



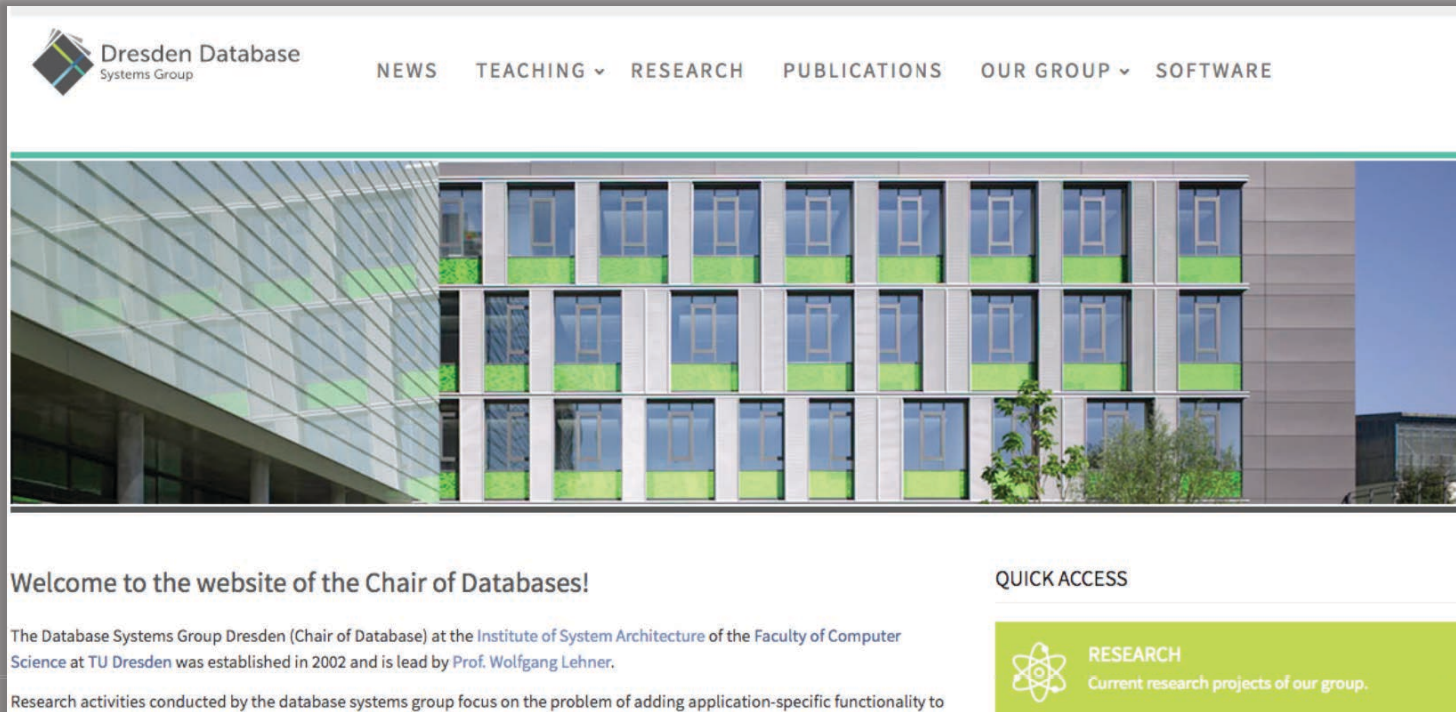
# Future Trends in Data Analytics

Hannes Voigt

# General Information

## LECTURE NOTES

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



The screenshot shows the website's header with the logo and navigation menu. Below the header is a large photograph of a modern building with a grid of windows. The main content area contains a welcome message and a quick access section.

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



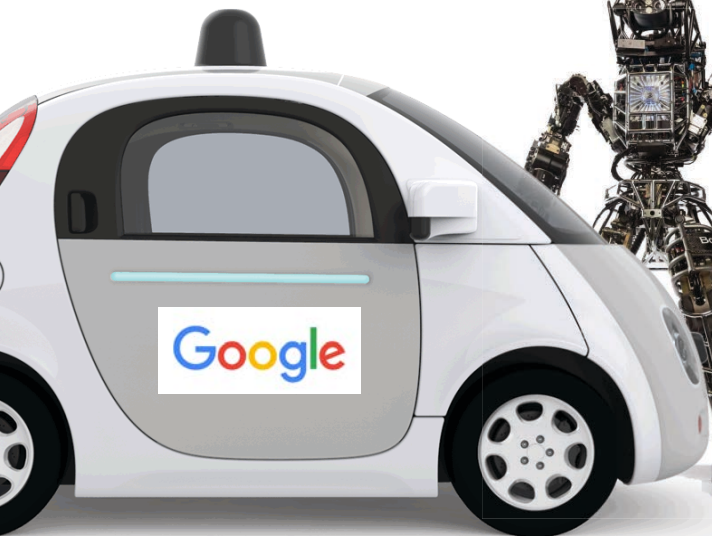
# Future is Now

JEOPARDY!

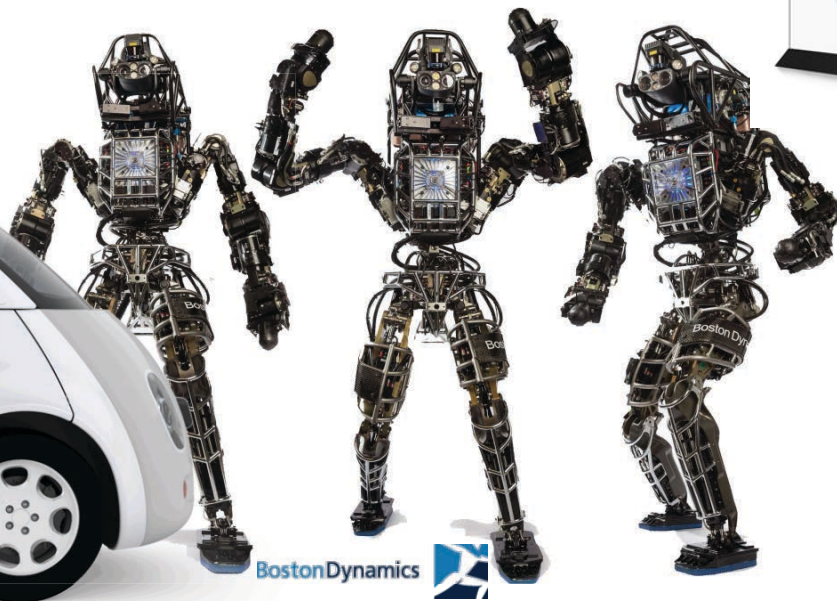
- ..., face detection & recognition, autobeam, Siri/Cortant/..., 3D printing, ...



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



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



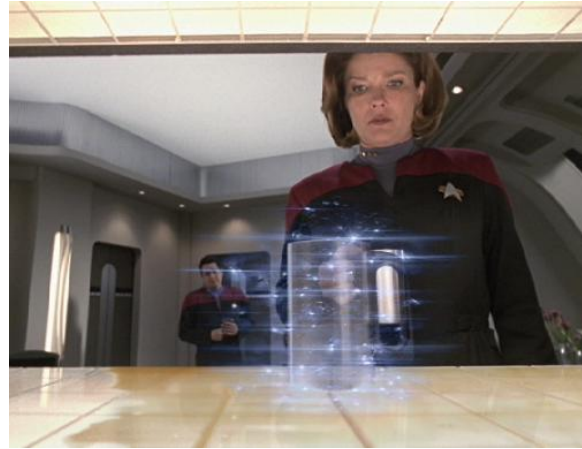
[<http://www.bostondynamics.com/>]



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

... tomorrow?

SOONER THEN YOU THINK!



When have we  
entered the future?

What was the threshold  
we crossed?

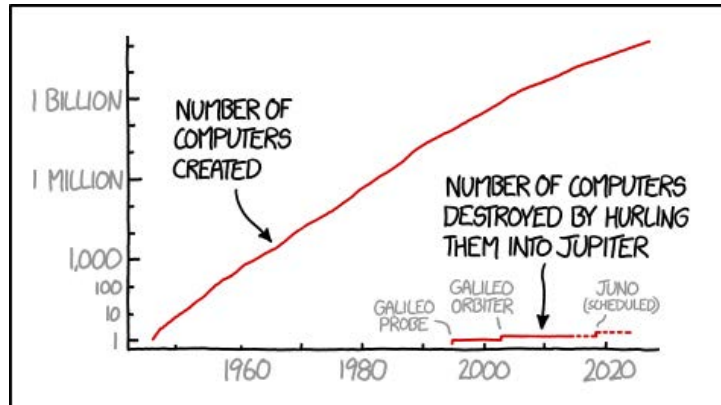


# Exponential Growth

# Exponential Growth

## OUTNUMBERS EVERYTHING ELSE QUICKLY

- Asymptotic advantage
- Quickly increasing add-on

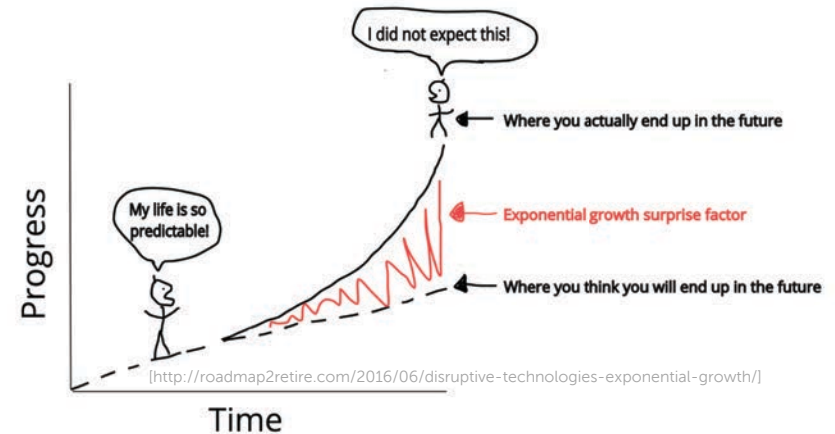


NASA NEEDS TO PICK UP THE PACE. IF THEY EVER WANT TO FINISH THE JOB.

