# Designing Engines For Data Analysis

Peter Boncz
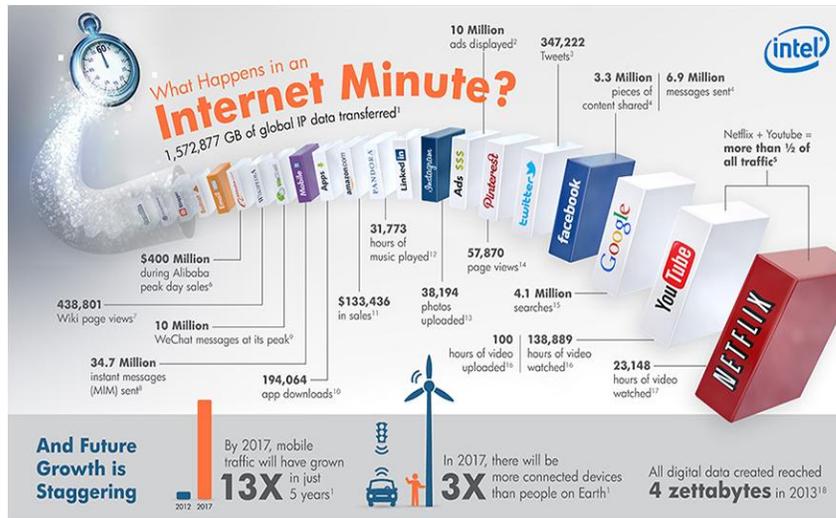
Inaugural Lecture for the Special Chair

"Large-Scale Analytical Data Management"

Vrije Universiteit Amsterdam

17 October 2014

15:45-16:30

# The age of Big Data

**The age of Big Data.** As the world turns increasingly digital, and all individuals and organizations interact more and more via the internet, for instance in social networks such as Facebook and Twitter, as they look for information online with Google, and watch movies or TV on YouTube and Netflix, the amount of information flying around is huge.

Slide 1 depicts what happens on the internet during one minute [1]. The numbers are staggering, and are expected to keep growing.  For instance, the amount of information being sent corresponds to a thousand full large disk drives, per minute. This is a hard disk drive stack of 20 meters high. All stored digital information together is estimated currently at 3 billion full disk drives, that is 4 zettabytes (when you start with a megabyte, a 1000 times that is a gigabyte, by multiplying again you get a petabyte, and then a zettabyte)

**Big Data.** All this data is being gathered, kept and analyzed. Not only by the NSA ;-), it is spread over many organizations. The question for all of them, including the NSA, is: how to make use of it?
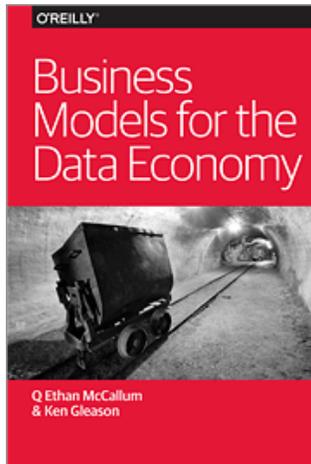
## "Big Data"

This problem is called Big Data, and is characterized by the "three V-s":

- The data comes in high **Volume**,
- The data comes at high **Velocity**,
- The data has high **Variety**, it consists of web pages, form, but also images, or any kind of text messages such as an email, Tweet or Facebook post. The data thus not only has variety, but may also be difficult to interpret (it is not always clear what people mean when they write a text message, due to missing context or slang, or simply weird abbreviations). Therefore, machine learning, information retrieval and data mining techniques to do this automatically are an important part of Big Data.

This Big Data with volume, velocity and variety requires new forms of processing to make better decisions, gain new insights and optimize processes

The Data Driven Economy

Frank van Harmelen

**The Data Driven Economy.** Across many industries, for instance health care, energy, pharmaceutical and medical research, logistics, the daily work processes, and activities in these processes increasingly revolve around data analysis. This data processing is becoming increasingly important for the functioning of organizations. It also ties into many important societal challenges, such as aging populations ("vergrijzing"), energy saving to preserve the environment, disease control, etc.

But not only inside organizations has data processing become essential, also between organizations, the processing, exchanging and enriching of data has become a business of their own. There is an economy emerging, consisting of data providers, data consumers and companies that turn this data into services, enriching it. By providing data in the open, for free, or in an open market, governments aim to stimulate new forms of growth. This has been dubbed the "Data Driven Economy" by the EU.

The picture on the right of Slide 3 depicts the EU commissioner for the digital agenda, Neelie Kroes, talking about how open government data can create value in this Data Driven Economy. The idea about open government data is that all

government data (on hospitals, roads, education etc) belongs to the public and should be publicly accessible on the internet. A relevant term here is Linked Open Data, used when data is expressed in a data format called "RDF" (the Resource Description Framework), that makes sharing data on the web easy. Thus it is a good format for open government data, because it is interlinkable in the semantic web with RDF. RDF is part of the Semantic Web Initiative, which is being driven by Tim Berners Lee for the past years. As you may know, Tim Berners Lee is the inventor of the internet as we know it, the World Wide Web.

At this point I would like to thank Frank van Harmelen. It is thanks to him that I am a professor here now, and also thanks to him that I know and understand much more about the semantic web. Frank runs a world-class research group on the semantic web at the VU to which I belong at VU.

**Disruptive Changes.** One example of this emerging Data Driven Economy is Uber, it has been in the news recently over taxi protests. You can simply see Uber as a taxi service with better software, software written for the smartphone age, so you see all taxis moving on a map and you know in advance who your taxi driver is, and do not need to pay cash, but its ambition is much wider, as it also has a

service that allows anyone (not only professional cab drivers) to share his or her car and perform taxi drives to make money.

