



IAE – Octobre 2019

Big Data - Résumé

**Christophe Menichetti**

Big Data and AI Lead Architect

[christophe.menichetti@fr.ibm.com](mailto:christophe.menichetti@fr.ibm.com)

IBM Montpellier Client Center (IBMCCMPL)

## 1 De la BI au Big Data : Introduction

0000000 000010101

&gt;&gt;&gt; 111111 010101

## 2 Usages et Technologies Big Data

0111100

## 3 Du Big Data vers la Data Science

## 4 Architectures BI, Big Data et AI

0000000 000010101

Introduction à l'analyse de la donnée  
Architecture et Solution Data Warehouse  
Les nouvelles problématiques de la data (3V)  
les objets Connectés (IoT)

Nouvelle approche de stockage/traitement donnée ---  
l'écosystème Hadoop  
Focus sur Map Reduce et Spark  
Les différentes offres du marché

La donnée : le nouveau pétrole  
Le Machine Learning et Intelligence artificielle  
Les différents Metiers du Big Data et rôles clés  
Focus sur le Data Scientist

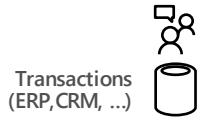
architecture On Premises  
Le Cloud et solutions Data dans le cloud  
Etudes de Cas

**Synthese de tout le cours**

10111010 1010 10101001

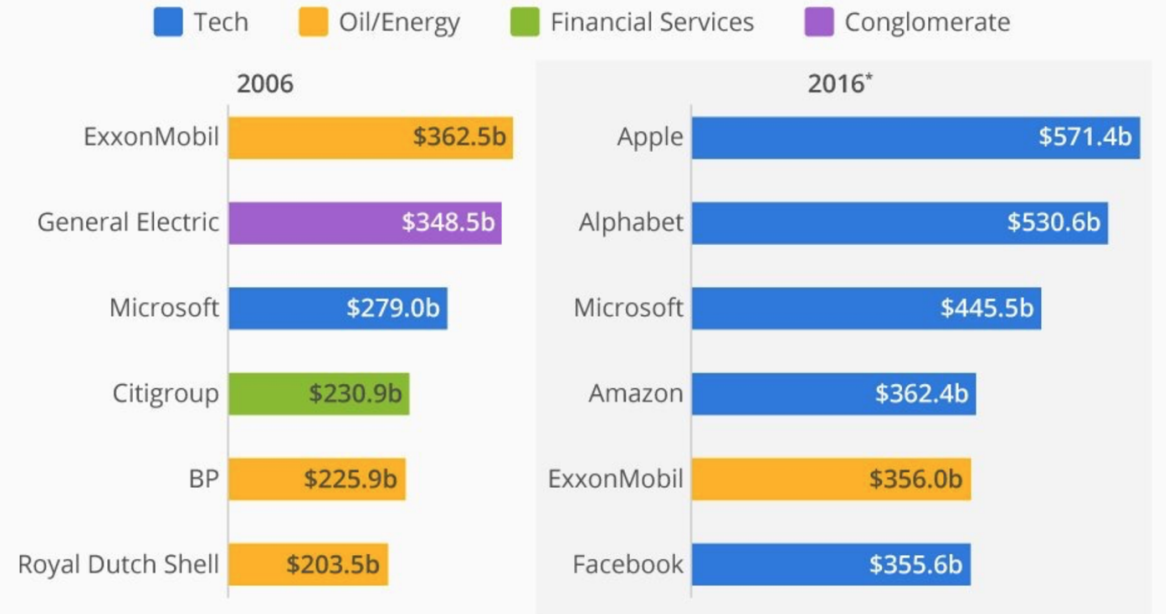
1010 1010100

# Data : the new oil



## The Age of Tech

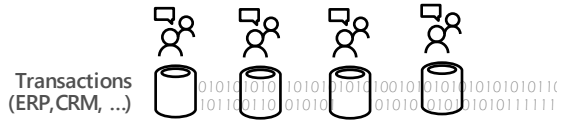
Market capitalization of the world's most valuable public companies



