# Detectors Data Flow and Processing: Present and Future in High Energy Physics Computing

## A. Zoccoli - INFN and University of Bologna

1

INFN

ALMA MATER STUDIORUM
UNIVERSITÀ DI BOLOGNA

# A Big (data) problem !!!

# Big Data definition

"Big data is data sets that are so voluminous and complex that traditional data processing application software are inadequate to deal with them. Big data challenges include capturing data, data storage, data analysis, search, sharing, transfer, visualization, querying, updating and information privacy. There are five dimensions to big data known as Volume, Variety, Velocity and the recently added Veracity and Value."
Wikipedia

A. Zoccoli Scuola F. Bonaudi - Cogne

"Big data is like teenage sex; everyone talks about it, nobody really knows how to do it, everyone thinks everyone else is doing it, so everyone claims they are doing it".

Dan Ariely, Duke University

# Outline

- Some history
- Computing for HEP experiments
- Some history: the GRID
- Future challenges and evolution
- High Performance Computing
- Deep Learning
- The European scenario
- The Italian Infrastructure
- Conclusion

Credits: S. Bassini , L. Benini, I. Bird, A. Cavalli, F. Fava., C. Grandi, D. Lucchesi, D. Salomoni, and many others. Thanks !!

# Introduction

# Scientific computing

Three different types of scientific computing :

1)   High Throughput Computing (HTC) for experimental physis (LHC @ CERN)

2)   High Performance Computing (HPC) for theoretical physics, materials studies, weather forecast, fluido-dynamics,deep learning ….

3)   "Traditional" Computing for small experiments or scientific initiatives

High Energy Physics developed a worldwide HTC computing infrastructure based on the GRID technology to analyse the data produced at the 4 LHC experiment at CERN.

A. Zoccoli Scuola F. Bonaudi - Cogne

# HEP computing

*HEP computing has different aspects*

> *For instance the characteristics of an accelerator-based experiment are different from those of an astro-particle experiment*
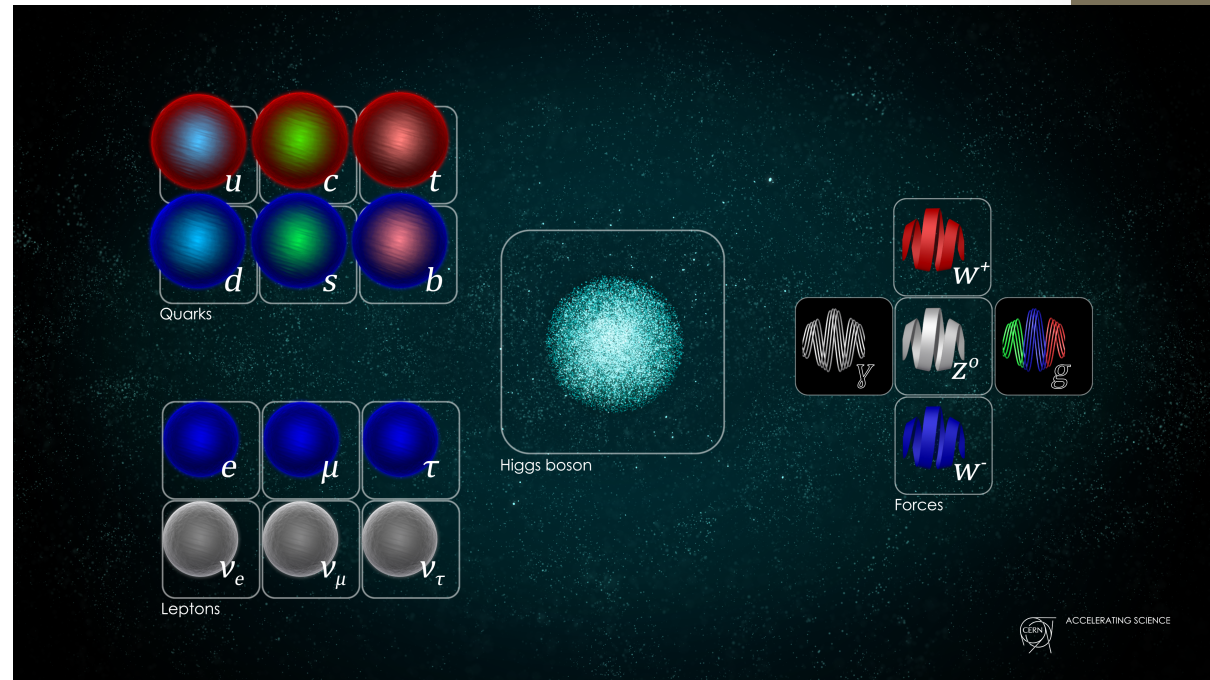
*The infrastructure built by the community is tailored on the needs of LHC that is the most demanding user at the moment (but it serves all the HEP community and more)*

# *What is HEP about?*

*High Energy Physics studies the fundamental constituents of matter and the forces that drive their interactions*
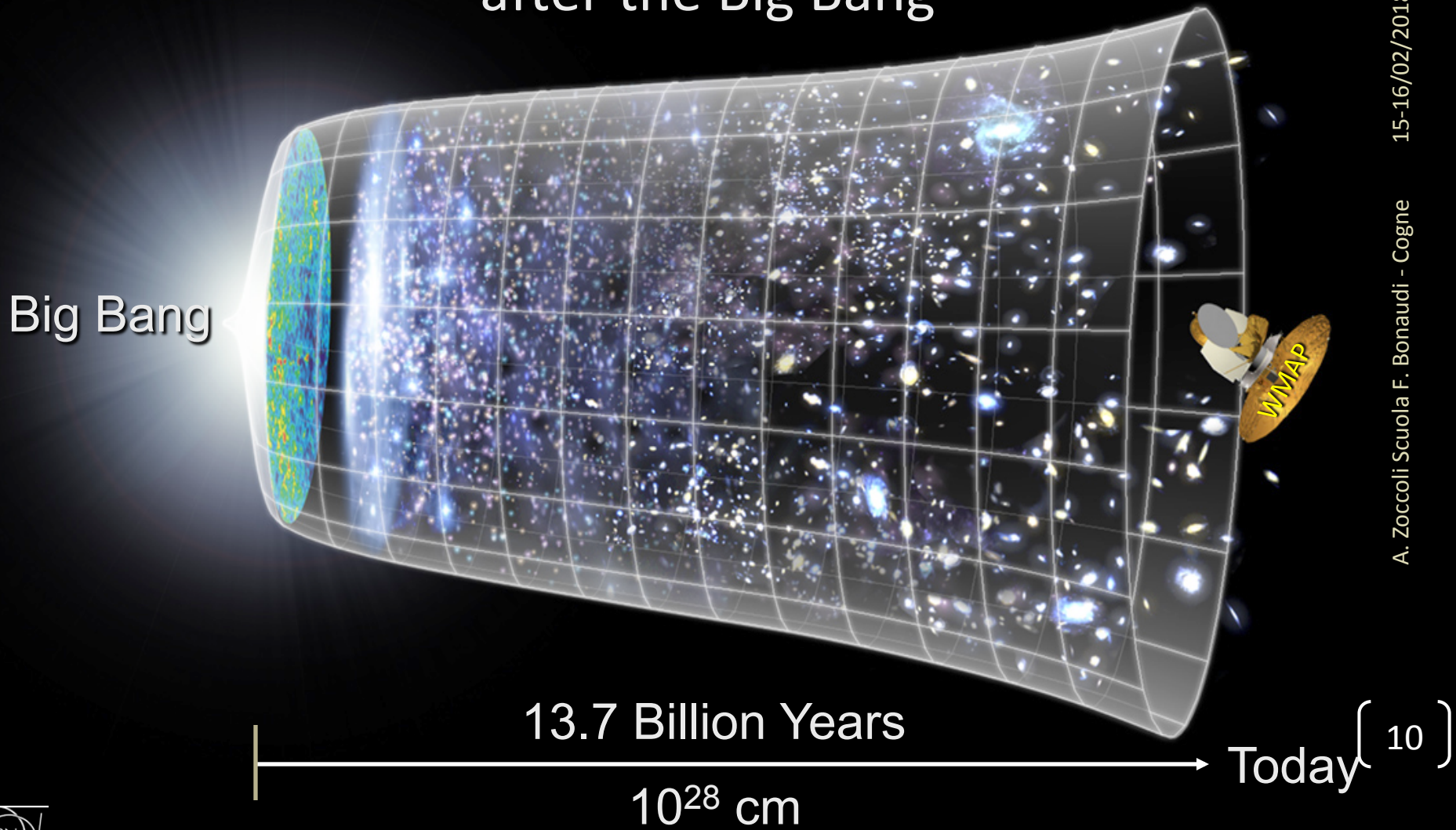
*One of the methods is to create very high energy densities*

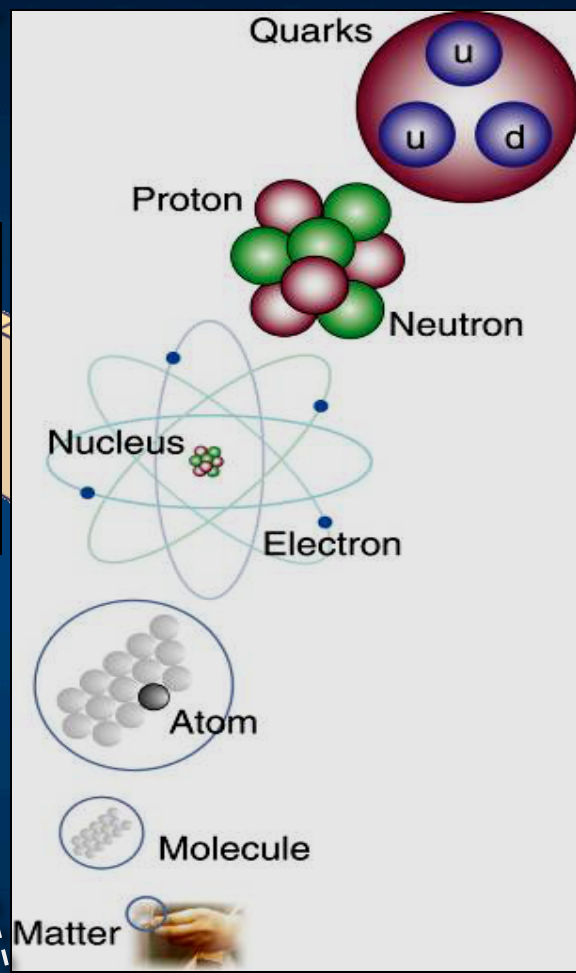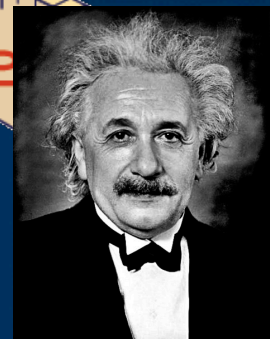*This reproduces the environmental conditions of the primordial universe*

# The Scientific Challenge:
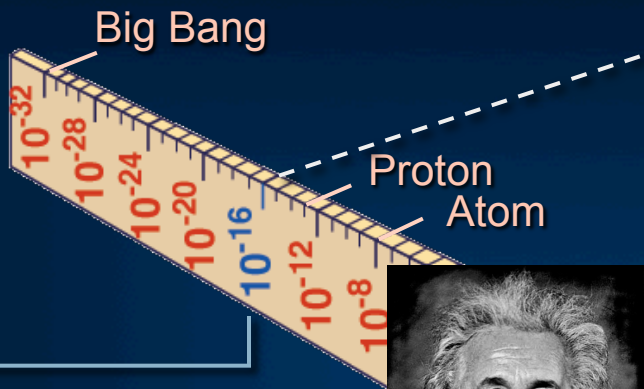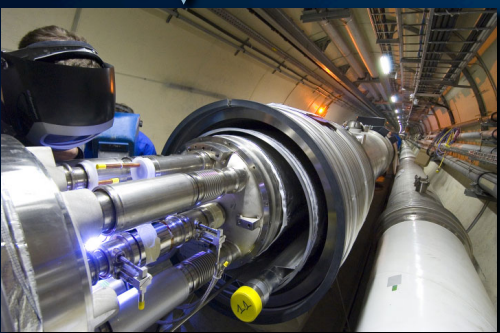## to understand the very first moments of our Universe after the Big Bang



Big Bang

WMAP

13.7 Billion Years

$10^{28}$ cm

Today

10

A. Zoccoli Scuola F. Bonaudi - Cogne

CERN

Big Bang

$10^{-32}$ $10^{-28}$ $10^{-24}$ $10^{-20}$ $10^{-16}$ $10^{-12}$ $10^{-8}$

Proton
Atom

CERN

LHC

Super-Microscope

Quarks

u

u    d

Proton

Neutron

Nucleus

Electron

Atom

Molecule

Matter

Radius of Galaxies

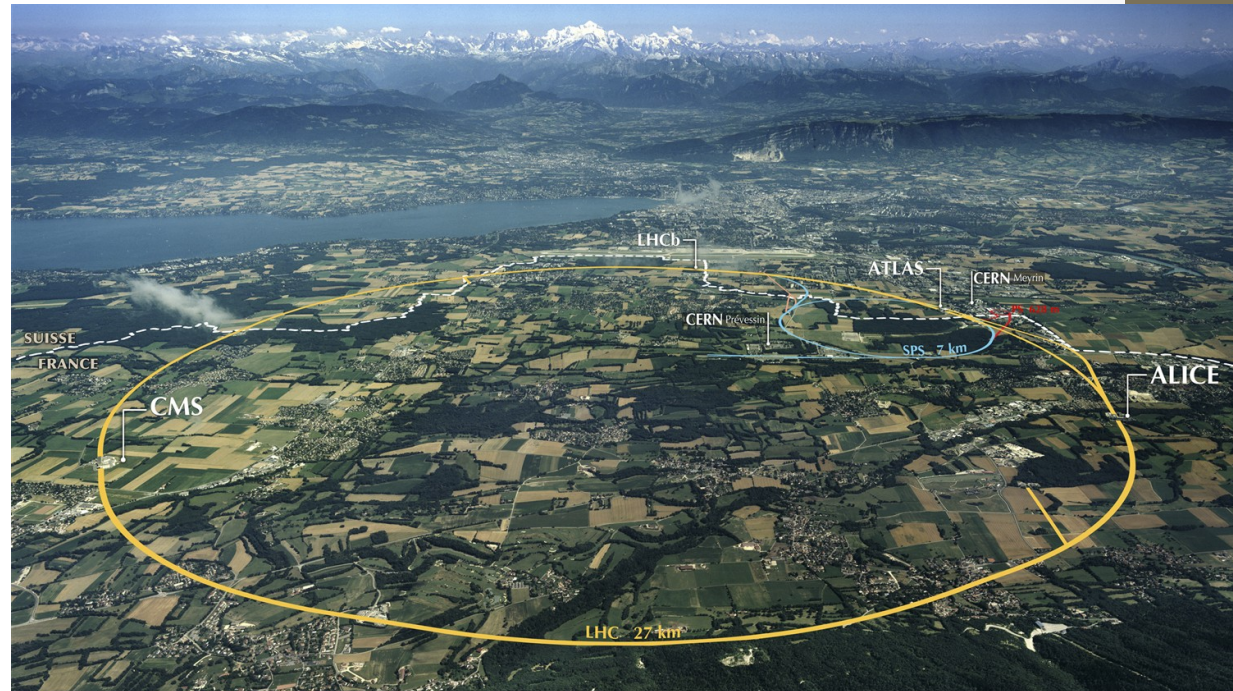Universe

$10^{24}$

$10^{28}$

cm

ALMA

Hubble

AMS

VLT

Study physics laws of first moments after Big Bang
increasing Symbiosis between Particle Physics,
Astrophysics and Cosmology

CERN

# *Particle accelerators*

*In order to create high energy densities we accelerate particles in opposite directions and make them collide one against the other*

*The CERN LHC accelerates protons. It has 27 km of circumference and is located in a tunnel about 100 m underground in the Geneva area*
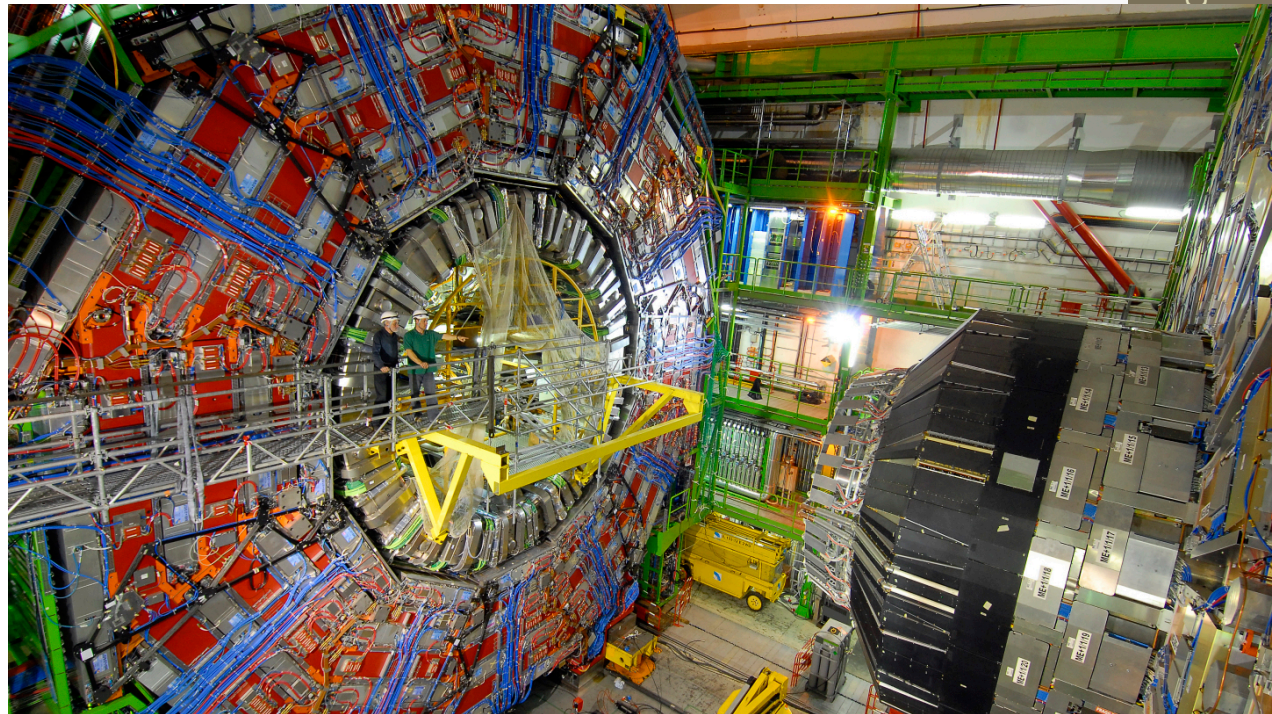
# *Particle detectors*

*Around collision points we have built particle detectors that can "see" the particle produced in the proton collision so that we can understand what happened.*

