

Big Data Management and Analytics

Lecture Notes

Winter semester 2015 / 2016
Ludwig-Maximilians-University Munich

© Prof. Dr. Matthias Renz 2015

Based on lectures by
Donald Kossmann (ETH Zürich), as well as
Jure Leskovec, Anand Rajaraman, and Jeff Ullman (Stanford University)

- Course website:
 - http://www.dbs.ifi.lmu.de/cms/Big_Data_Management_and_Analytics
 - Registration for this lecture is now open via Uniworx
 - Registration required to attend the exams!!!
- Organization:
 - Load: 3+2 hours weekly
 - Required: Lecture "Database Systems I" or equivalent
 - Beneficial: Lecture "Knowledge Discovery in Databases I" or equivalent
 - Lecture: Prof. Dr. Matthias Renz



- Assisting:
Klaus Arthur Schmid, Felix Borutta, Evgeniy Faermann, Christian Frey



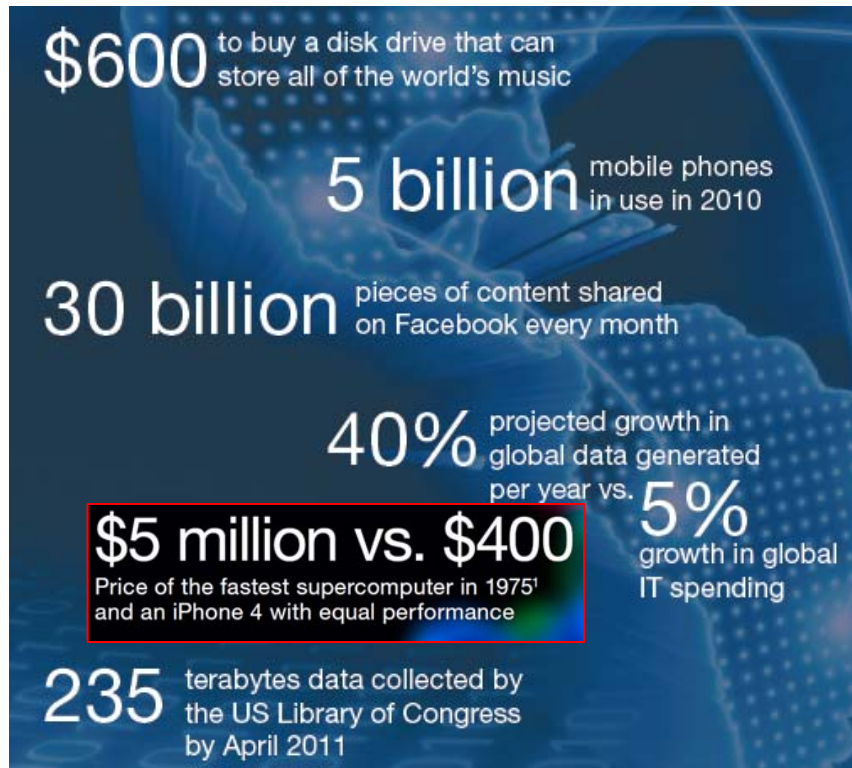
- Tutors: TBA

Why this course?

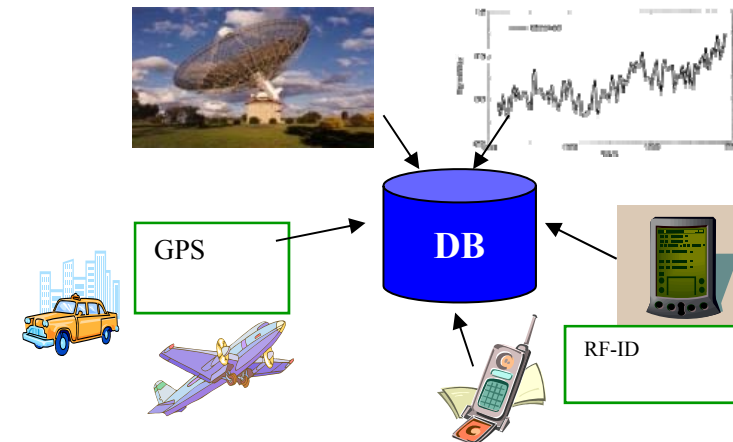
- **Big Data is big**
 - \$ and science: choose your poison

We are drowning in data ... but starving for information

- Exponential grows in data



J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets,
<http://www.mmids.org>



