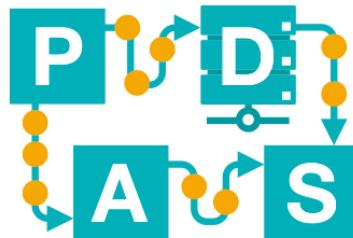


Big Data

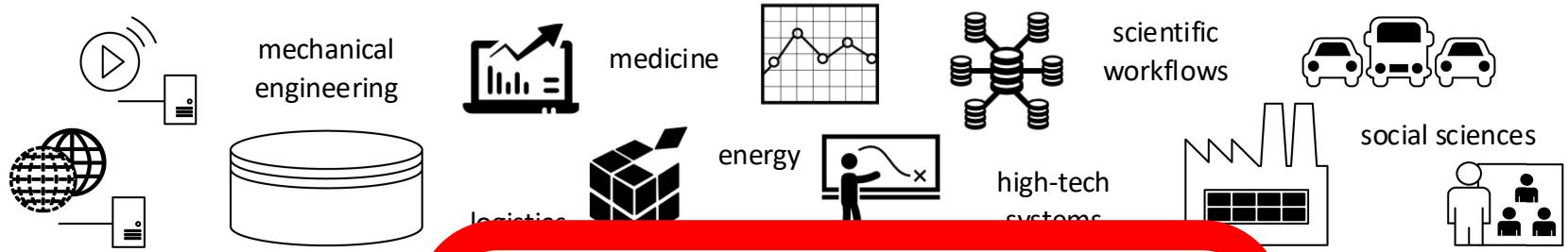
Lecture 22

IDS-L22



Chair of Process
and Data Science

RWTH AACHEN
UNIVERSITY



infrastructure

“volume and velocity”

analysis

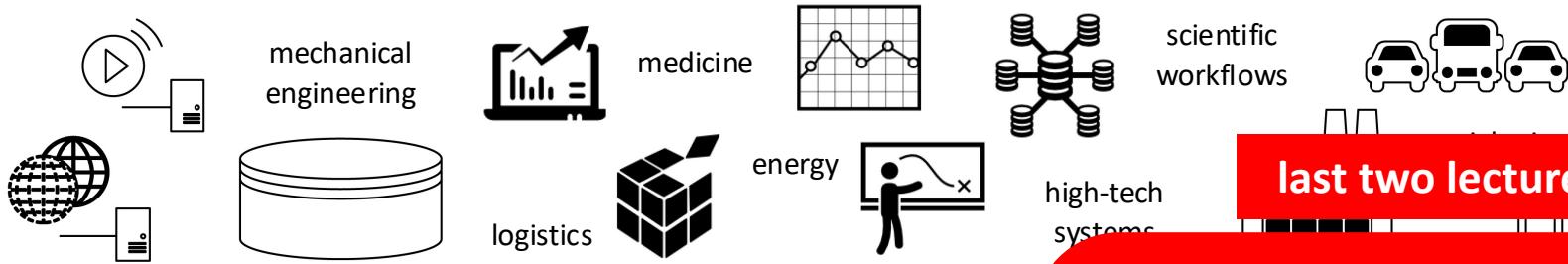
“extracting knowledge”

effect

“people, organizations, society”

- big data infrastructures
- distributed systems
- data engineering
- programming
- security
- ...
- statistics
- data/process mining
- machine learning
- artificial intelligence
- visualization
- ...
- ethics & privacy
- IT law
- operations management
- business models
- entrepreneurship
- ...





last two lectures

infrastructure

“volume and velocity”

analysis

“extracting knowledge”

effect

“people, organizations, society”

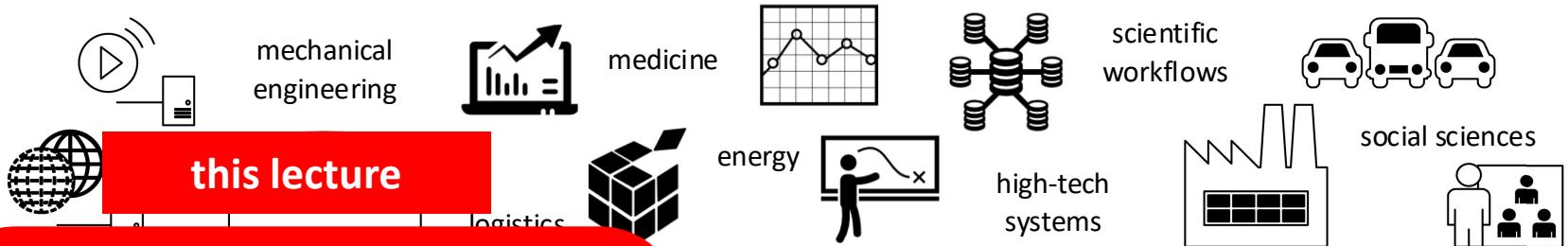
- big data infrastructures
- distributed systems
- data engineering
- programming
- security
- ...

- statistics
- data/process mining
- machine learning
- artificial intelligence
- visualization
- ...

- ethics & privacy
- IT law
- operations management
- business models
- entrepreneurship
- ...



Chair of Process
and Data Science



infrastructure

“volume and velocity”

analysis

“extracting knowledge”

effect

“people, organizations, society”

- big data infrastructures
- distributed systems
- data engineering
- programming
- security
- ...
- statistics
- data/process mining
- machine learning
- artificial intelligence
- visualization
- ...
- ethics & privacy
- IT law
- operations management
- business models
- entrepreneurship
- ...



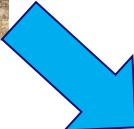
Outline of Today's Lecture

- **Big Data: Relevance and Applications**
- **Characteristics of Big Data: The 4 V's**
- **Big Data Infrastructures**
- **MapReduce**
- **Streaming**
- **Streaming k-means**

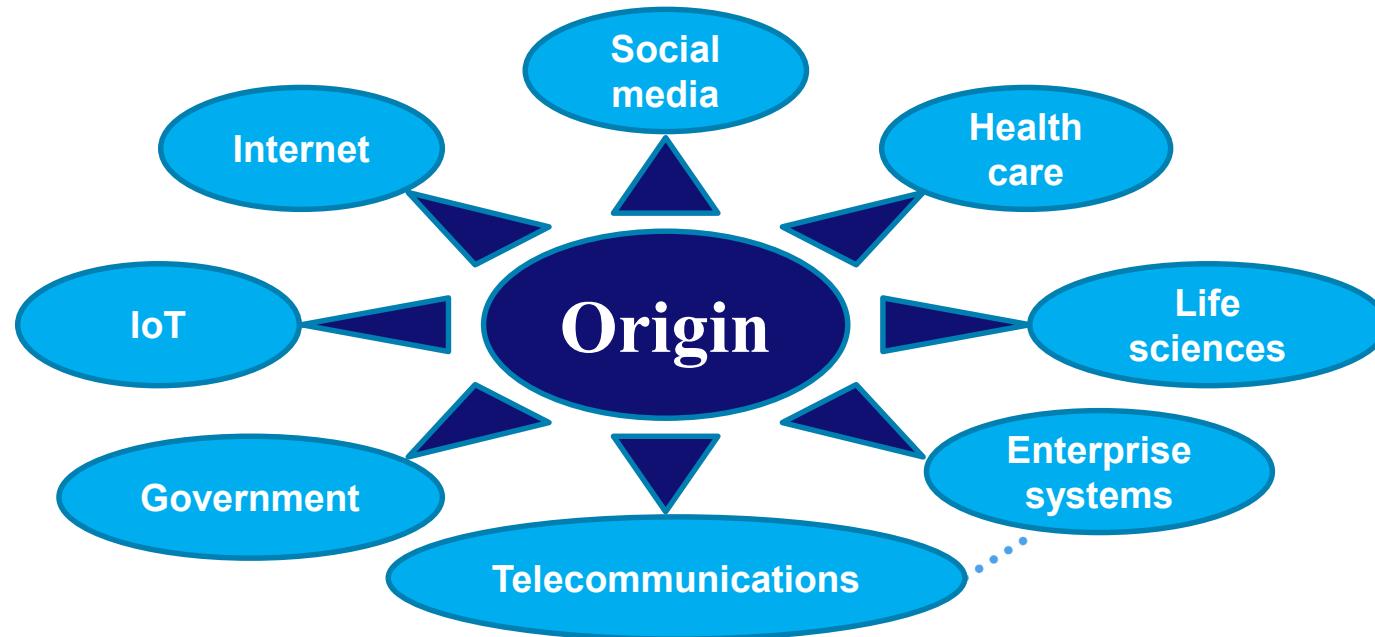
Big Data: Relevance and Applications



From Statistics to Big Data



Where does Big data come from?



Data Explosion

JAN
2018

DIGITAL AROUND THE WORLD IN 2018

KEY STATISTICAL INDICATORS FOR THE WORLD'S INTERNET, MOBILE, AND SOCIAL MEDIA USERS

TOTAL POPULATION



7.593
BILLION

URBANISATION:
55%

INTERNET USERS



4.021
BILLION

PENETRATION:
53%

ACTIVE SOCIAL MEDIA USERS



3.196
BILLION

PENETRATION:
42%

UNIQUE MOBILE USERS



5.135
BILLION

PENETRATION:
68%

ACTIVE MOBILE SOCIAL USERS



2.958
BILLION

PENETRATION:
39%

7

SOURCES: POPULATION: UNITED NATIONS; U.S. CENSUS BUREAU; INTERNET: INTERNET WORLD STATS; ITU: DIGITAL; INTERNET USERS: CIA WORLD FACTBOOK; MIDDLE EAST & MEA: ONG;

FACEBOOK: GOVERNMENT OFFICIALS; REGULATORY AUTHORITIES; REPUTABLE MEDIA; SOCIAL MEDIA AND MOBILE SOCIAL: WEB: FACEBOOK; TENCENT; VONTAKTE; KAKAO; NAVER; DING;

TECHRASA; SIMILARWEB; KEPPOS ANALYSIS; MOBILE: GSMA INTELLIGENCE; GOOGLE; ERICSSON; KEPPOS ANALYSIS. NOTE: PENETRATION FIGURES ARE FOR TOTAL POPULATION (ALL AGES).



Hootsuite™



- ◆ Around 8.4 billion connected things are in use worldwide in 2017.
- ◆ The number will reach 20.4 billion by 2020.
- ◆ The consumer segment is the largest user of connected things with 5.2 billion units in 2017.
- ◆ Regionally, Greater China, North America and Western Europe together represent 67 percent of overall internet of things in 2017.

Applications of Big Data

- **What benefits can the massive volume of available data bring?**

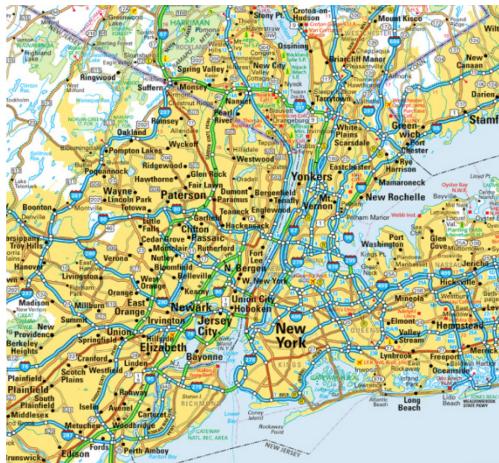


Applications of Big Data

- Employ data to predict crime hotspots



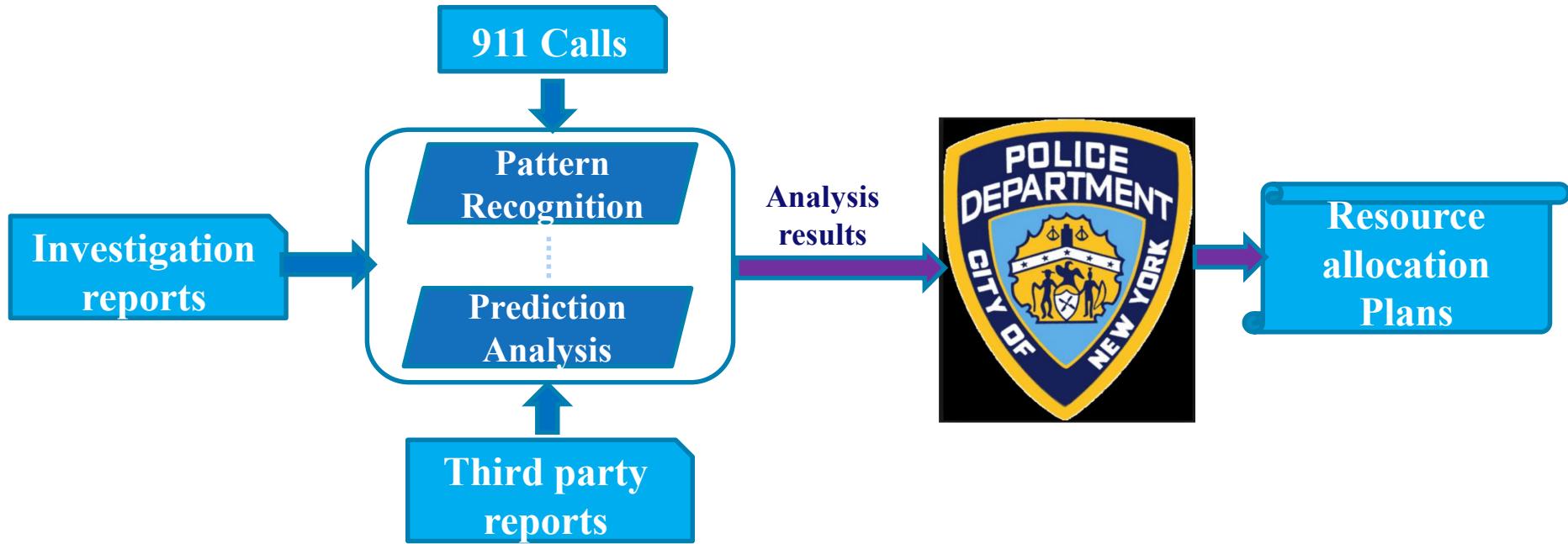
How to allocate
police resources
???????



1. Budgets are limited.
2. Hard to judge where crimes happen.

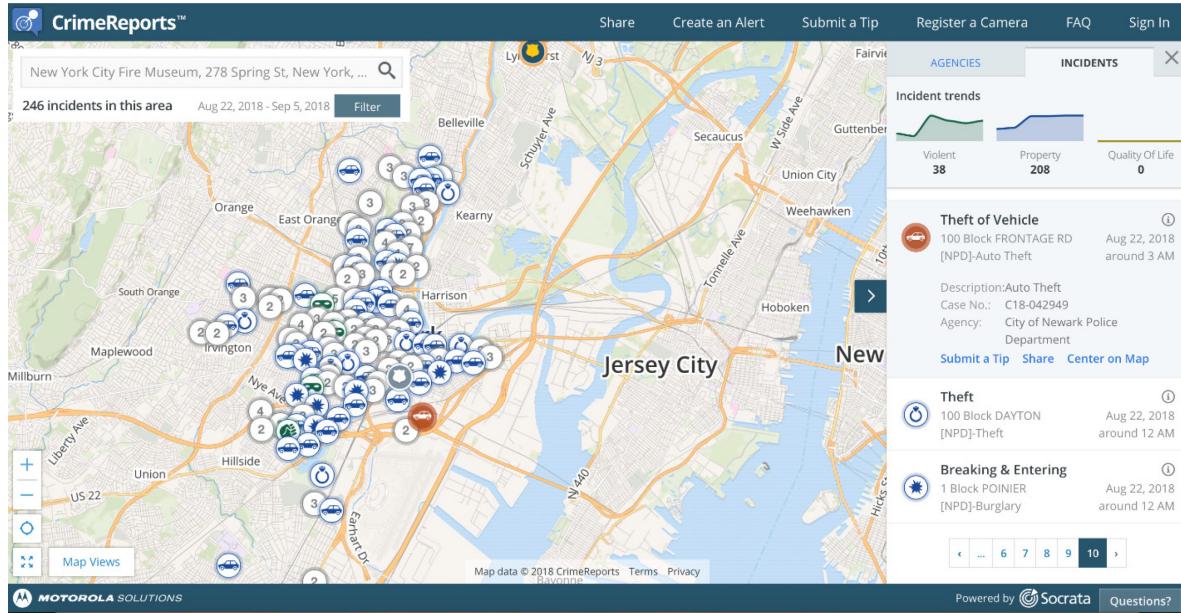
Applications of Big Data

- Employ data to predict crime hotspots



Applications of Big Data

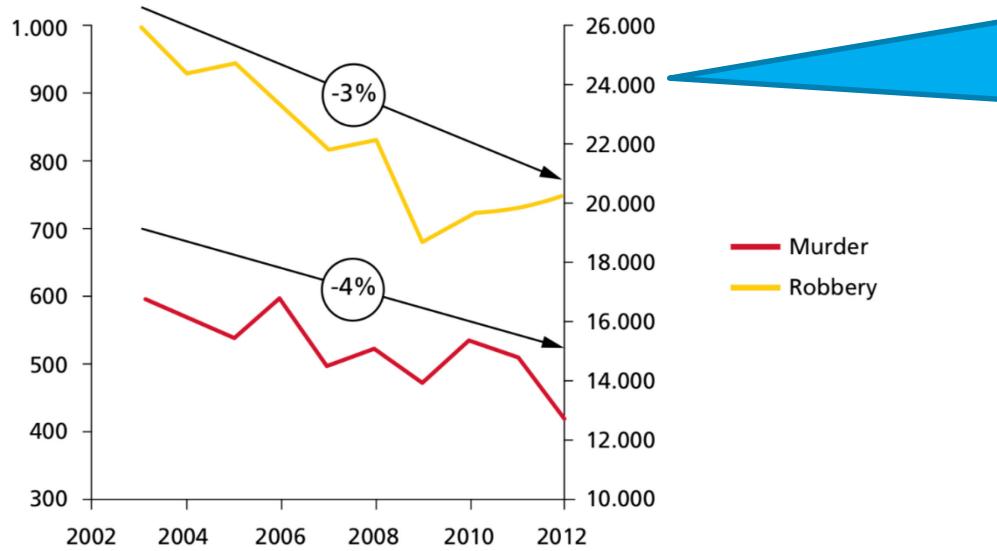
- Employ data to predict crime hotspots



- Information integration.
- Descriptive statistics.
- Data mining.
- Visualization.
- Visual analytics.

Applications of Big Data

- Employ data to predict crime hotspots



Using the correlated data sets, the NYPD is able to effectively target the deployment of manpower and resources

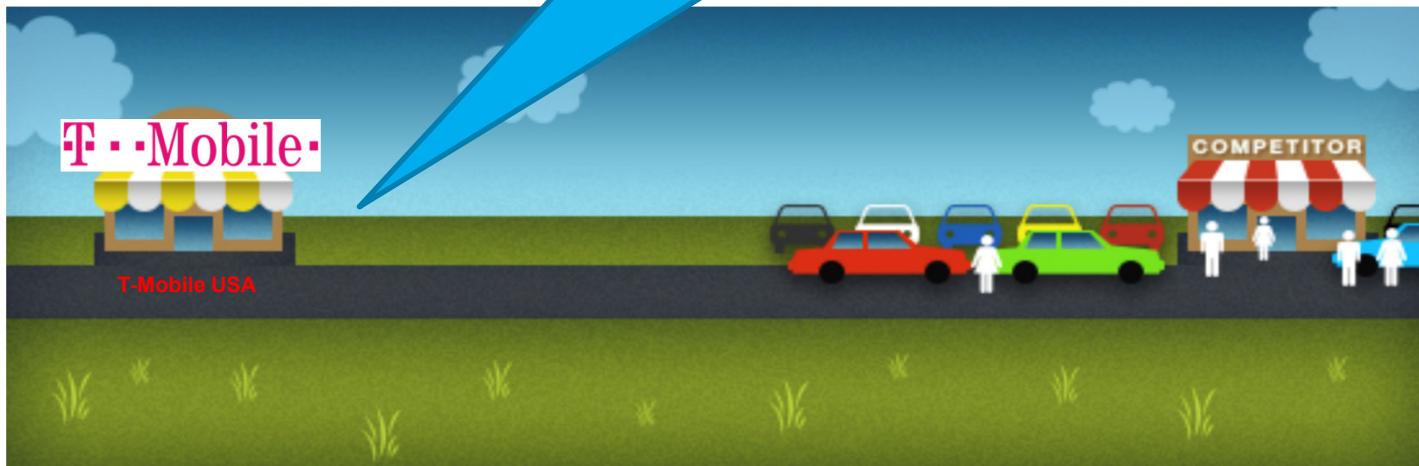
"By introducing statistical analysis and geographical mapping of crime spots, the NYPD has been able to create a "bigger picture" to guide resource deployment and patrol practice. Now the department can recognize crime patterns using computational analysis, and this delivers insights enabling each commanding officer to proactively identify hot spots of criminal activity."

Applications of Big Data

- Social network analysis for retaining customers

Facing the problem of
customer churn

(the loss of customers over a period of time)



Applications of Big Data



Communication network elements



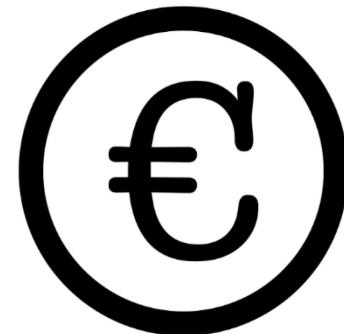
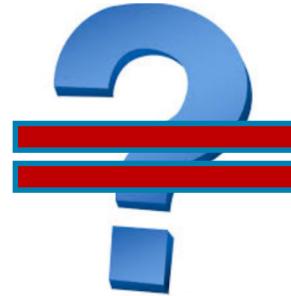
Billing system



Applications of Big Data

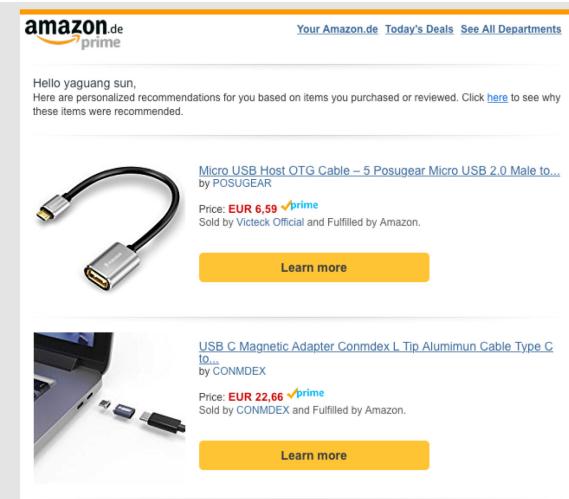
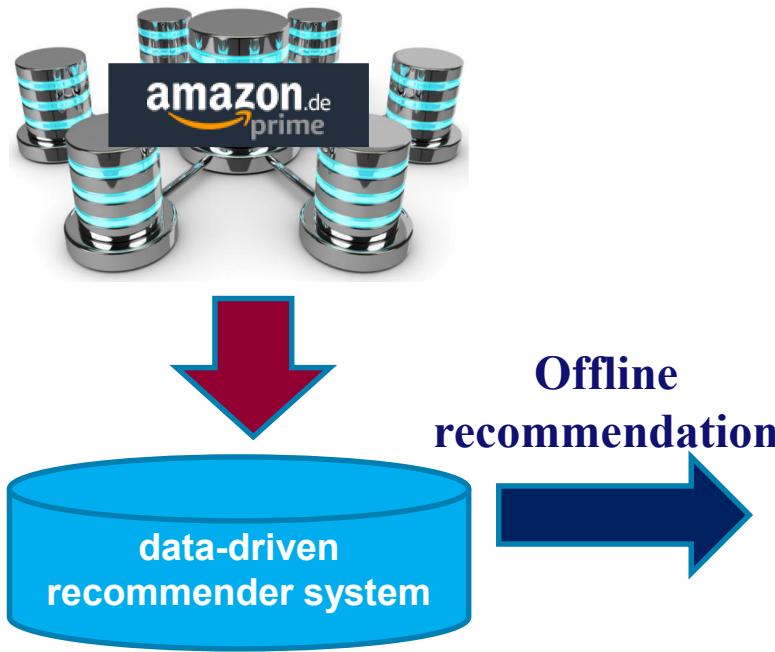
- **Data-driven marketing**

Amazon's database contains billions
of data points from tens of millions
of customers and products !



Applications of Big Data

- Data-driven marketing



Applications of Big Data



Recommendations for you in Kindle-Shop



Inspired by your shopping trends

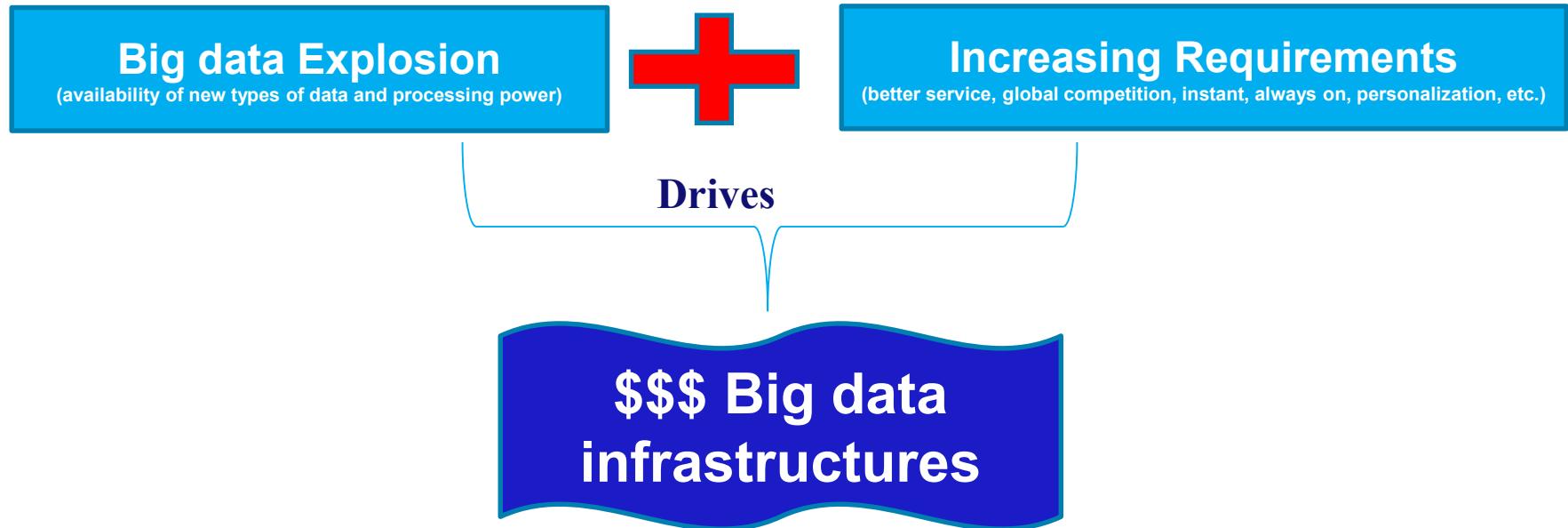


Applications of Big Data

35% of purchases from Amazon.com's customers are brought by its recommendation engine.

