# Chapter 1: What is big data?
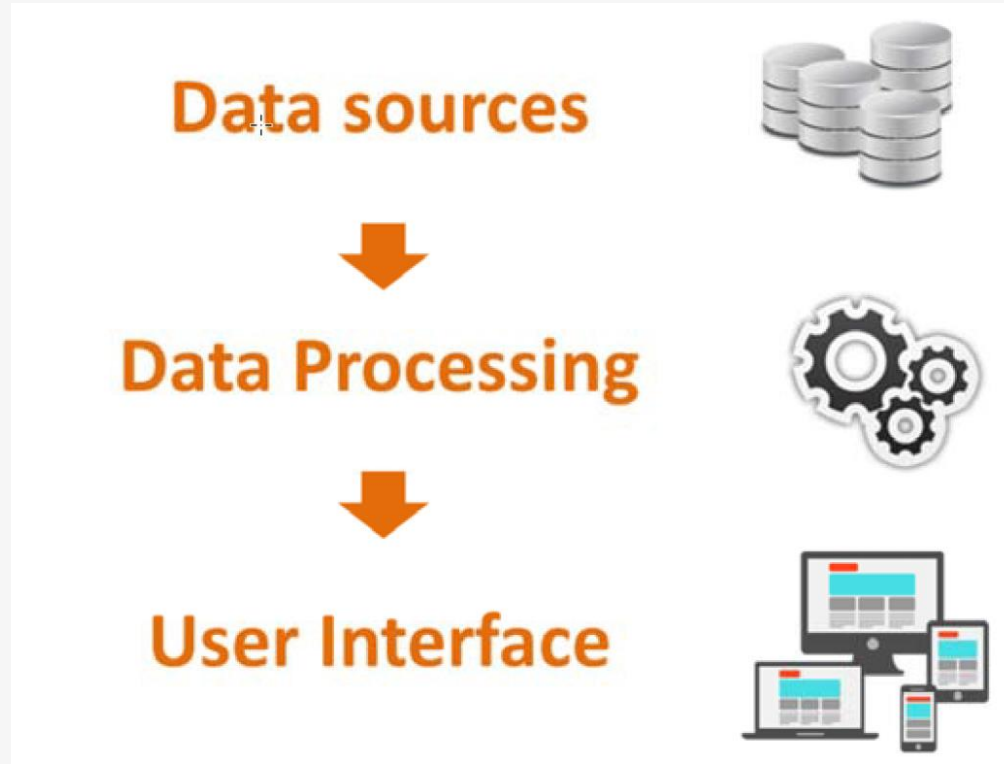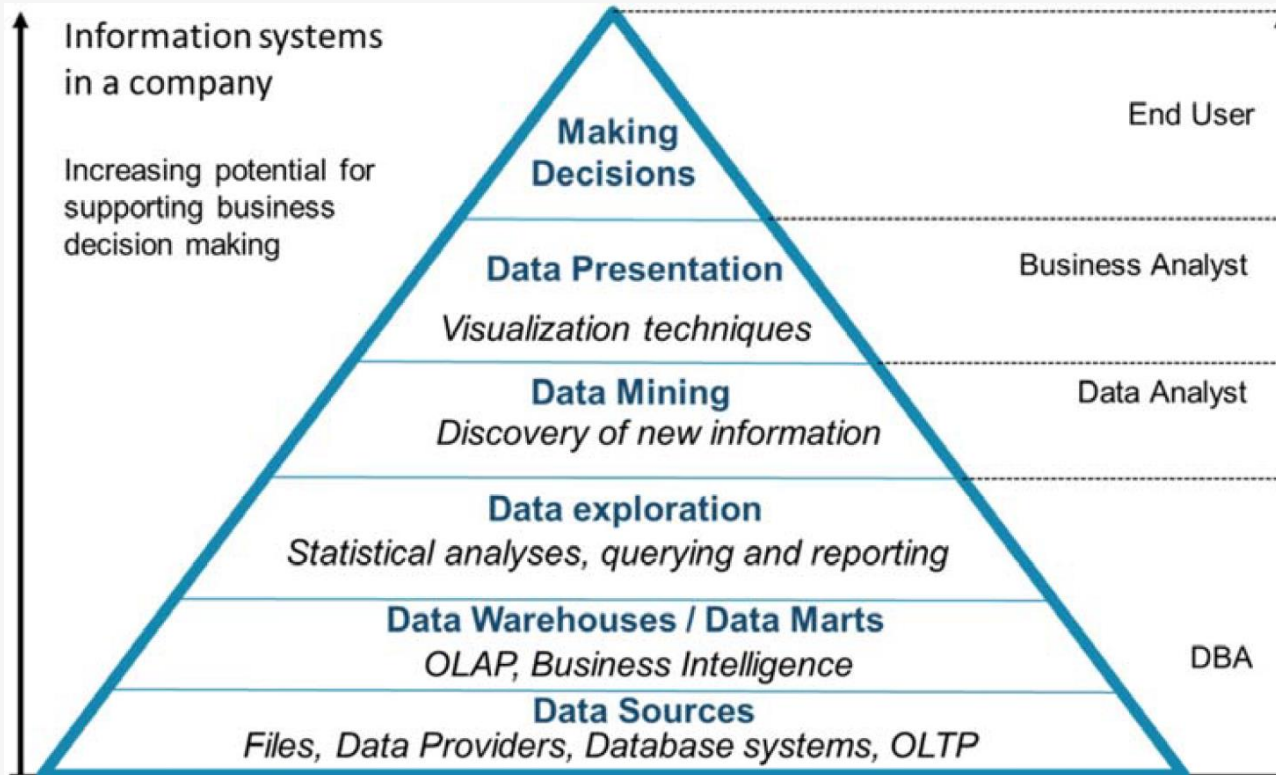
# Contents

