

## Introduction to Big Data

### Corso di Sistemi e Architetture per Big Data

A.A. 2021/22

Valeria Cardellini

Laurea Magistrale in Ingegneria Informatica

## Why Big Data?

How much data is created every single minute of the day?

Global Internet population in January 2022: 4.95 billion people (62.5% of world population)



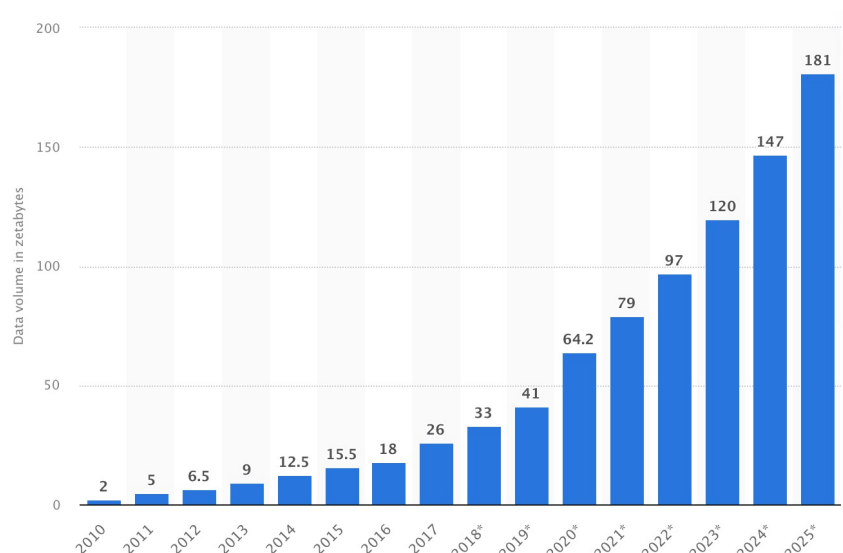
Source: <https://www.domo.com/learn/infographic/data-never-sleeps-9>

# How much data?

- Big data volume: from Terabytes to Zettabytes
  - How big is a Zettabyte? <http://bit.ly/2G15uVI>
  - $1 \text{ ZB} = 2^{70} \text{ B} = 2^{40} \text{ GB} \approx 10^{21} \text{ B}$ 
    - Remember that  $2^{10} = 1024 \approx 10^3$
- 79 Zettabytes of data generated by 2021
  - 79 Zettabytes ( $79 \times 2^{70} \approx 79 \times 10^{21}$ ) ...
  - $\approx 79,000$  Exabytes ( $79,000 \times 10^{18}$ ) ...
  - $\approx 79,000,000$  Petabytes ( $79,000,000 \times 10^{15}$ ) ...
  - $\approx 79,000,000,000$  Terabytes ( $79,000,000,000 \times 10^{12}$ ) ...
  - $\approx 79,000,000,000,000$  Gigabytes ( $79,000,000,000,000 \times 10^9$ ) ...
  - $\approx 79,000,000,000,000,000,000,000$  bytes!

# How much data?

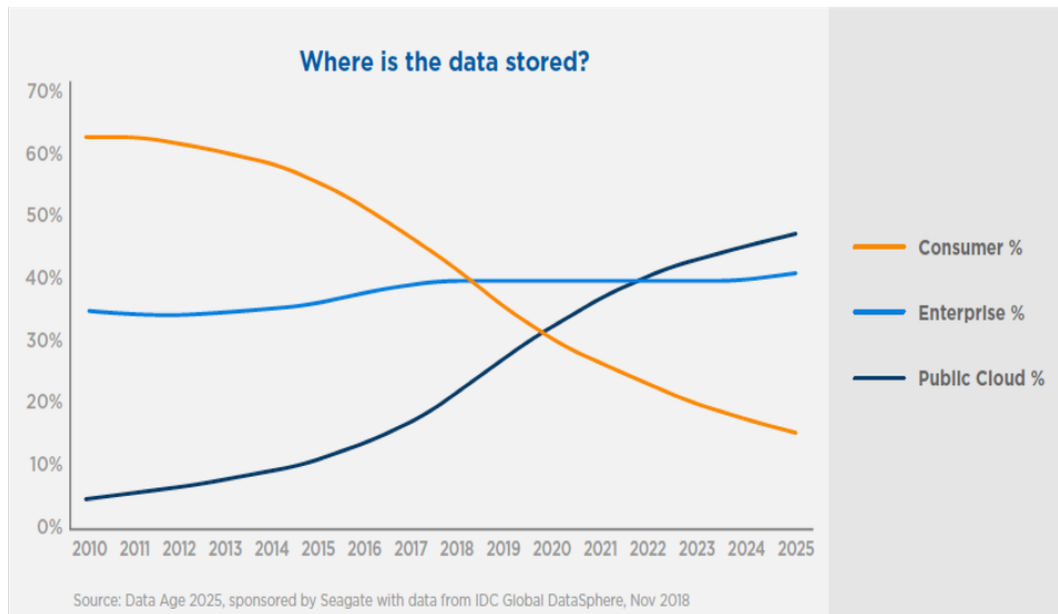
- Recent explosion in data volume
  - In 2013: 90% of all the data in the world was generated over the last two years
  - 30x growth from 2010 to 2020



# Where is data stored?

---

- Data will be increasingly stored in the public cloud



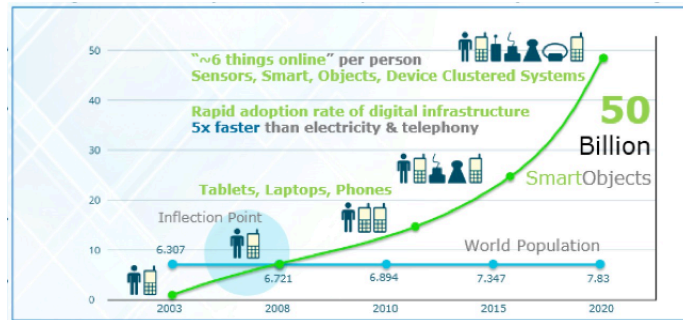
## Big data statistics and economic impact

