



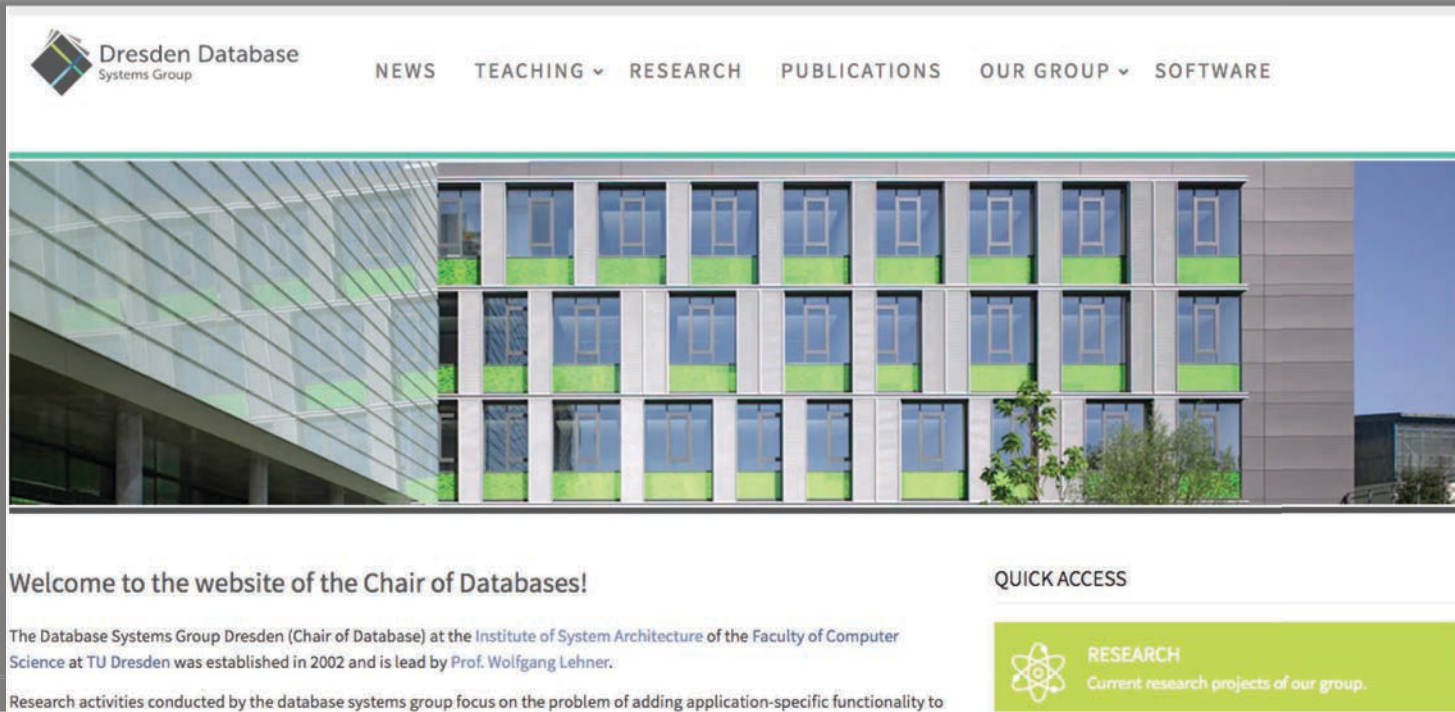
Future Trends in Data Analytics

Hannes Voigt

General Information

Lecture Notes

- Further information on our website <http://www.db.inf.tu-dresden.de>



The screenshot shows the homepage of the Dresden Database Systems Group. At the top is a navigation bar with the group's logo and the text "Dresden Database Systems Group". To the right of the logo are links for "NEWS", "TEACHING" (with a dropdown arrow), "RESEARCH", "PUBLICATIONS", "OUR GROUP" (with a dropdown arrow), and "SOFTWARE". Below the navigation bar is a large banner image of a modern building with a glass facade and green-tinted windows. Under the banner, the text "Welcome to the website of the Chair of Databases!" is displayed. To the right of this text is a "QUICK ACCESS" section with a green button labeled "RESEARCH" and the text "Current research projects of our group." below it. At the bottom left, there is a logo for "TECHNISCHE UNIVERSITÄT DRESDEN" and a paragraph of text: "The Database Systems Group Dresden (Chair of Database) at the Institute of System Architecture of the Faculty of Computer Science at TU Dresden was established in 2002 and is lead by Prof. Wolfgang Lehner." Below this paragraph is another line of text: "Research activities conducted by the database systems group focus on the problem of adding application-specific functionality to".

Dresden Database
Systems Group

NEWS TEACHING ▾ RESEARCH PUBLICATIONS OUR GROUP ▾ SOFTWARE

Welcome to the website of the Chair of Databases!

QUICK ACCESS

RESEARCH
Current research projects of our group.

TECHNISCHE
UNIVERSITÄT
DRESDEN

The Database Systems Group Dresden (Chair of Database) at the Institute of System Architecture of the Faculty of Computer Science at TU Dresden was established in 2002 and is lead by Prof. Wolfgang Lehner.

Research activities conducted by the database systems group focus on the problem of adding application-specific functionality to

Dresden Database System Group



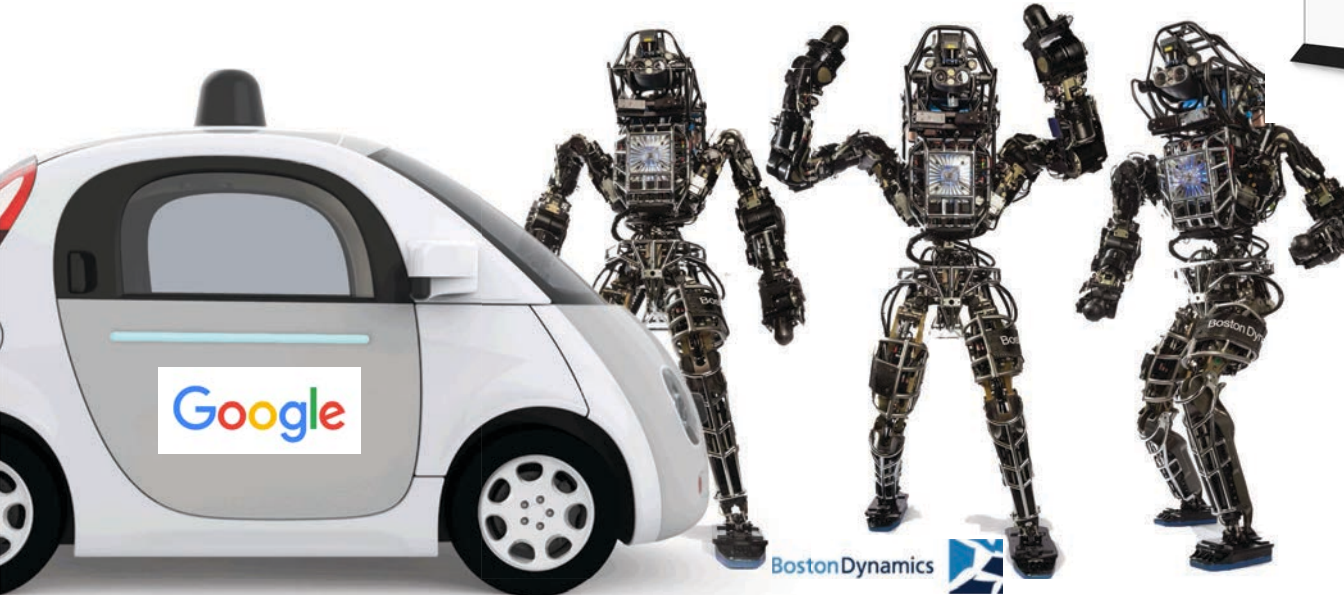
Future is Now

- ..., face detection & recognition, autobeam, Siri/Cortana/..., 3D printing, ...

JEOPARDY!



Dresden Database
Systems Group



[<https://www.google.com/selfdrivingcar/>]

[<http://www.bostondynamics.com/>]



[<http://www.idsc.ethz.ch/research-dandrea/research-projects/cubli.html>]

... tomorrow?

Sooner then you think!



Dresden Database
Systems Group





When have we
entered the future?

What was the threshold
we crossed?

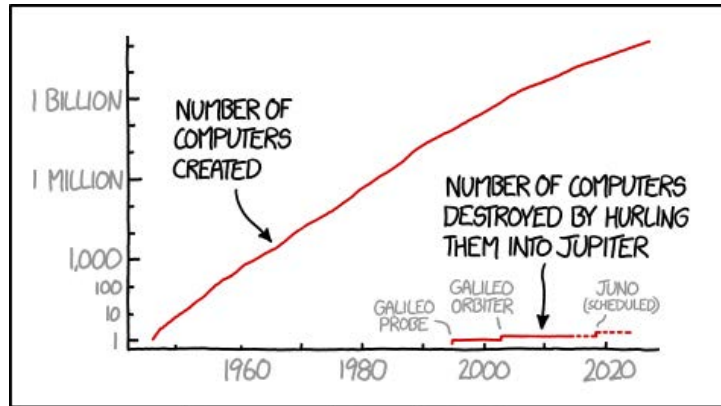


Exponential Growth

Exponential Growth

Outnumbers everything else quickly

- Asymptotic advantage
- Quickly increasing add-on



NASA NEEDS TO PICK UP THE PACE.
IF THEY EVER WANT TO FINISH THE JOB.

[<http://xkcd.com/1727/>]

Leads to surprising results

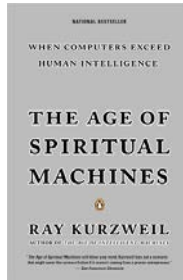
- Black swans in economic crisis
- Shooting stars in business, media, sport, etc.










Second Half of the Chessboard

When exponential growth really kicks in

- According to Ray Kurzweil
- Things start to get interesting in the second half of the chess board
- Beyond 4G numbers quickly go beyond human intuition
- What happens in the second half can hardly be foreseen

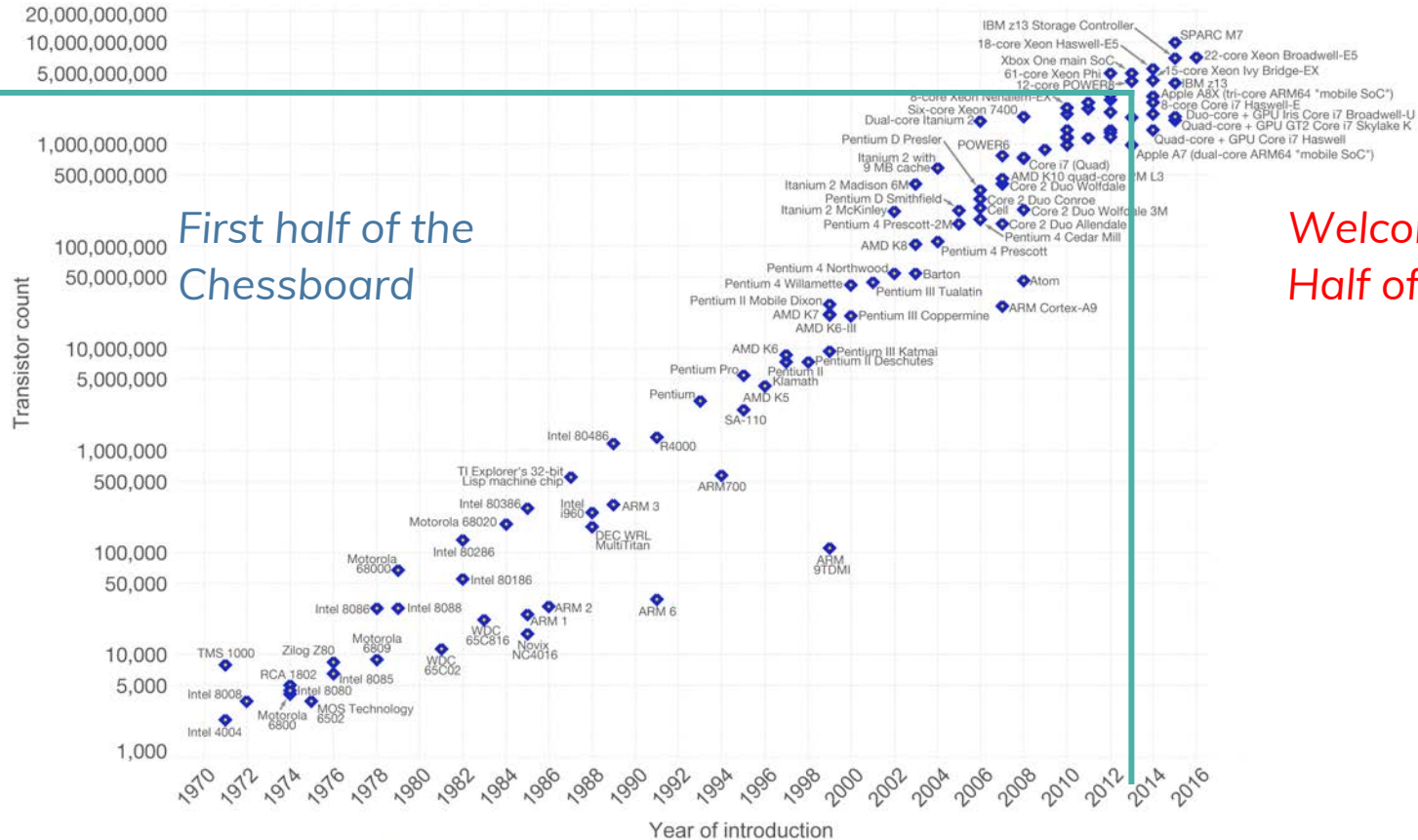


							128
1	2	4	8	16	32	64	128
256	512	1024	2048	4096	8192	16384	32768
256	512	1024	2048	4096	8192	16384	32768
65536	131K	262K	524K	1M	2M	4M	8M
65536	131K	262K	524K	1M	2M	4M	8M
16M	33M	67M	134M	268M	536M	1G	2G
16M	33M	67M	134M	268M	536M	1G	2G
4G	8G	17G	34G	68G	137G	274G	549G
4G	8G	17G	34G	68G	137G	274G	549G
1T	2T	4T	8T	17T	35T	70T	140T
1T	2T	4T	8T	17T	35T	70T	140T
281T	562T	1P	2P	4P	9P	18P	36P
281T	562T	1P	2P	4P	9P	18P	36P
72P	144P	288P	576P	1E	2E	4E	9E
72P	144P	288P	576P	1E	2E	4E	9E