Another example is AirBnB, also in the news, the leading website for renting private apartments, and letting your own, if you are away, also to make some money.
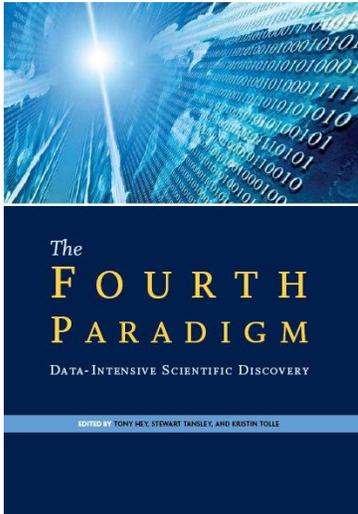
We see that those organizations that succeed at making better use of data, making themselves intermediaries between service providers and customers, can be disruptive. They are examples of the new Data Driven Economy.

Notice that these examples exploit the increasing online presence of people in the digital space, in social networks, since an important factor behind the quality of their service is in keeping track of online reputations. This can convince an AirBnB host that you will leave the apartment intact (because you have done so in the past, you have a good online reputation), and convince you in advance that a taxi-driver is punctual, and drives safely, choosing him over others.

I am not saying that such developments are unilaterally positive; taxi drivers and hotel owners disagree. And they may have certain good arguments. This should lead to democratic debate how to adjust to the new situation. But clearly the economy and society is going through disruptive changes due to this emerging Data Driven Economy.

**Data Driven Scientific Discoveries.** Science itself is also going through disruptive changes, due to large-scale data analysis. Here, I mean all sciences, not just my own, the computer science. The influential book *The Fourth Paradigm: Data-Intensive Scientific Discovery* [3] expands on the vision of pioneering computer scientist Jim Gray (more on him later). He envisioned a new, fourth paradigm of discovery based on data-intensive science and offers insights into how it can be fully realized.

## Data Disrupting Science

Scientific paragims:

1. Observing
2. Modeling
3. Simulating
4. Collecting and Analyzing Data

The FOURTH PARADIGM

*The* FOURTH PARADIGM

Data-Intensive Scientific Discovery

EDITED BY TONY HEY, STEWART TANSLEY, AND KRISTIN TOLLE

The four paradigms are summarized as follows:

1. Thousand years ago science was ***empirical***, describing natural phenomena

2. Last few hundred years, scientists developed ***models*** that explained as well as possible experimental results, allowing to understand natural phenomena, and predict the outcomes of future experiments.

3. In the last decades, the availability of compute power led scientists to further explore hypotheses by using **simulations**.

4. In the "Fourth Paradigm" scientists turn to **discovery through data analysis** making use of large datasets e.g. collected with sensors.

**Data Analysis Engines.** Let me then finally describe what I do. The title of this lecture already says it: designing engines for data analysis. In Slide 2 earlier we saw a so-called Big Data pipeline. Raw data with low information content (such as images or tweets) are being processed by algorithms on a massive scale, and massaged into high-density information or even knowledge, using machine learning and information retrieval techniques. I focus on the part in the middle, the engine.

# Engines for Data Analysis

The goal of a data processing engine is to make it easy to create such data processing pipelines. I do not typically do the analysis itself, I know little about astronomy, or pharmacy research. But what I know is to build software systems that can process large amounts of data fast.

Data Analysis Engines are building blocks for Big Data systems. Good examples are computational frameworks like MapReduce[4], graph programming frameworks like Pregel[5], and database systems like HBase[6] or Hive[7]. They shield problem owners from the complexities of large-scale system building, and make efficient use of hardware, ensuring that the Big Data task at hand is addressed by combining the computational power of multiple machines (clusters), that inside each machine, its parallel processing capabilities are leveraged (multi-core) and that data manipulation makes use of efficient data structures and algorithms, ideally in a self tuning system where the task at hand is automatically optimized into an efficient work plan.

Designing such engines for data analysis is the focus of my research.
The design of an engine is influenced by multiple factors. Important ones are the desired functionality (what should the engine accomplish?) and the capabilities of computer hardware.

# Is Designing Data Engines a Science?



▸ parallel with designing a house i.e. Architecture

The functionality can be compared to the wishes of the person wishing to build a house. The computer hardware capabilities can be compared with the available building materials, and their properties. From these, an architect must design a house. Can one compute the optimal house design? No. There is an aspect of art to architecture. This artistic part is also be perceived in some software designs. Database systems are of the most complex software systems, were:

- it is not possible to model the design space in a representative way.
- there is no such thing as one "optimal solution", the fitness of a system is highly functionality and hardware dependent (both drift over time).
- there are relevant characteristics of a system design that cannot be predicted without testing it (and thus implementing it).

System Architecture thus requires an **experimental approach** where we seek to apply the **scientific method**. In order to make progress, one must design the system **and implement it** to be able to test it in the field, or run "benchmarks". Benchmarks are an important concept for system architecture. A benchmark is a standardized test, which allows to quantify the properties of a system [8]. It allows to compare different systems, and thus to apply the scientific method to system architecture. I will come back to benchmarking later, when discussing future research.