*Detectors have about 100 million channels that are acquired at each collision*
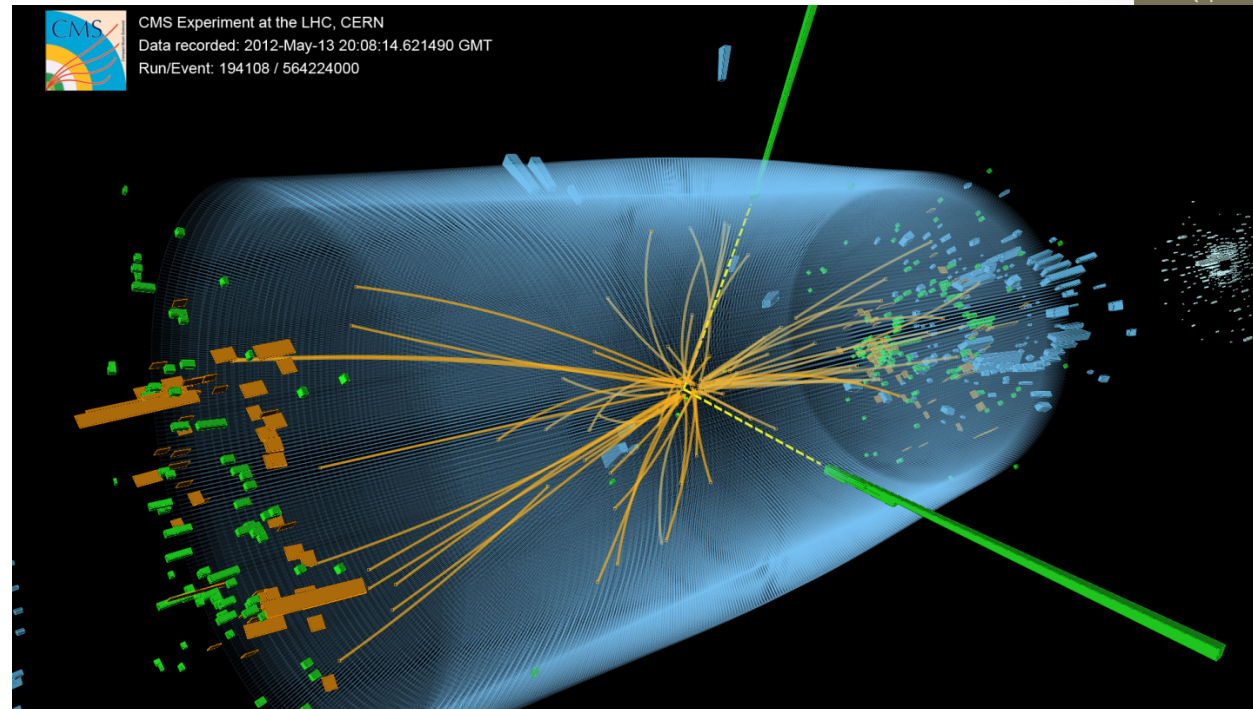
Cogne

# *Collision events*

*We call "event" a single crossing of the proton bunches in the detector area.*

*For each event we reconstruct the particles produced in the collisions.*

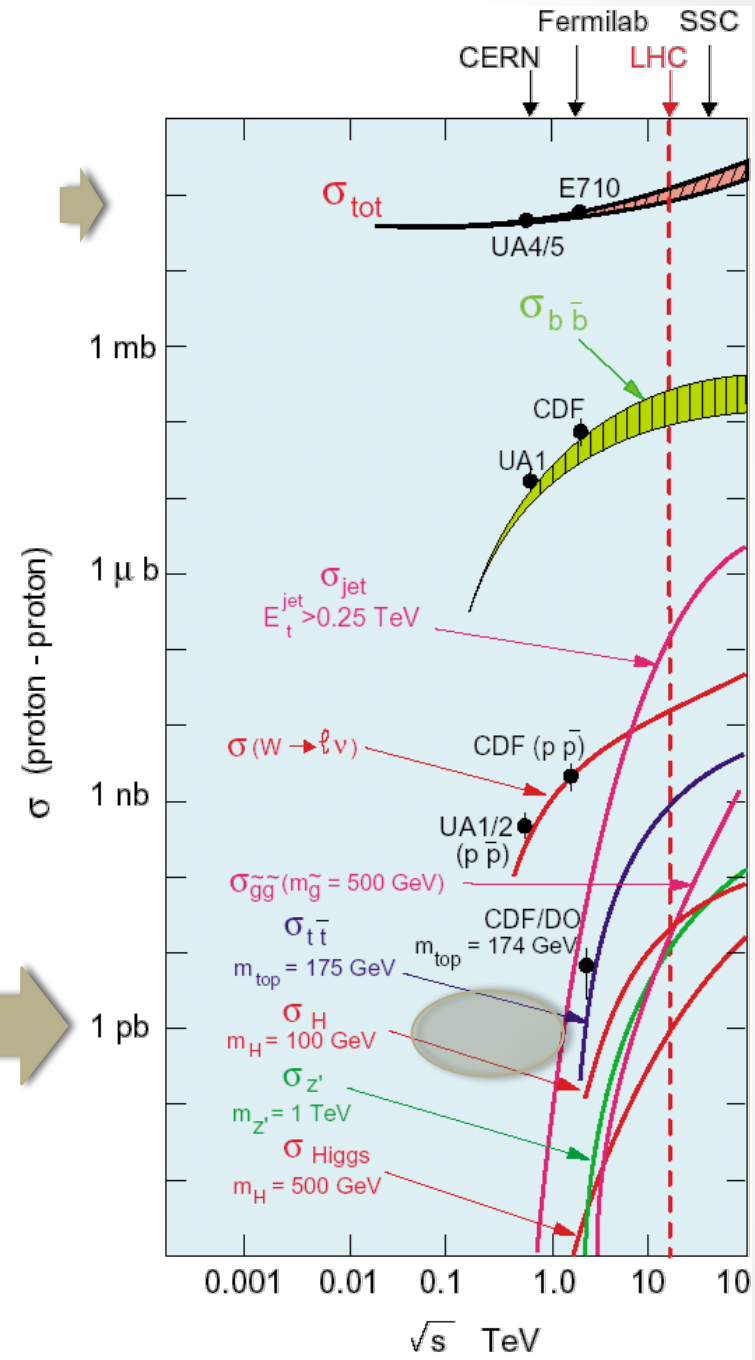*There are 40 millions crossings per second*



CMS Experiment at the LHC, CERN
Data recorded: 2012-May-13 20:08:14.621490 GMT
Run/Event: 194108 / 564224000

14

# LHC Physics

*The reason why in LHC we produce so many events is that experiments study rare events*

*For example the signal to noise ratio for Higgs events is ~ $10^{-13}$*

→ *Effective data reduction techniques are needed!*
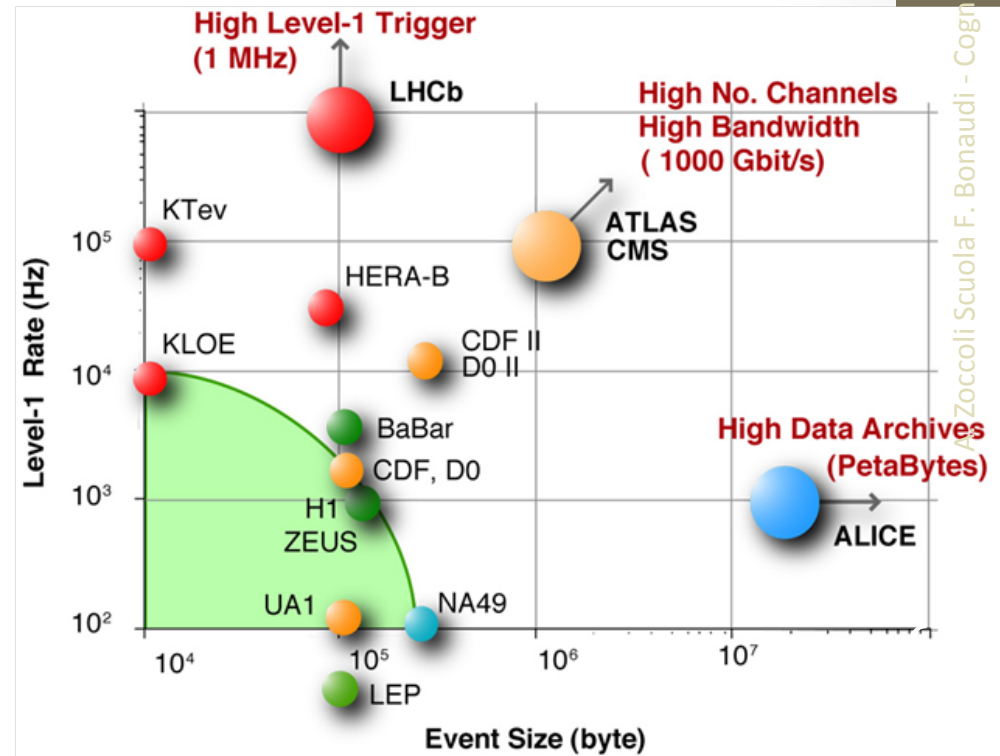
A. Zoccoli Scuola F. Bonaudi - Cogne

15

# LHC data

*In each LHC experiment there are 40 million bunch crossings per second. Every time 100 million channels are acquired (100 MB)*

➔ *40,000 EB/y (4x10$^{22}$ Byte)*

*Obviously it is not affordable!*

*The data reduction process brings to 1000 events per second each ~ 1 MB*
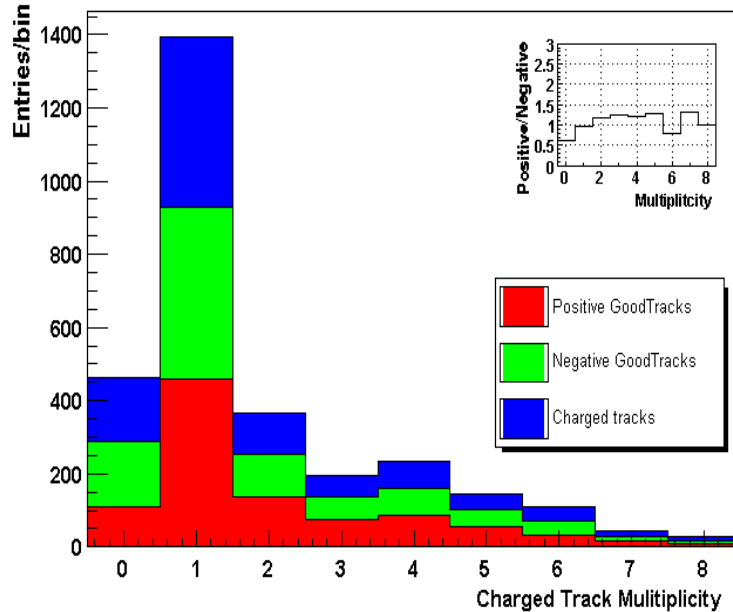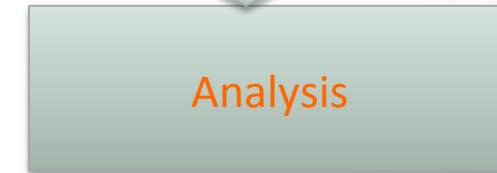
➔ *~10 PB/y (10$^{16}$ Byte)*

TRIGGER !!!!

A. Zoccoli Scuola F. Bonaudi - Cogne
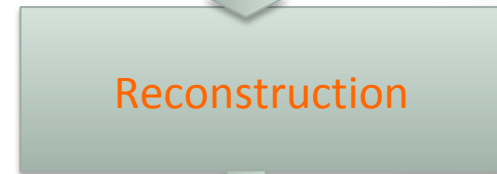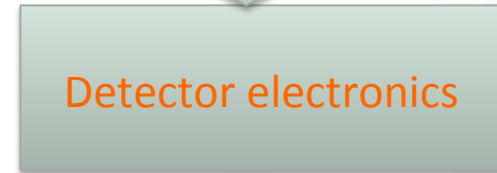
# *LHC Data processing*

- *In general physicists do not like to work on RAW data coming from the detector*

- *Typically they prefer to work with particles, jets, vertices, missing energy, etc…*

- *The process that interprets RAW data in terms of physics objects is the reconstruction*

- *Actually there are many reconstruction phases*

- *Physicists do analysis on reconstructed data*

# *LHC Real data*

LHC collisions → Decay of unstable particles →

ATLAS



Detector electronics

Trigger

Reconstruction

Analysis ←

# *LHC Simulation*

*Not just real data form detectors!*

*Since it is not possible to use analytical solutions of physic processes going from the proton interactions to the final state particles, we use simulations based on* Monte Carlo *techniques*

*Events are* generated *according to theoretical models and then* simulated *in order to reproduce the detector behaviour and then treated in the same way of the real data*

*The simulated data sample is 1 to 2 times the real data sample*

# *LHC Simulated data*

Theoretical model

Simulation of decays of unstable particles

Simulation of interactions particle-detector

Geant4, …

Pythia,…

Simulation of detector electronics

Trigger

Reconstruction

Analysis

# *Computing infrastructure*

*Management of different kinds of* <span style="color:orange">*data*</span> *(raw, reconstructed, simulated, analysis products) and of* <span style="color:orange">*processes*</span> *(different phases of* <span style="color:orange">*reconstruction*</span>*,* <span style="color:orange">*simulation*</span>*, end-user* <span style="color:orange">*analysis*</span>*)  is done on an infrastructure built by all countries participating to the LHC experiments*

*The project that coordinates the operations on the infrastructure is the*

<span style="color:blue">*World-wide LHC Computing Grid (WLCG)*</span>

# *Units used*

## *Storage*

*1 byte (B)= [0...255] = 8 bit*

*1 GB = $10^9$ B*

*1 PB = $10^{15}$ B*

*1 EB = $10^{18}$ B*

## Today: Hard Disk ~ 7 TB

## *Network*

*Gb/s = $2^{30}$ bit/s ~ 100 MB/s*

## Today: sites are connected at n x 10 Gb/s to n x 100 Gb/s

## *CPU*

*Using a unit specific for HEP: HepSpec06 (HS06)*

## *Today:*

*1 computing core ~> 10 HS06*

*1 CPU (~12 cores) ~> 100 HS06*

# *Some reference numbers*

- *600 million collisions every second*

- *Only 1 in a million collisions is of interest*

- *Fast electronic preselection passes 1 out of 10 000 events and stores them on computer memory*

- *100 GB/s transferred to the experiment computing farm*

- *15 000 processor cores select 1 out of 100 of the remaining events*

# Some history: How did we cope with the LHC requirements for data handling and processing ?

# The problem: LHC data handling

After filtering, LHC detectors select >200 interesting collisions per second.
Several MBs of data to be stored for each collision...

➡️ more than 25 Petabytes/year of data!

Global Collaboration

➡️ The GRID paradigm

8 Megabyte (8MB)
A digital photo

1 Gigabyte (1GB)
= 1000MB
A DVD movie

1 Terabyte (1TB)
= 1000GB
World annual book production

> 25 Petabytes (25PB)
= 25000TB
**Annual LHC data output**

# The solution: the GRID

At the end of 90's some important initiatives from the Physics community have provided the foundation for the GRID infrastructure:

1. The CERN action to prepare an EC project for the LHC computing

2. The launch of the eScience program in UK

3. The action of the INFN Management setting up the CNTC the Committee for the new IT technologies for LHC
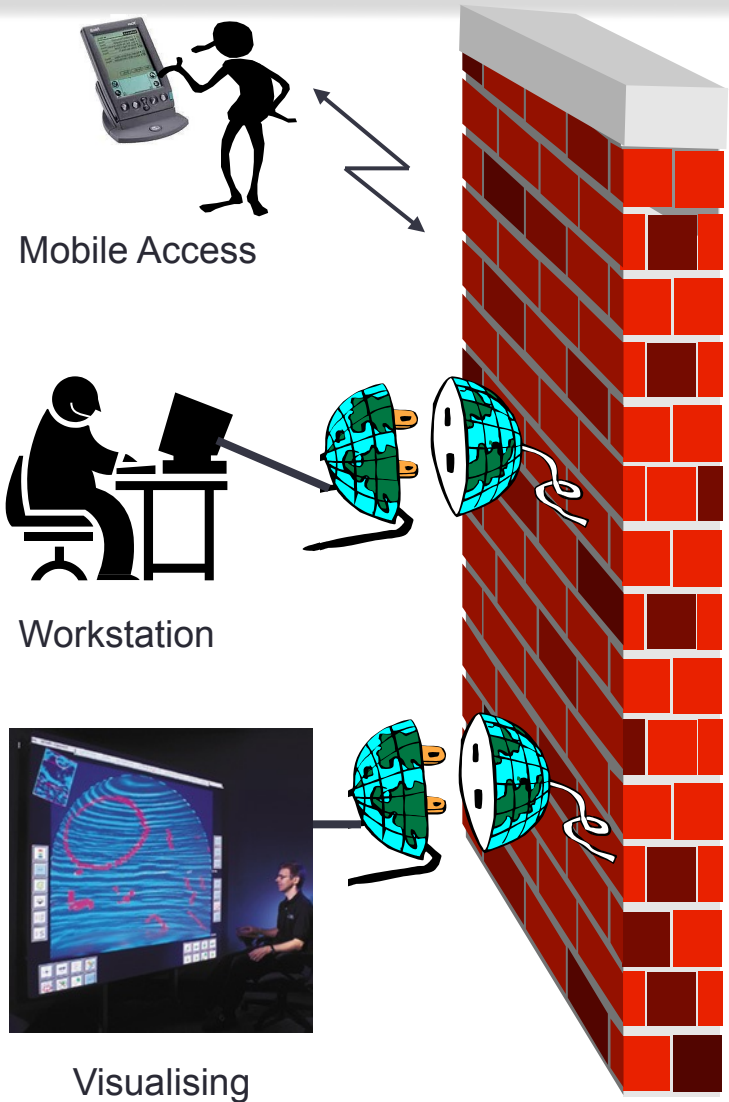
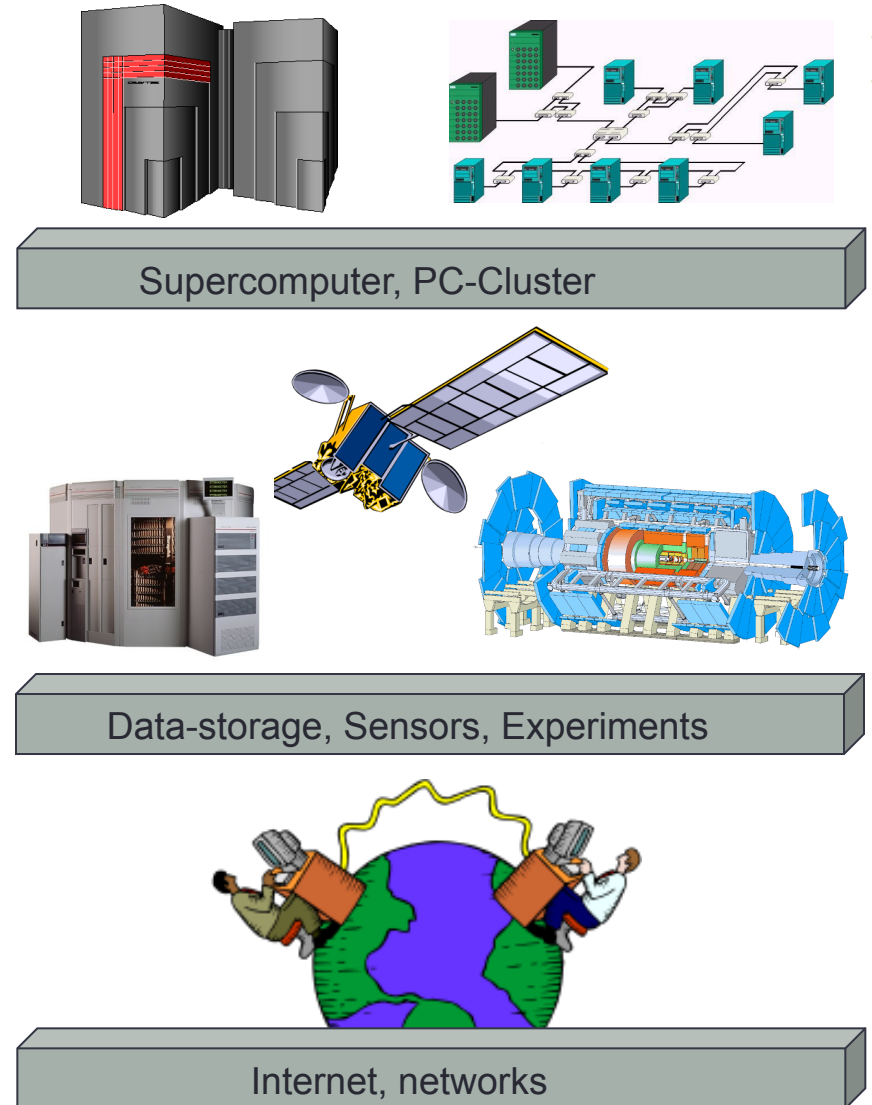4. ...........

