# Patrick Beaucamp

Founder of the Vanilla Project
Mail : Patrick.beaucamp@bpm-conseil.com

**Open Source Platforms to deploy Search and Maps
Visualization on top of a Big Data Database**

**II-SDV, Nice 15th April 2013**

**BPM-Conseil**
*the Company behind the Vanilla project*

II-SDV

# Presentation Agenda

**Big Data Use case**

    Various Application

**Big Data Landscape**

        Data Storage and Database Type
        PlayerMap - Hadoop
        Search & DataMining Tools
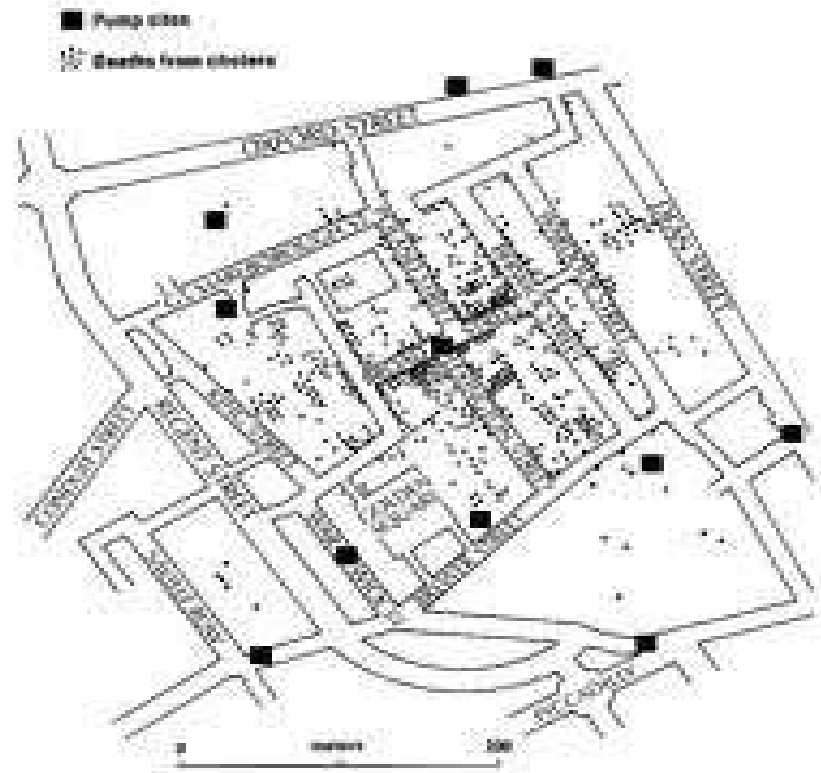        Data Visualisation

**Big Data Use case**

        Social Network
        Open Data

**BPM-Conseil**
*the Company behind the Vanilla project*

# Various Applications

**Discovery of Cholera – 1854 (John Snow)**

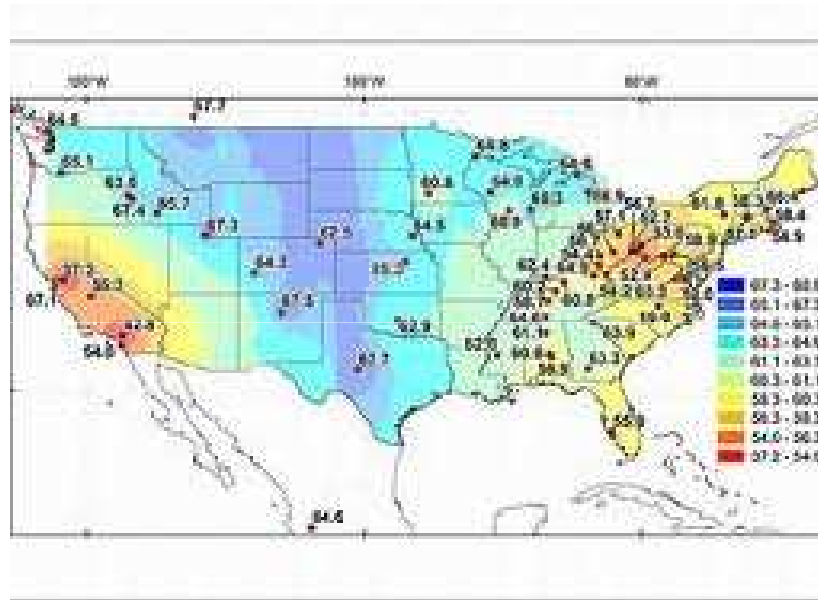http://en.wikipedia.org/wiki/1854_Broad_Street_cholera_outbreak

