

ERCIM



NEWS

www.ercim.eu

Special theme:

Machine Learning

Also in this issue:

Research and Society:

Open Access
Open Science

ERCIM News is the magazine of ERCIM. Published quarterly, it reports on joint actions of the ERCIM partners, and aims to reflect the contribution made by ERCIM to the European Community in Information Technology and Applied Mathematics. Through short articles and news items, it provides a forum for the exchange of information between the institutes and also with the wider scientific community. This issue has a circulation of about 6,000 printed copies and is also available online.

ERCIM News is published by ERCIM EEIG
BP 93, F-06902 Sophia Antipolis Cedex, France
Tel: +33 4 9238 5010, E-mail: contact@ercim.eu
Director: Jérôme Chailloux, ISSN 0926-4981

Contributions

Contributions should be submitted to the local editor of your country

Copyright notice

All authors, as identified in each article, retain copyright of their work. ERCIM News is licensed under a Creative Commons Attribution 4.0 International License (CC-BY).

Advertising

For current advertising rates and conditions, see <http://ercim-news.ercim.eu/> or contact peter.kunz@ercim.eu

ERCIM News online edition

<http://ercim-news.ercim.eu/>

Next issue

January 2017, Special theme: Computational Imaging

Subscription

Subscribe to ERCIM News by sending an email to en-subscriptions@ercim.eu or by filling out the form at the ERCIM News website: <http://ercim-news.ercim.eu/>

Editorial Board:

Central editor:

Peter Kunz, ERCIM office (peter.kunz@ercim.eu)

Local Editors:

Austria: Erwin Schoitsch (erwin.schoitsch@ait.ac.at)

Belgium: Benoît Michel (benoit.michel@uclouvain.be)

Cyprus: Ioannis Krikidis (krikidis.ioannis@ucy.ac.cy)

Czech Republic: Michal Haindl (haindl@utia.cas.cz)

France: Steve Kremer (steve.kremer@inria.fr)

Germany: Michael Krapp

(michael.krapp@scai.fraunhofer.de)

Greece: Eleni Orphanoudakis (eleni@ics.forth.gr), Artemios

Voyiatzis (bogart@isi.gr)

Hungary: Andras Benczur (benczur@info.ilab.sztaki.hu)

Italy: Carol Peters (carol.peters@isti.cnr.it)

Luxembourg: Thomas Tamisier (thomas.tamisier@list.lu)

Norway: Poul Heegaard (poul.heegaard@item.ntnu.no)

Poland: Hung Son Nguyen (son@mimuw.edu.pl)

Portugal: José Borbinha, Technical University of Lisbon
(jl@ist.utl.pt)

Spain: Silvia Abrahão (sabrahao@dsic.upv.es)

Sweden: Kersti Hedman (kersti@sics.se)

Switzerland: Harry Rudin (hrudin@smile.ch)

The Netherlands: Annette Kik (Annette.Kik@cwi.nl)

W3C: Marie-Claire Forgue (mcf@w3.org)

RESEARCH AND SOCIETY

The section “Research and Society” on “Open Access – Open Science” has been coordinated by Laurent Romary (Inria)

5 Open Science: Taking Our Destiny into Our Own Hands
by Laurent Romary (Inria)

6 ERCIM Goes to Open Access
by Jos Baeten (CWI) and Claude Kirchner (Inria)

7 Will Europe Liberate Knowledge through Content Mining?
by Peter Murray-Rust (University of Cambridge)

9 Roads to Open Access: The Good, the Bad and the Ugly
by Karim Ramdani (Inria)

10 Open-Access Repositories and the Open Science Challenge
by Leonardo Candela, Paolo Manghi, and Donatella Castelli (ISTI-CNR)

11 LIPIcs – an Open-Access Series for International Conference Proceedings
by Marc Herbstritt (Schloss Dagstuhl – Leibniz-Zentrum für Informatik) and Wolfgang Thomas (RWTH Aachen University)

13 Scientific Data and Preservation – Policy Issues for the Long-term Record
by Vera Sarkol (CWI)

14 Mathematics in Open Access – MathOA
by Johan Rooryck and Saskia de Vries

SPECIAL THEME

The special theme section “Machine Learnig” has been coordinated by Sander Bohte (CWI) and Hung Son Nguyen (University of Warsaw)

Introduction to the Special Theme

16 Modern Machine Learning: More with Less, Cheaper and Better
by Sander Bohte (CWI) and Hung Son Nguyen (University of Warsaw)

More with less

18 Micro-Data Learning: The Other End of the Spectrum
by Jean-Baptiste Mouret (Inria)

19 Making Learning Physical: Machine Intelligence and Quantum Resources
by Peter Wittek (ICFO-The Institute of Photonic Sciences and University of Borås)

20 Marrying Graphical Models with Deep Learning
by Max Welling (University of Amsterdam)

22 Privacy Aware Machine Learning and the “Right to be Forgotten”
by Bernd Malle, Peter Kieseberg (SBA Research), Sebastian Schrittwieser (JRC TARGET, St. Poelten University of Applied Sciences), and Andreas Holzinger (Graz University of Technology)

24 Robust and Adaptive Methods for Sequential Decision Making
by Wouter M. Koolen (CWI)

Research

25 Neural Random Access Machines
by Karol Kurach (University of Warsaw and Google), Marcin Andrychowicz and Ilya Sutskever (OpenAI (work done while at Google))

26 Mining Similarities and Concepts at Scale
by Olof Görnerup and Theodore Vasiloudis (SICS)

28 Fast Traversal of Large Ensembles of Regression Trees
by Claudio Lucchese, Franco Maria Nardini, Raffaele Perego, Nicola Tonello (ISTI-CNR), Salvatore Orlando (University of Venice) and Rossano Venturini (University of Pisa)

Massive data processing

29 Optimising Deep Learning for Infinite Applications in Text Analytics

by Mark Cieliebak (Zurich University of Applied Sciences)

31 Towards Streamlined Big Data Analytics

by András A. Benczúr, Róbert Pálovics (MTA SZTAKI), Márton Balassi (Cloudera), Volker Markl, Tilmann Rabl, Juan Soto (DFKI), Björn Hovstadius, Jim Dowling and Seif Haridi (SICS)

How does the brain do it?

32 Autonomous Machine Learning

by Frederic Alexandre (Inria)

34 Curiosity and Intrinsic Motivation for Autonomous Machine Learning

by Pierre-Yves Oudeyer, Manuel Lopes (Inria), Celeste Kidd (Univ. of Rochester) and Jacqueline Gottlieb (Univ. of Columbia)

Applications

35 Applied Data Science: Using Machine Learning for Alarm Verification

by Jan Stampfli and Kurt Stockinger (Zurich University of Applied Sciences)

37 Towards Predictive Pharmacogenomics Models

by George Potamias (FORTH)

38 Optimisation System for Cutting Continuous Flat Glass

by José Francisco García Cantos, Manuel Peinado, Miguel A. Salido and Federico Barber (AI2-UPV)

40 Online Learning for Aggregating Forecasts in Renewable Energy Systems

by Balázs Csanád Csáji, András Kovács and József Váncza (MTA SZTAKI)

42 Bonaparte: Bayesian Networks to Give Victims back their Names

by Bert Kappen and Wim Wiegerinck (University Nijmegen)

RESEARCH AND INNOVATION

This section features news about research activities and innovative developments from European research institutes

44 BASMATI – Cloud Brokerage Across Borders For Mobile Users And Applications

by Patrizio Dazzi (ISTI-CNR)

46 An Incident Management Tool for Cloud Provider Chains

by Martin Gilje Jaatun, Christian Frøystad and Inger Anne Tøndel (SINTEF ICT)

48 Predictive Modelling from Data Streams

by Olivier Parisot and Benoît Otjacques (Luxembourg Institute of Science and Technology)

49 Mandola: Monitoring and Detecting Online Hate Speech

by Marios Dikaiakos, George Pallis (University of Cyprus) and Evangelos Markatos (FORTH)

51 The BÆSE Testbed – Analytic Evaluation of IT Security Tools in Specified Network Environments

by Markus Wurzenberger and Florian Skopik (AIT Austrian Institute of Technology)

53 Behaviour-Based Security for Cyber-Physical Systems

by Dimitrios Serpanos (University of Patras and ISI), Howard Shrobe (CSAIL/MIT) and Muhammad Taimoor Khan (University of Klagenfurt)

54 The TISRIM-Telco Toolset – An IT Regulatory Framework to Support Security Compliance in the Telecommunications Sector

by Nicolas Mayer, Jocelyn Aubert, Hervé Cholez, Eric Grandry and Eric Dubois

56 Predicting the Extremely Low Frequency Magnetic Field Radiation Emitted from Laptops: A New Approach to Laptop Design

by Darko Brodić, Dejan Tanikić (University of Belgrade), and Alessia Amelio (University of Calabria)

57 Managing Security in Distributed Computing: Self-Protective Multi-Cloud Applications

by Erkuden Rios (Tecnalia), Massimiliano Rak (Second University of Naples) and Samuel Olaiya Afolaranmi (Tampere University of Technology)

EVENTS, IN BRIEF

Announcements

59 VaMoS 2017: 11th International Workshop on Variability Modelling of Software-intensive Systems

In Brief

59 2016 Internet Defense Prize for Quantum-safe Cryptography

ERCIM

Membership

After having successfully grown to become one of the most recognized ICT Societies in Europe, ERCIM has opened membership to multiple member institutes per country. By joining ERCIM, your research institution or university can directly participate in ERCIM's activities and contribute to the ERCIM members' common objectives playing a leading role in Information and Communication Technology in Europe:

- Building a Europe-wide, open network of centres of excellence in ICT and Applied Mathematics;
- Excelling in research and acting as a bridge for ICT applications;
- Being internationally recognised both as a major representative organisation in its field and as a portal giving access to all relevant ICT research groups in Europe;
- Liaising with other international organisations in its field;
- Promoting cooperation in research, technology transfer, innovation and training.

About ERCIM

ERCIM – the European Research Consortium for Informatics and Mathematics – aims to foster collaborative work within the European research community and to increase cooperation with European industry. Founded in 1989, ERCIM currently includes 21 leading research establishments from 18 European countries. Encompassing over 10,000 academics and researchers, ERCIM is able to undertake consultancy, development and educational projects on any subject related to its field of activity.

ERCIM members are centres of excellence across Europe. ERCIM is internationally recognized as a major representative organization in its field. ERCIM provides access to all major Information Communication Technology research groups in Europe and has established an extensive program in the fields of science, strategy, human capital and outreach. ERCIM publishes ERCIM News, a quarterly high quality magazine and delivers annually the Cor Baayen Award to outstanding young researchers in computer science or applied mathematics. ERCIM also hosts the European branch of the World Wide Web Consortium (W3C).

“Through a long history of successful research collaborations in projects and working groups and a highly-selective mobility programme, ERCIM has managed to become the premier network of ICT research institutions in Europe. ERCIM has a consistent presence in EU funded research programmes conducting and promoting high-end research with European and global impact. It has a strong position in advising at the research policy level and contributes significantly to the shaping of EC framework programmes. ERCIM provides a unique pool of research resources within Europe fostering both the career development of young researchers and the synergies among established groups. Membership is a privilege.”

Dimitris Plexousakis, ICS-FORTH, ERCIM AISBL Board

Benefits of Membership

As members of ERCIM AISBL, institutions benefit from:

- International recognition as a leading centre for ICT R&D, as member of the ERCIM European-wide network of centres of excellence;
- More influence on European and national government R&D strategy in ICT. ERCIM members team up to speak with a common voice and produce strategic reports to shape the European research agenda;
- Privileged access to standardisation bodies, such as the W3C which is hosted by ERCIM, and to other bodies with which ERCIM has also established strategic cooperation. These include ETSI, the European Mathematical Society and Informatics Europe;
- Invitations to join projects of strategic importance;
- Establishing personal contacts with executives of leading European research institutes during the bi-annual ERCIM meetings;
- Invitations to join committees and boards developing ICT strategy nationally and internationally;
- Excellent networking possibilities with more than 10,000 research colleagues across Europe. ERCIM's mobility activities, such as the fellowship programme, leverage scientific cooperation and excellence;
- Professional development of staff including international recognition;
- Publicity through the ERCIM website and ERCIM News, the widely read quarterly magazine.

How to Become a Member

- Prospective members must be outstanding research institutions (including universities) within their country;
- Applicants should address a request to the ERCIM Office. The application should include:
 - Name and address of the institution;
 - Short description of the institution's activities;
 - Staff (full time equivalent) relevant to ERCIM's fields of activity;
 - Number of European projects in which the institution is currently involved;
 - Name of the representative and a deputy.
- Membership applications will be reviewed by an internal board and may include an on-site visit;
- The decision on admission of new members is made by the General Assembly of the Association, in accordance with the procedure defined in the Bylaws (<http://kwz.me/U7>), and notified in writing by the Secretary to the applicant;
- Admission becomes effective upon payment of the appropriate membership fee in each year of membership;
- Membership is renewable as long as the criteria for excellence in research and an active participation in the ERCIM community, cooperating for excellence, are met.

Please contact the ERCIM Office: contact@ercim.eu

Open Science: Taking Our Destiny into Our Own Hands

by Laurent Romary (Inria)

There is currently a tug-of-war going on within the arena of scientific communication: scientists are exploring new, more efficient and affordable ways to disseminate research results, but at the same time, a web of private publishing companies (and even learned societies) are endeavouring to preserve their financial turnover on the basis of models from a previous era. This tension is echoed in the recent news relating to scholarly communication within Europe as a whole, and within individual countries:

- Various legislative initiatives have been launched to improve legal copyright settings with various degrees of success. Julia Reda's extremely ambitious proposal to reform the European copyright regulation does not seem to be reflected in the most recent drafts by the commission. On the contrary, new digital legislation is likely to be adopted in France in the coming weeks, with articles on both the freedom to deposit authors' manuscripts in publication repositories and data mining freedom for legally acquired material;
- Open science has been high on the agenda of the Dutch EU presidency during the first semester of 2016, and the final press release [L1] clearly states the objective that all scholarly papers should be freely available online by 2020. However, we have no defined strategy to guide us towards this ambitious goal and, at the same time, extremely conservative initiatives such as OA2020, riding on a tenuous connection to EU policy, are attempting to preserve the publishing landscape in its current state;
- There have been recent instances of large private publishing trusts acquiring other companies to enlarge the scope of their services and thus their grasp on our communication facilities. Elsevier has recently taken over SSRN, a publication repository in social sciences and Hivebench, a laboratory notebook platform, just a few months after acquiring Mendeley, a major online information management site.

Without providing a comprehensive overview of how this situation arose, we can identify a few milestones that may help to explain why many researchers and institutions have started to question the adequacy of the contemporary publication system.

Within Europe, the first real sign of a strong awareness of diverging interests between the scientific community and scholarly publishers dates back to 2006 when a petition [L2] of more than 28,000 signatures, including many higher education and research institutions, was fiercely answered by a

communiqué from the International Association of Scientific, Technical & Medical Publishers (STM) warning against the EU issuing any kind of open-access policy [L3]. Since then the EU has actually funded the OpenAire initiative and above all designed a mandatory open-access policy for all publications financed within its H2020 program.

The private sector has also taken up the open-access agenda and now presents itself as the key actor in the development of an economically viable solution with the author-pays model. Unfortunately, some countries have adopted this as a reference for the development of their public policies, as we have seen with the Finch report. Even the recent declaration by the League of European Research Universities LERU [L4] refers to the 'transition', a term that is inextricably linked to the dialectic of moving the subscription-based landscape to an author-pays scenario.

Finally, many new private actors are setting up online services related to communication (Academia, ResearchGate) or assessment (F1000, ScienceOpen, My ScienceWorks, peer.us) of scholarly content. It is alarming at times to see how much content is being redirected to such platforms, whose confidentiality and sustainability are far from guaranteed.

Given the current situation, shouldn't we be concerned about the relatively low level of involvement of research institutions in directing the evolution of science communication? Is it really wise to hand over the reins of publication repositories and associated services (surely a vital part of our research infrastructure) to private ventures?

It makes sense for the computer science and mathematics community to be at the forefront of any initiative that relies heavily on information technology. The question at hand for ERCIM is to determine how professionals within this community might play a leading role in designing new models for the dissemination of research results that could ensure a high level of scientific quality, appropriate rewarding of its authors, and be both affordable and sustainable in the long term.

To this end, we at Inria have designed and implemented an ambitious open-access policy based on two main pillars:

- a full-text deposit mandate on the French national repository HAL (<http://hal.inria.fr>) coupled with the annual reporting requirement of our institution;

- a cautious approach to the author-pays model, with the setting up of a central budget for a native open-access journal and a ban on ‘hybrid payment’, i.e. journals that are also based on subscriptions.

The success of our policy, which is similar to the one deployed at the Dutch ERCIM member CWI, has allowed us to reach very high levels of full text coverage but also to keep article-processing charges low over the last five years. We are also exploring the development of new publication models in collaboration with the CCSD (Centre pour la Communication Scientifique Directe) service unit in Lyon, with the launch of an overlay journal platform: Episciences.org, where we both launched new scientific journals in computer science and applied mathematics, but also migrated legacy publications such as LMCS (Logical Methods in Computer Science) or DMTCS (Discrete Mathematics & Theoretical Computer Science).

The contributions focussing on open access featured in this issue of ERCIM News reflect the variety of doubts and ambitions that have emerged within our community, but also more widely within European academic institutions. We start with a presentation by Jos Baeten and Claude Kirchner of the recommendations approved by the ERCIM board followed by a plea from Peter Murray-Rust for a systematisation of data mining services on scholarly content. Karim Ramdani makes a clear case for implementing a green open-access policy as opposed to models based on the payment of article processing charges, Leonardo Candela, Paolo Manghi and Donatella Castelli show how this requires an increase in service provision and connectivity for existing publication repositories. The role of public initiatives in setting up new publication standards is discussed by Johan Rooryck, in the domain of mathematics, and Marc Herbstritt and Wolfgang Thomas, who advocate for a ‘reconquista’ of our scientific communication means. Finally, Vera Sarkol extends the debate to scientific data, with a look ahead to the necessary infrastructures we have to put in place.

We all have a responsibility to make sure that our discoveries and results are widely available for our colleagues and the general public. It is time for all ERCIM members to take a clear position, but also for each of us, as researchers, to contribute to the debate and ensure we achieve a viable scientific communication scenario for the future.

Links:

[L1] <http://francais.eu2016.nl/documents/persberichten/2016/05/27/communiquede-presse---tous-les-articles-scientifiques-europeens-en-libre-acces-a-partir-de-2020>

[L2] <http://legacy.earlham.edu/~peters/fos/2007/02/20000-signatures-for-oa-presented-to-ec.html>

[L3] <http://legacy.earlham.edu/~peters/fos/2007/02/publishers-issue-brussels-declaration.html>

[L4] <http://data.consilium.europa.eu/doc/document/ST-9526-2016-INIT/en/pdf>

Please contact:

Laurent Romary, Inria, France
laurent.romary@inria.fr

ERCIM Goes to Open Access

by Jos Baeten (CWI) and Claude Kirchner (Inria)

At its October 2014 meeting, the EEIG ERCIM board installed a task group Boost Open Access Mastering (BOM), chaired by us, with the goal of facilitating the sharing of information and the strategies of ERCIM participants in regard to open access. The ensuing report [L1], a plea for author control, which was adopted by the board in October 2015, recommends an open-access strategy and identified tools shared or to be shared by several ERCIM members.

We need change

The current digital revolution is impacting the way science develops and the way we conduct research. The seminal vision of Jim Gray about big data as the fourth paradigm of science [L2] is an excellent entry point to understanding these phenomena, where the initial paradigms of theory building and experimentation are now completed or even replaced by digital simulation and data exploration.

In this profoundly renewed context, the role of scientific data is fundamental. Scientists of all disciplines are completely dependent on the data that allow them to understand, model, experiment, reproduce and communicate.

In the digital world, everything can be seen as source data: a text describing the results of a study, a computer program, a video, a picture, a sound, a MOOC, a lab book, a protocol, a data set captured by an instrument or generated by a computer, and so on. Secondary data or data generated from other data, like discussions, social network information or peer reviews are also crucial sources that may be relevant for further research.

Being in control of data is a matter of scientific sovereignty, and any restriction or hindrance in this respect will be to the detriment of science. Note that control is more than ownership, because ownership is transferable, and if something is sold you can no longer control it. ‘Control’ is used here in terms of ability to read, re-use, quote, analyse a common good. From this point of view, maintaining the sovereignty of scientific academic research is a crucial issue, which we need to preserve in the short as well as the long run.

The services that allow scientific data to be used are crucial. They include data mining, analysis and synthesis for scientific purposes as well as for societal, economic or industrial purposes. In particular they require access to the full texts of scientists’ contributions. Ideally, researchers would be able to make the most of the available data; this is an important goal that either scientists themselves, or public or private entities, should aim towards.

Recommendations

As a consequence, the BOM task group, consisting of J. Baeten, L. Candela, I. Fava, C. Kirchner, W. Mettrop, L. Romary, L. Schultze, makes the following recommendations

which could be adapted to the best practices of each scientific discipline as well as to local legislation, with the goal of making scientific sovereignty an unalterable reality by or before 2020.

Main principles

1. Scientists should maintain control over all their scholarly products (i.e., all the outcomes of their research activities, ranging from their publications — actually the full text — to the datasets they curated/contributed to);
2. The services that value scientific data should be open to competition.

Organisation principles

1. All research institutions should formulate and implement a strategic policy about the proper management of their scholarly outputs. Such policies should mandate scientists to deposit every scholarly product in a suitable open-access repository as soon as the product is produced. The policy should also mention the repositories trusted by the institution;
2. All research institutions should support the development of suitable publishing platforms for their research products (including open-access repositories and overlay journals). Such publishing platforms should be maintained as public infrastructure;
3. Scientists deserve proper credit for their scholarly products. Research institutions should promote and support the development of a comprehensive, scientific community-recognised and innovative set of scholarly products evaluation/assessment criteria.

ERCIM specifics

1. A network of repository and scientific information managers should be set up in order to share experience as well as develop better services related to the various institutions' open-access strategies;
2. ERCIM should be able to access reliable output figures from all institutions, which could then be shared between institutions;
3. A joint dashboard should be set up for sharing article processing charges (APC) across all ERCIM entities: the model suggested by University of Bielefeld could be used;
4. In the name of ERCIM and of each national research institution, address the recommendations of the BOM Report to the highest political level of the EU and of each country;
5. ERCIM should favour the re-use of publication facilities available among its members, such as repositories or overlay journals;
6. The involvement of ERCIM members into the emergence of open-access publication including overlay journals dedicated to data and software should be encouraged.

ERCIM has adopted these recommendations and is working further towards our goals.

Links:

[L1] <http://oai.cwi.nl/oai/asset/23589/23589B.pdf>

[L2] <http://kwz.me/VI>

Please contact:

Jos Baeten, CWI, Jos.Baeten@cwi.nl

Claude Kirchner, Inria, claud.kirchner@inria.fr

Will Europe Liberate Knowledge through Content Mining?

by Peter Murray-Rust (University of Cambridge)

Scholarly publications, especially science and medicine, have huge amounts of untapped knowledge, but it's a technical challenge to extract it and there's a political fight in Europe as to whether we can legally do it.

About three million peer-reviewed scholarly publications and technical reports, especially in life science and medicine, are published each year – one every 10 seconds. Many are filled with facts (species, diseases, drugs, countries, organisations) resulting from about one trillion USD of funded research. But they aren't properly used – for example the Zika outbreak was predicted 30 years ago [1] but in a scanned PDF behind a paywall so was never broadcast. Computers are essential to process this data – but there are major problems: the complexity of semi-structured information and socio-political conflict is being played out in Brussels even as I write this [2].

Scientists write in narrative text, with embedded data and images and no thought for computer processing. Most text is born digital (Word, TeX) and straightforward to process but turned into PDF. Data is collected from digital instruments, summarised to diagrams and turned into pixels (PNG, JPG) with total loss of data – even from summary diagrams (plots). We know of researchers who spend their whole time turning this back into computable information.

However, it's still possible to recover a vast amount of data with heuristics such as natural language processing and diagram processing. With Shuttleworth Foundation funding I've created ContentMine.org [3] to read the whole scientific literature and extract the facts on a daily basis. We've spent two years creating the code and pipeline and are now starting to process open and subscription articles automatically – and this can extend to theses and reports.

The pipeline covers many tasks including crawling and scraping websites, using APIs or papers already aggregated. Papers often need normalising from raw HTML to structured, annotated XHTML and marking up the sections ('Introduction', 'Methods', 'Results', etc) is an important way of reducing false positives. Captions for tables and figures are often the most important parts of some articles. We then search the text by discipline-specific plugins, most commonly using simple dictionaries enhanced by Wikidata. These often exist in current disciplines – e.g., ICD-10 for disease – and increasingly we can extract them directly from Wikidata. More complex tools are required for species and chemistry. And we have pioneered automatic methods for interpreting images of phylogenetic trees and constructed a supertree for 4,500 articles.

Among the sources are EuropePubMedCentral – over one million open articles on life science and medicine, converted

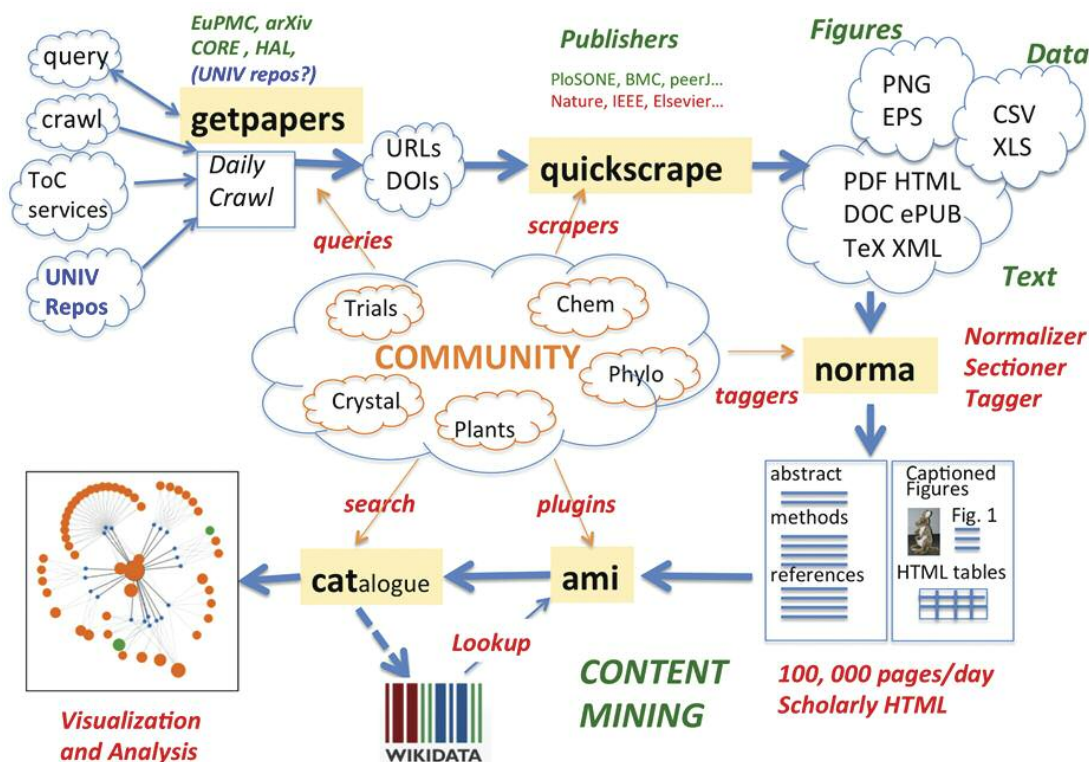


Figure 1: The ContentMine pipeline. Articles can be ingested by systematic crawling, push from a TableOfContents service, or user-initiated query (getpapers and/or quickscape). They are normalised to XHTML, annotated by section and facts extracted by ‘plugins’ and sent to public repositories (Wikidata, Zenodo). All components are Open Source/data.

into XML. Our getpapers tool directly uses EPMC’s search API and feeds text for conversion to scholarlyHTML. We can also get metadata from Crossref and scrape sites directly with per-publisher scrapers – it takes less than an hour to create a new one.

We see Wikidata as the future of much scientific fact, and cooperate with them by creating enhanced dictionaries for searching and also providing possible new entries. The Wikidata-enhanced facts will be stored in the public Zenodo database for all to use. Since facts are uncopyrightable we expect to extract over 100 million per year.

Text and data mining (or ContentMining) has been seen as a massive public good[5]. Sir Mark Walport, director of the Wellcome Trust, said "This is a complete no-brainer. This is scholarly research funded from the public purse, largely from taxpayer and philanthropic organisations. The taxpayer has the right to have maximum benefit extracted and that will only happen if there is maximum access to it."

But there’s huge politico-legal opposition, because the papers are copyrighted, normally by the publishers who see mining as a new revenue stream, even though they have not developed the technology. Innovative scientists carrying out mining risk their universities being cut off by publishers. The UK has pioneered reform to allow mining for non-commercial research, but it was strongly opposed by publishers and there’s little effective practical support. In 2013 organisations such as national libraries, funders, and academics were opposed by rightsholders (‘Licences for Europe’) leading to an impasse. The European parliament has tried to reform copyright, but recommendations have been heavily watered

down by the commission and leaks suggest that formalising the market for exploitation by publishers will be emphasised at the expense of innovation and freedom.

We desperately need open resources – content, dictionaries, software, infrastructure. The UK has led but not done enough. France is actively deciding on its future. Within two years decisions will become effectively irrevocable. Europe must choose whether it wants mining to be done by anyone, or controlled by corporations.

Links:

- [L1] http://www.nytimes.com/2015/04/08/opinion/yes-we-were-warned-about-ebola.html?_r=0
- [L2] <http://kluwercopyrightblog.com/2016/07/20/julia-reda-mep-discusses-harmonisation-copyright-law-ip-enforcement-brexit/>
- [L3] <http://contentmine.org>
- [L4] <http://www.statewatch.org/news/2016/aug/eu-com-copyright-draft.pdf>
- [L5] <https://www.jisc.ac.uk/news/text-mining-promises-huge-economic-and-research-benefit-but-barriers-limit-its-use-14-mar-2012>

Reference:

P Murray-Rust, J Molloy, and D Cabell: “Open content mining”, in Proc. of The First OpenForum Academy Conference, pp. 57-64. OpenForum Europe LTD, 2012.

Please contact:

Peter Murray-Rust, University of Cambridge, UK
+44 1223 336432, pm286@cam.ac.uk

Roads to Open Access: The Good, the Bad and the Ugly

by Karim Ramdani (Inria)

Promoting Open Access without specifying the road chosen to reach it makes no sense. The author-pays road (APC Gold Open Access) is without a doubt the worst option.

The Scientific Board of the French CNRS Institute for Mathematics (INSMI) has recently made the following recommendations to French mathematicians for their publications:

1. Do not choose the author-pays option for open access, especially for hybrid journals (a hybrid journal is a subscription-based journal, in which authors are given the option of paying publication fees (APC) to make their own article freely available);
2. Do not include in funding requests such publication fees (known as APC, author processing charges).

These recommendations perfectly illustrate the rejection of the author-pays model by French mathematicians, and more widely, by European ones [1].

Given that scientists are generally both authors and readers, the reader-pays model (the current dominant subscriptions based model) and the author-pays model (also known as APC Gold Open Access) might seem at first glance symmetrical, and hence equivalent. This is not the case for economic and ethical reasons.

Economic aspects

First, scholarly publishing costs in an author-pays model are higher than in the reader-pays model (whose costs are already unacceptably high). This statement is based on several projections made by French research institutions

(CNRS, INRA) and the data available for the UK [L1] [L2]. At the same time, publishers' costs decrease when moving to an open-access model (no printed versions, no managing fees for subscriptions and accesses rights). Second, the idea that universities will be able to control prices in an author-pays model by introducing competition between publishers is illusory. Indeed, most countries that started moving towards APC Gold Open Access have done so by signing contracts with big commercial publishers. Consequently, as with subscription negotiations today, universities will be in a weak position with no expected benefits from competition: it seems unlikely that any scientist will choose to pay €1,000 to publish with a small independent publisher, when Elsevier and Springer journals publish "for free" (the APC having already been paid at a national level). These economic arguments should disqualify any changeover towards author-pays models: either a partial one in which subscription costs coexist with APC costs (the ugly road to OA) or a complete one where only APC costs exist (the bad road to OA).