---

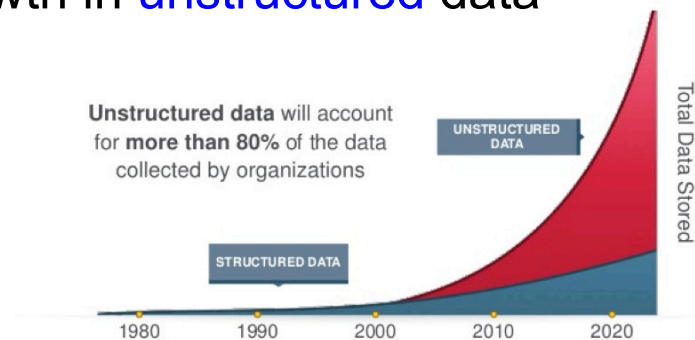
- In 2020, every person generated 1.7 MB in just a second
- Internet users generate about 2.5 EB of data each day
- Google, Facebook, Microsoft, and Amazon store at least 1,200 PB of information
- Big data and business analytics market is set to reach \$274 billion by 2022 (IDC source)
- 91% of organizations are investing in Big Data and AI
- Using Big Data, Netflix saves \$1 billion per year on customer retention

# Big data driving factors

- Big Data is growing fast
  - Smartphones
  - Social networks
  - Internet of Things (IoT)

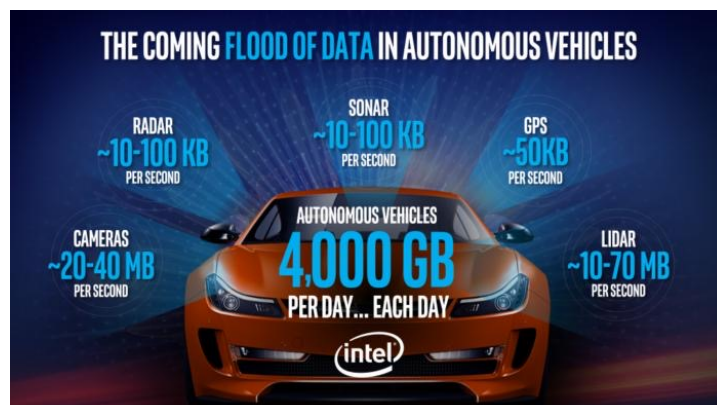


- Exponential growth in unstructured data



## How Big? IoT impact

- IoT is everywhere and largely contributes to increase Big Data challenges
  - Proliferation of data sources: by 2021 over 35 billion IoT devices installed worldwide
- Example: self-driving cars
  - Just one autonomous car will use 4 TB of data/day

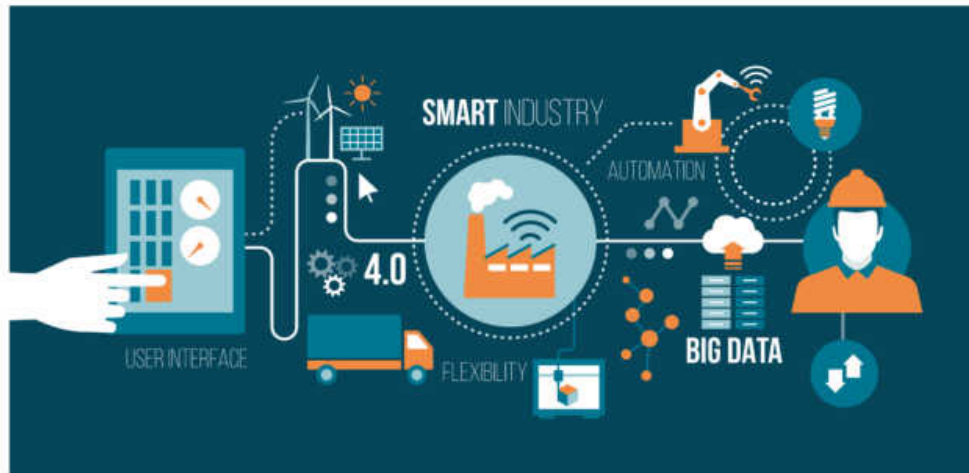




# IoT impact: Industrial IoT

---

- Industrial Internet of Things (IIoT) is a network of physical objects, systems, platforms and applications that contain embedded technology to communicate and share intelligence with each other, the external environment and with people



Valeria Cardellini - SABD 2021/22

8

## Big Data definitions

---

### Different definitions

- “Big data refers to data sets whose size is **beyond** the ability of typical database software tools to capture, store, manage and analyze.” *The McKinsey Global Institute, 2012*
- “Big data is a field that treats ways to analyze, systematically extract information from, or otherwise deal with data sets that are **too large or complex** to be dealt with by traditional data-processing application software.” *Wikipedia, 2020*
- “Big data is mostly about taking numbers and using those numbers to **make predictions about the future**. The bigger the data set you have, the more accurate the predictions about the future will be.” *Anthony Goldbloom, Kaggle’s founder*

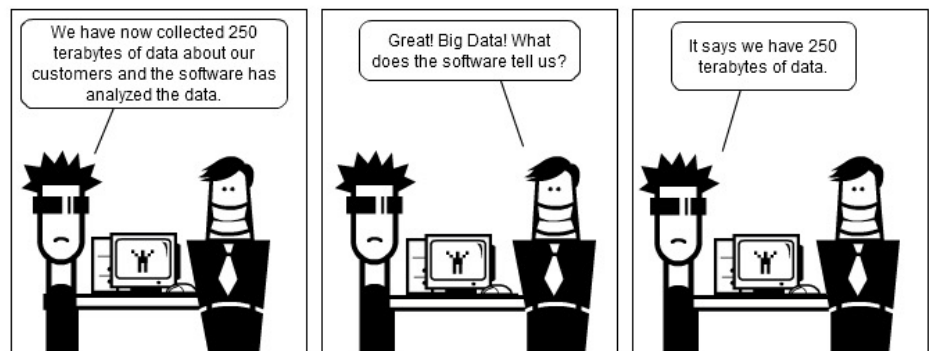
Valeria Cardellini - SABD 2021/22

9

## ... so, what is Big Data?

---

- “Big Data” is similar to “small data”, but bigger
- But bigger data requires different approaches: **scale changes everything!**
  - New methodologies, tools, architectures
- ...with an aim to solve new problems or old problems in a better way