26

# The GRID development and INFN

Active participation in all the development steps like for example:

1. DataGrid (2000) with CERN, INFN, CNRS, PPARC, NIKEHF

2. egee (2003)" Enabling GRID for E-Science in Europe": definition of the European Grid operation model

3.  WLCG (2004) designed for the LHC data analysis based on GRID

4. EGI (2010) "European Grid Infrastructure": probably we know what it is or should be..

# *The Grid paradigm*

Mobile Access

Workstation

Visualising

G R I D   M I D D L E W A R E

Supercomputer, PC-Cluster

Data-storage, Sensors, Experiments

Internet, networks

8

# *A distributed system*

*Advantages of a <span style="color:orange">distributed system</span> (w.r.t. a unique data centre)*
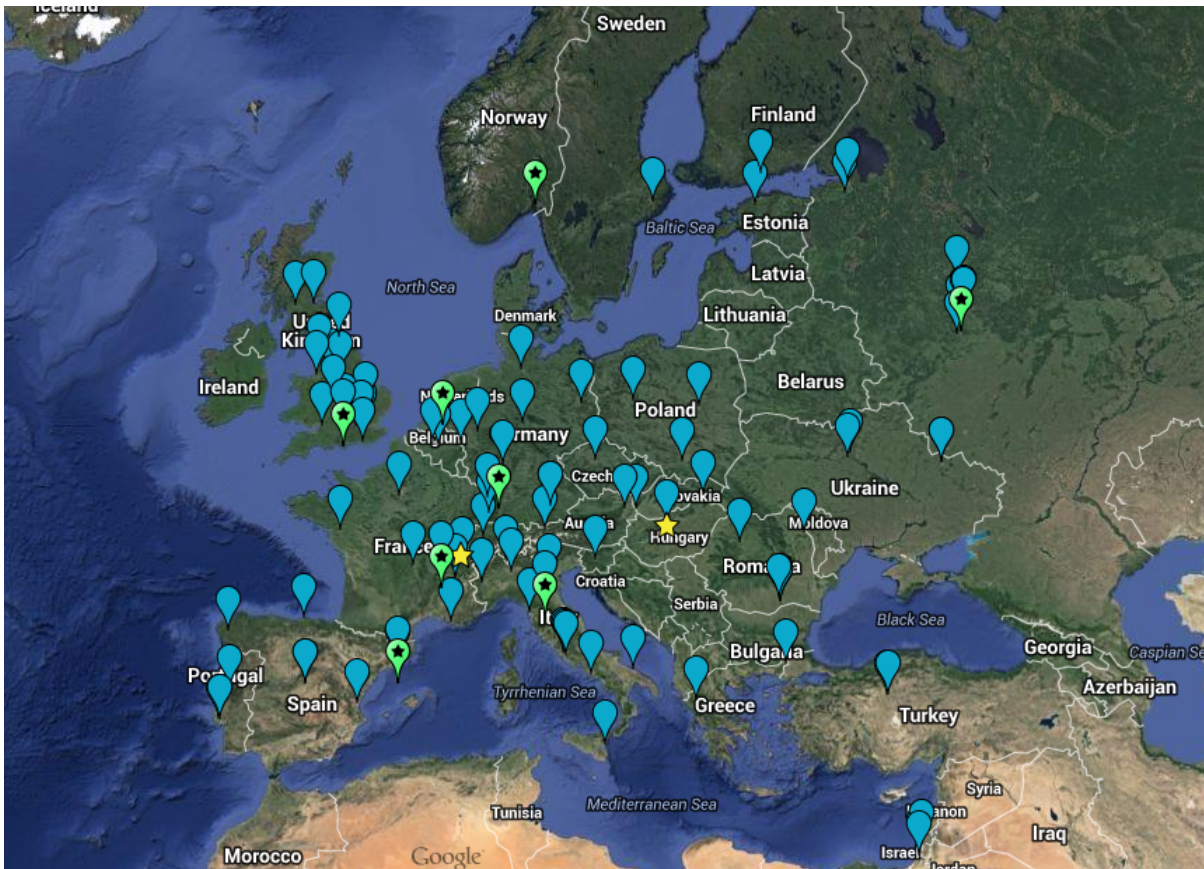
- *Avoid single point of failure*

- *Have access to local funding otherwise not provided by member states*

- *Investment on <span style="color:orange">manpower</span> available in different countries*

- *Build an adaptable system able to integrate external resources that are made available*

# The e-infrastructure

During the past years INFN + other FA + CERN with large contributions of EC has constructed and consolidate a large computing infrastructure



⭐ Tier-0    ⭐Tier-1    ⭐ Tier-2

Nearly 170 sites in 40 Counties:
~500.000 cores
~1000 PB storage
(400 disk + 600 Tape)
~ 2M jobs/day
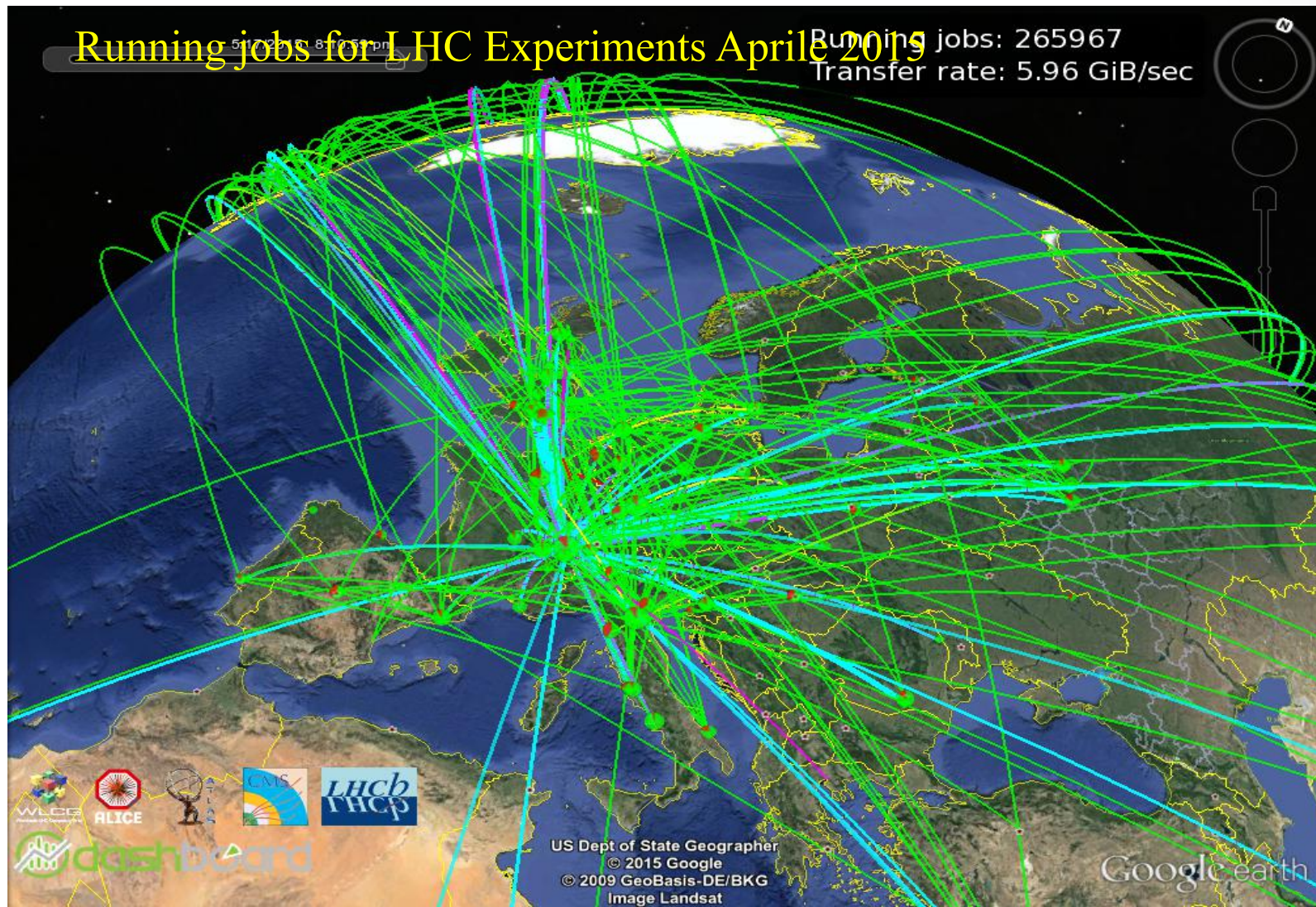in 2017

Half of the resources are sitting in Tier-1

# The GRID

Besides the computing centers the infrastructure rely on network which has modified the major experiments computer models



Running jobs for LHC Experiments Aprile 2015
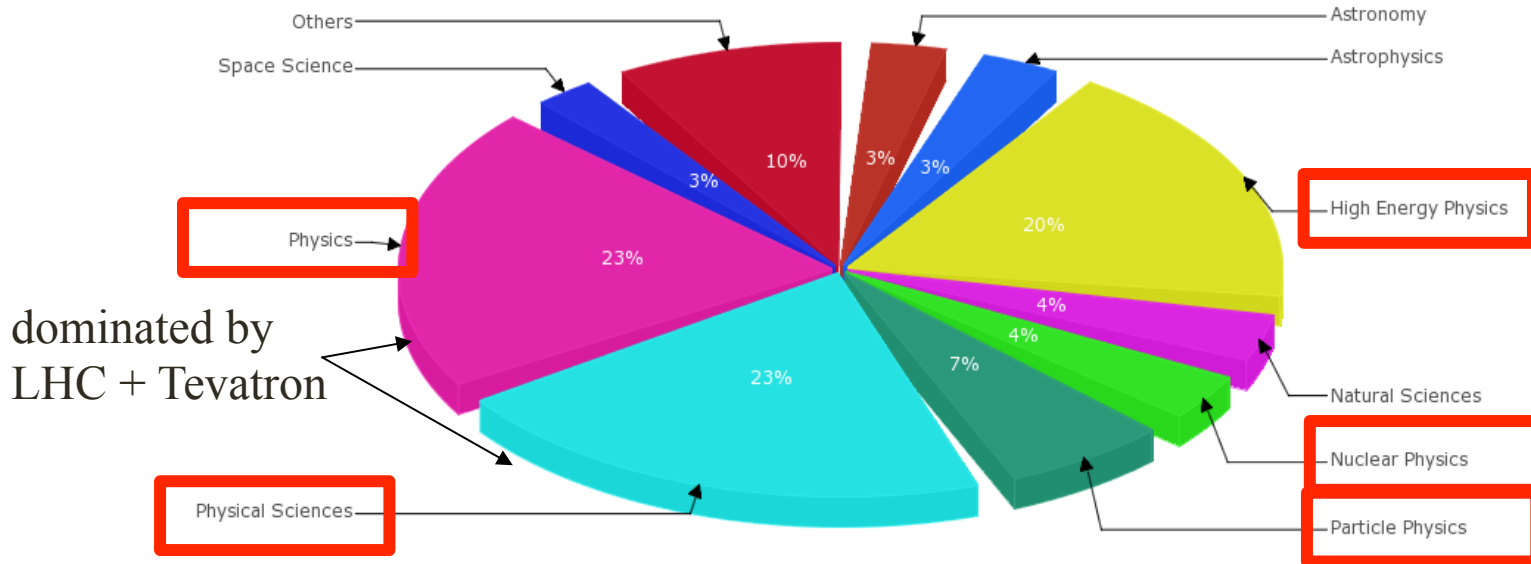
Running jobs: 265967
Transfer rate: 5.96 GiB/sec

# The Users

This infrastructure has served several disciplines but mainly high energy physics (> 75%)



VO_DISCIPLINE Normalised CPU time (HEPSPEC06) per DISCIPLINE

01-01-2011
01-01-2015

dominated by
LHC + Tevatron

Form EGI accounting portal

A. Zoccoli Scuola F. Bonaudi - Cogne     15-16/02/2018

*CERN Data Centre (Tier 0)*

> *100.000 processor cores*

*Data aggregation and initial data reconstruction*

*copy to long-term tape storage and distribute to other data centres*

*11 Tier 1 centres*

*Permanent storage, re-processing, analysis*

*140 Tier 2 centres*

*Simulation, ent-useer analysis*

*> 2 million jobs running every day*

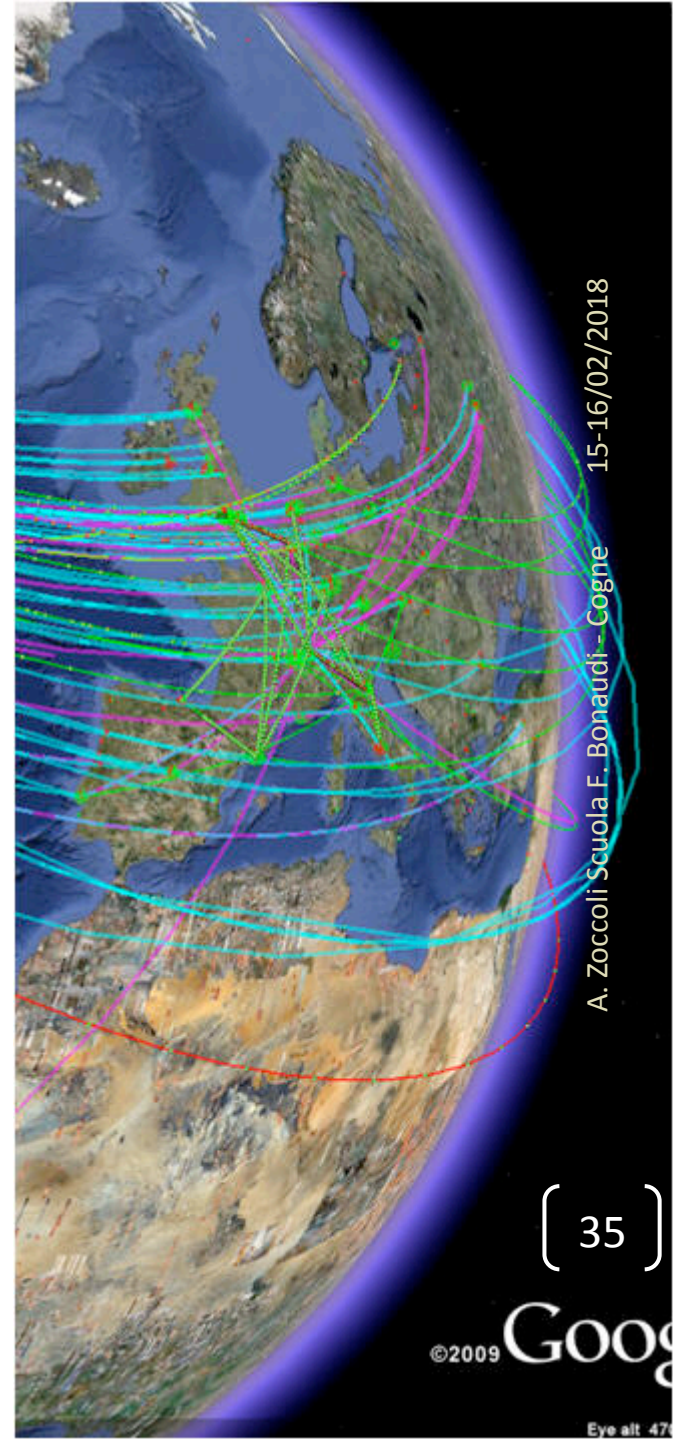*25GB/s global transfer rate*

34

# *...more numbers*

*Global resources for 2017 are:*

- *5,200,000 HS06 (~500.000 processor cores)*

- *395.000 TB disk*

- *590.000 TB tape*

- *Dedicated network connections (from multiples of 10 Gb/s to multiples of 100 Gb/s)*

*...and more available in collaborating institutes*

*More than 180 data centres in over 35 countries*

*More than 8000 analysts all over the world*
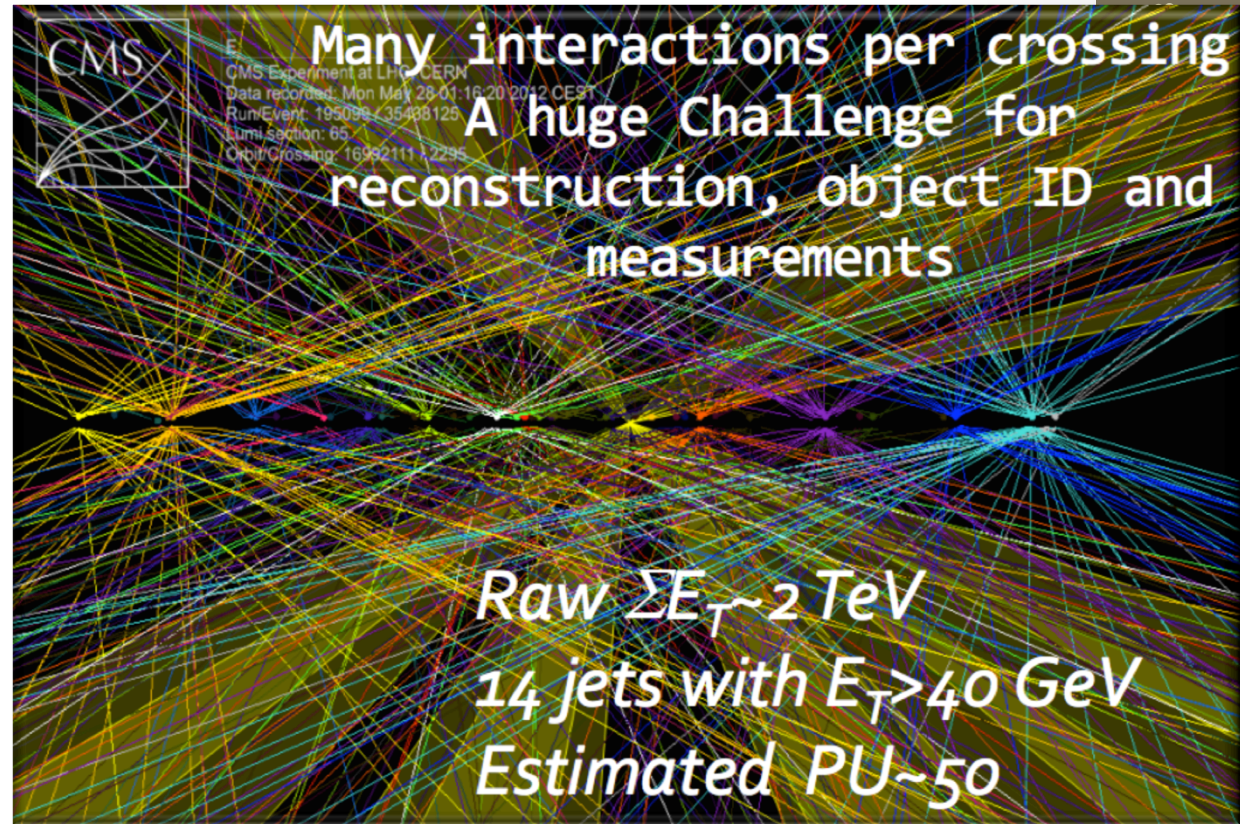
A. Zoccoli Scuola F. Bonaudi – Cogne