Scientific and ethical aspects

The author-pays model is unethical as well as costly. It introduces an unacceptable inequality in access to publishing between scientists (especially if APC expenses are not centralised at a national level). In such a system, only "rich" researchers will be able to publish in the "best" journals, often the most expensive ones (in the UK, the average APC by article was £1,575 in 2014 and £1,762 in 2015, with a maximum APC around £3,200). In return, this will increase their "visibility" and their ability to be funded. Besides introducing such discrimination, the author-pays model also carries ethical risks inherent in its philosophy: why would a journal refuse to publish a paper submitted for publication when its acceptance increases its profit? The answer is obvious, as shown by the emergence of several "predatory publishers" [L3] in recent years.

Good roads to Open Access

The above criticisms echo the recent joint statement on Open Access of UNESCO and COAR [L4], warning both governments and the research community against a large-scale shift from subscriptions to open access via APC. Refusing such a



shift, that will reinforce a historical oligopolistic situation, does not mean that the current situation is satisfactory. Many actions need to be undertaken:

- Denounce the obscene profits of big commercial publishers and protest against their business practices [L5].
- Cancel subscriptions when necessary [L6].
- Develop and promote good roads to OA:
 - green Open Access (articles are placed in a repository and can be freely accessed by all) with its institutional repositories,
 - fair Open Access with its sponsor-pays journals, like Discrete Analysis, Journal de l'École polytechnique or Epiga [L7].
- Create new economic models for scholarly publishing, free of charge for the author and the reader, for instance: using institutional support (Episciences [L8], SciELO [L9]), sale of premium services (e.g., OpenEdition [L10]), crowd-funding (e.g., OLH [L11]), or library subscriptions.
- Fight against the use and abuse of impact factors and bibliometrics and rethink the evaluation process.

Finally, perhaps the first battle we must fight is the one of words. For-profit publishers have appropriated the noble idea of open access to propose through APC Gold Open Access a model that preserves their commercial interests. We must denounce this openwashing [L12] that makes politicians think that all forms of open access are beneficial for scientists and taxpayers. Promoting open access without specifying the road chosen to reach it makes no sense. The author-pays road (APC Gold Open Access) is definitely the worst of them.

Links:

- [L1] <https://www.jisc.ac.uk/sites/default/files/apc-and-subscriptions-report.pdf>
- [L2] <http://www.rcuk.ac.uk/documents/documents/openaccessreport-pdf/>
- [L3] <https://scholarlyoa.com/2015/01/02/bealls-list-of-predatory-publishers-2015/>
- [L4] http://www.unesco.org/new/fileadmin/MULTIMEDIA/HQ/CI/CI/pdf/news/coar_unesco_oa_statement.pdf
- [L5] <http://thecostofknowledge.com/>
- [L6] <http://www.bib.umontreal.ca/communiqués/20160506-DC-annulation-springer-va.htm>
- [L7] <http://discreteanalysisjournal.com/>,
<http://jep.cedram.org/spip.php?article33&lang=en>,
<http://epiga.episciences.org/>
- [L8] <https://www.episciences.org/>
- [L9] <http://www.scielo.org/>
- [L10] <https://www.openedition.org/?lang=en>
- [L11] <https://www.openlibhums.org/>
- [L12] <https://twitter.com/audreywatters/status/184387170415558656>

Reference:

- [1] T. Pisanski: "Open Access – Who Pays?", Newsletter of the European Mathematical Society, June 2013, p. 54, <http://www.ems-ph.org/journals/newsletter/pdf/2013-06-88.pdf>

Please contact:

Karim Ramdani, Inria, France
karim.ramdani@inria.fr

Open-Access Repositories and the Open Science Challenge

by Leonardo Candela, Paolo Manghi, and Donatella Castelli (ISTI-CNR)

The open-access movement is promoting free-of-restriction access to, and use of, research outcomes. It is a key aspect of the open-science movement, which is pushing for the research community to go 'beyond papers'. This new paradigm calls for a new generation of repositories that are: (i) capable of smartly interfacing with the wealth of research infrastructure and services that scientists rely on, thus being able to intercept and publish research products, (ii) able to provide researchers with social networking tools for discovery, notification, sharing, discussion, and assessment of research products.

The landscape of scientific research has changed dramatically in the last few years. The forces driving the change include both new technology (namely ICT infrastructures and services) and the open-science movement that is supporting and encouraging an open-access-driven dissemination and exploitation of virtually every research product worth sharing; not only papers but datasets, software, notebooks and every computational object produced in the course of research.

However, the evolution is still underway. ICT infrastructures are quite diffuse among research communities and researchers, and the large majority of daily scientific activities relies on them, yet a gap remains between the 'places' where research is conducted and the 'places' where its dissemination and communication happen. This gap, which originates from the long tradition of paper-driven scientific communication that still characterises science, is one of the major barriers to overcome before open science becomes a reality. The traditional means of scientific communication are so ingrained that, when called upon to manage a new type of scientific product, i.e., the 'research data', the scientific community responded by proposing existing approaches such as specific journals, i.e., data journals [2], and/or repositories, i.e., data repositories [3]. Such approaches do not fit well with the entire spectrum of research products envisaged, for which effective interpretation, evaluation, and reuse can only be ensured if publishing includes the properties of 'within' the environment (and context) from which they originate and 'during' the research activity.

Motivated by these observations we envisioned a completely new kind of open access / science repository, SciRepo [1]. This is a sort of 'overlay repository' that is expected to sit on top of the research environment / infrastructure that researchers use to dynamically collect research artefacts (a) as soon as they are produced, (b) without needing to spend effort to repurpose them for publication purposes, and (c) fully equipped with their 'context', i.e., the wealth of information surrounding the artefact and key for its under-

standing. SciRepo's distinguishing features include: (a) hooks interfacing with ICT services to intercept the generation of products and to publish such products, i.e., to make them discoverable and accessible to other researchers; (b) provision of repository-like tools so that scientists can access and share research products generated during their research activities; (c) social networking based practices to modernise (scientific) communication both intra-community and inter-community, e.g., posting rather than deposition, 'like' and 'open discussions' for quality assessment, sharing rather than dissemination.

SciRepo repository-oriented facilities are largely based on the rich information graph characterising every published product. They include search and browse allowing search by product typology, but also permitting navigation from research activities to products and related products. Ingestion facilities are provided, allowing scientists to manually or semi-automatically upload 'external' products into the repository and associate them with a research activity, thus including them in the information graph. Ingestion allows scientists to complete the action of publishing a research activity with all products that are connected to it but generated out of the boundaries of the community. The way scientists or groups of scientists can interact with products (access and reuse them) is ruled by clear rights management functionalities. Rights are typically assigned when products are generated or ingested by scientists, but can vary over time.

SciRepo collaboration-oriented facilities include typical social networking facilities such as the option to subscribe to events that are relevant to research activities and products, and be promptly notified, e.g., the completion of a workflow execution, the generation of datasets that conform to a particular criteria. Users can reply to posts and, most importantly, can express opinions on the quality of products, e.g., 'like' actions or similar. SciRepo thus represents a step towards truly 'open' peer-review. More sophisticated assessment/peer-review functionalities (single/double blind) can be supported, in order to provide more traditional notions of quality. Interestingly, the posts themselves represent a special type of product of the research activity and are searchable and browsable in the information graph.

References:

- [1] M. Assante et al.: "Science 2.0 Repositories: Time for a Change in Scholarly Communication", *D-Lib Magazine*. 21 (1/2), (2015), doi: 10.1045/january2015-assante
- [2] L. Candela et al.: "Data Journals: A Survey", *Journal of the Association for Information Science and Technology*. 66 (1): 1747–1762, 2015), doi:10.1002/asi.23358
- [3] M. Assante et al.: "Are Scientific Data Repositories Coping with Research Data Publishing?", *Data Science Journal*. 15, 2016, doi:10.5334/dsj-2016-006/

Please contact:

Leonardo Candela, ISTI-CNR, Italy
leonardo.candela@isti.cnr.it

LIPIcs – an Open-Access Series for International Conference Proceedings

by Marc Herbstritt (Schloss Dagstuhl – Leibniz-Zentrum für Informatik) and Wolfgang Thomas (RWTH Aachen University)

The commercialisation of scientific publishing has resulted in a situation where more and more relevant literature is separated from the scientists by high pay walls; this has created an unacceptable impediment to scientific exchange. To illustrate how scientists can regain the essence of 'publishing' – namely to make research results public – we report on LIPIcs (Leibniz International Proceedings in Informatics), an open-access series for the proceedings of international conferences.

Background

With the advent of digital technologies, many tasks involved in scientific publishing have been facilitated enormously. This applies to scientific writing (using systems such as LaTeX) as well as the world-wide dissemination of literature via the internet. Somewhat paradoxically, at the same time the prices for accessing scientific literature have exploded, a development that was and is driven by commercial publishers and which imposes severe obstacles to scientific progress. It is not clear whether and how the world of science will be able to launch a "reconquista" of scientific publishing, taking it out of the hedgefunds and stock markets and making it more science-driven again.

We report here on an initiative, started ten years ago, that has the potential to be a successful chapter of this reconquista.

The Foundation of LIPIcs

Since the 1970s, a standard venue for proceedings of conferences in computer science was the series Lecture Notes in Computer Science (LNCS) published by Springer-Verlag. When the first editorial board of LNCS resigned in 2004, the number of published volumes drastically increased (to about two volumes a day) by inclusion of many workshop proceedings. At the same time, the price of the series increased significantly, resulting in many research institutions cancelling their subscriptions. LNCS was effectively alienating its readers and contributors.

Responding to this development, the steering committee of the renowned Symposium on Theoretical Aspects of Computer Science (STACS), together with the Asian conference Foundations of Software Technology and Theoretical Computer Science (FSTTCS), made the bold decision in 2007 to leave Springer-Verlag after more than 20 years. They elected instead to go open access with solely digital online proceedings. A strong and devoted partner was found in Reinhard Wilhelm, then scientific director of the Germany-based Leibniz Center of Informatics – Schloss Dagstuhl, which is well known in the community for hosting its 'Dagstuhl Seminars'. Together the open-access series Leibniz International Proceedings in Informatics (LIPIcs) [L1] was

founded in 2008. LIPIcs embodies two core principles (discussed further below): (i) gold open access while insisting on high scientific standards, and (ii) providing affordable, meticulously edited proceedings.

Editorial Board and Editorial Policy

The editorial board currently consists of nine members whose terms are limited to two periods of at most six years each. The task of the board is to ensure that conferences of high scientific standards are accepted for LIPIcs. The board must determine for instance: (i) whether there is evidence that a conference has a high reputation, (ii) whether there is a steering committee whose members are renowned scientists and change on regular terms, and (iii) whether the conference adequately represents its respective field.

Strict rules determine whether or not an application is successful: a secret vote is held which needs six positive votes (out of nine) for acceptance. Accepted conferences need to re-apply every five years. This policy has led to rejections of several conferences that could safely be considered solid. Such a rigorous process was essential, however, for LIPIcs to earn an excellent scientific reputation within a short time. The appeal and success of LIPIcs is evident, with 25 conferences having now been accepted. To date, for 2016, this amounts to about 1,000 conference papers which are published open-access.

Production of the Proceedings and Financial Matters

Clearly, considerable effort is needed to ensure high editorial quality beyond the scientific value of a paper. This involves more than just adopting some LaTeX style (which some authors tend to violate). It also means, for example, that the validity of citations must be checked. This tedious work is handled by the team of the LIPIcs editorial office, who managed, despite rather sparse resources, to deliver high-quality proceedings [L2] on a par with LNCS and other conference proceedings series.

LIPIcs has been charging an article-processing charge (APC) since 2010. Initially the APC was kept at a very low €15. In 2015, the funding agency of Schloss Dagstuhl, the German Federal Ministry of Education and Research, stipulated that general funds of Schloss Dagstuhl were no longer to be used to support the publishing activities of LIPIcs. Thus the APC had to be increased to €60 to cover the costs. This still compares favourably to the charges of commercial publishers for gold open access, which range from six to 12 times this amount. The APC will be increased incrementally, in three stages between now and 2019, using a generous donation that Schloss Dagstuhl – now under the scientific directorship of Raimund Seidel – received from the Heidelberg Institute for Theoretical Studies (HITS).

Perspectives

The open-access movement has gained a considerable boost in recent years. A complete switch to open-access publications now seems possible, and research organisations worldwide are working towards this goal (see, for example, the report by Schimmer et al. at <http://dx.doi.org/10.17617/1.3>).

In the area of computing research, LIPIcs is at the forefront of making relevant research results openly accessible. This is

underpinned by Schloss Dagstuhl's recently established Dagstuhl Artifacts Series (DARTS) [L3] which aims for persistent publication of research data and artifacts. DARTS was triggered by the needs of LIPIcs conferences and shows how science-driven publishing infrastructure can evolve.

There are also other not-for-profit open-access publishing services for proceedings that share similar goals as LIPIcs, for example, EPTCS [L4] and CEUR-WS [L5]. Not-for-profit publishing services of this kind rely on cooperative authors and editors to make gold open access for computer science conferences happen. The reconquista of scientific publication into the hands of science will only be successful if these services are not seen as a simple replacement for for-profit publishers but as collaborative academia-driven not-for-profit initiatives.

Links:

[L1] <http://www.dagstuhl.de/lipics>

[L2] <http://drops.dagstuhl.de/lipics>

[L3] <http://www.dagstuhl.de/darts>

[L4] <http://www.eptcs.org>

[L5] <http://ceur-ws.org>

Please contact:

Marc Herbstritt (head of LIPIcs editorial office)
Schloss Dagstuhl – Leibniz-Zentrum für Informatik,
Germany
+49 681 302 3849, marc.herbstritt@dagstuhl.de

Wolfgang Thomas (chair of editorial board of LIPIcs)
RWTH Aachen University, Germany
+49 241 8021701, thomas@cs.rwth-aachen.de

Scientific Data and Preservation – Policy Issues for the Long-term Record

by Vera Sarkol (CWI)

In order to keep open data accessible into the future, academics and librarians need to consider long-term preservation.

From open access of publications the trend is now expanding to open science, and, with that, open data. The progress of our communal knowledge is dependent on previously discovered truths, and therefore the data has to be openly available to the extent that others can find, understand and use it [1]. The concepts of ‘openness’ and ‘preservation’ are inextricably linked if we want to secure a continuous record of the path of discovery. The job of maintaining these records falls to the national or institutional libraries and repositories.

Many funders, such as the Netherlands Organisation for Scientific Research (NWO), are developing policies for data and software management which address openness and preservation. This puts some pressure on the issue, and it is the right place to raise the question of cost for documenting and depositing the artifacts, in terms of workload and resources. The most important challenges for long-term policy are selection, findability, and reusability.

Selecting what to preserve

Ideally we would preserve and make available every scientific artifact, but in reality this is neither feasible nor desirable [L1]. Constraints of size or legality will of course hinder preservation. Other constraints are the time it costs to properly document and describe datasets and software, and the environmental cost of storage. Therefore data that can easily be replicated or code that only serves to illustrate an algorithm does not necessarily need to be preserved. For now it is a good principle to preserve artifacts that underlie publications, but if in the future the boundaries of publications as the unit of scientific knowledge blur (e.g., if preprints and post-evaluation get integrated into the process), academics and librarians together will have to develop other criteria for selection.

Replication packages

Storing only data or software is no guarantee that a finding can be replicated if crucial information is missing. To avoid this problem NWO will soon make replication packages mandatory. This means that along with the dataset or program, the metadata, identifier and provenance information should be stored, as well as the software and hardware, or at the very least a description. However, even with that information, complex dependencies or outdated software packages may still prevent replication.

One project that provides a solution to this problem is being developed at CWI: Snakemake [2]. This is a text-based

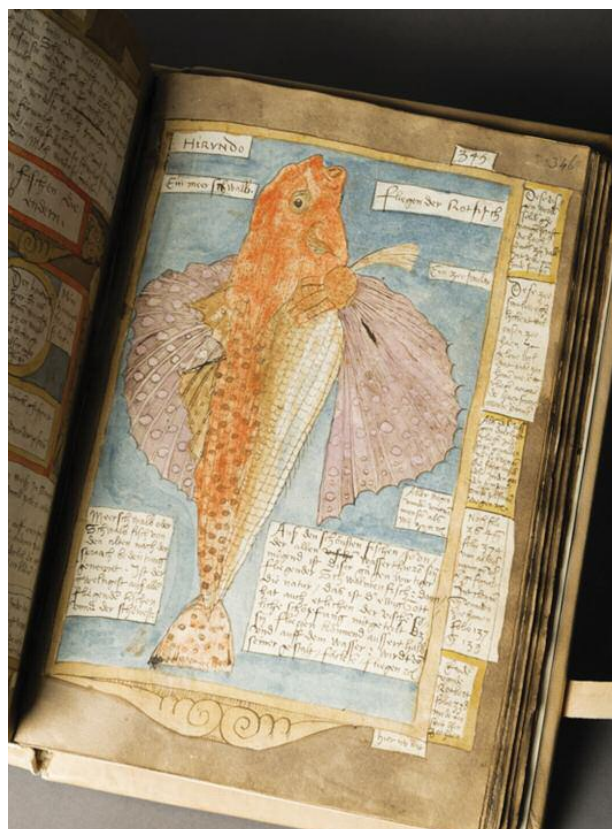


Figure 1: Preserved knowledge of fish, from Adriaen Coenensz' 'Visboeck', 1579, the National Library of the Netherlands. Location: KB Den Haag, KW 78 E 54 fol. 346r

workflow management system that was originally developed for bioinformatics but could be suitable for other fields of research as well. Using a domain specific language, Snakemake aims to formalise an analysis workflow, including a specification of used software packages. Upon execution of a workflow, software packages are deployed automatically so that an analysis is reproducible without extra work.

More drastic problems will occur when hardware becomes outdated. One possible way forward could be virtualisation [L2], where the old environment is emulated to access preserved scientific artifacts. However, at some point hardware may undergo such a large change that this too is no longer a viable option. It is necessary for the community to start thinking about what to do when that occurs.

Licences

A large variety of data and software licences are currently in use, sometimes prescribed by journals or repositories. When interoperability becomes more prominent these licences may not interact well with each other and this may lead to datasets being unable to recombine. Another problem is that not everyone has licences to proprietary (legacy) software and operating systems used. One solution is to fully commit to open-source. Another possibility is to licence everything to the public domain and that licenced legacy software is kept running centrally, for instance by national heritage institutions [L2] (the National Library in the Netherlands, for instance).

Findability of preserved software

For long term findability, citing an URL for a program in an article is not sufficient, since the content of an URL can easily change. More durable would be to give all scientific artifacts, including software, a persistent identifier which is a unique code given to an object by an organisation, irrespective of its location. The DOI has become the academic standard and thus may be expected to be maintained the longest. The version of a program that underlies a publication should be deposited in a repository and receive a DOI. For instance, Zenodo provides this service and is integrated with GitHub. Getting a DOI will make it easier to find the right version of the software with the publication, but it will also make it easier to find and cite the work for the broader (academic) community and funding agencies [3].

Conclusion

From a library's perspective the goal is to make the record of scientific knowledge as permanent as it was when there was only paper. While progress is being made, international consensus on a number of issues needs to be achieved. Consultation at a European level is necessary to establish guidelines for the long-term preservation of open data and software.

Links:

[L1] https://www.esciencecenter.nl/pdf/Software_Sustainability_DANS_NLeSC_2016.pdf

[L2] https://www.unesco.nl/sites/default/files/dossier/report_girona_session_persist.pdf

References:

- [1] M. D. Wilkinson et al.: "The FAIR Guiding Principles for scientific data management and stewardship", *Scientific Data* 3:160018, 2016.
<http://dx.doi.org/10.1038/sdata.2016.18>
- [2] J. Köster, S. Rahmann: "Snakemake – A scalable bioinformatics workflow engine", *Bioinformatics* 28(19): 2520-2522, 2012.
<http://dx.doi.org/10.1093/bioinformatics/bts480>
- [3] A.M. Smith et al.: "Software Citation Principles", *PeerJ Preprints*, 2016.
<http://dx.doi.org/10.7287/peerj.preprints.2169v2>

Please contact:

Vera Sarkol
CWI Information & Documentation
+31(0)205924051
vera.sarkol@cwi.nl

Mathematics in Open Access – MathOA

by Johan Rooryck and Saskia de Vries

The new project MathOA is a response to the EU Council call for a transition to open access by 2020. MathOA provides a large-scale passage to open access for mathematics research that addresses current market dysfunction by a uniquely sustainable and affordable transition based on Fair OA price pressure. Mathematics in Open Access builds further on the proven bottom up approach of Linguistics in Open Access (LingOA) that is discipline-based and editor-based. This Fair OA approach makes sure that no author pays individual article processing charges.

Background: LingOA and Fair Open Access

Open-access publishing is often said to be the future of academic journals, but the actual move from a subscription model to an open-access model is not easily achieved. Frequently, it only raises the total cost of access for libraries. In the meantime, researchers and libraries remain hostages of big publishers such as Elsevier, Wiley, Taylor & Francis, or Springer. These publishers make profits in excess of 35% or more on the public money most libraries use to pay for access to published research. Articles behind paywalls remain inaccessible not only for the taxpayers who paid for the research published in those articles, but also for scholars around the world who cannot afford expensive subscriptions.

Recently, the EU Competitiveness Council's Conclusion on Open Science [L1] stated that all scientific publications deriving from Horizon 2020 or other EC funding will have to be freely available by 2020. Carlos Moedas, the European Commissioner for Research, Innovation and Science, has called the move 'life-changing'. One of the routes towards this goal has started in the Netherlands with the 'OA big-deals', in which prices are recorded in licences. Another route lies in what's known as Fair Open Access, an alternative, researcher-driven path towards the same goal.

In 2015, the foundation Linguistics in Open Access (LingOA) [L2] was set up as a pilot project in the humanities to flip existing linguistics journals with an excellent reputation from subscription to open access. After this successful pilot, we were approached by a group of mathematicians who wished to replicate the LingOA incubation model. Furthermore, the board of the Conference of European Schools for Advanced Engineering, Education and Research (CESAER) has asked us to submit a proposal for their 50 universities to make this possible. We therefore are developing MathOA as a pilot project in the domain of hard sciences, thus providing an example for other hard science disciplines to follow suit. As a result, LingOA and MathOA will function as the two pioneering Fair Open-Access projects in their respective scientific domains.

Conditions of Fair Open Access

Under the LingOA Fair Open-Access model, reputed linguistics journals can join LingOA if their publisher agrees to comply with the following conditions of Fair Open Access:

- The editorial board or a learned society owns the title of the journals.
- Authors own the copyright of their articles, and a CC-BY license applies.
- All articles are published in a fully open-access mode (no subscriptions, no ‘hybrid model of both subscriptions and APCs a.k.a. ‘double dipping’).
- Article processing charges (APCs) are low, transparent, and in proportion to the cost of the work carried out by the publisher.
- Authors do not individually pay for APCs.

The journals *Laboratory Phonology* (De Gruyter), *Journal of Portuguese Linguistics* (University of Lisboa) and *Lingua* (Elsevier, now called *Glossa*) were the first ones to flip to a publisher who complies with these conditions of Fair Open Access. By the end of August 2016 it was already clear that the LingOA pilot was a success.

Fair Open Access: who pays for the APCs?

The Association of Dutch Universities (VSNU) and the Dutch Organization for Scientific Research (NWO) have provided LingOA with a five-year grant of 0.5 million euros to pay for the APCs of linguistics journals that move to Fair Open Access, as well as for legal and other advice and project management. As a result, authors submitting articles to journals that are members of LingOA do not pay for any APCs themselves. After the initial five years, the APCs of participating journals will be taken over by the Open Library of Humanities (OLH) [L4]. OLH is a charitable organisation dedicated to publishing open-access scholarship with no author-facing APCs. OLH is funded by an international consortium of 190+ prestigious libraries who make a contribution that covers the APCs of participating journals. Once again, this means that no linguist ever pays for APCs when they publish an article in a journal participating in LingOA. In this way, long term sustainable Fair Open Access is achieved for all participating linguistics journals. With MathOA, the OLH will be extended to an Open Library of Sciences.

MathOA partners

- MathOA will be founded and hosted by two prestigious mathematics societies in the Netherlands: (1) Centrum Wiskunde & Informatica (CWI). (2) The Royal Netherlands Mathematical Society (KWG).
- CESAER, the Conference of European Schools for Advanced Engineering Education and Research, is a non-profit international association of leading European universities of science and technology, technology and engineering schools/faculties at comprehensive universities and university colleges.
- Last but not least, a group of scientists around Sir Timothy Gowers (a Fields-medal winning mathematician) have

joined forces to convince editorial boards of important journals to flip to Fair Open Access.

Political momentum

MathOA is not only about flipping prestigious subscription journals to Fair Open Access, it is also about raising pressure on the commercial publishers to start providing their services on fair and transparent conditions. If CESAER decides to sponsor MathOA, the pilot that started with LingOA would be given enormous momentum, resulting in a tidal change in all sciences. We are confident that this would eventually lead to a reduction of the total cost of scientific communication, in line with the path Ralf Shimmer describes in the Max Planck white paper ‘Disrupting the subscription journals’ business model for the necessary large-scale transformation to open access’ [L3].

Links:

- [L1] <http://data.consilium.europa.eu/doc/document/ST-9526-2016-INIT/en/pdf>
- [L2] <http://www.lingOA.eu>
- [L3] <http://www.openlibhums.org>
- [L4] http://pubman.mpg.de/pubman/item/escidoc:2-148961:7/component/escidoc:2149096/MPDL_OA-Transition_White_Paper.pdf

Please contact:

Saskia de Vries
Sampan – academia & publishing, The Netherlands
s.c.j.devries@sampan.eu

Introduction to the Special Theme

Modern Machine Learning: More with Less, Cheaper and Better

by Sander Bohte and Hung Son Nguyen

While the discipline of machine learning is often conflated with the general field of AI, machine learning specifically is concerned with the question of how to program computers to automatically recognise complex patterns and make intelligent decisions based on data. This includes such diverse approaches as probability theory, logic, combinatorial optimisation, search, statistics, reinforcement learning and control theory. In this day and age with an abundance of sensors and computers, applications are ubiquitous, ranging from vision to language processing, forecasting, pattern recognition, games, data mining, expert systems and robotics.

Historically, rule-based programs like the Arthur Samuel checkers-playing program were developed alongside efforts to understand the computational principles underlying human learning, in the developing field of neural networks. In the '90s, statistical AI emerged as a third approach to machine learning, formulating machine learning problems in terms of probability measures. Since then, the emphasis has vacillated between statistical and probabilistic learning and progressively more competitive neural network approaches.

The breakthrough work by Krizhevsky, Sutskever & Hinton [1] on deep neural networks in 2012 has been a catalyst for AI research by demonstrating a step function in performance on the Imagenet computer vision competition. For this, they used a deep neural network trained exhaustively on 'GPUs': a garden-variety parallel computing hardware used for video-games. Similar advances were then quickly reported for speech recognition and later for machine translation and natural language processing. In short order, big companies like Google, Microsoft and Baidu established large machine learning groups, quickly followed by essentially all other big tech companies. Since then, with the combination of big data and big computers, rapid advances have been reported, including the use of machine learning for self-driving cars, and consumer-grade real-time speech-to-speech translation. Human performance has even been exceeded in some specialised domains. It is probably safe to say that at present, machine learning allows for many more applications than there are engineers capable of implementing them.

These rapid advances have also reached the general public, with often alarming implications: think tanks are declaring that up to 70% of all presently existing jobs will disappear in the near future, and serious attention is being given to potentially apocalyptic futures where AI capabilities exceed human intelligence. We believe, however, that it is safe to say that this will not happen in the next five years, as machine learning still faces some serious obstacles before reaching human levels of flexible intelligence.

Some of the current challenges in machine learning are reflected in the articles presented in this special issue: the much glorified deep learning approaches all rely on the availability of massive amounts of data, often needing millions of correctly

labelled examples. Many domains, however, including some important areas such as health care, will never have such massive labelled datasets. Similarly, robots cannot be trained for millions of trials, simply because they wear out long before. The question is thus how to learn more with less. Here, statistics and prior knowledge will likely play a big role, and some promising work is presented in this issue – see for example the articles by Mouret and by Welling. Some work is also examining whether quantum computing can help reduce the computational complexity of machine learning, as explained in the article by Wittek. At the same time, massive data streams generate problems of their own: Cieliebak and Benczur each present work on how to deal with huge torrents of streaming data.

Apart from these technical challenges, we, as a community, need to train the future experts in machine learning, such that the wider industry and society can benefit from what is currently already possible. Some of the exciting applications are laid out in this issue: Kappen for example showcases the Bonaparte Disaster Victim Identification system, which uses Bayesian statistical modelling to identify victims based on their DNA and that of next of kin. Potamias presents a wonderful application of machine learning in ‘Pharmagenomics’, where the challenge is to determine which genes interact with which drugs. This is a key determining factor in the efficacy of these drugs and central to the future of personalised medicine.

A separate line of ongoing research is the link between the one working example of intelligence, the brain, and learning principles. This relates to such diverse questions as ‘How can goals be selected in an autonomous fashion?’ and ‘How can we optimise over many different learning problems with one system?’, but also in reverse: ‘What can the success of deep neural networks tell us about the brain?’. Articles by Alexandre and Oudeyer cover current efforts on these topics.

The study of machine learning has thus grown from the efforts of a handful of computer engineers exploring whether computers could learn to play games and mimic the human brain, and a field of statistics that largely ignored computational considerations, to a booming discipline that is actively transforming the world in which we live.

Reference:

[1] Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E. Hinton. "Imagenet classification with deep convolutional neural networks." *Advances in neural information processing systems (NIPS-2012)*. Pages 1097-1105, 2012.

Please contact:

Sander Bohte, CWI, The Netherlands
sander.bohte@cwi.nl

Hung Son Nguyen, University of Warsaw, Poland
son@mimuw.edu.pl

Micro-Data Learning: The Other End of the Spectrum

by Jean-Baptiste Mouret (Inria)

Many fields are now snowed under with an avalanche of data, which raises considerable challenges for computer scientists. Meanwhile, robotics (among other fields) can often only use a few dozen data points because acquiring them involves a process that is expensive or time-consuming. How can an algorithm learn with only a few data points?

Watching a child learn reveals how well humans can learn: a child may only need a few examples of a concept to “learn it”. By contrast, the impressive results achieved with modern machine learning (in particular, by deep learning) are made possible largely by the use of huge datasets. For instance, the ImageNet database used in image recognition contains about 1.2 million labelled examples; DeepMinds's AlphaGo used more than 38 million positions to train their algorithm to play Go; and the same company used more than 38 days of play to train a neural network to play Atari 2600 games, such as Space Invaders or Breakout.

Like children, robots have to face the real world, in which trying something might take seconds, hours, or days. And seeing the consequence of this trial might take much more. When robots share our world, they are expected to learn like humans or animals, that is, in far fewer than a million trials. Robots are not alone to be cursed by the price of data: Any learning process that involves physical tests or precise simulations (e.g., computational fluid dynamics) comes up against the same issue. In short, while data might be abundant in the virtual world, it is often a scarce resource in the physical world. I refer to this challenge as “micro-data” learning (see Figure 1).

The first precept of micro-data learning is to choose as wisely as possible what to test next (active learning). Since computation tends to become cheaper every year, it is often effective to trade data resources for computational resources, that is, to employ computationally intensive algorithms to select the next data point to acquire. Bayesian optimisation [1] is such a data-efficient algorithm that has recently attracted a lot of interest in the machine learning community. Using the data acquired so far, this algorithm creates a probabilistic model of the function that needs to be optimised (e.g., the

walking speed of a robot or the lift generated by an airfoil); it then exploits this model to identify the most promising points of the search space. It can, for example, find good values for the gait of a quadruped robot (Sony Aibo / 15 parameters to learn) in just two hours of learning.