## Systems Architecture in Academia?

**Systems architecture research** in Academia is under pressure, as it does not match with the general incentive to maximize publication count, nor with the academic funding schemes, which are short-term. Systems architecture research requires a long-term (more than ten years) commitment and stable research focus, and a significant team of people, hence significant funding. As a result, many colleagues no longer work on what we call 'core topics'. They choose topics with shorter time-span (four years or less, to fit in the PhD track of a single person). They choose topics for which less software engineering investment is needed, so there is more time to write papers. I do not blame them, but it is not my choice. Should we just leave systems architecture research to companies? A danger is that academia then becomes less relevant for IT. We know that qualities we teach our PhD students in systems architecture research are in high demand. Also, large companies maximize profit, not efficiency, so innovation might slow down if we leave all systems design to big companies, such as Microsoft. I believe that systems architecture research has a high likelihood to lead to economic opportunities, like spin-off companies. Personally, I have been involved in three such spin-offs. Therefore, in my opinion the academic community should continue to look for ways in which to engage in systems architecture research, because there are important rewards to be reaped, both scientifically and economically.
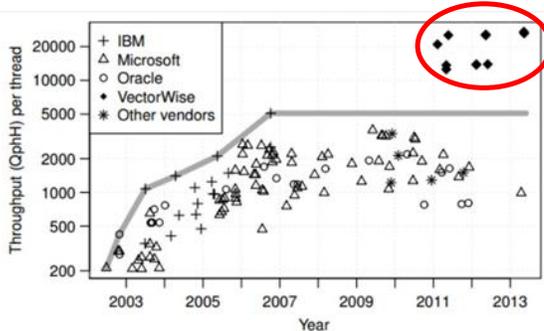
## MonetDB

**MonetDB.** My first main foray into database systems architecture research was MonetDB, this was my PhD project [9]. I am very much indebted to my mentor and advisor Martin Kersten, seen in this picture. The picture is from this summer, where he is just receiving the ACM SIGMOD Tedd Codd Award, for lifetime technical contributions to the database research community (Ted Codd invented the relational model – the basis for SQL-speaking database systems). This is a prestigious prize. That was the second major award for the MonetDB work, earlier came the VLDB 10 year best paper award in 2009 [10,11], in the small picture above. In MonetDB we rethought database architecture **for analysis,** as opposed to transaction processing. MonetDB, is called "the column-store pioneer". It has a radically new architecture in terms of data storage, query execution model, and processing algorithms. This led Michael Stonebraker - a pioneer in the field of database systems research (more on him later) - to proclaim: "the time one size fits all has gone" [12], meaning that the database systems market would split into specialized Transactional vs. Analytical DB systems, no longer offering engines that could serve all workloads, but rather offering specific engines systems for specific workloads. This has really materialized in the database industry in the past years. MonetDB is still going strong, it is a major group effort at CWI, the software is available for all in open source at http://monetdb.org.

VectorWise (Actian Vector)

▸ From the PhD thesis of Spyros Blanas (2013)

Marcin Zukowski

(b) Decision support performance per thread, measured in TPC-H queries answered per hour. VectorWise is a new database system that has been designed from scratch to better utilize modern hardware [80].

Figure 1.1: Performance per thread for transaction processing and decision support workloads. The thick gray line denotes peak performance per thread among the three established database software vendors.

**VectorWise**. After MonetDB, I started a new project, initially called X100, because it was supposed to be 100 times faster than anything else. Later we renamed the project VectorWise [13]. My first PhD student, Marcin Zukowski invented *Vectorized* Query execution. This is not the time and place to explain this in depth, but let me summarize that it introduced a series of new algorithms and techniques for optimizing data engines on modern hardware. The well-known Lucene information retrieval system borrows from it, specifically the new data compression schemes developed by Sándor Héman  that optimize not for compression ratio but for decompression speed. Vectorization aims to optimize use of deep CPU features such as caches, SIMD instructions and out-of-order execution. VectorWise can be applied also to data problems that are much larger than a computer memory (RAM), and includes new methods for use of precious disk bandwidth, where queries cooperate instead of compete. VectorWise makes use of multi-core execution and self-optimizes algorithms to modern hardware using micro-adaptivity [14].

The graph on Slide 11 I found in the recent PhD thesis of Spyros Blanas. It shows the evolution of the performance of the main commercial database engines over the past decade on one benchmark. It is quite flattering to see that VectorWise so significantly improved the state of the art.

the relational industry has been reshaped...

The ultimate form of flattery is imitation. In the past two years, the relational database industry leaders Oracle, SAP, IBM and Microsoft all have been introducing specialized analytical data engines. The first was SAP HANA (formerly T-REX[15]),a project that started by inviting our research group to SAP HQ in Waldorf to give talks about MonetDB. Then Microsoft introduced ColumnStore in SQLserver [16], IBM introduced DB2 BLU [17] and Oracle its in-memory option [18]. All three look a lot like VectorWise, but we have measured that VectorWise is in fact still faster.

The Spin-Off Company Experience

- 1996-2003 **Data Distilleries** — driving e-business personalization
- 2008- **vectorwise**
- 2013- **monetdb solutions**

It's a lot of fun! ☺

**Spin Off Companies.** Data management systems research is an applied field. A large data management industry has been created, and there is a thriving cooperation between the two. At the major data management research conferences (SIGMOD,VLDB,ICDE) there are scientific contributions not only from academia but also from industry. The close contact with industry allows for feedback from the field on challenges felt by data management practitioners. Social events are organized there by the main commercial sponsors, which are attracted also by a desire to hire the PhD students. A steady stream of spin-off companies have been formed by academic industry members, and careers of veterans in the field often span both academia and industry. Apart from the incentive of possible high financial reward, scientists who try to bring their inventions to market in a spin-off company are driven my multiple motives. One is to retain influence over the fate of the invention; another is driven by the desire to get feedback from real users and usage which is arguably more valuable than in-lab experimentation with benchmark. Further, such user feedback often provides new research ideas, and a final motive is the ability to address a problem with more engineering resources than is possible in a university.