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

## LEADS TO SURPRISING RESULTS

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

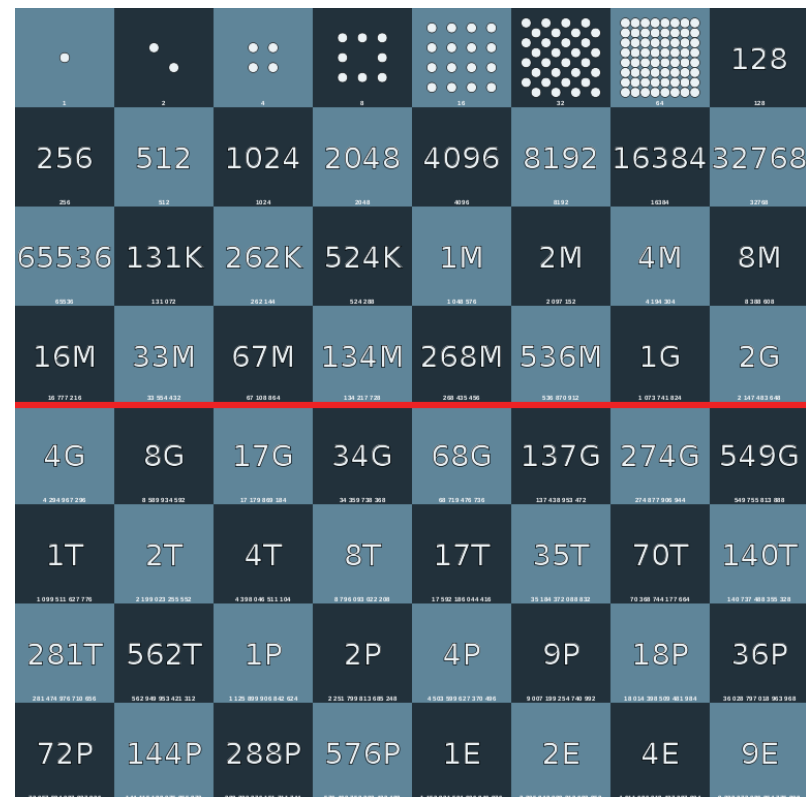
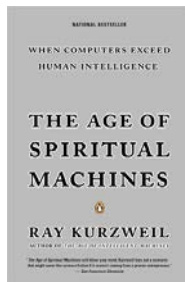




# Second Half of the Chessboard

## WHEN EXPONENTIAL GROWTH REALLY KICKS IN

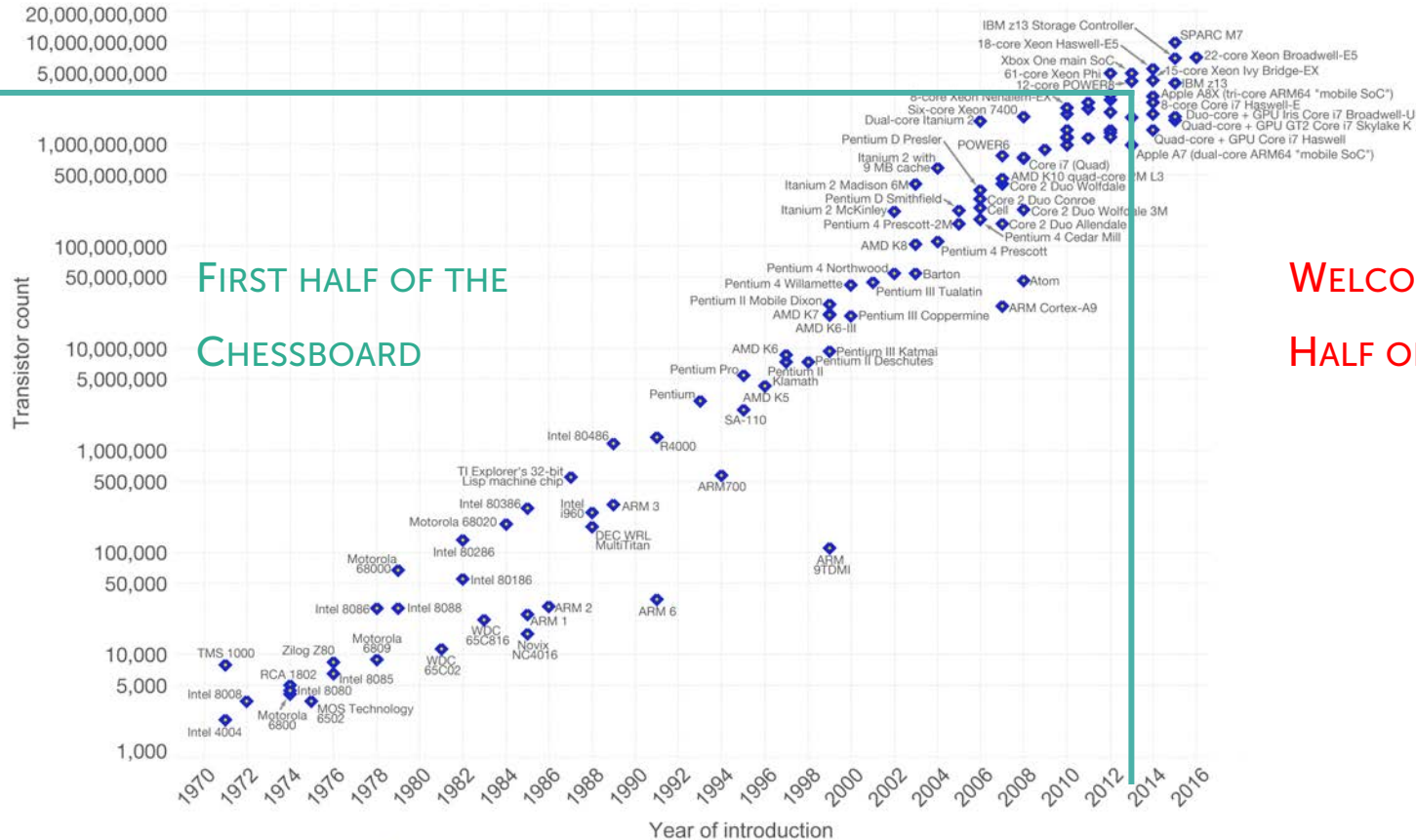
- According to Ray Kurzweil
- Things start to get interesting in the second half of the chess board
- Beyond 4G numbers quickly go beyond human intuition
- What happens in the second half can hardly be foreseen



[[https://en.wikipedia.org/wiki/File:Wheat\\_Chessboard\\_with\\_line.svg](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.

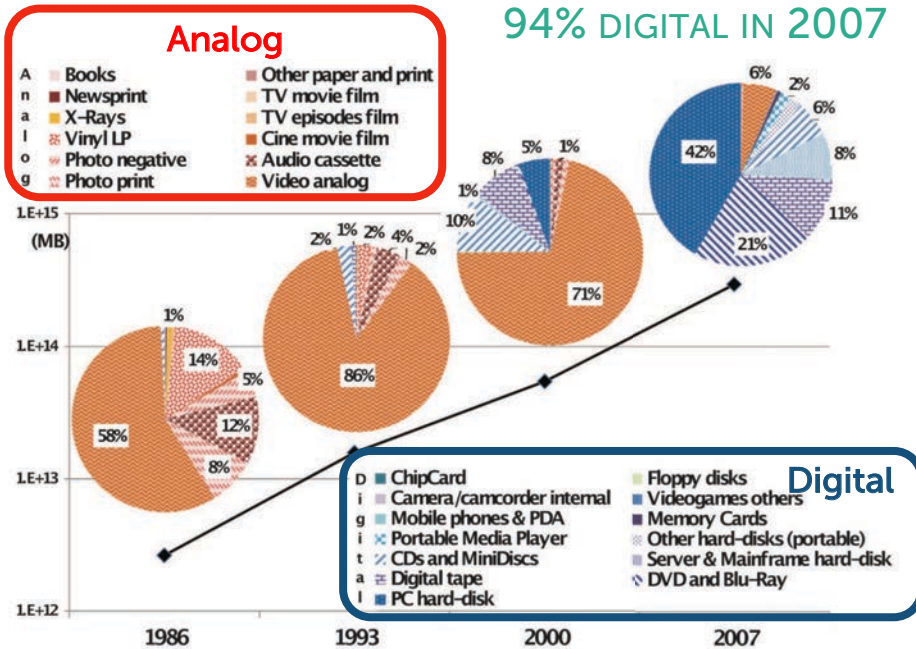


# Digitization

# Everything is Digital

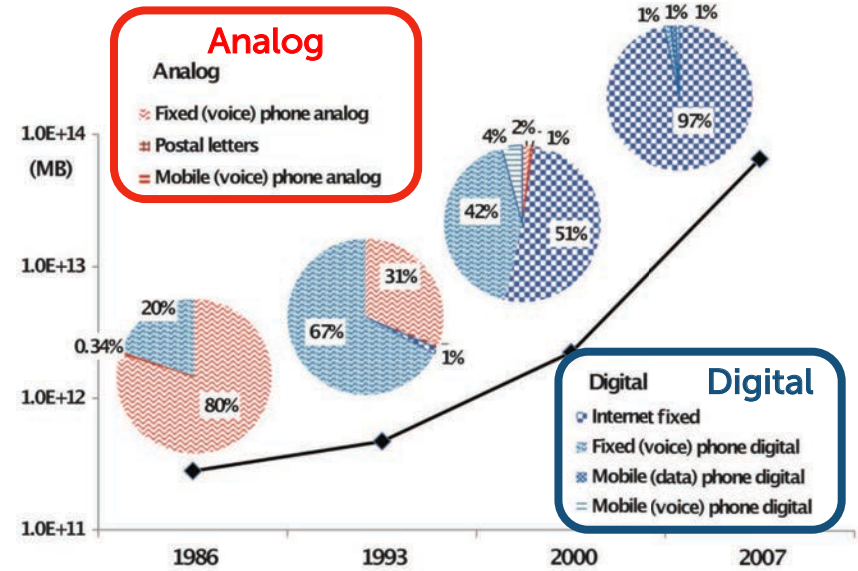
## WORLD'S CAPACITY TO STORE INFORMATION

94% DIGITAL IN 2007



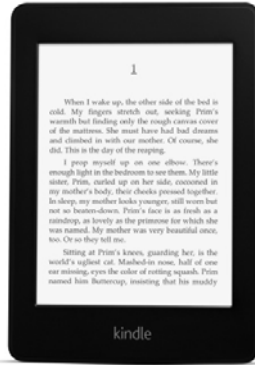
## WORLD'S CAPACITY TO TELECOMMUNICATE

99% DIGITAL IN 2007



[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



# 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



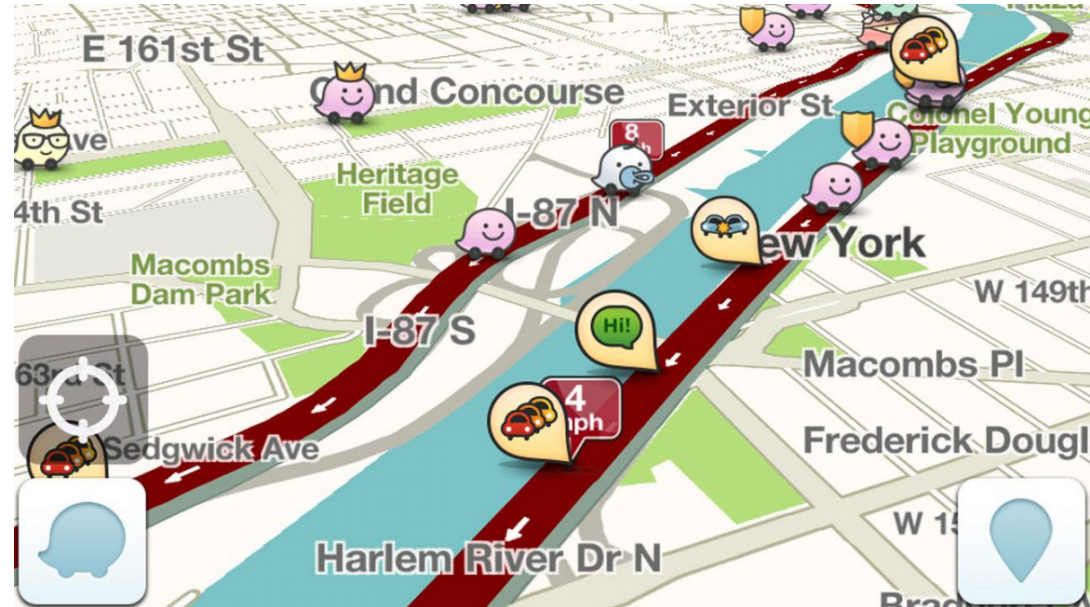
[<http://blog.acronis.com/posts/data-everything-8-noble-truths>]

# Recombinant Innovation



## CAR ROUTING APP THAT TURNS EVERY USER INTO A SENSOR

- Provides real-time traffic information
- Learns from users about traffic flow
- Crowd-sourcing of traffic information
- User can post alerts about
  - speed traps,
  - accidents,
  - traffic jams,
  - and even gas prices



