



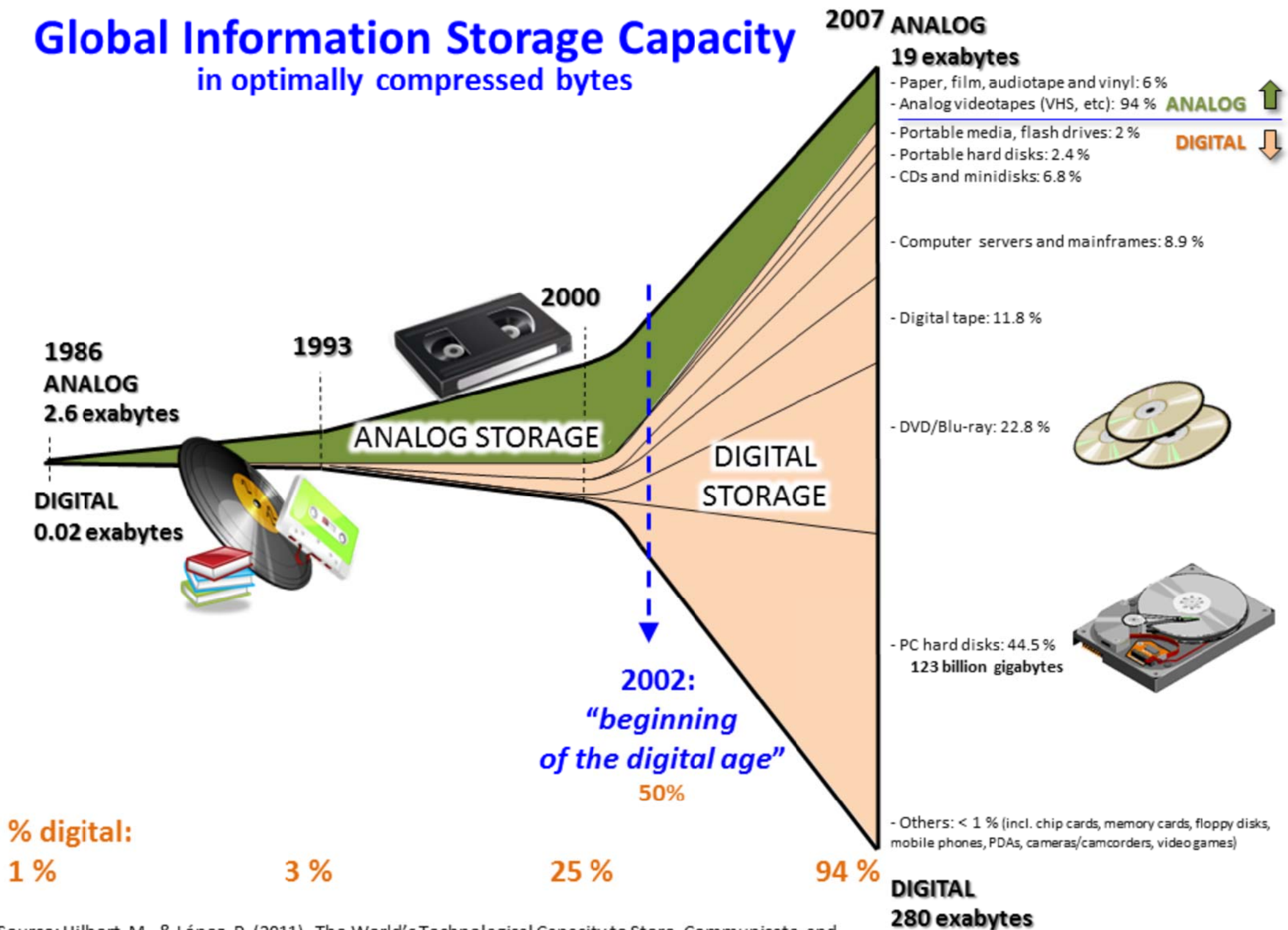
# Cursus

- DUT Informatique de Bayonne :) (2009-2011)
- INSA de Lyon 2011-2014)
- KAIST – Corée du Sud (2013)
  
- Amesys puis racheté par ATOS (2014-2016)
  - Wyplay → Ingénieur Middleware (C++/Python)
  - Orange → Ingénieur Big-Data (Java8, Hadoop)
- Criteo (Novembre 2016 - ...)
  - Ingénieur Développement et Opération (DevOps)