I have personally been involved in three spin-off companies from my research. In the first I was very young, during my PhD and other people, namely Marcel Holsheimer, Fred Kwakkel, Arno Siebes and Martin Kersten were at the steering wheel, however my PhD research was the engine of this company. Data Distilleries truly pioneered predictive analytics, by putting data mining into business processes even in the late 1990s. It grew up to more than 100 employees internationally, but after the internet bubble burst and 9/11 put the world economy into recession it retreated to Europe, still 40 person strong. It was acquired by SPSS in 2003.
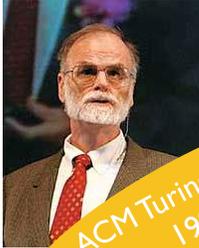
The second company was VectorWise. VectorWise created the leading analytical database engine out of the work of two PhD students, Marcin Zukowski and Sandor Heman, and Niels Nes, then a post-doc. It was set up by striking a deal with an existing database company (Actian Corporation), which helped us focus on the technical challenges, while they were doing the marketing and sales. This made financial life easier than with Data Distilleries, which had to continuously provide its own revenue stream.

Currently VectorWise is named "Actian Vector". In fact, Actian most prominently markets the technology as part of its Analytical Platform "Hadoop SQL Edition" (see www.actian.com). CWI and Actian have forged a research collaboration agreement after the sale of VectorWise. Out of that collaboration, a version of Vector that runs on Hadoop clusters has been developed ("Hadoop SQL Edition") which remains powered by the VectorWise system. The continued relationship is valuable for my research, because the development group in Amsterdam is a powerful base for new research by MSc and PhD students, and it is a valuable source of funding research at CWI.

The third spin-off company, MonetDB Solutions, is a small firm that works in symbiosis with the CWI research group, it acts when MonetDB open source users need serious commercial support.

# Database Research Pioneers

Jim Gray

Michael Stonebreaker

ACM Turing Award 1998

IEEE Von Neumann Medal 2004

Participating in these spin-off endeavors was a rewarding experience, as I learnt a lot of technical, business, and management skills.  It is great to see the research ideas take off in practice, and the user feedback one obtains while doing so gives great inspiration for new research work.

While some may see spin-off company activities as a diversion that unavoidably will diminish the scientific impact of a researcher, experience in the database research community tells a different story. Undisputedly the two most famous database systems researchers are Jim Gray and Michael Stonebraker – I mentioned both before.

Jim Gray started his database career as team member of the seminal SystemR project at IBM [19], the first SQL speaking database system based on the relational model. He subsequently changed to Tandem Corp. where he made further foundational contributions to database transaction processing [20] and database benchmarking [21]. Later in his career he moved to Microsoft, and founded his own lab (BARC – Bay Area Research Center) where he started to collaborate with scientists from diverse fields on their data management problems, specifically with astronomers he worked on making the Sloan Digital

Sky Survey (SDDS) fully database-driven [22]. From this work came his vision for a Fourth Paradigm of scientific discovery.

Michael Stonebraker is the research scientist with the longest and most successful track record of database startups. Ingres [23], Postgres [24], Illustra [25], StreamBase [26], Vertica [27], VoltDB [28] and SciDB [29] to name a few projects. Berkeley Ingres was one of the first relational database systems, Postgres and Illustra focsed on extensibility, whereas StreamBase is an engine for querying data streams. Vertica and VoltDB and specialized database systems targeting respectively analytical and transactional workloads ("one size does not fit all"). SciBase is a database system rich in matrix and array storage design to support science applications.

One would think that people who spent so much time in industry would not make top scientists. Yet these two people are the most highly decorated data management researchers, with a Tuning Award and the Von Neumann medal. The fact that they were able to prove the merits of their ideas an algorithms in systems that were used and conquered markets, gives their scientific papers about these more weight, credibility and recognition.
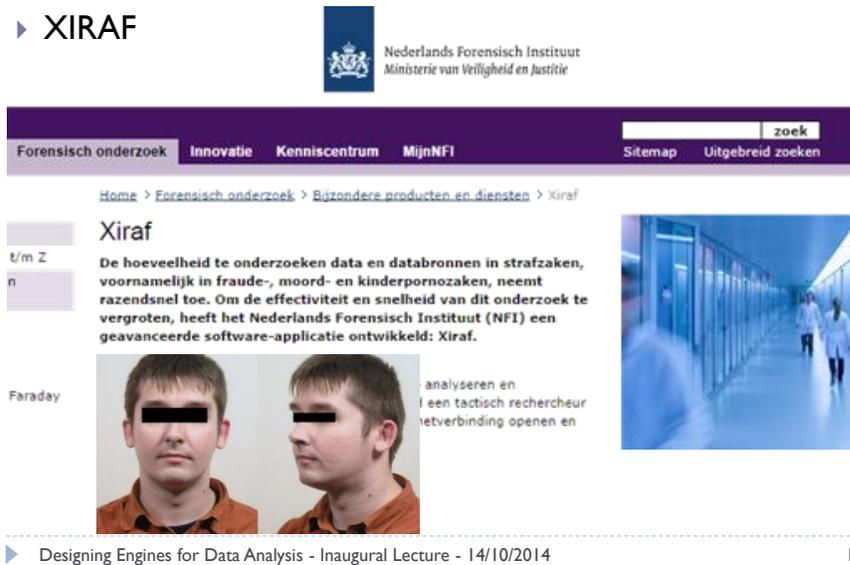
I cannot stand in the shoes of these two, but I have tried to continue in their direction.

The lesson I draw is that influential spin-off activities can in fact bolster a scientific career in a field where academia and industry closely cooperate. In my case, I might not have gotten a Humboldt Research Award for career-wide achievements if I had not proven the worth of my ideas by means of the VectorWise spin-off.