1. Information Management
2. What are big data?
3. The origin of big data
4. The 4 V's of big data
5. Some other V's
6. Bottlenecks
7. Examples of how big data can be used

**HO GENT**

# 1. Information Management

# Decision Support Systems

Information systems in a company

Increasing potential for supporting business decision making

**Making Decisions** — End User

**Data Presentation**
*Visualization techniques* — Business Analyst

**Data Mining**
*Discovery of new information* — Data Analyst

**Data exploration**
*Statistical analyses, querying and reporting*

**Data Warehouses / Data Marts**
*OLAP, Business Intelligence* — DBA

**Data Sources**
*Files, Data Providers, Database systems, OLTP*
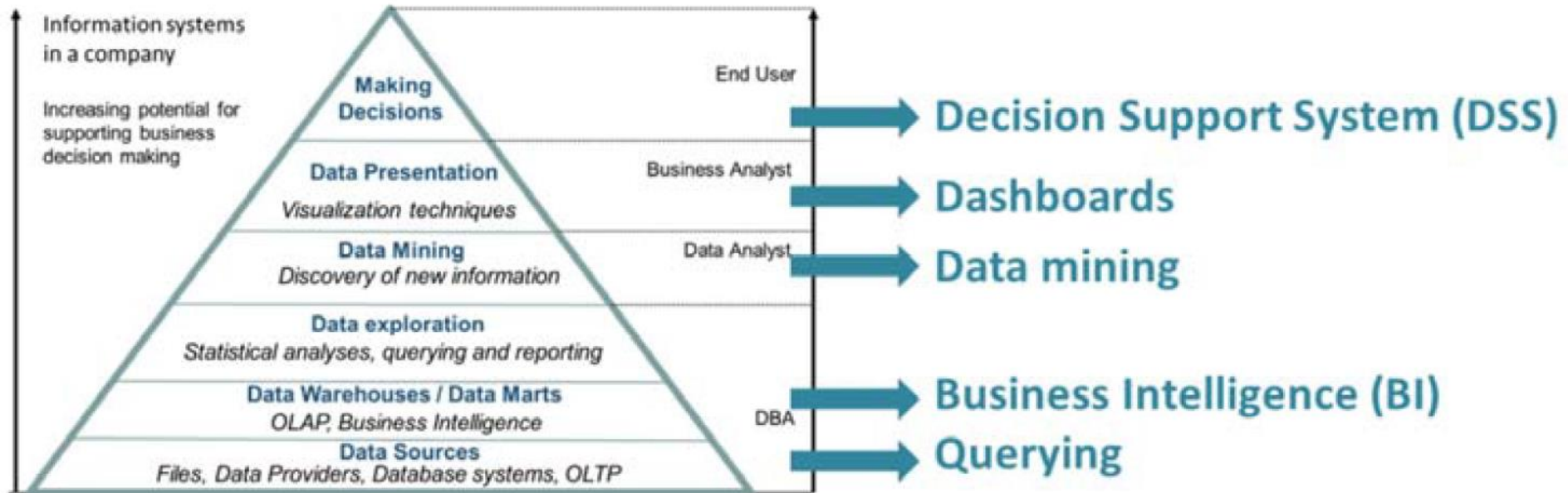
A decision support system (DSS) is a computer program application used to improve a company's decision-making capabilities. It analyzes large amounts of data and presents an organization with the best possible options available.