[https://en.wikipedia.org/wiki/File:Wheat_Chessboard_with_line.svg]

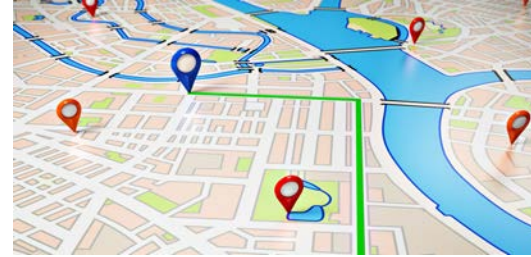
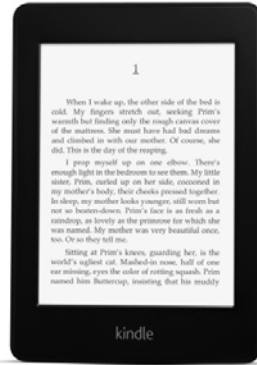
Moore's Law – The number of transistors on integrated circuit chips (1971-2016)

Moore's law describes the empirical regularity that the number of transistors on integrated circuits doubles approximately every two years. This advancement is important as other aspects of technological progress – such as processing speed or the price of electronic products – are strongly linked to Moore's law.



Digitization

Everything is Digital



Landscape has changed

From Islands of digital data ...

... to ponds of analog signals



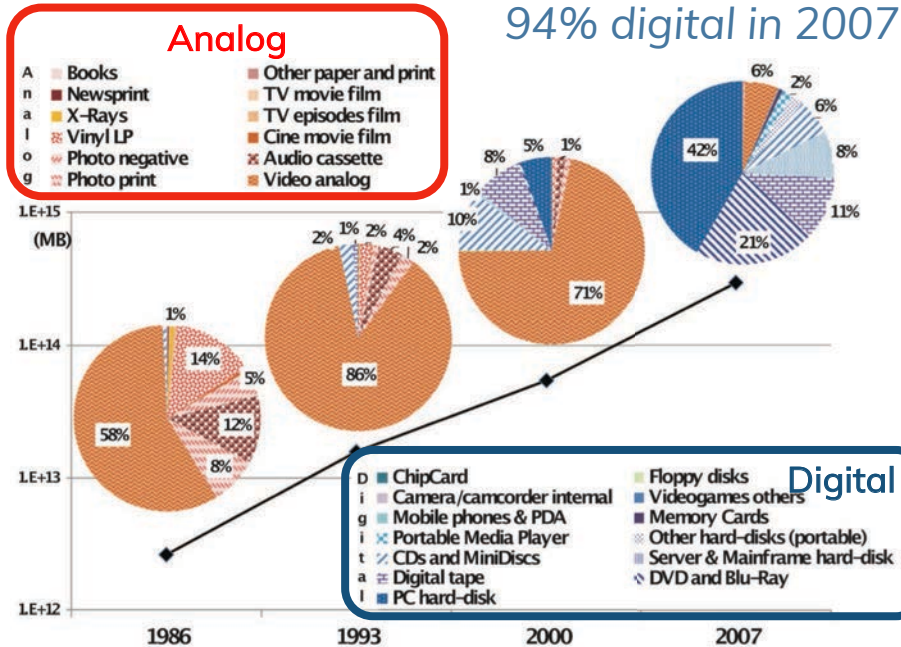
(Tuamotu Archipelago, French Polynesia)



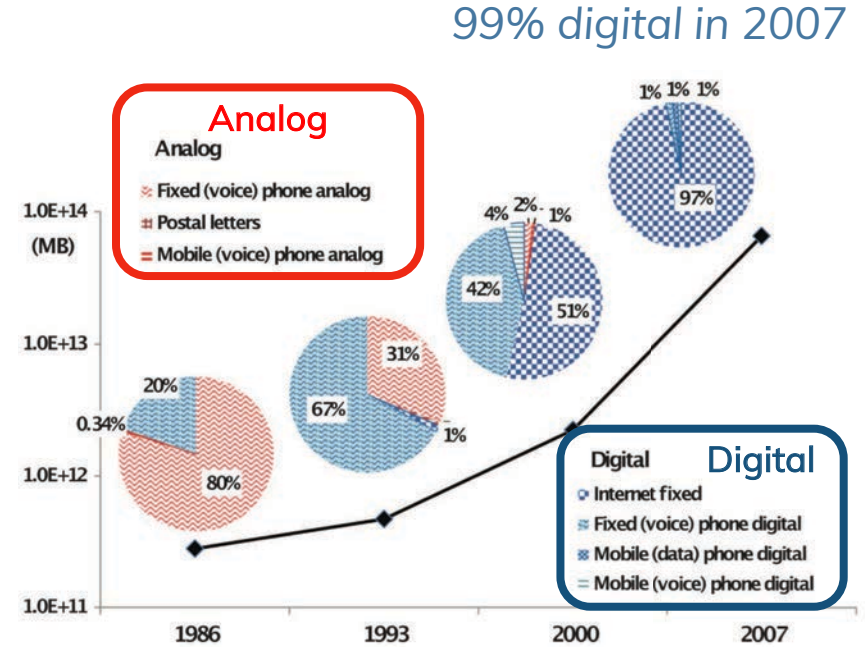
(Algonquin Provincial Park, Ontario, Canada)

Everything is Digital

World's capacity to store information



World's capacity to telecommunicate



[M. Hilbert and P. Lopez, The World's Technological Capacity to Store, Communicate, and Compute Information, Science, 332, April 2011, DOI: 10.1126/science.1200970]

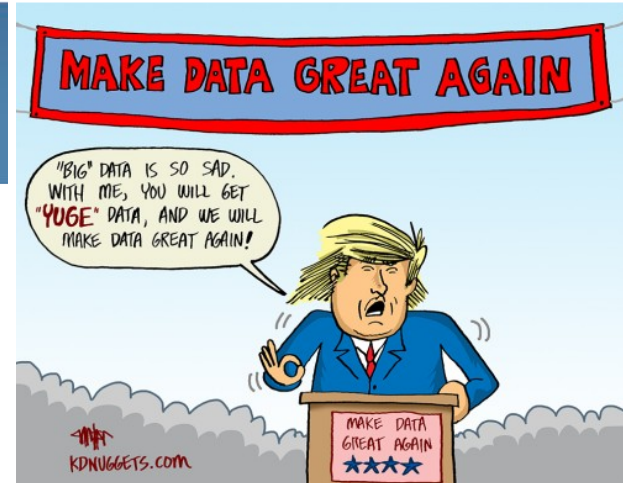
Everything is Digital



[<http://blog.acronis.com/posts/data-everything-8-noble-truths>]



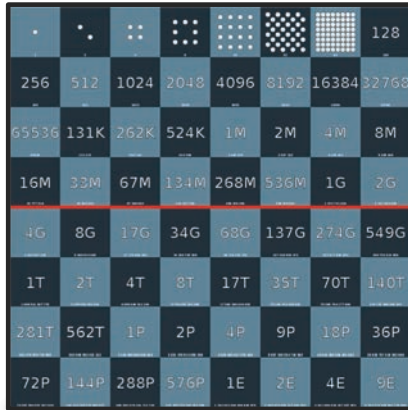
Big Data



Big Data term reflects the three drivers

Exponential growth

- beyond intuition
- second half of the chessboard



A chessboard illustrating exponential growth, showing the number of grains of rice on each square. The board is 8x8, and the number of grains doubles on each square, starting from 1 grain on the first square and reaching 128 grains on the last square.

1	2	4	8	16	32	64	128
256	512	1024	2048	4096	8192	16384	32768
65536	131K	262K	524K	1M	2M	4M	8M
16M	33M	67M	134M	268M	536M	1G	2G
4G	8G	17G	34G	68G	137G	274G	549G
1T	2T	4T	8T	17T	35T	70T	140T
281T	562T	1P	2P	4P	9P	18P	36P
72P	144P	288P	576P	1E	2E	4E	9E

Big

Digitization

- Everything is digital data
- Analog signals are not part of big data



Data

40 ZETTABYTES

[43 TRILLION GIGABYTES]
of data will be created by 2020, an increase of 300 times from 2005



Volume SCALE OF DATA

It's estimated that
2.5 QUINTILLION BYTES
[2.3 TRILLION GIGABYTES]
of data are created each day

Most companies in the U.S. have at least
100 TERABYTES
[100,000 GIGABYTES]
of data stored

The New York Stock Exchange captures
1 TB OF TRADE INFORMATION
during each trading session



Velocity ANALYSIS OF STREAMING DATA

By 2016, it is projected there will be
18.9 BILLION NETWORK CONNECTIONS
— almost 2.5 connections per person on earth



Modern cars have close to
100 SENSORS
that monitor items such as fuel level and tire pressure



The FOUR V's of Big Data

From traffic patterns and music downloads to web history and medical records, data is recorded, stored, and analyzed to enable the technology and services that the world relies on every day. But what exactly is big data, and how can these massive amounts of data be used?

As a leader in the sector, IBM data scientists break big data into four dimensions: **Volume, Velocity, Variety and Veracity**

Depending on the industry and organization, big data encompasses information from multiple internal and external sources such as transactions, social media, enterprise content, sensors and mobile devices. Companies can leverage data to adapt their products and services to better meet customer needs, optimize operations and infrastructure, and find new sources of revenue.

By 2015
4.4 MILLION IT JOBS
will be created globally to support big data, with 1.9 million in the United States



As of 2011, the global size of data in healthcare was estimated to be

150 EXABYTES
[161 BILLION GIGABYTES]



30 BILLION PIECES OF CONTENT
are shared on Facebook every month



Variety DIFFERENT FORMS OF DATA

By 2014, it's anticipated there will be

420 MILLION WEARABLE, WIRELESS HEALTH MONITORS

4 BILLION+ HOURS OF VIDEO
are watched on YouTube each month



400 MILLION TWEETS
are sent per day by about 200 million monthly active users



1 IN 3 BUSINESS LEADERS

don't trust the information they use to make decisions



Poor data quality costs the US economy around

\$3.1 TRILLION A YEAR



27% OF RESPONDENTS

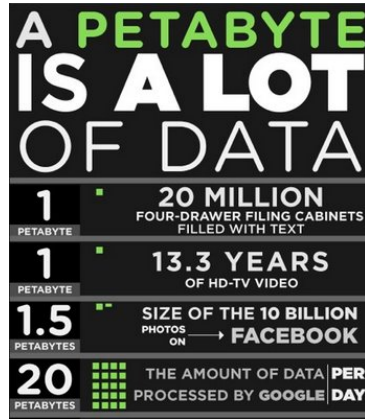
in one survey were unsure of how much of their data was inaccurate

Veracity UNCERTAINTY OF DATA

Data, data, everywhere... Volume

The Petabyte Age

- 2008



- Eric Schmidt (in 2010): Every 2 Days We Create As Much Information As We Did Up To 2003



