# Introduction to Big Data

Alexandru Costan

*alexandru.costan@insa-rennes.fr*

# Rennes capital of Brittany and city of over 400 000 inhabitants

**Rennes is the:**

- 11th largest city in France
- 2nd city in France for its student population

**Rennes is situated:**

- 45 minutes from the sea (Saint-Malo)
- 1,5 hours from Paris

**Rennes is a university and research hub of international importance in 3 sectors :**

- Health
- Digital Technology
- Eco-activities

# INSA Rennes,
## a public-funded graduate and post graduate engineering school member of the INSA group

# INSA Rennes, public-funded graduate and post graduate engineering school

Founded in 1966, INSA Rennes is classed among **the best graduate and post graduate engineering schools** in France.

In addition to being a research center, INSA Rennes gives professional training in engineering and research in **2 poles of excellence:**

- **Information and Communication Sciences and Technologies (ICST)**
- **Materials, Structures and Mechanics (MSM)**

# INSA Rennes, at the heart of a science and technology campus

## Equipped 17 hectare campus

- Accommodation and catering in situ
- Pedagogical, scientific and sports equipment

## Within the Le POOOL technology cluster

- More than 300 companies
- 80% engineering and technology companies

## Within the science campus of Rennes

- 70289 students
- 32 institutes of higher education and graduate schools

**Today**

**What is Big Data?**

- Recognize some of the main terminology

**What can we do with Big Data?**

- Realize the potential of Big Data

**Why is it difficult?**

- Understand why we need a different paradigm

**How to store and process Big Data?**

- Know the existing tools



HOW'S THE BIG DATA PROJECT COMING ALONG, HOSKINS?

© D.Fletcher for CloudTweaks.com

# Not focusing on

- ## machine learning

- ## data mining

- ## natural language processing

## although we will touch on these

INSTITUT NATIONAL DES SCIENCES APPLIQUÉES | RENNES

**Linked in Learning**

# The Skills Companies Need Most in 2020

Data Scientist
intersection
scientist is h

**MATH & STATIST**

☆ Machine lear
☆ Statistical m
☆ Experiment d
☆ Bayesian infe
☆ Supervised le
random fores
☆ Unsupervise
dimensionali
☆ Optimization
variants

**DOMAIN & SOFT S**

☆ Passionate a
☆ Curious abou
☆ Influence wit
☆ Hacker mind
☆ Problem solv
☆ Strategic, pr
innovative a

## Top 5 Soft Skills

1. Creativity — 
2. Persuasion — 
3. Collaboration — 
4. Adaptability — 
5. Emotional intelligence — new

## Top 10 Hard Skills

1. Blockchain — new
2. Cloud computing — −1
3. Analytical reasoning — 
4. Artificial intelligence — −2
5. UX design — 
6. Business analysis — +10
7. Affiliate marketing — new
8. Sales — 
9. Scientific computing — +3
10. Video production — −3

**RAMMING ABASE**

iter science fundamentals
ng language e.g. Python
cal computing packages, e.g., R
ases: SQL and NoSQL
onal algebra
I databases and parallel query
sing
duce concepts
p and Hive/Pig
n reducers
ence with xaaS like AWS

**MUNICATION UALIZATION**

o engage with senior
eement
elling skills
ate data-driven insights into
ons and actions
art design
kages like ggplot or lattice
edge of any of visualization
g. Flare, D3.js, Tableau

— means that it remains at the same spot as last year.

UEb

# Lectures outline

- **Big Data overview**
  - Data and Processing Models
  - Consistency

- **Programming with Big Data**
  - Google MapReduce
  - Apache Hadoop

# Readings

- Tony Hey, Stewart Tansley, Kristin Tolle „The Fourth Paradigm: Data-Intensive Scientific Discovery ", Microsoft Research

- Jeffrey Stanton, "Introduction to Data Science", Syracuse University Press

- Jeffrey Dean, Sanjay Ghemawat, "MapReduce: Simplified Data Processing on Large Clusters ", OSDI 2004

- Matei Zaharia, Mosharaf Chowdhury, Michael J. Franklin, Scott Shenker, Ion Stoica, "Spark: Cluster Computing with Working Sets", NSDI 2012