The second precept of micro-data learning is to exploit every bit of information from each test. For instance, when a robotic arm tries to reach a point

approach for learning control strategies in robotics; for example, the Pilco algorithm can learn to balance a non-actuated pole on an actuated moving cart in 15-20 seconds (about 3-5 trials) [2].

The third precept of micro-data learning is to use the “right” prior knowledge. Most problems are indeed simply too hard to be learned from scratch in a few trials, even with the best algorithms: The quick learning ability of humans and animals is due largely to their prior

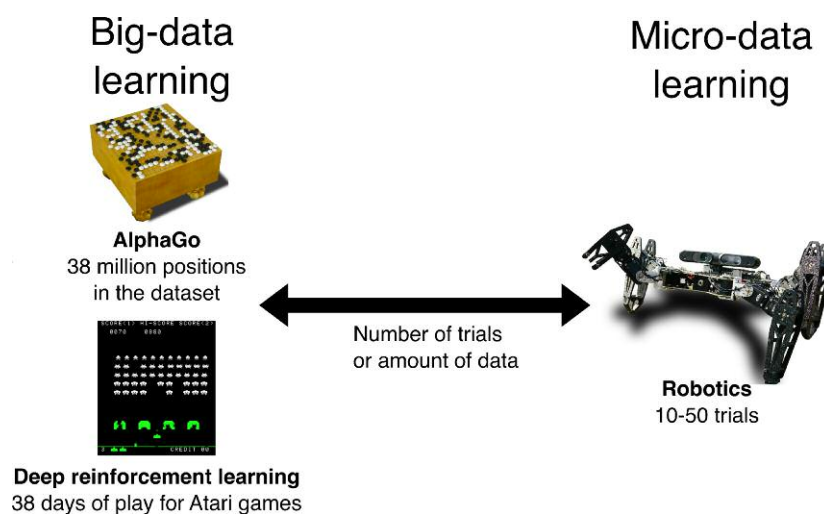


Figure 1: Modern machine learning (e.g., deep learning) is designed to work with a large amount of data. For example, the Go player AlphaGo by DeepMind used a dataset of 38 million positions, and the deep reinforcement learning experiments from the same team used the equivalent of 38 days to learn to play Atari 2600 video games. Robotics is at the opposite end of the spectrum: most of the time, it is difficult to perform more than a few dozen trials. Learning with such a small amount of data is what we term “Micro-data learning”.

in space, the learning algorithm can perform the movement, then, at the end of the trial, measure the distance to the target. In this case, each test corresponds to a single data point. However, the algorithm can also record the position of the “hand” every 10ms, thus getting thousands of data points from a single test. This is a very effective

knowledge about what could and could not work. When using priors, it is critical to make them as explicit as possible, and to make sure that the learning algorithm can question or even ignore them. In academic examples, it can also be challenging to distinguish between prior knowledge that is useful and prior knowledge that actually gives the solu-

tion to the algorithm, which leaves nothing to learn.

We focused on prior knowledge in our recent article about damage recovery in robotics [3, L1]. In this scenario, a six-legged walking robot needs to discover a new way to walk by trial-and-error because it is damaged. Before the mission, a novel algorithm explores a large search space with a simulation of the intact robot to identify the most promising solution of each "family". Metaphorically, this algorithm takes the needles out of a haystack to make a stack of needles. If the robot is damaged, the learning algorithm, which is a derivative of Bayesian optimisation [1], exploits this prior knowledge to choose the best trials. In our experiments, the

robot discovers compensatory gaits in less than two minutes and a dozen trials, for the five damage conditions that we tested [3].

In this learning approach, a data-efficient learning algorithm that works with the physical, damaged robot is guided by prior knowledge based on a simulation of the intact robot. This micro-data learning algorithm makes it possible to learn a complex task in only a few trials. The subsequent challenge is to exploit more knowledge from the trials [2] and select the next trials while taking the context into account (e.g., potential obstacles).

Link:
[L1] <http://www.resibots.eu>

References:

- [1] B. Shahriari, et al.: "Taking the human out of the loop: A review of bayesian optimization", Proc. of the IEEE, 2016.
- [2] M. P. Deisenroth, D. Fox, C. E. Rasmussen: "Gaussian processes for data-efficient learning in robotics and control", IEEE Trans. on Pattern Analysis and Machine Intelligence, 2016.
- [3] A. Cully, et al.: "Robots that can adapt like animals", Nature, 2015.

Please contact:

Jean-Baptiste Mouret
Inria, France
jean-baptiste.mouret@inria.fr

Making Learning Physical: Machine Intelligence and Quantum Resources

by Peter Wittek (ICFO-The Institute of Photonic Sciences and University of Borås)

It is not only machine learning that is advancing rapidly: quantum information processing has witnessed several breakthroughs in recent years. In theory, quantum protocols can offer an exponential speedup for certain learning algorithms, but even contemporary implementations show remarkable results – this new field is called quantum machine learning. The benefits work both ways: classical machine learning finds more and more applicability in problems in quantum computing.

After a history spanning over five decades, artificial general intelligence still remains out of reach. Machine learning has common roots with AI research, but focuses on more attainable goals and has achieved tremendous success in many application fields. Similarly, a universal quantum computer is still far ahead in the distant future: the criterion for this machine is to be able to simulate an arbitrary closed quantum system. Nevertheless, uses of quantum information processing are proliferating: two notable

examples are quantum key distribution systems and quantum random number generators.

Recently, there has been a surge of interest in the intersection of machine learning and quantum information processing. Combining ideas from these two fields leads to tremendous benefits for both. We are collaborating on several subjects in this domain between ICFO-The Institute of Photonic Sciences, the Autonomous University of Barcelona, the University of the

Basque Country, all in Spain, as well as the University of Calgary, Canada.

At the highest level, abstracting of the actual algorithms and focusing on the foundations of statistical learning theory, we can ask what it means to learn with quantum data and channels, what induction and transduction mean in this setting, how we can define figures of merit to quantify performance, and eventually establish bounds on generalisation performance using sample and model complexity. We studied supervised learning, and proved that in the asymptotic limit and under an assumption of exchangeability, quantum entanglement does not break our traditional notion of induction [L1]. This is an important stepping stone towards understanding generalisation properties of quantum learning protocols.

The next natural question to ask is that given a universal quantum computer, what kind of protocols can we use for

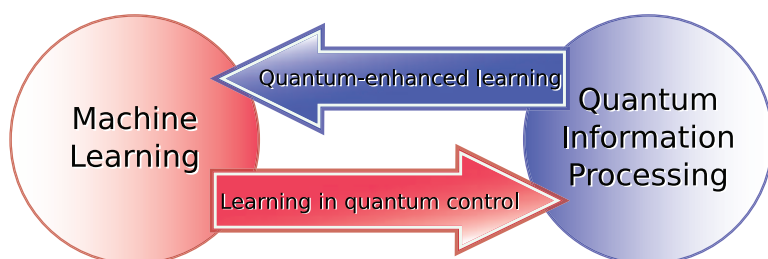


Figure 1: Overview of the interplay between quantum information processing and machine learning.

learning? Many proposals have been put forward, offering up to an exponential reduction in computational complexity [1], but it is worth looking at what is attainable with current technology. Quantum adiabatic optimisation leads the way in this regard; it uses quantum annealing – a physical process resembling the widely used global optimisation heuristic simulated annealing that uses both thermal fluctuations and quantum tunnelling to find the global energy minimum of a system. Scalable quantum annealers already exist, but they are not guaranteed to find a global optimum. On the other hand, the local optima retrieved from such machines closely follows a Gibbs distribution. For this reason, they have been used to train various configurations of Boltzmann machines [2]. Gibbs sampling, however, does not only occur in Boltzmann machines: probabilistic inference in Bayesian networks and Markov random fields makes extensive use of it. Since this a #P problem in general, currently we are looking at how

one can match these inference methods with existing implementations.

Using quantum resources in learning is only one side of the coin: we can also employ machine learning using classical computers in common problems that arise in quantum information processing, for instance, in quantum control. Quantum control steers quantum dynamics towards realising specific quantum states or operations, and thus it is an important component in constructing a universal quantum computer. We have been working on a reinforcement learning algorithm to control an adaptive quantum metrology scheme [3], improving noise tolerance and scalability, providing a tool that experimentalists can use [L2, L3].

Our overall vision is that this cross-fertilisation between machine learning and quantum information processing will continue, and that the future of artificial general intelligence and universal quantum computing are intertwined.

Links:

- [L1] <http://arxiv.org/abs/1605.07541>
- [L2] <http://arxiv.org/abs/1607.03428>
- [L3] https://panpalitta.github.io/phase_estimation/

References:

- [1] P. Rebentrost, M. Mohseni, S. Lloyd: “Quantum support vector machine for big feature and big data classification”, *Physical Review Letters*, 2014, 113, 130503.
- [2] S. H. Adachi, M. P. Henderson: “Application of Quantum Annealing to Training of Deep Neural Networks”, *arXiv:1510.06356*, 2015.
- [3] A. Hentschel, B. C. Sanders: “Machine Learning for Precise Quantum Measurement”, *Physical Review Letters*, 2010, 104, 063603.

Please contact:

Peter Wittek, ICFO-The Institute of Photonic Sciences, Barcelona, Spain
University of Borås, Borås, Sweden
+34935542237
<http://peterwittek.com/>

Marrying Graphical Models with Deep Learning

by Max Welling (University of Amsterdam)

In our research at the University of Amsterdam we have married two types of models into a single comprehensive framework which we have called “Variational Auto Encoders”. The two types of models are: 1) generative models where the data generation process is modelled, and 2) discriminative models, such as deep learning, where measurements are directly mapped to class labels.

Deep learning is particularly successful in learning powerful (e.g., predictive/discriminative) features from raw, unstructured sensor data. Deep neural networks can effectively turn raw data streams into new representations that represent abstract, disentangled and semantically meaningful concepts. Based on these, a simple linear classifier can achieve the state of the art. But to learn them one needs very large quantities of annotated data. They are flexible input-output mappings but do not incorporate a very sophisticated inductive bias about the world. An important question is how far will this take us?

If we are asked to analyse a scene depicted in an image we seek a story that can explain the things we see in the image. Yes, there is a fast feedforward pipeline that quickly segments out the objects and classifies them into object

classes. But when you need to truly understand a scene you will try to infer a story about which events caused other events, which in turn led to the image you are looking at. This causal story is also a powerful tool to predict how the events may unfold into the future.

So, to understand and reason about the world we need to find its causal atoms and their relationships. Now this is precisely what Bayesian networks [1] were intended to do. Each random variable connects to other random variables and their directed relations model their causal relationships. (Bayesian networks do not necessarily represent the causal relationships, but an extension called “structural equation models” does.) Another key advantage of interpretable models like Bayesian networks is that they can express our expert knowledge. If we know X causes Y then

we can simply hard-code that relation into the model. Relations that we do not know will need to be learned from the data. Incorporating expert knowledge (e.g., the laws of physics) into models is the everyday business of scientists. They build sophisticated simulators with relatively few unidentified parameters, for instance implemented as a collection of partial differential equations (PDEs).

Generative models can also be used for classification by inverting the relationship they model from class label to input features. When you have a lot of (labelled) data at your disposal this type of classifier will generally speaking not work as well as a direct mapping from input features to labels (such as a deep neural network). But when the amount of data is small relative to the complexity of the task, the opportunity to

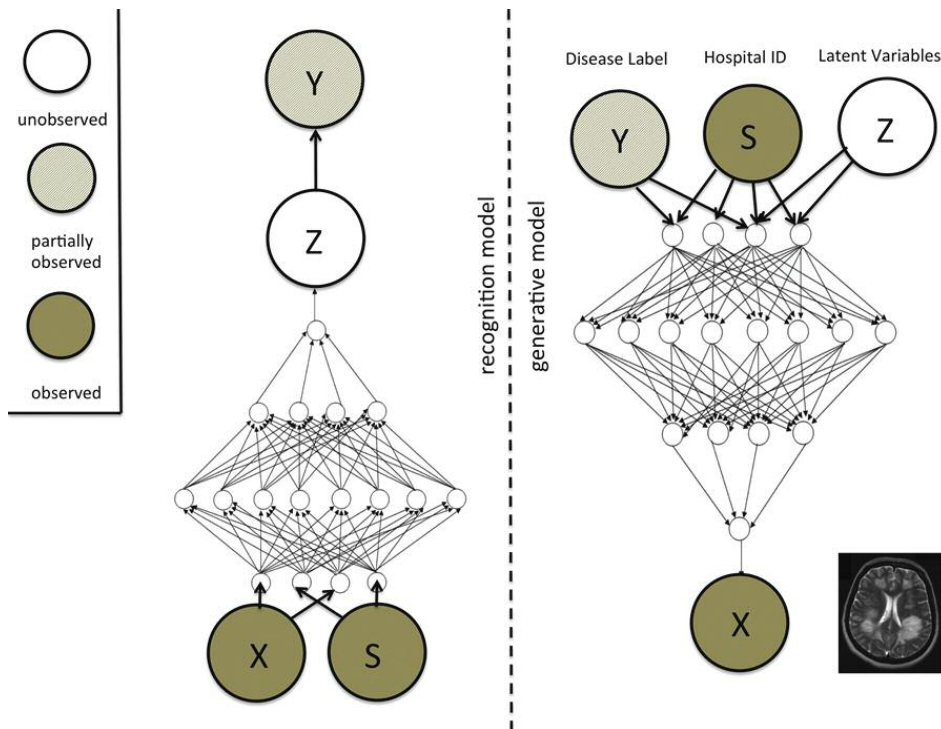


Figure 1: Example of a Variational Auto-encoder model. On the right we have the generative model with three groups of variables at the top: class labels Y , nuisance variables S and latent variables Z . In this example, we may think of Y as disease labels, S as hospital identifiers and Z as size, shape and other relevant properties of a brain. These variables are input to the generative process that generates a pseudo data case; in the example an MRI image of a brain. The discriminative or recognition model on the left takes all the observed variables as input and generates posterior distributions over the latent variables and possibly (if unobserved) the class labels. The models are trained jointly using the variational EM framework.

inject expert knowledge may pay back. Concluding, for very complex tasks (causal) generative models should in my opinion be part of the equation.

Can graphical models and deep neural networks be meaningfully combined into a more powerful framework? The variational auto-encoder (VAE) naturally combines generative models with discriminative models where the generative model can be a Bayesian network or a simulator and the discriminative model a deep neural network [2]. The discriminative model performs inference of the unobserved (latent) variables necessary to perform the (variational EM) learning updates. In this view the discriminative model approximately inverts the generative model. However, one can also interpret the VAE differently if we are more interested in the latent representation itself (based on which we can for instance perform classification). Now the generative model guides the discriminative model to learn interesting, semantically meaningful representations. They represent the fundamental sources of variation that are the input for the generative model. Thus, the

generative model may be viewed as an informed way to regularize the discriminative model.

VAEs as described above are a framework for unsupervised learning. However, they are easily extended to semi-supervised learning by incorporating labels in the generative model [3]. In this case the input data are generated by instantiating a label and some latent variables and sampling from the Bayesian network. In contrast, the discriminative model inverts this relationship and learns a mapping from input directly to class labels and latent variables (see Figure 1). Semi-supervised learning is now easy, because for the unlabelled examples, the label variable is treated as latent, while for a labelled data-case it is treated as observed.

In summary, marrying (discriminative) deep learning with causality and probabilistic reasoning in graphical models may be an important component in reaching the ambitious goals of Artificial General Intelligence. However, most likely completely new ideas are needed as well.

References:

- [1] Pearl, Judea: “Probabilistic reasoning in intelligent systems: networks of plausible inference” Morgan Kaufmann, 2014.
- [2] Kingma, P. Diederik, and M. Welling: “Auto-encoding variational bayes”, the International Conference on Learning Representations (ICLR), Banff, 2014.
- [3] Kingma, P. Diederik et al.: “Semi-supervised learning with deep generative models”, Advances in Neural Information Processing Systems 2014: 3581-3589. <https://www.youtube.com/watch?v=XNZIN7Jh3Sg> http://dpkingma.com/?page_id=277

Please contact:

Max Welling
University of Amsterdam (UvA)
+31 (0)20 525 8256
welling.max@gmail.com
<http://www.ics.uci.edu/~welling/staff.fnwi.uva.nl/m.welling/>

Privacy Aware Machine Learning and the “Right to be Forgotten”

by Bernd Malle, Peter Kieseberg (SBA Research), Sebastian Schrittwieser (JRC TARGET, St. Poelten University of Applied Sciences), and Andreas Holzinger (Graz University of Technology)

While machine learning is one of the fastest growing technologies in the area of computer science, the goal of analysing large amounts of data for information extraction collides with the privacy of individuals. Hence, in order to protect sensitive information, the effects of the right to be forgotten on machine learning algorithms need to be studied more extensively.

Data driven economy (and related concepts like Industry 4.0) and data driven science, as well as big data are the keywords most often heard in discussions on the future of high-profile industries and on the upcoming revolutions in the economic world. With the integration of modern information technology into “classical” industrial environments or services, many new opportunities can be envisioned, e.g., in the optimisation of supply chains or in on-demand production of specifically tailored goods, but even in governmental areas like health environments, where P4-medicine (predictive, preventive, personalised, participatory) is seen as a new paradigm that could revolutionise health care. With all these new opportunities, the challenges were traditionally located in the technical area, especially regarding technologies for enabling the efficient and correct analysis of the large amounts of

data produced by factories and large sensor networks. In recent years, the area of machine learning has seen a surge in new technologies developed and brought to the market. In combination with the ever increasing amount of computational power and storage that is available for a relatively reasonable price, many of these applications can now be applied in real life environments.

While many of the technological issues have been apparently solved, the legal aspects of the collection and processing of vast amounts of data using machine learning algorithms has been neglected (see [1]). The “right to be forgotten” has been recently discussed with a particular focus on removing personal data and sensitive (personal) information from automated analysis if requested by an individual. This brings up technical

as well as ethical questions. Especially in the European Union, the right to be forgotten has remained in political discussion, especially fed by a legal base for protection of personal information on the Internet by the European Commission (see [L1]), with the draft European Data Protection Regulation Article 17 (see [2]). While the main focus in the current discussion is leaning mostly towards the removal of information from search indexes of prominent search engines like Google, the underlying technological challenges run much deeper and touch fundamental aspects of machine learning and its application in the industry.

One of the major questions are the effects of removing information from knowledge bases on machine learning algorithms. This is especially important for algorithms or analytical systems that

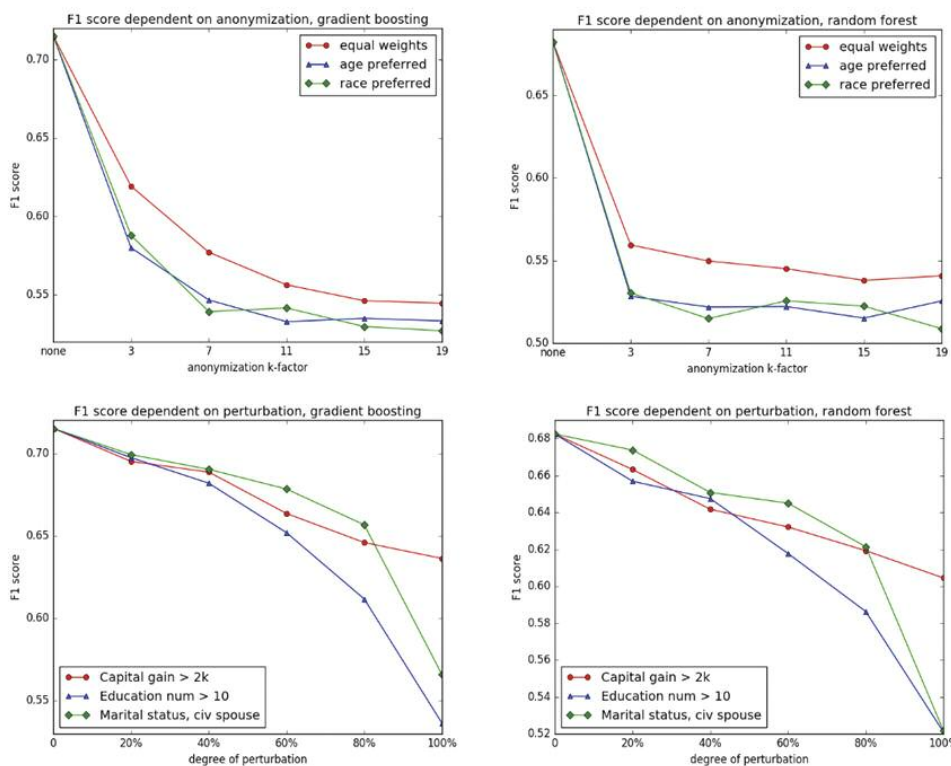


Figure 1: Effects of information removal on selected algorithms.

rely on a large knowledge base containing previously analysed data to learn from or which work by iteratively enhancing the global analytical result by continuously increasing the amount of information. Deletion of existing data thus tends to worsen the quality of the existing results, putting the organisations operating under such law at a serious disadvantage compared to organisations not subject to those restrictions. For instance, in developing clinical decision support software, companies not bound by the right to be forgotten can train their algorithms on a more comprehensive data set, giving them a worldwide advantage in marketing their product. In a project we studied the effects of selective deletion of valuable data items (in terms of their contribution to a classifier accuracy) as well as different levels of anonymization of a whole data set on machine learning algorithms [3]. In those experiments, which were based on data from the 1994 US-census with around 32,000 individual data records, the first phase tested four different classifiers (gradient boosting, linear SVC, logistic regression, random forest) with respect to precision, recall and F1-score. Subsequently, increasing fractions of valuable data were removed from the data set, which resulted in significant loss of classifier performance.

The other major research topic concerned the design of machine learning algorithms and applications that can cope with anonymized or otherwise generalised information. The fundamental idea is that the sensitive information is encoded by some privacy protecting means, analysed using machine learning algorithms and then prepared for inspection. Although several approaches exist for this strategy, they are currently not practical either due to their impact on the quality of the results, or due to the additional costs introduced:

- *Trusted environments.* While being the most popular strategy, the main issue of using a trusted environment lies in the large amount of resources that have to be reserved for this task, even when the analysis is only done very infrequently; e.g., using shared environments like the Cloud is not possible with this methodology. Furthermore, the environment needs to be set up with high security standards

and a lot of audit and control mechanisms, as mechanisms for thwarting insider attacks must be introduced, including fully trusted human operators.

- *Anonymization.* The data set is transformed into a derived set that blurs the sensitive attributes without actually removing them altogether. Popular techniques work by generalising records until a certain minimal amount of them form an equivalence group (are indistinguishable).
- *Pseudonymization.* Related to anonymization, Pseudonymization works by removing sensitive attributes and replacing them with a placeholder, while keeping the internal logical structure of the data set, i.e., records with the same sensitive attributes get assigned the same pseudonym.
- *Functional Encryption.* Functional encryption allows the calculation of certain mathematical operations on the encrypted values, e.g., let $F(x)$ be the encryption function of x and $F^{-1}(y)$ the decryption routing, then it holds true that $x+y=F^{-1}(F(x)+F(y))$. While this works well in theory, currently available algorithms are very slow and thus cannot be used on close to all real life scenarios.

Our experiments to date (see [3] and [L2]) have focused on the loss of classifier performance when applied to anonymised knowledge bases. We used the same data set as described above and anonymised it using the k-anonymity criterion. The classifiers were then re-applied on a series of increasingly anonymised data sets (by increasing the k-factor), again resulting in significant losses of classifier performance. A noteworthy difference between selective deletion and anonymisation of data is that classifier performance on reduced data decreased rather slowly at the beginning and became more drastic with increased fractions of data removed, whereas for anonymised data sets the greatest loss occurred instantly and subsequently mellowed with increasing factors of k-anonymity (see Figure 1).

In conclusion, we can see that the effects of introducing the right to be forgotten to machine learning predictably results in losses of algorithmic performance. However, our experiments so far have only considered classification. A

next logical step in our efforts would be the inclusion of predictors as well as unsupervised learning methods (clustering for automatic label provision, pattern / preference recognition for product design etc.). Even though in reality the effects might not be as drastic as produced by our initial experimental setting, even a few percentage points in ML performance could make a significant difference in crucial areas of application or highly competitive market environments. To sum up, we believe a lot of additional future research is needed in order to:

- fully understand the effects of the right to be forgotten on machine learning environments;
- be able to design algorithms more resilient to changes in the knowledge base;
- understand the effect of perturbing other forms of knowledge bases, e.g., graph based data sets, in which distance is derived from node feature vectors as well as associations.

Links:

- [L1] http://ec.europa.eu/justice/data-protection/document/review2012/com_2012_11_en.pdf
 [L2] <http://www.hci-kdd.org/>

References:

- [1] P. Kieseberg, H. Hobel, S. Schrittwieser et al: "Protecting Anonymity in Data-Driven Biomedical Science", pp. 301-316, 2014.
 [2] P. De Hert, V. Papakonstantinou : "The proposed data protection Regulation replacing Directive 95/46/EC: A sound system for the protection of individuals", Computer Law & Security Review 28, no. 2 (2012): 130-142, 2012.
 [3] B. Malle, P. Kieseberg, E. Weippl, A. Holzinger: "The Right to Be Forgotten: Towards Machine Learning on Perturbed Knowledge Bases", Workshop on Privacy Aware Machine Learning (PAML), August 2016

Please contact:

Peter Kieseberg
 SBA Research, Vienna, Austria
pkieseberg@sba-research.org

Robust and Adaptive Methods for Sequential Decision Making

by Wouter M. Koolen (CWI)

Machine learning systems for sequential decision making continuously improve the quality of their actions. Our new adaptive methods learn to improve as fast as possible in a wide range of applications.

Many practical problems can be cast as sequential decision making tasks, ranging from time series prediction, data compression and portfolio management to online versions of routing and ranking. Yet other problems may be reduced to sequential decision making, in particular batch learning from large data sets. The main challenge is the design of generic, efficient, robust and adaptive learning methods.

Following [1], we model the learning process as a series of interactions between the learner and its environment. Each round the learner picks an action. Then the environment reveals the quality of all available actions by means of a “loss function”. Based on this feedback the learner updates its internal state, with the goal of picking better future actions.

Robustness is traditionally achieved by taking a game-theoretic perspective. We interpret the learning problem as a sequential game, and regard the environment as an adversary that chooses the loss functions to maximally frustrate the learner. One of the accomplishments of the online learning community is the design of learners that are provably successful in such games. It is amazing that generic learning methods can be designed for such a rich class of problems.

Unfortunately, the robust game-theoretic methods do not always fare well in practice, in the sense that for some problems they are outperformed significantly by special-purpose methods. The reason is that, in practice, environments are not maximally evil, and may admit more aggressive learning techniques which are unsafe in general. This realisation has spurred the search for reasonable assumptions that describe the additional useful structure present in practical problems.

We studied two popular types of such assumptions. First, one may require that

the loss functions exhibit curvature, either in every direction (strong convexity) or along the gradient (exp-convexity). For either scenario with known curvature, methods have been developed that increase performance dramatically. Yet in practice the curvature often is a feature of the data, and hence unknowable a-priori. We aim to design methods that adapt to any exploitable curvature presented.

Second, one may consider the relation between loss functions issued by the environment over the course of multiple rounds. A simple (but rich and often

every round based on the newly observed loss function. To exploit the aforementioned structural properties of the environment (curvature and/or Bernstein exponent) it turns out that the magnitude of this update needs to be tuned carefully. The friendlier the problem, the larger the step size required (whence the title of the NIPS workshop series “Learning Faster from Easy Data”).

Tuning the step size, a single scalar, appropriately for some unknown parameters might sound like a simple learning problem in itself. Yet no

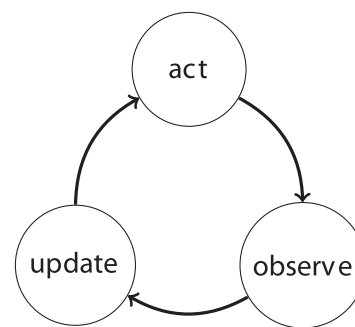


Figure 1: Sequential decision making.

practically reasonable) class of assumptions is obtained by modelling the sequence of loss functions as an independent and identically distributed stochastic process. We may then characterise the simplicity of the distribution by its “Bernstein exponent” (generalisation of the Tsybakov margin condition used in classification). For distributions with a known Bernstein exponent it is possible to design learning algorithms that perform optimally. But how does one deal with and adapt to an unknown Bernstein exponent?

Online learning methods maintain internal parameters that are updated

existing learning method would apply: the overhead for learning the step size would dwarf the benefit of setting it right. The key contribution of our new method, called MetaGrad, is a novel approach for learning the step size from the data.

We proved that MetaGrad has a certain performance guarantee of second-order form [2]. This bound implies in particular worst-case safety (robustness), adaptivity to curvature and adaptivity to the Bernstein exponent [3]. This establishes MetaGrad as the new state-of-the-art adaptive method for online convex optimisation. Since MetaGrad only has

modest computational overhead over earlier methods, we also expect it to be highly relevant in practical applications.

The research team consisted of Wouter M. Koolen (CWI), Tim van Erven (Leiden University) and Peter Grünwald (CWI and Leiden University). The author is funded by an NWO Veni grant. The results will be presented at the 30th NIPS conference in December 2016. Our reference implementation of MetaGrad is publicly available [L1].

Link:

[L1] <https://bitbucket.org/wmkoolen/metagrad>

References:

- [1] S. Shalev-Shwartz: “Online learning and online convex optimization”, *Foundations and Trends in Machine Learning*, 4(2):107–194, 2012.
- [2] T. van Erven, W. M. Koolen: “MetaGrad: Multiple learning rates in online learning”, accepted to *Advances in Neural Information Processing Systems (NIPS) 29*, 2016, Pre-print.

<https://arxiv.org/abs/1604.08740>.

[3] W. M. Koolen, P. D. Grünwald, T. van Erven: “Combining adversarial guarantees and stochastic fast rates in online learning”, accepted to *Advances in Neural Information Processing Systems (NIPS) 29*, 2016. Pre-print <https://arxiv.org/abs/1605.06439>.

Please contact:

Wouter M. Koolen, CWI, The Netherlands.
 +31 20 592 4244
 W.M.Koolen-Wijkstra@cwi.nl

Neural Random Access Machines

by Karol Kurach (University of Warsaw and Google), Marcin Andrychowicz and Ilya Sutskever (OpenAI (work done while at Google))

We propose “Neural Random Access Machine”, a new neural network architecture inspired by Neural Turing Machines. Our architecture can manipulate and dereference pointers to an external variable-size random-access memory. Our results show that the proposed model can learn to solve algorithmic tasks and is capable of discovering simple data structures like linked-lists and binary trees. For a subset of tasks, the learned solutions generalise to sequences of arbitrary length.

Recurrent Neural Networks (RNNs) have recently proven to be very successful in real-world tasks, like machine translation and computer vision. However, success has been achieved only on tasks which do not require a large memory to solve the problem, e.g., we can translate sentences using RNNs, but we cannot produce reasonable translations of really long pieces of text, like books. A high-capacity memory is a crucial component to deal with large-scale problems that contain multiple long-range dependencies.

Currently used RNNs do not scale well to larger memories, e.g., the number of

parameters in a popular LSTM architecture [1] grows quadratically with the size of the network’s memory. Ideally, the size of the memory would be independent of the number of model parameters. The first versatile and highly successful architecture with this property was the Neural Turing Machine (NTM) [2]. The main idea behind the NTM is to split the network into a trainable ‘controller’ and an ‘external’ variable-size memory.