75% of what people watch on Netflix comes from recommendations.

What drives the development of Big Data infrastructures?

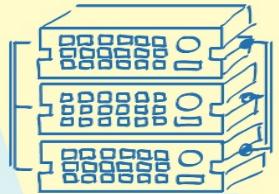


Characteristics of Big Data: The 4 V's



VOLUME

SCALE OF DATA



BIG DATA

ANALYSIS OF
DATA-FLOW

VELOCITY

VARIETY

FORMS OF DATA



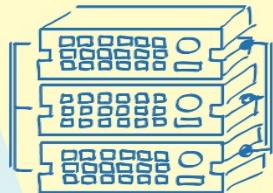
UNCERTAINTY
OF DATA



VERACITY

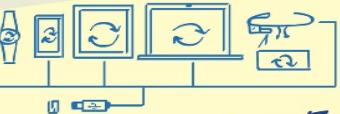
VOLUME

SCALE OF DATA



B

I G D A T A



ANALYSIS OF
DATA-FLOW

VELOCITY

VARIETY

FORMS OF DATA



UNCERTAINTY
OF DATA

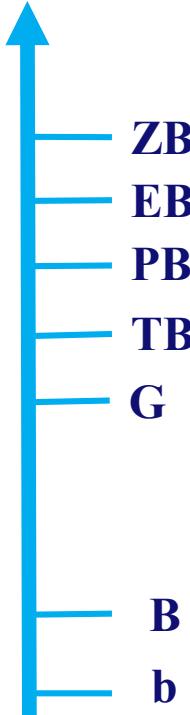
VERACITY

Characteristics of big data

- **Volume = Data Size**

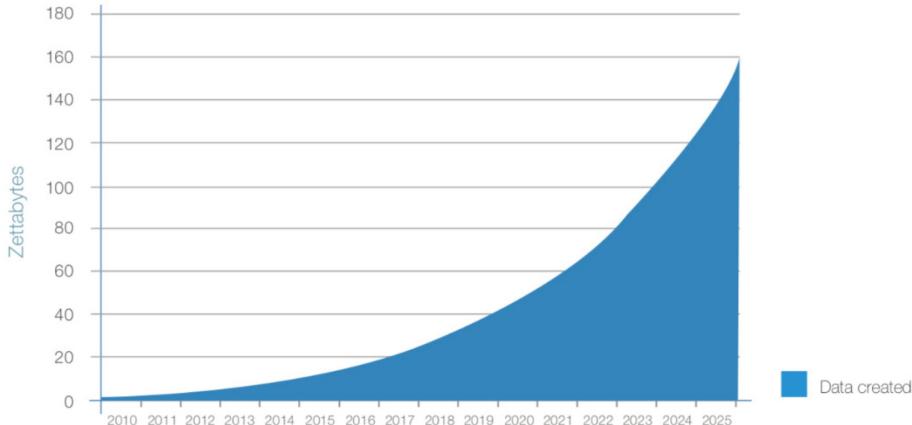
1	b	bit
$8 \times b$	B	Byte
$1024 \times B$	KB	Kilobyte
$1024^2 \times B$	MB	Megabyte
$1024^3 \times B$	GB	Gigabyte
$1024^4 \times B$	TB	Terabyte
$1024^5 \times B$	PB	Petabyte
$1024^6 \times B$	EB	Exabyte
$1024^7 \times B$	ZB	Zettabyte

Characteristics of big data

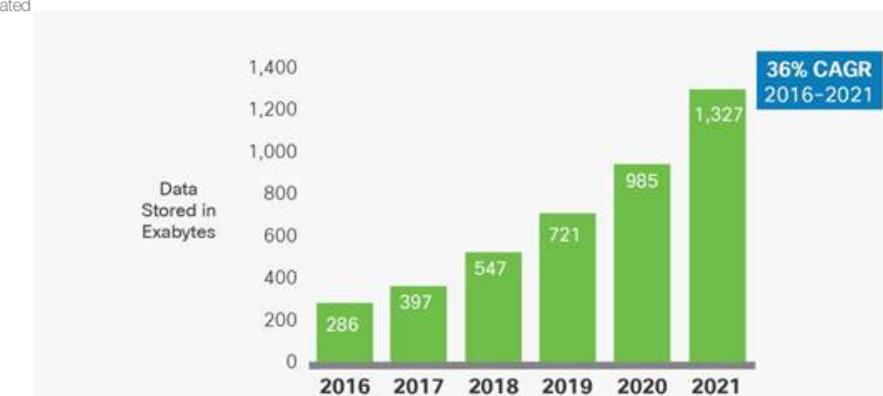


- 180 ZB ≈ estimated amount of digital data in 2025**
- 0.5 ZB ≈ entire World Wide Web in 2009**
- 5 EB ≈ all words ever spoken by human beings**
- 1 PB ≈ 2.7 billion facebook images**
- 1 TB ≈ 300 hours of high quality video**
- 1 GB ≈ 7 minutes of HDTV video at 19.39 Mbit/s**
- 1 Byte = one character**
- 1 bit = 0/1**

All point in the same direction



Source: IDC's Data Age 2025 study, sponsored by Seagate, April 2017



Source: Cisco Global Cloud Index, 2016-2021.

Examples

V1

Snapchat users share 527,760 photos every day

456,000 tweets are sent on Twitter every day

We send 16 million text messages every minute

Spotify adds 13 new songs every minute

Uber riders take 45,788 trips every minute

990,000 Tinder swipes every minute

600 pages are edited in Wikipedia every minute

154,200 calls on Skype every minute

More than 300 million photos are uploaded every day in Facebook

Every minute 510,000 comments posted and 293,000 statuses are updated in Facebook



Chair of Process
and Data Science

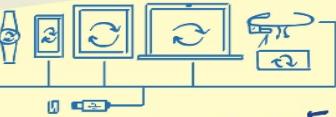
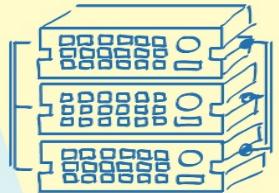
Key question

v1

**How to collect, store, and process all of
these data in an economical and reliable
manner, instantly?**

VOLUME

SCALE OF DATA



ANALYSIS OF
DATA-FLOW

VELOCITY

BIG DATA

VARIETY

FORMS OF DATA

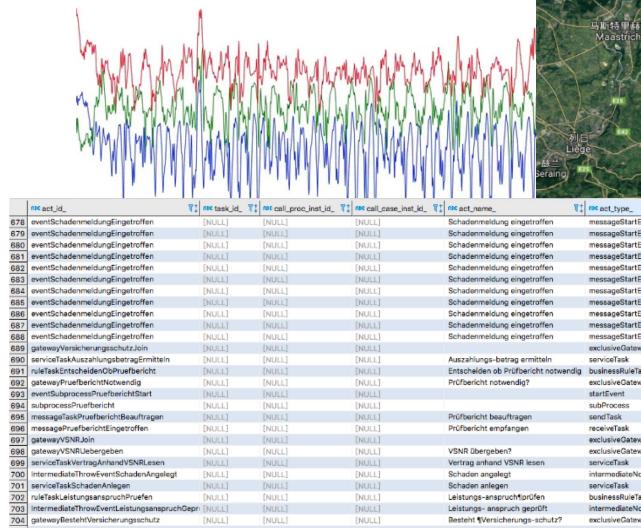


UNCERTAINTY
OF DATA

VERACITY

Characteristics of big data

- Data are more heterogeneous than before



be measured. One advantage of this technique is that it provides a full range of metrics for clustering traces. In [193] the authors point out that a full model (with high fitness) discovered may support a high variety of metrics that are not registered in event log, thus a real some significant structure may be concealed in the mined model. Such a problem can be dealt with considering the metric soundness [9] which measures the percentage of the mined model that are recorded in the log among all of the behaviour proposed by the model. An efficient technique is proposed in [193] where the whole process into a set of distinct sub-processes based on a tree which makes sure the further division of a process will lead to an equally sound sub-process. This method can also help solve the problem of complexity of the initial model. Context-aware trace clustering technique is proposed in [195] and [35]. In [195] the authors indicate that the basic idea is to calculate the edit distance between traces. The cost of edit operations is associated with the contexts of activitie. The calculated edit distance between traces is more accurate. The sequence technique based on first-order Markov chains is presented in [196] which learns a potential first-order Markov model for each cluster expectation-maximization algorithm. A sequence is assigned to a cluster if it is able to generate it with higher probability. The technique proposed paper also inherits the idea from sequence clustering, the difference is that it represents each cluster with a set of separate sequences (sign).

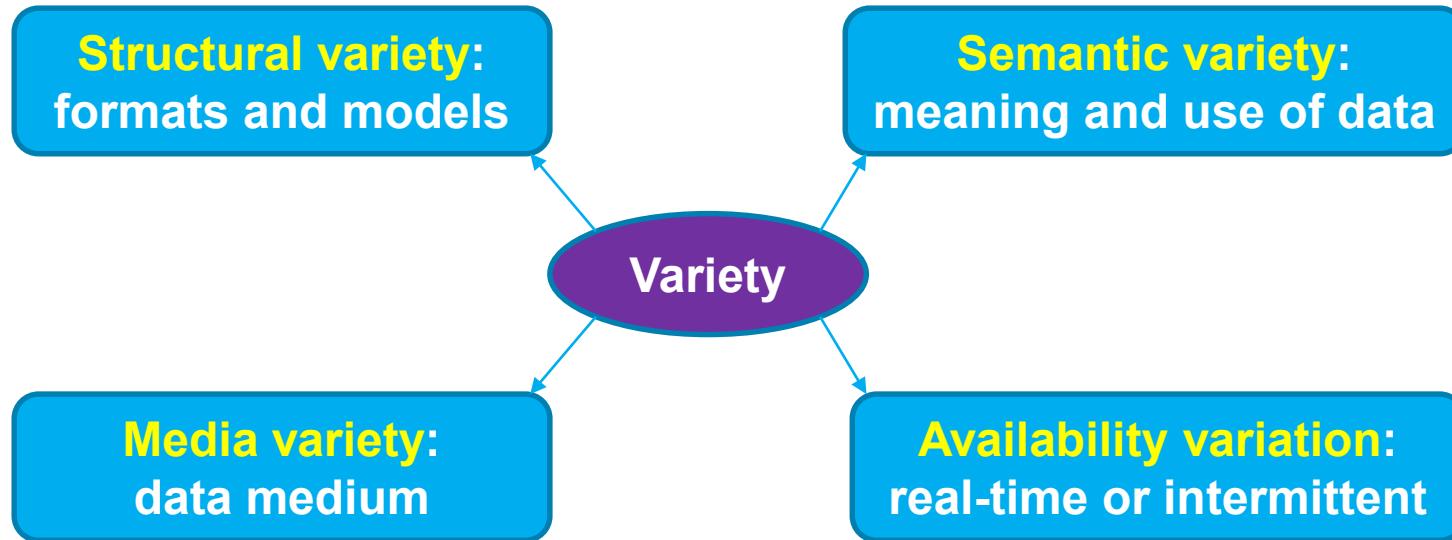
axigo High Performance PMS - Trailer Reaktion



The complexity of data storage and analysis is greatly increased with a variety of structured and unstructured data!!

Characteristics of big data

- Four types of variety



Characteristics of big data

V2

- Variety in social network

The screenshot shows a LinkedIn Learning course page. At the top, there's a navigation bar with icons for Home, People, Positions, and Messages. Below the bar, a course card for 'Learn how to overcome the most common hurdles when it comes to job seeking and career success with this course – free on LinkedIn Learning' is displayed. The card includes a play button and a photo of a woman with glasses. Below the card, there are comments from users like Gzim Siriniqi and Deon Botha.

LinkedIn 领英 搜索

Yaguang Sun
PADS in RWTH Aachen – Research Assistant

79 好友

添加好友

18 谁看过我

获取独家工具与洞察 试用高级帐号

LinkedIn (领英) 4,158,770 位关注者 广告

Learn how to overcome the most common hurdles when it comes to job seeking and career success with this course – free on LinkedIn Learning <https://lnkd.in/eRpYF2y>

Betty Liu Anchor & Founder of Radiate

赞 (14,542) · 评论 (216)

热门评论

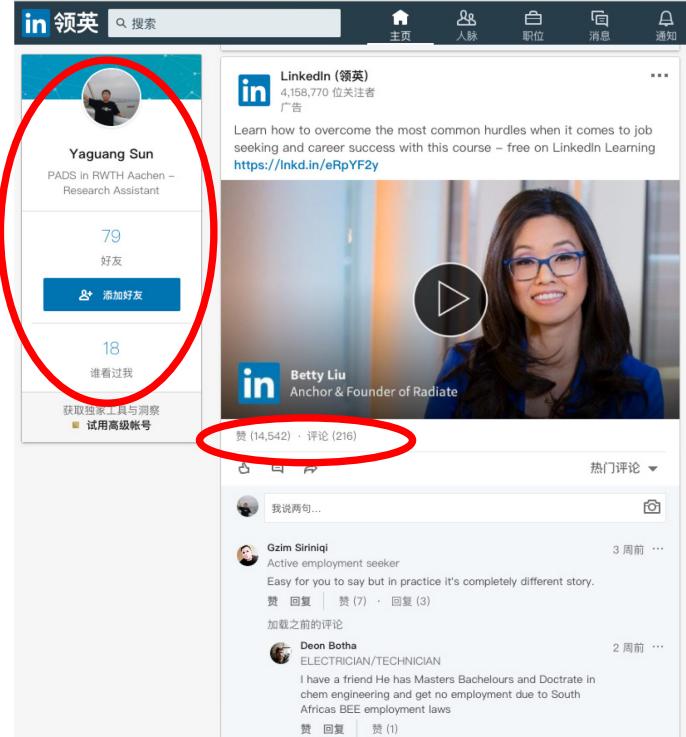
Gzim Siriniqi Active employment seeker 3周前 Easy for you to say but in practice it's completely different story.
赞 回复 | 赞 (7) · 回复 (3)
加载之前的评论

Deon Botha ELECTRICIAN/TECHNICIAN 2周前 I have a friend He has Masters Bachelours and Doctorate in chem engineering and get no employment due to South Africas BEE employment laws
赞 回复 | 赞 (1)

Characteristics of big data

V2

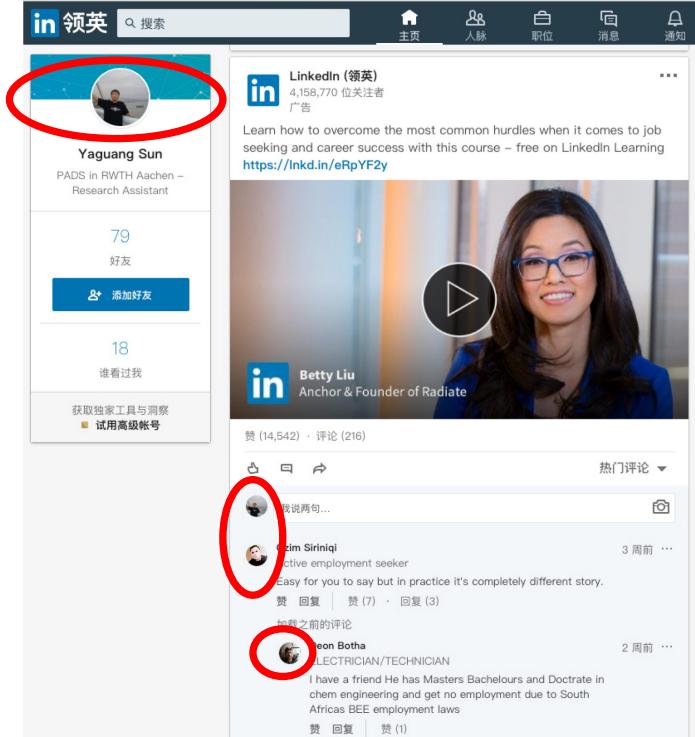
- Variety in social network
 - Structured data



Characteristics of big data

V2

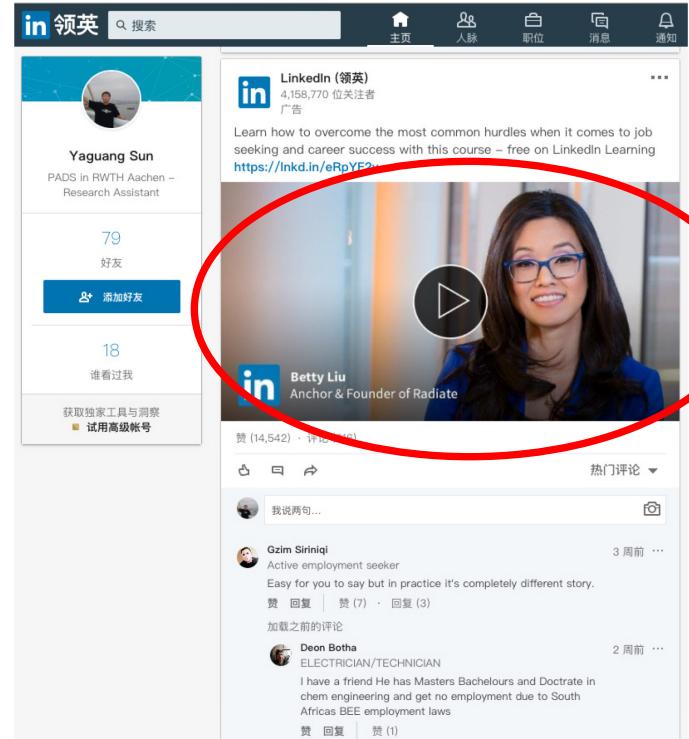
- Variety in social network
 - Structured data
 - Pictures



Characteristics of big data

V2

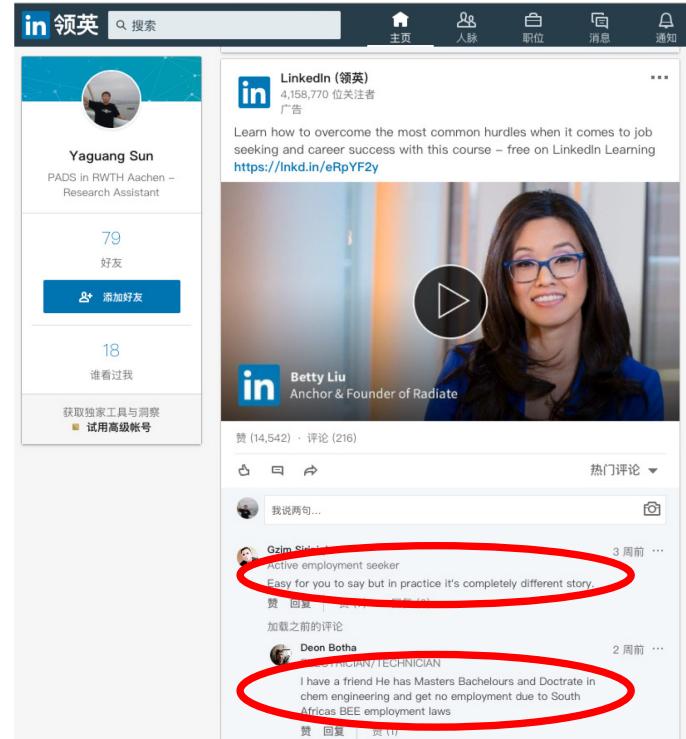
- Variety in social network
 - Structured data
 - Pictures
 - Video



Characteristics of big data

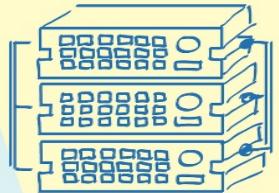
V2

- Variety in social network
 - Structured data
 - Pictures
 - Video
 - Unstructured text

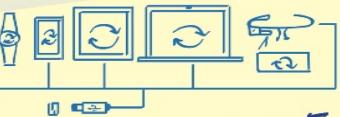


VOLUME

SCALE OF DATA



BIG DATA



ANALYSIS OF
DATA-FLOW

VELOCITY

VARIETY

FORMS OF DATA

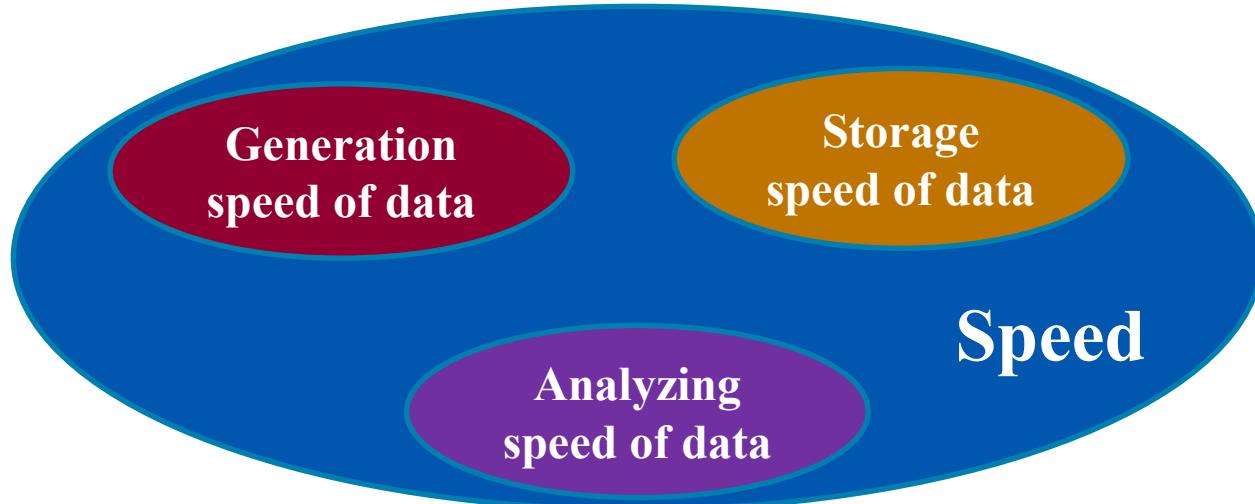


UNCERTAINTY
OF DATA

VERACITY

Characteristics of big data

- **Velocity = Speed**



Characteristics of big data

- Exponential growth of data

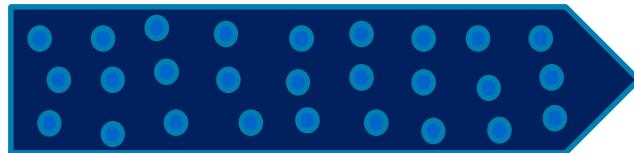


The Patt & Whitey's GTF engine for a Bombardier C Series jetliner contains 5000 sensors, which generate 10 GB of data per second.