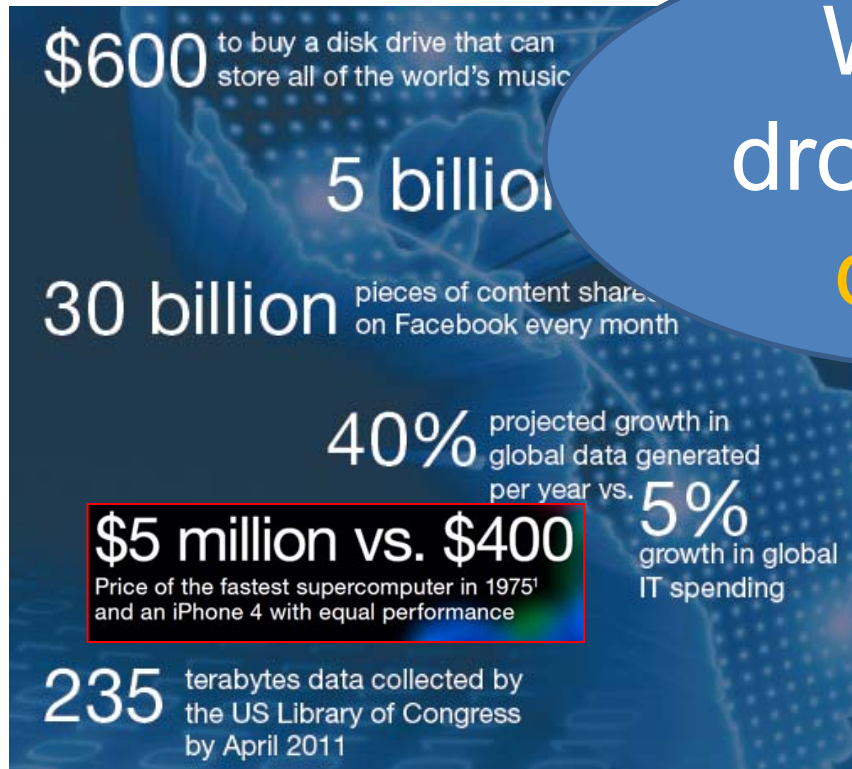
<http://www.popsoci.com/announcements/article/2011-10/november-2011-data-power>



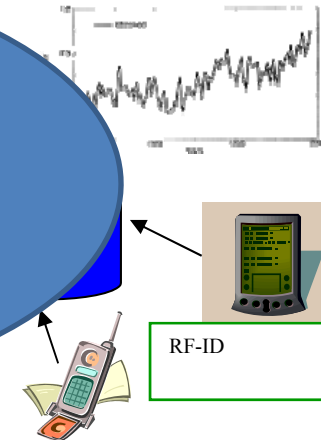
- Data contains value and knowledge

We are drowning in data ... but starving for information

- Exponential grows in data



We are drowning in data...



<http://www.popsoci.com/announcements/article/2011-10/november-2011-data-power>



J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets,
<http://www.mmids.org>

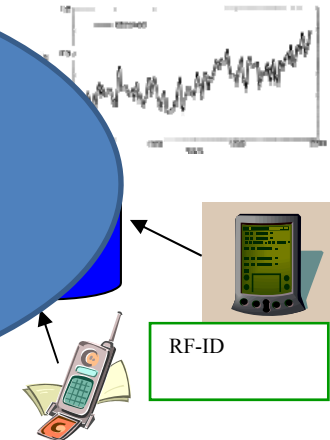
- Data contains value and knowledge

We are drowning in data ... but starving for information

- Exponential grows in data



We are drowning in data...



<http://www.popsoci.com/announcements/article/2011-10/november-2011-data-power>

...but starving for information

Datasets,

- Data can be used to create and knowledge

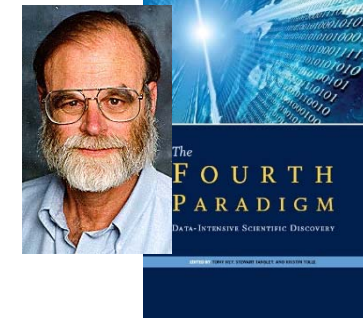


Why this course?

- **Big Data is big**
 - \$ and science: choose your poison
 - Big Data approaches required for Data Science
“move data from raw to relevant”

- **The Fourth Paradigm:**
Age of data driven exploration
→ **Data Science** (eScience / Industry 4.0)

[Informatik Pionier Jim Gray]



[Hey, Tansley, Tolle:
Fourth Paradigm, 2009]

