

9TH CONFERENCE OF PHD STUDENTS IN COMPUTER SCIENCE

Volume of extended abstracts

CS²

Organized by the Institute of Informatics of the University of Szeged



June 30 – July 2, 2014
Szeged, Hungary

National Development Agency
www.ujszecsenyiterv.gov.hu
06 40 638 638



The project is supported by the European Union
and co-financed by the European Social Fund.



Scientific Committee:

János Csirik (Co-Chair, SZTE)
Lajos Rónyai (Co-Chair, SZTAKI, BME)
András Benczúr (ELTE)
Hasszan Charaf (BME)
Tibor Csendes (SZTE)
László Cser (BCE)
János Demetrovics (SZTAKI, ELTE)
József Dombi (SZTE)
Ferenc Friedler (PE)
Zoltán Fülöp (SZTE)
Aurél Galántai (ÓE)
Zoltán Gingl (SZTE)
Tibor Gyimóthy (SZTE)
Zoltán Horváth (ELTE)
Csanád Imreh (SZTE)
Zoltán Kató (SZTE)
Zoltán Kása (Sapientia EMTE)
László Keviczky (SZIE)
László Kóczy (SZIE)
János Kormos (DE)
Erzsébet Kovács (Corvinus)
János Levendovszky (BME)
Dániel Marx (SZTAKI)
László Nyúl (SZTE)
Attila Pethő (DE)
András Recski (BME)
Jenő Szigeti (ME)
Tamás Szirányi (SZTAKI)
Péter Szolgay (PPKE)

Organizing Committee:

Rudolf Ferenc, Balázs Bánhelyi, Tamás Gergely, Zoltán Kincses

Address of the Organizing Committee

c/o. Rudolf Ferenc

University of Szeged, Institute of Informatics

H-6701 Szeged, P.O. Box 652, Hungary

Phone: +36 62 546 396, Fax: +36 62 546 397

E-mail: cscs@inf.u-szeged.hu

URL: <http://www.inf.u-szeged.hu/~cscs/>

Sponsors

Telemedicine Oriented Research in the Fields of Mathematics, Informatics and Medical Sciences

TÁMOP-4.2.2.A-11/1/KONV-2012-0073

Association of Hungarian PhD and DLA Students Scientific Section of Mathematics and Informatics

Evopro Group

Lufthansa Systems

University of Szeged, Institute of Informatics

Polygon Publisher

Preface

This conference is the ninth in a series. The organizers aimed to bring together PhD students working on any field of computer science and its applications to help them publishing one of their first abstracts and papers, and provide an opportunity to hold a scientific talk. As far as we know, this is one of the few such conferences. The aims of the scientific meeting were determined on the council meeting of the Hungarian PhD Schools in Informatics: it should

- provide a forum for PhD students in computer science to discuss their ideas and research results;
- give a possibility to have constructive criticism before they present the results at professional conferences;
- promote the publication of their results in the form of fully refereed journal articles; and finally,
- promote hopefully fruitful research collaboration among the participants.

The papers emerging from the presented talks will be forwarded to the Acta Cybernetica journal.

The organizers hope that the conference will be a valuable contribution to the research of the participants, and wish a pleasant stay in Szeged.

Szeged, June 2014

*Rudolf Ferenc
Balázs Bánhelyi
Tamás Gergely
Zoltán Kincses*

Contents

Preface	i
Contents	ii
Program	iv
Plenary talks	1
Zoltán Horváth, Zoltán Istenes, and Zsuzsa Várhalmi: <i>Research, Innovation and Education in Infocommunication Technology</i>	1
John Milton: <i>Transient stabilization of unstable states: why a stick balanced at the fingertip always falls?</i>	1
Pawel Pilarczyk: <i>Computational homology</i>	2
Abstracts	3
Elvira Antal and Tamás Vinkó: <i>Mathematical modeling of max–min fair bandwidth allocation in BitTorrent communities</i>	3
Viktor Árgilán, János Balogh, József Békési, Balázs Dávid, Miklós Krész: <i>Heuristic approach for the driver scheduling problem</i>	4
Viktor Árgilán, János Balogh, Balázs Dávid, Miklós Krész: <i>Integrating fast heuristics into a decision support system for vehicle rescheduling</i>	5
Gergő Balogh: <i>Validation of the City Metaphor in Software Visualization</i>	6
Áron Baráth and Zoltán Porkoláb: <i>Language Support for High-level Worst-case Execution Time Estimation</i>	7
Péter Bodnár and László G. Nyúl: <i>QR Code Localization Using Boosted Cascade of Weak Classifiers</i>	8
Zsolt Borsi: <i>A translation of interaction relationships into SMV language</i>	10
Csaba Faragó: <i>Variance of Source Code Quality Change Caused by Version Control Operations</i>	12
Bálint Fodor: <i>Sum of Gaussian Functions Based 3d Point Cloud Registration</i>	14
Roberto Giachetta and István Fekete: <i>A case study of advancing remotely sensed image processing</i>	16
Zoltán Gilián: <i>Generalized Haar systems toolbox for MATLAB</i>	18
Erika Griechisch and Gábor Németh: <i>Offline Signature Verification Using Similarity Measure for Skeletons</i>	19
László Hajdu, Miklós Krész, Attila Tóth: <i>Two-phase graph coloring heuristic for crew rostering</i>	20
Norbert Hantos and Péter Balázs: <i>Eliminating Switching Components in Binary Matrices</i>	21
Ferenc Horváth, Richárd Dévai, Tamás Gergely: <i>Structural Information Aided Automated Test Method for Magic 4GL</i>	22
Ferenc Horváth, Benjamin Mészáros, Tamás Gergely: <i>Usability Testing of Android Applications</i>	24
Péter Hudák and László Lengyel: <i>Test-driven verification of model transformations</i>	25
Gábor Imre and Gergely Mezei: <i>Graph Transformation-based Opinion Mining in Web Queries</i>	26
Máté Karácsony: <i>Comparison of Source Code Transformation Methods for Clang</i>	27
Melinda Katona and László G. Nyúl: <i>Fast recognition of natural feature identifiers by a mobile phone</i>	28
Sándor Kazi and Gábor Nagy: <i>A scalable parallel boosting scheme for bulk synchronous parallel environments</i>	29
György Kovács and László Tóth: <i>The Joint Optimization of Spectro-Temporal Features and Deep Neural Nets for Robust ASR</i>	31
László Kundra: <i>Methods for Feature Point Aggregation of Optical Flow for Orientation Estimation</i>	33
Gergely Ladányi: <i>Business process quality measurement using advances in static code analysis</i>	34
Balázs L. Lévai: <i>Automatic Failure Detection and Monitoring of Ventilation and Cooling Systems</i>	35
Balázs L. Lévai: <i>Automatic Design of LED Street Lights</i>	36

András London and Tamás Németh: <i>Application of graph based data mining techniques in administrative systems of education</i>	37
Gábor I. Nagy, Sándor Kazi, Győző Papp: <i>Distributed News Analytics Framework: Collecting News Feed Sources from social media</i>	39
Gábor I. Nagy and Sándor Kazi: <i>Filtering noise from stock related Twitter messages</i>	41
Tamás Dániel Nagy, Gergely Vadai, Zoltán Gingl, László Rudas, and Éva Zöllei: <i>Phase delay detecting methods to analyse the correlation between blood pressure and paced breathing</i>	43
István Nagy T.: <i>Automatic Detection of Multiword Expressions with Dependency Parsers on Different Languages</i>	45
Boldizsár Németh and Zoltán Kelemen: <i>Derivable Partial Locking for Algebraic Data Types</i>	47
Zoltán Németh: <i>Visualization and analysis of financial transaction networks - a multidimensional modeling approach</i>	49
Zoltán Ozsvár and Péter Balázs: <i>Reconstruction of $h\nu$-Convex Binary Matrices from Horizontal and Vertical Projections Based on Simulated Annealing</i>	50
Krisztián Pándi and Hassan Charaf: <i>Load Balancing Strategy in Mobile Resource Management</i>	51
Zsombor Paróczy: <i>x86 instruction reordering and split-stream compression benchmark</i>	52
Zsombor Paróczy, Bálint Fodor, Gábor Szűcs: <i>Diversification for Content-Based Visual Information Retrieval System for Video</i>	54
Attila Selmei and István Orosz: <i>Workflow processing using SAP Objects</i>	56
Attila Selmei and Tamás Gábor Orosz: <i>New Approaches of Storing ERP Data</i>	57
Dávid Szalóki: <i>Camera Placement Optimization in Object Localization Systems</i>	59
Péter Szűcs: <i>Challenges in Real-time Collaborative Editing</i>	60
Szabolcs Urbán, Antal Nagy, László Ruskó: <i>Tumor detection and segmentation on multimodal medical images</i>	62
Gyula Vörös: <i>Expanding Small Corpora to Aid People with Communication Impairments</i>	63
List of Authors	65

Program

Monday, June 30

08:00 – 09:00	Registration and Opening
09:00 – 11:00	Talks – Recognition (5x20 minutes)
11:00 – 12:00	Plenary talk
12:00 – 12:40	Talks – Medical Solutions (2x20 minutes)
12:40 – 14:00	Lunch
14:00 – 15:40	Talks – Software Quality (4x20 minutes)
15:40 – 17:00	Talks – Software Testing (3x20 minutes)
17:00 – 18:00	Talks – Mobile and ERP (3x20 minutes)
19:00	Reception at the Rector’s Building

Tuesday, July 1

09:00 – 11:00	Talks – Compilers and Languages (5x20 minutes)
11:00 – 11:50	Plenary talk
11:50 – 12:00	Association of Hungarian PhD and DLA Students
12:00 – 12:50	Plenary talk
12:50 – 14:00	Lunch
14:00 – 15:40	Talks – Image Processing (4x20 minutes)
15:40 – 16:40	Talks – Mathematical Problems/Models (3x20 minutes)
16:40	Social program
19:00	Gala Dinner at the Rector’s Building

Wednesday, July 2

09:00 – 11:00	Talks – Artificial Intelligence (5x20 minutes)
11:00 – 12:40	Talks – Optimization (5x20 minutes)
12:40 – 14:00	Lunch
14:00 – 15:20	Talks – Data Mining (4x20 minutes)
15:20	Closing

Detailed program

Monday, June 30

08:00	Registration
08:45	Opening
Session 1	Recognition
09:00	Szabolcs Urbán and Antal Nagy: <i>Tumor detection and segmentation on multimodal medical images</i>
09:20	Bálint Fodor: <i>Sum of Gaussian Functions Based 3d Point Cloud Registration</i>
09:40	Melinda Katona and László G. Nyúl: <i>Fast recognition of natural feature identifiers by a mobile phone</i>
10:00	Péter Bodnár and László G. Nyúl: <i>QR Code Localization Using Boosted Cascade of Weak Classifiers</i>
10:20	Erika Griechisch and Gábor Németh: <i>Offline Signature Verification Using Similarity Measure for Skeletons</i>
10:40	Break
11:00	Plenary Talk Prof. Pawel Pilarczyk: <i>Computational homology</i>
11:50	Break
Session 2	Medical Solutions
12:00	Gyula Vörös: <i>Expanding Small Corpora to Aid People with Communication Impairments</i>
12:20	Tamás Dániel Nagy, Gergely Vadai, Zoltán Gingl, László Rudas, Éva Zöllei: <i>Phase delay detecting methods to analyse the correlation between blood pressure and paced breathing</i>
12:40	Lunch
Session 3	Software Quality
14:00	Gergő Balogh: <i>Validation of the City Metaphor in Software Visualization</i>
14:20	Gergely Ladányi: <i>Business Process Quality Measurement Using Advances in Static Code Analysis</i>
14:40	Csaba Faragó: <i>Variance of Source Code Quality Change Caused by Version Control Operations</i>
15:00	Peter Szűcs: <i>Challenges in Real-time Collaborative Editing</i>
15:20	Break
Session 4	Software Testing
15:40	Péter Hudák, and László Lengyel: <i>Test-driven Verification of Model Transformations</i>
16:00	Ferenc Horváth, Richárd Dévai, and Tamás Gergely: <i>Structural Information Aided Automated Test Method for Magic 4GL</i>
16:20	Ferenc Horváth, Benjamin Mészáros, and Tamás Gergely: <i>Usability Testing of Android Applications</i>
16:40	Break
Session 5	Mobile and ERP
17:00	Krisztián Pándi, and Hassan Charaf: <i>Load Balancing Strategy in Mobile Resource Management</i>
17:20	Attila Selmeçi, and Tamás Gábor Orosz: <i>New Approaches of Storing ERP Data</i>
17:40	Attila Selmeçi, and István Orosz: <i>Workflow Processing Using SAP Objects</i>
18:00	Free program
19:00	Reception at the Rector's Building

Tuesday, July 1

Session 6	Compilers and Languages
09:00	Zsombor Paróczy: <i>x86 Instruction Reordering and Split-stream Compression Benchmark</i>
09:20	Máté Karácsony: <i>Comparison of Source Code Transformation Methods for Clang</i>
09:40	Áron Baráth, and Zoltán Porkoláb: <i>Language Support for High-level Worst-case Execution Time Estimation</i>
10:00	Boldizsár Németh, and Zoltán Kelemen: <i>Derivable Partial Locking for Algebraic Data Types</i>
10:20	Zsolt Borsi: <i>A Translation of Interaction Relationships into SMV Language</i>
10:40	Break
11:00	Plenary Talk Prof. Zoltán Horváth: <i>Research, Innovation and Education in Infocommunication Technology</i>
11:50	Introduction of Association of Hungarian PhD and DLA Students, Scientific Section of Mathematics and Informatics
12:00	Plenary Talk Prof. John Milton: <i>Transient stabilization of unstable states: why a stick balanced at the fingertip always falls?</i>
12:50	Lunch
Session 7	Image Processing
14:00	Roberto Giachetta and István Fekete: <i>A case study of advancing remotely sensed image processing</i>
14:20	László Kundra and Péter Ekler: <i>Methods for Feature Point Aggregation of Optical Flow for Orientation Estimation</i>
14:40	Norbert Hantos and Péter Balázs: <i>Eliminating Switching Components in Binary Matrices</i>
15:00	Zoltán Ozsvár and Péter Balázs: <i>Reconstruction of $h\nu$-Convex Binary Matrices from Horizontal and Vertical Projections Based on Simulated Annealing</i>
15:20	Break
Session 8	Mathematical Problems/Models
15:40	Zoltán Gilián: <i>Generalized Haar systems toolbox for MATLAB</i>
16:00	Elvira Antal and Tamás Vinkó: <i>Mathematical modeling of max-min fair bandwidth allocation in BitTorrent communities</i>
16:20	Zoltán Németh: <i>Visualization and analysis of financial transaction networks - a multidimensional modeling approach</i>
16:40	Social Program Visit the Informatorium at Szent-Györgyi Albert Agóra
19:00	Gala Dinner at the Rector's Building

Wednesday, July 2

Session 9	Artificial Intelligence
09:00	Gábor Nagy and Sándor Kazi: <i>Filtering Noise from Stock related Twitter Messages</i>
09:20	Sándor Kazi and Gábor Nagy: <i>A Scalable Parallel Boosting Scheme for Bulk Synchronous Parallel Environments</i>
09:40	Zsombor Paróczy, Bálint Fodor, and Gábor Szűcs: <i>Diversification for Content-Based Visual Information Retrieval System for Video</i>
10:00	Gábor Nagy, Sándor Kazi, and Győző Papp: <i>Distributed News Analytics Framework: Collecting News Feed Sources from Social Media</i>
10:20	Balázs L. Lévai: <i>Automatic Failure Detection and Monitoring of Ventilation and Cooling Systems</i>
10:40	Break

Session 10	Optimization
11:00	Balázs L. Lévai: <i>Automatic Design of LED Street Lights</i>
11:20	Dávid Szalóki: <i>Camera Placement Optimization in Object Localization Systems</i>
11:40	László Hajdu, Miklós Krész and Attila Tóth: <i>Two-phase graph coloring heuristic for crew rostering</i>
12:00	Viktor Árgilán, Balázs Dávid, János Balogh, József Békési and Miklós Krész: <i>Heuristic approach for the driver scheduling problem</i>
12:20	Balázs Dávid, János Balogh, Viktor Árgilán and Miklós Krész: <i>Integrating fast heuristics into a decision support system for vehicle rescheduling</i>

12:40	Lunch
-------	--------------

Session 11	Data mining
14:00	István T. Nagy: <i>Automatic Detection of Multiword Expressions with Dependency Parsers on Different Languages</i>
14:20	György Kovács and László Tóth: <i>The Joint Optimization of Spectro-Temporal Features and Deep Neural Nets for Robust ASR</i>
14:40	András London and Tamás Németh: <i>Application of graph based data mining techniques in administrative systems of education</i>
15:00	Gábor Imre and Gergely Mezei: <i>Graph Transformation-based Opinion Mining in Web Queries</i>

15:20	Closing
-------	----------------

PLENARY TALKS

Research, Innovation and Education in Infocommunication Technology

Zoltán Horváth, Zoltán Istenes, and Zsuzsa Várhalmi
Eötvös Loránd University, Hungary

The Eötvös Loránd University participates as the leader of the Budapest Associate Partner Group in the European Institute of Innovation and Technology (EIT) Information and Communications Technology Labs (ICT Labs). Researchers and also students from universities and industrial partners, more than two hundred researchers work on dozen innovative Research and Development projects. Their innovation potential is supported by national, European and EIT ICT Labs grants and funds. A joint project (EITKIC_12-1-2012-0001) of ELTE and BME – built on new basic research results – targeted the development of reliable technologies that help our daily lives in the increasingly developing world of ICT solutions and provides assistance for a wide range of users and IT enterprises, e.g. solutions to predict traffic jams, games that develop inductive reasoning, infrastructures to test mobile applications, solutions to defend critical systems, or an innovative speech synthesizer for future railways. The joint international Master School of EIT ICT Labs was launched with the participation of ELTE and BME; 10 students started their studies in two different Technical Majors at ELTE in September 2013. The EIT ICT Labs Budapest Doctoral Training Centre provides business development courses for PhD students of the ICT Doctoral Schools of ELTE and BME.

Transient stabilization of unstable states: why a stick balanced at the fingertip always falls?

John Milton
Claremont, California, US

The prevention of falls in the elderly, and their accompanying mortality and morbidity, are major challenges faced by aging industrialized societies. Increased risks of falling are associated with diseases of aging, such as diabetes, which enlarge sensory dead zones of peripherally located sensory receptors related to balance control. The sensory dead zone represents a range of the controlled variable over which no output is generated. Consequently mathematical models for balance control are posed as switching, or hybrid, type models in which the feedback control is turned on whenever the controlled variable exceeds a sensory threshold. From a dynamical point of view the presence of a dead zone represents a small-scale nonlinearity which does not affect large-scale linear stability, but can produce complex dynamics including limit cycle oscillations and micro-chaos. Here it is shown that for such systems it is possible that balance can be maintained for as long as minutes even though the feedback control is asymptotically unstable! Thus, for example, a stick balanced at the fingertip, an important laboratory paradigm for investigating the neural control of balance, can eventually fall no matter how skilled the stick balancer. These observations support the hypothesis that techniques which modify sensory dead zones, such as noisily vibrating insoles, may be useful to reduce the risk of falling in the elderly.

Computational homology

Pawel Pilarczyk
Institute of Science and Technology Austria

I am going to provide brief introduction to the subject of algorithmic homology computation of simplicial complexes and cubical sets, mainly limited to the theory and algorithms behind the Computational Homology Project software:

<http://chomp.rutgers.edu/software/>

In particular, I plan to discuss algorithms for the computation of homology of cubical sets and homomorphisms induced in homology by continuous maps, and to point out some reduction techniques that speed up the computations. I am also going to mention some applications of this machinery, and discuss a selection of related topics, such as persistent homology.

Mathematical modeling of max–min fair bandwidth allocation in BitTorrent communities

Elvira Antal and Tamás Vinkó

Data transfer in a BitTorrent community can be modeled as a network problem on a bipartite graph, where three types of nodes can be differentiated: the uploading peers (seeders and leechers), the downloading peers (leechers), and the content of the sharing sessions (torrents). The edges represent the peers' participation in the sessions, and bandwidth constraints are the weights of the edges.

The bandwidth allocation problem is setting feasible data flow values for every edge to achieve an optimal solution for a specific objective. In this case, we want to find a max–min fair bandwidth allocation. In other words, the minimum of the downloading data flows should be maximized. This type of allocation called “fair”, because no one's resources are allowed to be increased at the expense of decreasing the resources of any “weaker” peer. Max–min fairness could be a reasonable objective in direct media sharing, for example.

