# BIG DATA?

# THE GLOBAL IMBALANCE!

Stéphane Grumbach
INRIA

# The digital universe   2.7 Zettabytes

## Data deluge in all sectors of activity

**kilo** $10^3$
**mega** $10^6$
**giga** $10^9$
**tera** $10^{12}$
**peta** $10^{15}$
**exa** $10^{18}$
**zetta** $10^{21}$
**yotta** $10^{24}$

U.S. Library of Congress: 235 Terabytes of data

Walmart: 2.5 petabytes of data, 1 million customer transactions / hour

Facebook: 30 Petabytes of user data

Google: processing 20 petabytes a day (2008)

World: 5 billion people calling, tweeting, browsing on mobile phones

## Exponential increase

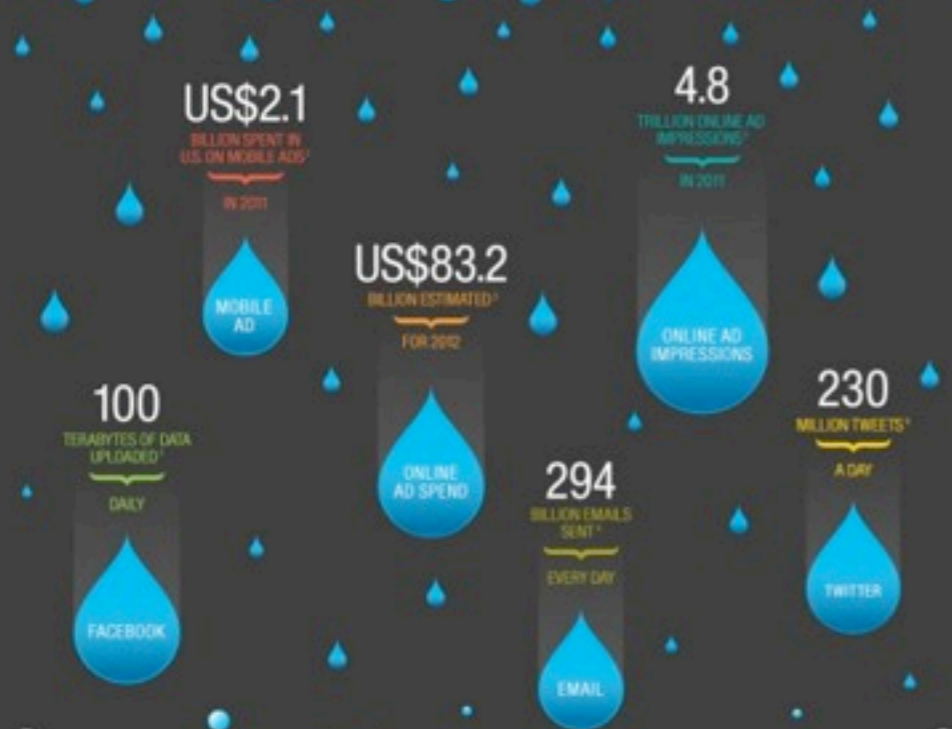doubles every two years   **35 zettabytes in 2020**

followed by the capacity to store, compute, and communicate

# The Big Data Industry

## **Advertising**

Capture users data

Generate users profiles

Target ads

# The Big Data Industry

beyond advertising

- $300 billion/year
  US health care
- €250 billion/year
  Europe public administration
  [McKinsey 2011]

Tremendous economic impact
Teraeuros (thousands billions)

4

# First challenge: Data Harvesting

**70% of the data produced by individuals**

directly produced by users:

email, photos, blogs, etc. (less than half)

indirectly digital shadow/footprint:

surveillance, web usage, transactions

## The free paradigm of the 2.0

Free services traded for private user data

Free exploitation of the accumulated data

# Second challenge: knowledge extraction

User profiles (business)

=> Ads target

Automatic discovery (science)

=> Google Flu
monitoring of flu related queries

*a search engine company knows everything*

=> Biological, sociological data...