- … additional bibliography specific to each lecture

# Sources of Big Data

# Big Data Units

| Unit | Size | What it means |
|---|---|---|
| Bit (b) | 1 or 0 | Short for "binary digit", after the binary code (1 or 0) computers use to store and process data |
| Byte (B) | 8 bits | Enough information to create an English letter or number in computer code. It is the basic unit of computing |
| Kilobyte (KB) | 1,000, or $2^{10}$, bytes | From "thousand" in Greek. One page of typed text is 2KB |
| Megabyte (MB) | 1,000KB; $2^{20}$ bytes | From "large" in Greek. The complete works of Shakespeare total 5MB. A typical pop song is about 4MB |
| Gigabyte (GB) | 1,000MB; $2^{30}$ bytes | From "giant" in Greek. A two-hour film can be compressed into 1-2GB |
| Terabyte (TB) | 1,000GB; $2^{40}$ bytes | From "monster" in Greek. All the catalogued books in America's Library of Congress total 15TB |
| Petabyte (PB) | 1,000TB; $2^{50}$ bytes | All letters delivered by America's postal service this year will amount to around 5PB. Google processes around 1PB every hour |
| Exabyte (EB) | 1,000PB; $2^{60}$ bytes | Equivalent to 10 billion copies of *The Economist* |
| Zettabyte (ZB) | 1,000EB; $2^{70}$ bytes | The total amount of information in existence this year is forecast to be around 1.2ZB |
| Yottabyte (YB) | 1,000ZB; $2^{80}$ bytes | Currently too big to imagine |

The prefixes are set by an intergovernmental group, the International Bureau of Weights and Measures. Yotta and Zetta were added in 1991; terms for larger amounts have yet to be established.
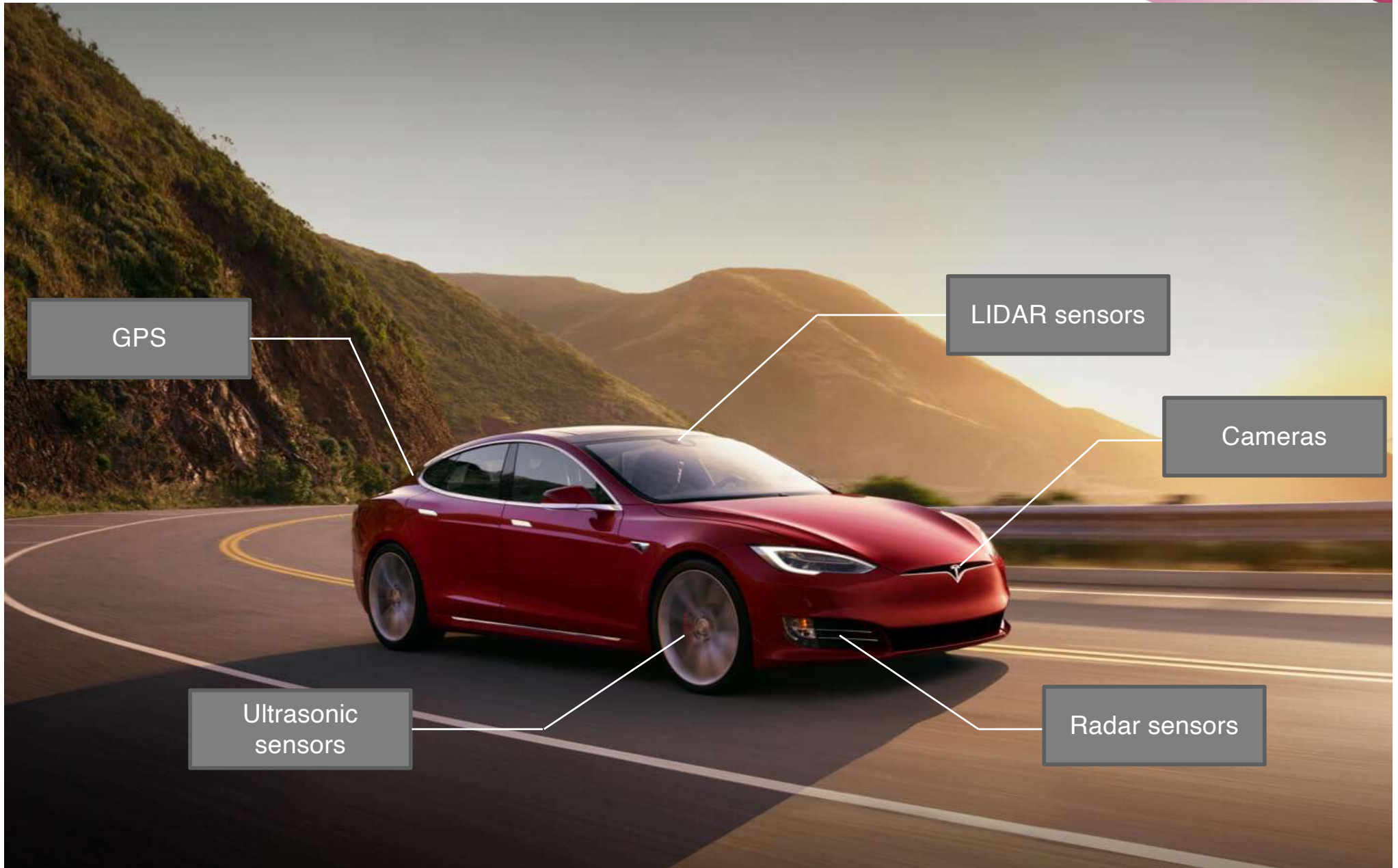
Source: *The Economist*

