

# Big-Data Tutorial

Marko Grobelnik  
[marko.grobelnik@ijs.si](mailto:marko.grobelnik@ijs.si)  
Jozef Stefan Institute  
Ljubljana, Slovenia

Stavanger, May 8<sup>th</sup> 2012

# Outline

- ▶ Introduction
  - What is Big data?
  - Why Big-Data?
  - When Big-Data is really a problem?
- ▶ Techniques
- ▶ Tools
- ▶ Applications
- ▶ Literature

# *Big data—a growing torrent*

**\$600** to buy a disk drive that can  
store all of the world's music

**5 billion** mobile phones  
in use in 2010

**30 billion** pieces of content shared  
on Facebook every month

**40%** projected growth in  
global data generated  
per year vs. **5%**  
growth in global  
IT spending

**235** terabytes data collected by  
the US Library of Congress  
by April 2011

**15 out of 17**  
sectors in the United States have  
more data stored per company  
than the US Library of Congress

# *Big data—capturing its value*

**\$300 billion**

potential annual value to US health care—more than double the total annual health care spending in Spain

**€250 billion**

potential annual value to Europe's public sector administration—more than GDP of Greece

**\$600 billion**

potential annual consumer surplus from using personal location data globally

**60%**

potential increase in retailers' operating margins possible with big data

**140,000–190,000**

more deep analytical talent positions, and

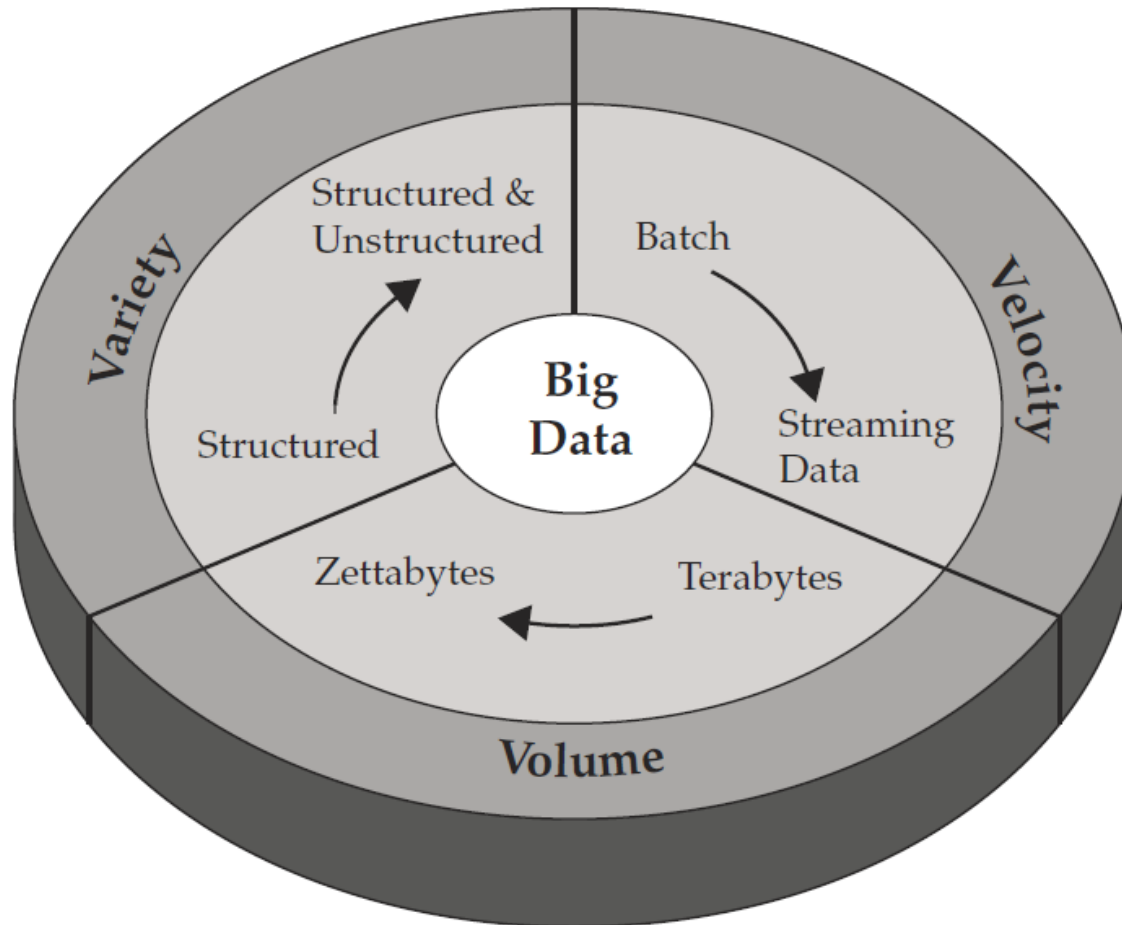
**1.5 million**

more data-savvy managers needed to take full advantage of big data in the United States

# What is Big-Data?

- ▶ ‘Big-data’ is similar to ‘Small-data’, but bigger
- ▶ ...but having data bigger consequently requires different approaches:
  - techniques, tools & architectures
- ▶ ...to solve:
  - New problems...
  - ...and old problems in a better way.

# Characterization of Big-Data: volume, velocity, variety (V3)





# Big-Data popularity on the Web

● big data ● data mining ● semantic web ● machine learning

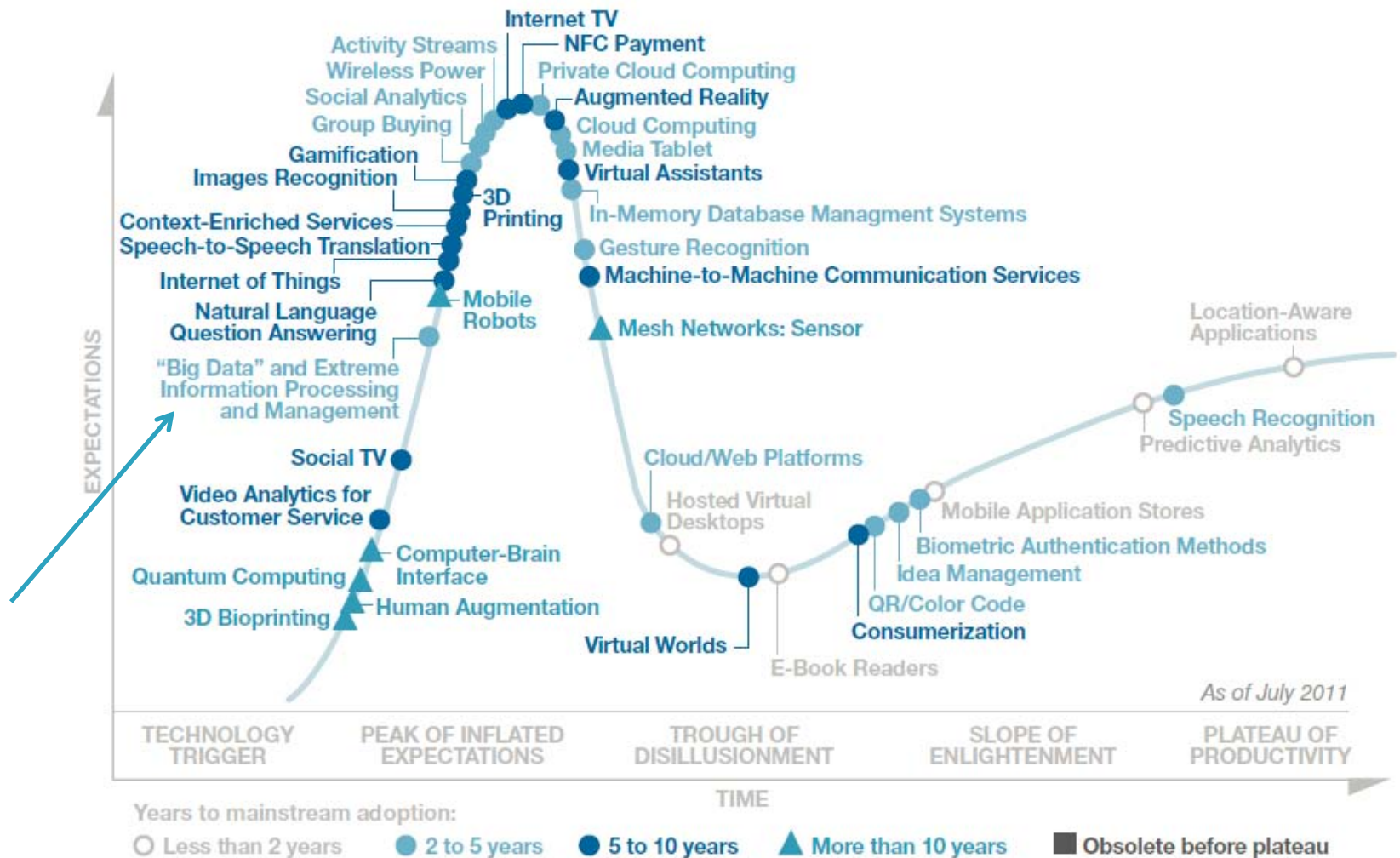


- A** [Spectra Logic Delivers ExaScale Storage for 'Big Data'; Announces Series of Products and Advancements and Unveils World's Highest Capacity Storage System](#)  
MarketWatch - Nov 1 2011
- B** [Webcast: Obama Goes Big on Big Data](#)  
Wired News - Mar 27 2012
- C** [Cisco Joins Forces with EMC to Advance IT Skills in Cloud, Big Data and Data Center Technologies](#)  
Justmeans - Apr 3 2012

- D** [Ferranti Unveils its MECOMS™ "Big Data" Strategy for Utility Meter Data Management and Real Time Billing](#)  
Victoria Times Colonist - Apr 10 2012
- E** [Deconstructing Big Data - BuildZoom Launches an Article Series that Reveals the Hype and Substance Behind Big Data](#)  
Houston Chronicle - Apr 17 2012
- F** [Harvard Releases Big Data for Books](#)  
New York Times - Apr 24 2012