Characteristics of big data

- Processing data in real-time to match its generation speed is very important in some applications

Alibaba's singles day shopping festival



Data flow



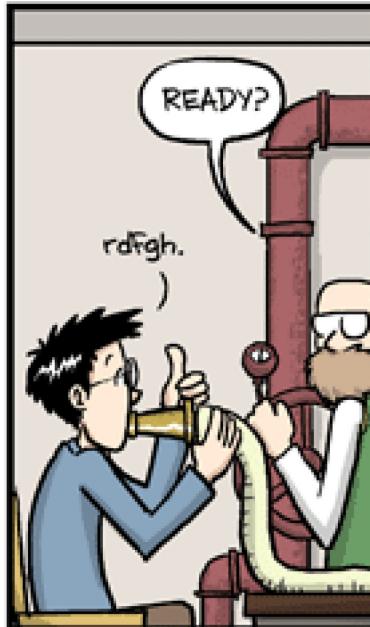
Peak volume of deals: 325k / s

Peak volume of payment: 256k / s

Peak volume in data base operation: 42000k / s

Streaming Data

Need to process data immediately while having strict memory constraints

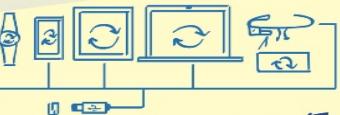
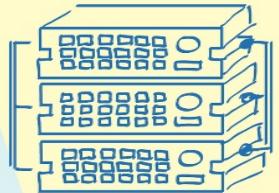


JORGE CHAM © 2007

WWW.PHDCOMICS.COM

VOLUME

SCALE OF DATA



ANALYSIS OF
DATA-FLOW

VELOCITY

VARIETY

FORMS OF DATA



BIG DATA

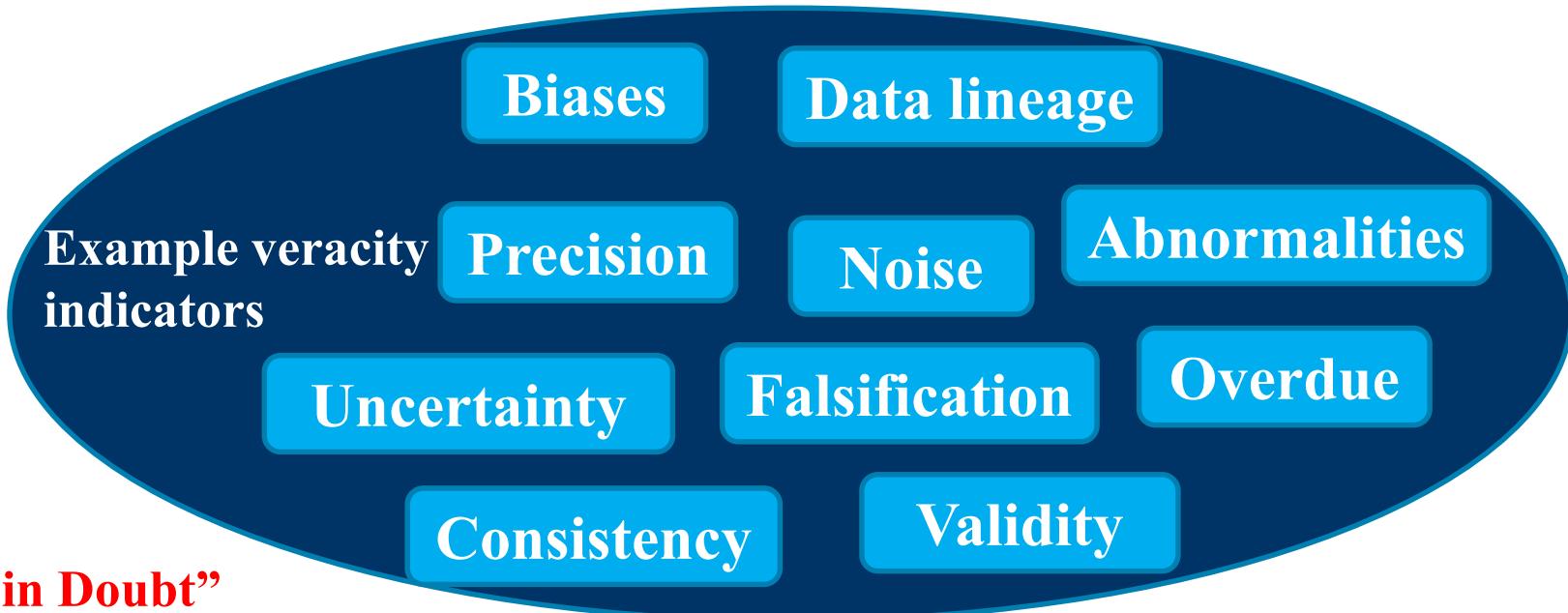
UNCERTAINTY
OF DATA



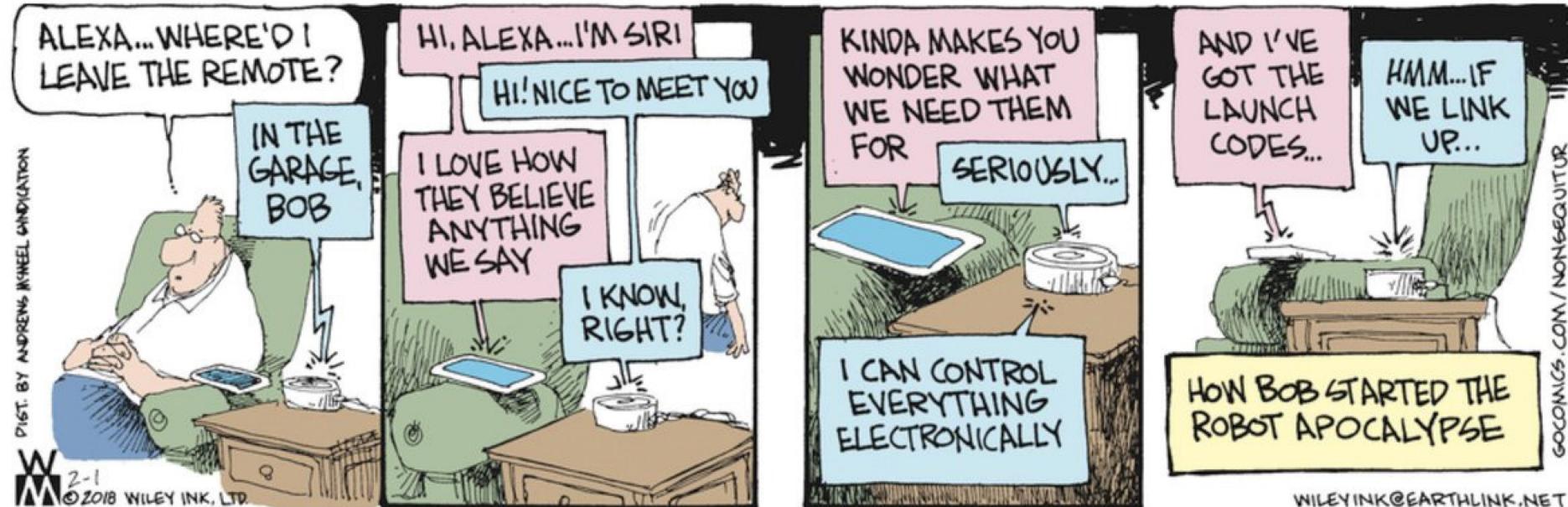
VERACITY

Characteristics of big data

- **Veracity = Uncertainty of Data**



Are events recorded correctly?



Accidental Alexa and Siri activations

Big Data Infrastructures



Big data infrastructures

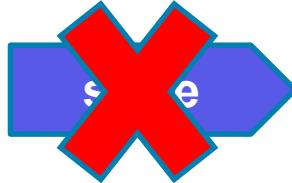
Motivations for the development of big data infrastructures



Big data infrastructures

We need to store and manage massive volumes of data (TB/EB/ZB)

Hundreds of petabytes of data



4 TB SSD storage

Hundreds of petabytes of data

store



Easy access and management of data

Big data infrastructures

We need to parallelize and distribute computation across hundreds or thousands of CPUs

Hundreds of petabytes of data



Hundreds of petabytes of data



1000 x



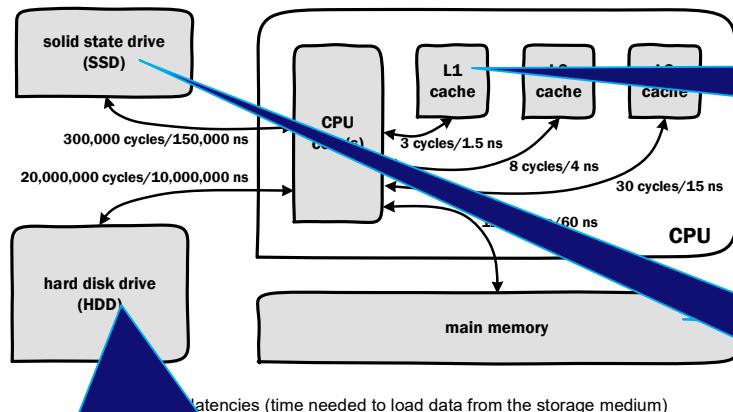
Chair of Process
and Data Science

Google Data Centers

cheap
commodity
hardware &
reliability by
redundancy



“Distances” in computing



getting a
Nespresso from
the kitchen

getting a coffee from the
Starbucks around the
corner

flying to Milano for a
coffee

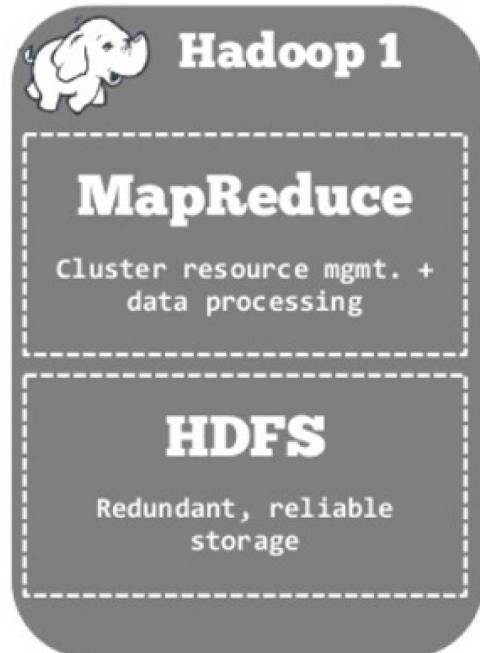
flying to Colombia, Ethiopia
and Kenya, process the beans
in Amsterdam, and then fly to
Rome

Three Important Trends

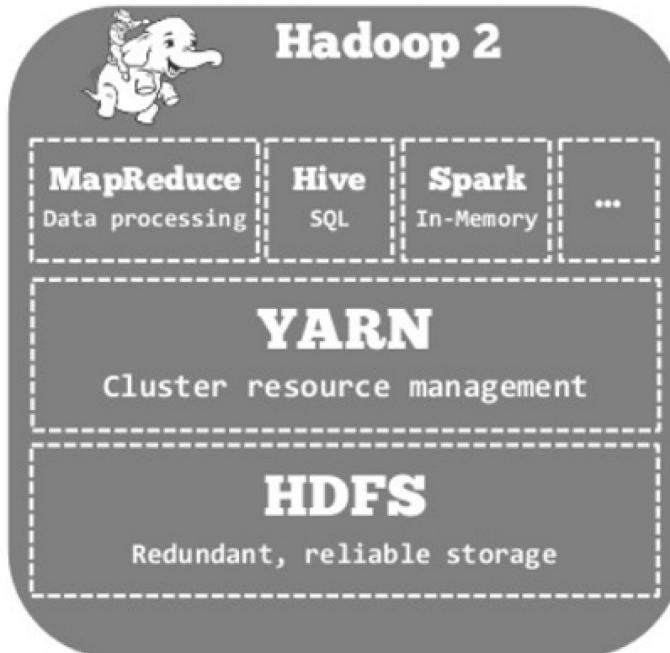
- More data in memory
- Distribution
 - Commodity hardware made reliable
 - Programming models like MapReduce
- Streaming

Hadoop History

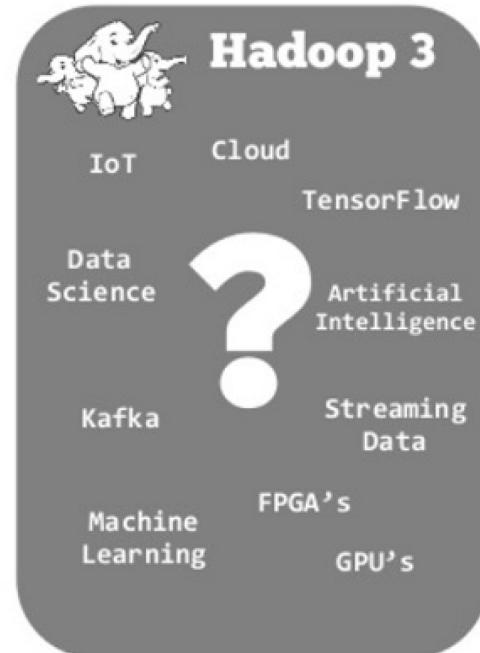
Era of Silicon Valley Hadoop



Era of Enterprise Hadoop



Era of ?



Apache Hadoop Ecosystem

(main goal is to hide the complexity of distribution from the user and be fast and reliable)

Management & Monitoring
(Ambari)

Scripting
(Pig)

Machine Learning
(Mahout)

Query
(Hive)

Coordination
(ZooKeeper)

Workflow & Scheduling
(Oozie)

Distributed Processing
(MapReduce)

Distributed Storage
(HDFS)

NoSQL Database
(HBase)

Data Integration
(Sqoop/REST/ODBC)

BIG DATA & AI LANDSCAPE 2018



INFRASTRUCTURE

HADOOP ON-PREMISE



HADOOP IN THE CLOUD



STREAMING / IN-MEMORY



ANALYTICS

DATA ANALYST PLATFORMS



DATA SCIENCE PLATFORMS



NoSQL DATABASES



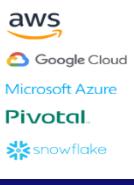
NewSQL DATABASES



GRAPH DBs



MPP DBs



CLOUD EDW



DATA TRANSFORMATION



DATA INTEG.



CLUSTER SVCS



APP DEV



CROWD-SOURCING

HARDWARE



GPU DBs



SEARCH



LOG ANALYTICS



SOCIAL ANALYTICS



WEB / MOBILE / COMMERCE ANALYTICS

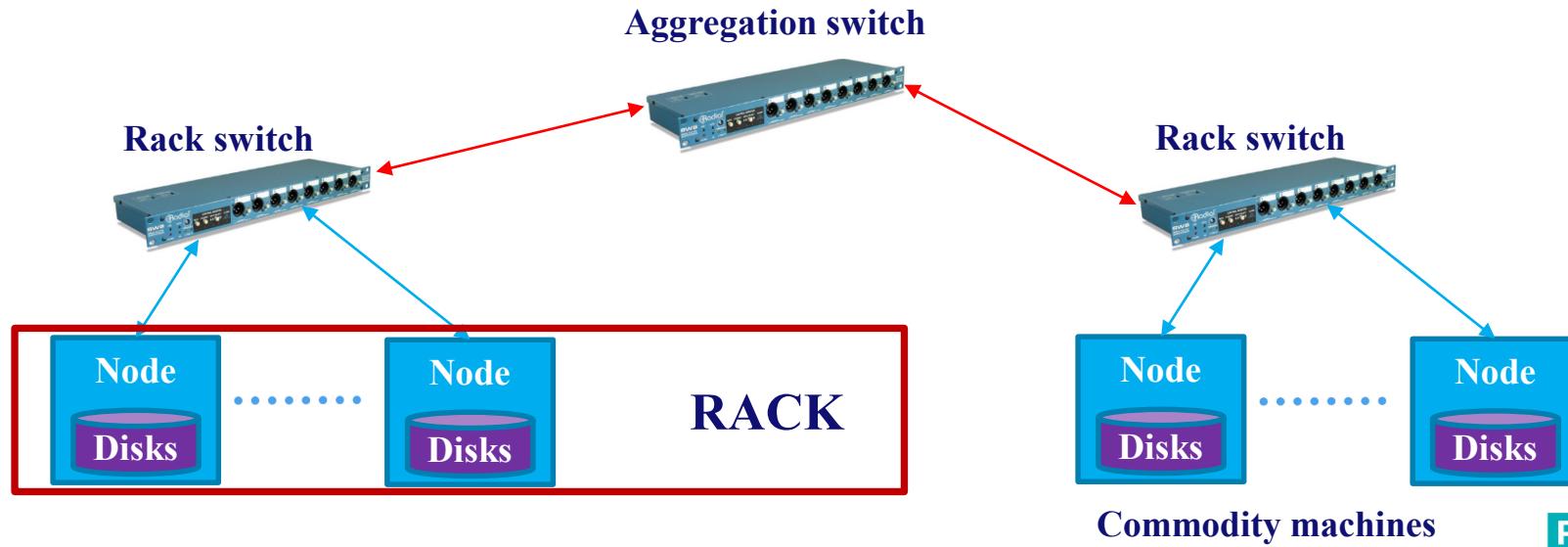


CROSS-INFRASTRUCTURE/ANALYTICS



Distributed File System (DFS)

Distributed File System (DFS): The core of any Big data infrastructure



Distributed File System (DFS)

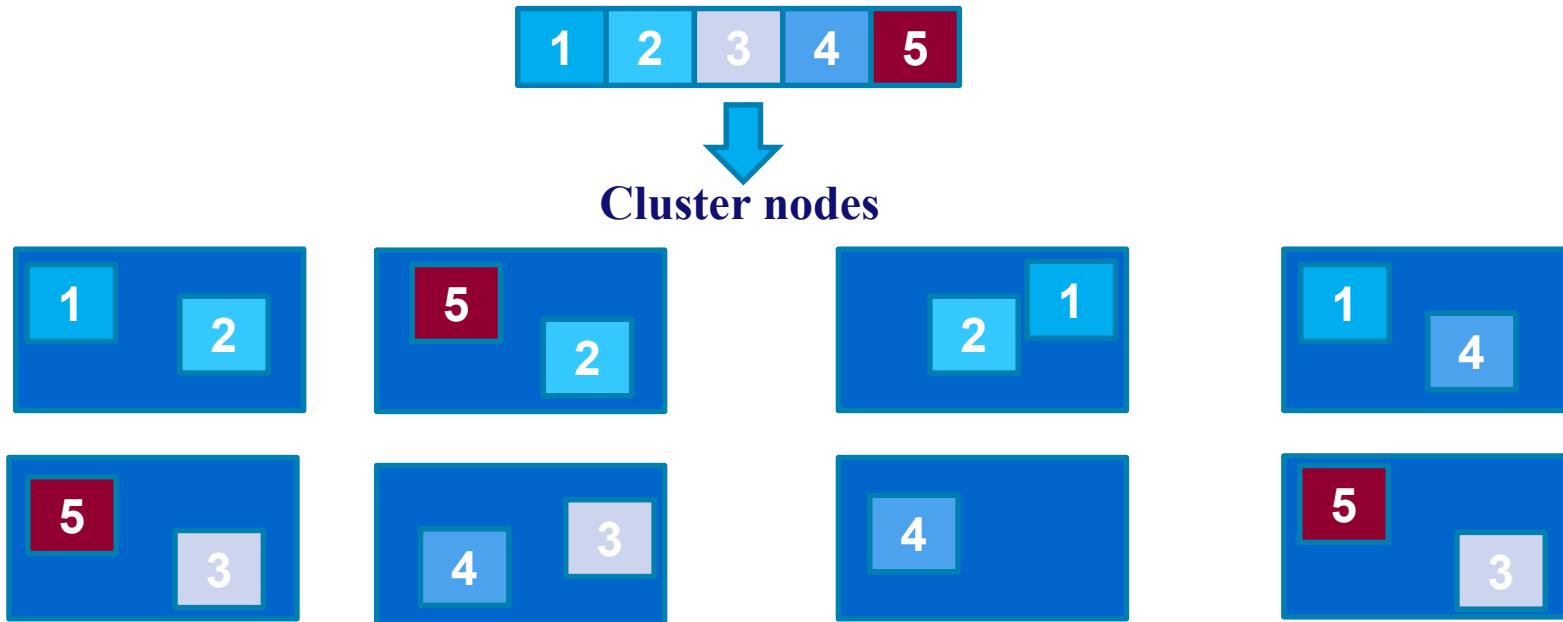
- **Distributed file system for big data**
 - **Each file written into DFS is split into data blocks**

One file is split
into five blocks



Distributed File System (DFS)

- **Distributed file system for big data**
 - Each block is stored on one or more nodes (3 copies / block)



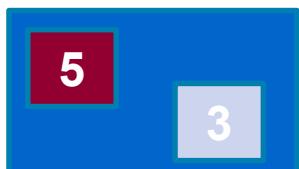
Distributed File System (DFS)

**data distribution
and replication**

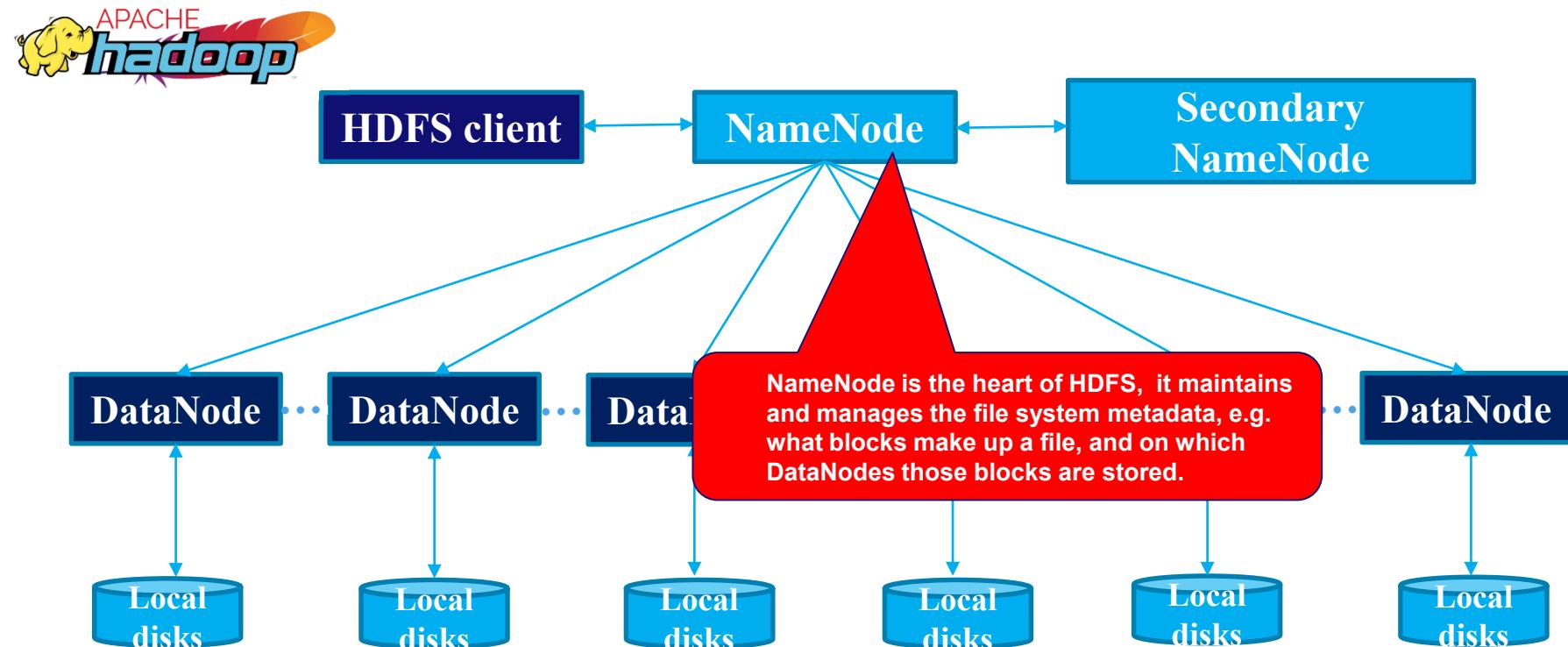


**high concurrency
scalability
fault tolerance**

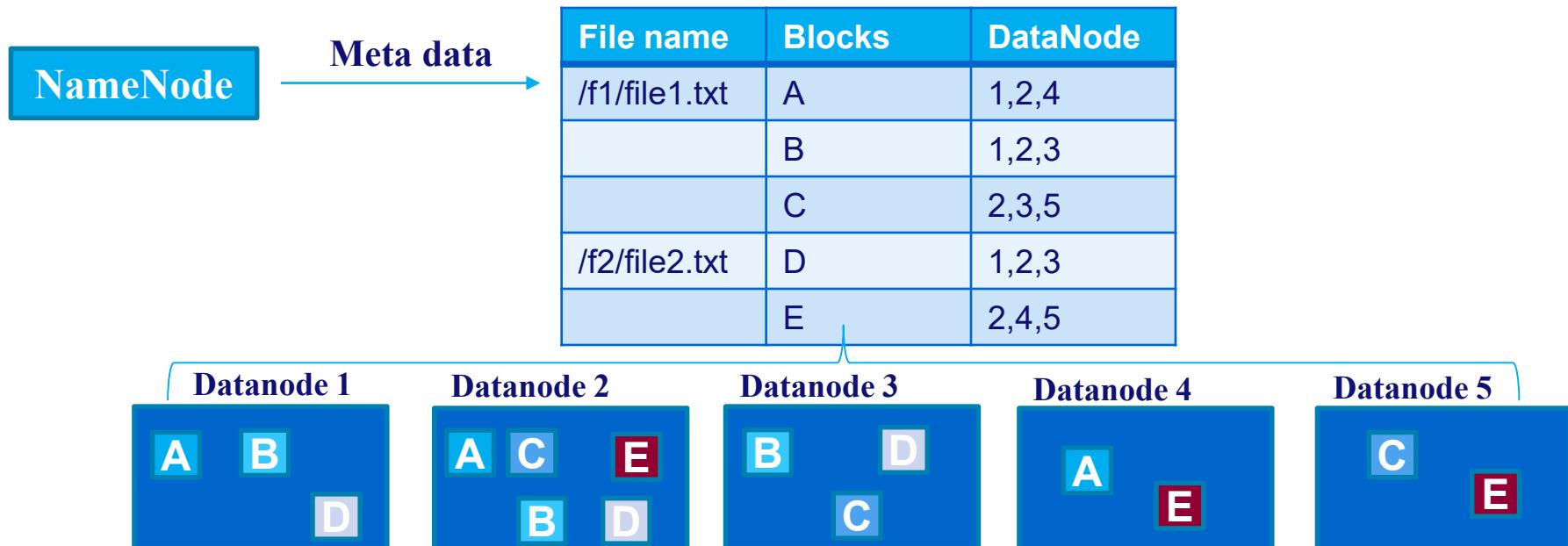
Cluster nodes



Hadoop Distributed File System (HDFS)

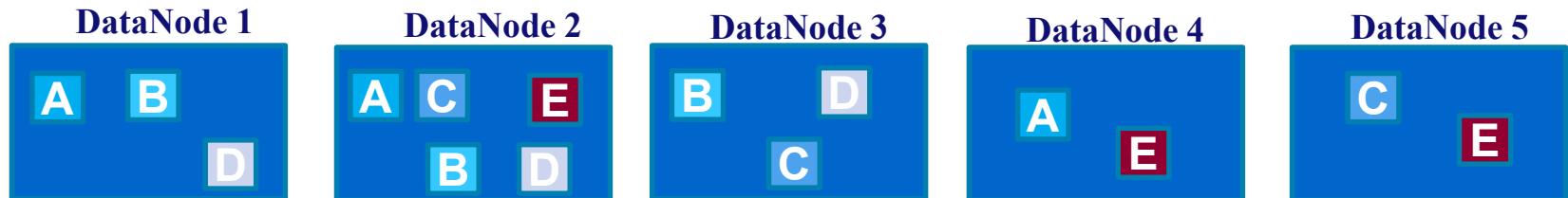


Hadoop Distributed File System (HDFS)



Hadoop Distributed File System (HDFS)

- **DataNode is where HDFS stores the actual data, there are usually many of such nodes.**



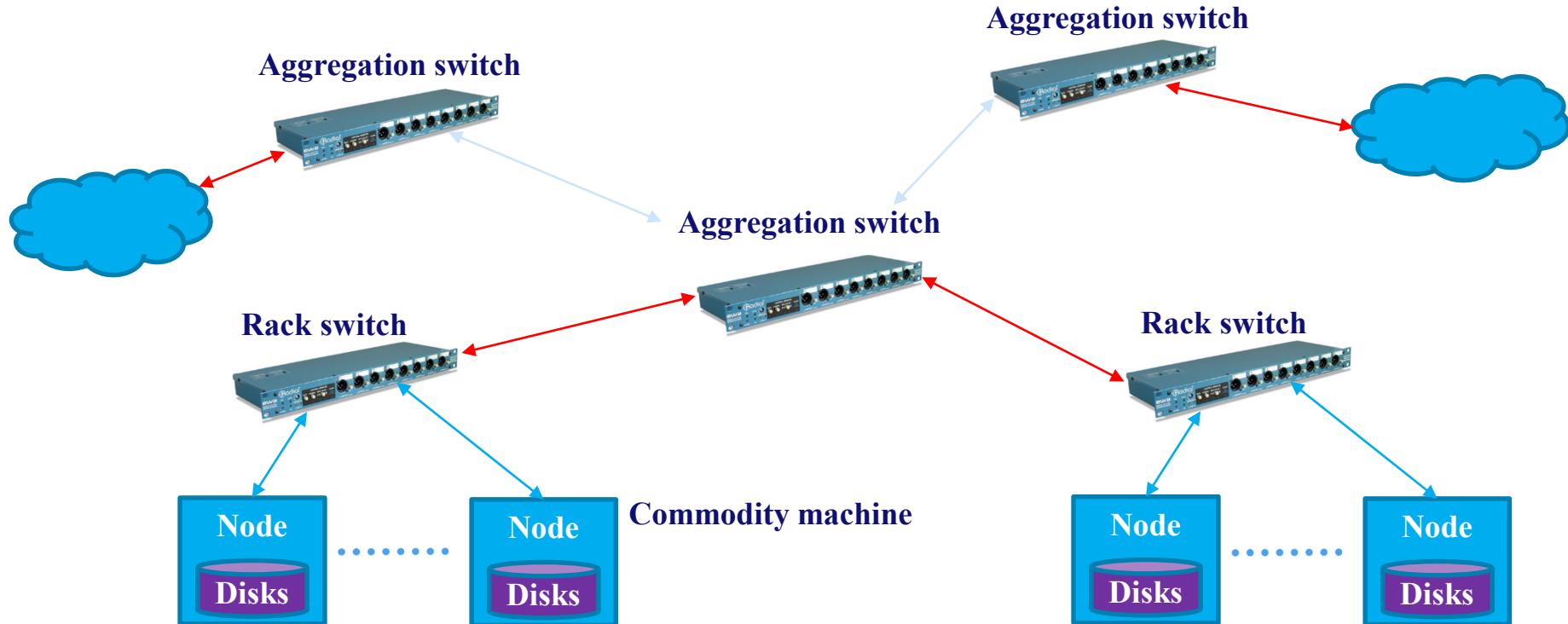
Hadoop Distributed File System (HDFS)

- HDFS is **fault tolerant** because the data are duplicated across multiple nodes to protect against machine failures.
- HDFS is highly **scalable**:
 - Data transfers happen directly with DataNodes so that read/write capacity scales well.
 - More capacity needed: just add more DataNodes and re-balance.

