universities?

Top level process mapping undertaken at the two case study universities studied to date (Misurata and Al-Mergab) suggests that seven main organizational processes can be identified (Fig.1). This process map can act as a framework for assessing how ICTs are supporting these processes, each one of which can be broken down into two or more sub-processes. These processes appear to be fairly standard across the two universities studied to date, and derive from interviews with key staff and assessment of university departmental structures and operations.

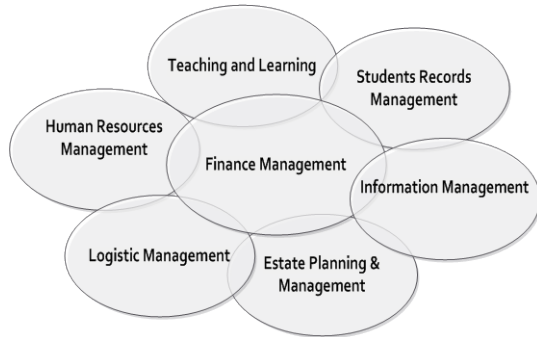


Fig. 1. Top level processes at Misurata and Al-Mergab Universities

The overall picture is of organizations in the early stages of ICT adoption. There is some use of standalone office automation packages (word processors and spreadsheets) and there are some in-house developed systems, particularly for supporting the student records management and financial processes and functions. However, there are only very basic hard-wired networks, although these are currently being developed. Thus, systems are standalone and the use of email is limited. Both universities have a website, but only for providing basic information to the outside world. Knowledge and awareness of the capabilities of ICT is limited to a few individuals and there are but a few ICT training programmes for staff.

Systems profiling, based on interview responses and first hand assessment of systems performance, suggests the majority of information systems now in place are in need of replacement or major upgrade. The problem is particularly acute in the student management process where the three sub-processes are supported by a mix of in-house systems, developed in Visual Basic and the Delphi programming language with the aid of external programmers over the past 12 years, and ad hoc developments in MS Excel and MS WORD.

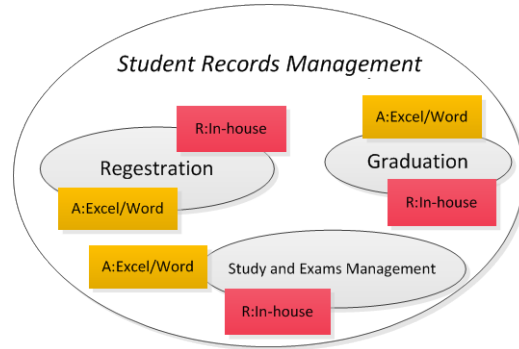


Fig. 2. Systems profiling in the student management process at Misurata University.

R = Red, in need of replacement

A = Amber, can possibly be retained, with upgrade/improvements

There is a basic local area network, which allows multi-user updates of these systems. In the engineering college, end-users have developed an in-house portal which allows students to edit and modify certain designated information regarding their courses and personal details. There is no electronic ID student card for identifying students as members of the university community. The in-house developed systems are not integrated nor interfaced with each other, leading to problems of data inconsistencies and multiple data updates. An assessment of these systems suggests they should be replaced, although the support spreadsheets and word-processed documents could remain (Fig. 2).

Similarly, the human resource management process is supported by an in-house developed system using the Visual Basic and Delphi programming languages, with an underpinning SQL database, introduced in 2006. This provided a great opportunity for keeping and reporting staff records, but it is a stand-alone system which is digitally isolated from the University's networks. Data is gathered manually and organized by administrative staff using MS WORD or Excel before entry into the in-house system.

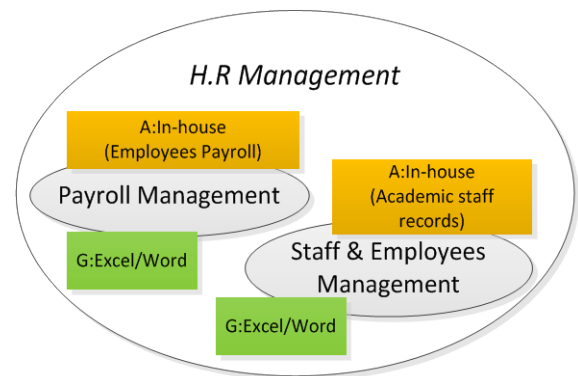


Fig. 3. Systems profiling in the human resources management process at Misurata University

A = Amber, can possibly be retained, with upgrade/improvements

G = Green, strategically sound technology deployment

This is a time consuming process and there remains a high degree of manual and semi-manual processes and a lack of information sharing in the human resource management process (Fig. 3). The financial management process is arguably the most automated, with the majority of activities and procedures supported by systems developed in house (again in Visual Basic and Delphi) and/or in MS WORD and Excel. In house systems are the students scholarship and grants system; the national ID system; staff salary management; and year end budgeting. WORD and Excel systems are used for the preparation of disbursements; procurement; financial audit records; receipt and collection of various university incomes; cash rewards; grants and other activities. The university financial process and sub-processes are illustrated in Fig. 4.

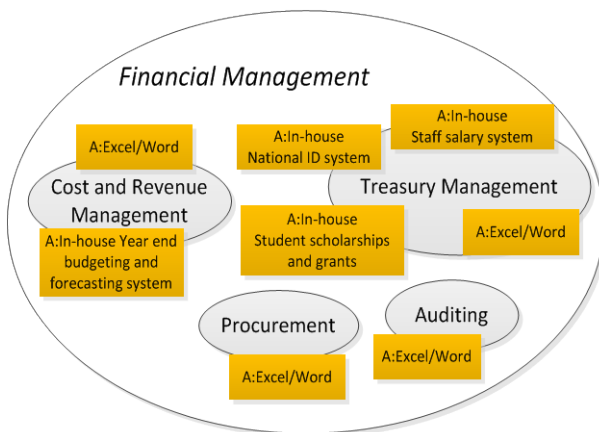


Fig. 4. Systems profiling in the financial management process at Misurata University
A = Amber, can possibly be retained, with upgrade/improvements

ICT provision in the other processes at Misurata University follows a similar pattern, but in general is less well developed. The current ICT provision is a range of bespoke systems, running either in standalone mode or across disjointed networks, supported by ad hoc use of WORD and Excel. Use of email is limited to a few users, and there are no managed server based systems. From a strategic perspective, although WORD and Excel can be retained as strategically sound office automation tools, the majority of these systems are likely to be in need of replacement in the mid-term (3-5 year horizon).

Let us now discuss how appropriate existing models are for assessing the current and potential use of ICTs in Libyan universities (RQ2). Other studies [13][14] have applied Nolan's model [21] to assess the overall status of the IT function in Libyan organizations. An assessment of the two universities studied to date indicates both universities are at or beyond the Contagion stage, but have not yet reached the Control stage, when a centralized IT department is established to plan, manage and control ICT implementation. There are however, some support functions within some of the process areas, notably the financial management process

at Misurata where systems support is combined with other job roles. The relatively limited impact of ICTs at process level and overall in the universities suggests that several of the models developed and applied in a developed world context may be inappropriate to gauge the current status of these universities. However, some of these models can be used as the basis for developing a new conceptual model with a better fit to the organizational and technology adoption situation in Libyan universities.

V. TOWARDS A NEW CONCEPTUAL MODEL FOR ICT UTILISATION IN LIBYAN UNIVERSITIES

This section addresses RQ3: Can a new conceptual and operational model be developed for ICT utilization in universities in Libya? Research to date suggests a stage model which scales down the levels of organizational development and ICT deployment would be of value in assessing and comparing universities' ICT utilization in Libya. A SCALE (start-connect-access-leverage-enterprise) model (Table II) is put forward as a useful provisional framework for further analysis. This model allows better differentiation of the use of ICTs at process level, compared with the application of other models, such as the CPIT model [11][16], which would show most processes registering relatively little progress.

Applying this model to Misurata University confirms that the three processes discussed above are the most advanced, with financial management being at the 'leverage' stage with the majority of staff working in this area using one or other of the deployed systems (Fig. 5). Learning and teaching, arguably the core university process where ICTs have the potential to radically improve the student learning experience, has some evidence of ICT utilisation in the classroom - computers, data shows, and related technologies. However, there is still no use of the internet (there is no connectivity) and no online course materials or lectures that can be accessed 24 hours a day, 7 days a week. In terms of the library resources, Misurata University is still in the early stage of ICT utilization - there is a library website, but access to learning resources is still reliant on printed books and other materials in the university site libraries. There is still a lack of electronic information on curriculum content and teaching materials that could be published and accessible via the Web, and there are neither e-learning applications nor educational portals in the University.

VI. CONCLUDING REMARKS

At the two Libyan universities studied so far, the utilization of ICTs remains limited and uncoordinated. There is a lack of electronic communication channels, such as email and web based media, networks are fragmented and systems are rudimentary and largely standalone. ICT knowledge and awareness is patchy amongst staff, and there has been a general lack of any overarching ICT strategy,

TABLE II. SCALE MODEL – STAGE CHARACTERISTICS

Stage	Characteristics
Start	A clear organisational structure, with clearly defined roles and procedures. A general awareness of the availability of ICTs for day to day running of an organisation One or two individuals using standalone technologies (e.g., a laptop or mobile phone/i-pad).
Connect	A sound electricity network in the majority of the main buildings. A partial in-house data communications network (wired or wireless), facilitating multi-point access. Connectivity with the Internet allowing ad hoc access to websites and inter organisation email exchange. A few standalone users of basic office systems (e.g., a word processor or spreadsheet).
Access	A wider take-up of office systems and the first use of transaction processing systems (often built in-house by end-users). The first servers appear allowing access to systems and applications from the organisation’s network. Wider use of the internet and email. Problems with systems integration, data transfer, back-up procedures and version control. The organization sets up its own email/web address and may have a basic website An awareness of the need for formal ICT support and cross-organisation standards and policies.
Leverage	An ICT support function is established, at either departmental level, or centralized within the organisation as a whole. Policies and/or standards for products and services are introduced, and access, back up and version control procedures are introduced. A range of multi-user systems accessing centrally held databases on a number of servers. Use of cloud computing is considered and may be pursued. The need for systems integration and/or interfacing is recognized and partially addressed. The organisation develops its website, and network capability is extended to allow wider access to systems from both within and outside of the organization.
Enterprise	ICTs and IS are in place to support all main processes/sub-processes, either using a range of integrated packages or in-house developments, or possibly via an off the shelf integrated system. Information availability plays a key role in supporting the organisation’s objectives and improving process operations; the end-user community includes information specialists. Data and information exchange within the organization and with business partners is the norm; external access to systems via the website is supported and controlled. Processes are improved and streamlined as new ICTs and IS are implemented and return on investment is monitored and controlled at organisational level.