# Big-Data in Gartner Hype-Cycle 2011

Hype Cycle for Emerging Technologies, 2011





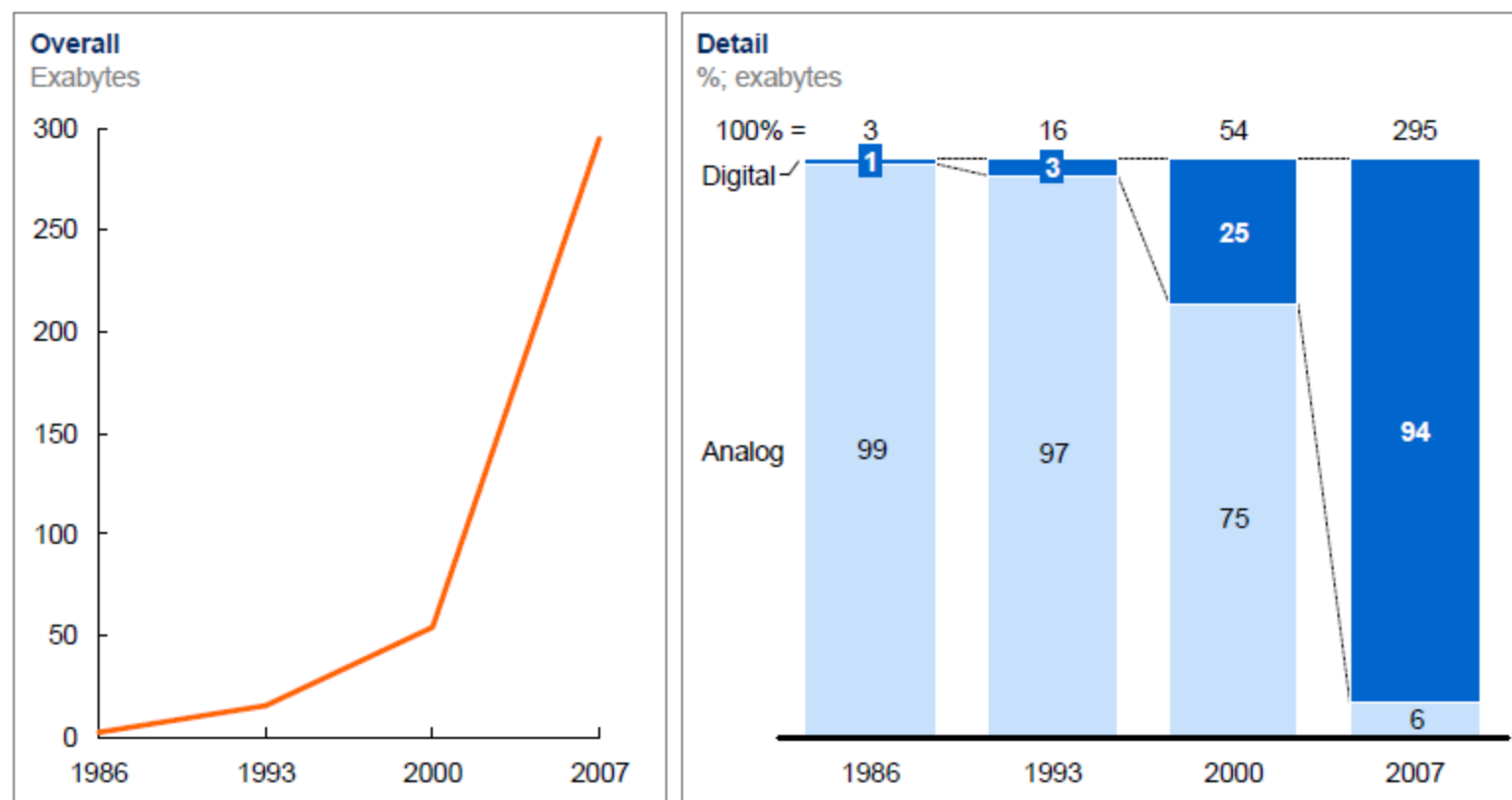
# Why Big-Data?

- ▶ Key enablers for the growth of “Big Data” are:
  - Increase of storage capacities
  - Increase of processing power
  - Availability of data

# Enabler: Data storage

**Data storage has grown significantly, shifting markedly from analog to digital after 2000**

Global installed, optimally compressed, storage



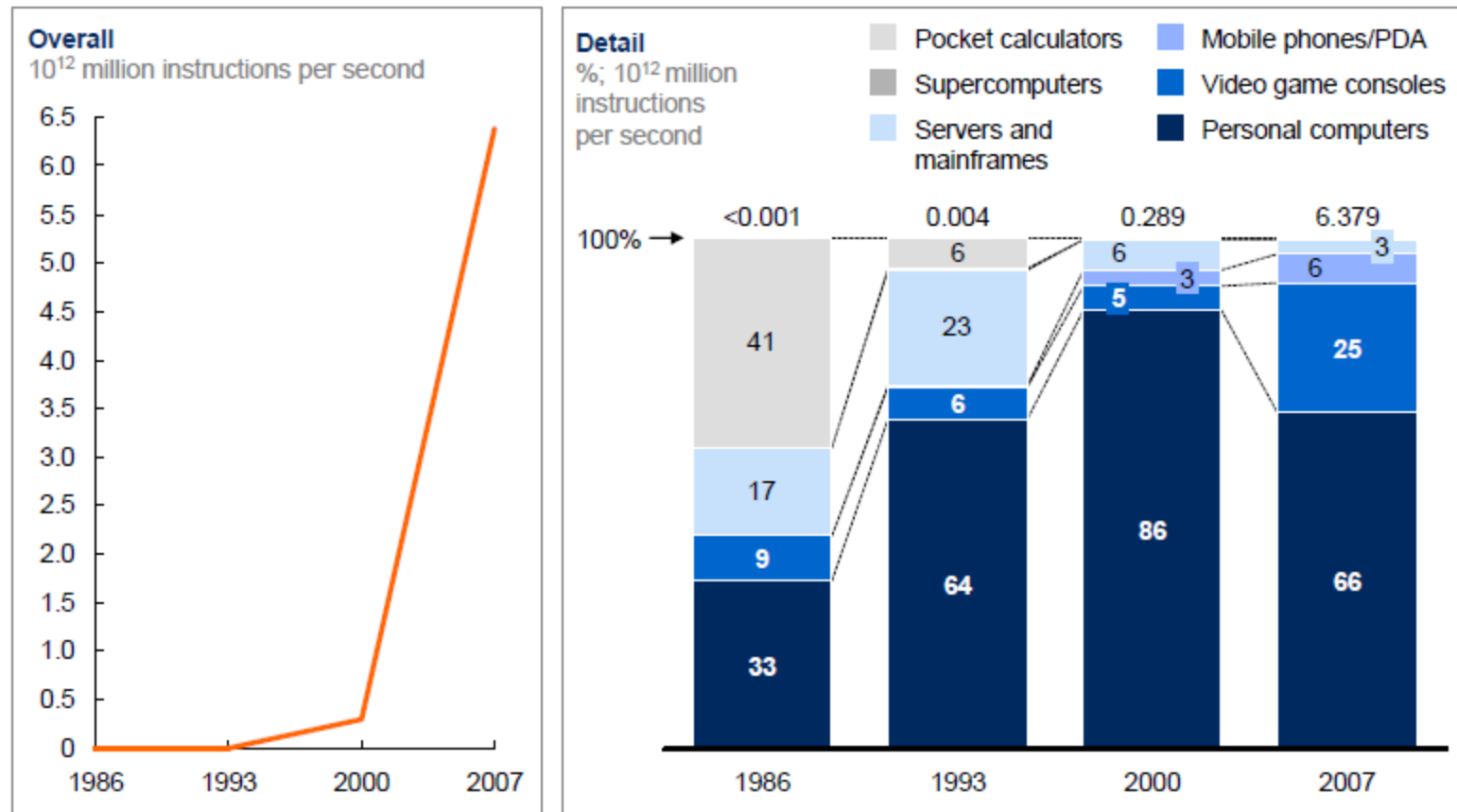
NOTE: Numbers may not sum due to rounding.

SOURCE: Hilbert and López, "The world's technological capacity to store, communicate, and compute information," *Science*, 2011

