



Introduction to Big Data

Daniel Hagimont
hagimont@enseiht.fr



The lecture starts with a general introduction to Big Data.

Context

- We generate more and more data
 - Individuals and companies
 - Kb → Mb → Gb → Tb → Pb → Eb → Zb → Yb → ???
- Few numbers
 - In 2013, Twitter generates 7 Tb per day and Facebook 10 Tb
 - The square kilometre array radio telescope
 - ◆ Products 7 Pb of raw data per second, 50 Tb of analyzed data per day
 - Airbus generates 40 Tb for each plane test
 - Created digital data worldwide
 - ◆ 2010 : 1,2 Zb / 2011 : 1,8 Zb / 2012 : 2,8 Zb / 2020 : 40 Zb
 - ◆ 90 % of data were created in the last 2 years



Today's digital applications are generating more and more data. These data may be generated by individuals (through social networks) or companies.

Examples are Twitter and Facebook who analyze the behavior of their users, large scientific equipments like the square kilometre array radio telescope.

Another example I have been working on comes from Airbus. When they have designed and built a new plane, during the test phase, the plane is equipped with a huge set of captors which monitor everything. Each test flight generates about 40 TB of data and they run several test flights per week over a year.

The amount of data generated worldwide becomes so large and it is exponential.

Context

What Happens in an Internet Minute?



3

This is a picture I grabbed on the net which illustrates this deluge of data.

Context

■ Many data sources

- Multiplication of computing devices and connected electronic equipments
- Geolocation, e-commerce, social networks, logs, internet of things ...

■ Many data formats

- Structured and unstructured data



4

It is important to remember that these data are coming from many sources, with the multiplication of devices (personal computers, phones, tablets) and the development of cloud applications.

These diverse sources generate data in many different formats. Such format can be structured (with XML, Json ..) or unstructured (like textual data).

Applications domains

- Scientific applications (biology, climate ...)
- E-commerce (recommendation)
- Equipment supervision (e.g. energy)
- Predictive maintenance (e.g. airlines)
- Espionage

The NSA has built an infrastructure that allows it to intercept almost everything. With this capability, the vast majority of human communications are automatically ingested without targeting. E Snowden

5

These data are generated and used (analyzed) in many application domains

Scientific application : biology (genomic) or climate (weather prediction)

E-commerce : profiling users for product recommendation

Equipment supervision : energy where they profile consumption

Predictive maintenance : for instance in the airline industry, they monitor each flight with captors embedded in planes. When a failure occurs, they analyze the monitored data before the failure to identify symptoms that preceded the failure. Therefore, by analyzing data from each flight, they can detect such symptoms before a failure occurs and anticipate a maintenance to prevent the failure.

Espionage : E. Snowden revealed that the NSA was analyzing data from so many sources in a wish to control everybody.

New jobs

■ Data Scientist

- IT specialist : know how to develop, parameterize, deploy tools
- HPC specialist : parallelism is key
- Statistician : know how to use mathematics to classify, group and analyze information
- Manager : know how to define objectives and identify the value of information

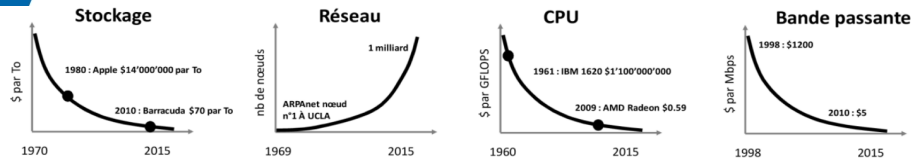
6

The field of big data is generating many jobs:

- manager who are dealing with high level objectives
- mathematicians who are analyzing and modeling data
- HPC specialists who know how to improve performance
- IT specialists who master the infrastructure (hardware and software)

Computing infrastructures

■ The reduced cost of infrastructures



- Main actors (Google, Facebook, Yahoo, Amazon ...) developed frameworks for storing and processing data
- We generally consider that we enter the Big Data world when processing cannot be performed with a single computer

7

Decades ago, data treatments were performed by a single computer, a multiprocessor machine if computing power was required.