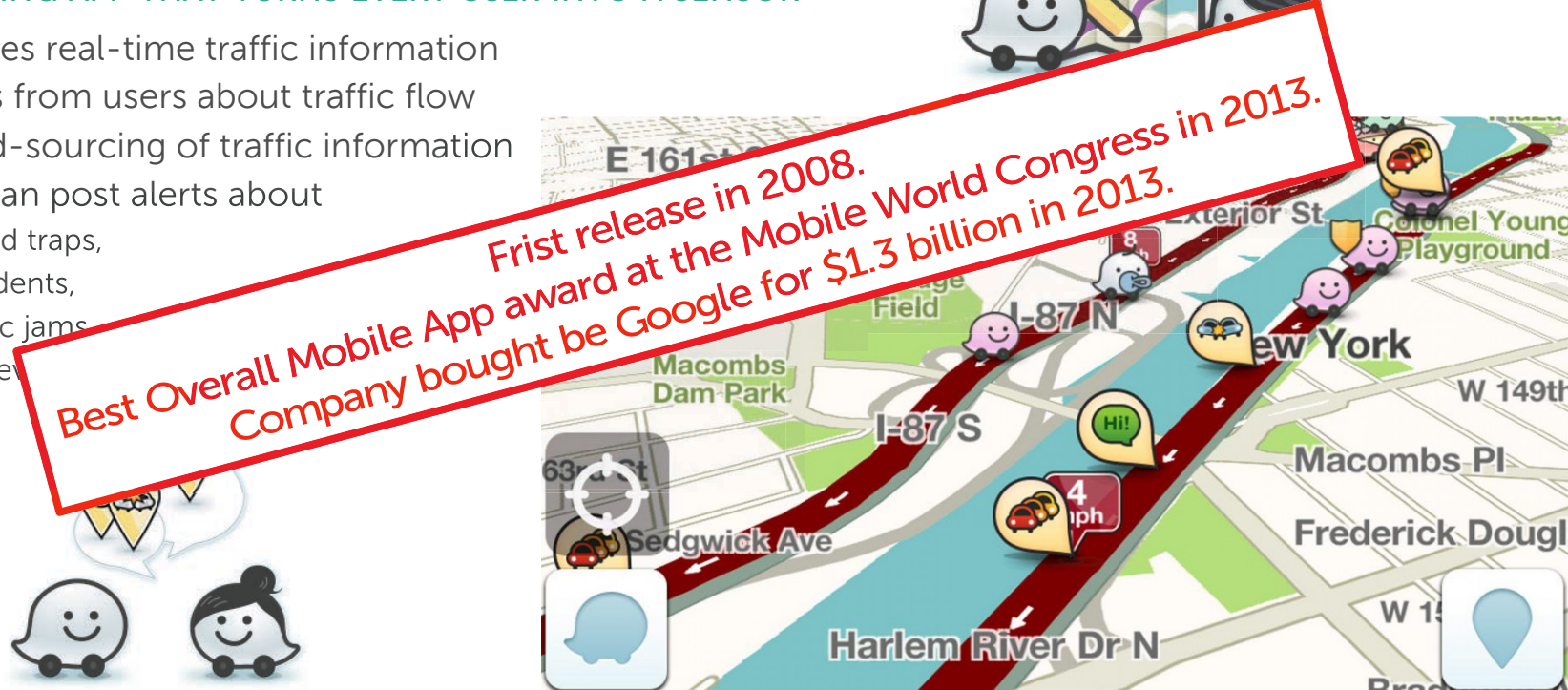
[<http://www.digitaltrends.com/cars/wazes-connected-citizens-program-turn-sunday-driver-ally/>]

Example:



## CAR ROUTING APP THAT TURNS EVERY USER INTO A SENSOR

- Provides real-time traffic information
- Learns from users about traffic flow
- Crowd-sourcing of traffic information
- User can post alerts about
  - speed traps,
  - accidents,
  - traffic jams
  - and e



[<http://www.digitaltrends.com/cars/wazes-connected-citizens-program-turn-sunday-driver-ally/>]

# Waze Ingredients

## ALL INGREDIENTS HAVE EXISTED LONG BEFORE

- Digital road maps
- Car navigation and routing
- GPS/GPS car navigation
- Mobile broadband
- PDAs/Smart Phones



TomTom Navigator, 2002

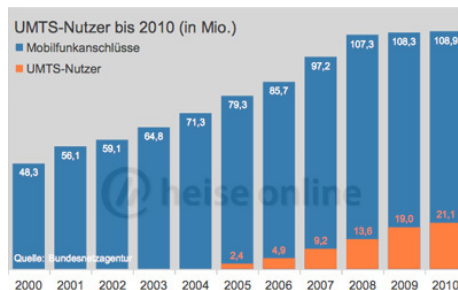


founded 1985



founded 1985  
as Karlin & Collins, Inc.

PSN-8 Manpack GPS Receiver  
used in Operation Desert Storm  
in 1991



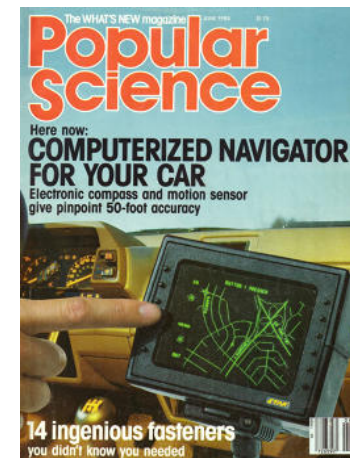
UMTS auction in 2000  
UMTS available in 2004



Newton, 1993



iPhone 1, 2007



Etak Navigator, 1985

# Waze Ingredients

## ALL INGREDIENTS HAVE EXISTED LONG BEFORE

- Digital road maps
- Car navigation and routing
- GPS/GPS car navigation
- Mobile broadband
- PDAs/Smart Phones



TomTom Navigator, 2002



founded 1985  
as Karlin & Collins, Inc.

**Although major break through in traffic organization,  
Waze is no major technological break through.  
Waze combines existing technology for innovation.**



Newton, 1993



iPhone 1, 2007



UMTS auction in 2000  
UMTS available in 2004



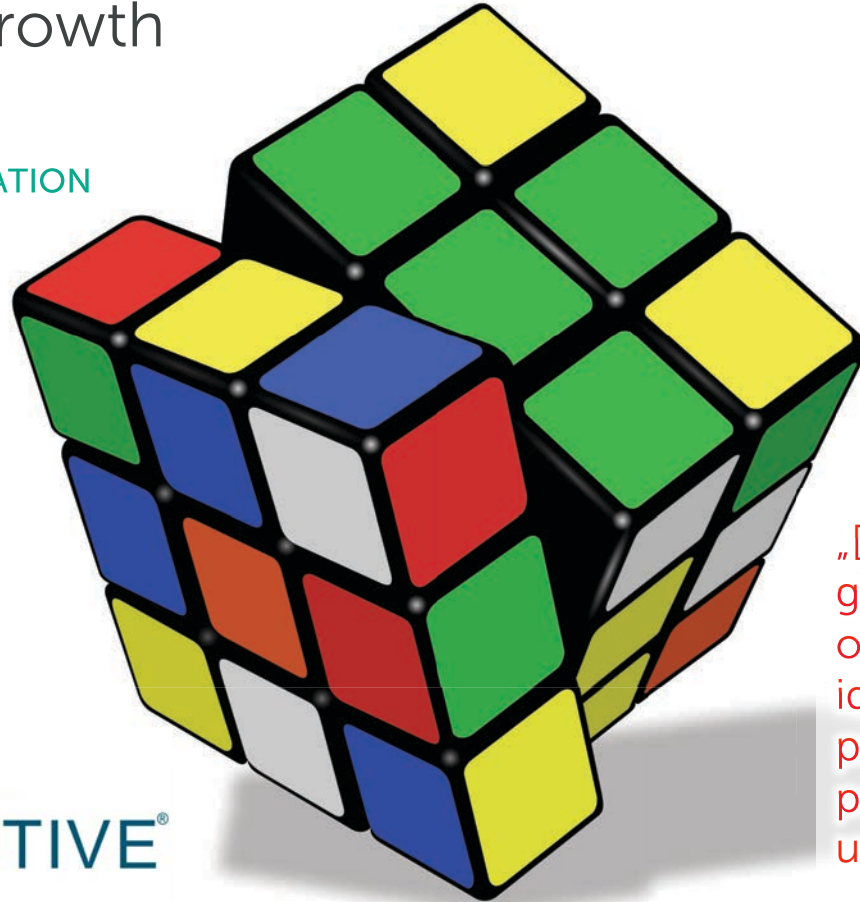
Etak Navigator, 1985

# Recombinant Growth

## INNOVATION IS RECOMBINATION

### OF EXISTING IDEAS

- Possibilities explode quickly
- But we got way better in processing ideas



Dresden Database  
Systems Group

## THE QUARTERLY JOURNAL OF ECONOMICS

Vol. CXIII

May 1998

Issue 2

### RECOMBINANT GROWTH\*

MARTIN L. WEITZMAN

This paper attempts to provide microfoundations for the knowledge production function in an idea-based growth model. Production of new ideas is made a function of newly reconfigured old ideas in the spirit of the way an agricultural research station develops improved plant varieties by cross-pollinating existing plant varieties. The model shows how knowledge can build upon itself in a combinatoric feedback process that may have significant implications for economic growth. The paper's main theme is that the ultimate limits to growth lie not so much in our ability to generate new ideas as in our ability to process an abundance of potentially new ideas into usable form.

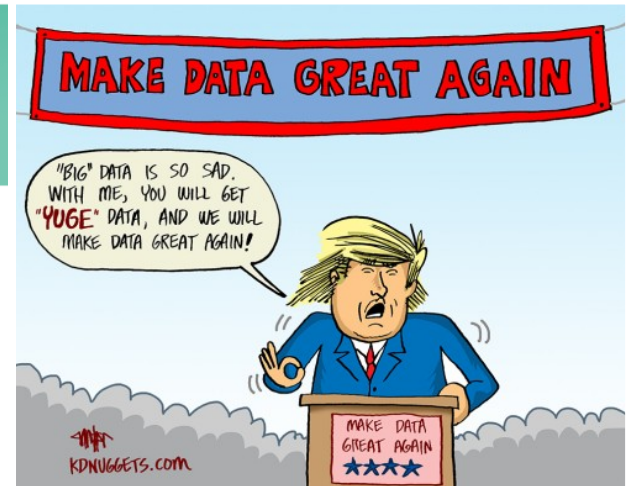
### I. INTRODUCTION

„[T]he ultimate limits to growth lie not so much in our ability to generate new ideas as in our ability to process an abundance of potentially new ideas into usable form.“

Harvard College and the Massachusetts Institute of Technology  
The Quarterly Journal of Economics, May 1998



# Big Data



# Big Data term reflects the three drivers

## EXPONENTIAL GROWTH

- beyond intuition
- second half of the chessboard

•	••	•••	••••	•••••	••••••	•••••••	••••••••	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

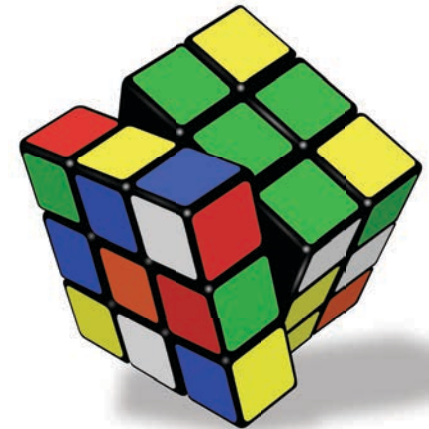
- Everything is digital data
- Analog signals are not part of big data