# 1 GB of data
# / person
# / day



# x 5,000,000,000 users

# The big picture (in 2020)

- **10PB** of Facebook data per day

- **500M** tweets **per day**

- Google processes **100PB a day**

GPS

LIDAR sensors

Cameras

Ultrasonic sensors

Radar sensors

| Sensor type | Quantity | Data generated |
|---|---|---|
| Radar | 4–6 | 0.1–15 Mbit/s |
| LIDAR | 1–5 | 20–100 Mbit/s |
| Camera | 6–12 | 500–3,500 Mbit/s |
| Ultrasonic | 8–16 | <0.01 Mbit/s |
| Vehicle motion, GNSS, IMU | - | <0.1 Mbit/s |

**TOTAL ESTIMATED BANDWIDTH**

3 Gbit/s (~1.4TB/h) to 40 Gbit/s (~19 TB/h)

# The levels of Autonomous Vehicles



**Autonomous Vehicles "Tipping Point"**
Transition from human drivers to vehicles driving

Our World Today

Near to Distant Future

| 0 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| **HUMAN ONLY** | **MODERN VEHICLE** | **MODERN PLUS** | **PARTIAL AUTONOMY** | **FULL AUTONOMY (+ HUMAN)** | **FULL AUTONOMY (NO HUMAN)** |
| The driver (human) controls everything: steering, brakes, throttle, power | Most functions are still controlled by a driver, but some (like braking) can be done automatically by the car | At least 2 functions are automated (like cruise control & lane-centering), but the driver must be ready to take control of the vehicle | Drivers are still necessary, but are not required to monitor the situation as with previous levels | Vehicles perform all safety-critical driving functions and monitor roadway conditions for an entire trip, with option for human driving | No option for human driving - no steering wheel or controls |

# Self-driving cars

## 10 light year away
The SKA will be so sensitive that it will be able to detect an airport radar on a planet at this distance

## 2'000'000 years
The data collected by the SKA in a single day would take nearly two million years to playback on an ipod

## 1'000'000+
of 500GB laptops can be filled with SKA data every year

## On two sites

### South Africa
*SKA1-MID*

**≈200** dishes

**5x** more sensitive than any other radio telescope

**33'000 m²** of total collecting area
*(=126 tennis courts)*

### Western Australia
*SKA1-LOW*

**8x** more sensitive than any other radio telescope

**≈130'000** antennas spread between 500 stations

**420'000 m²** of total collecting area
*(=58 football pitch)*

# 3D Rendering for Animations



## Avatar

- 40,000 processors handling 8 GB of data per second, 24 hours a day
- A final copy equated to 17 GB per minute of storage
- The sum of required computing power for the creation of Avatar reached 205 teraflops



## Monsters University

- 2,000 computers with more than 24,000 cores
- Still took 25-30 hours to render a single frame
- All in all, it took over 100 million hours of CPU time to render the entire movie

"Obama Signs Executive Order Banning The Pledge Of Allegiance In Schools Nationwide" ABCNews.com.co

2,177,000 Facebook shares, comments, and reactions

"Pope Francis Shocks World, Endorses Donald Trump for President, Releases Statement" Ending the Fed

961,000

"Trump Offering Free One-Way Tickets to Africa & Mexico for Those Who Wanna Leave America" tmzhiphop.com

802,000

"FBI Agent Suspected in Hillary Email Leaks Found Dead in Apparent Murder-Suicide" Denver Guardian

567,000

"RAGE AGAINST THE MACHINE To Reunite And Release Anti Donald Trump Album" heaviermetal.net

560,000

More than 50% of traffic to fake news sites comes from Facebook

Fake news spreads 6x faster than accurate news on Twitter, and falsehoods are 70% more likely to be retweeted. [MIT, 2018]

# Fighting fake news with Big Data

## Facebook to flag fake news stories

- Users report the story as bogus or Facebook's software detects something odd
- Facebook sends the story to some of the organizations that have signed on to provide free fact-checking (e.g. Snopes, Politifact)
- If two of those fact-checkers think it's bogus, the label goes on.

## Google presented an algorithm for fake news detection

- Websites are scored according to the accuracy of the of the facts presented



Betsy Marshall Barda shared Joe Redner's post.
March 2 at 5:56pm · 

OMG if this is true I will laugh soooo hard. He's right - we need to investigate the leaker!!! LOLOLOL

Joe Redner
March 2 at 4:39pm · 

Investigators from A.R.H. Intelligence and Z|13 Security believe that the unsecured Android device was most likely compromised by a suspicious animated GIF that was sent to President Trump via text message.