HO GENT

# Traditional systems



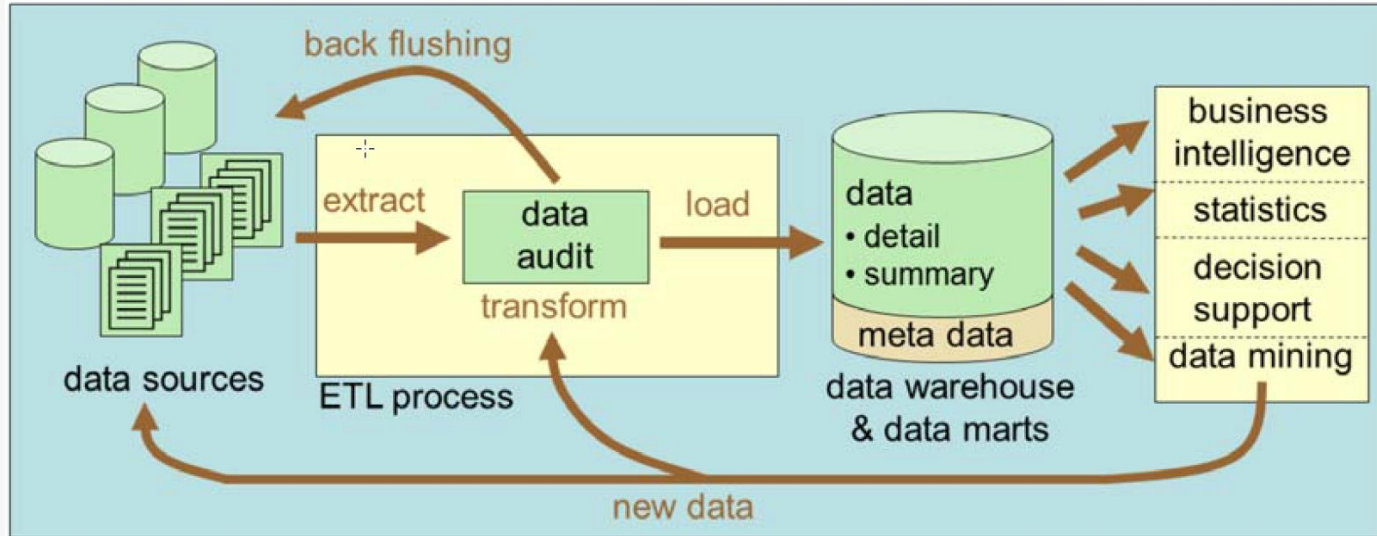OLTP: Online Transaction Systems, day-to-day operations
OLAP: Online Analytical Processing, midterm and longterm decision support
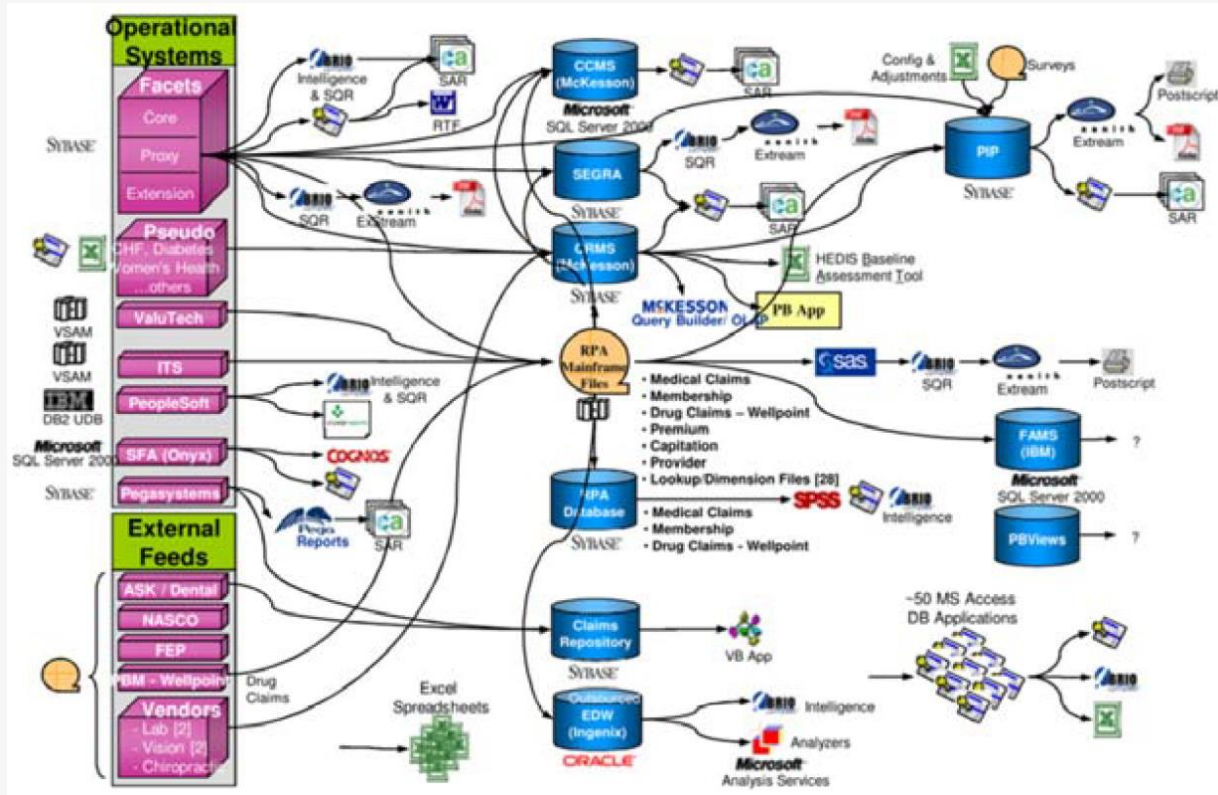
HO
GENT