**Data**

## RECOMBINANT INNOVATION

- explosion of ideas
- all somehow seem to be part of big data

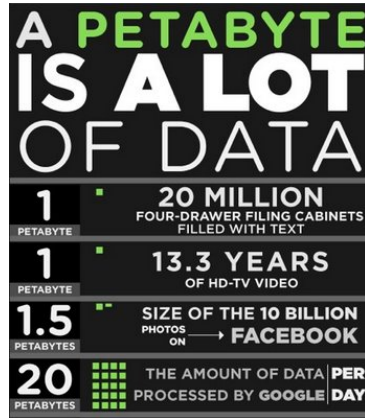


**Fuzziness of the term**

# 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



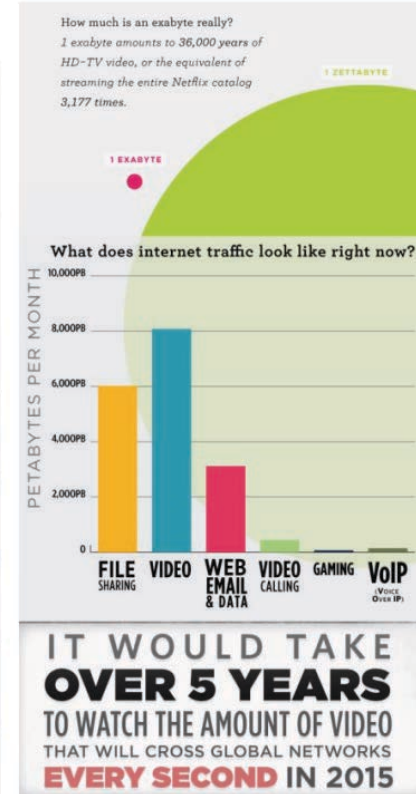
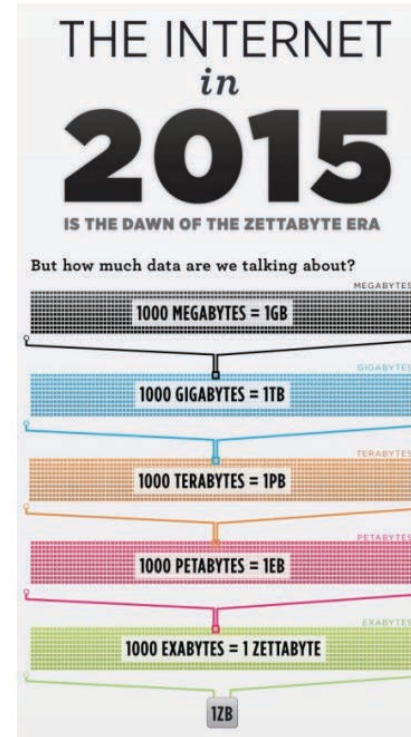
## THE ZETTABYTE AGE

- 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

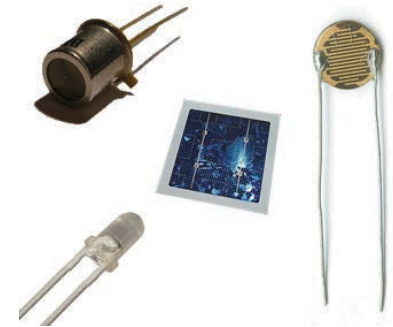
## HUMANE-PRODUCED DATA

- ~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



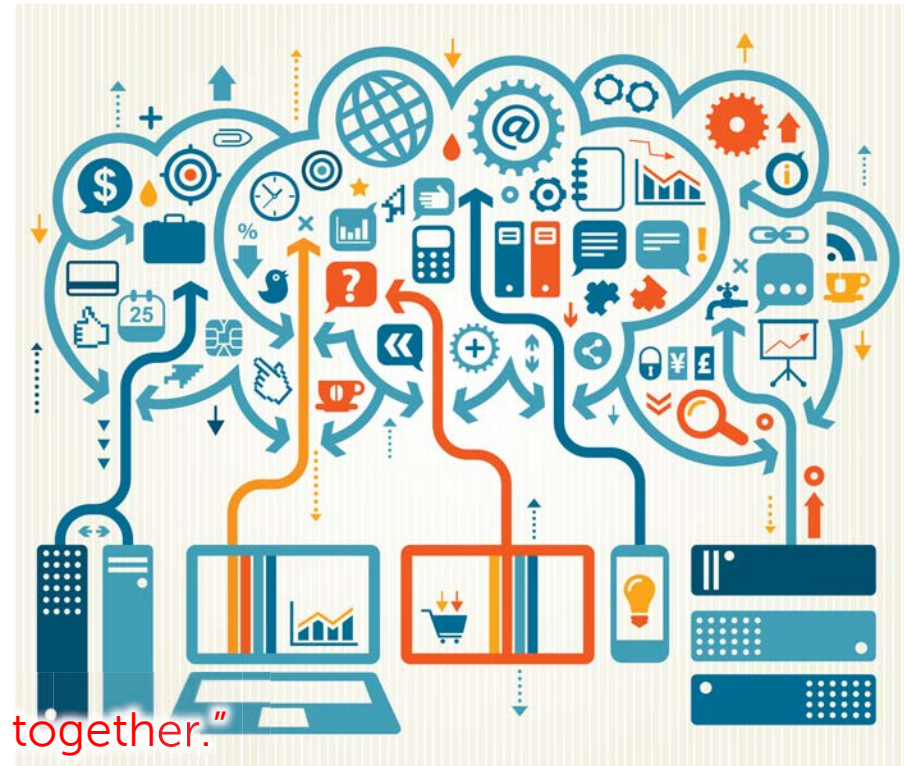
# Everything is data ► Variety

## DATA FROM ALL KINDS OF DIGITAL SOURCES

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

## SEMANTIC HETEROGENITY

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

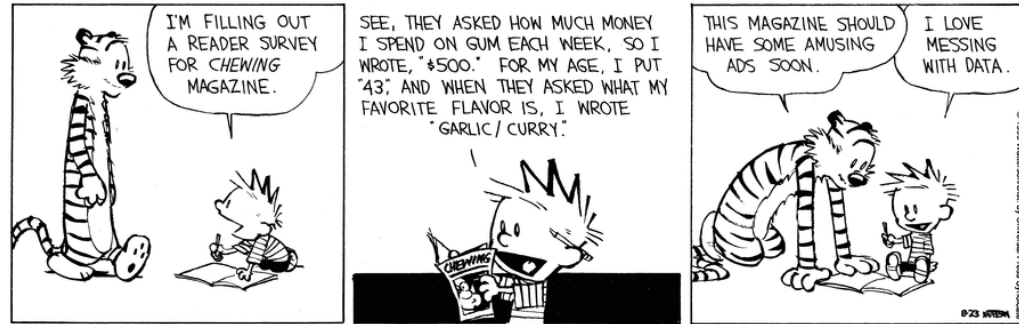


“A lot of Big Data is a lot of small data put together.”

# Data is messy ▶ Veracity

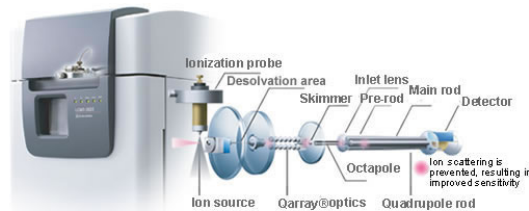
## QUALITY OF CAPTURED DATA VARIES GREATLY

- 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.