### Trump's Unsecured Android Device Source Of Recent White House Leaks

THESEATTLETRIBUNE.COM

⚠ Disputed by Snopes.com and PolitiFact

👍 Like   💬 Comment   ➤ Share   </> Embed

1 Comment

# Deep Fake

**Machine learning** is exceptionally good at learning how to **exploit human psychology**
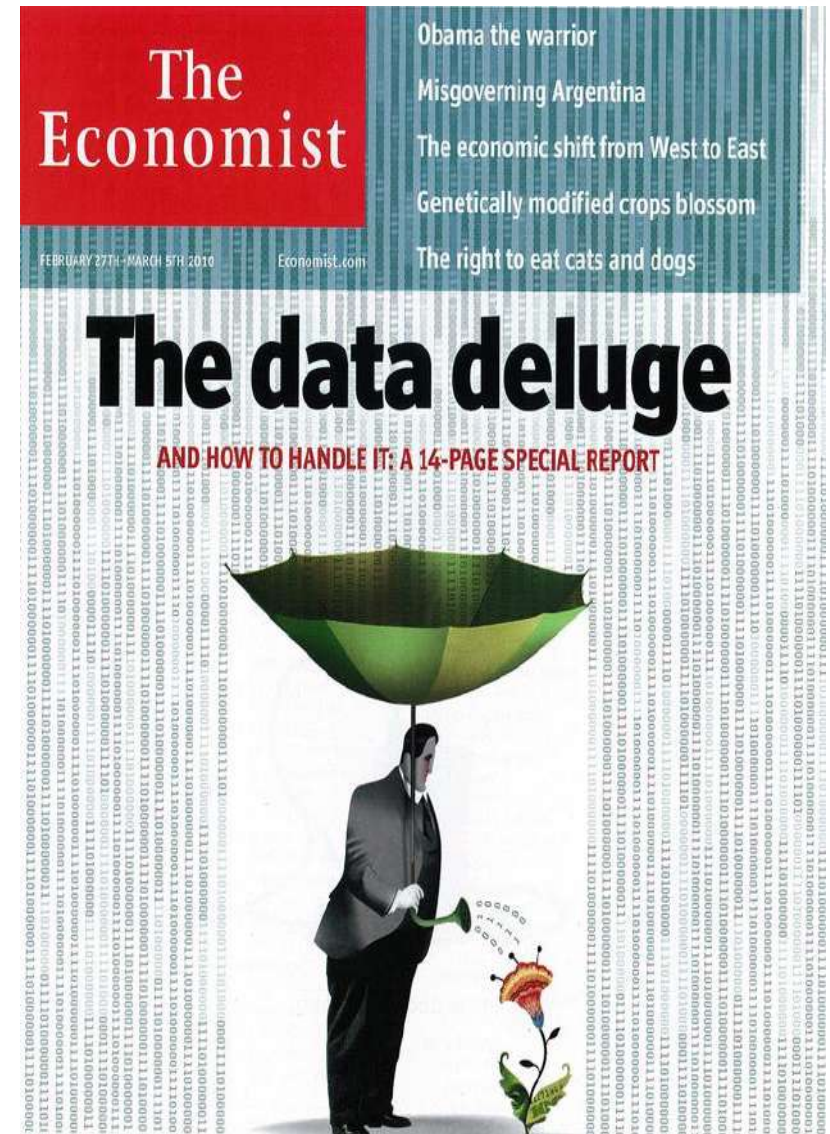
- because the internet provides a vast and fast **feedback loop** to learn what will reinforce and or break beliefs by demographic cohorts

- an AI engine that can generate messages and immediately test if the message is effective

**Generative adversarial networks** help create AI-generated images and deepfakes

- Two neural networks (a generator and a discriminator) work together to create the fictional image

# Total size of the Digital Universe

# 50 ZetaBytes
## in 2021

**Defining Big Data**

# What is Big Data ?

# Big Data Features: the 3 Vs



| Volume | Variety | Velocity |
|---|---|---|

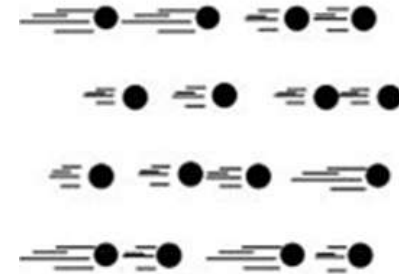Terabytes to Exabytes of existing data to process

Structured and unstructured data

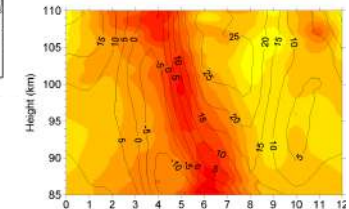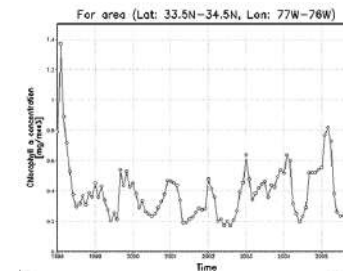Data requiring milliseconds to seconds to respond
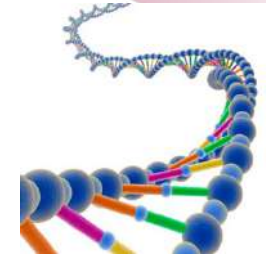
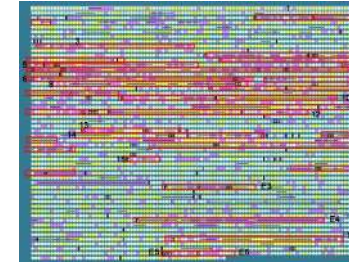**Structured Data**   **Semi-Structured Data**   **Unstructured Data**