# What do users want?



Users want a single logic database for their analyses: it looks like all distributed, heterogeneous data is available from one single database.
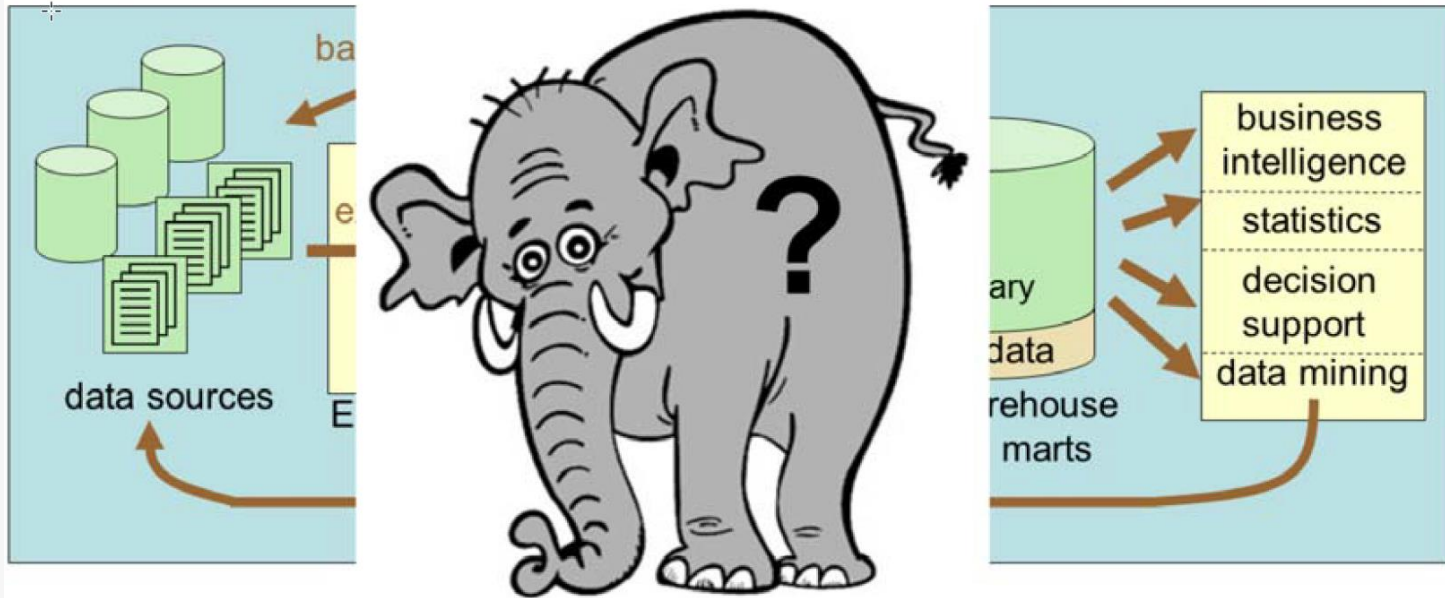
**HO GENT**

# What do we usually have? … in theory



This is the theoretical process of **E**xtracting data from different source systems, **T**ransforming and cleaning it to a suitable format and **L**oading it into a Datawarehouse for easy Business Intelligence.

HO GENT

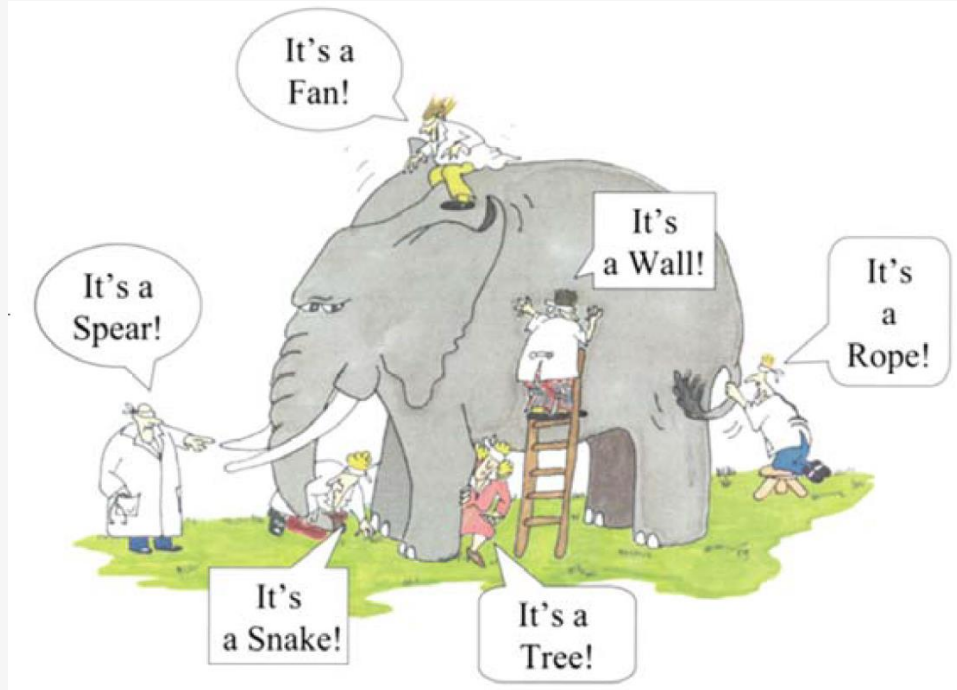# What do we usually have? … in practice

# How does Big Data fit in these pictures?



Furthermore, traditional databases, ETL processes, datawarehouses, data mining techniques etc. are not suited for processing data with "specific" characterics. Such data are called "big" data.