\* as of August 1, 2016

Sources: Yahoo! Finance, Forbes

# Data : but new challenges



**40 ZETTABYTES**  
of data will be created by 2020, an increase of 300 times from 2005



**6 BILLION PEOPLE**  
have cell phones  
WORLD POPULATION: 7 BILLION



The New York Stock Exchange captures **1TB OF TRADE INFORMATION** during each trading session



## THE 4 V'S OF BIG DATA

### Volume

SCALE OF DATA

**2.5 QUINTILLION BYTES**  
of data are created each day



Most companies in the U.S. have at least **100 TERABYTES** of data stored



Modern cars have close to **100 SENSORS** that monitor items such as fuel level and tire pressure



As of 2011, the global size of data in healthcare was estimated to be **150 EXABYTES**



**30 BILLION PIECES OF CONTENT** are shared on facebook every month



**1 IN 3 BUSINESS LEADERS** don't trust the information they use to make decisions



### Variety

DIFFERENT FORMS OF DATA

**4 BILLION + HOURS OF VIDEO** are watched on YouTube each month



**4 MILLION TWEETS** are sent per day by about 200 million monthly active users



### Veracity

UNCERTAINTY OF DATA

**27% OF RESPONDENTS** in one survey were unsure of how much of data was inaccurate



# Data : different types

Transactions  
(ERP, CRM, ...)



**STRUCTURED** =

Static schema, easy to insert into RDBMS and request with Structured Query Language

Graph



**UNSTRUCTURED** =

Dynamic schema, time and resources consuming (hard) to insert into Relation DB

Image



Social



IoT



Cloud



# Data : how to store it ?

**STORE**



Data Warehouse

Design : Transactions

: Consistent

: easy to request

Data Warehouse

**STRUCTURED**

Transactions  
(ERP, CRM, ...)



**Unstructured data**

The university has 5600 students. John's ID is number 1, he is 18 years old and already holds a B.Sc. degree. David's ID is number 2, he is 31 years old and holds a Ph.D. degree. Robert's ID is number 3, he is 51 years old and also holds the same degree as David, a Ph.D. degree.


**Semi-structured data**

```
<University>
  <Student ID="1">
    <Name>John</Name>
    <Age>18</Age>
    <Degree>B.Sc.</Degree>
  </Student>
  <Student ID="2">
    <Name>David</Name>
    <Age>31</Age>
    <Degree>Ph.D. </Degree>
  </Student>
  ....
</University>
```

**Structured data**

| ID | Name    | Age | Degree |
|----|---------|-----|--------|
| 1  | John    | 18  | B.Sc.  |
| 2  | David   | 31  | Ph.D.  |
| 3  | Robert  | 51  | Ph.D.  |
| 4  | Rick    | 26  | M.Sc.  |
| 5  | Michael | 19  | B.Sc.  |

**STORE**



Big data




**UNSTRUCTURED**

Hadoop/NoSQL (Big Data)

Design : Flexible

: Scalable

: cheap for perf

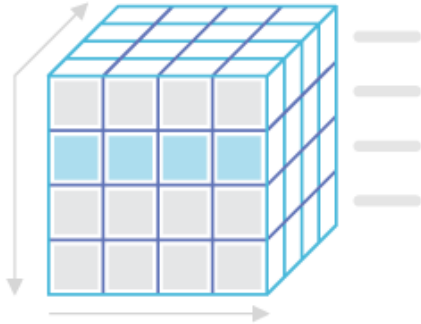
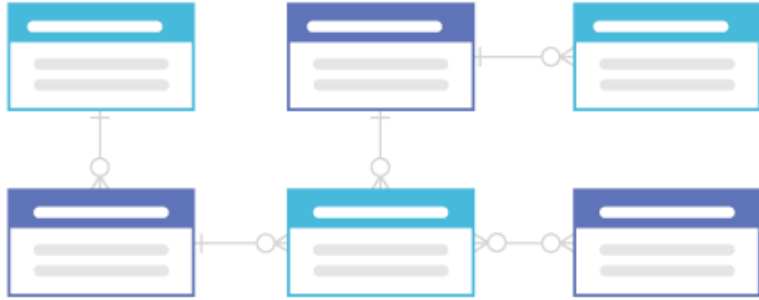
- Graph 
- Image 
- Social 
- IoT 
- Cloud 