Distributed Computing Using a DFS

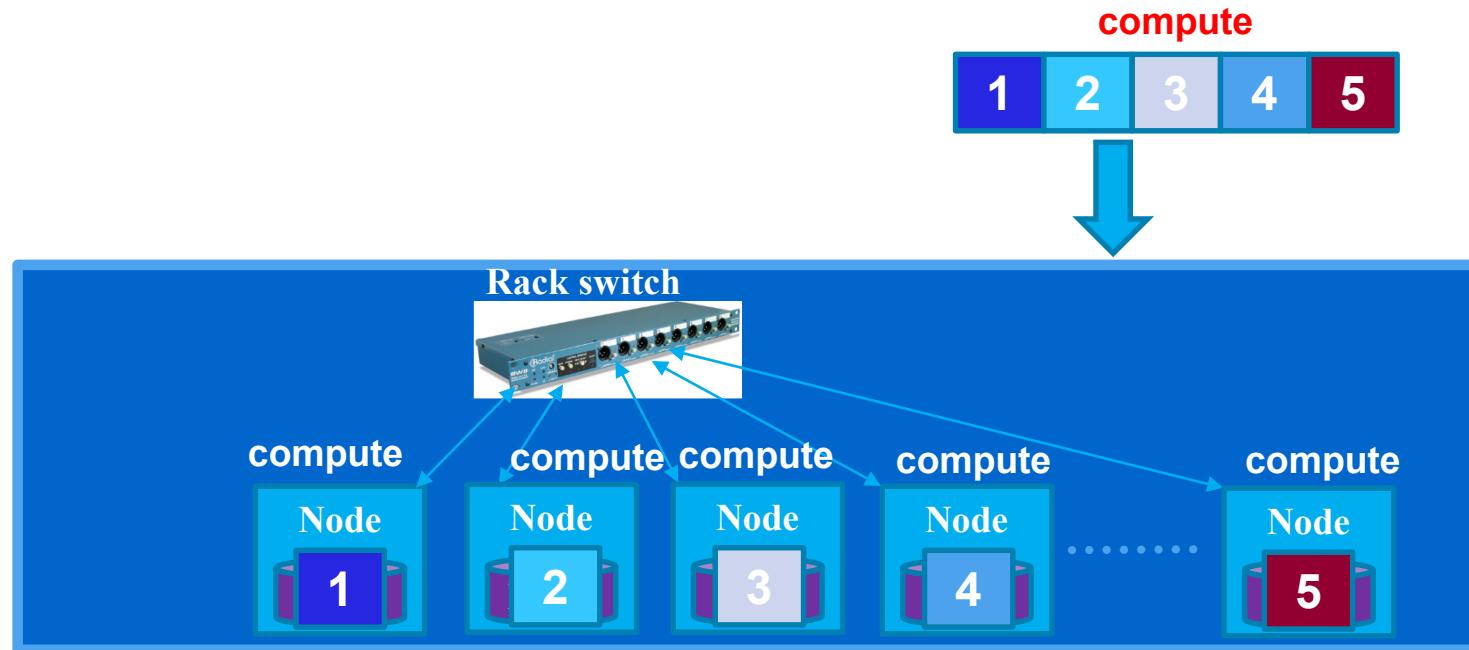
- **Distributed system is composed of components that are located on different networked computers.**
- **These communicate and coordinate their actions by passing messages to achieve a common goal.**

Distributed Computing Using a DFS



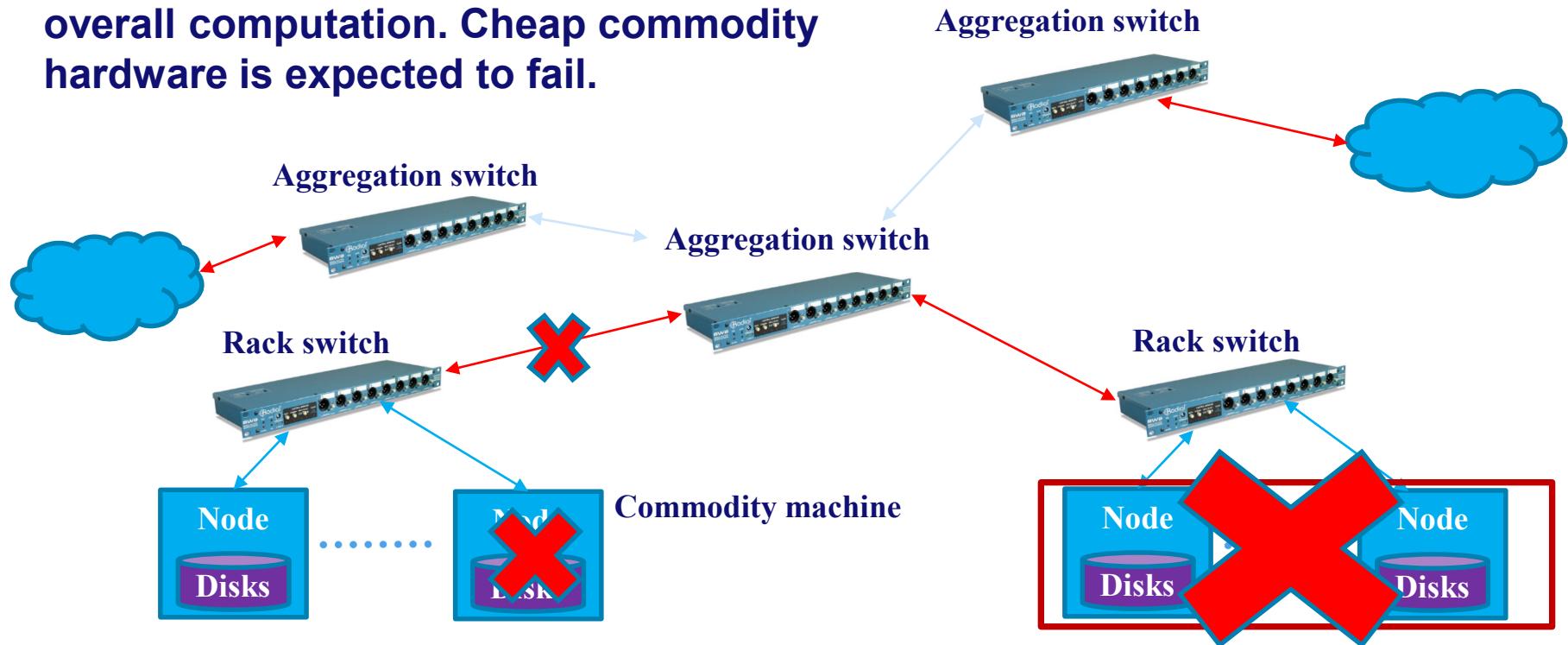
Distributed Computing Using a DFS

Data parallelism: Move computation to data!



Distributed Computing Using a DFS

Fault-tolerant: No need to restart the overall computation. Cheap commodity hardware is expected to fail.

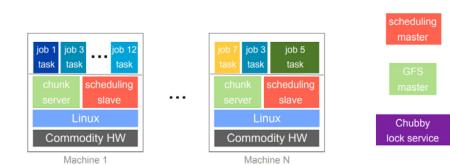


Some numbers from Jeff Dean (Google)

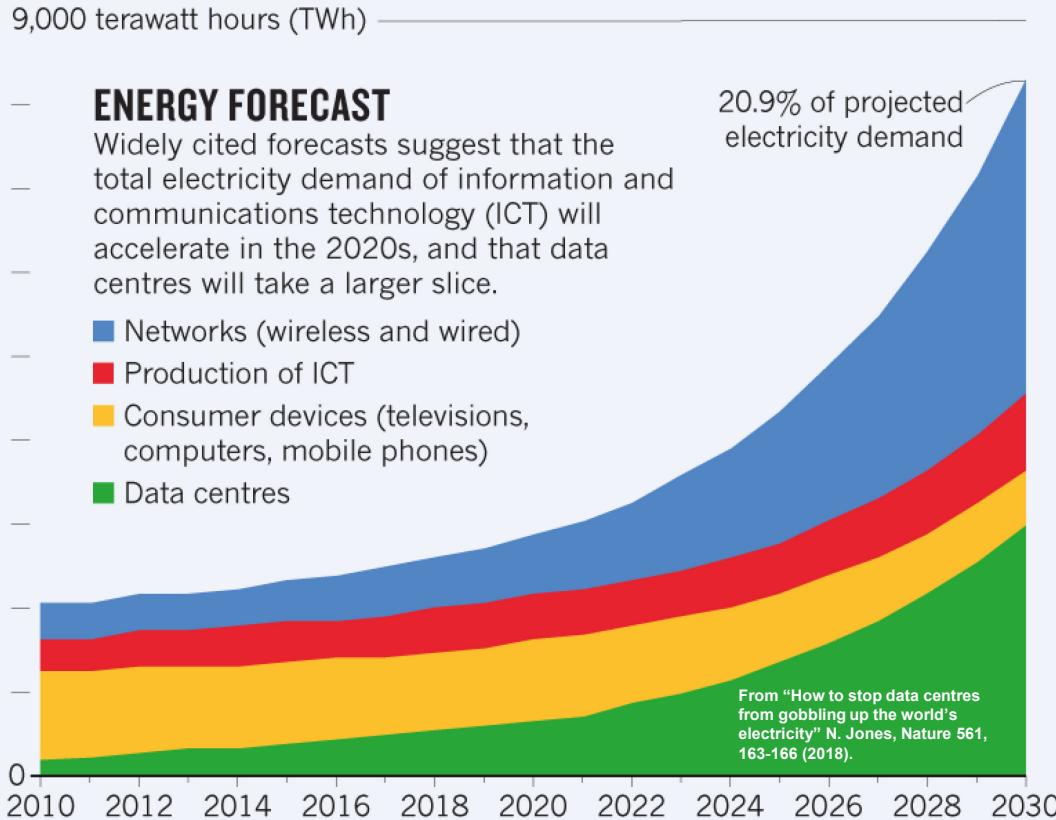
Typical first year for a new cluster:

- ~0.5 **overheating** (power down most machines in <5 mins, ~1-2 days to recover)
- ~1 **PDU failure** (~500-1000 machines suddenly disappear, ~6 hours to come back)
- ~1 **rack-move** (plenty of warning, ~500-1000 machines powered down, ~6 hours)
- ~1 **network rewiring** (rolling ~5% of machines down over 2-day span)
- ~20 **rack failures** (40-80 machines instantly disappear, 1-6 hours to get back)
- ~5 **racks go wonky** (40-80 machines see 50% packetloss)
- ~8 **network maintenances** (4 might cause ~30-minute random connectivity losses)
- ~12 **router reloads** (takes out DNS and external vips for a couple minutes)
- ~3 **router failures** (have to immediately pull traffic for an hour)
- ~dozens of minor **30-second blips for dns**
- ~1000 **individual machine failures**
- ~thousands of **hard drive failures**
- slow disks, bad memory, misconfigured machines, flaky machines, etc.

- Cluster is 1000s of machines, typically one or handful of configurations
- File system (GFS) + Cluster scheduling system are core services
- Typically 100s to 1000s of active jobs (some w/1 task, some w/1000s)
 - mix of batch and low-latency, user-facing production jobs



Energy consumption data centers



In this Google-owned data center, blue pipes supply cold water and red pipes return warm water to be cooled.

ENERGY SCALE

Global electricity demand

20,000 TWh

Electricity use by ICT

2,000 TWh

Data-centre electricity demand

200 TWh

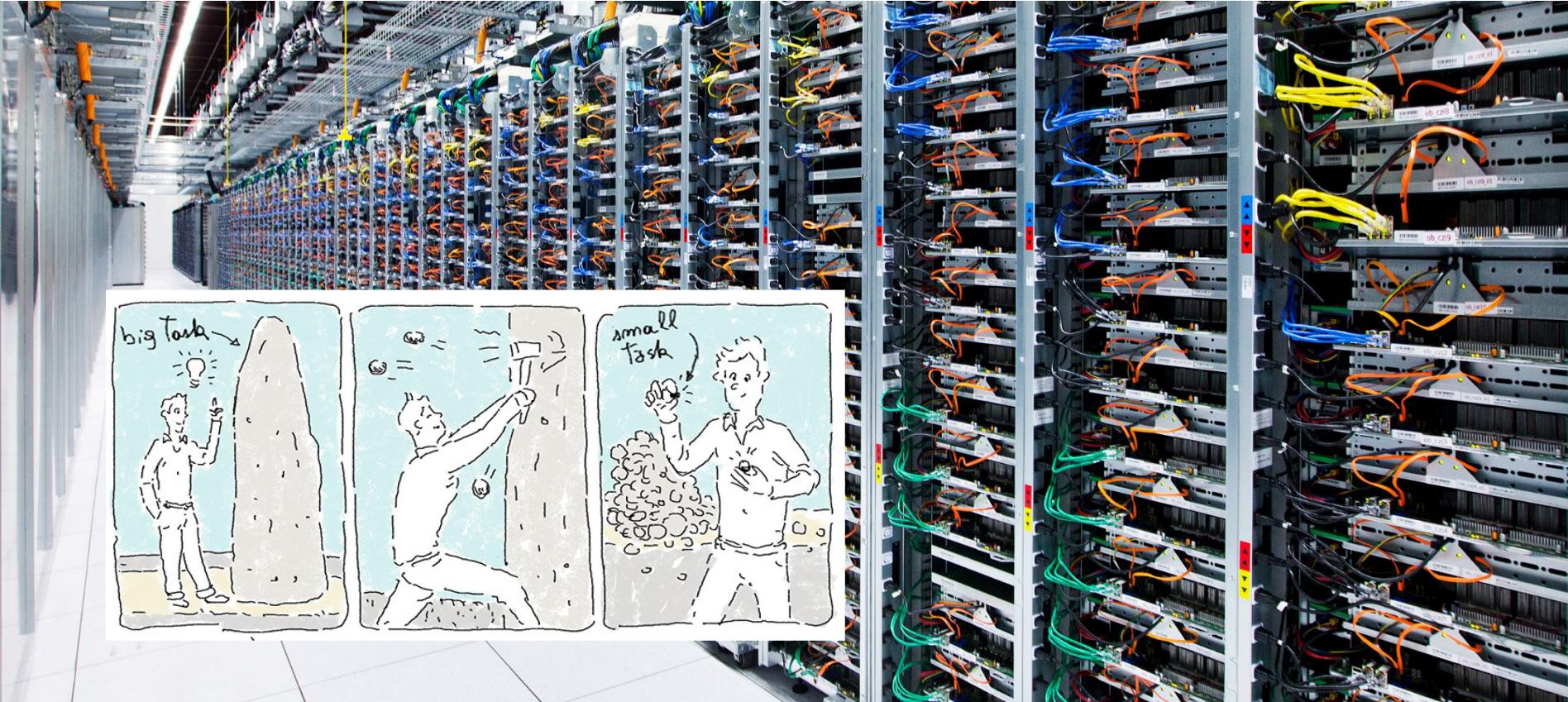
Bitcoin use by mid-2018

20 TWh

©nature

Figures are approximate.

How to exploit concurrency?



MapReduce



MapReduce (2004)

- MapReduce is a programming model.
- Speed-up through distribution.
- Part of Apache Hadoop.
- Google has the MapReduce patent (but not really new).

MapReduce: SIMPLIFIED DATA PROCESSING ON LARGE CLUSTERS

by Jeffrey Dean and Sanjay Ghemawat

1 Introduction

MapReduce is a programming model and an associated implementation for processing large datasets that is amenable to a broad variety of real-world tasks. Users specify the computation in terms of a map and reduce function, and the underlying runtime system automatically parallelizes the computation across large-scale clusters of machines, handles machine failures, and schedules inter-machine communication to make efficient use of the network and disks. Programmers find the system easy to use: more than ten thousand distinct MapReduce programs have been implemented internally at Google over the past four years, and an average of one hundred thousand MapReduce jobs are executed on Google's clusters every day, processing a total of more than twenty petabytes of data per day.

Prior to our development of MapReduce, the authors and many others at Google implemented a number of specialized purpose computations that required significant expertise to write. These included systems for generating request logs, etc., to compute various kinds of derived data, such as inverted indexes, various representations of graph structures, or Web page rank. In each case, the user had to program the parallel part of the system to handle the set of most frequent queries in a given day. Most such computations are conceptually straightforward. However, the input data is usually distributed across many machines, and the user must partition it among hundreds or thousands of machines in order to finish in a reasonable amount of time. The issues of how to parallelize the computation, distribute the data, and handle failures are complex and difficult to get right, especially when dealing with large amounts of complex code written with these issues.

As a reaction to this complexity, we designed a new abstraction that allows users to express their computation in a simple way, without having to worry about the details of parallelization, fault tolerance, data distribution and load balancing in a cluster. Our abstraction is inspired by the functional paradigm, and is similar to the well-known map and reduce idioms.

We realized that most of our computations involved applying a map operator to each logical record in an input record to compute a set of intermediate values, and then applying a reduce operator to all the intermediate values sharing the same key in order to combine the derived data appropriately. Our use of a functional model with user-specified map and reduce operators allows us to apply this language abstraction easily and to use recursion as the primary mechanism for fault tolerance.

2 Programming Model

The computation takes a set of input key-value pairs, and produces a set of output key-value pairs. The computation is specified by a library that expresses the computation as two functions: map and reduce.

Map writers tell the system, taking an input pair and producing a set of intermediate key-value pairs. Reducers take a set of intermediate key-value pairs and produce an output key-value pair. Reducers receive all intermediate values associated with the same intermediate key I and pass them to the reduce function.

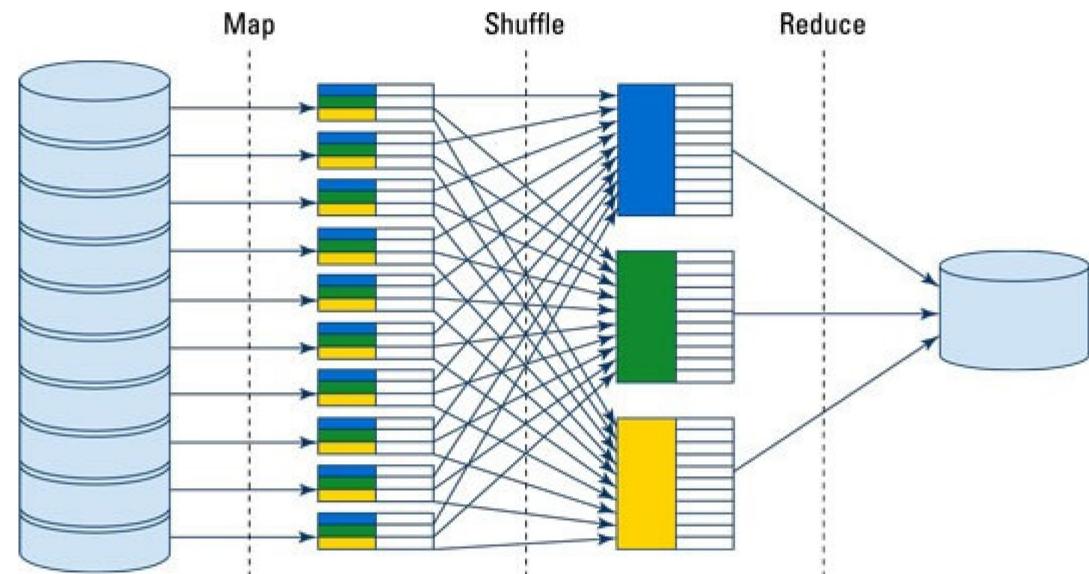
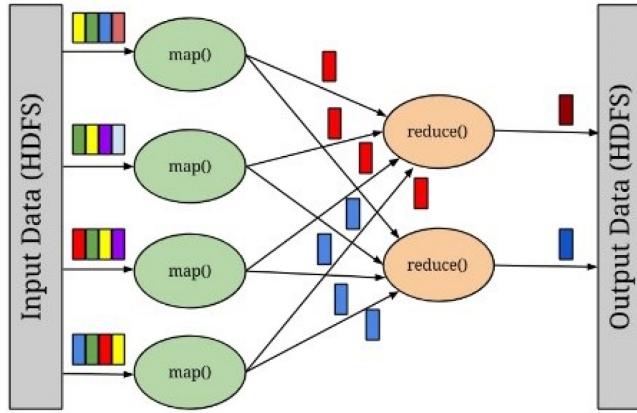
Map and reduce functions are pure, i.e., they do not modify the state of the system. They are also stateless, i.e., they do not remember previous inputs. The user, accepting an intermediate key I and a set of values for that key, it merges these values together to form a possibly smaller set of values. Typically just zero or one value is produced for each intermediate key. If multiple intermediate values are applied to the user's reduce function via an iterator, this allows us to handle lists of values that are too large to fit in memory.

3 Example

Consider the problem of counting the number of occurrences of each word in a large collection of documents. The user would write code similar to the following pseudocode:

Biographies
Jeffrey Dean currently works at Google Fellow and is a member of the Google Research team. He has worked on a large variety of large-scale distributed systems at Google, Microsoft Research, and Stanford University. He received his Ph.D. from the University of California, Berkeley.
Sanjay Ghemawat (sanja@google.com) is a Google Fellow and works on the distributed computing infrastructure used by most of the company products. He is based at Google Mountain View, CA, working

Main idea



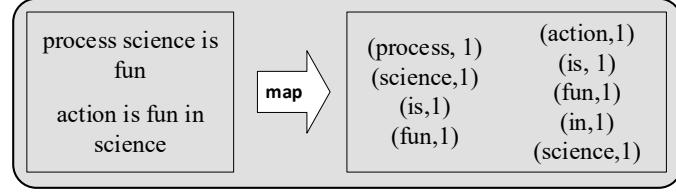
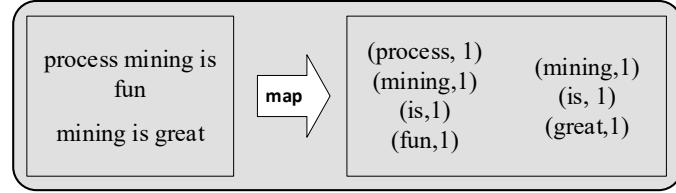
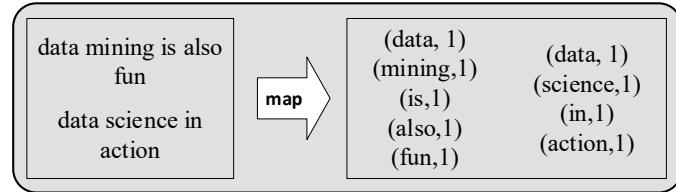
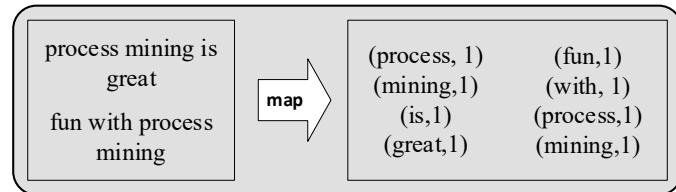
MapReduce

- K_1, K_2 are key domains.
- V_1, V_2, V_3 are value domains.
- $\text{map} \in K_1 \times V_1 \rightarrow (K_2 \times V_2)^*$ is the map function.
 - The map function is applied in parallel to every pair (keyed by K_1) in the input dataset.
 - This produces a list of pairs (keyed by K_2) for each call.
- $\text{reduce} \in K_2 \times (V_2)^* \rightarrow (V_3)^*$ is the reduce function.
 - The reduce function is applied to the grouped values (keyed by K_2).
 - Often the result is just a single value (keyed by K_2).
- The infrastructure takes care of the shuffle connecting sets of $(K_2 \times V_2)^*$ values to sets of $K_2 \times (V_2)^*$ values in a distributed environment.



