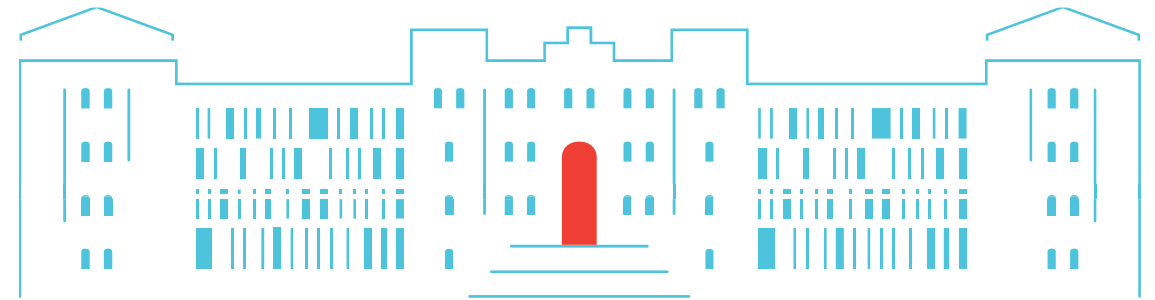
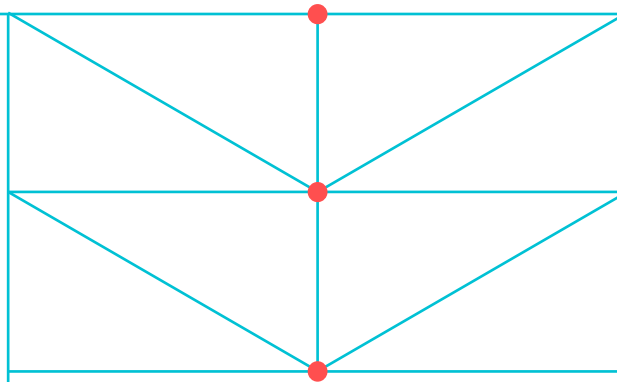


Lecture 01: Introduction Big Data

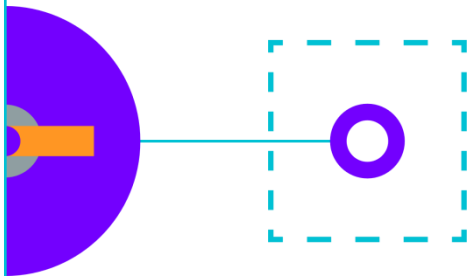
TUHH
Hamburg
University of
Technology



Goals of this Course	TUHH
<ul style="list-style-type: none">▪ Describe basic concepts and characteristics of the “Big Data Stack”, including MapReduce, NoSQL Databases, Stream Processing, as well as according technologies.▪ Discuss and assess the benefits and drawbacks of Big Data technologies▪ Select and apply Big Data software for particular application areas▪ Design and develop practical solutions	
	2

Agenda for Today

1. Lecture and Lab Overview and Organizational Information
2. Overview of the Topics
 - Course Content Overview
 - Big Data: Motivation and Basic Characteristics
 - Data Engineering



Lecturers

TUHH

Team:

- Lecture: Stefan Schulte
- Lab: Nisal Hemadasa



Teaching Approach

TUHH

- Language: Entirely in English
- In general, on-site lectures
 - No on-site lecture on June 27th (instead: lecture recording)
- Recordings of the lectures (but not of the exercises)
 - These may go wrong!
- Utilization of classroom response tools from time to time (Particify and/or EduVote)
- Brief break (2–3 minutes) after 45 minutes
- Interrupt me whenever you think that's a good idea

5

Lecture Overview (Tentative)

TUHH

- Lecture 1 (Today): General Introduction
- Lecture 2: Lab kickoff, Introduction to Big Data Processing
- Lectures 3–5: MapReduce
- Lectures 5–9: Key-Value Stores / NoSQL
- Lectures 9–11: Stream Processing
- Lectures 12–13: Still open

Literature

TUHH

- Unfortunately, there is no single textbook that covers the whole lecture
- M. Kleppmann, *Designing Data-intensive Applications*, O'Reilly Media, Inc., 2017.
- T. Akidau, S. Chernyak, and R. Lax, *Streaming Systems*, O'Reilly Media, Inc., 2018.
- Literature will be provided throughout the semester

Dates and Rooms

TUHH

- Lectures: Thursday, 16:45–18:15, A-1.15
- Tutor Consultation Hours: Thursday, 11:30–13:00 (D-2.022) or 13:15–14:45 (D-2.022), starting on April 18th
- Fixed consultation hours per team, i.e., either the 11:30 or 13:15 slot on Thursday
- There are some exceptions, so please check the schedule in Stud.IP. If there are changes in the semester, there will also be an announcement.