SQL

NO ONE FIT ALL

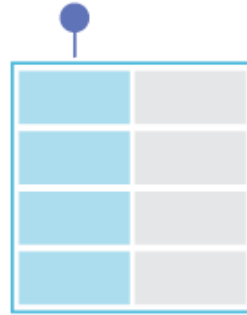
NoSQL

*Relational Database Management Systems (RDBMS)*

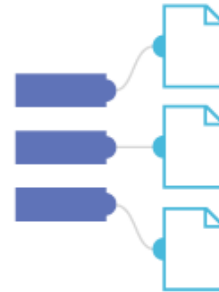
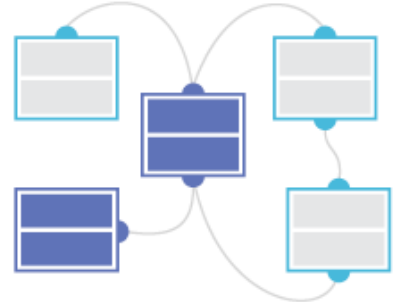


*Online Analytical Processing (OLAP) Cube*

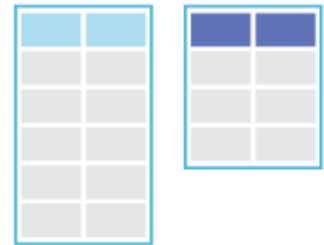
*Key-Value*



*Graph*



*Document*



*Column store*

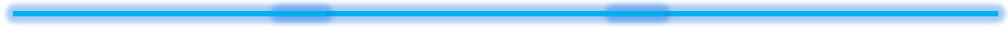
# > Les offres du marché BI (non exhaustif)

Analytics  
Appli & Tools

Advanced & Predictive Analysis      Corporate Performance Management      Planning, Budgeting & Forecasting



IBM (Cognos) Oracle SAP (BO) MicroStrategy QlikTech Tableau .....



Predictive, Data Mining      Query & Reporting      Dash boards, Scorecards, Visualization

Data  
Warehouse



IBM (DB2) Oracle, HP (Vertica) SAP (HANA) Teradata, .....



ETL

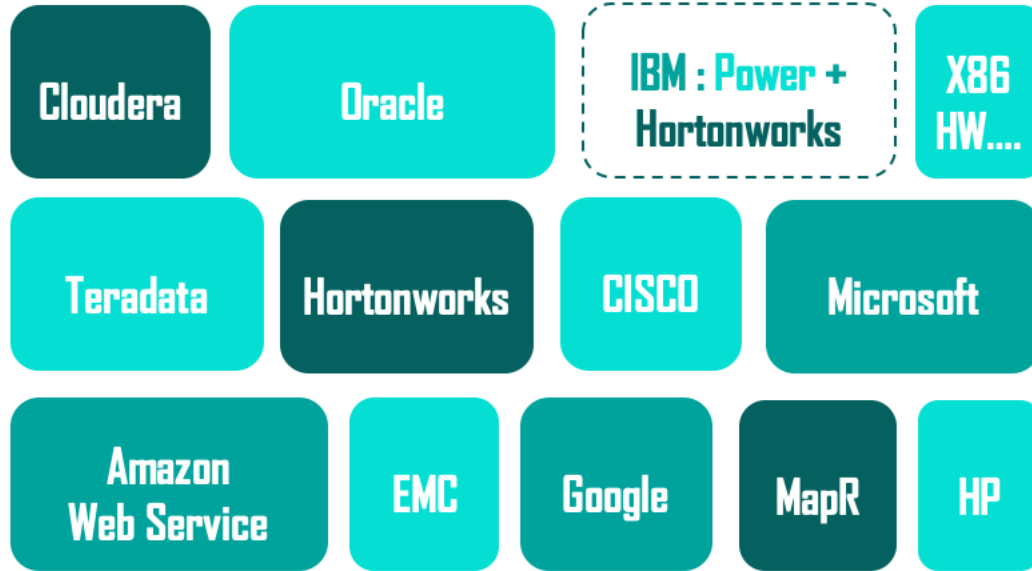


IBM (Infosphere Server) Oracle (Oracle Data Integrator) SAP (BO Data Services)  
Informatica