Valeria Cardellini - SABD 2021/22

10

## Gartner's Big data definition

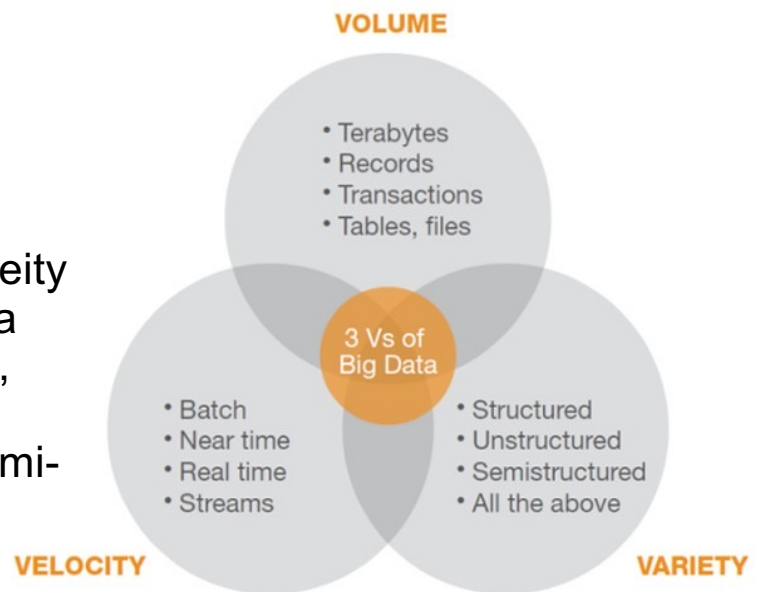
---

- The most-frequently used and perhaps, somewhat abused definition (revised version by Gartner, 2012)  
*Big data is **high volume**, **high velocity**, and/or **high variety** information assets that require **new forms of processing** to enable enhanced decision making, insight discovery and process optimization.*

# 3V model for Big Data

---

1. **Volume**: data size challenging to store and process (how to index, retrieve)
  2. **Variety**: data heterogeneity because of different data types (text, audio, video, record) and degree of structure (structured, semi-structured, unstructured data)
  3. **Velocity**: data generation rate and analysis rate
- Defined in 2001 by D. Laney



## The extended (3+n)V model

---

4. **Value**: Big data can generate huge competitive advantages
  - “Big data technologies describe a new generation of technologies and architectures, designed to economically extract value from very large volumes of a wide variety of data, by enabling high-velocity capture, discovery, and/or analysis.” (IDC, 2011)
  - “The bigger the data set you have, the more accurate the predictions about the future will be” (A. Goldbloom)
5. **Veracity**: uncertainty of accuracy and authenticity of data
6. **Variability**: data flows can be highly inconsistent with periodic peaks
7. **Visualization**

# Big Data visualization

---

- Presentation of data in a pictorial and graphical format
- Why? Our brain processes images 60,000x faster than text
- Some examples



Valeria Cardellini - SABD 2021/22

14

# Big Data visualization

---

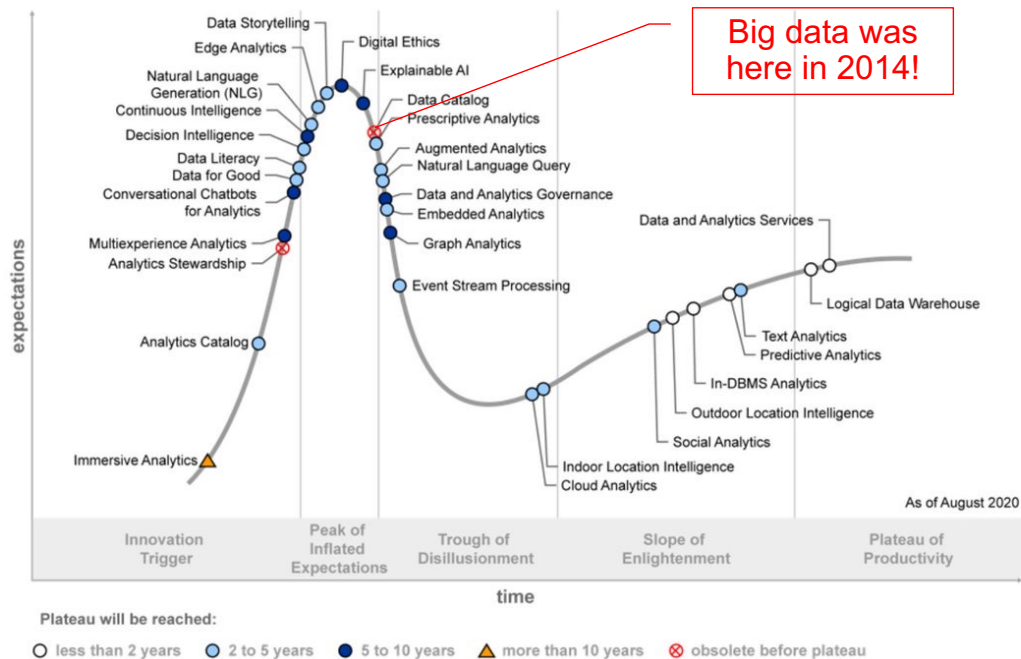
- Some examples
  - Flight patterns in US <http://bit.ly/1rEKMIR>
  - Pollution map <https://www.pollution.org>
  - Ocean surface currents  
<https://www.nasa.gov/topics/earth/features/perpetual-ocean.html>
  - World tweet map  
<https://www.omnisci.com/demos/tweetmap>

Valeria Cardellini - SABD 2021/22

15

# Gartner's 2020 hype cycle for analytics and business intelligence

## Hype Cycle for Analytics and Business Intelligence, 2020



Valeria Cardellini - SABD 2021/22

16

## Why now?

- Because we have data
  - Data born already in digital form
  - 40% of data growth per year
- Because we can
  - 400\$ for a drive in which to store all the music of the world
  - More than 40 years of Moore's law: we have large computing resources

Valeria Cardellini - SABD 2021/22

17

## Examples of Big Data applications in very diverse sectors

---

- Customer analytics in retail industry
  - E.g., to increase customer retention and loyalty
- Predictive maintenance for Industry 4.0
  - E.g., detecting anomalous machine states to reduce maintenance costs
- Crime prevention
  - To analyze crime patterns and trends
- Health care
  - E.g., to diagnose and treat genetic diseases
- Finance
  - To anticipate customer behaviors and create strategies for banks and financial institutions

## Examples of Big Data applications in very diverse sectors

---

- Government sector, e.g. using Open Data  
<https://www.europeandataportal.eu/>
- Education
  - E.g., to improve the learning process, to design a new course
- Space science
  - E.g., astronomical discoveries



# Batch vs. real-time analytics

---

- **Batch analytics**: analysis of set of data collected over a period of time and that has already been stored
  - We will study **batch processing engines**
- **Real-time analytics**: analysis of **high-velocity, continuous** data streams as soon as they are ingested without (or before) storing them
  - Goal: get insights **immediately** (or very rapidly after) data enters the system
  - We will study **stream processing engines**

## Examples of real-time analytics

---

- Grand Challenge at DEBS conferences  
<https://debs.org/grand-challenges/>
  - Over high volume sensor data: analysis of energy consumption measurements (DEBS 2014)
  - Over high volume geospatial data streams: analysis of taxi trips based on a stream of trip reports from New York City (DEBS 2015)
  - Over social network: to identify posts that trigger the most activity and large communities that are involved in a topic (DEBS 2016)
  - Over maritime transportation data: to predict destinations and arrival times of ships (DEBS 2018)
  - Technical analysis of market data: to compute specific trend indicators and detect patterns resembling those used by traders to decide on buying or selling (DEBS 2022)