# 2. What are big data?



Many definitions of big data exist, some are leading to confusion and unjustified criticism …
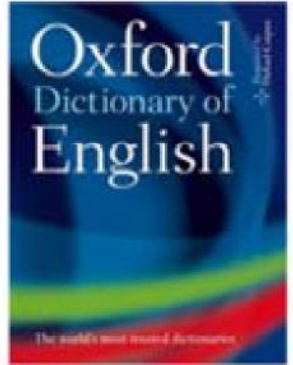
# Big Data: Just the next hype or reality?



Source: Google Trends

HO
GENT

# Definitions of big data

*Oxford English Dictionary:*
"data of a very large size, typically to the extent that its manipulation and management present significant logistical challenges."

*Wikipedia:*
"an all-encompassing term for any collection of data sets so large and complex that it becomes difficult to process using on-hand data management tools or traditional data processing applications."
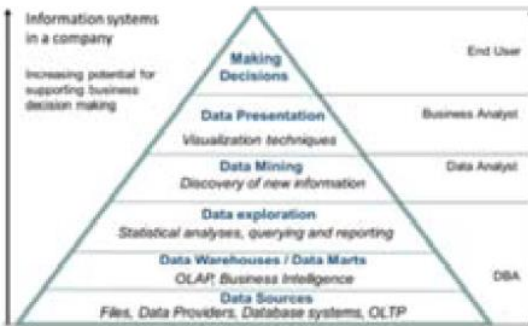
# Definitions of big data

My personal definition (Johan Decorte):

*"Data that is not (consciously) entered by a user but that arise, often spontaneously, as a by-product of other processes and that are (usually) used for purposes for which they were not originally intended".*

# 3. The origin of 'big' data



➡ New data characteristics
➡ New data challenges

HO
GENT

The origin of 'big' data

# **Internet Of Things**



VIDEO

## What is the Internet of Things (IoT)?

in LinkedIn Learning · May 2019 · From the course: Tech Sense

2m 21s

https://www.linkedin.com/learning/tech-sense/what-is-the-internet-of-things-iot?autoplay=true&u=77666441

(watch for free with your HOGENT-account)

# 4. The 4 V's of big data

- **Volume**: the amount of data, also referred to the data "at rest"
- **Variety**: the range of data types and sources that are used, data in its "many forms"
- **Velocity**: the speed at which data comes in and goes out, data "in motion", streaming data
- **Veracity**: the uncertainty of the data, data "in doubt"

HO
GENT

# **Big Data Challenges**

- **Volume**: BIG data → horizontal scaling/distributed data
- **Variety**: varied data → schemaless databases
- **Velocity**: fast data → NoACID databases
- **Veracity**: bad data → <span style="color:red">Data veracity handling</span>

HO
GENT

The 4 V's of big data

# BIG Data



WILL BE MEASURED IN **TERABYTES** 1TB = 1,000GB

WILL BE MEASURED IN **PETABYTES** 1PB = 1,000TB

WILL BE MEASURED IN **EXABYTES** 1EB = 1,000PB

VOLUME OF INFORMATION

LARGE

SMALL

1990's (RDMBS, DATA WAREHOUSE, ETC.)

2000's (CONTENT & DIGITAL ASSET MANAGEMENT)

2010's (NO-SQL, KEY/VALUE, ETC.)

- Don't decide too fast you have BIG data: tradition SQL databases can easily store several terabytes of data. Relational databases and SQL are still the preferred technology for most data related applications
- Automated and extremely fast generation and (near) real-time data processing of unstructured data require new technologies
  - NoSQL databases
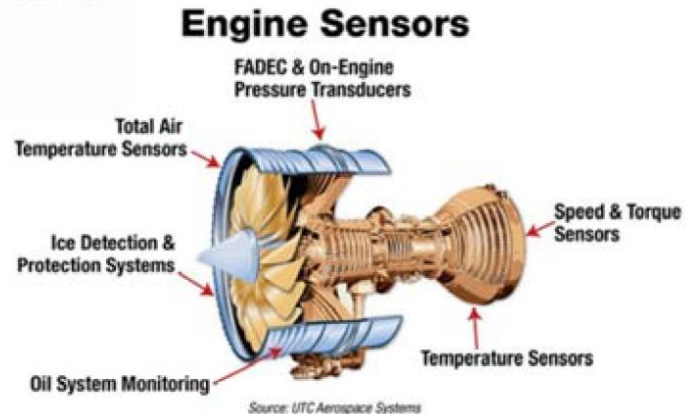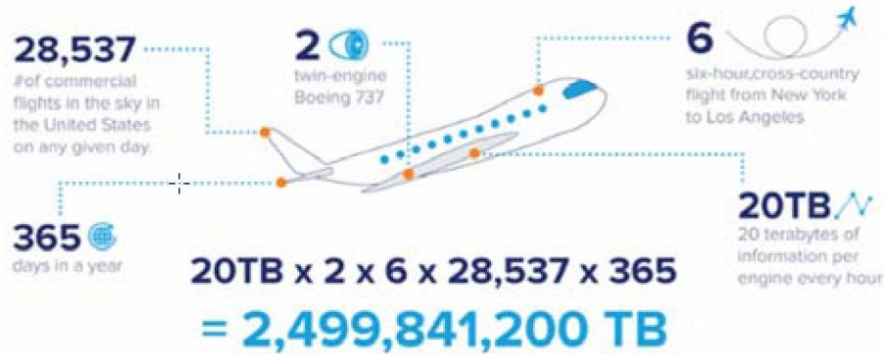  - Big data analysis via Python and R

