www.bsc.es

**Barcelona Supercomputing Center**
*Centro Nacional de Supercomputación*

EXCELENCIA SEVERO OCHOA
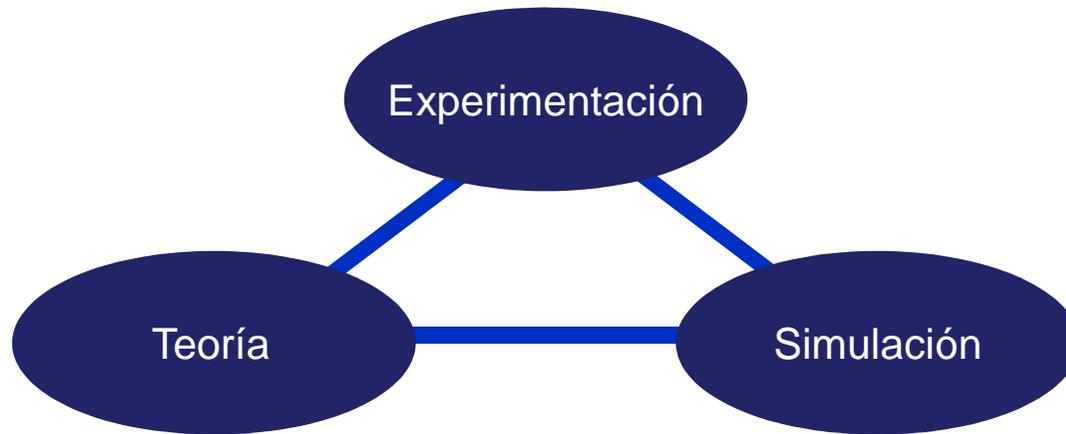
# Supercomputación, Big Data y Computación Cognitiva
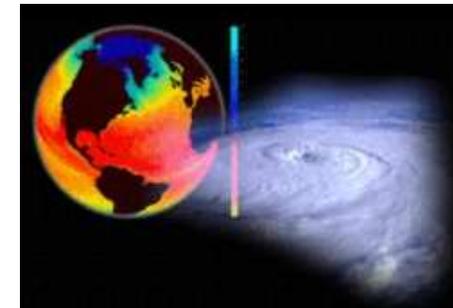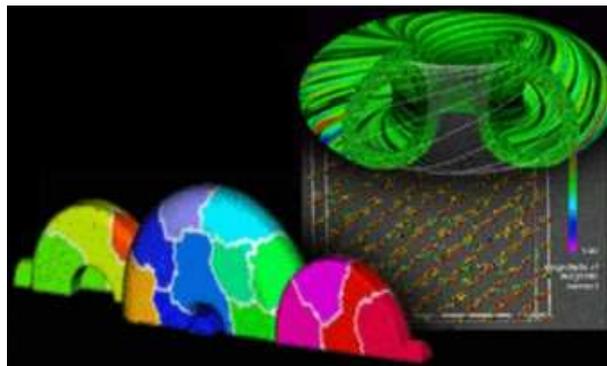
## Prof. Mateo Valero

BSC Director

erc

Granada, Abril, 2015

# ¿Cómo avanza la ciencia hoy?



Experimentación

Teoría

Simulación

Simulación = Calcular las formulas de la teoría



CARO



PELIGROSO



IMPOSIBLE

# Simulation on Supercomputers helps to solve scientific, industrial and societal challenges



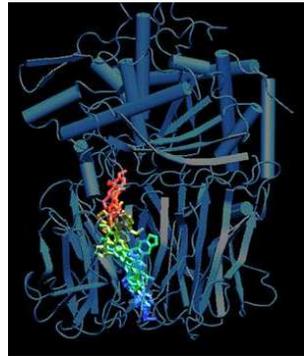**EXCELENCIA SEVERO OCHOA**
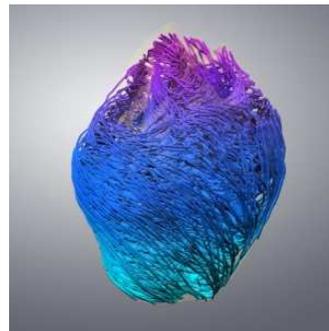
## Environment

Climate prediction

Air quality

## Life Sciences

Personalised medicine

The virtual patient

## Engineering

Industrial process improvement

Virtual prototyping

## Energy

Oil exploration

**REPSOL**

Wind farm design

**IBERDROLA**

Fusion (ITER)

**Barcelona Supercomputing Center**
*Centro Nacional de Supercomputación*

# Computers are now an essential part of almost all research

October 11, 2013
**Science: Beyond the God particle**
By Clive Cookson

The Nobel prizes in chemistry and physics show
how computing is changing every field of research

# Top10

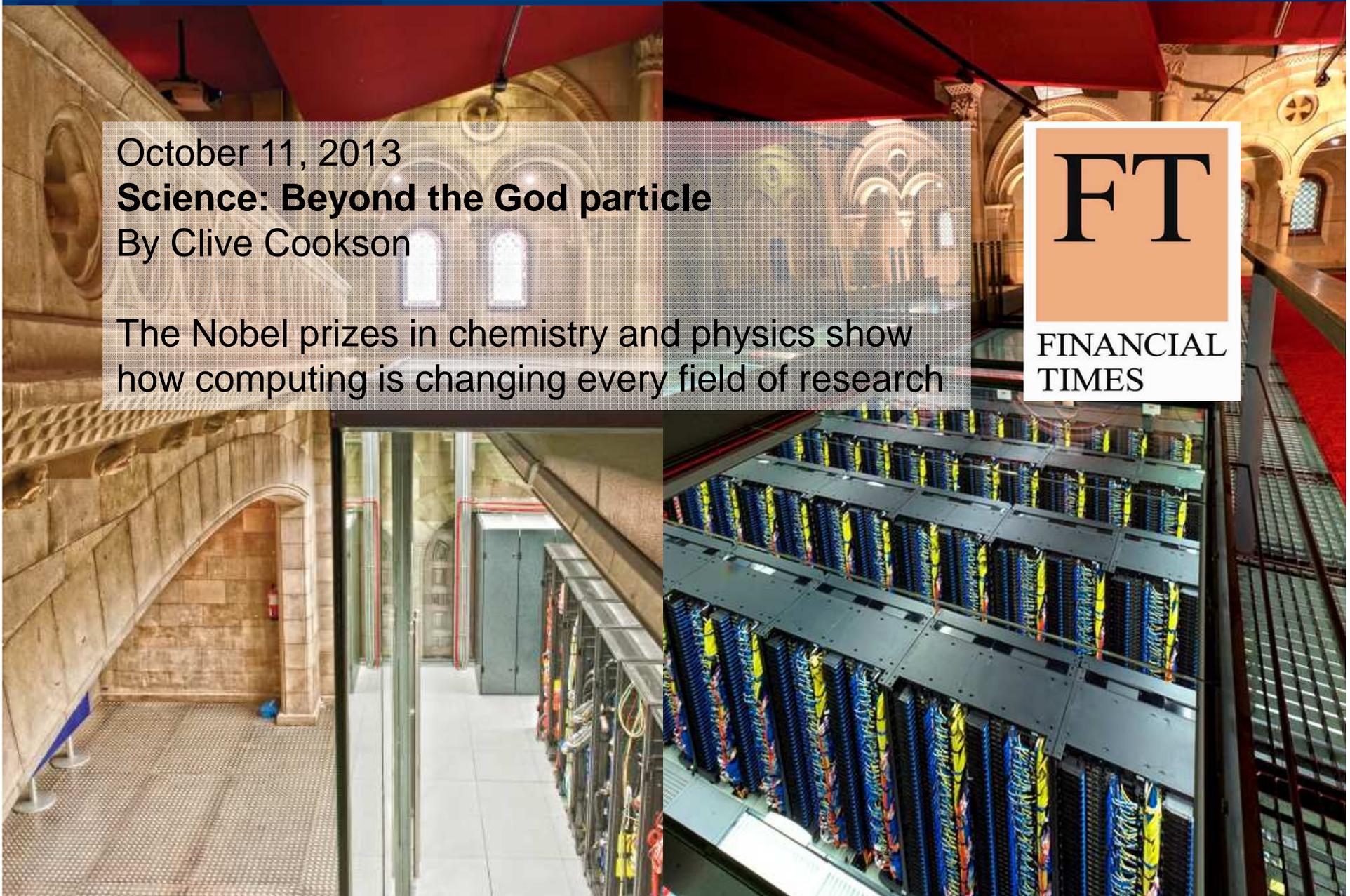| Rank | Site | Computer | Procs | Rmax | Rpeak | Power | GFlops/Watt | Name |
|------|------|----------|-------|------|-------|-------|-------------|------|
| 1 | National Super Computer Center in Guangzhou | TH-IVB-FEP Cluster, Intel Xeon E5-2692 12C 2.200GHz, TH Express-2, Intel Xeon Phi 31S1P | 3120000 2736000 | 33,86 | 54,90 | 17,8 | 1,90 | Tianhe-2 (MilkyWay-2) |
| 2 | DOE/SC/OAK Ridge National Lab | CRAY XK7, Opteron 6274 16C, 2.20 GHz, Cray Gemini interconnect, NVIDIA K20x | 560640 261632 | 17,59 | 27,11 | 8,21 | 2,14 | Titan |
| 3 | DOE/NNSA/LLNL | BlueGene/Q, Power BQC 16C 1.60 GHz, Custom | 1572864 | 17,17 | 20,13 | 7,89 | 2,18 | Sequoia |
| 4 | RIKEN Advanced Institute for Computational Science (AICS) | Fujitsu, K computer, SPARC64 VIIIfx 2.0GHz, Tofu interconnect | 705024 | 10,51 | 11,28 | 12,65 | 0,83 | K |
| 5 | DOE/SC/Argonne National Laboratory | BlueGene/Q, Power BQC 16C 1.60GHz, Custom | 786432 | 8,58 | 10,06 | 3,94 | 2,18 | Mira |
| 6 | CSCS | Cray XC30, Xeon E5-2670 8C 2.600GHz, Aries interconnect , NVIDIA K20x | 115984 73808 | 6,27 | 7,79 | 2,32 | 2,70 | Piz Daint |
| 7 | Texas Advanced Computing Center | PowerEdge C8220, Xeon E5-2680 8C 2.700GHz, Infiniband FDR, Intel Xeon Phi | 462462 366366 | 5,17 | 8,52 | 4,51 | 1,14 | Stampede |
| 8 | Forschungszentrum Juelich (FZJ) | BlueGene/Q, Power BQC 16C 1.60GHz, Custom | 458752 | 5,00 | 5,87 | 2,30 | 2,18 | JUQUEEN |
| 9 | DOE/NNSA/LLNL | BlueGene/Q, Power BQC 16C 1.60 GHz, Custom | 393216 | 4,29 | 5,03 | 1,97 | 2,18 | Vulcan |
| 10 | Government | Cray XC30, Intel Xeon E5-2660v2 10C 2.2GHz, Aries, NVIDIA K40 | 72800 62400 | 3,57 | 6,13 | 1,50 | 2,39 | |

# Evolution of the computing power of Supercomputers

**FLOP/second** (operaciones sobre números reales 64 bits)

$$\overset{E}{1}\,000\overset{P}{0}00\overset{T}{0}00\overset{G}{0}00_{000000}$$

~2018
? ($1 \times 10^7$ processadors

2008
Cray XT5 (15000 processadors)

1988
Cray Y-MP (8 processadors)

1998
Cray T3E (1024 processadors)

**FLOP/segon** (operaciones sobre números reales de 64 bits)

E       P      T     G

$1000000000000_{000000}$

#1 (55 PF)

Tianhe @ National University of Defense Technology, 54.9 PFlops

Samsung Exynos > 50 Gflops
Prototips MontBlanc @ BSC

#1 Espanya (1PF)
MN3 @ BSC

#1 EU (5PF)

JUQUEEN @ Forschungszentrum Jülich

**Barcelona**
**Supercomputing**
**Center**
Centro Nacional de Supercomputación

# Barcelona Supercomputing Center
# Centro Nacional de Supercomputación


EXCELENCIA SEVERO OCHOA

**«** **BSC-CNS objectives:**

– R&D in Computer, Life, Earth and Engineering Sciences

– Supercomputing services and support to Spanish and European researchers



**«** **BSC-CNS is a consortium that includes:**

– Spanish Government                                    51%

– Catalonian Government                              37%
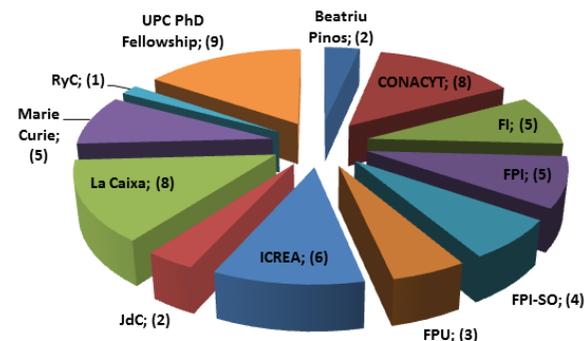
– Universitat Politècnica de Catalunya (UPC)    12%



**«** **425 people, 41 countries**



**BSC Staff Funding 2014 (425)**

Ordinary Budget; (80)

Personnel Grants; (58)

Competitive Projects; (287)



**Staff with Personnel Grants (58)**

UPC PhD Fellowship; (9)

Beatriu Pinos; (2)

RyC; (1)

CONACYT; (8)

FI; (5)

Marie Curie; (5)

La Caixa; (8)

FPI; (5)

ICREA; (6)

FPI-SO; (4)

JdC; (2)

FPU; (3)




