



# Trends in computing for the INFN experiments

*Claudio Grandi*  
*INFN Bologna*





# *Computing in INFN*







Claudio Grandi

WWW.GARR.IT



onale di Fisica Nucleare



# Activities

*Management of all INFN Operational Units (Sections, Laboratories, Centres)*

- *User support (network, e-mail, administration, ...)*

*Management of resources for scientific computing*

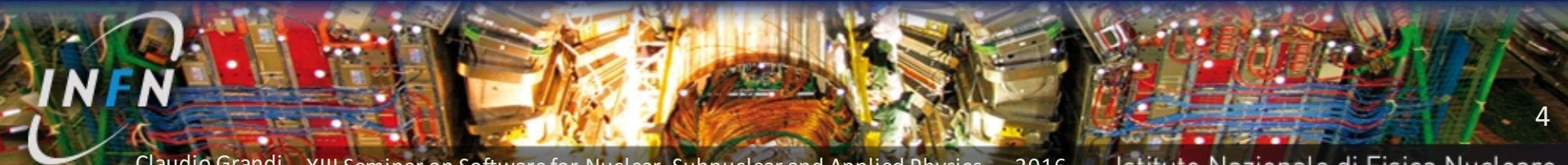
- *Management of **computing centres** and of the distributed infrastructure*

*Development of **technologies** useful for the different activities*

- *Distributed computing, data preservation and open access*
- *Data acquisition and real-time, low-latency computing*
- *Data processing*