Capotă et al. [1] have given an iterative procedure for computing max–min fair bandwidth allocation, however, they have used a complicated and computationally intensive filtering phase in their algorithm. Our aim was to give an intuitive, clear mathematical model for the same problem, to help better understanding of the problem, verify the results of the earlier algorithm, and conduce to the development of a more efficient algorithm, which would serve to analyze the effects of structural alterations in large-scale BitTorrent communities.

Acknowledgements

This work was partially supported by the European Union and the European Social Fund through project FuturICT.hu (grant no.: TAMOP-4.2.2.C-11/1/KONV-2012-0013). T. Vinkó was supported by the Bolyai Scholarship of the Hungarian Academy of Sciences.

References

- [1] M. Capotă, N. Andrade, T. Vinkó, F. Santos, J. Pouwelse, and D. Epema. Inter-swarm resource allocation in BitTorrent communities, *In Proceedings of IEEE International Conference on Peer-to-Peer Computing (P2P 2011)*, pp. 300–309, 2011

Heuristic approach for the driver scheduling problem

Viktor Árgilán, János Balogh, József Békési, Balázs Dávid, Miklós Krész

A central problem of public transportation companies is to optimize their operational process. Since the minimization of the overall operational cost is a very complex task, the arising subproblems are considered as separated optimization problems. Vehicle scheduling assigns the trips of the input to vehicles that execute them, while driver scheduling creates the daily duties of the drivers based on the vehicle schedules.

In this presentation, we introduce a heuristic for the driver scheduling problem. Both vehicle- and driver scheduling are NP hard [2, 3]. An exact solution for our problem can be found in [1], where the problem is presented as a set-covering problem, and is solved by column generation. Using this approach, optimal solution for a real-life instance of a middle-sized city (approximately 5000 trips) can take up to several weeks, as the number of possible driver duty combinations generated by the pricing problem is too large. This way, the method cannot be applied in practice, so we developed a heuristic to speed up the generation process. The method partitions the trips of the input into two sets based on their starting time. Trips starting "too late" are selected, and the driver scheduling problem is solved using these as an input. Suitable driver schedules of the solution are saved, while trips of the other schedules are joined with the remaining partition. Trips of this resulting partition are divided using time windows, and the problem is solved for each individual time window separately in a sequential manner. Each resulting schedule of a phase is considered as a single trip in the future, and joined with trips belonging to the next time window. With the decrease of the number of trips, this approach decreases the size of our problem significantly, and allows the exact method to solve the smaller sub-problems in an acceptable running time.

Our new heuristic approach has been tested on real-life input. The test results show that running time significantly decreases, while the cost of the solutions remains close to the optimum.

Acknowledgements

This work was partially supported by the European Union and co-funded by the European Social Fund through project HPC (grant no.: TÁMOP-4.2.2.C-11/1/KONV-2012-0010).

References

- [1] Balogh, J., Békési, J. Driver scheduling for vehicle schedules using a set covering approach: a case study. *Proceedings of the 10th International Conference on Applied Informatics*, Accepted, 2014.
- [2] Bertossi, A.A., Carraresi, P., Gallo, G. On Some Matching Problems Arising in Vehicle Scheduling Models. *Networks* 17, pages 271-281, 1987.
- [3] Fischetti, M., Lodi, A., Martello, S., Toth, P. The fixed job schedule problem with working-times constraint. *Operations Research* 37, pages 395-403, 1989.

Integrating fast heuristics into a decision support system for vehicle rescheduling

Viktor Árgilán, János Balogh, Balázs Dávid, Miklós Krész

Transportation companies create their schedules in advance for a longer planning period. However, several unforeseen events can occur that render these pre-planned schedules infeasible. Companies have to deal with the disruptions almost immediately, and propose a new schedule that executes every disrupted event in a feasible manner once again. This process is called the vehicle rescheduling problem (VRP). Exact solution of this problem is not an option, because creating an optimal vehicle schedule with multiple different vehicle types is NP hard [1]. Because of these facts, fast heuristics have to be introduced.

Operators of a transportation company use their past experience in the case of such disruptions. Solutions given by fast methods for the VRP can be used as suggestions for the operators to help their planning process. In this talk, we introduce a decision support system for the VRP that provides multiple solutions with good quality in "quasi-real time". The system works independently of the solution methods, and because of this, new heuristics can be easily integrated into the framework. The best solutions of the methods are selected and given to the operator as suggestions. We also present some of our proposed heuristics from [2], and show their integration in such a system.

Acknowledgements

This work was partially supported by the European Union and co-funded by the European Social Fund through project HPC (grant no.: TÁMOP-4.2.2.C-11/1/KONV-2012-0010).

References

- [1] Bertossi, A.A., Carraresi, P., Gallo, G. On Some Matching Problems Arising in Vehicle Scheduling Models. *Networks* 17, pages 271-281, 1987.
- [2] Dávid, B., Krész, M. A model and fast heuristics for the multiple depot bus rescheduling problem. *10th International Conference on the Practice and Theory of Automated Timetabling (PATAT)*, 2014, Accepted.

Validation of the City Metaphor in Software Visualization

Gergő Balogh

The rapid evolution of computers made it possible to handle a large amount of information. New algorithms were invented to process data and new ways emerged to store their results.

However, the final recipients of these are still the users themselves, so we have to present the information in such a way that the human brain could understand it. One of the many possibilities is to convert the data into a graphical representation. This conversion is called visualization. Various kinds of method exist, beginning with simple charts through compound curves and splines to complex three-dimensional scene rendering. However, there is one point in common; all of these methods use some underlying model, a sort of language to express its content.

The increased performance of graphical units and processors made it possible and the data-processing technologies made it necessary to renew and to reinvent these visualization methods. In this research, we focus on the so called city-metaphor which displays information as buildings, districts, and streets.

Our main goal is to find a way to map the data to the entities in the fictional city. To allow the users to navigate freely in the artificial environment and to perceive the meaning of the objects, we have to find the difference between a realistic and an unrealistic city. To do this, we have to measure how much it is truth to reality or the city-likeness of our virtual creations. In this paper, we present four computable metrics which express various features of a city. These are compactness for measuring space consumption, eccentricity for describing the shape of the city, connectivity for showing the low level coherence between the buildings, and homogeneity for expressing the smoothness of the landscape. These metrics will be defined in a formal and an informal way, illustrated with examples. The connection among the high level city-likeness and these low level metrics will be analysed. Our preliminary assumptions about these metrics will be compared to the intuitions of users collected by an online survey. Finally, we will summarise our results and propose a way to compute the city-likeness metric.

References

- [1] Wettel, Richard and Lanza, Michele. CodeCity: 3D visualization of large-scale software, *Companion of the 30th international conference on Software engineering*, 2008
- [2] Balogh, Gergő and Beszedes, Arpad. CodeMetropolis-code visualisation in Minecraft, *Source Code Analysis and Manipulation (SCAM), 2013 IEEE 13th ICSM*, 2013

Language Support for High-level Worst-case Execution Time Estimation

Áron Baráth and Zoltán Porkoláb

There are many ways to implement one specific task, but some of them are faster than the others. In general we can specify an acceptable complexity that we can validate based on the source code or the object file. However, the validation requires external tool, which cause undesired dependency from a third-party utility. A better solution when the compiler provides language support for checking the expected worst-case execution time for each functions.

Validating the expected worst-case execution time should be an aspect of the correctness of the program, because it guarantees high-level correctness at design time. Therefore, the complexity of program will not exceed an expected value. It is important to note that, our goal is to give an expected upper limit for the complexity at *design time*.

In this paper, we first define the method to estimate the worst-case execution time of a function. This is a structural algorithm, which analyze e.g. loops, function calls and assignments, and detect correlations between elements to provide as accurate estimation as possible. We introduce our experimental language, which can perform worst-case execution time analysis, and the programmer can annotate the function with the expected time-complexity. This is a C-like language with strict syntax, strong static type-system, an powerful compile-time checks. The compiler can validate the given value with the estimated worst-case execution time, and gives an error when the complexity exceeds the expected value. Furthermore, the accepted values for the time-complexity will be used at runtime in the dynamic program loader. It is used to ensure that, the proper function is loaded with the correct computed worst-case execution time value.

References

- [1] Wilhelm, Reinhard, et al. "The worst-case execution-time problem – overview of methods and survey of tools." ACM Transactions on Embedded Computing Systems (TECS) 7.3 (2008): 36.
- [2] Huynh, Bach Khoa, Lei Ju, and Abhik Roychoudhury. "Scope-aware data cache analysis for WCET estimation." Real-Time and Embedded Technology and Applications Symposium (RTAS), 2011 17th IEEE. IEEE, 2011.

QR Code Localization Using Boosted Cascade of Weak Classifiers

Péter Bodnár and László G. Nyúl

QR code is a common type of visual code format that is used at various industrial setups and private projects as well. Its structure is well-defined and makes automatic reading available by computers and embedded systems. The recognition process consists of two steps, localization and decoding.

Belussi et al. [1] built an algorithm around the Viola-Jones framework [2], which proved that, even though the framework was originally designed for face detection, it is also suitable for QR code localization, even on low resolutions. The authors used a cascade of weak classifiers with Haar-like features, trained on the finder patterns (FIP) of the QR code. We extend their original idea, and propose improvements for both choosing the feature type of classifiers and their training target.

While Haar-like feature based classifiers are the state of the art in face detection, the training process is more difficult on FIPs. In order to increase the strong features of the object intended to detect, we propose training of a classifier for the whole code area. Even though QR codes have high variability on the data region, they contain data density patterns, a fourth, smaller FIP that can be perfectly covered with the center-type Haar-feature, furthermore, they contain the three discussed FIPs at more prominent location in the ROI.

Instead of Haar-like features, we also propose Local Binary Patterns (LBP) and Histograms of Oriented Gradients (HOG) for the feature evaluation. LBP and HOG based classifiers also can be trained both to FIPs and whole code areas, and since they are also considered as fast and accurate general purpose object detectors, evaluation of their performance on code localization is highly motivated. Furthermore, LBP can be more suitable than Haar classifiers, since it is not restricted to a pre-selected set of patterns, while HOG can also be efficient due to the strict visual structure and limited number of distinct gradient directions of the QR code.

We trained a total number of six classifiers, based on Haar-like features, LBP and HOG, both for FIPs and full code objects. For the FIPs, feature symmetry is also recommended to speed up the training process, while usage of the rotated features of Lienhart et al. [3] is not very useful, since these classifiers are not flexible enough to detect QR codes of any orientation. However, this issue can be solved by training two classifiers, for codes with orientation of 0° and 45° , respectively. We used a 32×32 sample size, which is larger than the one of the reference method, since training to the whole code object requires finer sample resolution. We decided cascade topology for the classifier instead of tree, since it showed higher overall recall in [1], and left required recall and false positive rate at the default values for each stage, with a total number of 10 stages. We trained our classifiers on a synthetic database consisting of 10,000 images, divided into 4:1 proportion for training and testing. Images of the database are artificially generated QR codes, each containing a permutation of all lower- and uppercase letters and numerals. QR codes generated this way are rendered onto images not having QR codes, with perspective distortion.

For the classifiers trained to FIPs, post-processing is needed to reduce the amount of false detections. Belussi et al. proposes searching through the set of FIP candidates for triplets that can form QR code, using geometrical constraints. Since real-life images of QR codes also suffer perspective distortion, it is obligatory to give tolerance values for positive triplet response. Contrary to this construction, our proposed classifiers trained to the whole code area need no post-processing.

Cascade classifiers are general purpose tools for object detection, and the discussed approach can be adapted to other two-dimensional code types as well. According to our experiments, cascade classifiers seem to be a decent option for QR code localization, especially a classifier using LBP for features and trained for the whole code object.

References

- [1] Luiz F. F. Belussi and Nina S. T. Hirata. Fast QR code detection in arbitrarily acquired images. In *Graphics, Patterns and Images (Sibgrapi), 2011 24th SIBGRAPI Conference on*, pages 281–288, 2011.
- [2] Paul Viola and Michael Jones. Rapid object detection using a boosted cascade of simple features. In *Computer Vision and Pattern Recognition, 2001. CVPR 2001. Proceedings of the 2001 IEEE Computer Society Conference on*, volume 1, pages I–511–I–518 vol.1, 2001.
- [3] Rainer Lienhart, Alexander Kuranov, and Vadim Pisarevsky. Empirical analysis of detection cascades of boosted classifiers for rapid object detection. In Bernd Michaelis and Gerald Krell, editors, *Pattern Recognition*, volume 2781 of *Lecture Notes in Computer Science*, pages 297–304. Springer Berlin Heidelberg, 2003.

A translation of interaction relationships into SMV language

Zsolt Borsi

In this paper a translation of a particular scenario-based model into SMV language will be presented. SMV is the input language of the NuSMV model-checker tool[4]. Model checkers in general provide a verification way to prove that a given system meets its specification. By using model checking, errors of the system can be detected even in the very early phases of the software development process.

The Unified Modeling Language (UML) provides diagrams to describe the same system from different aspects. The notion of Interaction Overview Diagram (IOD) was introduced in the second version of UML for specifying the relationships between UML interaction diagrams and the control flow passing between them. The notion of IOD is based on activity diagrams. An activity diagram is a directed graph, consisting of nodes and edges. Each node in an IOD is a reference to an interaction and the edges between activity nodes allow the definition of relationships between interactions. This paper does not want to take into account all the constructs included in IODs. On the other hand, some authors (including Whittle in [1]) found useful some additional construct, which are not part of the IODs. The construct included in this extended version of Interaction Overview Diagram (namely EIOD) will be considered in this paper.

The algorithm presented in this paper will translate a hierarchical construct containing EIODs into SMV. In the given construct there is an EIOD at the top level and each node of the top-level EIOD is refined at the underneath level by an EIOD. The algorithm is inspired by two previous works[2, 3]. The first one gives a translation of state charts into SMV. The notation of state charts are similar to IODs in the sense that any state in a statechart may contain whole statecharts like the nodes of EIODs are EIODs. The second one presents an algorithm for translating activity diagrams into SMV language.

The module concept of SMV provides the means for representing every EIOD in a separate module. SMV modules operate parallel to each other. The top level diagram will be represented by the main module in SMV. The model specification in SMV consists of three parts. First, the possible values of variables should be given to determine the space of states. Then the initial values of the variables and the transition relation should be defined. The state space of our model consists of boolean variables for each node of the top-level EIOD indicating whether the corresponding node is active. All nodes of the respective EIOD are instantiated in that module as well. The specification of these nodes are modelled in separate modules. The transition relation describes the mechanism that passes control along the hierarchical structure.

The idea of using model checking for verification is not new. The novelty of this paper is, that the base of the translation are IODs. Moreover, the transition takes into account additional constructs which are not part of UML, but are used by various authors.

References

- [1] J. Whittle, P.K. Jayaraman. Synthesizing hierarchical state machines from expressive scenario descriptions, *ACM Transactions on Software Engineering and Methodology (TOSEM)*, v.19 n.3, pp.1-45, 2010.
- [2] E.M. Clark, W. Heinle. Modular Translation of Statecharts to SMV, *Technical Report CMU-CS-00-XXX*, School of Computer Science, Carnegie Mellon University, Pittsburgh, PA 15213., 2000.
- [3] R. Eshuis. Symbolic model checking of UML activity diagrams, *ACM Transactions on Software Engineering and Methodology (TOSEM)*, v.15 n.1, pp1-38, 2006.

[4] NuSMV Model Checker Home Page, <http://nusmv.fbk.eu>

Variance of Source Code Quality Change Caused by Version Control Operations

Csaba Faragó

Motivation

Maintainability is one of the six sub-characteristics of software quality, as defined by the ISO/IEC 9126 standard [5]. Software maintenance consumes huge efforts: based on experiences, about half of the total amount of software development costs are spent on this activity. As maintainability is in direct connection with maintenance costs [2], we investigated the effect of a particular development process task – performing version control operations – on maintainability. Our goal was to explore typical patterns causing similar results in software quality change, which could either help to avoid software erosion, or provide information how to better allocate efforts spent on improving software quality.

Rich data about the developer actions can be found in the version control systems. Number of operations is the most obvious information we can gain for every commit, therefore we decided that we first examine the impact of these data on maintainability. We already performed research in this topic. First we found that there is a strong connection between the version control operations and the maintainability change of the source code [4]. Afterwards, we investigated the impact of the version control operations on maintainability change [3]. We showed that file additions have rather positive, file updates have rather negative effect on maintainability, while a clear effect of file deletions was not identified.

In this current paper we examine the variance of the maintainability changes caused by version control operations. We decided to check this for several reasons. First of all, if the net effect of one commit set is similar to the other's one, the difference in amplitudes could be important. The limited amount of efforts which could be spent on source code quality improvements could be better allocated to those commits which statistically cause higher amplitude. Eliminating the drastic maintainability decreases will result in net maintainability increase (similarly to the greenhouse effect). Second, discovering other dimension of the connections between version control operation and quality change could help in fine-tuning the results of the long-term research. We were especially interested in the variance caused by file deletions, as we did not identify clear impact of this operation on the maintainability previously. Finally, by discovering new connections other questions may raise. These potential new questions result in new research which might bring us closer and closer to the final goal: to create the formula of the developer interactions' impact on the quality of the source code.

Methodology

We took all available revisions of the source code of one industrial and three open-source projects. For each revision we recorded the following values: *number of each version control operation* (Add, Update, Delete), and *maintainability change* caused by that commit.

We estimated the maintainability of each revision by employing the ColumbusQM probabilistic software quality model [1]. Then we calculated the absolute maintainability change based on these values with transformations.

The variance tests are generally executed on two sets of numbers, with the null hypothesis that their variances are the same. Therefore we created two disjoint subsets of commits based on the number of version control operations in several ways. We examined all three operations one-by-one, and defined 7 combinations of divisions for every operation based on the existence, and absolute and relative medians. Then we considered the maintainability change values in

both subsets of every division as the input of variance tests. We executed the test itself with help of `var.test()` function of the R statistical program [6].

In the study we asked the following research question: *What is the impact of each operation (Add, Update, Delete) on the variance of maintainability change?*

Results

We executed the variance tests on one industrial (Gremon) and three open-source projects (Ant, Struts2, Tomcat). As result, we found clear connection between version control operations and the variance of the maintainability change. File Additions and Deletions caused significantly higher variance of maintainability change, compared to file Updates. Commits containing higher number of operations – regardless which operation it was – caused higher variance of maintainability change than those commits containing lower number of operations.

These results help us for better allocating the efforts spent on improving code quality. It is recommended to pay special attention on operations which could cause drastic change on maintainability. These are file additions and file deletions. There is nothing to do with deletions: if something needs to be removed, then it should be removed. On the other hand, it is recommended to pay special attention on new code. For example, it is recommended to mandate code review at least in case of new code development; or, if the code review is mandatory anyway, then it is recommended to do this in these cases with more strict rules. That the number of commits containing file Additions is relatively low considering all the commits, and this is especially true for commits containing file Additions exclusively.

Acknowledgements

This work was supported by Lufthansa Systems Hungary. The author would like to thank the help provided by Rudolf Ferenc and Péter Hegedűs in providing advices to this article.

References

- [1] Bakota, T., Hegedűs, P., Körtvélyesi, P., Ferenc, R., and Gyimóthy, T. A Probabilistic Software Quality Model. In *Proceedings of the 27th IEEE International Conference on Software Maintenance (ICSM 2011)*, pages 368–377, Williamsburg, VA, USA, 2011. IEEE Computer Society.
- [2] Bakota, T., Hegedűs, P., Ladányi, G., Körtvélyesi, P., Ferenc, R., and Gyimóthy, T. A Cost Model Based on Software Maintainability. In *Proceedings of the 28th IEEE International Conference on Software Maintenance (ICSM 2012)*, pages 316–325, Riva del Garda, Italy, 2012. IEEE Computer Society.
- [3] Faragó, Cs., Hegedűs, P., and Ferenc, R. The impact of version control operations on the quality change of the source code. In *14th International Conference on Computational Science and Applications (ICCSA 2014)*. Springer Verlag (accepted, to appear), June 2014.
- [4] Faragó, Cs., Hegedűs, P., Végh, Á. Z., and Ferenc, R. Connection Between Version Control Operations and Quality Change of the Source Code. *Acta Cybernetica (submitted)*.
- [5] ISO/IEC. *ISO/IEC 9126. Software Engineering – Product quality 6.5*. ISO/IEC, 2001.
- [6] R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2013.

Sum of Gaussian Functions Based 3d Point Cloud Registration

Bálint Fodor

Acquiring real-world 3d point cloud data is getting more and more accessible as RGB-D cameras (like Microsoft Kinect), Time-of-Flight cameras and other range imaging devices get cheaper. These products can also provide the depth information for the pixels of a captured frame. The 3d point set of a video frame is positioned in the coordinate system of the camera, which means that the range images acquired while moving the device will be represented in separate coordinate systems.

