



# Distributed Data Management Foundations

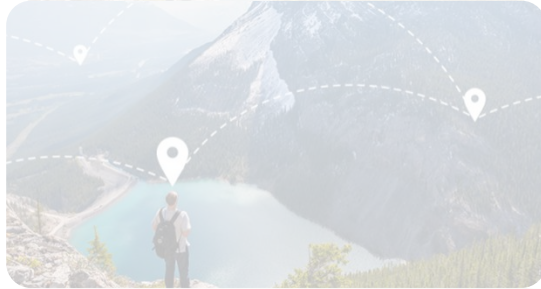
Thorsten Papenbrock

F-2.04, Campus II  
Hasso Plattner Institut

## Big Data



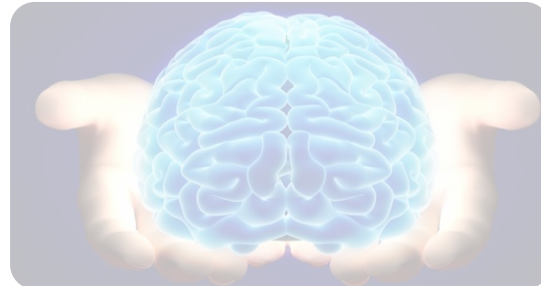
## Data-Intensive Applications



## Consistency Models



## Distributed Computing



## Distributed Data Management

Foundations

ThorstenPapenbrock  
Slide 2

used to refer to the *study and applications of* ~~big~~ ~~large~~ or complex data sets that are so ~~large~~ or complex that traditional ~~database management tools or data processing software~~ are inadequate to deal with them. *data-processing application software*

(Wikipedia ~~2017~~ 2018)

- The challenges include data ...
  - capturing
  - storage
  - extraction
  - curation
  - analysis
  - search
  - sharing
  - transfer
  - visualization
  - querying
  - updating
  - privacy

**Distributed Data  
Management**

Foundations

If data is **too big**, **too fast**, or **too hard** for existing tools to process,  
it is Big Data.

ThorstenPapenbrock  
Slide 3

# Properties of Big Data – Gartner’s 3 V’s

## Volume

- 12 terabytes of Tweets (calculate sentiment analysis)
- 350 billion annual meter readings (predict power consumption)

## Velocity

- 5 million daily trade events (identify potential fraud)
- 500 million daily call detail records (predict customer churn faster)

## Variety

- 100’s of live video feeds from surveillance cameras (find persons)
- 80% data growth in images, videos and documents (improve customer satisfaction)



**Gartner’s 3 V’s:** M. Beyer: Gartner Says Solving „Big Data“ Challenge Involves More Than Just Managing Volumes of Data, [www.gartner.com/it/page.jsp](http://www.gartner.com/it/page.jsp)

**Examples for V’s:** [www.ibm.com/software/data/bigdata](http://www.ibm.com/software/data/bigdata)

## Veracity (Wahrhaftigkeit)

- Trust in correctness and completeness of the data

## Viscosity

- Integration and dataflow friction

## Venue

- Different locations that require different access & extraction methods

## Vocabulary

- Different language and vocabulary

## Value

- Added-value of data to organization and use-case

## Virality

- Speed of dispersal among community

## Variability

- Data, formats, schema, semantics change

### Big Data can be very small:

- Example: streaming data from aircraft sensors
  - A sensor produces an eight byte reading every second (8 byte/sec)
  - Hundred thousand sensors on an aircraft
  - About 2.7 GB of data in an hour of flying (100,000 sensors x 60 min/hour x 60 sec/min x 8 bytes/sec)
  - Difficult to process due to strong real-time requirements and on plane!

### Not all large datasets are “big”:

- Example: video streams plus metadata
  - A live TV stream sends about twenty megabyte per second (20 MB/sec)
  - About 70 GB of data in an hour of streaming (60 min/hour x 60 sec/min x 20 MB/sec)
  - Easy to parse and process, because content is well structured

**Distributed Data  
Management**

Foundations

ThorstenPapenbrock  
Slide 6

➤ The task at hand makes data “big”

## Amazon.com

- Millions of back-end operations every day
- Catalog, searches, clicks, wish lists, shopping carts, third-party sellers, ...



## Walmart

- > 1 million customer transactions per hour
- 2.5 petabytes (2560 terabytes)



## Facebook

- 250 PB, 600TB added daily (2013)
- 1 billion photos on one day (Halloween)



## FICO Credit Card Fraud Detection

- Protects 2.1 billion active accounts



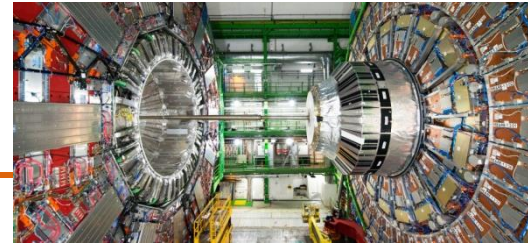
**Distributed Data  
Management**

Foundations

Thorsten Papenbrock  
Slide 7

# Big Data

## Big Data in Use – Science

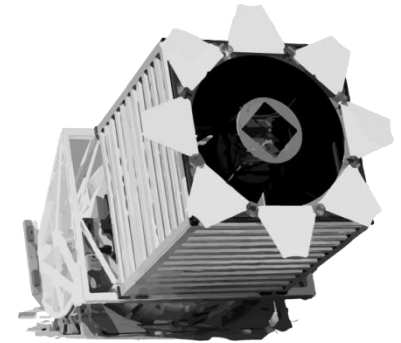


### Large Hadron Collider

- 150 million sensors: 40 million deliveries and 600 million collisions per sec
- Theoretically: 500 exabytes per day (500 quintillion bytes)
- Filtering: 100 collisions of interest per second (→ 99.999% reduction rate)
- 200 petabytes annual rate

### Sloan Digital Sky Survey (SDSS)

- Began collecting astronomical data in 2000
- 200 gigabyte per night; 140 terabytes overall (more data in first few weeks than all data in the history of astronomy)
- Large Synoptic Survey Telescope, successor to SDSS since 2016
  - Acquires that amount of data every five days!

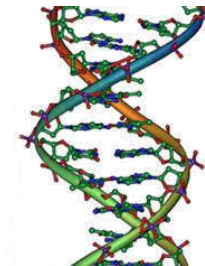


### Distributed Data Management

Foundations

### Human Genome Project

- Human genome: 3,234.83 Mb
- Processing one genome originally took 10 years; now less than a day



ThorstenPapenbrock  
Slide 8



## Correlation

- Correlation describes a linear **statistical relationship** of two random variables (or bivariate data), i.e., the values of both variables change synchronously.

## Causation

- Causation describes a directed, **semantic dependence** of one variable (= cause) to another variable (= effect) such that a change in the first variable always causes a corresponding change in the second variable.
- Correlating variables might share the same causal variable.

