# Big Data – Retos y Oportunidades

Raúl Ramos Pollán

**Líder Big Data & Large Scale Machine Learning**

**Laboratorio de Supercomputación y Cálculo Científico**

**Universidad Industrial de Santander – Bucaramanga, Colombia**

**http://sc3.uis.edu.co**

# Avalancha de datos – *Data Deluge*

In the same way that past Federal investments in information-technology R&D led to dramatic advances in supercopmuting and the creation of the Internet, the initiative we are launching today promises to transform our ability to use **Big Data** for scientific discovery, environmental and biomedical research, education, and national security

**Dr. J Holdren,**
**Director of the White House Office of Science and Technology ddd2012**

The World Economic Forum convened in Switzerland in January, 2012 highlighted **Big Data** as a new economic asset comparable to currency and gold

**Rethinking Personal Data: Strengthening Trus**
**World Economic Forum. May, 2012**

Today, we are experiencing a major shift in decision making driven by several factors: **unprecedented amounts of data** from a variety of sources […] on a variety of socio-economic, technological, and ecological systems, our vastly increasing ability to store and perform computation over very large data sets …
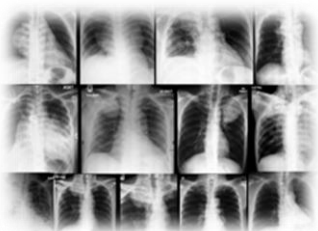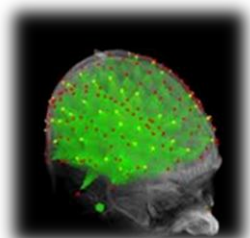
EU Digital Agenda, May 2013

Universidad Industrial de Santander

# Sources for Big Data

consumer logs

mobile devices

social networks

industry

patient registries

document repositories

web logs

bank transactions

government

sensors, instrumentation

simulations

science

image acquisition

biological data

Universidad Industrial de Santander

# Data Science / Big Data



**Bussiness Intelligence** → managing existing data for monitoring management concerns

**Data Science** → discovering knowledge

# Big Data

**datasets size >> traditional DB capacity**

**datasets groth >> tecnology evolution**

**new models & techniques for processing, analyzing, visualizing**

**BIG DATA MARKET FORECAST (US BILLIONS)**

$5.1  $10.2  $16.8  $32.1  $48.0  $53.4

DIGITAL DATA IN THE WORLD

2012: **2.75** ZETTABYTES

2015: **8** ZETTABYTES

Universidad
Industrial de
Santander

# Ley de Amdahl

# Big Data

Coarse grained
data parallelism

Trivial scalability

Variety
Velociy
Volume

# Big Data

**Coarse grained data parallelism**

**Trivial scalability**

Variety
Velociy
Volume

# Learning representations
## Latent semantics

# Learning representations
## Bag of Features / Latent semantics

## How to crunch 1PB?

Lots of disks spinning all the time
Redundancy, since disks die
Lots of CPU cores, working all the time
Retry, since network errors happen

## Design Qualities

Scalable – many servers with cores and disks
Reliable – redundant storage
Fault-tolerant – auto retry, self healing

## Computation to Data

Very simple computing model → Map-Reduce
Each computing node is also a storage node
HDFS → on top of ext3, fixed 64MB file blocks
write once, read many

# NoSQL

- Expressivity SQL vs. Scalability

- Simpler data model (key, values)
- Simpler operations

    Scan/access per key, basic transactions (check&put)

    No joins, no SQL language

- Simple failover and scale up

- Big table, Hbase, DynamoDB, Azure, Cassandra, etc.

→ more work for programmers!!!

# NoSQL Eventual Consistency

# Big Data in Marketing

Spencer Stuart survey
171 US marketing executives

**In which of the following areas are big data analytics currently having the largest impact on the way marketing is executed or how decisions are made in your organization?**