*Not just batch data analysis but also statistical analysis and visualization*

- *Theoretical computations and High Performance Computing*





# *Computing for the experiments*

## *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)*

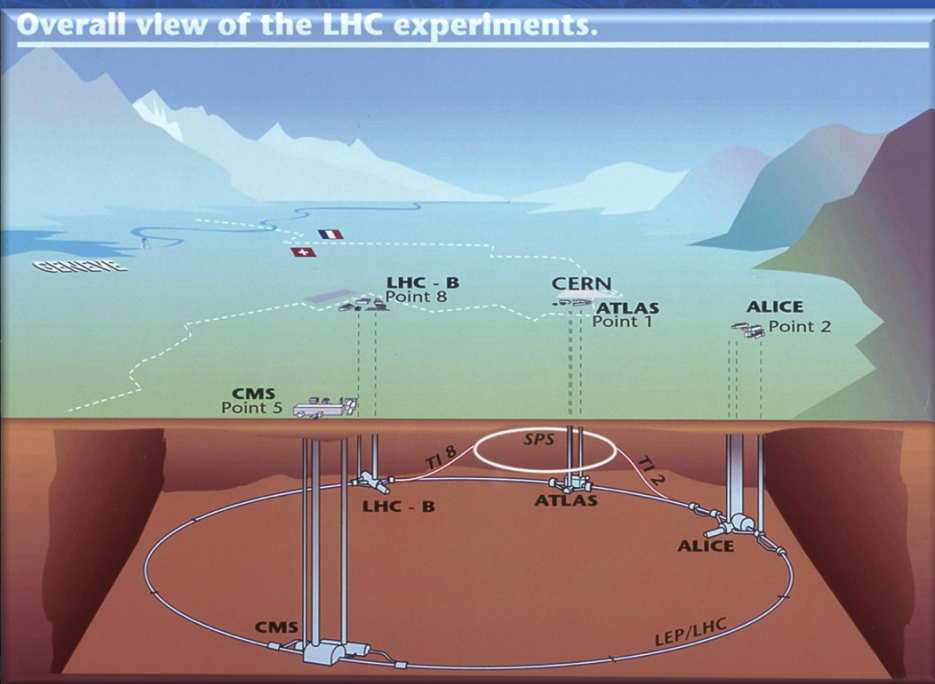


# LHC

*27 km proton-proton collider, ~100 m underground*

*13 TeV c.m. energy,  $10^{34} \text{ cm}^{-2}\text{s}^{-1}$  luminosity*

*40 MHz bunch crossing rate in each of the four experiments*







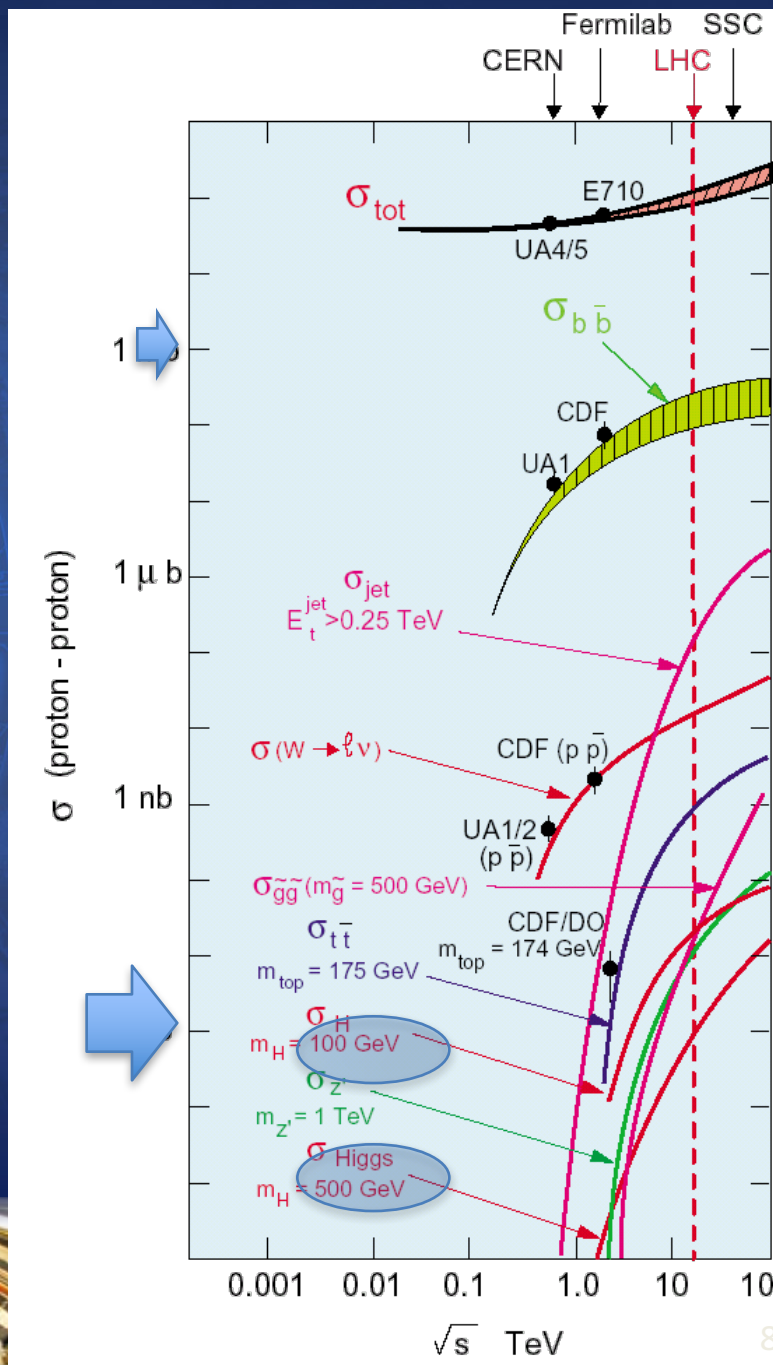
# LHC

Not going to talk about physics and detectors. From the computing point of view:

LHC experiments study rare events!

Signal to noise ratio  $\sim 10^{-13}$

Effective data reduction techniques are needed!





# LHC

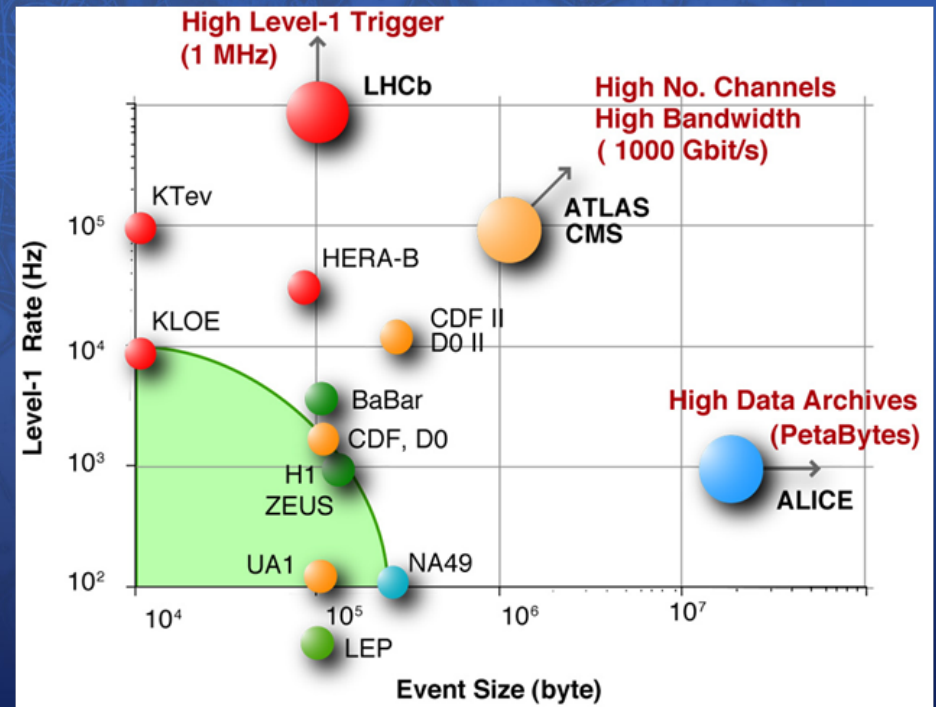
In each LHC experiment there are 40 million bunch crossings per second. Every time 100 million channels are acquired (100 MB)

→ 40,000 EB/y ( $4 \times 10^{22}$  Byte)

Obviously it is not affordable!