HO GENT

# Varied Data



To store business cards in a traditional database system you have to create a uniform database structure with fixed database schema which requires expensive data transformations

→ Use case for NoSQL databases

HO
GENT

# Fast Data (IoT)



Sensor data from a cross-country flight

**28,537** # of commercial flights in the sky in the United States on any given day.

**2** twin-engine Boeing 737

**6** six-hour, cross-country flight from New York to Los Angeles

**365** days in a year

**20TB** 20 terabytes of information per engine every hour

20TB x 2 x 6 x 28,537 x 365
= 2,499,841,200 TB

**Engine Sensors**

FADEC & On-Engine Pressure Transducers

Total Air Temperature Sensors

Ice Detection & Protection Systems

Speed & Torque Sensors

Temperature Sensors

Oil System Monitoring

Source: UTC Aerospace Systems
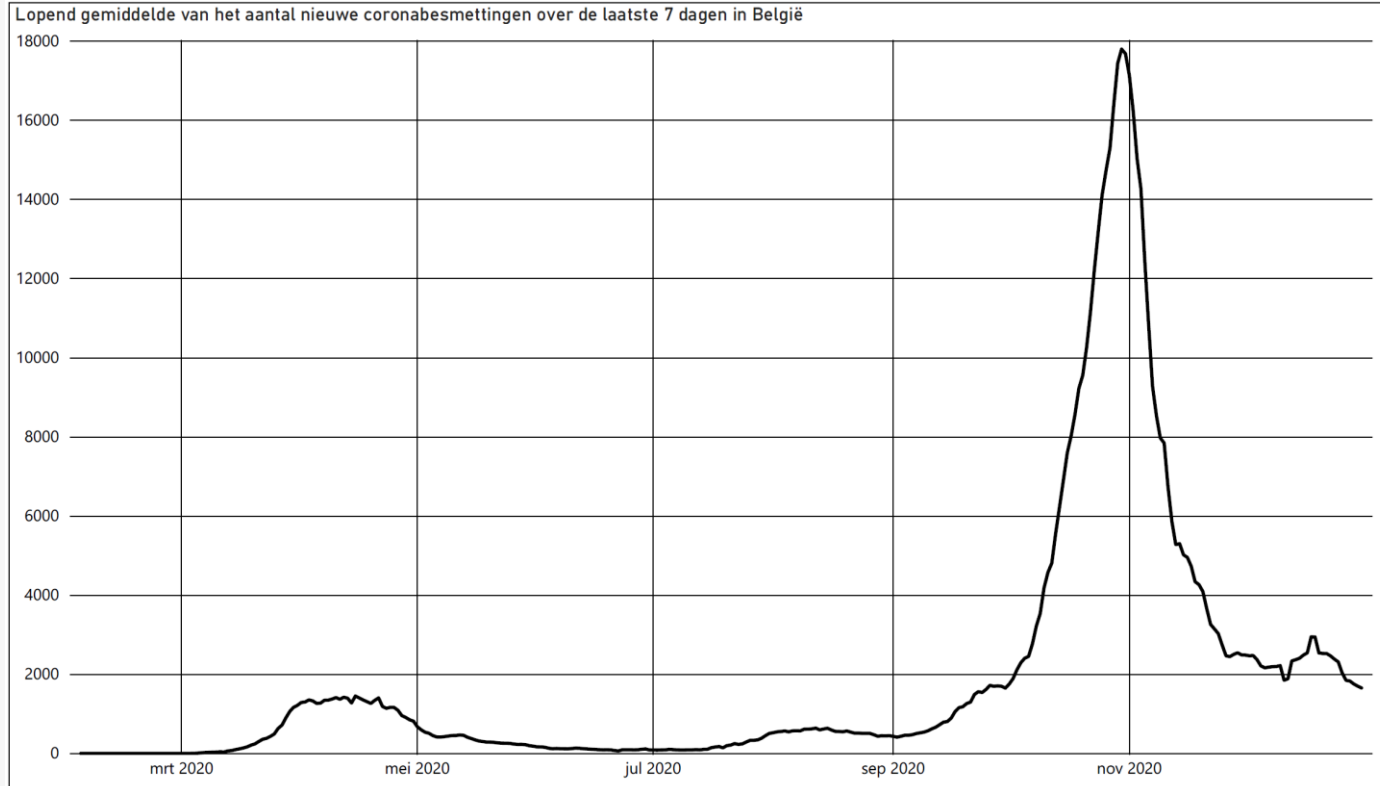
HO GENT

# Voluminous, Fast and Varied  Data (social media)

# Bad Data (bad veracity)

- Imprecise data
- Vague data
- Uncertain data
- Incomplete (missing) data
- Inconsistent data

→ Large volumes of varied data that need to be processed very fast are susceptive to bad data quality (failing sensors, networks, ...)