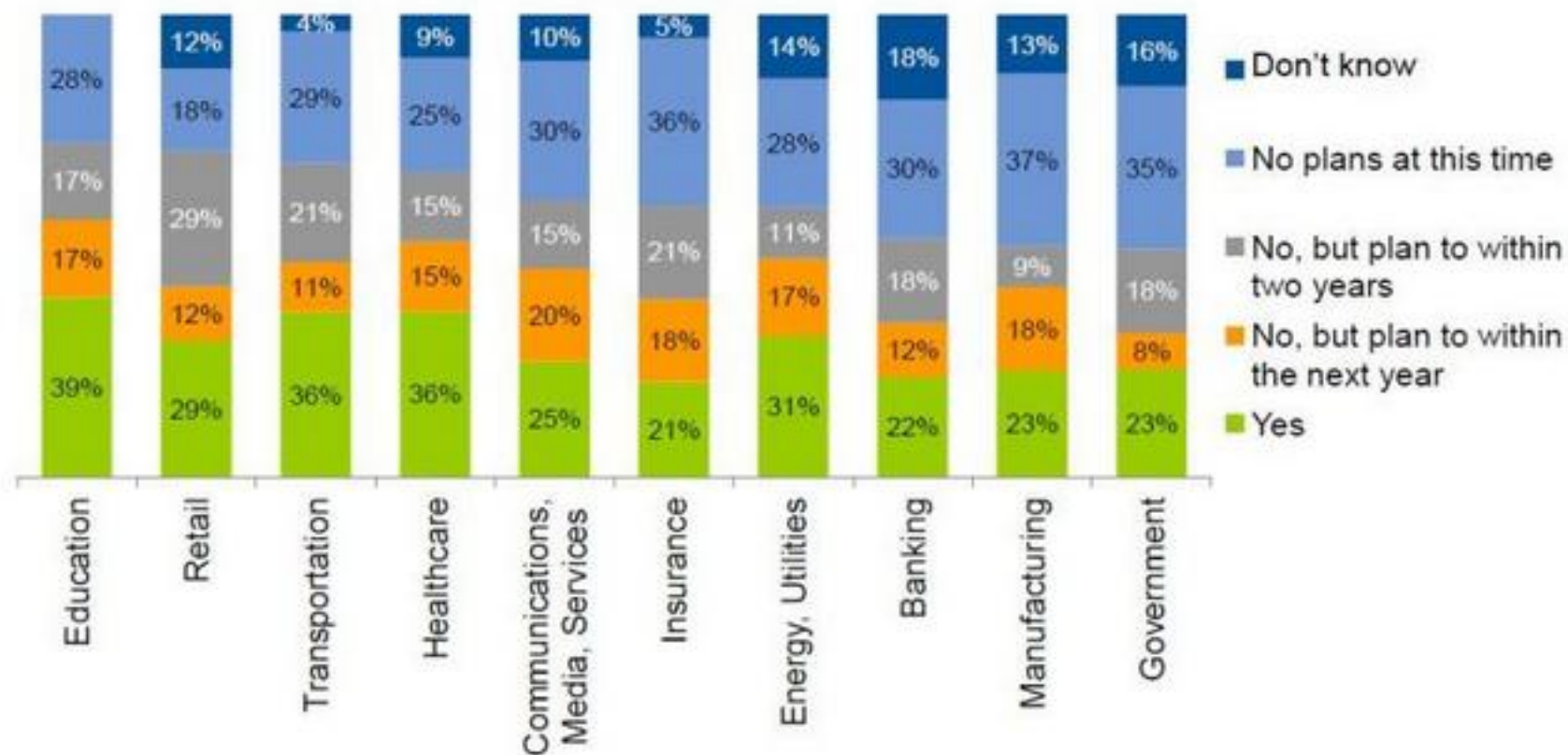| Area | Percentage |
| --- | --- |
| Non-digital advertising | 14% |
| SEO/SEM/email/SMS marketing | 58% |
| Social media outreach | 18% |
| Social media analysis | 35% |
| Marketing strategy | 41% |
| Brand strategy | 22% |
| Loyalty/retention programs | 36% |
| Customer segmentation | 49% |
| Customer service | 13% |
| Product/service development | 14% |
| Public relations/crisis awareness | 5% |
| Other | 5% |

Respondents were allowed to choose multiple responses.

Santander

Number of processed transactions at Eurex Exchange & response times

# Key points for Big Data

→ Integration of different data sources

→ Continuous prediction models from historical + real time data

→ (Semi)automatic knowledge discovery is now possible

→ Computing infrastructure

→ Multidisciplinary teams (in house, w/ academy)

# Multidisciplinary teams

The world of computing is flat, and anyone can do it. What will distinguish us from the rest of the world is **our ability to do it better and to exploit new architectures** we develop before those architectures become ubiquitous.

There is a clear and urgent need for a **new, modern approach to educating and training the next generation of researchers** in high performance computing specifically, and in modeling and simulation generally, for scientific discovery and engineering innovation.

Inadequate **education and training of the next generation of computational scientists** threatens global as well as U.S. growth of SBE&S […] unless we prepare researchers to develop and use the next generation of algorithms and computer architectures, we will not be able to exploit their game-changing capabilities.

There are clear and urgent opportunities for **industry-driven partnerships with universities** and national laboratories to hardwire scientific discovery and engineering innovation through SBE&S.

**www.wtec.org/sbes    2009**

# Challenges → Focus at UIS-SC3

- Identify data assets within your organization
- Identify diversity and sources of data
- Assess data links to support bandwidth
- Assess data quality and private data
- Understand gap data collection-interpretation
- Building specialized skills
- Building multidisciplinary teams

Universidad Industrial de Santander
# Supercomputación y Cálculo Científico

**Guane I**
60 TFlops
GPU Powered
**Launched 4/2012**
Doubling soon

**Gas Oil Energy**
**Big Data**
**Agua Recursos**
**Medio Ambiente**
**Astrofísica**
**Biología**

# http://sc3.uis.edu.co