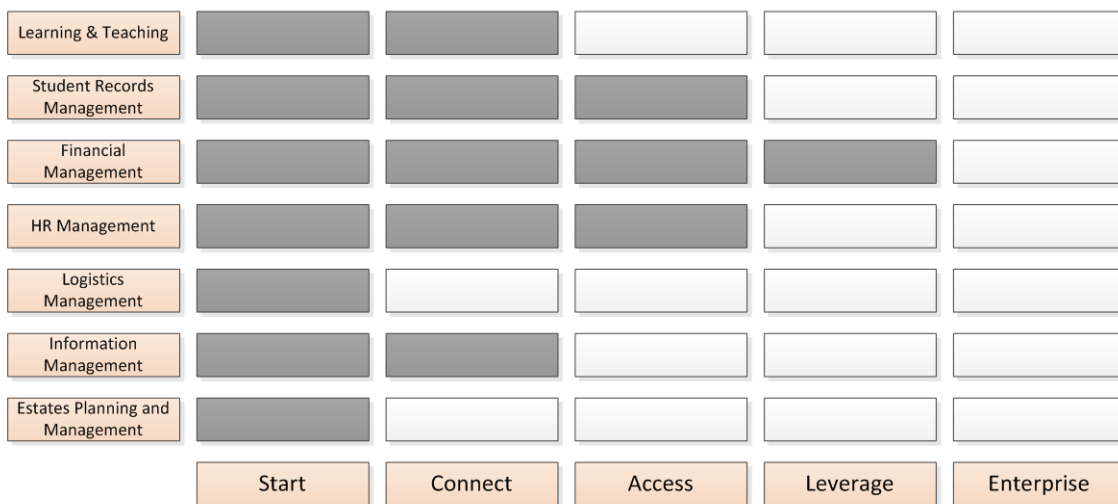


Fig. 5. SCALE model applied to Misurata University processes
The darker shading indicates the level or stage of ICT utilization in each process area.

with ICT spending and policies thus being *ad hoc*, short-term and limited in perspective.

As noted at the outset, there is a general lack of literature and particularly a lack of evaluation studies in developing countries. “Until very recently, the entire literature on IS and developing countries would struggle to fill a single bookshelf... those who have the will to evaluate, such as academics, often lack the resources and capacity” [22]. This research has developed an innovative framework for assessing ICT utilization in Libyan universities, building upon existing models designed for, and applied to, western technology environments. Research case studies in developing countries “help to elaborate such models and, in so doing, contribute to underlying theory... cases from developing countries therefore provide fertile ground for helping understand the complex interplay of action and context that underlies all organizational change” [17]. This research will further develop the framework in a number of ways. First, the existing model will be applied and refined in several other Libyan universities. Second, the other dimensions will be explored and incorporated into the model. This will include an evaluation and classification of the barriers that universities have to overcome to progress their utilization of ICTs, and it will also look at how processes and people skills must change, if these universities are to better exploit the benefits of modern information and communications technologies.

REFERENCES

- [1] C. Yang and S. Lee, “Entry barrier’s difference between ICT and non-ICT industries”, *Industrial Management & Data Systems*, 113, (3), 2012, pp. 461-480.
- [2] E.Turban, E.McLean, and J.Wetherbe, *Information Technology for Management - Improving Quality and Productivity*. New York: Wiley, 1996.
- [3] Free Online Dictionary of Computing (FOLDOC), cited in Wikipedia.org
Available: http://en.wikipedia.org/wiki/Information_and_communications_technology. Retrieved: February, 2014.
- [4] O. Espinoza and L. Gonzalez “Accreditation in higher education in Chile”, *Quality Assurance in Education*, 21, (1), 2012, pp. 20-38.
- [5] A. Elzawi and J. Underwood, “How Higher Engineering Researchers in Libya Perceive the Use of Internet Technology”. In: *Proceedings of The International Arab Conf. on Information Technology*, University of Garyounis, Benghazi, Libya, December 2010, pp. 89-98.
- [6] ITU, *ICT adoption and prospects in the Arab region*, Geneva: ITU Press office, 2012.
- [7] A. Kumar and M. Arteimi, “Potential Opportunities, Barriers and Enablers to Use E-Learning within Libyan Medical Educational Institutions”, *The New York Times*, 20th March, 2009.
- [8] A. Rhema and I. Miliszewska, “Towards E-Learning in Higher Education in Libya”, *Issues in Informing Science and Information Technology*, 7, 2010, pp. 423-437.
- [9] Monitor Group, *National Economic Strategy*, 9th February 2006. United Kingdom: pp. 140-147.
- [10] M. Wynn and O. Olubanjo, “Demand-supply chain management: systems implications in an SME packaging business in the UK”, *International Journal of Manufacturing Research*, vol. 7, no. 2, 2012, pp. 198-212.
- [11] M. Wynn, P. Turner, and E. Lau, “E-business and process change in the UK SME sector”, *Journal of Small Business and Enterprise Development*, vol. 20, iss. 4, 2013, pp. 913-933.
- [12] Libyan Ministry of Higher Education (LMOHE), Libyan Universities. Available at: http://highereducation.gov.ly/?page_id=108 [Accessed: November, 2012].
- [13] O. Sharkasi and M. Wynn, “Deployment evaluation of accounting information systems in Libyan commercial banks”, *The African Journal of Information Systems*, vol. 3, iss. 3, 2011, pp. 87-106.
- [14] H. Akeel, M. Wynn, and S. Zhang, “Information systems deployment in Libyan oil companies: two case studies”, *Electronic Journal of Information Systems in Developing Countries*, (EJISDC), vol. 59, iss. 4, 2013, pp. 1-18.
- [15] Department of Trade and Industry (DTI), *Business in the information age: International Benchmarking Study 2001*, London: Booz Allen Hamilton, 2001.
- [16] M. Taylor and A. Murphy, “SMEs and e-business”, *Journal of Small Business and Enterprise Development*, 11, (3), 2004, pp. 280-289.
- [17] R. Heeks, “Information Systems and Developing Countries: Failure, Success, and Local Improvisations,” *The Information Society*, 18, 2002, pp.101–112.
- [18] M. Saunders, P. Lewis, and A. Thornhill, *Research methods for business students*, 5th edn., 2009, England: Pearson Education Limited.
- [19] R. K. Yin, *Applications of Case Study Research*. 3rd edn, 2012, London: SAGE Publications, Inc.
- [20] H. Al-Mobaideen, “ICT diffusion in Jordanian universities”. In: *Proceedings of the European and Mediterranean Conference on Information Systems*, Izmir, Turkey, July 2009, pp. 13-14.
- [21] R. L. Nolan, “Managing the crisis in data processing”, *Harvard Business Review*, 57 (2), Mar-April, 1979, pp. 115-126.
- [22] R. Gomez and S. Pather, “ICT Evaluation: are we asking the right questions”, *Electronic Journal of Information Systems in Developing Countries* (EJISDC), 50, 5, 2012, pp.1-14.

Mutation Testing: Guideline and Mutation Operator Classification

Lorena Gutiérrez-Madroñal

Juan José Domínguez-Jiménez

Inmaculada Medina-Bulo

UCASE research group
University of Cádiz
Spain, Cádiz

UCASE research group
University of Cádiz
Spain, Cádiz

UCASE research group
University of Cádiz
Spain, Cádiz

Email: lorena.gutierrez@uca.es Email: juanjose.dominguez@uca.es Email: inmaculada.medina@uca.es

Abstract—Mutation testing has been found to be effective to assess test suites quality and also to generate new test cases. In fact, it has been applied to many languages. Unfortunately there is no research work which focuses its attention on detect deficiencies in the mutation testing studies already done. Moreover, the mutation operators classification has not been tackled at all. When has a mutation testing study an enough grade of maturity? It has not been found an study which determines the grade of maturity of a mutation testing study. We propose a classification of mutation operators which lets us know if it is necessary to define more mutation operators, and also lets us determine if the mutation operators are good enough to consider mature the mutation testing study. The evaluation results show the benefits of the proposed classification. This classification lets the developer do a good evaluation of the programming language under mutation testing, as well as its defined mutation operators. The mutation testing process described in this paper and the mutation operators classification have been developed and analyzed with real cases. New mutation operators for Event Processing Language have been defined and used as examples to understand the proposed mutation operators classification.

Keywords—Mutation testing; mutation testing guideline, mutation operators classification.

I. INTRODUCTION

With the success of the application of mutation testing to a number of common implementation languages [1]–[8], this paper presents the elaboration of a formal guide which lets the developer follow the right process to do a mutation testing study (MTS) of a new language, as well as a mutation operators classification to evaluate the grade of maturity of the MTS.

Mutation testing [9]–[11] is a fault-based testing technique providing a test criterion: *the mutation score*, which can be used to measure the effectiveness of a test suite in terms of its ability to detect faults. Mutants will be generated from applying *mutation operators* to the program under test. These mutation operators introduce slight syntactical changes into the program that should be detected by a high-quality test suite.

The majority of the mutant generation systems generate all the possible mutants, and normally include a vast array of mutation operators. Each mutation operator generates a large number of mutants which need to be run against the test suite to determine whether they can be told apart from the original program in some of its test cases (that means whether they are *killed* by the test suite or not). The entire process can take a long time for nontrivial programs.

Why are we proposing a guide to do an MTS? Nowadays, programming languages are adapted or created to cover new functionalities, solve problems or for new systems. Moreover mutation-based testing is a well-know and effective testing approach to reveal code-level vulnerabilities [9], [12]–[16] and has been applied for many traditional implementation languages, in which some deficiencies in mutation testing process have been found. The quality and efficiency of the programs, implemented or part-implemented with these program languages, could be proved and improved after a mutation testing study. For these reasons, we present a formal pattern for a MTS.

Each programming language under MTS should have a mutation operators classification. However, no research work which covers or studies mutation operators classifications has been found. Some definitions of the mutation operators are based on another ones because they have a similar behavior, in other cases the definitions are the same because these could be applied to languages with the same nature or the mutation operators definitions are about general changes that can be done whatever the nature of the programming language. In [17] a similar idea is proposed comparing mutation operators for Fortran, C and Java. One of Marthur's conclusions is "A *basic understanding of mutation operators, and experience in using them, helps a tester understand the purpose and strength of each operator. Further, such understanding opens the door for the construction of a completely new class of operators should one need to test an application written in a language for which no mutation operators have ever been designed!*". That means that it is necessary to study the already defined mutation operators to avoid the duplication and define new mutation operators for new programming languages. The proposed classification helps, not only to organize the mutation operators, but also to determine if there exist mutation operators definitions with a similar behavior, and as consequence if it is necessary to define more operators to cover specific characteristics of the programming language.

It is not an easy task to determine the *grade of maturity* of an MTS. The mutation operators definitions are the key to evaluate the maturity of an MTS. A grade of maturity definition for mutation testing studies is proposed. That definition will make use of the classification presented also in this paper.