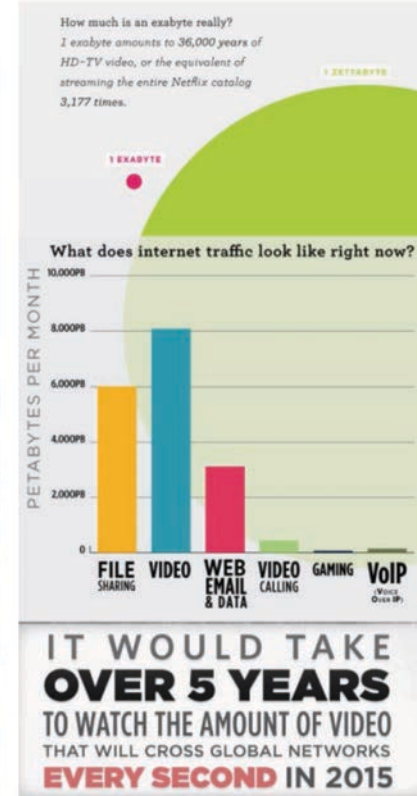
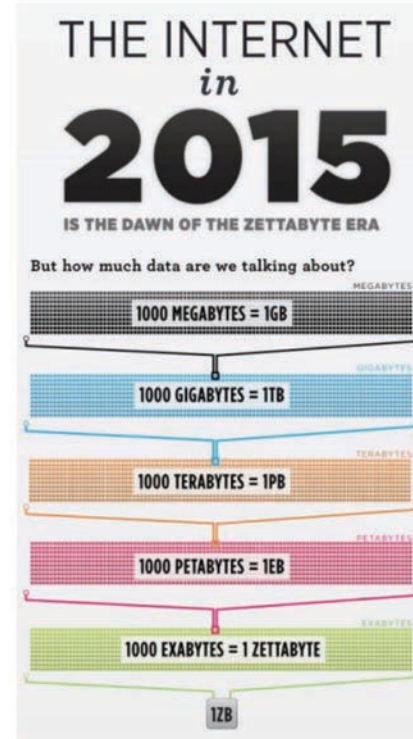
The Zettabyte Age

- 2015
- One Zettabyte = Stack of books from Earth to Pluto 20 times



The internet in 2020

- ~26.3 billion networked devices
- 25.1 GB average traffic per capita per month
- 2.3 Zettabytes annual IP-Traffic



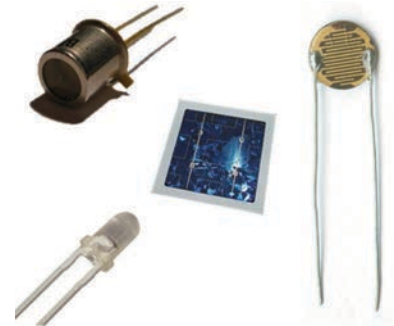
Data is produced continuously Velocity

Humane-produced data

- ~300 million email sent/received per minute in 2016
- >0.4 million tweets per minute in 2016
- ~138 million google searches per minute in 2016
- >400 hours of video was uploaded to YouTube per minute in 2016

Machine produced data

- IoT will be boost data velocity greatly
- Sensors become ubiquitous
- Sensors for sound, images, position, motion, temperature, pressure, etc ...
- Resolution (in time and space) is continuously increasing
- Assume Waze like cars collecting 20 double values every second
- With one million driving cars that is almost 10 TB every minute



- Structure vs. unstructured
- Text vs. image
- Curated vs. automatically collected
- Raw vs. edited vs. refined

- Decentralized content generation
- Multiple perspectives (conceptualizations) of the reality
- Ambiguity, vagueness, inconsistency

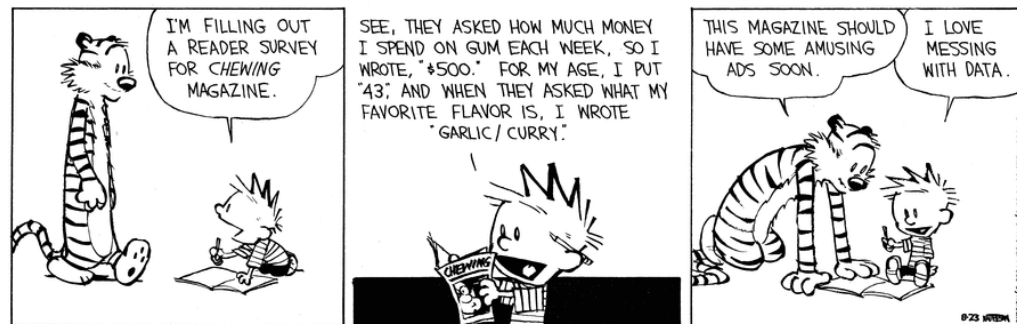


“A lot of Big Data is a lot of small data put together.”

Data is messy Veracity

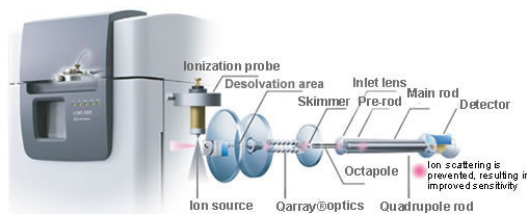
Quality of captured data varies greatly

- Sensor inaccuracy
- Human mistakes
- Incompleteness
- Untrusted sources
- Deterministic processes
- etc.

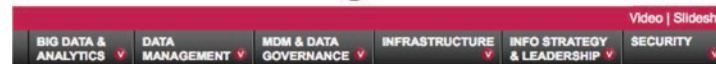


“Next to the analytes, we see everything in the results, from the perfume of the lab assistant to the softener in the new machine sitting next.”

—About Liquid chromatography–mass spectrometry at IPB Halle



information management



GET BREAKING NEWS TO YOUR INBOX PLUS MORE EXCLUSIVE BENEFITS! BECOME A REGISTERED MEMBER

NEWS

Messy Big Data Overwhelms Data Scientists

by BOB VIOLINO
FEB 20, 2015 2:00pm ET

Print

Email

Reprints

Comments (2)

Data scientists see messy, disorganized data as a major hurdle preventing them from doing what they find most interesting in their jobs: predictive analysis and data mining for behavioral patterns and future trends, according to a new report from CrowdFlower, a data enrichment platform provider.

A majority of the 153 CrowdFlower online research panel members surveyed (80%) also acknowledged the skills shortage within their field. The respondents work for companies of varied sizes and sectors, mostly in the U.S. All respondents have the term "data scientist" in their job title or job description on LinkedIn, CrowdFlower says.

Value - the fifth
V of Big Data

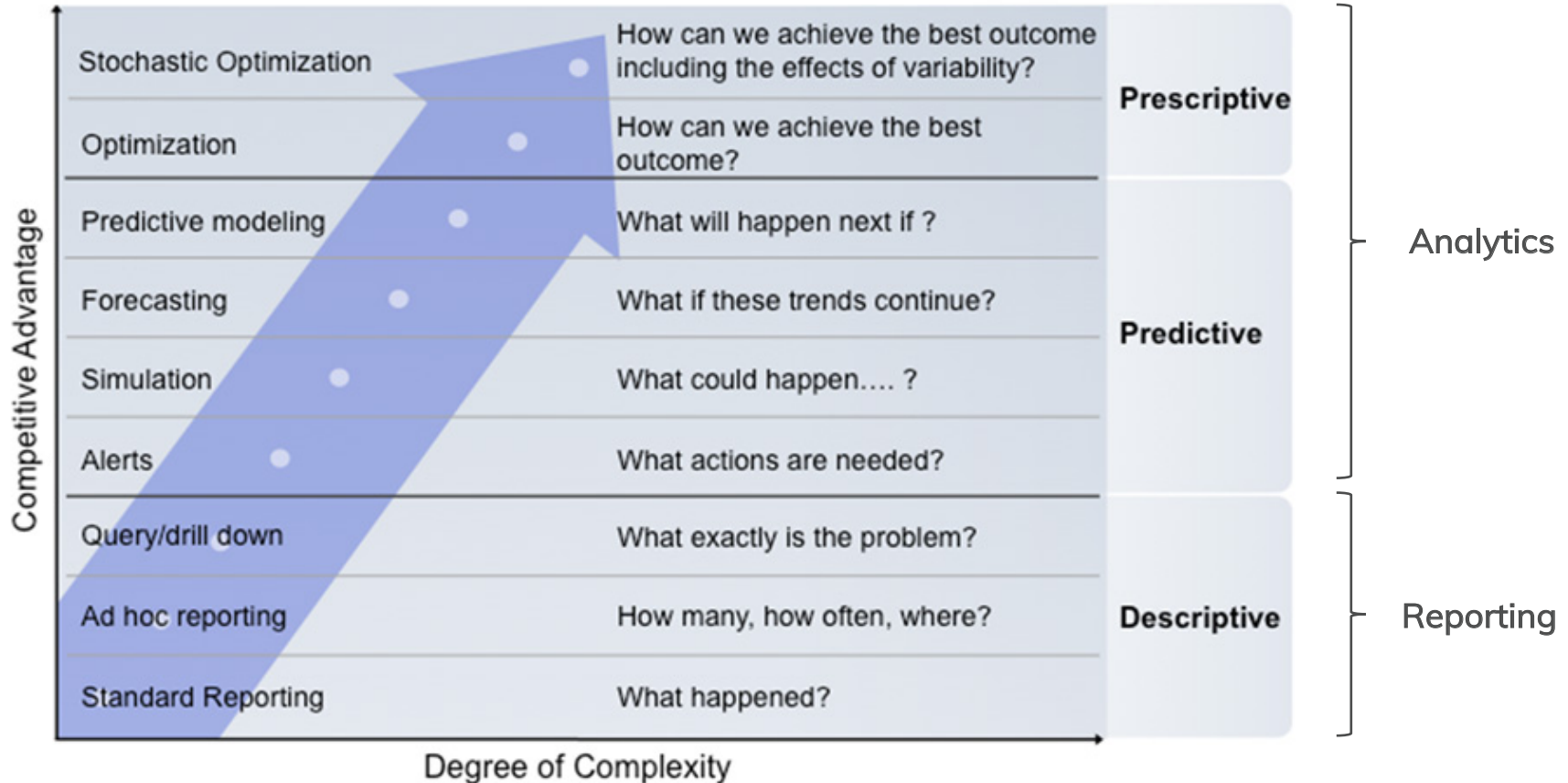


Data Science/Data Analysis

... or how to turn raw data into something valuable?

“Data is the new oil. It’s valuable, but if unrefined it cannot really be used. It has to be changed into gas, plastic, chemicals, etc. to create a valuable entity that drives profitable activity; **so must data be broken down, analyzed for it to have value.**” –Clive Humby

Levels of Analysis



Descriptive Analytics

How have I done? (and why?)

- Simplest class of analytics
- Condense big data into smaller, more useful bits of information
- Summary of what happened
- 70-80% penetration

Examples

- Database aggregation queries
- Business reporting (e.g. Sales figures)
- Market survey (e.g. GFK)
- (classical) business intelligence, dashboards, scorecards
- Google Analytics



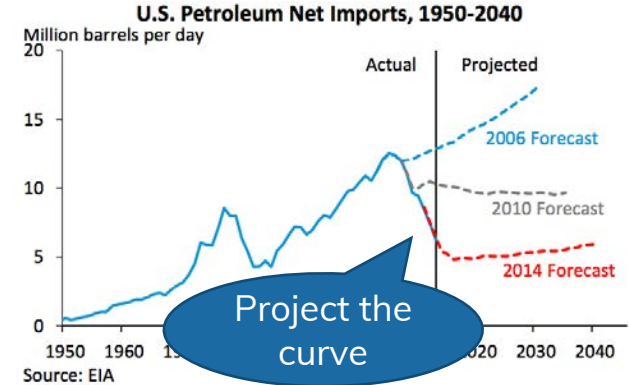
Predictive Analytics

How will I do?

- Next step up in data reduction
- Studies recent and historical data
- Utilizes a variety of statistical, modeling, data mining and machine learning techniques
- Allows (potential inaccurate) predictions about the future
- Use data you have, to create data you do not have
- 15-25% penetration

Examples

- Market developments
- Stock developments
- Movie/product recommendations on netflix/amazon
- Energy demand and supply forecasting
- Predictive policing (e.g. precobs, predpol)



			
Alice	?	★★★★★	★★
Michael	★★★★★★	?	★★★★★

Fill in the blanks

Prescriptive Analytics



ORION –
On-Road Integrated
Optimization and Navigation



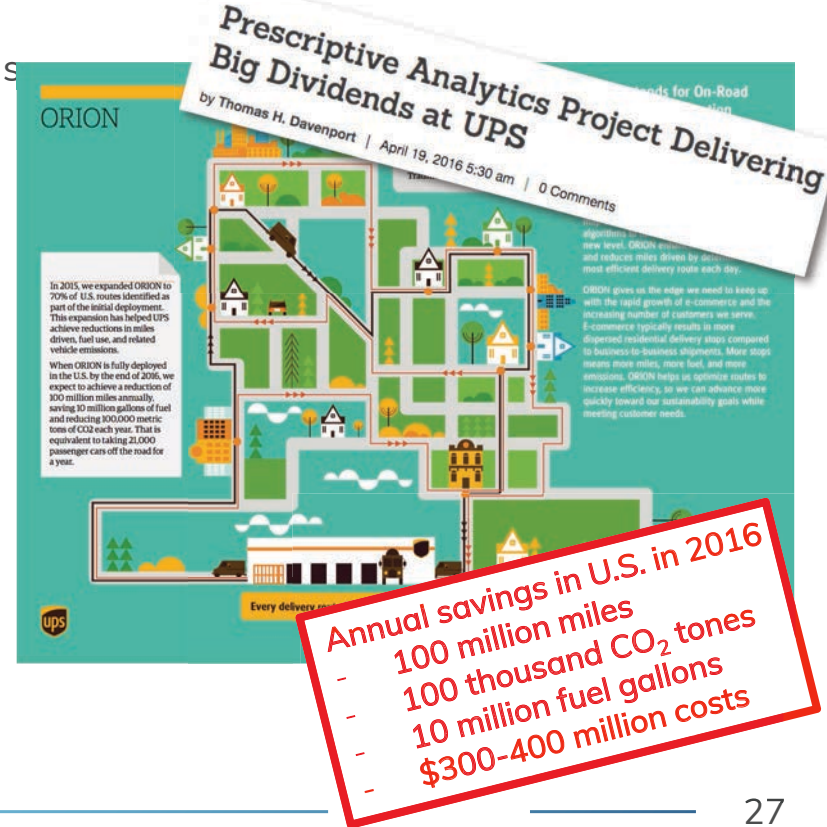
Dresden Database
Systems Group

What should I do?