# Various Applications

Bicycle Accident in Street : who is taking care of trafic management
Example in Boston :

http://www.boston.com/bostonglobe/editorial_opinion/blogs/the_angle/2010/12/bike_crash_map.html

BPM-Conseil
the Company behind the Vanilla project
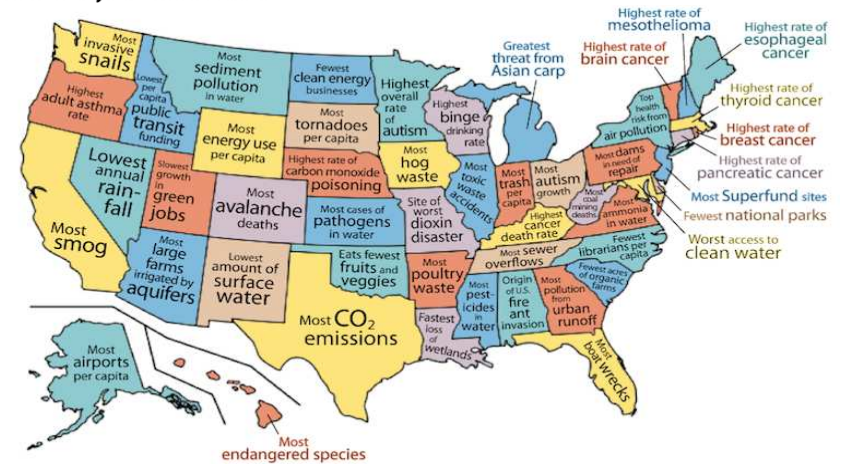
# Various Applications

Pollution data with electoral data : who is taking care of pollution
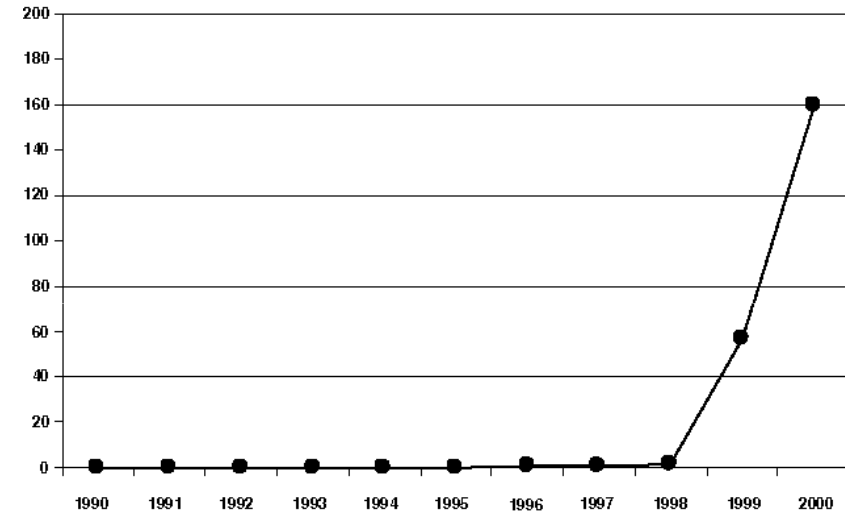
# 1 - Big Data Landscape - Storage