**Barcelona Supercomputing Center**
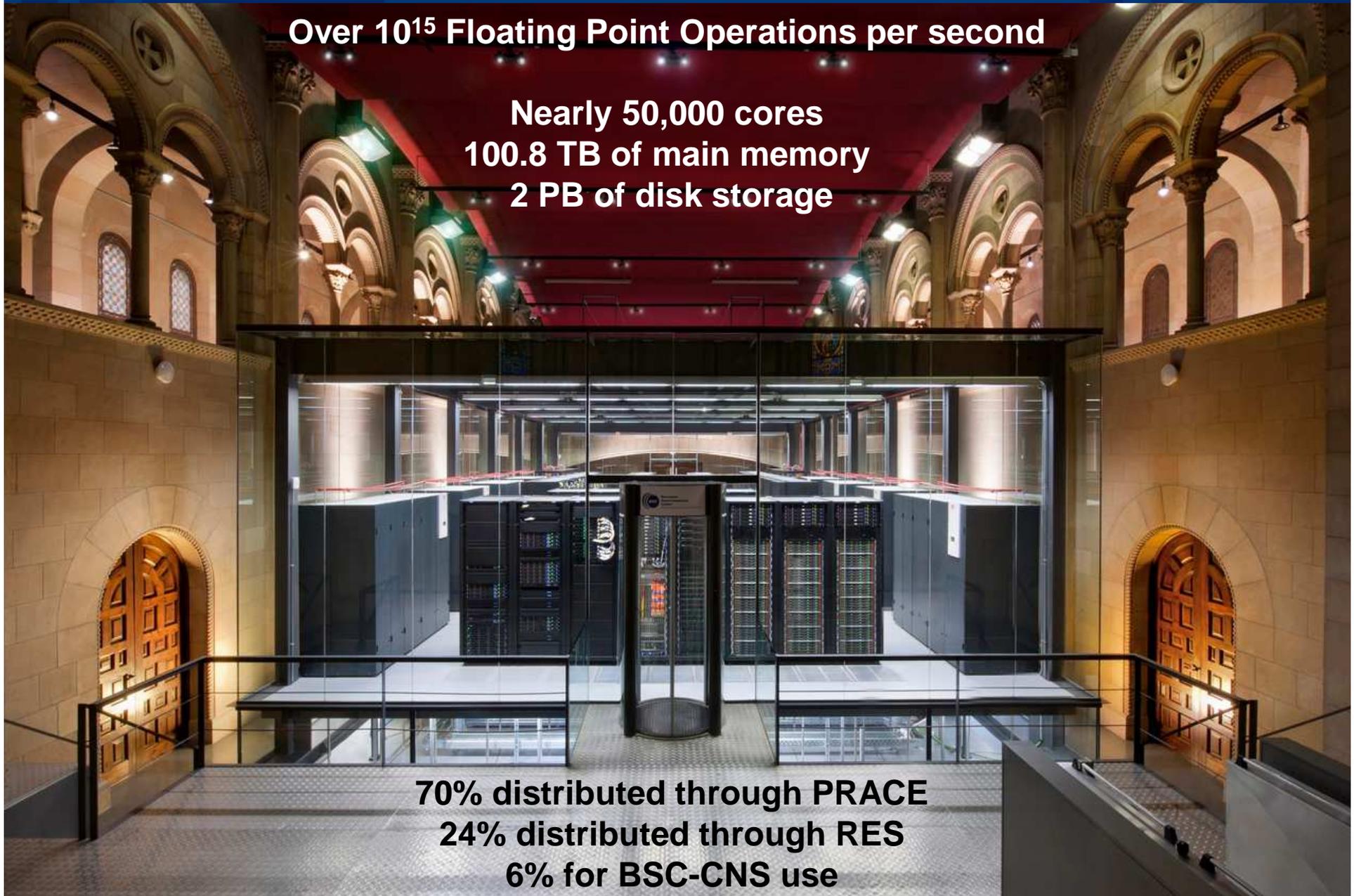Centro Nacional de Supercomputación

# The MareNostrum 3 Supercomputer

**Over $10^{15}$ Floating Point Operations per second**

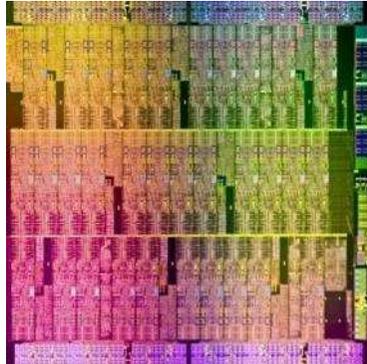**Nearly 50,000 cores**
**100.8 TB of main memory**
**2 PB of disk storage**

**70% distributed through PRACE**
**24% distributed through RES**
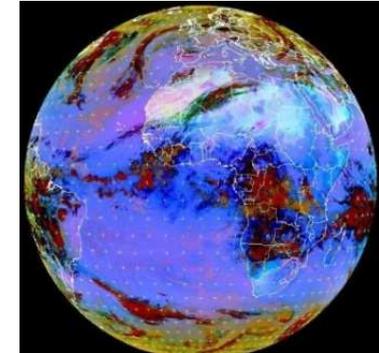**6% for BSC-CNS use**

EXCELENCIA
SEVERO
OCHOA

## COMPUTER SCIENCES

To influence the way machines are built, programmed and used: programming models, performance tools, Big Data, computer architecture, energy efficiency.
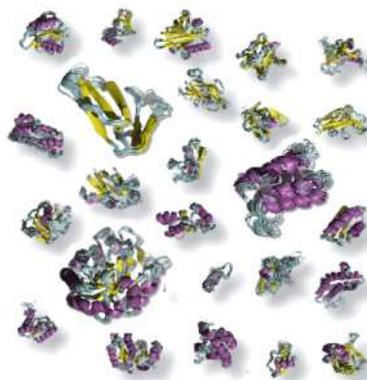


## EARTH SCIENCES

To develop and implement global and regional state-of-the-art models for short-term air quality forecast and long-term climate applications.
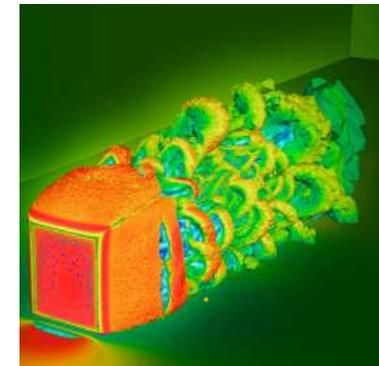


## LIFE SCIENCES

To understand living organisms by means of theoretical and computational methods (molecular modeling, genomics, proteomics).
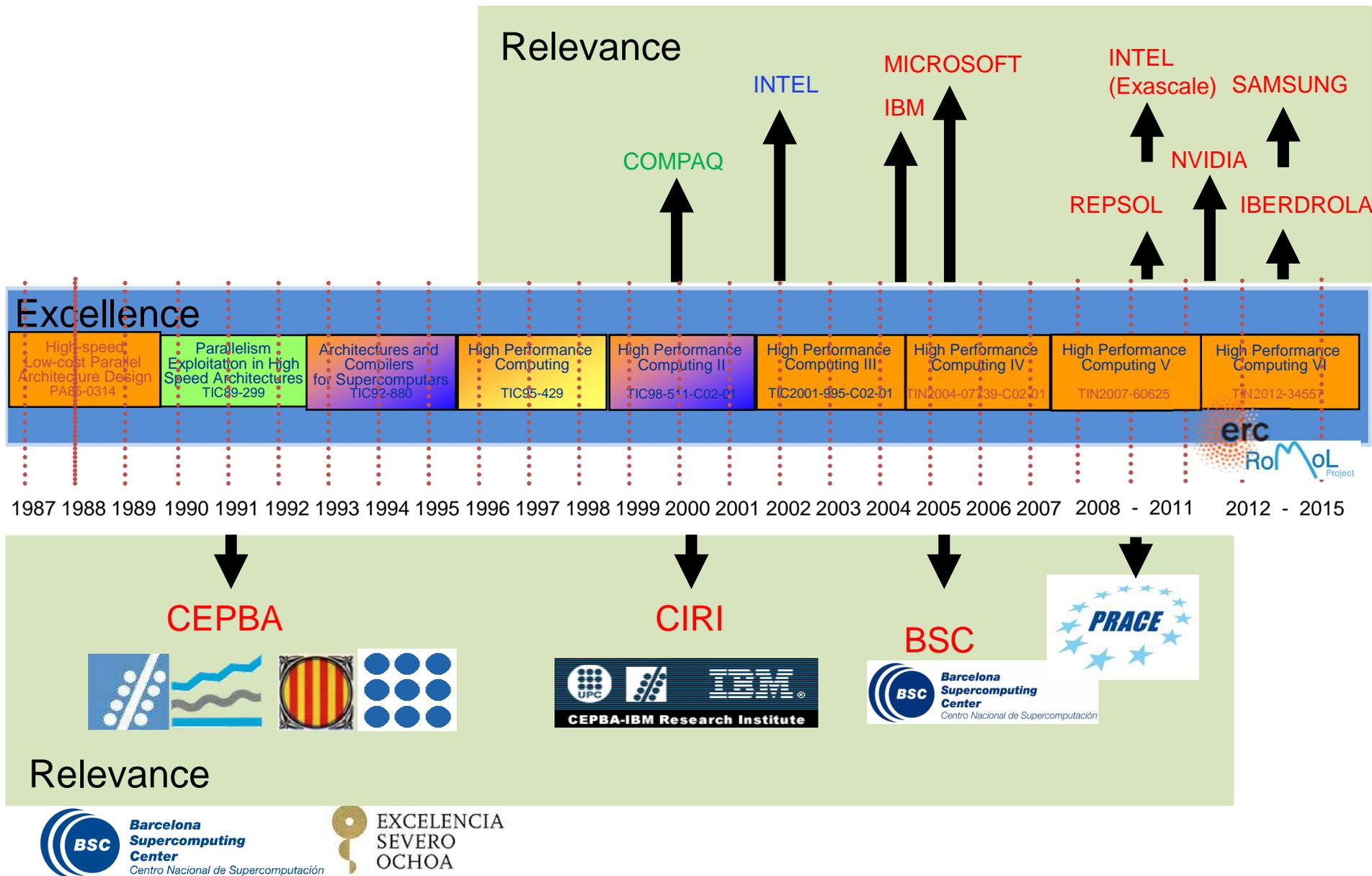


## CASE

To develop scientific and engineering software to efficiently exploit super-computing capabilities (biomedical, geophysics, atmospheric, energy, social and economic simulations).
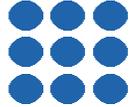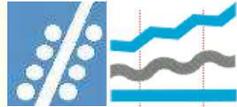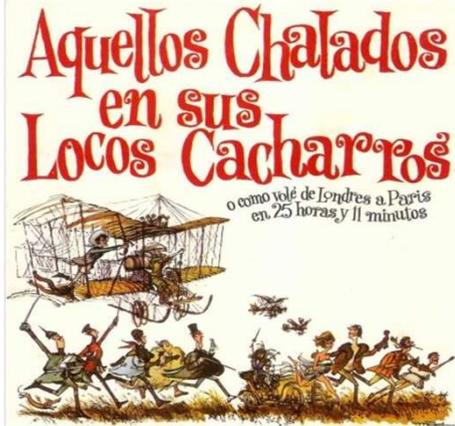
# Venimos de muy lejos …

Aquellos Chalados en sus Locos Cacharros
o como volé de Londres a Paris en 25 horas y 11 minutos

Transputer cluster

Research prototypes

Parsys Multiprocessor

Convex C3800

Connection Machine CM-200
0,64 Gflop/s

Parsytec CCi-8D
4.45 Gflop/s

SGI Origin 2000
32 Gflop/s

Compaq GS-140
12.5 Gflop/s

Compaq GS-160
23.4 Gflop/s

BULL NovaScale 5160
48 Gflop/s

IBM RS-6000 SP & IBM p630
192+144 Gflop/s

Maricel
14.4 Tflops, 20 KW

SGI Altix 4700
819.2 Gflops

SL8500
6 Petabytes

IBM PP970 / Myrinet
MareNostrum
42.35, 94.21 Tflop/s

CEPBA-IBM Research Institute

Barcelona Supercomputing Center
Centro Nacional de Supercomputación

1985  1986  1987  1988  1989  1990  1991  1992  1993  1994  1995  1996  1997  1998  1999  2000  2001  2002  2003  2004  2005  2006  2007  2008  2009  2010

Barcelona Supercomputing Center
Centro Nacional de Supercomputación

# Venimos de muy lejos …

**Ayto. Barcelona**
**Uitesa**
**UPC-EIO**

**AMES, CIMNE**

**TGI**
**UPM-DATSI**

**Hesperia**
**Neosystem**
**s**
**UPG-EIO**

**INDO, CEPBA-UPC**

**Soler y Palau**
**CIMNE**
**CEPBA-UPC**

**Torres**
**Soft Greenhouse**
**CEPBA-UPC**

**CEPBA**
**CESCA**
**UMA**
**UNICAN**
**UPM**

**Metodos Cuantitativos**
**Gonfiesa**
**CESCA, CESGA**

**Tecnatom, UMA**

**Iberdrola, Uitesa, UPV**

**AZTI**
**UPC-LIM**

**BCN COSIVER**
**Mides**
**UPC-EIO**

**CASA**
**Envision**
**GTD**
**Intespace**
**RUS**

**Ospedali Galliera**
**Le Molinette**
**Parsytec**
**PAC**
**EDS**

**ENEL**
**EDF**
**CSR4**
**Reiter**
**Kenijoki**

**Ferrari, Genias, P3C**

**ST Mecanica**
**DERBI**
**AUSA**
**CEPBA-UPC**

**Italeco**
**Geospace**
**Intecs**
**Univ. Leiden**

**Intera SP**
**Intera UK**
**UPC-DIT**
**CEPBA-UPC**

**Volkswagen**
**Ricardo**
**PAC**

**CANDEMAT**
**CIMNE**
**CEPBA-UPC**

**SENER**
**CIC**
**UNICAN**

**CEBAL-ENTEC**
**NEOSYSTEMS**

**Cari Verona**
**AIS**
**PAC**
**Univ.**
**Cat. Milan**

**Cristaleria Española**
**UNICAN**
**CEPBA-UPC**

**Inisel Espacio**
**Infocarto**
**UPC-TSC**
**CEPBA-UPC**

**Iberdrola**
**SAGE**
**CEPBA-UPC**

**Barcelona**
**Supercomputing**
**Center**
Centro Nacional de Supercomputación

# Scientific mission

## What

Influence the way machines are …
     … built …
     … programmed …
     … and used