Following the study by Gutiérrez-Madroñal et al. [5], some new mutation operators for Event Processing Language (EPL) [18] have been defined in this paper. Their definitions

not only are included in the list of EPL mutation operators but also they help to explain the proposed mutant operators classification.

The structure of the rest of the paper is as follows: Section 2 introduces mutation testing and a description of the main steps involved in the process. Section 3 describes the process to follow to do an MTS, its goals and the relation between them, as well as the proposed classification of the mutation operators (explained with real and new mutation operators). In Section 4 are defined the new EPL mutant operators which are classified with the proposed classification. This classification is also applied and analyzed in Section 5 with mutation operators defined in different mutation testing studies. This section also includes the grade of maturity definition for mutation testing studies as well as the characteristic of its definition. And finally, in Section 6, the conclusions and future work are presented.

II. MUTATION TESTING BACKGROUND

Mutation testing [9]–[11] is a fault-based testing technique that introduces simple syntactic changes in the original program by applying mutation operators. Unlike other fault-based strategies that directly inject artificial faults into the program, the mutation method generates syntactic variations, *mutants*, of the original program by applying mutation operators. Each mutation operator represents “typical” programming errors, that the developer could make. For example, a mutation system replaces a relational operator (say $>$; i.e. $a > 26$) in the original program with other operators (such as $<$, $=$, $>=$, $<=$ and $<>$; i.e., $a < 26$), which is intended to represent a wrong instruction typed by the programmer. If a test case is able to distinguish between the original program and the mutant, it is said that this test case *kills* the mutant. On the contrary, if no test case in the test suite is able to distinguish between the mutant and the original program, it is said that the mutant stays *alive*. An *equivalent mutant* always produces the same output as the original program, hence it cannot be told apart from the original program. The general problem of determining if a mutant is equivalent to the original program is undecidable [19]. There are so many factors to take into account in order to find out if a mutant is equivalent. The code length, the mutant itself, and the variables used in the test case are good examples of it.

Mutation testing can be used in two ways. The first use is to measure the quality of a test suite with its *mutation score*, which is defined as the percentage of killed mutants. A formal definition:

$$ms = \frac{Km}{Tm - Em} \quad (1)$$

where Km is the number of killed mutants, Tm is the total number of mutants, and Em is the number of equivalent mutants. Normally the value of Em is not known, which is a problem because it is necessary to manually inspect the mutants to identify those that are equivalent. The second use for mutation testing is to generate new test cases in order to kill the surviving mutants, and thus improve the quality of the initial test suite. In an ideal case the mutation score reaches

100% which indicate that the test suite is adequate to detect all the faults modeled by the mutants.

One of the main drawbacks of mutation testing is the high computational cost involved. Commonly there is a large number of mutation operators that generate a wide numbers of mutants, each of them must be executed against the test suite. Under certain conditions described in an empirical study by Offut et al. [20], the number of mutants has quadratic complexity in program size.

From a theoretical point of view, the mutation testing process is divided in four main steps: analysis of the program, generation of mutants, execution of the mutants against the test suite and outputs analysis (Figure 1). Each process requires the previous results to be developed. And all of them, with the exception of the last one, output analysis, need a mutant system.

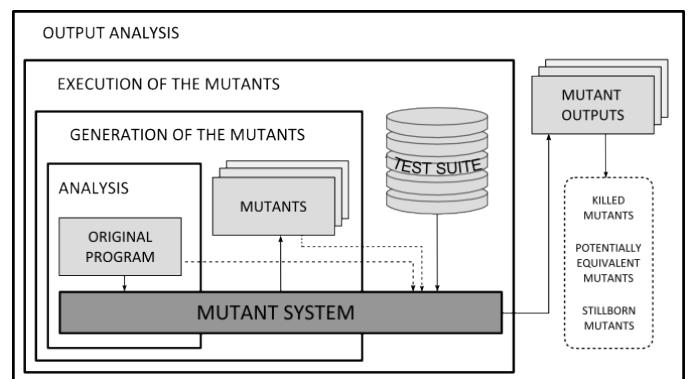


Figure 1. Phases of the Mutation Testing Process

The *analysis* can be defined as the part of the process where the programming language is studied. The analysis result is the number of occurrences of each mutation operator in the *program*. At this point it is necessary to clarify that program is used to denote the software under test, which could be a complete program or some smaller unit, such as a query. With this information, the mutant system moves to the next step, *generation of the mutants*. The mutant system does slight changes in the original program depending on the occurrences of each mutation operator and their definitions. For example, if the program under test were the next line of code “ $a + 6$ ”, the analysis result will say that there is one occurrence of an arithmetic operator. And if the mutation operator definition says that each occurrence has to be changed by one of the following set of elements: $\{-, /, *, \%\}$, the number of generated mutants based on the previous conditions will be four. The *execution of the mutants* is the last task of the mutant system. At this point the mutants are executed against the test suite and just if their outputs can be distinguished between the mutant and the original program, the mutant can be classified as a killed mutant. The stillborn mutants also are automatically classified, because if the mutant fails to be run against the test suite or violates some static constraint defined by the language, the mutant system will notify it. The stillborn mutants fail to execute because of syntax errors or whatsoever conditions that the mutation operators producing them may have introduced. Finally in the *output analysis* the mutants are classified in order to check if the 100% in the mutation score has been reached.

In other case the test suite has to be improved.

III. MUTATION TESTING PROCESS

We propose a guide which describes a formal process to apply an MTS to a new language. Furthermore one of its steps includes the proposed classification for the mutation operators. Figure 2 shows the subtasks involved in each main step of the mutation testing process, our guideline describes them.

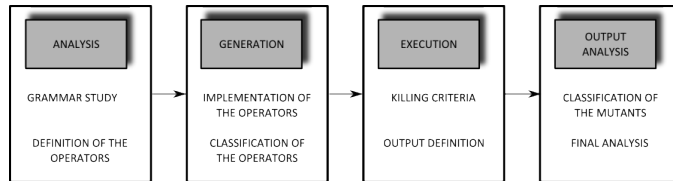


Figure 2. Subtasks of each Mutation Testing Step

A. Language selection

Depending on the programming language, some considerations should be taken into account. If the chosen programming language is based on another programming language, they should be compared (i.e., Event Processing Language (EPL) [18], an SQL-like language and SQL [5]), especially if this base programming language has an MTS. Otherwise, it should be recommendable start the MTS with the base programming language. If there is an MTS for the chosen programming language new contributions can be significant, i.e., some contributions may be done if the new MTS was focused on the changes of the latest version of the programming language. If the chosen programming language is not involved in any of the previous situations, it has to be in consideration its availability as well as its popularity. This means that it is a common programming language and it is easy to find open programs for a further study. Its availability and its popularity are important, because both will influence the MTS, i.e., it is a private programming language which is used by a private company its grammar could be not accessible.

B. Grammar study

The programming language has to suffer an exhaustive survey. Each code line has to be considered and covered with all the possibilities. This part of the progress is linked with the next one because based on the grammar, the mutation operator definitions will be made. Studying the grammar let us know the right changes (mutations) that could be done according to the context. For example in a SQL query:

```
update table_A set field_A =
field_B * (select count (*) from table_B)
```

The first * correspond to an arithmetic operator, which could be replace by another +, -, / and %. On the contrary the last one is a wildcard that can not be replace by any of the last arithmetic operators. If this kind of situations were not considered, the resulting mutation programs would be syntactically wrong as well as the mutations would not represent common errors that a developer could make.

C. Definition of the operators

Before doing a formal definition, it has to be checked if there is a definition in other language that could be directly applied, or applied with a slight modification. After studying the grammar, the mutation operators can be defined. The operators have to represent typical programming errors such as: change the variable name for another, type a wrong (logical, arithmetic) operator, add or subtract one unit to a number, date or time, duplicate or forget a part of code and so on. The mutation operator definition has to explain how the mutation has to be done, and also each special situation that could happen. Sometimes, it is preferable to introduce an example for a visual understanding.

If the mutation operator has a similar behavior as another one that is already defined, it is preferable to use the same name but with a slight modification. For example, the SQL mutation operator "JOI, JOIN Clause", which definition says [21]:

JOI; JOIN Clause - Each occurrence of a join-type keyword (INNER JOIN, LEFT OUTER JOIN, RIGHT OUTER JOIN, FULL OUTER JOIN, CROSS JOIN) is replaced by each of the others. When a join-type is replaced by CROSS JOIN, the search-conditions under the ON keyword are removed. When CROSS JOIN is replaced by another join-type, an ON clause is added and its corresponding join-condition is created based on the primary keys of the joined tables.

This definition is the base for "EJOI" definition, a mutation operator of EPL (Event Processing Language) [5]:

EJOI; JOIN clause - It is applicable for INNER JOIN, LEFT OUTER JOIN, RIGHT OUTER JOIN, and FULL OUTER JOIN. Each one of them is replaced by the others.

The name has been slightly modified because is a new definition but based on an existing one. In this case, the key word CROSS JOIN is not in EPL grammar.

D. Implementation of the operators

Once the operators have been defined, the next step is their implementation in the mutation system. These implementations have two elements. The first one is for detecting the operators in the program, which is included in the analysis step of the mutation testing process. Each mutation operator has a function which localizes each occurrence of the mutation operator in the program. This code is influenced by the studied grammar because it determines if the mutation operator is well located. The analysis output is the number of occurrences of each mutation operator. If in an analysis output the number of occurrences of a mutation operator is zero, this mutation operator is not applied.

The second element in the mutation operator implementation is the generation of the mutants. This code is influenced by the mutation operator definition, which says what have to be done to do the mutation. Every change generates a mutant, so if a definition of a mutation operator says that you can make five possible changes to the original instruction, five mutants will be generated. If the goal is to generate all the possible mutants, there has to be a control in order to generate the correct mutants as well as the exact number of mutants. Its output is a set of generated mutants.

E. Classification of the operators and SoMO classification

This part of the process has to be done after the definitions of mutation operators, but it would be recommendable to do it after their implementation as their behavior could no be the expected one.

The authors usually present their mutation operators grouping them by the nature of the change proposed; for example if the mutation operator changes values in the code, they are group in the category "Value mutation operators", if the changes are replacements, its category is "Replacement mutation operators" and so on. In this paper the approach goes a step further by proposing the SoMO (Sets of Mutation Operators) classification for the mutation operators which divides them in the following sets: *traditional*, *nature* and *specific*. The following lines describe each set and show an example for a better comprehension.

The first set of mutation operators is *traditional*. The mutation operators of this set can be found in the majority of the MTSs of any language. It is easy to discover the following mutation operators types: "Relational mutation operator", "Arithmetic mutation operator", "Logical mutation operator" or any other with a similar name. In some studies this set of mutation operators is called *traditional*, for this reason the SoMO classification uses the same name for this first set.