In our paper Neural Random-Access Machines [3] we propose a neural archi-

itecture inspired by the NTM. The Neural Random-Access Machine (NRAM) is a computationally-universal model employing an external memory, whose size does not depend on how many parameters the model has. It has, as primitive operations, the ability to manipulate, store in memory, and dereference pointers into its working memory. The pointers in our architectures are represented as distribution over all memory cells. That is, the memory consisting of M cells is an $M \times M$ matrix, where each row fulfils probability conditions (all elements are non-negative and sum to 1). This trick

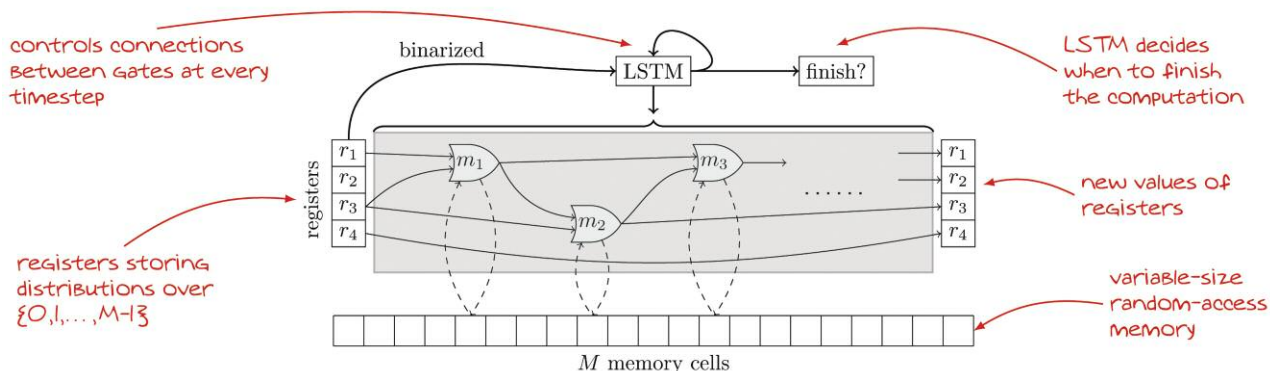


Figure 1: One timestep of the NRAM architecture with 4 registers. The weights of the solid thin connections are outputted by the controller. The weights of the solid thick connections are trainable parameters of the model. Some of the modules (i.e. READ and WRITE) may interact with the memory tape (dashed connections).

allows our model to be fully differentiable and the standard backpropagation algorithm can be used to train it.

Figure 1 gives an overview of the model. NRAM consists of a neural network controller, memory, registers and a set of built-in operations (such as number addition and comparison). An operation can either read (a subset of) contents from the memory, write content to the memory or perform an arithmetic operation on either input registers or outputs from other operations. The controller runs for a fixed number of time steps. At every step, the model selects both the operation to be executed and its inputs. To make this step differentiable, we are using soft attention – each operation is given a linear combination of the inputs, where weights of this linear combination can be controlled by the network. Among other novel techniques, our model employs a differentiable mechanism

for deciding when to stop the computation.

By providing our model with dereferencing as a primitive, it becomes possible to train it on problems whose solutions require pointer manipulation and chasing. It has learned to solve algorithmic tasks and is capable of learning the concept of data structures that require pointers, like linked-lists and binary trees. For a subset of tasks, we show that the found solution can generalise to sequences of arbitrary length. Moreover, memory access during inference can be done in a constant time under some assumptions.

Finally, one may well ask: why train a network to solve a task for which people already know the optimal solution? One reason is that a powerful neural architecture, capable of learning sophisticated algorithms, should be also be able to learn solutions (or approximations) for

complex tasks for which we do not yet know algorithms. We also believe that algorithmic reasoning is one of the necessary (and missing) components in the design of systems able to solve a wide range of real-world problems. The presented model is a step in this direction.

References:

- [1] S. Hochreiter, J. Schmidhuber “Long short-term memory”, *Neural Computation*, Vol 9 Issue 8, pp 1735-1780, MIT Press Cambridge, MA, USA, 1997
- [2] A. Graves et al.: “Neural Turing Machines”, *CoRR*, 2014 <http://arxiv.org/abs/1410.5401>
- [3] K. Kurach et al. “Neural Random-Access Machines”, *ICLR* 2016, <http://arxiv.org/abs/1511.06392>

Please contact:

Karol Kurach, University of Warsaw and Google, Poland
kkurach@gmail.com

Mining Similarities and Concepts at Scale

by Olof Görnerup and Theodore Vasiloudis (SICS)

In machine learning, similarities and abstractions are fundamental for understanding and efficiently representing data. At SICS Swedish ICT, we have developed a domain-agnostic, data-driven and scalable approach for finding intrinsic similarities and concepts in large datasets. This approach enables us to discover semantic classes in text, musical genres in playlists, the genetic code from biomolecular processes and much more.

What is similarity? This fundamental, almost philosophical, question is seldom asked in machine learning in all but specific contexts. And although similarities are ubiquitous in the field, they are often limited to specific domains or applications; calculated between users to make recommendations, between websites to improve web searches, or between proteins to study diseases, for example.

To enable similarity and concept mining applicable in a broad range of areas, we have proposed a generalisation of the distributional hypothesis commonly used in natural language processing: Firstly, the context of an object (a word, artist or molecule, for instance) is its correlations (co-occurrences, play order, interactions etc) to other objects. Secondly, similar objects have similar contexts. And thirdly, a group of inter-similar objects form an abstraction – a concept.

In this way, we essentially base the notion of similarity on exchangeability. In the case of language, for example, it should be possible to remove a word from a sentence and replace it with a similar one (consequently belonging to the same concept) and the sentence should still make sense. And the same should be true regardless of the dataset and types of objects, whether it is words, items a user interacts with, people in a social network, or something else.

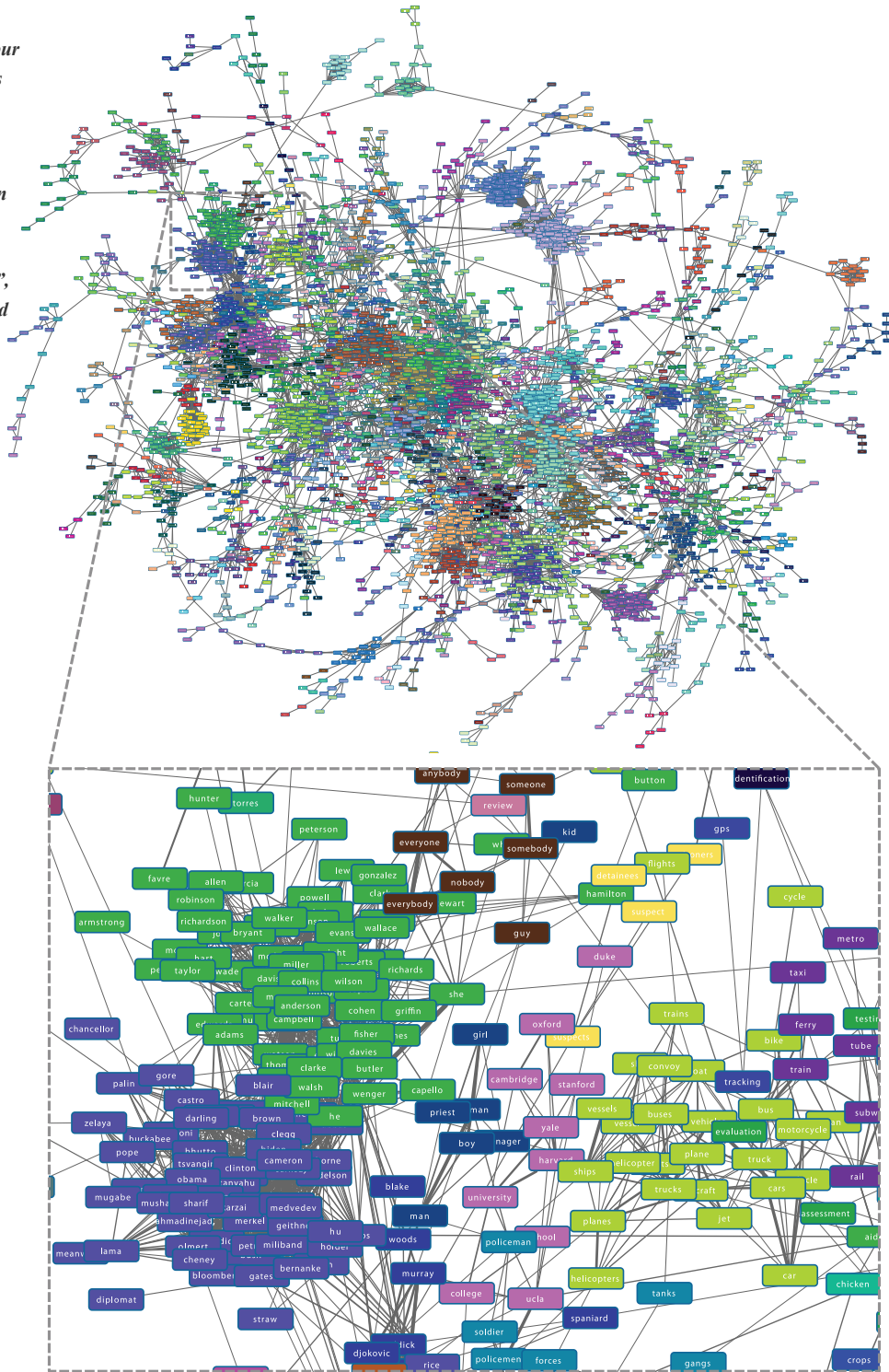
Using these principles, we have developed a method that allows us to find all relevant similarities among a set of objects and their correlations, as well as concepts with respect to these similarities [1,2]. From a bird’s eye view, we achieve this by representing objects and their correlations as a graph, where objects constitute vertices with edges weighed by correlations. We transform

this graph to a similarity graph, where edges are instead weighted by similarities. Concepts are then revealed as clusters of objects in the similarity graph.

This approach is scalable and can therefore be applied to very large datasets, such as the whole of Wikipedia or Google Books (4% of all books ever published). It is also completely data-driven and does not require any human-curated data. At the same time, the method is transparent, so we can interpret the similarities and understand their meaning.

We have demonstrated our approach for three different types of data – text, playlist data, and biomolecular process-data – where similarities and concepts correspond to very different things. In the text case (see Figure 1), objects are words, and words that have

Figure 1: A similarity graph colour coded by concept. The graph was generated using data from the Billion Word corpus. After the similarity transformation was applied to the original correlation graph, community detection was performed to uncover concepts such as “politician”, “university”, or “vehicle” shown in the zoomed in portion.



similar co-occurrence patterns form concepts that correspond to semantic classes, including colours, days of the week, newspapers, vehicles etc. In the playlist case, artists are conceptualised to subgenres and genres. And in the biomolecular case, where codons (triplets of nucleotides in DNA) are related by mutation rates, concepts crystallise as groups of codons that code to the same amino acids, and in effect the method recreates the standard genetic code.

Several improvements and extensions of the method are in the pipeline, including supporting object polysemy (where an object can have several different meanings) and, in particular, mining of higher-order concepts, since an object may also be a concept.

The code for the algorithm is open-source licensed and available at [L1].

Link:
[L1] <https://github.com/sics-dna/concepts>

- References:**
- [1] O. Görnerup, D. Gillblad, T. Vasiloudis: “Knowing an object by the company it keeps: A domain-agnostic scheme for similarity discovery”, in ICDM, 2015, pp. 121–130.
 - [2] O. Görnerup, D. Gillblad, T. Vasiloudis: “Domain-agnostic discovery of similarities and concepts at scale”, in Knowl. Inf. Syst., 2016, pp. 1–30.

Please contact:
Olof Görnerup, SICS, Sweden
+46 70 252 10 62, olof@sics.se

Fast Traversal of Large Ensembles of Regression Trees

by Claudio Lucchese, Franco Maria Nardini, Raffaele Perego, Nicola Tonellotto (ISTI-CNR), Salvatore Orlando (Ca' Foscari University of Venice) and Rossano Venturini (University of Pisa)

The complexity of tree-based, machine-learned models and their widespread use in web-scale systems requires novel algorithmic solutions to make the models fast and scalable, both in the learning phase and in the real-world.

Machine-learned models based on additive ensembles of regression trees have been shown to be very effective in several classification, regression, and ranking tasks. These ensemble models, generated by boosting meta-algorithms that iteratively learn and combine thousands of simple decision trees, are very demanding from a computational point of view. In fact, all the trees of the ensemble have to be traversed for each item to which the model is applied in order to compute their additive contribution to the final score.

This high computational cost becomes a challenging issue in the case of large-scale applications. Consider, for example, the problem of ranking query results in a web-scale information retrieval system: the time budget available to rank the possibly huge number of candidate results is limited due to the incoming rate of queries and user expectations of quality-of-service. On the other hand, effective and complex rankers with thousands of trees have to be exploited to return precise and accurate results [1].

To improve the efficiency of these systems, in collaboration with Tiscali Italia

S.p.A, we recently proposed QuickScorer (QS), a solution that remarkably improves the performance of the scoring process by dealing with features and characteristics of modern CPUs and memory hierarchies [2]. QS adopts a novel bit-vector representation of the tree-based model, and performs the traversal of the ensemble by means of simple logical bitwise operations. The traversal is not performed by QS one tree after another, as one would expect, but is instead interleaved, feature by feature, over the whole tree ensemble. Due to its cache-aware approach, both in terms of data layout and access patterns, and to a control flow that entails very low branch misprediction rates, the QS performance is impressive, resulting in speedups of up to 6.5x over state-of-the-art competitors.

An ensemble model includes thousands of binary decision trees, each composed of a set of internal nodes and a set of leaves. Each item to be scored is in turn represented by a real-valued vector x of features. As shown in Figure 1, the internal nodes of all the trees in the ensemble are associated with a Boolean test over a specific feature of the input

vector (e.g., $x[4] \leq \gamma_2$). Each leaf node stores the potential contribution of the specific tree to the final score of the item. The scoring process of each item requires the traversing of all the trees in the ensemble, starting at their root nodes, until a leaf node is reached, where the value of the prediction is considered. Once all the trees in the ensemble have been visited, the final score for the item is given by the sum of the partial contributions of all the trees.

One important result of QS is that to compute the final score, we only need to identify, in any order, all the internal nodes of the tree ensemble for which the Boolean tests fail, hereinafter false nodes. To perform this task efficiently, QS relies on a bit-vector representation of the trees. Each node is represented by a compact binary mask identifying the leaves of the current tree that are unreachable when the corresponding node test evaluates to false. Whenever a false node is found, the set of unreachable leaves, represented as a bit-vector, is updated through a logical AND bitwise operation. Eventually, the position of the leaf storing the correct contribution for each tree is identified. Moreover, in order to find all the false

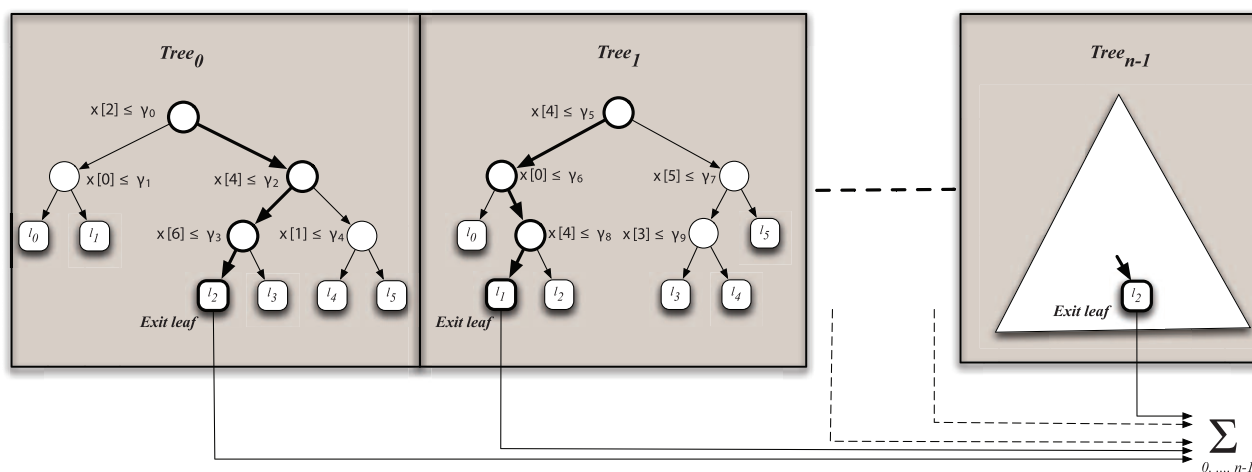


Figure 1: An ensemble of binary decision trees.

nodes for the scored item efficiently, QS processes the nodes of all the trees feature by feature. Specifically, for each feature $x[i]$, QS builds the list of all the nodes of the ensemble where $x[i]$ is tested, and sorts this list in ascending order of the associated threshold γk . During the scoring process for feature $x[i]$, as soon as the first test in the list evaluating to true is encountered, i.e., $x[i] \leq \gamma k$, the subsequent tests also evaluate to true, and their evaluation can be safely skipped and the next feature $x[i+1]$ considered.

This organisation allows QS to actually visit a consistently lower number of nodes than in traditional methods, which recursively visit the small and unbalanced trees of the ensemble from the root to the exit leaf. In addition, QS exploits only linear arrays to store the tree ensemble and mostly performs cache-friendly access patterns to these data structures.

Considering that in most application scenarios the same tree-based model is applied to a multitude of items, we recently introduced further optimisations in QS. In particular, we introduced vQS [3], a parallelised version of QS that exploits the SIMD capabilities of mainstream CPUs to score multiple items in parallel. Streaming SIMD Extensions (SSE) and Advanced Vector Extensions (AVX) are sets of instructions exploiting wide registers of 128 and 256 bits that allow parallel operations to be performed on simple data types. Using SSE and AVX, vQS can process up to eight items in parallel, resulting in a further performance improvement up to a factor of 2.4x over QS. In the same line of research we are finalising the porting of QS to GPUs, which, preliminary tests indicate, allows impressive speedups to be achieved.

More information on QS and vQS can be found in [2] and [3].

References:

- [1] G. Capannini, et al.: “Quality versus efficiency in document scoring with learning-to-rank models”, *Information Processing & Management*, Elsevier, 2016, <http://dx.doi.org/10.1016/j.ipm.2016.05.004>.
- [2] C. Lucchese et al.: “QuickScorer: A Fast Algorithm to Rank Documents with Additive Ensembles of Regression Trees”, *ACM SIGIR 2015*: 73-82 [best paper award].
- [3] Cl. Lucchese, et al.: “Exploiting CPU SIMD Extensions to Speed-up Document Scoring with Tree Ensembles”, *ACM SIGIR 2016*: 833-836.

Please contact:

Raffaele Perego
ISTI-CNR, Pisa, Italy
+39 (0)50 3152993
raffaele.perego@isti.cnr.it

Optimising Deep Learning for Infinite Applications in Text Analytics

by Mark Cieliebak (Zurich University of Applied Sciences)

Deep Neural Networks (DNN) can achieve excellent results in text analytics tasks such as sentiment analysis, topic detection and entity extraction. In many cases they even come close to human performance. To achieve this, however, they are highly-optimised for one specific task, and a huge amount of human effort is usually needed to design a DNN for a new task. With DeepText, we will develop a software pipeline that can solve arbitrary text analytics tasks with DNNs with minimal human input.

Assume you want to build a software for automatic sentiment analysis: given a text such as a Twitter message, the tool should decide whether the text is positive, negative, or neutral. Until recently, typical solutions used a feature-based approach with classical machine learning algorithms (e.g., SVMs). Typical features were number of positive/negative words, n-grams, text length, negation words, part-of-speech tags etc. Over the last two decades a huge amount of research has been invested in designing and optimising these features, and new features had to be developed for each new task.

With the advent of deep learning, the situation has changed: now the computer is able to learn relevant features from the texts by itself, given enough

training data. Solving a task like sentiment analysis now requires three major steps: define the architecture of the deep neural network; aggregate enough training data (labelled and unlabelled); and train and optimise the parameters of the network.

For instance, Figure 1 shows the architecture of a system that won Task 4 of SemEval 2016, an international competition for sentiment analysis on Twitter [1]. This system uses a combination of established techniques in deep learning: word embedding and convolutional neural networks. Its success is primarily based on three factors: a proper architecture, a huge amount of training data (literally billions of tweets), and a huge amount of computational power to optimise its parameters. Live demos of various

deep learning technologies are available at [2].

Goal of DeepText

In DeepText, we will automate the three steps above as far as possible. The ultimate goal is a software pipeline that works as follows (see Figure 2):

1. The user uploads his or her training data in a standard format. The data can consist of unlabelled texts (for pre-training) and labelled texts, and the labels implicitly define the task to solve.
2. The system defines several DNNs to solve the task. Here, different fundamental architectures will be used, such as convolutional or recurrent neural networks.
3. The system then trains these DNNs and optimises their parameters.

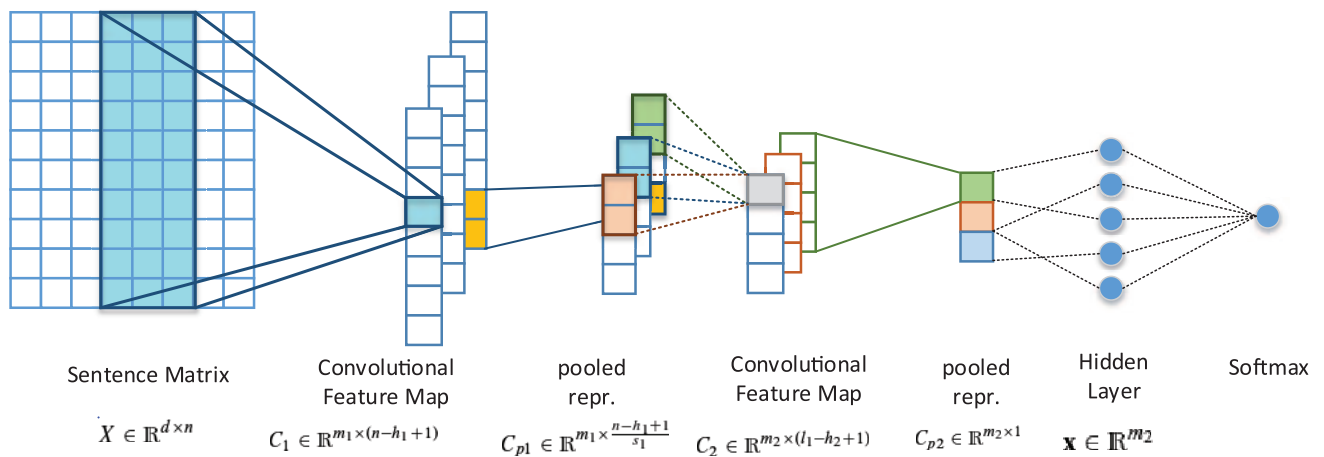


Figure 1: Deep nNeural nNetwork for sSentiment aAnalysis [1].

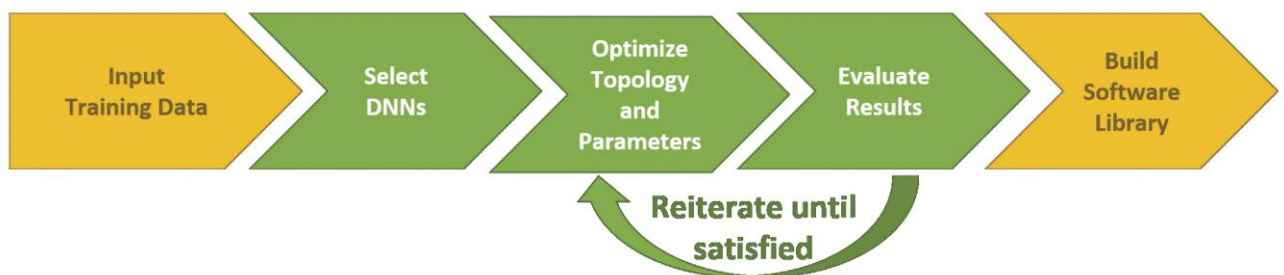


Figure 2: Generating a software library for arbitrary text understanding tasks.

4. Performance of each DNN is measured, e.g. in terms of F1-score, and the best DNN is selected.
5. Finally, the system wraps the “winning” DNN into a software library with a simple interface. This library is ready-to-use in production.

In principle, only the first step – collecting and labelling the training data – needs to be done by humans, since this step defines which task should be solved, and how. For instance, for sentiment analysis on Twitter, each text is labelled with positive, negative, or neutral; on the other hand, if we want to detect companies or persons in text (“entity recognition”), then the proper position of each occurrence of an entity within the text needs to be labelled.

The last three steps in the process above are straightforward, and basically require substantial computational resources and appropriate skills in software engineering.

Challenge: Find a good DNN architecture

The most challenging part is Step 2: to come up with “appropriate” DNNs for

the task at hand. There exist several established DNN architectures for text analytics, such as convolutional neural networks (CNNs) or recurrent neural networks (RNNs). For each architecture, there exist various parameters: in the case of CNNs, this is the number of convolutional layers, size and number of filters, number and type of pooling layers, ordering of the layers etc. In theory, each configuration of a DNN could be used, but this would lead to an explosion of DNNs to evaluate.

For this reason, we will develop several template DNNs for different types of text analytics tasks: classification, topic detection, information extraction etc. Based on these templates, the system will run a pre-training where each template is applied to the task at hand and evaluated. Only the most promising DNNs will then be used for parameter tuning and optimisation.

Our goal is that, given the training data, the system will generate a suitable software library within three days.

About the Project

Deep Text is an applied research project of Zurich University of Applied Sciences (ZHAW) and SpinningBytes AG, a Swiss startup for data analytics. It started in 2016 and is funded by the Commission for Technology and Innovation (CTI) in Switzerland (No. 18832.1 PFES-ES).

Link:

[L1] <http://spinningbytes.com/demos/>

Reference:

[1] J. Deriu et al.: “SwissCheese at SemEval 2016 Task 4”, SemEval (2016).

Please contact:

Mark Cieliebak
 School of Engineering, Zurich
 University of Applied Sciences (ZHAW), Switzerland
 +41 58 934 72 39
 ciel@zhaw.ch

Towards Streamlined Big Data Analytics

by András A. Benczúr, Róbert Pálovics (MTA SZTAKI) , Márton Balassi (Cloudera) , Volker Markl, Tilmann Rabl, Juan Soto (DFKI) , Björn Hovstadius, Jim Dowling and Seif Haridi (SICS)

Big data analytics promise to deliver valuable business insights. However, this will be difficult to realise using today's state-of-the-art technologies, given the flood of data generated from various sources. The European STREAMLINE project develops scalable, fast reacting, and high accuracy machine learning techniques for the needs of European online media companies.

Big data analytics promise to deliver valuable business insights. However, this will be difficult to realise using today's state-of-the-art technologies, given the flood of data generated from various sources. The European STREAMLINE project [L1] develops scalable, fast reacting, and high accuracy machine learning techniques for the needs of European online media companies.

A few years ago, the term “fast data” arose to capture the idea that streams of data are generated at very high rates, and that these need to be analysed quickly in order to arrive at actionable intelligence.

To this end, the EU Horizon 2020 funded STREAMLINE project aims to address the aforementioned technical and business challenges. The STREAMLINE consortium's three research partners, MTA SZTAKI (Hungary), SICS (Sweden), and DFKI (Germany) and four industry partners, Rovio (Finland), Portugal Telecom, NMusic (Portugal), and Internet Memory Research (France) will jointly tackle problems routinely faced by online media companies, including customer retention, personalised recommendation, and web-scale data extraction. Collectively, these four industrial partners serve over 100 million users, offer services that produce billions of events yielding over 10 TB of data daily, and possess over a PB of data-at-rest.

Today's big data analytics systems face two types of latency bottlenecks, namely, system latency and human latency. System latency issues arise due to the absence of appropriate (data) stream-oriented analytics tools and more importantly the added complexity, cost, and burden associated with simultaneously supporting analytics for both data-at-rest and data-in-motion. Human latency results from the heterogeneity of existing tools and the low-level programming languages required for

product or service development that rely on an inordinate number of boilerplate codes are system specific (e.g., Hadoop, SolR, Esper, Storm, and relational database management systems) and demand a plethora of scripts to glue systems together.

Developing analytics that that are well-suited for both data-at-rest and data-in-motion is non-trivial. Prior to our development, even the simplest case had no integral solution when we train a pre-

tion with the items (e.g., clicks, listening, view). In this top-k recommendation task, we have to provide a list of the best k items for a given user. In this task, we contrasted batch and online learning methods. Batch machine learning repeatedly reads all training data multiple times, e.g., via stochastic gradient descent, which uses records multiple times in random order, or via elaborate optimisation procedures, e.g., involving SVMs. The common belief is that these methods are more accurate

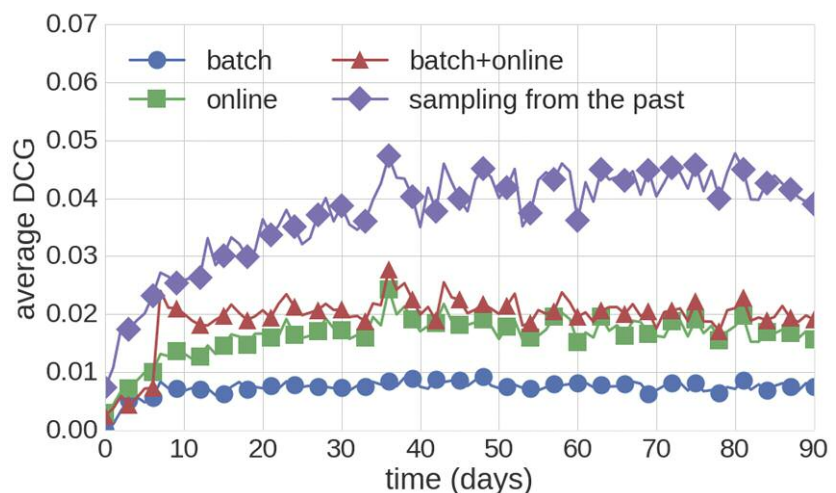


Figure 1: Combined batch and online method system performance.

dictor on historic data and use the streaming API to give predictions with the trained model. In our solution using Flink's unified runtime environment [L2], the model's input for the fitting is drawn from the batch environment and the unlabelled data for the predictor is drawn from the streaming environment without ever needing to explicitly store models or switch between different runtime environments.

As another machine learning example, let us consider a typical recommendation problem of our partners. Their tasks are implicit as users give no ratings: we only have information on their interac-

and easier to implement than online machine learning.

Online learning recommenders perform updates immediately, in the data-streaming model by reading events only once. These models adapt fast to the actual mood of the user, as seen in our experiments over the 30M music listening dataset, crawled by the CrowdRec [L3] team. As shown in Figure 1, online learning outperforms batch trained models in terms of Discounted Cumulative Gain (DCG). We may only gain slight improvements by a loose integration of a linear combination of batch and streaming predic-

tions. The strongest method, on the other hand, requires tight integration by injecting past samples into a data stream matrix factorisation method.

Although Apache Spark is currently a clear leader in the next generation open source big data platform scene both in terms of market penetration and community support, Flink is gaining solid momentum to rival it. We aim to determine whether Apache Flink as a data processing platform has the potential to become a leading player of the open source Big Data market. For example, in Figure 2 we have repeated the benchmarks of Data Artisans [L4] with the latest version of Flink and Spark to depict Alternating Least Squares performance. Another important benchmark that emphasises Flink's low latency capabilities is performed by Yahoo [L5].

In addition to evaluating system latency in comparison with existing alternatives (Spark, Storm), we also aim to evaluate human latency, defined as the effort and time spent on setup, preparation, and development. To do that we will conduct some test subjects to implement data analysis problems, including (i)

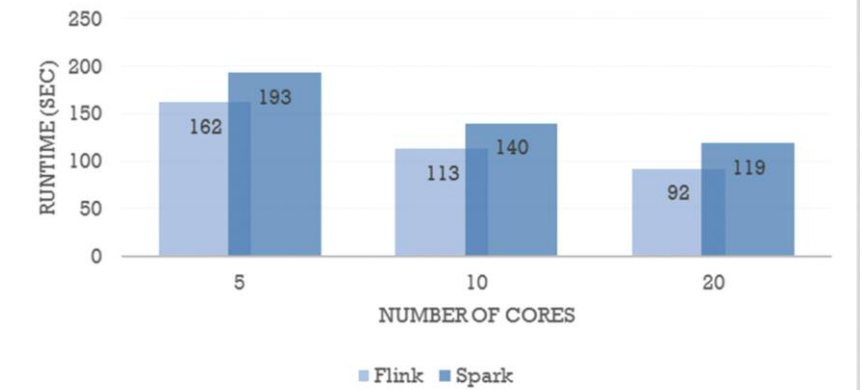


Figure 2: Matrix Factorization, 1 billion entries.