Two main evolutions modified that:

- new applications which require huge computing power
- the reduced cost of hardware (mainly servers and networks), leading to the development of clusters (a huge set of interconnected servers).

A cluster is a set of interconnected servers which are exploited by a distributed operating system in order to give the illusion of a giant computer.

The main actors of the field developed such distributed operating system frameworks to implement a sort of giant computer for storing and processing data.

We generally consider that we enter the Big Data world when processing cannot be performed with a single computer (even a large one) and when a cluster is required.

Definition of Big Data

■ Definition

- Rapid treatment of large data volumes, that could hardly be handled with traditional techniques and tools

■ The three V of Big Data

- Volume
- Velocity
- Variety
- Two additional V
 - ◆ Veracity
 - ◆ Value

8

One (among many) definition of Big Data is the rapid treatment of large data volumes that could hardly be handled with traditional techniques and tools (e.g. a database system on a single computer).

In the field, people are talking about the 3 V which characterize Big Data:

- Volume : dealing with very large data volumes
- Velocity : treatment should be fast enough
- Variety : data may be of very different types

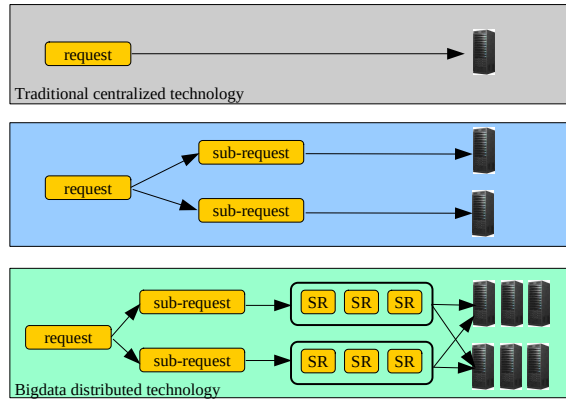
Sometimes, they add 2 additional V :

- Veracity : big data may also consider the trust we have into the handled data
- Value : it may also consider the value of the data

And may be at the time of today, they added again other Vs.

General approach

- Main principle : divide and conquer
 - Distribute IO and computing between several devices



9

The general approach used to efficiently treat very large datasets is to apply a well known principle called "divide and conquer".

The principle is to divide a task into several sub-tasks which can be distributed on distributed computers, therefore benefiting from parallel IO (reads from different disks) and/or parallel computation (on different processors).

Figure top : in a traditional centralized setting, a request or task is executed on one computer.

Figure middle : a request can be divided into 2 sub requests executed in parallel on 2 separate computers.

Figure bottom : the request can be divided into many sub-requests scheduled on many computers. This can be done if the initial request is sufficiently large.

Notice here that in this general principle, we don't precise how data are moved to the computers where sub-requests are executed.

Solutions

■ Two main families of solution

- Processing in batch mode (e.g. Hadoop)
 - ◆ Data are initially stored in the cluster
 - ◆ Various requests are executed on these data
 - ◆ Data don't change / requests change
- Processing in streaming mode (e.g. Storm)
 - ◆ Data are continuously arriving in streaming mode
 - ◆ Treatments are executed on the fly on these data
 - ◆ Data change / Requests don't change

10

The previous general principle leads to 2 main families of solution.

Processing in batch mode (Hadoop is an example). Data are initially stored on the computers' disks. For instance a very large dataset is divided into blocks and the blocks are distributed on the machines. Then, various requests can be issued to analyze the data. Such a request is divided into sub-requests which will handle the blocks on the different machines (in parallel). What is important here is that the data are installed in the cluster (installed means here that data don't change, they are here to be read and analyzed, not modified) and that many requests can be issued on the same dataset.

Processing in streaming mode (Storm is an example). Data are continuously arriving in streaming mode. Treatments are decomposed in tasks that are deployed in the cluster