To get a comprehensive 3d model in the same coordinate system the frames should be transformed based on the pose information of the camera. In most of the cases however we do not have direct knowledge about the relative, nor the absolute translation and orientation of the camera.

We assume that subsequent frames contain overlapping parts of the same 3d object. Finding these parts the transformation between the range images can be estimated. The main problem is finding good matchings between points of the overlapping parts, while obtaining such a matching the rigid transformation between the parts can be easily calculated with well studied, closed form solutions[1].

Point cloud registration is one of the earliest problems risen in the 3d computer vision communities. Several approaches and formulations have been investigated, but due to the under-determined nature of the problem the topic has still relevance. A detailed categorization of the point registration problem can be found in the inspiring work of Tam et. al.[2].

The main idea of our proposed method is turning the point sets into continuous 3d vector-scalar functions. The score of a transformation between the functions is calculated by the space integral of the multiplication of the two functions. Of course the higher the score the better the transformation is. So the matching and transform estimation problems are together treated as a function fitting task.

It is important to note that we are not seeking for a global maximum during the fitting task, but aiming to get stuck in a local maximum where two overlapping regions are merging. Thus a fair initial transformation guess is essential.

For the two point sets we form corresponding continuous functions. It is done by replacing all points in a set with a 3d Gaussian function. The center of the Gaussian will be the original point. The resulting function for a point set will be a sum of Gaussian functions. Thus the points of a set $\{a_i\}$ are replaced by the function:

$$f_a(x) = \sum_{i=1}^n e^{-\frac{(x-a_i)^2}{w}} \quad (1)$$

where x is a 3d variable and w is a parameter of the Gaussian function that controls the 'slimness' of the function.

We define the rigid transformation as a rotation, translation and scale of the 3d point set x and denote with $t(x, p)$, where p is the parameter vector of the transformation. The score function for the parameters p is defined by the integral of the function multiplication:

$$S(p) = \int_{\mathbb{R}^3} f_a(x) \cdot t(f_b(x), p) dx \quad (2)$$

where f_a and f_b are the corresponding functions to the point sets $\{a_i\}$ and $\{b_j\}$.

Now the problem is to find a local maximum of $S(p)$ given an initial guess parameter p_0 . Fortunately, the equation can be simplified due to special properties of the Gaussian functions.

The proposed method was tested on synthetic data. Figure 1a. shows the same object in two different poses. The blue model is the original one and the green one is to be transformed. Both the models were treated as point clouds by taking the vertices of the objects. Each object consists of 505 vertices.

Figure 1b. shows the resulting transformation after 100 iterations of the steepest gradient optimization. The initial guess for the transformation was set to identity, so no translation, rotation or scale was considered. The test took 33s to run on a 2.4 GHz Intel Core i5 CPU without utilizing parallel computation.

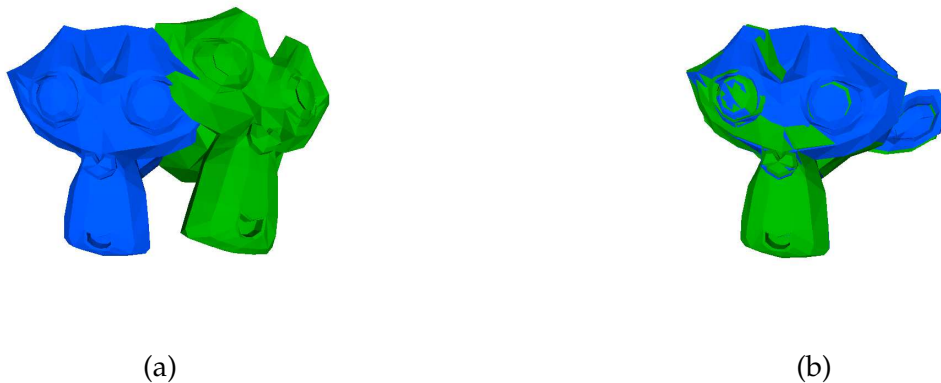


Figure 1: (a) initial poses of the models (b) final transformation

References

- [1] BKP Horn, HM Hilden, and S Negahdaripour. Closed-form solution of absolute orientation using orthonormal matrices. *JOSA A*, 5(July), 1988.
- [2] Gary K L Tam, Zhi-Quan Cheng, Yu-Kun Lai, Frank C Langbein, Yonghuai Liu, David Marshall, Ralph R Martin, Xian-Fang Sun, and Paul L Rosin. Registration of 3D point clouds and meshes: a survey from rigid to nonrigid. *IEEE transactions on visualization and computer graphics*, 19(7):1199–217, July 2013.

A case study of advancing remotely sensed image processing

Roberto Giachetta and István Fekete

Today's information systems present virtually unlimited opportunities for data analysis due to advances in algorithmic capabilities, evolution of parallel and distributed processing, leading to the increasing popularity of cloud computing. Due to data sets becoming increasingly large and complex (usually noted as *Big data*), these new paradigms cannot be ignored. In recent years many systems have been developed for data processing in the cloud, most notably the *MapReduce* model [1], and its open-source implementation, the *Apache Hadoop* framework [2]. However to take advantage of today's state of the art computing, previous data processing methodologies and workflows have to be revisited and redesigned.

In case of geographical information systems (GIS) and remotely sensed image analysis, the new paradigms have already been successfully applied to several problems [3], and systems have been developed to support processing of geographical or remotely sensed data in the cloud (e.g. *Hadoop-GIS* [4]). However in the case of complex processing operations the paradigm shift involves several considerations affecting both design and implementation. In contrary, most spatial data analysis processes performed at organizations such as the *Institute of Geodesy, Cartography and Remote Sensing (FÖMI)* have their evolved workflows using multiple (proprietary and open-source) software and GIS expertise.

For example, in the case of *waterlogging and flood detection*, the current practice involves a workflow with multiple steps with different software and including several parameters usually specified manually by a remote sensing expert. This is primarily due to the nature of waterlogging phenomena depending on the environment. For this problem, some specific algorithms have been introduced with environmental circumstances in mind [5, 6]. The process also involves multiple datasets including high-resolution remotely sensed imagery and digital elevation model (DEM). Thus the size of input is usually too large to be handled by a single machine, making the process worthwhile to be moved into cloud architecture. For this purpose, not only the distributed execution, but also complete automation is required. Hence, the advancement of the workflow to the state of the art requires both architectural, and algorithmic considerations.

In our paper, waterlogging and flood detection serves as a case study of shifting remotely sensed image processing from its current, semi-automatic, multi software environment to an integrated and fully automated distributed system based on the Hadoop framework.

The algorithm itself is enhanced by including segmentation based image analysis techniques based on reference data with auto tuning [7]. This approach enables the exclusion of intervention by the expert during the process, as the parameters are automatically tuned towards best result (with respect to reference data). The application of segmentation is not only an option, but a necessity in the processing of very high resolution images, as their pixels usually cannot be interpreted individually. Using reference data and automatic parameter calibration, the manual steps of the process can be eliminated, leading to a fully automated workflow.

To enable the distributed execution of the process, we rely on the proven Hadoop implementation, but extend its functionalities to enable a more flexible and controllable execution of the workflow and a more advanced management of resources. Instead of the direct approach to implement a MapReduce versions of the used algorithms, the distribution is accomplished by providing an environment, in which the process is automatically transformed to the MapReduce paradigm. For this purpose, we extend the capabilities of the *AEGIS spatio-temporal framework* [8], as it contains a generally defined, flexible processing system enabling operations to be executed in a variety of environments without any modification. This approach enables other operations to be executed in the Hadoop environment, not just our case study. However this technique also comes with the challenge of enhancing the data management possibilities of

Hadoop distributed file system, as the basic data splitting methodology does not allow execution of complex operations (such as segmentation and clustering) on the distributed data in a uniform manner.

In conclusion our approach advances the process of waterlogging and flood detection in both automation and efficiency using state of the art technology. It enables the replacement of multiple software to a single, generic framework, and it also paves the way for applying the paradigm shift to other GIS workflows.

References

- [1] Dean, J., Ghemawat, S.: MapReduce: simplified data processing on large clusters. *Communications of the ACM*, 51 (1), 107-113, 2008.
- [2] Agrawal, D., Das, S., Abbadi, A. E.: Big data and cloud computing: current state and future opportunities. *Proceedings of the 14th International Conference on Extending Database Technology (EDBT/ICDT '11)*, 530-533, 2011.
- [3] Almeer, M. H.: Hadoop Mapreduce for Remote Sensing Image Analysis. *International Journal of Emerging Technology and Advanced Engineering*, 2 (4), 443-451, 2012.
- [4] Aji, A. et. al.: HadoopGIS: A High Performance Spatial Data Warehousing System over MapReduce. *Proceedings of the VLDB Endowment*, 6 (11), 1009-1020, 2013.
- [5] Sreenivas, K., Dwivedi, R. S., Singh, A. N., Raviprakash, S.: Detection of sub-surface waterlogging using Terra-1 MODIS data. *Journal of the Indian Society of Remote Sensing*, 38 (1), 119-132, 2010.
- [6] El Bastawesy, M., Ali, R. R., Deocampo, D. M., Al Baroudi, M. S.: Detection and Assessment of the Waterlogging in the Dryland Drainage Basins Using Remote Sensing and GIS Techniques. *IEEE Journal of Selected Topics in Applied Earth Observation and Remote Sensing*, 5 (5), 1564-1571, 2012.
- [7] Dezső, B., Fekete, I., Gera, D., Giachetta, R., László, I.: Object-based image analysis in remote sensing applications using various segmentation techniques. *Annales Universitatis Scientiarum Budapest, Sectio Computatorica*, 37 (1), 103-120, 2012.
- [8] Giachetta, R.: AEGIS - A state-of-the-art spatio-temporal framework for education and research. *OSGeo Journal*, 13, 68-77, 2014.

Generalized Haar systems toolbox for MATLAB

Zoltán Gilián

The Haar system or Haar wavelet is a well-known mathematical tool proposed by the Hungarian mathematician Alfréd Haar with an abundance of applications in the field of signal and image processing. The system forms an orthonormal function basis in $L^p[0, 1)$ ($1 \leq p < \infty$), and the Fourier series of any continuous function with respect to this system converges uniformly to it. Furthermore the structural simplicity of the Haar wavelet makes it ideal for application in computer software.

Considering these notable properties it is desirable to generalize the Haar wavelet and construct different systems with similar properties [1]. A key concept in this regard is the notion of product systems. A well-known example is the relationship of the Haar, Rademacher and Walsh functions. Using this concept, it is possible to construct orthogonal systems of rational functions analogous to the Haar wavelet using the Malmquist-Takenaka system as a starting point [2] (see Figure 2). An advantage of rational functions over the Haar wavelet is that the former are analytic, while the latter is not even continuous.

In this paper we present a MATLAB toolbox for working with generalized Haar systems. The toolbox contains efficient implementations of signal transformation and reconstruction algorithms using various function systems. Two dimensional variants for image processing are also presented. The library is based on the RAIT rational approximation and interpolation toolbox [3].

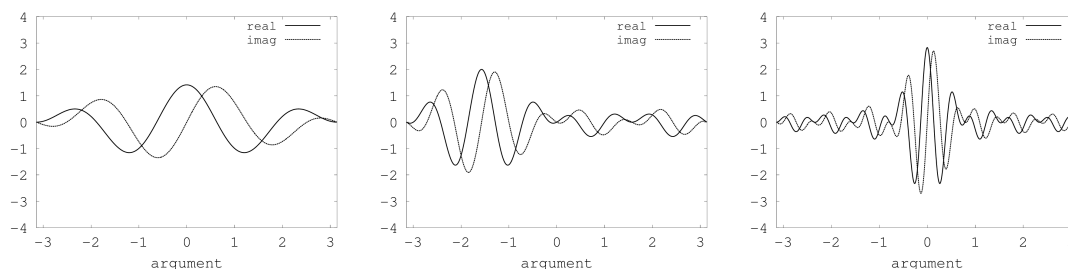


Figure 2: Rational Haar-like functions on the complex unit circle.

References

- [1] F. Schipp. On a generalization of the Haar system, *Acta Math. Acad. Sci. Hung.*, 33 (1979), 183-188.
- [2] F. Schipp. Rational Haar systems and fractals on the hyperbolic plan, *Sacks Memorial Conference*, Szentgotthárd, 2003, Oskar Kiadó.
- [3] P. Kovacs and L. Locsi. RAIT: The rational approximation and interpolation toolbox for Matlab, *Telecommunications and Signal Processing (TSP) 2012 35th International Conference*

Offline Signature Verification Using Similarity Measure for Skeletons

Erika Griechisch and Gábor Németh

Digital signature becomes more common for author identification, however handwritten signature still plays an important role in several aspects of life, e.g., business and bank sector. Offline signature verification methods analyse the images and shapes of the signatures. The main challenge in offline (and online) signature verification is the intra-class variability of the signatures: it is a known fact that every individual has variations in his or her own signatures. It means if someone writes two signatures after each other they will never be the same. Over a longer period of time the writing style may vary even more.

In this paper we present an offline signature verification method. Signature verification methods consist of preprocessing, feature extraction, and classification steps. As preprocessing segmentation was applied on each signature images to produce binarized signature image. In the feature extraction step the centerline was extracted for each signature images.

For pairwise comparison centerlines were normalized and registered to each other. As similarity measure we used one that Lee, Lam, and Suen proposed for quantitative comparison of 2D centerlines [1]. This measure has a value in range $[0, 1]$, where higher value represents higher similarity.

In the classification phase for each author the reference signatures are compared pairwise. One reference distance \mathcal{R}_i is calculated for author i , based on statistical values (i.e., mean and variance) of the pairwise distances. Then, each questioned signature of author i is compared to each reference signature. A questioned signature is accepted as genuine if its similarity is above $c \cdot \mathcal{R}_i$, otherwise it is rejected (considered as forged). For evaluation of the proposed method the multiplier c is varied in a range and false acceptance and false rejection rates are calculated as a function of this multiplier.

The proposed method has been evaluated on the publicly available SigComp2011 database and the results are competitive to the systems submitted to the concerning competition.

Acknowledgements

This work was supported by the MTA-SZTE Research Group on Artificial Intelligence.

References

- [1] S.-W. Lee, L. Lam, and C. Y. Suen, "A systematic evaluation of skeletonization algorithms," *Thinning Methodologies for Pattern Recognition*, vol. 8, pp. 239–261, 1994.

Two-phase graph coloring heuristic for crew rostering

László Hajdu, Miklós Krész, Attila Tóth

Nowadays, the companies and institutions have numerous employees, therefore the crew rostering problem became increasingly important. In most cases, some of the shifts take more time than others, which means that the employees don't spend the same amount of time working. These differences produced an "overtime cost" which is added to the basic salary. However, the companies and institutions must guarantee the basic salary for everyone, even if the employee does not spend the normal amount of time at work. This causes an additional cost for the companies.

The objective is to assign the crew members to shifts, meeting the constraints, and optimize the overall cost in such a way that the sum of the guaranteed basic salaries and the induced overtime cost is minimized. We improved a two-phase graph coloring method for the crew rostering. In the first step, a graph is built and colored, and in the second step, the graph is recolored with the tabu search method by our algorithm. The lower bound of the algorithm depends on the number of the employees and the working time. Our method has been tested with artificially generated and real inputs. For moderate size problems, the results of the new algorithm have been compared to the solutions of the appropriate integer programming model.

We obtained that our algorithm is able to handle relatively large inputs, and in the majority of the test cases, it has reached the theoretical lower bound with producing a satisfactory running time.

Acknowledgements

This work was partially supported by the European Union and co-funded by the European Social Fund through project HPC (grant no.: TÁMOP-4.2.2.C-11/1/KONV-2012-0010).

Eliminating Switching Components in Binary Matrices

Norbert Hantos and Péter Balázs

Let $A = (a_{ij})_{m \times n}$ be a binary matrix of size $m \times n$. The indices $1 \leq i_1 < i_2 \leq m$ and $1 \leq j_1 < j_2 \leq n$ form a *switching component* in A , if either $a_{i_1 j_1} = a_{i_2 j_2} = 1$ and $a_{i_1 j_2} = a_{i_2 j_1} = 0$, or $a_{i_1 j_1} = a_{i_2 j_2} = 0$ and $a_{i_1 j_2} = a_{i_2 j_1} = 1$. In other words, a switching component is a 2×2 sub-matrix of A with exactly two 1-s in its diagonal and two 0-s in its antidiagonal, or vice versa.

Switching components play an important role in image reconstruction and lossless image compression. The absence of switching components in the matrix is a necessary and sufficient condition for the unique reconstruction of the matrix from its horizontal and vertical projections, i.e., from the row and column sums of the matrix. Therefore, in that case the binary image represented by the binary matrix can be stored in a (lossless) compressed form by those two projections. However, if the matrix contains switching components, there is still a chance to reconstruct the matrix uniquely, if properly chosen elements of the matrix are stored as well. One can store, e.g., the positions of 0-s which need to be inverted into 1-s in order to make the matrix switching component free. These positions are called *0-1 flips*. Then, the aim is to find the minimal number of flips needed to achieve uniqueness. Unfortunately, the problem is generally NP-complete, thus there is no efficient exact algorithm to solve it, unless $P = NP$.

Switching components are also important in biogeography, where matrices represent the presence or absence of certain species (rows) on certain locations (columns). Here, the so-called nestedness is a relevant measurement of the matrix, which has a strong connection to the 0-1 flips.

In this paper we show that the minimal number of 0-1 flips can be found by determining the proper ordering of the columns regarding a certain filling function, rather than searching through matrix elements and switching components. Based on theoretical results, we develop two deterministic, polynomial-time heuristics to minimize the number of 0-1 flips. We compare those algorithms to another well-known ones in the literature, both on artificial random binary matrices and real-life biogeographical matrices. We conclude that the algorithms searching for proper column permutations perform better, both in the number of 0-1 flips and running time, especially on sparse matrices.

Acknowledgements

The work of Péter Balázs was supported by the OTKA PD100950 grant of the National Scientific Research Fund and by the European Union and the State of Hungary, co-financed by the European Social Fund in the framework of TÁMOP 4.2.4.A/2-11-1-2012-0001 'National Excellence Program'. The research of Norbert Hantos was supported by the Research Group on Artificial Intelligence (RGAI), belonging to the Hungarian Academy of Sciences (MTA) and University of Szeged (SZTE) Informatics Departments.

Structural Information Aided Automated Test Method for Magic 4GL

Ferenc Horváth, Richárd Dévai, Tamás Gergely

Nowadays testing data intensive, GUI enhanced applications (like those designed with Magic 4GL) properly on an easily maintainable way has become a more crucial part of the application life cycle. For 3GL applications there are many evolving technologies to support automatized GUI testing for web applications and also for their desktop counterparts.

The first generation GUI tester tools were layout dependent and more like simple recorder and player tools which could replay the users' earlier actions. The main drawbacks of this solution were the poor verifiability and the sensitiveness to even the smallest modifications of the layout. Also, these tools needed to be controlled by user interactions.

In the next generation the tools got more sophisticated by supporting the identification of GUI elements in layout independent ways. This identification allowed the record or generation of test cases, as the interactions on the GUI elements took place based on a general model instead of layout positions. Based on previous works [2] and [3] we aimed to create a new test generator system. We created a new path and test script generator tool which can provide a sophisticated and simple way to create interpretable test cases for the used test execution tool.

Path generation is a very complex task because one has to face obstacles like the exponential relation between the number of branches and the number of the possible paths, or the presence of loops in every real program which induces the number of paths to converge to infinity [1]. Therefore robust algorithms are needed.

Our path generation strategy is based on well known algorithms like breadth-first and depth-first search. Considering the time and space complexities, these algorithms perform well enough to allow us to build upon them while developing our advanced strategies. However, we have to introduce some modifications to meet the requirements of our domain. These modifications affect the generation process in numerous ways. For example, we created an interface that can be used to build a business model for the program being tested. This model defines constraints which will be used by the generator algorithm as a form of selection method to narrow down the set of used program subcomponents (e.g. tasks which represent main functional components in Magic 4GL). This way the set of the generated paths can be reduced significantly and contains only the most relevant ones.

Path generation strategy is a crucial part of the test script generation. After the selection of the appropriate path set, we have to convert the information we gained into interpretable test scripts. The number of the generated paths has a strong influence on the final number of test scripts as well as on the strategy we use to determine how we select and permute values among different variables.