- Predicts “multiple futures” based on the potential actions
- Recommends the best course of action for any pre-specified outcome
- Typically involves a feedback system to track outcome produced by the action taken
- Utilizes predictive methods + optimization techniques
- 1-5% penetration

Examples

- Energy load balancing by flexoffer scheduling
- Inventory optimization in supply chains
- Targeted marketing campaign optimization
- Focus treatment of clinical obesity in health care
- Waze-like car navigation



Roles in Big Data Projects

Data scientist

- Data science is a systematic method dedicated to knowledge discovery via data analysis
- In business, optimize organizational processes for efficiency
- In science, analyze experimental/observational data to derive results
- Typical skills
 - Statistics + (mathematics) background
 - Computer science: Programming, e.g.: R, (SAS,) Java, Scala, Python; Machine learning
 - Some domain knowledge for the problem to solve

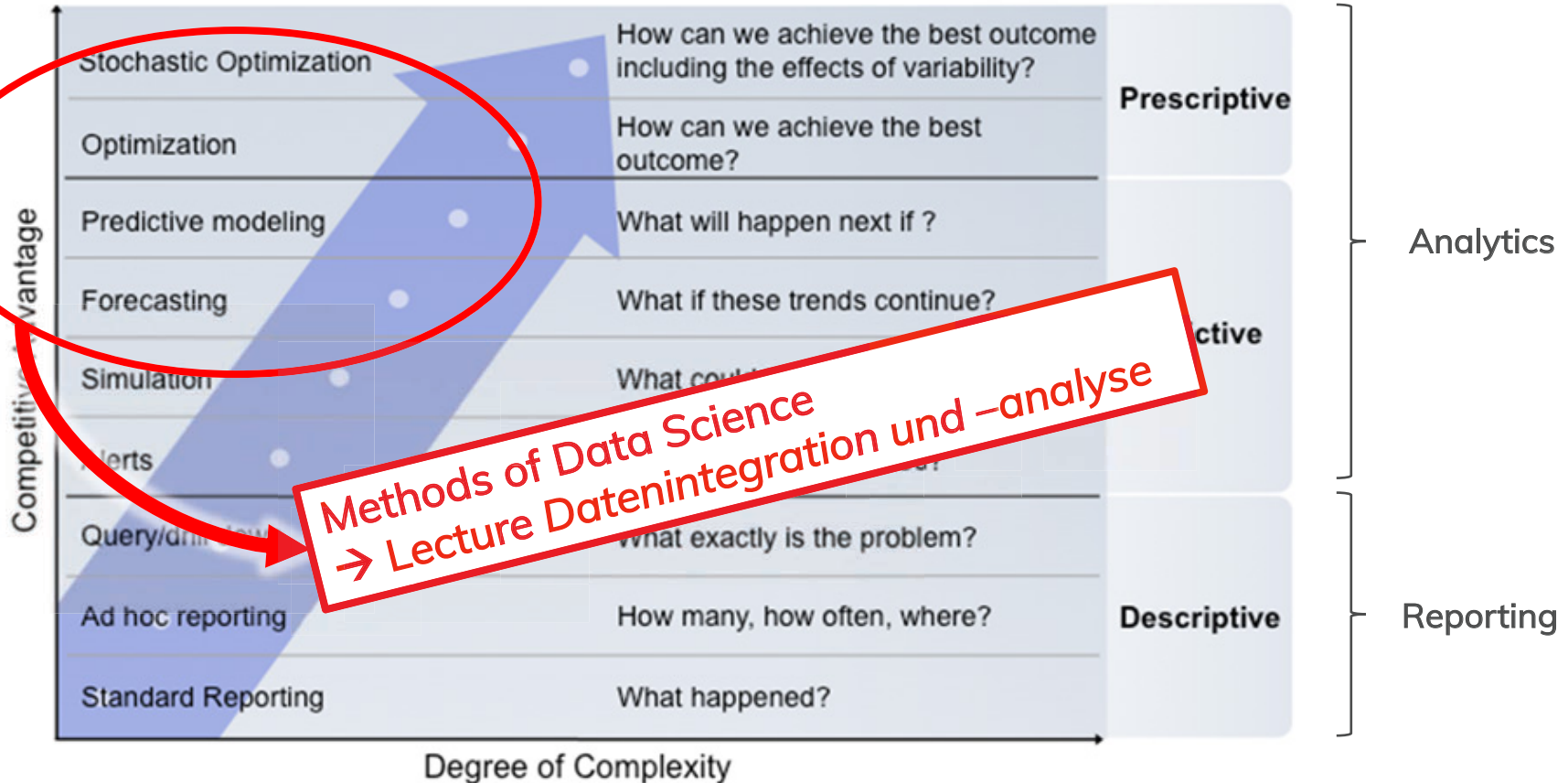
Data engineer

- Data engineering is the domain that develops and provides systems for managing and analyzing big data
- Build modular and scalable data platforms for data scientists
- Deploy big data solutions
- Typical skills
 - Computer science background
 - Databases
 - Software engineering
 - Massively parallel processing
 - Real-time processing
 - Languages: C++, Java, (Scala,) Python
 - Understand performance factors and limitations of systems



Our Focus in Teaching and Research

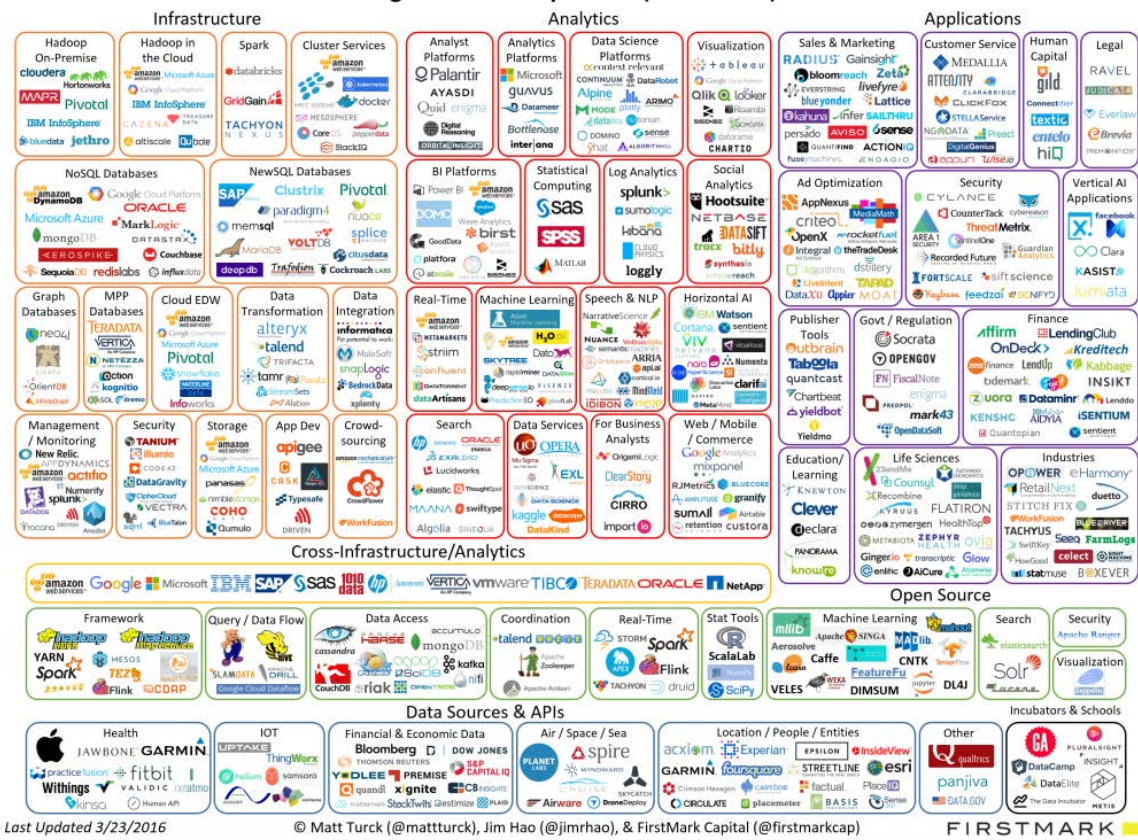
Levels of Analysis



Big Data Landscape(s)



Big Data Landscape 2016 (Version 3.0)



Last Updated 3/23/2016

© Matt Turck (@mattturck), Jim Hao (@jimhao), & FirstMark Capital (@firstmarkcap)

FIRSTMARK

Data Systems

data systems are in the middle of all this

a data system...

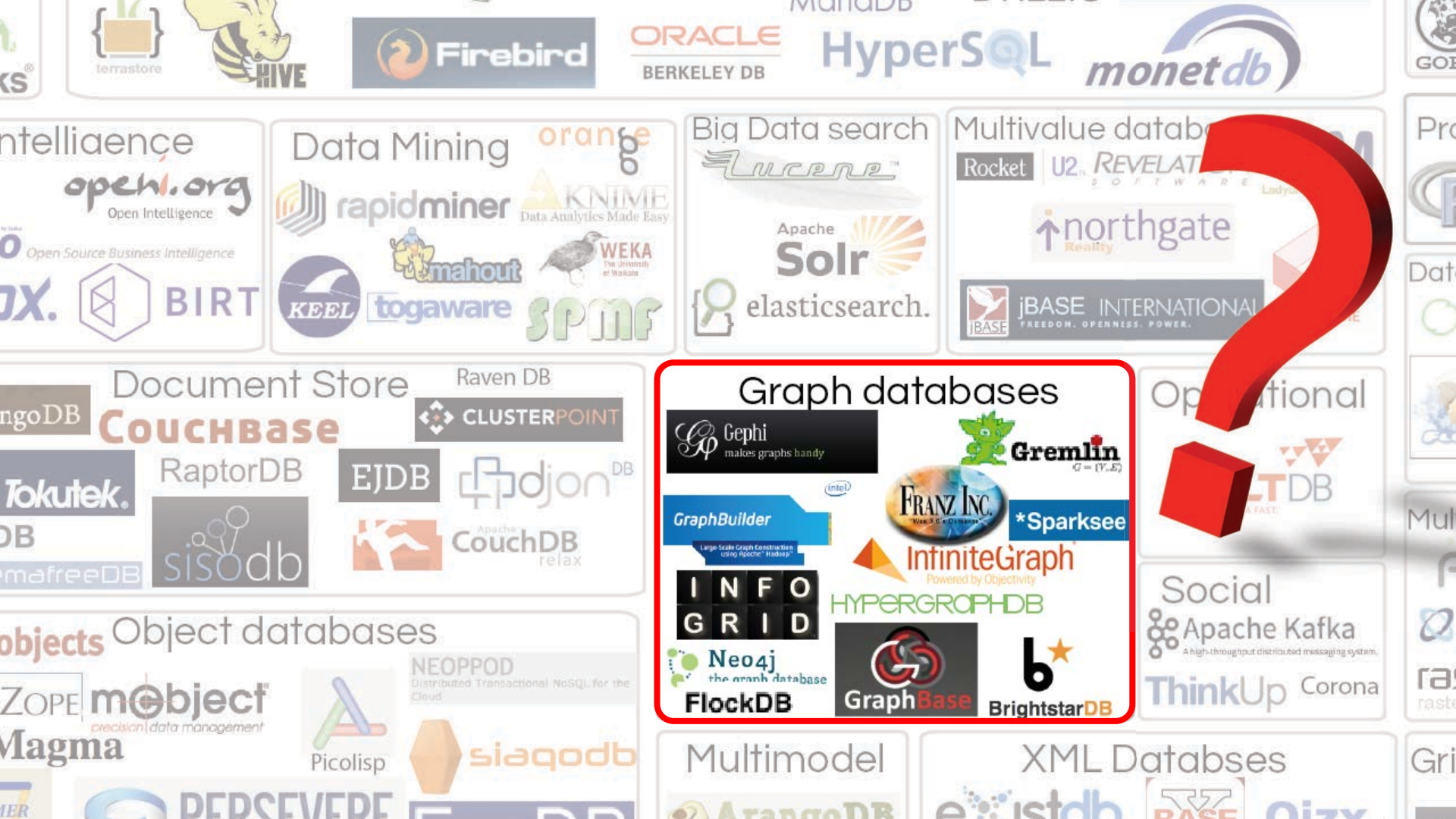
- ...stores data...
- ...provides access to data...
- ...and (ideally) makes data analysis easy

Different data systems use different data models



(Big) Data System Landscape(s)





Intelligence

openi.org
Open Intelligence

Open Source Business Intelligence

BIRT

Data Mining

orange

rapidminer

KNIME
Data Analytics Made Easy

mahout

WEKA
The University of Waikato

KEEL

togaware

SPMF

Big Data search

Lucene

Apache Solr

elasticsearch.