## How

Through ideas, …
     … demonstration, …
     … cooperation with manufacturers, …
     … and "products"

**Performance, productivity, power/energy and reliability**

## Why

Our strength …
     … critical mass of people …
     … holistic/vertical vision/background …
     … stable and exploratory paths …
     … and co-design approach

| Programming environments | Performance tools |
|---|---|
| Resource management | Architecture, from single core HPC to other areas |

Supercomputing

BigData

Multicores real-time

Mobile and embedded

# Objectives

## What

Perform research in Earth sciences for the development, implementation and refinement of global and regional state-of-the-art models for short-term air quality forecast and climate applications

## Why

Our strength …
 … operations …
 … research …
 … service …
 … high resolution …

## How

Bringing together knowledge in atmospheric dynamics, natural and anthropogenic emissions, improvement of air quality forecasting, transport and dispersion of pollutants in complex terrain, urban air quality, aerosol optical properties, aerosol radiative effects and the feedback between meteorology and air pollution with the advances in the parallelization of air quality model codes

## What

The understanding of living organisms by means of theoretical methods.

## How

Use computational methods to get information and simulate biological systems …

… with the final goal of explaining biological systems from the basic rules of physics and chemistry

**Vision**

**Research**

Main Stream

Integration

Focus

## What

Develop relevant simulation software
…
   … science
   … engineering

## How

Close contact with industry …
   … innovative manufacturing
   … energy
   … pharmaceutical

## Why

Our strength …
   … multidisciplinary background
   … access to hardware
   … co-design approach

| ALYA | BSIT |
|------|------|
| FALL3D | PANDORA |

EXCELENCIA SEVERO OCHOA

## Over 106M€ in grants and contracts

BSC achieved the 9th greatest return from FP7 of all Spanish institutions according to CDTI report 11/11/2013

Fourteen H2020 projects approved or under negotiation (Grant value: 5.97M€)

EU
50.3M€
47%

National
23.5M€
22%

Companies
32.4M€
31%

**Barcelona Supercomputing Center**
Centro Nacional de Supercomputación

- Includes personnel grants and future income for approved projects 21/01/2015
- National income includes ICTS calls

**EXCELENCIA SEVERO OCHOA**

**BSC-IBM Technology Center for Supercomputing**

*Future challenges for supercomputers including power efficiency and scalability, new programming models, and tools for analysis and optimization of applications*

**BSC-NVIDIA CUDA Center of Excellence**

*Training in Parallel Programming using CUDA and StarSs Optimising management of execution resources in multi-GPU environments with GMAC*

**BSC-Microsoft Research Centre**

*Analysis of Hadoop workload performance under different software parameters and hardware configurations. Results available online*

**Intel-BSC Exascale Lab**

*Multi-year agreement focussing on optimising efficiency through research into Programming Models, Performance Tools and Applications*

*Agreement on memory performance in HPC systems with*

**SAMSUNG**

# Key Spanish Industrial Partners

## Repsol-BSC Research Center

Research into advanced technologies for the exploration of hydrocarbons, subterranean and subsea reserve modelling and fluid flows

## Iberdrola Renovables

IBERDROLA

Design and optimization of wind farms

JUAN YACHT DESIGN
Juan Kouyoumdjian · Naval architecture

ANAXOMICS
Systems biology solutions

INTELLIGENT PHARMA

TF
Termo Fluids

Telefónica
Telefónica
Investigación y Desarrollo

FUNDACIÓN BOTÍN

Mind the Gap: 500K €
Nostrum Drug Discovery Spin off

erc

Proof of Concept: 150K€ PELE

**Patents**
- dpEDMD: method for drug discovery
- SMUFIN: Somatic Mutations Finder
- miRNA markers for Morbid obesity diseases

# Some Strategic Projects

**Severo Ochoa**

*A multidisciplinary research program to address the complex challenges in the path towards Exascale. A set of key strategic scientific projects and improvements in HR management, training, mobility and communication.*
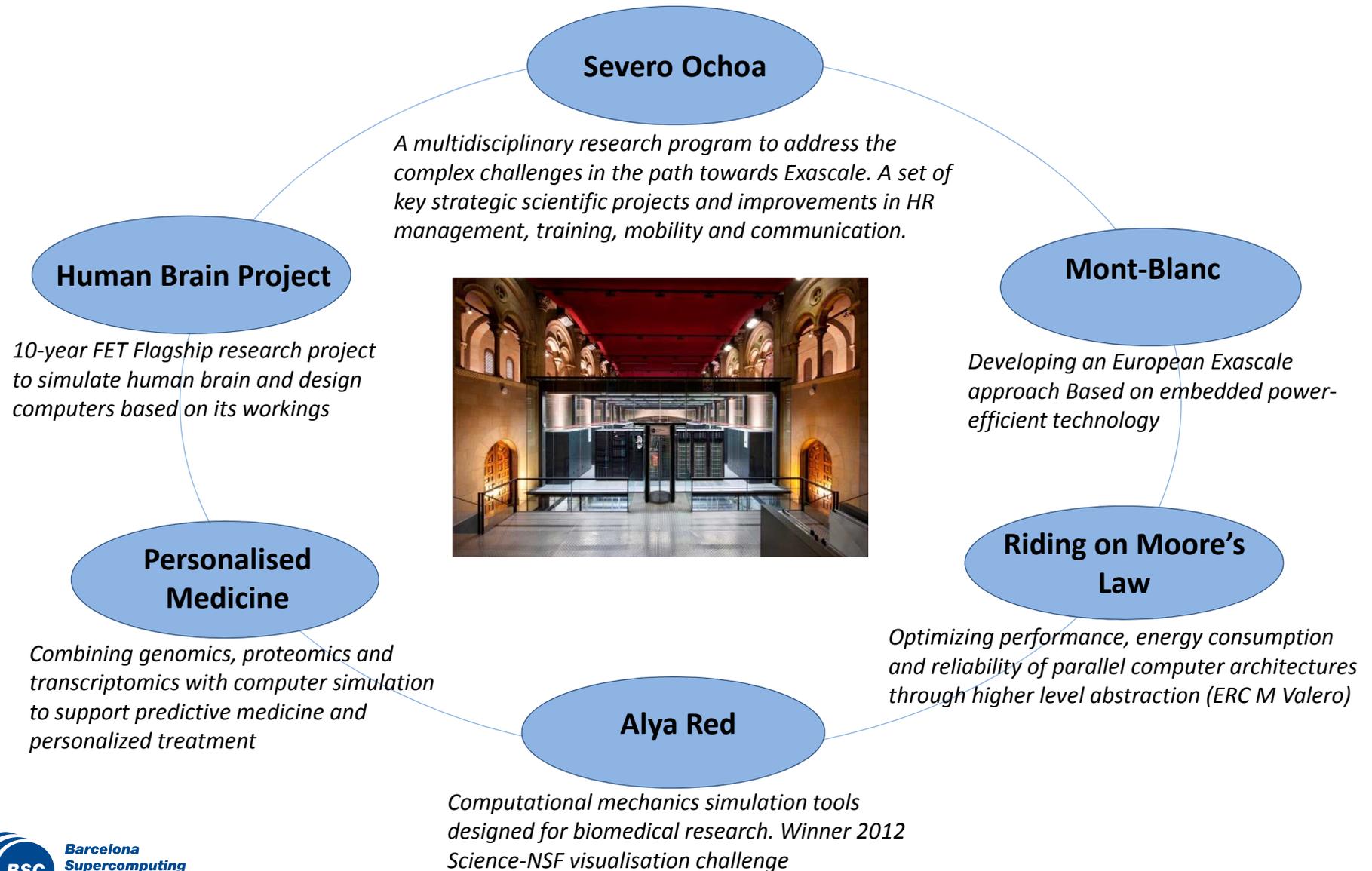
**Human Brain Project**

*10-year FET Flagship research project to simulate human brain and design computers based on its workings*

**Mont-Blanc**

*Developing an European Exascale approach Based on embedded power-efficient technology*

**Personalised Medicine**

*Combining genomics, proteomics and transcriptomics with computer simulation to support predictive medicine and personalized treatment*

**Riding on Moore's Law**

*Optimizing performance, energy consumption and reliability of parallel computer architectures through higher level abstraction (ERC M Valero)*

**Alya Red**

*Computational mechanics simulation tools designed for biomedical research. Winner 2012 Science-NSF visualisation challenge*

**Barcelona Supercomputing Center**
Centro Nacional de Supercomputación

BSC

21

# Internationalisation

### HPC Infrastructure for Europe's Best Scientists



### Over 90 EC Framework Projects



### Part of Future Spanish node of ICT-Labs



### Bridge between EU and Latin America

# Helping to define the future of global HPC

## Enabling the Data Revolution



## International Roadmapping



## Leadership in Exascale
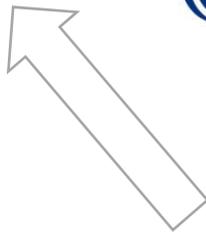


## Contributing to Standardisation



Barcelona
Supercomputing
Center
Centro Nacional de Supercomputación

**Barcelona Supercomputing Center**
Centro Nacional de Supercomputación

BSC-CNS, PRACE Hosting Member

BSC-CNS, ETP Founding Member
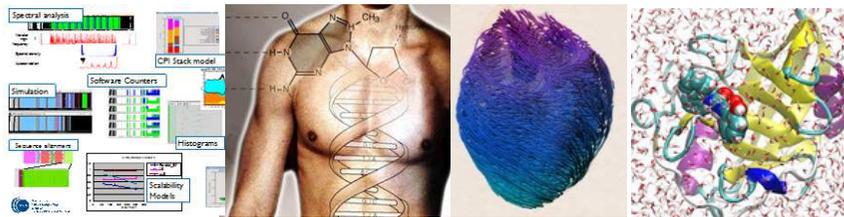
**PRACE**

**ETP 4 HPC**

THE EUROPEAN TECHNOLOGY PLATFORM FOR HIGH PERFORMANCE COMPUTING

Access to **best HPC resources** for industry and academia

**Autonomous** EU development of **Exascale technologies**

BSC-CNS leads two CoE proposals, participates in others

Centers of Excellence in **HPC applications**
SME Competence Centers

**Barcelona Supercomputing Center**
Centro Nacional de Supercomputación

## Maintain leadership and visibility

– On programming models and performance analytics

– More platforms, more intelligence

– More Apps, engage with communities
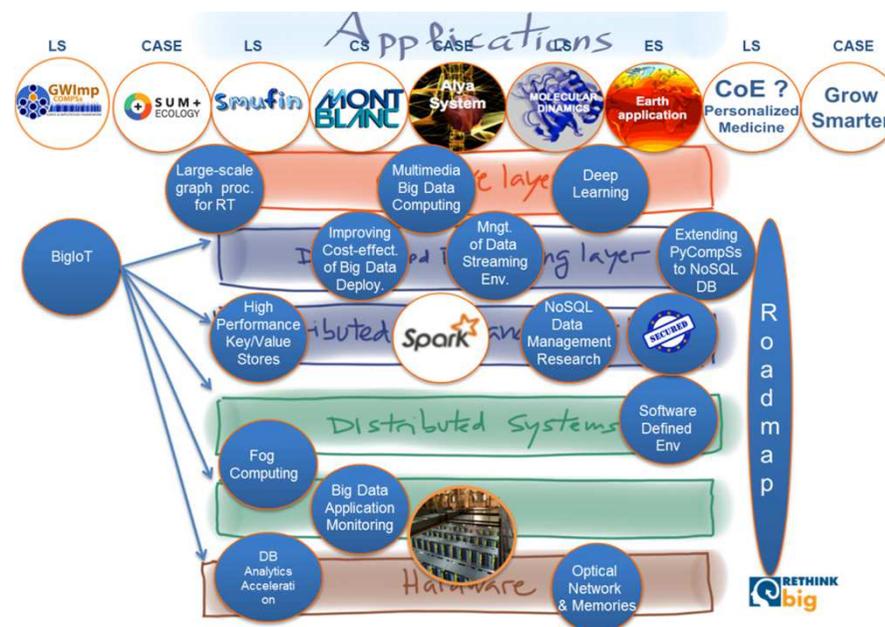
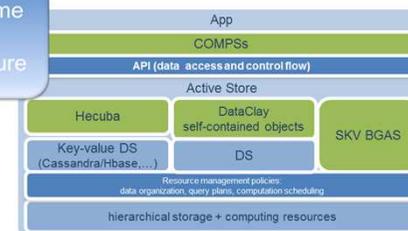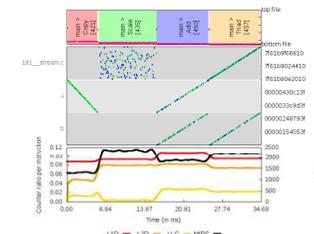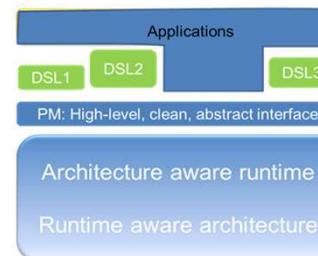– Influence standards