# Sommaire

- Enjeux
- Compétences
- Outils
- Questions

# Global Information Storage Capacity in optimally compressed bytes



# Enjeux (2/3)

## TOUS LES 2 ANS

On crée autant d'informations que depuis la nuit des temps à 2003

## CHAQUE SECONDE

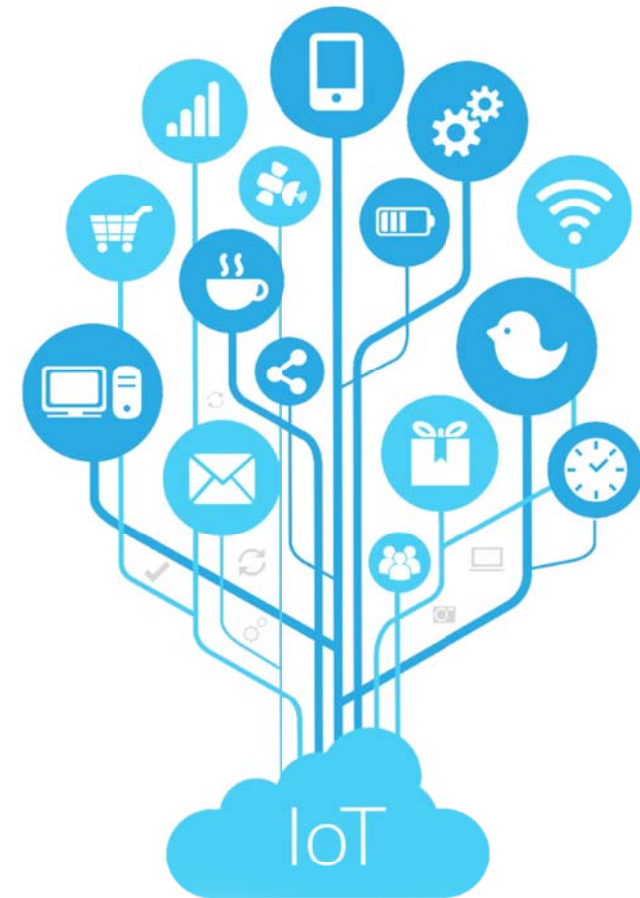
En 2020, chaque **humain** va générer **1.7 megaoctets** d'informations

## Maitriser le changement

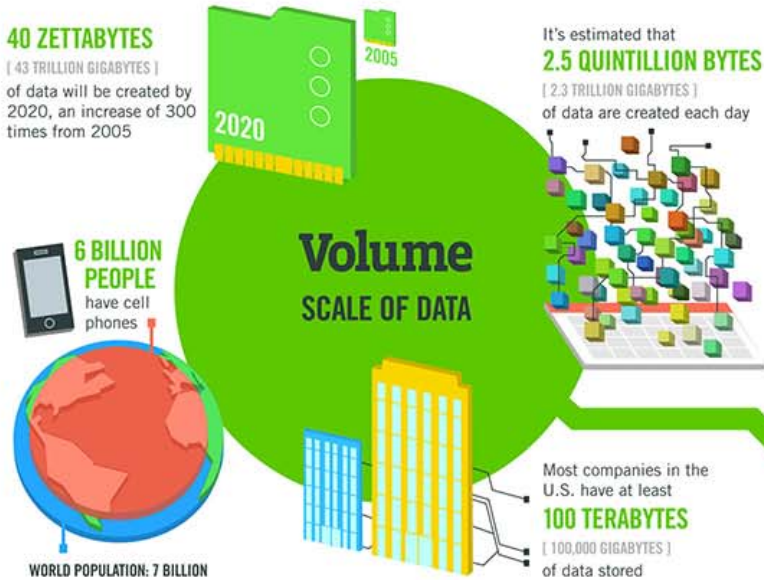
Comment se donner les moyens de traiter et de réagir à cette masse d'informations ?