# Variety
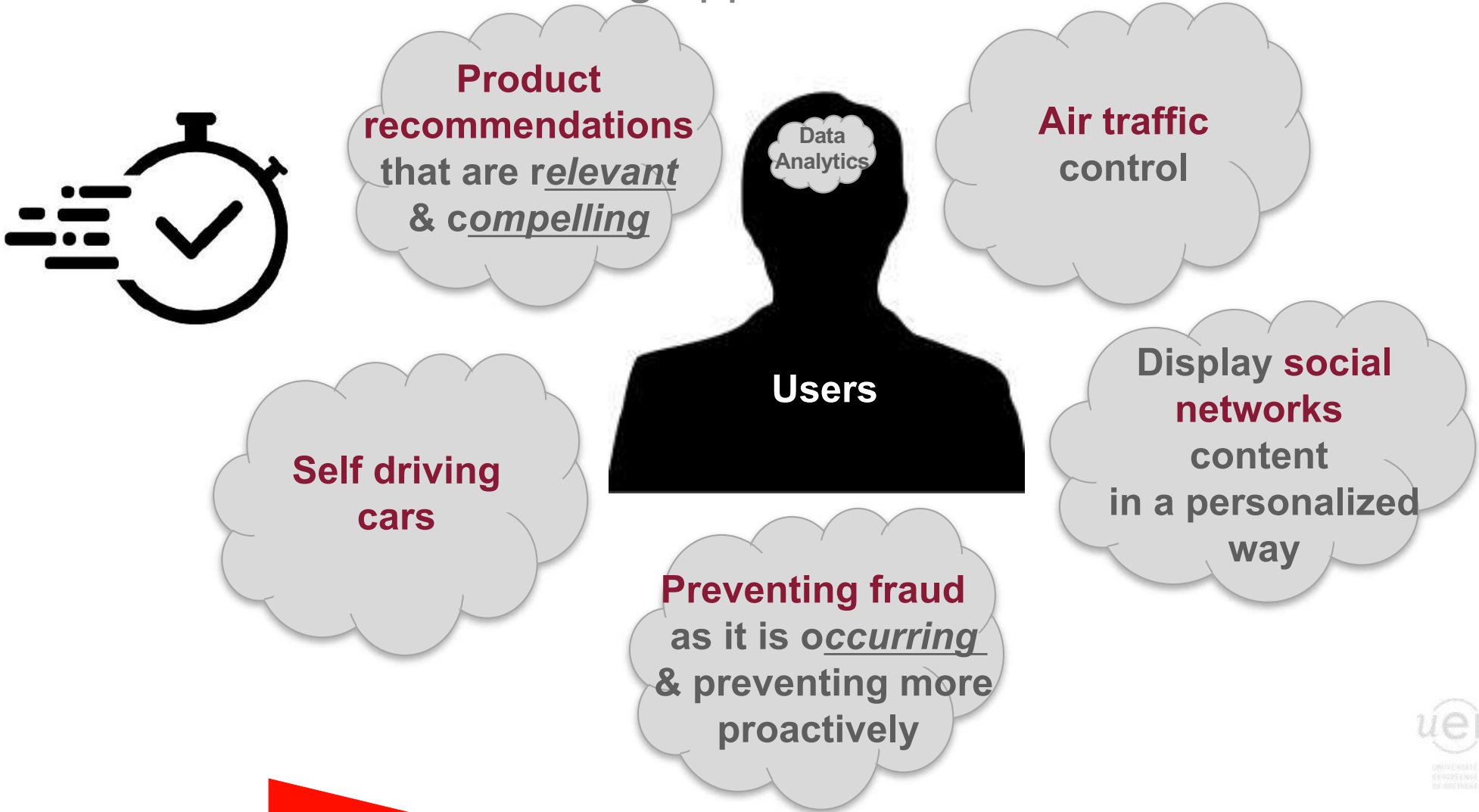
- **Structured**
  - Relational data: tables, transactions, legacy data

- **Semi-structured**
  - Text: web, documents
  - XML, JSON
  - Large graphs: social networks, semantic web (RDF)

- **Unstructured**
  - Streaming: You can only scan the data once

- A single application can process many types of data

To extract knowledge➔ all these types of data need to be linked together

# Velocity

- **Real-time / fast data:** data is being generated fast and need to be processed fast