## Storage … what for ?



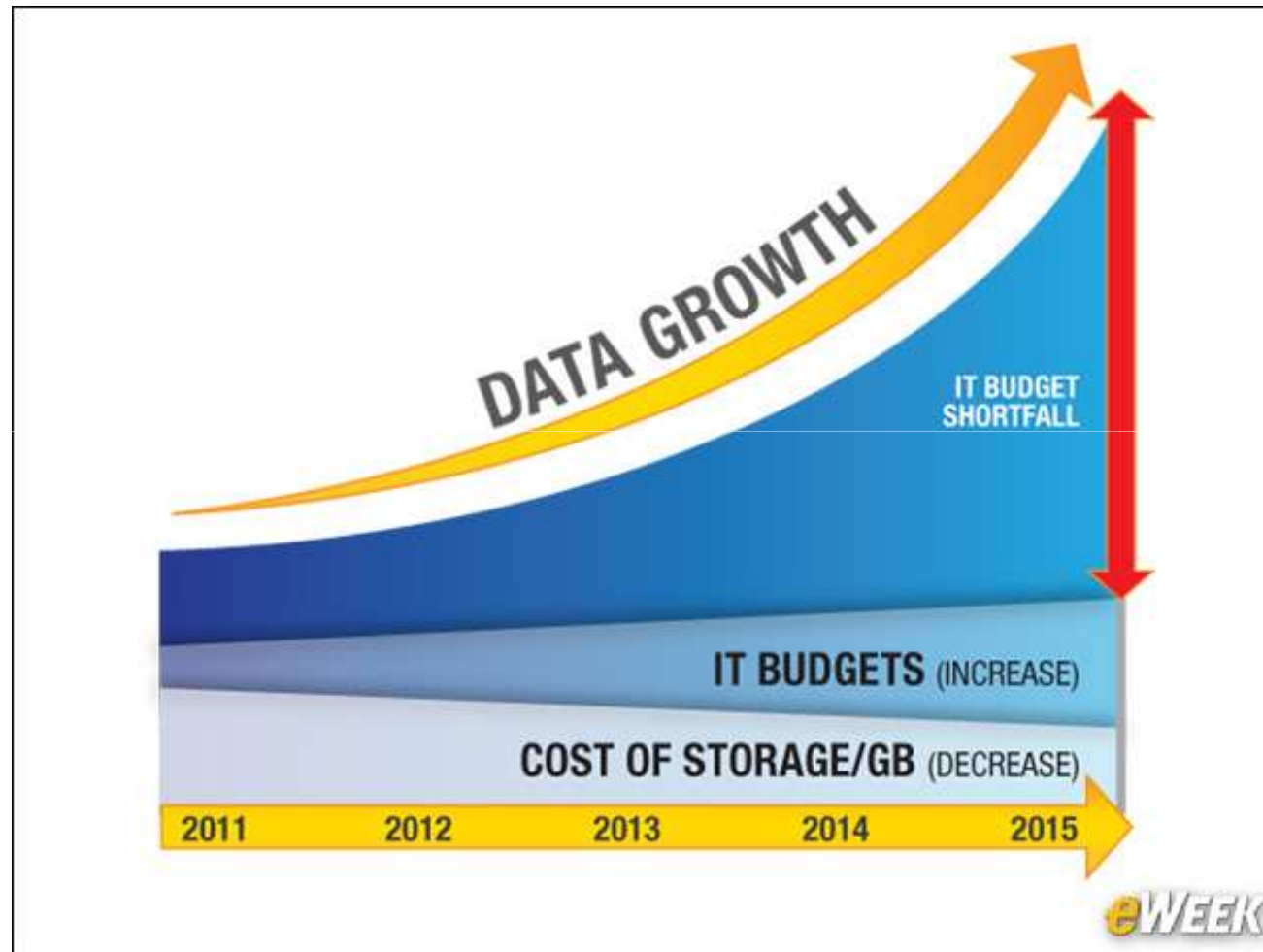The red button in a IBM 3380 cabinet is as big as three MicroSD cards.