8

Lab Overview

TUHH

- Learning by doing
- Projects from the field of Big Data
 - This semester, focus on data stream processing
 - Topic descriptions will be released on Monday (April 8th)
- Team size: 4 students
- Expected outcomes:
 - Project solution (software prototype)
 - Two presentations

9

Lab Overview II

TUHH

- The lab is the “PBL” part of the module, and therefore makes up 50% of the overall grade, i.e., you cannot pass the module without doing the lab
 - The other half of the grade is coming from the end-of-semester written exam
- But of course, the main aspect is to learn how to work with modern Big Data technologies
- Decide until the kickoff meeting (April 11th) if you want to do it or not
 - Check the topic descriptions provided in Stud.IP
 - But please do not drop out during the semester

10

Lab Overview III

TUHH

- If you really want to do the lab exercise: Sign up for a consultation hour in Stud.IP
 - Start: April 7th, 9:00
 - End: April 9th, 23:59
 - Topics will be distributed at the kickoff meeting next week
- Lab meetings:
 - Kickoff meeting: Next week, as part of the lecture; if there are any things to be clarified, we'll do it there
 - Interim presentations: As part of the consultation hours. **Attendance is mandatory for members of the presenting teams.**
 - Final presentations: July 11th, 10:00 to 15:00. **Attendance is mandatory for all!**

General Information: Stud.IP

TUHH

Content:

- Announcements
- Slides (“Files”)
- Dates and rooms (“Schedule”)
- Content (especially the lab topics), exam regulations and grading (to be added), references, ...
- Lecture recordings (“Media Links”)
- Forum

Questions/requests/comments:

- If the answer could be interesting to more students, please use the forum in Stud.IP!
- If there are further questions w.r.t. the lab, nisal.hemadasa@tuhh.de
- If there are further questions w.r.t. the lecture, contact stefan.schulte@tuhh.de

12

Code of Conduct

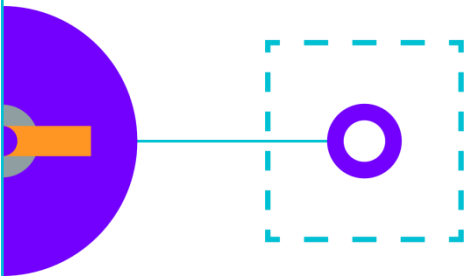
TUHH

- Asking questions is greatly encouraged
 - Discuss questions with each other (except in the exams)
 - Submit lab exercises individually (i.e., per team), but feel free to discuss with other teams
- The limits of collaboration
 - Do not just share solutions with each other – explain your solutions
 - Plagiarism, copying, or other forms of dishonesty **will result in failing the course**
- Communication
 - Write professional, polite emails
 - Use netiquette in forum, email, chats, etc.