## Push forward

– Architectures for real-time

– Convergence of HPC and BigData

– Runtime-aware multicore architecture

- RoMoL, Mont-Blanc 3

## Further exploration

– Convergence of embedded and HPC, IoT

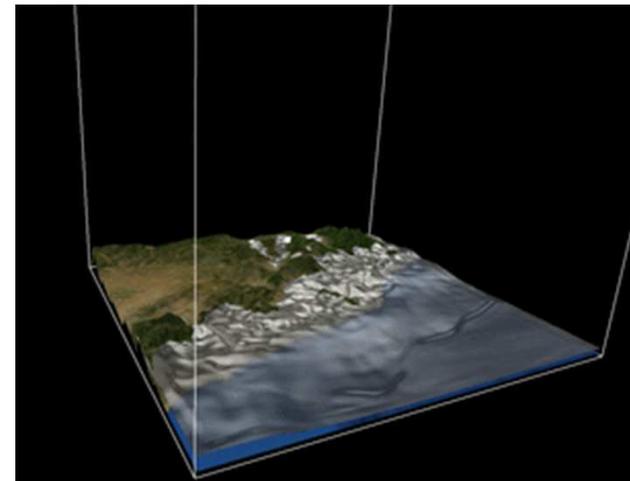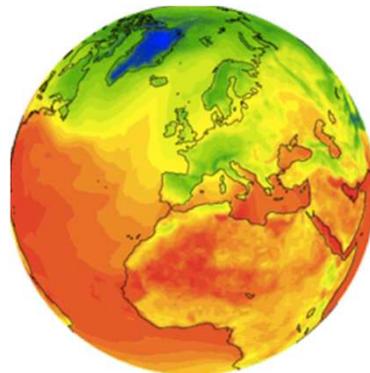– Algorithmic development in different domains

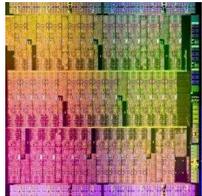– Cognitive computing

OpenMP
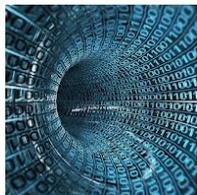Forerunner

# Environmental Forecasting for Services

- Development of the best knowledge of plausible predictions of the direction, scope, speed and intensity of a set of environmental changes.

- Bringing together the human activity and the environment through the most efficient combination of monitoring and modeling.

- Focused on the development of modeling solutions to provide weather, climate and air quality information at a global scale for the public and private sectors, with a special interest in the Mediterranean, African, Arctic and South American regions.
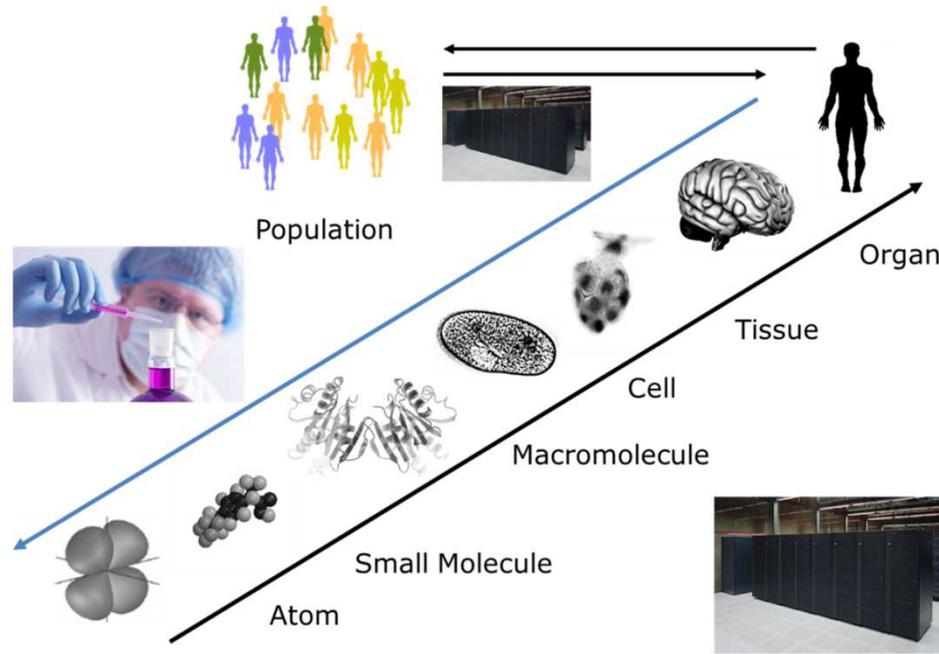
**Barcelona Supercomputing Center**
Centro Nacional de Supercomputación

EXCELENCIA SEVERO OCHOA

## Computer Sciences

Architecture

Big Data

Software

## Personalized Medicine Model



Population

Organ

Tissue

Cell

Macromolecule

Small Molecule

Atom

## Life Sciences
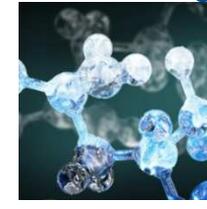
Genomics

Modeling

Chemistry

## Solutions for Biotech, Pharma and Clinical

| GENOTYPING | GWAS | DIAGNOSIS/ PROGNOSIS | LEAD GENERATION | DRUG REPOSITION | DRUG RESISTANCE |

- Integrate Personalized Medicine into the HPC ecosystem
- Implement an HPC platform serving the needs of the Personalized Medicine community

**Barcelona Supercomputing Center**
Centro Nacional de Supercomputación

# Alya software consolidation as standard for HPC modelling

Alya is the BSC in-house simulation code:
- Coupled multi-scale and multi-physics
- Complex simulation scenarios
- Parallel efficiency in supercomputers

In 2014 it becomes the first in its class, simulating three complex multi-physics problems in 100.000 cores ("2014 Top Supercomputing Achievement" by HPCWire and announced in SC).

The target is 1.000.000 cores in 2016: CODEX project

# HPC & Data Analytics service convergence

- Massive storage capabilities.

- Workflows definition.

- Adequate User Support skills.

**Barcelona
Supercomputing
Center**
*Centro Nacional de Supercomputación*

# THE BIG DATA ERA:
# DATA GENERATION EXPLOSION

Last year one of the most computer-intensive scientific experiments ever undertaken confirmed Peter Higgs and François Englert's theory by making the Higgs boson – the so-called "God particle" – in an $8bn atom smasher, the Large Hadron Collider at Cern outside Geneva.

" the LHC produces 600 TB/sec… and after filtering needs to store 25 PB/year"… 15 million sensors….

**TECHNOLOGY FEATURE**

# THE BIG CHALLENGES OF BIG DATA

As they grapple with increasingly large data sets, biologists and computer scientists uncork new bottlenecks.

« High resolution imaging

« Clinical records

« Simulations

« Omics

# Sequencing Costs



Source: **National Human Genome Research Institute** (NHGRI)
http://www.genome.gov/sequencingcosts/

(1)    "Cost per Megabase of DNA Sequence" — the cost of determining one megabase (Mb; a million bases) of DNA sequence of a specified quality

(2)    "Cost per Genome" - the cost of sequencing a human-sized genome. For each, a graph is provided showing the data since 2001

In both graphs, the data from 2001 through October 2007 represent the costs of generating DNA sequence using Sanger-based chemistries and capillary-based instruments ('first generation' sequencing platforms). Beginning in January 2008, the data represent the costs of generating DNA sequence using 'second-generation' (or 'next-generation') sequencing platforms. The change in instruments represents the rapid evolution of DNA sequencing technologies that has occurred in recent years.

**Prof. Mateo Valero – Big Data**

**Operations on SKA-2 expected to start by 2024** ← SKA2

UNIVERSITY OF CAMBRIDGE

~ 250 Dense Aperture Array Stations 300-1400MHz

~ 2700 Dishes

3-Core Central Region

~250 Sparse Aperture Array Stations 70-450MHz

Wide Band Single Pixel Feeds

Phased Array Feeds 800 MHz – 2GHz

Artist renditions from Swinburne Astronomy Productions

**Prof. Mateo Valero – Big Data**

SKA₂ wide area data flow

Prof. Mateo Valero – Big Data

## Current infrastructure sagging under its own weight

**2013**

**By 2020**

**Internet of Things**

**98,000** tweets

**23,148** apps downloaded

**400,710** ad requests

**In 60 sec today**

**2000** lyrics played on Tunewiki

**1,500** pings sent on PingMe

**208,333** minutes Angry Birds played

**30 Billion** (1)
**Devices**

**40 Trillion GB** (2)
**DATA**

**Mobile Apps** **10 Million** (3)

**… for 8 Billion** (4)

Pervasive Connectivity | Smart Device Expansion | Explosion of Information

**Barcelona Supercomputing Center**
Centro Nacional de Supercomputación

**EXCELENCIA SEVERO OCHOA**

**Prof. Mateo Valero – Big Data**

RoMoL Project

# Challenges of data generation

## Volume
### SCALE OF DATA

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

2005

2020

**6 BILLION PEOPLE**
have cell phones

WORLD POPULATION: 7 BILLION

It's estimated that
**2.5 QUINTILLION BYTES**
[ 2.3 TRILLION GIGABYTES ]
of data are created each day

Most companies in the U.S. have at least
**100 TERABYTES**
[ 100,000 GIGABYTES ]
of data stored

## 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

## Variety
### DIFFERENT FORMS OF DATA

As of 2011, the global size of data in healthcare was estimated to be
**150 EXABYTES**
[ 161 BILLION GIGABYTES ]

**30 BILLION PIECES OF CONTENT**
are shared on Facebook every month

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

**400 MILLION TWEETS**
are sent per day by about 200 million monthly active users

## Velocity
### ANALYSIS OF STREAMING DATA

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

By 2016, it is projected there will be
**18.9 BILLION NETWORK CONNECTIONS**
– almost 2.5 connections per person on earth

## Veracity
### UNCERTAINTY OF DATA

**1 IN 3 BUSINESS LEADERS**
don't trust the information they use to make decisions

**27% OF RESPONDENTS**
in one survey were unsure of how much of their data was inaccurate

Poor data quality costs the US economy around
**$3.1 TRILLION A YEAR**

Source: *http://www-01.ibm.com/software/data/bigdata/*

IBM

**Barcelona Supercomputing Center**
Centro Nacional de Supercomputación

**EXCELENCIA SEVERO OCHOA**

RoMoL Project

# The Data Deluge



**2005**
0.1ZB

**2010**
1.2ZB

**2012**
2.8ZB

**2015**
8.5ZB

**2020**
**40ZB***
(figure exceeds prior forecasts
by 5 ZBs)

**Barcelona Supercomputing Center**
Centro Nacional de Supercomputación

EXCELENCIA SEVERO OCHOA

**Prof. Mateo Valero – Big Data**

RoMoL Project

# How big is big?

Saganbyte, Jotabyte,…

**Geopbyte**

$10^{30}$

**This will take us beyond our decimal system**

**Brontobyte**

$10^{27}$

**This will be our digital universe tomorrow…**

**Yottabyte**

$10^{24}$

**This is our digital universe today** = 250 trillion of DVDs

$10^{21}$

Zettabyte

**1.3 ZB** of network traffic by **2016**

$10^{18}$

Exabyte

**1 EB** of data is created on the internet each day = **250 million DVDs** worth of information. The proposed **Square Kilometer Array telescope** will generated an **EB** of data per day

$10^{15}$

Petabyte

The **CERN Large Hadron Collider** generates **1PB** per second

$10^{12}$

Terabyte

**500TB** of new data per day are ingested in **Facebook** databases

$10^{9}$

Gigabyte

$10^{6}$

Megabyte

2025
2020
2015
2010
2005

# The data explosion is transforming science



| Experiments | Simulations | Archives | Literature | Consumer |

Petabytes
Doubling & Doubling

**((** Every area of science is now engaged in data-intensive research

**((** Researchers need

– Technology to publish and share data in the cloud

– Data analytics tools to explore massive data collections

– A sustainable economic model for scientific analysis, collaboration and data curation

**Barcelona**
**Supercomputing**
**Center**
*Centro Nacional de Supercomputación*

# BIG DATA TECHNOLOGIES:

## (( Big Data Technologies

– Storing data

– Processing data

– Where do we place data?

– Managing Big Data

# Magnetic Tape Memory

« Invented by Eckert & Mauchly for the UNIVAC I, March, 21,1951

– Model UNISERVO

– 224 KB of data

– 1/2 inches of diameter,1200 feets, 128 characters per inch

– Speed: 100 inches per second, equivalent to 12800 characters per second

« Storage Tech, 2013, T-10000D,

– 8.5 Terabytes (40.000.000 increase)

– EBW: 250 Mbytes/second (20.000 increase)

– Load Time: 10 seconds

**May 2014**

IBM and Fujifilm have demoed a **154TB** LTO-size tape cartridge which could come to market in **10 years' time.**

The Sony development involved a **148Gbit/in2** tape density and its own tape design to achieve a **185TB** uncompressed capacity.



Figure 40: Areal Density of Hard Disk and Tape Laboratory Demonstrations and Products [71].