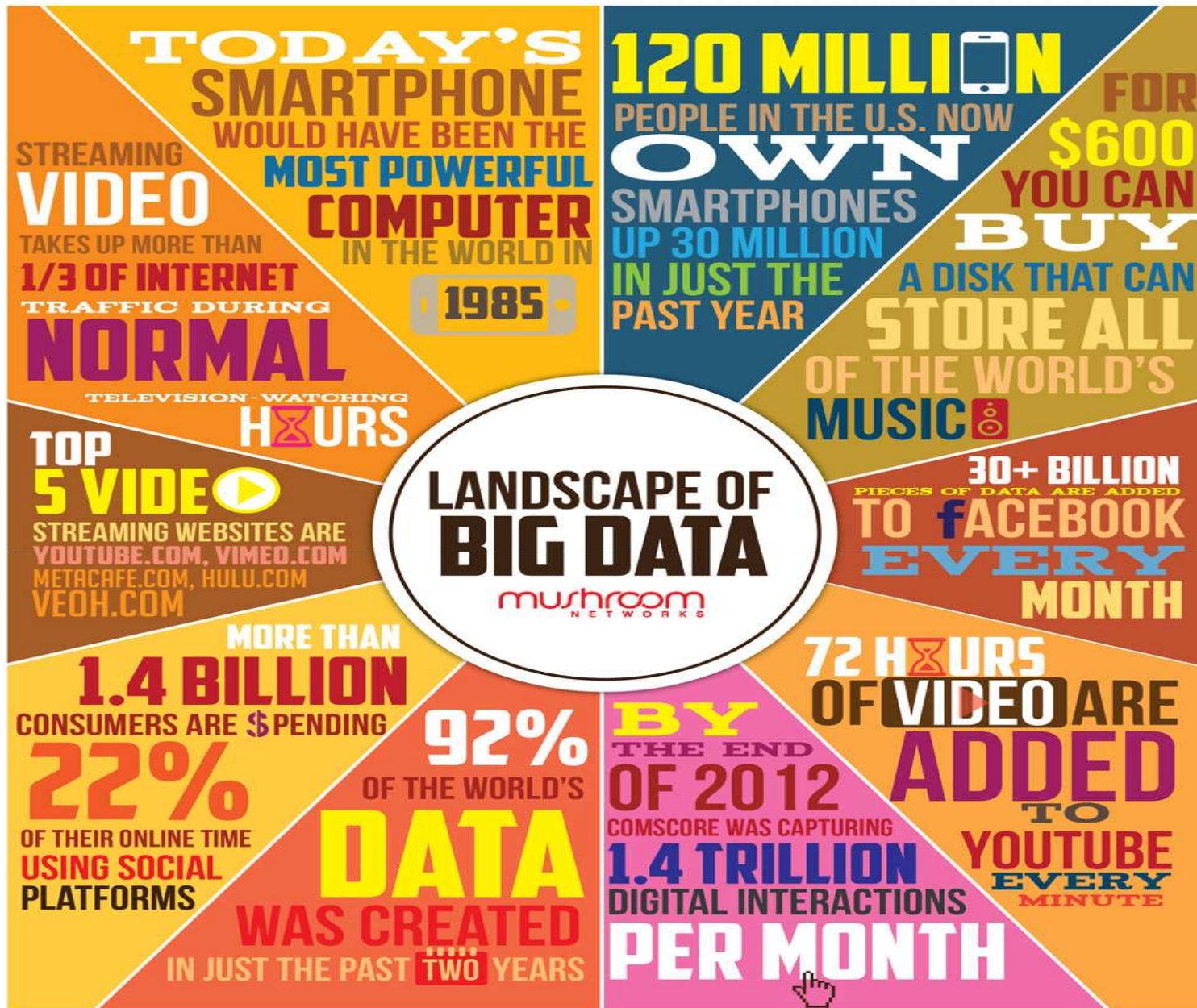1980 — Eight 2.5GB IBM 3380 Disk Systems: 20GB
Estimated value: $648,000 - $1,137,600
Weight: 2,000,000 grams (4,400 pounds)

2010 — One MicroSD Card: 32GB
Estimated value: $100 - $150
Weight: 0.5 grams (0.001 pounds)