The co-domain for each variable is based on its type and control flow context. After we extract the influencing data properties (D/U information), we have to narrow down the co-domain of a variable based on the influencing expressions and statements which appear on the corresponding path. Value selection on co-domains is also a crucial part of the generation as it can cause an enormous growth of the number of test scripts. Choosing a good strategy for solving the problems we listed above is essential in order to get an accurate and usable test script set as the output of the described script generator algorithm. Only a small enough test set size can be manageable even if the pre-configuration of the Magic xpa program and the execution of the scripts are totally automated.

We present a schematic model of the system we built to test Magic xpa applications in Figure 3. This system is based on the work of Dévai et al. [2] to gain usable behavioral information about the program under test and Fritsi et al. [3] to run our assembled scripts. We connected

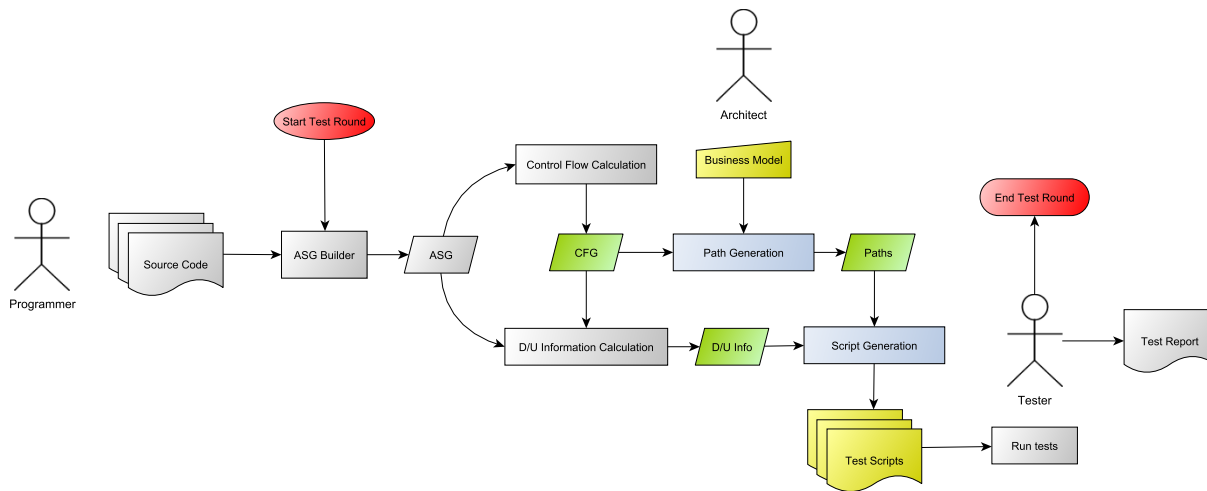


Figure 3: Flowchart of Magic test process

this two components by developing a path and a script generator that process behavioral information and assemble executable scripts for the test runner.

In conclusion, we would like to present our path and script generator algorithms and their implementations for Magic xpa applications. In addition, we intend to demonstrate how these components cooperate with the existent solutions and as a possible use of the results we introduce a test method that has been completed by the application of our path and script generator tools.

Acknowledgements

This research was supported by the Hungarian national grant GOP-1.1.1-11-2011-0039.

References

- [1] T. Bakota, Á Beszedes, T. Gergely, M. Gyalai, T. Gyimóthy, and D. Füleki. Semi-automatic test case generation from business process models. In *Proceedings of the 11th Symposium on Programming Languages and Software Tools (SPLST'09) and 7th Nordic Workshop on Model Driven Software Engineering (NW-MODE'09)*, pages 5–18, Tampere, Finland, August 26-28 2009.
- [2] Richárd Dévai, Judit Jász, Csaba Nagy, and Rudolf Ferenc. Designing and implementing control flow graph for Magic 4th generation language. In Ákos Kiss, editor, *Proceedings of the 13th Symposium on Programming Languages and Software Tools (SPLST'13)*, pages 200–214, Szeged, Hungary, 2013. University of Szeged.
- [3] Dániel Fritsi, Csaba Nagy, Rudolf Ferenc, and Tibor Gyimóthy. A layout independent GUI test automation tool for applications developed in Magic/uniPaaS. In *Proceedings of the 12th Symposium on Programming Languages and Software Tools (SPLST'11)*, pages 248–259, 2011.

Usability Testing of Android Applications

Ferenc Horváth, Benjamin Mészáros, Tamás Gergely

Nowadays the quality of a company is determined mainly by the quality of the software it uses, thus it has become more important to obtain high quality software. Therefore quality assurance plays a central role in the software industry.

In addition, if we observe the importance of the certification steps in Microsoft's Windows Store application publishing process [1], or we notice the fact that because of various quality concerns Google has permanently restricted the downloading of about a hundred and fifty thousand apps which were available in the official Android Play Store in mid-2013, then we can see that quality assurance is not only essential for the desktop applications, but it is more and more crucial in the constantly growing segment of the portable devices.

Software quality is a complex attribute. It has many aspects so numerous qualitative and quantitative characteristics should be taken into consideration during its calculation. Landauer [2] states that a significant part of software bugs are related to some kind of usability problems, so we can suppose that measuring applications from a usability point of view will give us an insight into the overall quality of the actual software.

In general there are several solutions which help to create qualitative applications and to qualify our existing applications, but in case of Android the situation is far from being optimal. Although there are some automatized tools that can help to determine the goodness of Android applications, there are hardly any software that can provide a comprehensive assessment regarding the usability of an application.

In this paper we would like to present a method that can help to solve the problem described above. We created a framework that helps to analyse the usability of android applications.

To assess usability we had to be able to log the relevant interactions of the users on the graphical user interface, thus, we developed a tool which injects logging methods into the applications by instrumenting them, but without modifying the original functionality of the programs. For example, it partially modifies the inheritance tree and the event handling mechanism, as also it replaces some of the UI descriptions, but preserves the original behaviour.

We applied this tool on a few Android apps and for each of them we defined use case scenarios that respects the features of the actually tested application. After we executed these use cases we got several log files that served as a base for the further investigations. We analysed these logs and based on our observations and the experiences of our usability experts we defined metrics (e.g., navigation anomaly, misclick, etc.) and a model that can provide sophisticated information about usability.

Acknowledgements

This research was supported by the Hungarian national grant GOP 1.1.1-11-2011-0006.

References

- [1] Microsoft Inc. Overview of publishing an app to the windows store, May 2014.
- [2] Thomas K. Landauer. *The Trouble With Computers: Usefulness, Usability, and Productivity*. MIT Press, Cambridge, Massachusetts, 1995.

Test-driven verification of model transformations

Péter Hudák and László Lengyel

Model-driven engineering provides a numerous number of fast and efficient solutions on the field of software engineering. The focused problem is decomposed to different abstraction levels. This approach helps understanding and modeling the problem to be solved. As a result, software models represent the essence of the actual problem and show the critical aspects of it.

Different planning and development phases of a software project produces models belonging to different abstraction levels. Model processors are dedicated to map the different abstraction levels. A set of rules what are applied in a defined order on source models and result target models are called transformation. Model-driven engineering utilizes this method widely. One of the current significant fields is the software development for cyber-physical systems (next generation of embedded systems).

Both syntactically and semantically appropriate outputs should be generated by model processors. This means that the resulted software artifact has to contain the correct semantic information. This is the point where test-driven verification of model transformations gives proper feedback to check the correctness of model processors. Testing of model transformations is a challenging and complex problem. Design errors are usually the main source of defects in model transformations. Complexity of software artifacts and models make searching for bugs and errors even harder. The goal of test-driven method is to test model transformations by automatically generating right input models. Transformations execute these input models. Semantic correctness verification is made after the execution via comparing the input and output models. The method is hard to be fully automatized because additional information about the models are usually required by the process. Models tend to change continuously during software development so model transformation tests should adapt to these changes. Model transformations should conform to predefinitions so coherent test cases could be generated in this way. This set of tests is handled then as unit tests.

Visual Modeling and Transformation System (VMTS) is a domain-specific modeling and metamodeling framework, which supports defining domain-specific languages, editing models, furthermore, defining and performing model-to-model and model-to-code transformations. The latest version is currently under development, we are working on the extension of the model transformation engine. A new automatic test generation solution will be added that supports test-driven verification of model transformations. This framework generates a set of input models. These input models function as unit tests therefore unit tests could be generated for model transformations which fill predefinitions. In this paper we give an overview about the test-driven verification methods of model transformations, furthermore, we demonstrate the test generation engine of VMTS 4.

Graph Transformation-based Opinion Mining in Web Queries

Gábor Imre and Gergely Mezei

In our research we evaluate the possible uses of opinion mining algorithms in a web-based graph transformation system. Text analysis algorithms, such as opinion mining [1], can be applied within graph transformations to provide effective analyzing and processing capabilities for interpreting public web data.

The web provides tremendous amounts of data available in various formats. Some data providers publish their data in a semi-structured format, where plain text and formally defined data is mixed. A typical example is the Stack Exchange API¹, where huge amounts of Q&A forum data is available. The data is semi-structured: it has a formal structure for tags, scores, connections between questions and answers, etc.; yet the text of the questions and answers is available as plain text. Text mining algorithms can be used to process the plain text, while graph transformations are strong in processing structured data. We examine the possibilities and requirements of applying opinion mining algorithms as parts of graph transformations.

Graph transformations are widely used in various industrial and scientific tasks in Model-Driven Architecture [2]. Graph transformations over typed attributed graphs [3] are especially strong in practical representation of object oriented structures, in addition they have strong mathematical background. Combining graph transformations with text analysis algorithms can lead to powerful analysis applications.

We have constructed a case study that demonstrates the usage of opinion mining through graph transformations over semi-structured data. We analyze the questions and answers available at Stack Overflow², the world's greatest programming Q&A site. Some of the most important Linux distributions are compared based on the question and answer plain texts. The sentiment of a question or answer is extracted by a simple, dictionary-based opinion mining algorithm. If the user opinions are generally more positive, then the regarding distribution is considered more comfortable. This simple heuristics can be used to compare similar technologies based on the general user opinions.

Acknowledgements

This work was partially supported by the European Union and the European Social Fund through project FuturICT.hu (grant no.: TAMOP-4.2.2.C-11/1/KONV-2012-0013) organized by VIKING Zrt. Balatonfüred; and by the Hungarian Government, managed by the National Development Agency, and financed by the Research and Technology Innovation Fund (grant no.: KMR 12-1-2012-0441).

References

- [1] B. Pang and L. Lee, "Opinion mining and sentiment analysis," *Foundations and trends in information retrieval*, vol. 2, pp. 1-135, 2008.
- [2] "The Model-Driven Architecture", <http://www.omg.org/mda/>, OMG, Needham, MA, 2002
- [3] H. Ehrig, U. Prange and G. Taentzer, *Fundamental theory for typed attributed graph transformation*, Springer Berlin Heidelberg, 2004.

¹<https://api.stackexchange.com/docs/>

²<https://stackoverflow.com>

Comparison of Source Code Transformation Methods for Clang

Máté Karácsony

Clang is a robust, industrial strength C, C++ and Objective-C compiler. Its library based architecture made it a common platform for custom tools which are analyzing and transforming programs. As it provides a rich abstract syntax tree representation, it is widely used especially for source-to-source transformations, like vectorization of loops [1], or translation between different language dialects [2]. Despite the presence of these tools, Clang currently does not have a complete and uniform library to support specifically source-to-source transformations.

Conventional transformation engines are working by manipulating the abstract syntax tree (AST) of the input source code, which is generated by a parser. A particular transformation is usually described by a set of rules. Application of a single rule is performed by replacing or mutating matching AST nodes. These rules are applied according to a strategy to create a modified AST. The output of the transformation process is computed by pretty-printing this structure into its source code representation [3].

Clang provides different methods to traverse and match specific parts of AST, but it does not provide an interface to define and apply transformation rules in a straightforward way. While its AST is designed to be immutable, classical AST-based transformations are simply not eligible. Moreover, pretty-printing the whole transformed syntax tree is not practical in the case of preprocessed languages, since all macros will be expanded in the output, which impedes further code maintenance. For these reasons, Clang incorporates a simple source code rewriting facility, that is able to replace source text in the given ranges. This low-level API can be used to safely change small text fragments in input source files, but it is not able to handle complex cases.

As these built-in transformation capabilities are limited, every tool implementation created its own procedures and utilities for this purpose. In this paper we will present and compare various existing transformation techniques, including the built-in options and custom solutions found in Clang-based tools.

Three key aspects of these methods will be examined: expressiveness, implementation cost and reusability. Expressiveness shows what kind of transformations can be described, and how these are represented. The second aspect measures the effort needed to implement the given technique. Finally, reusability expresses how easy is to apply the specific method to achieve different kinds of transformations.

References

- [1] Krzikalla, Olaf, et al. "Scout: a source-to-source transformator for SIMD-Optimizations." Euro-Par 2011: Parallel Processing Workshops. Springer Berlin Heidelberg, 2012.
- [2] Martinez, Gabriel, Mark Gardner, and Wu-chun Feng. "CU2CL: A CUDA-to-OpenCL translator for multi-and many-core architectures." Parallel and Distributed Systems (ICPADS), 2011 IEEE 17th International Conference on. IEEE, 2011.
- [3] Van Wijngaarden, Jonne, and Eelco Visser. "Program transformation mechanics. a classification of mechanisms for program transformation with a survey of existing transformation systems." Technical report UU-CS 2003-048 (2003).

Fast recognition of natural feature identifiers by a mobile phone

Melinda Katona and László G. Nyúl

As privacy issues are becoming increasingly important, concealing the identities of individual persons or objects is essential for the analysis of the mass amount of data captured by cameras or other means about our world, that is full of artificial and natural identifiers. For reliable anonymization or de-identification, techniques are required to automatically recognize identifying features, markers, patterns, and to manipulate the data so that the content can still be used for the intended purposes without any privacy issues.

When coming to the automatic identification/recognition of objects, algorithms are needed to automatically locate and decode the identifiers attached to the objects. Very different techniques are required for natural biometric IDs, such as fingerprints, iris or retina patterns, and for artificial IDs, such as barcodes or QR codes [1, 2].

Barcode technology is the pillar of automatic identification, that is used in a wide range of real-time applications with various types of codes. The different types of codes and applications impose special problems, so there is a continuous need for solutions with improved effectiveness. Barcode localization methods have two objectives, speed and accuracy. For industrial environment, accuracy is crucial since undetected (missed) codes may lead to loss of profit. Processing speed is a secondary desired property of the detectors. On smartphones, the accuracy is not so critical, since the device interacts with the user and re-shoting is easily possible, but a fast (and reasonably accurate) barcode detection is desirable.

In this paper, we focus on the automatic localization and recognition of a kind of natural feature identifier (NFI). We present an algorithm that successfully locates NFI code region in an image taken by a camera, extracts features of the NFI that can be the basis for recognition or matching. We show our preliminary experimental results using a moderate set of labels and images.

Acknowledgements

Research of Melinda Katona was supported by the European Union and the State of Hungary, co-financed by the European Social Fund in the framework of TÁMOP 4.2.4.A/2-11-1-2012-0001 'National Excellence Program'.

The research reported here was financed by InterSoft Hungary Ltd. within an R&D project and all rights to commercial use of the resulting technology have been transferred to the sponsoring firm.

References

- [1] L. Belussi and N. S. T. Hirata, "Fast qr code detection in arbitrarily acquired images," in *Graphics, Patterns and Images (Sibgrapi)*, 2011, pp. 281-288
- [2] I. Szentandrási, A. Herout, and M. Dubská, "Fast detection and recognition of qr codes in high-resolution images," in *Proceedings of the 28th Spring Conference on Computer Graphics*, 2013, pp. 129-136

A scalable parallel boosting scheme for bulk synchronous parallel environments

Sándor Kazi and Gábor Nagy

Under the pressure of the several V-s (volume, velocity, variety, etc.) of BigData [1] a common approach is to use distributed computation frameworks on top of distributed filesystems. The open source flag-carrier of this approach is Hadoop which has transformed from an open-source project to a widely applied business solution in the last few years. Machine learning possibilities on top of Hadoop is basically identified by the somewhat limited capabilities of Mahout, a machine learning library for Hadoop, and the programming capabilities of a user. However MapReduce has disadvantages when it comes to iterative task execution [2]. Therefore a large set of machine learning algorithms call for a different approach. A possible way to implement distributed iterative machine learning algorithms on top of the Hadoop infrastructure (overcoming the limitations mentioned before) is to use the BSP computation model (designed by Leslie Valiant [3] in the early '90s and revisited by Google in 2010 [4]). This model is supported by two Hadoop packages now: Apache Giraph (mainly for graph processing) and Apache Hama. If any dimension of the data (including its velocity, etc.) is too much for one node to handle these distributed frameworks can provide scalable parallel implementation possibilities.

We hereby present a scheme to create distributed versions of a boosting algorithm using the BSP model. The idea of boosting comes from the task of training a set of weak learners to form a strong learner, a popular representative of this meta modeler group is Gradient Boosting which uses a gradient based method to calculate the new labels for each weak learner to use for training. One of the most (or the most) common weak learners used with gradient boosting are Decision Trees, the abbreviation GBDT refers to this construction [5, 6]. To introduce the scheme we use GBDT in our demonstrations, although the approach can be similarly applied for some other boosting algorithms.

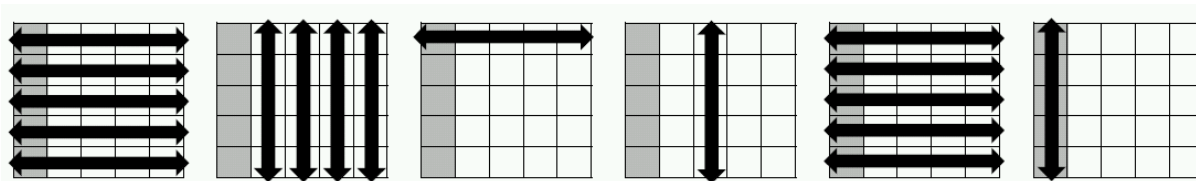


Figure 4: A schematic representation of the parallel GBDT training process.

The data is presumed to be distributed among multiple nodes, several factors are considered to ensure efficient execution of the GBDT training. Some of the nodes should have one of the two special roles, the role setting can be mapped to the data distribution setting. A schematic representation of the parallel GBDT training process is presented on figure 4: each row/column represents nodes having parts of the same data row/column. Residuals are distributed in the first step, statistics are gathered in each column in the second, the best split for each column is broadcasted and the global best split is selected afterwards. Then the leaf is updated using the new split description at the nodes having that column. If the tree is not completed yet we restart from statistics gathering (but for the new leaves only) otherwise the model is updated with the current tree (with appropriate multipliers calculated for each leaf). Every new iteration of GBDT starts with a new residual calculation step. Most of these can be executed in parallel and can be efficiently balanced to minimize waiting times by (even dynamic) distribution of the data.

The efficiency of the parallel setting can depend on the properties of the dataset, the distribution setting and the loss function. Communication costs in this scheme can be and should be minimized, this can be done in several ways. Some of these methods can guarantee the model to be the same as it would be in a non-distributed approach, some others (like the histogram setting of Ben-Haim and Tom-Tov [7]) are approximations but operate without reasonable degradation in efficiency [8].

Acknowledgements

This work is supported by the grant: FUTURICT, TÁMOP-4.2.2.C-111KONV, "Financial Systems" subproject.

References

- [1] D. Laney (2001-02-06). *The Importance of 'Big Data': A Definition*. Gartner.
- [2] K. Lee, Y-J. Lee, H. Choi, Y. D. Chung, and B. Moon. Parallel data processing with mapreduce: a survey. *SIGMOD Rec.*, 40(4):11–20, January 2012.
- [3] L. G. Valiant, "A bridging model for parallel computation," *Commun. ACM*, vol. 33, no. 8, pp. 103–111
- [4] G. Malewicz, M. H. Austern, A. J. C. Bik, J. C. Dehnert, I. Horn, N. Leiser, and G. Czajkowski, "Pregel: a system for large-scale graph processing," in *Proceedings of the 2010 ACM SIGMOD International Conference on Management of data*, ser. SIGMOD '10. New York, NY, USA: ACM, 2010, pp. 135–146.
- [5] J. H. Friedman, "Stochastic gradient boosting," *Computational Statistics and Data Analysis*, vol. 38, pp. 367–378, 1999.
- [6] J. H. Friedman, "Greedy function approximation: A gradient boosting machine," *The Annals of Statistics*, vol. 29, no. 5, pp. 1189–1232, 10 2001.
- [7] Y. Ben-Haim, E. Tom-Tov, "A Streaming Parallel Decision Tree Algorithm," *The Journal of Machine Learning Research*, vol. 11, 3/1/2010, pp. 849–872
- [8] S. Tyree, K. Q. Weinberger, K. Agrawal, and J. Paykin, "Parallel boosted regression trees for web search ranking," in *Proceedings of the 20th International Conference on World Wide Web*, ser. WWW '11. New York, NY, USA: ACM, 2011, pp. 387–396.