35

©2009 Goog

Eye alt 47

# *Pile-up*

*If you're wondering why a bunch crossing rate of 40 MHz produces 600 collisions per second:*

*Every bunch crossing (event) there are on average 15 p-p collisions (AKA pileup)*



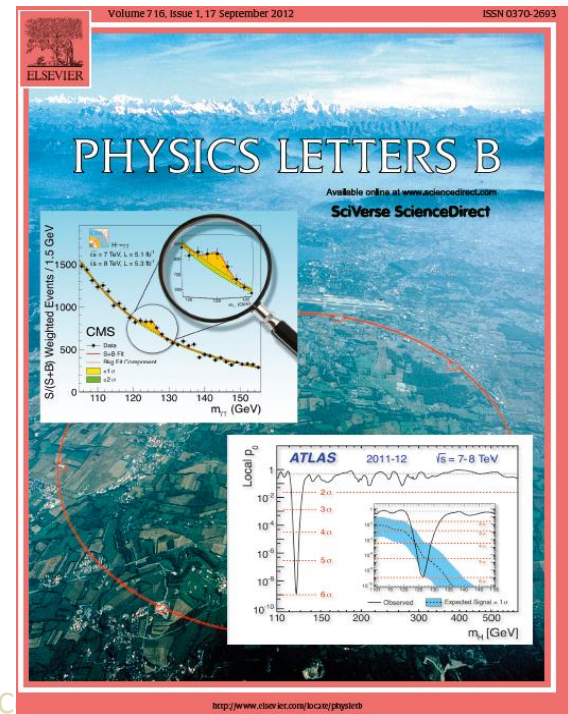*Pileup is increased in 2017 to 50 and eventually to more than 150 in HL-LHC*

# A successful GRID story

On 4 July 2012 both of the
CERN experiments ATLAS
and CMS announced
they had independently made
the same discovery



In summary

We have observed a new
boson with a mass of
**125.3 ± 0.6 GeV**
at
**4.9 σ** significance

# A successful GRID story lasting an year

On 14 March 2013 CERN confirmed that:

"CMS and ATLAS have compared a number of options for the spin-parity of this particle, and these all prefer no spin and positive parity [two fundamental criteria of a Higgs boson consistent with the Standard Model]. This, coupled with the measured interactions of the new particle with other particles, strongly indicates that it is a Higgs boson." (wikipedia)
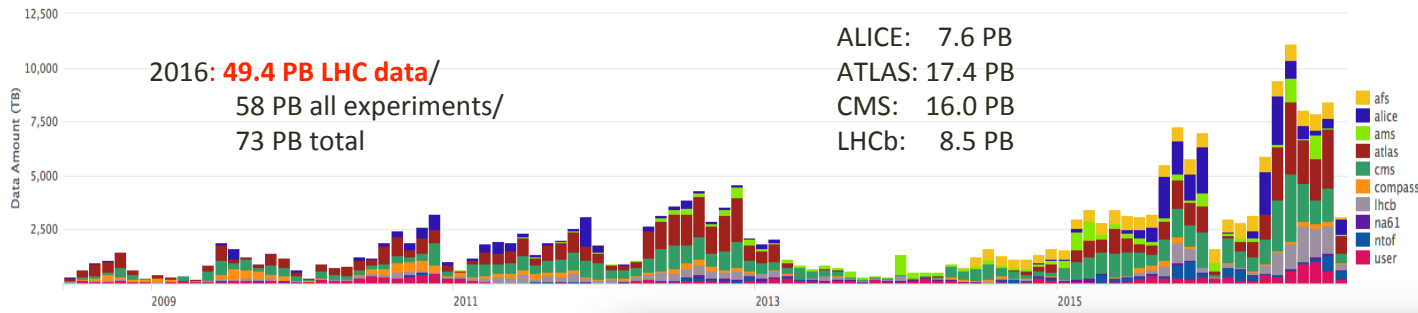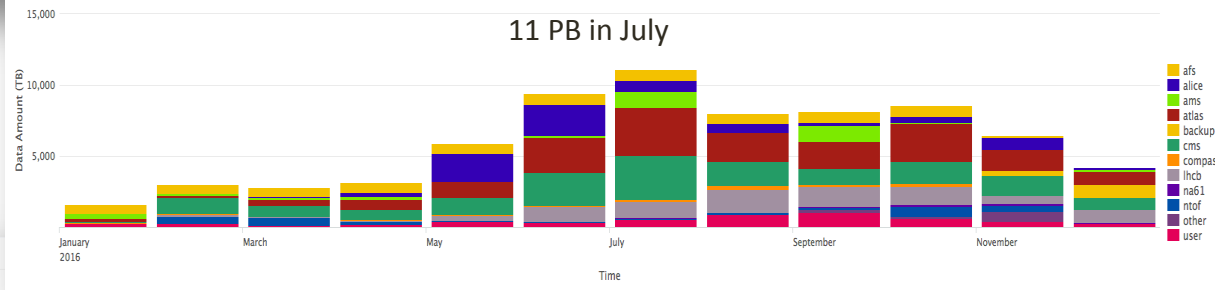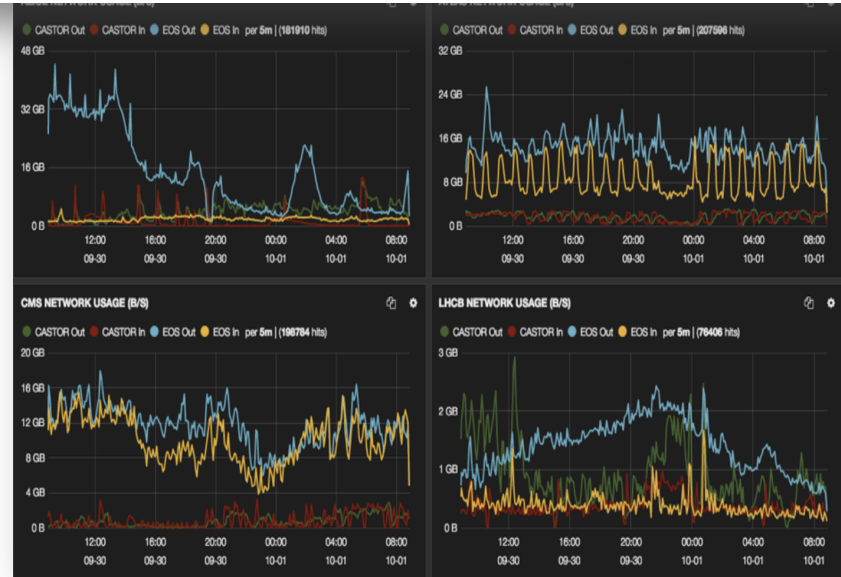
# Some more details…..

# Data in 2016 - updated
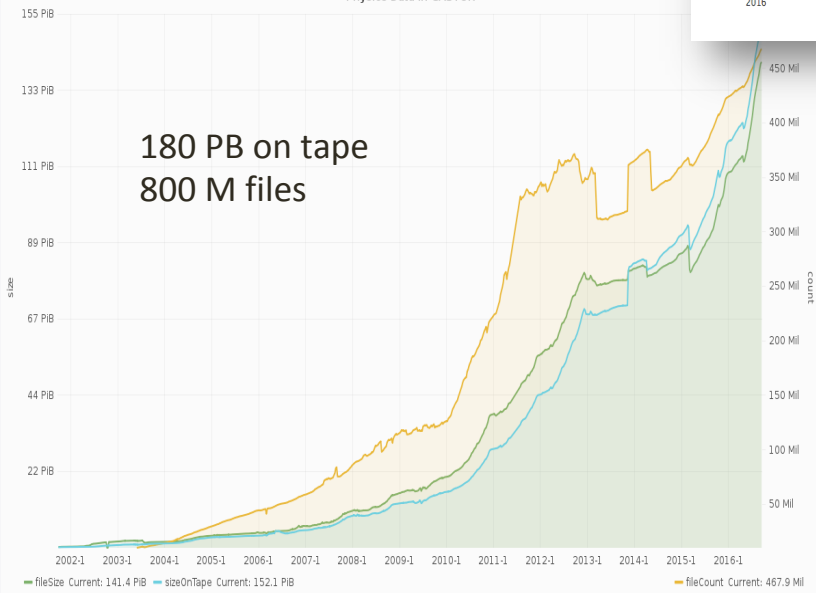


**Transfered Data Amount per Virtual Organization for WRITE Requests**

2016: **49.4 PB LHC data**/
58 PB all experiments/
73 PB total

ALICE:  7.6 PB
ATLAS: 17.4 PB
CMS:   16.0 PB
LHCb:   8.5 PB

**Transfered Data Amount per Virtual Organization for WRITE Requests**

11 PB in July

Physics Data in CASTOR

180 PB on tape
800 M files

# Data transfers

# CPU Delivered



CPU Delivered: HS06-days/month

New peak: ~180 M HS06-days/month
~ 600 k cores continuous

# *Grid Security management*

- Authentication based on x.509 certificates

- Authorization based on *attribute certificates* (VOMS)

- *Policy management* system (ARGUS)

# *Grid Computing management*

*Access is based on batch jobs: asynchronous execution*

*Dedicated interfaces allow to manage remote submissions as if local*

*Interactive processing is limited and based on local resources or on systems able to manage part of the load in batch mode (e.g. PoD)*

User workspace
PoD UI
PoD Server
Workers
A Local cluster

A Grid Site
Firewall / Sub-net
Workers

Server validates PoD WNs, which if validated will become your PROOF WNs.

# The "pilot" model

*Separation of resource allocation and job management*

# Storage Federations

Starts from the possibility to have remote data access

Clients always ask the closest location for files

If the file is not available, the request is forwarded to a *hierarchy of redirectors* until it is satisfied (or fails globally)



46

# The future challenges

# The evolution of scientific computing

The evolution of scientific computing is mainly driven by the necessity to process unprecedented data samples and by the needs from different disciplines (astroparticle, biology, science of materials, medicine, industry etc. etc.)

→ New computing models for LHC experiments

→ Fast networking

→ New concepts for e-infrastructures

→ Cloud computing

48

# LHC resources for High-Lumi



By assuming:
- Trigger rate 5-10 kHz, factor10 wrt today
- Expected Pileup ~140-200 ⇒ complexity of events increased by a factor 10

→Back on the envelope evaluation of resources needs = 100x wrt today !!

# LHC resources for High-Lumi - 2

Expected mitigation effects:
- Technological evolution for CPU and storage: +20%/y performances → 6x in 10 years
- Computing Models Optimization (CPU: exploitation of parallel architectures) and of the trigger scheme (raw data vs AOD)

## Current estimated gap:

CPU: 6x (CMS) e 4x (ATLAS)

Disk space: 4x (CMS) e 7x(ATLAS)

# *Not only LHC…*

# Future Computing Resources



Future Astroparticle experiments (SKA, CTA, Euclid, …) will produce unprecedented amount of data.

# Future Computing Resources

Huge resources demand resulting in requests of HPC and Big Data management will come from many different research fields in the next years:

- HEP and astroparticle physics
- Human brain
- Biology
- Personalized medicine & genomics
- Weather predictions
- Climate change studies
- Material studies
- Smart manufacturing & Industry 4.0
- IoT
- SmartCity

Impact on private sector

# Does the technological evolution help?

# *CPU power*

*Moore's law* (CPU performance doubles every 18 months at the same cost) does not hold any more



Moore's Law - 2005



Intel Processor features

Source: Andrzej Nowak – CERN OpenLab

*We may reasonably expect a 20% increase per year but we need to cope with multi-core systems*

55

# Disk



ASTC Technology Roadmap

*Extrapolation is more difficult for disk because there are technology changes foreseen*

*It is safe to assume that disk size in 2023 will be around 40 TB*



*The number of disks may not need to increase*

A. Zoccoli Scuola F. Bonaudi - Cogne

# *Electrical power*

*Disk power consumption does not depend on size in first approx.*



*CPU power to electrical power ratio increasing linearly. In 2023 foreseen 2 HS06/W*
*→ Low power architectures?*

*Total power (including services) in 2023 is foreseen to be ~ 1 MW*

57

# *Costs*

- *Provisioning of CPU, disk and tape*

- *Electrical power for IT*

- *Electrical power for cooling*

    *~60% of power for IT at CNAF (PUE 1.5 to 1.7 depending on the season)*

- *Infrastructure maintenance*

→ *Far from a "flat budget" hypothesis for Run3*

    *And Run4 is even worse!*

*Need to change models and exploit new technologies*

# More actions….

A. Zoccoli Scuola F. Bonaudi - Cogne

# New computing models in HEP

- Computing Models are not static. Continuous evolution
  - since the beginning of the data taking, the "ideal" CMs have been replaced by realistic ones exploiting the technology and infrastructure improvements
- In Run-1 the LHC experiments have been able to cope with an unforeseen amount of data transferred and analysed

*HEP is not different from the rest of the world*

*We can try to follow what others are doing*

*Even though Google, Facebook, & C. are making money out of investments while we have budget restrictions*

*We can also try to exploit resources that others may make available to science in opportunistic mode*

→ *From GRID to Cloud*

# Cloud Computing: definition

- The canonical definition comes from the US National Institute of Standards and Technology (NIST) (http://goo.gl/eBGBk)

- In a nutshell, Cloud Computing deals with:

"Grid Computing" vs. "Cloud Computing" Google lookup stats (from Google Trends)



**1** Supplying

**2** information and communication technologies

**3** as a service

# Why the Cloud?
# (or, what was missing?)

1. Ease of access to IT resources for small as well as big companies and [scientific] communities.

2. Software and financial sustainability.

3. Robustness (mitigation of vulnerabilities).

4. Modular and scalable architecture (or, flexibility).

5. Open source software, vendor independence.

6. Clear business model(s).

# *From the Grid...*

*The "factory" harvests job slots*

# ...to the Cloud

*The "factory" harvests machines (or containers)*



Client

*Resource allocation*

Glidein Factory
- schedd
- collector

Cloud interface (EC2)

CMS WM tools
- schedd

Glidein Frontend
- Central Manager
  collector
  negotiator

VM

(Glidein Startup)
- startd
- CMS job

*Job management*

# *Hybrid Cloud model*

*The use of standard cloud interface will allow to exploit private and commercial clouds at the same time*

## Helix Nebula Hybrid Cloud Model

### Users

Big Science   Small al Medium Scale Science   Other market sectors

Hybrid cloud platform

Network Commercial/GÉANT

Front end   Front end   Front end   Front end   Front end   Front end

Publicly funded   Commercial   EGI   PRACE

# A possible computing evolution in HEP

- Evolution toward a federated infrastructure with larger installation with integration of HPC and commercial clouds:
  - Economies of scale and improved efficiency
  - Reduction of operating costs
- Evolve the current e-infrastructure and the computing models:
  - serve HEP at large and other sciences (astrophysics, astronomy, photon science, chemistry, biology, medicine …)
  - flexible to technological development on datastorage solutions and HPC
- exploit different budget lines

# Model for future HEP computing infrastructure



ATLAS HLT

CMS HLT

LHCb HLT

ALICE O2

turbo

Reconstruction facility (calibration, alignment, reconstruction)

Archive at "AOD" level

**Distribute**

Simulation resources

**LHC Data cloud**
**Storage and compute**

DC

Compute

1-10 Tb/s

DC

Compute

DC

Cloud users: Analysis

**Datalake**

# Build a "data cloud"

- Few – O(5-10) - large centres
  - Multi-Tb private (SDN) network between them
  - Treat as a single "virtual data centre"
    - Policy replicates data inside for security and performance
    - Think of RAID across data centres
  - Store all of the "AOD" data here; Do not replicate data to global physics institutes (major cost)
- Pluggable compute capacity:
  - HEP resources at these centres & other large centres
  - Commercial compute
- Model allows commercial data centres
  - For storage – enough redundancy that a commercial centre could unplug
  - For compute
  - Relies on networking and GEANT/Esnet etc. connections to commercial entities, policy
- ☞ Users access data in this cloud remotely
  - Eventually download "ntuples" – or equivalent
  - All organised processing is done in this model
- ☞ Enables new analysis models: all data can be seen as colocated
  - Get away from the "event-loop" → queries, machine-learning, etc.

This idea has been discussed in the WLCG community (e.g. see I. Fisk CHEP plenary)

❑ Hybrid model:
  ▪ HEP-resources at a level we guarantee to fill → cost-effective
  ▪ Commercial resources for "elasticity"
❑ Needs new funding models

A. Zoccoli Scuola F. Bonaudi - Cogne

69

# Are Commercial Providers the solution ?

# Some Cloud-related risks

- **Security and privacy**
- **Lock-in**
- **Isolation failure**
- **Management interface compromise**
- **Insecure or incomplete data deletion**

- Some examples taken from a very popular public Cloud follow...

# Non-exclusive rights

- **You lose ownership**

  - Amazon (for instance) could develop products directly competing with what you yourself develop on AWS, adopt technologies that you are using...

  - ... or assist somebody else in developing products competing with yours.

amazon web services

AWS Products & Solutions ▼

| Legal |
| --- |
| ▪ AWS Acceptable Use Policy |
| ▪ **AWS Customer Agreement** |

**AWS Customer Agreement**

**Last updated March 15, 2012**
*(current AWS customers: See What's Changed)*

Example of a typical ToC (Amazon)

**13.3 Independent Contractors; Non-Exclusive Rights.** We and you are independent contractors, and neither party, nor any of their respective affiliates, is an agent of the other for any purpose or has the authority to bind the other. Both parties reserve the right (a) to develop or have developed for it products, services, concepts, systems, or techniques that are similar to or compete with the products, services, concepts, systems, or techniques developed or contemplated by the other party and (b) to assist third party developers or systems integrators who may offer products or services which compete with the other party's products or services.

# Unavailability

- **Limitations of liability** in case of unavailability of data or services.

  - Due e.g. to power outages, system failures, or to any other service interruption.

  - Or due to unauthorized access, alteration, loss, or anything else of data or any other content stored in AWS.