# Enabler: Computation capacity

## Computation capacity has also risen sharply

Global installed computation to handle information

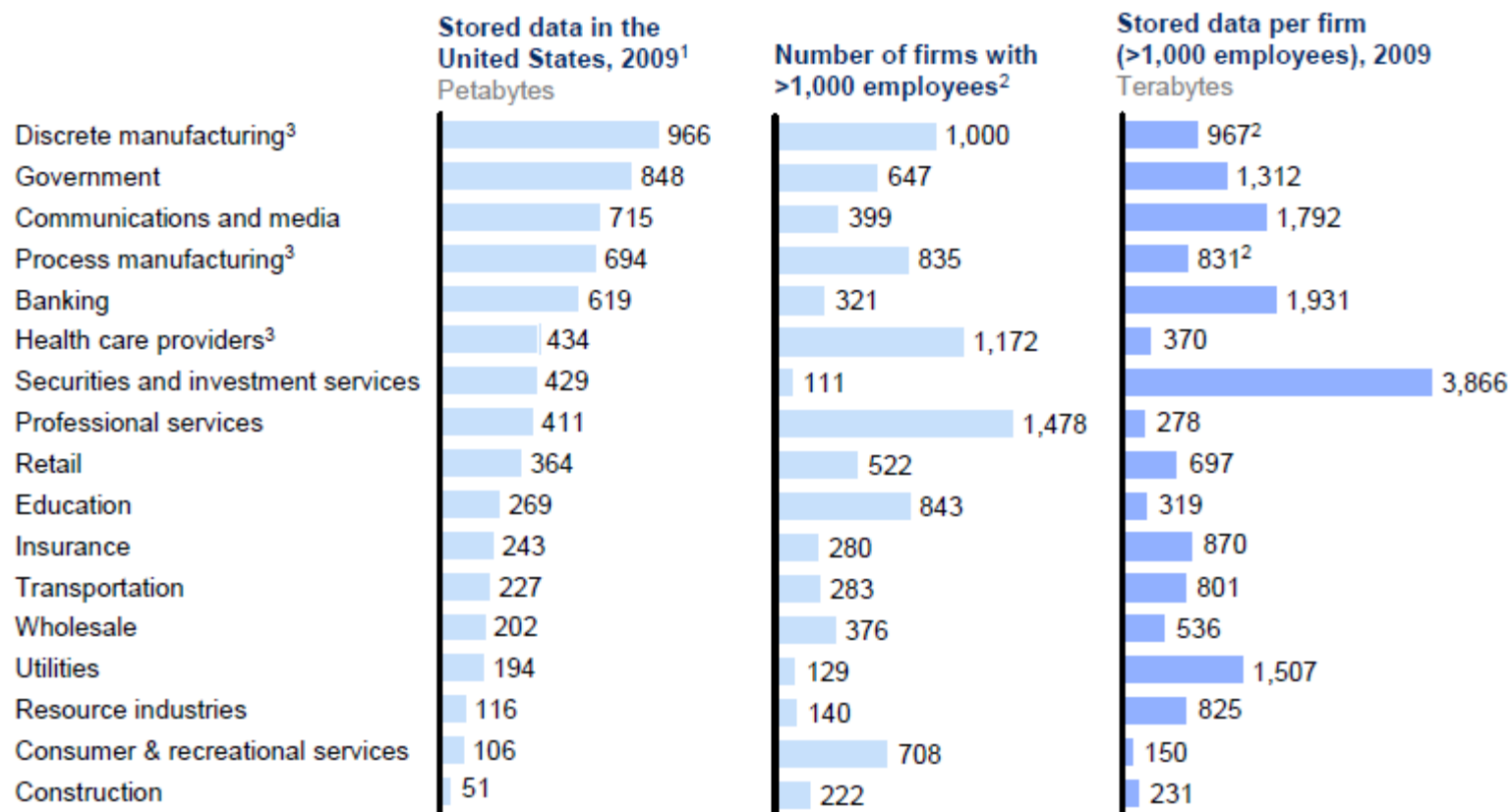


NOTE: Numbers may not sum due to rounding.

SOURCE: Hilbert and López, "The world's technological capacity to store, communicate, and compute information," *Science*, 2011

# Enabler: Data availability

**Companies in all sectors have at least 100 terabytes of stored data in the United States; many have more than 1 petabyte**



1 Storage data by sector derived from IDC.

2 Firm data split into sectors, when needed, using employment

3 The particularly large number of firms in manufacturing and health care provider sectors make the available storage per company much smaller.

SOURCE: IDC; US Bureau of Labor Statistics; McKinsey Global Institute analysis

# Type of available data

The type of data generated and stored varies by sector<sup>1</sup>

|                                       | Video  | Image  | Audio  | Text/<br>numbers |
|---------------------------------------|--------|--------|--------|------------------|
| Banking                               | Medium | Medium | Medium | High             |
| Insurance                             | Low    | Low    | Low    | High             |
| Securities and investment services    | Low    | Low    | Low    | High             |
| Discrete manufacturing                | Medium | Medium | Low    | High             |
| Process manufacturing                 | Medium | Medium | Low    | High             |
| Retail                                | Medium | Low    | Low    | High             |
| Wholesale                             | Low    | Low    | Low    | High             |
| Professional services                 | Medium | Medium | Medium | High             |
| Consumer and recreational services    | Medium | Low    | Medium | Medium           |
| Health care                           | Low    | High   | Low    | High             |
| Transportation                        | Medium | Low    | Low    | High             |
| Communications and media <sup>2</sup> | High   | Medium | High   | High             |
| Utilities                             | Medium | Medium | Low    | High             |
| Construction                          | Low    | High   | Low    | Medium           |
| Resource industries                   | Medium | Medium | Low    | High             |
| Government                            | High   | Medium | High   | High             |
| Education                             | High   | Medium | High   | Medium           |

**Penetration**

High  
 Medium  
 Low

<sup>1</sup> We compiled this heat map using units of data (in files or minutes of video) rather than bytes.

<sup>2</sup> Video and audio are high in some subsectors.

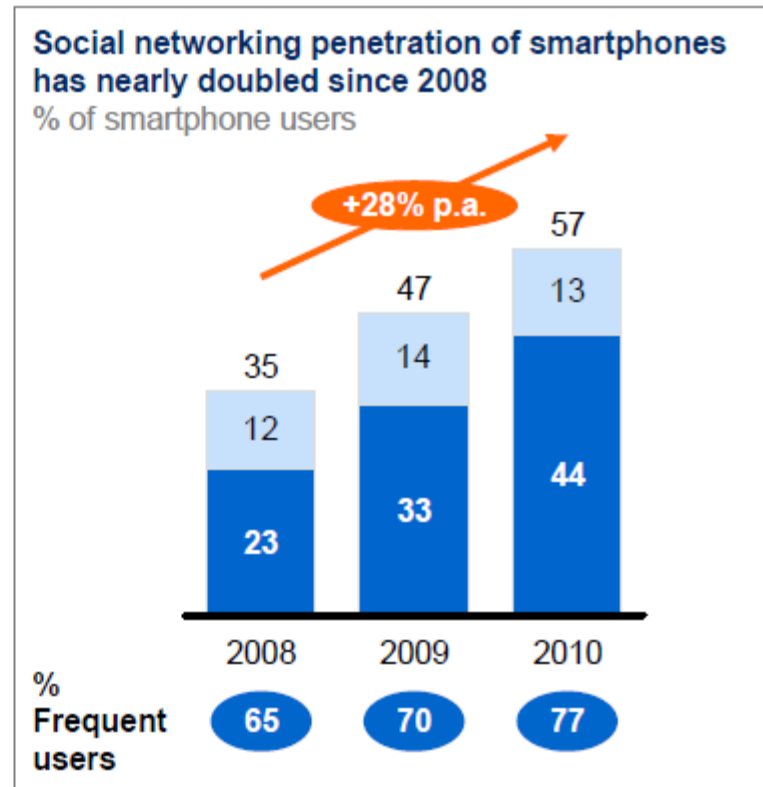
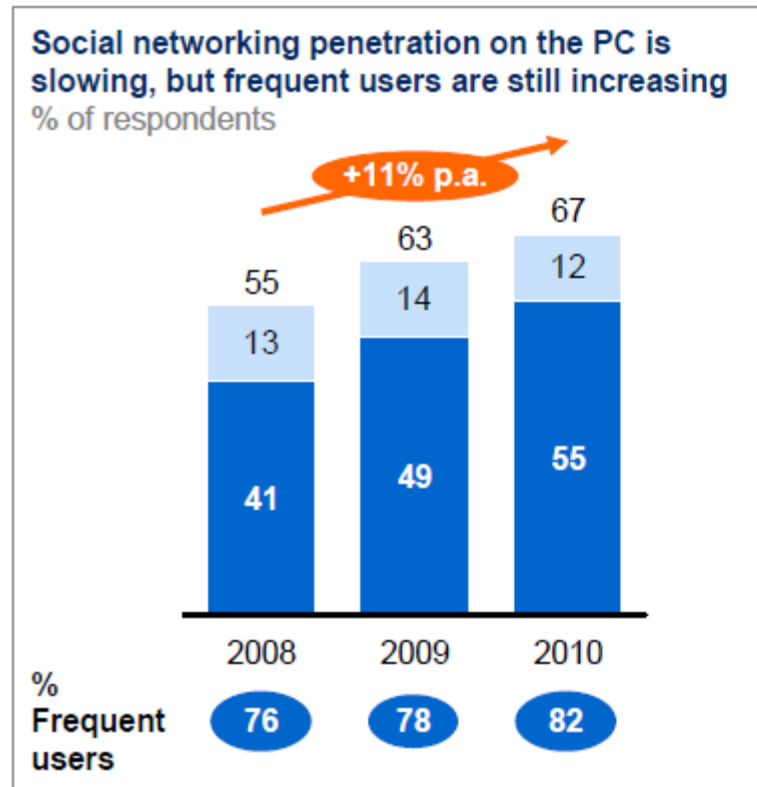
SOURCE: McKinsey Global Institute analysis



# Data available from social networks and mobile devices

The penetration of social networks is increasing online and on smartphones; frequent users are increasing as a share of total users<sup>1</sup>

■ Frequent user<sup>2</sup>



- 1 Based on penetration of users who browse social network sites. For consistency, we exclude Twitter-specific questions (added to survey in 2009) and location-based mobile social networks (e.g., Foursquare, added to survey in 2010).
- 2 Frequent users defined as those that use social networking at least once a week.
- SOURCE: McKinsey iConsumer Survey