## ... other example of real-time analytics in very diverse sectors

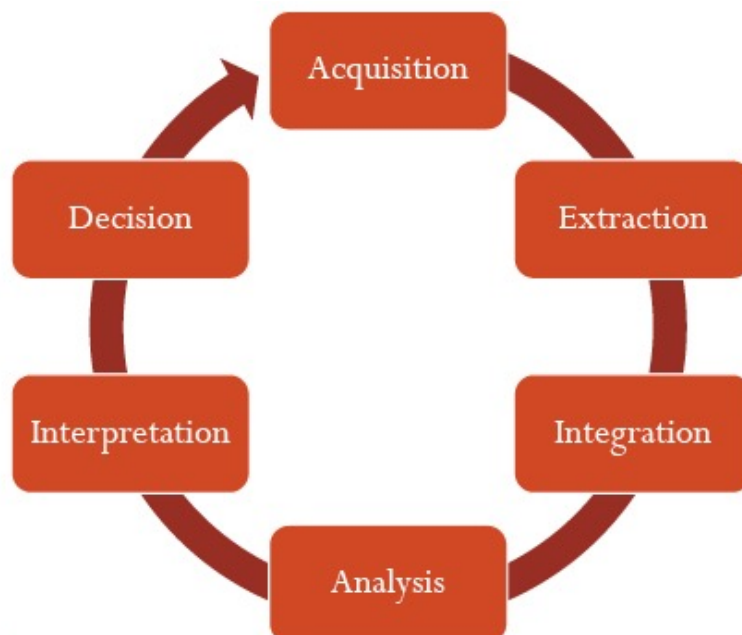
---

- Medicine
  - To track epidemic diseases, to prevent diseases through wearable health care technologies
- Security
  - To detect frauds or DDOS attacks, to recognize behavioral patterns
- Urban traffic management
  - To address traffic congestion and lack of parking, to plan public transportation

## The Big Data process

---

- 6 stages of the Big data analytics lifecycle



# The Big Data process

---

- Acquisition

- Requires:

- Selecting data
    - Filtering data
    - Generating metadata
    - Managing data provenance
      - E.g., GDPR compliance



# The Big Data process

---

- Extraction

- To transform data into a format that can be used by Big data processing frameworks

- Requires:

- Data transformation
    - Data normalization
      - E.g., avoid duplication
    - Data cleaning
      - Detect and correct (or remove) corrupted or inaccurate data
    - Data aggregation
      - E.g., from multiple sources



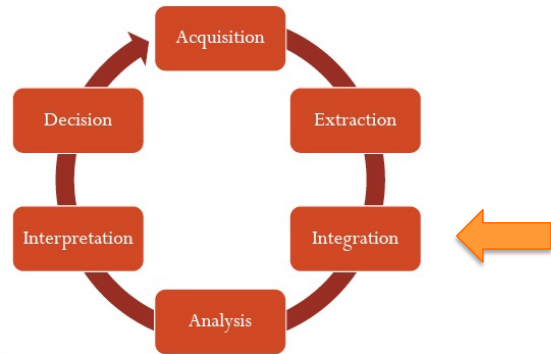
# The Big Data process

---

- Integration

- Requires:

- Standardization
    - Conflict management
    - Reconciliation
    - Mapping definition



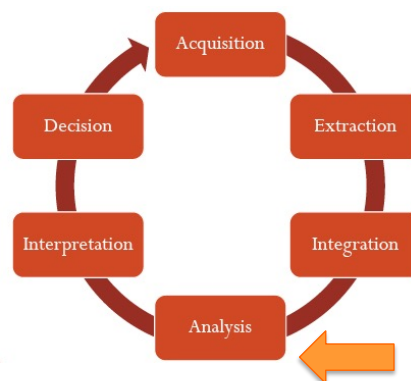
# The Big Data process

---

- Analysis

- Requires:

- Data analytics techniques
      - Exploration
      - Data mining
      - Machine learning
      - Visualization



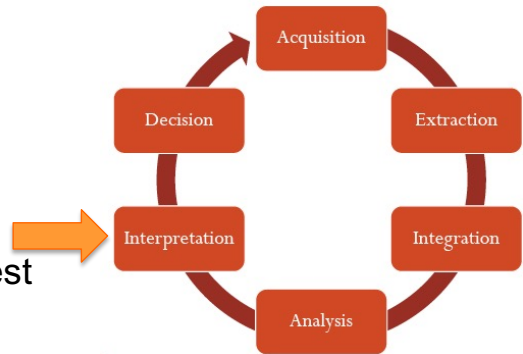
# The Big Data process

---

- Interpretation

- Requires:

- Knowledge of domain
    - Knowledge of data provenance
    - Identification of patterns of interest
    - Flexibility of the process