➤ Correlation  $\neq$  Causation

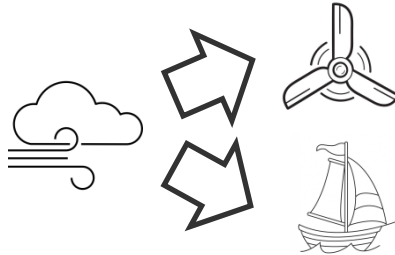
## Correlation

- “energy production of wind turbines” and “top-speed of sailing boats”



## Causation

- “wind speed” causes “energy production of wind turbines”
- “wind speed” causes “top-speed of sailing boats”

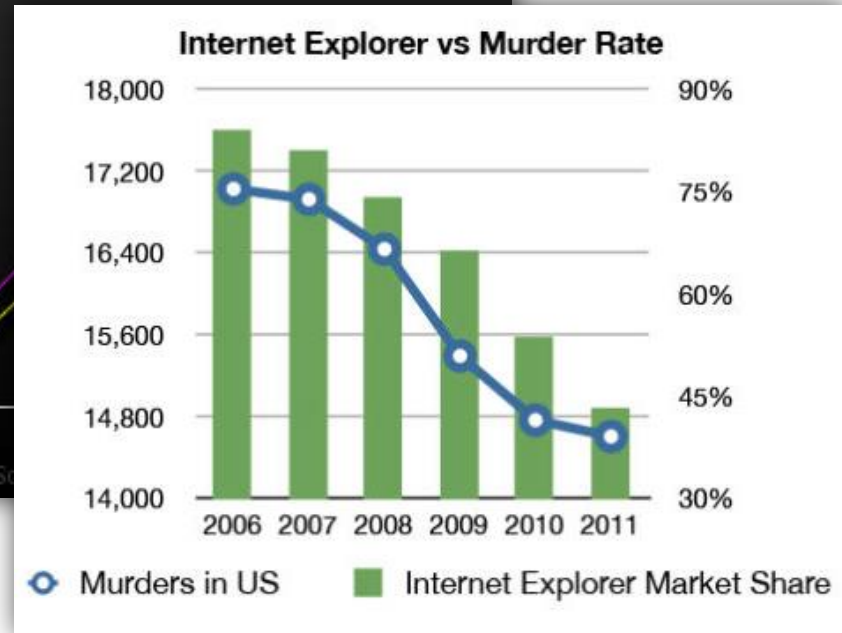
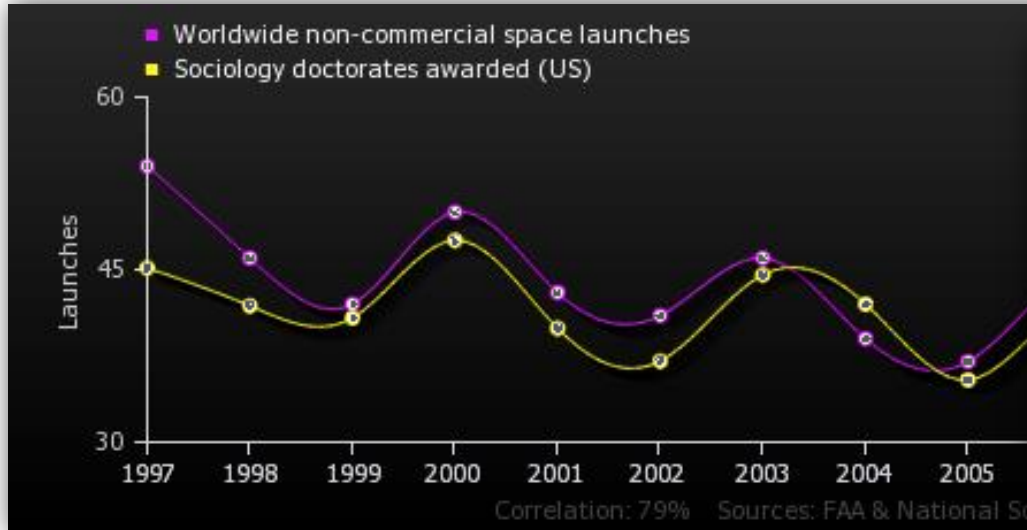


➤ Correlation  $\neq$  Causation

# Big Data

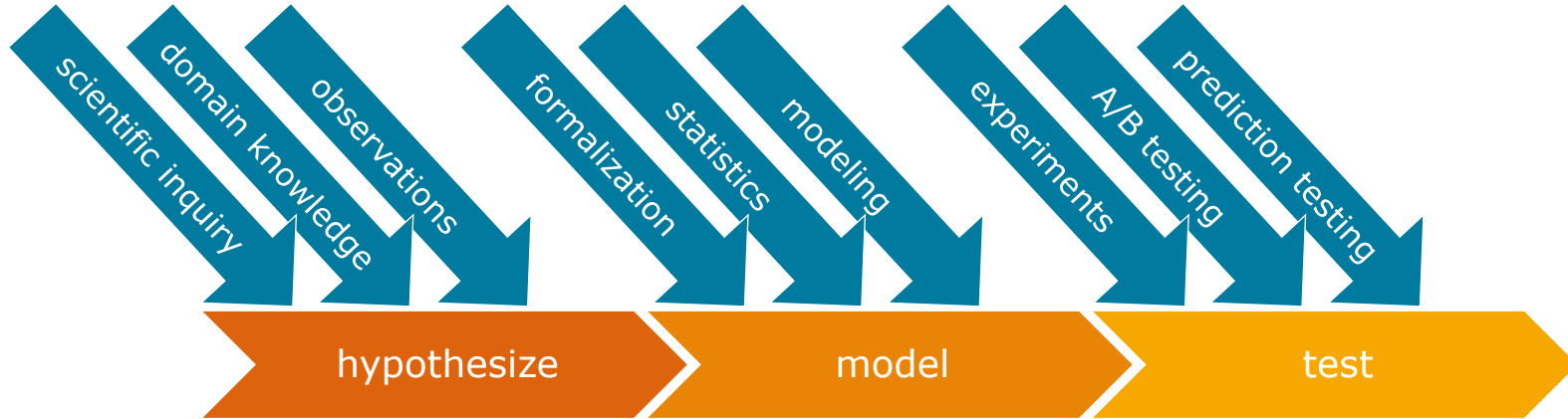
## Correlation vs. Causation

- Correlation  $\neq$  Causation
  - Examples:



# Correlation vs. Causation

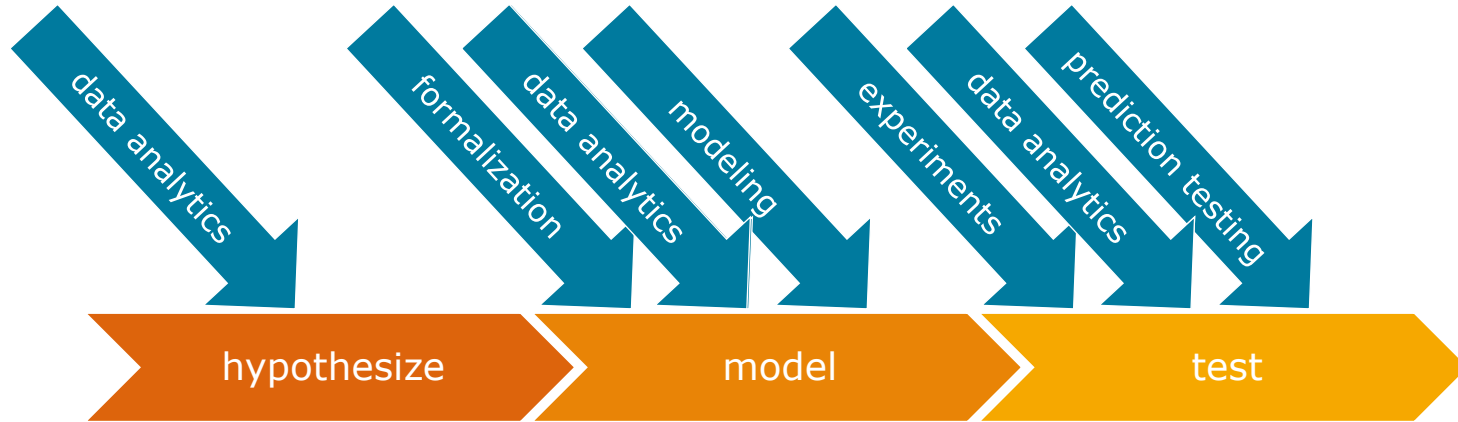
- Correlation  $\neq$  Causation
  - Good science **before** Big Data:



# Big Data

## Correlation vs. Causation

- Correlation  $\neq$  Causation
  - Good science **with** Big Data:

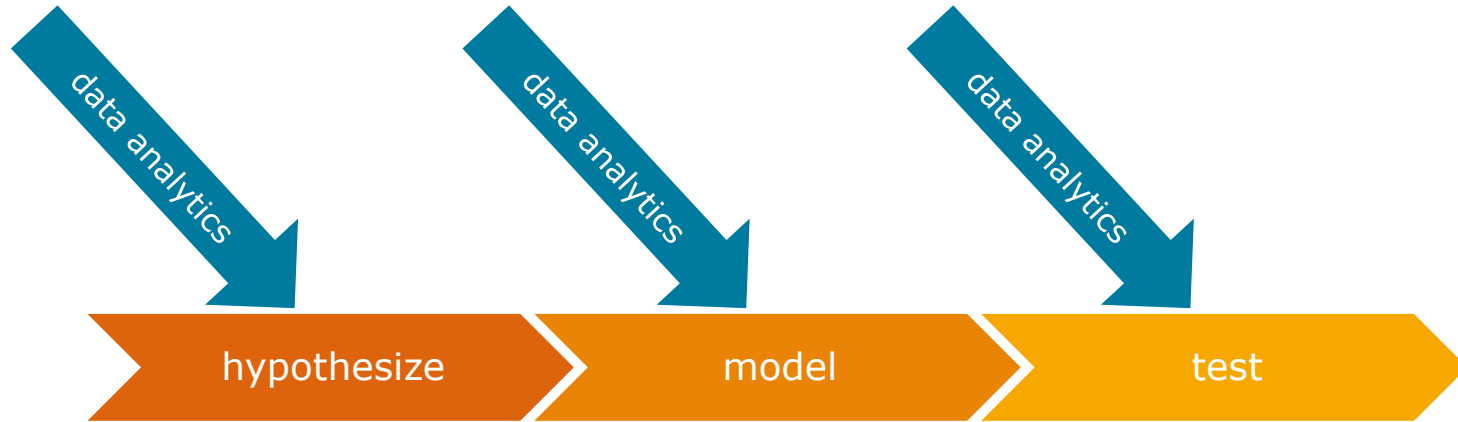


- Hypothesizing is hard: Use discovered correlations to formulate them!
- Modeling is hard: Use automatically trained models!
- Testing is hard: Use Big Data to verify your model!

# Big Data

## Correlation vs. Causation

- Correlation  $\neq$  Causation
  - Good science **with** Big Data:



- If correlation holds **for very large data sets**, it's likely a causation.
  - Big Data Analytics: find correlations  $\rightarrow$  derive causations

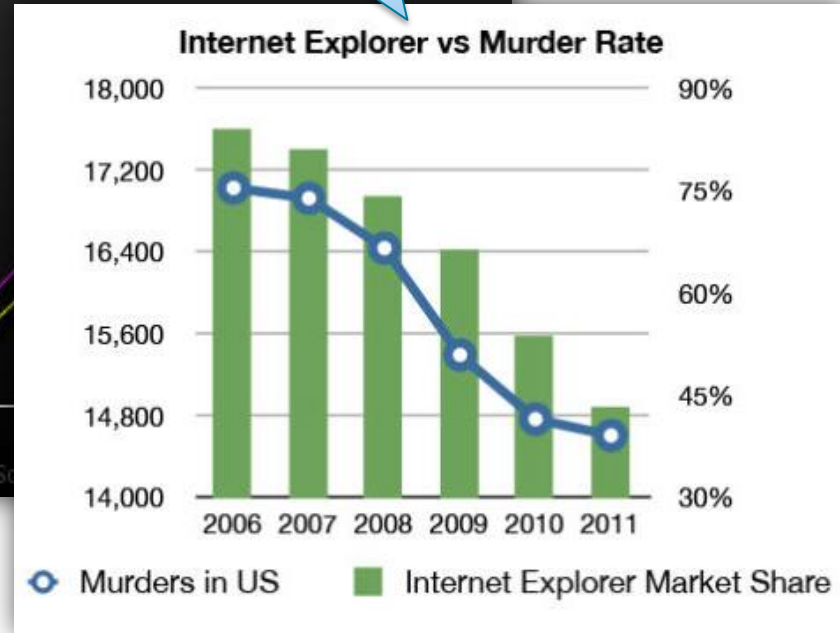
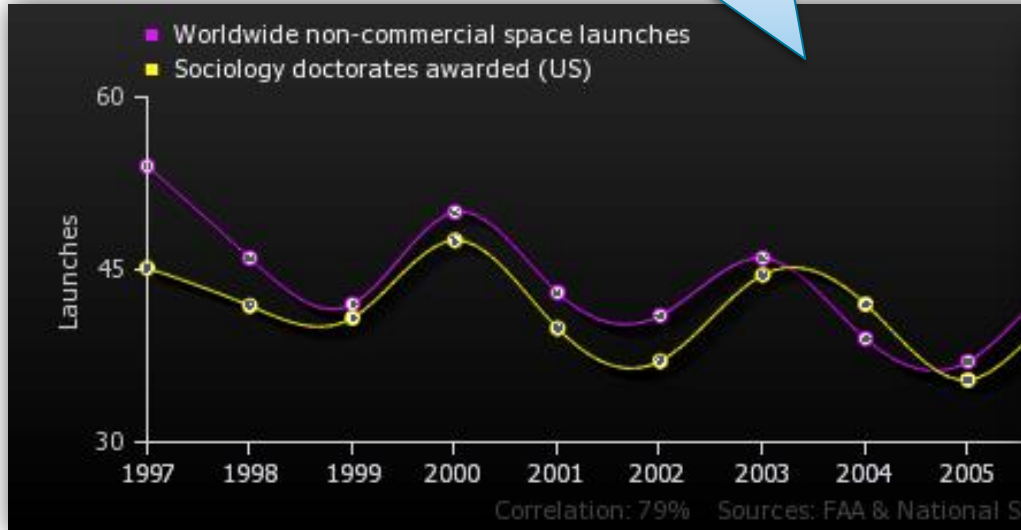
# Big Data