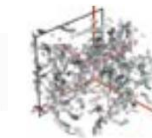
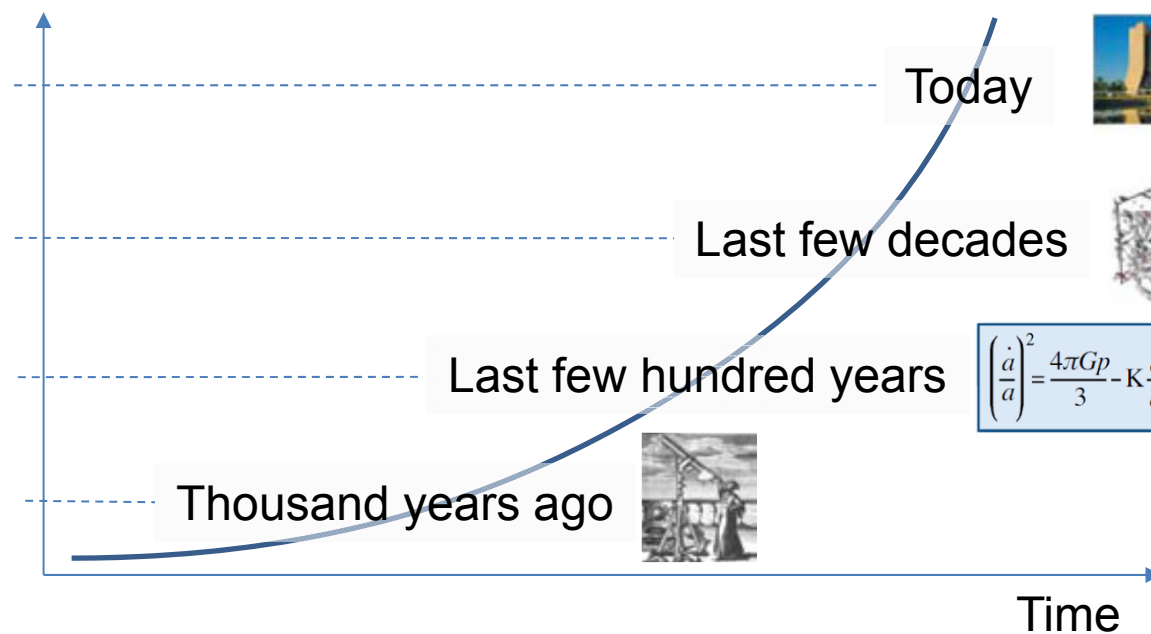
- **Science Paradigms**

Data driven –
Data Science
unify theory,
experiment,
and simulation

Computational –
simulating complex
phenomena

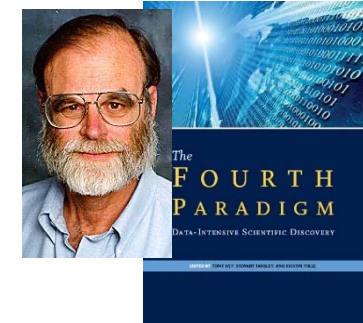
Theoretical –
using models,
generalizations

Empirical -
describing natural
phenomena



- **The Fourth Paradigm:**
Age of data driven exploration
→ **Data Science** (eScience / Industry 4.0)

[Informatik Pionier Jim Gray]

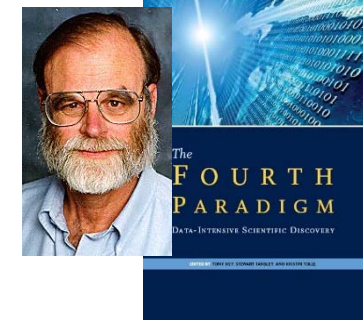


[Hey, Tansley, Tolle:
Fourth Paradigm, 2009]

- **Data Science**
 - Data captured by instruments or generated by simulator
 - Processed by software
 - Information/knowledge stored in computer
 - Scientist/Analyst analyzes database / files using data management and statistics

- **The Fourth Paradigm:**
Age of data driven exploration
→ **Data Science** (eScience / Industry 4.0)

[Informatik Pionier Jim Gray]



[Hey, Tansley, Tolle:
Fourth Paradigm, 2009]

- **Data Science**

- Data gene
Proc
- “Modern science increasingly relies on integrated information technologies and computation to **collect, process, and analyze complex data.**”*

[Hey, Tansley, Tolle: Fourth Paradigm, 2009]

- Information/knowledge stored in computer
- Scientist/Analyst analyzes database / files using data management and statistics

Why this course?

- **Big Data is big**
 - \$ and science: choose your poison
 - Big Data approaches required for Data Science
“move data from raw to relevant”
- **Big Data is exciting**
 - gives a new twist to almost everything
 - allows you to reinvent the wheel

Why this course?

- **Big Data is big**
 - \$ and science: choose your poison
 - Big Data approaches required for Data Science
“move data from raw to relevant”
- **Big Data is exciting**
 - gives a new twist to almost everything
 - allows you to reinvent the wheel
- **Big data is old**
 - opportunity to teach you some fundamental technology

- Introduction (Motivation and Overview)
- Introduction to Big Data — the four V's
- NoSQL
- Hadoop / HDFS / MapReduce & Applications
- Spark
- Data Stream Processing & Applications & Algorithms
- Text Processing
- High-Dimensional Data
- Graph Data Processing
(Link Analysis, Page Rank, Community Detection)
- Uncertain Data Processing
(Concepts of probabilistic query processing and mining)

- This course is mainly based on a mixture of existing external lectures, Surveys, Papers and Reports on Big Data
- There is NO, or better, I'm not aware of a single book or script that is equivalent to this course (and addresses all issues discussed in this course)
- Since Big Data is a quite new and hot topic, standards and basic concepts are quite dynamic => The Web is a very appropriate source of relevant information
- External lectures basically used for this course:
 - Big Data: Donald Kossmann & Nesime Tatbul, Systems Group ETH Zurich - <http://www.systems.ethz.ch/node/217>
 - Mining of Massive Datasets: Jure Leskovec, Anand Rajaraman, Jeff Ullman, Stanford University - <http://www.mmds.org>
- Further material will appear at our web page (check for updates during the course / open to further suggestions!)