13

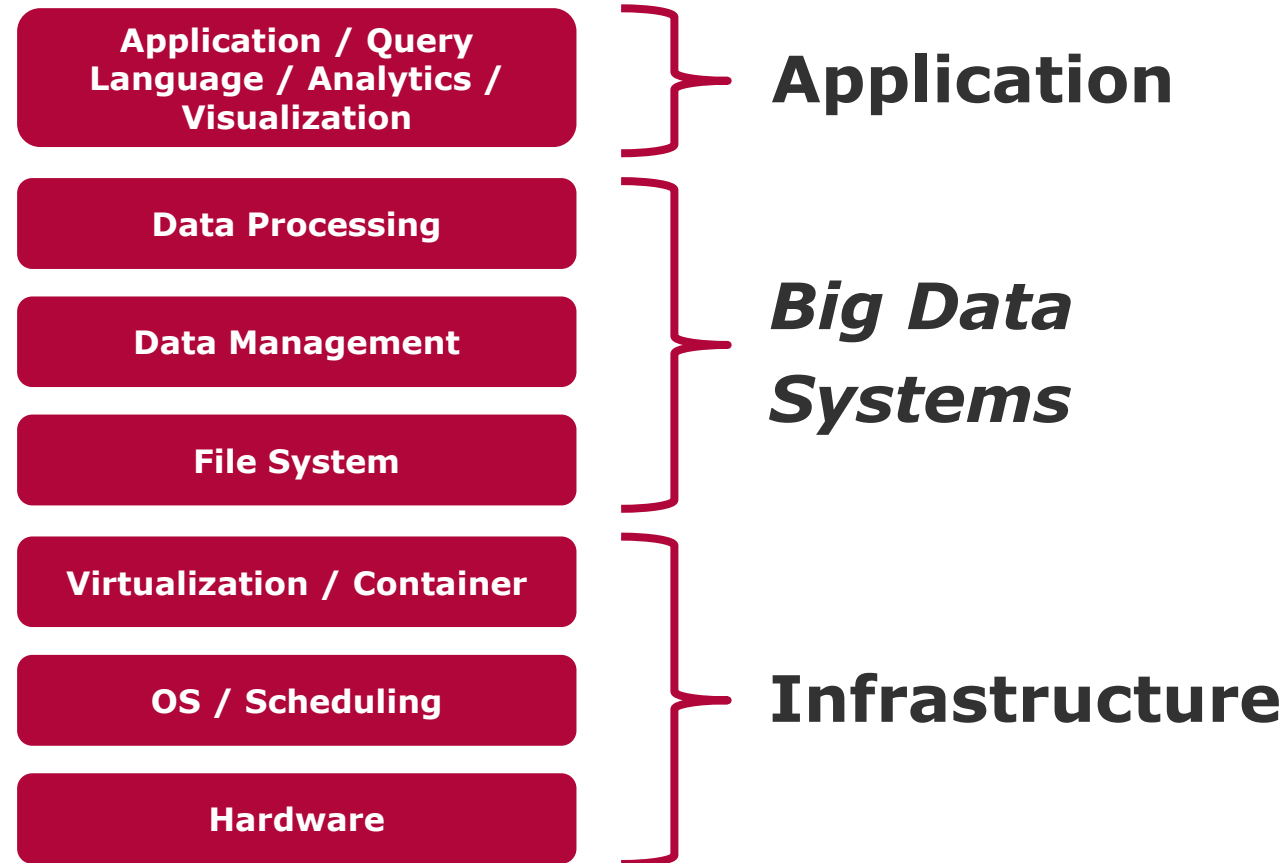
Agenda for Today

1. Lecture and Lab Overview and Organizational Information
2. Overview of the Topics
 - Course Content Overview
 - Big Data: Motivation and Basic Characteristics
 - Data Engineering



Course Contents Overview

TUHH



Big Data Stack

15

Big Data

Data is growing

Messages, tweets, social networks (statuses, check-ins, **shared content**),
blogs, click streams, various logs, ...

- *Facebook: > 3B monthly active users, > 100B messages/day* (Source: Statista, Meta)
- *Instagram: > 2B monthly active users, 66K photos and videos shared per minute* (Source: Statista, LocaliQ)

The value of data is decreasing with its age!



What is Big Data?

TUHH

- **Big data** is an *accumulation* of data that cannot be processed / handled using traditional data management processes / tools, e.g., (object-) relational database management systems
- A *big data management infrastructure* should ensure that the underlying hardware, software, and architecture have *the ability to enable learning (from data) using analytics*.

17

Big Data vs. Small Data

TUHH

Small Data:

- Small enough for human inference
- Accumulated slowly
- Relatively consistent, and (semi-) structured
- Mostly located in storage systems within organizations and data centers

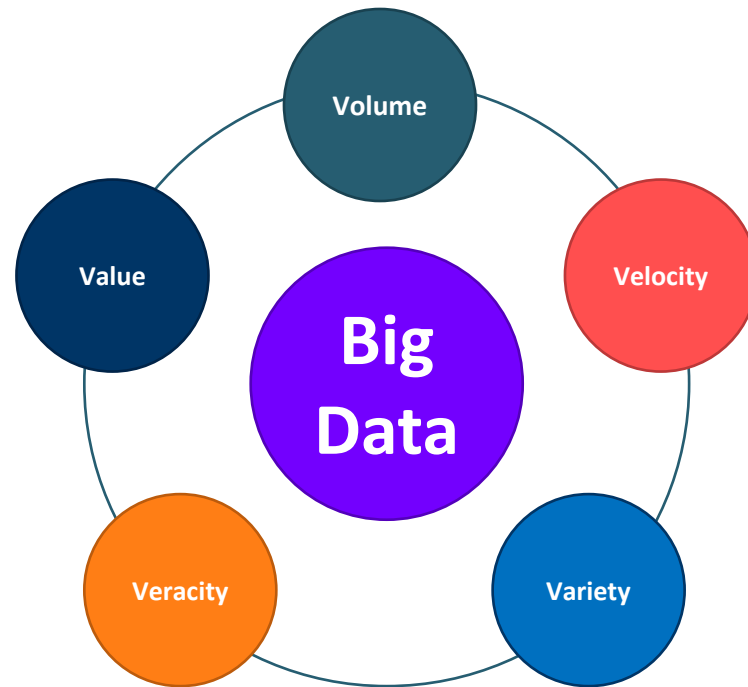
Big Data:

- Generated in huge volumes, and could be structured, semi-structured, or unstructured
- Needs processing to generate insights for human consumption
- Arrives continuously at enormous speed from multiple sources
- Comprises any form of data including video, photos, and more
- Distributed in the cloud or a server farm

18

Big Data: 5 Vs

TUHH



19

Velocity

TUHH

- Data that is generated very fast
 - Data generation rate
 - Data generation never stops
 - Data arrival rate
 - How quick do we have to act on the data?
- Drivers:
 - Improved connectivity and hardware
 - Rapid response times needed in many settings

20

Volume

TUHH

- Scale of the data
- Increased amount of stored data: With respect to the number of observations (size of the data) and the number of variables (dimensionality of the data).
- Drivers:
 - Increase in data sources
 - Higher resolution sensors
 - Scalable hardware infrastructure

21

Variety

TUHH

- Data that comes from machines, people, and processes
- “Data in many forms”: Structured, semi-structured, and unstructured data
- Drivers:
 - Mobile technologies
 - Scalable infrastructure
 - Efficient storage and retrieval

22

Veracity

TUHH

- “Data in doubt”
- Quality, origin, and conformity of facts
- Accuracy of data: Processing errors, noise
- Drivers:
 - Robust ingestion
 - ETL mechanisms (Extract, transform, load)

23

Value

TUHH



Ultimate goal of organizations: Produce value in the form of:

- Faster and smarter (business) decisions
- Increase efficient use of resources
- Discover new (market) opportunities

24

Impact of Big Data on Your Daily Life

Für dich



Text übersetzen
32 Sprachen

Dateien übersetzen
.pdf, .docx, .pptx

DeepL Write BETA
KI-Schreibassistent

Deutsch (erkannt) ▾

↔ Englisch (US) ▾

Glossar

Hey DeepL, was ist ein künstliches neuronales Netz?

Hey DeepL, what is an artificial neural network?

👍👎✎📄🔗

ChatGPT 4 ▾

M You

What would you do on a weekend trip to Philadelphia?

🌀 ChatGPT

A weekend trip to Philadelphia offers a blend of historical exploration, cultural experiences, culinary adventures, and some unexpected discoveries. Here's a mix of activities that could make your weekend memorable:

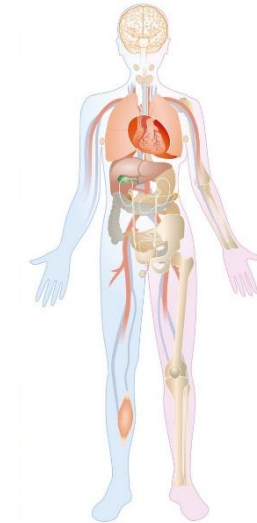
Day 1: Historical and Cultural Exploration

- **Morning: Visit Independence National Historical Park.** Start your day early to beat the crowds at the Liberty Bell and Independence Hall, where the Declaration of Independence and the U.S. Constitution were debated and adopted. Don't miss the chance to see Congress Hall and the National ●

And there is so much more...