The Joint Optimization of Spectro-Temporal Features and Deep Neural Nets for Robust ASR

György Kovács and László Tóth

One of the biggest challenges facing automatic speech recognition is to get an acceptable performance even in an adverse environment, e.g speech with background noise. One way of increasing the robustness of ASR systems is to apply spectro-temporal processing on the speech signal [1]. In this approach, the features for the subsequent classification are got by processing small, spectrally and temporally localized patches of the spectrogram. Hence, unlike in traditional processing methods, some of the features may be unaffected by the noise component, making the whole feature set more robust. The spectro-temporal processing of the patches can be performed by using the two-dimensional discrete cosine transform (2D DCT) or Gabor filters. Good recognition results were reported with both of these approaches earlier [1, 2, 3].

After the initial step of feature extraction, the features are passed on to a machine learning algorithm – in most cases a Hidden Markov Model (HMM) or an Artificial Neural Net (ANN). Normally, the feature extraction and recognition steps are performed separately. In our previous paper, we showed that the spectro-temporal feature extraction step and the ANN-based recognition step can be integrated, and in this way the parameters of the two phases can be trained together [4]. Our solution was based on the observation that the spectro-temporal filters can be treated as special types of neurons, and so the standard backpropagation training algorithm of ANNs can be extended to the feature extraction step as well. We experimented with neurons that simulated three types of spectro-temporal feature extraction methods. In the first case, a set of 2D DCT filters was applied; in the second we applied Gabor filters; while for the third configuration we used randomly generated filters. What we found was that in each case, our integrated method enhanced the performance of the filter sets by extending the scope of the backpropagation algorithm to the neurons simulating them.

In this study, we improve our system further by incorporating recent advances in neural networks into it. In the standard ANN implementations there are three layers, namely an input layer, an output layer (applying the softmax nonlinearity), and in between a hidden layer that uses a sigmoid activation function. Recently, it has been shown that significant improvements in performance can be achieved by increasing the number of hidden layers [5]. Unfortunately, training these ‘deep’ networks with many (three or more) hidden layers using the classic backpropagation algorithm has certain problems associated with it. A solution to these problems was given by Hinton et al., leading to a renaissance of ANN-based technologies in speech [5]. An even simpler solution was later given with the introduction of rectifier neural networks [6]. Here, we apply the latter in combination with deep neural nets to the model introduced in our earlier paper [4]. By evaluating our system in phone recognition tasks on the widely used TIMIT speech database, we will show that these techniques allow us to further improve the model performance in the case of clean speech and also noise contaminated speech. It is important to note that the training and testing was performed on different portions of the database, and that to obtain phone recognition results in noisy environment, we did not do any additional training of the models, but used the models that were trained on clean speech data.

Acknowledgements

This publication is supported by the European Union and co-funded by the European Social Fund. Project title: Telemedicine-focused research activities in the fields of mathematics, informatics and medical sciences. Project number: TÁMOP-4.2.2.A-11/1/KONV-2012-0073.

References

- [1] Kleinschmidt M. Robust Speech Recognition Based on Spectro-Temporal Processing. PhD thesis at Carl von Ossietzky University, Oldenburg, Germany, 2002.
- [2] Bouvrie, J., Ezzat, T., and Poggio, T. Localized spectro-temporal cepstral analysis of speech. *Proceedings of ICASSP* pp. 4733–4736, 2008.
- [3] Kovács, G., Tóth, L. Phone Recognition Experiments with 2D DCT Spectro-Temporal Features. *Proceedings of SACI*, pp. 143–146, 2011.
- [4] Kovács, G., Tóth, L. The Joint Optimization of Spectro-Temporal Features and Neural Net Classifiers. *Proceedings of TSD*, pp. 552–559, 2013.
- [5] Hinton, G. et al. Deep Neural Networks for Acoustic Modeling in Speech Recognition. *in IEEE Signal Processing Magazine*, pp. 82–97, 2012.
- [6] Tóth, L. Phone Recognition with Deep Sparse Rectifier Neural Networks. *Proceedings of ICASSP*, pp. 6985–6969, 2013.

Methods for Feature Point Aggregation of Optical Flow for Orientation Estimation

László Kundra

Estimation of precise orientation using only inertial sensors is a basic problem of the field. Several algorithms have been developed and extensively discussed regarding bias compensation, sensor fusion with complementary and Kalman filters, and so. Although there are impressive results for specific application domains, their restrictions that make them feasible are not generalizable. In this research, inertial measurement unit (IMU) was aided with optical flow based feature displacements from an on-board camera, to calculate a robust and reliable orientation. Such situations when a camera is available is quite frequent. (Mobile phones, unmanned aerial vehicles (UAVs), etc.) However, calculating optical flow on whole images will result a large amount of vectors, representing the addition of physical displacement and rotation. Using algorithms that are beyond the scope of this paper, these vectors can be assigned a reliability value, and physical displacement can also be compensated. Having the orientation change as a set of vectors, an appropriate method can aggregate them into one representative angle. Using these angles and angular velocity in combination with IMU data, a more robust and reliable orientation can be estimated.

The purpose of this paper is to introduce and compare different methods for aggregation on both simulation and real world data. First, the most straightforward ways are examined to aggregate, like mean and median of the vectors. To avoid enormous errors brought into the calculation by outlier points, when calculating the mean these points are rejected. Median by its definition is insensitive for such perturbations, however the central element is not always a good selection. For this, weighted medians has also been considered. Regular weighted median assigns a weight to each item, and by doing so, the central position may point to an item different from the median. An other inspected modification to the median is a method using coefficients to sum the ordered list. In this case the sum of the coefficients should be 1. An other approach to find the most plausible angular rate is to find where the density of these vectors is maximal. In this case, we suppose that there are some vectors that have nearly the same representation of the rotation, thus lying densely to each other, while deceiving vectors are spread in the range. Finding the location of this maximum can be achieved in various way. One method is to count nearby elements (in an ϵ range) for each input point. Where the density is high, this function will have a large output. An other technique is to take the sum of Gaussian functions fit to each point. The density can be calculated the same way as for the previous method. These methods all need some parameters, or parameter vectors, thus these can be calculated based on the nature of the application field.

Measurements were performed on mobile devices (Samsung Galaxy S2 and S4, with different camera resolution and feature point count). Using measurement data, we performed several comparison and tuning of parameters, and for orientation estimation we could produce very low angular rate bias, that led to negligible drift over time.

Acknowledgements

This work was partially supported by the European Union and the European Social Fund through project FuturICT.hu (grant no.: TAMOP-4.2.2.C-11/1/KONV-2012-0013) organized by VIKING Zrt. Balatonfured. This work was partially supported by the Hungarian Government, managed by the National Development Agency, and financed by the Research and Technology Innovation Fund (grant no.: KMR-12-1-2012-0441).

Business process quality measurement using advances in static code analysis

Gergely Ladányi

Business process is a collection of activities that takes one or more kinds of input and creates an output that is of value to the customer[1]. Since it has direct effect on the customer, it is really important to have an up-to-date knowledge about its weak and strong points. Using flow charts like Event-driven process chain (EPC) is a widespread solution for modeling, analyzing, and redesigning business processes. EPC is commonly used because models described with it are flexible, easy to manage and understand. More precisely it is an ordered graph of events and functions with various connectors (AND, OR, XOR) that allow alternative and parallel execution of processes. In this paper we focus on the quality of an EPC and then on the extension of this property to the entire process group hierarchy.

Several other papers[2] showed that the results of software quality measurement can be adopted to the problem of measuring the quality of business processes because an EPC is actually a simple program, which can be characterized with static code metrics like McCabe's cyclomatic complexity. With our contributions the problems of measuring the quality of business processes and software systems became even more similar.

- If a task behind a function is managed by a software system, it is a reasonable assumption that the quality or test coverage of this software has a serious effect on the function and the EPC itself. First we approximated the execution probability of each function and then we used these values as weights for aggregating the quality values of the employed software systems to the EPC level.
- In software quality measurement very often a quality model uses code metrics as predictors. Then we aggregate these predictors to higher level characteristics by comparing them with a benchmark. We also created a quality model for business processes based on the ISO/IEC 25010 quality standard. As predictors we used simple metrics calculated to the EPC and the previously aggregated quality metrics from the functions.
- Since we did not stop at the level of the EPC, and extended the measurement to the entire business process hierarchy we have quality information about the whole business configuration. The calculation was started from the most atomic elements, but the interpretation of the results is done from top to down. Meaning that besides measuring quality of business configuration the approach also provides the explanation.

Acknowledgements

This research was supported by the hungarian GOP-111-11-2011-0038 grant.

References

- [1] Michael Hammer and James Champy, Reengineering the Corporation: A Manifesto for Business Revolution, Harper Business, 1993
- [2] Khlif, Wiem and Makni, Lobna and Zaaboub, Nahla and Ben-Abdallah, Hanene, Quality Metrics for Business Process Modeling. In: *Proceedings of the 9th WSEAS International Conference on Applied Computer Science ACS'09*, Genova: WSEAS. pp. 195–200.

Automatic Failure Detection and Monitoring of Ventilation and Cooling Systems

Balázs L. Lévai

Nowadays, due to the general market conditions and the ever increasing competition and rivalry, industrial companies are pushed more than ever to optimize their operational costs, and improve the product and service related business processes. In addition to the regular requirements, service continuity and availability are also expected. Maintaining a high level quality requires the detection of occasional failures of the included service elements as soon as possible. Consequently, the automatization of this aspect of support could be essential to minimize the time, personnel, and financial costs in a company's everyday operation.

A local service company, which maintains and installs ventilation and cooling systems for a wide range of customers, including hospitals, office buildings, and private homes too, asked us to develop a computer system for failure detection purposes in cooperation with the engineering department of our institute.

Considering the operation of cooling and ventilation systems, one may see or even hear that malfunctions affect and alter the normal vibration caused by the moving parts, mainly the engine and the fans, therefore it is a sound way to monitor such systems by motion sensors. These regularly measure the vibration over a fixed-length of time interval in terms of signed one dimensional acceleration. The moving parts work on specific frequency which means that the collected data of each sensor is a periodic one dimensional time series, or signal in other words.

Our task is the classification of the current condition of the monitored system based on the information of the last measurement using predefined malfunction classes. We present this industrial sample classification problem, the structure and operation of the developed solution including the basic theoretical background [1, 2]. In addition, we discuss the possibilities of reducing the number of monitoring devices with an acceptable loss of detection rate. The test results of our technique, implemented in MATLAB [3], and its future applicability are also presented in detail.

References

- [1] BISHOP, C. M. , *Neural Networks for Pattern Recognition, Oxford University Press, Inc. New York, NY, USA (1995).*
- [2] SMITH, S. W. , *The Scientist and Engineer's Guide to Digital Signal Processing, California Technical Publishing, San Diego, California (1997).*
- [3] THE MATHWORKS, INC. , MATLAB,
<http://www.mathworks.com/products/matlab>.

Automatic Design of LED Street Lights

Balázs L. Lévai

An interesting side effect of industrial civilization - becoming an emphatic problem of our modern age - is light pollution, also known as photopollution or luminous pollution. This phenomenon is the unwanted or unneeded lighting of outdoor areas by artificial light. It obscures the stars in the night sky making impossible the observation of certain astronomical objects. Moreover, light pollution raises questions about energy conservation, a subject having an increasing importance as the world's energy consumption grows [1].

A light-emitting diode, or simply a LED, is a semiconductor light source. The colour of the emitted light can be in the infrared, visible light or ultraviolet wavelength range. Nowadays, the application of LED technology is available for street lighting purposes. It has many advantages over the commonly used incandescent light bulbs such as the lower energy consumption, the longer lifetime, or the dimmability.

Assume, that we would like to light a rectangle shaped street section by a LED street light. As a single LED can only light a small area, we have to use multiple LEDs in the lamp in order to cover the entire street section. Consequently, we need to set the direction and type of every LED to create a configuration. Rigorous regulations specify the minimum, maximum value and other characteristics of public lighting in protection of the motorists, therefore the proper design of a LED configuration is quite a challenge even for an experienced engineer. If we add the obvious expectation of economic optimality to the list, the manual design is certainly impossible.

We developed a software solution to automatize this process. Considering the lighted area of a single LED as the intersection of its sphere shaped lighting characteristic and the street section to be lighted, we are facing an interesting circle covering problem. Since we try to cover a rectangle with circles whose extent of contribution to the coverage is not just different comparing any two of them, but it is also changing within the same circle, one may undoubtedly realize that our problem is far from ordinary. To make things more complicated, we also have to take into account the effect of the neighbouring street sections based on the symmetries, hence calculating the coverage of a single configuration is not trivial too.

We approached this task as a global optimization problem. We used a genetic algorithm to find a suitable configuration. Due to the fact that the bulk of the computations are related to the determination of the actual lightings, we made parallel the independently executable parts of our calculations using NVIDIA's CUDA technology [2]. As a result, the average runtime of 3 hours has dropped to 15 minutes although we remark that the non parallelized version also provided an acceptable solution after 20-30 minutes.

Acknowledgements

We were asked by Wemont Ltd. to develop a LED street light designer software. This work has been partially supported by the Grant K-2010-GOP-1.1.1-09/1.-0240762/129.

References

- [1] GEREFFI, G., DUBAY, K., LOWE, M., Manufacturing Climate Solutions Carbon Reducing Technologies and U.S. Jobs Center on Globalization, Governance & Competitiveness, Duke University, USA (2008).
- [2] NVIDIA CORPORATION, CUDA technology, <http://www.nvidia.com/object/about-nvidia.html> (visited in 2014).

Application of graph based data mining techniques in administrative systems of education

András London and Tamás Németh

Graph based data mining and network analysis have become widespread in the last decade since these tools have been proved to be extremely applicable and useful in a wide range of areas including biology, economy and social sciences, among others. The large amount of available data allowed us to study of such large-scale systems that appears in the mentioned areas. Usually, these complex systems can be represented by graphs (or networks), where vertices (or nodes) stand for individuals, while edges (or links) represent the interaction between pairs of these individuals (for an excellent review, see *e.g.* [1]). The network approach is not only useful for simplifying and visualising this enormous amount of data, but it also very effective to pick up the most important elements and find their most important interactions. Besides, several techniques have been developed to explore the deeper topological features of a network, such as community structure [2], core-periphery structure [3] or small-world property [4] and scale-freeness [5]. Ranking individuals based on their position in the model interaction graph has also become an important direction of studies in the last decade. Random walk based ranking algorithms (such as the widely-known PageRank algorithm by Google), that were originally developed for ranking web pages, have been used recently for such different purposes like citation network analysis [6], ranking in sports *e.g.* in [7, 8] or evaluating the quality of wines and skills of the tasters [9], etc.

In this work, by following the complex network approach, we introduce a novel example of a real social system taken from the world of public education. Since a huge amount of data (that is more accurate as well) is produced by a complex administrative software system of educational institutes, new type of data processing methods are required to handle with it in order to information extraction, instead of the classical statistical analysis. The maintainers, the leaders and the teachers of the institute, the students and their parents would have asked new type of questions about the educational work and quality of the institute, the teachers and students. Such questions, among others, can be the following:

- Can we say something more useful about the efficiency of the teachers' work by using this data than by using the classical questionnaire system?
- Can we make spectacular statements that are easier to understand, *e.g.* with data visualization techniques?
- By applying new type of models, can we get results that are modeling the reality better, than the simple statistical statements, *e.g.* for comparing the achievements of the students in the same year?
- Can we "measure" the improvement of the students and by using these results, can we detect accidentally occurring problems, like drug use, alcoholic problems, crisis in the family, etc.?

We define several suitable network representations of the system with the goal of answering or partially answering these questions. Depending on the construction of the underlying graph we consider three different network models. First we construct an undirected, weighted graph based on various similarity measures between students (in the same school and year) and investigate the topological features, especially the community structure of it. Second we define a directed and weighted graph based on the same data set and apply a ranking algorithm to this network in order to quantify the achievements of the students. Third we construct a bipartite

graph of students and teachers that can be used to compare both the students and teachers with each other according to different aspects. Although we use different algorithms to handle with the different representations, we try to analyze the results simultaneously to get a clear and detailed picture about the achievement and quality of students and teachers and to answer such questions mentioned above.

Keywords Data mining, Educational evaluation, PageRank, Modularity

Acknowledgements

This work was partially supported by the European Union and the European Social Fund through project FuturICT.hu (grant no.: TAMOP-4.2.2.C-11/1/KONV-2012-0013).

András London was supported by the **European Union** and the **State of Hungary, co-financed by the European Social Fund** in the framework of TAMOP-4.2.4.A/2-11-1-2012-0001 'National Excellence Program'.

References

- [1] M.E.J Newman. The structure and function of complex networks. *SIAM review*, 45(2), 167-256, 2003
- [2] M.E.J Newman and M. Girvan. Finding and evaluating community structure in networks. *Physical review E*, 69(2), 026113, 2004
- [3] P. Csermely, A. London, L.Y. Wu and B. Uzzi. Structure and dynamics of core/periphery networks. *Journal of Complex Networks*, 1(2), 93-123, 2013
- [4] D.J. Watts and S.H. Strogatz. Collective dynamics of 'small-world' networks. *Nature*, 393(6684), 440-442, 1998
- [5] A.L. Barabási. Scale-free networks: a decade and beyond. *Science* 325(5939), 412-413, 2009
- [6] P. Chen, H. Xie, S. Maslov and S. Redner. Finding scientific gems with Google's PageRank algorithm. *Journal of Informetrics*, 1(1), 8-15, 2007
- [7] S. Motegi and N. Masuda. A network-based dynamical ranking system for competitive sports. *Scientific reports*, 2, 904, 2012
- [8] F. Radicchi. Who is the best player ever? A complex network analysis of the history of professional tennis. *PLoS One*, 6(2), e17249, 2011
- [9] A. London and T. Csendes. HITS based network algorithm for evaluating the professional skills of wine tasters. *Proc. of the 8th International Symposium on Applied Computational Intelligence and Informatics*, 2013, pp. 197-200.

Distributed News Analytics Framework: Collecting News Feed Sources from social media

Gábor I. Nagy, Sándor Kazi, Győző Papp

Text mining in news articles and social media in the financial context became a vibrant research topic in the past years [1]. Distributed computing helps overcome the difficulties of processing vast amounts of various, heterogenic textual data. The implementation and extension of an earlier version of such a system is discussed in this paper: a distributed adaptive news analytics framework that gathers, stores, curates and process vast amounts of textual data from conventional news feeds and social media. The main components of the system is shown on Figure 5.

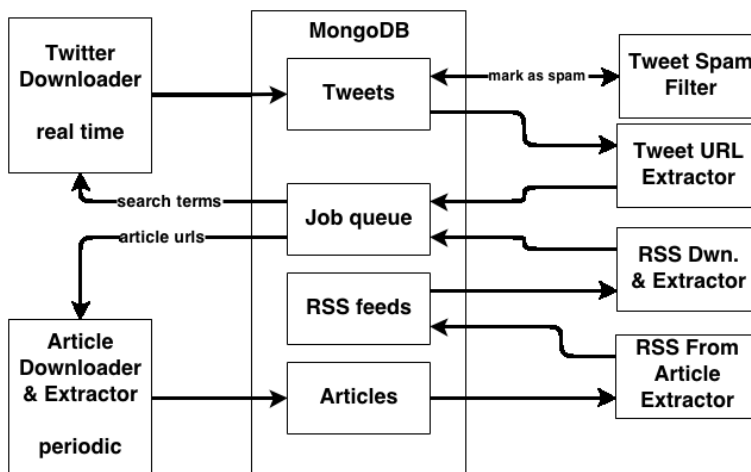


Figure 5: Distributed real-time content monitoring service - system components

The Twitter Downloader component uses Twitters Streaming API to listen to the conversation in social media and stores messages in MongoDB. The API is real-time: tweets are pushed to the data collection server as soon as they are published by Twitter. However one can only search for a limited number of terms (be it ticker symbol, hashtag or user) and can have updates on these terms only. There will be search terms that generate most of the content stored in our database, on the other hand there will be obsolete search terms with little or no messages as topics grow older and interest falls. These obsolete search terms can be replaced by popular terms and ticker symbols found in recent news articles, extracted from either URLs found in Twitter posts. Twitter content, and specifically URLs posted can be monitored for interesting new news RSS feeds, that are currently not used by the Article Downloader & Extractor.