**BPM-Conseil**
the Company behind the Vanilla project

# 1 - Big Data Landscape - Storage

LANDSCAPE OF **BIG DATA**
mushroom NETWORKS

TODAY'S SMARTPHONE WOULD HAVE BEEN THE MOST POWERFUL COMPUTER IN THE WORLD IN 1985

120 MILLION PEOPLE IN THE U.S. NOW OWN SMARTPHONES UP 30 MILLION IN JUST THE PAST YEAR

FOR $600 YOU CAN BUY A DISK THAT CAN STORE ALL OF THE WORLD'S MUSIC

STREAMING VIDEO TAKES UP MORE THAN 1/3 OF INTERNET TRAFFIC DURING NORMAL TELEVISION-WATCHING HOURS

TOP 5 VIDEO STREAMING WEBSITES ARE YOUTUBE.COM, VIMEO.COM METACAFE.COM, HULU.COM VEOH.COM

30+ BILLION PIECES OF DATA ARE ADDED TO FACEBOOK EVERY MONTH

MORE THAN 1.4 BILLION CONSUMERS ARE $PENDING 22% OF THEIR ONLINE TIME USING SOCIAL PLATFORMS

92% OF THE WORLD'S DATA WAS CREATED IN JUST THE PAST TWO YEARS

BY THE END OF 2012 COMSCORE WAS CAPTURING 1.4 TRILLION DIGITAL INTERACTIONS PER MONTH

72 HOURS OF VIDEO ARE ADDED TO YOUTUBE EVERY MINUTE

**BPM-Conseil**
the Company behind the Vanilla project

# 1 - State of art of different database technologies (1/6)

**Sql Database**

Database that complies with SQL language and its sub-definition like pl/sql, transact-sql …, where data are structured in tables and columns

Usage : Erp, transactional application (Oltp), any kind of application with structured datas

Famous one : Oracle, Sybase, SqlServer, PostgreSql, Mysql

# 1 - State of art of different database technologies (2/6)

**Column Database (vertical database)**

Database that complies with SQL language where index strategy along with data compression are the key subjects

Usage : Bi Applications – Log Analysis application

Famous one : SybaseIQ, Vertica, Infobright, VectorWise, Exadata

# 1 - State of art of different database technologies (3/6)

**Olap Database (Cubes)**

Database that complies with MDX language (and its Xmla implementation) where data is structured with dimensions and measures

Usage : Bi Applications (Olap)

Famous one : Hyperion Essbase, Cognos PowerPlay, Sap Bw

# 1 - State of art of different database technologies (4/6)

**NoSql Database (Hadoop world, BigData)**

Database that provides support for any kind of datastorage for structured data, unstructured data and any kind of document (image, message). Development language specific to read/write data (Api)

Usage : Large unstructured Web application (Twitter, Facebook : could be view like table with millions of columns)

Famous one : Cassandra, Hbase, MongoDb, Google BigTable

# 1 - State of art of different database technologies (5/6)

**Application Database (Hidden structure)**

Database that provides access & exchange to their data through public Api, keeping hidden or too-complex the database structure.

Usage : Any kind of application, usually large Erp and large Web based cloud application

Warning : Api licence & usage (twitter example)

Famous one : GoogleAnalytics, Sap-Bapi (even Sql behind), twitter (Database & Application)

# 1 - State of art of different database technologies (6/6)



Twitter Unlinks From LinkedIn

forbes.com - Social media addicts may be dismayed to learn today that Twitter updates will no longer automatically sync with LinkedIn accounts. In a blog post earlier today, Ryan Roslansky, Head of Content Products at LinkedIn, noted that...

Trending within the following companies

LinkedIn          IBM

**BPM-Conseil**
*the Company behind the Vanilla project*

# Big Data Landscape

# 2 - Database proposal in big data : nosql+column database, Hadoop (1/6)

Hadoop is the natural platform for BigData, providing the fundation and the Framework.

Hadoop software library is a framework that allows for the distributed processing of large data sets across clusters of computers using a simple programming model

# 2 - Database proposal in big data : nosql+column database, Hadoop (2/6)

Column Database position is difficult, as Hadoop Hbase provides column indexation services. Column database editor trying to push « BigData offer » with some Hadoop components

3 majors databases for NoSql : Cassandra, Hbase and MongoDb, 2 of the 3 (Cassandra, Hbase) under the Apache fundation umbrella

# 2 - Database proposal in big data : nosql+column database, Hadoop (3/6)

**Why Hadoop is unique**



It's an Open Source Framework, without any commercial competitor at that level

New companies endorse Hadoop to build stable distribution with major Hadoop components : HortonWorks, MapR, Cloudera



Hadoop architecture target no single point of failure (multi master) , with usefull component (Zookeeper, MapReduce)

BPM-Conseil
the Company behind the Vanilla project

# 2 - Database proposal in big data : nosql+column database, Hadoop (4/6)

## Hadoop Framework to manage big data set on the Cloud

# 2 - Database proposal in big data : nosql+column database, Hadoop (5/6)

**Hadoop Components**

It was inevitable that Hadoop would soon evolve
to adopt some of the characteristics of a database

HBase, stepped in to partially fill the gap. It implements a
column-oriented data store modeled on Google's Bigtable on
top of Hadoop and HDFS

Pig, originally developed at Yahoo Research, is a high-level
language for building MapReduce programs for Hadoop,
thus simplifying the use of MapReduce. It is a data flow
language that provides high-level commands.