TUHH

- Autonomous Driving
 - Requires rich navigation info
 - Rich data sensor readings
 - 1GB data per minute per car (all sensors)¹
- E-health
 - 3.2 billion base pairs of DNA (genomics)
 - 10 million proteins in a person (proteomics)
- Preprocessing of sensor data
 - CERN experiments generate ~1PB of measurements per second
 - Unfeasible to store or process directly, fast preprocessing is a must



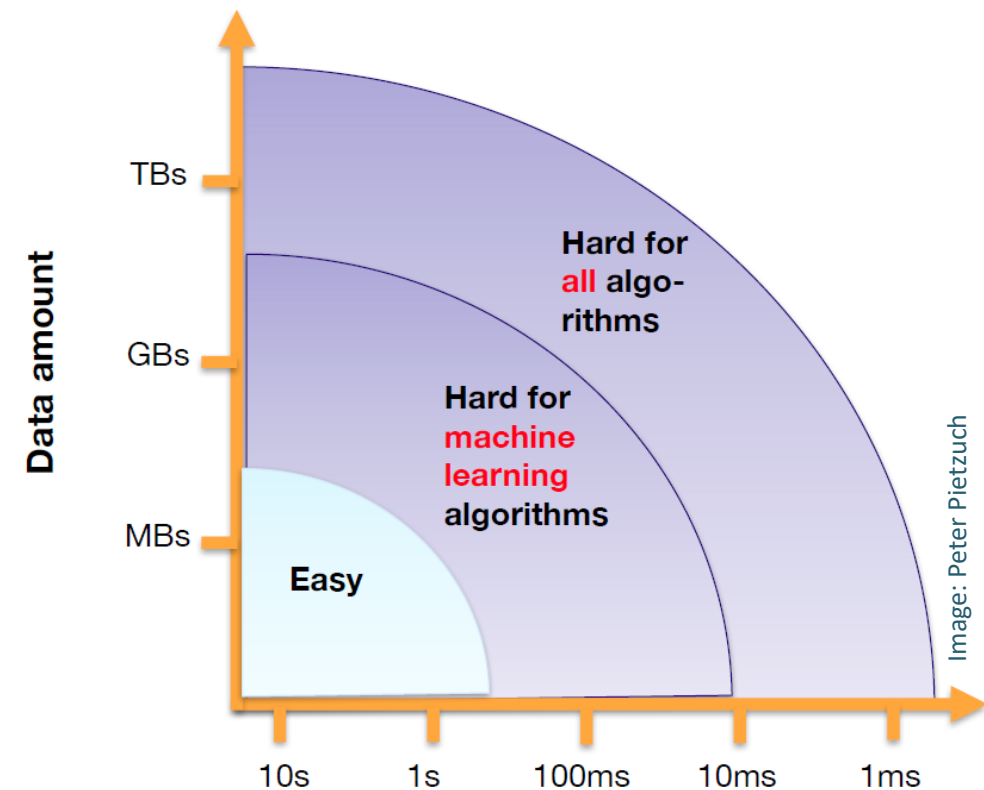
26

¹Cobb: <http://www.hybridcars.com/tech-experts-put-the-brakes-on-autonomous-cars/>

However – Basic Considerations First

TUHH

- Do I have a big data problem?
 - My data fits in memory -> probably no



27

Key Takeaway Messages

TUHH

- Lectures will be recorded
- Register for a lab topic (more precisely: for a certain group)
- Stud.IP should have all the necessary information – if something is missing, please let me know
- The impact of “Big Data”
- The 5 Vs: Velocity, Volume, Variety, Veracity, Value

28

Further Readings / Acknowledgements

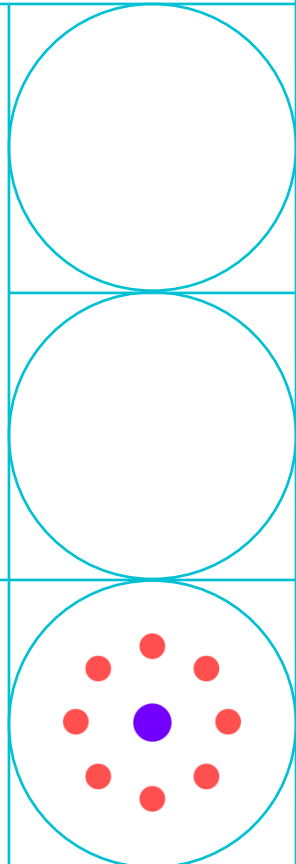
TUHH

- A lot of these slides are (heavily) based on input by Prof. Tilman Rabl, Hasso-Plattner-Institut, Potsdam

Thank You Very Much!

Hamburg University of Technology
Stefan Schulte
Institute for Data Engineering
Blohmstr. 15
21079 Hamburg
Tel.: +49 40 42878-3027
stefan.schulte@tuhh.de

<https://www.tuhh.de/ide/>



TUHH
Hamburg
University of
Technology