Other set of mutation operators to consider is the one which covers the mutation operators that do the changes in line/s of code according to the nature of the programming language. If the programming language is an object oriented one, is likely that the object oriented mutation operators will be found, such as "Heritage mutation operators" or "Polymorphism mutation operators" and so on. In the case of a query language, "Query mutation operators" will exist, for example mutation operators which affect *select clauses* or *where conditions*. The SoMO classification includes this type of mutation operators in the *nature* set.

And finally, the last set of mutation operators is the one which involves the mutation operators whose changes are based on the finality of the programming language. These types of mutation operators are common in new programming languages which have been created *recently* and/or have been created based on a *traditional* programming language. For example, the base or origin of query languages is SEQUEL (Structured English Query Language) [22], now known as SQL. The majority of query languages are based on SQL, in fact they are defined as SQL-like languages. Some of them are created for specific tools or functionalities like GQL (Google Query Language) created for developing tools which use the Google App Engine or the Google Cloud Datastore [23], [24], or EPL (Event Processing Language) for dealing with high frequency time-based event data [18], YQL (Yahoo! Query Language) to query, filter, and combine data from different sources across the Internet [25], and so on. If a study of mutation testing of these query programming languages is done, it is likely to find mutation operators which would be included in the first set as well as in the second set of mutation operators previously mentioned. But some of them can not be included in any of them because the changes are done in a specific part of code which can be only found in this programming language. These mutation operators are included

in SoMO classification in *specific* set.

So the different sets of mutation operators (SoMO) are:

- *Traditional mutation operators*: They could be applied to any programming language no matter its nature.
- *Nature mutation operators*: They just could be applied according to the nature of the programming language. These operators are defined according to syntactic-like language faults.
- *Specific mutation operators*: They can not be applied to any other programming language. These operators are defined according to non syntactic-like language faults. That means that the mutations are done in the part of the language which differs from the rest of the languages.

What is the main contribution of the SoMO classification? SoMO classification is focused on the vulgarity or exclusivity of the mutation operators.

Lets define for each set of SoMO classification a GQL mutant operator (for the App Engine of Google). First set, *traditional mutation operators*:

GROR; *Google relational operator replacement* - Each occurrence of one of the relational operator {<, <=, >, >=, =, !=} is replaced by each of the other operators.

Relational operators can be found in several programming languages (their symbols could be different), no matter their nature.

For the second set of SoMO, *nature mutation operators*:

GSEL; *Google SELECT clause* - Each occurrence of one of the SELECT or SELECT DISTINCT keywords is replaced by the other.

This definition is the same that SEL mutation operator definition of SQL [21]. This reaffirms the necessity of SoMO classification, in this case it is not necessary the GSEL definition because there is a programming language with the same nature that already have that definition.

And for the third set of SoMO, *specific mutation operators*:

FRM; *From clause* - The optional FROM clause limits the result set to those entities of the given kind. A query without a FROM clause is called a kindless query and cannot include filters on properties. When this happens, the mutation consists in removing the FROM clause.

This mutation operator can just be applied to GQL because of its grammar and finality, so it is included in the *specific mutation operators* set.

What programming languages can be classified using SoMO classification? All the programming languages which are under an MTS can be classified by SoMO. If it is a base or traditional programming language such as Fortran, COBOL or Pascal, the *specific mutation operators* set will not have mutation operators. On the other hand, if it is a new programming language such as Ruby, Objective C, GQL, the *specific mutation operators* set will have mutation operators.

F. Killing criteria and output definition

This part of the study it is very important. If it is not done, there is just a list of definitions and generated mutants without any goal. In this section is where the conditions to consider when a mutant is killed are explained. There are two steps to develop the killing criteria:

- 1) According to the mutation operator definition, the mutant expected output should be determined.
- 2) Execute the original program and the mutant program and check if the mutant output is the expected one. If the mutant output is different, it is necessary to study the behavior of the mutant to determine the killing criteria.

The criteria to take a mutant apart depend on the output definition. For example, if EPL is the programming language for which the killing criteria wants to be done, it is necessary to take into account its nature. Let us consider the difference of the number of events as the EPL output definition. If the number of events obtained by the original program and the mutants are checked, the mutant outputs which differ in the time of response are equivalent according to the output definition, and this is fault. On the other hand, if the killing criteria is the latency of events between original and mutants, the mutant outputs which differ in the number of events are equivalent. So for EPL, the output definition must cover not only the number of events but also the latency between the original and the mutants.

It has to be defined a killing criteria for each mutation operator. Following with the EPL example, for the mutation operators which affect the time of response, their killing criteria are the latency between original and mutant, and for the ones which affect the amount of events in their output, their killing criteria are the differences in the number of events.

G. Mutant Classification

Once the killing criteria as well as the output definition are completed, it is time to execute all the mutants against a test suite which covers all the mutation operators defined before. Depending on the number of mutation operators defined and the generated mutants, the computational cost will be affected. When the execution is finished, the mutant outputs can be analyzed. Part of this analysis can be automated, in particular to determine if the mutant is killed, alive or stillborn. But it has also a manual part that is focused on the mutants that are still alive. This part consists in to discern whether these mutants are equivalent or it is necessary to improve the test suites. In [26] is described a technique called *constraint-based test data generation* which overcomes partially the equivalent mutant problem.

In the literature there are some definitions which will help the manual analysis in which the equivalent mutants have to be separated:

Equivalent mutant: A mutant is equivalent when there is not a test case which after applying it, the original program output and mutant program output can be differentiated. This means that there is not a single test case can kill the mutant.

Stubborn non-equivalent mutant: A mutant is stubborn non-equivalent when there is no test case in a test suite which

after applying it, the original program output and mutant program output can be differentiated yet. This means that there is not an adequate test suite which detect the mutant, but it is not equivalent.

Equivalent mutants should not be confused with stubborn non-equivalent mutants, the set of equivalent and stubborn non-equivalent mutants is called *potentially equivalent mutants*. So if there are stubborn non-equivalent mutants in the study, the test suite has to be revised and a new test case which solves the stubborn non-equivalent mutants issue must be defined.

The automatic analysis can differentiate if the mutants are alive, killed or invalid. In the set of the killed mutants other definitions about their resistance have to be considered:

Weak mutant: A mutant is weak when is killed by every test case.

Resistant mutant: A mutant is resistant when is killed by a single test case.

Hard to kill mutant: A mutant is hard to kill when it is a resistant mutant, and the test case which kills it, just kills the resistant mutant.

This classification can be done automatically if in a matrix M (a row for each generated mutant R and a column for each test case of the test suite T), the status of the mutants after their execution are saved. The status can be considered as follows: 0 in $m_{i,j}$ means that the mutant in the i -row is alive for the test case in the j -column, 1 is use killed mutants and 2 for invalid mutants.

So it is easy to identify a weak mutant, because the value sum of its row is equal to the number of test cases in the test suite T used in the MTS:

$$\sum_{j=1}^T m_{i,j} = j$$

For a resistant mutant the sum of the elements of its row is equal to 1:

$$\sum_{j=1}^T m_{i,j} = 1$$

And finally a hard to kill mutant can be detected if the sum of the elements of its row and the sum of the elements of the column of the test case which kills it, are both equal to 1:

$$\sum_{j=1}^T m_{i,j} = 1 \quad \text{and} \quad \sum_{i=1}^R m_{i,j} = 1$$

H. Final Analysis

The previous mutant classification helps to discard some mutants that are not important in the MTS such as equivalent mutants and weak mutants. Detecting stubborn non-equivalent is a hard task because it has to be done manually. It is necessary to discover the path to the mutation and check the different modules and code lines which are affected because of the change, so in order to kill the mutant a new test case should be developed. To determine if a mutant is equivalent it is also necessary to do a manual analysis in the code, which has to be so exhaustive as the one for stubborn non-equivalent mutants because many stubborn non-equivalent mutants could be discarded.

IV. ANALYSIS OF REAL STUDY CASES

In [5] a set of mutation operators for EPL is presented, to obtain these results the six first steps of the described guideline were followed. The classification of the operators just covers their categories, so a SoMO classification have to be applied to differentiate the traditional, nature and specific mutation operators.

In this section, we first show a briefly background about EPL, then a discussion of new EPL mutation operators divided in categories is proposed. After their definition, mutation killing criteria for the proposed operators are presented. And finally, a SoMO classification, recovering the EPL mutation operators presented in [5].

A. Event Processing Language

Event Processing Language (EPL) is a SQL like query language. However, unlike SQL that operates on tables, EPL operates on continuous stream of events. As a result, a row from a table in SQL is analogous to an event present in an event stream. Example applications for EPL queries can be found in business process management and automation, finance, network and application monitoring and sensor network systems. These systems require processing of events in real-time.

Despite the fact that EPL is a SQL like language, considerable dissimilarities have been notice between EPL and SQL. Thus, in [5] were developed mutation operators due to specific features of EPL (pattern expression, sliding window of length and time, batch processing). Other point presented in [5] uses mutation testing to reveal vulnerabilities in event processing queries to assess the quality of input event streams. EPL was chosen as our case of study to develop mutation operators and killing criteria to generate high quality event streams and malicious inputs. In order to complete that study as well as to explain and to show the benefits of SoMO classification, new EPL mutation operators are proposed.

B. EPL Mutation Operators

Extension View Set Operators

SWVI; *Sorted Window View increase* - This view (ext:sort) sorts by values returned by the specified expression or list of expressions and keeps only the top (or bottom) events up to the given size. This view retains all events in the stream that fall into the sort range. The mutation consists in increasing the window size by one, so this will increase in one the specified events in the view. See an example in Table I.

SWVD; *Sorted Window View decrease* - The mutation consists in decreasing by one the number of the specified events in the view, so this will decreased by one the window size. See Table I.

SWVR; *Sorted Window View Replacement* - An expression may be followed by the optional *asc* or *desc* keywords to indicate that the values returned by that expression are sorted in ascending or descending sort order. The mutation consists in replacing each specified sort keyword by the other, or remove it. Table I shows an example of SWVR.

RWVI; *Ranked Window View increase* - This view (ext:rank) retains only the most recent among events having the same

TABLE I. SORTED WINDOW VIEW OPERATORS EXAMPLES

Original		SELECT sum(price) FROM StockEvent.ext:sort(10, price desc)
Mutant	SWVI	SELECT sum(price) FROM StockEvent.ext:sort(11, price desc)
	SWVD	SELECT sum(price) FROM StockEvent.ext:sort(9, price desc)
	SWVR	SELECT sum(price) FROM StockEvent.ext:sort(10, price)

value for the criteria expression(s), sorted by sort criteria expressions and keeps only the top events up to the given size. This view is similar to the sorted window in that it keeps only the top (or bottom) events up to the given size, however the view also retains only the most recent among events having the same value(s) for the specified uniqueness expression(s). The mutation increases by one the number of specified events. An example in Table II.