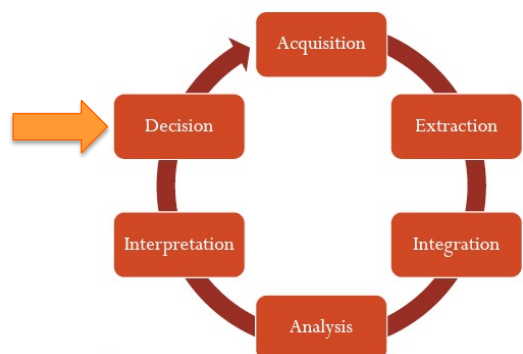
# The Big Data process

---

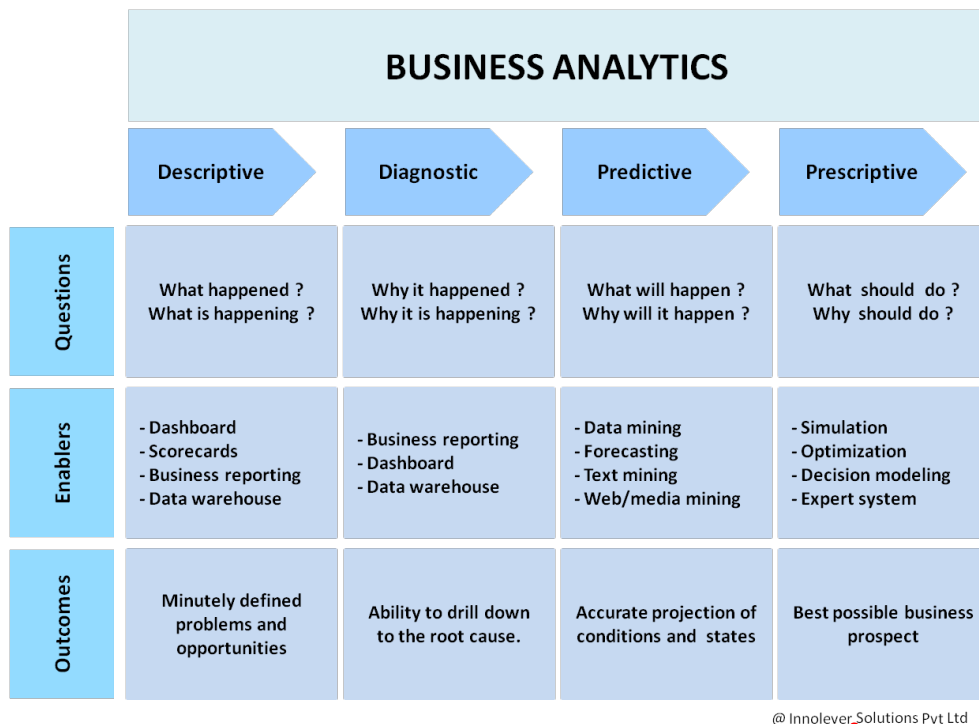
- Decision

- Requires:

- Managerial skills
    - Continuous improvement of the process (loop)



# Some techniques for Big Data analytics



Added value, complexity

Valeria Cardellini - SABD 2021/22

30

# Some techniques for Big Data analytics

- **Data mining**: anomaly detection, association rule mining, classification, clustering, regression, summarization
- **Machine learning**: supervised learning, unsupervised learning, reinforcement learning
- **Crowdsourcing**
  - Outsourcing human-intelligence tasks to a large group of unspecified people via Internet

In this course we focus on systems and architectures for Big Data, not on data analysis techniques

Valeria Cardellini - SABD 2021/22

31

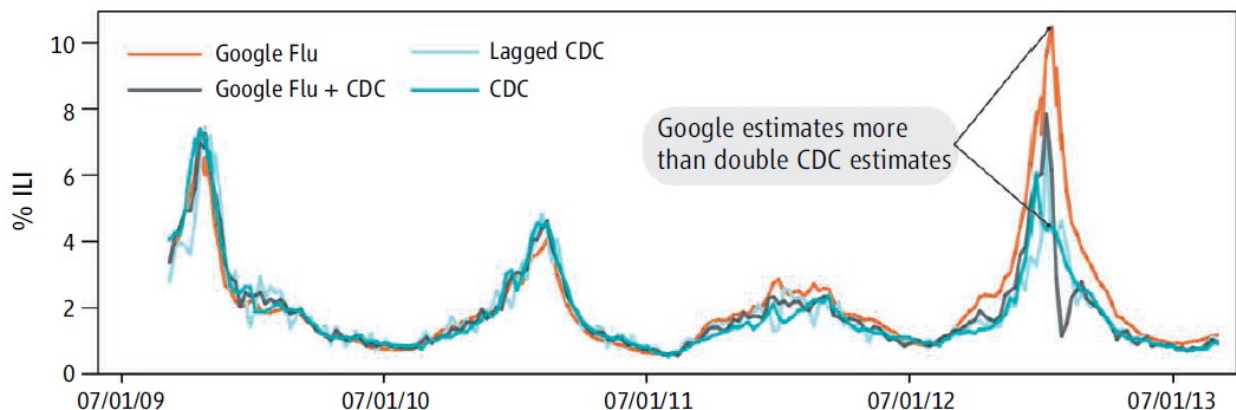


# Risks and challenges of Big Data

- Effectiveness of data analysis
- Performance
  - Efficiency
  - Scalability and elasticity
    - Scale linearly as workloads and data volumes grow
  - Fault tolerance
  - Sustainability
    - Data grows faster than energy on chip
- Heterogeneity
  - Data, processing environment, network latencies, ...
- Flexibility
- Privacy and security
- Costs

## Effectiveness of Big data analysis

- A famous example of inaccurate analysis
- Google Flu Trends' predictions
  - Sometimes very inaccurate: over the interval 2011-2013, when it consistently overestimated flu prevalence and over one interval in the 2012-2013 flu season predicted twice as many doctors' visits as those recorded



Lazer et al., "The Parable of Google Flu: Traps in Big Data Analysis". *Science*. 343 (6176): 1203–1205. doi:10.1126/science.1248506

# Taming performance: distribution and replication

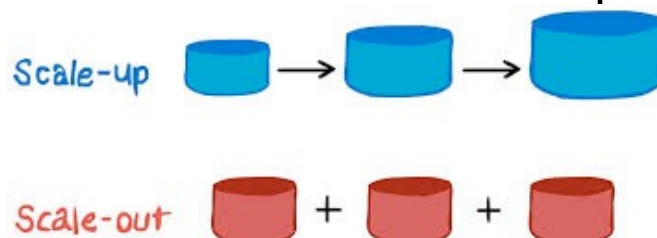
---

- Distributed architecture
  - Common architectural solution for Big Data processing: cluster of commodity hardware resources, also in Cloud
  - **Scale out** (horizontally), not up (vertically)!
  - Challenges: *elasticity* and data processing at the *network edges*
- Distributed processing
  - **Shared-nothing** model
  - New programming paradigms, e.g., functional programming
- Resource replication
  - The well-known solution to achieve fault tolerance
  - Eventual consistency (CAP theorem!)