## 40 ZETTABYTES

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

6 BILLION PEOPLE  
have cell phones



WORLD POPULATION: 7 BILLION

## 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



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

## 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



# 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**

## Veracity UNCERTAINTY OF DATA

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

Value - the fifth V  
of Big Data

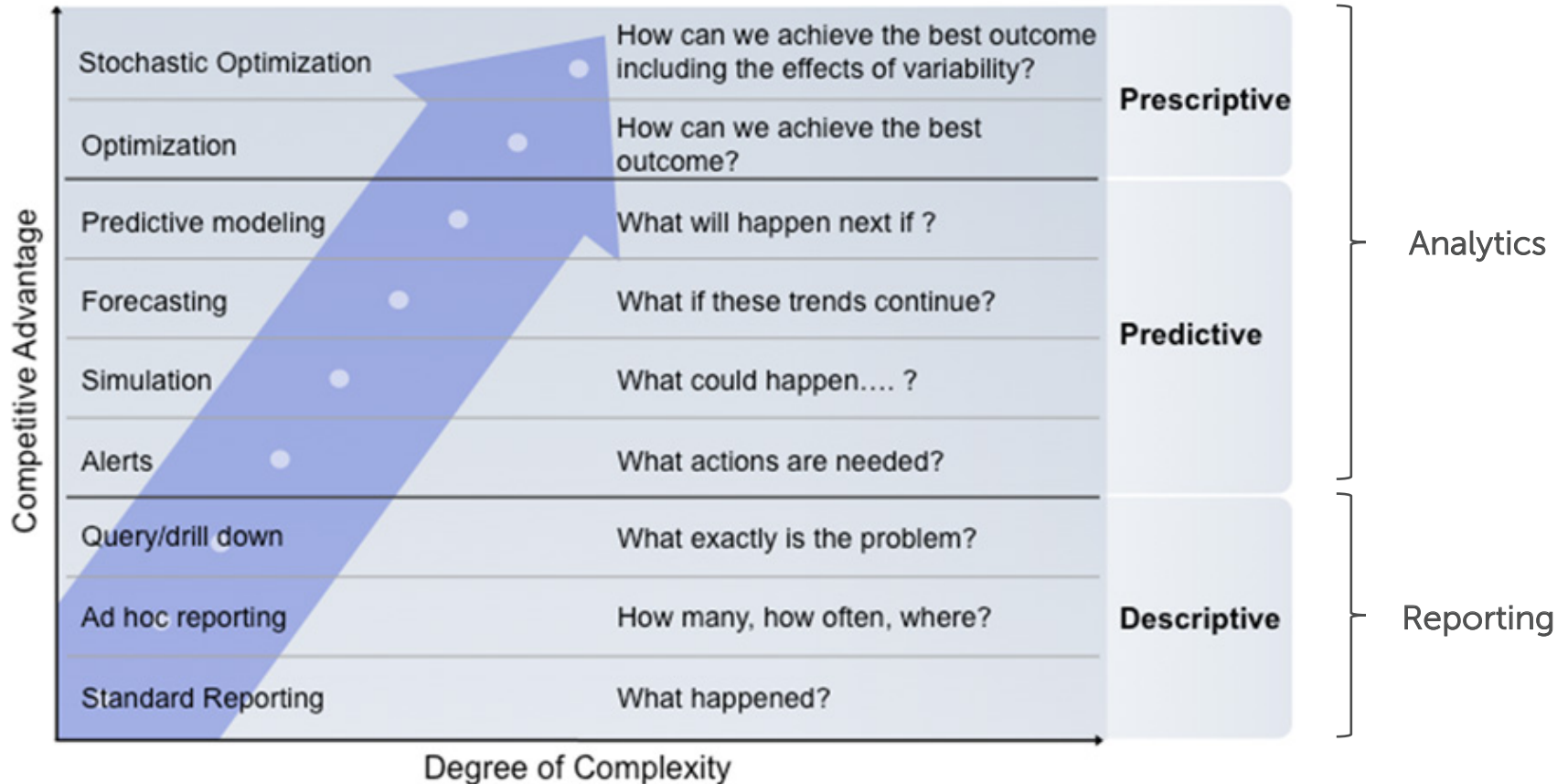


## Data Science/Data Analysis

... 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



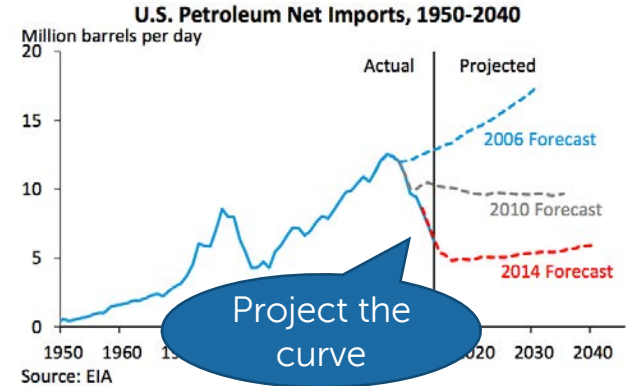
# Predictive Analytics

## 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. precops, predpol)



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

Fill in the blanks



# Prescriptive Analytics



ORION –  
On-Road Integrated  
Optimization and Navigation

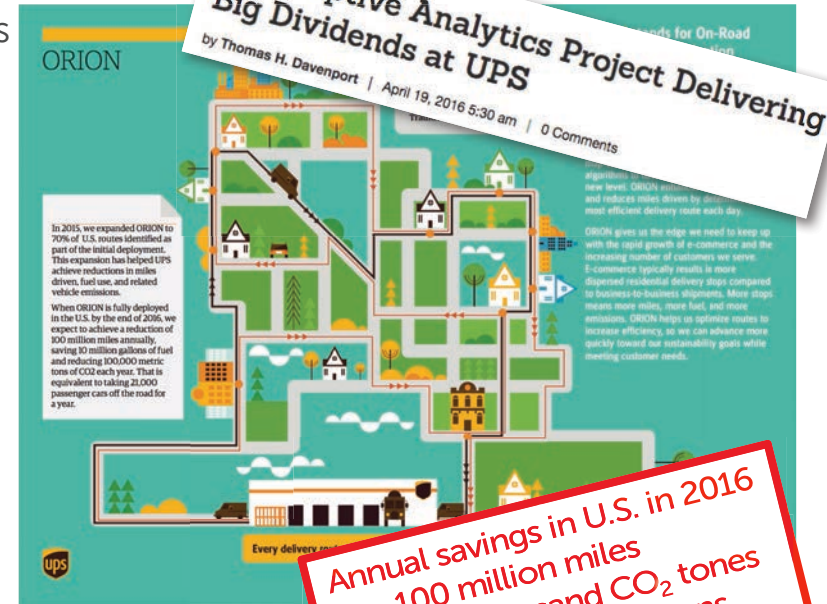


## WHAT SHOULD I DO?

- 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



- Annual savings in U.S. in 2016
- 100 million miles
- 100 thousand CO<sub>2</sub> tones
- 10 million fuel gallons
- \$300-400 million costs

Same same but different

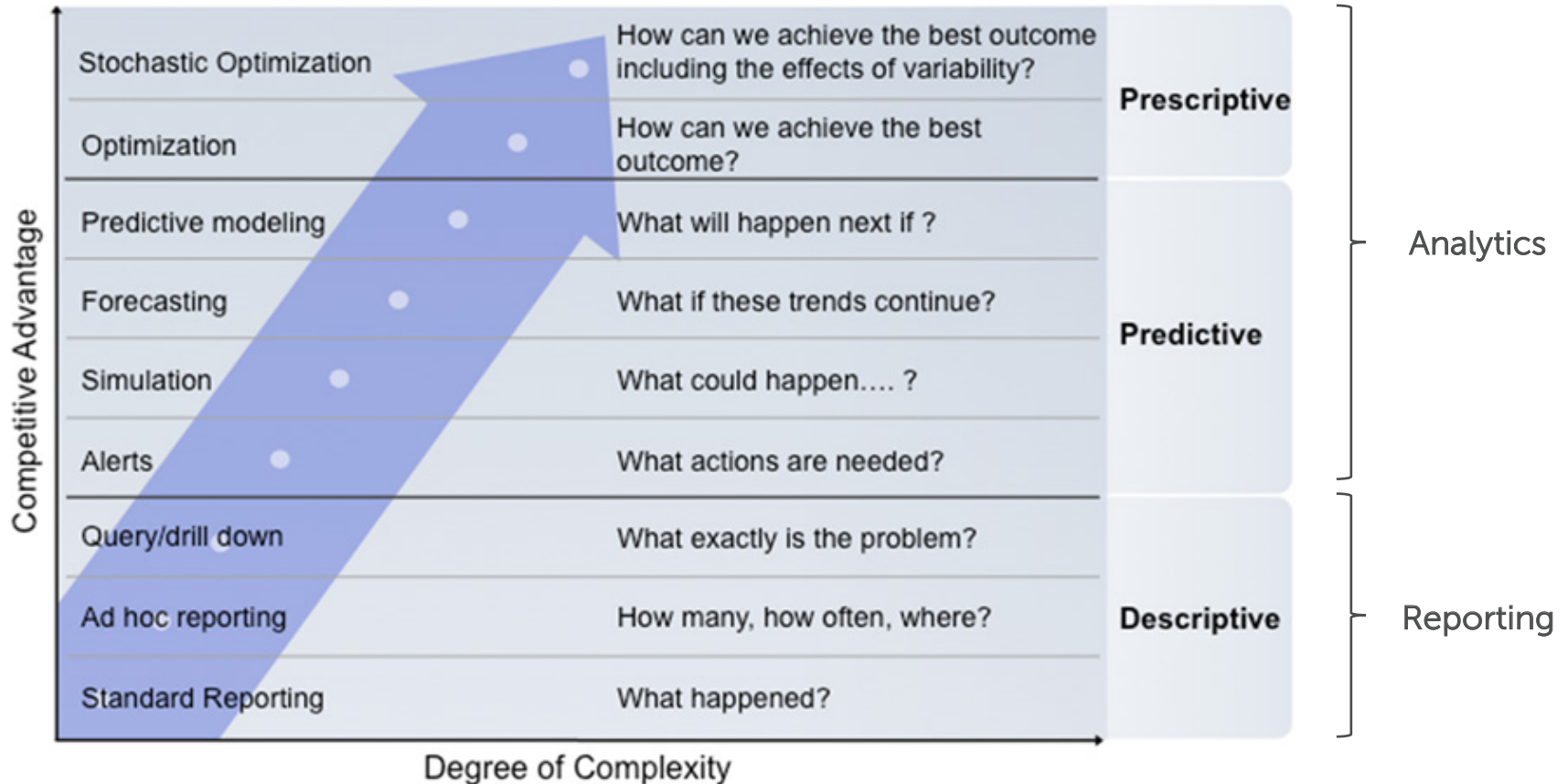
## CHARACTERISTICS OF BUSINESS INTELLIGENCE (BI)

- Provides pre-created dashboards for management
- Repeated visualization of well known analysis steps
- Deals with structured data
- Typically data is generated within the organization
- Central data storage (vs. multiple data silos)
- Handled well by specialized database techniques

## TYPICAL TYPES OF QUESTIONS AND INSIGHT

- Customer service data: "what business causes customer wait times"
- Sales and marketing data: "which marketing is most effective"
- Operational data: "efficiency of the help desk"
- Employee performance data: "who is most/least productive"

# Levels of Analysis



# From Data Warehouse to Data Lake

## WITH CHEAP STORAGE COSTS, PEOPLE PROMOTE THE CONCEPT OF THE DATA LAKE

- Combines data from many sources and of any type
- Allows for conducting future analysis and not miss any opportunity

## COLLECT EVERYTHING

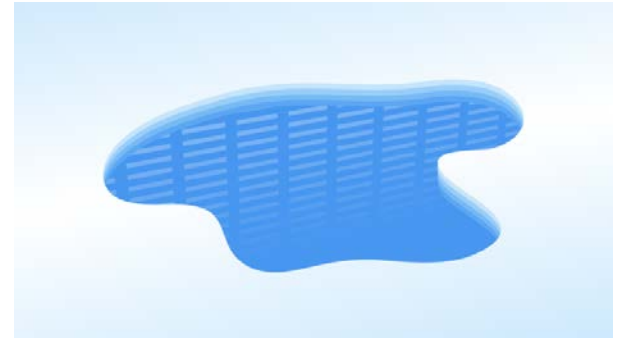
- All data, both raw sources over extended periods of time as well as any processed data
- Decide during analysis which data is important, e.g., no “schema” until read

## DIVE IN ANYWHERE

- Enable users across multiple business units to refine, explore and enrich data on their terms

## FLEXIBLE ACCESS

- Enable multiple data access patterns across a shared infrastructure: batch, interactive, online, search, and others



# Roles in Big Data Projects

## DATA SCIENTIST

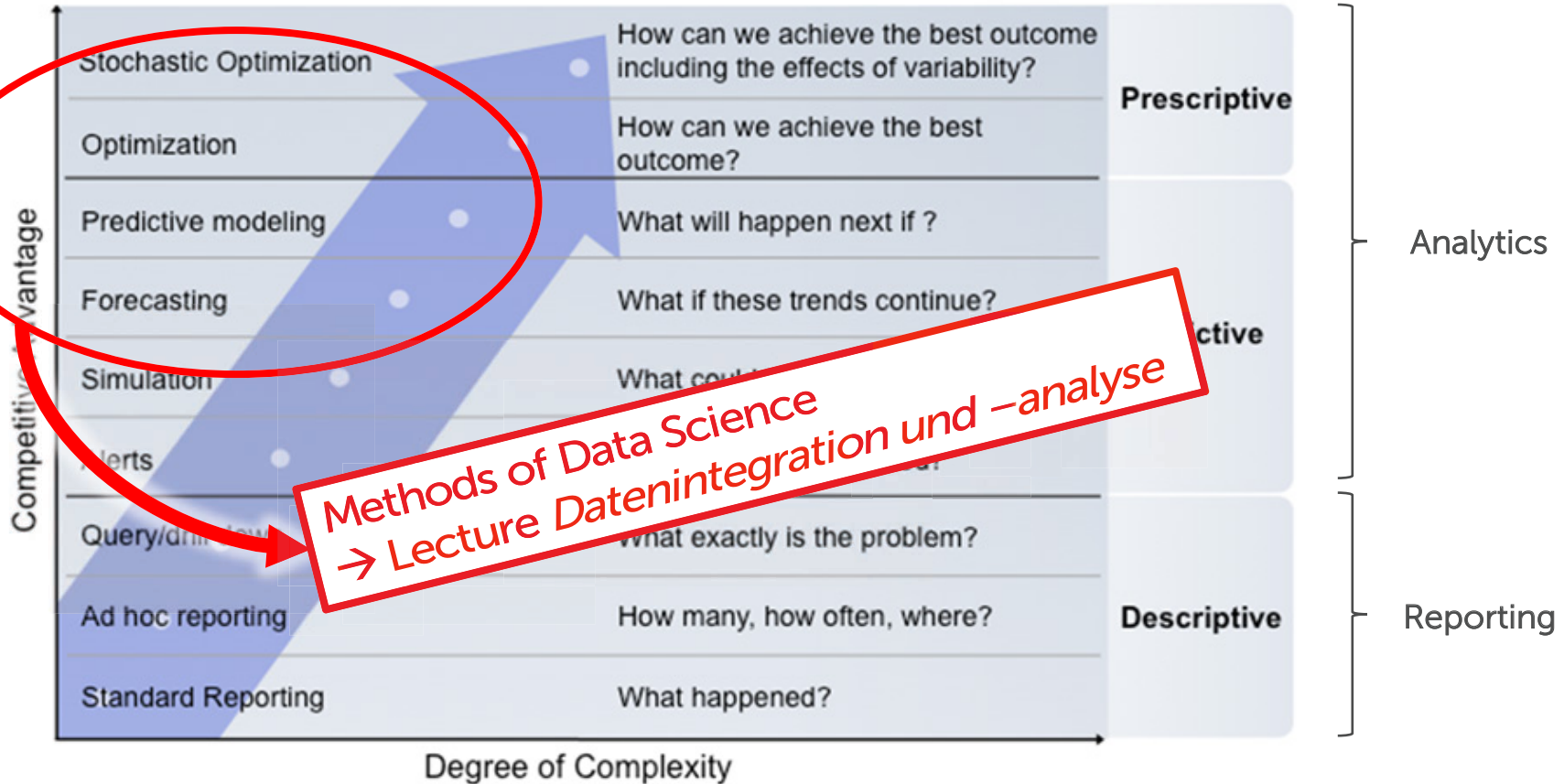
- Data science is a systematic method dedicated to knowledge discovery via data analysis
- In business, optimize organizational processes for efficiency
- In science, analyze experimental/observational data to derive results
- Typical skills
  - Statistics + (mathematics) background
  - Computer science: 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