148 Gbit/in2

3.1 Gbit/in2 (T10000C)

http://www.insic.org/news/2012Roadmap/12index.html



Image capture of resulting magnetic particles (top layer)

Comparison of magnetic particle size (image)

7.7nm — Sony's results

Several tens of nanometers — Currently

# HDD: Hard Drives Disk

- 1956, IBM 305 RAMAC
- 4 MB, 50x24" disks, 1200 rpm, 100 bits/track
- Intertracks: 0.1 inches, Density: 1000 bits/in2
- 100 ms access , Tubes, 35k$/y rent

- Year 2013: 4 Terabytes (1.000.000 increase)
- Average access time: few milliseconds (40 to 1)
- Areal: doubling in average every 2/4 years, but not now

- Predicted: 14 Terabytes in 2020 at the cost of $40

# How We Increase 10x and Beyond…

**《 Seagate Step up to SAS; SAS Roadmap - Source: SCSI Trade Association**

**Capacity**  **Drive Technologies**  **Bandwidth and Connectivity**



Example features for capacity driven usage

# Areal Density

Seagate Storage Effect Web Site

| Shingled Magnetic Recording (SMR) | Heat-Assisted Magnetic Recording(HAMR) | Future Candidate(s) Single Molecule Magnets |
|---|---|---|
| **Projected 25%** Areal Density Increase | **Projected 55%** Areal Density Increase | **Envisioned 1,000x** Areal Density Increase |
| **Enables +10TB and above** | **Theoretical up to 30TB to 60TB Hard Drive** | **Theoretical up to 3000TB Hard Drive** |

**((** Can *"trade"* memory capacity for other metrics of interest – e.g. bandwidth?

**((** Packaging
  – 2.5D stacking
  – 3D stacking

**((** Technologies
  – Wide I/O
  – Hybrid memory cube

# Emerging (non volatile) Memories

| | Traditional Technologies | | | | Emerging Technologies | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | | Improved Flash | | | | | | |
| | DRAM | SRAM | NOR | NAND | FeRAM | MRAM | PCM | Memristor | NEMS |
| Cell Elements | 1T1C | 6T | 1T | | 1T1C | 1T1R | 1T1R | 1M | 1T1N |
| Half pitch (F) (nm) | 50 | 65 | 90 | 90 | 180 | 130 | 65 | 3-10 | 10 |
| Smallest cell area ($F^2$) | 6 | 140 | 10 | 5 | 22 | 45 | 16 | 4 | 36 |
| Read time (ns) | < 1 | < 0.3 | < 10 | < 50 | < 45 | < 20 | < 60 | < 50 | 0 |
| Write/Erase time (ns) | < 0.5 | < 0.3 | $10^5$ | $10^6$ | 10 | 20 | 60 | < 250 | 1ns(140ps-5ns) |
| Retention time (years) | seconds | N/A | > 10 | > 10 | > 10 | > 10 | > 10 | > 10 | > 10 |
| Write op. Voltage (V) | 2.5 | 1 | 12 | 15 | 0.9-3.3 | 1.5 | 3 | < 3 | < 1 |
| Read op. Voltage (V) | 1.8 | 1 | 2 | 2 | 0.9-3.3 | 1.5 | 3 | < 3 | < 1 |
| Write endurance | $10^{16}$ | $10^{16}$ | $10^5$ | $10^5$ | $10^{14}$ | $10^{16}$ | $10^9$ | $10^{15}$ | $10^{11}$ |
| Write energy (fJ/bit) | 5 | 0.7 | 10 | 10 | 30 | $1.5{\times}10^5$ | $6{\times}10^3$ | < 50 | < 0.7 |
| Density (Gbit/cm$^2$) | 6.67 | 0.17 | 1.23 | 2.47 | 0.14 | 0.13 | 1.48 | 250 | 48 |
| Voltage scaling | Fairly scalable | | | | | no | poor | promising | promising |
| Highly scalable | Major technological barriers | | | | poor | | promising | promising | promising |

Sources: Chong et al ICCAD09, Eshraghian TVLSI11, ITRS13

## )) Big Data Technologies

- Storing data
- Processing data
- Where do we place data?
- Managing Big Data

## Moore's Law + Memory Wall + Power Wall

### Chip MultiProcessors (CMPs)



POWER4 (2001)

Intel Xeon 7100 (2006)

UltraSPARC T2 (2007)

Legend:
- Transistors (Thousands)
- Frequency (MHz)
- Power (W)
- Cores

**Prof. Mateo Valero – Big Data**

RoMoL Project

# Fujitsu SPARC64 XIfx

- 32 computing cores (single threaded) + 2 assistant cores
- 24MB L2 sector cache
- 256-bit wide SIMD
- 20nm, 3.75M transistors
- 2.2GHz frequency
- 1.1TFlops peak performance
- High BW interconnects
  - HMC (240GB/s x 2 in/out)
  - Tofu2 (125GB/s x 2 in/out)

## Knights Corner (2011)

- Coprocessor, 61x86 cores, 22nm, AVX-512, 4 HTs
- 1.2TFLOPS (DP), 300W TDP, 4 GFLOPS/W
- 512KB/core L2 coherent
- Int Netw: Ring
- Mem BW: 352GB/s

## Knights Landing (exp 2015)

- Coprocessor or **host processor**
- 72 Atom cores, 14nm, AVX512 per core, 4 HTs
- Up to 16GB of DRAM 3D stacked on-package, 384GB GDDR
- 3TFLOPS (DP), 200W TDP, 15GFLOPS/W

- **DP Performance:**
  1.43 Tflop

- **Mem BW (ECC off):**
  288 GB/s

- **Memory size (GDDR5):**
  12 GB

- **15 SMX units**
  - 192 single-precision CUDA cores
  - 64 double-precision units
  - 32 special function units
  - 32 load/store units

- **Six 64-bit memory controllers**

# Graph500 vs. Top500

- Top500 defined a benchmark (Linpack) to rate HPC machines upon performance. This benchmark is not suitable to address the characteristics of Big Data applications.

  - Linpack:
    - computation bounded
    - focused on floating-point operations
    - Bulk-Synchronous-Parallel model: behavior based on big computation and short communication bursts
    - dense data structures highly organized and coalesced (spatial locality)

  - Graph500:
    - communication bounded
    - focused on integer operations
    - asynchronous spatial uniform communication interleaved with computation
    - larger sparse datasets (very low spatial and temporal locality)

- Green Graph 500 list:
  - Collects performance-per-watt metrics
  - To compare the energy consumption of data intensive computing workloads.

  Graph500: graph500.org  Top500: top500.org

❰❰ **Aimed to fix data movement bottleneck**

❰❰ **Based on NVlinks**

- chip-to-chip communication approach

- comprised of bi-directional 8-lane links

- provide between 80 and 200 GB/s of bandwidth

❰❰ **This approach is expected to provide 4x speedups w. r. t. current GPU-based designs**

Source: Bob Broderson, Berkeley Wireless group

**Prof. Mateo Valero – Big Data**

# Outline

## ❰❰ Big Data Technologies

- Storing data
- Processing data
- Where do we place data?
- Managing Big Data

## 100.000 subjects suffering different diseases



**1 PB of compressed data**

**If we ever had a 1PB disk (100MB/s)**

**scanning 1 Petabyte:**

**3,000 hours / 125 days**

Supercomputing is about doing things FAST…

What if we want to process the data in
1 hour ?

Prof. Mateo Valero – Big Data

Evolution of storage architecture for Big Data

**Compute Network**
**(40Gbps Node Adapter)**

**Compute Nodes**

**Storage Network**
**(1Gbps Node Adapter)**

**11hrs for 1 PB**

**Storage Racks**

**1PB**

# Solid State Hybrid (SSHD) Technology

**Reduces Costs**
**Now practical to employ enterprise-class SLC NAND flash**

**SSHD: Capacity, performance, and value**

**HDD large capacity + SSD high speed**

**Adaptive Memory™ technology**

- Self-learning software algorithms make HDD SLC NAND flash work together.
- Enables SSD-like performance when accessing most frequently-used files
- Reduces workload and increases reliability of SLC NAND flash

- More energy to move data than to compute on it
  - Computation almost feels "free" relative to communication
  - Time will make this worse
- There are two long poles in the communication energy tent
  - Memory
  - Storage
- This is a predicate for this talk

| Operation | Energy (pJ) |
|---|---|
| 64-bit integer operation | 1 |
| 64-bit floating-point operation | 20 |
| 256 bit on-die SRAM access | 50 |
| 256 bit bus transfer (short) | 26 |
| 256 bit bus transfer (1/2 die) | 256 |
| Off-die link (efficient) | 500 |
| 256 bit bus transfer(across die) | 1,000 |
| DRAM read/write | 16,000 |
| HDD read/write | $O(10^6)$ |

Source: Processors and sockets: What's next? Greg Astfalk / Salishan Conference / April 25, 2013

**Old**

**Compute-centric Model**

**New**

**Data-centric Model**

Manycore    Accelerators



Massive Parallelism
Persistent Memory

Flash    Storage Class Mem

Source: Heiko Joerg http://www.slideshare.net/schihei/petascale-analytics-the-world-of-big-data-requires-big-analytics

**Barcelona Supercomputing Center**
Centro Nacional de Supercomputación

EXCELENCIA SEVERO OCHOA

**Prof. Mateo Valero – Big Data**

RoMoL Project

Data Centric Deep Computing

Trends in Transaction Processing
(example is current Neteeza)

Traditional Computing

Blue Gene Active Storage (BGAS) Concept

IBM

"How to" guide:
- Remove 512 of 1024 BG/Q compute nodes in rack – to make room for solid state storage
- Integrate 512 Solid State (Flash+) Storage Cards in BG/Q compute node form factor

Standard BG/Q Compute Fabric

BQC Compute Card

Node card 16 BQC + 16 PCI Flash cards

PCIe Flash Board

10GbE    FPGA    PCIe

| Flash Storage | 2012 Targets |
|---|---|
| Capacity | 2 TB |
| I/O Bandwidth | 2 GB/s |
| IOPS | 200 K |

System Software Environment
- Linux OS enabling storage + embedded compute
- OFED RDMA & TCP/IP over BG/Q Torus – failure resilient
- Standard middleware – GPFS, DB2, MapReduce, Streams

Active Storage Target Applications
- Parallel File and Object Storage Systems
- Graph, Join, Sort, order-by, group-by, MR, aggregation
- Application specific storage interface

512 Hs4 Cards

BGAS Rack Targets

| Nodes | 512 |
|---|---|
| Storage Cap | 1 PB |
| I/O Bandwidth | 1 TB/s |
| Random IOPS | 100 Million |
| Compute Power | 104 TF |
| Network Bisect. | 512 GB/s |
| External 10GbE | 512 |

... scale it like BG/Q.

Key architectural balance point: All-to-all throughput roughly equivalent to Flash throughput

© 2013 IBM Corporation

14

Blue Gene Active Storage

© 2013 IBM Corporation

# Microsoft Catapult

Doug Burger From Microsoft Research
Talks About Project Catapult Which Will
Make Bing Twice As Fast
(*Jun 17, 2014*)

http://microsoft-news.com/doug-burger-from-microsoft-research-talks-about-project-catapult-which-will-make-bing-twice-as-fast-video/

*Microsoft is planning to replace traditional CPUs in data centers with field-programmable arrays, or FPGAs, processors that Microsoft could modify specifically for use with its own software. These FPGAs are already available in the market and Microsoft is sourcing it from a company called Altera. The FPGAs are 40 times faster than a CPU at processing Bing's custom algorithms.*

*June 11, 2014*

**hp** next

*HP unveils "The Machine"*

*It uses clusters of special-purpose cores, rather than a few generalized cores; photonics link everything instead of slow, energy-hungry copper wires; memristors give it unified memory that's as fast as RAM yet stores data permanently, like a flash drive.*

*A Machine server could address 160 petabytes of data in 250 nanoseconds; HP says its hardware should be about six times more powerful than an existing server, even as it consumes 80 times less energy. Ditching older technology like copper also encourages non-traditional, three-dimensional computing shapes, since you're not bound by the usual distance limits.*



**Performance estimates – transaction speed**
What could you do with 168 GUPS?

160 GUPS
28.8 GUPS
Performance

12,600 kW
160 kW
Power

- The Machine
- Fujitsu K

Fujitsu K – current record holder
- 73,000 SPARC nodes

HP New Architecture
- 8 racks
- 64 byte packet
- 256 SoCs per rack
- 24 cores per SoC, 122K total
- 2 GHz cores
- 256 GB NVM per SoC, 1.3PB total
- 256 NICs per rack
- 2*100 Gbps links per NIC
- Utilization < 70%
- Network interface bandwidth 140M 64-byte transactions per second
- Unrestricted bisection bandwidth



**The Machine**

Special purpose cores — Photonics — Massive memory pool



**Future History**

**2015**
- Memristors begin sampling
- Physical infrastructure of Core prototypes established
- Silicon Scale Integration discipline established across HP and supply chain
- Open Source Machine OS SDK released
- ISV Partner collaborations begin

**2017**
- Edge devices begin sampling
- Machine OS enters public beta

**2019**
- Core devices ship
- Machine available as product, service and as a business process transformation

**2016**
Memristor DIMMs launched
ISV partners collaborations ship as services

**2020**
Distributed mesh compute goes mainstream

**2014**
- Memristor media controller, protocols and standards established
- SoC Partners selected for co-development
- Machine OS development begins

**2018**
- Edge devices ship in volume
- Core Machine POCs at scale
- Machine OS released

# Outline

**《 Big Data Technologies**

- – Storing data
- – Processing data
- – Where do we place data?
- – Managing Big Data

NoSQL for non-structured data, eventual consistency

- **To meet the challenges: MapReduce**

  - **Programming Model** introduced by Google in early 2000s to support **distributed computing** (special emphasis in fault-tolerance)

    - Ecosystem of big data processing tools
      - open source, distributed, and run on commodity hardware.