Bachelor and Master students, as well as data scientists from our industrial partners and participants at Flink Hackathons.

The goal of our evaluation under various evaluation criteria and use cases is to determine whether Flink is technologically suited to eventually replace Spark or is destined to remain a niche solution for streaming analytics.

This project has received funding from the European Union's Horizon 2020 research and innovation program under grant agreement No 688191.

Links:

[L1] <http://streamline.sics.se/>

[L2] <http://flink.apache.org/>

[L3] <http://crowdrec.eu/>

[L4] <http://data-artisans.com/computing-recommendations-at-extreme-scale-with-apache-flink/#more-81>

[L5] <http://yahoeng.tumblr.com/post/135321837876/benchmarking-streaming-computation-engines-at>

Please contact:

Björn Hovstadius, STREAMLINE Project Coordinator, SICS, Sweden
bjornh@sics.se

Autonomous Machine Learning

by Frederic Alexandre (Inria)

Inspiration from human learning sets the focus on one essential but poorly studied characteristic of learning: Autonomy.

One remarkable characteristic of human learning is that, although we may not excel in any specific domain, we are quite good in most of them, and able to adapt when a new problem appears. We are versatile and adaptable, which are critical properties for autonomous learning: we can learn in a changing and uncertain world. With neither explicit labels, nor data preprocessing or segmentation, we are able to pay attention to important information and neglect noise. We define by ourselves our goals and the means to reach them, self-evaluate our performances and apply previously learned knowledge and strategies in different contexts. In contrast, recent advances in machine learning exhibit impressive results, with powerful algorithms surpassing human performance in some

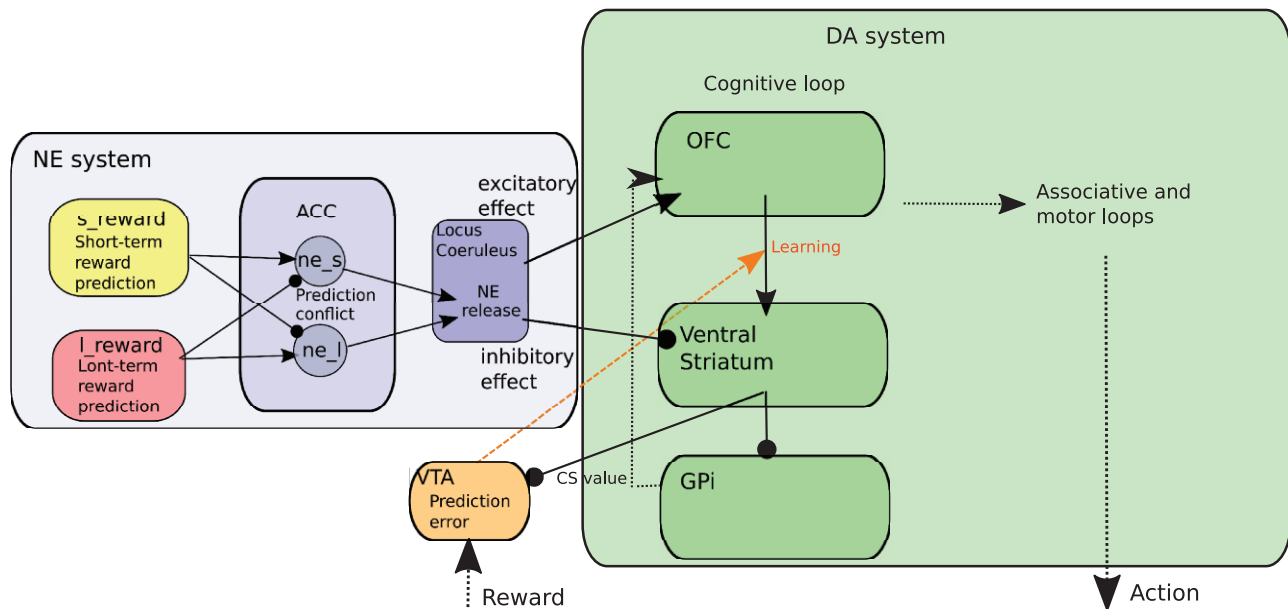
very specific domains of expertise, but these models still have very poor autonomy.

Our Mnemosyne Inria project-team is working in the Bordeaux Neurocampus with medical and neuroscientist teams to develop systemic models in computational neuroscience, focusing on these original characteristics of human learning. Our primary goal is to develop models of the different kinds of memory in the brain and of their interactions, with the objective to exploit them to study neurodegenerative diseases, and another important outcome of our work is to propose original models in machine learning, integrating some of these important characteristics.

We believe that important steps toward autonomous learning can be made along the following lines of research:

Developing an interacting system of memories

Specific circuits in the brain are mobilised to learn explicit knowledge and others to learn procedures. In addition to modelling these circuits, studying their interactions is crucial to understanding how one system can supervise another, resulting in a more autonomous way of learning. In the domain of perceptual learning in the medial temporal lobe, we model episodic memories storing important events in one trial, and forming later, by consolidation in other circuits, new semantic categories. In the domain of decision-making in the loops between



Example of a large scale model described in [3], studying the interactions between two systems in the brain, respectively influenced by noradrenaline (NE, released by Locus Coeruleus) and dopamine (DA, released by VTA) and involving different regions of the loops between the prefrontal cortex (ACC, OFC) and the basal ganglia (Ventral Striatum, GPI). The NE system evaluates the level of non-stationarity of sensory input and modifies accordingly the level of attention on sensory cues, resulting in a shift between exploitation of previously learned rules and exploration of new rules in the DA system performing action selection.

the prefrontal cortex and the basal ganglia, we model cerebral mechanisms by which goal-directed behaviour relying on explicit evaluation of expected rewards can later become habits, automatically triggered with less flexibility but increased effectiveness.

Coping with uncertainty

We learn the rules that govern the world and consider it uncertain for two main reasons: it can be predictable up to a certain level (stochastic rules) or non-stationary (changing rules). Whereas standard probabilistic models are rather good at tackling the first kind of uncertainty, non-stationarity in a dynamic world raises more difficult problems. We are studying how regions of the medial prefrontal cortex detect and evaluate the kind and the level of uncertainty by monitoring recent history of performance at managing correctly incoming events. These regions are also reported to activate the release of neuro-modulators like monoamines, known to play a central role in adaptation to uncertainties [1]. In a nutshell, instead of developing large sets of circuits to manage uncertainty as stable rules in various contexts, the cerebral system has developed a general-purpose system adaptable to uncertainty with hyperparameters sensitive to meta-learning by neuromodulation, which is what we are currently trying to understand more precisely.

Embodiment for emotional learning

One important source of autonomy is our body itself that tells us what is good or bad for us; what must be sought out or avoided. Pavlovian learning is modelled to detect and learn to predict biologically-significant aversive and appetitive (emotional) stimuli which are key targets for attentional processing and for the organisation of behaviour. This learning can be done autonomously if the model of the cerebral system is associated with a substrate corresponding to the body, including sensors for pain and pleasure. We take this a step further, extending the study of the Pavlovian rules to integrate the effects of Pavlovian responses on the body and the neuromodulatory system.

From motivation to self-evaluation

Considering the brain and the body also introduces physiological needs, fundamental to introducing internal goals in addition to the external goals evoked above. This is the basis for renewed approaches regarding reinforcement learning, defining criteria more complex than a simple scalar representing an abstract reward. In humans, another important source of information for learning autonomously is based on self-evaluation of performance. It is noticeable that both motivation and self-evaluation processing are central in cognitive control [2] and reported to be

located in the anterior part of the prefrontal cortex, as we endeavour to integrate in our models.

In addition to developing models to explore each of these mechanisms in interaction with neuroscience and medicine, we also integrate them in a common platform defining the adaptive characteristics of an autonomous agent exploring an unknown virtual world together with the characteristics of its artificial body. Beyond machine learning, this numerical testbed is also a valuable simulation tool for our medical and neuroscientist colleagues.

References:

- [1] A.J. Yu, P. Dayan; "Uncertainty, Neuromodulation and Attention", *Neuron* 46(4), 2005
- [2] M. Carrere, F. Alexandre: "Modeling the sensory roles of noradrenaline in action selection", the Sixth Joint IEEE International Conference Developmental Learning and Epigenetic Robotics, 2016
- [3] E. Koehlin, C. Ody, F. Kouneither: "The Architecture of Cognitive Control in the Human Prefrontal Cortex", *Science*, 302(5648):1181–1185, 2003.

Please contact:

Frederic Alexandre
Inria Bordeaux Sud-Ouest, France
frederic.alexandre@inria.fr
<https://team.inria.fr/mnemosyne/>

Curiosity and Intrinsic Motivation for Autonomous Machine Learning

by Pierre-Yves Oudeyer, Manuel Lopes (Inria), Celeste Kidd (Univ. of Rochester) and Jacqueline Gottlieb (Univ. of Columbia)

Autonomous lifelong multitask learning is a grand challenge of artificial intelligence and robotics. Recent interdisciplinary research has been investigating a key ingredient to reach this goal: curiosity-driven exploration and intrinsic motivation.

A major difference between human learning and most current machine learning systems is that humans are capable of autonomously learning an open-ended repertoire of skills, often from very little data that they actively

for every new task the machines are given.

One of the major components that enables autonomous, open learning in humans is curiosity, a form of intrinsic

aims at pushing the frontiers of what we know about human active learning and how it can be built into machines.

Various strands of work in developmental robotics, AI and machine



Figure 1: Curiosity-driven learning in humans and robots (left: photo by Adam Fenster/Univ. Rochester; right: Milo Keller/ECAL).

collect themselves. Humans show an extraordinary capacity to adapt incrementally to new situations and new tasks. They proactively seek, select, and explore new information to develop skills before they are actually needed.

On the contrary, typical machine learning systems—including those associated with recent advances in deep (reinforcement) learning—learn to solve finite sets of tasks that are predefined by the engineer, and only by access to very large databases of examples. As a consequence, such machines require a new dedicated reward/cost function to be programmed by an engineer and time to reprocess millions of learning examples

motivation that pushes us to actively seek out information and practice new skills for the mere pleasure of learning and mastering them (as opposed to practicing them for extrinsic rewards such as money or social recognition). In the context of the interdisciplinary HFSP project “Curiosity”, Flowers team [L1] at Inria (France), Gottlieb Lab [L2] at Columbia University (US) and Kidd Lab [L2] at University of Rochester (US) are joining forces to study the mechanisms of curiosity-driven active learning in children, adults and monkeys and how they can be modelled and applied with machine learning systems. Mixing artificial intelligence, machine learning, psychology and neuroscience, this project

learning have begun to explore formal models of curiosity and intrinsic motivation (see [1] for a review), providing theoretical tools used in this project. In these models, curiosity is typically operationalised as a mechanism that selects which action to experiment or which (sub-)goals to pursue, based on various information-theoretic measures of their “interestingness”. Many such measures have already been studied with machines and robots—e.g., Bayesian surprise, uncertainty, information gain, learning progress or empowerment—and are often optimised within the reinforcement learning framework, where they are used as intrinsic rewards.

Such algorithmic systems were recently shown to allow machines to learn how to solve efficiently difficult tasks in which extrinsic rewards are rare or deceptive, precluding an easy solution through traditional reinforcement learning methods [1]. These systems were shown to allow robots to efficiently learn multiple fields of parameterised high-dimensional continuous action policies [1]. They also allow robots to self-organise their own learning curriculum, self-generating and self-selecting their own goals, showing a progressive development of new skills with stages that reproduce fundamental properties of human development, for example, in vocal development or tool use [2].

However, many open questions remain. For example, what are the features of interestingness that stimulate the curiosity of human brains? Can current computational models account for them, or be improved by taking inspiration from the heuristics used by humans? Are these mechanisms of

curiosity hardwired or adapted during lifelong learning? As curiosity is a form of guidance for exploration and data collection for autonomous machines, it is also possible to investigate how it can be combined with other forms of guidance used by human-like imitation. For example, in recent robotics experiments, curiosity-driven robots learn repertoires of skills by actively seeking help from human teachers [2].

Finally, curiosity has long been known to be key in fostering efficient education. Computational models of these mechanisms open the possibility for new kinds of educational technologies that could foster intrinsically motivated learning. In recent work, Clement et al. [3] showed one way by presenting active teaching algorithms that were capable of personalising sequences of pedagogical exercises (e.g., math exercises for primary school children), through the dynamic selection of exercises that maximise informational quantities such as learning progress.

Links:

- [L1] <https://flowers.inria.fr>
- [L2] <http://www.gottliebblab.com>
- [L3] <http://www.bcs.rochester.edu/people/ckidd/>

References:

- [1] P-Y. Oudeyer, J. Gottlieb, M. Lopes: "Intrinsic motivation, curiosity and learning: theory and applications in educational technologies", *Progress in Brain Research*, 2016.
- [2] P-Y. Oudeyer, L. Smith: "How Evolution may work through Curiosity-driven Developmental Process", *Topics in Cognitive Science*, 1-11, 2016.
- [3] B. Clement, D. Roy, P-Y. Oudeyer, M. Lopes: "Multi-Armed Bandits for Intelligent Tutoring Systems", *Journal of Educational Data Mining (JEDM)*, Vol 7, No 2, 2015.

Please contact:

Pierre-Yves Oudeyer
Inria and Ensta ParisTech
pierre-yves.oudeyer@inria.fr

Applied Data Science: Using Machine Learning for Alarm Verification

by Jan Stampfli and Kurt Stockinger (Zurich University of Applied Sciences)

A novel alarm verification service applying various machine learning algorithms can identify false alarms.

False alarms triggered by sensors of alarm systems are a frequent and costly inconvenience for the emergency services and owners of alarm systems. Around 90% of false alarms are caused by either technical failures such as network downtimes or human error.

To remedy this problem, we develop a novel alarm verification service by leveraging the power of an alarm data warehouse. In addition, we apply various machine learning algorithms to identify false alarms. The goal of our system is to help human responders in their decision about whether or not to trigger costly intervention forces.

Approach

We are working with a security company that is a major player in secure alarm transmission with the aim of sig-

nificantly reducing the number of false alarms. Alarms can be triggered by devices installed at banks, jewellery stores, private homes, etc.

The problem of alarm prediction is conceptually similar to anomaly detection [1] or prediction of failures [2]. Hence, we can borrow some ideas from these fields and take advantage of the latest progress in deep learning [3]. We have chosen four machine learning approaches to predict false alarms:

- Random forests
- Support vector machines
- Logistic regression
- Deep neural networks.

We based our experiments on a set of more than 300,000 alarms. For each alarm we had multiple features such as device ID, device location, type of

alarm, alarm trigger time, etc. We used these features as input for our machine learning approaches. An overview of the system architecture is given in Figure 1.

For security reasons we do not have direct information about whether an alarm is a false alarm or not. However, we have indirect information that we can use as labels for our machine learning approach. In particular, for each alarm we know when it was triggered and when it was reset again. Our hypothesis is, that if the time difference Δt between triggering and resetting an alarm is small, there is a high chance that the alarm is false.

Consider the case of an alarm system deployed at a private home. Further assume that kids or a pet triggered an

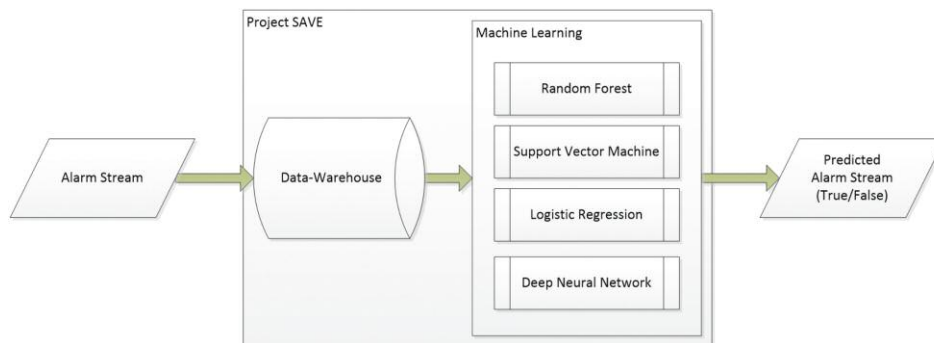


Figure 1: Overview of our alarm prediction system: alarm streams are stored in a data-warehouse and each alarm is evaluated as true or false.

alarm. In this case, the parents or pet owners might call the alarm receiving centre and inform them about a false alarm. In any case, the alarm receiving centre will reset the alarm after verification of the alarm system responsible or the caller identity. In this case, the delta t, i.e., the time between triggering an alarm and resetting it is very small. Details about the four experiments are beyond the scope of this article but can be found in our technical report [L1].

In summary, the goal of our machine learning approach is to learn whether the delta t is below or above a certain threshold value.

Results

In order to evaluate the effectiveness of our machine learning approach, we experimented with various values for delta t ranging between 1 and 10 minutes. The goals of our evaluation were as follows:

- Evaluate the accuracy of four different machine learning algorithms
- Study the impact of various deltas t on the prediction accuracy.

For the training and evaluation of our machine learning approach, we split the data into two sets containing about 170,000 alarms each. Apart from support vector machines, we observe that the performance of the algorithms is not affected by the deltas t (see Figure 2). Random forest and deep neural networks show the best performance with a prediction accuracy of up to 92%. These results show that our system is reliable even if we replace our hypothetical labels when we get access to the real ground truth in the future.

Conclusions

The results demonstrate that our machine learning approaches are very effective for predicting false alarms with an accuracy of up to 92%. These results can be directly used by typical alarm receiving centres for prioritising alarms and thus have a large potential to significantly reduce the costs of dispatching intervention forces.

As part of future work we will integrate this machine learning approach into our alarm data warehouse to enable stream

and batch processing. The idea is to apply the machine learning algorithms in real time on alarm streams and correlate the results with the alarm response to potentially further increase the accuracy of false alarm prediction.

Link:

[L1] <http://pd.zhaw.ch/publikation/upload/210931.pdf>

References:

[1] V. Chandola, A. Banerjee, V. Kumar: “Anomaly detection: A survey”, ACM Comput. Surv. 41, 3, Article 15, 2009.
 [2] F. Salfner, M. Lenk, M. Malek: “A survey of online failure prediction methods”, ACM Comput. Surv. 42, 3, Article 10, 2010.
 [3] Y. LeCun, Y. Bengio, G. Hinton: “Deep learning”, Nature, 521(7553), 436-444, 2015.

Please contact:

Kurt Stockinger
 Zurich University of Applied Sciences,
 Switzerland
Kurt.Stockinger@zhaw.ch

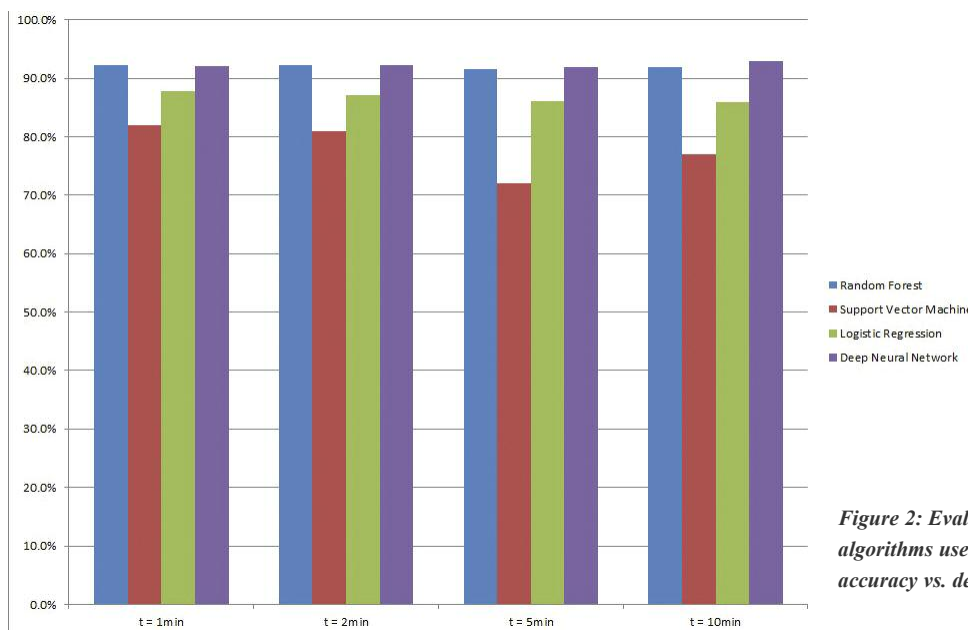


Figure 2: Evaluation of the machine learning algorithms used to predict alarm validity: prediction accuracy vs. delta t.

Towards Predictive Pharmacogenomics Models

by George Potamias (FORTH)

Within the framework of the Greek eMoDiA (electronic Pharmacogenomics Assistant) project, we developed an electronic pharmacogenomics assistance platform using machine learning techniques.

Pharmacogenomics (PGx) has revolutionised drug therapy and unearthed hundreds of associations between genes and drug response, with genome wide association studies (GWAS) and next generation sequencing (NGS) approaches to boost the engaged association discovery process [1]. In order to serve the needs of the different PGx communities – from biomedical researchers to clinical decision-makers and therapy planners – we designed and implemented an electronic PGx Assistance (ePGA) platform [L1], [2,3]. The whole endeavour was conducted in the context of the eMoDiA (Greek funded) project [L2].

ePGA offers two basic services: (i) *explore* – a browser-based service to search through established PGx associations and their accompanying metabolic phenotypes (i.e., extensive, intermediate, poor and ultra-rapid); (ii) *translate* (see Figure 1(a) – an algorithmic process that infers PGx metabolic phenotypes (referred also as a metaboliser statuses) from individual genotype profiles. In the heart of the translation process are the haplotype tables that define PGx haplotypes with reference to gene-variants. PharmGKB (a state-of-the-art PGx knowledge base, www.phar-

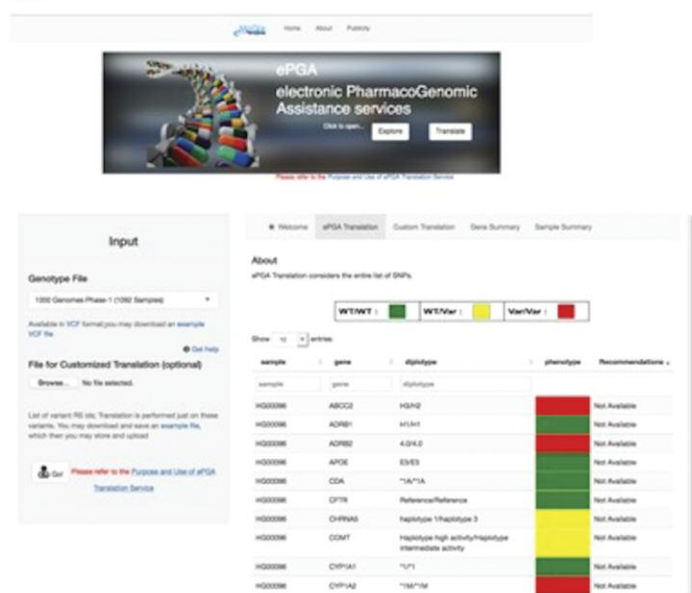
mgkb.org), provides haplotype tables for a total of 69 pharmacogenes (genes for which respective haplotypes are related with metabolising response profiles) that engage a total of 786 variants.

The ePGA translation process was applied on the phase-I 1000 Genomes (1kG) project samples (1092 in total, covering fourteen populations, www.1000genomes.org). For each gene, a sample is assigned to a PGx phenotype status, i.e., the class target variable, which may take one of the following three (phenotype status) values: “RR” (reference/wild-type), “RV” (reference-variant), or “VV” (variant-variant). The variants (rsIDs) represent the descriptor features that take values over all the possible allele combinations, e.g., the variant “rs12475068: C>G” may take the values CC, CG (or GC) and GG. So, we have at our disposal a genotype matrix with rows the variants, columns the samples, and cell values the genotypes of the respective samples. The target is to form patterns of variants that could model, and in a way explain the assignment (by the translation process) of phenotype class to the samples – a classical machine-learning problem, which we tackle with the utilisation of decision-tree induction techniques. We

restricted our experiments to 28 genes for which at least 20% of the samples are assigned to one phenotypic class. For each gene the respective genotype data (restricted to the rsIDs present in the original gene haplotype table) were used as input to a decision-tree induction algorithm (the Weka J48 tool was utilised). For each gene a respective decision-tree was induced. The results are presented in Figure 1(b).

The high performance figures, LOOCV: 99.7%, Sensitivity: 99.8%, Specificity: 99.6% and ROC/AUC: 0.997 are indicative for the validity of the whole approach. The highly predictive results are achieved with a reduced set of variants – the percentage of common variants is 49.6% resulting in a 50.4% reduction in the number of utilised variants (column “redund.V%”). So, we can safely state that: using a representative set of genotyped samples we are able to identify a reduced number of “critical” and “informative” variants, and form respective predictive models (PGx-decision-trees) that are able to accurately infer the PGx phenotype phenotypic metaboliser status of the sample cases. We also applied the same approach on genotype profiles that include the whole set of variants present

1.a



1.b

Gene	#Cases	1KG %	LOOCV	SE	SP	AUC	table#V	tree#V	redund.V%
ABCB1	533	48.8%	100.0%	100.0%	100.0%	1.000	6	1	83.3%
ADRB1	870	79.7%	100.0%	100.0%	100.0%	1.000	2	2	0.0%
APOE	1091	99.9%	99.0%	99.9%	100.0%	1.000	2	2	0.0%
BRCA1	553	50.6%	99.1%	99.1%	99.4%	0.991	9	2	77.8%
CHRNA5	910	83.3%	100.0%	100.0%	100.0%	1.000	2	2	0.0%
COMT	581	53.2%	100.0%	100.0%	100.0%	1.000	4	1	75.0%
CYP1A1	928	85.0%	99.4%	99.4%	99.6%	0.997	10	4	60.0%
CYP1A2	561	51.4%	98.9%	98.9%	98.4%	0.995	25	5	80.0%
CYP1B1	790	72.3%	99.9%	99.9%	100.0%	0.999	18	3	83.3%
CYP2C9	999	91.5%	98.8%	98.8%	98.2%	0.984	20	6	70.0%
CYP2E1	913	83.6%	99.8%	99.8%	99.0%	0.998	6	3	50.0%
CYP3A4	1092	100.0%	100.0%	100.0%	100.0%	1.000	1	1	0.0%
CYP3A5	974	89.2%	99.5%	99.5%	99.3%	0.998	12	4	66.7%
CYP3A7	1002	100.0%	100.0%	100.0%	100.0%	1.000	1	1	0.0%
CYP4B1	1044	95.6%	99.5%	99.5%	99.6%	0.996	6	3	50.0%
CYP4F2	1039	95.1%	100.0%	100.0%	100.0%	1.000	2	2	0.0%
DDC	415	38.0%	100.0%	100.0%	100.0%	1.000	5	1	80.0%
LDLR	541	49.5%	100.0%	100.0%	100.0%	1.000	5	4	20.0%
NAT1	981	89.8%	99.3%	99.3%	99.5%	0.993	24	2	91.7%
PK3CA	1020	93.4%	100.0%	100.0%	100.0%	1.000	5	1	80.0%
SCN1A	862	78.9%	100.0%	100.0%	100.0%	1.000	2	1	50.0%
SCNN1B	795	72.8%	100.0%	100.0%	100.0%	1.000	4	1	75.0%
SLC25A27	687	62.9%	100.0%	100.0%	100.0%	1.000	4	2	50.0%
SULT1A2	929	85.1%	100.0%	100.0%	100.0%	1.000	3	2	33.3%
SULT2A1	1088	99.6%	99.8%	99.8%	96.7%	0.965	3	2	33.3%
TPMT	259	23.7%	100.0%	100.0%	100.0%	1.000	30	1	96.7%
UGT2B15	897	82.1%	99.7%	99.7%	99.5%	0.999	4	3	25.0%
VKORC1	562	51.5%	100.0%	100.0%	100.0%	1.000	10	2	80.0%
Average	821.6	75.2%	99.7%	99.8%	99.6%	0.997	8.0	2.3	50.4%

Figure 1: (a) The ePGA translation service; and (b) Results of PGx decision-tree prediction models (using just the variants in haplotype tables).

in the region of the respective genes. The high performance figures, LOOCV: 99.6%, SE: 96.1%, SP: 99.5% and AUC: 0.996, achieved with a common variants' percentage of 60%, also confirm the validity of our decision-tree induction approach.

Our on-going and future R&D work in the field includes: (a) investigation of the ways to form new haplotype tables utilising the set of variants present in the respective PGx decision-trees, and (b) experimentation with more sets of population representative genotyped samples in order to achieve a more extended and deeper validation of the approach.

Links:

- [L1] <http://www.epga.gr>
- [L2] <http://www.emodia.gr>

References:

- [1] G. Potamias K.Lakiotaki, T. Katsila, et al.: "Deciphering next-generation pharmacogenomics: an information technology perspective", Open Biol. 4(7), 2014, doi: 10.1098/rsob.140071.
- [2] K. Lakiotaki, G. Patrinos, G. Potamias: "Information Technology meets Pharmacogenomics: Design Specifications of an Integrated Personalized Pharmacogenomics Information System", IEEE-EMBS Int.

- Conf. Biomed. Heal. Informatics, pp. 13–16, 2014.
- [3] K. Lakiotaki, et al.: "ePGA: a web-based Information System for Translational Pharmacogenomics," PLOS One 11(9):e0162801. doi:10.1371/journal.pone.0162801, 2016.

Please contact:

George Potamias, ICS-FORTH, Greece
 potamias@ics.forth.gr

Optimisation System for Cutting Continuous Flat Glass

by José Francisco García Cantos, Manuel Peinado, Miguel A. Salido and Federico Barber (AI2-UPV)

Metaheuristic techniques can help to optimise various industrial processes. These techniques can be used, for example, to help plan cutting processes in continuous float glass manufacturing, resulting in reduced losses due to material defects.

The processes within industry have been evolving over the years. In the early days, workers were necessary to run the machines, but today most of the work is being automated, resulting in cost savings but also creating difficult optimisation problems. This can be illustrated with the problem of cutting stock in general and glassmaking in particular. In this manufacturing process, cuts are generated automatically, but culling due to existing defects should be minimised.

Float glass is a sheet of glass made by decanting molten glass on a bed of molten metal, typically tin. The manufacturing process starts with a mixture of materials in an oven at a high temperature (Figure 1). The molten mixture passes through a pool of liquid tin where the glass surface is formed. An array of scanners detects the type and position of defects that may appear in the manufacturing process. According to the production orders (number, size and quality of glass sheets) and found defects, cut planning is carried out ensuring the required quality of each order and minimising losses of glass (cullet). These losses occur when, due to existing defects in a zone of the glass ribbon, the required quality cannot be achieved in any order.

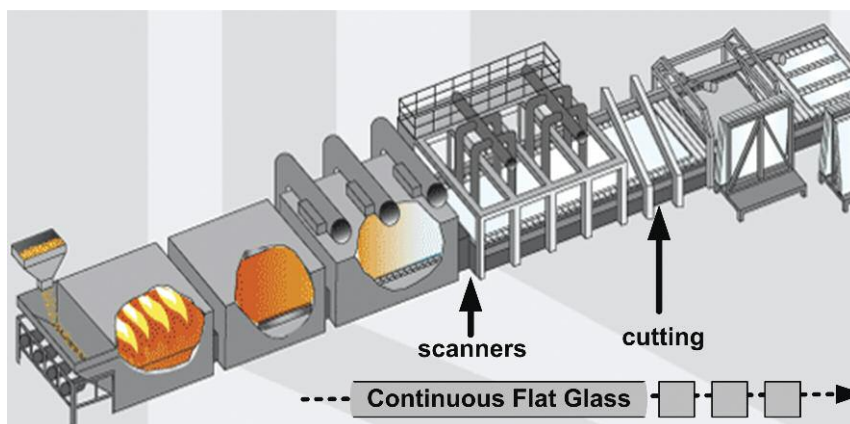


Figure 1: Flat glass production process. For more details see [L1].