Multivalue database

Rocket

U2

REVELAT
SOFTWARE

northgate
Reality

jBASE INTERNATIONAL
FREEDOM. OPENNESS. POWER.

Document Store

CouchBase

Raven DB

CLUSTERPOINT

RaptorDB

EJDB

djon DB

CouchDB relax

sisodb

Graph databases

Gephi
makes graphs handy

Gremlin
G = {V, E}

GraphBuilder
Large Scale Graph Construction using Hadoop, MapReduce

FRANZ INC.
Mark 3 G4 GraphDB

***Sparksee**

InfiniteGraph
Powered by Objectivity

HYPERGRAPHDB

INFO GRID

Neo4j
the graph database

FlockDB

GraphBase

BrightstarDB

Object databases

NEOPPOD
Distributed Transactional NoSQL for the Cloud

siaqodb

Picolisip

PERCEVEDE

Multimodel

ArangoDB

XML Databases

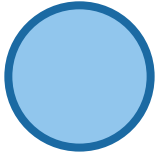
istdb

BASE

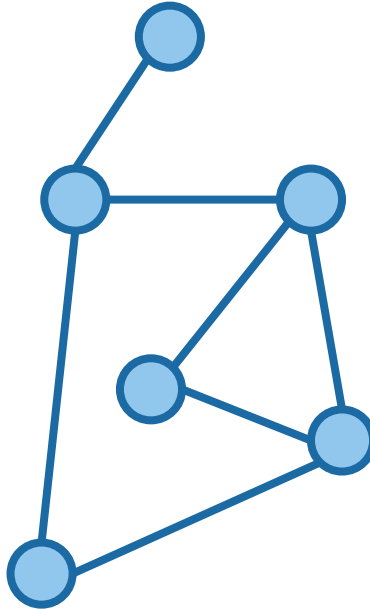
oizy

Graph Building Blocks

Nodes (Dots)



- Like an entity in ER
- Exist on their own
- Have object identity

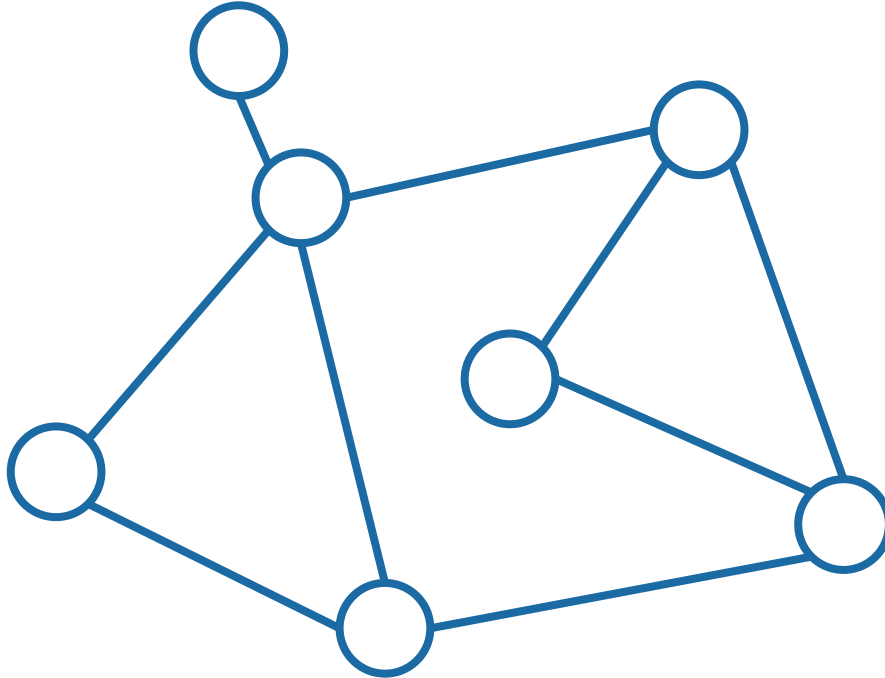


Edges (Lines)

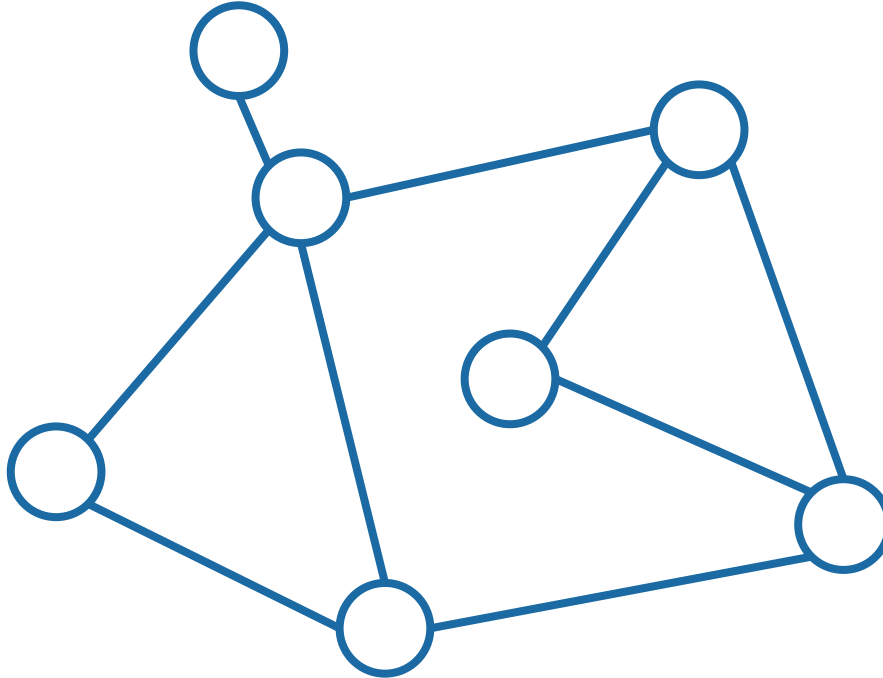


- Like a relationship in ER
- Exist only between nodes
- Identity depends on the nodes they connect

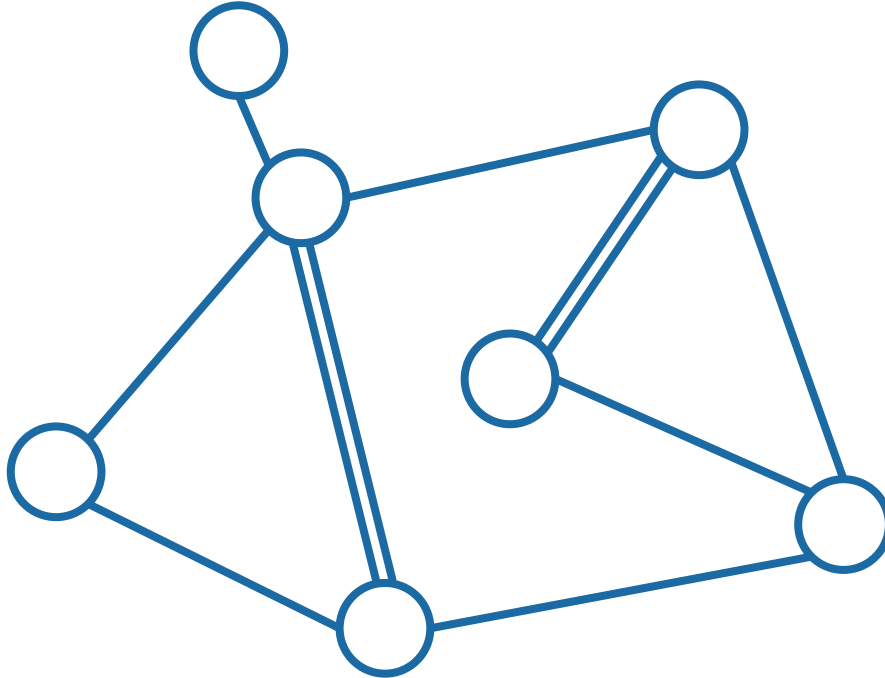
Graph Data – Social Network



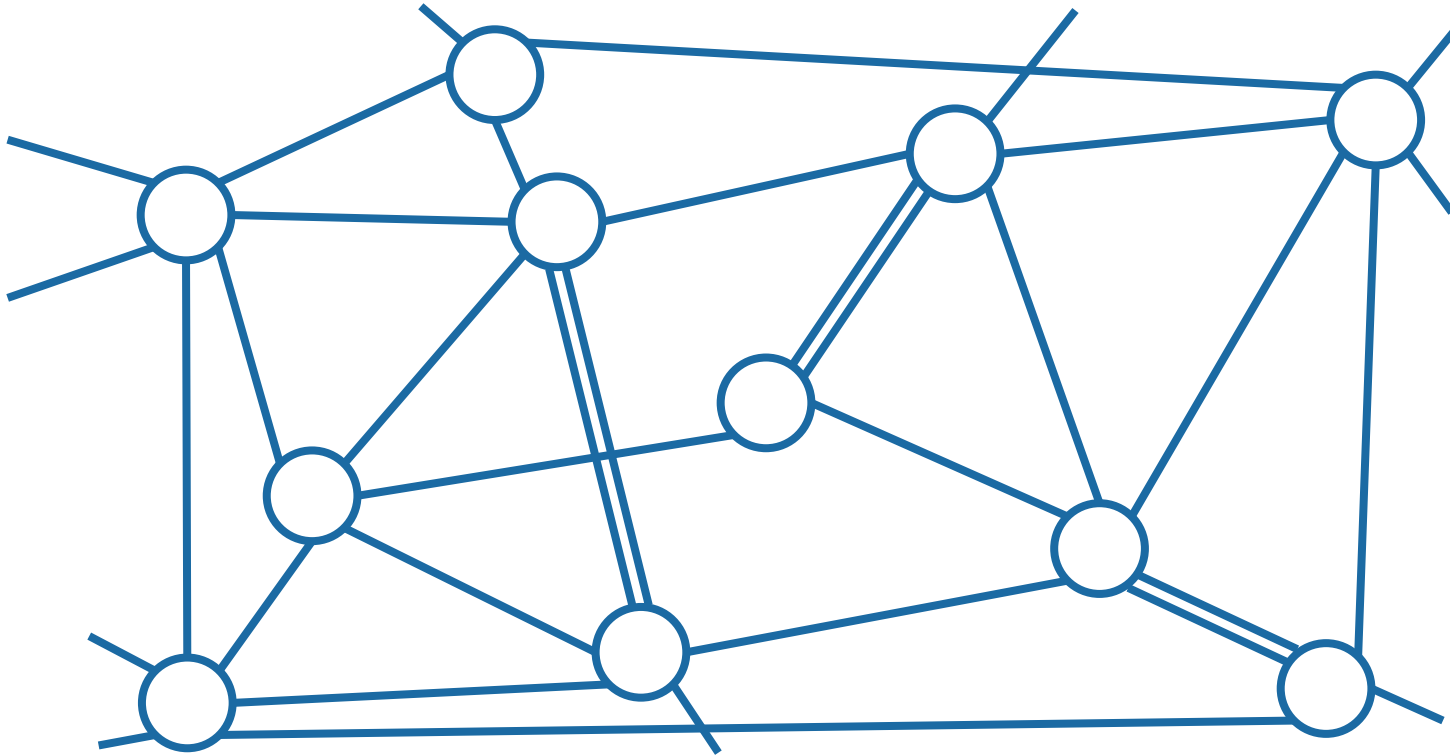
Graph Data – Social Network



Graph Data – Social Network



Graph Data



Social Graphs

Facebook

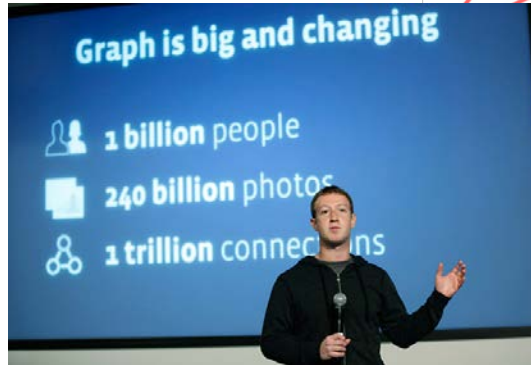
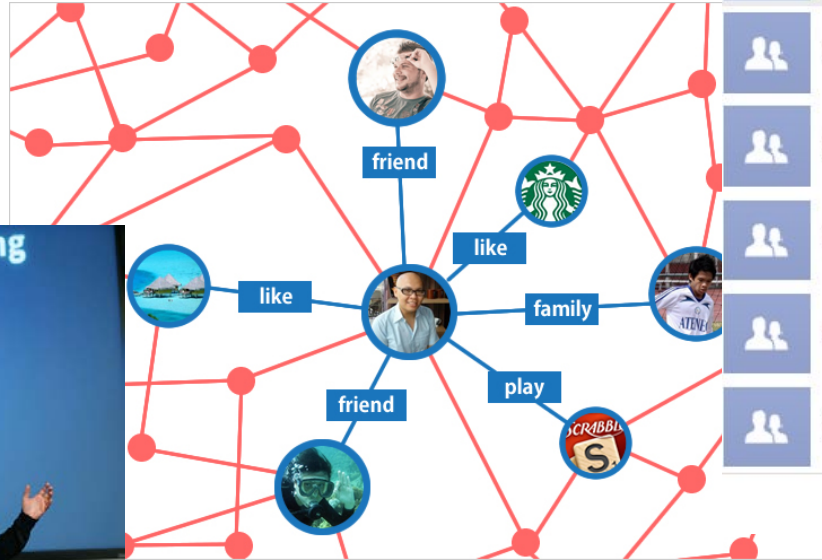
- May 2013