# L'offre du marché hadoop (soft + hardware)



 Hadoop related - Cloud

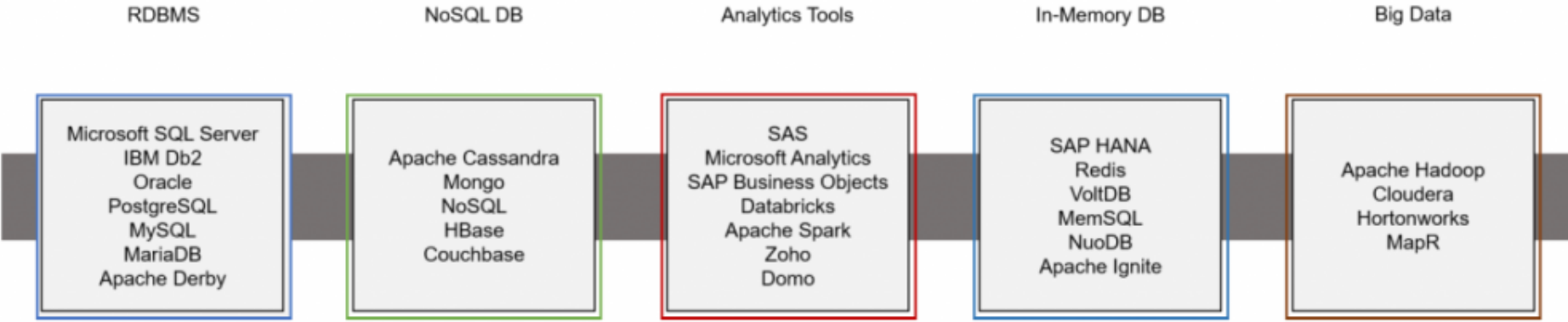
 Hadoop distributions

 Hadoop related - Hardware