**Application Examples.** The MonetDB and VectorWise technologies have by now been deployed in hundreds of cases. Most we know little about, but in some of these, we are actively involved.

## Application Examples

▸ XIRAF

**Nederlands Forensisch Instituut**
*Ministerie van Veiligheid en Justitie*

| Forensisch onderzoek | Innovatie | Kenniscentrum | MijnNFI | | Sitemap | Uitgebreid zoeken | zoek |

Home > Forensisch onderzoek > Bijzondere producten en diensten > Xiraf

### Xiraf

t/m Z
n

De hoeveelheid te onderzoeken data en databronnen in strafzaken, voornamelijk in fraude-, moord- en kinderpornozaken, neemt razendsnel toe. Om de effectiviteit en snelheid van dit onderzoek te vergroten, heeft het Nederlands Forensisch Instituut (NFI) een geavanceerde software-applicatie ontwikkeld: Xiraf.

Faraday

... analyseren en
... een tactisch rechercheur
... netverbinding openen en

For instance, in the case of LOFAR, the digital telescope I talked about earlier, Bart Scheers works between the MonetDB group and the Pannekoek Institute on a project to detect quick events in the sky. In realtime, the sky is being monitored and software must distinguish between light sources that are always there and those who are transient. The latter ones, are the interesting ones [30]. Clearly, a fast database system is a necessary ingredient for this real-time monitoring. Vectorwise has more than a hundred customers, from all walks of life. An inspiring example is the Oxford University Clinical Trial Service Unit which uses it for Cancer and Heart desease research in the BioBank project [31].

Slide 15 describes the XIRAF system [32], designed to help analyze digital evidence at the Dutch Forensic Institute (NFI), also a Big Data problem. I stood at the cradle of this system, co-advising the Msc project of Wouter Alink, who later joined the NFI to develop XIRAF into a usable product, while I cooperated in a joint research project. Since then this software has been significantly evolved by the team at NFI, originally led by Raoul Bhoedjang. I know XIRAF was used to unravel the pedophilia network through wich Roberts M. (infamous in The Netherlands) disseminated his images and videos. I am proud of the achievements of the XIRAF team at NFI.

unethically

▸ Selling  financial products

high risk

SPAARBELEG

100% DOCHTER VAN AEGON

Dat is makkelijk verdiend.

Stichting Woekerpolisproces

De sprintplanclaim succesvol tegen Aegon

Sprintplanclaim is de naam waaronder de Vereniging Consument en Geldzaken een rechtszaak voert tegen verzekeraar Aegon. Sprintplan was de naam van de woekerpolis van Aegon die meer dan honderdduizend verzekerden heeft gedupeerd. Deze gedupeerden liepen zo een schade op van om en nabij de zevenhonderd miljoen euro. Kenmerkend aan het Sprintplan was dat het werd vormgegeven als een aandelenleaseplan. Zo werd dit echter geenszins naar de consument toe gepresenteerd. Zo dwaalden
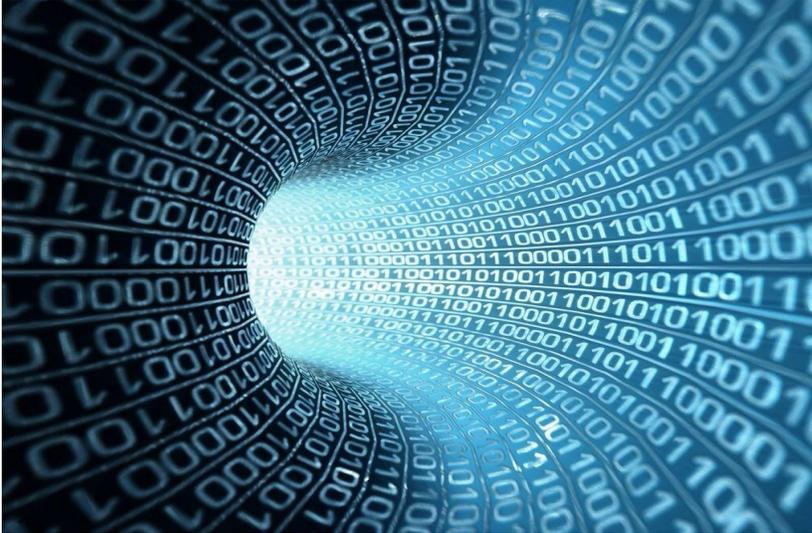
One of the early customers of Data Distilleries was the well known financial firm Aegon, specifically its Spaarbeleg service. We helped Aegon sell more of its financial products, like SprintPlan, by using data mining and predictive analytics to improve marketing efficiency. Regrettably, we were not aware that these products were very shaky and could leave the investors with a debt, and that Spaarbeleg was not honest about that in its sales message. As such, unwillingly we had some part in the so-called "woekerpolis affair" [33].

**Ethical Aspects.** This story brings me to the more reflexive section of this lecture. Clearly, data analysis and data analysis engines can be used both for good and bad purposes. There are both ethical and political questions around this. I am not a philosopher nor a sociologist or a political scientist, but I probably have a broader view than most on what is happening with data.

Let me thus share some thoughts on:

▸ Who guarantees or arbitrates on issues of Fairness of the digital playing field in this new Data Driven Economy?

▸ Similarly, for Preserving Privacy and Reputation?

▸ How should risks be shared between, e.g. banks and their customers?