Chair of Process
and Data Science

Counting words: Map step

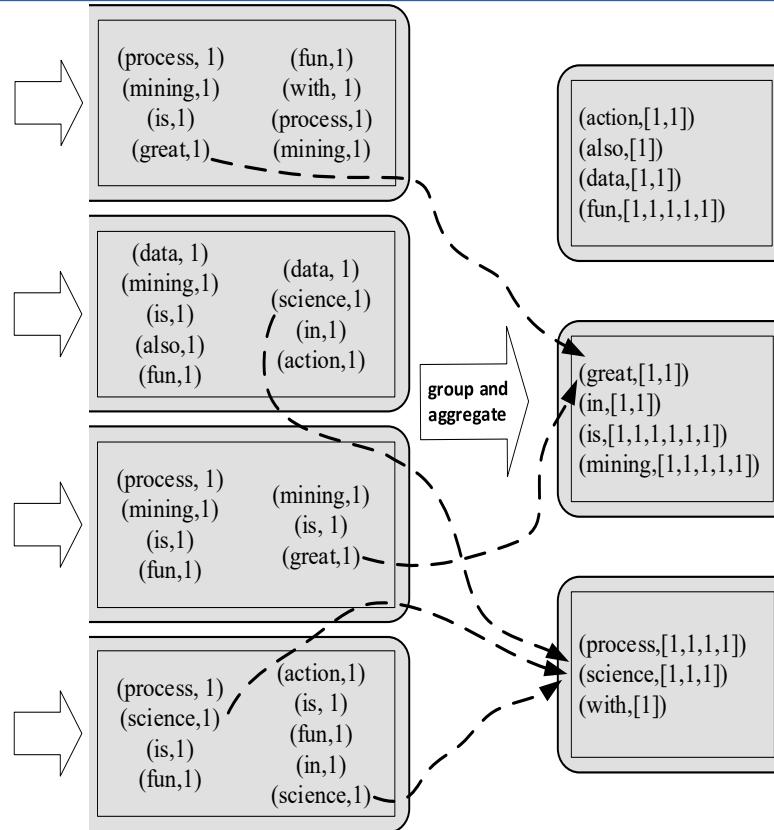


$$\text{map} \in K_1 \times V_1 \rightarrow (K_2 \times V_2)^*$$

Distributed!

Counting words: Shuffle step

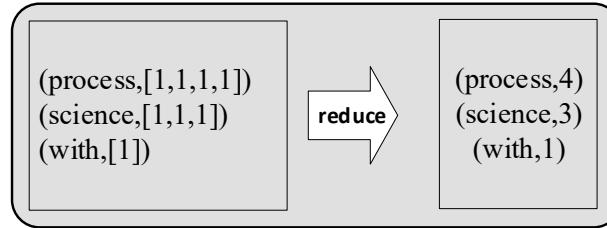
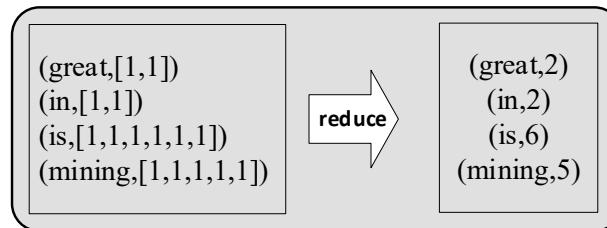
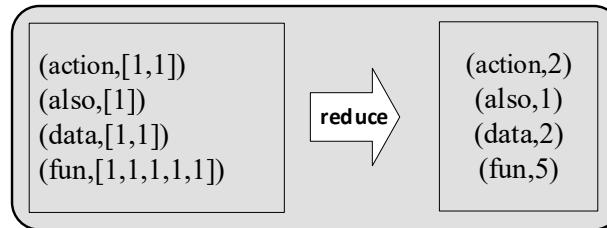
Result of
Map step



Input for
Reduce step

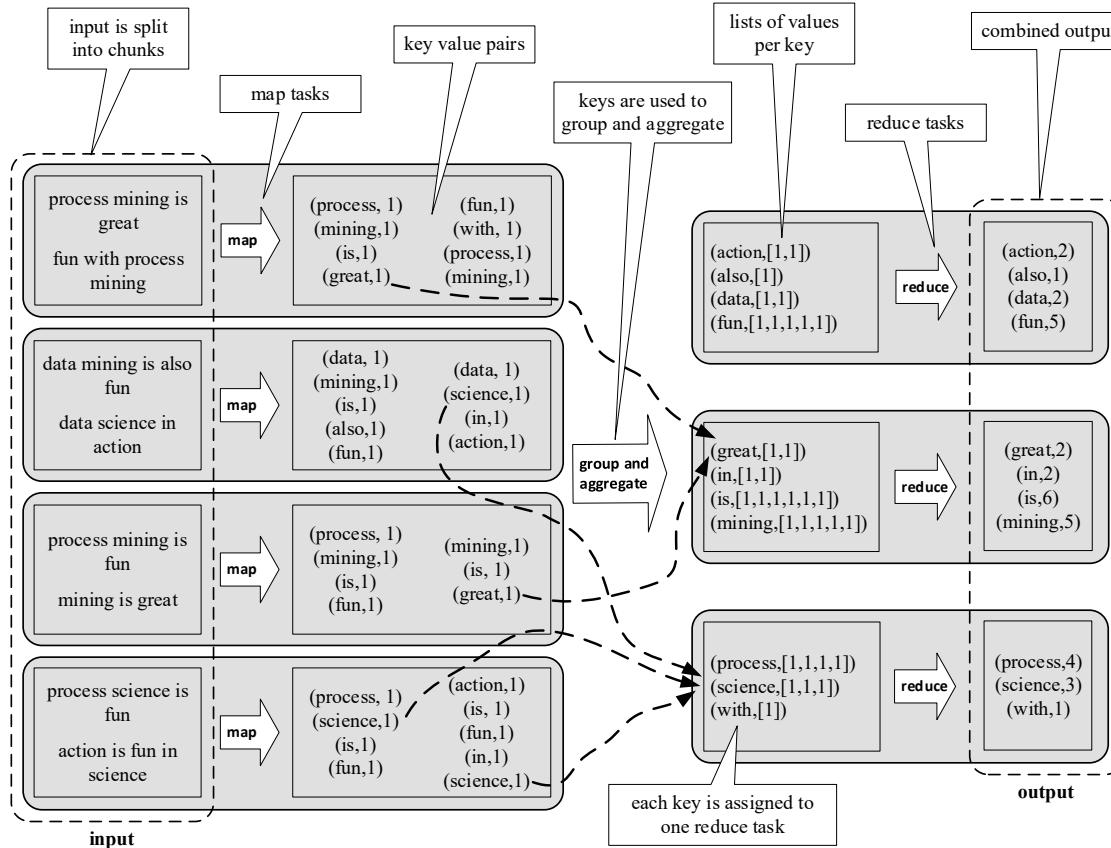
Counting words: Reduce step

reduce $\in K_2 \times (V_2)^* \rightarrow (V_3)^*$

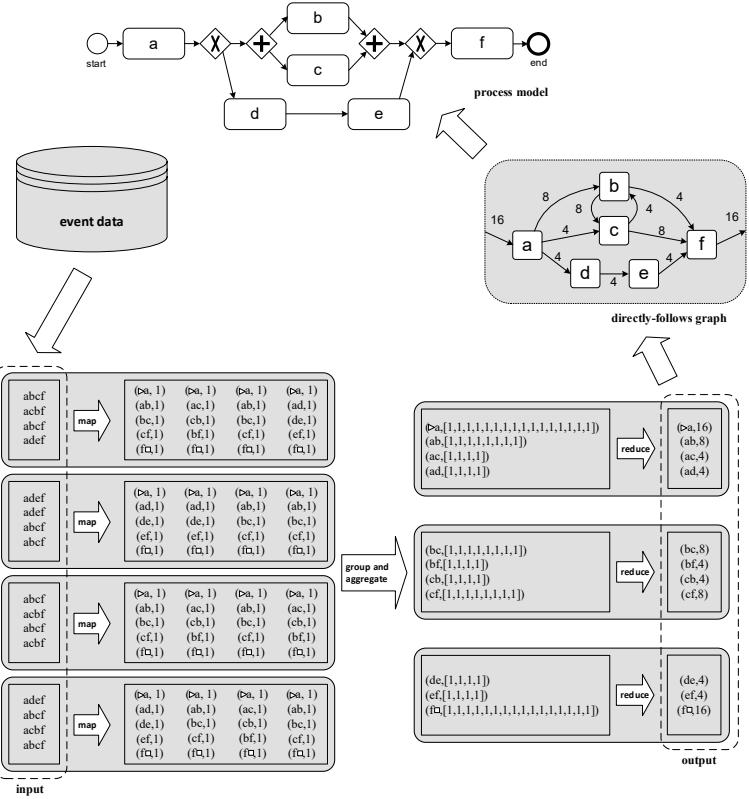


Distributed!

Overview

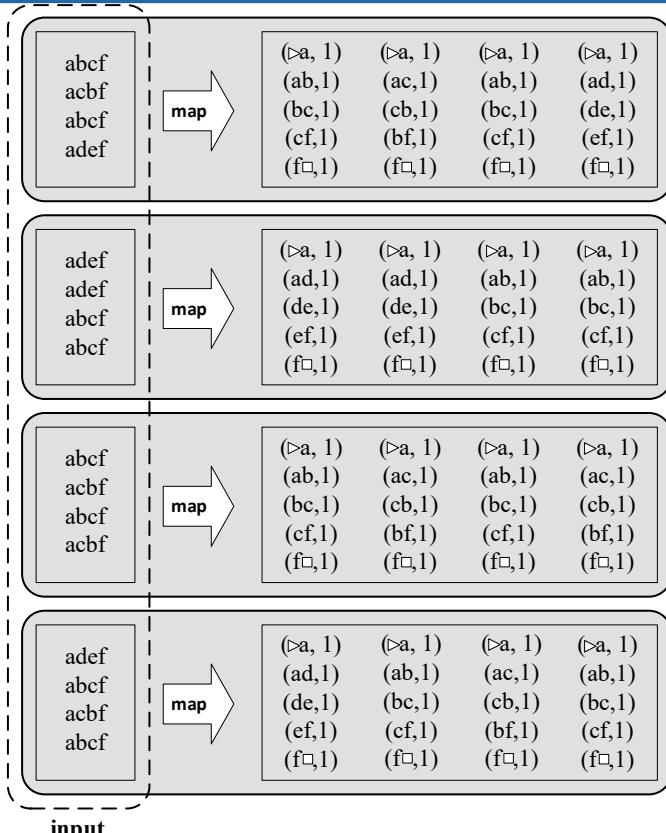


MapReduce applied to process discovery



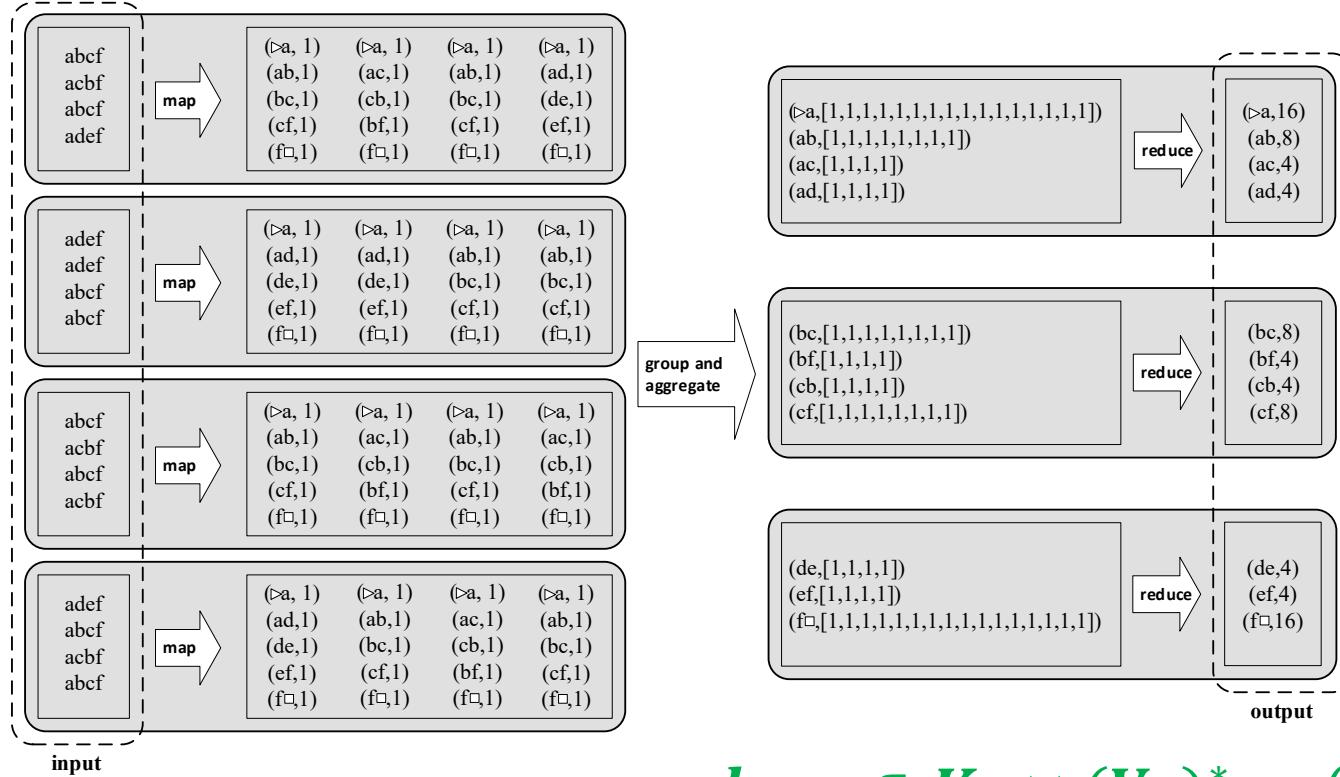
- We need to count how many times activity x is followed by y in an event log with billions of events.
- Assume traces are split over multiple nodes.

MapReduce applied to process discovery



$$map \in K_1 \times V_1 \rightarrow (K_2 \times V_2)^*$$

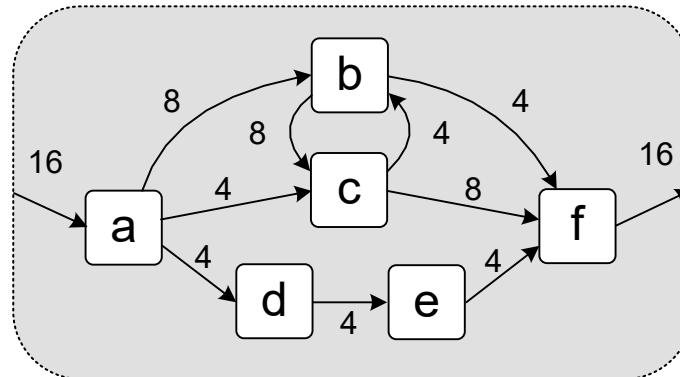
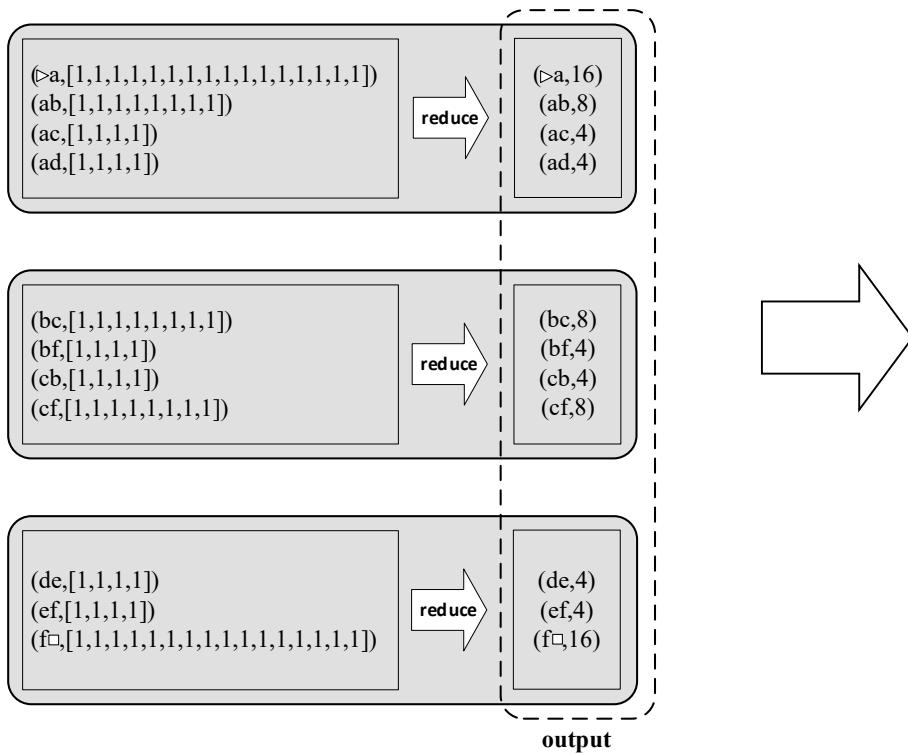
MapReduce applied to process discovery



$$\text{reduce} \in K_2 \times (V_2)^* \rightarrow (V_3)^*$$

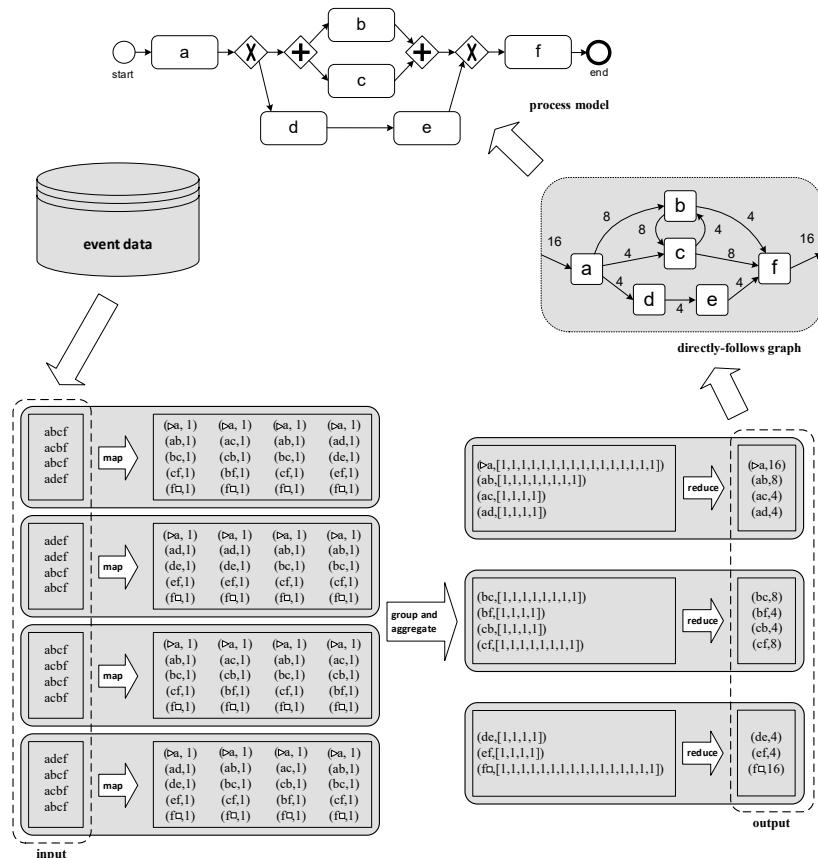


MapReduce applied to process discovery



directly-follows graph

MapReduce applied to process discovery



This way process discovery is scalable at the scale of Google (provided you have enough hardware).



Chair of Process
and Data Science

Streaming

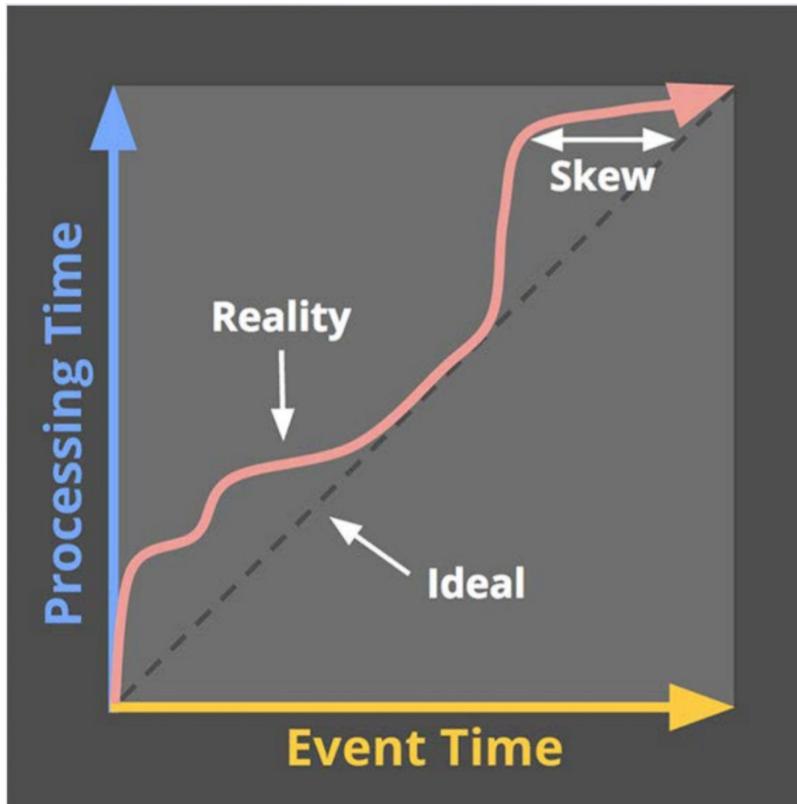


Streaming



streaming data

Streaming: Respond “immediately”



- Event time: the time at which the data object occurred
- Processing time: the time at which the data object is observed in streaming system



Streaming: Limited capacity



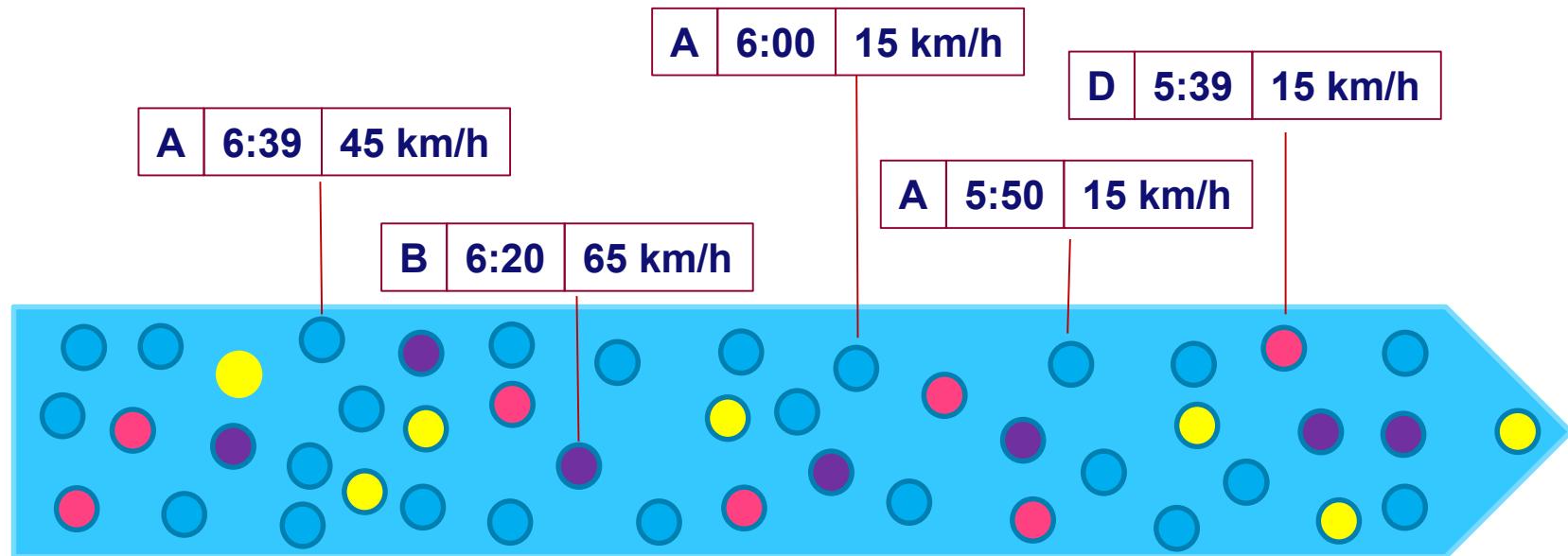
- Limited capacity
- When “new data” get stored, “old data” need to be removed.

Streaming Examples

- Stock market surveillance (e.g., predict the trend of stock).
- Identification of life-threatening conditions in health care (e.g., continuously monitor the physical signs of patients and alert when anomaly happens).
- Computer system and network monitoring (e.g., monitor and analyze the amount of packets sent from the IP addresses that access web server).
- Dealing with data stream from various kinds of sensors (e.g., continuously analyze the sensors in aircrafts for anomalies).