**11. Limitations of Liability.**

WE AND OUR AFFILIATES OR LICENSORS WILL NOT BE LIABLE TO YOU FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, CONSEQUENTIAL OR EXEMPLARY DAMAGES (INCLUDING DAMAGES FOR LOSS OF PROFITS, GOODWILL, USE, OR DATA), EVEN IF A PARTY HAS BEEN ADVISED OF THE POSSIBILITY OF SUCH DAMAGES. FURTHER, NEITHER WE NOR ANY OF OUR AFFILIATES OR LICENSORS WILL BE RESPONSIBLE FOR ANY COMPENSATION, REIMBURSEMENT, OR DAMAGES ARISING IN CONNECTION WITH: (A) YOUR INABILITY TO USE THE SERVICES, INCLUDING AS A RESULT OF ANY (I) TERMINATION OR SUSPENSION OF THIS AGREEMENT OR YOUR USE OF OR ACCESS TO THE SERVICE OFFERINGS, (II) OUR DISCONTINUATION OF ANY OR ALL OF THE SERVICE OFFERINGS, OR, (III) WITHOUT LIMITING ANY OBLIGATIONS UNDER THE SLAS, ANY UNANTICIPATED OR UNSCHEDULED DOWNTIME OF ALL OR A PORTION OF THE SERVICES FOR ANY REASON, INCLUDING AS A RESULT OF POWER OUTAGES, SYSTEM FAILURES OR OTHER INTERRUPTIONS; (B) THE COST OF PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; (c) ANY INVESTMENTS, EXPENDITURES, OR COMMITMENTS BY YOU IN CONNECTION WITH THIS AGREEMENT OR YOUR USE OF OR ACCESS TO THE SERVICE OFFERINGS; OR (D) ANY UNAUTHORIZED ACCESS TO, ALTERATION OF, OR THE DELETION, DESTRUCTION, DAMAGE, LOSS OR FAILURE TO STORE ANY OF YOUR CONTENT OR OTHER DATA. IN ANY CASE, OUR AND OUR AFFILIATES' AND LICENSORS' AGGREGATE LIABILITY UNDER THIS AGREEMENT WILL BE LIMITED TO THE AMOUNT YOU ACTUALLY PAY US UNDER THIS AGREEMENT FOR THE SERVICE THAT GAVE RISE TO THE CLAIM DURING THE 12 MONTHS PRECEDING THE CLAIM.

# No guarantees

- **Disclaimers = no guarantee** that the service will be "uninterrupted, error free or free of harmful components".
  - Or that what is stored in AWS is safe, is not lost, or damaged.
  - What if I decided to use AWS to store my scientific data (some tens of PB maybe...)

## 10. Disclaimers.

THE SERVICE OFFERINGS ARE PROVIDED "AS IS." WE AND OUR AFFILIATES AND LICENSORS MAKE NO REPRESENTATIONS OR WARRANTIES OF ANY KIND, WHETHER EXPRESS, IMPLIED, STATUTORY OR OTHERWISE REGARDING THE SERVICE OFFERINGS OR THE THIRD PARTY CONTENT, INCLUDING ANY WARRANTY THAT THE SERVICE OFFERINGS OR THIRD PARTY CONTENT WILL BE UNINTERRUPTED, ERROR FREE OR FREE OF HARMFUL COMPONENTS, OR THAT ANY CONTENT, INCLUDING YOUR CONTENT OR THE THIRD PARTY CONTENT, WILL BE SECURE OR NOT OTHERWISE LOST OR DAMAGED. EXCEPT TO THE EXTENT PROHIBITED BY LAW, WE AND OUR AFFILIATES AND LICENSORS DISCLAIM ALL WARRANTIES, INCLUDING ANY IMPLIED WARRANTIES OF MERCHANTABILITY, SATISFACTORY QUALITY, FITNESS FOR A PARTICULAR PURPOSE, NON-INFRINGEMENT, OR QUIET ENJOYMENT, AND ANY WARRANTIES ARISING OUT OF ANY COURSE OF DEALING OR USAGE OF TRADE.

# But you are responsible

- **You are responsible to make sure your data, code, etc. is safe, protected from unauthorized access, and *you are responsible for your own backup*** (again – with if it's in the order of several PB?)

**4.2 Other Security and Backup.** You are responsible for properly configuring and using the Service Offerings and taking your own steps to maintain appropriate security, protection and backup of Your Content, which may include the use of encryption technology to protect Your Content from unauthorized access and routine archiving Your Content. AWS log-in credentials and private keys generated by the Services are for your internal use only and you may not sell, transfer or sublicense them to any other entity or person, except that you may disclose your private key to your agents and subcontractors performing work on your behalf.

# Data property / privacy?

- When a contract with a Cloud provider gets cancelled, how can we make sure that **all our data is removed**?

- And how can I avoid *vendor lock-in*?

- But where is my data? How about *tapping*?

**Edward Snowden** @Snowden
Follow

The New York Times: @FBI's war on #Apple will aid China. nytimes.com/2016/02/18/tec...

> China is watching the dispute closely. Analysts say that the Chinese government does take cues from the United States when it comes to encryption regulations, and that it would most likely demand that multinational companies provide accommodations similar to those in the United States.
>
> Last year, Beijing backed off several proposals that would have mandated that foreign firms provide encryption keys for devices sold in China after heavy pressure from foreign trade groups...
>
> "...a push from American law enforcement agencies to unlock iPhones would embolden Beijing to demand the same."

RETWEETS 3,002    LIKES 2,299

1:43 PM - 17 Feb 2016

3K        2.3K

**Microsoft admits Patriot Act can access EU-based cloud data**

Microsoft's U.K. head admitted today that no cloud data is safe from the Patriot Act, and the company can be forced to hand EU-stored data over to U.S. authorities.

By Zack Whittaker for iGeneration | June 28, 2011 -- 08:10 GMT (09:10 BST) | Topic: Government : US

**NSA infiltrates links to Yahoo, Google data centers worldwide, Snowden documents say**

Some of the sites "impossible to be deleted from" (the ones in **black**):
- Backblaze (Cloud backup)
- Blogger, Picasa (and other Google services)
- eDreams
- Evernote
- Netflix
- OpenShift
- Pastebin
- PlayStation Network
- Slashdot
- Udacity
- Wikipedia
- Wordpress.com

# justdelete.me

A directory of direct links to delete your account from web services.

Chrome Extension    Fork on GitHub    Tweet JDM

POPULAR    A - Z    DIFFICULTY    RESET

| 4shared | 500px | 9GAG | Abload |
|---|---|---|---|
| EASY | EASY | EASY | EASY |
| NO INFO AVAILABLE | NO INFO AVAILABLE | SHOW INFO... | NO INFO AVAILABLE |

| About.me | Adobe | Affero | Airbnb |
|---|---|---|---|
| EASY | HARD | EASY | EASY |
| NO INFO AVAILABLE | SHOW INFO... | SHOW INFO... | NO INFO AVAILABLE |

| Album Reminder | Alibaba | Alvanista | Amara |
|---|---|---|---|
| EASY | HARD | EASY | EASY |
| NO INFO AVAILABLE | SHOW INFO... | SHOW INFO... | SHOW INFO... |

| Amazon | Amazon AWS | Animal Crossing Community | AOL / Instant Messenger |
|---|---|---|---|
| HARD | EASY | IMPOSSIBLE | EASY |
| SHOW INFO... | SHOW INFO... | SHOW INFO... | NO INFO AVAILABLE |

| App.net | AppFog | Argyle Social | ArmorGames |
|---|---|---|---|
| EASY | HARD | IMPOSSIBLE | EASY |
| NO INFO AVAILABLE | SHOW INFO... | SHOW INFO... | NO INFO AVAILABLE |

| Artsy | Ask.fm | Asos | Assembla |
|---|---|---|---|
| EASY | EASY | HARD | EASY |
| NO INFO AVAILABLE | NO INFO AVAILABLE | SHOW INFO... | NO INFO AVAILABLE |

| Audiomack | AutoScout24 | Avast! | Awdio.com |
|---|---|---|---|
| EASY | EASY | EASY | HARD |
| NO INFO AVAILABLE | NO INFO AVAILABLE | SHOW INFO... | SHOW INFO... |

# Last but not least, the big misunderstanding

- **Capacity is not infinite** (although this is one of the postulates of Cloud computing). *Nor are credit card limits*.

  - Hence, resources might not be available when we need them; or, if available, they might not have the characteristics we need.

  - Unless maybe we are willing to pay some hefty over-provisioning costs.

# High Performance Computing (HPC)

# New architectures

*Up to now HEP computing is based on a single architecture (x86-64)*

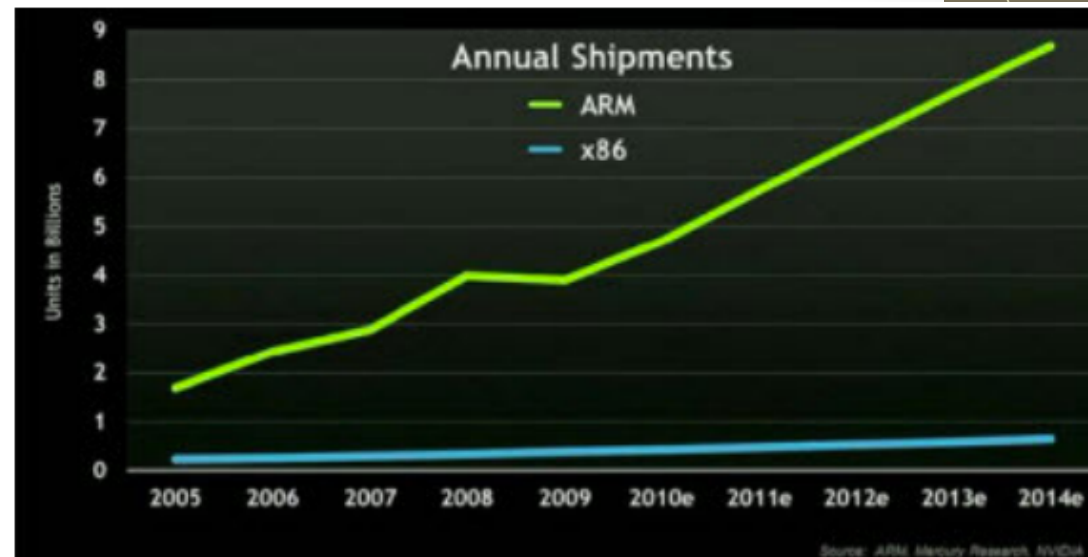*→ Follow the market mainstream*

*→ Use highly available architectures*
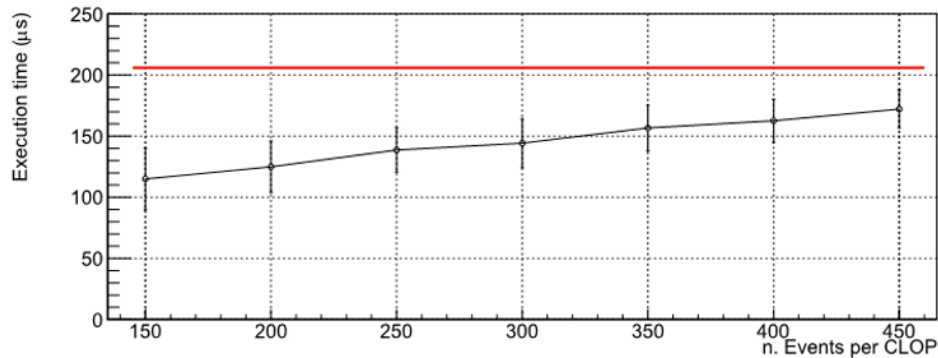
*ARM, …*

*→ Exploit parallelization*

*Multi/many-core,*

*GPGPU, …*

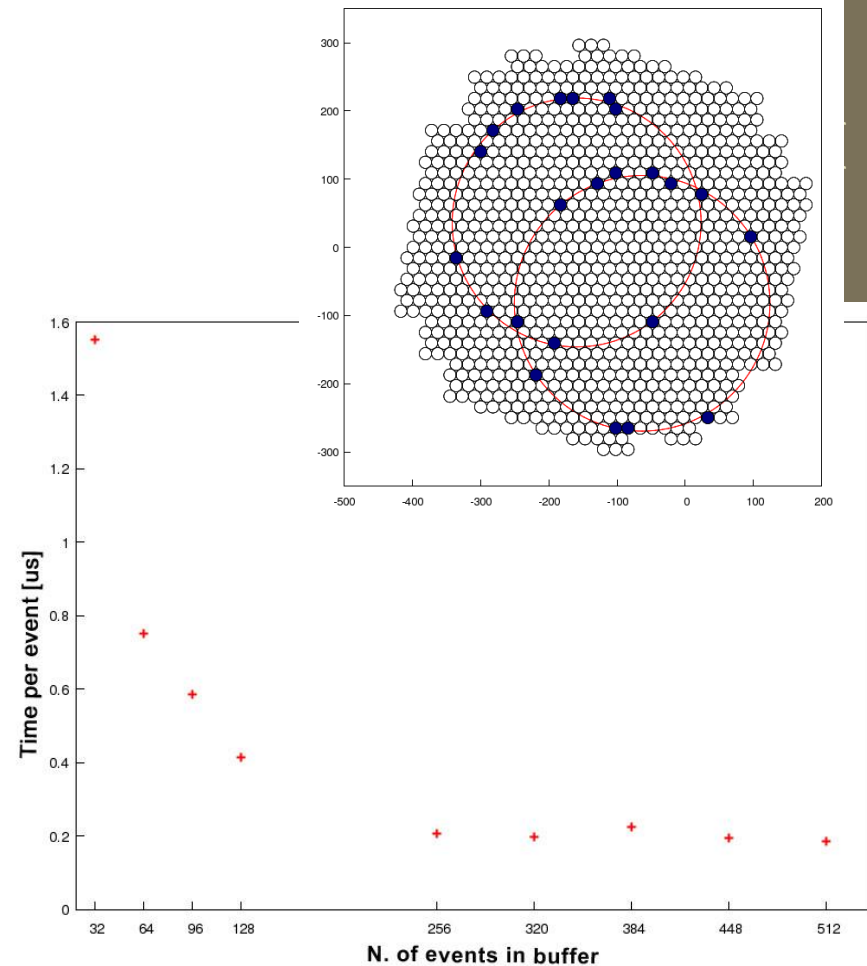*→ Use low-power architectures*

# *Exploit hardware  capabilities*



*Algorithm parallelization* *in*

*NA64 RICH pattern recognition*

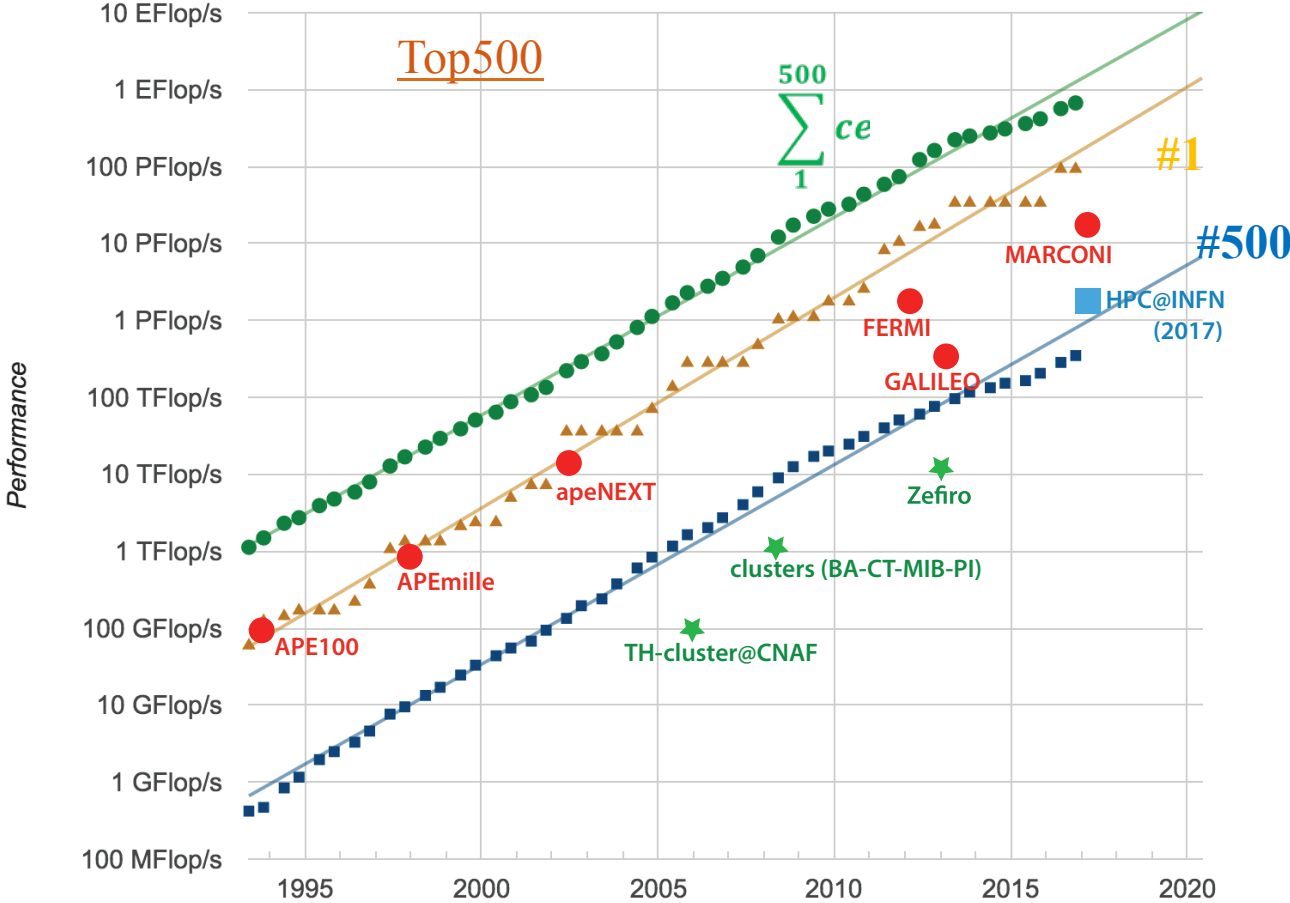*Execution on NaNet-10*

    *Based on GPU the Tesla K20c*
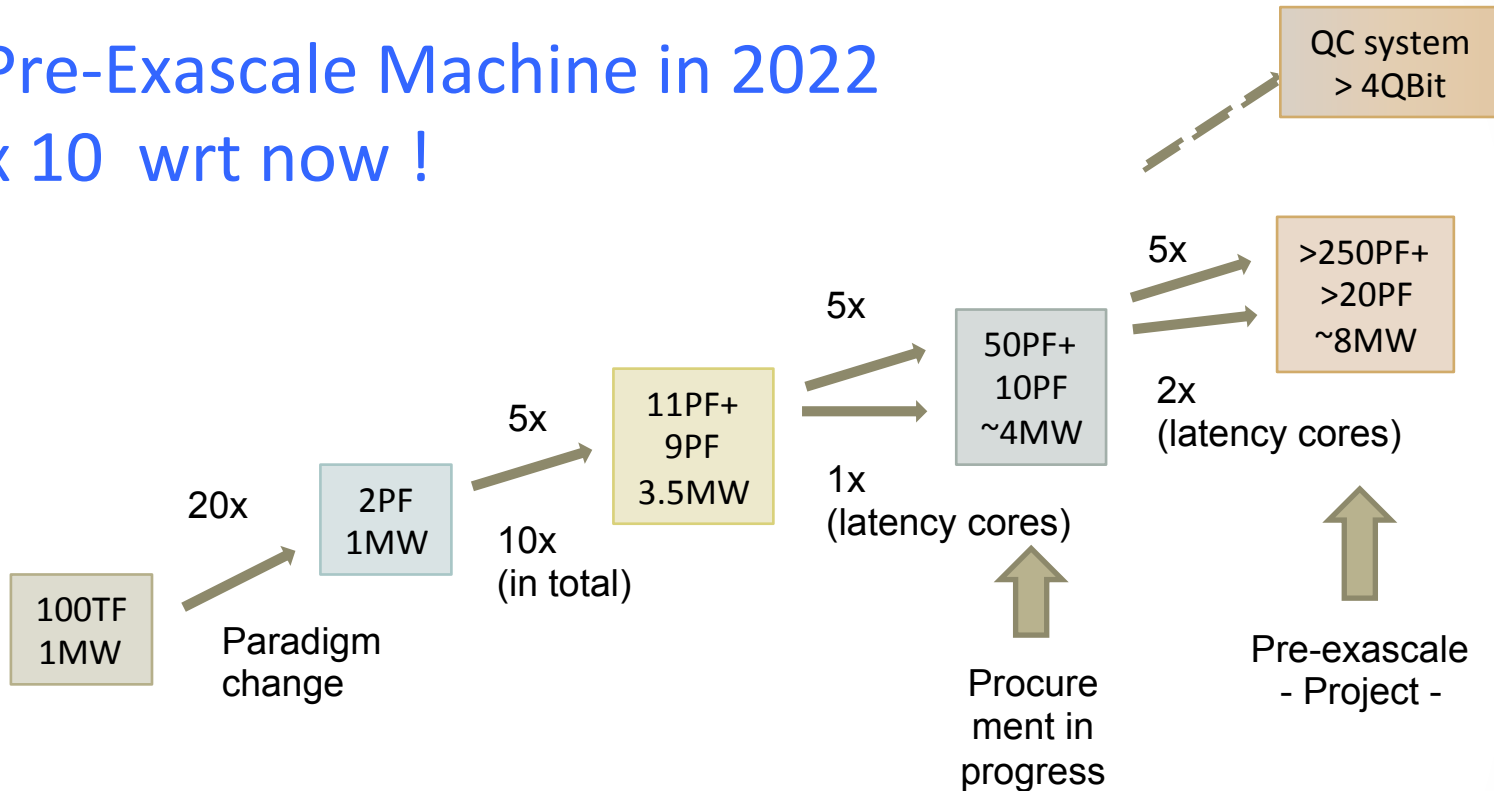
# High performance computing