- Late decisions: missing opportunities

**Product recommendations that are r*elevant* & c*ompelling*

**Data Analytics**

**Air traffic control**

**Users**

**Display social networks content in a personalized way**

**Self driving cars**

**Preventing fraud as it is o*ccurring* & preventing more proactively**

Processing
Big Data

# What to do with all these data ?

- Stored data = **Costs**

- Information from data = **Profit**

**Goal:** Extract valuable information / added value from these huge data

# Hardware: distributed infrastructures

– Cloud computing allows to lease computing and storage resources



# Software: new programming models

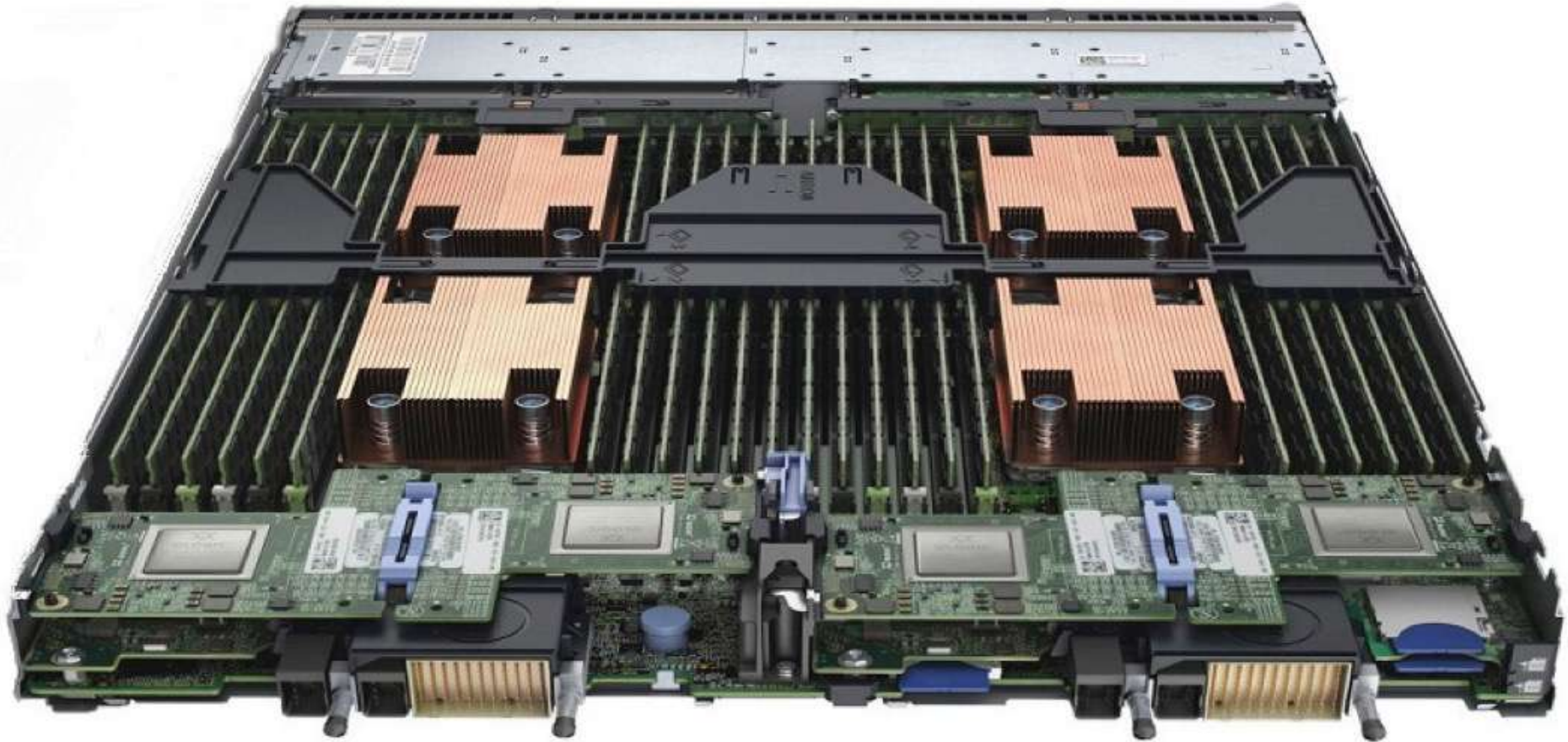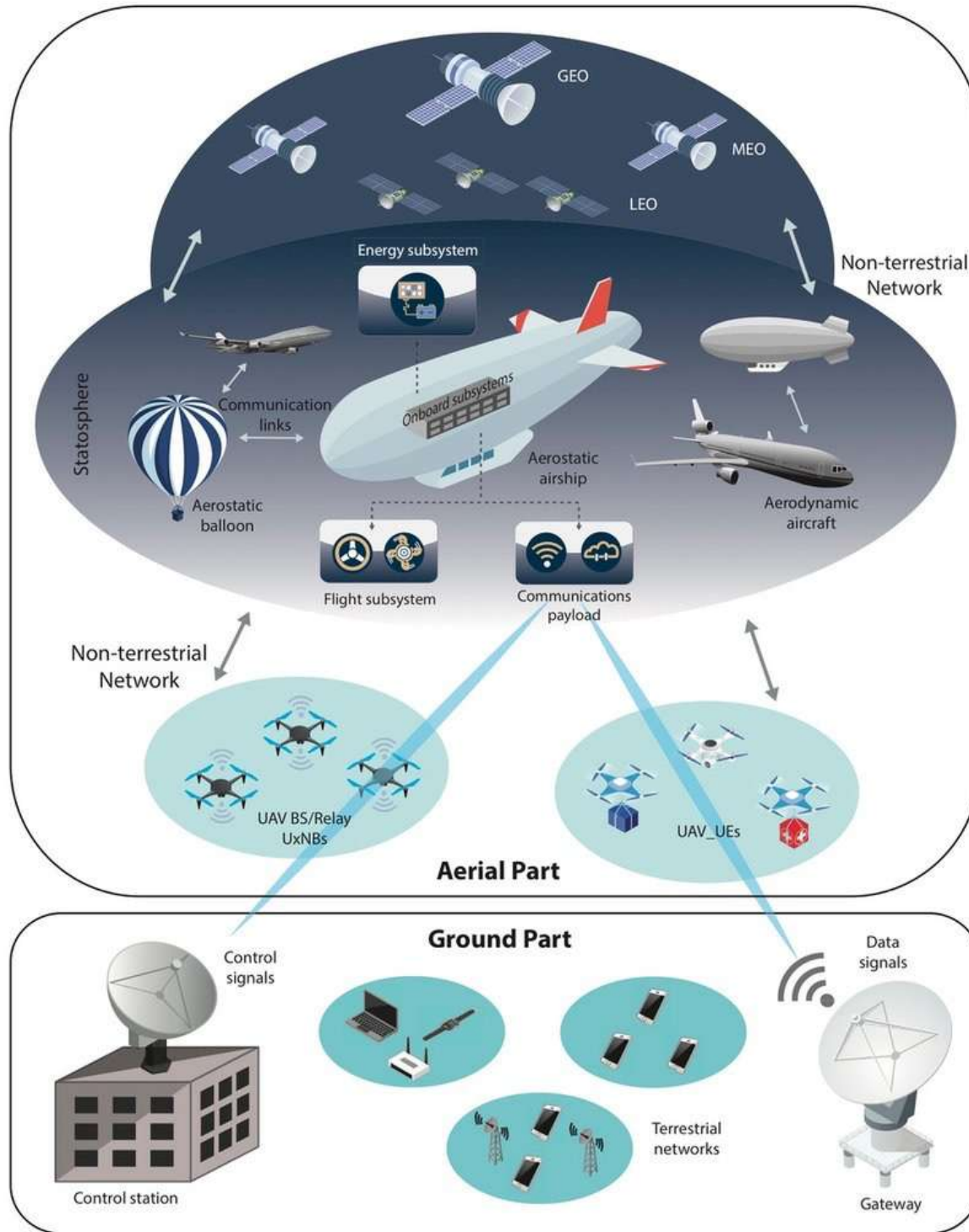– MapReduce: simple yet scalable model for Big Data processing