→ Handling bad or unknown quality is a big challenge.

→ Such data is also called big data

HO
GENT

# Example of inconsistent data



Lopend gemiddelde van het aantal nieuwe coronabesmettingen over de laatste 7 dagen in België

In the first COVID wave in Belgium, the laboratories detected an estimated one-thirtieth of the actual corona infections, in the second wave an estimated one-third. This gives the wrong impression that the second wave was much worse than the first in terms of number of infections.

*Data source: https://ourworldindata.org/coronavirus-source-data. Visualisation software used: Microsoft Power BI.*

HO GENT
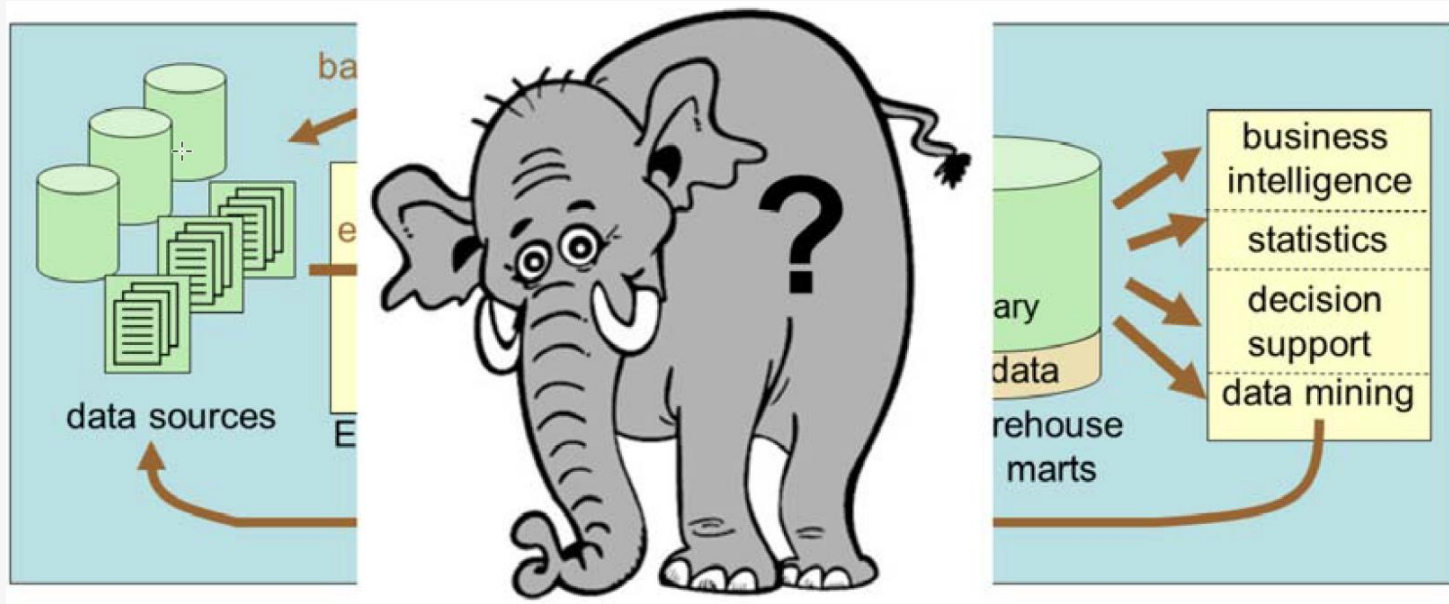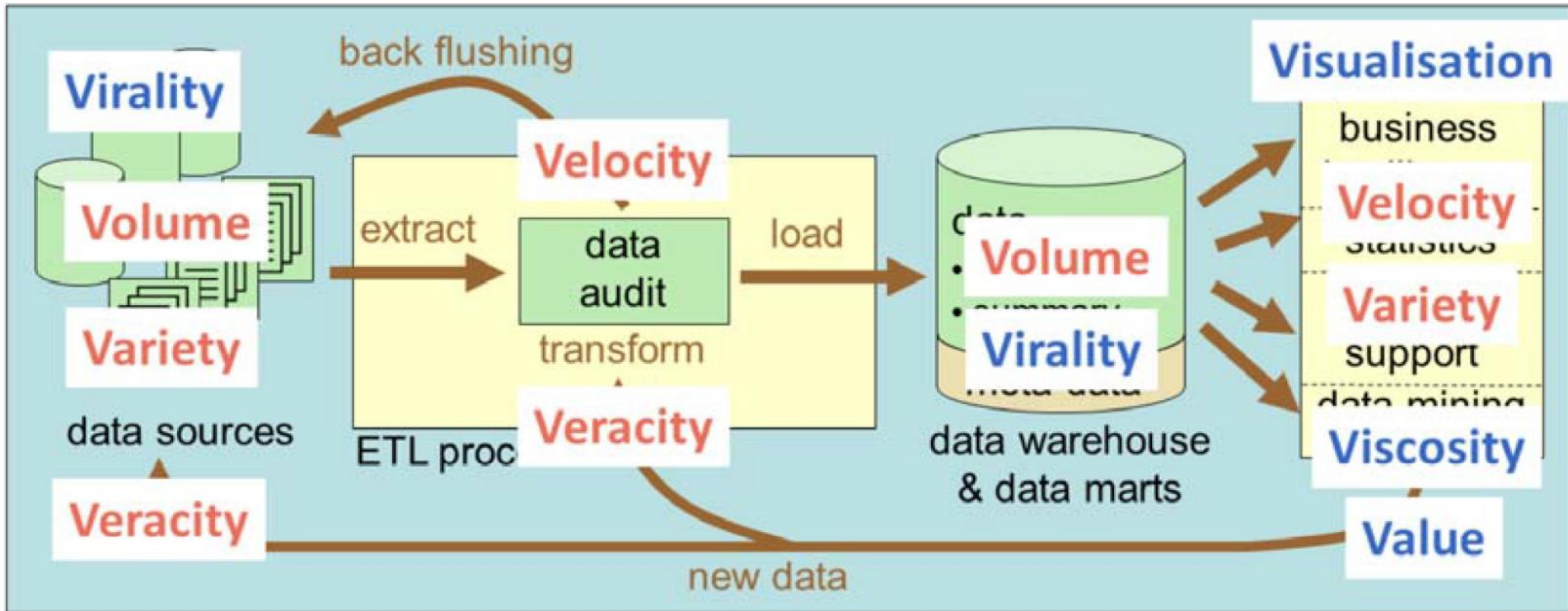
# 5. Some other V's