RWVD; *Ranked Window View decrease* - The mutation decreases by one the number of the specified events in the view, so the window size is decreased by one. See Table II.

RWVR; *Ranked Window View Replacement* - The sort criteria expressions may be followed by the optional *asc* or *desc* keywords to indicate that the values returned by that expression are sorted in ascending or descending sort order. The mutation consists in replacing each specified sort keyword by the other, or remove it. See Table II for an example.

TABLE II. RANKED WINDOW VIEW OPERATORS EXAMPLES

Original		SELECT sum(price) FROM StockEvent.ext:rank (symbol, 8, price desc)
Mutant	RWVI	SELECT sum(price) FROM StockEvent.ext:rank (symbol, 9, price desc)
	RWVD	SELECT sum(price) FROM StockEvent.ext:rank (symbol, 7, price desc)
	RWVR	SELECT sum(price) FROM StockEvent.ext:rank (symbol, 8, price asc)

TOVI; *Time-Order View increase* - This view (ext:time_order) orders events that arrive out-of-order, using timestamp-values provided by an expression, and by comparing that timestamp value to engine system time. The mutation increases the specified timestamp by one second.

TOVD; *Time-Order View decrease* - The mutation decreases the specified timestamp by one.

Table III shows examples of the previous mutation operators.

TABLE III. TIME-ORDER VIEW OPERATORS EXAMPLES

Original		SELECT * FROM TimeEvent.ext:time_order (arrivalTime, 8 sec)
Mutant	TOVI	SELECT * FROM TimeEvent.ext:time_order (arrivalTime, 9 sec)
	TOVD	SELECT * FROM TimeEvent.ext:time_order (arrivalTime, 7 sec)

Pattern Expression Operators

NRK; *Not in regexp Keyword* - The regexp keyword is a form of pattern matching based on regular expressions. The result of a regexp expression is of type Boolean. If the input value matches the regular expression, the result value is true. Otherwise, the result value is false. The mutation removes or inserts the keyword 'not' before the keyword 'regexp'. In the Table IV there is an example.

Operator replacement

TABLE IV. NRK OPERATOR EXAMPLE

Original	SELECT * FROM PersonEvent WHERE name not regexp '.*Jack.*'
Mutant	SELECT * FROM PersonEvent WHERE name regexp '.*Jack.*'

FRR; Filter Ranges Replacement - Ranges come in the following 4 varieties. The use of round () or square [] bracket dictates whether an endpoint is included or excluded. The low point and the high-point of the range are separated by the colon : character. The mutation consists in changing a round bracket by a square one and viceversa, the possible outputs: {(), (], [), [)}. Table V shows an example of FRR.

TABLE V. FRR OPERATOR EXAMPLE

Original	mypackage.Event(x not in [0:100])
Mutant	mypackage.Event(x not in [0:100])

SQL Injection Attack Operators

LCR; Limit Clause Remove - The limit clause is used to limit the query results to those that fall within a specified range. This operator removes the limit clause in the query. The Table VI has a LCR example.

TABLE VI. LCR OPERATOR EXAMPLE

Original	select age, count(*) from PersonEvent group by age order by count(*) desc limit 10
Mutant	select age, count(*) from PersonEvent group by age order by count(*) desc

OPR; Offset Parameter Remove - The limit clause has an optional offset parameter which specifies the number of rows that should be skipped at the beginning of the result set. This operator removes the offset parameter in the query. Table VII shows an example of the OPR mutation operator.

TABLE VII. OPR OPERATOR EXAMPLE

Original	select age, count(*) from PersonEvent group by age order by count(*) desc limit 10 offset 2
Mutant	select age, count(*) from PersonEvent group by age order by count(*) desc limit 10

Killing Criteria

The Table VIII shows the list of operators and their corresponding killing criteria. The total number of events reported by an original query (O) and a mutated query (M) are compared in order to decide if a mutant is killed or not. A mutant is killed if the number of reported events is not equal to the original query. Some operators (SWVR, RWVR) require defining a different killing criterion; in these cases a mutant is killed if the order of the reported events by the original and the mutated query is different.

C. EPL SoMO classification

The EPL SoMO classification is showed in Table IX (including the new EPL mutation operators and the defined in [5]). The first column shows the sets, the second column shows the mutation operators which are included in each set and the third shows the percentage of operators involved.

TABLE VIII. KILLING CRITERIA

Category	Operator	Killing criteria
EVS	SWVI, SWVD, RWVI, RWVD, TOVI, TOVD	Event count between O and M.
	SWVR, RWVR	Event sort between O and M.
PEP	NRK	Event count between O and M.
SQJ	LCR, OPR	Event count between O and M.

TABLE IX. EPL SoMO CLASSIFICATION

Set	Mutation Op.	% Op.
Traditional	RLOP, OEDIP, OEDDP, LINC, LDEC, TINC, TDEC, TRUN, SQNC, SQUP, EABS, EAOR, EROR, EUOI, ENLF, ENLI, ENLO, TOVI, TOVD, RWVI, RWVD, SWVI, SWVD, NRK, FRR	41'67
Nature	SQRC, SQFD, EAGR, ESEL, EGRUE, EGRUA, EJOI, EORDE, EORDK, EORDS, ESUBRI, ESUBRII, ESUBRIII, ESUBIBII, ESUBIBIII, ESUBIIBI, EBTW, ELKEWC, ELKEWR, ELKECA, ELKECB, ELKEAE, ELKEAB, LCR, OPR	43'3
Specific	RREP, CEOP, ESIRF, ESIRL, RGEP, BATL, BATT, SWVR, RWVR	15

The EPL SoMO classification includes more than 40% of the mutation operators in both the nature and the traditional set and a 15 percent in the specific set.

V. SoMO CLASSIFICATION ANALYSIS FOR OTHER STUDIES

Lets do the SoMO classification with the following mutation testing studies: XML Schema [27], [28], HTML [29], JSP [29], WS-BPEL 2.0 [30], [31], SOAP [32], PHP [33] and XACML [34]. All of them have something in common with markup languages: are markup languages, are like a markup language or are written in markup languages.

Table X shows the percentage of mutation operators classified in the different sets. The first column is the name of the domain in which the study has been done, the second, third and fourth are the percentage of mutation operators in each set respectively and finally the total number of mutation operators.

TABLE X. SoMO CLASSIFICATION FOR OTHER STUDIES

Domain	% Set			Total
	Traditional	Nature	Specific	
XML	17'9	82'1	0	28
HTML	0	100	0	8
JSP	0	0	100	3
WS-BPEL	38'24	26'47	35'3	34
SOAP	25	0	75	7
PHP	86'4	13'6	0	44
XACML	0	36'4	63'6	11

As it was expected the traditional and base markup languages XML and HTML have not any mutation operators in their specific set. The newest programming languages have the highest percentage of mutation operators in the specific set, with the exception of PHP which has the majority of its mutation operators in the traditional set. According to the PHP study it is an ongoing work and more mutation operators will be added over the time, so it seems to be a huge MTS. Another point to highlight on the PHP study is the 13'6% value in its nature set, this percentage is not because the markup language nature, it is because the influence from Java, C, C++, Perl and

Python. PHP MTS needs definitions for mutation operators of the specific set because these are the important and interesting ones, in other case this study would be incomplete given the fact that the mutation operators already defined do not contribute to a PHP MTS.

JSP is a programming language based on a markup language and SOAP is a protocol written using a markup language, but as we can see, they have no mutation operator in their nature set. The reason of the non-existence of mutation operators in the nature set could be because their authors considered them already defined, or maybe, the study was not focused on defining a set of mutation operators. However, in contrast to JSP and SOAP, WS-BPEL and XACML studies have mutation operators in their nature set. The ones which belong to this set are those whose changes are done in a part of the syntax-like markup language (i.e., remove targets, change the targets order).

The WS-BPEL study is the only one which contemplates mutation operators in all the sets, this is because this study is a complete MTS and WS-BPEL is a new programming language. As it was said in section 3.5, if a programming language is new, its specific set will have mutation operators, in other case (it is not a new programming language) its specific set will have no mutation operators.

In the majority of the cases, it is very difficult to determine if all the mutation operators in an MTS have been defined, although the interesting thing to know is if the already defined mutation operators are good enough to consider the MTS mature. The number of defined mutation operators and the SoMO classification let us define the grade of maturity of an MTS. In order to give a grade of maturity definition, some considerations have to be taken into account. Due to WS-BPEL MTS has an exhaustive, huge and tested mutation operators list and MuBPEL [4], a mutation tool for WS-BPEL language which has been used to verify this study, WS-BPEL study can be considered a finished and mature MTS. Considering WS-BPEL a mature MTS and after studying its SoMO classification, a grade of maturity definition of an MTS can be explained:

Grade of maturity of an MTS: The grade of maturity of an MTS is measured with its SoMO sets. When the percentage of its traditional set is about 40%, and the sum of the percentages of its nature and specific sets is about 60%, it can be said that the defined mutation operators are good enough to consider the MTS mature.

Why these percentages? When an MTS is started, it is normal to begin with traditional mutation operators definitions. These mutation operators are easy to identify and they are already considered as common mistakes that a developer could make. On the other hand, nature and specific sets can not be separated because, depending on the similarity of the programming language to the one in which is based on, these percentages can change. If a programming language has a high percentage of mutation operators in the nature set that means that the original programming language and the new one are very similar. But if the new programming language syntax does not allow those changes, it means that the languages are not so similar.

According to the grade of maturity definition for an MTS,

the previous mutation operators definitions for the EPL MTS described in this paper, we can consider it as a mature MTS (see the percentages in Table IX). The EPL MTS has a vast mutation operators list and satisfies the grade of maturity definition of an MTS (About 40% of mutation operators belong to the traditional set and the rest of them belong to the nature and specific sets). In accordance with the grade of maturity definition, the XML study needs mutation operators of the traditional set to be a mature MTS. On the other hand, the PHP study needs mutation operators of the nature and specific sets (as it was previously indicated). The rest of MTSs have a low number of mutation operators definitions so this definition can not be applied to them.

If a comparison is made between the SoMO classifications of WS-BPEL and EPL, as the percentage of the nature set in EPL is higher than in WS-BPEL, it can be stated that EPL programming language is more similar to SQL than WS-BPEL to XML.

VI. CONCLUSION AND FUTURE WORK

Doing an MTS is an exhaustive task and it is necessary to know the steps to follow and their goals. Due to the computational cost which involves the execution of the mutants, it is interesting to have no errors along the mutation testing process which can affect the results. This paper proposes a mutation testing guideline, with real examples in each step (including mutation operators definitions for EPL), and their goals. Moreover, the SoMO classification divides the mutation operators in traditional, nature and specific sets, and it also determines if the MTS has a enough grade of maturity. The percentage of each SoMO set determines not only what kind of programming language is under study, but also if it is an ongoing study. This helps us to determine what set of mutation operators need to be completed.