Projected Performance Development

HPC computing power vs years

# Cineca (Italy) Roadmap

**Pre-Exascale Machine in 2022**
**x 10  wrt now !**



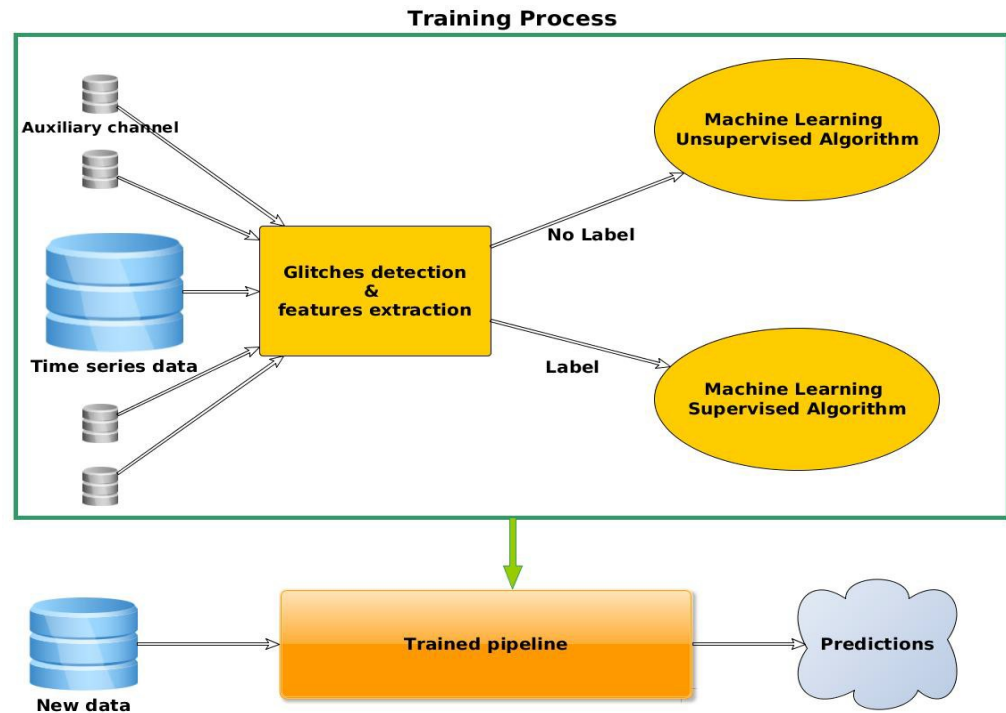| 2009 | 2012 | 2016 | 2019 | 2021/2022 |
|------|------|------|------|-----------|
| IBM SP6 Power6 | Fermi IBM BGQ PowerA2 | Marconi Lenovo Xeon+KNL | *To be defined* Scalar + Vector and / or Accelerator | *To be defined* Scalar + Vector and / or Accelerator |

# Machine/Deep Learning

*Starting adopting Machine Learning & Deep Learning techniques for data processing*

*Example:*

*Glitches detection in Gravitational Waves searches*
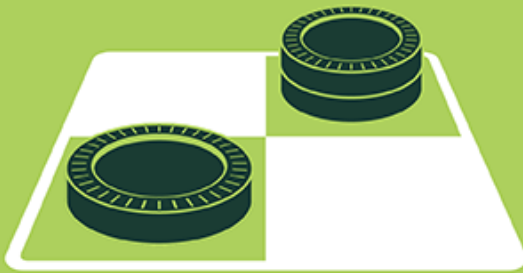
GlitchesClassificationStrategy

# Why machine learning

- It is very hard to write programs that solve problems like recognizing a three-dimensional object from a novel viewpoint in new lighting conditions in a cluttered scene.

- It is hard to write a program to compute the probability that a credit card transaction is fraudulent, that an object (car, animal, ..)  is present in a picture etc.

# ARTIFICIAL INTELLIGENCE

Early artificial intelligence stirs excitement.

# MACHINE LEARNING

Machine learning begins to flourish.

# DEEP LEARNING

Deep learning breakthroughs drive AI boom.

1950's   1960's   1970's   1980's   1990's   2000's   2010's

Since an early flush of optimism in the 1950s, smaller subsets of artificial intelligence – first machine learning, then deep learning, a subset of machine learning – have created ever larger disruptions.

# The machine learning approach

- Definition Machine Learning is a field of study that gives computers the ability to learn without being explicitly programmed [Arthur Samuel,1959]

- Instead of writing a program by hand for each specific task, we collect lots of examples that specify the correct output for a given input.

- A machine learning algorithm then takes these examples and produces a program that does the job.

- Massive amounts of computation are now cheaper than paying someone to write a task-specific program.

# Some examples of tasks best solved by learning

- Recognizing patterns:
  - Objects in real scenes
  - Facial identities or facial expressions
  - Spoken words
- Recognizing anomalies:
  - Unusual sequences of credit card transactions
  - Unusual patterns of sensor readings in a nuclear power plant
- Prediction:
  - Future stock prices or currency exchange rates
  - Which movies will a person like?

# Types of learning tasks

- Supervised learning
  - Learn to predict an output when given an input vector.
  - Each training example consists of an input vector x and a target output t.

- Unsupervised learning
  - Discover a good internal representation of the input

So, 1. **what exactly is deep learning** ?

And, 2. **why is it generally better** than other methods
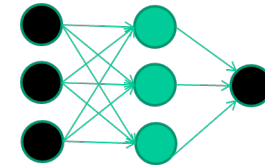on image, speech and certain other types of data?

**The short answers**

1. **'Deep Learning' means** using a **neural network**
with **several layers of nodes** between input and output

2. the series of layers between input & output do
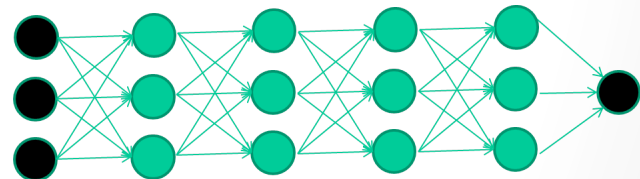feature identification and processing in a series of stages,
just as our brains seem to.

hmmm… OK, but:

3. **multilayer neural networks have been around for 25 years. What's actually new?**

we have always had good algorithms for learning the weights in networks with 1 hidden layer

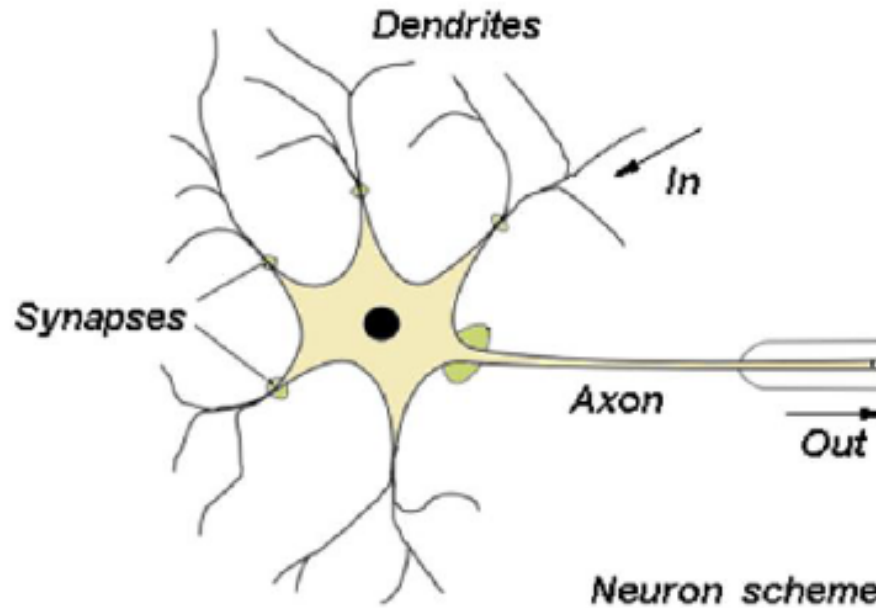but these algorithms are not good at learning the weights for networks with more hidden layers
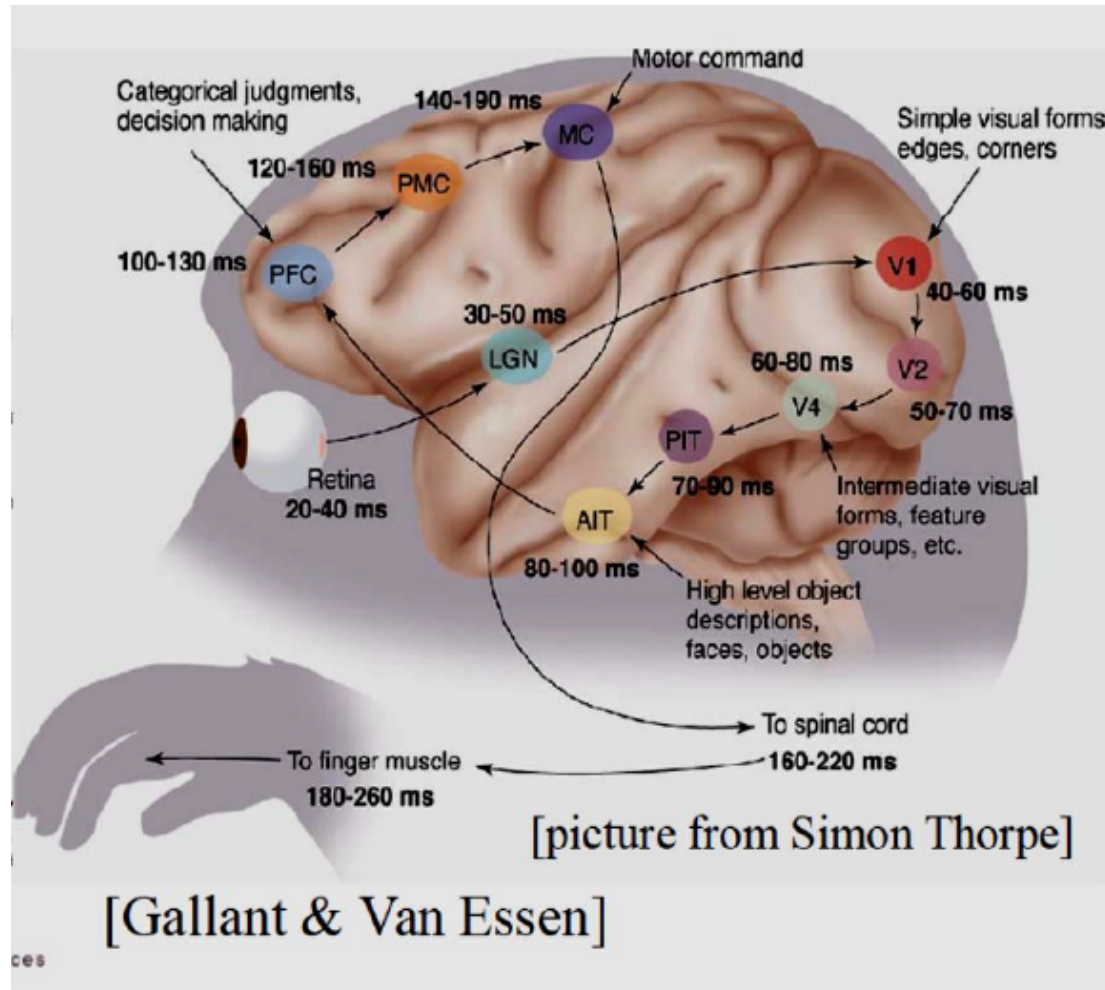
what's new is: algorithms for training many-later networks

# Deep learning: mimic the brain

Neuron scheme

- Each neuron receives inputs from other neurons
- The effect of each input line on the neuron is controlled by a synaptic weight
- The synaptic weights adapt so that the whole network learns to perform useful computations
- There are about 10^11 neurons each with about 10^4 weights.

# Deep learning: mimic the brain

The ventral (recognition) pathway in the visual cortex has multiple stages



[picture from Simon Thorpe]

[Gallant & Van Essen]