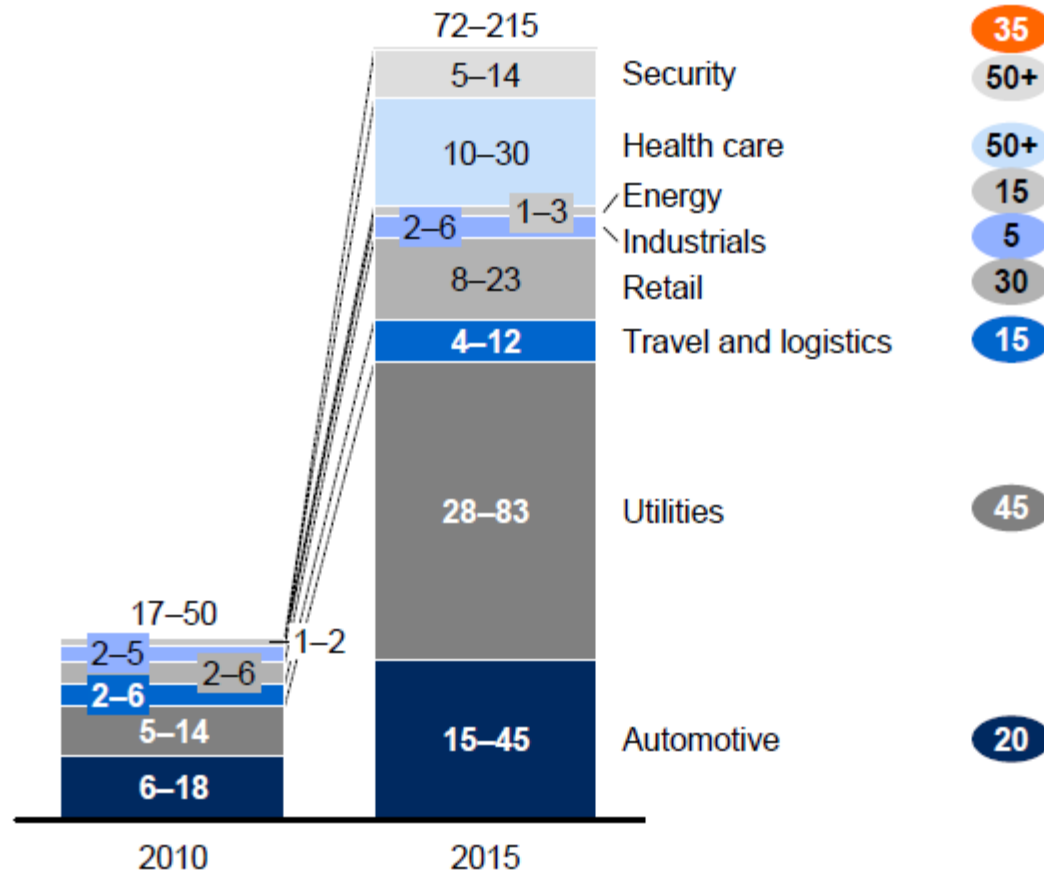
# Data available from “Internet of Things”

**Data generated from the Internet of Things will grow exponentially as the number of connected nodes increases**

Estimated number of connected nodes

Million

Compound annual growth rate 2010–15, %



NOTE: Numbers may not sum due to rounding.

SOURCE: Analyst interviews; McKinsey Global Institute analysis

# Big-data value chain

## Big data constituencies

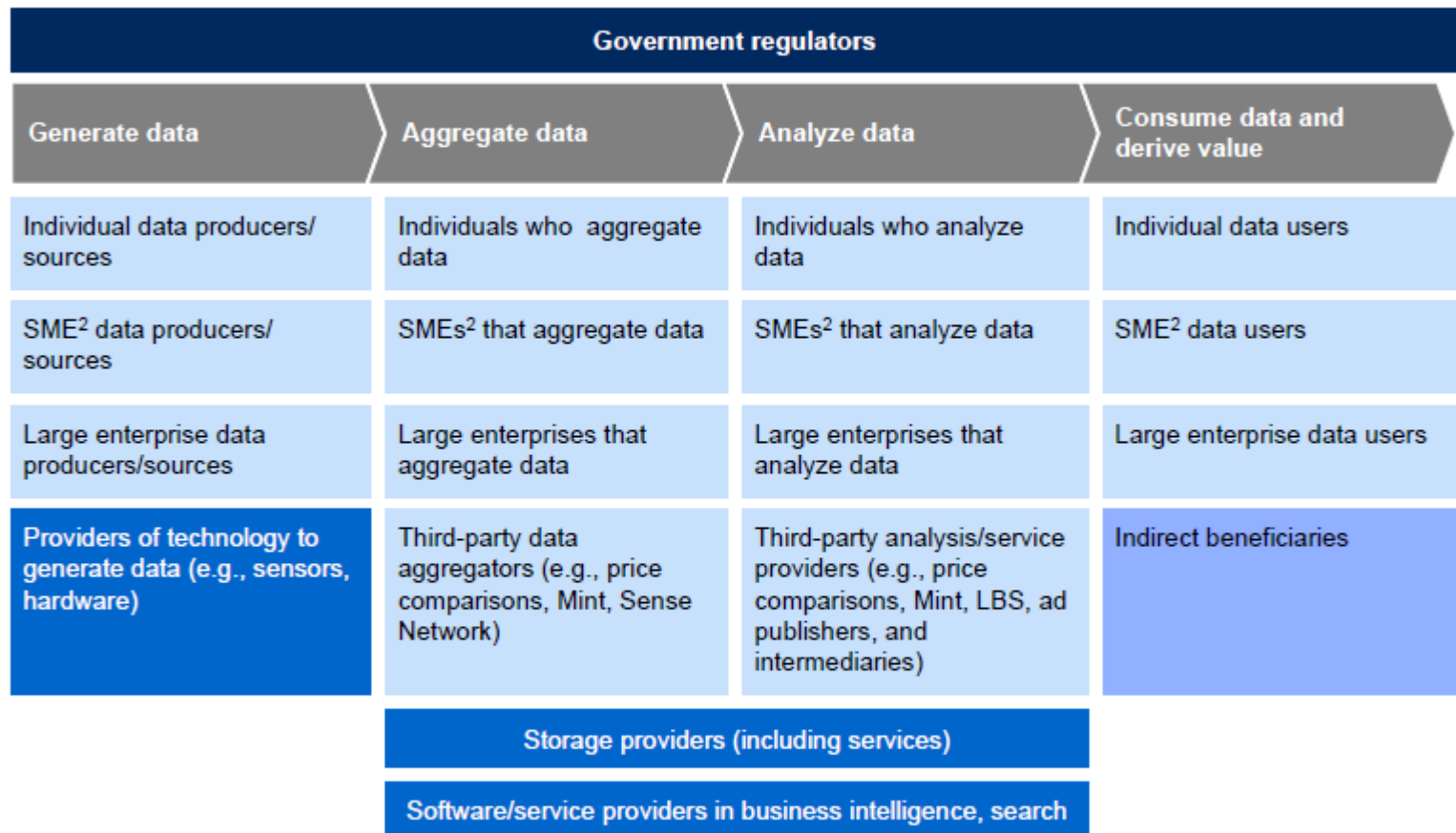
### Big data activity/value chain

Individuals/organizations using data<sup>1</sup>

Indirect beneficiaries

Providers of technology

Government regulators



<sup>1</sup> Individuals/organizations generating, aggregating, analyzing, or consuming data.

<sup>2</sup> Small and medium-sized enterprises.

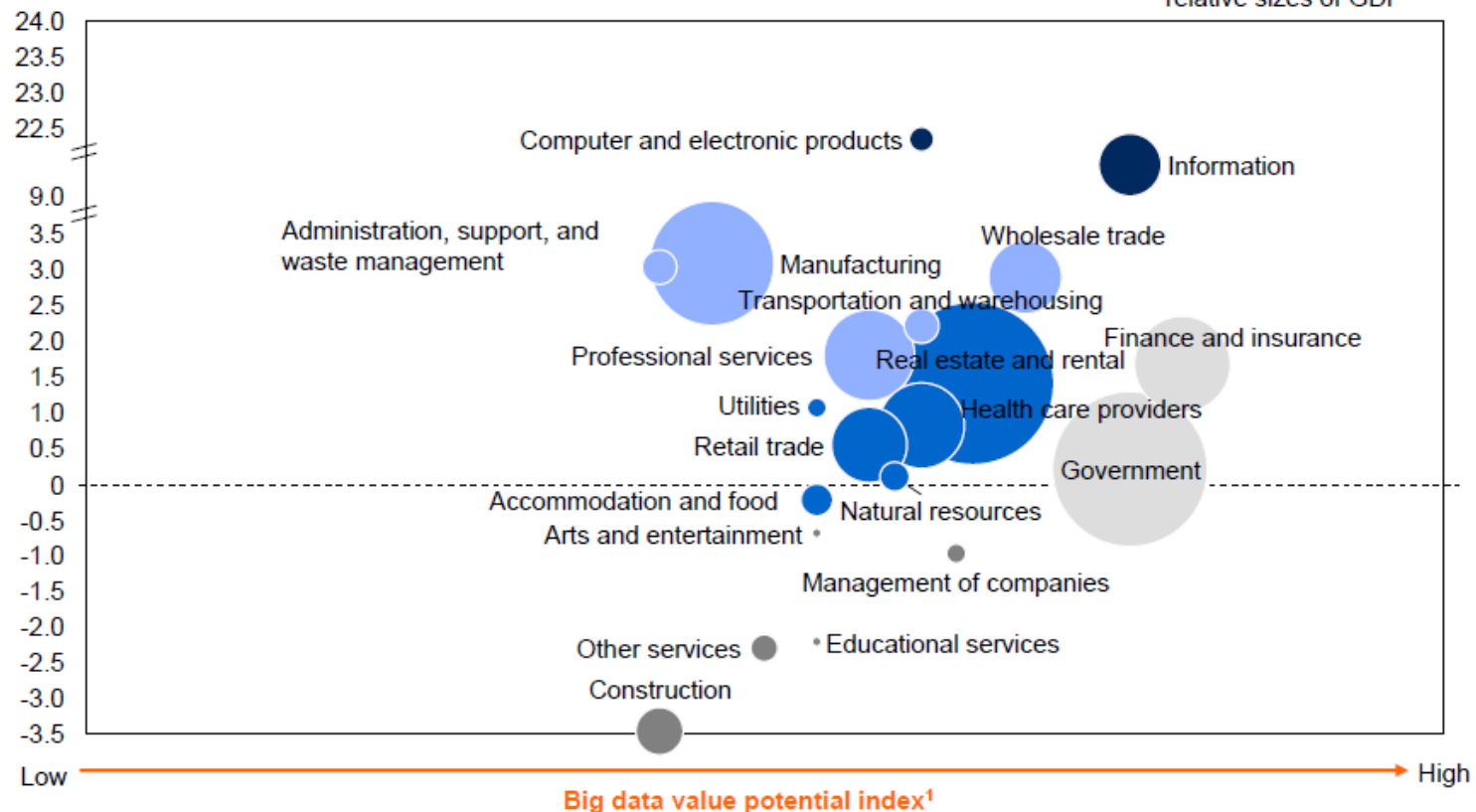
SOURCE: McKinsey Global Institute analysis