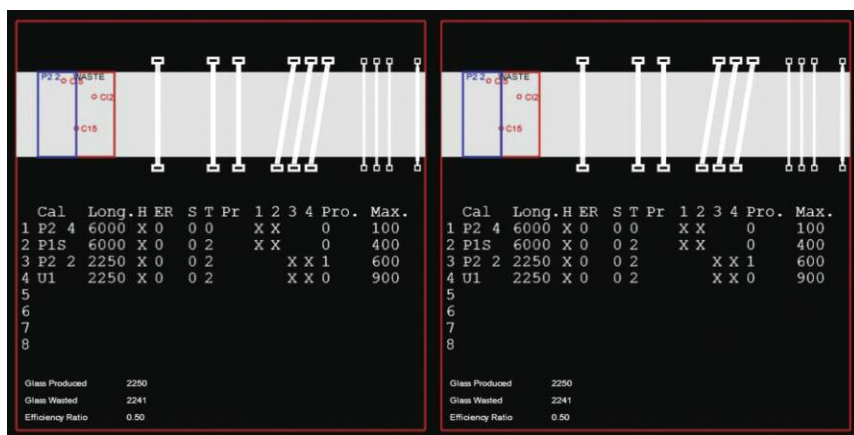


Figure 2: Cut planning process [L1].

After cutting, glass sheets are allocated to a loader, so there cannot be more orders in production than available loaders.

The project developed at the Instituto Universitario de Automática e Informática Industrial (AI2) of the Universitat Politècnica de València (UPV) in the framework of the TETRACOM action “OPTIGLASS: Artificial intelligence-based techniques for optimizing the continuous glass cutting problems” aims to optimise the cutting process planning in order to reduce losses [1]. The action was realised in collaboration with AGC Glass Europe, one of the most important global glass manufacturers.

The company was already planning cuts automatically, although the cut planning process was not fully optimised. A greedy algorithm that chooses the highest priority order that meets the quality requirements to produce a glass sheet was applied. Therefore, if no order fulfilled the quality requirements, a loss was generated. Consequently, this

method allows glass to be produced following a priority, but it does not minimise cullet.

Metaheuristic techniques are higher-level searching procedures, which guide subordinate heuristics, allowing combining exploring and exploiting of the search space in combinatorial optimization problems and finding efficiently near-optimal solutions. Our metaheuristic-based method starts by obtaining all cutting patterns that exist every 16 metres (the distance between the guillotine and the scanners) choosing the production of the sheets that generates less cullet. Afterwards, it searches for the best pattern, taking into account the objective of reducing cullet at an affordable computational cost (Figure 2).

The second problem was to decide the best sequence of orders to be introduced in the system, since the maximum number of simultaneous orders is limited by the number of available loaders. This sequence was decided by production experts, but this problem can be

also automated. The planning process cut was complemented by another heuristic process for determining the best set of orders, which is the best combination of required quality and size that minimises cullet, according to the defects found in the glass ribbon.

Once we had solved the optimisation problems encountered in flat glass manufacturing, a cutting glass simulator was developed. In this simulator, a series of orders are introduced, as well as randomly generated defects according to historical data of the company. The simulator obtains the best cutting positions and the insertion of new orders in the system (Figure 3).

A set of tests of the developed system were carried out in different scenarios including different sets of orders and various glass configurations and defects. In these tests, over 90% of glass was used (and in most cases over 95%) with the new planning system, therefore losses are low, below 10% (Figure 4). As figures show, a reduction of losses greater than 20% can be obtained with respect to a non-optimised planning.

The project demonstrates the applicability of optimisation systems for material cost savings, as well as a clear energy saving in manufacturing processes of flat glass. As recognised by the company management: “The project has been very profitable. The optimisation reduces losses and increases the competitiveness of the company in the market. This kind of relationship and technology transfer should be encouraged” [L3].

Links:

[L1] <http://www.metroglass.co.nz/catalogue/009.aspx>

[L2] <http://gps.blogs.upv.es/software-transferencia/>

[L3] <http://www.tetracom.eu/impact>

Reference:

[1] “OPTIGLASS: Artificial Intelligence-based techniques for optimizing the continuous Glass Cutting Problems”. FP7-ICT-2013-10-Nº 609491, TEchnology TRAnsfer in COMputing systems, TETRACOM-U.E.

Please contact:

Federico Barber, Ai2, UPV, Spain
fbarber@dsic.upv.es

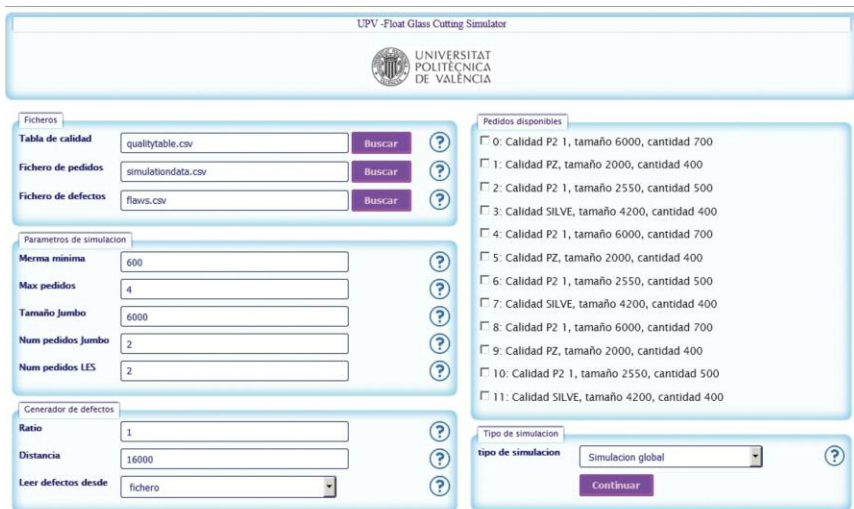


Figure 3: Web-based simulator [L2].

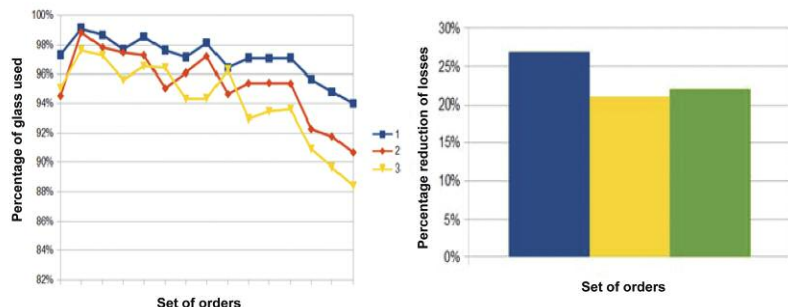


Figure 4: Ribbon glass exploitation depending on combination of orders and obtained improvement.

Online Learning for Aggregating Forecasts in Renewable Energy Systems

by Balázs Csanád Csáji, András Kovács and József Váncza (MTA SZTAKI)

One of the key problems in renewable energy systems is how to model and forecast the energy flow. At MTA SZTAKI we investigated various stochastic times-series models to predict energy production and consumption, and suggested an online learning method which can adaptively aggregate different forecasts while also taking side information into account. The approach was demonstrated on data coming from a prototype public lighting microgrid containing photovoltaic panels and LED luminaries.

The presented research was motivated by the E+grid project which aims at building an energy-positive public lighting microgrid using photovoltaic panels, LED luminaries that regulate their lighting levels based on motion sensor signals, energy storage, various sensors and smart meters, wireless communication and a central controller (see Figure 1). A prototype system was developed by an industry-academy consortium formed by GE Hungary, the Budapest University of Technology and Economics, and two institutes of the Hungarian Academy of Sciences (MFA and SZTAKI). The physical prototype, containing 191 luminaries and 152 m² of PV panels, is located in Budapest at the campus of the MFA Institute of the Hungarian Academy of Sciences [1].

Stochastic Models of Energy Flow

A crucial problem in renewable energy systems is to model energy flow. It is a challenging task, as both energy production and energy consumption are affected by various external factors, and

hence, highly uncertain and dynamically changing. On the other hand, such models are needed to generate forecasts and to build efficient controllers. Several models have been suggested for this in the past, including clear-sky models (i.e., an estimate of the terrestrial solar irradiance under the assumption of a cloudless sky based on astronomical calculations), persistence approaches, autoregressive models, neural networks, fuzzy and hybrid models [2].

During the E+grid project we experimented with a number of time-series models and, after suitable preprocessing (such as removing outliers, noise reduction and normalisation), fitted separate dynamic models to the energy production and consumption processes. We used discrete-time stochastic models with one hour as the time step. The applied models can be classified in two groups: linear and nonlinear. The linear models were: FIR (finite impulse response), AR (autoregressive), ARX

(autoregressive exogenous), ARMA (autoregressive moving average), BJ (Box-Jenkins) and state space, while the nonlinear models were: HW (Hammerstein-Wiener), Wavelet, MLP (Multilayer Perceptron), MLPX (MLP with exogenous inputs), SVR (Support Vector Regression) and SVRX (SVR with exogenous inputs) [2].

For the models with exogenous components, we supplied side information as the inputs to help, for example, to cope with the quasi-periodic nature of the problem as well as to provide the available background knowledge on the modelled phenomenon. Side information included the clear-sky prediction for the case of photovoltaic energy production, while it was the typical consumption pattern (based on historical data) for the specific hour of the day, in case of consumption.

After the models were identified, the innovation (noise) sequences driving the processes were estimated. Based on the

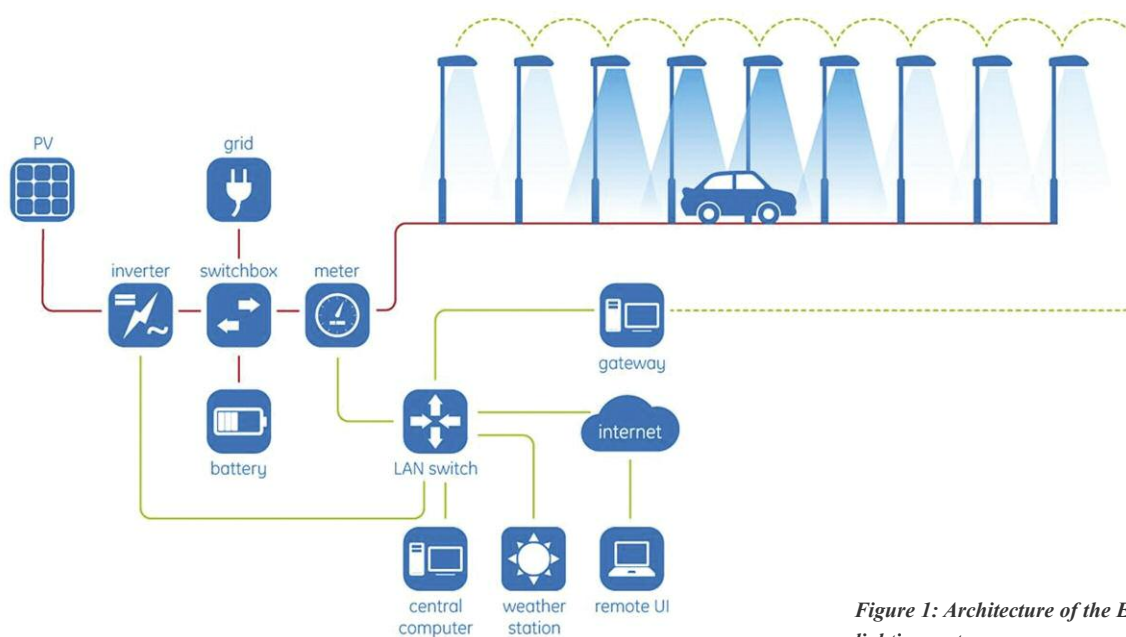
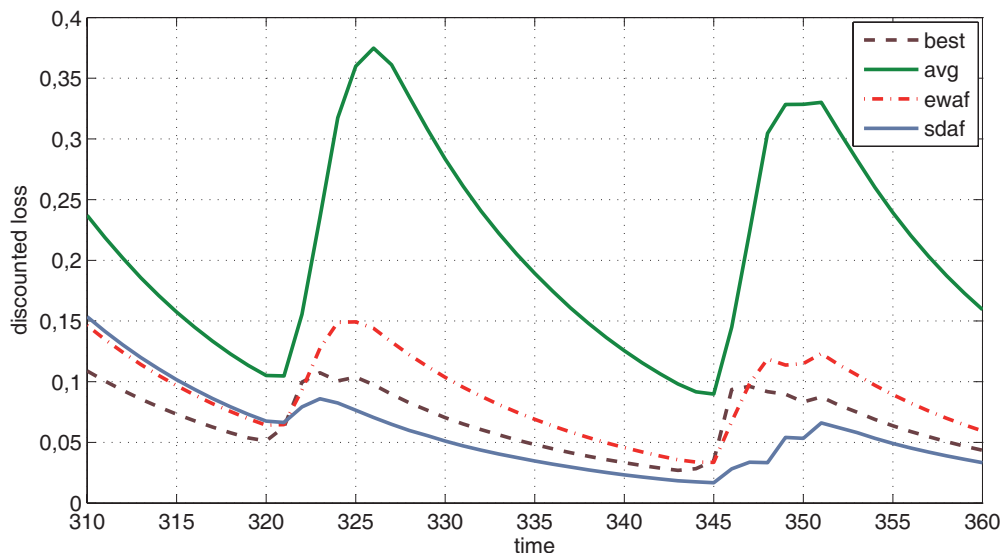


Figure 1: Architecture of the E+grid lighting system.

Figure 2. Discounted cumulative loss of predicting photovoltaic energy production for the best, the average, the exponentially weighted average and the state-dependent exponentially weighted average forecaster.



process and the noise models, forecasts were made by Monte Carlo methods. In the E+grid system, 24-hour forecasts were generated hourly and were used by a receding horizon controller [1].

Online Learning for Context Dependent Forecast Aggregation

While experimenting with various stochastic models we observed that there was no uniformly best model; some models performed better in some situations but worse in others. Since generating forecasts with the already estimated models is computationally cheap, we decided to use all of the models and aggregate their predictions online, based on their past performances in similar situations.

For online learning the best forecasts, we applied the framework of prediction with expert advice. In this framework a learner sequentially faces the problem of predicting an unknown environment based on the predictions and past performances of a pool of experts. The learner aims at minimising its regret, i.e., the difference of its cumulative loss compared to that of the best performing expert so far. The loss is typically defined as the distance between the predicted and actual outcomes of particular variables in the environment. A standard and widely used aggregation rule to combine the predictions of the experts based on their past losses is the exponentially weighted average forecaster (EWA) [3].

In our case the experts were the estimated time-series models based on which the forecasts were generated. We

refined the standard framework by taking contextual information into account as well; namely the losses were weighted by a suitably defined similarity kernel which described how similar the current situation was to the past one in which the expert (the stochastic model) incurred the loss. We also applied discounting to help focus on recent events (e.g., losses incurred a long time ago had lower weights). In addition to the similarity and temporal weighting, our approach, called the state dependent average forecaster (SDAF) was similar to EWA, e.g., exponential weighting was applied [2].

Experimental Results and Conclusions

Several numerical experiments were performed on the energy production and consumption data coming from the prototype E+grid system [2]. Our results indicate that the applied time-series models, especially the ones using side information, can be efficiently applied to forecast the energy flow in the system. They also demonstrate (see Figure 2) that aggregated approaches can provide better forecasts than single time-series models in themselves. Furthermore, they show that our context dependent aggregation approach (SDAF) outperforms the standard context independent EWA for this kind of prediction problems.

This work has been supported by the Hungarian Scientific Research Fund (OTKA), projects 113038 and 111797. B. Cs. Csáji acknowledges the support of the János Bolyai Research Fellowship No. BO/00217/16/6.

References:

- [1] A. Kovács, R. Bártai, B. Cs. Csáji, P. Dudás, B. Háy, G. Pedone, T. Révész, and J. Váncza: "Intelligent Control for Energy-Positive Street Lighting", *Energy*, Elsevier, Vol. 114, 2016, pp. 40–51.
- [2] B. Cs. Csáji, A. Kovács, and J. Váncza: "Adaptive Aggregated Predictions for Renewable Energy Systems", in *Proceedings of the 2014 IEEE Symposium on Adaptive Dynamic Programming and Reinforcement Learning (ADPRL 2014)*, Orlando, Florida, Dec. 9-12, 2014, pp. 132-139.
- [3] N. Cesa-Bianchi, G. Lugosi: "Prediction, Learning, and Games", Cambridge University Press, 2006.

Please contact:

Balázs Csanád Csáji
MTA SZTAKI, Hungary
+ 36 1 279 6231
balazs.csaji@sztaki.mta.hu

Bonaparte: Bayesian Networks to Give Victims back their Names

by Bert Kappen and Wim Wiegerinck (University Nijmegen)

Bonaparte is a software application developed by Smart Research, a subsidiary of SNN at Radboud University Nijmegen, which uses Bayesian Networks to quickly identify a large number of disaster victims.

When a mass disaster occurs – for example, an airplane crash or a terrorist attack – there are many important issues to be addressed, including the quick recovery and identification of the remains of the victims. This is of great importance for the families of the victims, so that they can start the grieving process. Bayesian networks can help forensic scientists perform this complex task promptly and reliably.

Disaster victim identification (DVI) is greatly facilitated by modern DNA technology. Forensic laboratories can extract and record DNA profiles from miniscule samples. The probability that two individuals (except identical twins) have the same DNA profile is extremely small. This makes DNA excellent for identification.

Matching the remains of the victims with reported missing persons, however, is complicated by the absence of comparative direct DNA samples from the missing person. In this situation, DNA samples are usually obtained from indirect sources (family members). Linking victims with their closest relatives instead of their own DNA is far more difficult since they share some, but not all, of their DNA. In a mass disaster this is further complicated by the scale of the problem.

In order to be well prepared for a mass disaster incident, The Netherlands Forensic Institute (NFI) needed a computer assisted disaster victim identification system. The demands were high from the start: it had to be designed for large scale incidents where large numbers of samples would have to be matched as promptly and reliably as possible.

This has resulted in the development of Bonaparte, a software application that does the matchmaking in minutes instead of months. The software was developed by Smart Research, a sub-

siary of SNN at Radboud University Nijmegen by order of and in close collaboration with the NFI.

Bonaparte is named after emperor Napoleon Bonaparte, who introduced the register office in the Netherlands and thus gave every Dutch person a name. The aim of the Bonaparte software is to give the unknown victims back their names.

Bonaparte's core computation task is the following. Given a family with a missing person (MP), for whom there is no DNA sample, and given the avail-

these two hypotheses. If LR is sufficiently high, the forensic expert will decide that there is a match.

To perform the LR computations, Bonaparte uses Bayesian networks. Bayesian networks are a generic class of multivariate probabilistic models. The Bayesian networks in Bonaparte model the inheritance of DNA from parents to child. In addition, they contain an observation model, which, for instance, can handle allele drop-out. The structure of the Bayesian network can be derived from the structure of the family tree. When a family tree is



Convoy of hearses with MH-17 victims on the Dutch highway.

Source: Ministry of Defence, The Netherlands.

ability of DNA profiles both of other family members and an unknown individual (UI), is UI equal to MP? Bonaparte computes the likelihood ratio LR of the two hypotheses. The first hypothesis is that UI=MP. The second hypothesis is that UI is a random person from the population, unrelated to the family. The LR is defined as the ratio of probabilities of finding UI's actual DNA profile under

entered in Bonaparte, the Bayesian network is automatically generated. The strength of the links between individuals follow the rules of Mendelian inheritance. With the resulting model and generic Bayesian network inference algorithms, combined with a method called value abstraction to reduce memory size, all the probabilities required for the LR can be computed.

Bonaparte has been in use by NFI since 2010. The software played an important role in the identification of the victims of the 2010 air disaster in Tripoli, and more recently, in the identification of the victims of Malaysia Airlines flight MH17 in the Ukraine in 2014. Bonaparte is also used for familial searching: in the case of a serious crime, NFI is permitted to search in the national criminal DNA database for relatives of the (unknown) donor of a DNA trace at the crime scene. In 2014 this led to the conviction of a serial rapist who was found via his brother.

Bonaparte has received international interest. Interpol will use Bonaparte for their international missing persons' program. The Australian government bought the software to improve their DNA case work. Finally, the Vietnam government will use Bonaparte for the identification of victims of the Vietnam War, which is expected to be the largest DNA identification project ever conducted.

Bonaparte is being continually refined, for example, by extending the types of DNA profiles that can be entered into the system according to the latest forensic developments.

Link:

<http://www.bonaparte-dvi.com>

References:

[1] Heckerman D. & Wellman, M.P. (1995) Bayesian networks. *Commun. ACM* 38, 3, 27-30.

[2] Bruijning-van Dongen, C. J., Slooten, K., Burgers, W., & Wiegerinck, W. (2009). Bayesian networks for victim identification on the basis of DNA profiles. *Forensic Science International: Genetics Supplement Series*, 2(1), 466-468.

Please contact:

Bert Kappen
Radboud University Nijmegen,
The Netherlands
+31 243614241
b.kappen@science.ru.nl

NIPS 2016 – 30th Annual Conference on Neural Information Processing Systems

Barcelona Spain, 5-10 December 2016

NIPS is a multi-track machine learning and computational neuroscience conference that includes invited talks, demonstrations, symposia and oral and poster presentations of refereed papers. Following the conference, there are workshops which provide a less formal setting.

Invited Speakers:

Yann LeCun (Facebook), Susan Holmes (Stanford), Kyle Cranmer (NYU), Saket Navlakha (Salk Institute), Drew Purves (Deep Mind), Marc Raibert (Boston Dynamics), Irina Rish (IBM)

More information:

<https://nips.cc/Conferences/2016>

ICLR 2017 - 5th International Conference on Learning Representations

Toulon, France, 24-26 April 2017

The performance of machine learning methods is heavily dependent on the choice of data representation (or features) on which they are applied. The rapidly developing field of representation learning is concerned with questions surrounding how we can best learn meaningful and useful representations of data. The conference takes a broad view of the field and includes topics such as deep learning and feature learning, metric learning, compositional modeling, structured prediction, reinforcement learning, and issues regarding large-scale learning and non-convex optimization. The range of domains to which these techniques apply is also very broad, from vision to speech recognition, text understanding, gaming, music, etc. The program will include keynote presentations from invited speakers, oral presentations, and posters.

Topics

A non-exhaustive list of relevant topics:

- Unsupervised, semi-supervised, and supervised representation learning
- Representation learning for planning and reinforcement learning
- Metric learning and kernel learning
- Sparse coding and dimensionality expansion
- Hierarchical models
- Optimization for representation learning
- Learning representations of outputs or states
- Implementation issues, parallelization, software platforms, hardware
- Applications in vision, audio, speech, natural language processing, robotics, neuroscience, or any other field.

Paper submission deadline: 4 November 2016

More information:

<http://www.iclr.cc/doku.php?id=ICLR2017:main>

European Research and Innovation

BASMATI – Cloud Brokerage Across Borders For Mobile Users And Applications

by Patrizio Dazzi (ISTI-CNR)

Cloud Computing and mobile applications are key drivers for innovation. However, mobile device limitations still hinder today's mobile applications from reaching their full potential. The joint South-Korea and EU Horizon 2020 project BASMATI is developing an integrated brokerage platform that targets federated clouds and supports the dynamic needs of mobile applications and users.

Computational clouds and mobile applications have been two of the most relevant drivers of innovation in the last decade. The combination of their particular features has fostered the development of more and more pervasive services with fast roll-out and great scalability, and with low costs for initial assets, thus strengthening industrial competitiveness and promoting economic growth. However, current technological and social landscapes are now calling for a shift towards the introduction of a more flexible, hybrid computing paradigm. Cloud platforms need dynamic, automatic management mechanisms to enable service migration and scalability that is sufficiently efficient and autonomous to cope with the requirements of very large mobile applications serving a mass of nomadic users.

The BASMATI project [L1] aims at providing a fully featured ecosystem able to integrate widely deployed cloud federations along with a range of smaller computing devices, including mobile devices, targeting crowds of users that access their data and applications while moving across national borders. BASMATI will deal with computational and storage resources localized at the edge of the network, addressing challenges related to resource heterogeneity, ultra-scalable provisioning, computation offloading, context- and situation identification, quality of service and privacy / security guarantees.

To achieve these objectives, BASMATI aims at the delivery of an integrated brokerage platform targeting federated clouds with heterogeneous resources and supporting the efficient, cost-effective execution of mobile cloud applications, in a transparent and ubiquitous manner. The design of the BASMATI platform focuses on four main axes: (i) enabling of mobile cloud services, (ii) federation of cloud infrastructures, (iii) scalable infrastructure management, and (iv) brokerage and offloading.

With respect to service enablement, BASMATI plans to support and model user mobility patterns and behaviour, and classify mobile applications in terms of their different static and dynamic characteristics (e.g., structural topology, resource access patterns, performance issues and bottlenecks, to name but a few). Furthermore, BASMATI will study functional and non-functional properties that provide

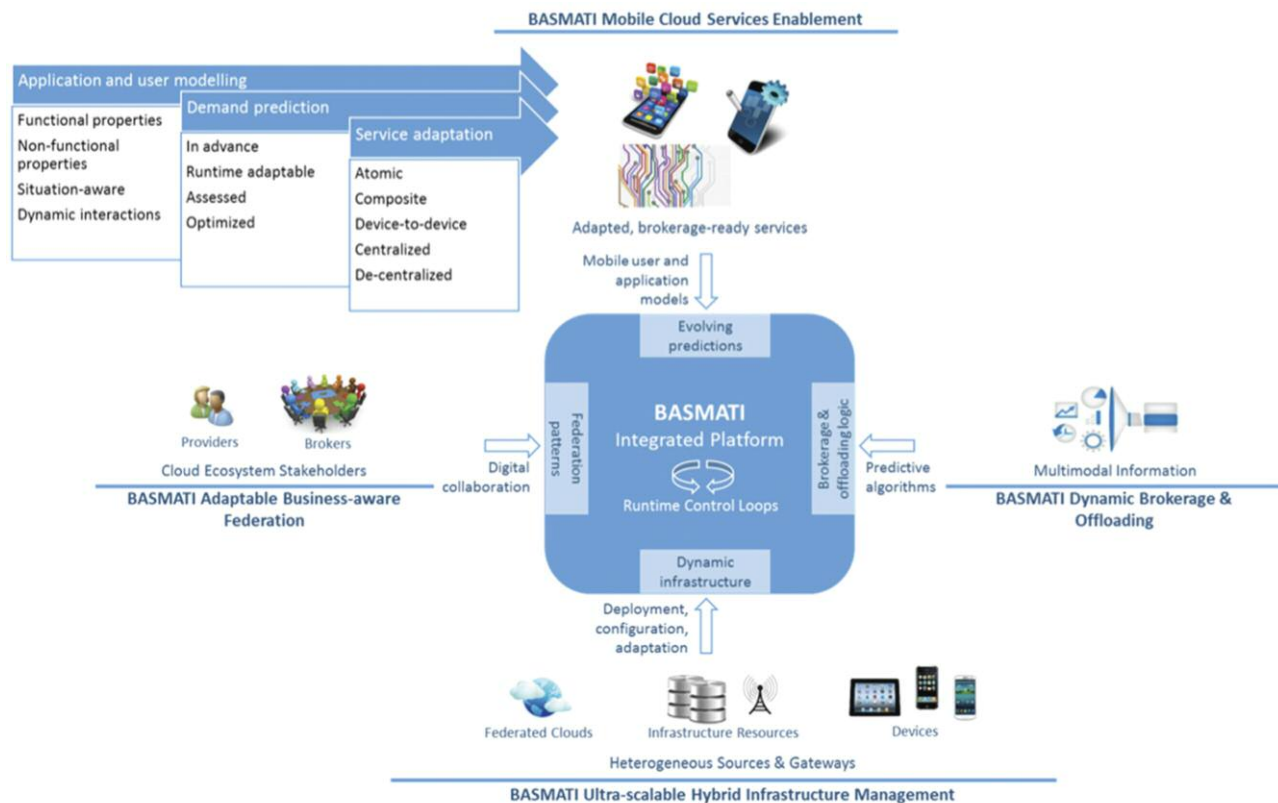


Figure 1: BASMATI platform architecture.

(some degree of) a-priori knowledge of the resource demands, interactions between atomic service components of applications and the situational/contextual aspects that provide additional insight to the mobile user and application space.

The second axis focuses on the cloud federation. The BASMATI infrastructure will blend together existing and new technological solutions into a platform that can spread globally and leverage information from many sources, including complex resource management and business aspects, as well as user and application modeling.

The goal of the third main axis is to facilitate the development of a dynamic infrastructure consisting of different, heterogeneous resources like clouds, data and networking resources, and other kinds of edge devices. Deployment and configuration patterns will be dynamic according to the emerging requirements of mobile cloud applications and users.

The fourth axis focuses on dynamic and automated brokerage and offloading, implementing innovative algorithms (e.g., economic models and machine-learning based solutions). The ultimate aim is to leverage the information from the models to build advanced management algorithms that can drive the efficient exploitation of the Basmati infrastructure despite the many constraints: resource heterogeneity and location, user QoS requirements, user dynamicity and geographic dispersion.

BASMATI is designed to support both existing and new applications, with a focus on three real-world use cases: (1) Large events management, in which BASMATI will support the management of dynamic information about the event

(including user-generated items) in a time- and context-aware fashion; (2) Virtual Mobile Desktop for highly nomadic users, enabling seamless and efficient access to the user environment by reacting to the different conditions occurring when users cross national and international borders; (3) Trip building [L2] of personalized time-budgeted tours of cities, with dynamic definition and detection of points of interests, users communities and visit trajectories, as well as context-aware and personalized pathways to enhance the user experience. The validation of the entire BASMATI platform is planned during the 2018 edition of Das Fest [L3], a music festival that expects more than 250.000 participants over three days.

BASMATI was launched on 1 June, 2016. The project consortium is composed of five European and three Korean partners. On the European side, the project members are: the National Technical University of Athens (NTUA), the overall project coordinator; the National Research Council of Italy (CNR), the scientific and technical coordinator; ATOS Barcelona, CAS Software AG and Amenesik. On the Korean side, the participants are: the Electronics and Telecommunications Research Institute (ETRI), serving as coordinator; Seoul National University and Innogrid.

Links:

- [L1] <http://www.basmati.cloud>
- [L2] <http://tripbuilder.isti.cnr.it/>
- [L3] <http://www.dasfest.de>

Please contact:

Patrizio Dazzi, BASMATI scientific coordinator
 ISTI-CNR, Italy
patrizio.dazzi@isti.cnr.it

An Incident Management Tool for Cloud Provider Chains

by Martin Gilje Jaatun, Christian Frøystad and Inger Anne Tøndel (SINTEF ICT)

The complex provider landscape in cloud computing makes incident handling difficult, as cloud service providers (CSPs) with end-user customers do not necessarily get sufficient information about incidents that occur at upstream CSPs. As part of the FP7 project 'A4Cloud', we have developed an incident management tool that can embed standard representation formats for incidents in notification messages, and a web-based dashboard for handling the incident workflow.

New tools, procedures and guidelines are needed to help cloud service providers be accountable to their customers. An accountable organisation must commit to responsible stewardship of other people's data, which implies that it must define what it does, perform what it defined, monitor how it acts, remedy any discrepancies between the definition of what should occur and what is actually occurring, and finally must explain and justify any and all actions that are per-

formed [1]. Simply put, being accountable means 'doing the right thing'. The objective of the Accountability for Cloud and other Future Internet Services project (A4Cloud) [L1] is to develop tools, guidelines and procedures to make being accountable a business advantage.

The Incident Management Tool (IMT) is a tool targeted at organisations and teams that handle computer security incidents – in practice any organisation that provides or consumes an internet service. The targeted audience of IMT is not the end user, but rather professional incident handlers and privacy officers. The contribution of IMT is a simplified incident format and a simplified incident exchange – making the solution usable for smaller companies as well. A dashboard (see Figure 1) is developed to demonstrate and visualise how such a tool can be useful for IM. Through the integration with the A4Cloud toolset, the incident handler is able to send notifications directly to the affected end users.

A problem experienced by incident handlers in the context of cloud computing, is the lack of access to sufficient incident information throughout the cloud provider chain [2]. Furthermore, complicated cloud provider chains with multiple participants increase the need for more automated sharing of incident information – potentially allowing for automation of response actions, such as notification of availability-related SLA breaches to cloud providers and end users.