# Enjeux (3/3)

- Encore plus de croissance à prévoir grâce aux
  - Objets connectés
    - Puces RFID pour suivre les objets
    - Capteurs en tous genres
  - Médecine
    - Bilan de santé
    - Analyses biologiques/ADN
  - Auto alimentation
    - + on automatise == + de données encore



# Compétences



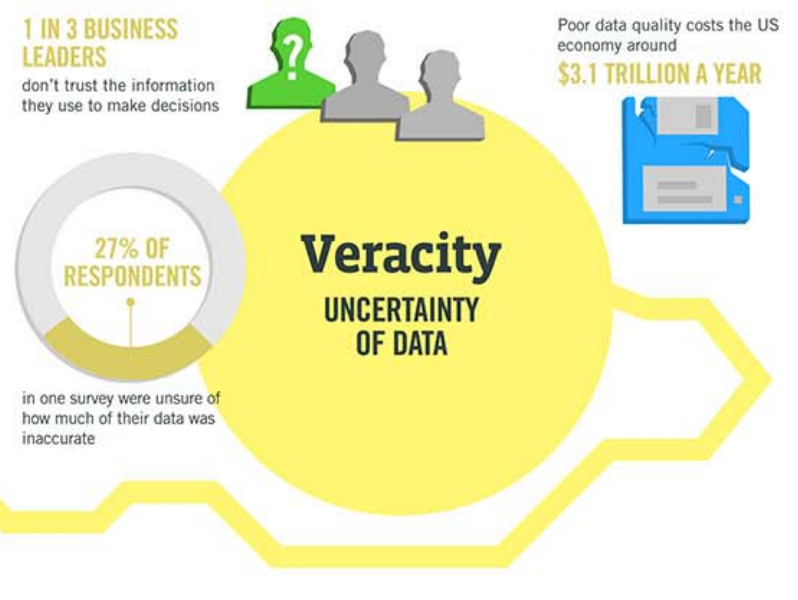
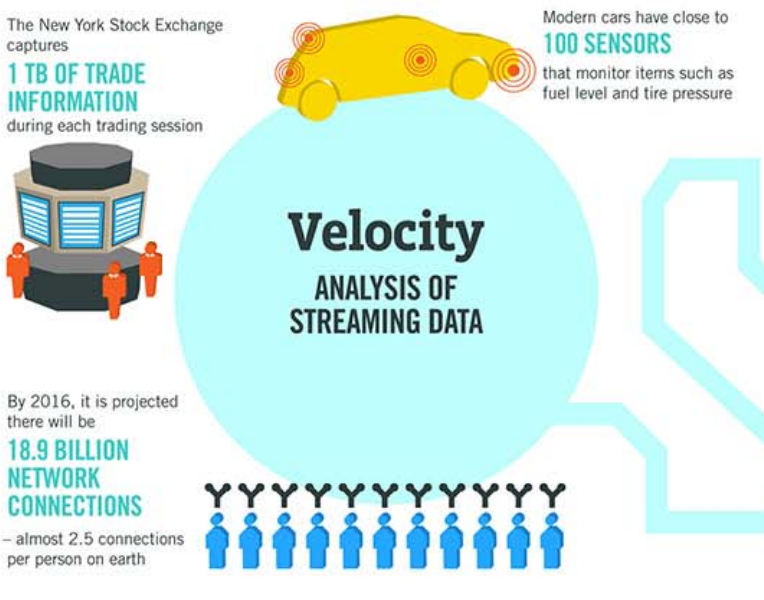
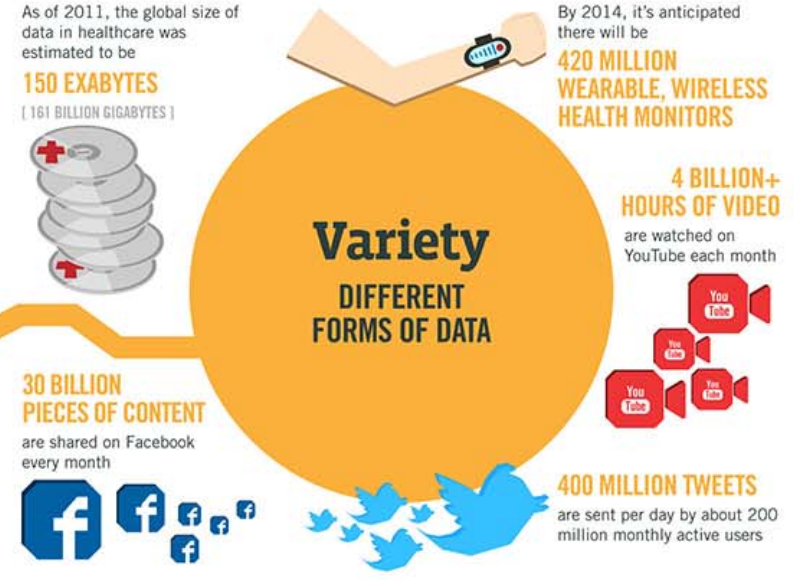
## The FOUR V's of Big Data

From traffic patterns and music downloads to web history and medical records, data is recorded, stored, and analyzed to enable the technology and services that the world relies on every day. But what exactly is big data, and how can these massive amounts of data be used?