*Non exhaustive !!!!*

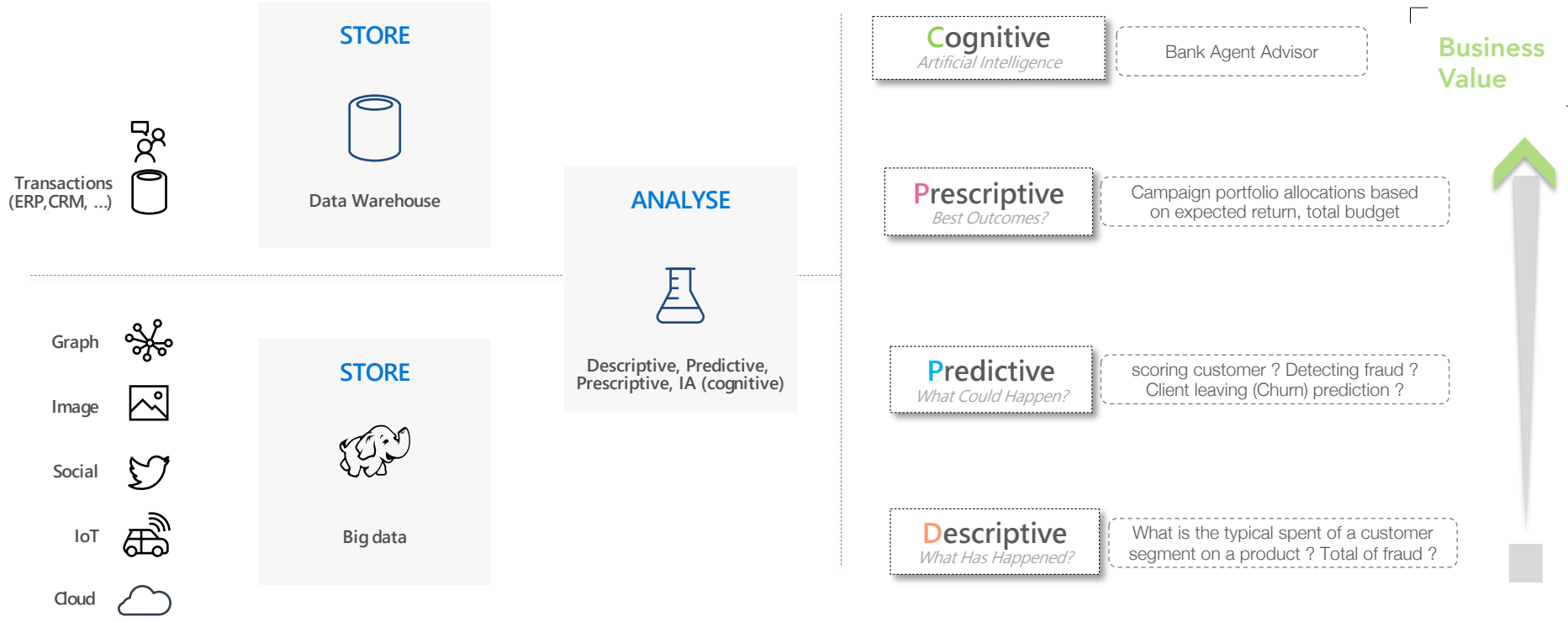
# L'offre DATA STORE : vue globale

## FIGURE 2 – THE DATA MANAGEMENT LANDSCAPE