Structured Search

Example: Facebook Graph Search

- Finding subgraph structures
- Very natural way of formulating queries



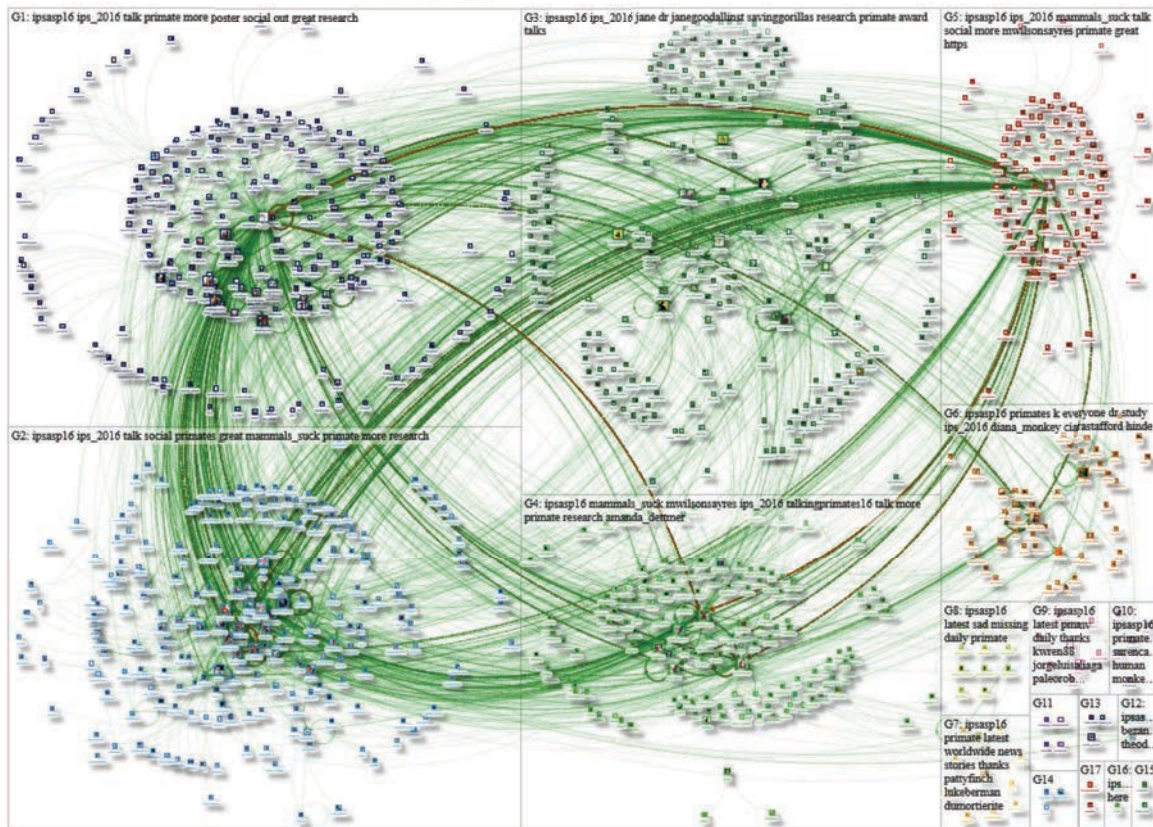
[<http://socialnewsdaily.com/15865/facebook-social-graph-search-a-great-way-to-find-working-professionals-in-your-network/>]

Social Network Analysis



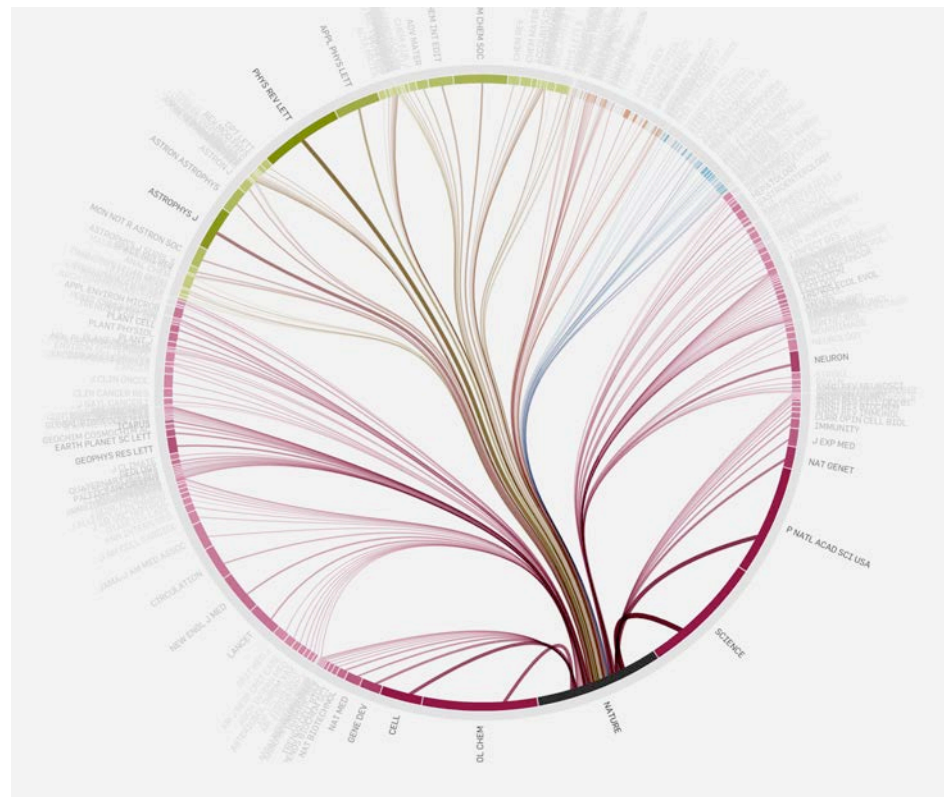
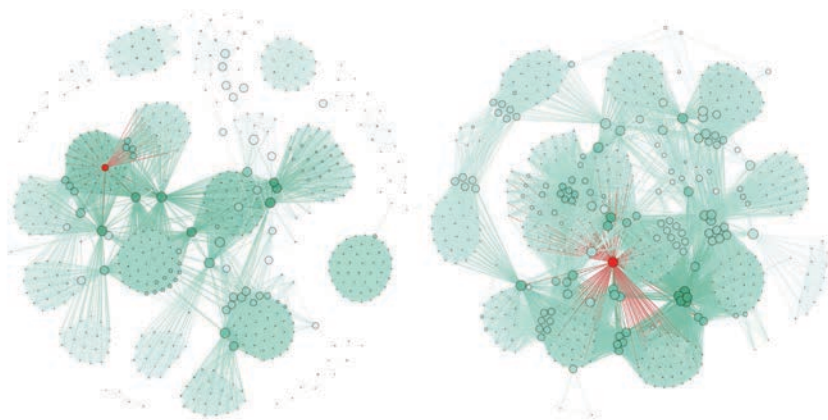
Example: Twitter communication

- Users tweeting on a specific topic
- Others reply of retweet
- Users can be grouped based on communication topology (-> graph clustering)
- Analysis reveals user groups and dominant communication patterns



[<https://nodexlgraphgallery.org/Pages/Graph.aspx?graphID=76277>]

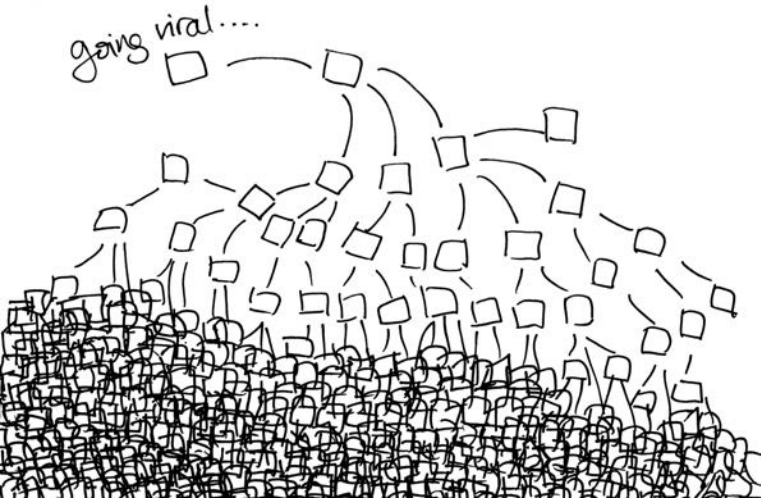
- Open bibliographic information on major computer science journals and proceedings
- >3.4 million publication
- >7000 new publication per month
- >1.7 million authors



Viral Marketing

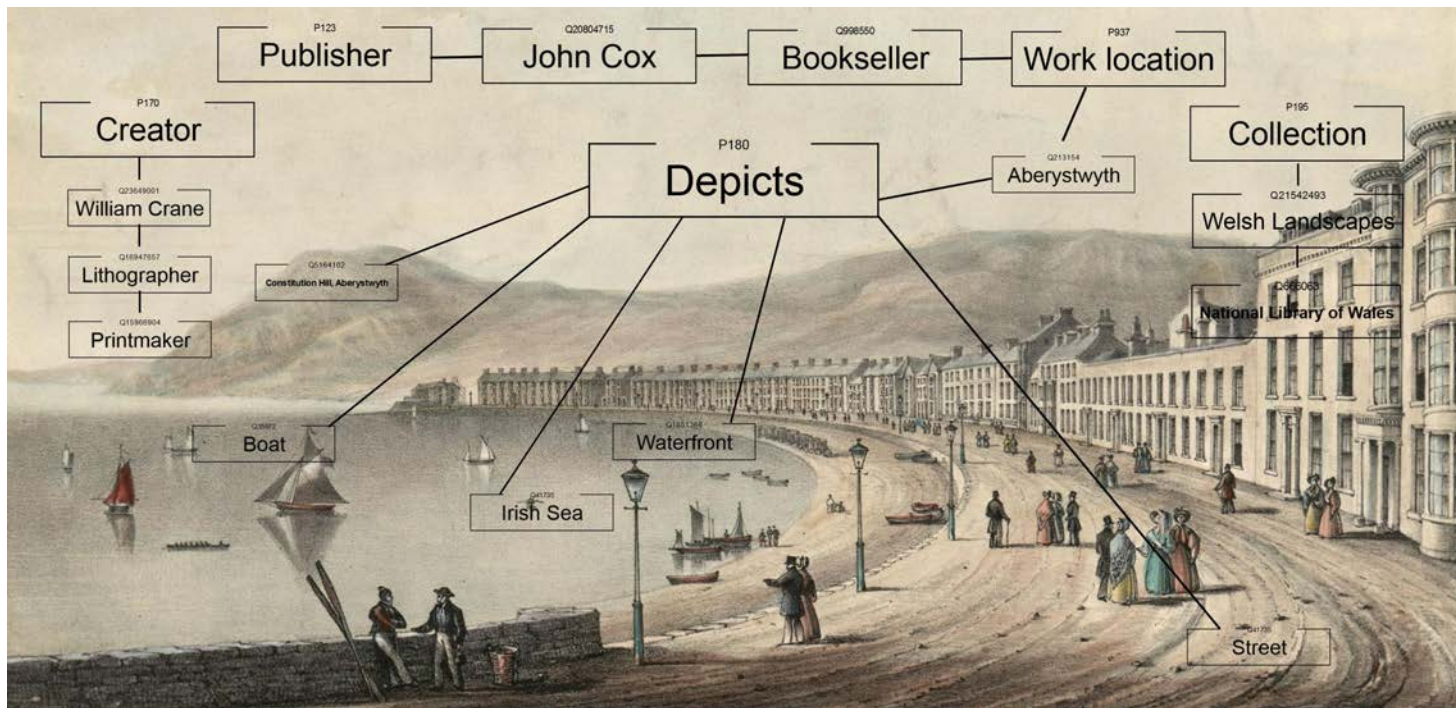
Viral Marketing

- spreading content to one person so that more than one person engaging with the content
- Techniques
 - Influence estimation
 - Influence maximization