- **Virality**: How long do we need to keep data, when does it become outdated?
- **Viscosity**: Do we have enough data to perform (statistically) relevant analyses?
- **Visualisation**: can the results easily be presented?
- **Value**: what the value of our data? "Data is the new gold".

HO
GENT

# 6. Bottlenecks



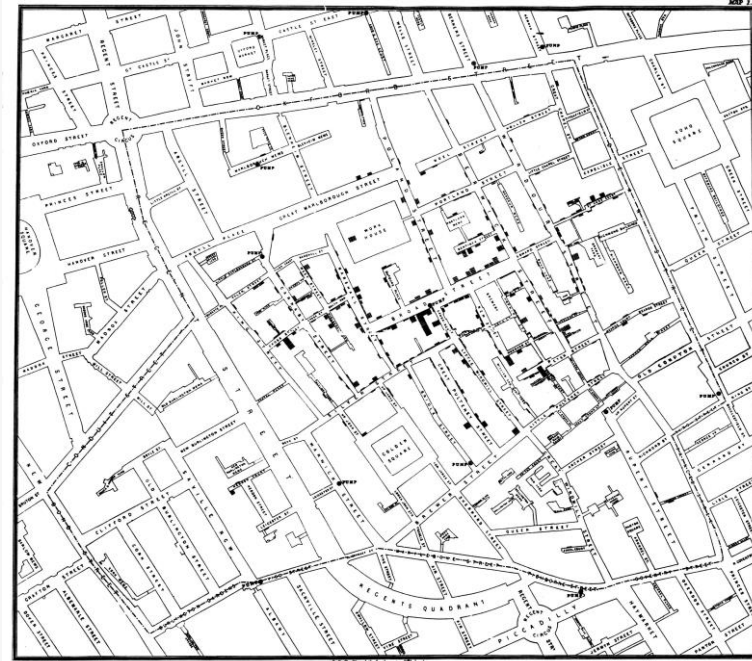Where do we find these challenges in the (theoretical) information processing?

# Bottlenecks

# 7. **Examples of how big data can be used**

- London (Soho) cholera outbreak in 1854 →
 (deaths per household,
  centralized around infected water pump)

- Counting the number of visitors of the 'Gentse Feesten':
https://www.nieuwsblad.be/cnt/dmf20180722_03626134

- Better predictions using Big Data:
see next video (dutch):

https://www.youtube.com/watch?v=m9YcrXIc2IM



HO GENT

# If you want to now more ...



Lannoo

Koopjes tot -70% | Agenda | Inspiratie | Vacatures | Contact          0 | 🛒  Inloggen

Kies een categorie

Zoek op auteur, trefwoord, titel

Home > Mens & Maatschappij > Big Data

## Big Data
### Johan Decorte

*Big data: een revolutie ontrafeld*

Elke twee jaar verdubbelt de hoeveelheid beschikbare gegevens in de wereld. Met krachtige en betaalbare computers en toegankelijke algoritmes kunnen we die data omvormen tot bruikbare inzichten over mensen, machines en processen. Het is niet alleen een nieuw businessmodel voor heel wat bedrijven: ook onze gezondheid en het klimaat kunnen er wel bij varen. Maar elke technologie kan misbruikt worden, en big data is daarin geen uitzondering. Het risico bestaat dat onze privacy en ethische waarden in het gedrang komen.

Johan Decorte legt uit hoe datawetenschappers orde scheppen in de immense massa aan gegevens en belicht enkele van de belangrijkste toepassingsgebieden.

€ 12,⁵⁰

**Reserveer nu**

Gratis verzending vanaf € 20 (Benelux)

Binnen 1-2 werkdagen in huis

Beschikbaarheid: Binnenkort leverbaar

• Boekhandel

HO GENT