## Correlation vs. Causation

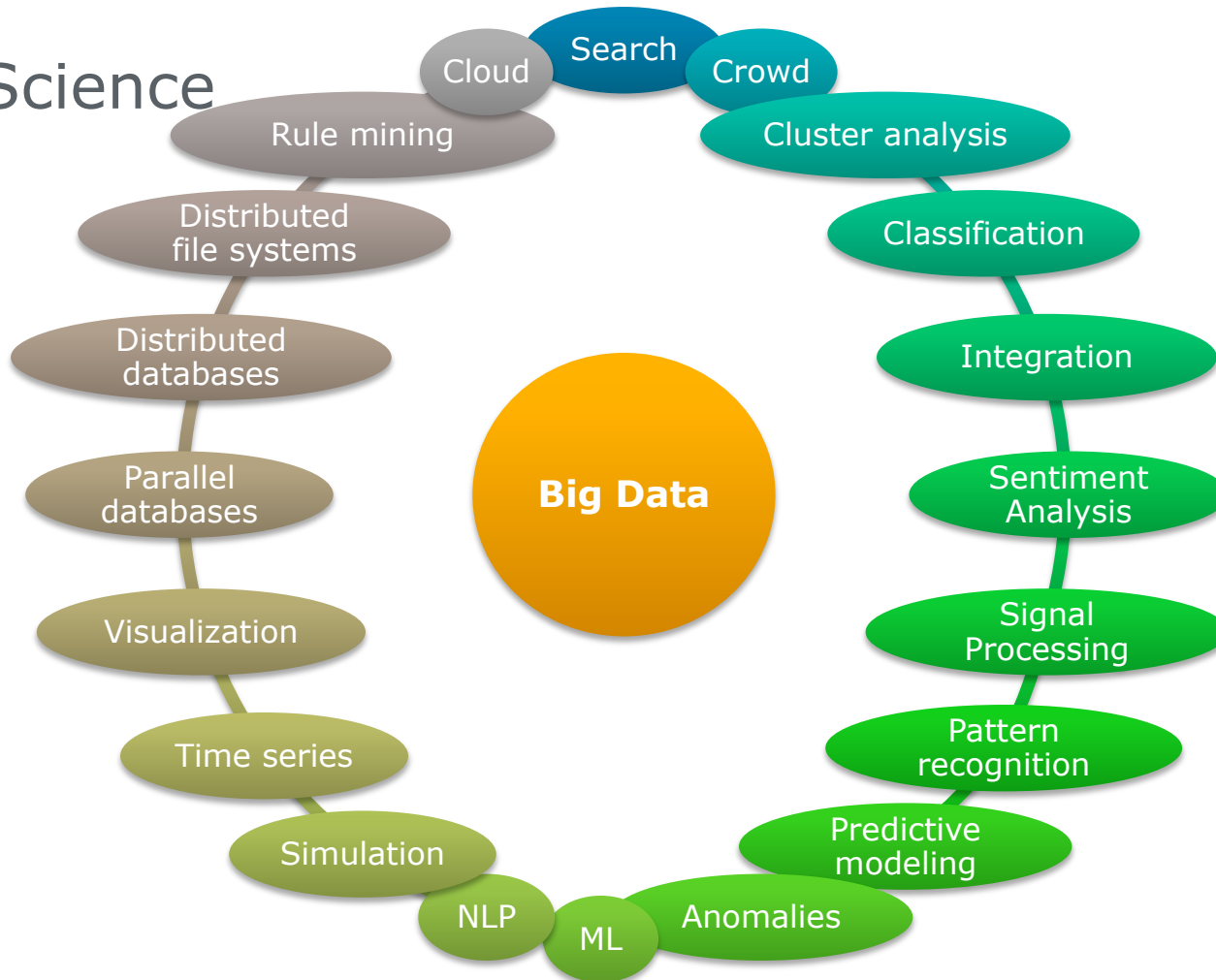
- Correlation = Causation ?
  - Why did it fail here?

13 measurements

6 measurements



# Big Data Data Science

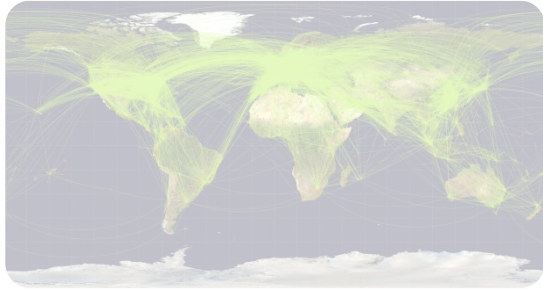


## Distributed Data Management

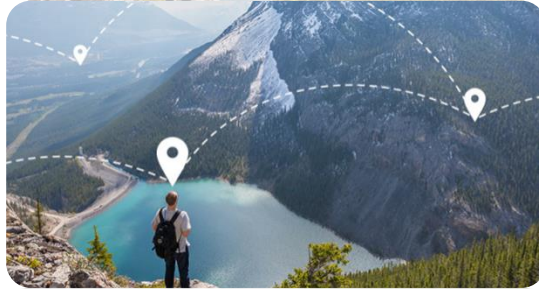
Foundations



## Big Data



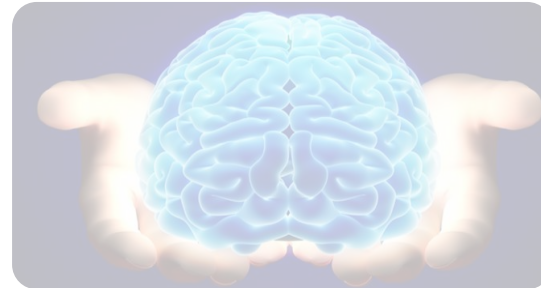
## Data-Intensive Applications



## Consistency Models



## Distributed Computing



## Distributed Data Management

Foundations

ThorstenPapenbrock  
Slide 17

### Databases

- Data storage and persistence

### Search indexes

- Keyword search and filtering

### Caches

- Optimization of expensive and re-occurring queries

### Visualization

- Presentation of data and control options to human users

### Batch processing

- Processing of large amounts of accumulated data (transform, analyze)

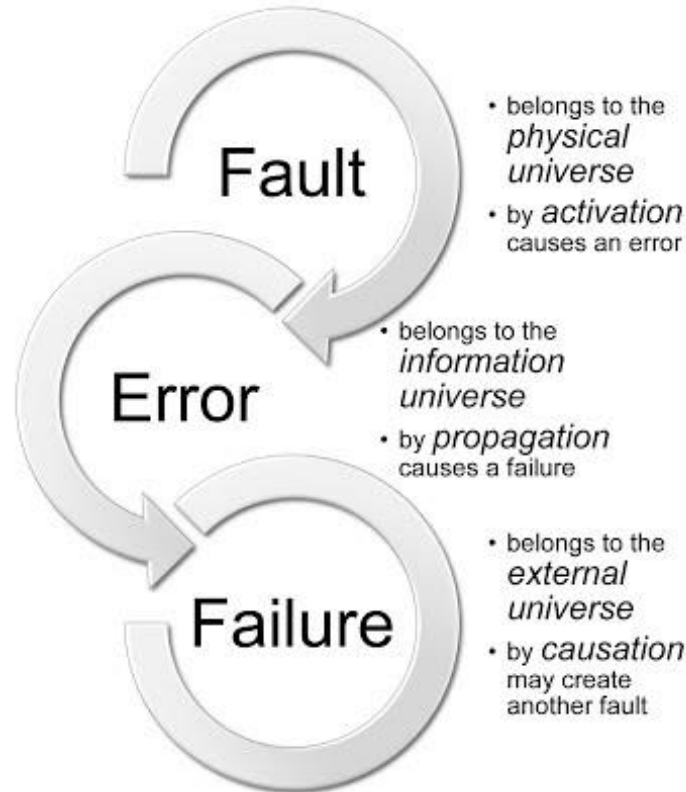
### Stream processing

- Processing of continuous data flows (operate, analyze, store)

#### Design Concerns

1. Reliability
2. Scalability
3. Maintainability

### Reliability



**Distributed Data Management**

Foundations

Thorsten Papenbrock  
Slide **19**

### Reliability

- “The system continues to work **correctly** (= correct functionality at the desired level of performance) even in the face of **adversity** (= hardware or software faults; human faults).”

- = *fault-tolerance*:



- Techniques to ensure Reliability:
  - **Careful design** (clear interfaces, decoupling of code, ...)
  - **Testing** (fault-injection, unit/integration/system/random tests, ...)
  - ! ▪ **Redundancy** (RAID systems, failover systems, backups, ...)
  - ! ▪ **Process isolation** (allowing processes to crash and restart)
  - **Measuring, monitoring, and analyzing** system behavior in production

### Scalability

- “The system supports **growths** (in data volume, traffic volume, or complexity) with reasonable ways of dealing with it (e.g. more resources).”
- **Load:**
  - = *measure to quantify scalability*
  - E.g.: requests per second (= throughput), cache hit rate, read/write ratio to disk, ...
- **Performance:**
  - = *load a system can handle*
  - Usually calculated as the mean, median, or x-percentile of load measurements
- **Reasoning:**
  - a) How does an increasing load with fixed resources affect performance?
  - b) How much must the resources be increase when the load increases and the performance should be fix?