# Gains from Big-Data per sector

**Some sectors are positioned for greater gains from the use of big data**

Historical productivity growth in the United States, 2000–08

%



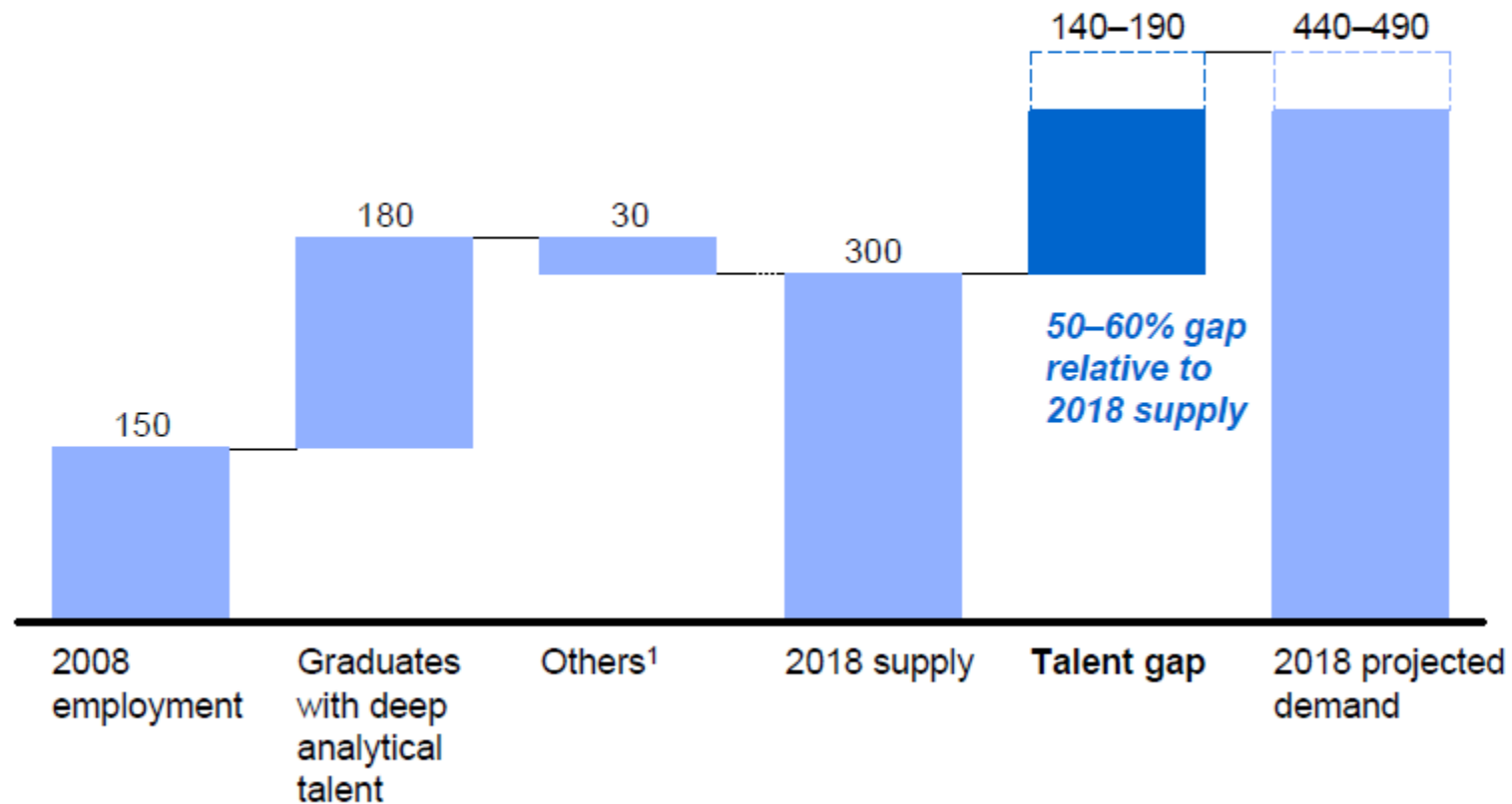
1 See appendix for detailed definitions and metrics used for value potential index.  
SOURCE: US Bureau of Labor Statistics; McKinsey Global Institute analysis

# Predicted lack of talent for Big-Data related technologies

**Demand for deep analytical talent in the United States could be 50 to 60 percent greater than its projected supply by 2018**

Supply and demand of deep analytical talent by 2018

Thousand people



<sup>1</sup> Other supply drivers include attrition (-), immigration (+), and reemploying previously unemployed deep analytical talent (+).



Tools

# Tools typically used in Big-Data scenarios

- ▶ NoSQL
  - Databases MongoDB, CouchDB, Cassandra, Redis, BigTable, Hbase, Hypertable, Voldemort, Riak, ZooKeeper
- ▶ MapReduce
  - Hadoop, Hive, Pig, Cascading, Cascalog, mrjob, Caffeine, S4, MapR, Acunu, Flume, Kafka, Azkaban, Oozie, Greenplum
- ▶ Storage
  - S3, Hadoop Distributed File System
- ▶ Servers
  - EC2, Google App Engine, Elastic, Beanstalk, Heroku
- ▶ Processing
  - R, Yahoo! Pipes, Mechanical Turk, Solr/Lucene, ElasticSearch, Datameer, BigSheets, Tinkerpop

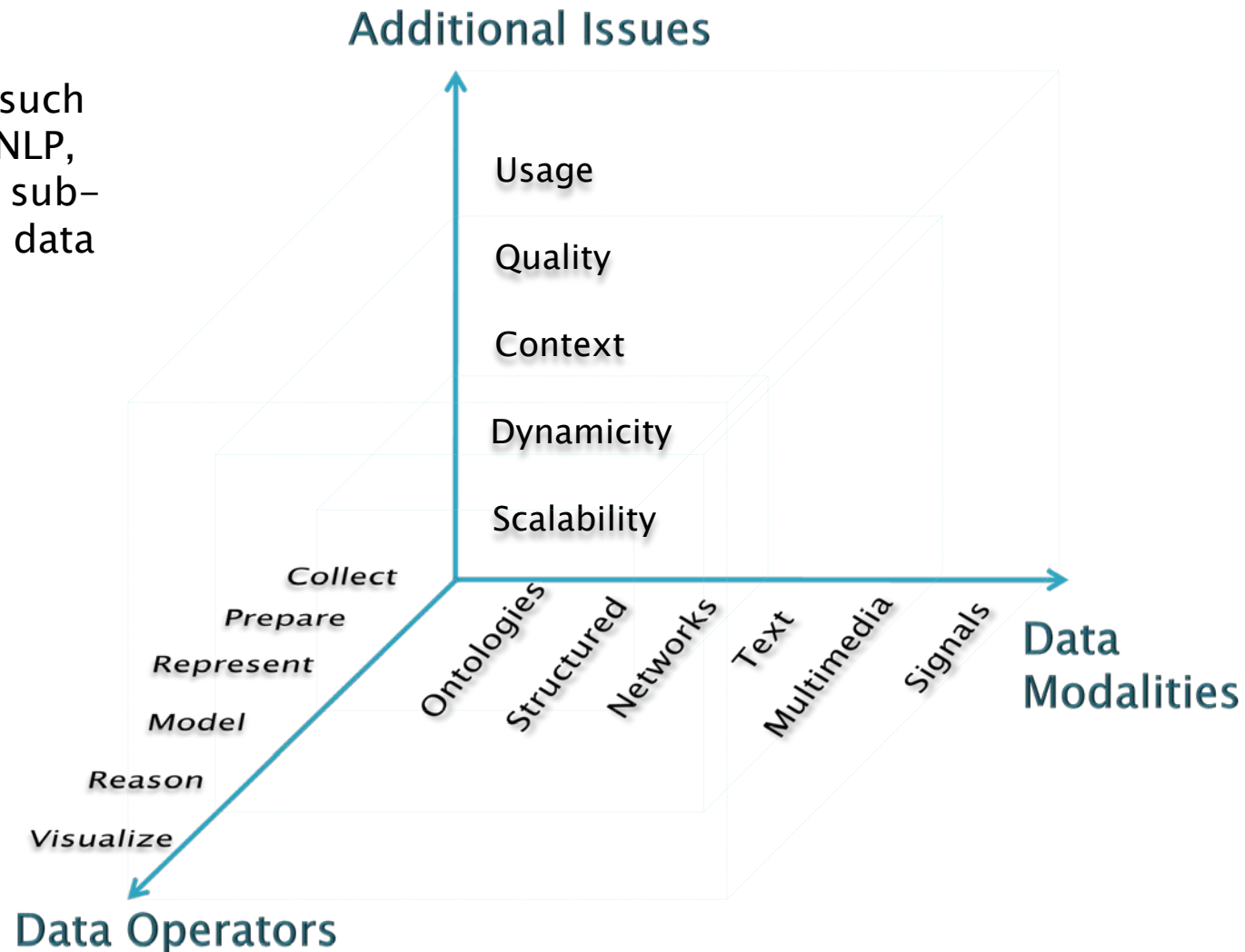
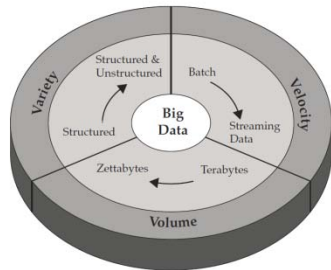
# Techniques