- # Everything as a Service

- The **delivery of computing** as a service rather than a product

- Shared resources, software, and information are provided as a **metered service** over a network (typically the Internet)

# What is the cloud? One datancenter?
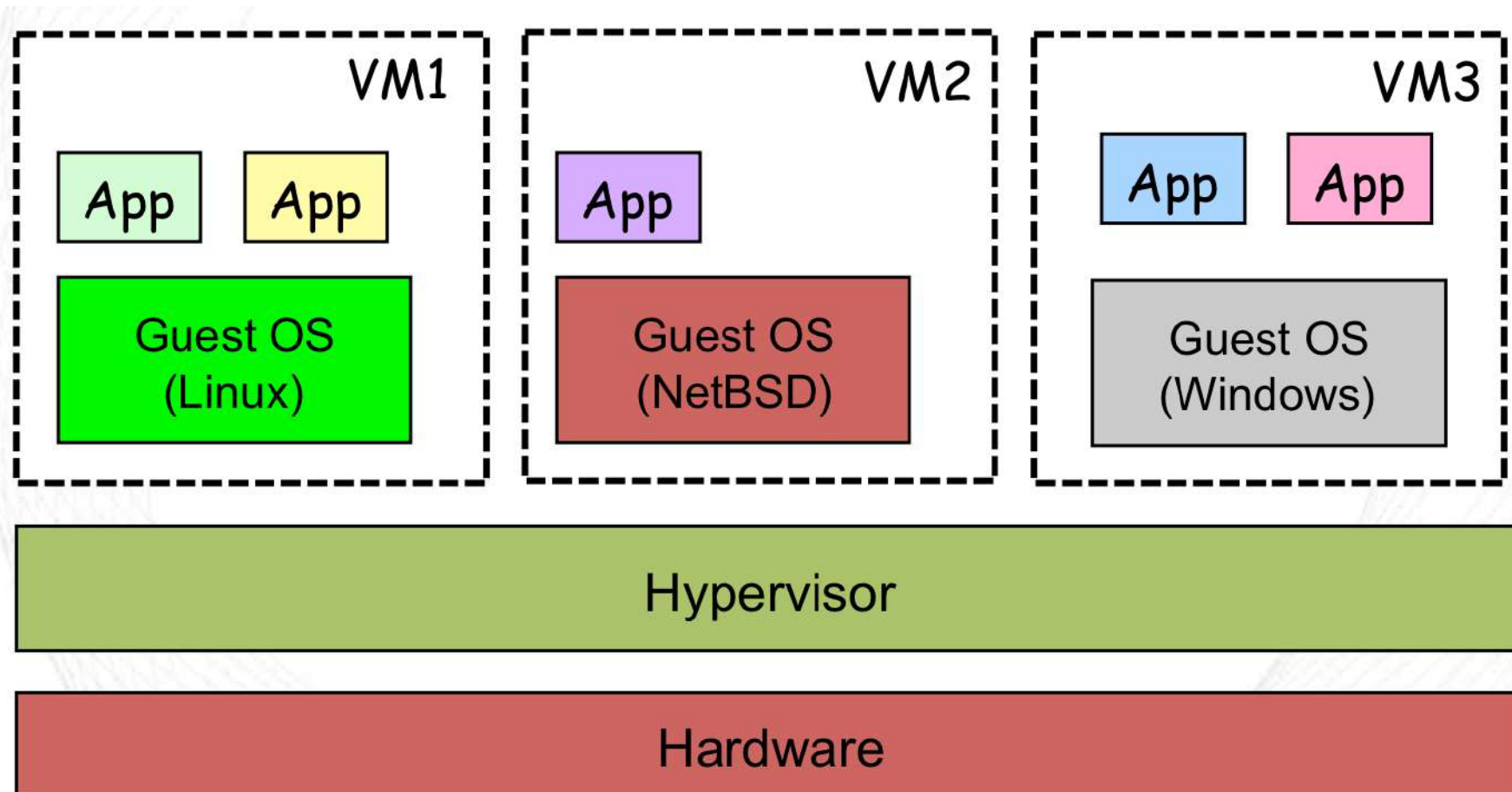
# What is the cloud? One datancenter?

# Many geographically distributed datacenters

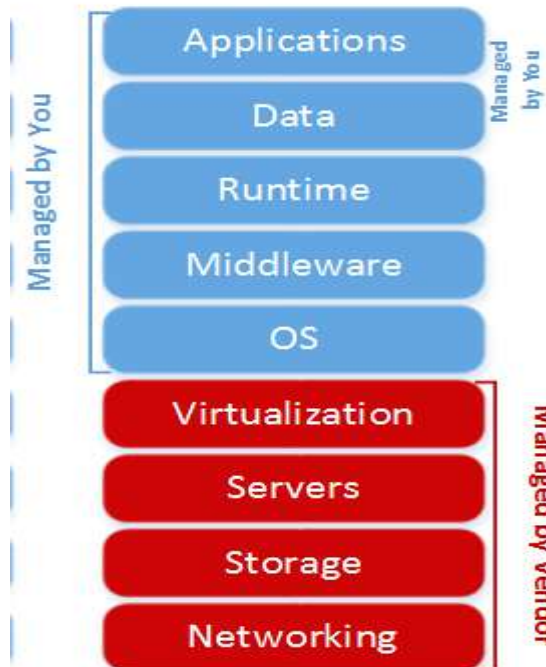# Enabling technology: virtualisation

- Allows multiple virtual machines to run on a single physical machine
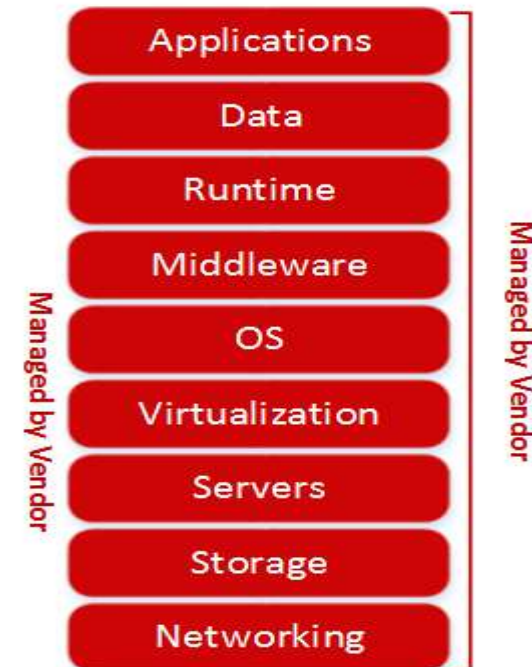
# Types of clouds



**IaaS:** Infrastructure as a Service

**PaaS:** Platform as a Service

**SaaS:** Software as a Service

Bare **hardware**: storage and compute

Development **environments** to create applications

Business **applications** (e.g. Dropbox, Office 365)