The screenshot shows the IMT dashboard interface. On the left is a dark sidebar with navigation options: Dashboard, Incidents (Browse Incidents, Add Incident, Browse Incident Types, Add Incident Types), Subscribers, Providers, and Settings. The main content area is titled 'Incidents Details' and shows an incident titled 'Service 1 was moved to US region'. The incident details include: Origin (DataSpacer), TLP (EHSAC, Amber), Status (Unresolved), Impact (Low-high), Type (Unauthorized government access to data), Language (English), and Description (Due to a system failure in the EU-based data center, some services were moved to the US in order to ensure their operability. Service 1 was one of these services. This is a direct violation of the policy specifying that these services should only be operated from within the EU in order to comply with privacy laws). It also shows 'Occurred at' (Nov. 2, 2015, 9:10 a.m.) and 'Detected at' (Nov. 6, 2015, 6:50 p.m.). Below this is 'Additional Information' with resources ['creditcard'] and users ['hOlson', 'ofordmann', 'jacob', 'tsh123']. The 'Attachments' section states 'This incident has no attachments' with an 'Add attachment' button. On the right side, there is a profile for the 'INCIDENT LEAD' Ola Nordmann, contact details for the 'Liaison' Andrew Smith, a red box indicating 'SUBSCRIBERS HAVE NOT BEEN NOTIFIED' with a 'Notify Subscribers' button, a list of 'End user notifications' with two entries, and 'Actions' buttons for 'Update Incident' and 'Derive Incident'.

Figure 1: IMT dashboard description: In the upper right corner, the current incident handler ('Ola Nordmann') is indicated, with icons indicating active alerts and pending messages. Below that, there is the lead handler for the current incident. Below the lead handler, there are contact details for the liaison at the originating provider. The next box down indicates whether downstream subscribers have been notified of the incident. At the bottom, there are action buttons to derive an incident or update the incident with more information.

IMT operates in the direct context of multiple tools from the A4Cloud toolkit, namely DTMT, AAS and A-PPLE. IMT receives detected data transfer and audit incidents from DTMT and AAS, and utilises A-PPLE to notify end users about incidents that are relevant for them. When a notification of end users is to be performed, IMT sends a notification to A-PPLE, A-PPLE provides this information to Transparency Log (TL), and Data Track fetches this information from TL in order to inform the end user about the incident. IMT could also be used outside the context of A4Cloud tools as a way for organisations to communicate incident information and have this information propagate the cloud service provision chain.

The IMT interacts with other instances of IMT and other tools by a simple, extensible incident format and a publish-subscribe based API. The integration with A4Cloud tools allows for easy notification of end users. The solution supports incidents propagating through the Cloud Service Provision Chain while preserving traceability. The IMT user interface targeting humans consists of a dashboard in which incident handlers and privacy officers can manage subscriptions, incidents and notifications. The notifications can be directed both to other instances of IMT and to A-PPLE instances capable of notifying end users, as appropriate.

In IMT, a human is involved in making the decision on whether or not to notify subscribers and end users. This is because few or no companies would agree to send their incidents directly to the end users or customers upon happening. Thus, the company can decide when to notify their subscribers and end users. A potential problem with this approach could be that the company might decide not to notify about some incidents, but this should be prevented by maintaining an audit trail.

The A4Cloud project has reached its conclusion, but the project partners will continue to develop the various tools in new research opportunities. For IMT, the natural next step would be a large-scale pilot implementation in a real cloud provider chain, and we are currently exploring different options in this regard.

The Accountability for Cloud and other Future Internet Services project (A4Cloud) was led by Hewlett Packard Labs (Bristol, UK), with partners SINTEF, SAP, ATC, Cloud Security Alliance EMEA, Eurecom, Ecole des Mines de Nantes, University of Stavanger, Furtwangen University, Karlstad University, Queen Mary University of London, Tilburg University and University of Malaga. The project started in October 2012, and had its final review in May 2016.

Tools developed in the A4Cloud project:

- Incident Management Tool (IMT): described in this article
- Cloud Offerings Advisory Tool (COAT): an online tool where a prospective cloud customer can select possible cloud providers based on a set of criteria
- Data Protection Impact Assessment Tool (DPIAT): a checklist-based tool to assist a cloud customer in performing a data protection impact assessment in accordance with the EU General Data Protection Regulation

- Data Track tool (DT): a client-based tool for end-users to keep track of what kind of personal data they have disclosed to various cloud providers
- Accountability PrimeLife Policy Language Engine (APPL-E): a server application running at each cloud provider to ensure that accountability and privacy policies are adhered to
- Audit Agent System (AAS): a distributed server tool running at each cloud provider that is capable of performing continuous audit monitoring of various processes
- Data Transfer Monitoring Tool (DTMT): a server tool running at each cloud provider, monitoring data transfers between physical storage locations for possible violations of customer policies regarding data location.

Link:

[L1] <http://a4cloud.eu>

References:

- [1] M. G. Jaatun, et al.: "Enhancing Accountability in the Cloud", to appear in *International Journal of Information Management*, DOI 10.1016/j.ijinfomgt.2016.03.004, 2016
- [2] C. Frøystad, et al.: "Security Incident Information Exchange for Cloud Services", in *Proc. of International Conference on IoT and Big Data, Rome, 2016*

Please contact:

Martin Gilje Jaatun, SINTEF ICT, Norway
+47 900 26 921
Martin.G.Jaatun@sintef.no

Predictive Modelling from Data Streams

by Olivier Parisot and Benoît Otjacques (Luxembourg Institute of Science and Technology)

In order to support knowledge extraction from data streams, we propose a visual platform for quickly identifying main features and for computing predictive models in real time. To this end, we have adapted state-of-the-art algorithms in stream learning and visualisation.

Useful information may be extracted from the high frequency data streams that are common in various domains. For example, textual data from social media such as Twitter and Facebook can be studied to extract the hot topics and anticipate trends. In a different scenario, numerical data collected by a network of environmental sensors can be inspected in order to capture events that could precede potential disaster like floods, storms or pollution peaks.

Therefore, numerous data mining techniques have recently been proposed in order to extract predictive models from data streams [1]. On the one hand, classical analytics techniques can be applied on streams by using a certain pool of observations (by using a sliding window, for example). On the other hand, specific online/incremental methods can be applied to dynamically refresh results. A clever data obsolescence strategy is necessary to consider both significant and up-to-date data windows and allow efficient methods (without accessing too much historical data).

In order to improve stream analytics, we have developed a JAVA platform to inspect data streams on-the-fly and to apply the leading predictive models. Various specific third-parties components can be integrated into the software such as WEKA for static data mining or MOA for specific stream processing.

The platform was designed to support two kinds of data sources:

- Remote streams (i.e., available through web APIs): processed on-the-fly.
- Local streams (i.e., obtained from potentially huge files): iteratively processed in a single-pass, without accessing the previous values.

The user interface was designed to be reactive (by plotting on-the-fly the continuously arriving values) and interactive (by providing a real control to the end-user like play/pause/stop the data stream or select the processing speed).

Additionally, various analytics modules were developed in order to continuously inspect the considered data streams.

First, we have implemented a “Feature similarity” module to extract the meaningful characteristics from data. More precisely, we have designed an innovative real time multidimensional scaling 2D projection dedicated to time series, in order to show the correlations (respectively inverse correlations) for the recent history. As an example, this module could help to determine if the IBM and ORACLE stocks quotes are following the same pattern.

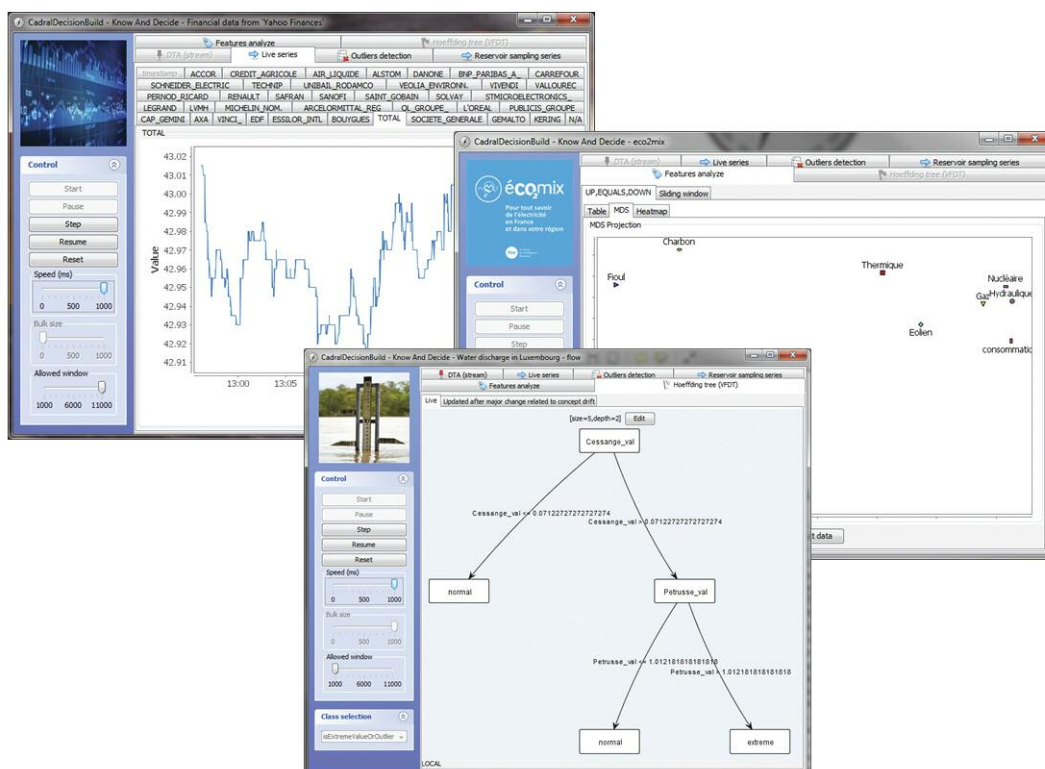


Figure 1: Live visualisation of quotes (Yahoo finance API), feature similarity analysis on the French energy consumption data (RTE – eco2mix) and extreme flooding prediction using hydrological time series from Luxembourg [3].

Second, we have developed a “Predictive modelling” component to create and refresh models, which continuously takes into account recent history. A multitude of techniques exists for predictive analytics, and a critical issue for the data scientist is to select the appropriate technique according to the data characteristics (completion, linear/non-linear relationships, noise, etc.) and the tasks to be carried out. Our aim is to make it easier for the user to understand the predictive models. Therefore, decision trees were selected because they allow a model to be built that is both efficient and easy to interpret. On the one hand, we have applied VFDT, the reference method for classification tree induction. On the other hand, we have used model trees (i.e., decision trees combined to linear regressions) with the recent FIMT-DD algorithm [2] to predict numerical values.

The platform was applied on various real-world data streams (Figure 1). Initially, we tested our approach on the live stocks quotes from the Yahoo Finances website (CAC40 index – one record per second): it helped us to check the “Feature Similarity” module with real life settings. Then, we processed the data from the French electricity transmission system operator (RTE) in order to analyse energy consumption in France (oil, coal, gas, nuclear, wind, solar, bioenergy, hydraulic and pumping – 15-min time series). In this case, we processed heterogeneous values with different scales (for instance: how to use both gas and oil consumption in order to produce meaningful predictions?)

Finally, we inspected hydrological data obtained from the hydrometric stations in Luxembourg: the considered sensor network is composed of 24 stations and continuously produces 15-min time-series [3]. Owing to the poor quality of the data, we had to apply techniques that are robust to noise and missing values in sensor data: for example, the platform was successfully used to fill data gaps in hydrological time series [3].

We plan to extend the software in order to help data scientists quickly identify and eliminate bad data that pollute predictive models. To this end, we are implementing real-time modules for extreme value detection, missing data imputation and live clustering.

Links:

<http://www.list.lu/en/erin/>

<http://www.list.lu/en/erin/news/le-list-effectue-des-recherches-sur-les-precipitations-et-les-crues-extremes/>

References:

[1] H. Nguyen et al.: “A survey on data stream clustering and classification“, Knowledge and Information Systems, Springer, 12/2015

[2] E. Ikonomovska, J. Gama: “Learning model trees from data streams”, Discovery Science, 10/2008

[3] L. Giustarini et al.: “A user-driven case-based reasoning tool for infilling missing values in daily mean river flow records”, Environmental Modelling and Software, 8/2016

Please contact:

Olivier Parisot, Benoît Otjacques

Luxembourg Institute of Science and Technology
olivier.parisot@list.lu, benoit.otjacques@list.lu

Mandola: Monitoring and Detecting Online Hate Speech

by Marios Dikaiakos, George Pallis (University of Cyprus) and Evangelos Markatos (FORTH)

MANDOLA wants to make a bold step towards improving our understanding of the prevalence and spread of online hate speech and towards empowering ordinary citizens and policy makers to monitor and report hate speech.

In recent years an ominous picture about online hate speech has started to materialise within cyberspace. Indeed, recent polls suggest that as many as four out of five respondents have encountered hate speech online and two out of five have personally felt attacked or threatened [L1]. Although it is difficult to get accurate statistics about the spread of hate speech in cyberspace, the picture is becoming increasingly clear: the internet is alarmingly effective at spreading hate speech – so much so that most internet users have encountered it at some point. To make matters worse, hate speech usually targets the most vulnerable groups within society: children, minorities, and immigrants – groups that by definition have little capacity to protect themselves, both in the online and the physical worlds. There are two major difficulties in dealing with online hate speech: (i) Lack of reliable data that can show detailed online hate speech trends. (ii) Poor awareness about how to deal with the issue since there is a fine line between hate speech and freedom of speech: the boundaries between legal and illegal hate speech are sometimes blurred, and may vary between territories.

MANDOLA [L2] plans to contribute towards filling this gap by: (i) monitoring the spread and penetration of online hate-related speech in Europe and in member states using big data approaches; (ii) providing policy makers with actionable information that can be used to promote policies that mitigate the spread of online hate speech; (iii) providing ordinary citizens with useful tools that can help them deal with online hate speech; (iv) transferring best practices among member states; (v) setting up a reporting infrastructure that will connect concerned citizens with the police and appropriate abuse desks and which will enable the reporting of hate-related speech and dangerous speech.

These goals will be achieved by: (i) developing a multi-lingual monitoring dashboard that will offer reliable information about online hate speech enabling users to focus on their geographic region ranging from their city to their country to the entire European Union. This monitoring dashboard will distinguish, if found feasible, between potentially illegal and potentially legal content. The dashboard will use Twitter and Google as sources of possible hate-related online content. Previous works [1, 2, 3] have demonstrated that Twitter can be used to monitor and detect trends in online hate speech; (ii) developing a smartphone app which will (a) spread awareness to users about online hate speech enabling users to understand and isolate such contents, and which will (b)

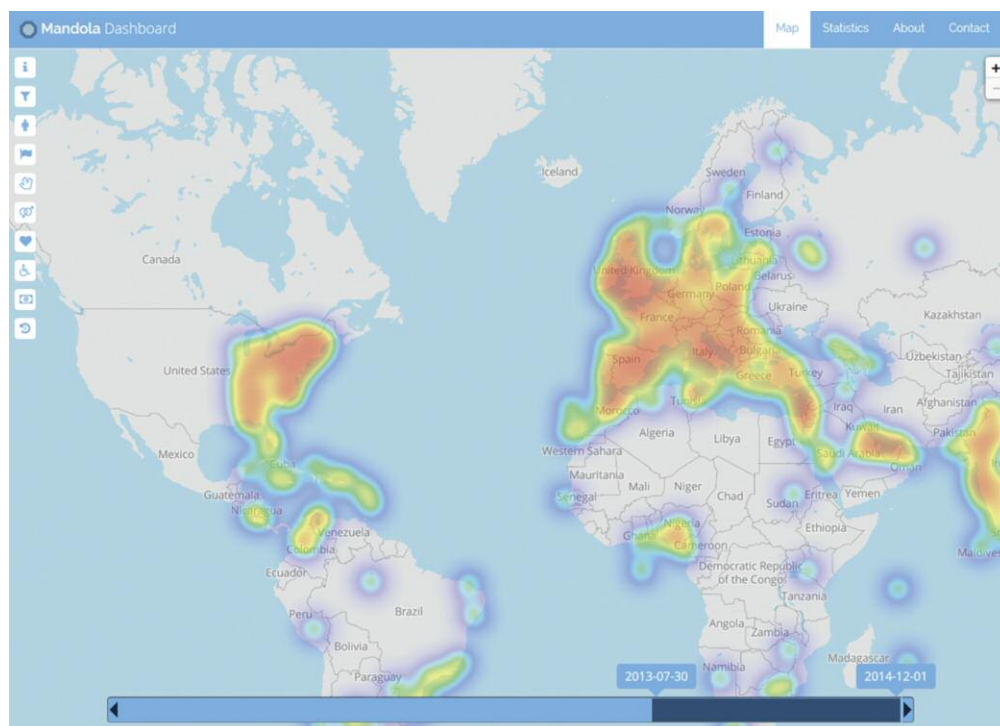


Figure 1: Mandola Monitoring Dashboard.

enable illegal online hate speech to be reported when encountered; (iii) developing a network of National Liaison Officers (NLOs) who will act as the main contact point for their member state; (iv) delivering a Frequently Asked Questions Manual on “Responding to online Hate Speech”; (v) conducting a study of the definition of illegal hatred throughout the European Union, which will enable us to clarify the precise kind of contents to be targeted as well as a study of the legal framework surrounding hate-related speech monitoring and reporting.

The main innovative aspect of MANDOLA is the use of a clear technology-based big data approach to monitoring and reporting online hate speech. This approach is based on the state of the entire visible internet (both Google and Twitter) and can show a clear and representative state of hate speech online. This technology-based approach, which takes care of ethical and legal aspects, gives us a full picture of hate speech on the entire internet allowing us to zoom in to any place and time interval needed for each specific purpose.

The target groups of MANDOLA are:

- Ordinary citizens (i) will have a better understanding of what online hate speech is and how it evolves, (ii) will be able to recognise it when they see it and understand when freedom of speech crosses the boundary into illegal hate speech and (iii) will know how to cope with illegal online hate speech when they encounter it.
- Policy makers will have online hate speech-related information at their fingertips – actionable information that can be used to make decisions.
- Witnesses of online hate speech incidents will have an opportunity to report hate speech anonymously.

MANDOLA is a two year project (started 1st October 2015) and is co-funded by the Rights, Equality and Citizenship (REC) Programme of the European Commission. MAN-

DOLA consortium consists of seven partners from six countries: FORTH (Foundation for Research and Technology – Hellas), Aconite Internet Solutions (Ireland), the International Cyber Investigation Training Academy (Bulgaria), Inthemis (France), the Autonomous University of Madrid (Spain), the University of Cyprus (Cyprus) and the University of Montpellier (France). The project is led by FORTH. This article was written with the financial support of the Rights Equality and Citizenship (REC) programme of the European Union. The contents of this publication are the sole responsibility of the authors and can in no way be taken to reflect the views of the European Commission.

Links:

- [L1] <http://eeagrants.org/News/2012/Countering-hate-speech-online>
 [L2] Mandola Project <http://mandola-project.eu/>

References:

- [1] P. Burnap et al.: “Detecting tension in online communities with computational Twitter analysis”, *Technol Forecast Soc Change* 95:96-108, 2015.
 [2] P. Burnap, M. L. Williams: “Us and Them: Identifying Cyber Hate on Twitter across Multiple Protected Characteristics”, *EPJ Data Science*, 5:11, 2016, <http://link.springer.com/article/10.1140/epjds/s13688-016-0072-6>.
 [3] L. Silva, et al.: “Analyzing the Targets of Hate in Online Social Media”, in *Proc. of the International AAAI Conference on Weblogs and Social (ICWSM’16)*, Cologne, Germany. 2016.

Please contact:

George Pallis, University of Cyprus
gpallis@cs.ucy.ac.cy
<http://www.cs.ucy.ac.cy/~gpallis/>

The BÆSE Testbed – Analytic Evaluation of IT Security Tools in Specified Network Environments

by Markus Wurzenberger and Florian Skopik (AIT Austrian Institute of Technology)

Recent years have seen a dramatic increase in damage caused by cyber-criminals. Although there are many IT security tools on the market, there is currently no way to test, compare and evaluate them without actually running them in real systems. The BÆSE testbed offers a solution to challenge and benchmark IT security solutions for dedicated network environments under attack conditions, without putting real systems in danger.

The demand for interconnected digital services has been growing rapidly in recent years. With the introduction of low-cost devices, which are even used in industrial environments now, the emergence of cyber physical systems (CPS) spanning widely spread network components has begun. While this broad connectivity improves society's productivity and optimises industrial production through automation, it makes organisations and private life vulnerable to cyber-attacks. Modern attacks are sophisticated, tailored to specific purposes and often use customised tools. Consequently, prompt detection is difficult and victims are often oblivious to the attack. Attacks typically result in data breaches, high financial losses (actual reports estimate \$500 billion per year [L1]) and damage to reputation.

Actual Situation and Problem Statement

There are various types of IT security tools in use, including: anomaly detection systems, intrusion detection systems (IDS), antivirus scanner and security information event management (SIEM) tools. These systems aim to promptly detect attacks, but they are often inadequate when it comes to

recognising sophisticated and tailored intrusions. While the market for security is growing (projected to rise from \$75 billion in 2015 to \$170 billion in 2020 [L2]) and vendors are entering the market with new products in ever shorter cycles, the financial losses are still growing.

While many great IT security solutions have been developed, the problem of how to rate, compare and evaluate them to facilitate their optimal configuration and application in a specific organisational context remains unsolved [1]. This is because there is no effective way of testing solutions prior to their deployment. Challenging IT security solutions with attacks under realistic circumstances is mandatory to rate their detection capabilities. However, this needs to be performed in the target environment – or at least in an environment which simulates the target environment as closely as possible to rate the individual detection capabilities, since the structure of a network, and the way individuals use it, differs from one organisation to another. Unfortunately, this is a highly non-trivial task.

Goals and Innovations

In the BÆSE (Benchmarking and Analytic Evaluation of IDSs in Specified Environments) project, we invent the BÆSE testbed, which allows vendor-independent evaluation and comparison of IT security tools for user-specified environments. There are two steps to the BÆSE testbed approach:

- (i) Generate semi-synthetic test data based on the properties and characteristics of a specified network environment (see Figure 1).
- (ii) Feed the generated test data into various detective IT security tools for comparing, rating and evaluating their capabilities with respect to different configurations using the BÆSE Testbed (see Figure 2).

Part (i) of the model takes a small part of captured logging, netflow, or packet data from a real network environment as input. This input is processed with machine learning methods to obtain the properties and characteristics of the considered network. For example in the case of log data, clustering can be used to group the log lines. Based on the results of the machine learning algorithms, stable and variable parts of the

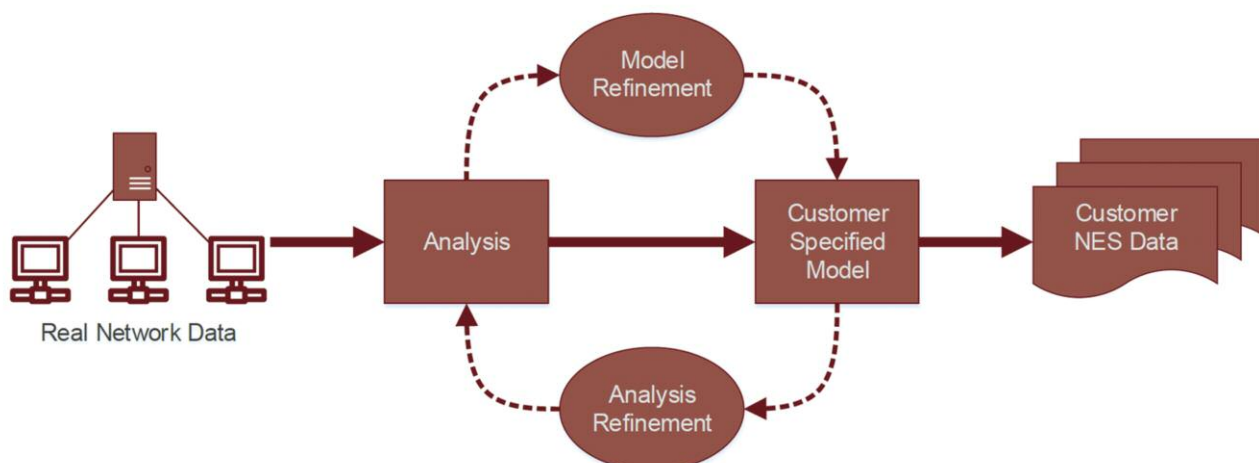


Figure 1: NES data generation process.

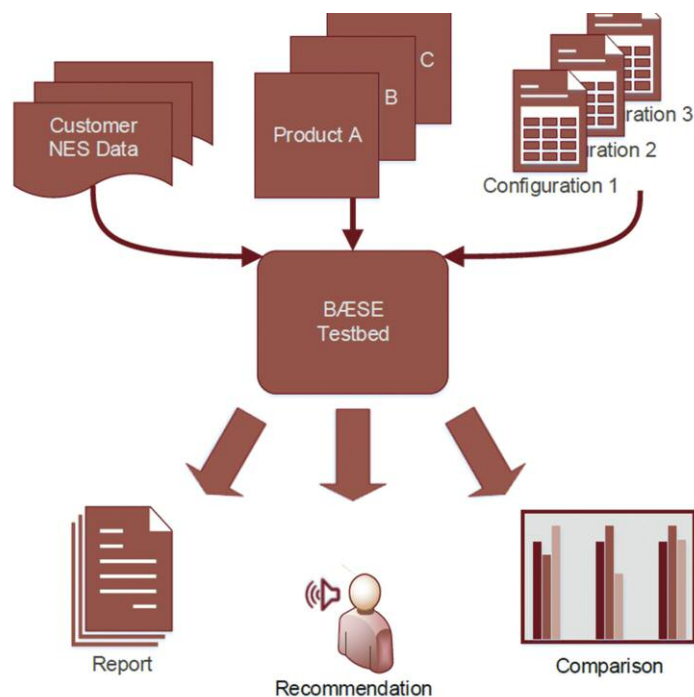


Figure 2: BÆSE testbed.

data within a cluster are determined. Furthermore, the transition probabilities between log lines representing various sorts of events are stored as well as the temporal distribution of these events. This information then allows the application of a Markov chain simulation to generate a semi-synthetic test data file of any size, which in Figure 1 is named Network Event Sequence (NES) since our approach is not limited to log data and can also be used, for example, with netflow and packet data. Iterative and interactive refinement of the analysis and the model allows NES data to be generated with varying degrees of detail. A detailed description of the approach can be found in [2].

Part (ii) describes the main building block of the BÆSE testbed. The testbed takes as input different sets of NES data, and a choice of IT security tools, which are evaluated with different configurations. The BÆSE testbed then compares and rates the tools and configurations, which are considered for evaluation. The outputs of the BÆSE testbed are reports, recommendation and statistics. Furthermore we plan to develop a concept to simulate realistic attacks in the NES data to stress the IT security tools in an appropriate way. This will allow an evaluation of security solutions under most realistic circumstances for a specified network environment, without exposing the real network infrastructure to risky situations. The objectives of BÆSE are to:

- find the optimal security solution
- optimise the usage of security tools
- find the most efficient configurations of security tools
- raise the detection capability of security tools.

The BÆSE Project

The BÆSE project is financially supported by the Austrian Research Promotion Agency FFG under grant number 852301, and carried out in the course of an industry-related PhD thesis. Project stakeholders are the Austrian Institute of Technology (AIT) as coordinator, the Vienna University of

Technologies as academic partner, and T-Systems Austria. AIT has extensive experience in anomaly detection and is developing its own IDS named Automatic Event Correlation for Incident Detection (AECID).

Links:

- [L1] <http://www.forbes.com/sites/stevemorgan/2016/01/17/cyber-crime-costs-projected-to-reach-2-trillion-by-2019/#16549f753bb0>
- [L2] <http://cybersecurityventures.com/cybersecurity-market-report/>

References:

- [1] ISO/IEC 27039, Information Technology – Security techniques – Selection, deployment and operations of intrusion detection systems.
- [2] M. Wurzenberger, F. Skopik, G. Settanni, and W. Scherrer, “Complex log file synthesis for rapid sandbox-benchmarking of security- and computer network analysis tools,” *Information Systems*, vol. 60, pp. 13–33, Aug. 2016.

Please contact:

Markus Wurzenberger
AIT Austrian Institute of Technology, Austria
+43 664 8157937
markus.wurzenberger@ait.ac.at

Florian Skopik

AIT Austrian Institute of Technology, Austria
+43 664 8251495
florian.skopik@ait.ac.at

Behaviour-Based Security for Cyber-Physical Systems

by Dimitrios Serpanos (University of Patras and ISI), Howard Shrobe (CSAIL/MIT) and Muhammad Taimoor Khan (University of Klagenfurt)

Behaviour-based security enables industrial control systems to address a wider range of risks than classical IT security, providing an integrated approach to detecting security attacks, dependability failures and violation of safety properties.

Cyber-physical systems are increasingly employed for the control and management of processes, ranging from cars and traffic lights to medical devices; from industrial floors to nuclear power generation and distribution. Importantly, industrial control systems (ICS), a large class of cyber-physical systems composed of specialised industrial computers and networks, are employed to control and manage critical national infrastructure, such as power, transportation and health infrastructure. ICS are now vital to the welfare of nations and individuals, and consequently these systems are increasingly becoming targets of cyber-attacks as evidenced by a number of recent sophisticated incidents, such as Stuxnet and the hacking of cars and health devices.

Industrial control systems differ from traditional IT systems in several ways, from their purpose to ownership and maintenance; they support processes instead of people, they interface with physical systems, they are typically owned by engineers and their requirements for real-time continuous operation necessitate different approaches to upgrade and maintenance. These characteristics have led them to be known as OT (Operational Technology), instead of IT, and they are often targeted by novel attacks that are different from the traditional IT attacks. For example, in addition to the typical

computational or networking attacks, where malicious software is inserted into a system or transmitted to it through a network, a new attack that so far has only inflicted OT systems is ‘false data injection’ (FDI): instead of attacking the computers or the networks of sensor-based systems, an attacker attacks the sensors and inserts false data (measurements) in order to lead the system to a wrong decision, although no malicious software is running and network packets are not manipulated in any way. FDI attacks can be quite powerful and can lead to catastrophic results as simple experiments easily demonstrate (for water networks, power networks, etc.). A second class of novel attacks involve timing disruptions in which the control system is prevented from issuing new commands within the time constant of the system under control; these attacks can be achieved through subtle disruption of the networks that carry sensor data to the controller or by the introduction of parasitic computations on the host whose only goal is to slow down the controller.

Traditional computer and network security has been addressing the problem of fortifying systems and networks for IT, but has been quite limited in addressing systems with real-time requirements and has not considered FDI or timing disruption attacks. The traditional approach to detect malicious processes and/or traffic is through two basic methods: static and dynamic. Static analyses focus on static characteristics of processes and packets, such as signatures, while dynamic ones focus on process and traffic behaviour. Traditionally, behaviour is defined as patterns of resource usage, such as memory and i/o in a computing system, network ports, source and destination addresses, connections, etc. This approach leads to pattern-based behaviour definition, where several directions exist: one can define bad behaviour patterns, monitor and try to detect when one (or more) occur, or define good behaviour patterns, monitor and try to detect deviation from them.