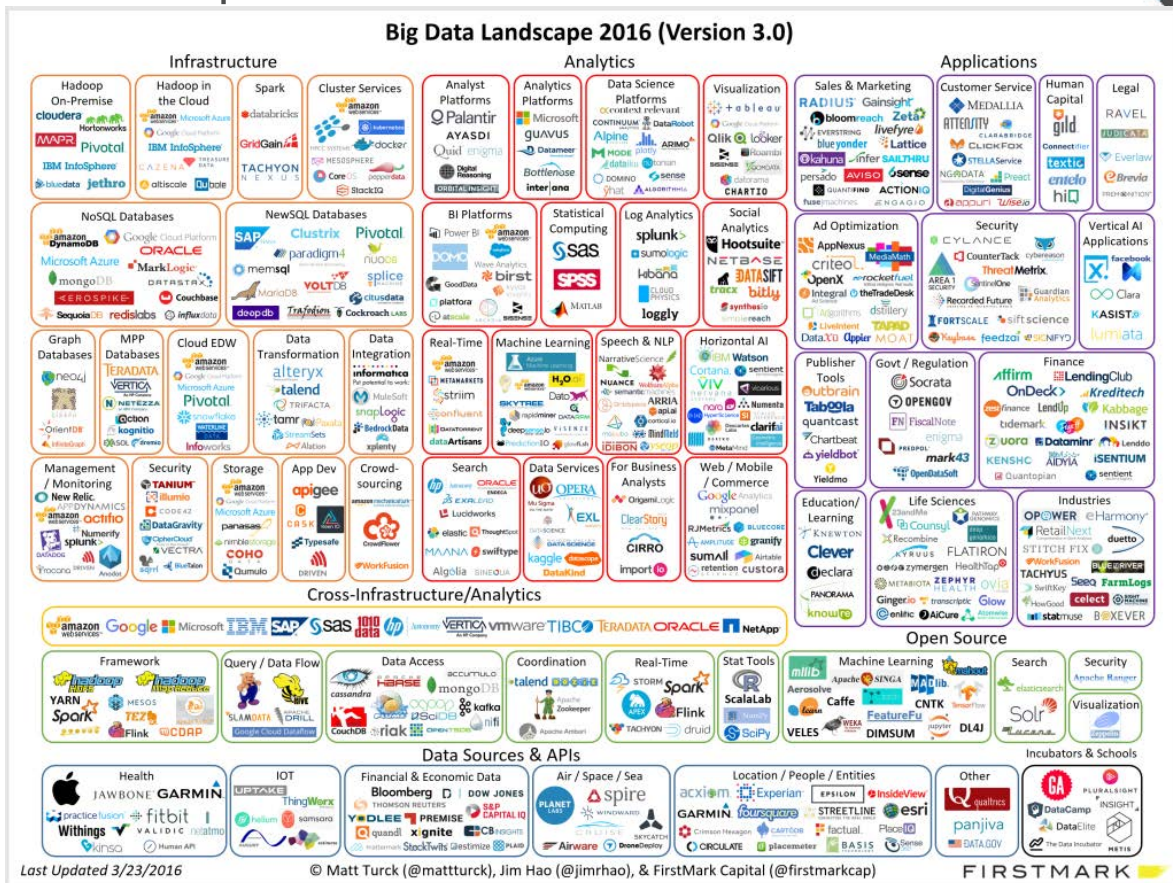
# Our Focus in Teaching in Research

# Levels of Analysis





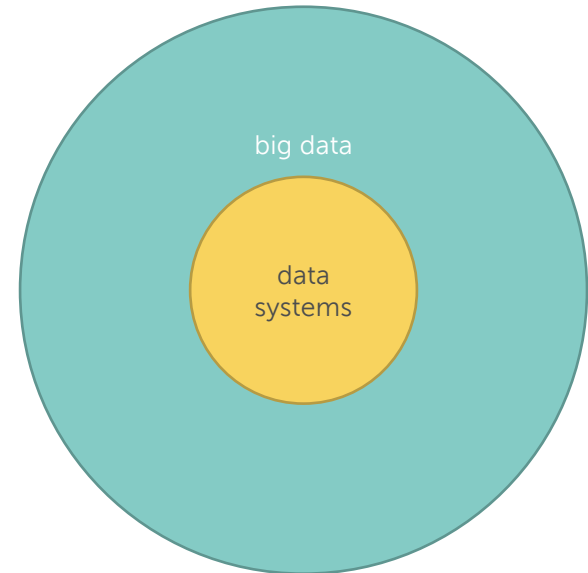
# Big Data Landscape(s)



## DATA SYSTEMS ARE IN THE MIDDLE OF ALL THIS

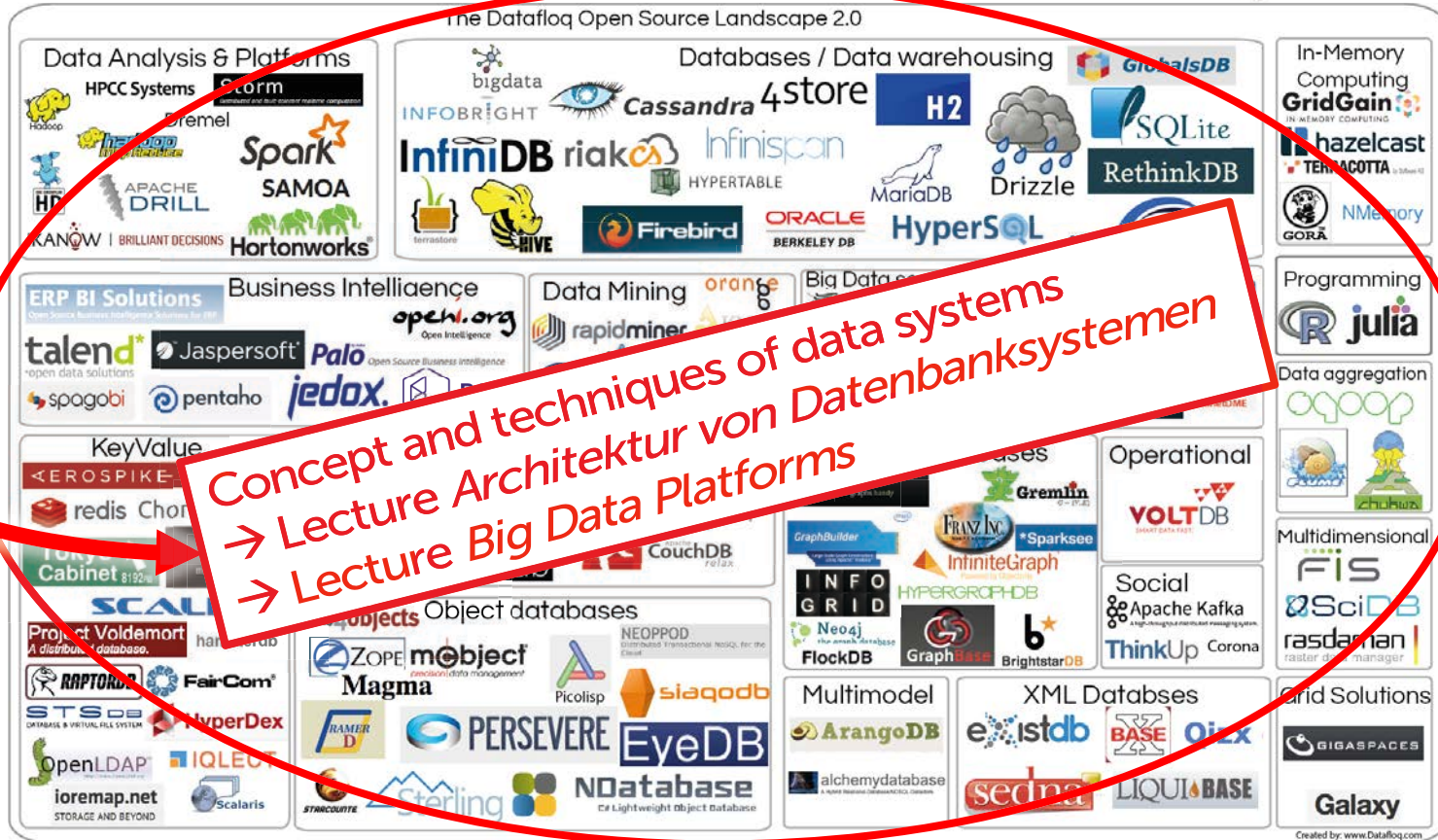
### A DATA SYSTEM...

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



## DIFFERENT DATA SYSTEMS USE DIFFERENT DATA MODELS

# (Big) Data System Landscape(s)



Concept and techniques of data systems  
→ Lecture Architektur von Datenbanksystemen  
→ Lecture Big Data Platforms