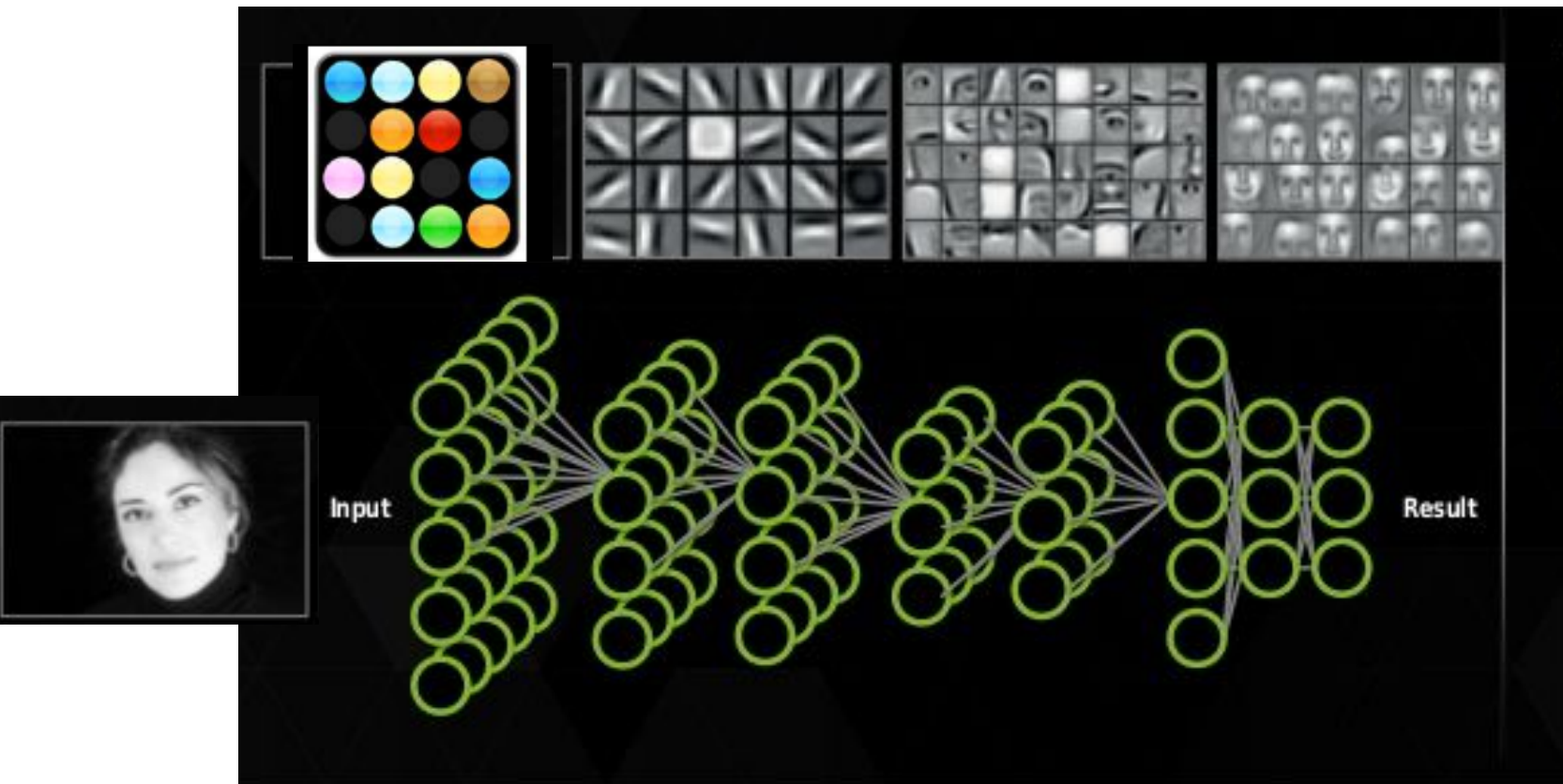# Deep Learning: multilayer approach

Unsupervised: Need to train the network with many input data

**pixels**        **edges**        **object parts**        **object models**

# Computing resources
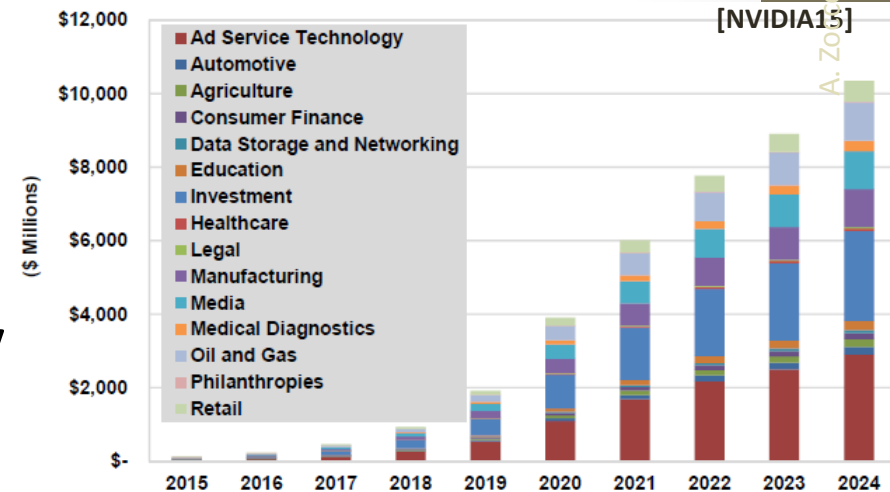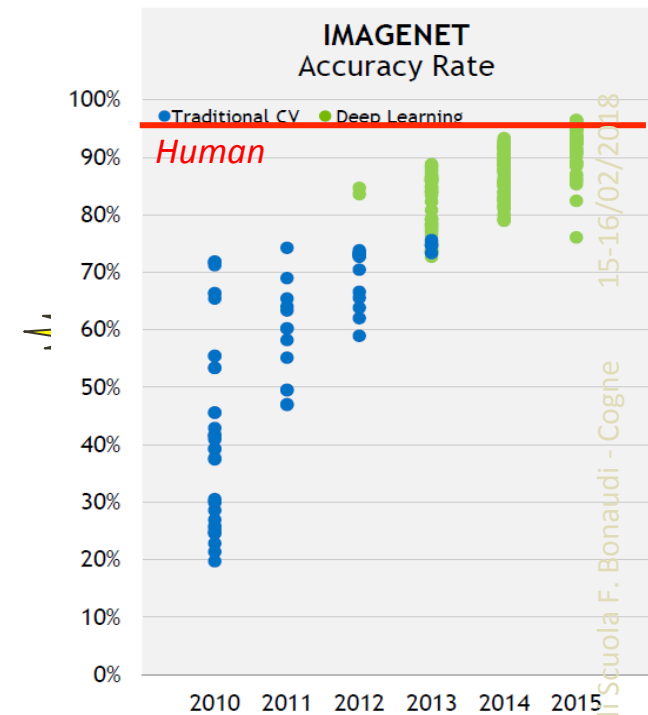


Today's Largest Networks

~10 layers
1B parameters
10M images
~30 Exaflops
~30 GPU days

Human brain has trillions of parameters - only 1,000 more.
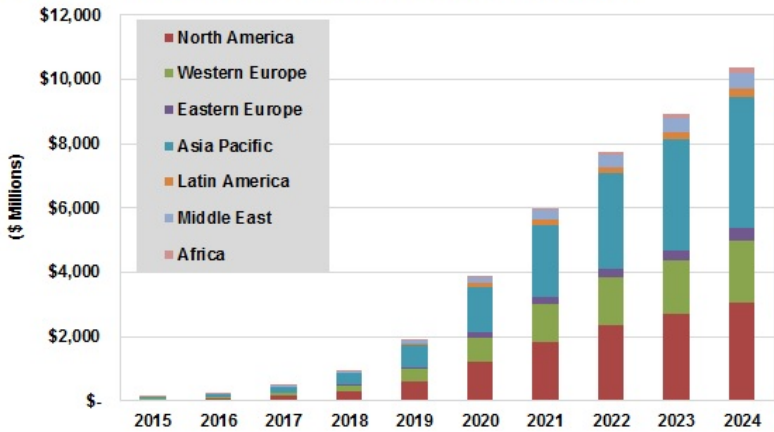
# Deep Learning applications

- It's a Big Data / Supercomputing Problem:

  - 1 Exaflop to train a SOA deep network [Baidu]

  - Need fast turnaround time for idea→test→code

  - Unsupervised learning is orders of magnitude more demanding than supervised learning

- Supply chains Pilot

  - **Automatic vision-based quality inspection**



IMAGENET
Accuracy Rate

[NVIDIA15]

# Computer Vision and Deep Learning Market



Deep Learning



Computer Vision

The market for **computer vision** technologies will grow from $5.7 billion in 2014 to $33.3 billion by 2019, representing CAGR of 42%

**The machine vision market** size is estimated to grow from USD 8.08 billion in 2015 to USD 12.49 billion by 2020, at an estimated CAGR of 9.1%.

**3D Machine Vision Market** Global Forecast to 2020 says, the market is expected to grow at a CAGR of 10.53% during the forecast period between 2015 and 2020
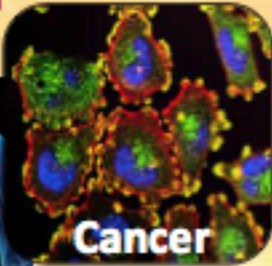
In "**Automated Guided Vehicle Market**", the total market is expected to reach USD 2.81 Billion by 2022, at a CAGR of 10.2%

97

# The european strategy for computing

# Future: Big data & HPC

# Future: Big data & HPC

# Future: Big data & HPC

# Facts and Figures

- **EU: ~33%** *of GDP and* **~60 million jobs**
- **USA: ~55%** *of GDP*

**For 97%** of industry using HPC: it is key for innovating

## Very high Return on Investment (RoI)

- *In Europe*: every 1€ invested in HPC generated ~€870 in revenues for businesses and €69 in profits
- *In Japan*: every $ invested by projects run on the K-supercomputer yielded a RoI of $571 in revenue

# The European scenario

# Eureopean Cloud Initiative (EOSC)



19 April 2016
Blueprint of
EC

# European Open Science Cloud (EOSC)

**Carlos Moedas – Commissioner for Research, Science and Innovation**

Submitted by alim on 18 Sep 2015



European Commission - Speech - [Check Against Delivery]

First, we are preparing a call for European Science Cloud Project in order to identify the possibility of creating a cloud for our scientists. We need more open access to research results and the underlying data. Open access publication is already a requirement under Horizon 2020, but we now need to look seriously at open data.

Launched directly by Commission
Formed an High Level Expert Group
October 11 released a first report
HLEG Report

**H2020 call Infradev-04-2016: Proposal EOSCPilot on European Open Science Cloud for Research**

# European Data Infrastructure

**Pilot project**
*Important Project of Common European Interest* **IPCEI**

# IPCEI *on* **High Performance Computing and Big Data Enabled Applications**

**IMPORTANT PROJECT OF COMMON EUROPEAN INTEREST (IPCEI)**

ON
HIGH PERFORMANCE COMPUTING
AND
BIG DATA ENABLED APPLICATIONS
(IPCEI-HPC-BDA)

**European Strategic Positioning Paper**

1. **European exa-scale technology**
   prove the capacity of the European industry to answer the challenges of building an exa-scale machine by 2023 via a prototype to be ready by 2020.

2. **Large test beds and applications**
   Establish Centers of Excellence, starting from research, to develop and test HPC-enabled and big data based applications in specific and strategic sectors at regional, national and pan-European scale;
   Deploy application test beds on:
   Personalized Medicine, Smart Space, Industry 4.0 and Smart Manufacturing, New advanced Materials,Fintech, Smart Agrifood and Smart City Applications.

A. Zoccoli Scuola

# EuroHPC: Declaration
# Cooperation Framework on High Performance Computing



DECLARATION

Cooperation framework on High Performance Computing

Bundesrepublik Deutschland

and

República Portuguesa

and

République française

and

Reino de España

and

Repubblica Italiana

and

Grand–Duché de Luxembourg

and

the Netherlands



Agree to work together towards making available across EU an integrated world class computing (HPC) infrastructure which in combination with European data and network infrastructures would upraise Europe Scientific capabilities and industrial competitiveness.

# Member states

➤ Agree to work towards the establishment of a cooperation framework – EuroHPC- for acquiring and deploying an integrated exascale supercomputer

➤ Agree to work together and with EC to prepare by the end of 2017 a roadmap to address:
- Procurements process for acquisition of 2 world-wide pre-exascale computer in 2019-2020 and 2 exascale computers by 2023
- Development of high quality competitive European technology and its optimization through co-design approach
- Development of test-beds for HPC and big data applications for scientific, public administration and industrial purposes
- ……………………

# Italian computing infrastructure

# L'Istituto Nazionale di Fisica Nucleare

The National Institute for Nuclear Physics (INFN) is the Italian research agency dedicated to the study of the fundamental constituents of matter and the laws that govern them, under the supervision of the Ministry of Education, Universities and Research (MIUR). It conducts theoretical and experimental research in the fields of subnuclear, nuclear and astroparticle physics.

# The Italian Network for Research, Education …

- 100 Universities, Conservatories and art Academies

- 350 Research Institutes and Laboratories

- 60 Biomedical Research Institutes

- 65 Libraries, Museum and Cultural Institutions

- More than 300 schools

# Current INFN e-Infrastructure

**Tier1**:
- 200 kHS06 CPU power equivalent to 20000 cores
- 19 PB disk space
- 57 PB tape library

**9 Tier2**:
- 240 kHS06 CPU power equivalent to 24000 cores
- 18 PB disk space

**Network** provided by GARR Consortium (10-100 Gbps), under upgrading



Major part of the computing centers are multidisciplinary (HEP, astroparticle, neutrino, etc. )

# Resources distribution

INFN



SUM Wallclock Work HS06 Hours (hours) by Country

- United States of America (23.34%)
- Other (12.81%)
- Canada (2.08%)
- France (7.89%)
- Germany (8.96%)
- Italy (7.97%)
- Netherlands (2.55%)
- Russia (3.46%)
- Spain (2.45%)
- Switzerland (18.54%)
- United Kingdom (9.94%)

Italian contribution 8%

In the world → WLCG

Half of the resources @ Tier1

In Italy

Very effective infrastructure for LHC data analysis !



Normalized Elapsed time (HEPSPEC06) * Number of Processors (hours) by Resource Centre

- INFN-T1 (47.04%)
- Other (7.20%)
- INFN-BARI (9.43%)
- INFN-CATANIA (2.62%)
- INFN-COSENZA (3.37%)
- INFN-FRASCATI (3.19%)
- INFN-LNL-2 (10.52%)
- INFN-MILANO-ATLASC (2.08%)
- INFN-NAPOLI-ATLAS (6.01%)
- INFN-PISA (5.66%)
- INFN-ROMA1 (2.87%)

# Costs & manpower

**Budget**

| Item | Cost (M€) |
|---|---|
| CPU, disk and tape | 4 |
| HPC@CINECA | 0.5 |
| Electric power | 2 |
| Network(@GARR) | 5 |
| Total/year | 11.5 |

| Infrastructure | Middleware and software |
|---|---|
| 50 FTE | 30 FTE |

**Personnel**

# *Italian resources in 2017*

*Let's take CNAF, the Italian Tier-1, as an example to understand what changes:*

|  | CPU (kHS06) | Disk (PB) | Tape (PB) |
|---|---|---|---|
| *All WLCG* | *5200* | *340* | *590* |
| *INFN Tier-1 & 2* | *520* | *38* | *57* |
| *% INFN* | *10* | *11* | *10* |

*From: https://wlcg-rebus.cern.ch/*

# CNAF evolution - LHC Run 1 & 2



*Run2 is ok with the flat budget hypothesis:*

- *CPU + 20 - 30%*
- *Disk + 15 - 25%*
- *Tape + 30% - 60%*

A. Zoccoli Scuola F. Bonaudi - Cogne

# *CNAF evolution up to LHC Run 3*

# HPC in Italy, a crumb of history

*In 1984 Nicola Cabibbo, Giorgio Parisi and young students participated at a workshop on Lattice Gauge Theory.*

*The formal conclusions of the workshop were that even thinking to build a computer for Lattice Gauge Theory was insane.*

*Despite the official conclusions of the meeting, after a few hours of discussion the basic ideas of APE (Array Processor Experiment) were sketched. In the following months it was quickly organized a scientific collaboration led by Nicola Cabibbo and Giorgio Parisi involving INFN of Padua, Pisa and Rome, and the CNAF.*

**The Processor of the first APE**



Nicola Cabibbo — The APE experience — 8/2/2007 — 6 / 21

**I Progetti APE — i primi 10 anni**

## I Progetti APE

**1984-1989: APE1**

- 16 nodi di calcolo, 1 Gflops
- Software "primitivo"
- Prodotti alcuni prototipi



**1990-1995: APE100**

- Modulare, 2048 nodi "custom", 100 GFlops
- Sviluppo di un linguaggio dedicato (TAO)
- Ambiente software "user friendly"
- Affidabilita' alta, 300 GFlops installati

Nicola Cabibbo — Il progetto APE — 8/2/2006 — 4 / 20

121

# *HPC italian infrastructure*

Cineca - SuperComputing Applications & Innovations

# CINECA (Bologna)

| Logical Name | Tier 0 - FERMI (June 2012) | Tier 1 - GALILEO (December 2014) | Big data - PICO (October 2014) |
|---|---|---|---|
| Peak Performance | ~ 2 Pflops; ~ 5 PByte | ~ 0,5 PFlops | ~ 0,3 Pflops; ~ 15 Pbyte |

| Logical Name | MARCONI (2016 / 2017) | BIG DATA (2016) |
|---|---|---|
| Peak Performance | ~ 20 Pflops; ~ 15 Pbyte | ~1 Pflops; ~ 20 Pbyte |

| Logical Name | Tier 0 BIG DATA (2019 – 2020) |
|---|---|
| Peak Performance | < 50 Pflops; < 50 Pbyte on line storage; < 50 Pbyte repository |

123

# HPC today in Europe and Italy

| System Name | Hosting Center | Architecture | Capacity |
|---|---|---|---|
| CURIE | GENCI@CEA | Bull x86 | 2 PFlop/s |
| MARCONI | CINECA | Intel Broadwell | 20 PFlop/s |
| HAZEL HEN | HLRS | Cray XC40 | 7.42 PFlop/s |
| JUQUEEN | GCS@FZJ | IBM BlueGene/Q | 5.87 PFlop/s |
| MareNostrum | BSC | IBM iDataPlex | 1 PFlop/s |
| Piz Daint | CSCS | Cray XC30 | 7.8 Pflop/s |



HPC has already European dimension

- INFN has an agreement with CINECA renewed every three years
- INFN participates to H2020 HPC project:
    - ExaNeSt, to study network and storage for exa-scale facilities
    - EuroExa, to build a prototype of exa-scale machine with new network

# Cineca (Italy) Roadmap

**Pre-Exascale Machine in 2022**
**x 10 wrt now !**

QC system
> 4QBit

>250PF+
>20PF
~8MW

5x

2x
(latency cores)

Pre-exascale
- Project -

50PF+
10PF
~4MW

5x

1x
(latency cores)

Procure
ment in
progress

11PF+
9PF
3.5MW

5x

10x
(in total)

2PF
1MW

20x

Paradigm
change

100TF
1MW

| 2009 | 2012 | 2016 | 2019 | 2021/2022 |
|------|------|------|------|-----------|
| IBM SP6 Power6 | Fermi IBM BGQ PowerA2 | Marconi Lenovo Xeon+KNL | *To be defined* Scalar + Vector and / or Accelerator | *To be defined* Scalar + Vector and / or Accelerator |

# Where are we going ?

# New italian integrated e-infrastructure

Integration of CINECA-HPC and INFN-HTC computing infrastructure and progressive inclusion of the others qualified systems owned by national entities (CMCC for weather, climate and ocean predictions and simulations; ENEA for ITER, INAF for SKA, INGV for EPOS…):

- Institutional basic and applied research
- Enabling for Public administrations
- Proof of concept and innovation for private organizations and industries

# INFN-CNAF (HTC) connection to CINECA (HPC)

# Integration Model



Research

Users Users · · Users · · Users Users

Private & Public Sectors

Service Layer

| HPC & Big Data | Data Infrastructure | HTC & Big Data |
|---|---|---|

Network & Transport layer

# Opportunities

◆Building-up an integrated super-computing center, probably the most powerful and flexible in Europe !

◆Complementarity of HPC and HTC resources

◆Easy user accessibility

◆Powerful infrastructure for research in HEP and other disciplines

◆… but also for private companies (Sme's etc. )

# Bologna Big Data Technopole



Conference and Education Center

Big Data Foundation

«Ballette» innovation and creativity center

ECMWF Data Center

CINECA & INFN Supercomputing center

Competence center Industry 4.0

Protezione civile and regional agency for development and innovation

Enea center

IOR bio-banks Regenerative medicine center

Academic Doctoral School

Agenzia Nazionale Meteo

131

# Bologna Big Data Technopole



The Emilia Romagna Region (ERR) is the National and European hub for big data: about the 70% of the Italian research data is stored / processed in research centers sited in this region.

CINECA and INFN have their top tier computing facilities and data centers in ERR, integrating infrastructures of the top universities and National Research institutions.