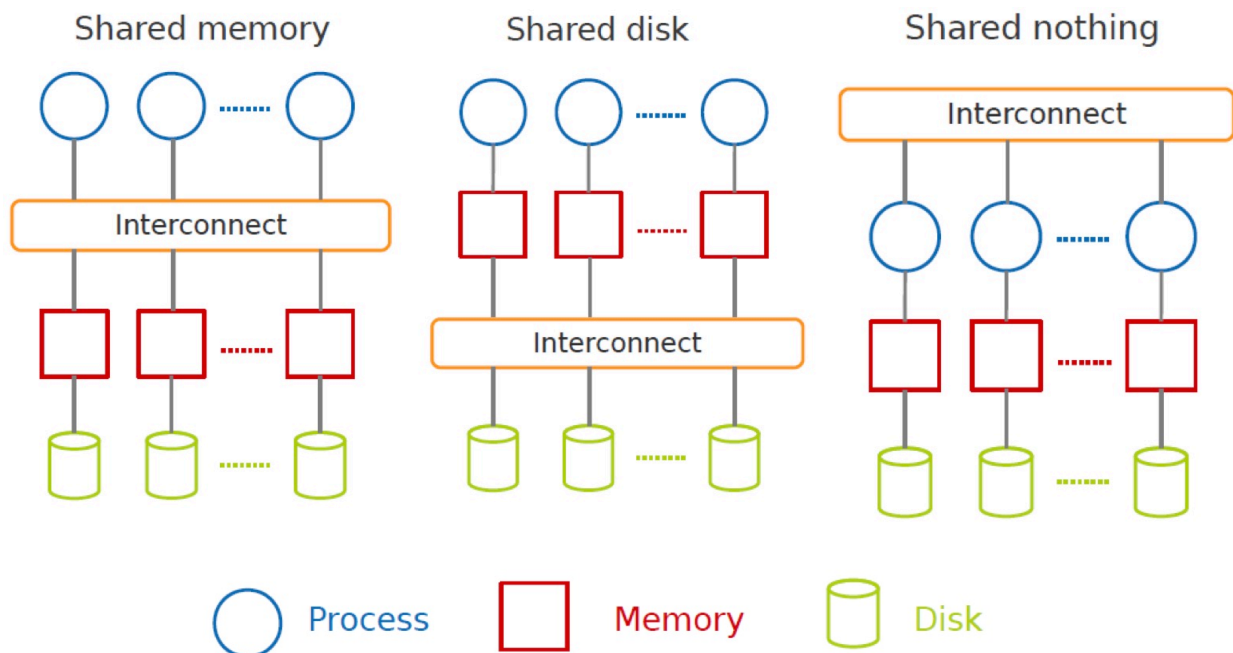
## Scaling out vs. scaling up

---

- Two different ways of addressing the need for more processor capacity, memory and other resources
- **Scaling up** (or vertical scalability) refers to purchasing and installing a more powerful server
  - E.g., with more processing capacity and RAM
- **Scaling out** (or horizontal scalability) means adding other lower-performance servers to collectively do the work of a much more powerful one



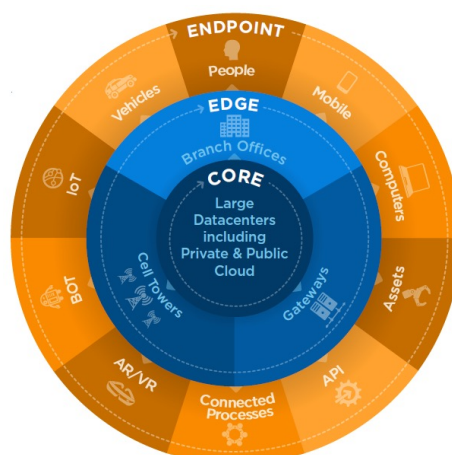
# Shared nothing vs. other parallel architectures



D. DeWitt and J. Gray, "Parallel database systems: the future of high performance database systems", *ACM Communications*, 1992

## Big Data architectures

- Ingest data
- Process data
- Analyze data
- Store data
- Where?



## Where to process Big Data

---

- The traditional way: using a **cluster of servers** on premises
  - Compute nodes are stored on racks
    - 8-64 compute nodes on a rack
  - There can be many racks of compute nodes
    - The nodes on a single rack are connected by a network, typically gigabit Ethernet
    - Racks are connected by another level of network or a switch
    - The bandwidth of intra-rack communication is usually much greater than that of inter-rack communication
- Cons:
  - Need to manage hardware infrastructure and processing platforms (acquire, install, configure, ...)

## Where to process Big Data

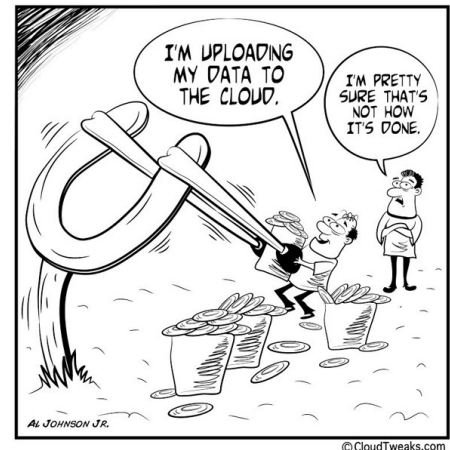
---

- The **Cloud** way: using a Cloud analytics service
- Some examples
  - Amazon EMR and Google Dataproc: Hadoop and Spark clusters (plus high-level frameworks) in the Cloud
- Pros:
  - Gain Cloud scalability and elasticity
  - Do not need to manage and provision the infrastructure and the platform

# Where to process Big Data

---

- But Cloud data centers are located in the network core
- Main challenges:
  - Move data to the Cloud
    - Latency is not zero (because of speed of light)!
    - Minor issue: network bandwidth
  - Data security and privacy



Valeria Cardellini - SABD 2021/22

40

# Where to process Big Data

---

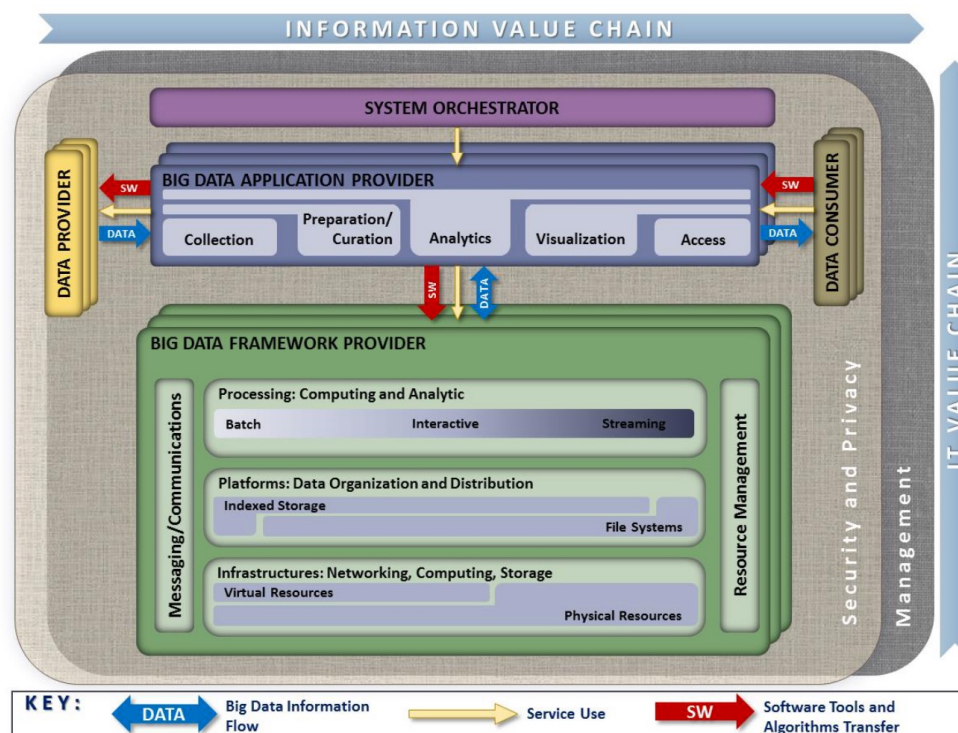
- The new scenario: **edge/Fog computing**
  - “The cloud close to the ground”: many micro data centers located at the network edge
  - Move data processing close to data producers and data consumers