  - The key innovation of MapReduce is
    - the ability to **take a query** over a data set, **divide it**, and **run it** in parallel over many nodes.

  - Two phases
    - **Map phase**
    - **Reduce phase**

Input Data

Mappers

Reducers

Output Data

**2006**
The Apache Hadoop Project brings an opensource MapReduce Implementation, with the team of the Apache Nutch Crawler and the support of Yahoo!

**Batch only**

**Graph Processing**

Pregel: A System for Large-Scale Graph Processing

Paper by: Grzegorz Malewicz, Matthew Austern, Aart Bik, James Dehnert, Ilan Horn, Naty Leiser, Grzegorz Czajkowski (Google, Inc.)

APACHE GIRAPH

**2004**
introduces
**MapReduce**

**Batch & Stream**

Announcing Cloud Dataflow

**2010**
Google announces Pregel used for building incremental reverse indexes instead of MapReduce
Phylosophy: '*think like a vertex*" (event-driven)

**June 25, 2014**
Google announces Cloud Dataflow: write code once, run it in batch or stream mode

**Cloud Dataflow** is a managed service for creating data pipelines that ingest, transform and analyze data in **both batch and streaming modes**.

**Prof. Mateo Valero – Big Data**

**Barcelona Supercomputing Center**
*Centro Nacional de Supercomputación*

# RESEARCH IN BIG DATA AT BSC

# Abstraction of computer middleware



Applications — e.g. Genomics

Distributed Processing layer — e.g. MapReduce

Distributed Data Management layer — e.g. NoSQL DB

Distributed Systems — e.g. Cloud

OS — e.g. Linux

Hardware — e.g. FPGAs, GPUs

# Urban semantics

IBM  **B**
**EC**OLOGIA
**B**
**N**

UN@HABITAT
FOR A BETTER URBAN FUTURE

HORIZON **2020**
LE PROGRAMME DE RECHERCHE ET
D'INNOVATION DE L'UNION EUROPÉENNE

**+ S U M +**
ECOLOGY

KPIs for urban resilience

**Grow Smarter**

Models a set of KPI on city mobility, sustainability, energy efficiency,…

Extending an ontology that integrates urban concepts

Stockholm, Köln, Barcelona cities. Urban Pollution model: to reduce transport emissions by 60%. Reducing energy consumption in buildings. Smart lighting and better communication facilities. Create 1500 jobs.

Distributed Processing layer

Distributed ~~nagement layer~~

Urban concepts Ontology

Distributed Systems

OS

Hardware

Ajuntament de **Barcelona**

Stadt Köln

**Barcelona Supercomputing Center**
Centro Nacional de Supercomputación

BSC

SEVERO OCHOA
SEVERO OCHOA

RoMoL Project

79

**RETHINK big**

Bring together the key European hardware, networking, and system architects with the key producers and consumers of Big Data to identify the industry coordination points that will maximize European competitiveness

Series of international workshops with participation of worldwide experts to address the challenges of both extreme scale computing and big data focusing on its convergence

**AXLE**

Focus on automatic scaling of complex analytics, while addressing the full requirements of real data sets.

Roadmap

Convergence

Distributed Processing layer

Distributed Data Management

Distributed Systems

Hardware

Tech transfer

DB Analytics Acceleration

Optical Network & Memories

**BIGDATACoE** BARCELONA

Technology transfer to SMEs in the Barcelona region

New architectures developed using optical interconnections and optical memory

Barcelona Supercomputing Center
Centro Nacional de Supercomputación

SEVERO OCHOA

EXCELENCIA SEVERO OCHOA Project

**FOG Computing**

Fog distributed platform capable of processing the data close to the end-points, becoming an excellent platform for Internet of Things (IoT)

**erc**
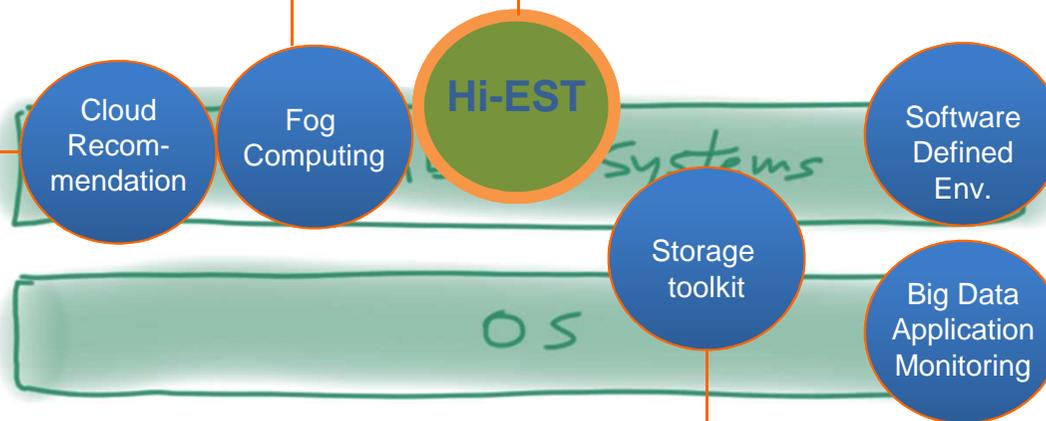
Holistic Integration of Emerging Supercomputing Technologies Automatic Optimization of Programmable Data Centers

**ca** technologies

Enhance a decision support system for the selection of the right set of Clouds
Trade-off between cost, reliability, risks and quality impacts

**IBM**

Building cloud environments embracing hardware and network heterogeneity to host a variety of workloads (JSA-SDN)

- Cloud Recom-mendation
- Fog Computing
- **Hi-EST**
- Software Defined Env.
- Storage toolkit
- Big Data Application Monitoring

Systems

OS

Creation of a Software-defined Storage toolkit for Big Data on top of the OpenStack platform

Design and implement a set of tools to monitor and to trace network traffic for hadoop applications
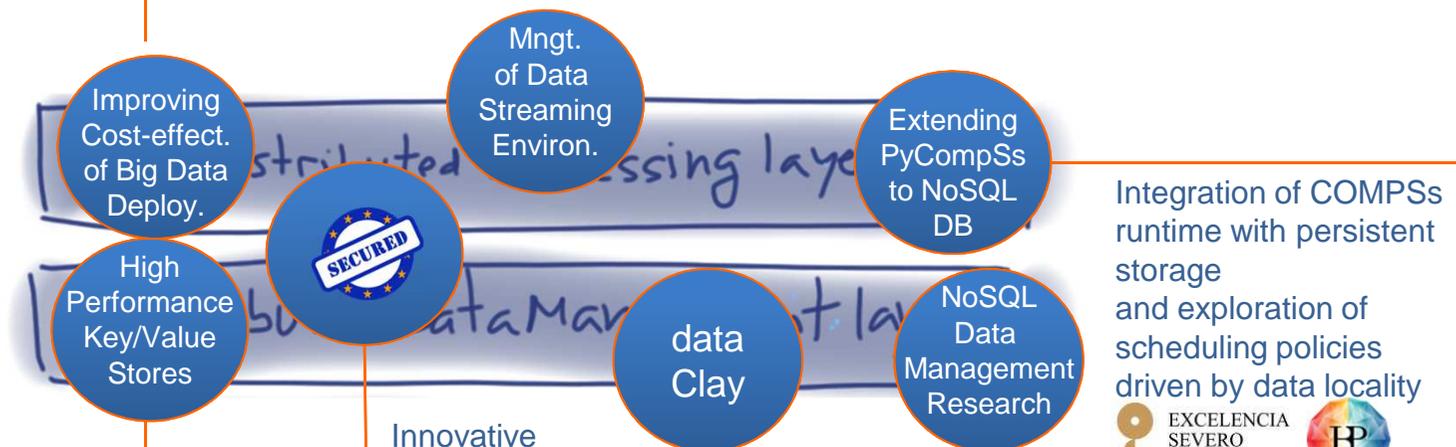
**IOSTACK**

**Lightness**

**RoMoL** Project

**Barcelona Supercomputing Center** Centro Nacional de Supercomputación

**EXCELENCIA SEVERO OCHOA**

**Microsoft**

Automated characterization of cost-effectiveness of Hadoop deployments (runtime performance vs software and hardware configuration choices)

**compose**

Explore novel architectures of the emerging **IoT stream processing platforms**, that provide the capabilities of data stream composition, transformation and filtering in real time

Improving Cost-effect. of Big Data Deploy.

Mngt. of Data Streaming Environ.

Extending PyCompSs to NoSQL DB

Integration of COMPSs runtime with persistent storage and exploration of scheduling policies driven by data locality

EXCELENCIA SEVERO OCHOA

High Performance Key/Value Stores

SECURED

data Clay

NoSQL Data Management Research

**IBM**

Explore the use of high performance key/value databases for fast persistent memory technologies

Innovative architecture to achieve protection from Internet threats by offloading execution of security applications into programmable devices at the edge of the network

Self-contained objects library

EXCELENCIA SEVERO OCHOA

Design and implement a software layer to enable NoSQL databases to decouple data organization from data model and provide NoSQL databases with efficient multi-level indexing

**Barcelona Supercomputing Center**
Centro Nacional de Supercomputación

EXCELENCIA SEVERO OCHOA

support

**RoMoL** Project

82

# Ongoing Big Data projects - V

**GWImp COMPSs**
GWAS & IMPUTATION FRAMEWORK

Large-scale Imputation and Genome-wide Association studies.

**Smufin**
Somatic Mutation Finder on cancer genomes.

Large-scale Genome analysis.

Large-scale Input data analysis

Use of statistical imputation theory to infer unknown genetic variants.
Large input data analysis and efficient parallel computation.
Big-data analytics

Applications

Optimized DNA compression algorithms, such as BWA, Suffix Arrays, FM-Indexing, etc., to compare multi-patients DNA information .
Big-data analytics

Distributed Processing layer

Distributed Data Management layer

Efficient Big-data genomic workflows deployment in current and future: HPC Clusters, Clouds, etc.

CLOUD

OS

Hardware

New architectures development for Computational Genomics tasks

**Barcelona Supercomputing Center**
Centro Nacional de Supercomputación

OCHOA

RoMoL Project

83

# Ongoing Big Data projects :Global picture

LS    CASE    LS    CS    CASE    LS    ES    LS    CASE

GWImp COMPSs — GWAS & IMPUTATION FRAMEWORK

SUM + ECOLOGY

Smufin

MONT BLANC

Alya System

MOLECULAR DINAMICS

Earth application

CoE ? Personalized Medicine

Grow Smarter

Applications

BigIoT

Improving Cost-effect. of Big Data Deploy.

Mngt. of Data Streaming Env.

Extending PyCompSs to NoSQL DB

High Performance Key/Value Stores

NoSQL Data Management Research

SECURED

Software Defined Env

Distributed Systems

Fog Computing

Big Data Application Monitoring

OS

DB Analytics Acceleration

Hardware

Optical Network & Memories

Roadmap

Barcelona Supercomputing Center
Centro Nacional de Super

RETHINK big VoL Project

84

**Barcelona
Supercomputing
Center**
*Centro Nacional de Supercomputación*

# BIG DATA ANALYTICS

**《** the **technology is often so subtle** that consumers have no idea that big data is actually helping make their lives easier

Source: www.DesignedInBarcelona.com

# Example: Online Shopping

**((** **Amazon's recommendation engine** uses big data and its database of around 250 million customers to suggest products by looking at previous purchases and other variables.



**((** Amazon are also developing a new technology which predicts what items you might want and sends it to your nearest delivery hub, (meaning faster deliveries for you)

**《** In 2012 identified a pregnant teenager before her family knew about her condition.

- The company can analyze customers' purchasing habits by monitoring credit card data, coupon usage, customer help lines, and emails for specific activities associated with pregnancy.

- They've identified 25 items that, when purchased in a particular order.

**❰❰ To do so, they are using <u>predictive models</u>**

**a collection of mathematical and programming techniques used to determine the probability of future events, analyzing historic and current data to create a model to predict future outcomes.**



Predictive Modeling

**Today, predictive models form the basis of many of the things that we do online: search engines, computer translation, voice recognition systems, etc.**

Source: http://bigsonata.com/wp-content/uploads/2014/07/PredictiveModeling.jpg

ability of computer systems to improve their performance by exposure to data without the need to follow explicitly programmed instructions.

# Computing waves

**((  We are now at a turning point of the history of computing**



The **first wave** of computing made **numbers** computable

The **second wave** has **made text and rich media** computable and accessible digitally

# Computing waves

**((** We are in the **next, third wave** that will also make **context computable**



Systems that embed predictive capabilities, providing the right functionality and content at the right time, for the right application, by continuously learning about them and predicting what they will need.

For example identify and extract context features such as hour, location, task, history or profile to present an information set that is appropriate for a person at a specific time and place.

# New self-learning systems are required

**《** Today computers require programming, and by definition programming does not allow for alternate scenarios that have not been programmed

**《** To allow alternating outcomes would require going up a level, creating a self-learning systems



The general idea is that instead of instructing a computer what to do, we are going to simply throw data at the problem and **tell the computer to figure it out itself**.

# Cognitive?

**❰❰** For this purpose the computer software takes functions from the brain like: inference, prediction, correlation, abstraction, … giving to the systems to possibility to do this by themselves.