▸ How to better enforce ethical data usage?
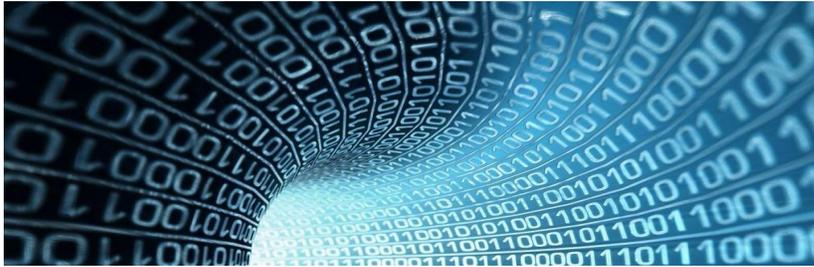
Ethics and Policies in the Data Driven Society

**Reputation Protection Rights.** It has been shown that it is easy to manipulate the Google Autocomplete feature. It displays you the most often used next words in a search query, to help you enter the query faster. However, if enough mean people type the search query "Peter Boncz pedophile" , and you can cheaply buy such manual labor in sites like Mechanical Turk, then after typing "Peter Boncz" in the search box, Google will suggest "Peter Boncz pedophile" first, so everyone looking for me will notice that from then on. Because of this association and the thought that there where there is smoke, there must be some fire, my reputation is easily tarnished. The bad part here, is that there is no (good) way for arbitration or to complain at Google when faced with such problems. Google just says (or at least used to say) the algorithm is right. This is the cheap way out, arbitration is overhead and costs money. As search engines and social networks start dominating the conversational space and the reputation space, internet companies can earn a lot of advertising money. While they have been quick to seize advertizing markets, these companies have been slow though to recognize their responsibility and role to effectively deal with ethical dilemmas such as maintaining free speech while also protecting reputations. Nation states have mostly stayed out of this space, so this is too much of a no-mans land right now.

# Ethics and Policies in the Data Driven Society

**Banken stellen nieuwe regels voor internetbankieren**
zondag 24 november 2013, 09:55 door **Redactie**, 105 **reacties**
Laatst bijgewerkt: 24-11-2013, 14:30

Nederlandse banken hebben nieuwe regels voor internetbankieren opgesteld waar klanten aan moeten voldoen als ze in het geval van fraude hun geld terug willen krijgen. Zo mag er geen illegale software zijn geïnstalleerd, moet de computer up-to-date zijn en moet de rekening regelmatig worden gecontroleerd.

Dat laat de Nederlandse Vereniging van Banken (NVB) weten. Tot nu toe hanteerde iedere bank zijn eigen veiligheidsvoorschriften. Banken bepalen zelf wanneer ze de nieuwe uniforme regels aan hun klanten communiceren. Tot dat moment zullen banken bij de behandeling van nieuwe claims van klanten vanaf 1 januari 2014 in de geest van de nieuwe regels handelen.

▷ Designing Engines for Data Analysis - Inaugural Lecture - 14/10/2014                18

**Shifting Fraud Liability.** Dutch banks announced at the start of this year, that in case of theft by electronic means from a bank account, they would no longer guarantee to cover the damages, if your PC (and its anti-virus) was not up-to-date or if the PC had illegal software on it. My reaction was to ditch our Microsoft PC, because in my perception these computers run more significant security risks than others like Apple or Linux, and with these kinds of conditions from the bank it is very unclear whether your case would hold, in case of trouble. And trouble is quite likely, given the amount of computer viruses that go undetected.

There is no easy technological fix to the pluriform potential problems associated with Big Data and the Data Driven Economy.  There is a strong need for debate and democratic checks and balances on the now digital aspects of our society. In my opinion, technology is not questioned enough.  Compare the current situation with a hundred years ago when cities had become very noisy after the industrial revolution, car engine motors were loud, trams turned corners screeching and factories made a lot of industrial noise. In response, "Anti-Noise Leagues" were formed and started campaigning against the noise. In the end, train suspension reduced screeching and exhausts dampers were fitted to cars. Technology

proponents back then had been saying that noise just was part of "progress" and people had to live with it. Thanks to the Anti-Noise Leagues, things took a different turn. We should similarly campaign against negative side effects of Big Data, and hope for privacy and reputation exhaust damper technology. To change the data society for the better, pressure by the public and government may lead the involved parties to self-regulate and take more responsibility.

**Algorthmic Auditing.** Solely depending on voluntary action of companies is probably not enough, so nation states (or the EU, here) should under circumstances regulate certain aspects of the Data Driven Economy. Such regulation needs to be enforced. Just like companies get inspected on the truthfulness of their financial reporting by a neutral accountant or auditor who must keep the detailed information confidential, a similar thing could be done for the algorithms data processing systems. One can see this as an advanced form of EDP (Electronic Data Processing) Auditing. An Algorithmic Auditor would be a Big Data savy person who can ask the right questions to see whether the data policies and algorithms deployed conform to regulations. It would be a natural extension of normal accounting practices that could be relevant for companies of certain size and focus.

## Education: Data Science



talent

Amsterdam Data Science develops multi-disciplinary curricula where students experiment with realistic big data sets throughout the program. The data science program aims to attract talented students who will be educated in all topics underlying data science. From understanding and configuring large scale infrastructures to deep knowledge in data analytics and beyond to presentation of and interaction with data.

▸ Feb2014, new course, MSc level
   "**Large-Scale Data Engineering**"

**Education.** Currently we are observing an acute shortage of Big Data savy personnel, even without Accounts firms hiring Algorithmic Auditors. The only qualified people right now tend to be people with a PhD in some data processing related topic. There are very few of those going round. To fill this gap, there is an educational demand that I hope to help address as part of my professor appointment. This I will do at VU in the context of Amsterdam Data Science (ADS), which is an initiative from VU, UvA,CWI and HvA, since 2014. Here, "Data Science" is the study of the generalizable extraction of knowledge from data. It is closely related to Big Data. In Data Science the data does not necessarily have to be huge, though it often is.