This adaptive behaviour of search terms allows monitoring larger amounts of content than one could monitor with setting constant search terms in Twitter Stream API. However one should be careful with posted content, as Twitter is a prime target of content spam. [2] [3] The Twitter Spam Filter allows ranking user messages as being spam. The Spam Filter is discussed in details in [4]. The Article Downloader & Extractor downloads conventional news articles given by their URL. Extracts various features from these articles: links on the page, raw text extract of the news for text minign, links to rss feeds, or number of ads on a given page from the downloaded HTML. Content is stored in MongoDB. RSS Downloader & Extractor downloads feeds given by their URL and dispatches their links to the Job Queue for to the Article Downloader. The Tweet URL Extractor periodically extracts the URLs from the latest tweets and dispatches them as jobs for the Article Downloader. This module of the system was

described in an earlier work [5]. The module named RSS From Article Extractor copies the interesting new entries from extracted documents originating from Twitter, that containing RSS feed URLs. The purpose of the Job Queue is to be a centralized temporary storage for tasks and messages of the processes. Performance tests of the proposed system is the core focus in the paper together with implementation details.

Acknowledgements

This work is supported by the grant: FUTURICT, TÁMOP-4.2.2.C-11/1/KONV, Financial Systems subproject.

References

- [1] J. Bollen, H. Mao, X. Zeng. Twitter mood predicts the stock market, *Journal of Computational Science*, Vol. 2, 2011, 1-8
- [2] K. Thomas, C. Grier, D. Song, V. Paxson. Suspended accounts in retrospect: an analysis of twitter spam, *IMC '11 Proceedings of the 2011 ACM SIGCOMM conference on Internet measurement conference*, 243-258
- [3] A. H. Wang. Detecting Spam Bots in Online Social Networking Sites: A Machine Learning Approach *Data and Applications Security and Privacy XXIV Lecture Notes in Computer Science*, Volume 6166, 2010, 335-342.
- [4] G. I. Nagy, S. Kazi. Filtering noise from stock related Twitter messages. $(CS)^2$ - *The 9th Conference of PhD Students in Computer Science*. Submitted.
- [5] G. I. Nagy, S. Kazi. Distributed News Analysis Framework for Text Mining. *IEEE, CogInfo-Com 2013 - Workshop on Future Internet Science and Engineering*, 2013.

Filtering noise from stock related Twitter messages

Gábor I. Nagy and Sándor Kazi

Twitter is a micro-blogging service that is used by millions of people to publish very short messages and broadcast it to their followers. This real-time service allows users to generate, read and discover interesting content. Parallel to its increasing popularity researchers also took interest in the service. Various research topics include text mining, sentiment analysis, topic discovery, social network analysis. The massive user base - measured in millions of people - also made the service valuable for marketers, who target consumer with advertisements promoting a particular website or service. An increasing volume of these messages hold little or no information regarding research of regular user content, and this commercial noise bias research and erode the performance of data mining algorithms that extract topics or model sentiment from user generated content. To avoid performance degradation a distinction should be made whether a message is relevant in the context of research or it is just spam.

To tap into the social conversation Twitter allows members to use its API. A prevalent use-case in a data mining is that researchers use the Twitter Stream API. This API helps to gather all the tweets real-time that contain particular terms, stocks (referred to as symbols), hashtags or user entities. For example if a researcher wishes to follow the conversation regarding stocks from S&P500, one can set up the query containing only the symbols of interest. The stream will inevitably contain tweets with a commercial incentive.

In this paper we explore methods for filtering relevant content from Twitter messages related to stocks, indexes and currencies represented by their symbols in tweets. We use the Twitter Streaming API to gather messages real-time for 249 assets. Majority of related work rely on crowdsourced, annotated data to filter relevant messages, user characteristics [1] [2] or textual features [3]. Our approach is comprised of three steps:

1. Label a number of spammers and non spammers and build a classifier to rank users. Features of the classifier include posting behaviour, vocabulary and general characteristics of the user (eg. number of statuses so far, number of friends, number of followers, account age, etc.).
2. Use classifier and rank unseen users. Annotate messages from users with high probability of spamming behaviour. Ranking users and tweets this way helps annotation.
3. Build classifier on messages textual features and evaluate performance.

Evaluation is based on standard accuracy measures (accuracy, precision, recall, F-measure) as well as AUC-ROC. Preliminary results show that ensemble methods, such as Gradient Boosting Decision Tree Classifiers (GBDT) [4] and Random Forest Classifiers (RF) [5] work best on the annotated dataset. The best GBDT model has the following error characteristics based on 5-fold cross validation: mean accuracy 81.6%, mean recall 86.6%, mean f1-measure 82.7% and a mean AUC score of 91.1%. The main characteristics of a spammer is that it uses a lot of different hashtags and stock symbols in its messages with large number of URLs. The account age is found to be less important which is in accordance with the findings in [5], that a lot of spamming accounts are compromised ones. Findings are in line with results in related work experiments. Further refinement of attributes and classifier parameters are needed, with a larger number of annotated users for better classification.

Acknowledgements

This work is supported by the grant: FUTURICT, TÁMOP-4.2.2.C-11/1/KONV, Financial Systems subproject.

References

- [1] F. Benevenuto, G. Magno, T. Rodrigues, and V. Almeida. Detecting Spammers on Twitter, *CEAS 2010 - Seventh annual Collaboration, Electronic messaging, Anti-Abuse and Spam Conference*, July 13-14, 2010, Redmond, Washington, US.
- [2] M. McCord, M. Chuah. Spam Detection on Twitter Using Traditional Classifiers, *Autonomic and Trusted Computing*, Lecture Notes in Computer Science Volume 6906, 2011, 175-186.
- [3] Juan Martinez-Romo, Lourdes Araujo. Detecting malicious tweets in trending topics using a statistical analysis of language. *Expert Systems with Applications*, Volume 40, Issue 8, 15 June 2013, 2992-3000.
- [4] Jerome H. Friedman. Greedy function approximation: A gradient boosting machine. *The Annals of Statistics*, Volume 29, Number 5 (2001), 1189-1232.
- [5] Leo Breiman. Random Forests. *Machine Learning*, Volume 45, Issue 1, 5-32.

Phase delay detecting methods to analyse the correlation between blood pressure and paced breathing

Tamás Dániel Nagy, Gergely Vadai, Zoltán Gingl, László Rudas, and Éva Zöllei

The analysis of the temporal relationship between blood pressure and respiration is useful to understand the underlying physiological processes and in diagnostic applications. Our work is based on a clinical experiment, where ECG, blood pressure and capnograph signals were measured during paced breathing in normal and simulated hypovolemic states.

The *pulse pressure* defined as the difference of the blood pressures maximum (*systolic*) and minimum (*diastolic*) values in a cardiac cycle. The temporal fluctuations of these three quantities show correlation with respiration, furthermore the phase delay between these signals and respiration seems to be relevant in diagnostics. For this reason, the detection of this delay is important and useful, unfortunately in practical cases it can be difficult. Only a few heart beats occur during one respiratory cycle, therefore the sampling frequency of the mentioned signals is relative low, furthermore because of the heart rate fluctuations, the signals are unevenly sampled. Moreover, the inaccuracy of the respiration frequency and further physiological processes can distort the sinusoidal-like shape of the signals.

Several clinical and physiological studies used the fundamental methods of phase delay detection in time or frequency domain. Our addition to this problem is a comparison of these methods' efficiency in the case of differently distorted signals. In this work, we present a systematic analysis of the methods' reliability and accuracy on the measured time series and numerically simulated signals.

The heartbeats occur at different time instants in a respiratory cycle, therefore after synchronizing these cycles in a long-time registration using the square wave-like capnograph signal, we can calculate an averaged pulse pressure (or systolic/diastolic blood pressure) period between two inspirations/expiration [1, 2]. Then we can compute the time delay between the pulse pressure peaks/valleys and the falling/rising edges of the capnograph signal. Unfortunately the detection of these peaks can be rather difficult in most cases.

We have examined the usage of spectral-domain based detection methods also. These methods compute the phase delay using the relevant phase values from the signals' spectra or from the interpolated signals' cross spectra [3, 4].

The mentioned detecting methods have been applied on the data series measured in the clinical experiment. The spectral methods seemed to be more reliable at the cases where the peaks of the averaged pulse pressure cycle were hardly detectable. Furthermore, the accuracy and reliability of the methods have been tested on numerically simulated signals with different distortion and variance of breathing frequency and heart rate.

Acknowledgements

The publication/presentation is supported by the European Union and co-funded by the European Social Fund. Project title: "Telemedicine-focused research activities on the field of Mathematics, Informatics and Medical sciences" Project number: TÁMOP-4.2.2.A-11/1/KONV-2012-0073.

References

- [1] P. Y. W. Sin, D. C. Galletly, Y. C. Tzeng, Influence of breathing frequency on the pattern of respiratory sinus arrhythmia and blood pressure: old questions revisited, *American Journal of Physiology - Heart and Circulatory Physiology*, 298(5) pp.1588-1599, 2010.

- [2] P. Sundblad, D. Linnarsson, Relationship between breath-synchronous arterial pressure and heart rate variations during orthostatic stress, *Clinical Physiology and Functional Imaging*, 23(2), pp.103-109, 2003
- [3] J. K. Triedman J. P. Saul, Blood pressure modulation by central venous pressure and respiration. Buffering effects of the heart rate reflexes, *Circulation*, 89(1), pp. 169-179, 1994.
- [4] M. Orini, R. Bailón, P. Laguna, L. T. Mainardi, R. Barbieri, A multivariate time-frequency method to characterize the influence of respiration over heart period and arterial pressure, *EURASIP Journal on Advances in Signal Processing*, 2012(1), pp. 1-17, 2012.

Automatic Detection of Multiword Expressions with Dependency Parsers on Different Languages

István Nagy T.

Here, we present how different types of MWEs can be identified by dependency parsers in different languages. In our investigations, we focus on English verb-particle constructions (VPCs), Hungarian light verb constructions (LVCs) and German light verb constructions. In our experiments, we exploit the fact that some treebanks contain MWE-aware annotations, i.e. there are MWE-specific morphological or syntactic tags in them. For instance, the French Treebank contains explicit annotations for MWEs [1] and different version of the Turkish Treebank are also annotated for MWEs [2]. Here, we make use of the Penn Treebank [3], which contains annotation for VPCs, the TIGER corpus [4] and the Szeged Dependency Treebank [5], both of which contain annotation for LVCs. In these treebanks, the special relation of the two components of the MWE is distinctively marked by a certain syntactic label. This entails that if a data-driven syntactic parser is trained on a dataset annotated with extra information for MWEs, it will be able to assign such tags as well, in other words, the syntactic parser itself will be able to identify MWEs in texts. In our experiments, we investigate the performance of such dependency parsers for three languages and two different MWE types.

English VPCs

The special relation of the verb and particle within a VPC is distinctively marked in the Penn Treebank, the particle is assigned a specific part of speech tag (`RP`) and it also has a specific syntactic label (`PRT`). Thus, parsers trained on the Penn Treebank are able to identify VPCs in texts. We experimented with two dependency parsers, namely the Stanford parser [6] and the Bohnet parser [7] and examined how they can perform on the Wiki50 corpus [8]. This corpus contains English Wikipedia articles which are annotated for several types of MWEs – thus for VPCs as well – and named entities. We parsed the texts of Wiki50 with the two parsers, using their default settings and if the parser correctly identified a `PRT` label, we considered it as a true positive. For evaluation, we employed the metrics precision, recall and F-measure interpreted on VPCs. The two parsers obtained the following results: the Stanford Parser achieved 91.09 (precision), 52.57 (recall) and 66.67 (F-measure) and the Bohnet Parser achieved 89.04 (precision), 58.16 (recall) and 70.36 (F-measure). Thus, precision values are rather high but recall values are lower, which suggests that the sets of VPCs found in the Penn Treebank and Wiki50 may differ significantly.

Hungarian LVCs

The Szeged Dependency Treebank contains manual annotation for light verb constructions [9]. Dependency relations were enhanced with LVC-specific relations that can be found between the two members of the constructions. For instance, the relation `OBJ-LVC` can be found between the words *döntést* (`decision-ACC`) and *hoz* “bring”, members of the LVC *döntést hoz* “to make a decision”.

We used the Bohnet dependency parser to identify LVCs in the legal subdomain of the corpus. We applied 10-fold cross validation here and got the following values: 86.60 (precision), 67.12 (recall), 75.63 (F-measure). According to the results and error analysis, the main advantages of the system are the high precision value on the one hand and the adequate treatment of non-contiguous LVCs on the other hand [5].

German LVCs

In the TIGER corpus, LVCs that consist of a verb and a prepositional phrase are annotated with the relation `CVC`. The German model of the Bohnet parser trained on the Tiger corpus is able to assign such a label, so we used it in our experiments with its default settings. For evaluation, we selected a subset of the German part of the JRC-Acquis corpus, which has recently been annotated for LVCs [10]. If the parser correctly identified a `CVC` label, we considered it as a true positive. We obtained a result of 84.81 (precision), 60.91 (recall) and 70.90 (F-measure), which indicates that similar to English VPCs, the set of LVCs in the test corpus may just partly overlap with the set of LVCs in the TIGER corpus.

Acknowledgements

István Nagy T. was partially funded by the State of a Hungary, co-financed by the European Social Fund in the framework of TAMOP-4.2.4.A/ 2-11/1-2012-0001 “National Excellence Program”.

References

- [1] Abeillé, A., Clément, L., Toussnel, F.: Building a Treebank for French. In: *Treebanks : Building and Using Parsed Corpora*. Springer (2003) 165–188
- [2] Eryiğit, G., Ilbay, T., Can, O.A.: Multiword expressions in statistical dependency parsing. In: *Proceedings of the Second Workshop on Statistical Parsing of Morphologically Rich Languages*, Dublin, Ireland, Association for Computational Linguistics (October 2011) 45–55
- [3] Marcus, M.P., Santorini, B., Marcinkiewicz, M.A.: Building a Large Annotated Corpus of English: The Penn Treebank. *Computational Linguistics* **19**(2) (1993) 313–331
- [4] Brants, S., Dipper, S., Eisenberg, P., Hansen-Schirra, S., König, E., Lezius, W., Rohrer, C., Smith, G., Uszkoreit, H.: TIGER: Linguistic interpretation of a German corpus. *Research on Language and Computation* **2**(4) (2004) 597–620
- [5] Vincze, V., Zsibrita, J., Nagy T., I.: Dependency Parsing for Identifying Hungarian Light Verb Constructions. In: *Proceedings of the Sixth International Joint Conference on Natural Language Processing*, Nagoya, Japan, Asian Federation of Natural Language Processing (October 2013) 207–215
- [6] Klein, D., Manning, C.D.: Accurate unlexicalized parsing. In: *Annual Meeting of the ACL*. Volume 41. (2003) 423–430
- [7] Bohnet, B.: Top accuracy and fast dependency parsing is not a contradiction. In: *Proceedings of Coling 2010*. (2010) 89–97
- [8] Vincze, V., Nagy T., I., Berend, G.: Multiword Expressions and Named Entities in the Wiki50 Corpus. In: *Proceedings of RANLP 2011*, Hissar, Bulgaria, RANLP 2011 Organising Committee (September 2011) 289–295
- [9] Vincze, V., Csirik, J.: Hungarian corpus of light verb constructions. In: *Proceedings of Coling 2010*, Beijing, China, Coling 2010 Organizing Committee (August 2010) 1110–1118
- [10] Rácz, A., Nagy T., I., Vincze, V.: 4FX: Light Verb Constructions in a Multilingual Parallel Corpus. In: *Proceedings of LREC*. (2014)

Derivable Partial Locking for Algebraic Data Types

Boldizsár Németh and Zoltán Kelemen

Concurrency is one of the most actively researched fields of Computer Science. Writing concurrent programs is challenging. The causes are the need for synchronization and solving possible race conditions and deadlocks while avoiding to unnecessary waiting and overhead. This can be very difficult when a transaction needs to lock multiple data elements, even with using previously defined concurrent data structures [1].

Algebraic Data Types are composite data structures that naturally support pattern matching. When the type is parametric then the ADTs support writing generic algorithms without additional complexity.

The integrity of the program data can be archived by providing locks for a data structure or using concurrent data structures. Central locking is a safe strategy to avoid race conditions, but it has a high cost because globally used objects are frequently accessed by different threads. Hierarchic locking allows transactions to lock exactly the data elements that they need [2]. This makes small transactions take small locks and work in parallel without waiting for each other.

This article focuses on a method that helps the implementation of thread-safe programs with ADTs. By transforming the data model of the application to thread-safe data structures with a built-in locking mechanism, a programmer can focus on the business logic of his application when writing the program [3].

First, we need to transform the original ADT to the thread-safe version. For this, the programmer can configure what parts of the data structure should be locked. Too many locks cause performance loss. Second, we define a frontend to access the locked parts of the data structure. This must be done with minimal syntactical noise, to ease the development of the application.

We implement our solution to this problem in Haskell. We use the Haskell concurrency primitive `mvar` [4] to create a concurrent version of the data model. Two tools are inspected that are capable of transforming the data structure to a thread-safe version. The generics of GHC provide a way to derive type class instances for ADTs, like how it can be done with built-in type classes. This can be done by decomposing the structure of the ADT and defining the meaning of the functions on these primitive blocks. The second tool is the Template Haskell compiler extension. With TH the internal representation of the program can be inspected and modified. This grants more freedom to transparently change the program, but there is the danger of confusing users and committing errors while transforming the program.

Acknowledgements

Supported by EITKIC 12-1-2012-0001.

References

- [1] Ohad Shacham, Nathan Bronson, Alex Aiken, Mooly Sagiv, Martin Vechev, and Eran Yahav. 2011. Testing Atomicity of Composed Concurrent Operations. In Proceedings of the 2011 ACM international conference on Object oriented programming systems languages and applications (OOPSLA '11). (51-64).
- [2] Goetz Graefe. Hierarchical locking in B-tree indexes.
- [3] Peter Hawkins, Alex Aiken, Kathleen Fisher, Martin Rinard, and Mooly Sagiv. 2012. Concurrent Data Representation Synthesis. Proceedings of the 33rd ACM SIGPLAN conference on Programming Language Design and Implementation. Pages 417-428

- [4] Simon Peyton Jones, Andrew Gordon, Sigbjorn Finne. 1996. Concurrent Haskell. 23rd ACM Symposium on Principles of Programming Languages.

Visualization and analysis of financial transaction networks - a multidimensional modeling approach

Zoltán Németh

In recent years network analysis and visualization has become an important concept in the growing field of business intelligence researches and applications. Through data mining methods, modeling customer value and behavior based on an abundance of data were already widely used during the last decade. Interest in customer relationships has increased particularly in telecommunications and banking industry since then, as efficient methods provided by the emerging network science paired with sufficient computing capacity has proven to be able to handle transactional data from a network perspective.

The operation of a particular business area represented by different entities and relations may be visualized and analyzed for palpable business benefits [1]. In this paper we focus mainly on financial transactions on an institutional level, using local data warehouses, data marts as standard data sources. Assuming a sales, risk or churn analysis perspective, relying on local business knowledge and additional steps of ensuring a proper level of data quality are to be considered in order to model and process transactional relationships resulting in money transfers. Types of the network units are defined aiming for relatively simple structures and flexibility, always considering both the benefits and limits of the end-users' ability of perceiving and processing visual information.

Core aims of this research are expanding and enhancing application development aspects, on the basis of our past research and development experience in displaying a browseable part of a network, and leaning on the analysis toolbox of network science, recognizing the importance of environmental statistics or path and pattern search algorithms. Apart from a constant need for improved performance, programability and scalability in such applications, added functionality can be produced by investigating and incorporating best practices in existing widely popular BI solutions such as on-line analytical processing. While table-based, normalized storage of nodes and edges represent big challenge for the network approach, data modeling has evolved through the years with new technologies and tools. Recognizing the fallbacks and limited effectiveness of prior attempts a different approach is proposed: applying multidimensional models in back-end components that are highly optimized to query performance and have flexible ways of defining business context. As we are mainly focusing on cost-effective solutions, designing an effective caching strategy and aggregation scheme is an essential step in creating such a data layer, while in-memory engine benefits and drawbacks need to be evaluated, and recent improvements in tabular modelling options are also examined. Furthermore, possible parallel solutions, filtering and ranking options are discussed regarding the business logic layer algorithms.

Keywords: network visualisation, financial transaction, multidimensional data model

References

- [1] Cser, L. (ed.) (2013) Business value in an ocean of data, Budapest: Alinea.