# When Big-Data is really a hard problem?

- ▶ ...when the operations on data are complex:
  - ...e.g. simple counting is not a complex problem
  - Modeling and reasoning with data of different kinds can get extremely complex
- ▶ Good news about big-data:
  - Often, because of vast amount of data, modeling techniques can get simpler (e.g. smart counting can replace complex model based analytics)...
  - ...as long as we deal with the scale

# What matters when dealing with data?

- ▶ Research areas (such as IR, KDD, ML, NLP, SemWeb, ...) are sub-cubes within the data cube

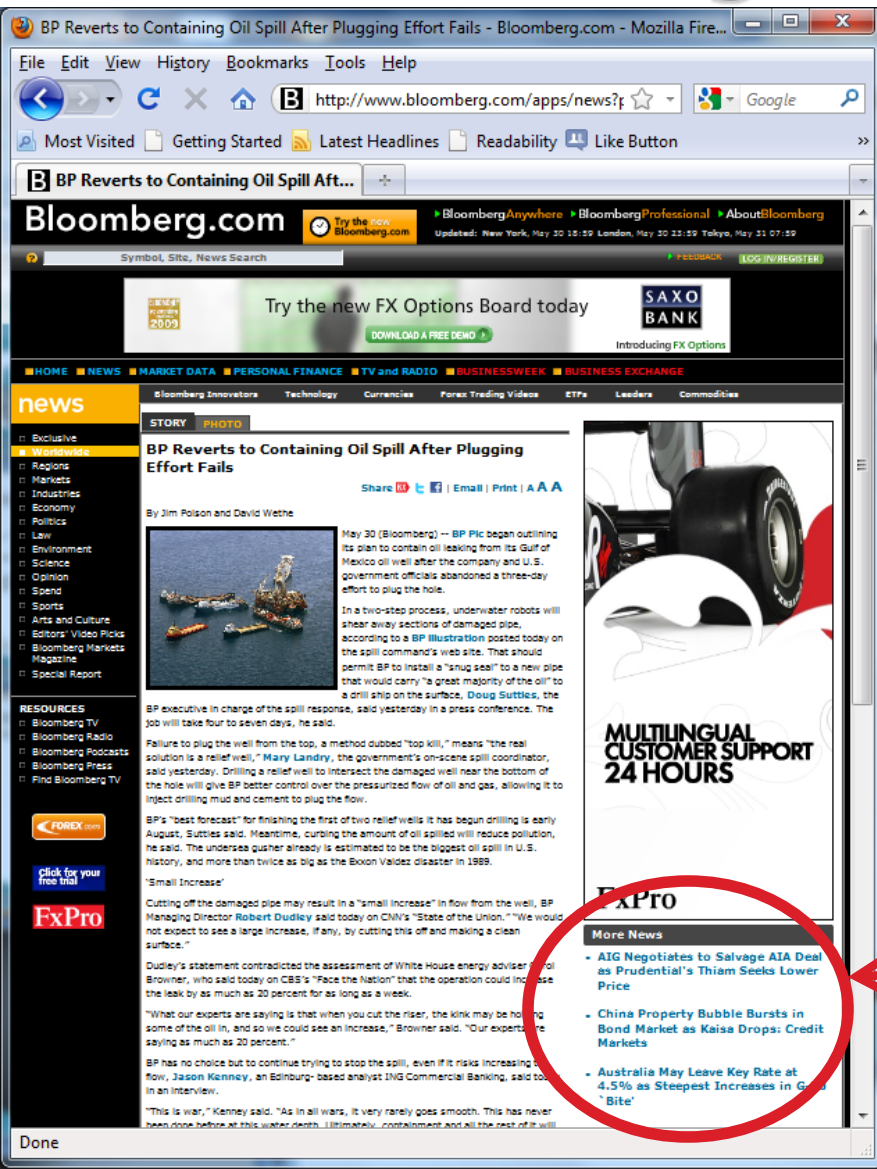




# Applications

# Recommendation

# ...an example: recommendation @Bloomberg.com



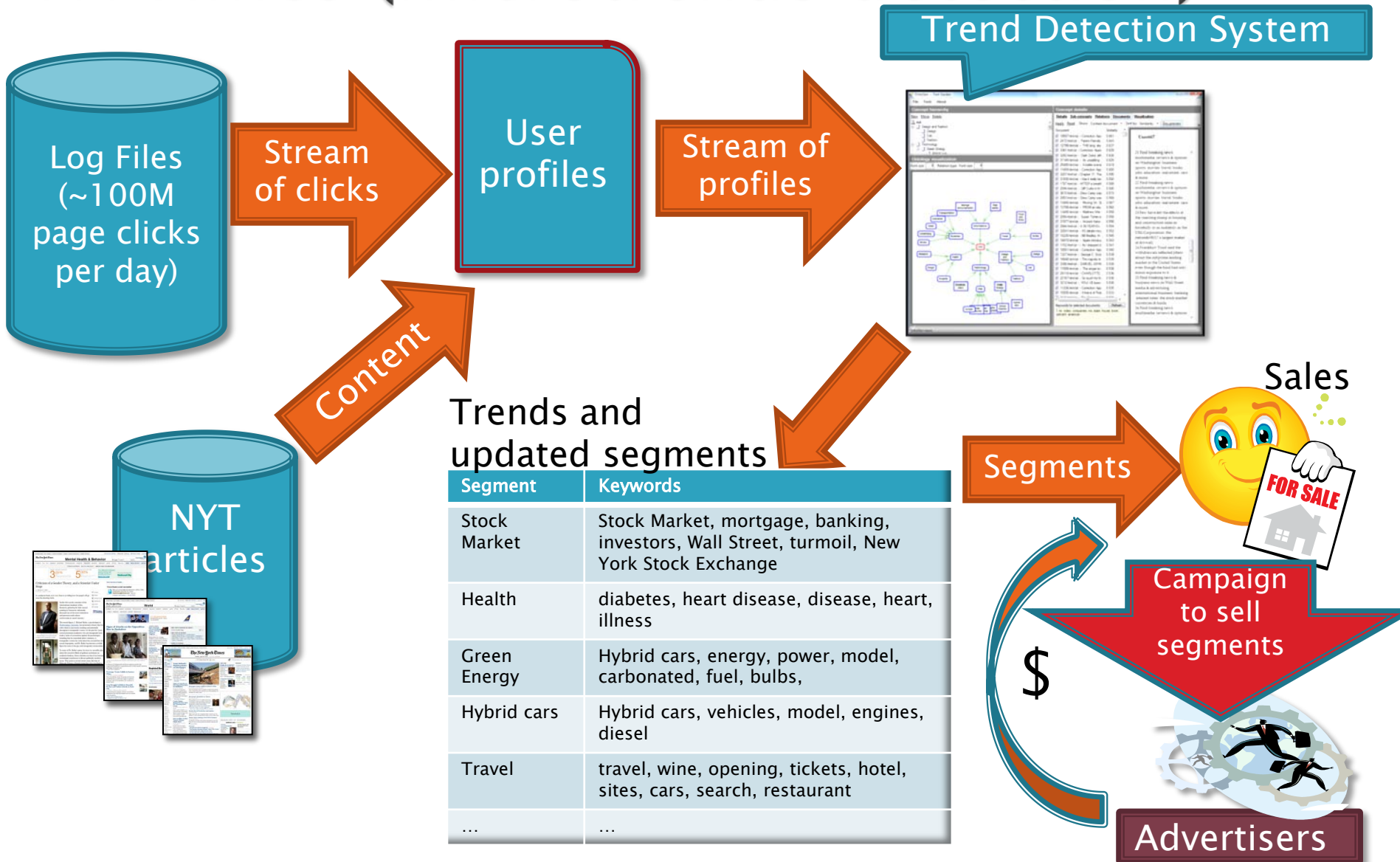
- ▶ Good recommendations can make a big difference when keeping a user on a web site
  - ...the key is how rich context model a system is using to select information for a user
  - Bad recommendations <1% users, good ones >5% users click

Contextual  
personalized  
recommendations  
generated in ~20ms

# Each click on the web site is enriched and indexed using:

- ▶ Domain
- ▶ Sub-domain
- ▶ Page URL
- ▶ URL sub-directories
- ▶ Page Meta Tags
- ▶ Page Title
- ▶ Page Content
- ▶ Named Entities
- ▶ Has Query
- ▶ Referrer Query
- ▶ Referring Domain
- ▶ Referring URL
- ▶ Outgoing URL
- ▶ GeolP Country
- ▶ GeolP State
- ▶ GeolP City
- ▶ Absolute Date
- ▶ Day of the Week
- ▶ Day period
- ▶ Hour of the day
- ▶ User Agent
- ▶ Zip Code
- ▶ State
- ▶ Income
- ▶ Age
- ▶ Gender
- ▶ Country
- ▶ Job Title
- ▶ Job Industry