Amsterdam Data Science  is developing an Education Program, and with help of Hannes Mühleisen I will develop and teach a new master course at VUA called "Large-Scale Data Engineering" early 2015. This course will teach the fundamentals of efficient data processing algorithms and its practical will train students in using tools for data processing in clusters. As such, the course should an important technical piece of the Data Science education program.
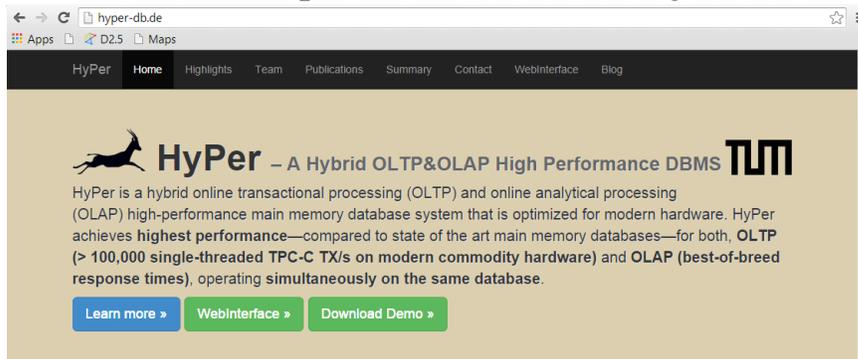
Research: Graph Analytics

**Future Research: Graph Analysis Systems.** One of my current research interests and research cooperations at VU (with Claudio Martella and Spyros Voulgaris and students) is analyzing not just data in simple tables, but data that has complex shapes. The shape, for instance, of a social network, which in mathematical terms is a graph. Such search problems that focus on the connections, the nearness, or the clustering of data occur in very many Big Data scenarios.

However, the field of graph management systems, being it graph database systems such as NEO4J[34] or graph programming frameworks such as Pregel, is still not very well understood. There are very strong systems architecture challenges there, due to this unclear functionality and purpose on the one hand, and the high complexity of graph problems on the other hand. In order to apply the scientific method, to be able to quantitatively compare approaches, we need standard tests or benchmarks. So this is where we are focusing on in the LDBC project. LDBC stands for Linked Data Benchmark Council. More information can be found on http://ldbcouncil.org

# Research: Compilers & Database Systems

Thomas Neumann

Alfons Kemper

**Future Research: Database Compiler Integration.** My specific interest in compilers further has two causes. One is the inability in Big Data systems to express all query needs in SQL. Very often, program code must be integrated. Running these together however is error prone and often slow. Compiler techniques promise to bring data engines and user code together in a seamless way, also allowing better correctness checking. The second reason is hardware related. Right now, the hardware market is splitting into a mobile and a server market. Until now, the market has been fully driven by consumer devices, so far these were PCs and later laptops. Nowadays, however, mobile phone hardware dominates and such hardware must run with low power, for better battery life. As consumer migrated to the lower end, the use of power-hungry hardware is no longer dominated by the consumer market, but by the server market. Also, the server hardware buyers get more influential, since a few of them (Amazon, Google, Apple, Oracle,RackSpace..) buy hardware by the tens of thousands, destined for a life in huge compute centers serving the cloud. This means that server feature wishes are likely to influence hardware designs. In the past, database systems architects would be forced to work with hardware features designed for consumer use cases (multi-media instructions such as SSE [35], even graphic cards [36]). In the future, likely we are going to see dedicated server

hardware. Oracle certainly is trying with its RAPID project [37]. Also, there are so many transistors on a chip now, that more and more non-frequently used features are still viable to put in hardware. The functions one does not need, just get switched off. This is called "dark silicon": CPU features that are most often switched off.

Now, one must understand that hardware has already evolved and is very diverse. Modeling performance is already impossible, due to the many complex interactions and limits that CPU architects introduce and programmers cannot know about. CPU performance depends on room temperature, because CPU chips routinely speed up or slow down depending on how hot they run. The unpredictable nature of hardware is hard for data management systems, which need cost models and optimizers in order to self-tune tasks. The further trend of diversification of server hardware is going to make this variability even larger, and the cost modeling problem even harder. The complexities of hardware are such and its diversity is such, that we cannot ship a single binary to an end-user. Rather, we need systems that dynamically adapt themselves to heterogeneous hardware, when and where the system is run. The most likely way forward to me indeed seems to be to exploit compiler techniques in the future, because just-in-time compilation of data management systems allows adjusting the program to the specific properties of the specific hardware it is compiled and run on.

During the past year I have stayed half-time at TU Munich, graciously hosted by professors Thomas Neumann and Alfons Kemper. During this stay, I have become very interested in combining two fields, namely programming language compilers and data processing systems. Their Hyper system [38] uses just-in-time compilation in new ways and the result is amazing, in fact their technology could be capable of re-joining the specialized analytical and transactional engines from the split we caused with MonetDB and VectorWise, since using such just-in-time compilation (and a slew of other intelligent techniques) both kinds of workloads can be made efficient in a single system. Hyper certainly is an inspiration for future research activities.