# 2 - Database proposal in big data : nosql+column database, Hadoop (6/6)

**Some Hadoop Components**

Hive, initially a sub-project of Hadoop, evolved to provide a formal query capability. In effect, Hive turns Hadoop into a data warehouse-like system, allowing data summarization, ad hoc queries and the analysis of data stored by Hadoop

Hive holds metadata describing the contents of files and allows queries in HiveQL, a SQL-like language. It also allows MapReduce programmers to get around the limitations of HiveQL by plugging in MapReduce routines.

# 3 - Big Data, Search & Datamining

-Work on your data with Dataming component such as

-R : http://www.revolutionanalytics.com/
- Weka : http://www.cs.waikato.ac.nz/ml/weka/
- Apache Mahout : http://mahout.apache.org/



- Manage different kind of data, using Hadoop Cassandra database

-Use existing Search Infrastructure like Solr/Lucene (Vanilla certified)
http://www.lucidworks.com/

# 3 - Big Data, Search & Vanilla

*Response Time is stable with Cassandra – Impressive !!!*

# 4 - Big Data & Data Visualisation

*History of data visualisation*

# 4 - Big Data & Data Visualisation

# 4 – Clustering & Data Mining



done

Big Data

**Adjecent Nodes:**
- root
- David Milward
- Patrick Beaucamp
- Renaud Garat
- Roger Bradford
- Steve Kearns

root
Big Data
David Milward
Patrick Beaucamp
Renaud Garat
Roger Bradford
Steve Kearns

## Règle de Bayes

$$\Pr(k / x) = \frac{\pi_k L_k(x)}{L(x)}$$

probabilité a posteriori

posteriori F
posteriori H

Région de F

Région de H

BPM-Conseil
the Company behind the Vanilla project

# 4 - Big Data & Data Visualisation

http://www.gapminder.org

# 1 - Social Network

- Facebook – Linked – Twitter
  Your Private Life – Your Professional Life – What you think

  Intensive usage of BigData database (Cassandra) and Platform (Hadoop)
  Unstructured Data (chat) – Image and Video (face recognition)

  Just rethink what they do with your data

  http://www.yesprofile.com

# 2 - Open Data Global Site

http://www.opengovdata.org/

*The 8 Principles of Open Government Data*

*1. Data Must Be Complete*
*2. Data Must Be Primary*
*3. Data Must Be Timely*
*4. Data Must Be Accessible*
*5. Data Must Be Machine processable*
*6. Access Must Be Non-Discriminatory*
*7. Data Formats Must Be Non-Proprietary*
*8. Data Must Be License-free*

# 2 - Open Data Government Site (1/4)

French Portal : http://www.data.gouv.fr/ - 352'431 dataset

# 2 - Open Data Government Site (2/4)

US Portal : http://www.data.gov/ - 420'894 dataset

# 2 - Open Data Government Site (3/4)

UK Portal : http://data.gov.uk/ - 8'610 dataset

# 2 - Open Data Government Site (4/4)

Russian Portal : http://opengovdata.ru/ - ?? dataset

# 2 - Open Data Government - Bonus

European Decision on Open Data for Government (April 10)

  https://ec.europa.eu/digital-agenda/en/public-sector-information-raw-data-new-services-and-products

  http://gigaom.com/2013/04/10/european-governments-agree-to-open-up-public-data/

II-SDV, Nice

European Decision on Mine (Diamond, Gold, Uranium), Gaz, Oil (April 13)

  http://eurodad.org/wp-content/uploads/2011/11/CBC-report_french-21.pdf

**BPM-Conseil**
the Company behind the Vanilla project