# Application: Online Advertising for NYTimes (microtrends detection)



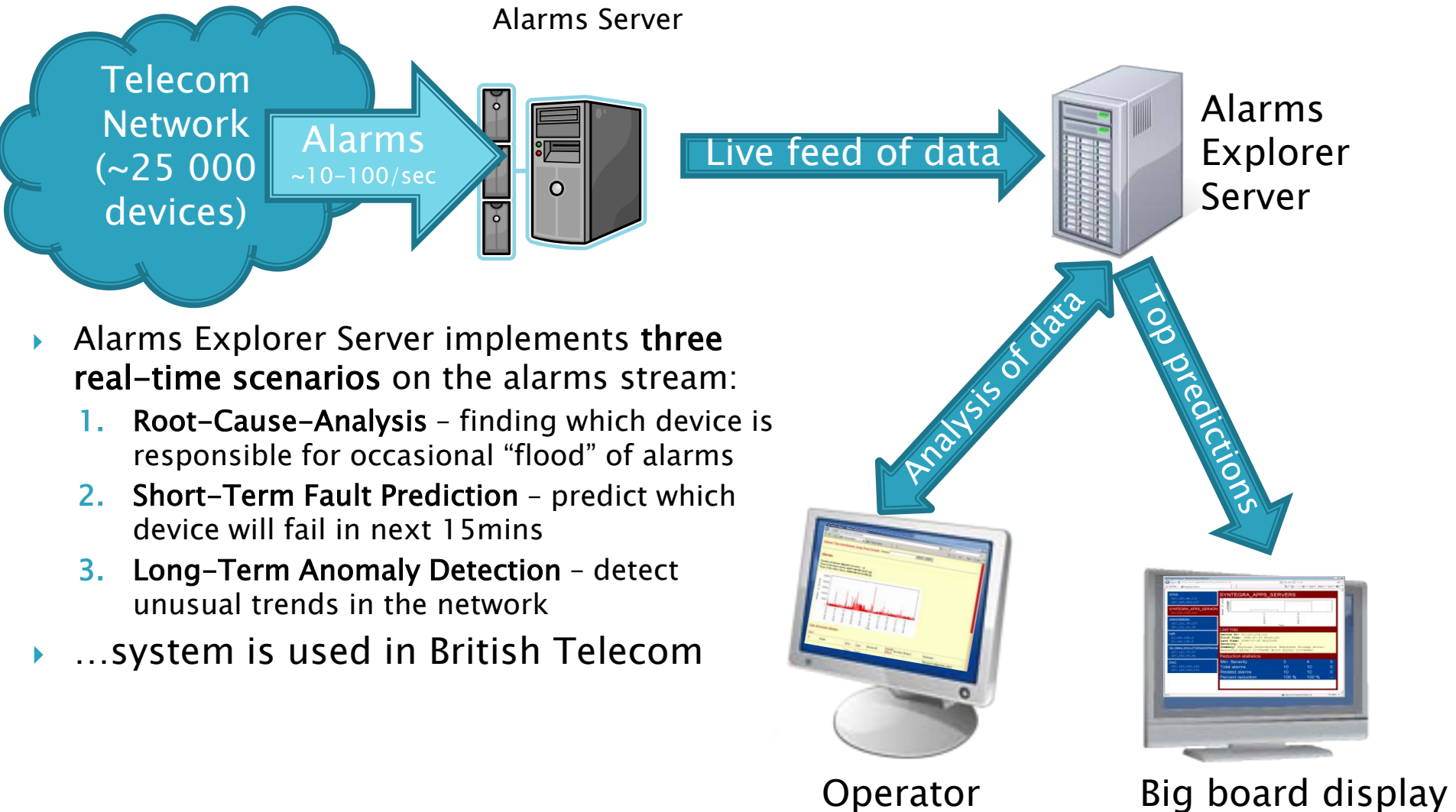


# Figures for one day of NYTimes

- ▶ 50Gb of uncompressed log files
- ▶ 10Gb of compressed log files
- ▶ 0.5Gb of processed log files
- ▶ 50–100M clicks
- ▶ 4–6M unique users
- ▶ 7000 unique pages with more than 100 hits
- ▶ Index size 2Gb
- ▶ Pre-processing & indexing time
  - ~10min on workstation (4 cores & 32Gb)
  - ~1 hour on EC2 (2 cores & 16Gb)

# Root-cause analysis

# Applications: Telecommunication Network Monitoring



# Analysis of MSN–Messenger Social–network

- ▶ Presented in “Planetary–Scale Views on a Large Instant–Messaging Network” by Jure Leskovec and Eric Horvitz WWW2008

# Instant Messenger – Phenomena at a planetary scale

- ▶ Observe social and communication phenomena at a *planetary* scale
- ▶ Largest social network analyzed to date

## Research questions:

- ▶ How does communication change with user demographics (age, sex, language, country)?
- ▶ How does geography affect communication?
- ▶ What is the structure of the communication network?

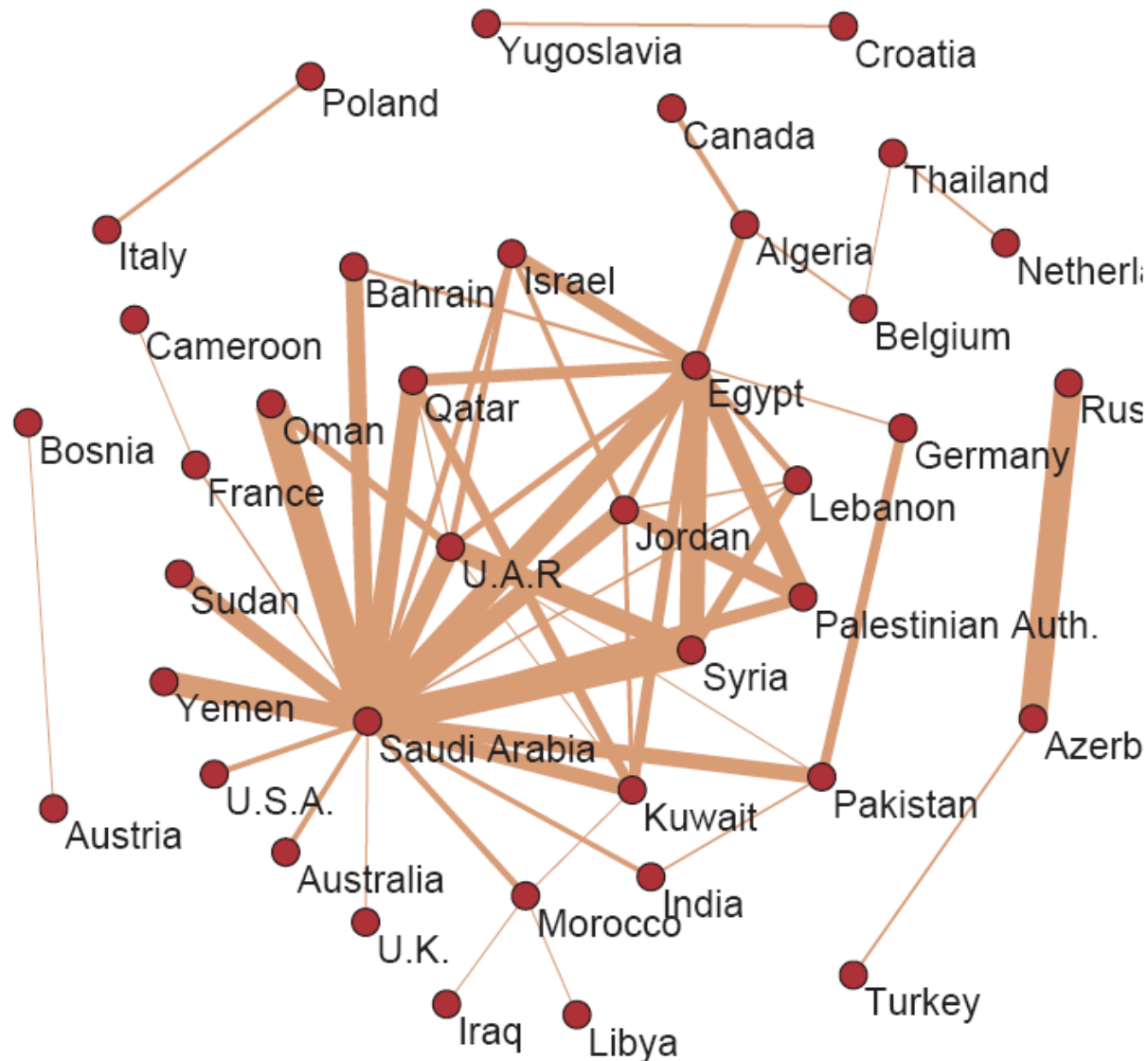
# Data statistics: Total activity

- ▶ We collected the data for **June 2006**
- ▶ Log size:
  - 150Gb/day (compressed)**
- ▶ Total: 1 month of communication data:
  - 4.5Tb of compressed data**
- ▶ **Activity over June 2006 (30 days)**
  - 245 million users logged in
  - 180 million users engaged in conversations
  - 17,5 million new accounts activated
  - More than 30 billion conversations
  - More than 255 billion exchanged messages