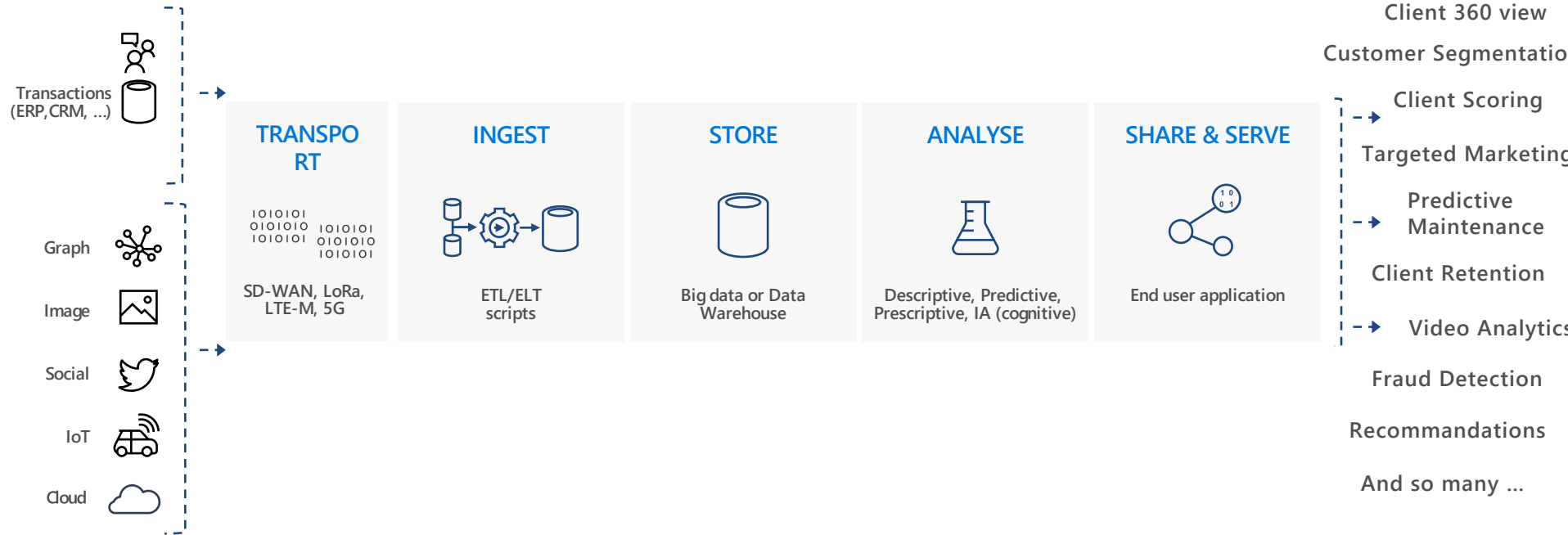


Source: Moor Insights & Strategy

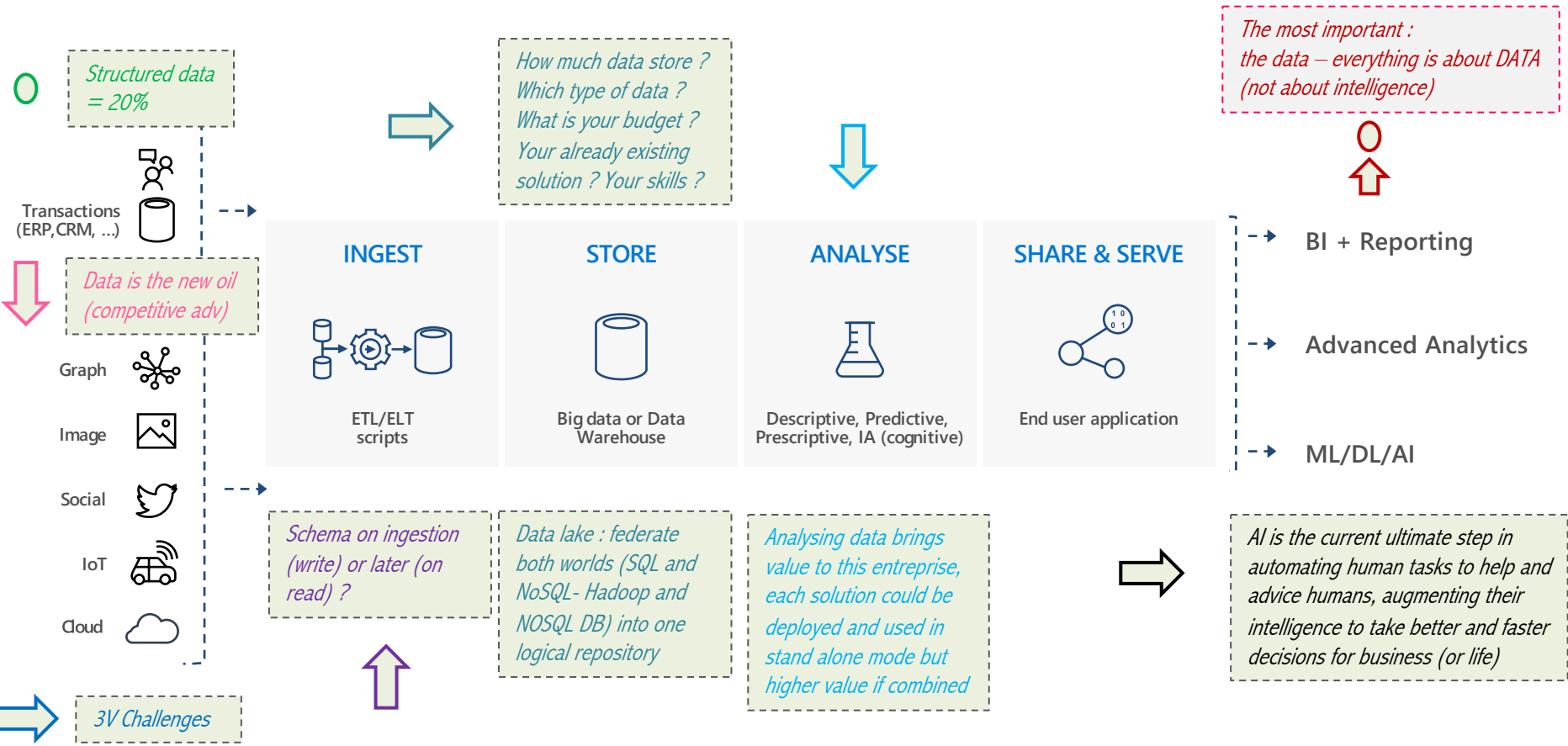
# Data : different analysis for different values