# Data-Intensive Applications

## Design Concerns

### Scalability (cont.)

- Approaches to cope with load:
  - Vertical scaling** (scale up)
    - Add CPUs, RAM, Disk
    - Replace entire machine
  - Horizontal scaling** (scale out)
    - Add additional machines
- Scalable software design:
  - Manual scaling (a human scales the system resources manually)
  - Elastic scaling (the system automatically adds resources if the load increases)

➤ Easier for programmers

➤ More expensive

scale up



scale out

The default strategy for the past 40 years.

Became increasingly important in the past years; probably the future default.

**Distributed Data Management**

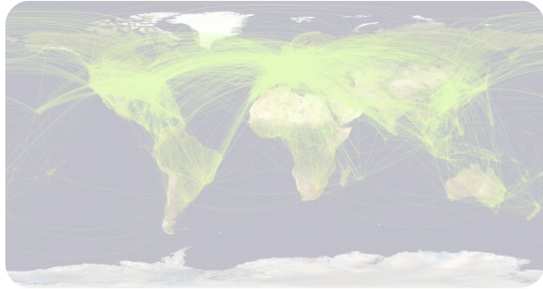
Foundations

ThorstenPapenbrock  
Slide 22

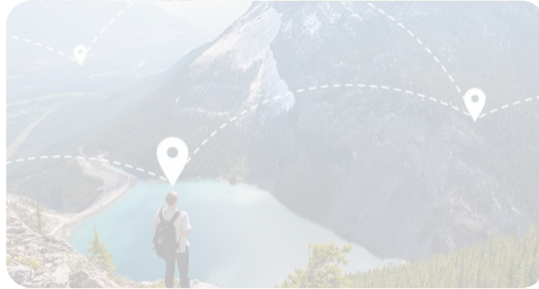
### Maintainability

- “The system allows its productive, further **development** by different engineers at different times in its operation.”
- Design principles to achieve maintainability:
  - **Operability**: Make it easy for operators to **keep the system running**.
    - Monitoring, documentation, testing, design patterns, ...
  - **Simplicity**: Make it easy for engineers to **understand the system**.
    - Clear interfaces, abstraction layers, no over-engineering, ...
  - **Evolvability**: Make it easy for engineers to **change the system**.
    - Agile techniques, test-driven development, pair programming, ...
- See lectures “Software-Architecture” and “Software-Technique” for details!
- See also: “Spotify Engineering Culture”  
<https://labs.spotify.com/2014/03/27/spotify-engineering-culture-part-1/>

## Big Data



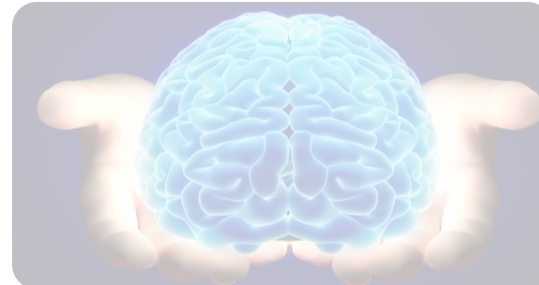
## Data-Intensive Applications



## Consistency Models



## Distributed Computing



## Distributed Data Management

Foundations

Thorsten Papenbrock  
Slide 24



### ACID

- The ACID consistency model stands for the following four guarantees:
  - **Atomicity**: All operations in a transaction succeed or every operation is rolled back.
  - **Consistency**: Before the start and after the completion of a transaction, the database is structurally sound.
  - **Isolation**: Transactions do not contend with one another. Contentious access to data is moderated by the database so that transactions appear to run sequentially.
  - **Durability**: The results of applying a transaction are permanent, even in the presence of failures.
- Requires moderated data access, locks, and failover protection
- Ensures a safe and reliable data storage environment for applications

### CAP Theorem

- It is impossible for a distributed data store to simultaneously provide more than two out of the following three guarantees:
  - **Consistency**: Every read receives the most recent write or an error. This condition includes consistency from ACID, i.e., consistent transaction processing, but also widens the scope from an individual node's data consistency to cluster-wide data consistency.
  - **Availability**: Every request receives a (non-error) response – without guarantee that it contains the most recent write. Server crashes, query congestion, or resource overload may deny service availability.
  - **Partition tolerance**: The system continues to operate despite an arbitrary number of messages being dropped (or delayed) by the network between nodes. Only total network failure might cause the system to respond incorrectly.

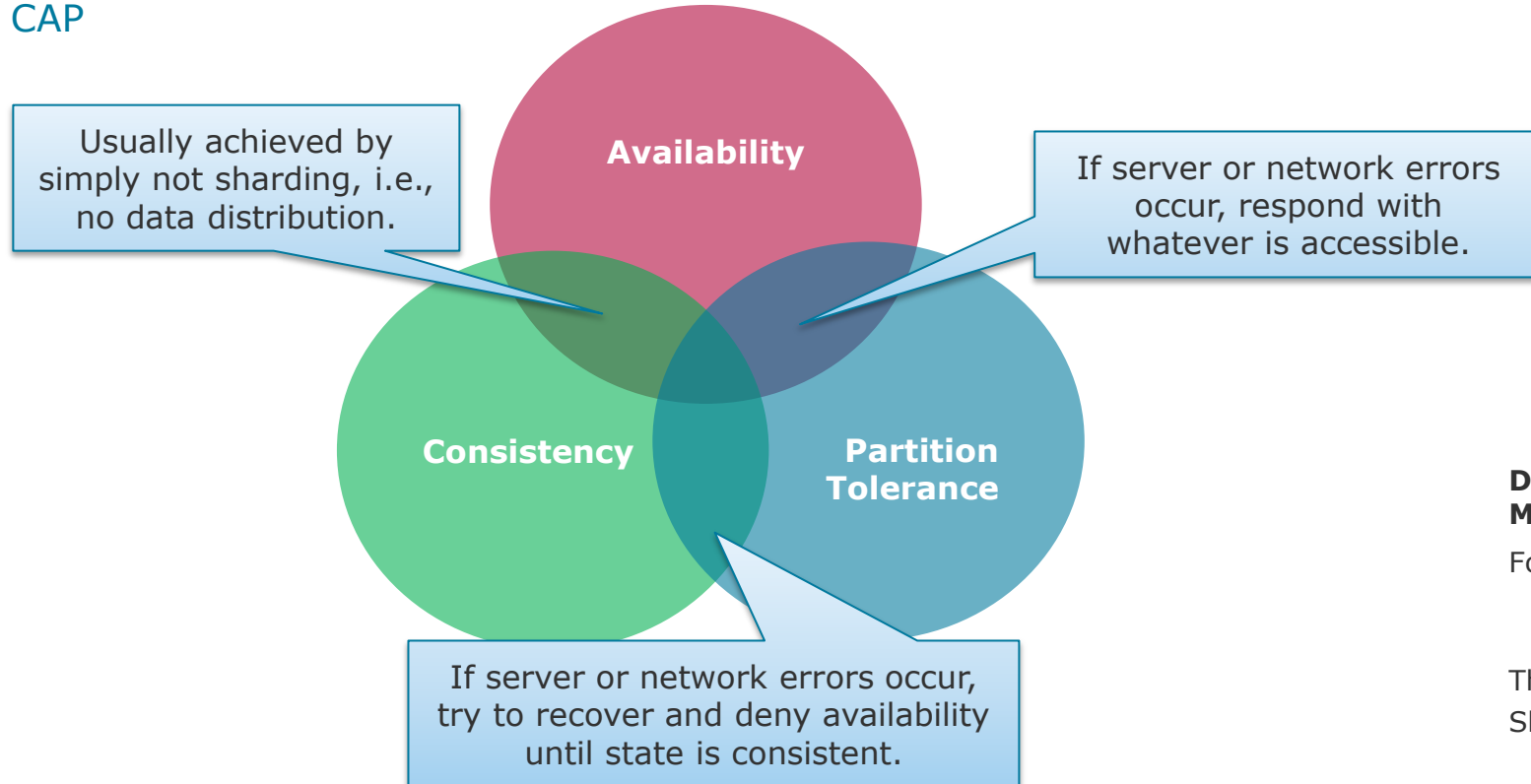
Usually stores achieve all three, but they must drop one dimension **if they are distributed and errors occur.**