**Acknowledgements – Dankwoord.** Ik wil graag mijn ouders Hanny and Árpád bedanken, zonder hun opvoeding had ik hier niet gestaan. Ik bedank mijn vrienden voor alle steun, Hylke,Jacqueline,Reinier,Annemarie,Pilar, Frank en alle anderen die ik niet zo snel kan opnoemen. A huge thank-you goes to Marcin for flying over for just two days.I would like to thank my colleagues at CWI, and Spinque; great to work with you and have you here. I would also like to thank my colleagues at VU, specifically Maarten van Steen and Frank van Harmelen, thanks for hosting me in the department and the group and for guiding me to this point. I am looking forward to more with all VU colleagues. Specific thanks go to all professors in the Cortege. Sehr vielen dank Thomas, Sabine und Alfons das ihr aus München gekommen seid!! Und auch vielen Dank führ das verganene Jahr, es war sehr inspirierend. I would like to thank my colleagues from the Vectorwise group of Actian. Very special thanks to Emma and Seetha for flying from the US. It is great to be advising such a talented group of people, many of you are my ex-students and I am very proud of you. Thanks to everyone else who is here. Bedankt iedereen voor het komen. Ik bedank ook het bestuur van de stichting Vu-VUmc, het College van Bestuur en het bestuur van de Faculteit Exacte Wetenschappen, voor het mogelijk maken van mijn leerstoel. Tot slot dank voor mijn familie, en met name mijn gezin.  Laura stond gisteren nog voor 180 mensen te zingen en is minder zenuwachtig dan ik. Matias, bedankt voor het luisteren, sorry dat het allemaal Engels was. Cecilia, muchisimas gracias por ayudarme en tantas maneras, te quiero mucho.  Ik heb gezegd.

## Bibliography

[1] www.intel.com/content/www/us/en/communications/internet-minute-infographic.html

[2] ec.europa.eu/digital-agenda/en/towards-thriving-data-driven-economy

[3] A Hey. The fourth paradigm: data-intensive scientific discovery. *Microsoft Research*, 2009.

[4] J Dean, et al. MapReduce: simplified data processing on large clusters. *CACM* 51.1,2008.

[5] G Malewicz, et al. Pregel: a system for large-scale graph processing, 2010.

[6] L George. *HBase: the definitive guide*.  O'Reilly Media, Inc., 2011.

[7] A Thusoo, et al. Hive: a warehousing solution over a map-reduce framework. *Proceedings of the VLDB Endowment* 2.2, 2009.

[8] J Gray. *Benchmark handbook: for database and transaction processing systems*. Morgan Kaufmann Publishers Inc., 1992.

[9] P Boncz. *Monet; a next-Generation DBMS Kernel For Query-Intensive Applications*. UvA, 2002.

[10] P Boncz, S Manegold, and M Kersten. Database architecture optimized for the new bottleneck: Memory access. *VLDB*. 1999.

[11 P Boncz, M Kersten, S Manegold. Breaking the memory wall in MonetDB. *CACM* 51.12, 2008.

[12] M Stonebraker, U Cetintemel. One size fits all: an idea whose time has come and gone. *Data Engineering, 2005. ICDE 2005*

[13] P Boncz, M Zukowski. Vectorwise: Beyond Column Stores. *DEBULL* 35.1, 2012.

[14] B Răducanu, P Boncz, and M Zukowski. Micro adaptivity in vectorwise. *SIGMOD,* 2013.

[15] F Färber, et al. SAP HANA database: data management for modern business applications. *ACM SIGMOD Record* 40.4,2012.

[16] P-A Larson, et. al. Columnar Storage in SQL Server 2012. *DEBULL*.35.1, 2012.

[17] V Raman, et al. DB2 with BLU Acceleration: So much more than just a column store. *PVLDB* 6.11, 2013.

[18] ORACLE DATABASE 12 C IN-MEMORY OPTION – company whitepaper 2014.

[19] M Astrahan, et al. System R: relational approach to database management. *TODS* 1.2, 1976.

[20] J Gray, A Reuter. *Transaction processing*. Morgan Kaufíann Publishers, 1993.

[21] J Gray, *Benchmark handbook: for database and transaction processing systems*. Morgan Kaufmann Publishers Inc., 1992.

[22] C Stoughton, et al. Sloan digital sky survey. *The Astronomical Journal* 123.1, 2002.

[23] M Stonebraker, et al. The design and implementation of INGRES. *TODS* 1.3, 1976.

[24] M Stonebraker, L. Rowe. *The design of Postgres*. ACM, 1986.

[25] M Stonebraker,D Moore. *Object Relational DBMSs: The Next Great Wave*. Morgan Kaufmann Publishers Inc., 1995.

[26] D Abadi, et al. Aurora: a new model and architecture for data stream management. *VLDB Journal* 12.2, 2003.

[27] M Stonebraker, et al. C-store: a column-oriented DBMS. *PVLDB*, 2005.

[28] M Stonebraker, et al. TheVoltDB Main Memory DBMS.*DEBULL*. 36.2, 2013.

[29] M Stonebraker, et al. Requirements for Science Data Bases and SciDB. *CIDR*.. 2009.

[30] Y Zhang, et al. Astronomical Data Processing Using SciQL, an SQL Based Query Language for Array Data. *ASP Conference Series*. 2012.

[31] www.ukbiobank.ac.uk/

[32] W Alink et al. XIRAF–XML-based indexing and querying for digital forensics. *Digital investigation* 3, 2006.

[33] www.rijksoverheid.nl/verzekeren

[34] J Webber. A programmatic introduction to Neo4j. *Proceedings of the 3rd annual conference on Systems, programming, and applications: software for humanity*. ACM, 2012.

[35] J Zhou, K Ross. Implementing database operations using SIMD instructions.*SIGMOD* 2002.

[36] B. He, et al. Relational query coprocessing on graphics processors. *TODS* 34.4, 2009.

[37] T Neumann Efficiently compiling efficient query plans for modern hardware. *PVLDB* 4.9, 2011.