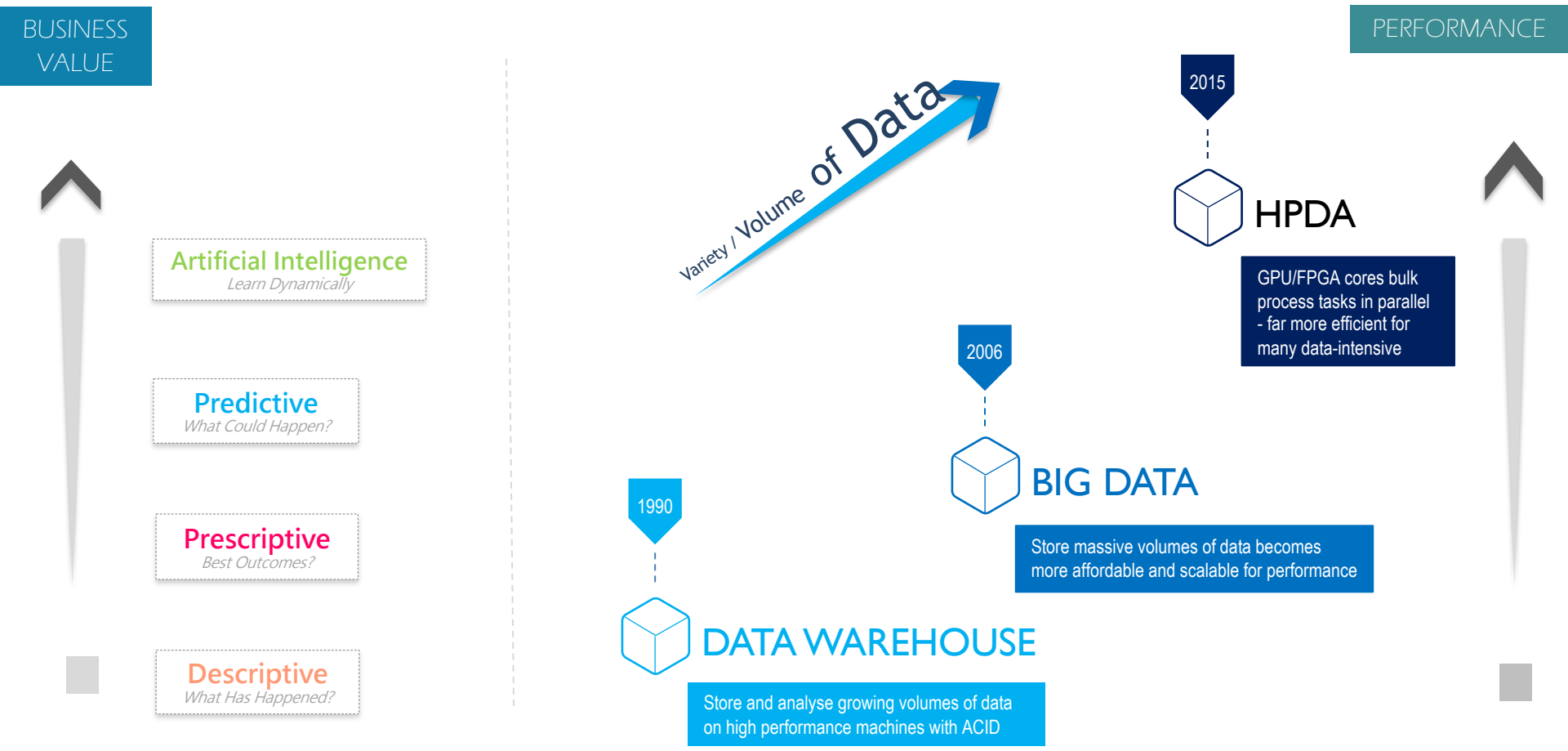
# The Data's Journey



# Le voyage de la donnée : TOUT RESUMER



# L'analyse de la donnée : enjeux business et IT



# Exemple de question

- Qu'est-ce qu'un entrepôt de données ?
- Quels sont les grands principes de la modélisation multidimensionnelle ?
- A quoi sert un ETL ?
- C'est quoi les 3V ?
- Qu'est-ce qu'une base noSQL ?
- Qu'est-ce que Hadoop ?
- Qu'est-ce que MapReduce ?
- Qu'est-ce qu'une appliance Big Data ?