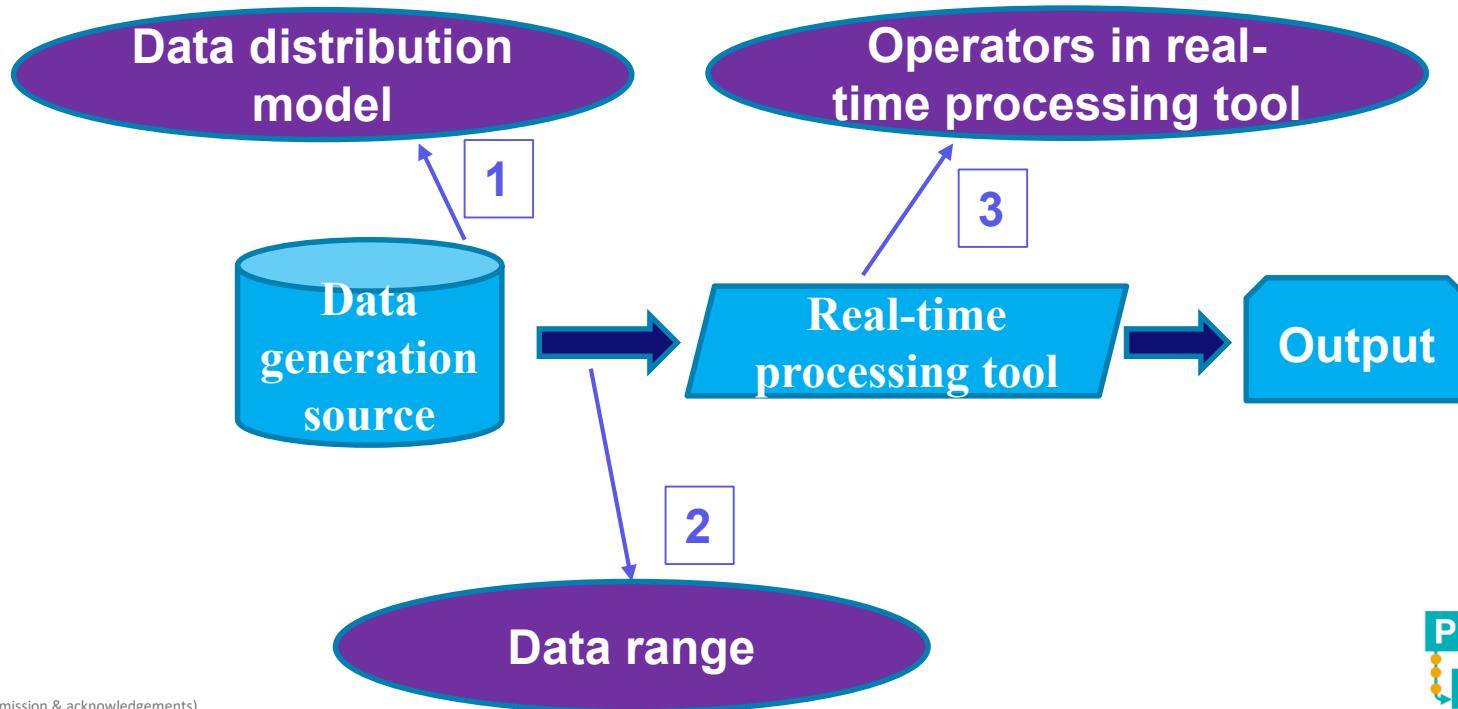
Streaming

- An example data stream coming from cars that are monitored



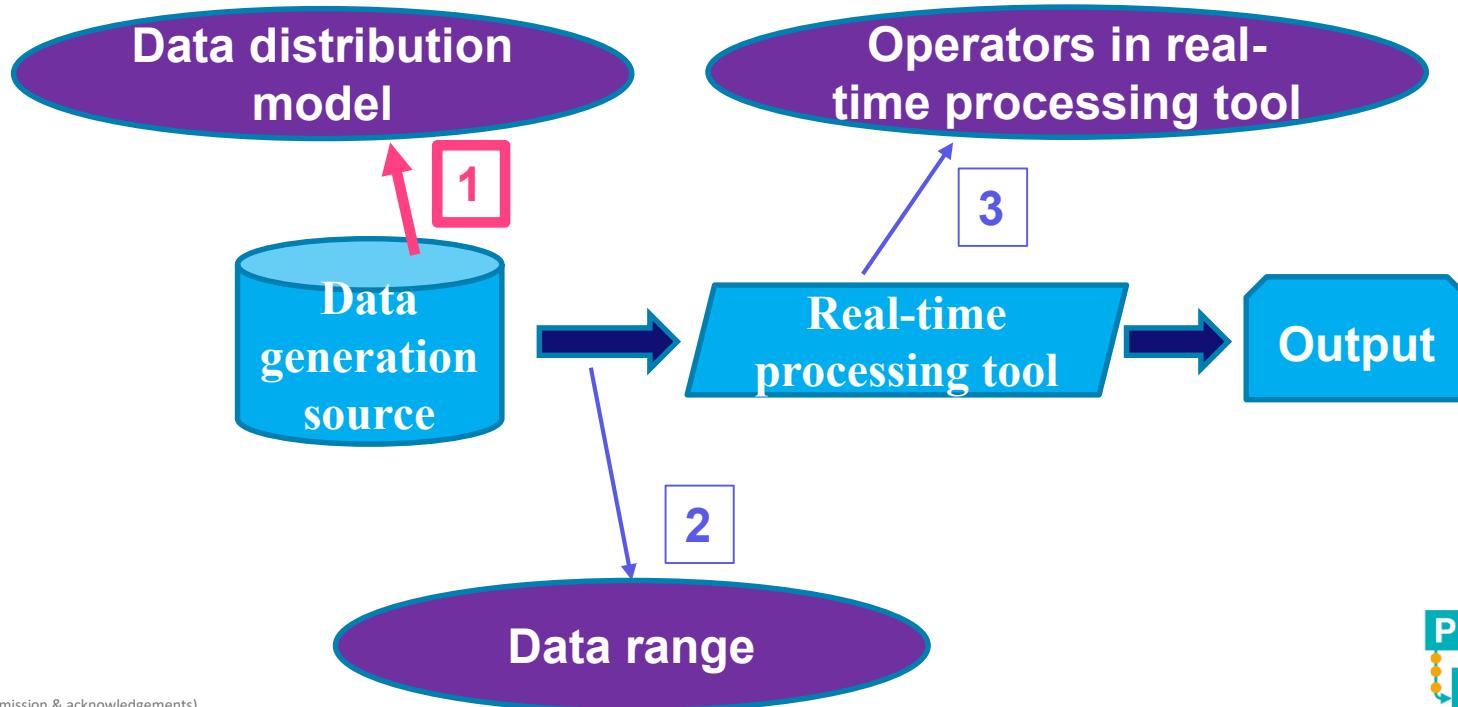
Streaming

- Important elements in streaming



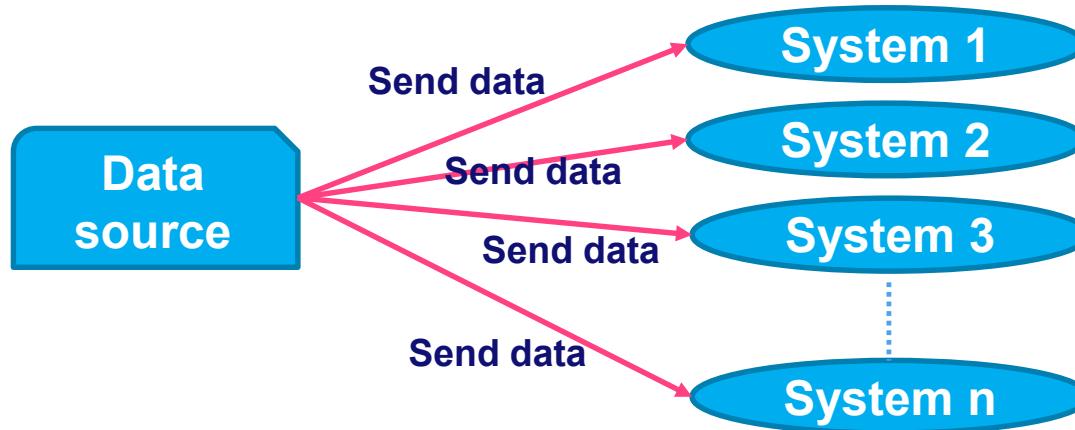
Streaming

- Important elements in streaming



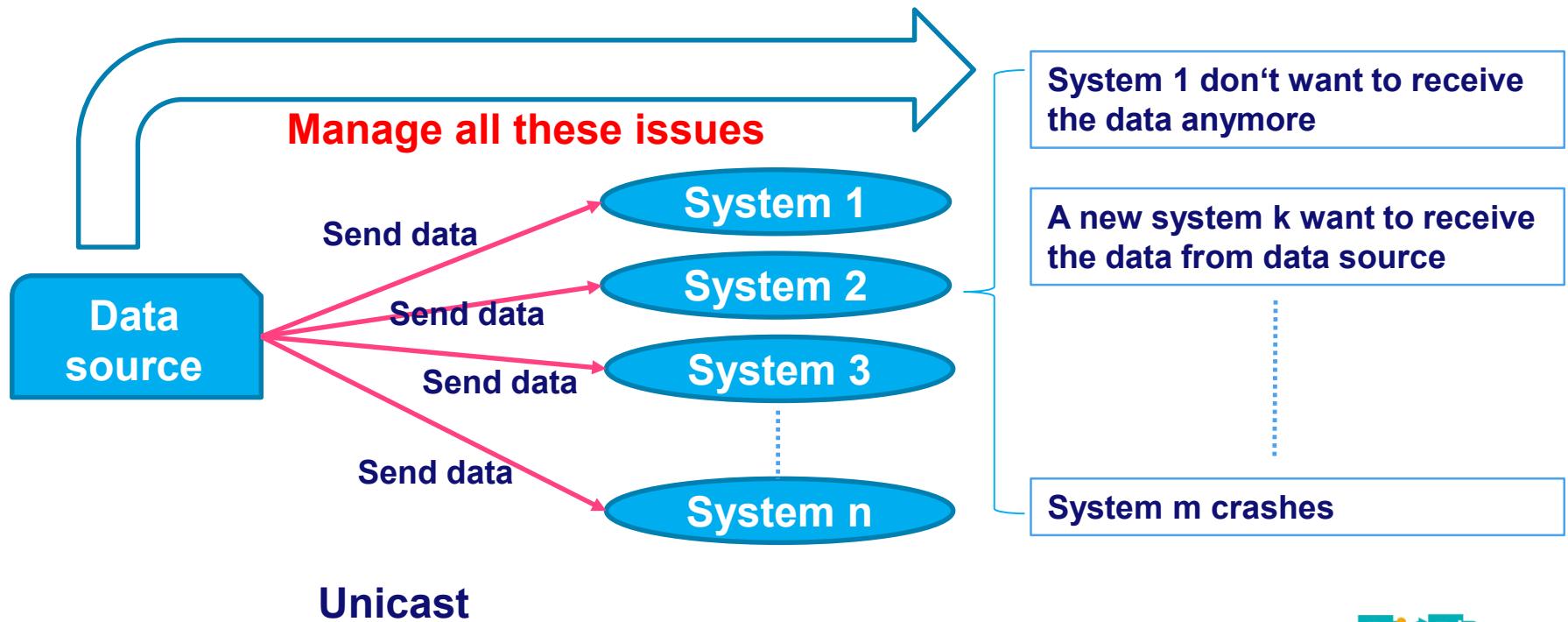
Streaming

- **Data distribution model**
 - Why not unicast?



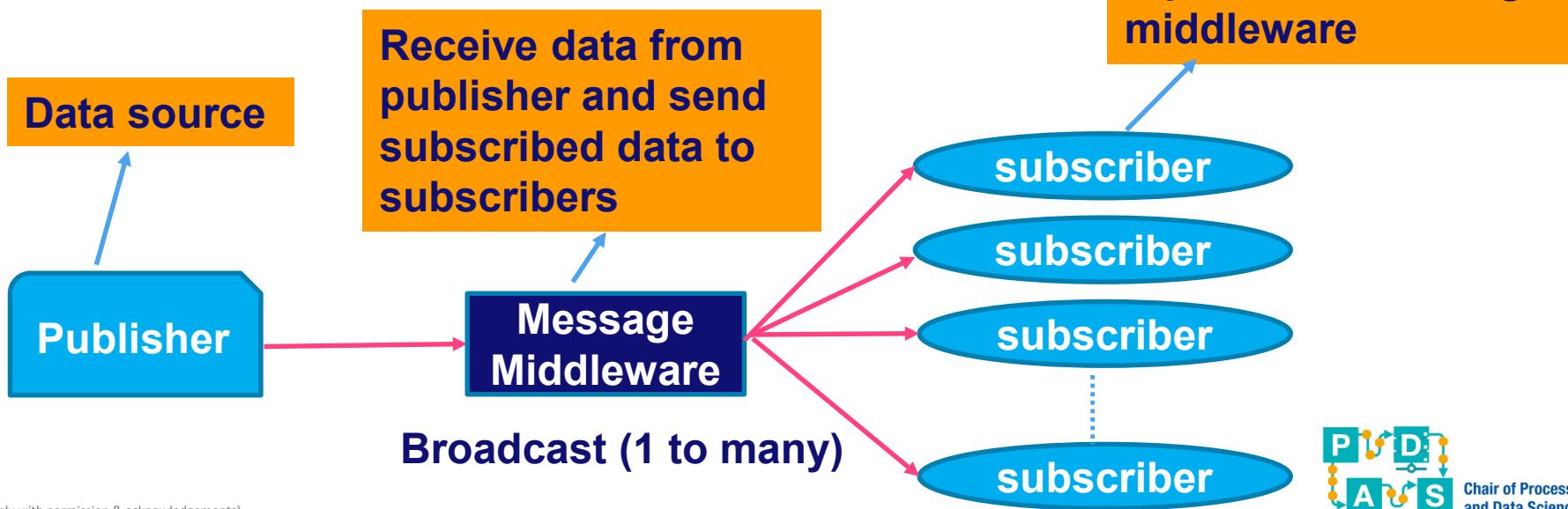
Unicast (one-to-one transmission)

Streaming



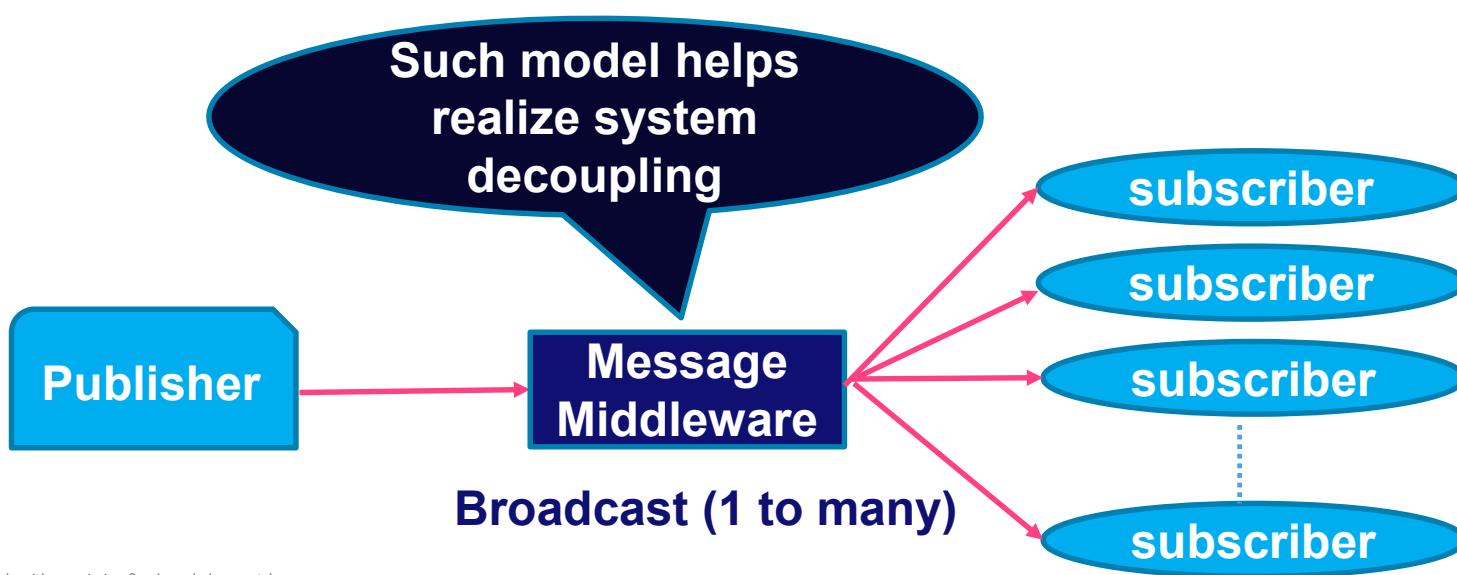
Streaming

- Data distribution model
 - Publish/Subscribe messaging model



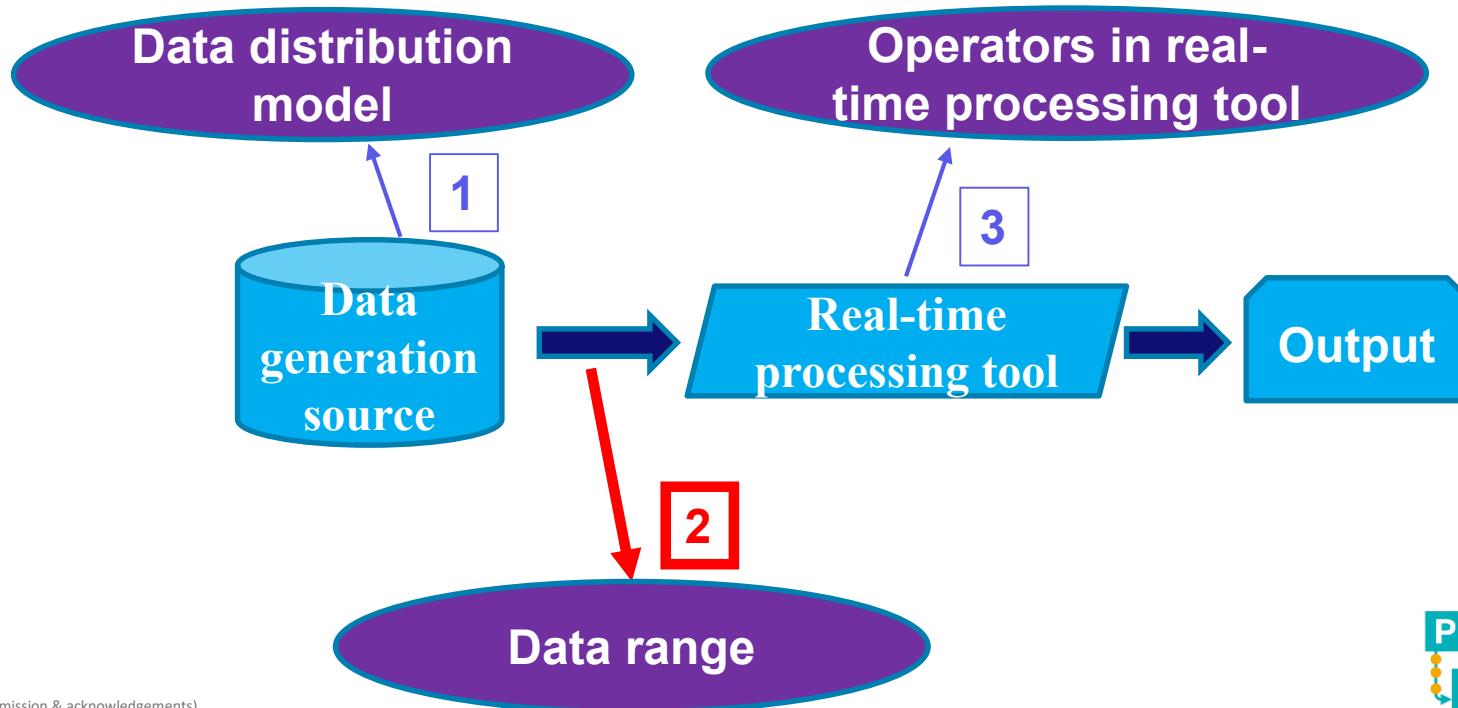
Streaming

- Data distribution model
 - Publish/Subscribe messaging model



Streaming

- Important elements in streaming

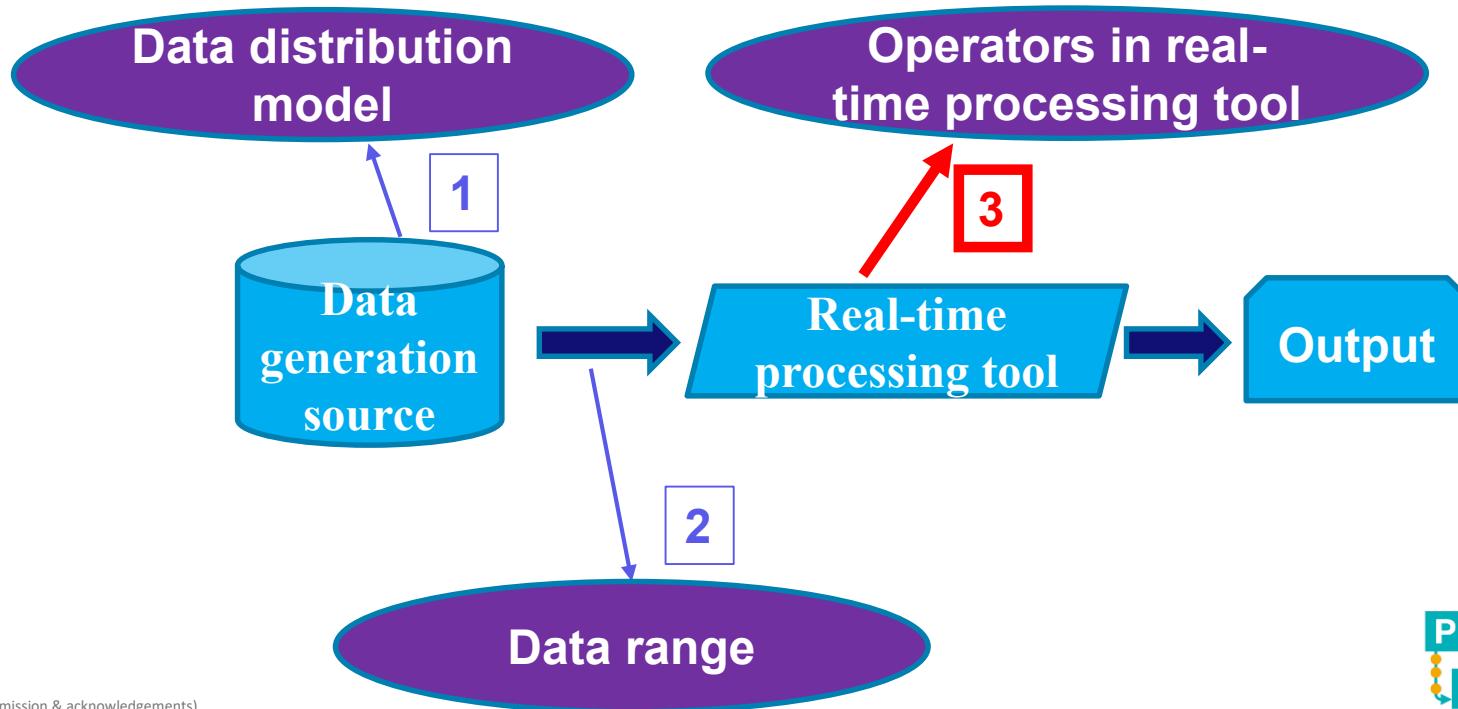


Data range

- What data to consider?
 - **Data stream model:** Consider all the data objects from a specific time point until now (e.g., for calculating the average speed of a car, consider all its speed data received until now).
 - **Sliding window model:** compute over a fixed size of window in the stream (e.g., compute the average speed of a car, only consider the speed data received during the last 10 days).
- Due to capacity constraints not everything can be stored or processed!

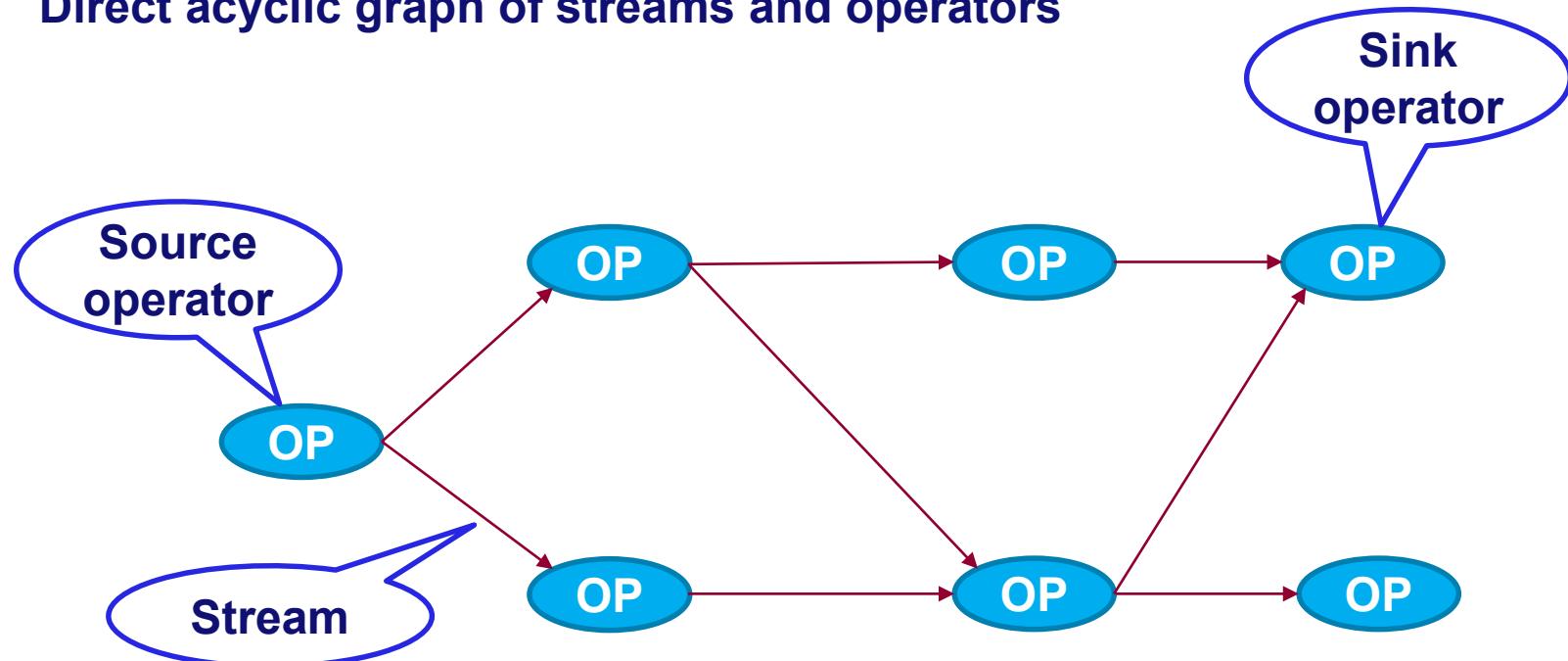
Streaming

- Important elements in streaming



Operator in real-time processing

Direct acyclic graph of streams and operators

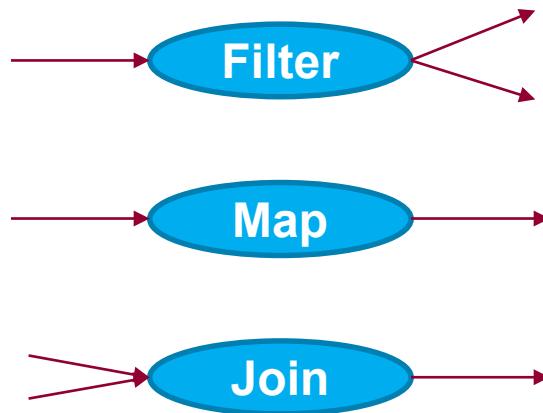


Two types of operators in real-time processing

- **Stateless operators:**
 - Don't record historical data object in memory
 - One-by-one processing possible
- **Stateful operators:**
 - Need to record historical data object in memory
 - Outputs depend on multiple input data objects

Streaming

- **Stateless operators**



Filter data objects according to one or more conditions.