The definitions of new EPL mutation operators have been used to explain one of the steps of the mutation testing process (mutation operator classification), as well as to clarify how SoMO classification has to be applied. Adding these operators to the EPL MTS, this MTS meets the grade of maturity definition and indicates that it is a mature MTS.

In future, we plan to finalize the EPL MTS in which an EPL tool will be developed. This tool will help to determine if mutation-based testing reveals code-level vulnerabilities not only in traditional implementation languages, but also in event processing queries. This process assesses the quality of input event streams and generates event streams that can reveal implementation bugs in queries.

The presented grade of maturity definition for an MTS needs to be checked with other real mutation testing studies to verify and adjust (in case to be necessary) the proposed percentages for SoMO classification. This definition will be polished applying the SoMO classification in other MTSs which are focused in different programming paradigms, such as query languages or object oriented languages. This will help to accurate the proposed percentages of the grade of maturity definition.

ACKNOWLEDGMENT

This work has been funded by the Ministry of Science and Innovation (Spain) under the National Program for Research, Development and Innovation, Project MoDSOA TIN2011-27242.

REFERENCES

- [1] W. K. Chan, S. C. Cheung, and T. H. Tse, "Fault-based testing of database application programs with conceptual data model," in Proceedings of the Fifth International Conference on Quality Software, ser. QSIK '05. Washington, DC, USA: IEEE Computer Society, 2005, pp. 187–196. [Online]. Available: <http://dx.doi.org/10.1109/QSIK.2005.27>
- [2] M. E. Delamaro and J. C. Maldonado, "Proteum - a tool for the assessment of test adequacy for C programs: User's guide," in PCS'96: Conference on Performance in Computing Systems, 1996, pp. 79–95.
- [3] S. Hussain, "Mutation clustering," Master's thesis, King's College London, 2008.
- [4] A. García-Domínguez, A. Estero-Botaro, J.-J. Domínguez-Jiménez, I. Medina-Bulo, and F. Palomo-Lozano, "Mubpel: una herramienta de mutación firme para ws-bpel 2.0," in Actas de las XVII Jornadas de Ingeniería del Software y Bases de Datos, A. Ruiz and L. Iribarne, Eds., Almería, Spain, 2012. [Online]. Available: <http://sistedes2012.ual.es/sistedes/jisbd>
- [5] L. Gutiérrez-Madroñal, H. Shahriar, M. Zulkernine, J. Dominguez-Jimenez, and I. Medina-Bulo, "Mutation testing of event processing queries," in Software Reliability Engineering (ISSRE), 2012 IEEE 23rd International Symposium on, 2012, pp. 21–30.
- [6] K. N. King and A. J. Offutt, "A FORTRAN language system for mutation-based software testing," *Software - Practice and Experience*, vol. 21, no. 7, 1991, pp. 685–718.
- [7] Y. Ma, J. Offutt, and Y. R. Kwon, "MuJava: An automated class mutation system," *Software Testing, Verification and Reliability*, vol. 15, no. 2, 2005, pp. 97–133.
- [8] J. Tuya, M. Suarez-Cabal, and C. de la Riva, "Sqlmutation: A tool to generate mutants of sql database queries," in Mutation Analysis, 2006. Second Workshop on, 2006, pp. 1–1.
- [9] R. DeMillo, R. Lipton, and F. Sayward, "Hints on test data selection: Help for the practicing programmer," *Computer*, vol. 11, no. 4, 1978, pp. 34–41.
- [10] R. Hamlet, "Testing programs with the aid of a compiler," *Software Engineering, IEEE Transactions on*, vol. SE-3, no. 4, 1977, pp. 279–290.
- [11] M. Woodward, "Mutation testing - its origin and evolution," *Information and Software Technology*, vol. 35, no. 3, 1993, pp. 163 – 169. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/0950584993900536>
- [12] Y. Jia and M. Harman, "Higher order mutation testing," *Information and Software Technology*, vol. 51, no. 10, 2009, pp. 1379–1393.
- [13] J. Domínguez-Jiménez, A. Estero-Botaro, A. García-Domínguez, and I. Medina-Bulo, "Evolutionary mutation testing," *Information and Software Technology*, vol. 53, no. 10, 2011, p. 1108–1123. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S095058491100084X>
- [14] A. Estero-Botaro, J. Domínguez-Jiménez, L. Gutiérrez-Madroñal, and I. Medina-Bulo, "Evaluación de la calidad de los mutantes en la prueba de mutaciones," in Actas de las XVI Jornadas de Ingeniería del Software y Bases de Datos, A Coruña, Spain, 2011.
- [15] M. Polo Usaola and P. Reales Mateo, "Mutation testing cost reduction techniques: A survey," *Software, IEEE*, vol. 27, no. 3, 2010, pp. 80–86.
- [16] J. Andrews, L. Briand, and Y. Labiche, "Is mutation an appropriate tool for testing experiments? [software testing]," in Software Engineering, 2005. ICSE 2005. Proceedings. 27th International Conference on, 2005, pp. 402–411.
- [17] A. P. Mathur, *Foundations of Software Testing*, 1st ed. Addison-Wesley Professional, 2008.
- [18] EsperTech, "Event processing with wsper and nesper," Last access nov, 2013. [Online]. Available: <http://esper.codehaus.org/>
- [19] H. Zhu, P. A. V. Hall, and J. H. R. May, "Software unit test coverage and adequacy," *ACM Comput. Surv.*, vol. 29, no. 4, Dec. 1997, pp. 366–427. [Online]. Available: <http://doi.acm.org/10.1145/267580.267590>
- [20] A. Offutt, G. Rothermel, and C. Zapf, "An experimental evaluation of selective mutation," in Software Engineering, 1993. Proceedings., 15th International Conference on, 1993, pp. 100–107.
- [21] J. Tuya, M. Suarez-Cabal, and C. de la Riva, "Mutating database queries," *Information and Software Technology*, vol. 49, no. 4, 2007, pp. 398 – 417. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0950584906000814>
- [22] D. D. Chamberlin and R. F. Boyce, "Sequel: A structured english query language," in Proceedings of the 1974 ACM SIGFIDET (Now SIGMOD) Workshop on Data Description, Access and Control, ser. SIGFIDET '74. New York, NY, USA: ACM, 1974, pp. 249–264. [Online]. Available: <http://doi.acm.org/10.1145/800296.811515>
- [23] Google, "Google query language in google app engine," Last access nov, 2013. [Online]. Available: <https://developers.google.com/appengine/docs/python/datastore/gqlreference>
- [24] —, "Google query language in google cloud datastore," Last access nov, 2013. [Online]. Available: <https://developers.google.com/datastore/docs/concepts/gql>
- [25] Yahoo!, "Yahoo query language," Last access nov, 2013. [Online]. Available: <http://developer.yahoo.com/yql/>
- [26] A. Offutt and J. Pan, "Detecting equivalent mutants and the feasible path problem," in Computer Assurance, 1996. COMPASS '96, Systems Integrity. Software Safety. Process Security. Proceedings of the Eleventh Annual Conference on, 1996, pp. 224–236.
- [27] J. B. Li and J. Miller, "Testing the semantics of w3c xml schema," in Computer Software and Applications Conference, 2005. COMPSAC 2005. 29th Annual International, vol. 1, 2005, pp. 443–448 Vol. 2.
- [28] L. Franzotte and S. R. Vergilio, "Applying mutation testing in xml schemas," in SEKE, K. Zhang, G. Spanoudakis, and G. Visaggio, Eds., 2006, pp. 511–516.
- [29] U. Praphamontripong and J. Offutt, "Applying mutation testing to web applications," in Software Testing, Verification, and Validation Workshops (ICSTW), 2010 Third International Conference on, 2010, pp. 132–141.
- [30] A. Estero Botaro, F. Palomo Lozano, and I. Medina Bulo, "Mutation operators for WS-BPEL 2.0," in ICSSEA 2008: Proceedings of the 21th International Conference on Software & Systems Engineering and their Applications, 2008. [Online]. Available: <http://neptuno.uca.es/redmine/wiki/gamera>
- [31] A. Estero-Botaro, J. Boubeta-Puig, V. Liñeiro-Barea, and I. Medina-Bulo, "Operadores de mutación de cobertura para ws-bpel 2.0," in XVII Jornadas de Ingeniería del Software y Bases de Datos (JISBD - SISTEDES 2012), Almería, Spain, 2012.
- [32] R. Wang and N. Huang, "Requirement model-based mutation testing for web service," in Next Generation Web Services Practices, 2008. NWESP '08. 4th International Conference on, 2008, pp. 71–76.
- [33] padraic, "Mutagenesis," Last access dec, 2013. [Online]. Available: <https://github.com/padraic/mutagenesis>
- [34] E. Martin and T. Xie, "A fault model and mutation testing of access control policies," in Proceedings of the 16th International Conference on World Wide Web, ser. WWW '07. New York, NY, USA: ACM, 2007, pp. 667–676. [Online]. Available: <http://doi.acm.org/10.1145/1242572.1242663>

MeterGoat: A Low Cost Hardware Platform for Teaching Smart Meter Security

Jefferson Capovilla, Nelson Uto
GTSIC – Information Security Department
CPqD
Campinas, São Paulo, Brazil
{jrodrigo, uto}@cpqd.com.br

Danilo Suiama, Jose Resende
Management of Metering, Losses and Technology
ELEKTRO
Campinas, São Paulo, Brazil
{danilo.suiama, jose.resende}@elektro.com.br

Abstract—Smart meters play an important role in smart grid architectures by enabling best operational efficiency, enhanced usage monitoring and variable pricing structure, among other advantages. On the other hand, meters have often been deployed with several security vulnerabilities that can compromise the mentioned benefits and result in cyber attacks. Therefore, teams involved in the development of smart meters should be trained in security aspects related to the area. In this context, this paper describes the work in progress project called *MeterGoat*, whose main objective is to develop an open source and low cost hardware platform for teaching smart meter security. In order to achieve this purpose, it deliberately implements meter's functionalities in a vulnerable way, providing a flexible framework for smart meter pentesting including hardware and firmware weaknesses.

Keywords- Smart meter pentesting; training platform; cyber attacks; vulnerability.

I. INTRODUCTION

Smart grids have the objective of transforming today's power grid by providing an extra level of grid status control, energy consumption profiles, powerful control mechanisms, and flexible billing processes. However, these advantages increase the overall complexity of the system, since it is necessary to use more advanced components, in order to maintain and transmit the information required by the utility companies. A major drawback of this new scenario is the high number of vulnerabilities reported on smart meters and the underlying infrastructure, as can be seen in the papers of Skopik et al. [1] and Carpenter et al. [2].

Clearly, there is a gap between the energy sector and information technology field with respect to information security. The main objective of this work in progress paper is then to fill this gap by introducing a low cost and open source hardware platform for teaching security aspects that must be considered in the implementation of a smart meter.