Giving computers a greater ability to understand information, and to learn, to reason, and act upon it



**❰❰** And here it comes the use of cognitive word to describe this new computing!

# Augment our reasoning capabilities



*These new systems will raise the potential to augment our reasoning capabilities and empower us to make better informed decisions in order to address complex situations that are characterized by ambiguity and uncertainty.*

# Cognitive Computing

**《 Its meaning is not clear yet …**

The term "cognitive computing" remains a bit confusing since it covers systems that use different analytic approaches.



(*) Others use Smart Computing, Intelligent Computing, …

# Why Now?



1. Along the explosion of data …

now algorithms can be "trained" by exposing them to large data sets that were previously unavailable.



2. And the computing power necessary to implement these algorithms are now available

# Cognitive Computing: Foundational Building Blocks

## Foundational Building Blocks

1. HPC resources
2. Big Data  Technologies
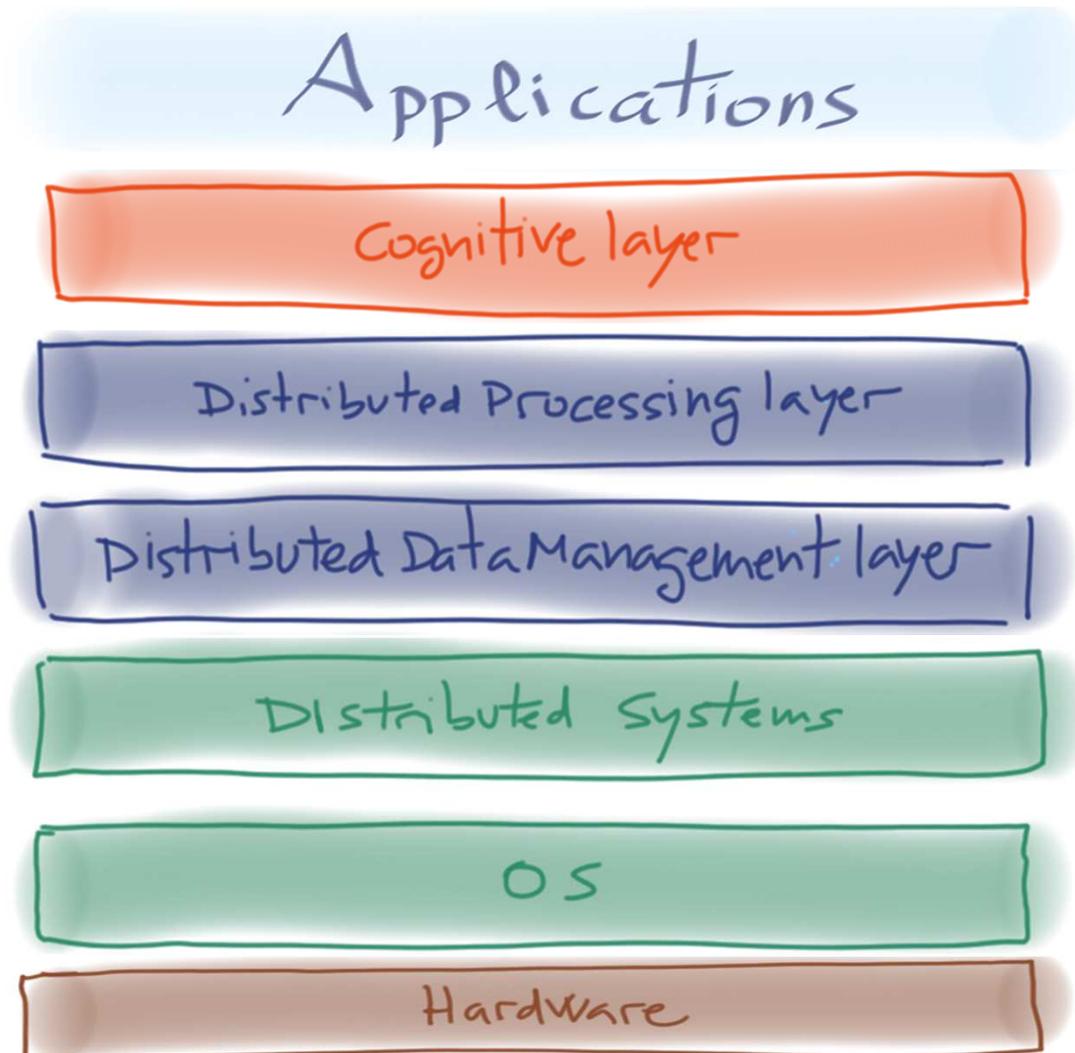3. Cognitive Layer: ML & AI Tech

*For us "Cognitive Computing" refers to the continuous development of supercomputing systems enabling the convergence of advanced analytic algorithms and big data technologies driving new insights based on the massive amounts of available data*



Advanced Analytic Algorithms

DATA

Big Data Technologies

Supercomputers Research

**((** Systems will have a new cognitive abstraction layer in the software stack

→ Cognitive Layer

*offering learning tools, but at the same time, abstracting lower layers to simplify the big data software stack.*



Applications

Cognitive layer

Distributed Processing layer

Distributed Data Management layer

Distributed Systems

OS

Hardware

# Cognitive Layer includes

**《 Machine Learning algorithms**

– (Neural Networks, SVM, Bayesian methods, …)

**《 Statistics**

– (regressions, general linear models, decision trees, …)

**《 Technologies enabled by Artificial Intelligence as**

Computer
Vision

Speech
Recognition

Natural Language
Processing

# Example: Cognitive Computing is already in business

**((** In 2011  IBM Watson computer defeated two of Jeopardy's greatest champions



*Since then, Watson supercomputer has become 24 times faster and smarter, 90% smaller, with a 2,400% improvement in performance*

*Watson Group has collaborated with partners to build 6,000 apps*

# Example: Cognitive Computing is already in business

**《** The project <u>Viv</u> built by Siri's creators



Booking a flights …

"I want a flight to MWC Barcelona with a return five days later via London."

Just closed on $12.5 M in venture capital funding.

# HOW BIG DATA IS ALREADY INFLUENCING OUR EVERYDAY LIVES

The solution analyzed 70,000 scientific articles on p53

## #Research



Baylor College of Medicine (Houston, Texas)

# Privacy & Big Data

❰❰ Even when real names and other personal information are stripped from big data sets, it is often possible to use just a few pieces of the information to identify a specific person.

❰❰ Example:
- Data: credit card transactions made by 1.1 million people in 10,000 stores over a three-month period.
- Results: knowing just four random pieces of information was enough to reidentify 90 percent of the shoppers as unique individuals and to uncover their records.
- And that uniqueness of behavior combined with publicly available information, like Instagram or Twitter posts, could make it possible to reidentify people's records by name.
- Source: http://www.sciencemag.org/content/347/6221/536.abstract

Barcelona
Supercomputing
Center
Centro Nacional de Supercomputación

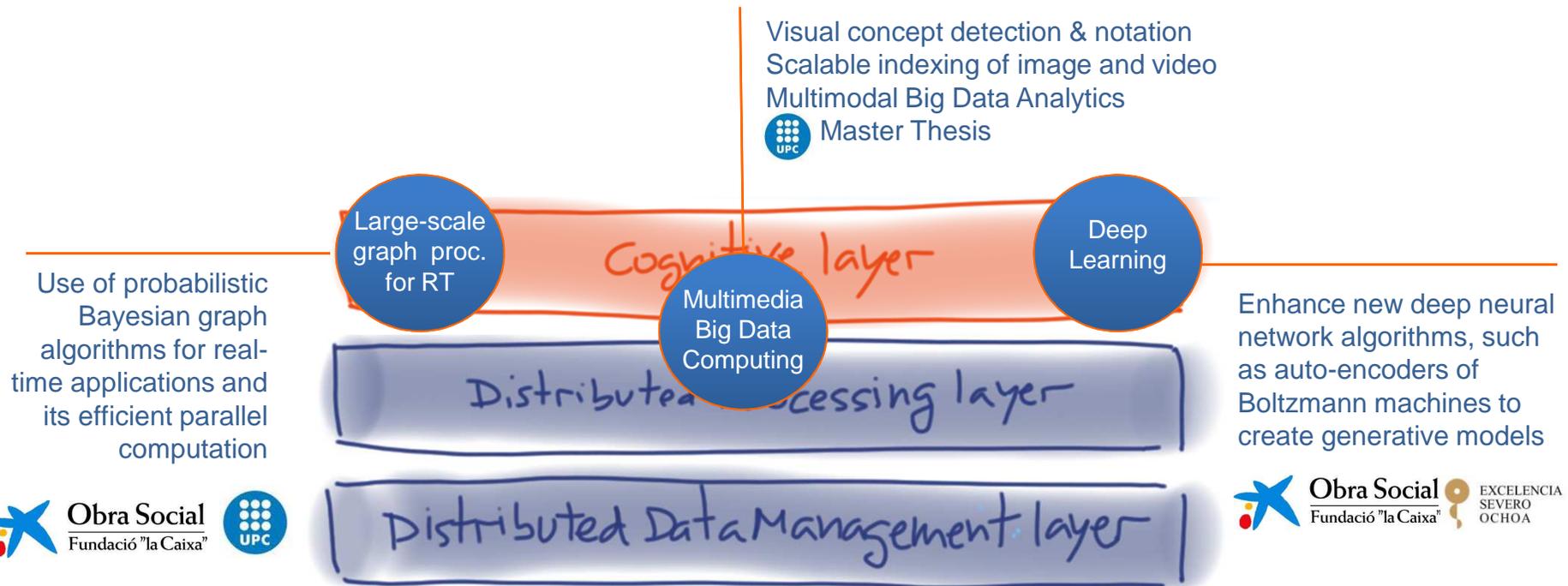EXCELENCIA
SEVERO
OCHOA

# Privacy & Big Data

**((** The old model of anonymity doesn't seem to be the right model when we are talking about large-scale metadata.

**((** Example:



– *Data: 15 months of data from 1.5 million people*

– *4 points (approximate places and times) are enough to identify 95% of individuals in a mobility database.*

– *human behavior puts fundamental natural constraints to the privacy of individuals and these constraints hold even when the resolution of the dataset is low; even coarse datasets provide little anonymity.*

– *Source:* M.I.T. Media Lab
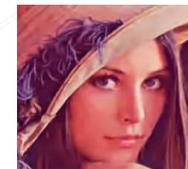 www.demontjoye.com/projects.html

Visual concept detection & notation
Scalable indexing of image and video
Multimodal Big Data Analytics
UPC Master Thesis

Large-scale graph proc. for RT

Deep Learning

Multimedia Big Data Computing

Use of probabilistic Bayesian graph algorithms for real-time applications and its efficient parallel computation

Enhance new deep neural network algorithms, such as auto-encoders of Boltzmann machines to create generative models

Obra Social
Fundació "la Caixa"   UPC

Obra Social
Fundació "la Caixa"   EXCELENCIA SEVERO OCHOA

**((** The challenge is to work with three kinds of data, at the same time:

– METADATA: Mainly geolocation, time and user defined tags. Also short descriptions, titles, surrounding text (twitter).

– SOCIAL NETWORK: Graphs of followers, likes and comments

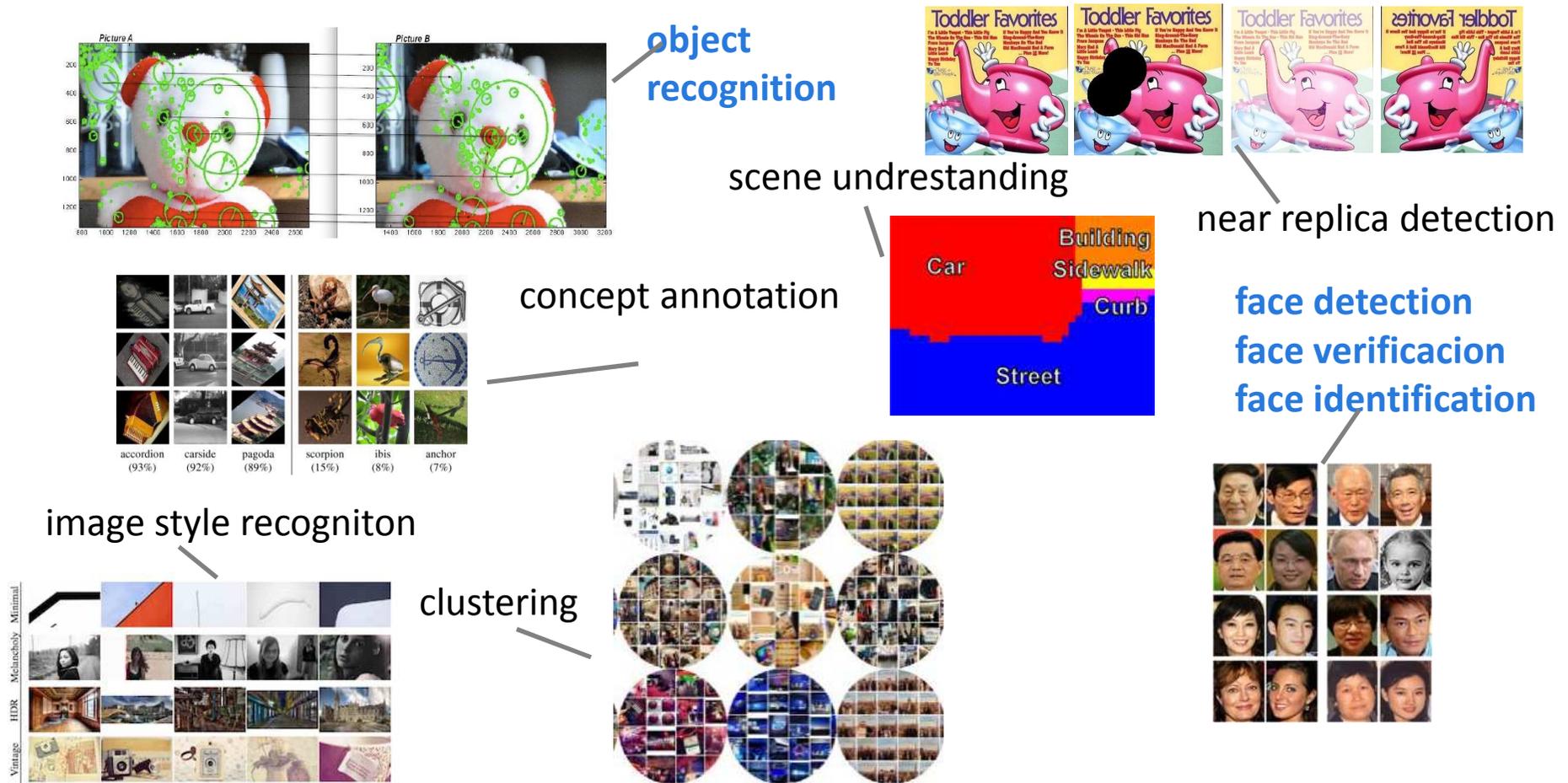– AUDIOVISUAL: We are focusing on still images.