Knowledge Graphs

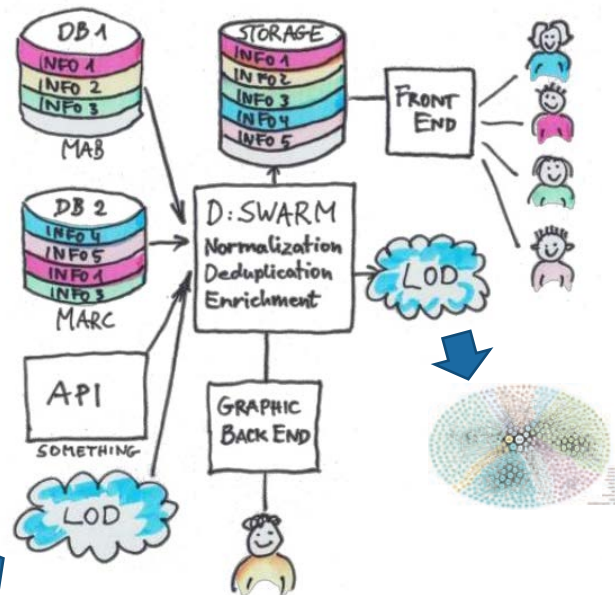
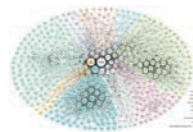
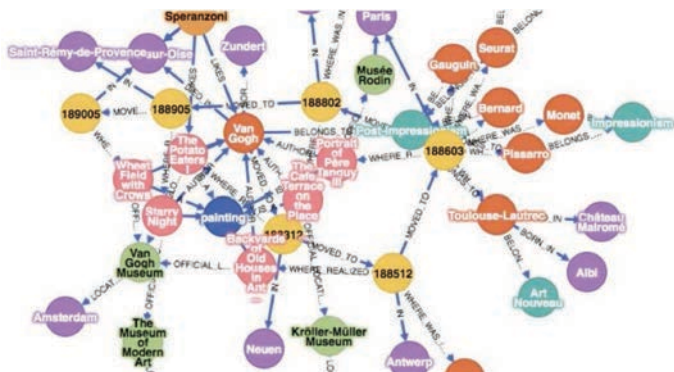
Knowledge graph of a picture



[The National Library of Wales, <https://www.llgc.org.uk/blog/?p=11246>]

Sächsische Landesbibliothek – Staats- und Universitätsbibliothek Dresden (SLUB)

- Adds semantics search to library online catalog
- Utilizing multi-lingual knowledge data from Wikipedia
- Significant improvements in search quality for library users



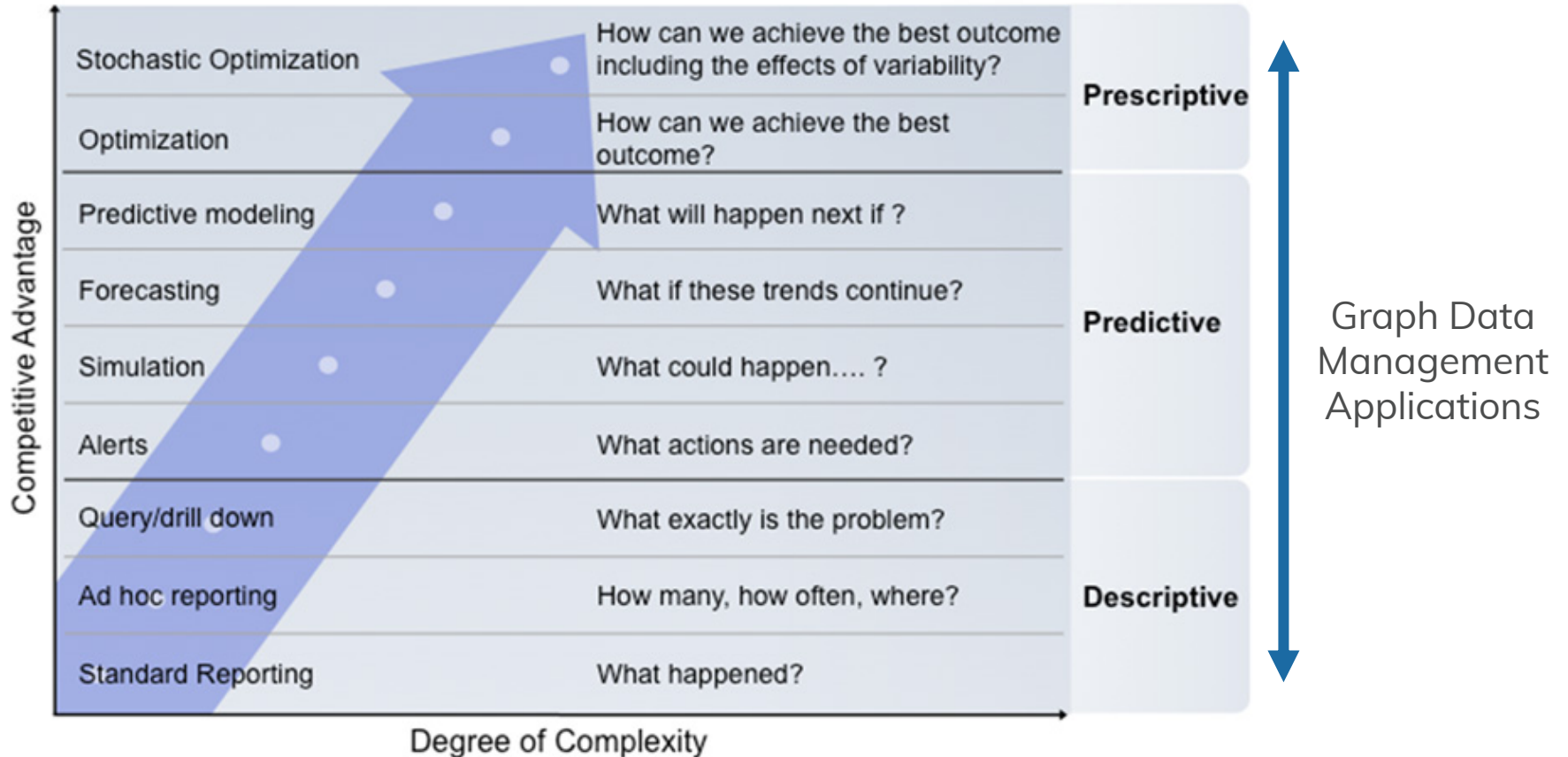
Supply Chain Management



Supply Chain Optimization

- Customer A: 10-30% reduction in inventory
- Customer B: 8% reduction in transportation costs

Level of Analytics



The Power of Networks

CognitiveMedia • Manuel Lima, The Power of Networks, London 8.12.2011

Lecture given at the RSA, London on the 8th December, 2011. Animated May, 2012
Manuel Lima explores the way in which we are all connected to the world around us, and how this notion has adapted over the years. From trees of knowledge to the Web of Life to modern social networking, Manuel unravels all the ways life is linked in this RSA Animate.



Watch on YouTube: <https://youtu.be/nJmGrNdJSGw>

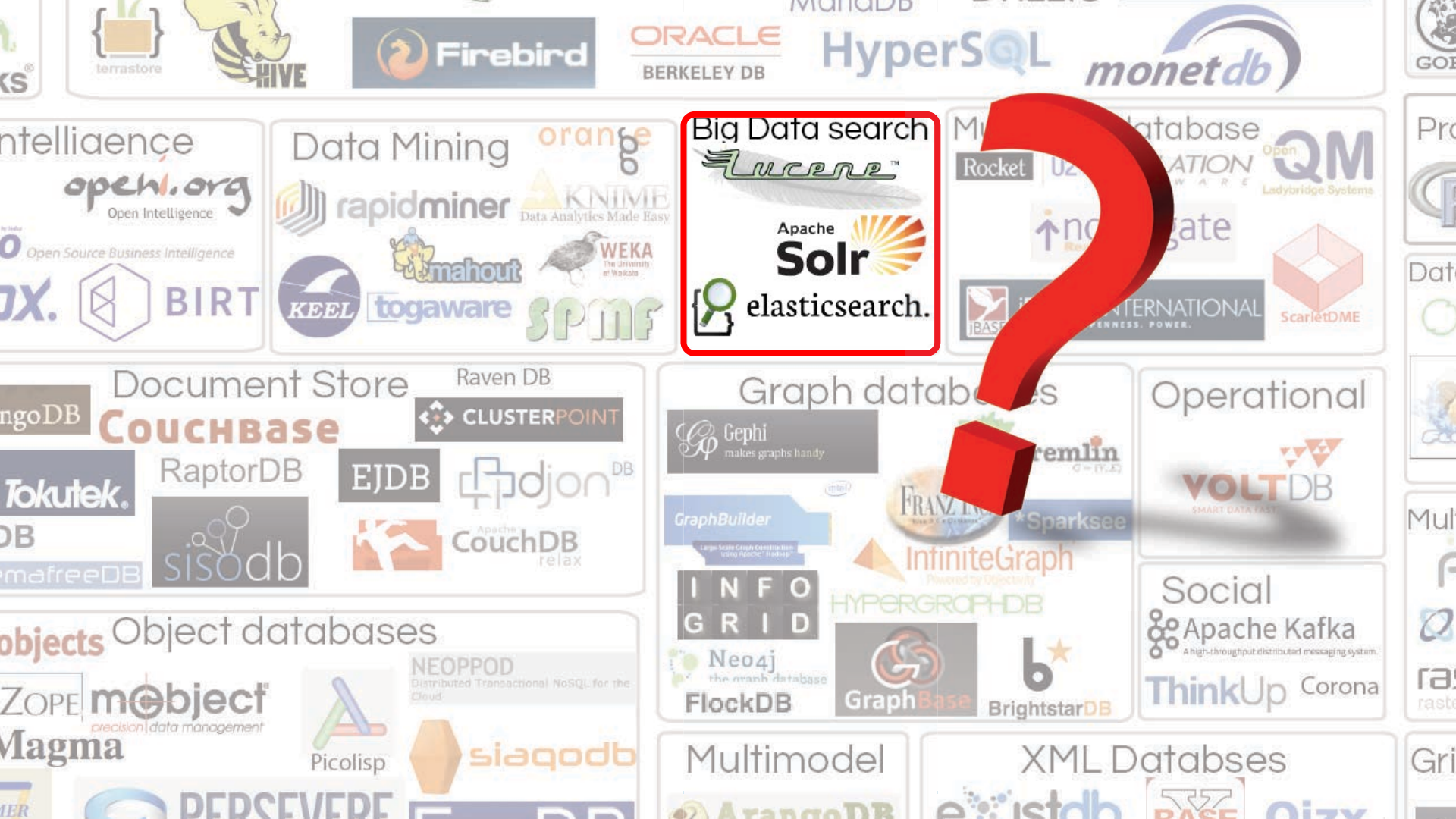
Slide: +46(0)1803 253689 contact@cognitivemedia.co.uk www.cognitivemedia.co.uk

→ Lecture Graph Data Management and Analytics

Graph databases

The collage includes the following logos and text:

- Gephi: makes graphs handy
- Gremlin: $G = (V, E)$
- GraphBuilder: Large Scale Graph Construction using Hadoop, MapReduce
- Franz Inc. (Mark 3 G4 GraphDB)
- *Sparksee
- InfiniteGraph: Powered by Objectivity
- InfoGrid
- HypergraphDB
- Neo4j: the graph database
- FlockDB
- GraphBase
- BrightstarDB



Use Cases of Search

Searching Textual Contents

+Ich Gmail Bilder  [Anmelden](#)




Google-Suche


Auf gut Glück!



Use Cases of Search

Stack Overflow → Programming Q & A site

 **stackoverflow** Questions Tags Users Badges Unanswered Ask Question

Suppress install outputs in R

 **CAREERS 2.0**
by stackoverflow


 +  Have projects on Google Code?
Import them easily to your profile

▲ This is really starting to bug me...I have tried a few methods and none seem to work

6 I am running an install from a function which generates a lot of unnecessary messages that I would like to suppress, but all of the methods I tried to do this have not worked.