#### Distributed Data Management

Foundations

ThorstenPapenbrock  
Slide **26**

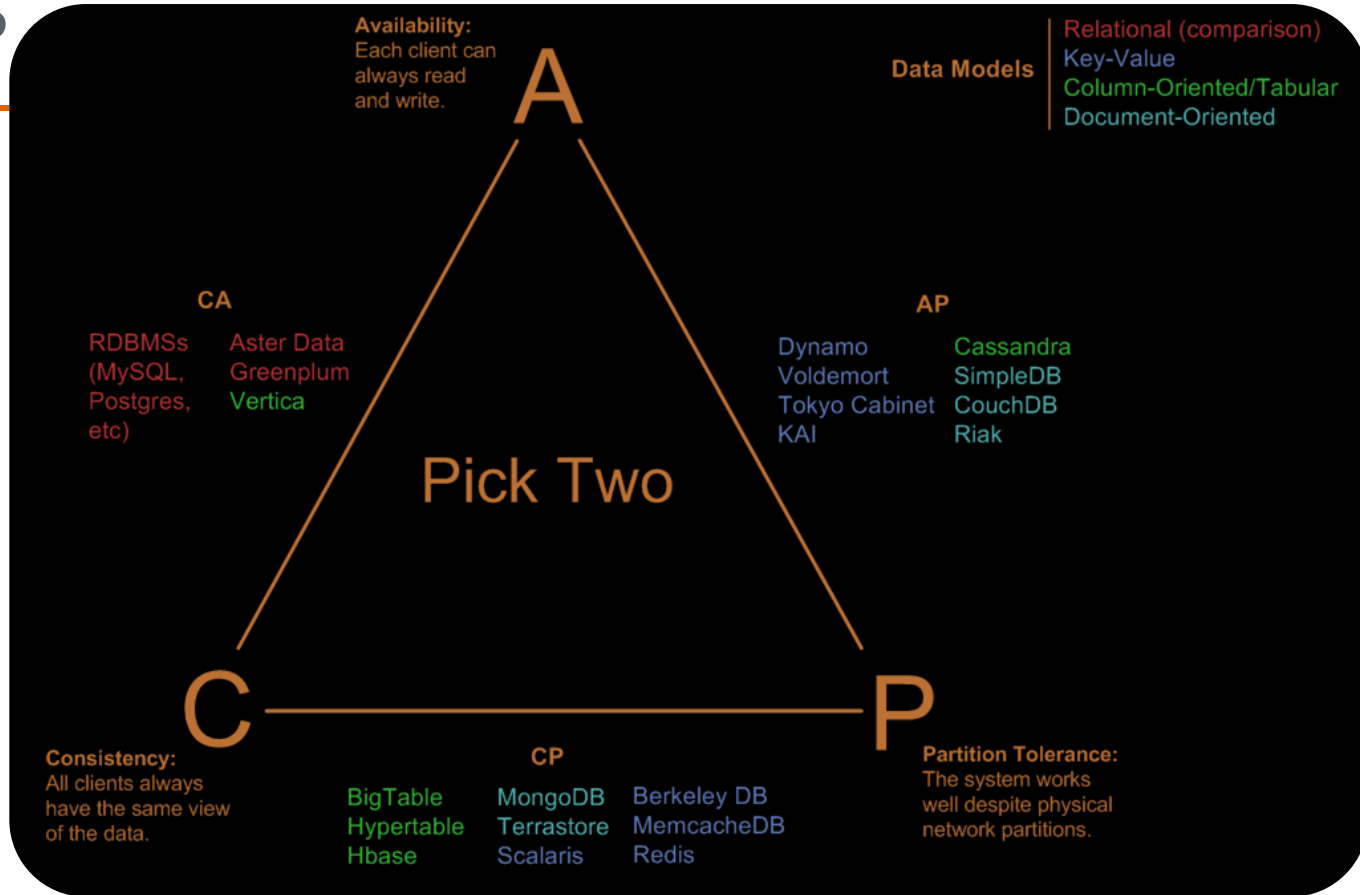
### CAP



# Consistency Models

## CAP

CAP



### Distributed Data Management

Foundations

Thorsten Papenbrock  
Slide 28

### BASE

BASE = “not (fully) ACID”

- The BASE consistency model relaxes CAP dimensions:
  - **Basic Availability**: The database appears to work most of the time.
    - Availability might be less than 100%
    - “Most of the time” is often quantified as lower bound, e.g., 90%
  - **Soft-state**: Stores don’t have to be write-consistent, nor do different replicas have to be mutually consistent all the time.
    - Stored data might be inconsistent, but the store can derive consistent states
  - **Eventual consistency**: Stores exhibit consistency at some later point (e.g., lazily at read time).
    - Usually consistent within milliseconds
    - Does not mean “no-consistency”, which would be fatal for a store

### BASE

- In comparison to ACID **often** means:

ACID	BASE
Transactions	Programmer managed
Strong consistency	Weak consistency
Isolation	Last write wins
Robust database	Simple database
Simpler application code	Harder application code
Conservative (pessimistic)	Aggressive (optimistic)

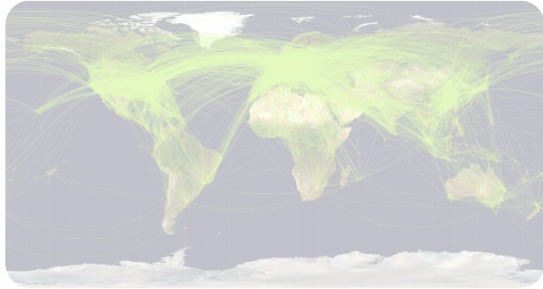
**Distributed Data Management**

Foundations

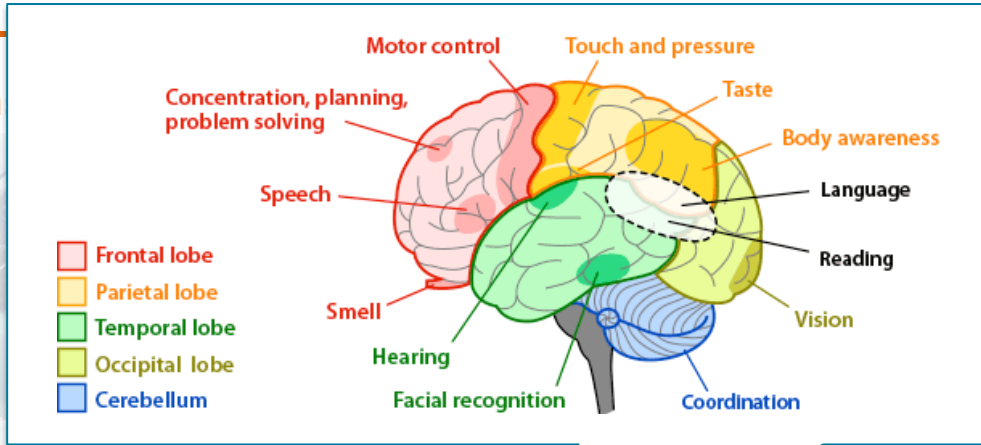
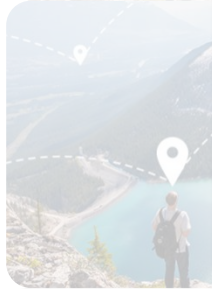
Thorsten Papenbrock  
Slide **30**