Convert each data object.

Merge multiple streams into one stream.

Streaming

- **Stateful operators**



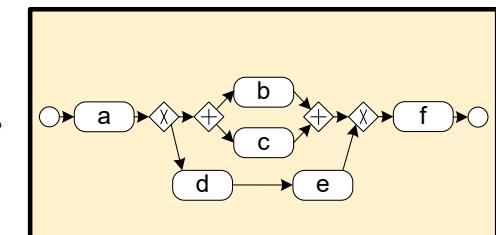
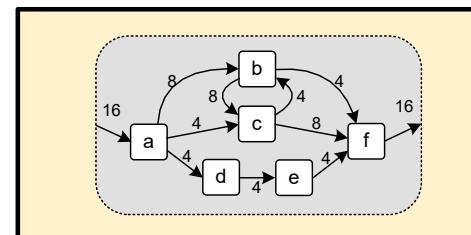
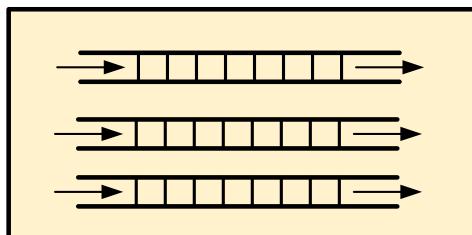
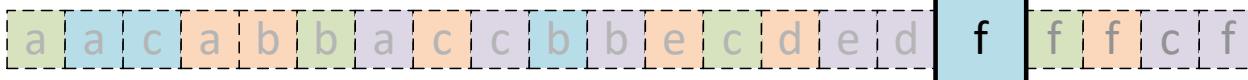
Aggregate multiple data objects.
(e.g., compute average speed of
data objects in last five minutes)



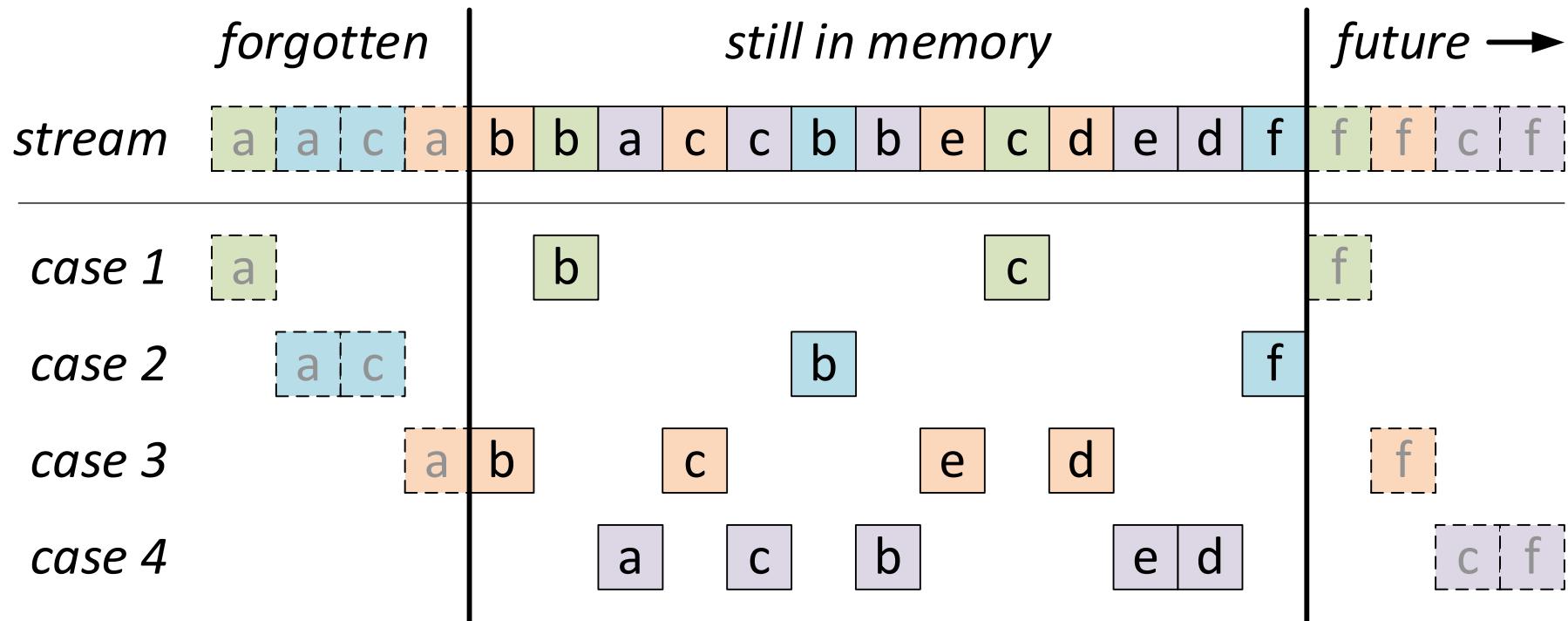
Join data objects from two streams given
a specific predicate.
(e.g., for the last 3 data objects from each
stream, join every pair with the same ID)

Example: Streaming process mining

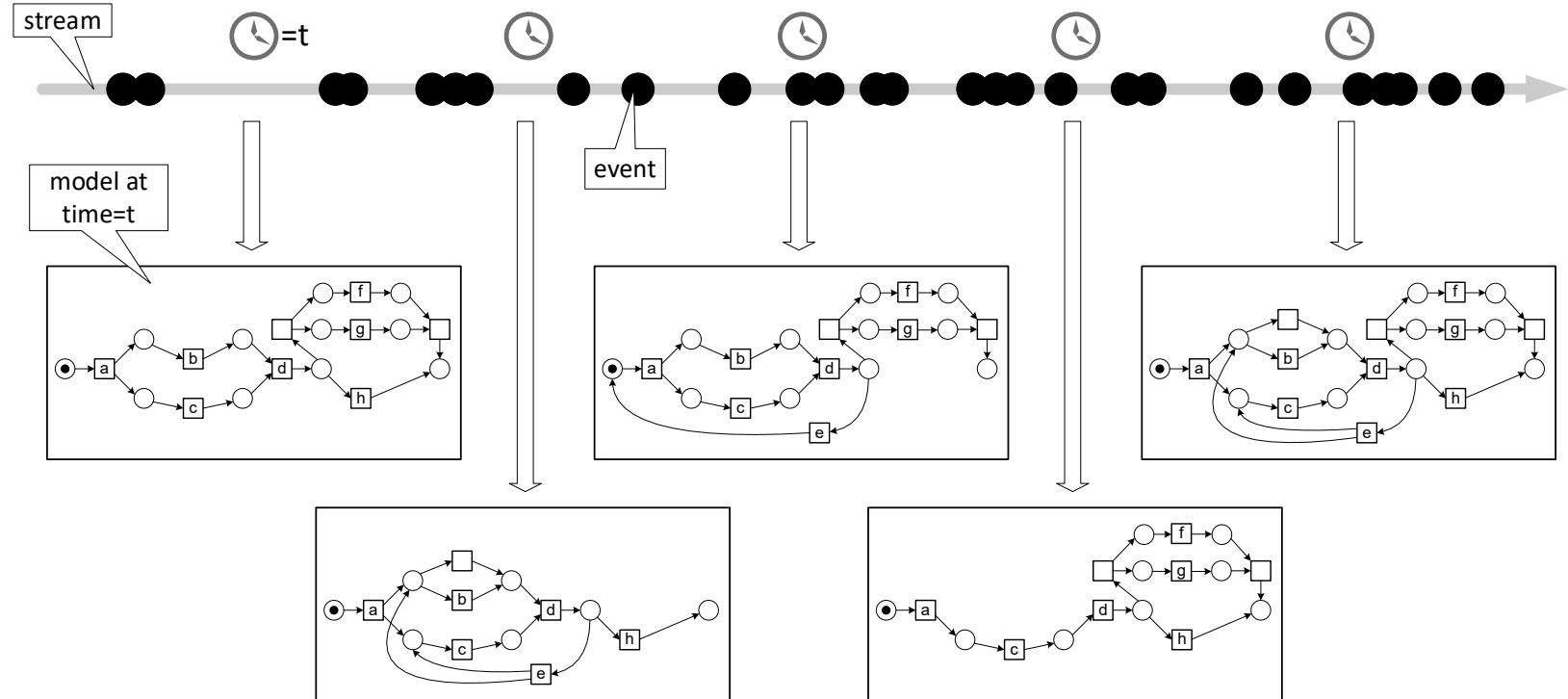
stream



Example: Streaming process mining



Example: Streaming process mining



Relates to concept drift (see later).



Streaming k-means



Initial non-streaming k-means algorithm

Algorithm: *k*-means. The *k*-means algorithm for partitioning, where each cluster's center is represented by the mean value of the objects in the cluster.

Input:

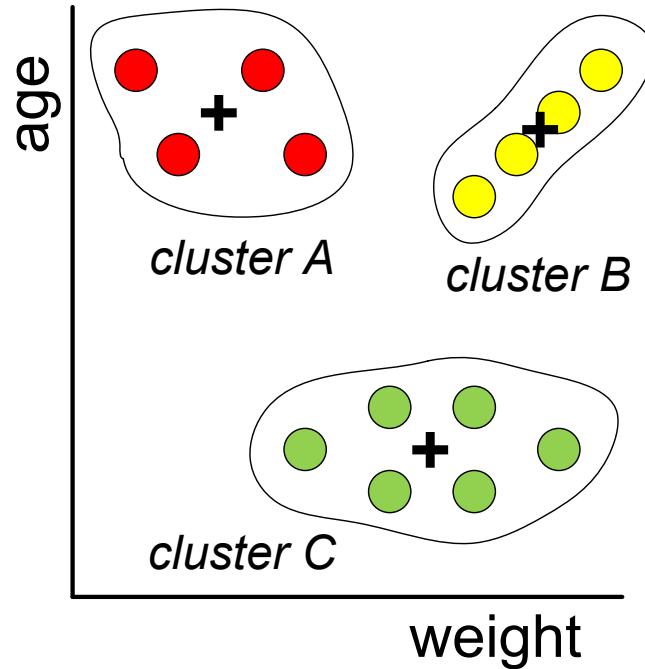
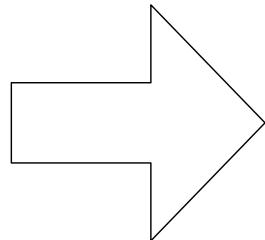
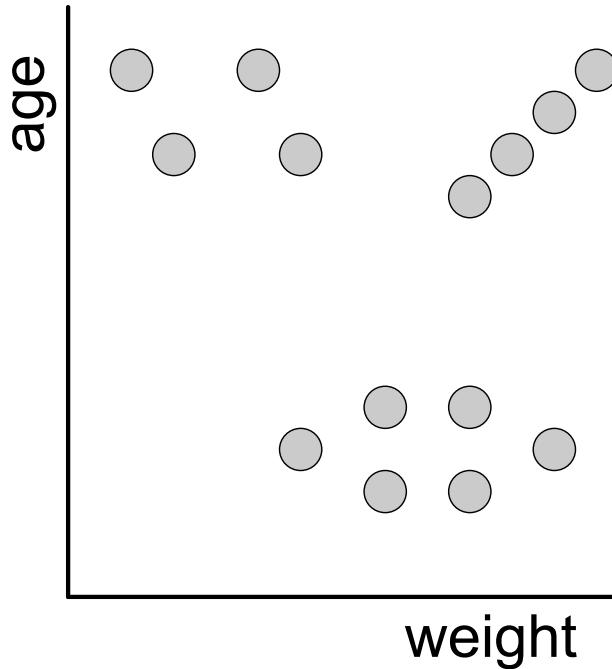
- k : the number of clusters,
- D : a data set containing n objects.

Output: A set of k clusters.

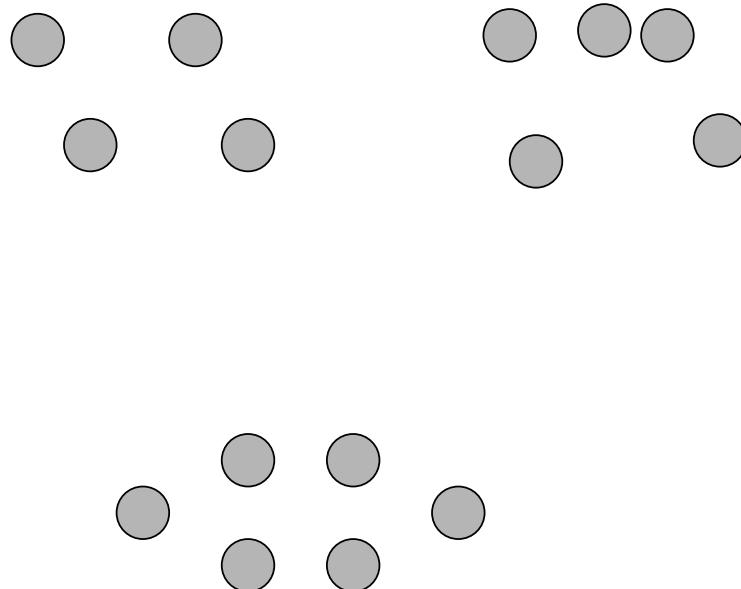
Method:

- (1) arbitrarily choose k objects from D as the initial cluster centers;
- (2) **repeat**
- (3) (re)assign each object to the cluster to which the object is the most similar,
 based on the mean value of the objects in the cluster;
- (4) update the cluster means, that is, calculate the mean value of the objects for
 each cluster;
- (5) **until** no change;

Clustering



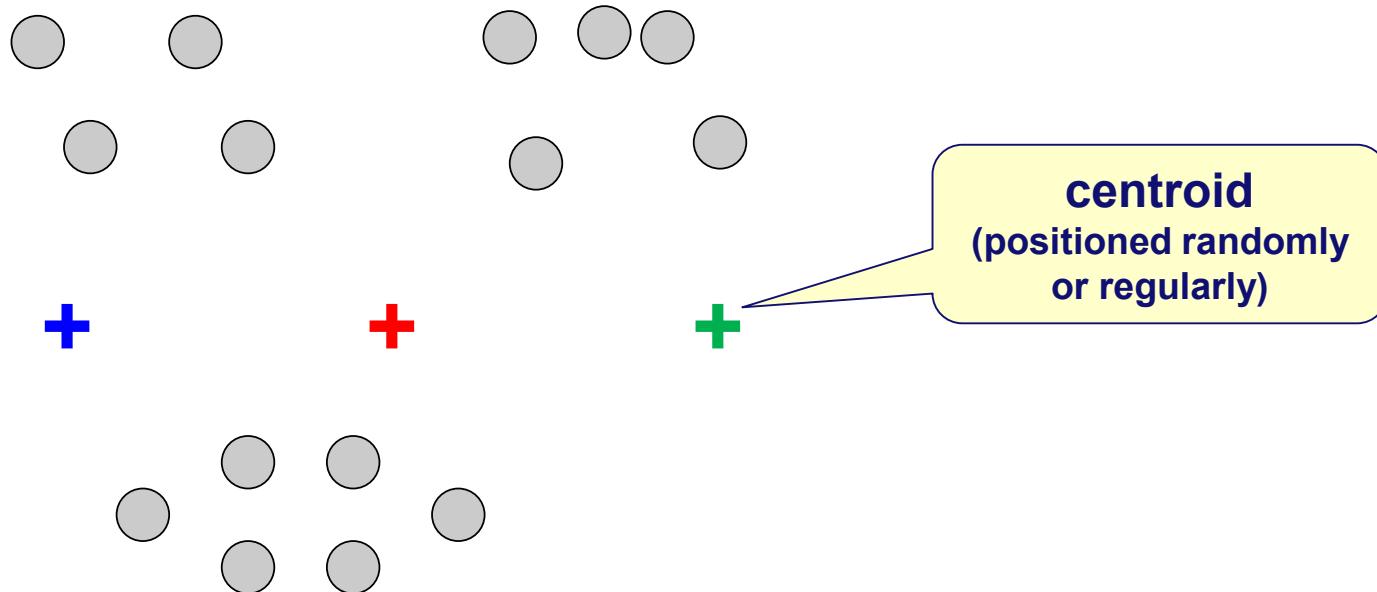
k -means clustering



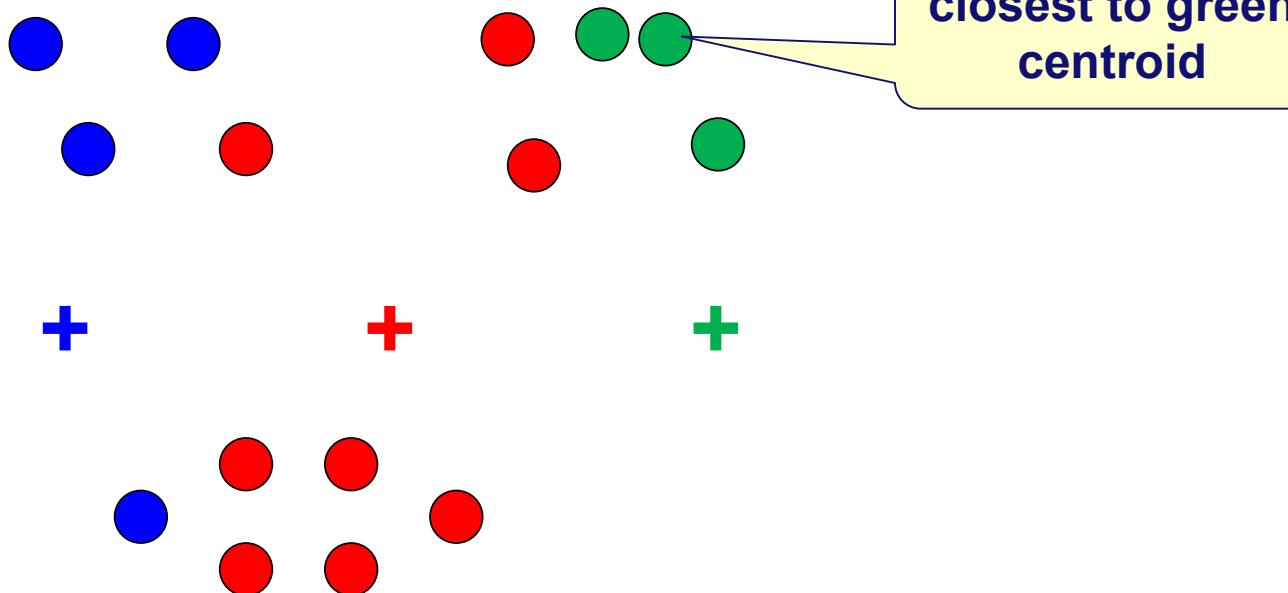
assume $k=3$

**Just an illustration.
May be misleading:
often many
dimensions!**

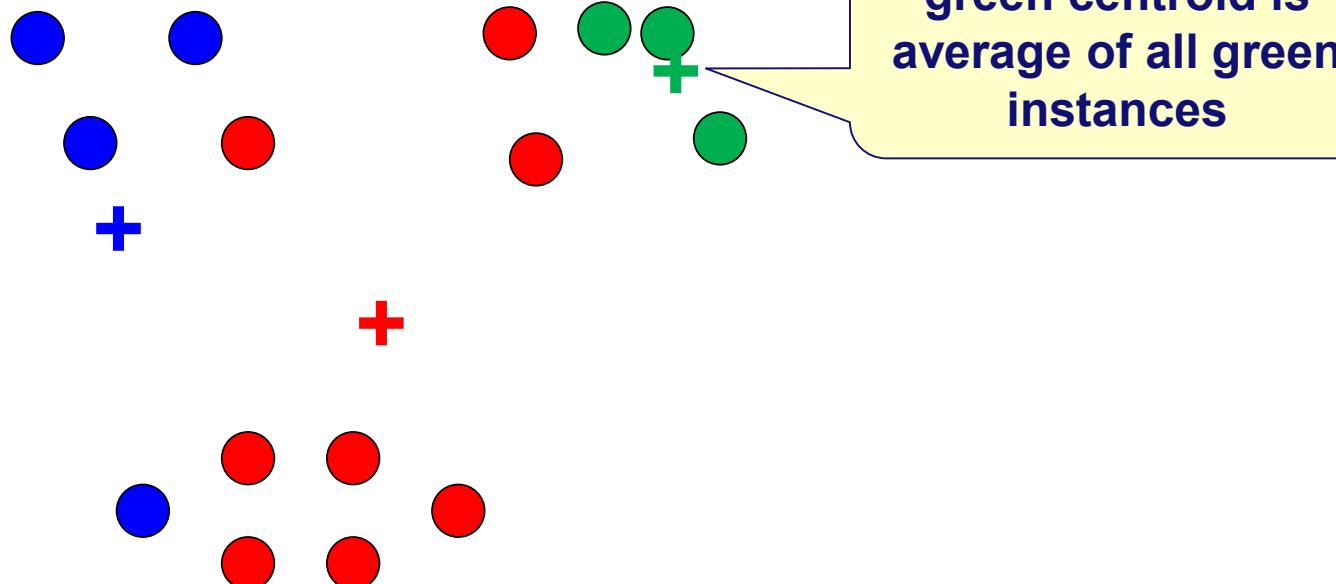
Generate $k=3$ centroids (e.g., random)



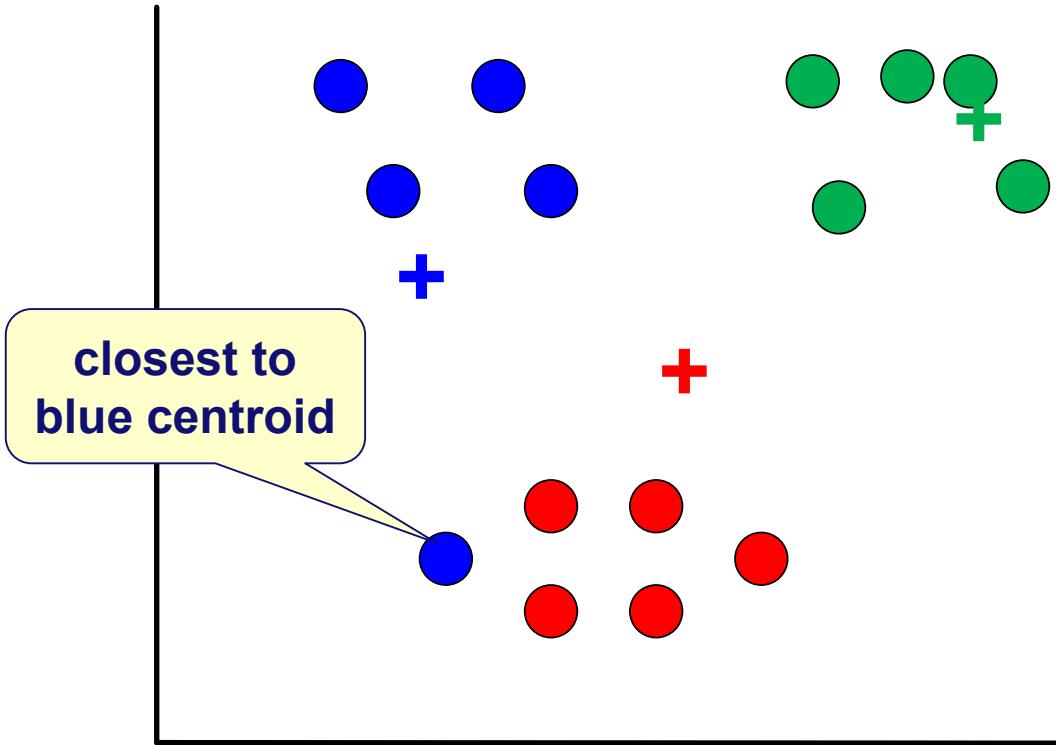
Assign instances to closest centroid



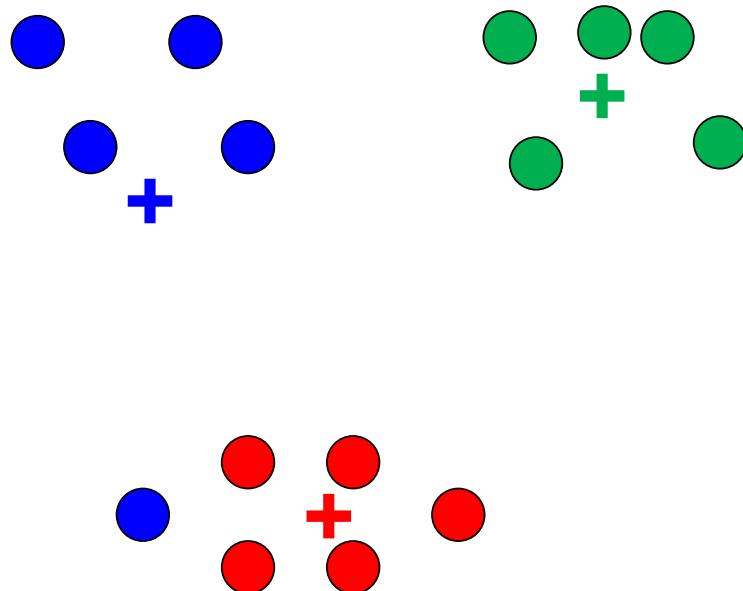
Recompute centroids based on assigned instances