★ The bit of code I am trying to suppress is : `install_github('ROAuth', 'duncantl')`, it requires the package `devtools` to be loaded beforehand.

1 Anyway, I tried `invisible`, `capture.output` and `sink`, none of which work...or perhaps I am not using them correctly... either way...any ideas?

 `suppressmessage`

share | improve this question


asked Sep 13 at 23:09

 **h.l.m.**
358 • 7
70% accept rate

3 perhaps `suppressMessages()` or `suppressPackageStartupMessages()` is what you want? – Chase Sep 13 at 23:11

1 @Chase is right. Your function in the other question is a bit convoluted and should not be calling `install_github()` every time. See my answer there. – MaiaSaura Sep 13 at 23:50

tagged

 × 17244

`suppressmessage` × 13

asked 8 days ago


viewed 146 times

active 2 days ago

Community Bulletin

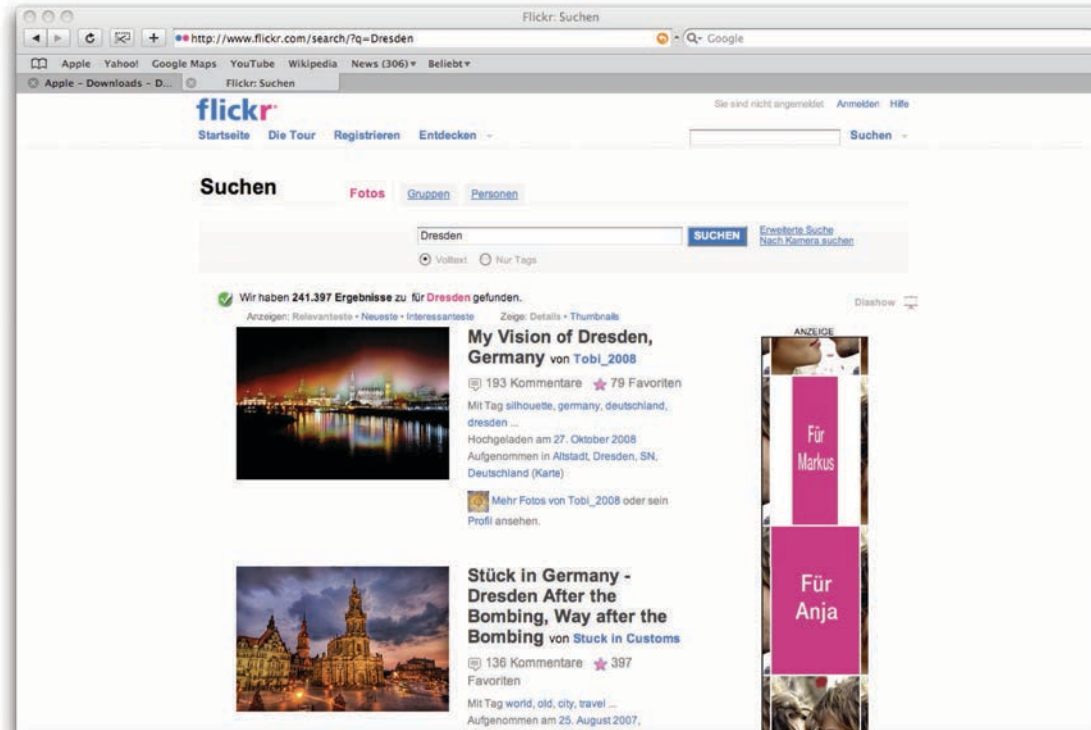
blog [AskPatents.com: A Stack Exchange To Prevent Bad Patents](#)

Better Jobs.
Better Pay.

 **CAREERS 2.0**
by stackoverflow

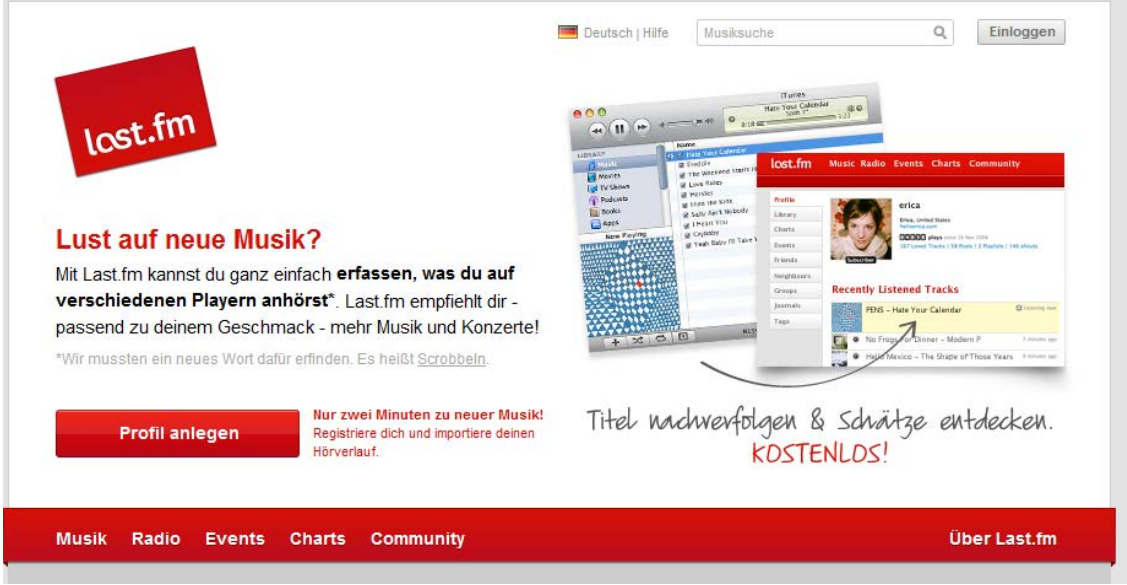
Use Cases of Search

Image Search



Use Cases of Search

Last.fm → Music Search



lost.fm

Deutsch | Hilfe Musiksuche

Lust auf neue Musik?

Mit Last.fm kannst du ganz einfach **erfassen, was du auf verschiedenen Playern anhörst***. Last.fm empfiehlt dir - passend zu deinem Geschmack - mehr Musik und Konzerte!

*Wir mussten ein neues Wort dafür erfinden. Es heißt Scrobblen.

Profil anlegen Nur zwei Minuten zu neuer Musik!
Registriere dich und importiere deinen Hörverlauf.

Titel nachverfolgen & Schätze entdecken.
KOSTENLOS!

Musik Radio Events Charts Community Über Last.fm

Example: Library of Congress

FACTS AT A GLANCE

In fiscal year 2008, the Library of Congress:

- Welcomed more than 1.6 million onsite visitors
 - 1,207,776 moving images
- Provided reference services to 545,084 individuals in person, by telephone and through written and electronic correspondence
 - 12,536,764 photographs
 - 98,288 posters
- Recorded a total of 141,847,810 items in the collections:
 - 545,347 prints and drawings
- 21,218,408 cataloged books in the Library of Congress classification system
- 11,599,606 books in large type and raised characters, incunabula (books printed before 1501), monographs and serials, music, bound newspapers, pamphlets, technical reports and other printed material
- Circulated more than 22 million disc, cassette and braille items to more than 500,000 blind and physically handicapped patrons
- Registered 232,907 claims to copyright
- Completed 871,287 research assignments for Congress through the Congressional Research Service
- Prepared 1,529 legal research reports for Congress and other federal agencies through the Law Library
- Recorded more than 85 million visits and 610 million page views on the Library's website. At year's end, the Library's online historical collections contained 15.3 million digital files
- Employed a permanent staff of 3,637 employees
- Operated with a total fiscal 2008 appropriation of \$613,496,414, including the authority to spend \$50,447,565 in receipts
- 109,029,796 items in the nonclassified (special) collections, including:
 - 3,005,028 audio materials, such as discs, tapes, talking books and other recorded formats
 - 62,778,118 manuscripts
 - 5,357,385 maps
 - 16,086,572 microforms
 - 5,674,956 pieces of printed sheet music
 - 14,388,175 visual materials, as follows:

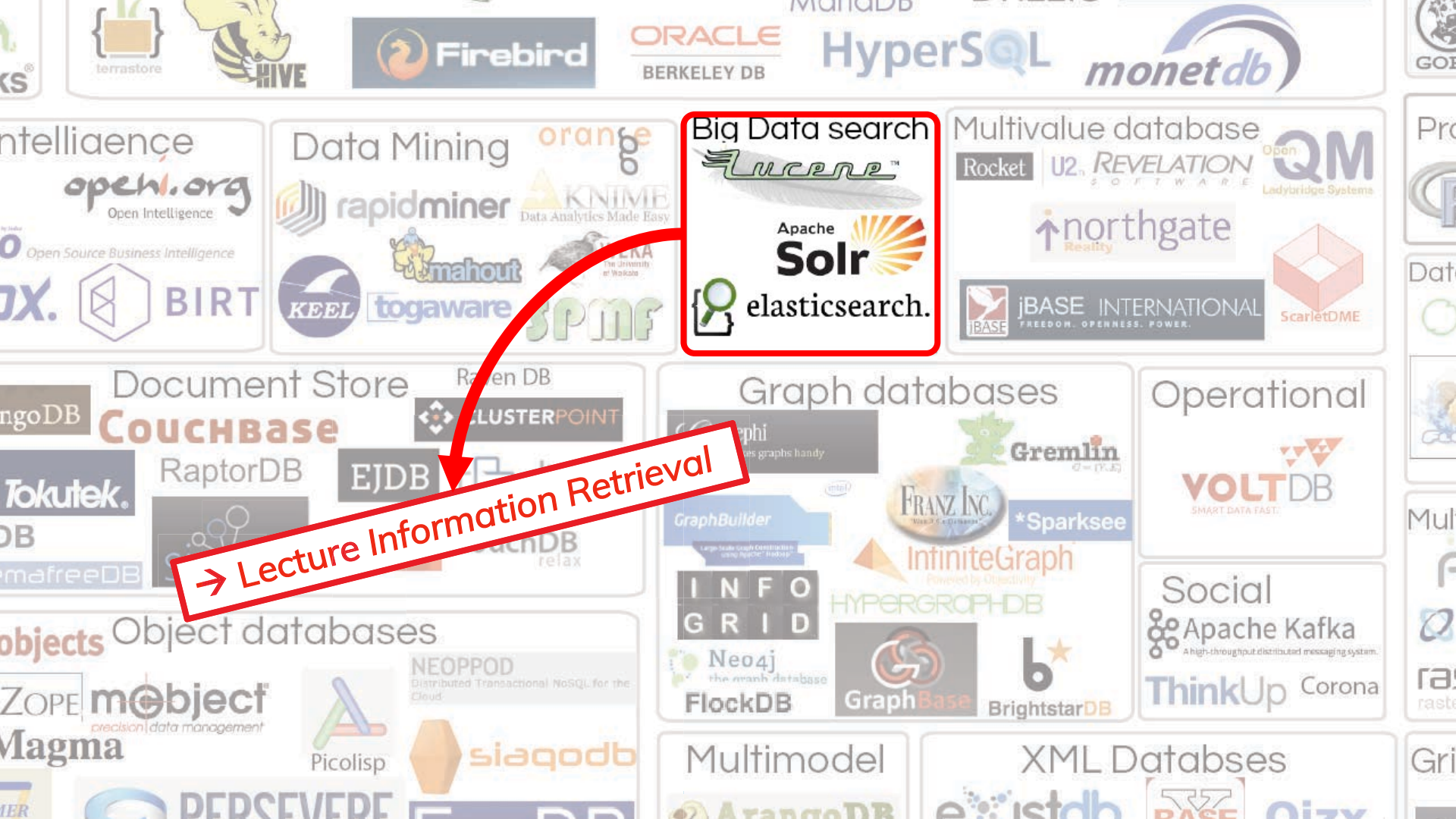


*The Great Hall
of the
Thomas Jefferson Building,
Library of Congress*

Databases versus Information Retrieval

	Database	Information Retrieval
Matching	Exact Match	Partial Match, Best Match
Model	Deterministic	Probabilistic
Query Language	Structured / Formal	Natural
Query Specification	Complete	Incomplete
Queried Objects	Matching	Relevant
Error Sensitivity	Sensitive	Insensitive

- Hard to formulate Queries
- Iterative workflow base on reponses
- Tons of results, but only a few are relevant
- Ranking of results (instead of set of results)
- Representation of document content often inadequate / inaccurate



→ Lecture Information Retrieval

Summary

Big Data

- Crossing thresholds in exponential growth & digitization
- Technical challenges in volume, velocity, variety, veracity, value

Related research and lectures

- Data Science → Datenintegration und -analyse
- Data Systems → Architektur von Datenbanksystemen
→ Big Data Platforms
- Graph Data → Graph Data Management and Analytics
- Search → Information Retrieval