# Overview Foundations

## Big Data



## Data-Inten



## Consistency Models



## Distributed Computing



### Distributed Data Management

Foundations

ThorstenPapenbrock  
Slide 31

# Distributed Computing Definition

## What is a distributed system?



**One machine**

**Data** in multiple caches, in memory, on disk ...

**Control-flow** over multiple cores, CPUs, GPUs, ...



**One big machine**

**Specialized racks** with shared infrastructure ...



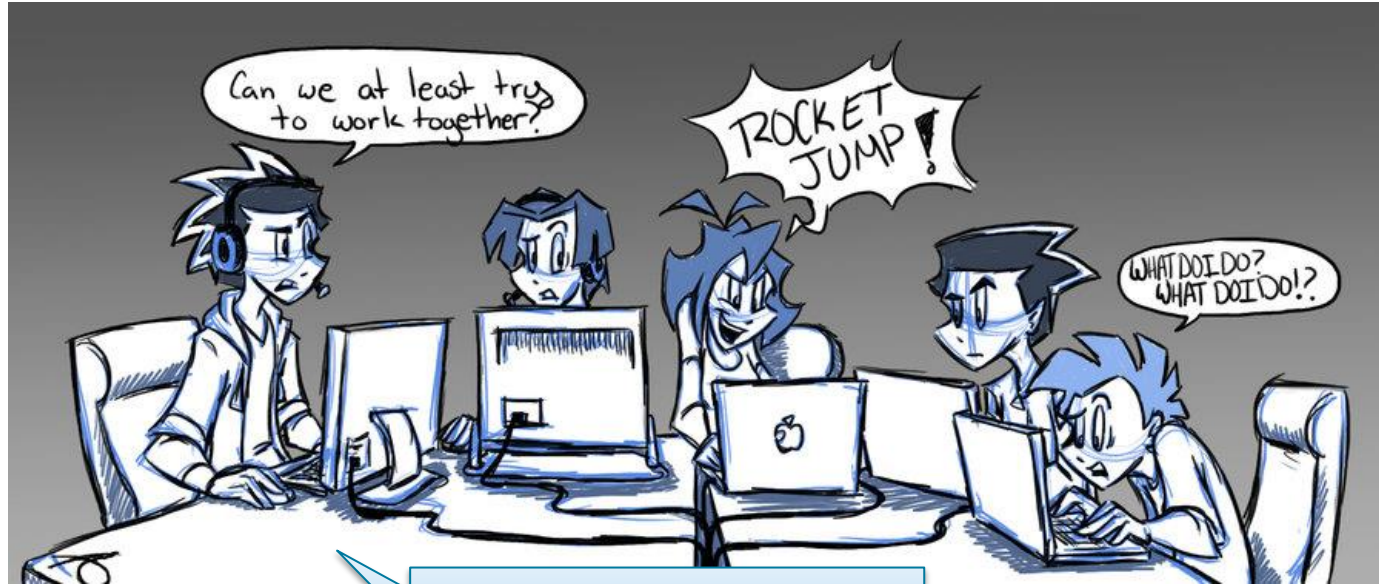
**Independent systems** connected via network

**Multiple, connected machines**



# Distributed Computing Definition

## What is a distributed system?



Do the system components  
need to work together?

**Distributed Data  
Management**

Foundations

ThorstenPapenbrock  
Slide **33**

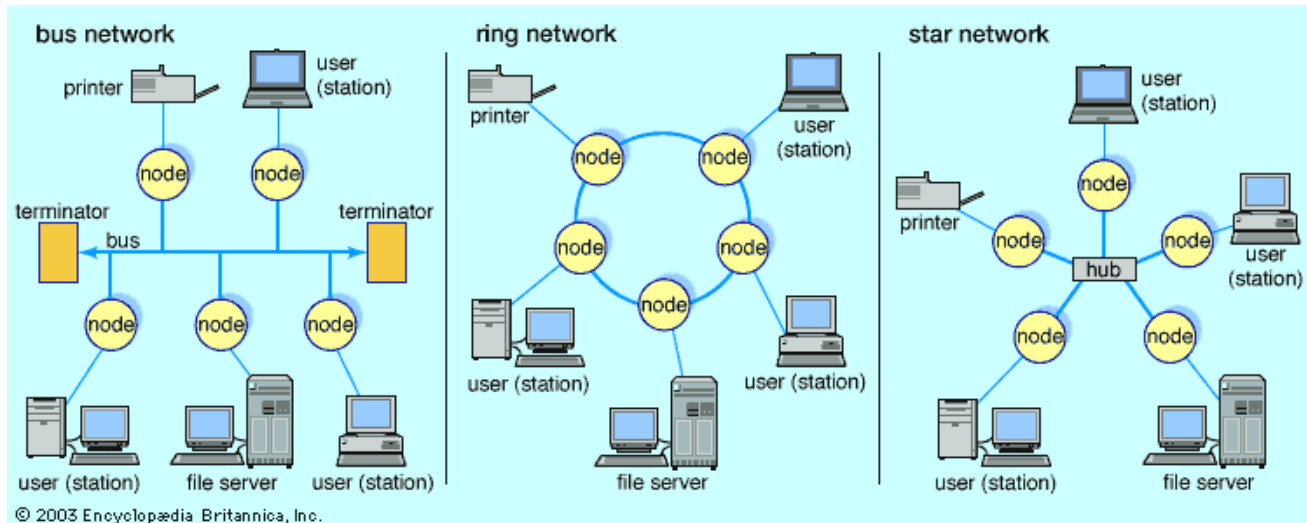
# Distributed Computing Definition

shared-nothing systems

## Practical Definition:

“A *distributed computing system* [...] is a number of **autonomous processing elements** (not necessarily homogeneous) that are interconnected by a **computer network** and that **cooperate** in performing their assigned task.”

(M. Tamer Özsu, Patrick Valduriez: “Principles of Distributed Database Systems”)



**Distributed Data  
Management**

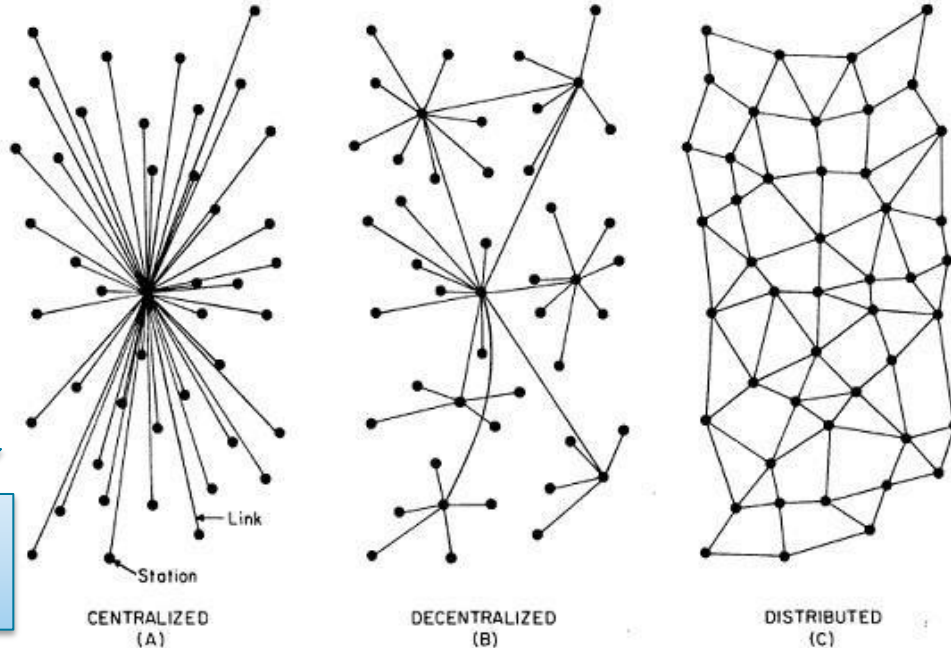
Foundations

ThorstenPapenbrock  
Slide **34**

# Distributed Computing Definition

## Topological Definition:

“A *distributed computing system* is a (fully) decentralized network of computing elements/stations, i.e., one that has multiple roots.”