**Graph databases**

Gephi makes graphs handy

Gremlin

FRANZ INC. Max 3 G + Oracle

GraphBuilder Large Scale Graph Construction using Apache Hadoop

\*Sparksee

InfiniteGraph Powered by Objectivity

HYPERGRAPHDB

Neo4j the graph database

FlockDB

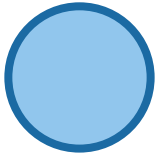
GraphBase

BrightstarDB

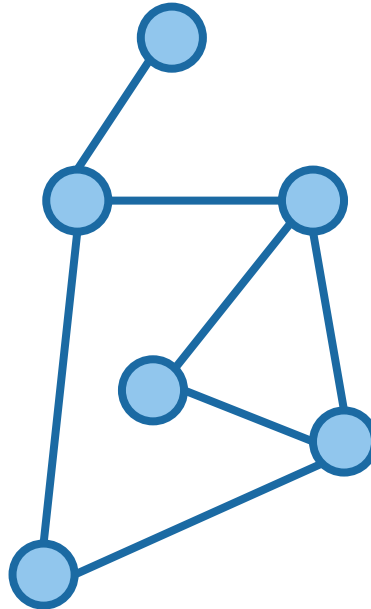


# Graph Building Blocks

## NODES (DOTS)



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

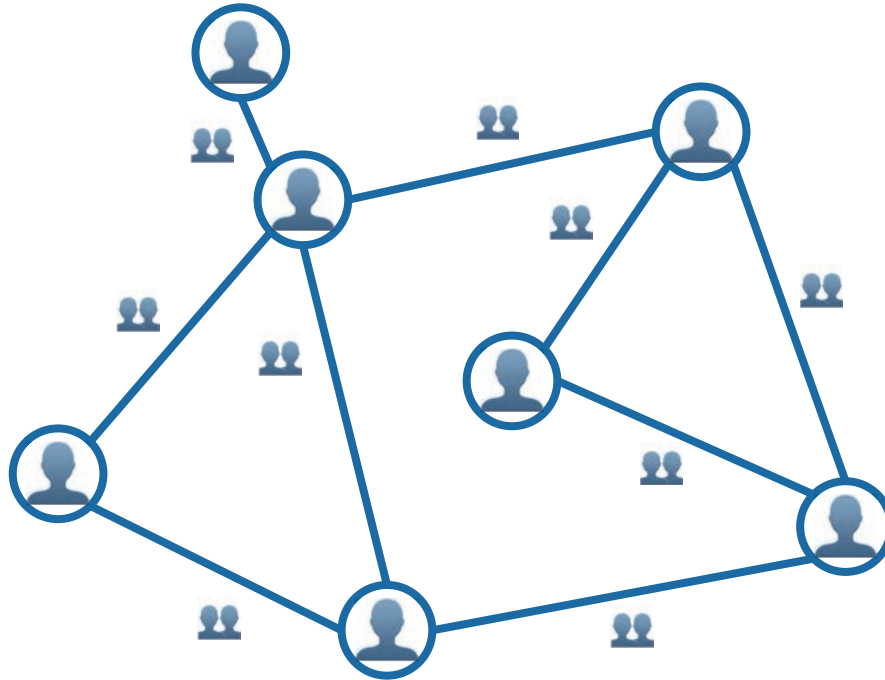


## EDGES (LINES)

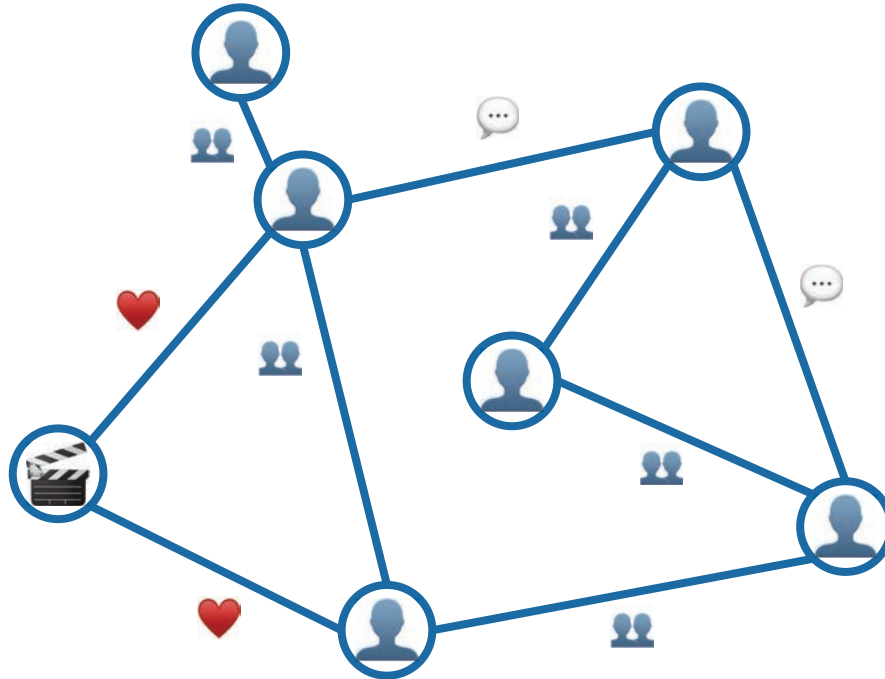


- Like a relationship in ER
- Exist only between nodes
- Identity depends on the nodes they connect

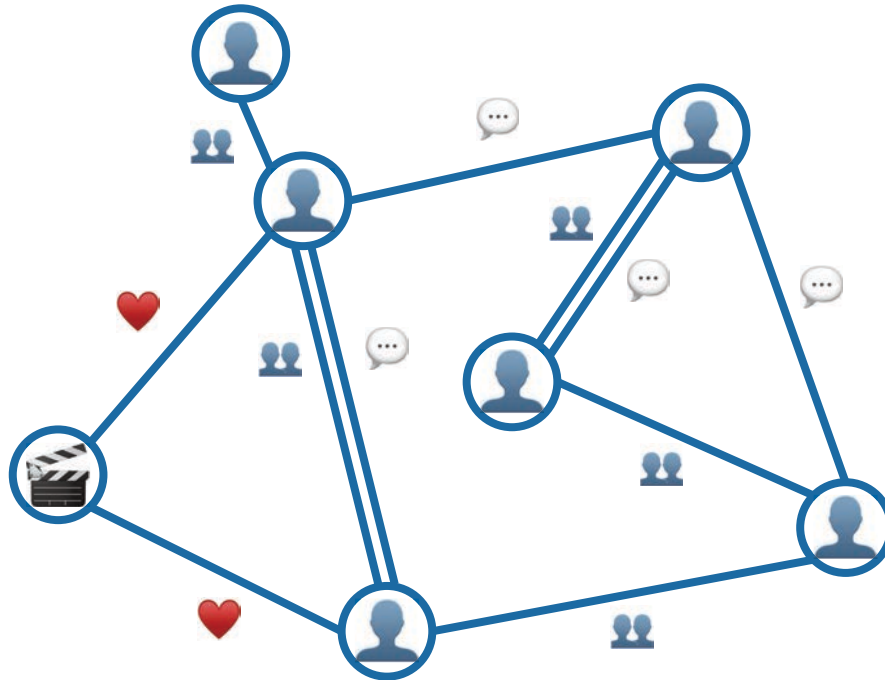
# Graph Data – Social Network



# Graph Data – Social Network



# Graph Data – Social Network







# Social Graphs

## FACEBOOK

- May 2013



As of 12/2014: "1.39 billion active users with more than 400 billion edges"





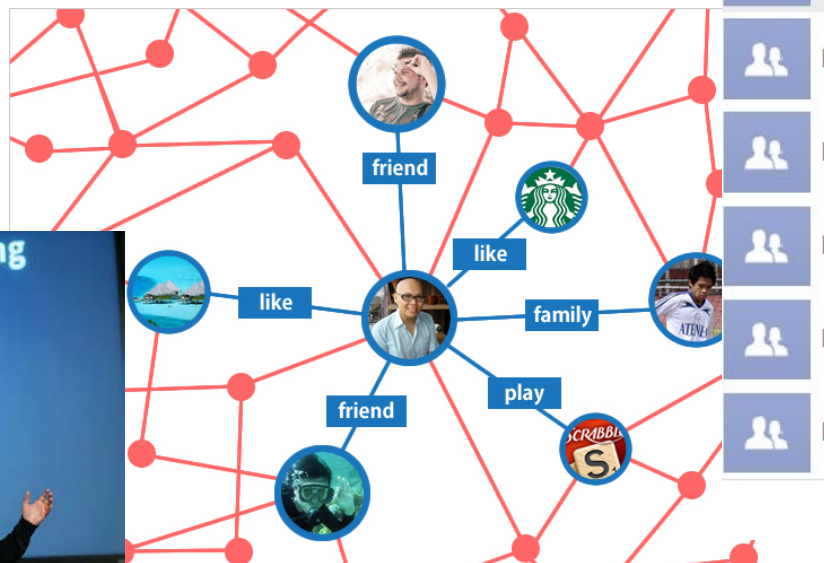




# Structured Search

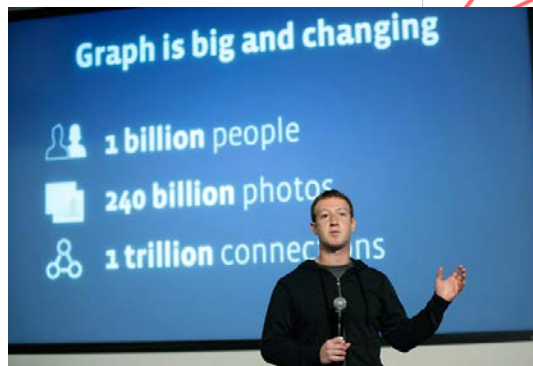
## EXAMPLE: FACEBOOK GRAPH SEARCH

- Finding subgraph structures
- Very natural way of formulating queries



Search: People who like Michael Jackson

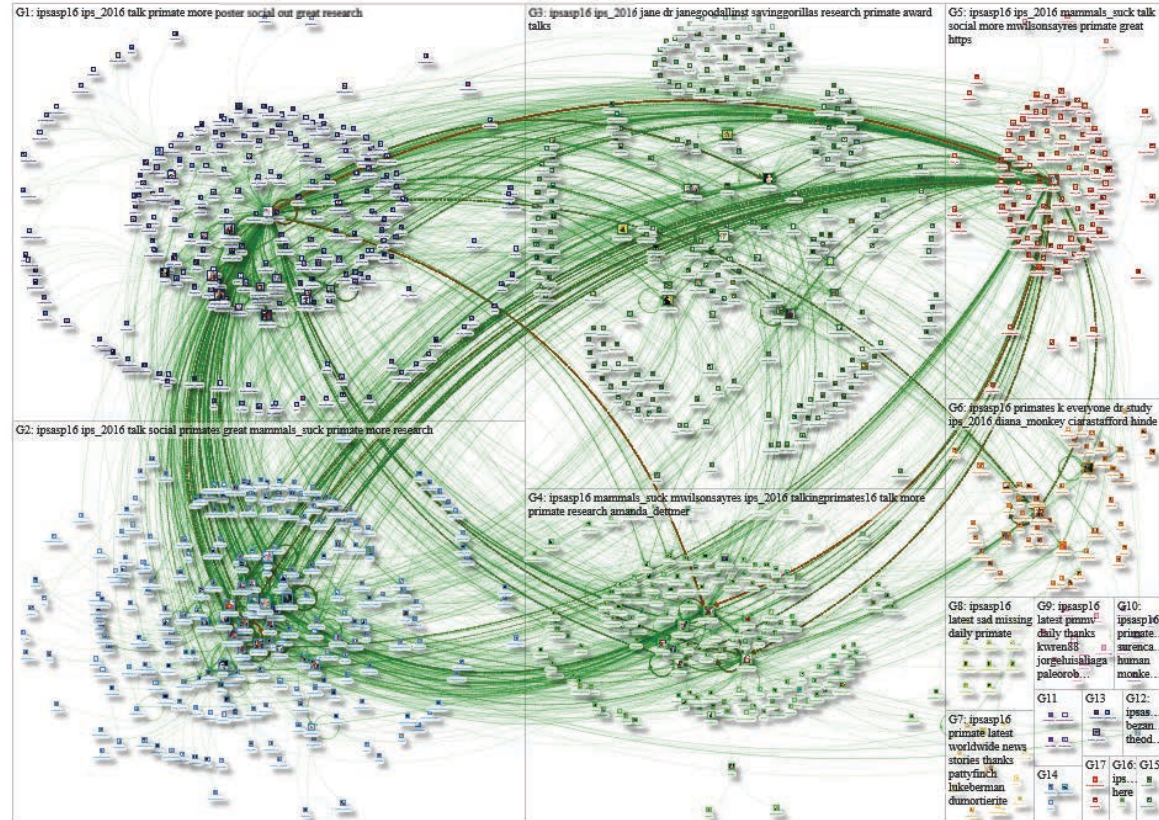
- People who like **Michael Jackson** (Musician/Band · 55,100)
- People who like **Michael Jackson** (Interest)
- People who like **Michael Jackson's This Is It**
- People who like **Michael Jackson The Experience**
- People who like **Michael Jackson Legend Never Die**
- People who like **Michael Jackson** and live in **Pune, M**