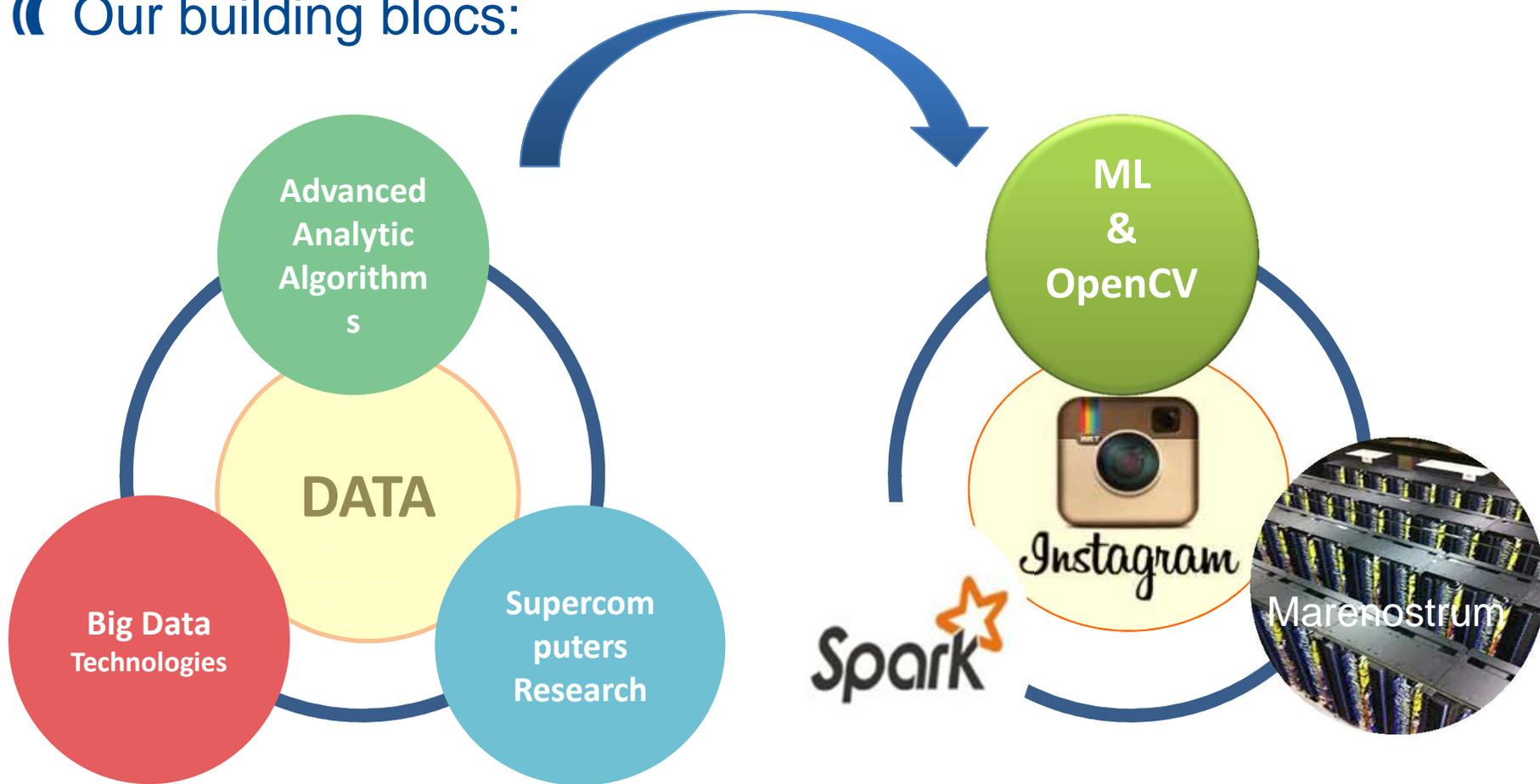| metadata | social network relationships | audiovisual content |
|---|---|---|

# Example of research at BSC: Multimodal Data Computing

**((** We are analyzing the "CONTENT" of the images



**object recognition**

scene undrestanding

near replica detection

concept annotation

**face detection**
**face verificacion**
**face identification**
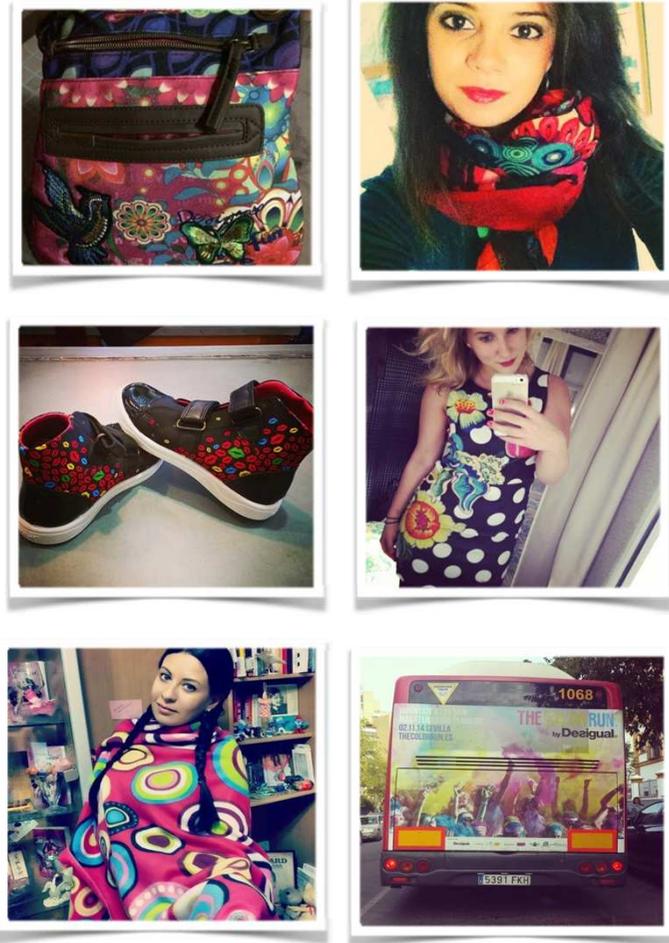
image style recogniton

clustering

**(( Our building blocs:**

# Case Study: Desigual

**((** Multimodal Data Analytics systems can aid Desigual in better understanding their customers and potential customers through the analysis of social media data sources

**#desigual #lavidaeschula**
**#mydesigual**

followers



30.000 photos

100 photos x 2K followers = 200K
Photos (100 GB)

**E.g. Predicting Desigual Followers**

AGE
GENDER
HOME LOCATION
TRAVEL PATTERNS
LIFESTYLE/CONSUMPTION PATTERNS
...

**Barcelona Supercomputing Center**
Centro Nacional de Supercomputación

# BIG DATA IMPLICATIONS

# Implications of Big Data Analytics

1. **We may be controlled by algorithms** that are likely to predict what we are about to do.

   - Privacy was the central challenge in the second wave era. In the next wave of Cognitive Computing, the challenge will be safeguarding free will.

   - Example: Last week at MWC, Qualcomm announced their latest chip, the Snapdragon 820. The new platform, called Zeroth, is said to anticipate the users actions in advance (deep learning devices)

# Implications of Big Data Analytics

2. Big Data Analytics is going to **challenge white collar**, professional knowledge work in the 21st century in the same way that factory automation and the assembly line challenged blue collar labor in the 20th century.

*Researchers at Oxford published a study estimating that 47 percent of total US employment is "at risk" due to the automation of cognitive tasks.*

### THE FUTURE OF EMPLOYMENT: HOW SUSCEPTIBLE ARE JOBS TO COMPUTERISATION?*

Carl Benedikt Frey[†] and Michael A. Osborne[‡]

September 17, 2013

**Abstract**

We examine how susceptible jobs are to computerisation. To assess this, we begin by implementing a novel methodology to estimate the probability of computerisation for 702 detailed occupations, using a Gaussian process classifier. Based on these estimates, we examine expected impacts of future computerisation on US labour market outcomes, with the primary objective of analysing the number of jobs at risk and the relationship between an occupation's probability of computerisation, wages and educational attainment. According to our estimates, about 47 percent of total US employment is at risk. We further provide evidence that wages and educational attainment exhibit a strong negative relationship with an occupation's probability of computerisation.

**Keywords:** Occupational Choice, Technological Change, Wage Inequality, Employment, Skill Demand

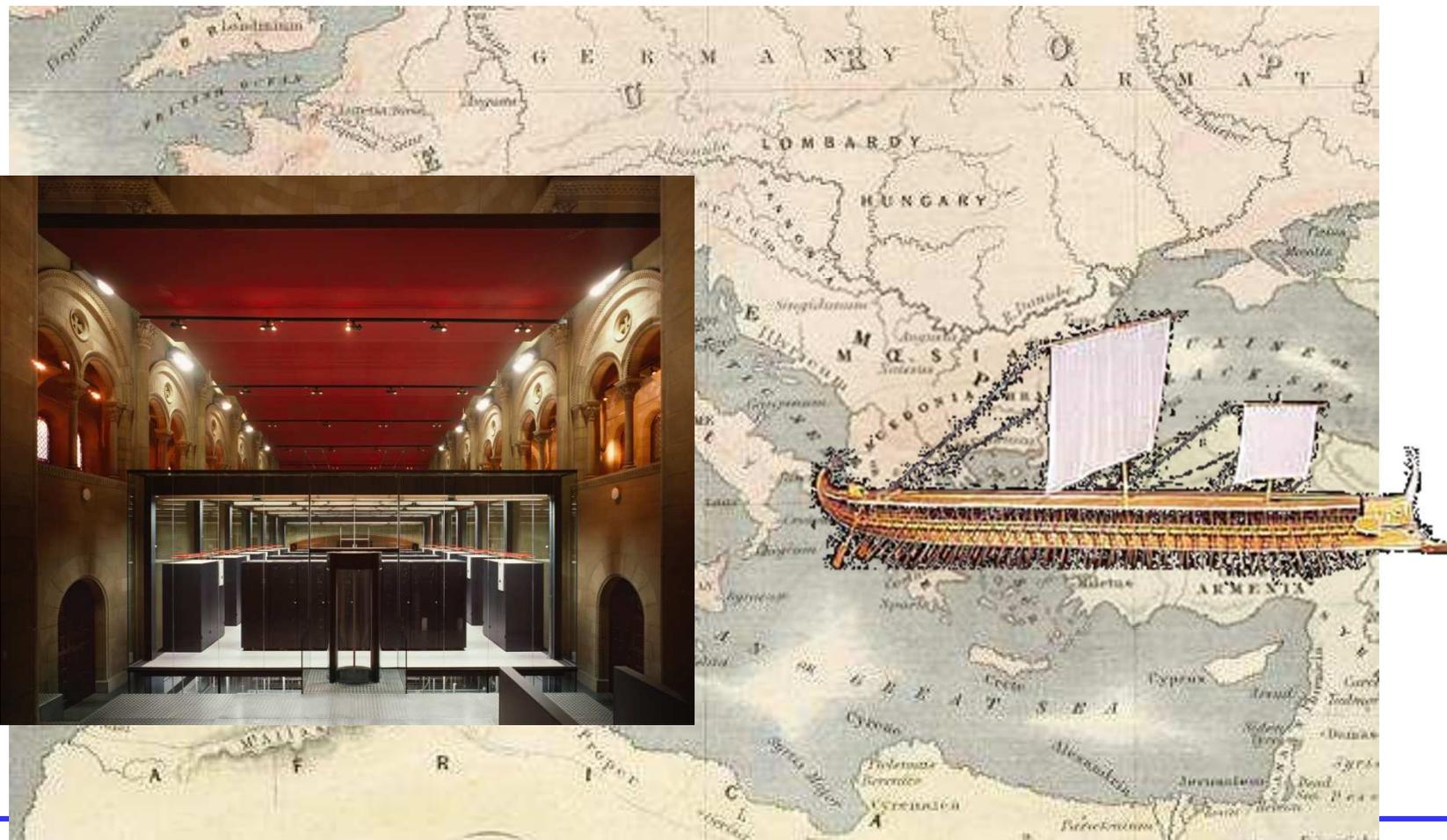**JEL Classification:** E24, J24, J31, J62, O33.

1

Barcelona Supercomputing Center
Centro Nacional de Supercomputación

EXCELENCIA SEVERO OCHOA

RoMoL Project

www.bsc.es

**Barcelona**
**Supercomputing**
**Center**
*Centro Nacional de Supercomputación*

**Thank you!**