We named it *MeterGoat*, after the Open Web Application Security Project (OWASP) *WebGoat* project [3], which has the same purpose, but in the scope of web applications. *MeterGoat* will then implement the most relevant smart meter's functionalities in a vulnerable way, allowing the students to perform real attacks without affecting live environments. In this way, they can fully understand why the exploits are possible and avoid making the same mistakes in real products they develop.

The remaining part of this paper is structured as follows: Section II presents a few works related to security training platforms and smart meter penetration testing. Section III describes the *MeterGoat* project, outlining the vulnerabilities that are being implemented and the hardware components we are using. In Section IV, we list a few tools to compose the toolkit for testing smart meter security. Section V gives an overview of training scenarios, and, finally, Section VI concludes the paper and discusses future work.

II. RELATED WORK

The recent demand of countries in substituting standard energy meters by the smart version contributes to foster research in Advanced Metering Infrastructure (AMI) security, such as the AMI Penetration Test Plan [4], by the National Electric Sector Cybersecurity Organization Resource (NESCOR). That document proposes that the security evaluation of smart meters should cover four main areas: embedded devices, network communications, server Operating System (OS), and server applications.

The work of Grand [5] complements the aforementioned guideline, since, besides testing, it also covers the concept of designing secure hardware for embedded systems, introducing, in this way, security in the earlier phases of development lifecycle. According to the author, in order to have a more secure product, one needs, at least, to avoid vulnerabilities in the enclosure, the circuit board, and the firmware. In our project, we are using several examples of insecure design given by Grand.

Regarding security training, one can find several open source projects, which implement insecure web and mobile applications, such as *WebGoat*, *iGoat* [6], and *GoatDroid* [7]. However, to the best of the author's knowledge, there is no testbed related to smart meters, making it hard for those interested in the area to learn how to test this type of device. Thus, we believe that *MeterGoat* can help to fill this gap.

As we already mentioned, *Webgoat* is a deliberately insecure web application, designed to teach penetration testing. It contains more than 30 lessons that emulate vulnerabilities commonly found in real applications. The *iGoat* and the *GoatDroid* are quite similar projects, presenting themselves as fully functional and self-contained security training environments, for iOS and Android, respectively. All three projects are maintained by OWASP.

III. THE METERGOAT PROJECT

This section presents the functional requirements of the *MeterGoat* Project, as well as preliminary project decisions.

A. Hardware Architecture of Commercial Smart Meters

The hardware architecture of *MeterGoat* is based on commercial smart meters that, in general, are composed of the elements showed in Fig. 1: current and voltage sensors, Analog-to-Digital Converters (ADC), Central Processing Unit (CPU) for metrology calculations, Random-Access Memory (RAM) for volatile data and firmware execution, non-volatile storage, e.g. flash and Electrically-Erasable Programmable Read-Only Memory (EEPROM), for firmware and data storage, and communication peripherals for maintenance and update procedures.

The majority of the microcontrollers contain integrated memory in the System on Chip (SoC) itself. However, if system functionalities require more memory than the default amount, external memory chips can be used to increase the overall capacity. Architectures that use microcontrollers generally comply with low cost, low power, and compact design requirements. On the other hand, systems that demand high performance employ dedicated processors and components, resulting in a more expensive device. Since smart meters are produced in high volumes, the first, cheaper, architecture is commonly adopted.

Anti-tampering mechanisms comprise a layer of defense against a possible integrity violation of an equipment. In this way, they must cope with any attempt, physical or electronic, of adulteration, such as opening the device's case, accessing internal memory, and replacing components. These controls, explained in detail in [5], are usually categorized into four groups: resistance, evidence, detection, and response.

The Joint Test Action Group (JTAG) Debug, illustrated in Fig. 1, provides a direct connection to most of the components of a SoC. One can use it, through the JTAG interface, to perform memory programming, boundary scanning, and on-chip debugging. Although the security best practice is to disable this interface, by physically damaging the Printed Circuit Board (PCB) connection or by blowing the JTAG fuse, several meters do not implement this recommendation.

It is important to note that the architectural diversity, presented in this section, determines the types of hardware

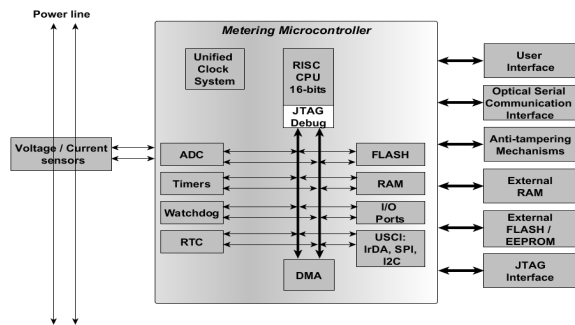


Figure 1 - Smart Meter High Level Architecture

vulnerabilities that one may encounter on these devices, since it provides different levels of access to metrological information and code.

B. Hardware Architecture of MeterGoat

The *MeterGoat* architecture, depicted in Fig. 2, derives from the models discussed in the previous section. We made some simplifications in order to reduce the project final cost and to obtain a modularization of the training platform, while still implementing the main hardware and software security weaknesses commonly found on real devices [8]. We briefly explain the main components of our project below:

1) Microcontroller MSP430F5438 and Development Kit

Based on the architecture illustrated in Fig. 2, we selected the Texas Instruments (TI) MSP-EXP430F5438 [9] development kit and the corresponding microcontroller, MSP430F5438 [10], for building *MeterGoat*. The later has, as main features, a 256KB internal flash memory, 16KB internal Static RAM (SRAM) memory, two Serial Peripheral Interface (SPI) interfaces, and 87 Input/Output (I/O) pins.

The microcontroller above contains a hardware multiplier that performs signed and unsigned operations with 8, 16, 24, and 32-bit operands. This feature allows the optimization of some metrologic operations, and, for this reason together with the aforementioned characteristics, several smart meters are built over this specific model or similar members from the same family of microcontrollers.

The development kit also comes with a 138x110 grayscale Liquid Crystal Display (LCD) for measurement display, a JTAG interface, a 5-position joystick, and two push buttons for user interface control. Finally, it can be powered in three different ways: Universal Serial Bus (USB), Flash Emulation Tool (FET), and two AA batteries.

2) External Components

Since the objective of the *MeterGoat* project is not to build a real smart meter, we are not using a high-priced energy metering Integrated Circuit (IC) such as the ADE7758 [11]. Instead, we are employing a programmable pulse generator, in order to emulate the signal from the voltage/current sensor, responsible for measuring energy consumption, and sample it using a digital I/O pin. One of the ways to generate such signal consists in assembling

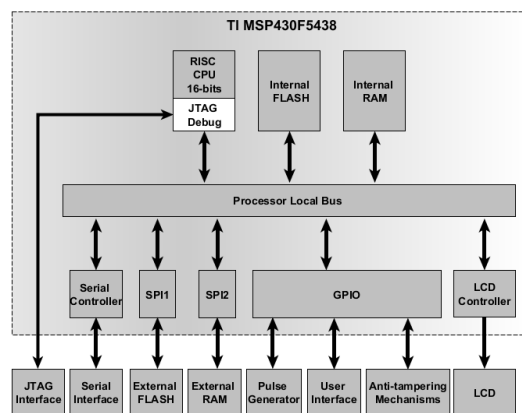


Figure 2 - *MeterGoat* High Level Architecture

discrete components and a LM555 timer [12] in a breadboard or using an Arduino Uno microcontroller board [13].

We also use a breadboard to assemble and power the external Flash and RAM memories. These components communicate with the microcontroller using an SPI interface.

3) *Anti-tampering Mechanisms*

We are implementing detection and response anti-tampering mechanisms in *MeterGoat*. The former includes: (i) anti-tampering switches, which one normally places on the meter's case so they trigger when someone attempts to open the chassis; (ii) temperature gradient measurement, against attacks that freeze components in order to take advantage of data remanence. This countermeasure will use the internal temperature sensor included in MSP430F5438.

With regard to response mechanisms against tampering, *MeterGoat* will implement passive zeroization, which involves disconnecting power from volatile memory so that content is lost, and active zeroization, which securely erases sensitive data used by the system.

C. *Exposed Interfaces*

In order to perform an attack, it is fundamental to initially identify the target's surface attack. In the case of smart meters, one must consider external configuration ports, network interfaces, buses, and electronic component pins, as explained in the present section.

1) *JTAG Interface*

JTAG is a standard interface composed at minimum by four pins, Test Mode Select (TMS), Test Data In (TDI), Test Data Out (TDO), and Test Clock (TCK), that is commonly used for interfacing with circuit boards, microprocessors, and several other peripherals for debugging purposes. In microcontrollers, the JTAG is also used to upload and store the firmware in internal flash memory. As showed in Fig. 2, by having access to an enabled JTAG interface, it is possible to perform attacks such as memory dumps and firmware extraction. For these reasons, the JTAG interface must always be disabled in final products, otherwise an attacker can dump and modify the component's contents.

2) *External Flash/EEPROM Memory*

These are discrete chips with the exclusive function of non-volatile storage, usually providing no protection whatsoever. When a smart meter architecture includes this type of element in the project, it is possible to dump or modify the entire memory content, that may contain measurement data, configuration parameters, or the firmware. This task can be performed with the chip soldered on the PCB, if it uses the Inter-Integrated Circuit (I²C) communication protocol (accepting multi-master on the bus). Otherwise, the chip needs to be desoldered for direct pin manipulation using the tools presented in Section IV.

3) *External RAM*

Consists in discrete chips with the exclusive function of volatile storage, usually providing no native protection against an attacker with physical access to them. When a smart meter architecture includes this kind of element in the project, it is possible to dump the memory addresses

accessed by the microcontroller through bus snooping. This task can be done using the tools presented in Section IV.

4) *Optical Interface*

The optical port is often used by a smart meter for configuration purposes. It is based on short range communication, making vendors wrongly assume it is less prone to sniffing and adulteration. A weak access control system may allow an attacker to send commands to the meter by this type of interface.

5) *Network Interface*

A smart meter can use a wired or a wireless network to communicate with the central system or aggregators. If a secure communication protocol is not employed, it is possible to capture and alter information in transit as well as inject packets.

6) *Communication Interface*

A smart meter can have several communication interfaces such as USB, I²C, Universal Asynchronous Receiver/Transmitter (UART), Controller Area Network (CAN), SPI, etc., which do not provide any tamper protection. Therefore, an attacker can use them to analyze information being transmitted, dump accessible data, and disturb the communication by injecting invalid packets or random data.