ERR is the headquarters of world-leading companies in the sectors of automotive, mechatronics, digital production, bio-medical, e-commerce, agri-food

ERR is the Italian leading region in term of economic growth rate and one of the most dynamic regions in Europe

# Bologna Big Data Technopole

**Private and public foundations and Institutions:**

Opificio Golinelli for the Doctoral School in *Data Science and Computation*, Mast Cultural and Philanthropic Foundation, Marco Biagi Foundation in Labor Law Research, Bologna Business School, Foundation of Religious Science John XXIII, Nomisma Economic Research, Prometeia Consulting on risk, wealth and performance, European Food Safety Authority, Consortium for the future in Research, …

**Leading IT companies:**
IBM, Yoox, CRIF, Engineering, Dedagruop, Eon reality, T3Lab, …



Consiglio Nazionale delle Ricerche
www.cnr.it

Inter-University Consortium for supercomputing
www.cineca.it

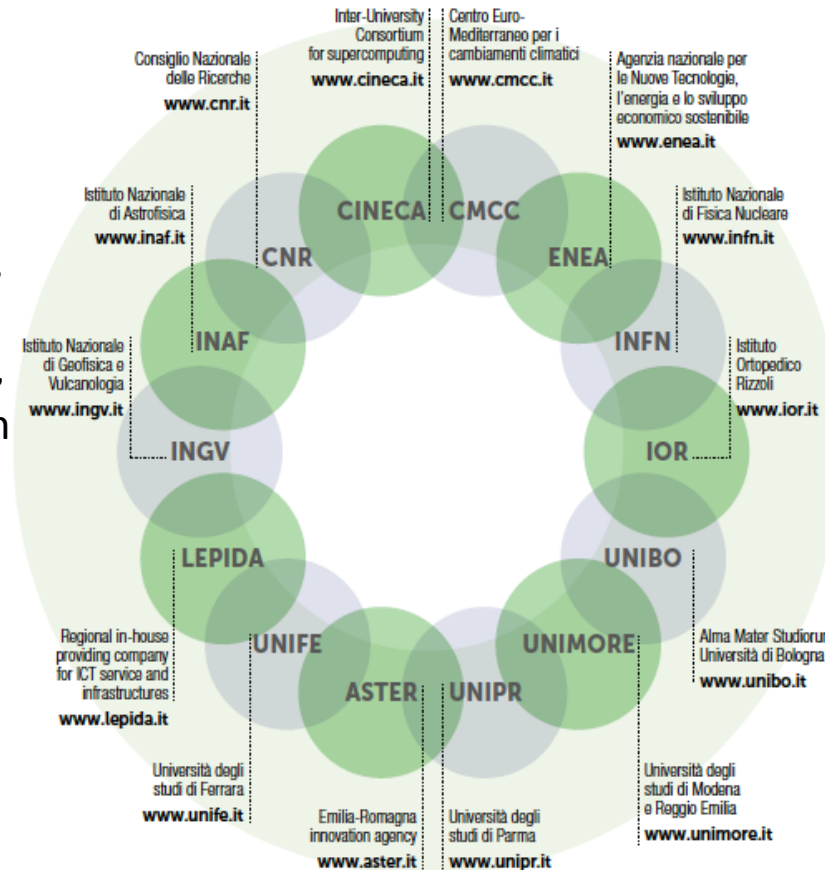Centro Euro-Mediterraneo per i cambiamenti climatici
www.cmcc.it

Agenzia nazionale per le Nuove Tecnologie, l'energia e lo sviluppo economico sostenibile
www.enea.it

Istituto Nazionale di Astrofisica
www.inaf.it

Istituto Nazionale di Fisica Nucleare
www.infn.it

Istituto Nazionale di Geofisica e Vulcanologia
www.ingv.it

Istituto Ortopedico Rizzoli
www.ior.it

Regional in-house providing company for ICT service and infrastructures
www.lepida.it

Alma Mater Studiorum Università di Bologna
www.unibo.it

Università degli studi di Ferrara
www.unife.it

Emilia-Romagna innovation agency
www.aster.it

Università degli studi di Parma
www.unipr.it

Università degli studi di Modena e Reggio Emilia
www.unimore.it

**Strong International connections** among universities, companies, institutions

**Connectivity**
LEPIDA, GARR, GEANT

**HW Infrastructure**
CINECA, INFN

**SW Infrastructure**
HPC Center of excellence in Material Design, Operational work flow for Climate Services, Visual information technology lab, Genomic and bioinformatics data processing, …

**End users National and European**
University system, Research Institutions and Agencies, Public administrations, Private and industrial organizations

15-16/02/2018

A. Zoccoli Scuola F. Bonaudi - Cogne

133

Fundamental analytics and applications.
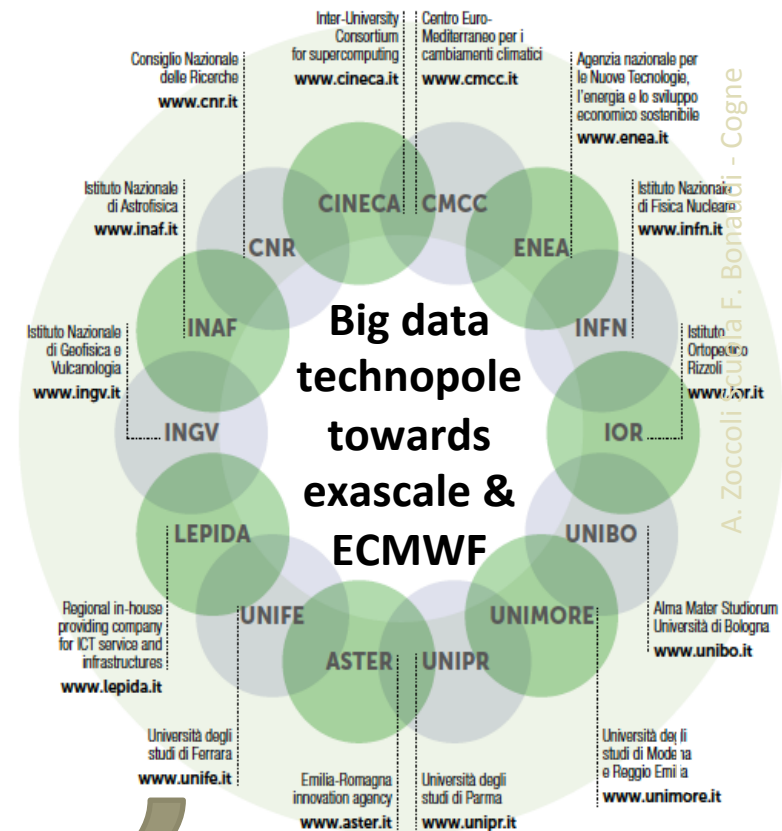Scientific, Industrial and societal challenges.

- Climate change and services
- Welfare, health and aging
- Production and digital transformation (e.g. mechatronics, automotive, agri-food, biomedical and pharma sectors, etc)
- Cultural heritage, humanities and society
- Sustainable cities
- Security, Cyber security and artificial intelligence
- Education and skills

A. Zoccoli ...... F. Bonaccorsi – Cogne

# The final goal

- provide a common infrastructure to the different research communities (physics, astrophysics, Biology, Medicine, engennering, …)
- but also to public and private sectors (test beds)
- attract National, regional and European funds

| Research | Industry4.0 Vertical test bed | Health Vertical test bed | Material Science test bed |
|---|---|---|---|

INFN + CINECA Federated HW Infrastructure

# Summary

❑ *HEP computing is continuously evolving*

❑ *Experiment requests impose an evolution of the model in order to comply with the (flat) budget*

❑ *Need to understand and exploit new technologies*

❑ *Software is the key to scalability and sustainability*

❑ *Flexibility, cloud and exascale computers will play a key role*

❑ *There is room for new ideas and innovative projects!*

❑ *……. and a final remark….*

"Big data is like teenage sex;
everyone talks about it,
nobody really knows how to do it,
everyone thinks everyone else is doing it,
so everyone claims they are doing it".

Dan Ariely, Duke University

… we know what we are talking about….

# Backup slides

# New middleware: from Grid to Cloud

- **Grid**:
  - Long, steep learning curve.
  - Difficult to use for real-time analysis, visualization, provisioning of complex virtual environments.
  - Storage management normally at the file/block level, not as distributed objects.

- **Cloud**:
  - Provide new services;  In addition to grid interface
  - Site Virtualisation, for efficiency, service provision, etc
  - Access also to academic infrastructures
  - Possibility to use commercial clouds

# Big Data for Personalized Medicine

## Scientific Objectives

- To collect genomic and clinical data from diseased and healthy patients
- To define best practice for the creation of databases in a standard and exploitable fashion
- To define clear policy for privacy and to create an ethical and transparent program based on consent
- To create reliable correlations between clinical and genomic data
- To exploit genomic and clinical data for patients stratification and clinical trial design

141

# Big Data for Personalized Medicine

**Worldwide Situation**

- Dozens of projects have been launched worldwide from Asia (China and Japan) to the US through Europe

- 100,000 Genomes Project in UK, announced in July 2013 aimed and completing the sequencing of 100,000 genes by end 2017

- Obama's Precision Medicine Initiative launched in 2015 with the objective of sequencing million patients

**Biobanks: constant increase of biological collections:**

- end of 2012 about 600 million pieces in USA [1]

**Biobanks: global Market:**

- Some billion dollars volume [2]
- Technavio report [2] estimates a stable yearly increase of the global market of about 8% in the period 2016–2020, mainly in the nord european area
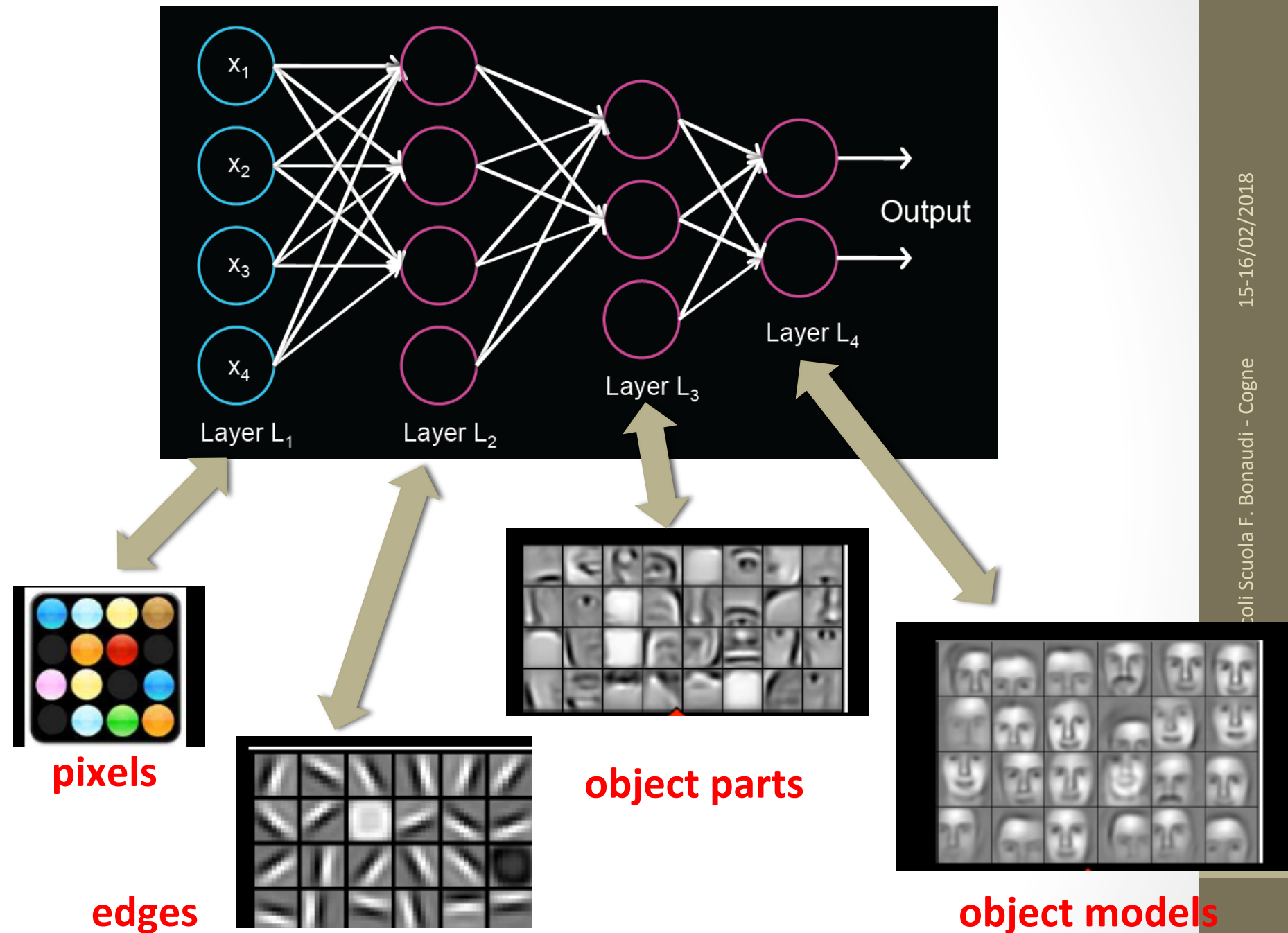
**Evolution of the field due to:**

- new technologies for the samples conservation
- new platform to store and analyse the data
- new efficient methods for data analysis

**The players:**

- public and private companies
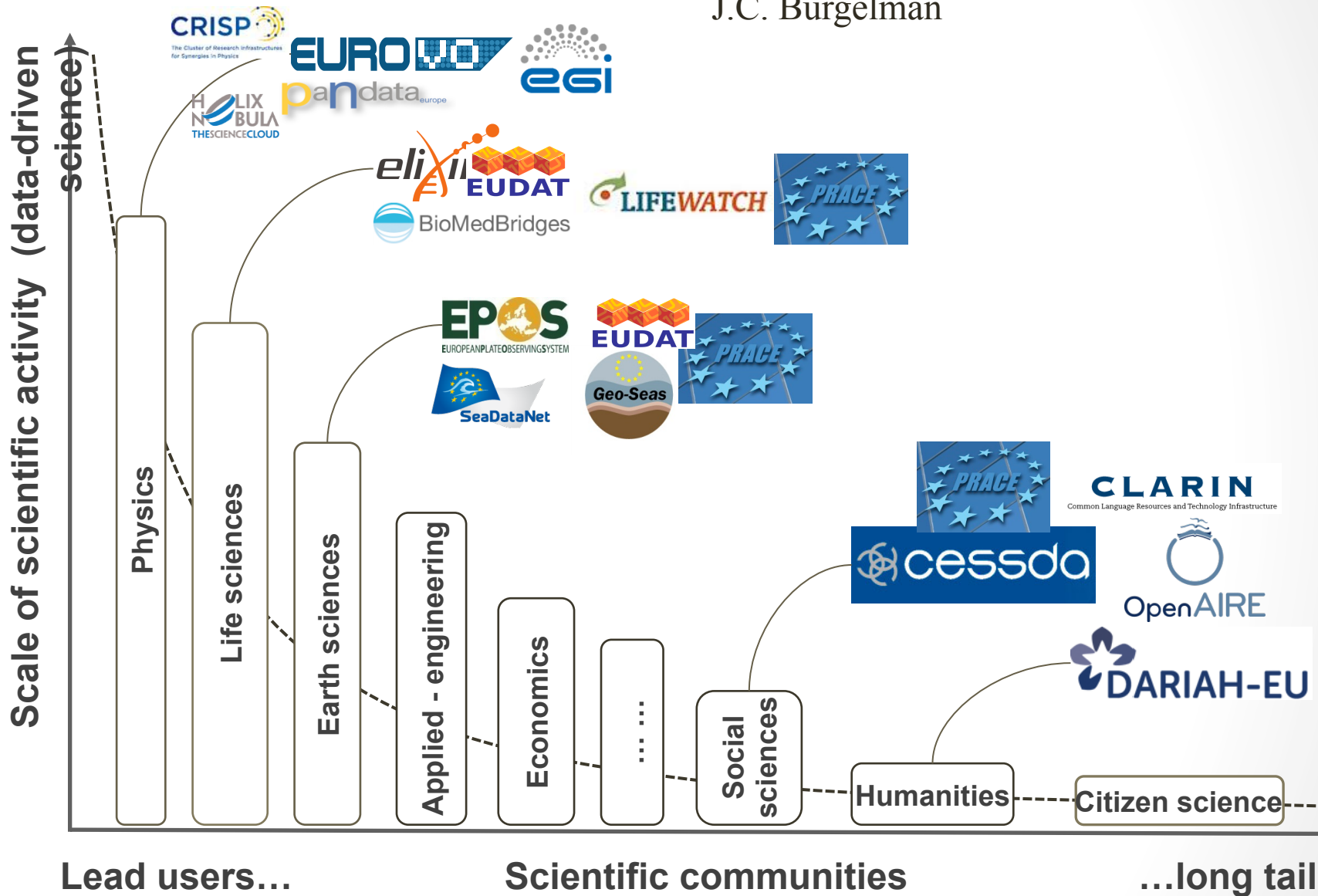- research communities

[1] Monya Baker, "Building better biobanks", Nature 46 (2012), 141–146

[2] Global Biobanking Market 2016-2020, Technavio Division of Cambridge Healthtech Institute, Tec. Rep. IRTNTR7852 (Nov. 2015), 1–145.

A. Zoccoli Scuola F. Bonaudi - Cogne

Output

Layer $L_1$

Layer $L_2$

Layer $L_3$

Layer $L_4$

pixels

edges

object parts

object models

# European Computing Landscape

J.C. Burgelman

# European Computing e-infrastructures

**EGI** is a federated e-Infrastructure set up to provide advanced computing services for research and innovation.
The EGI e-infrastructure is publicly-funded and comprises over 300 data centres and cloud providers spread across Europe and worldwide. The federation is coordinated by the EGI Foundation (also known as EGI.eu)
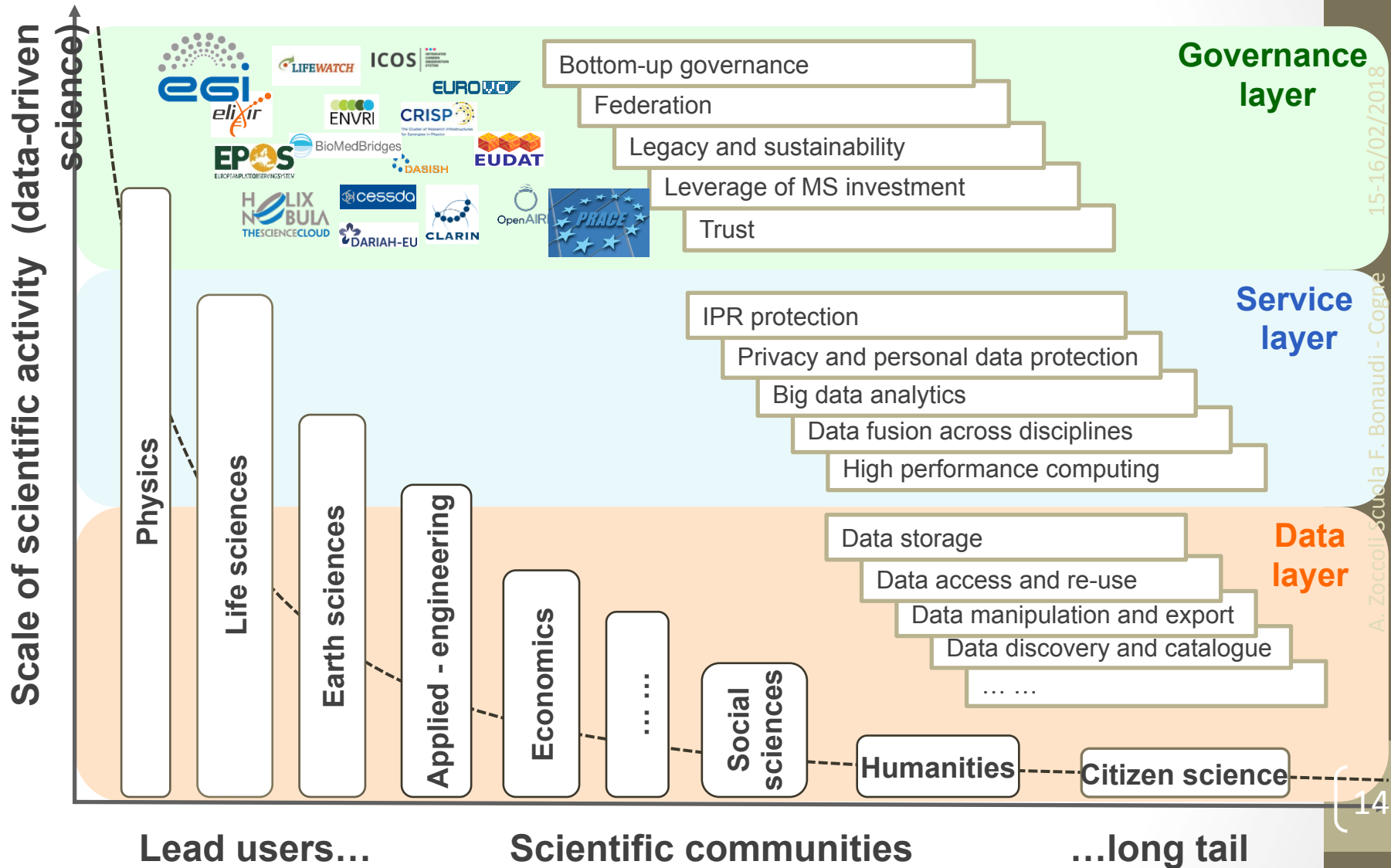
**PRACE** Partnership for Advanced Computing in Europe, an international not-for-profit association with 24 member countries. A pan-European supercomputing infrastructure, providing access to computing and data management resources and services for large-scale scientific and engineering applications at the highest performance level.

**EUDAT** Collaborative Data Infrastructure consists of a European e-infrastructure of integrated data services and resources to support research.

# European Computing Landscape Evolution

# European Cloud Initiative Business Model