Reconstruction of hv-Convex Binary Matrices from Horizontal and Vertical Projections Based on Simulated Annealing

Zoltán Ozsvár and Péter Balázs

Tomography is a method of producing a three-dimensional image of the internal structure of an object from its projections, without damaging it. This is usually achieved by reconstructing 2D slices from the projections and then assembling them. In Binary Tomography we assume, that the examined object is homogeneous to reduce the number of projections needed for the reconstruction. In case of just two projections the task is usually extremely underdetermined. To overcome that problem we further assume that the reconstructed image satisfies certain geometrical conditions, such as hv-convexity. The reconstruction of hv-convex binary matrices from their horizontal and vertical projections is proved to be NP-hard. In this work we design and study two different algorithms based on simulated annealing and compared them with a formerly published method. The first algorithm based on horizontal convex strips, and the second inverts pixels during the reconstruction process. We use a large set of test data, with different size and number of components. We study two different methods to estimate the convexity of a binary image, the directional convexity measure and another one which calculates the number of the neighbouring object pixels. We deduce that the directional convexity measure is more efficient to analyse experimentally the performance of two approaches.

Acknowledgements

Péter Balázs was supported by the OTKA PD100950 grant of the National Scientific Research Fund, and the European Union and the State of Hungary co-financed by the European Social Fund under the grant agreement TÁMOP 4.2.4.A/2-11-1-2012-0001 ('National Excellence Program').

Load Balancing Strategy in Mobile Resource Management

Krisztián Pándi and Hassan Charaf

This paper is dealing with mobile resource management that uses cloud computing resources. Load balancing is an important part of mobile resource management. In order to be able to solve load balancing issues, we must examine its methods and their effectiveness. Further issues such as solution of load balancing problem in mobile and cloud computing environments also need further examination. In the following, we investigate load balancing procedures, methods and their customization, with a particular attention on mobile and cloud computing requirements. As a result, we expect that important design aspects will become apparent.

Mobile devices have become a part of everyday life. Due to their small size, being always at hand and having relatively high calculating capacity, offering wide variety of applications with very different resource need.

Cloud computing promises to provide high performance, flexible and low cost on-demand computing services. Their resources compared to mobile devices are significantly larger and more scalable.

The use of cloud resources in mobile device seems tangible. In our previous article we suggested such architecture of mobile resource management, which can utilize benefits of cloud computing, expanded with smart using of available network interface parallel. The goal of the mechanism is to decide where the optimal place is for a certain service/application to run; on the mobile terminal itself or on public cloud computing server. One of the most interesting question is the load balancing strategy of suggested resource management layer. In this article load balancing topic will be investigated; the way it can be inserted into current resource management architecture, which if this kind of strategy is more forward-looking and most advantageous. Load balancing strategy has a key role in resource management architecture; it must meet several parallel requirements. The main goal of load balancing is to optimize resource usage. In proposed resource management architecture currently two resources are available; mobile device and cloud computing environment. Well-chosen load balancing strategy may benefit from the extra resource, and can lead to increased performance and reliability.

In the current article software load balancing open questions were investigated, which were aroused during resource management architecture planning and realization.

Acknowledgements

This work was partially supported by the European Union and the European Social Fund through project FuturICT.hu (grant no.: TAMOP-4.2.2.C-11/1/KONV-2012-0013) organized by VIKING Zrt. Balatonfured. This work was partially supported by the Hungarian Government, managed by the National Development Agency, and financed by the Research and Technology Innovation Fund (grant no.: KMR 12-1-2012-0441).

x86 instruction reordering and split-stream compression benchmark

Zsombor Paróczy

Binary executables are the result of the compilation process, which produces machine code, initialized and uninitialized data and operation system specific headers from source code.

Runtime executable compression is a method which uses standard data compression methods and binary machine code transformations to achieve smaller file size, yet maintaining the ability to execute the compressed file as a regular executable. In our work we only focus on Intel x86 instruction set with 32bit registers, which mainly used in personal computers. In this paper we perform a benchmark on different compressions using various binary transformations, the main contribution to the field is (i) testing Zopfli on binaries and (ii) benchmark the instruction reordering method combined with split-stream algorithm.

Method	Based on	Homepage
Aplib	LZ-based	http://www.ibsensoftware.com/
Gzip	deflate (LZ + Huffman)	http://www.gzip.org/
LZMA	Lempel-Ziv-Markov chain algorithm	http://www.7-zip.org/sdk.html
Zopfli	deflate based	https://code.google.com/p/zopfli/

Table 1: Compression methods

For our benchmark we selected four different compression methods, they can be seen in Table 1. Besides of the common algorithms widely used in Windows executable compression programs [1], we also added Zopfli to the benchmark. [3]. All of the data compression algorithms are lossless. In our benchmark we tested these compression method on various sample executables (details in Table 3) in four test cases (i) without modification (ii) using instruction reordering only (iii) using split-stream only (iv) the combination of both method. The combination of these two methods are: executing the instruction reordering algorithm than using the split-stream method before the compression. The (i) test case serves as a baseline for evaluation the improvements, the other three uses binary machine code transformation methods.

Instruction reordering is a method, which modifies the order of the instructions inside a basic blocks using data flow constraints, this reordering can improve code compression without changing the behavior of the code [5].

Split-stream is an algorithm that initially partitions a program into a large number of sub-streams with high auto-correlation and then, heuristically merges certain sub-streams in order to achieve the benefits provided by classical split-stream. It can reduce the increase in compression ratio which typically occurs when a PPM-like algorithm compresses small amounts of data [4]. The actual implementation was developed by Fabian Giesen for the compression program kkrunchy [2].

The implementations didn't use parallel processing methods, each compression method / machine code transformation method was run in a single thread mode on the same machine. We analyzed both the compression time and the result size, in Table 2 you can see each methods' runtime on the nodejs binary. The running time of the compression is influenced by two major factors: the compression method performance and the permutation count for the instruction reordering. Decompression is not effected by the permutation count of instruction reordering, the compression has to be done once for each binary.

	instr. reordering	split-stream	combined
Aplib	2877	3	3224
Gzip	1445	<1	1497
LZMA	2047	2	3620
Zopfli	7426	15	29413

Table 2: Compression times in seconds

Name	OS	Compiler	Source
libc-2.13.so (-0ubuntu13.1)	Ubuntu	gcc	http://packages.ubuntu.com/natty/libc6-i386
unzip (6.0-4)	Debian	gcc	http://packages.debian.org/squeeze/unzip
libconfig.dll (1.4.8)	Windows	VS2008	http://www.hyperrealm.com/libconfig/
node.js (0.8.8)	Mac	llvm	http://nodejs.org/dist/latest/node-v0.8.8-darwin-x86.tar.gz

Table 3: Binaries in the benchmark

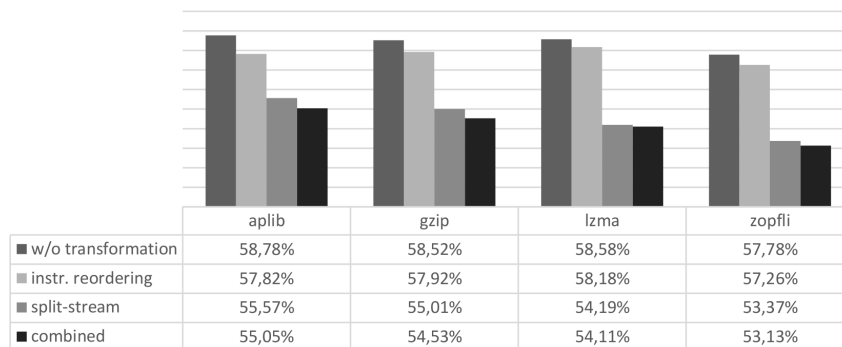


Figure 6: Compressed binary size compared to the uncompressed binary

On Figure 6 you can see the compressed binary sizes compared to the original binary size for the unzip binary. Each method produces significant improvement over the compressed original data.

The combined method gives an average of 9.46% smaller compressed file, than without any binary transformation. The gain by compression method: 6.8% in Aplib, 9.7% in Gzip, 8.7% in LZMA, 12.64% in Zopfli.

References

- [1] Árpád Beszédes, Rudolf Ferenc, Tibor Gyimóthy, André Dolenc, and Konsta Karsisto. Survey of code-size reduction methods. *ACM Comput. Surv.*, 35(3):223–267, September 2003.
- [2] Fabian Giesen. Working with compression. Online. Last accessed: 2013-01-15, 2006. Breakpoint 2006.
- [3] Lode Vandevenne M.Sc. Google Inc. Jyrki Alakuijala, Ph.D. Data compression using zopfli. February 2013.
- [4] Steven Lucco. Split-stream dictionary program compression. *SIGPLAN Not.*, 35(5):27–34, May 2000.
- [5] Zsombor Paroczi. x86 instruction reordering for code compression. *Acta Cybernetica*, 21(1):177–190, 2013.

Diversification for Content-Based Visual Information Retrieval System for Video

Zsombor Paróczy, Bálint Fodor, Gábor Szűcs

The number of user-generated content has grown exponentially due to the media recording and sharing methods becoming easier and cheaper for everyone. In recent years computer vision and data mining communities devoted significant attention to user submitted media analysis and retrieval, this field is usually referred to as content-based visual information retrieval (CBVIR). Due to the huge number of media content, without an efficient search system we could not find the relevant information. Existing retrieval technologies focus on the precision of the results that often provides the user with redundant answer list with near duplicated items. The users would like to retrieve not only relevant items, but also diverse results.

Our work focuses on a potential tourist, who tries to find more information about a famous place, only knowing the name of the place. The literature calls this type of attitude the surfer, who has a moderate goal, has a starting point but may want slightly different outcome in the same context. The goal besides finding the relevant items is to filter the duplicated or similar content. In this special case a diverse result for a location name based search should include shots from various weather conditions and seasons, day/night shots from the same place. We only focus on search results of YouTube, but the same method should work on any video based service.

CBVIR systems usually have three distinct steps: (i) content retrieval, (ii) data organization and indexing and (iii) data-driven applications [2]. In the first, the data is extracted from the source, in our case we used a self written video extractor which helped us to overcome the YouTube location based suggestion system. In this preprocessing phase every video was cut into multiple shots along the time domain. Every shot is intended to contain a coherent scene of the video, this gave us ability for fine-grained analysis of the video content. The dataset we worked on was cut into shots manually to ensure a quality input for the following phases. The raw data is obviously too big for direct analysis, so in the data organization and indexing step the extracted data is transformed to a smaller, but still distinct and representative dataset - in our case we created feature vectors for each shot, using a wide variety of video and image processing algorithms (including edge detection, color histograms, histograms of gradients). In the data-driven application step our goal was to reorder the existing results, we used a combination of machine learning and data mining algorithms.

The relevance of a shot was determined by a pre-trained neural network classifier. The training data was the 20% of the whole dataset. For improving the early diversity of the resulting shot order, first we clustered the initial ordering. The new result list was formed by picking relevant shots from each cluster. While the relevant and diverse shots are preferred to take the first places of the new ordering the algorithm tries to preserve the initial ordering as much as possible.

Our evaluation considers both the relevancy and the diversity factors, this is why the CBVIR algorithms are either evaluated by these two factors separately or using some kind of jointly optimization metric [1]. In our evaluation we used average precision at N seconds ($P@N$), cluster recall at N seconds ($CR@N$) (measure of how many of the existing clusters are represented in the final refinement, so this is the diversity) and harmonic mean of them, the F1-measure at N seconds ($F1@N$).

We used the same dataset, including the shot annotations in our previous work [3]. In that we used a smaller feature vector, and the reordering was done using K-means algorithm and cluster diameters for the reordering. In that work the evaluation was done in a shot based scale, we managed to improve the original reordering by 9-11% (measured at the 10th shot, 20th shot

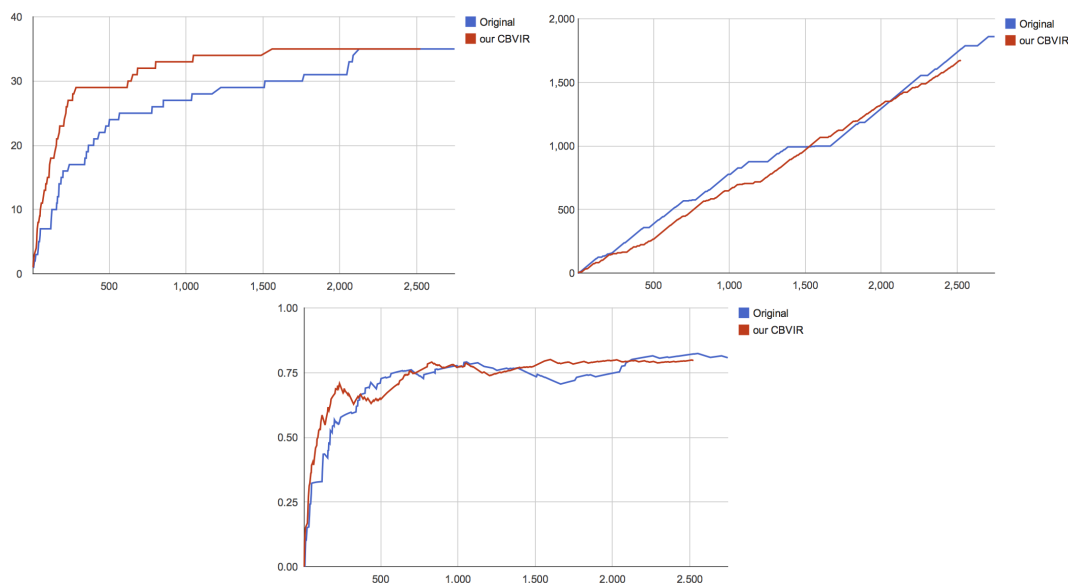


Figure 7: Results from left to right: CR@N, P@N, F1@N

and 50th shot). In this paper we extended the feature vector set with more metrics, and used neural network with support vector clustering. Due to the shot length variation the evaluation was done using the time dimension in seconds, which gives us a much more precise result.

Figure 7 shows the diversity, the relevance and the F1-measure for the “Acropolis of Athens” search result as a function of the time in seconds. The YouTube ordering (Original) is compared to the re-ranked ordering (Our CBVIR). The evaluation shows that our approach improves the original ordering by 14-18% at the first result shots.

Acknowledgments

The publication was supported by the TÁMOP-4.2.2.C-11/1/KONV-2012-0001 project. The project has been supported by the European Union, co-financed by the European Social Fund.

References

- [1] Thomas Deselaers, Tobias Gass, Philippe Dreuw, and Hermann Ney. Jointly optimising relevance and diversity in image retrieval. In *Proceedings of the ACM International Conference on Image and Video Retrieval, CIVR '09*, pages 39:1–39:8, New York, NY, USA, 2009. ACM.
- [2] Shi-Min Hu, Tao Chen, Kun Xu, Ming-Ming Cheng, and RalphR. Martin. Internet visual media processing: a survey with graphics and vision applications. *The Visual Computer*, 29(5):393–405, 2013.
- [3] Zsombor Paróczy and Bálint Fodor. Video content retrieval diversification system. In *Tavaszi Szél Konferencia*, march 2014.

Workflow processing using SAP Objects

Attila Selmeçi and István Orosz

The SAP, the developer of the leading Enterprise Resources Planning System in the World, has been providing business solutions in different areas required by the market for more than 40 years. The basic development environment of the SAP applications is smoothly implemented into the system. All business applications are developed in the internal ABAP (Advanced Business Application Programming) language. This language has its origin from the old style COBOL, PL1, but in the last 40 years the company had enough time to implement features from new technologies, programming languages. In the last ten years SAP implemented and uses for its own applications the object oriented programming paradigm, which is an extension of the original ABAP.

The business logic is implemented as reports, so-called transactions and function modules. Reports generate lists, transactions help to store, change business data in the system through many input and output screens defining a single logical unit of work. The function modules are special modularization units executing a determined functionality (like booking). Grouped function modules can share data through the main program of the group, which cannot be touched directly only via the connected function modules. This approach is similar to the static methods and attributes of an OOP class.

The Business Framework Architecture provides a special environment for internal and external usage of the system services by collecting them into Business Components and offering the real business entities as so-called Business Objects. The Business Objects (a determined entity from an entity type) are instances of the Business Object types. This Business Object layer extends the original report, transaction and function module layer with an object-oriented view, which is usable from outside as well. The Business Object Types contain as subcomponents attributes, methods, and events as well. There are special attributes, the key fields, which refer to the real, underlying data model. The other attributes can be simple fields, structures, arrays or even references to other Business Object Types. The events are information to the "world" about status changes of the Business Object. For example an employee can be hired or fired, an invoice can be approved, rejected, booked, parked as well. If a status change would be interesting for the system it can be publicized by an event. The methods execute the status changes of the business objects, like book, approve for an invoice. A method (as an attribute) can be instance specific or static (class level), and it can have parameters as well. There are special methods in the Business Object world, like CreateFromData, GetList, etc. These can get information from the available entities of the entity type or create a new instance, entity from scratch in the database. There is a unique meaning of public methods, because it means released for external availability. Only these methods can be called from outside of the system. These public methods of Business Objects are the BAPIs (Business Application Programming Interfaces).

The main constituent of process control is an object, which can store its current state; the state can be modified from outside, and the status changes can be publicized. The Business Objects are good candidate to be process control components. In our study we uncover the options to build up process flows or workflows within SAP systems and Microsoft Dynamics, and check the potential of creating outside driven workflows. We try to stretch the borders of the study to show the possibilities of creating and using system-wide workflows or processes by help of Business Objects.

Acknowledgements

This work was supported by TAMOP-4.2.2/B-10/1-2010-0012 Programme.

New Approaches of Storing ERP Data

Attila Selmei and Tamás Gábor Orosz

To install an earlier ERP system we had to design the disk layout accurately keeping in mind the behavior, functionality, and mechanism of many different hardware and software components. Some of these components should be mentioned, like the ERP solution itself, the underlying database management system solution and the operating system. The business data is stored generally always in database tables with primary keys. In some cases the relation between the tables are defined in the database level, but there are ERP solutions, where the table correlations are defined on application level having own dictionary for that. In both cases the data is read from (and stored in) tables, which are not always accessed by unique key attribute values, but several times by not key attributes. To have an acceptable response time from the database layer indexes should be created for the frequently searched attribute sets. The first database system based ERP systems required special attention to isolate table data and index parts to separate disk areas. The main reason was the speed, because the database can read the index and data parts simultaneously in that case. The data store of ERP systems as technique and method can be divided into three main areas:

- Business data representation in the ERP system (logical view)
- Database storage of the internal representation (technical view)
- Data or disk storage technology (physical view)

These different approaches cannot be handled as separated, unequally replaceable methods, but they are rather layers, which not in any case are really distinguishable. In the hardware world the physical layer was improved in the last two decades. The two main technical goals were generally the reliability (e.g. redundancy, high availability) and the response time, speed (e.g. load distribution, bandwidth). The early data stores used different RAID (redundant array of independent disks) technologies to approximate both goals together. The capacity and the speed of the disks increased and it was available to build and use larger data stores. The first RAID controllers were hardware components only, but afterwards some software configuration tools were offered to manage, design the disk distribution in a server. This direction founded the early storage system and solutions having specially connected disks with special-purpose hardware and an operating system offering many different features and services.

According to our study the storage vendors did not really follow the new hardware developments, options, but only the palette of services grew in a wide range. Nowadays the storage systems offer several redundancy, high availability, data handling (like de-duplication), backup, fast copy, read cache and other options. These features are very useful and required in some cases, but the costs are much more higher and the value that they provide. On the other hand the speed as a goal is not in the main focus of the vendors, because they advice only faster disks or even flash drives (SSD-s), and of course more disks to distribute the load among the spindles. These directions are correct but with having many features, special connectivity layers, like Fiber channels and switches, the disks, the real data storing surfaces are far away from the database server or even application server CPU-s. The newer and faster hardware elements made easier to use bigger, faster and reliable 'local' disks in a server. The local disks or even PCIe storage cards provide data in a much higher speed limit, which are sometimes necessary. Our measurements proved that less CPU power is enough if the data storage layer can faster provide data for the ERP system (in case of commercial RDBMS it could decries license amount as well).

The top layer of the above mentioned three data handling approaches determine the mechanism of data management and data store in lower layers as well. Each ERP system contains

several business modules and functions, which handles business objects, like accounting documents, employees, sales or purchase orders, etc. The business objects can be defined as object oriented entities with status (defined via properties), object modification functions (declared through class methods), and of course published status modifications (managed by events) to inform other objects about status change (like employees were hired or even fired). The business objects on technical level in an ERP system are not always defined as real object oriented entities, but only logical correspondence and relationships are defined between data or even tables. Many ERP environments can manage them as objects in higher level, but each time these objects are stored in database tables separating the current status to several tables connected in the database (or as described above in ERP system) level via foreign key relationships. The object as a unit or even entity cannot be recognized on database level.