D. *Vulnerabilities on MeterGoat*

In the first version of *MeterGoat*, we are implementing all the vulnerabilities explained below:

- **Unprotected interfaces** – the exposed interfaces described in Section III.C allows an attacker to dump information from the meter as well as to send commands to be performed by the device.
- **Broken cryptographic algorithms** – this class of vulnerability includes the use of home-made algorithms and classical cryptosystems, such as the shift cipher and Vigenère's, for which cryptanalysis is possible.
- **Incorrect use of cryptograph** – includes Electronic CodeBook (ECB) block cipher mode of operation for large messages, binary additive stream cipher with key reuse, and use of a cryptographic mechanism for a purpose different than the originally intended [14].
- **Insecure cryptographic key management** – consists in using predictable keys or known weak keys such as those for DES [14], embedding cryptographic keys in code, storing cryptographic keys in cleartext files, the lack of use of a secure static memory for key protection, and the absence of a functionality for key substitution.
- **Unprotected data at rest** – consists in the storage of sensitive information in cleartext form and the use of an encoding mechanism such as Base64 instead of a cipher.
- **Insecure communication** – encompasses cleartext communication between client software and the meter or through a protocol with known vulnerabilities such as the one described in [15].

- **Broken user authentication mechanism** – includes weak or absent password policy, flawed authentication protocol, and insecure storage of passwords.
- **Authorization flaws** – this type of problem arises when the concept of a reference monitor is not properly implemented, resulting, for instance, in privilege escalation attacks and direct accesses to resources.
- **Lack of integrity mechanism** – the absence of such a security mechanism allows an attacker to replace the firmware, altering or inserting functionalities, and to tamper with metrology data.
- **Firmware implementation flaws** – involve classical attacks such as buffer overflow, format string attack, and integer overflow.
- **Flawed anti-tampering mechanism** – although *MeterGoat* will not have a case at all, the idea of this exercise is to show how a bad casing could be explored to bypass anti-tampering controls.

IV. TRAINING TOOLKIT

The training toolkit comprises a selection of tools that can be used in the evaluation of smart meter security. One should not assume, however, ours is the only possible list, since there are similar tools available for choosing. If one decides to build its own set, one should select those that cover as many features as possible.

We divided our list in three groups: *interface identifier*, *interface manipulator*, and *software*. The former aids in identifying pins and test points, verifying pin voltage, and monitoring signal transitions. In this set, we recommend the tools below:

- **Multimeter** - is an electronic instrument that combines several measurement functions in a single unit. Since it can measure a wide range of voltages without being damaged, it is a robust equipment for measuring pins and test points for the first time. It is also used to find the connections between components by using the continuity test feature. When two points are electrically connected, a tone is emitted.
- **Oscilloscope** - is used to observe the change of an electrical signal over time. For reverse engineer, the purpose is to verify signal transitions to check the communication between the microcontroller and other components for later bus sniffing.
- **JTAGulator** - is an open source hardware tool that assists in identifying JTAG pin (TMS, TDI, TDO, TCK, TRST) connections from test points, vias, or component pads on a target device [16].

The interface manipulator tools aid in obtaining information stored in memory or transmitted through interfaces and buses. In this category, we suggest the following tools:

- **Bus Pirate** - performs serial bus manipulation. Supports many serial protocols at 0-5.5 volts such as 1-Wire, I²C, SPI, and asynchronous serial [17].

- **Logic Analyzer** - performs bus sniffing by recording, viewing, and measuring digital signals in transit between components. It can understand different protocols including serial, I²C, SPI, and CAN.
- **GoodFET** - is an open source JTAG adapter, used for TI MSP430 as a debugger and flash emulation tool [18].
- **MSP-FET430UIF** - is the official equipment provided by TI to MSP430 for JTAG debugging and flash emulation tool [19].

Finally, in order to perform firmware analysis, one needs a disassembler and a debugger, which support the platforms employed by the smart meter. One of the best tools in this category is the powerful IDA, but it has the disadvantage of being relatively expensive. Unfortunately, no open source counterpart works with the very specific processors used by meters. Moreover, there is no disassembler at all for some architectures, and, then, one would have to be built.

V. TRAINING SCENARIOS

The training scenarios we propose in this section are based on references [4] and [5] and on security weaknesses found in tests we performed against commercial smart meters. Since security is a dynamic area, with new vulnerabilities being discovered every day, the list of lessons may be updated in future versions.

In order to make the most of the training, the student must have basic knowledge of the following topics: electrical circuits, communication protocols (e.g., I²C, SPI, and serial), embedded systems, assembly and Python languages, software reverse engineering, and cryptography.

Each lesson of the course covers the analysis of a different part of a smart meter, resulting in four groups [4]: (1) electronic components; (2) field technician interface; (3) binary firmware; and (4) cryptographic mechanisms.

In the first part of the training, the students learn how to identify the main components and to read the corresponding datasheets, in order to map pins and understand their functionalities. After that, we proceed to the reverse engineer of the PCB, by identifying, with the help of a multimeter and an oscilloscope, connections among components, operating voltages, and signal transitions. Techniques to identify and bypass a simple tamper detection mechanism are taught using a multimeter. The idea is to show how to detect a logic '1' pin and to employ an extra wire, with the purpose of keeping the signal high for the port responsible for tamper monitoring.

The lecture ends with dumping of data from the non-volatile memories and snooping of the bus connecting the memory modules to the microcontroller, respectively, by using Bus Pirate and Logic Analyzer. *MeterGoat* will store and transmit information protected by a multitude of mechanisms, so the students can practice string analysis, entropy analysis, data decoding, and systematic key search.

For the second part of the course, *MeterGoat* provides a serial port as a field technician interface. The instructor will

teach how it can be used to interact with the equipment, using standard protocols, and how to implement these with the Python language. This will be the base for teaching protocol fuzzing and vulnerability exploitation.

The third group of lessons comprises binary firmware extraction and analysis, through the manipulation of the JTAG interface. The students will use GoodFET or MSP-FET430UIF, in order to dump the simple firmware stored in microcontroller's internal memory of *MeterGoat*. Since reverse engineering requires very advanced knowledge, we will cover only the disassembly and analysis of a short piece of authentication code.

Although advanced cryptanalysis is beyond the scope of a course about smart meter security, one needs to know how to identify basic mistakes in the implementation of cryptographic mechanisms. For this reason, the lessons of the fourth group are related to the most common vulnerabilities in this area, mainly those resulting from improper key management.

VI. CONCLUSIONS AND FUTURE WORK

This paper presented the work in progress project called *MeterGoat*, whose main objective is to provide a low cost platform for smart meter security training. The platform will provide most of the functionalities and interfaces of real smart meters, all implemented in a vulnerable way. It is important to mention that we do not intend to build a real meter. Thus, some of the functions will just be emulated, but in a way that allows the student to learn about a given security vulnerability we want to stress. The framework we defined is very flexible and one can easily extend it to include new types of vulnerabilities or variations of old weaknesses.

Currently, we have already specified the full platform, selected the list of components, and started the hardware assembly and coding of the chosen firmware vulnerabilities. We expect to spend six more months in this task, and, once finished, we intend to provide the schematics and firmware as an open source project. In this way, engineers and security analysts will be able to build and use *MeterGoat*, totally free of charge, for personal use, provided the user license be respected.

The estimated cost for building the platform lies under US\$200, which is reasonable for this type of equipment. On the other hand, the training toolkit is more expensive, and one can expect to spend about US\$1700, without a disassembler/debugger.

ACKNOWLEDGMENTS

This research project has been carried out in partnership between CPqD (independent institution focused on R&D and innovation) and ELEKTRO (the eighth largest power supply organization in Brazil) and relies on R&D financial resources provided and managed by ANEEL.

We would like to give special thanks to Rafael de Simone Cividanes, CPqD's consultant, for coordinating the project and reviewing this paper.

REFERENCES

- [1] F. Skopik, Z. Ma, T. Bleier, and H. Gruneis, "A survey on threats and vulnerabilities in smart metering infrastructures", *International Journal of Smart Grid and Clean Energy*, vol. 1, no. 1, Sep. 2012, pp. 22-28.
- [2] M. Carpenter, T. Goodspeed, B. Singletary, J. Searle, E. Skoudis, and J. Wright, "Advanced metering infrastructure attack methodology", Mar. 2011, InGuardians, Inc. 2.0.
- [3] WebGoat. Available from: https://www.owasp.org/index.php/Webgoat_WebGoat_Project. Accessed: Feb. 26th, 2014.
- [4] J. Searle, G. Rasche, A. Wright, and S. Dinnage, "AMI Penetration Test Plan". Available from: <http://www.smartgrid.epri.com/doc/AMI-Penetration-Test-Plan-1-0-RC3.pdf>. Accessed: Apr. 11th, 2014.
- [5] J. Grand, "Practical Secure Hardware Design for Embedded Systems", *Proc. 2004 Embedded Systems Conference (ESC 04)*, Mar. 2004, pp. 1-25.
- [6] iGoat. Available from: https://www.owasp.org/index.php/OWASP_iGoat_Project. Accessed: Apr. 11th, 2014.
- [7] GoatDroid. Available from: <https://github.com/jackMannino/OWASP-GoatDroid-Project/wiki>. Accessed: Apr. 11th, 2014.
- [8] InGuardians, Inc., "Advanced metering infrastructure attack methodology", Mar. 2011.
- [9] MSP430F5438 Experimenter board. Available from: <http://www.ti.com/tool/msp-exp430f5438>. Accessed: Feb. 26th, 2014.
- [10] MSP430F5438 Features. Available from: <http://in.embeddeddeveloper.com/processors/3252/Texas-Instruments/MSP430F5438.htm>. Accessed: Feb. 26th, 2014.
- [11] Analog Devices, "Poly phase multifunction energy metering IC with per phase information", *ADE7758 datasheet*, Oct. 2011.
- [12] Texas Instruments, "LM555 Timer", Mar. 2013.
- [13] Arduino Uno. Available from: <http://arduino.cc/en/Main/arduinoBoardUno>. Accessed: Feb. 26th, 2014.
- [14] A. Menezes, P. C. van Oorschot, and S. A. Vanstone, "Handbook of applied cryptography". CRC Press. Aug. 2001.
- [15] N. AlFardam, D. Bernstein, K. Paterson, and J. Schuldt, "On the security of RC4 in TLS", *Proc. of 22nd USENIX Security Symposium*. Aug. 2013, pp. 305-320.
- [16] JTAGulator - open source hardware for OCD identification. Available from: <http://www.grandideastudio.com/portfolio/jtagulator/>. Accessed: Feb. 26th, 2014.
- [17] Bus Pirate - open source hacker multi-tool. Available from: <http://dangerousprototypes.com/bus-pirate-manual/>. Accessed: Feb. 26th, 2014.
- [18] GoodFET - JTAG debugger for TI MSP430. Available from: <http://goodfet.sourceforge.net/>. Accessed: Feb. 26th, 2014.
- [19] MSP-FET430UIF - Official JTAG debugger for TI MSP430. Available from: <http://www.ti.com/tool/msp-fet430uif>. Accessed: Feb. 26th, 2014.