Valeria Cardellini - SABD 2021/22

41

# NIST Big Data reference architecture



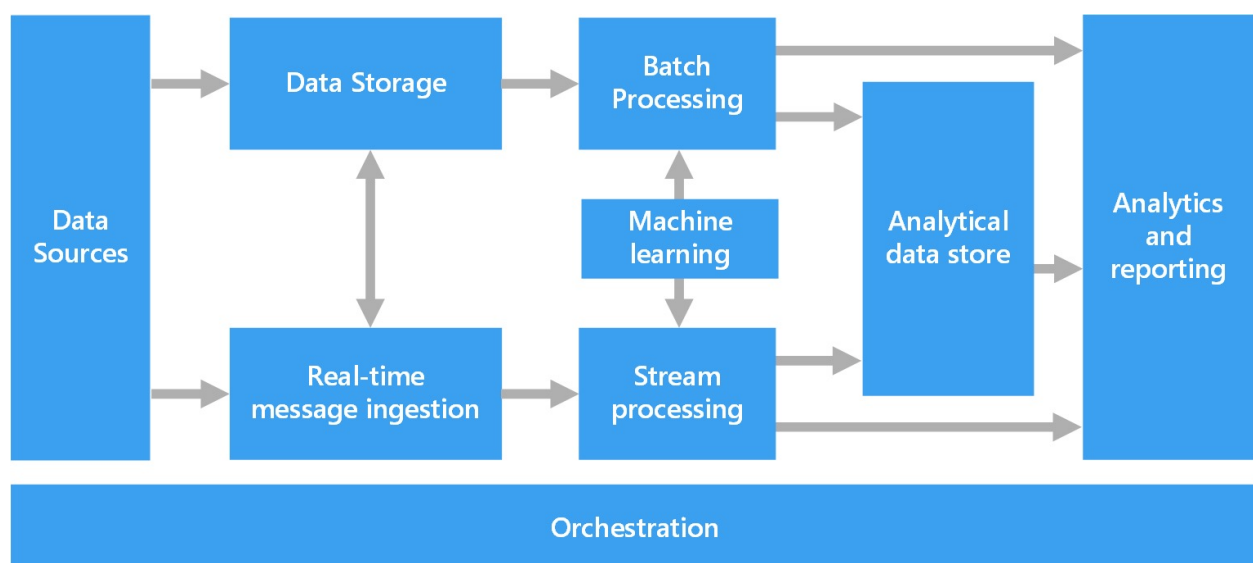
<https://doi.org/10.6028/NIST.SP.1500-6r2>