# Who talks to whom: Number of conversations

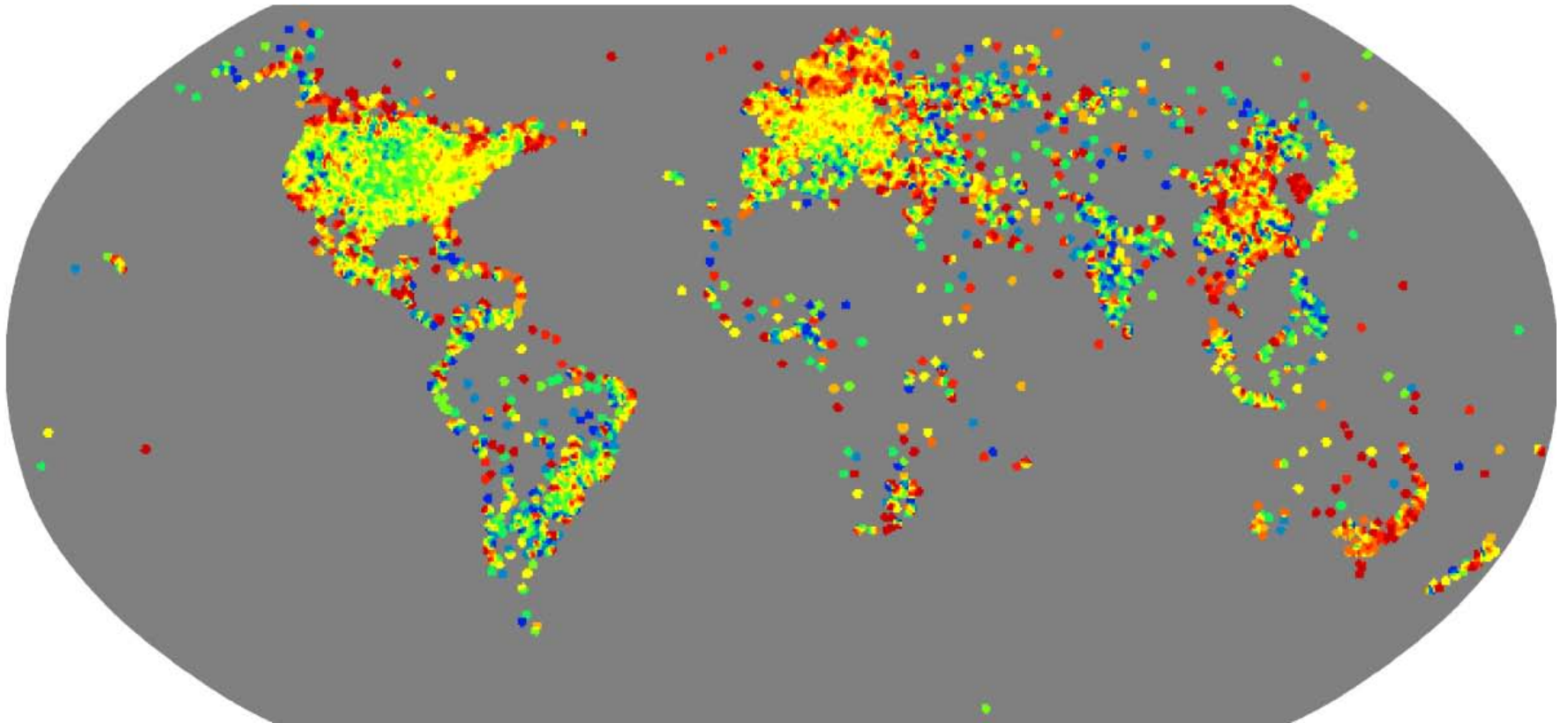


# Who talks to whom: Conversation duration



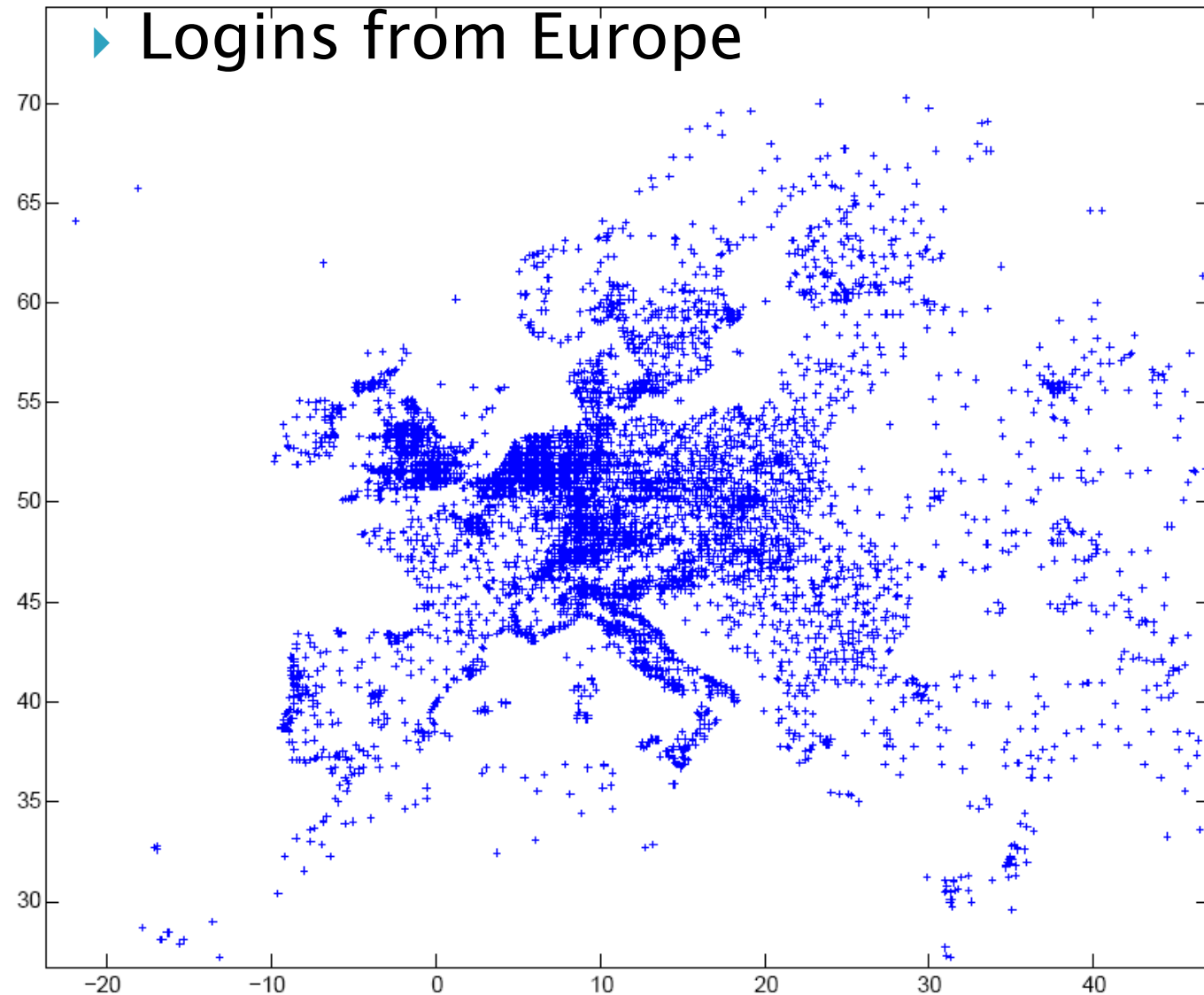


# Geography and communication

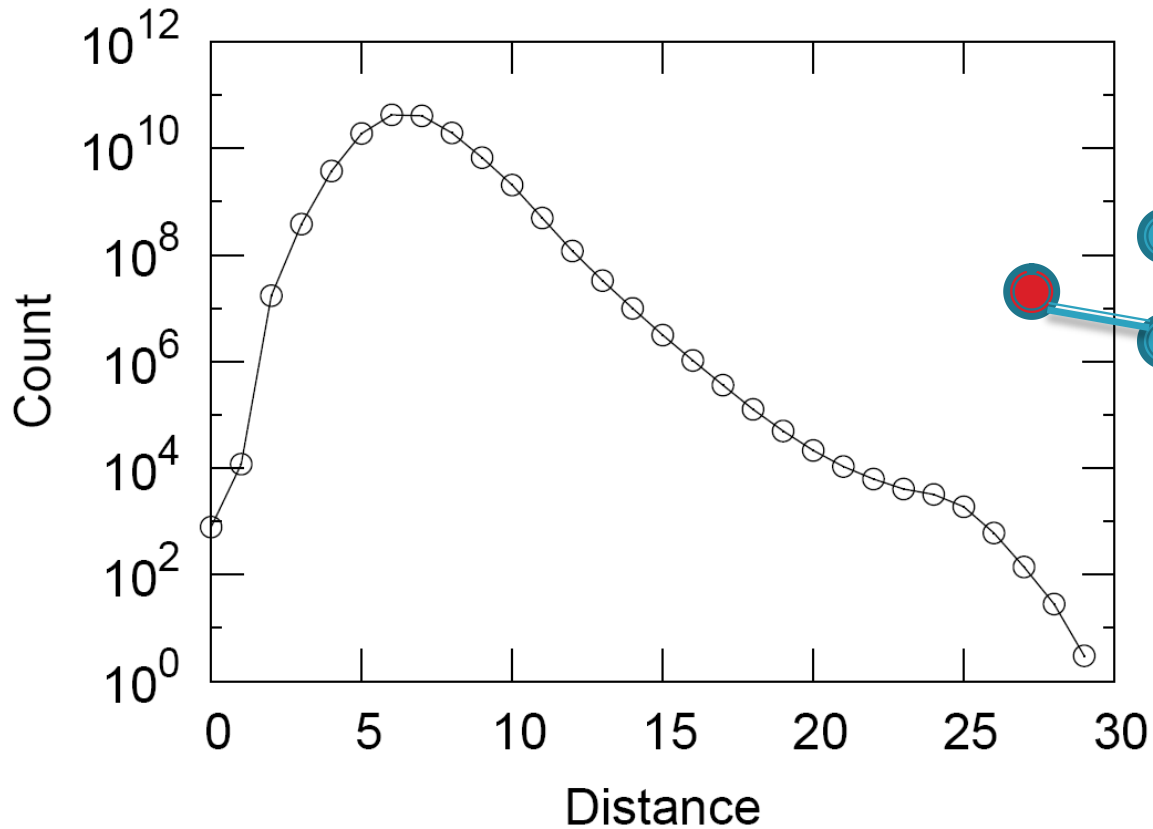


- ▶ Count the number of users logging in from particular location on the earth

# How is Europe talking



# Network: Small-world



- ▶ 6 degrees of separation [Milgram '60s]
- ▶ Average distance between two random users is 6.6
- ▶ 90% of nodes can be reached in  $< 8$  hops

| Hops | Nodes    |
|------|----------|
| 1    | 10       |
| 2    | 78       |
| 3    | 396      |
| 4    | 8648     |
| 5    | 3299252  |
| 6    | 28395849 |
| 7    | 79059497 |
| 8    | 52995778 |
| 9    | 10321008 |
| 10   | 1955007  |
| 11   | 518410   |
| 12   | 149945   |
| 13   | 44616    |
| 14   | 13740    |
| 15   | 4476     |
| 16   | 1542     |
| 17   | 536      |
| 18   | 167      |
| 19   | 71       |
| 20   | 29       |
| 21   | 16       |
| 22   | 10       |
| 23   | 3        |
| 24   | 2        |
| 25   | 3        |

# Web-of-Things

# Literature on Big-Data