single-client or  
single-master  
systems

peer-to-peer  
systems

**Distributed Data  
Management**

Foundations

ThorstenPapenbrock  
Slide 35

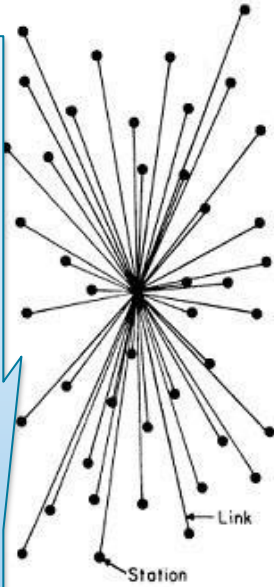
# Distributed Computing Definition

## Topological Definition:

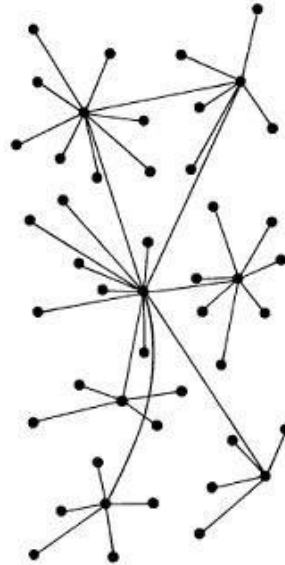
“A *distributed computing system* is a (fully) decentralized network of computing elements/stations, i.e., one that has multiple roots.”

### Examples:

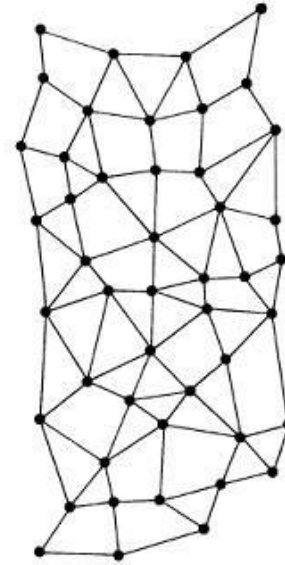
- Weather stations and their central control station
- Human workers and the central MTurk web service in Amazon Mechanical Turk



CENTRALIZED  
(A)



DECENTRALIZED  
(B)



DISTRIBUTED  
(C)

### Examples:

- BitTorrent file sharing clients
- Bitcoin miner networks
- InterPlanetary File System (IPFS) that connects arbitrary computers to a DFS storing hypermedia

Distributed Data  
Management

Foundations

StenPapenbrock

36

# Distributed Computing

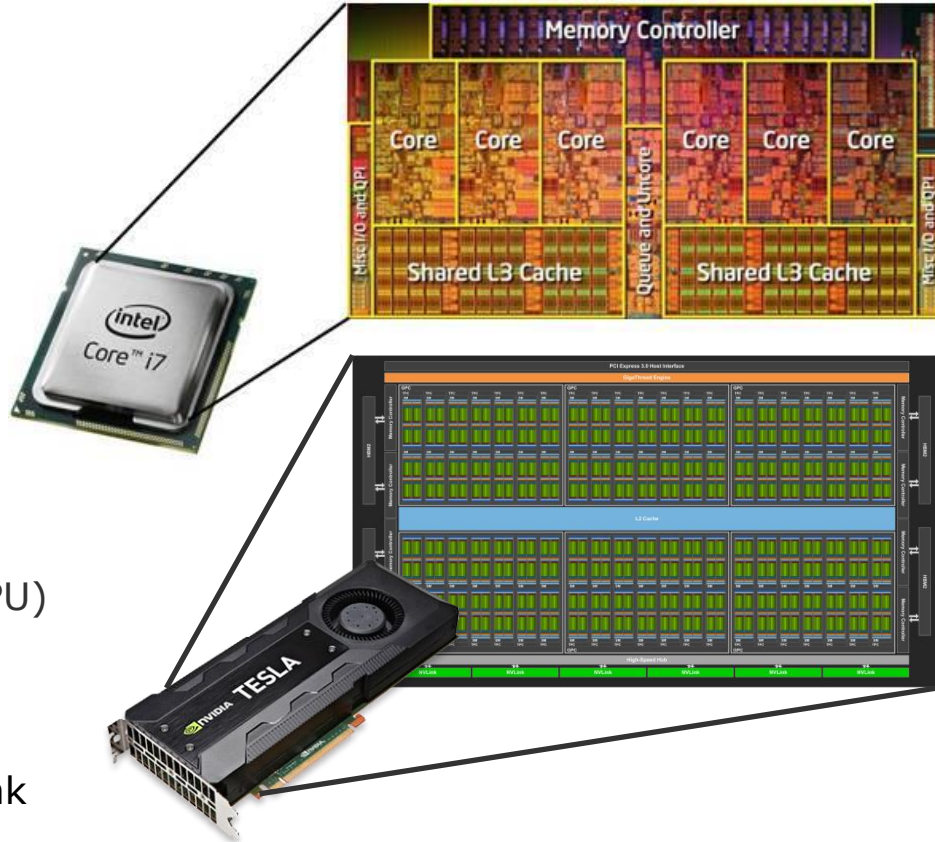
## Parallel Computing

### Parallelization

- Multiple processing units perform work simultaneously, i.e., in parallel
- Long tradition in databases
- One approach to address Big Data issues

### Trends

- **Multicore CPUs**
  - E.g. `java.util.concurrent` or `pthread`
- **General-purpose computing on GPUs (GPGPU)**
  - E.g. OpenCL or CUDA
- **Cluster frameworks**
  - E.g. Hadoop MapReduce, Spark, or Flink



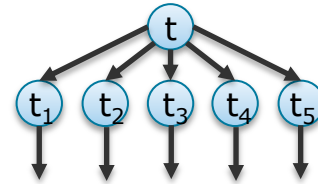
## Distributed computing vs. multi-threading:

- **Shared nothing:**
  - Communication and data sharing only via messaging
  - No shared memory, shared process resources, shared error handling, shared garbage collection, ...
- **Autonomous systems:**
  - Synchronization only via messaging
  - No mutexes, semaphores, atomic counters, lock-free data structures, blocking queues, ...
- **More constricted parallelism:**
  - A distributed algorithm can run parallel on one machine but a multi-threaded algorithm (usually) cannot run on many machines.

### Approaches

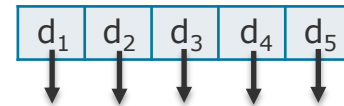
- **Task parallelism:**

- Breaks the task into sub-tasks that are processed in parallel
- Each processing unit performs a different subtask
  - Usually OLTP: Akka, RabbitMQ, Kafka, ...



- **Data parallelism:**

- Breaks the data of a task into packages that are processed in parallel
- Each processing unit performs the same task on different data
  - Usually OLAP: MapReduce, Spark, Flink, ...



- **Instruction-level parallelism:**

- Breaks the task into instructions that are processed in parallel
- One processing element performs multiple instructions simultaneously
  - In hardware: instruction pipelining, superscalar, branch prediction, ...

**Distributed Data Management**

Foundations

Thorsten Papenbrock

Slide 39