As a leader in the sector, IBM data scientists break big data into four dimensions: **Volume, Velocity, Variety and Veracity**

Depending on the industry and organization, big data encompasses information from multiple internal and external sources such as transactions, social media, enterprise content, sensors and mobile devices. Companies can leverage data to adapt their products and services to better meet customer needs, optimize operations and infrastructure, and find new sources of revenue.

By 2015 **4.4 MILLION IT JOBS** will be created globally to support big data, with 1.9 million in the United States



# 40 ZETTABYTES

[ 43 TRILLION GIGABYTES ]  
of data will be created by 2020, an increase of 300 times from 2005



# 2.5 QUINTILLION BYTES

[ 2.3 TRILLION GIGABYTES ]  
of data are created each day



**6 BILLION PEOPLE**  
have cell phones



**WORLD POPULATION: 7 BILLION**

# Volume SCALE OF DATA



Most companies in the U.S. have at least

# 100 TERABYTES

[ 100,000 GIGABYTES ]  
of data stored



The New York Stock Exchange captures

**1 TB OF TRADE INFORMATION**

during each trading session



Modern cars have close to

**100 SENSORS**

that monitor items such as fuel level and tire pressure

# Velocity

ANALYSIS OF STREAMING DATA

By 2016, it is projected there will be

**18.9 BILLION NETWORK CONNECTIONS**

– almost 2.5 connections per person on earth



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

**150 EXABYTES**

[ 161 BILLION GIGABYTES ]



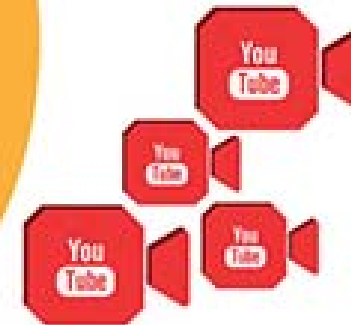
By 2014, it's anticipated there will be

**420 MILLION WEARABLE, WIRELESS HEALTH MONITORS**



**4 BILLION+ HOURS OF VIDEO**

are watched on YouTube each month



# Variety

## DIFFERENT FORMS OF DATA

**30 BILLION PIECES OF CONTENT**

are shared on Facebook every month



**400 MILLION TWEETS**

are sent per day by about 200 million monthly active users



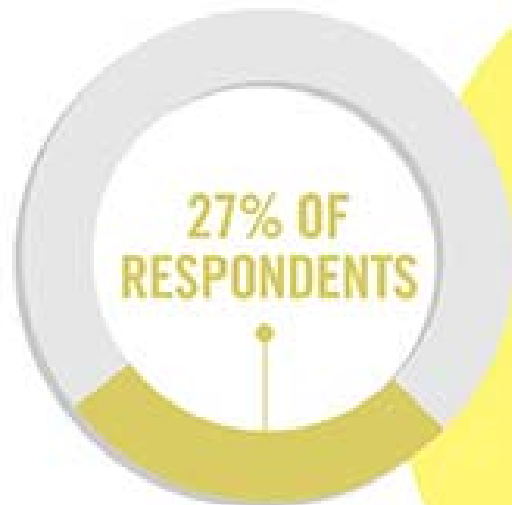
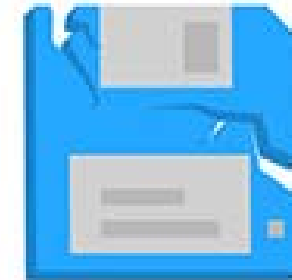
## 1 IN 3 BUSINESS LEADERS

don't trust the information they use to make decisions



Poor data quality costs the US economy around

**\$3.1 TRILLION A YEAR**



in one survey were unsure of how much of their data was inaccurate

# Veracity

## UNCERTAINTY OF DATA

# Outils (1/2)

- Système de fichiers distribués



- Système de calculs distribués



- Système de traitements évènementiels



# Outils (2/2)

- Système de croisement des données

-  

- Visualisation

- D3.js, Tableau, Apache Zeppelin



Questions ?



# BIG DATA



VOLUME

DATA SIZE



VELOCITY

SPEED OF CHANGE



VARIETY

DIFFERENT FORMS  
OF DATA SOURCES



VERACITY

UNCERTAINTY OF  
DATA