Valeria Cardellini - SABD 2021/22

42

## Components of a big data architecture

- Lambda architecture: both batch and stream processing

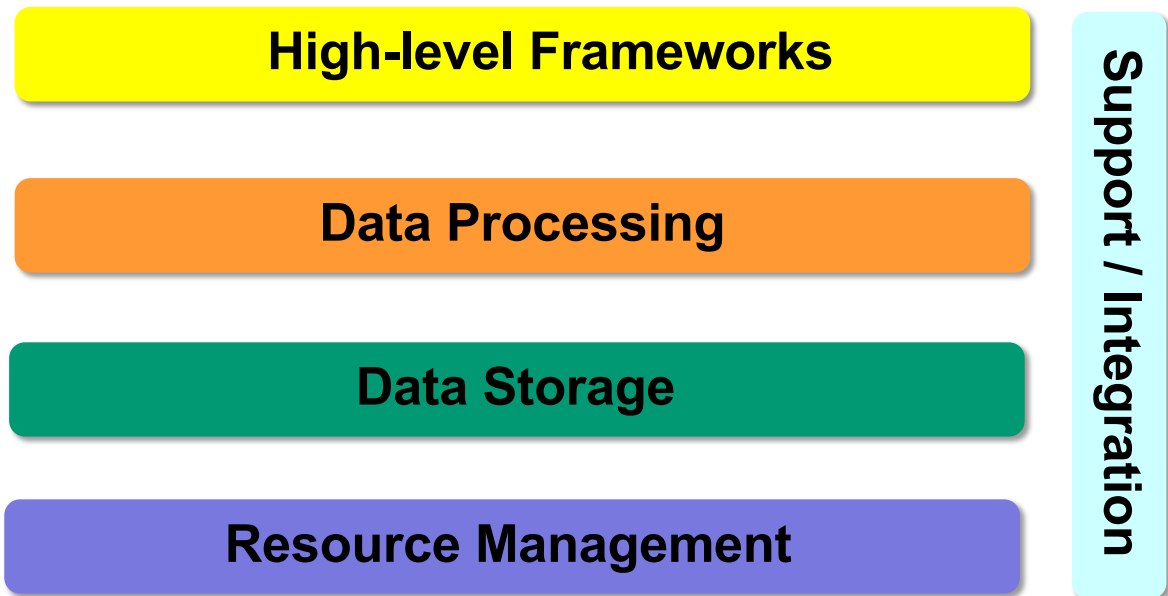


Valeria Cardellini - SABD 2021/22

43



# Our Big Data stack



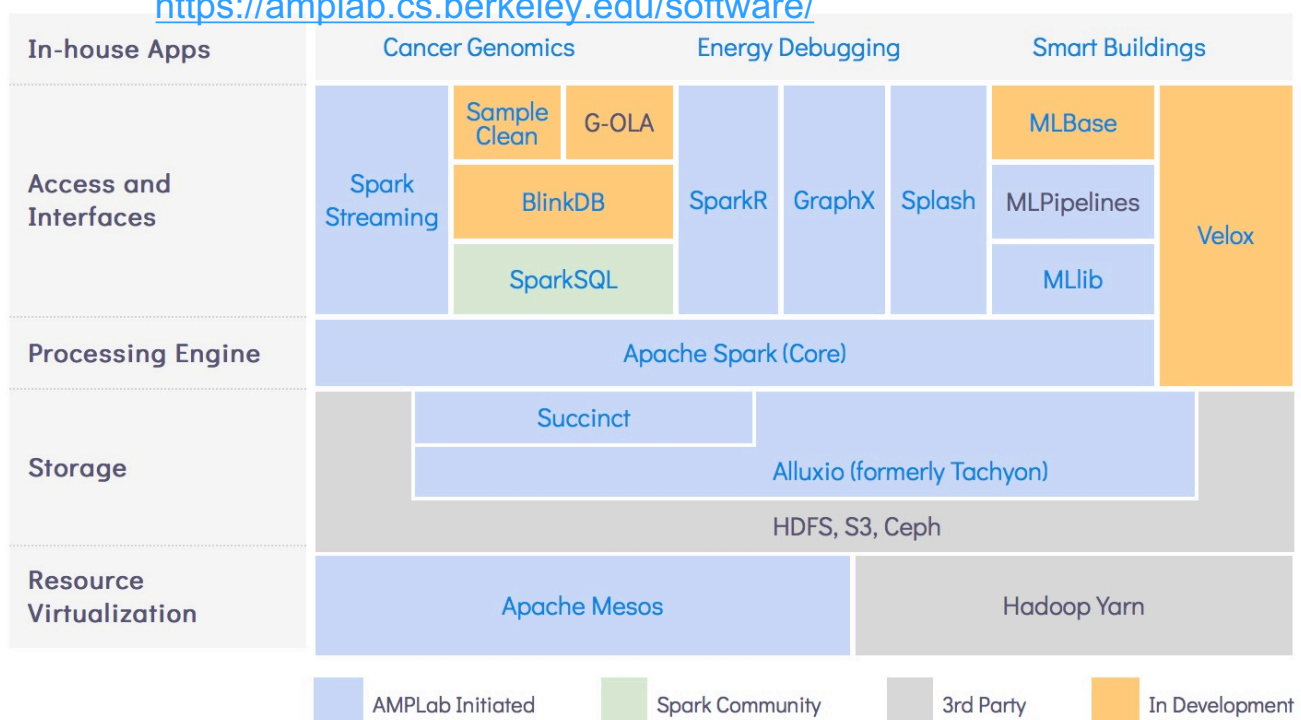
Valeria Cardellini - SABD 2021/22

44

## Example of Big Data stack: BDAS

- BDAS: the Berkeley Data Analytics Stack

<https://amplab.cs.berkeley.edu/software/>

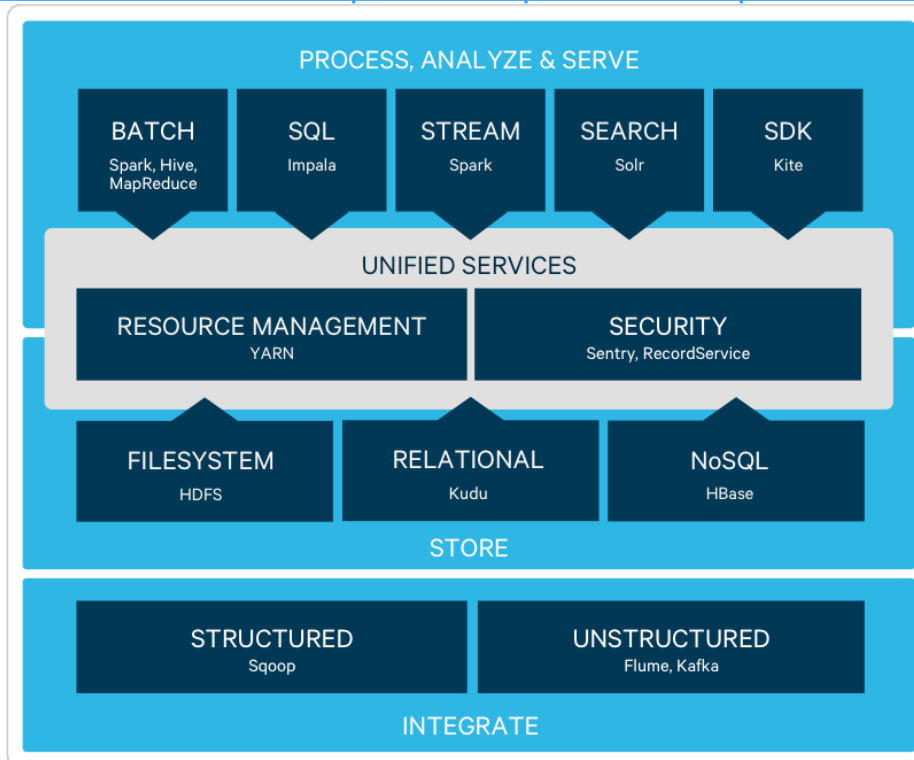


Valeria Cardellini - SABD 2021/22

45

# Example of Big Data stack: Cloudera

<https://www.cloudera.com/products/open-source/apache-hadoop.html>



Valeria Cardellini - SABD 2021/22

46

## Data lake

- Method of **storing data** within a system or repository, in its **natural format**, that facilitates the collocation of data in various schemata and structural forms, usually object blobs or files
- Designed for **quickly changing data**

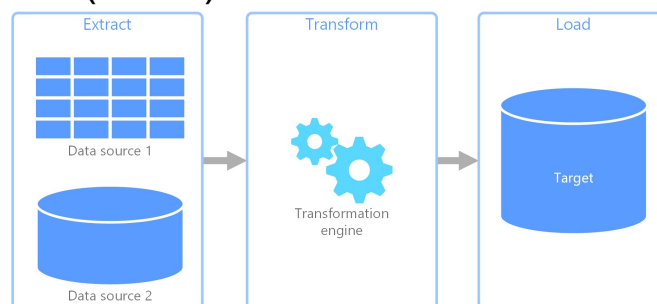


Valeria Cardellini - SABD 2021/22

47

## Paradigm shift in the data pipeline

- From the traditional way: Extract, Transform, and Load (**ETL**)
  - Also known as Structure, Ingest and Analyze
  - *Extract* data from multiple sources
  - *Transform* data into the proper format (or structure) for the purposes of storing
  - *Load* data into the target system, i.e., database or **data warehouse** (DWH)



## Paradigm shift in the data pipeline

- ... to the new way: Extract, Load, and Transform (**ELT**)
  - Also known as Ingest, Analyze, and Structure
  - *Extract* data from multiple sources
  - *Load* data into a **data lake**, where data is held in original format
  - *Transform* data using the processing capabilities of target system
- Advantages:
  - No need for separate transformation engine
  - Data transformation and loading happen in parallel
  - More effective when speed is critical
  - Works well when target system is powerful enough to transform data efficiently