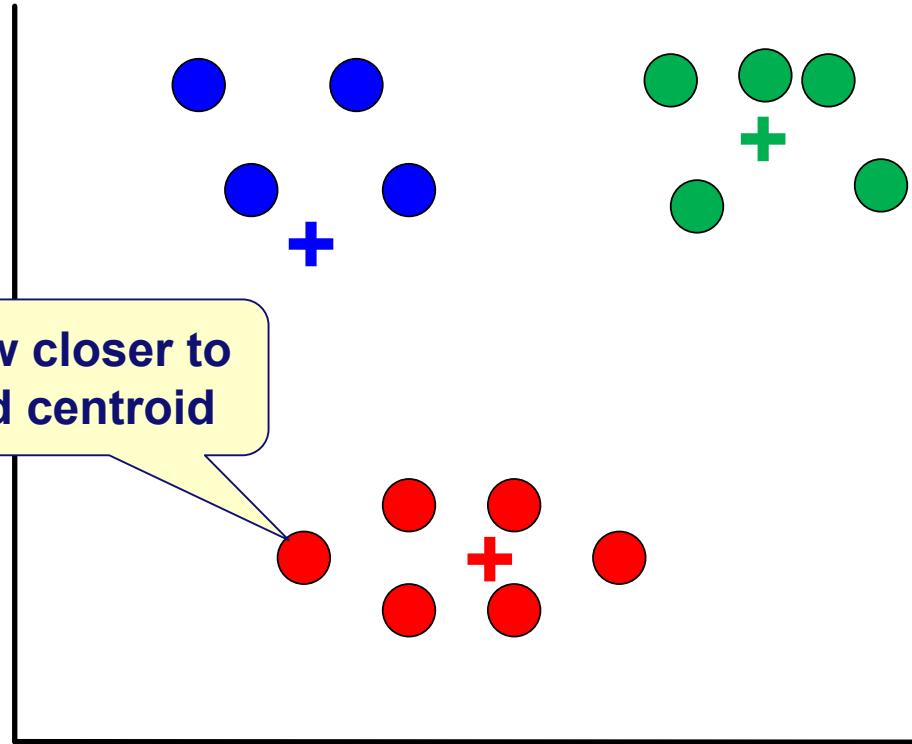
Assign instances to closest centroid



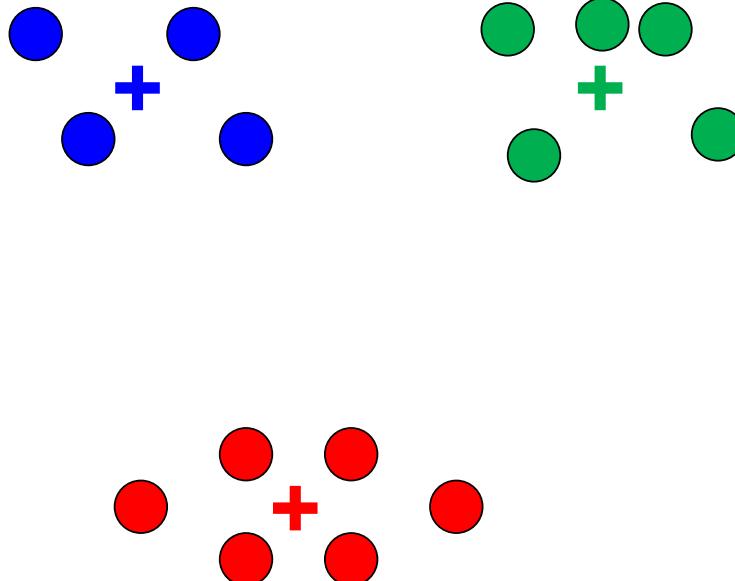
Recompute centroids based on assigned instances



Assign instances to closest centroid



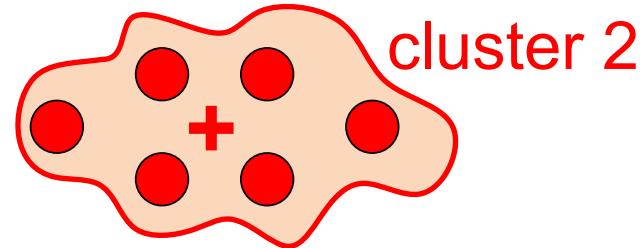
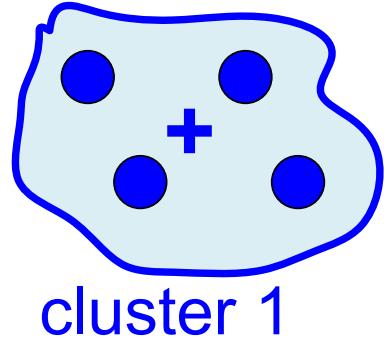
Recompute centroids based on assigned instances



Fixpoint has been reached: Nothing changes anymore.

Due to non-deterministic nature (random initial centroids), experiment may be repeated multiple times. Select "best clustering" at end.

Clusters returned



Main idea:
The instances in a cluster are more similar to each other than to those in other clusters.



**Assume now that we have limited
storage and are too slow to access old
events and do computations over them!**

Streaming k-means algorithm

- Assume we **cannot store all instances**.
- New instances arrive in **batches** and are processed in **batch** (and then forgotten).
- Behavior is expected to change over time (**concept drift**). The clusters should reflect **both past and current** instances.
- Many variations are possible, we just describe one variant.

Some notation (1/2)

- At any point in time there are k clusters $\{1, 2, \dots, k\}$.
- c_i^t is the **centroid** of cluster $i \in \{1, 2, \dots, k\}$ before the t -th batch B^t arrives.
- Hence, $\{c_1^t, c_2^t, \dots, c_k^t\}$ are the centroids at time t .
- $\{c_1^1, c_2^1, \dots, c_k^1\}$ are the initial centroids.

Some notation (2/2)

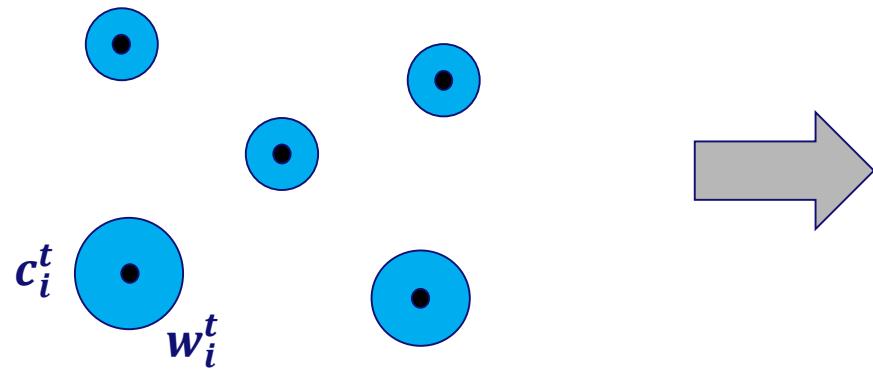
- At any point in time there are k clusters $\{1, 2, \dots, k\}$.
- w_i^t is the **weight** of cluster i before the t -th batch B^t arrives.
- Hence, $\{w_1^t, w_2^t, \dots, w_k^t\}$ are the **weights** at time t .
- $\{w_1^1, w_2^1, \dots, w_k^1\}$ are the **initial weights**.

Main idea

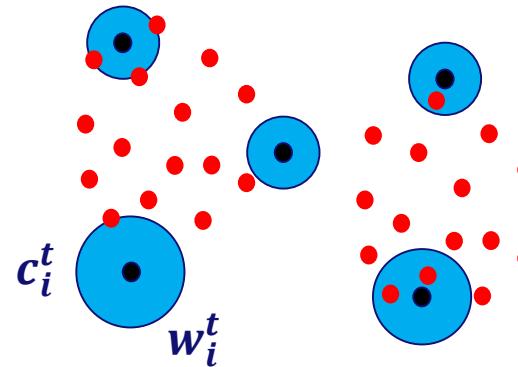
- For each new batch arrival B^t :
 - Assign each instance to the nearest cluster centroid.
 - Update the cluster centroids and weights.
- If the weight of a cluster i drops below a preset threshold δ , then cluster i is removed and the largest cluster j is split into two clusters, one taking over the role of i .
- I.e., the old cluster j is split into the new clusters i and j .

Example: Update clusters

clusters at time t

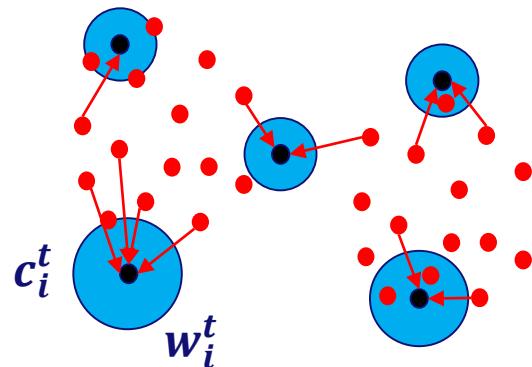


new batch

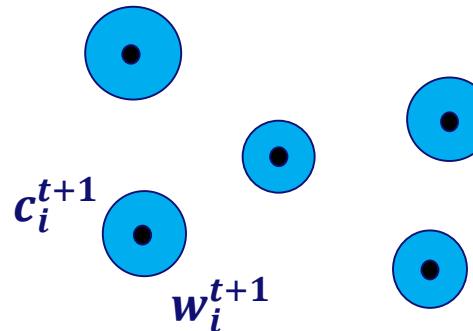


Example: Update clusters

clusters at time t and new batch

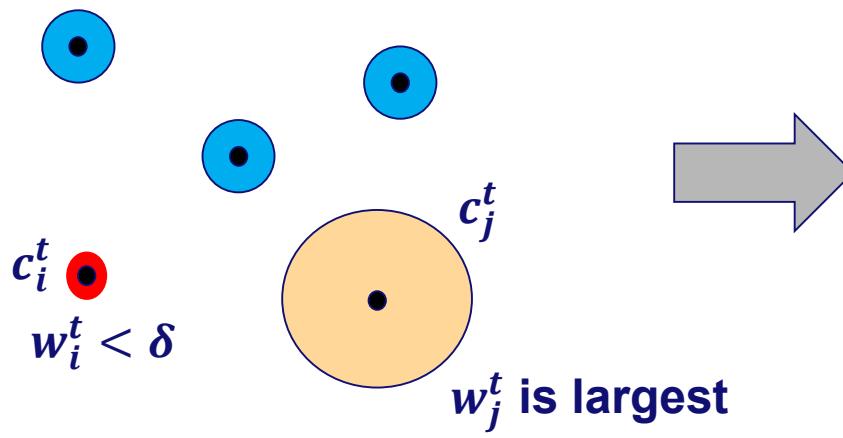


clusters at time $t + 1$

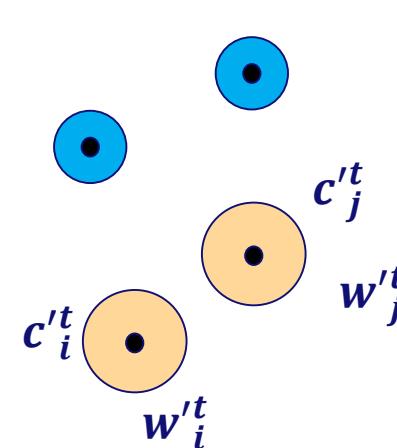


Example: Splitting

subcritical cluster



largest cluster is split



Different ways of splitting c_j^t (outside of scope).

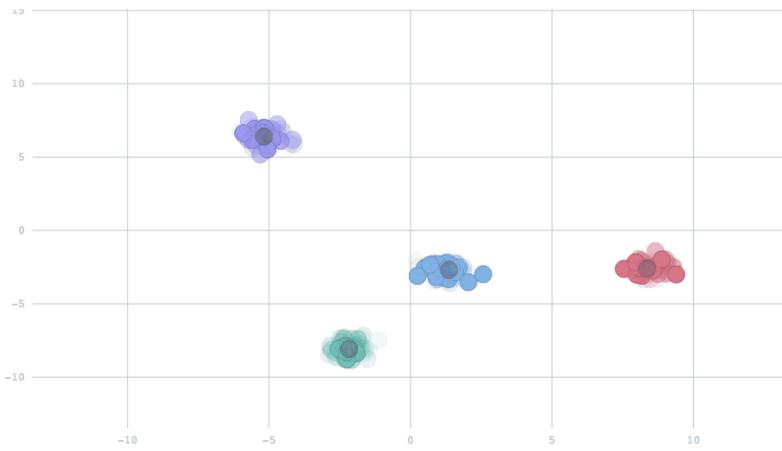
Batch arrival $B^t = B_1^t \cup B_2^t, \dots, B_k^t$

- $\{c_1^t, c_2^t, \dots, c_k^t\}$ and $\{w_1^t, w_2^t, \dots, w_k^t\}$ describe the centroids at time t when a batch B^t of instances arrives.
- $B^t = B_1^t \cup B_2^t, \dots, B_k^t$ are the instances assigned to the nearest centroid.
- $\{bc_1^t, bc_2^t, \dots, bc_k^t\}$ and $\{bw_1^t, bw_2^t, \dots, bw_k^t\}$ are the centroids and weights of the partitioned batches where $bw_i^t = |B_i^t|$ and bc_i^t is computed as usual (average vector).

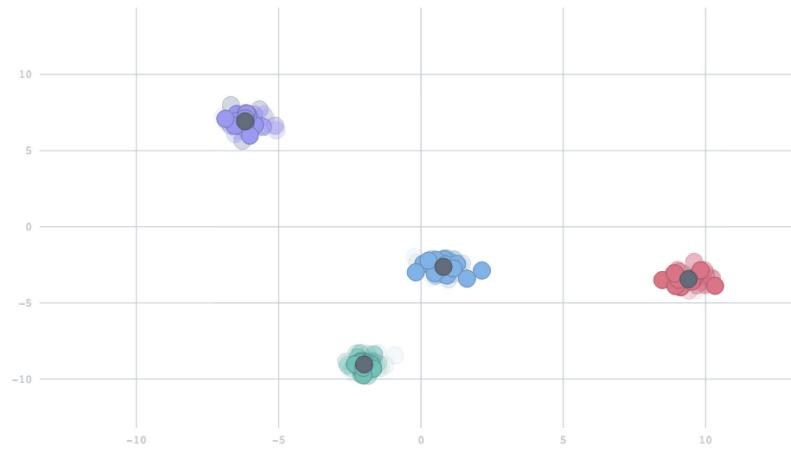
Decay factor $\alpha \in [0, 1]$

- $\alpha \in [0, 1]$ is the decay factor. The smaller α , the bigger the decay.
- $c_i^{t+1} = \frac{\alpha \cdot w_i^t \cdot c_i^t + bc_i^t}{\alpha \cdot w_i^t + bw_i^t}$ is the new centroid after adding B_i^t .
Note that c_i^t and bc_i^t are vectors.
- $w_i^{t+1} = \alpha \cdot w_i^t + bw_i^t$ is the new centroid after adding B_i^t .
- Remove cluster i if $w_i^{t+1} < \delta$ and split the largest cluster into two.

The effect of the decay factor

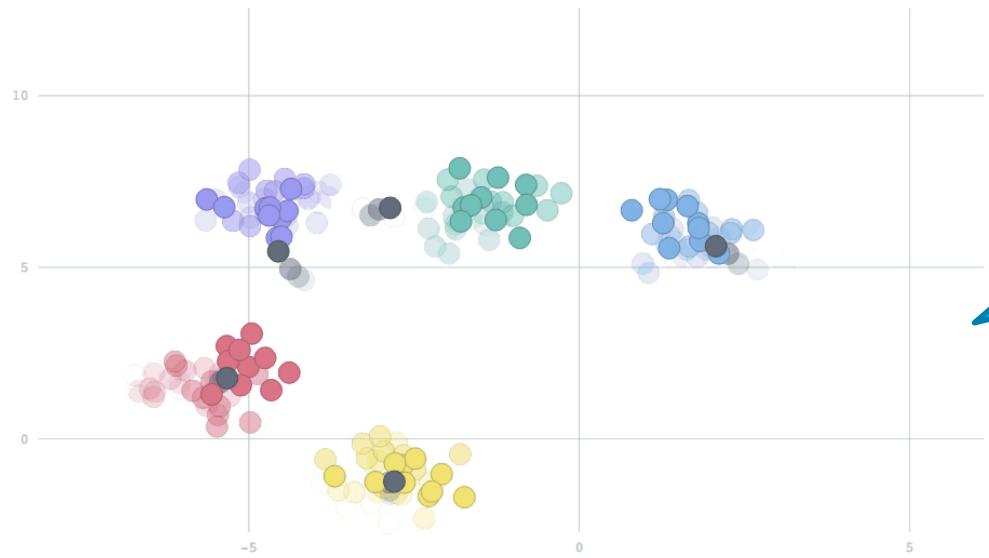


Higher value of α
(slower adjustment)



Lower value of α
(faster adjustment)

The effect of the decay factor

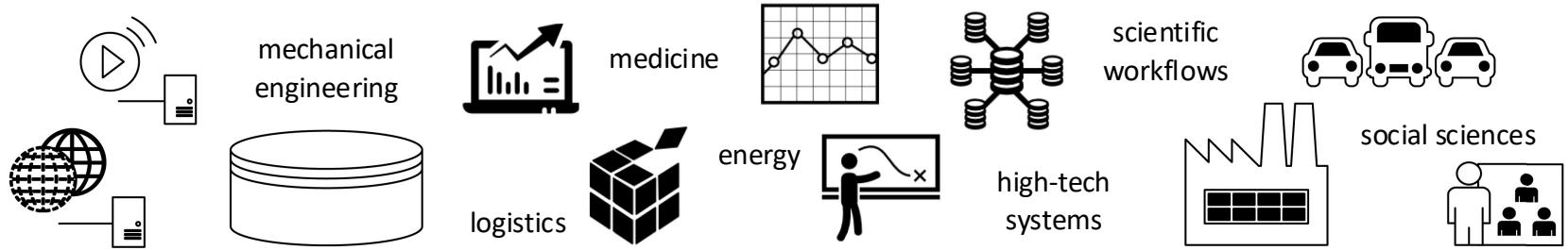


Concept drift

- Different types of drift:
 - Sudden drifts (e.g., due to Brexit).
 - Gradual drifts (e.g., due to improving economy).
 - Seasonal drifts (e.g., due to weather).
 - Intermittent drifts (e.g., due to system failures).
- What should the model describe?
 - Recent behaviors?
 - Long-term behaviors?

Conclusion





infrastructure

“volume and velocity”

analysis

“extracting knowledge”

effect

“people, organizations, society”

- big data infrastructures
- distributed systems
- data engineering
- programming
- security
- ...
- statistics
- data/process mining
- machine learning
- artificial intelligence
- visualization
- ...
- ethics & privacy
- IT law
- operations management
- business models
- entrepreneurship
- ...

Summary

- **Big Data: Relevance and Applications**
- **Characteristics of Big Data: The 4 V's**
- **Big Data Infrastructures**
- **MapReduce**
- **Streaming**
- **Streaming k-means**

Relevant Literature

- Viktor Mazer-Schönberger. “Big Data: A Revolution That Will Transform How We Live, Work, and Think”. 04.03.2014.
- Apache Hadoop: MapReduce Tutorial.
https://hadoop.apache.org/docs/r1.2.1/mapred_tutorial.html
- Streaming 101: The world beyond batch. <https://www.oreilly.com/ideas/the-world-beyond-batch-streaming-101>
- Streaming 102: <https://www.oreilly.com/ideas/the-world-beyond-batch-streaming-102>
- van der Aalst, Wil M.P. “Process Mining in the Large: A Tutorial”. Business Intelligence, pp 33-76.
- van der Aalst, Wil M.P. “Decomposing Petri nets for process mining: A generic approach”. Distributed and Parallel Databases, Volume 31, pp 471-507.
- Nada Elgendz, Ahmed Elragal. “Big Data Analytics: A Literature Review Paper”. Industrial Conference on Data Mining. pp 214-227.
- Tom White. “Hadoop: The Definitive Guide”. O'Reillz Media, April 2015.



Instruction on Friday

- It is essential that you have Hadoop installed before Friday. It does not make sense to try an install it during the instruction. See posting Yaguang Sun.
- Tomorrow (Thursday lecture), Yaguang Sun will demo this, summarize the installation, and help people that did not succeed.
- During the instruction you will learn how to implement a MapReduce algorithm and run your code on Hadoop.

#	Lecture	date	day
	Lecture 1	Introduction	10/10/2018 Wednesday
	Lecture 2	Crash Course in Python	
Instruction 1	Lecture 3	Basic data visualisation	Lecture 22 Big data (1/2)
	Lecture 4	Decision trees	Lecture 23 Big data (2/2)
Instruction 2	Decision trees	Instruction 11	Big data
	Lecture 5	Regression	Lecture 24 Closing
	Lecture 6	Support vector	backup
Instruction 3	Regression and	Instruction 12	Example exam questions
	Lecture 7	Neural network	backup
Instruction 4	Neural network	Instruction 13	extra Question hour
	Lecture 8	Neural network	Association rules
	Lecture 9	Evaluation of systems	Sequence mining
Instruction 5	Neural network	Clustering, frequent items sets, association rules	Clustering, frequent items sets, association rules
	Lecture 10	Clustering	Process mining (unsupervised)
	Lecture 11	Frequent items	Process mining (supervised)
	Lecture 12	Association rules	Process mining and sequence mining
Instruction 6	Lecture 13	Sequence mining	Text mining (1/2)
	Lecture 14	Clustering, frequent items sets, association rules	Text mining and process mining
	Lecture 15	Process mining (unsupervised)	Text mining (2/2)
Instruction 7	Lecture 16	Process mining (supervised)	Data preprocessing, data quality, binning, etc.
	Lecture 17	Text mining and sequence mining	Visual analytics & information visualization
Instruction 8	Lecture 18	Text mining (1/2)	backup
	Lecture 19	Text mining and process mining	Text
	Lecture 20	Text mining (2/2)	Res
Instruction 9	Lecture 21	Data preprocessing, data quality, binning, etc.	Res
	Lecture 22	Visual analytics & information visualization	Res
Instruction 10	Lecture 23	backup	Res
	Lecture 24	Big data	Res
Instruction 11	Big data	Big data	Res
	Closing	Question hour	Wednesday
	extra	Question hour	16/01/2019 Wednesday
	backup		17/01/2019 Thursday
	extra		18/01/2019 Friday
	backup		23/01/2019 Wednesday
	extra		24/01/2019 Thursday
	backup		25/01/2018 Friday
	extra		30/01/2019 Wednesday
	backup		31/01/2019 Thursday
	extra	Question hour	01/02/2019 Friday
	backup		
Instruction 12	Example exam questions		
	extra		
	backup		
	extra		

Solutions for two trial exams are given.

You can ask questions to the instructors and get two trial exams

#	Lecture	date	day
	Lecture 1	Introduction	10/10/2018 Wednesday
	Lecture 2	Crash Course in Python	
Instruction 1	Lecture 3	Basic data visualisation	Lecture 22 Big data (1/2)
	Lecture 4	Decision trees	Lecture 23 Big data (2/2)
Instruction 2	Decision trees	Instruction 11	Big data
	Lecture 5	Regression	Lecture 24 Closing
	Lecture 6	Support vector	backup
Instruction 3	Regression and	Instruction 12	Example exam questions
	Lecture 7	Neural network	backup
Instruction 4	Neural network	Instruction 13	extra
	Lecture 8	Neural network	Question hour
	Lecture 9	Evaluation of systems	
Instruction 5	Neural network	Backup	
	Lecture 10	Clustering	
	Lecture 11	Frequent items	
	Lecture 12	Association rules	21/11/2018 Tuesday
	Lecture 13	Sequence mining	22/11/2018 Wednesday
Instruction 6	Clustering, frequent items sets, association rules		23/11/2018 Thursday
	Lecture 14	Process mining (unsupervised)	28/11/2018 Tuesday
	Lecture 15	Process mining (supervised)	29/11/2018 Wednesday
Instruction 7	Process mining and sequence mining		30/11/2018 Thursday
	Lecture 16	Text mining (1/2)	05/12/2018 Tuesday
Instruction 8	Text mining and process mining		06/12/2018 Wednesday
	Lecture 17	Text mining (2/2)	12/12/2018 Tuesday
	Lecture 18	Data preprocessing, data quality, binning, etc.	
	Lecture 19	Visual analytics & information visualization	
	Backup		
Instruction 9	Text mining, preprocessing and visualization		
	Lecture 20	Responsible data science (1/2)	
	Lecture 21	Responsible data science (2/2)	
Instruction 10	Responsible data science		
	Lecture 22	Big data (1/2)	
	Lecture 23	Big data (2/2)	
Instruction 11	Big data		
	Lecture 24	Closing	23/01/2019 Wednesday
	Backup		24/01/2019 Thursday
Instruction 12	Example exam questions		25/01/2018 Friday
	Backup		30/01/2019 Wednesday
	Backup		31/01/2019 Thursday
	extra	Question hour	01/02/2019 Friday

Prepare questions and
send them before
Monday 21/1/2019 12.00.

Yaguang Sun provides support on 17/1/2019 for those that could not install Hadoop. Instruction assumes Hadoop is up and running (no regular lecture)