We follow a different approach to define behaviour and build secure industrial control systems. We define behavior as the (executable) specification of an application. Based on this

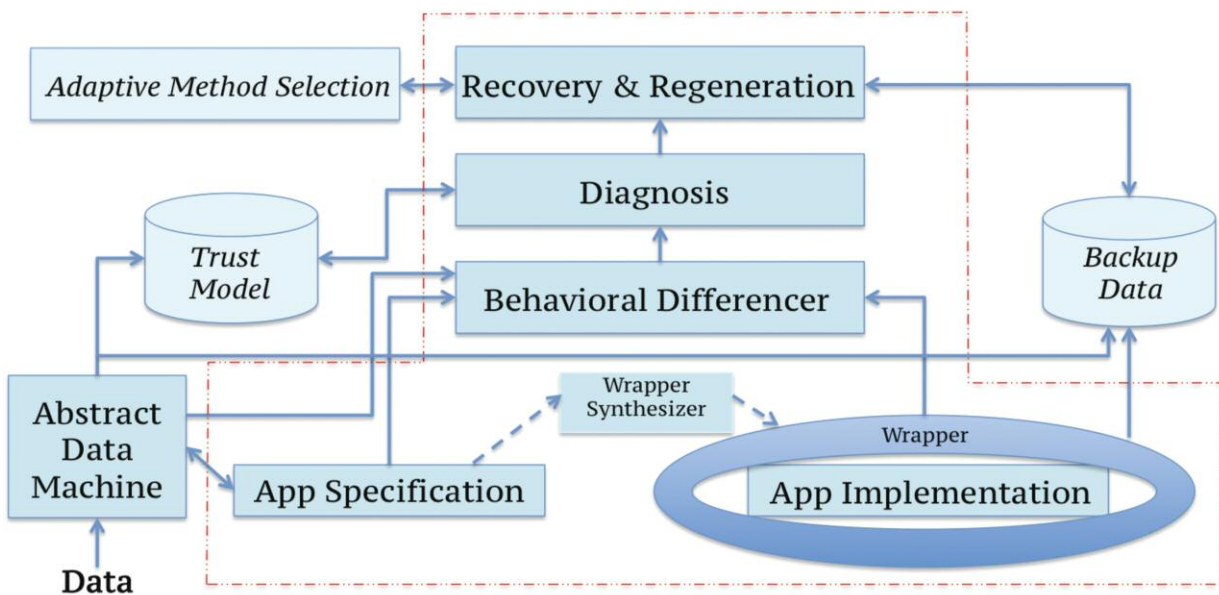


Figure 1: ARMET, a middleware system for ICS.

definition, we develop ICS that execute the (executable) specification in parallel with the application code, we monitor the application code execution and identify deviations between specification execution, which produces predictions, and the application code execution, which produces observations. Importantly, this approach to behaviour definition enables systems to detect all types of behavioural deviations from the specification, independently of motive, i.e., both malicious and accidental, integrating security with fault tolerance in the same approach.

Exploiting this approach, we have been developing ARMET, a middleware system for ICS shown in Figure 1, working in three main security research directions:

- Safe code derivation from specifications, based on the FIAT approach [1], in order to develop safe application code (App Implementation in Figure 1) with specified security properties from the application specification (App Specification)
- Monitor development (Behavioural Differencer) that accurately detects (without false positives or negatives) deviations between application specification and execution [2], and
- Methods to identify vulnerabilities for false data injection attacks and encode them in the Abstract Data Machine, so that the monitor can protect against them. Our work has already provided promising results for power (smartgrid) systems [3].

This approach is practical today, capitalising on the significant recent advances in software verification and formal methods that enable analyses of large programs, and especially in cyber-physical and ICS, which implement specific processes or applications ('plants' in control system terminology). Advantages include automated development of safe code and design and implementation of robust monitors that can identify when security properties are violated.

References:

- [1] B. Delaware, et al.: "Fiat: Deductive Synthesis of Abstract Data Types in a Proof Assistant", in Proc. of POPL'15. Jan. 2015.
- [2] M.T. Khan, D. Serpanos, H. Shrobe: "Sound and Complete Runtime Security Monitor for Application Software", arXiv:1601.04263 [cs.CR].
- [3] S. Gao et al.: "Automated Vulnerability Analysis of AC State Estimation under Constrained False Data Injection in Electric Power Systems", in Proc. of IEEE CDC'15. Dec. 2015.

Please contact:

Dimitrios Serpanos,
University of Patras and ISI, Greece,
Tel: +30 261 091 0299
serpanos@isi.gr

The TISRIM-Telco Toolset – An IT Regulatory Framework to Support Security Compliance in the Telecommunications Sector

by Nicolas Mayer, Jocelyn Aubert, Hervé Cholez, Eric Grandry and Eric Dubois

The objective of our project is to adapt and facilitate Information System (IS) security risk management in the telecommunications sector. To this end, we have developed: first, a model-based approach and a tool to support the adoption of IS security risk management by Luxembourg's telecommunications service providers (TSPs); and second, a framework to analyse the data collected by Luxembourg's National Regulation Authority (NRA).

There is currently a strong emphasis on the security of information systems (IS) and the management of information security risks. Numerous regulations are emerging that impose a risk-based approach for IS security on entire economic sectors. In the telecommunications sector, the EU Directive 2009/140/EC introduces Article 13a about security and integrity of networks and services. This article states that member states shall ensure that providers of public communications networks 'take appropriate technical and organisational measures to appropriately manage the risks posed to security of networks and services'. To harmonise the implementation of this at a national level, the European Network and Information Security Agency (ENISA), as the centre of network and information security expertise for the European Union, published in December 2011 a document entitled "Technical Guideline for Minimum Security Measures" [L1].

As part of the adoption of this directive at the national level in Luxembourg, we have developed a project that aims to adapt and facilitate IS security risk management in the telecommunications sector. To this end, the project is composed of two parts. In the first part we have developed a model-based approach and a tool to support the adoption of this regulation by telecommunications service providers (TSPs) at the national level [1]. The second part involves developing a framework to analyse the data collected by the NRA through this standard approach [2].

For the first part of this project, the starting point of our analysis is that the different TSPs in Luxembourg have very different levels of expertise in security risk management. Thus, letting them report to the NRA without strong guidance would have resulted in very different types of reports, with various quality levels. In order to build a harmonised reporting approach and to meet the needs of the users' (i.e., TSPs in Luxembourg), we decided to define both the methodology and its associated tool in collaboration with the TSPs. Furthermore, we established shared business and architecture models supporting the methodology. Regarding the definition

of these sector-specific models, the first task consisted of defining the different processes composing each regulated telecommunications service. Process reference models such as Business Process Framework (“eTOM”) of the TMForum [L2] or the Telecommunications Process Classification Framework of the American Productivity and Quality Center [L3] were used as input. Then the second task was to describe the IS supporting each telecommunications service. The works of The Open Group and TMForum have been specifically analysed and confronted with the state of practice of the national TSPs. Finally, we defined for each telecommunications service the most relevant threats and vulnerabilities, based on the reference IS architecture

previously defined, and the most relevant impacts, based on the business processes previously defined. Finally, we have represented the resulting knowledge in ArchiMate [L4]: an Enterprise Architecture modelling language. We have extended ArchiMate with the appropriate concepts from the risk management domain [3]. We then integrated all of the different models into a software tool. This task was performed by adapting TISRIM, a risk management tool developed in-house, that was initially released in 2009. TISRIM is currently the tool recommended to the TSPs by our national NRA to comply with the regulation.

After having defined and implemented a method to support the adoption of the regulation by TSPs, there was also a strong need to develop a platform in order to manage the reports received annually by the NRA, and to be able to efficiently analyse their contents. The purpose was therefore to define a set of measurements depicting the trust the NRA can have in the security of TSPs, as well as in the whole telecommunications sector. The outcome for the NRA is to be able to provide recommendations to the TSPs and to facilitate policy-making. The first task when defining the measurement framework was to establish a template for the measurement constructs, inspired by the state of the art, and in particular the template proposed in ISO/IEC 27004. Then, once the measurement template was established, two types of measurements were defined: compliance measurements, measuring the compliance to requirements imposed by legislation and performance measurements, measuring the effectiveness of IS security. The final set obtained is composed of 10 measurements defined for TSPs and 11 measurements defined for the whole telecommunications sector. Finally, the measurements were implemented in a tool named TISRIMonitor, which is currently under evaluation by the NRA.

Our objective is now to extend this approach to other critical and regulated sectors, such as the health and finance sectors, or the privacy regulator. All of these approaches will be man-

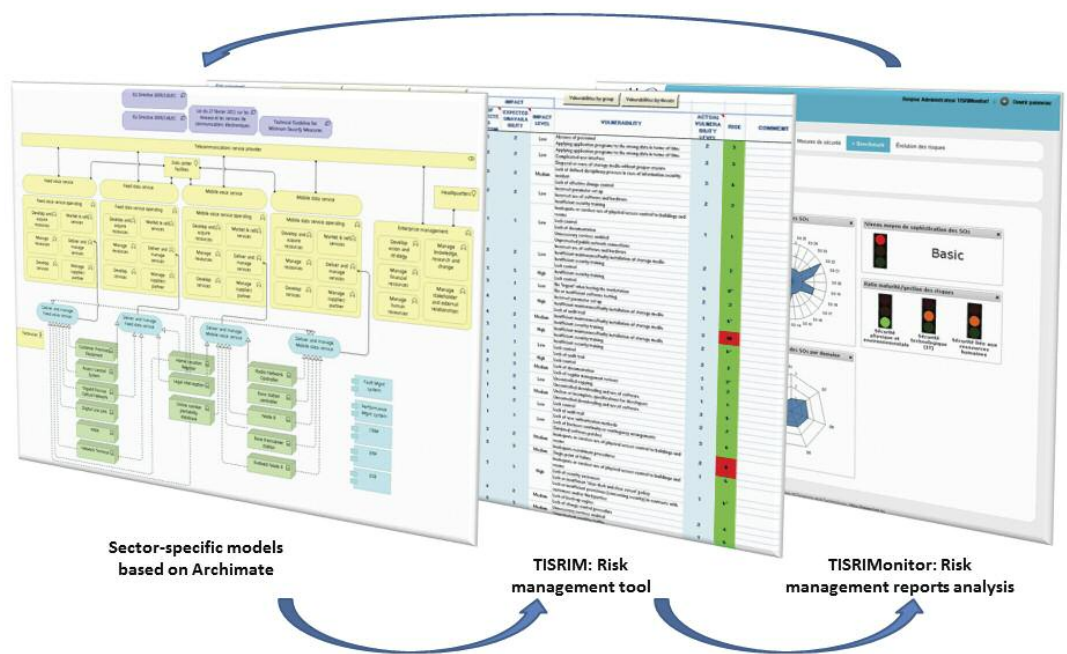


Figure 1: Overview of the TISRIM-Telco toolset.

aged and linked in a regulatory framework called ‘RegTech platform’ – a technological platform to conduct regulation activities. Another follow-up of this project is to adopt a more holistic approach. The objective is to extend the scope of risk management to networked enterprises, and to lead to systemic risk management, i.e., risk management on the entire system formed by networked enterprises, to avoid perturbations of the ecosystem due to local, individual decision-making.

Links:

- [L1] <https://www.enisa.europa.eu/topics/incident-reporting-for-telcos/guidelines/technical-guideline-on-minimum-security-measures>
- [L2] <https://www.tmforum.org/business-process-framework/>
- [L3] <https://www.apqc.org/knowledge-base/documents/apqc-process-classification-framework-pcf-telecommunications-pdf-version-50>
- [L4] <http://www.opengroup.org/subjectareas/enterprise/archimate>

References:

- [1] N. Mayer, J. Aubert, H. Cholez, E. Grandry: “Sector-Based Improvement of the Information Security Risk Management Process in the Context of Telecommunications Regulation”, EuroSPI 2013.
- [2] Y. Le Bray, N. Mayer, J. Aubert: “Defining Measurements for Analyzing Information Security Risk Reports in the Telecommunications Sector”, SAC 2016.
- [3] E. Grandry, C. Feltus, and E. Dubois: “Conceptual Integration of Enterprise Architecture Management and Security Risk Management”, EDOCW 2013.

Please contact:

Nicolas Mayer
 Luxembourg Institute of Science and Technology (LIST)
 Tel: +352 275 888 1, Nicolas.Mayer@list.lu

Predicting the Extremely Low Frequency Magnetic Field Radiation Emitted from Laptops: A New Approach to Laptop Design

by Darko Brodić, Dejan Tanikić (University of Belgrade), and Alessia Amelio (University of Calabria)

Known laptop characteristics can be used to create a model that predicts the extremely low frequency magnetic field radiation emitted at the top and underside of laptops.

Owing to their user-friendly characteristics, particularly their portability and the option of being powered by current or battery, laptops play an important role in many individuals' lives, with many users becoming quite dependant on them. Typically, when working with a laptop, the user's body is in constant close contact with the body of the device; either the laptop is on a desk or on the knees of the user. Either way, the laptop is implicitly in contact with the skin, the lymph and the bones of the user.

Because of the current flowing through the laptop's electronic components, an extremely low frequency (ELF) magnetic field radiation up to 300 Hz is generated at the top and at the bottom parts of the laptop, potentially posing a signifi-

cant risk to the user's health. In particular, the electric current densities up to 483% higher than the reference level of the International Commission for the Non-ionized Radiation Protection (ICNIRP) are caused by the laptop's power supply [1]. Also, different studies have analysed the correlation between exposure to ELF magnetic field radiation and the occurrence of serious illnesses, like leukemia and brain cancer.

To eliminate the high risk to users that is associated with everyday ELF magnetic field exposure, we propose a new model for predicting the ELF magnetic field emission from laptops, based on known laptop characteristics during their normal working conditions including popular office programs like Word, Excel and Internet browsing [2]. The model is based on artificial neural networks (ANN), which are software systems inspired by the human central nervous system. In this model, artificial neurons are connected to each other to create the network. The input of the network is a set of three parameters: (i) Passmark, which is the measure that estimates the processor calculation power, (ii) CPU total dissipation (CPU TD) which represents the maximum dissipation that the CPU achieves to process instructions and data, and (iii) Maximum power consumption (MPC), which is the maximum consumption of the laptop in its normal use as well as for charging its battery. The output of the network is the ELF magnetic field radiation that the laptop is likely to emit at its nine positions at the top and nine positions at the bottom. Figure 1 shows the typical measuring positions on a laptop.

The model is based on the high correlation between three input parameters and the ELF magnetic field emission pro-

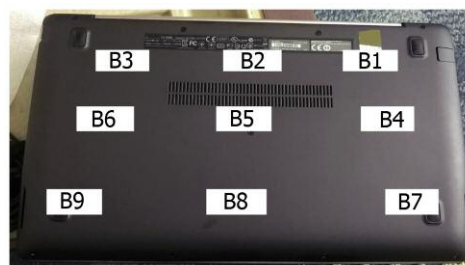
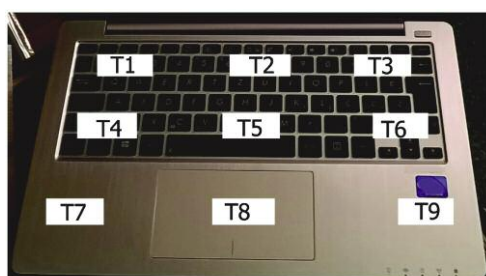


Figure 1: Typical positions at the top (left) and bottom (right) of the laptop.

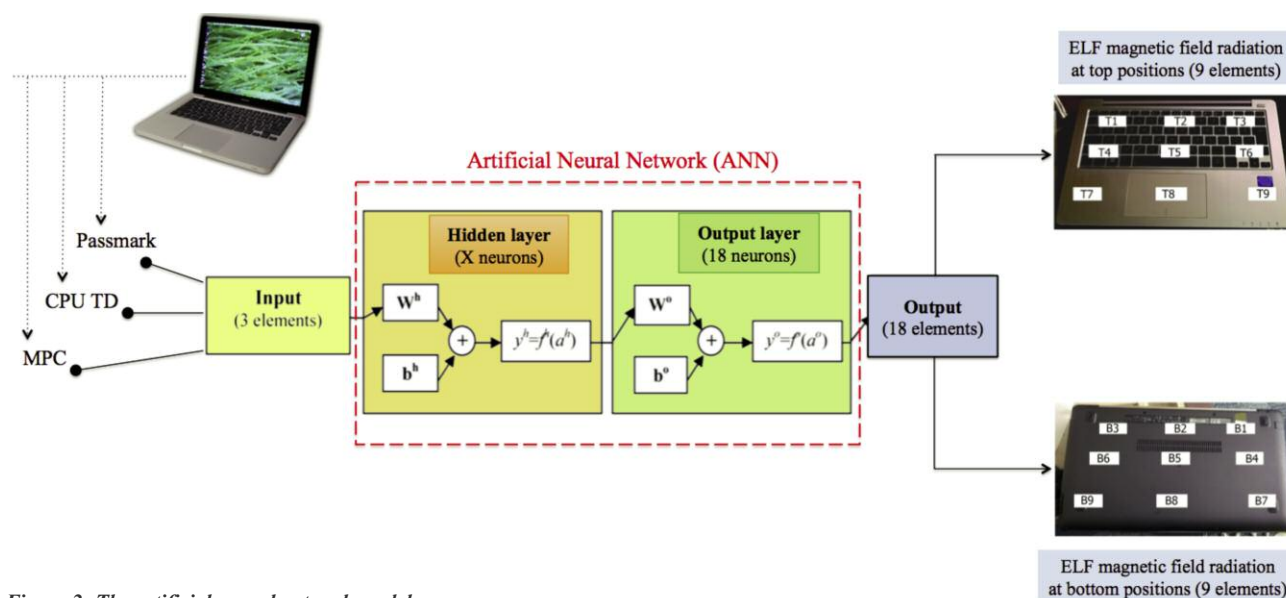


Figure 2: The artificial neural network model.

duced by the laptop, which has been formally demonstrated. Then, we classify the obtained ELF magnetic field emission into three or four dangerousness classes based on the possible effects of the emission on the human body, according to the reference limit extracted from the current safety standards [3]. Figure 2 illustrates the ANN model.

The proposed model has been tested on 10 laptops of different manufacturers, each characterised by its Passmark, CPU TD and MPC values. First of all, the ELF magnetic field radiation has been measured at the nine positions at the top and at the bottom parts of the laptops. Then, the model has been employed to predict the ELF magnetic field values of the laptops at all 18 positions, given the Passmark, CPU TD and MPC values as input of the model. Comparison between the results obtained from the model and real measured values also demonstrates the accuracy and the efficacy of the model in predicting the highest peaks in magnetic field emission. One of the strengths of the model is that it predicts with great accuracy the dangerousness classes associated with a particular laptop.

In conclusion, this model can play an important role in predicting the distribution, the emission levels and the dangerous levels associated with a given laptop in order to suggest safety rules for working with the device. The model could also provide useful information to manufacturers, helping them to prevent ELF magnetic field laptop emissions by tailoring the inner components that influence laptop characteristics. Consequently, the model could favour the design of low emission laptops.

This study is part of ongoing research involving the Technical Faculty in Bor, University of Belgrade, Serbia and the Department of Computer Science Engineering, Modelling, Electronics and Systems, University of Calabria, Italy. Future work will investigate the efficacy of the model when a laptop operates under conditions of stress.

This work was partially supported by the Ministry of Education, Science and Technological Development of the Republic of Serbia TR33037.

References:

- [1] C. V. Bellieni, et al.: "Exposure to Electromagnetic Fields From Laptop Use of "Laptop" Computers", *Archives of Environmental & Occupational Health*, 67(1):31-36, 2012.
- [2] D. Brodić, D. Tanikić, A. Amelio: "An approach to evaluation of the extremely low-frequency magnetic field radiation in the laptop computer neighborhood by artificial neural networks", *Neural Computing and Applications*, Springer, 1-13, 2016.
- [3] D. Brodić, A. Amelio: "Detecting of the Extremely Low Frequency Magnetic Field Ranges for Laptop in Normal Operating Condition or Under Stress", *Measurement*, Elsevier, 91:318-341, 2016.

Please contact:

Darko Brodić
Technical Faculty in Bor, University of Belgrade, Serbia
dbrodic@tf.bor.ac.rs

Managing Security in Distributed Computing: Self-Protective Multi-Cloud Applications

by Erkuden Rios (Tecnalia), Massimiliano Rak (Second University of Naples) and Samuel Olaiya Afolaranmi (Tampere University of Technology)

MUSA (Multi-cloud Secure Applications) is an EU H2020 funded research project which is aimed at ensuring security in multi-cloud environments. The main goal of MUSA is to support the lifecycle of applications with strict security requirements over heterogenous cloud resources. MUSA will result in a security framework that includes security-by-design mechanisms as well as runtime security monitoring and enforcement to mitigate security incidents.

Multi-cloud applications rely on the adoption of cloud services of different capability types (i.e. infrastructure, platform or software as a service) from different Cloud Service Providers (CSPs). Multi-cloud follows the concept of distributed computing in which the components are dispersed but communicate in an integrated manner to achieve the desired goal. This model offers the opportunity to select the best CSPs that satisfy both application and component level requirements. However, the distributed model makes security management even more complex as the need arises to tackle it at different levels: individual components, component-to-component communication and overall application. This calls for approaching security in a holistic manner. MUSA aims to address this need by providing the MUSA framework which considers security throughout the multi-cloud application lifecycle (i.e design, deployment and runtime) relying on security-by-design and integrated security assurance to allow application self-protection at runtime.

MUSA framework

The MUSA framework offers methods and tools to support the integration of the security within the multi-cloud application lifecycle phases, as follows:

- Design phase: the MUSA IDE, which helps in both specifying the end user security requirements and integrating such requirements in the application design. The IDE includes two main tools; the SLA Generator for the creation of the needed Security Service Level Agreements (SLAs) (see details below) and the Modeller which allows the creation of the architecture model of the application, i.e. the specification of the multi-cloud application requirements with respect to component interfaces, cloud deployment needs, etc. The MUSA IDE will allow embedding security agents in the application components for self-protection, i.e. they will enable the activation of security monitors and controls at runtime without modifying the programming model.
- Deployment phase: MUSA offers a Decision Support tool and a Distributed Deployment tool, helping in the choice of the CSPs to use (according to not only their functional

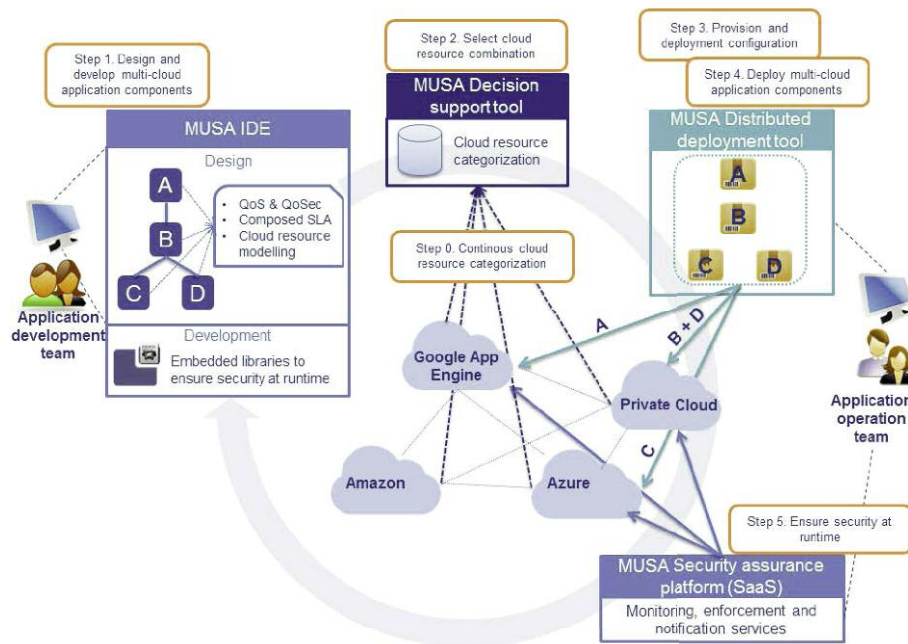


Figure 1: MUSA framework.

but also security features) and in the deployment of the multi-cloud application, respectively.

- Execution phase: The MUSA Security Assurance Platform (provided in form of a SaaS) supports the monitoring, notification and enforcement of correction actions to grant the security features in the Security SLA. The MUSA monitoring supports the collection of operation and security metrics of both the components of the application and the cloud resources provisioned. The approach relies on the use of standard APIs (when they are used by the CSPs), cloud interoperability frameworks such as jclouds, or measures provided by MUSA security embedded libraries.

SLA-driven Security

In order to consider the security features of the overall multi-cloud application, MUSA framework adopts the concept of Security SLA (i.e. contract between customer and provider that states the security terms granted to each other). The MUSA framework proposes a tool, SLA generator, which can be used to identify the Security SLA that each application component must grant.

The SLA Generation relies on a simplified risk analysis process that enable developers to identify major threats to the components and, according to them, build up the Security SLA describing countermeasures in terms of security controls (according to standards like NIST SP-800-53 [1] or to frameworks like CSA Cloud Control Matrix [2]) as well as offering Service Level Objectives, expressed with respect to measurable security metrics that demonstrate the correct application of the offered security controls. The MUSA framework enriches the process with the decision support tools to evaluate the services offered by real CSPs and to outline the feasibility of the Security SLA, suggesting development improvements in order to satisfy the security requirements. The final result of this process will be a multi-cloud application enriched with a set of Security SLAs granted to application components and application customers. The MUSA security assurance platform will provide application components monitoring and apply corrective actions needed to respect the agreed SLA.

Application validation and conclusion

In order to demonstrate the MUSA framework feasibility and effectiveness, two case studies are being implemented:

- NetLine/Sched flight scheduling application by Lufthansa Systems Germany. MUSA will support data integrity, confidentiality, localization and access control in this multi-cloud application which is used nowadays by 55 airlines around the world.
- Smart mobility services by Tampere University of Technology Finland. This open data based multi-cloud application optimizes urban travel experience in Tampere city. MUSA will facilitate the design and deployment of the needed privacy and protection for citizen's mobility data.

The initial validation of MUSA is planned for the end of this year and will serve to improve the framework tools and their integration towards fully fulfilling customers' requirements. The MUSA Project started in January 2015 and will run until December 2017. It receives funding from the EU's H2020 Research and Innovation programme under grant agreement No 644429. The project is coordinated by Fundación TECNALIA Research & Innovation (Spain). The MUSA consortium consists of academia and industry partners from six countries: Spain, Finland, Italy, England, France and Germany.

Links:

<http://musa-project.eu/>

References:

- [1] NIST Special Publication 800-53 Revision 4. Available at: <http://nvlpubs.nist.gov/nistpubs/SpecialPublications/NIST.SP.800-53r4.pdf>
- [2] The CSA Cloud Control Matrix v3.0.1, Available at: <https://cloudsecurityalliance.org/download/cloud-controls-matrix-v3-0-1/>

Please contact:

Erkuden Rios, TecNALIA (Project Coordinator)
erkuden.rios@tecnalia.com



CALL FOR PAPERS

VaMoS 2017: 11th International Workshop on Variability Modelling of Software-intensive Systems

Eindhoven, The Netherlands, 1-3 February 2017

The VaMoS workshop series aims to bring together researchers from different areas dedicated to mastering variability in order to discuss advantages, drawbacks, and complementarities of various approaches, and to present new results for mastering variability throughout the life cycle of systems, system families, and (software) product lines.

VaMoS 2017 particularly welcomes approaches that deal with the measurement, prediction, and modelling of non-functional features and properties, as well as approaches that address the wider spectrum of variability management, like requirements, architecture, implementation, verification, evolution and refactoring.

Important Dates

Submission: 4 November 2016

Notification: 2 December 2016

Camera ready: 16 December 2016

Submission and Publication

We welcome submissions related to the topics of VaMoS:

- Research papers describing novel contributions
- Problem statements describing open issues of theoretical or practical nature
- Reports on positive or negative experiences with techniques and tools
- Surveys and comparative studies investigating pros, cons and complementarities
- Research-in-progress including research results at a premature stage
- Papers presenting interesting and important data sets to the community
- Case studies and empirical studies
- Tool papers or demonstrations
- Vision papers.

Submissions must be in English, and between four and eight pages in length in the ACM proceedings format (<http://www.acm.org/sigs/publications/proceedings-templates>). The proceedings will be published in ACM's International Conference Series.

More information: <https://vamos2017.wordpress.com/>

2016 Internet Defense Prize for Quantum-safe Cryptography

Cryptographer Léo Ducas from CWI has won the 2016 Internet Defense Prize. He was awarded the prize with his co-authors Erdem Alkim (Ege University, Turkey), Thomas Pöppelmann (Infineon Technologies AG, Germany) and Peter Schwabe (Radboud University, the Netherlands) for their paper 'Post-Quantum Key Exchange – A New Hope'.



Some of the 2016 Internet Defense Prize winners: second from the left is Thomas Pöppelmann; Peter Schwabe is standing next to him (Alkim and Ducas not pictured). Source: USENIX.

The prize was awarded on 10 August 2016 at the 25th USENIX Security Symposium in Austin, Texas. Facebook created the Internet Defense Prize in 2014 through a partnership with USENIX. It consists of 100,000 dollars.

The winning team proposed an improved cryptosystem, called 'NewHope', that is designed to resist attacks by future quantum computers. Such quantum computers would have a devastating impact on the security of our current protocols. NewHope can, for example, be integrated into TLS and HTTPS, security protocols used by web-browsers. This was recently done by Google as an experiment toward post-quantum security. The research has been partly funded by an NWO Free Competition Grant and by a Public-Private Partnership between CWI and NXP Semiconductors.

More information on the prize:

<https://internetdefenseprize.org/>.

Full news item: <https://www.cwi.nl/news/2016/scientists-netherlands-win-2016-internet-defense-prize-newhope>



ERCIM is the European Host of the World Wide Web Consortium.



Consiglio Nazionale delle Ricerche
Area della Ricerca CNR di Pisa
Via G. Moruzzi 1, 56124 Pisa, Italy
<http://www.iit.cnr.it/>



Norwegian University of Science and Technology
Faculty of Information Technology, Mathematics and Electrical Engineering, N 7491 Trondheim, Norway
<http://www.ntnu.no/>



Centrum Wiskunde & Informatica

Centrum Wiskunde & Informatica
Science Park 123,
NL-1098 XG Amsterdam, The Netherlands
<http://www.cwi.nl/>



SBA Research gGmbH
Favoritenstraße 16, 1040 Wien
<http://www.sba-research.org/>



Fonds National de la
Recherche Luxembourg

Fonds National de la Recherche
6, rue Antoine de Saint-Exupéry, B.P. 1777
L-1017 Luxembourg-Kirchberg
<http://www.fnrl.lu/>



SICS Swedish ICT
Box 1263,
SE-164 29 Kista, Sweden
<http://www.sics.se/>



FWO
Egmontstraat 5
B-1000 Brussels, Belgium
<http://www.fwo.be/>

F.R.S.-FNRS
rue d'Egmont 5
B-1000 Brussels, Belgium
<http://www.fnrs.be/>



Spanish Research Consortium for Informatics and Mathematics
D3301, Facultad de Informática, Universidad Politécnica de Madrid
28660 Boadilla del Monte, Madrid, Spain,
<http://www.sparcim.es/>



Foundation for Research and Technology – Hellas
Institute of Computer Science
P.O. Box 1385, GR-71110 Heraklion, Crete, Greece
<http://www.ics.forth.gr/>



Magyar Tudományos Akadémia
Számítástechnikai és Automatizálási Kutató Intézet
P.O. Box 63, H-1518 Budapest, Hungary
<http://www.sztaki.hu/>



University of Cyprus
P.O. Box 20537
1678 Nicosia, Cyprus
<http://www.cs.ucy.ac.cy/>



Fraunhofer ICT Group
Anna-Louisa-Karsch-Str. 2
10178 Berlin, Germany
<http://www.iuk.fraunhofer.de/>



University of Southampton
University Road
Southampton SO17 1BJ, United Kingdom
<http://www.southampton.ac.uk/>



INESC
c/o INESC Porto, Campus da FEUP,
Rua Dr. Roberto Frias, nº 378,
4200-465 Porto, Portugal



University of Warsaw
Faculty of Mathematics, Informatics and Mechanics
Banacha 2, 02-097 Warsaw, Poland
<http://www.mimuw.edu.pl/>



Institut National de Recherche en Informatique
et en Automatique
B.P. 105, F-78153 Le Chesnay, France
<http://www.inria.fr/>



University of Wrocław
Institute of Computer Science
Joliot-Curie 15, 50-383 Wrocław, Poland
<http://www.ii.uni.wroc.pl/>



I.S.I. – Industrial Systems Institute
Patras Science Park building
Platani, Patras, Greece, GR-26504
<http://www.isi.gr/>



VTT Technical Research Centre of Finland Ltd
PO Box 1000
FIN-02044 VTT, Finland
<http://www.vttresearch.com>