According to our study for ERP and other kind of OLTP (Online Transaction Processing) systems the object based data handling on logical and technical level would be more effective than the currently widely used table based representation. This direction goes forward to the database layer, where the object oriented database solution would handle better the ERP or OLTP requirements than the relational database management systems.

Our paper collects the availability of the offered technologies, which could be used for ERP systems and architectures as future environment and technical solution. We try to refer to data-warehouse systems as well pointing out the differences in architecture and requirements positioning the possible directions. We figure out the advantages and disadvantages of the different technologies, approaches and try to sketch a way, which leads to new opportunities of data store for ERP systems.

Acknowledgements

This work was supported by TAMOP-4.2.2/B-10/1-2010-0012 Programme.

Camera Placement Optimization in Object Localization Systems

Dávid Szalóki

This paper discusses the placement of cameras in order to achieve the highest possible localization accuracy. It is reached by using several cameras with redundant fields of views. A camera model is introduced and the components which cause the localization errors are identified. The localization accuracy measure is defined for one and for multiple cameras too. The problem of adding a new camera to the system in order to improve the accuracy is formulated. The method for finding the optimal placement of this new camera is presented. Some features are applied for getting an advanced method for optimizing the placement of multiple cameras. In this area two optimization algorithms are introduced and examined.

Assume that we have a camera system containing fixed and movable cameras. We would like to track an object in the world as accurately as possible. The movable cameras can be placed by the system optimally so that the localization accuracy gets it's highest possible value. The position and orientation constraints of the movable cameras are the limitations. We would like to suggest a method for placing the movable cameras in order to get the highest possible localization accuracy. The localization accuracy is defined as the largest eigenvalue of the inverse of the resulting covariance matrix.

At first a camera model has to be formulated. The components which cause the localization inaccuracy can be identified. The localization accuracy has to be defined and a measure has to be chosen. The nature of these has to be examined for better understanding their behavior. As a first step they can be formulated in 2D for one observed point. Later this model can be generalized into 3D and for an observed area or volume instead of one observed point. The optimal placement of one camera can be formulated. The optimization of multiple cameras together can be deduced from this one camera case. Finally, two algorithms are formulated and examined.

Acknowledgements

This work was partially supported by the European Union and the European Social Fund through project FuturICT.hu (grant no.: TAMOP-4.2.2.C-11/1/KONV-2012-0013) organized by VIKING Zrt. Balatonfüred.

This work was partially supported by the Hungarian Government, managed by the National Development Agency, and financed by the Research and Technology Innovation Fund (grant no.: KMR_12-1-2012-0441).

References

- [1] Dávid Szalóki and Norbert Koszó and Kristóf Csorba and Gábor Tevesz, Marker localization with a multi-camera system, 2013 IEEE International Conference on System Science and Engineering (ICSSE)
- [2] Kristóf Csorba and Dávid Szalóki and Norbert Koszó, Towards a Visual Sensor Network Built from Smartphones, 3rd Eastern European Regional Conference on the Engineering of Computer Based Systems (ECBS-EERC)
- [3] Dávid Szalóki and Norbert Koszó and Kristóf Csorba and Gábor Tevesz, Optimizing Camera placement for localization Accuracy, 14th IEEE International Symposium on Computational Intelligence and Informatics (CINTI)

Challenges in Real-time Collaborative Editing

Péter Szűcs

Collaborative software or groupware is a software that is designed to support a group of people in achieving a common goal [6]. This task typically includes viewing and editing shared media, which can be basically anything from a text document to a CAD model provided collaborative editing makes sense. A real-time collaborative editor allows its users who are connected by some communication network, like the internet, to view and edit this shared media in a parallel fashion regardless of their geographical locations. Consistency maintenance proved to be the most significant challenge in designing and implementing these kinds of systems. Due to this, collaborative editing has been research topic for over twenty years. This paper aims to go through some of the milestones of the past two decades, giving a brief introduction to a popular technique used in present-day implementations, called Operational Transformation, and also, to outline the concept of a framework that could support such systems.

The term Operational Transformation was introduced by S.J. Gibbs and C.A. Ellis in 1989 and the technique was first used in the GROVE (GRoup Outline Viewer Editor) system which is an outline editor that enables its users to view and edit textual outlines simultaneously [1]. Since it was first introduced numerous research groups have contributed to the technique and created their own implementations. One of them is the Jupiter system [4] (developed at Xerox-PARC) which later lead to Google's wave protocol [5].

In Operational Transformation we take an edit-based approach, which means that we capture user actions and mirror them across the network to other users [2]. It's important note that in order to achieve a consistent state, all actions must be captured. The richness of modern user interfaces can make this rather problematic, since besides basic operations like textual insertion and deletion there are more complex actions that are much more difficult to handle, like automatic completion or drag & drop [2]. However, for simplicity, in this paper we will stick to basic textual editing that only involves insertion and deletion of text.

The basic idea is that changes to the document can be modeled as *operations*. When the user changes the document the changes are usually immediately applied locally, for fast response then an operation that represents this change is then sent to the other users. Since the others might have changed the document as well the incoming operation may be out of context. For example, if the incoming operation would insert character 'a' at position 2, but we deleted a character at position 0, the operation should be changed to refer to position 1 instead. In such cases, the incoming operation needs to be *transformed* first and then executed [3].

While it may seem simple, there are several difficulties to be handled in order to achieve a consistent document state at every user. In a complex scenario operations may arrive in different orders for every user, still we have to ensure that when no operations occur, the state of the edited document is the same for every user and that this state is consistent with what the users wanted to achieve. This isn't a trivial task even when we only deal with simple operations like textual insertion and deletion, not to mention more complex operations like drag & drop.

In this paper, we will see the basic concepts of Operational Transformation, some of the difficulties regarding consistency maintenance, solutions to them, as well as possible alternatives. In the contributions section we will also see a concept of a framework to support such systems. The framework is based on the Jupiter system and also uses ideas from the Google Wave white papers.

References

- [1] C.A. Ellis, S.J. Gibbs. *Concurrency Control in Groupware Systems*, Proceedings of the 1989 ACM SIGMOD international conference on Management of data, 1989

- [2] Neil Fraser. *Differential Synchronization*, Proceedings of the 9th ACM symposium on Document engineering, 2009
- [3] Chengzheng Sun, Clearence Ellis. *Operational Transformation in Real-Time Group Editors: Issues, Algorithms, and Achievements*, Proceedings of the 1998 ACM conference on Computer, 1998
- [4] David A. Nichols, Pavel Curtis, Michael Dixon, John Lamping. *High-Latency, Low-Bandwidth Windowing in the Jupiter Collaboration System*, Proceedings of the 8th annual ACM symposium on User interface and software technology, 1995
- [5] David Wang, Alex Mah, Soren Lassen. *Google Wave Operational Transformation*, 2010
- [6] *Collaborative software*, Wikipedia, the free encyclopedia, http://en.wikipedia.org/wiki/Collaborative_software, 2014

Tumor detection and segmentation on multimodal medical images

Szabolcs Urbán, Antal Nagy, László Ruskó

The number of patients living with cancer is increasing from year to year. Therefore tumor detection and segmentation are important tasks in daily clinical practice. Due to large variety of tumors in localization, shape, size and heterogeneity the manual detection and quantification of tumors are hard and also time-consuming even for experienced physicians. Computer-assisted methods can improve these challenging tasks to help the physician's work in tumor detection, quantification and staging. For these purposes there are many published algorithms, each one's has advantages as well as limitations. Most of the methods are using only one image series however in clinical practice several image series and even more modalities are used simultaneously.

In this work we focused on using more than one image series and/or modalities to produce high accuracy tumor detection and segmentation. The automated detection is based on image series which gives functional information about the tumor (e.g. DWI, PET). The detection incorporates intensities and symmetrical analysis of the abnormal regions. The segmentation is based on image series which provide detailed anatomical and soft tissue information (e.g. MRI, CT). Due to large variety of tumor shapes a new boundary based method is used.

The proposed method enables fast and accurate detection and segmentation when more than one image series are incorporated. The generated result can be manually adjusted by the user based on the information from the input images. According to our preliminary tests involving few clinical cases higher segmentation accuracy obtainable using the complementary information. In the future more evaluations are needed incorporating gold standard segmentation defined manually by physicians.

Expanding Small Corpora to Aid People with Communication Impairments

Gyula Vörös

There are people with various movement and cognitive disorders who are unable to speak or write. Many of them are able to communicate with pictorial symbols, but this still narrows the circle of possible communication partners to those who are familiar with this method.

The effects of this problem can be diminished by translating symbol sequences to sentences in natural language. For example, a sequence of symbols that stand for 'I', 'have' and 'sandwich', should be translated to something like 'I would like to have a sandwich'. An approach that is satisfactory in practice is to define all possible sentences manually, for example, by storing them in a text corpus. The construction of a such a corpus would require a lot a work.

In this paper, it is shown that writing all the sentences is not necessary; only a small initial seed corpus is required, which can be expanded automatically. A method is proposed for this expansion, based on replacing the nouns of the input sentences with other nouns. The resulting candidate sentences are then filtered using n-gram statistics from a much larger corpus. This is similar to the 5gramSum method described in [1]. The sentences in the expanded corpus will be similar to the initial ones in structure, but different in content. The method is language-independent, although using it in agglutinative languages such as Hungarian would require morphological processing.

The error rate of the proposed method was evaluated on a small corpus, which was collected from an English language learning website, and contained dialogues. The size of the corpus was doubled by introducing new sentences using our method. For example, the sentence 'I would like to buy some bread' was generated from the original 'I would like to buy some beef'. Samples of 100-100 sentences were randomly selected from the original and the new corpora. Additionally, as a baseline, 100 sentences were generated from the original corpus by replacing nouns randomly, without filtering them using n-gram statistics. Two annotators evaluated each sentence from the samples. The results show that the majority of the introduced sentences were potentially useful. The method produced 3-4 times as many good sentences as the baseline. This indicates that it may be possible to define a small corpus, extend it automatically, and use the resulting set of sentences for communication.

The practical applicability of the method was demonstrated by implementing a sentence production prototype software for alternative communication. A symbol set was defined that enables communication in a food buying situation. A small corpus of appropriate sentences was defined manually, then expanded automatically by including other words from the symbol set. The system was able to produce new meaningful sentences. This way, the amount of manual work necessary to create a communication aid was reduced considerably.

Acknowledgements

This work was carried out as part of the EITKIC 12-1-2012-0001 project, which is supported by the Hungarian Government, managed by the National Development Agency, financed by the Research and Technology Innovation Fund and was performed in cooperation with the EIT ICT Labs Budapest Associate Partner Group.

I would like to thank the work of the people involved with the project, especially András Lőrincz, András Sárkány, Anita Verő, Balázs Pintér and Brigitta Miksztai-Réthey.

References

- [1] Sinha, R. and Mihalcea, R.: Combining lexical resources for contextual synonym expansion, in *Proceedings of the International Conference RANLP*, pp. 404–410 (2009).

LIST OF AUTHORS

- Antal, Elvira:** University of Szeged, Institute of Informatics, Hungary,
E-mail: antale@inf.u-szeged.hu
- Asztalos, Márk:** Budapest University of Technology and Economics, Department of Automation and Applied Informatics, Hungary, E-mail: Asztalos.Mark@aut.bme.hu
- Árgilán, Viktor:** University of Szeged, Juhász Gyula Faculty of Education, Hungary,
E-mail: gilan@jgyppk.u-szeged.hu
- Balogh, Gergő:** University of Szeged, Institute of Informatics, Hungary,
E-mail: geryxyz@inf.u-szeged.hu
- Balogh, János:** University of Szeged, Juhász Gyula Faculty of Education, Hungary,
E-mail: balogh@jgyppk.u-szeged.hu
- Balázs, Péter:** University of Szeged, Institute of Informatics, Hungary,
E-mail: pbalazs@inf.u-szeged.hu
- Barath, Aron:** Eötvös Loránd University, Faculty of Informatics, Hungary,
E-mail: baratharon@caesar.elte.hu
- Bodnár, Péter:** University of Szeged, Institute of Informatics, Hungary,
E-mail: bodnaar@inf.u-szeged.hu
- Borsi, Zsolt:** Eötvös Loránd University, Faculty of Informatics, Hungary,
E-mail: bzsor@inf.elte.hu
- Békési, József:** University of Szeged, Juhász Gyula Faculty of Education, Hungary,
E-mail: bekesi@jgyppk.u-szeged.hu
- Charaf, Hassan:** Budapest University of Technology and Economics, Department of Automation and Applied Informatics, Hungary, E-mail: hassan@aut.bme.hu
- Dávid, Balázs:** University of Szeged, Juhász Gyula Faculty of Education, Hungary,
E-mail: davidb@jgyppk.u-szeged.hu
- Dévai, Richárd:** University of Szeged, Institute of Informatics, Hungary,
E-mail: devai@frontendart.com
- Ekler, Péter:** Budapest University of Technology and Economics, Department of Automation and Applied Informatics, Hungary, E-mail: peter.ekler@aut.bme.hu
- Faragó, Csaba:** University of Szeged, Institute of Informatics, Hungary,
E-mail: farago@inf.u-szeged.hu
- Fekete, István:** Eötvös Loránd University, Faculty of Informatics, Hungary,
E-mail: feketef@inf.elte.hu
- Fodor, Bálint:** Budapest University of Technology and Economics, Department of Automation and Applied Informatics, Hungary, E-mail: balint.fodor@gmail.com
- Gergely, Tamás:** University of Szeged, Institute of Informatics, Hungary,
E-mail: gertom@inf.u-szeged.hu

Giachetta, Roberto: Eötvös Loránd University, Faculty of Informatics, Hungary,
E-mail: groberto@inf.elte.hu

Gilián, Zoltán: Eötvös Loránd University, Faculty of Informatics, Hungary,
E-mail: zoltan.gilian@gmail.com

Gingl, Zoltán: University of Szeged, Institute of Informatics, Hungary,
E-mail: gingl@inf.u-szeged.hu

Griechisch, Erika: University of Szeged, Institute of Informatics, Hungary,
E-mail: grerika@inf.u-szeged.hu

Hajdu, László: University of Szeged, Juhász Gyula Faculty of Education, Hungary,
E-mail: hajdul@jgypk.u-szeged.hu

Hantos, Norbert: University of Szeged, Institute of Informatics, Hungary,
E-mail: nhantos@inf.u-szeged.hu

Horváth, Ferenc: University of Szeged, Institute of Informatics, Hungary,
E-mail: hferenc@inf.u-szeged.hu

Hudák, Péter: Budapest University of Technology and Economics, Department of Automation and Applied Informatics, Hungary, E-mail: hudakpe@gmail.com

Imre, Gábor: Budapest University of Technology and Economics, Department of Automation and Applied Informatics, Hungary, E-mail: imre.gabesz@aut.bme.hu

Karácsony, Máté: Eötvös Loránd University, Faculty of Informatics, Hungary,
E-mail: k_mate@inf.elte.hu

Katona, Melinda: University of Szeged, Institute of Informatics, Hungary,
E-mail: katona.melinda@stud.u-szeged.hu

Kazi, Sándor: Budapest University of Technology, Department of Telecommunication and Media Informatics, Hungary, E-mail: kazi@tmit.bme.hu

Kelemen, Zoltán: Eötvös Loránd University, Faculty of Informatics, Hungary,
E-mail: kelemzol@elte.hu

Kovács, György: MTA-SzTE Research Group on Artificial Intelligence, Hungary,
E-mail: wirth6@gmail.com

Krész, Miklós: University of Szeged, Juhász Gyula Faculty of Education, Hungary,
E-mail: kresz@jgypk.u-szeged.hu

Kundra, László: Budapest University of Technology and Economics, Department of Automation and Applied Informatics, Hungary, E-mail: laszlo.kundra@aut.bme.hu

Ladányi, Gergely: University of Szeged, Institute of Informatics, Hungary,
E-mail: lgergely@inf.u-szeged.hu

Lengyel, László: Budapest University of Technology and Economics, Department of Automation and Applied Informatics, Hungary, E-mail: Lengyel.Laszlo@aut.bme.hu

London, András: University of Szeged, Institute of Informatics, Hungary,
E-mail: london@inf.u-szeged.hu

Lévai, Balázs L.: University of Szeged, Institute of Informatics, Hungary,
E-mail: levaib@inf.u-szeged.hu

Mezei, Gergely: Budapest University of Technology and Economics, Department of Automation and Applied Informatics, Hungary, E-mail: gmezei@aut.bme.hu

Mészáros, Benjamin: University of Szeged, Institute of Informatics, Hungary,
E-mail: meszaros.benjamin@stud.u-szeged.hu

Nagy, Antal: University of Szeged, Institute of Informatics, Hungary,
E-mail: nagya@inf.u-szeged.hu

Nagy, Gábor: Budapest University of Technology and Economics, Department of Telecommunications and Media Informatics, Hungary, E-mail: nagyg@tmit.bme.hu

Nagy, Tamás Dániel: University of Szeged, Institute of Informatics, Hungary,
E-mail: nag.tams@gmail.com

Nyúl, László G.: University of Szeged, Institute of Informatics, Hungary,
E-mail: nyul@inf.u-szeged.hu

Németh, Boldizsár: Eötvös Loránd University, Faculty of Informatics, Hungary,
E-mail: nboldi@caesar.elte.hu

Németh, Gábor: University of Szeged, Institute of Informatics, Hungary,
E-mail: gnemeth@inf.u-szeged.hu

Németh, Zoltán: Corvinus University of Budapest, Hungary,
E-mail: zoltan.nemeth4@uni-corvinus.hu

Orosz, István: Óbudai Univesity, Hungary, E-mail: istvan.orosz@sznet.hu

Orosz, Tamás: Óbudai Univesity, Hungary, E-mail: orosz.tamas@arek.uni-obuda.hu

Ozsvár, Zoltán: University of Szeged, Institute of Informatics, Hungary,
E-mail: ozsvar.zoltan@gmail.com

Papp, Győző: Budapest University of Technology and Economics, Department of Telecommunications and Media Informatics, Hungary, E-mail: pgyozo@tmit.bme.hu

Paróczy, Zsombor: Budapest University of Technology and Economics, Department of Telecommunications and Media Informatics, Hungary, E-mail: paroczi@tmit.bme.hu

Porkoláb, Zoltán: Eötvös Loránd University, Faculty of Informatics, Hungary,
E-mail: gsd@elte.hu

Pándi, Krisztián: Budapest University of Technology and Economics, Department of Automation and Applied Informatics, Hungary, E-mail: pandi.krisztian@aut.bme.hu

Rudas, László: University of Szeged, Faculty of Medicine, Hungary

Selmeci, Attila: Óbudai Univesity, Hungary, E-mail: selmeci.attila@arek.uni-obuda.hu

Szalóki, Dávid: Budapest University of Technology and Economics, Department of Automation and Applied Informatics, Hungary, E-mail: Szaloki.David@aut.bme.hu

Szűcs, Gábor: Budapest University of Technology and Economics, Department of Telecommunications and Media Informatics, Hungary, E-mail: szucs@tmit.bme.hu

Szűcs, Péter: Budapest University of Technology and Economics, Hungary,

E-mail: pszucs@outlook.com

T. Nagy, István: University of Szeged, Institute of Informatics, Hungary,

E-mail: nistvan@inf.u-szeged.hu

Tóth, Attila: University of Szeged, Juhász Gyula Faculty of Education, Hungary,

E-mail: attila@jgypk.u-szeged.hu

Tóth, László: MTA-SzTE Research Group on Artificial Intelligence, Hungary,

E-mail: tothl@inf.u-szeged.hu

Urbán, Szabolcs: University of Szeged, Institute of Informatics, Hungary,

E-mail: urbansz@inf.u-szeged.hu

Vadai, Gergely: University of Szeged, Institute of Informatics, Hungary,

E-mail: vadaig@inf.u-szeged.hu

Vinkó, Tamás: University of Szeged, Institute of Informatics, Hungary,

E-mail: tvinko@inf.u-szeged.hu

Vörös, Gyula: Eötvös Loránd University, Faculty of Informatics, Hungary,

E-mail: vorosgy@gmail.com

Zöllei, Éva: University of Szeged, Faculty of Medicine, Hungary

NOTES

CSCS²

Supported by Telemedicine Oriented Research in the Fields of Mathematics, Informatics and
Medical Sciences T ÁMOP-4.2.2.A-11/1/KONV-2012-0073.

National Development Agency
www.ujszchenyiterv.gov.hu
06 40 638 638



HUNGARY'S RENEWAL



The project is supported by the European Union
and co-financed by the European Social Fund.