NSA (security)

=> Ambition to handle yottabytes ($10^{24}$) !!!

# Data: raw material of the 21st century (much like crude oil)



Major trade movements
Trade flows worldwide (million tonnes)

extraction from natural reservoirs

transport

refining

consumption at users

data analytics

accumulation in large repositories

Internet

production of data at users

244.2
126.1
25.1
50.1
208.4
102.0
70.8
36.4
154.3
51.8
19.7
90.6
333.7
34.0
28.8
81.5
120.9
35.3
36.0

USA
Canada
Mexico
S. & Cent. America
Europe & Eurasia
Middle East
Africa
Asia Pacific

# Where are these data?

Huge concentration of data

85% of data handled by (large) corporations
Virtualization/dematerialization of infrastructures
Social networks, Cloud, ...

Most of the prominent corporations based in the USA
Google, Facebook, Amazon, Twitter, ...
Storage capacity of Europe = 70% USA [McKinsey 2011]

1/3 of world data stored in the cloud by 2020

# Geopolitics of big data

Data from the Web 2.0

    produced by users everywhere in the world

    but accumulated by corporations most often abroad

Percentage of national web corporations among top 25 by country

- USA: 100%
- China: 92% (only Google makes it in the top 25)
- France: 36% (but mostly marginal sites, not data intensive)
  *leboncoin, Orange, Free, commentcamarche, lemonde, lequipe, lefigaro, pagesjaunes, sfr*

Alexa.com

# Geopolitics of big data

## The Top 50 websites worldwide

- USA: 72 %
- China: 16 % (Baidu: 5; QQ: 8; Taobao: 13; Sina:17; 163: 28; Soso:29; Sina weibo:31; Sohu:43)
- Russia: 6 % (Yandex: 21; kontakte:30; Mail: 33; )
- Israel: 2 % (Babylon: 22)
- UK: 2 % (BBC: 46)
- Netherland: 2 % (AVG: 47)

# Geopolitics of big data

Diversity of search engines

- USA: Google: 65 % ; Bing: 15% ; Yahoo: 15%
- China: Baidu: 78% ; Google: 16%
- Russia: Yandex: 60% ; Google: 25%
- UK: Google: 91 % ; Bing: 5%

- France: Google: 92 % ; Bing: 3%

In France,
- Google has a de facto monopoly
- Google knows more about France than INSEE

# The global imbalance

→ Information asymmetry

"Since asymmetries of information give rise to market power, and perfect competition is required if markets are to be efficient, it is perhaps not surprising that markets with information asymmetries and other information imperfections are far from efficient."

JOSEPH E. STIGLITZ

# Impact of the global imbalance

## Regulation

What legislations over a dematerialized global industry?

Aren't the rules defined by those who have the control?

## Business

How to face monopolistic positions?

How to handle the information asymmetry?

## Security

Data at the core of nations independence

# The power of data



Map Ecological Footprint
http://www.csa.com/discoveryguides/china/review.php

14

# What's at stake in Europe?

Suspicion (fear?) regarding data
  concern for privacy protection high in Europe
  active legislative work
  historical reasons?

Weak industrial/innovation environment
  no strong corporation emerging

But essential dependence on foreign systems

# Are there alternatives?

dominant (centralized) model

unclear privacy

lost property

active (centralized) business

little share of business capacity



decentralized 'utopian' model

high privacy

real ownership

little business

Faroo, Yacy

Diaspora

dominant (centralized) model

unclear privacy
lost property
active (centralized) business
little share of business capacity

Google  f

an alternative path ?

active (competitive) business

symmetry of information

ownership & privacy

anti monopoly

decentralized 'utopian' model

high privacy
real ownership
little business

Faroo, Yacy

Diaspora

# An alternative path for Europe?

**The information society**

   it is only emerging

   it will continue to evolve

   it will impact political systems

   new business models, new equilibrium will appear

Europe should embrace the future

谢谢