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

## EXAMPLE: TWITTER COMMUNICATION

- Users tweeting on a specific topic
- Others reply or 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>





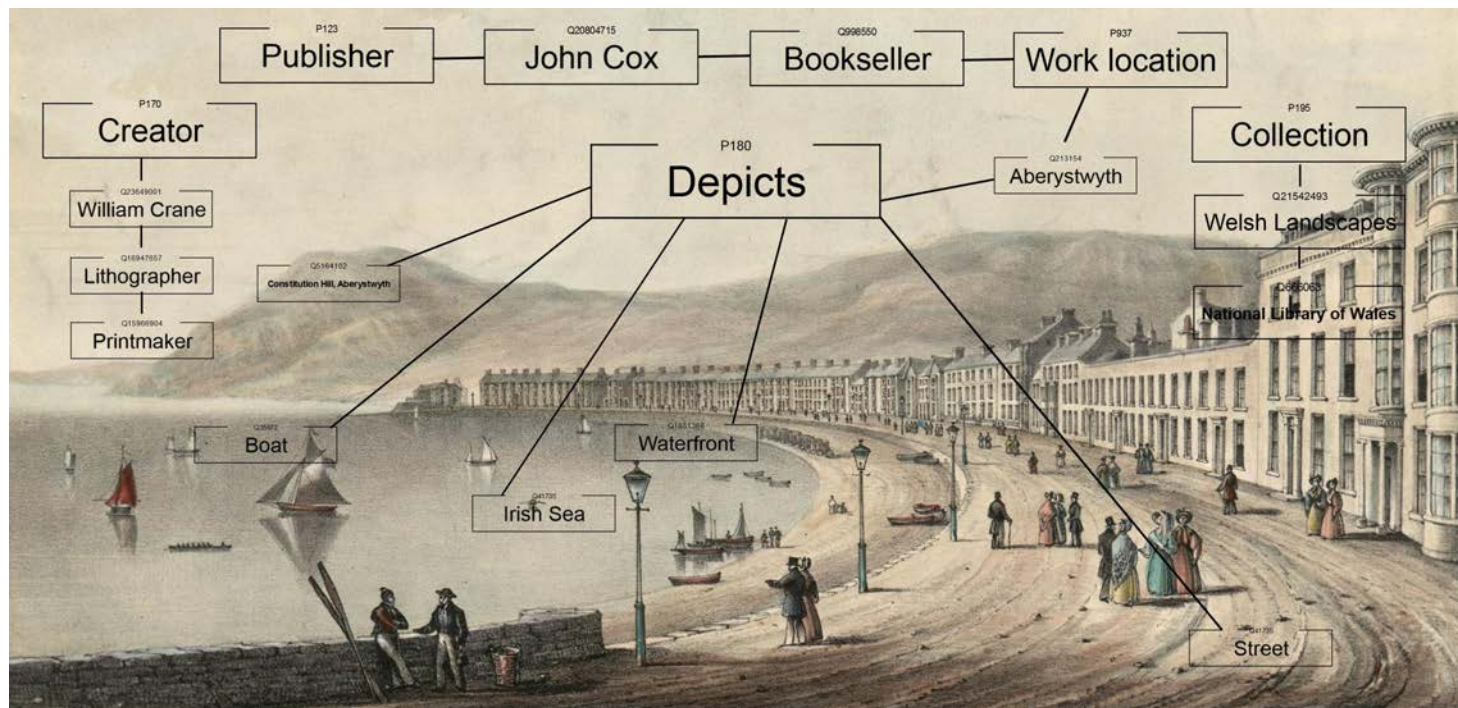
# Viral Marketing

## VIRAL MARKETING

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



## KNOWLEDGE GRAPH OF A PICTURE

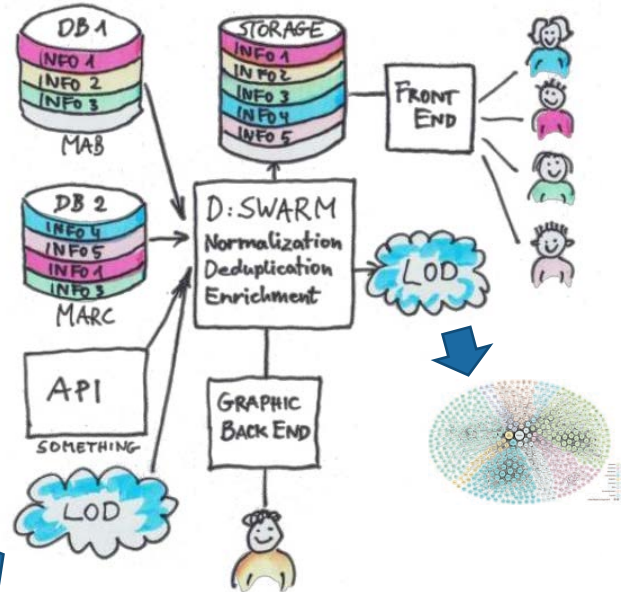
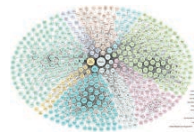
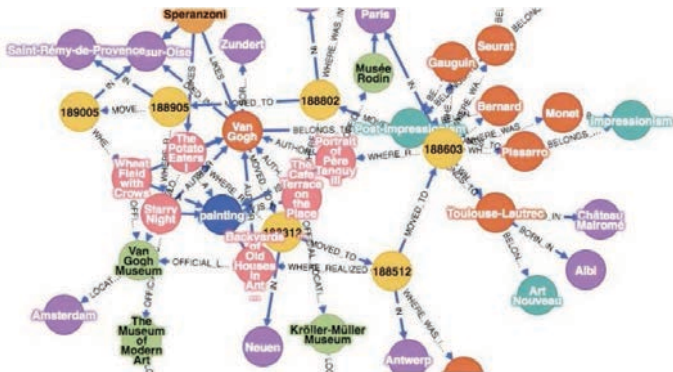


[The National Library of Wales, <https://www.llgc.org.uk/blog/?p=11246>]

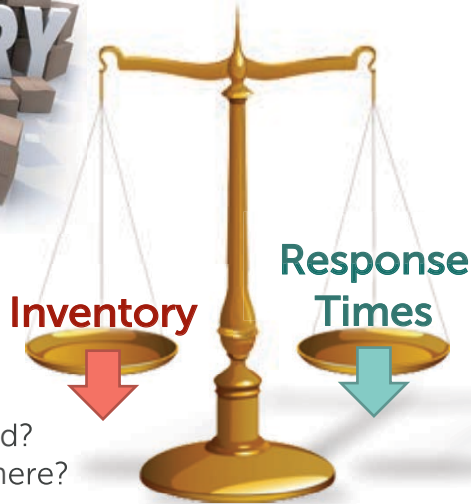


## SÄCHSISCHE LANDESBIBLIOTHEK – STAATS- UND UNIVERSITÄTSBIBLIOTHEK DRESDEN (SLUB)

- 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



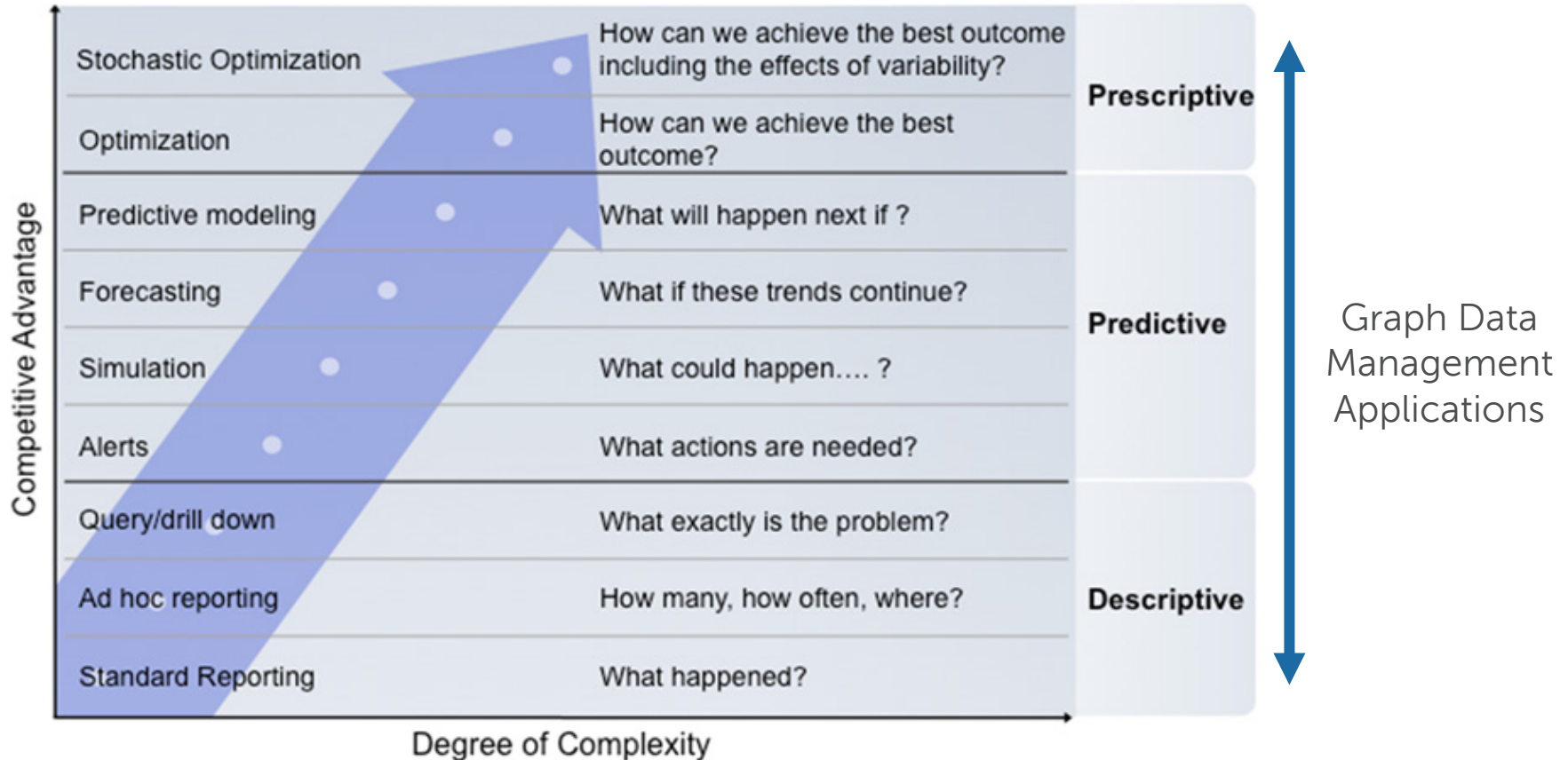
When to send?  
How much to send?  
From where to where?



Supply Chain Optimization

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

# Level of Analytics



→ Lecture Graph Data Management and Analytics

### Graph databases

- Gephi: makes graphs handy
- Gremlin:  $G = (V, E)$
- Franz Inc. (Mark 3.0 + Oracle)
- Sparksee
- InfiniteGraph: Powered by Objectivity
- InfoGrid
- HypergraphDB
- Neo4j: the graph database
- FlockDB
- GraphBase
- BrightstarDB

Logos visible in the collage include: Firebird, Oracle Berkeley DB, HyperSQL, monetdb, Couchbase, RaptorDB, EJDB, CouchDB, Neo4j, FlockDB, GraphBase, BrightstarDB, Apache Kafka, ThinkUp, Corona, and many others.

# Summary

## BIG DATA

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

## RELATED RESEARCH AND LECTURES

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