The data reduction process brings to 1000 events per second each  $\sim 1$  MB

→  $\sim 10$  PB/y ( $10^{16}$  Byte)





# *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*





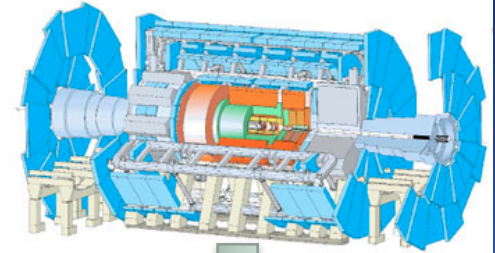


# LHC Real data

LHC collisions

Decay of unstable particles

ATLAS

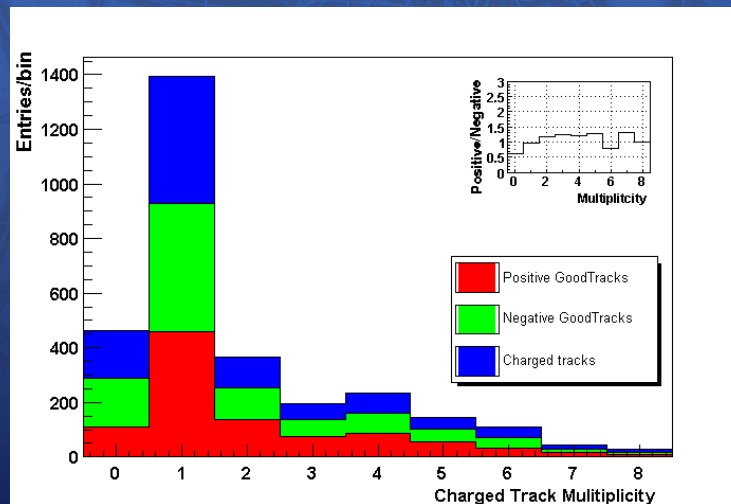


Detector electronics

Trigger

Reconstruction

Analysis



# LHC Simulation

*Not just real data from detectors!*

*Since it is not possible to use analytical solutions of physics 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,  
...

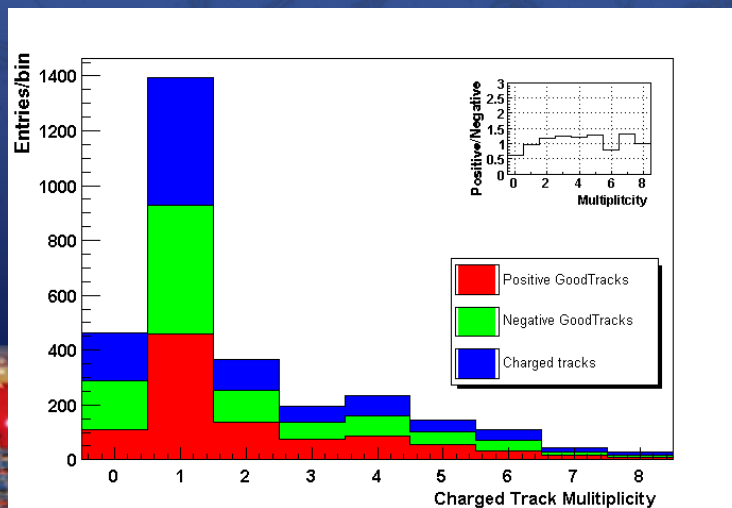
Pythia,  
...

Simulation of detector electronics

Trigger

Reconstruction

Analysis





# Computing infrastructure

Management different kinds of *data* (raw, reconstructed, simulated, analysis products) and of *processes* (different phases of *reconstruction*, *simulation*, end-user *analysis*) 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

*World-wide LHC Computing Grid (WLCG)*





# 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  $n \times 10$  Gb/s to  $n \times 100$  Gb/s

## CPU

Using a unit specific for HEP:  
HepSpec06 (HS06)

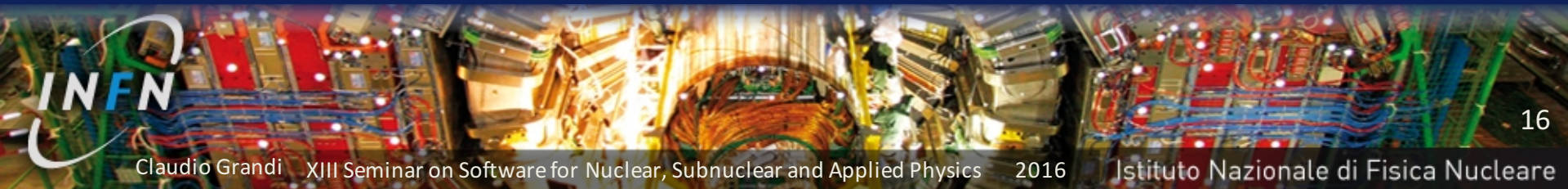
Today:

1 computing core ~ 10 HS06

1 CPU (~12 cores) ~> 100 HS06

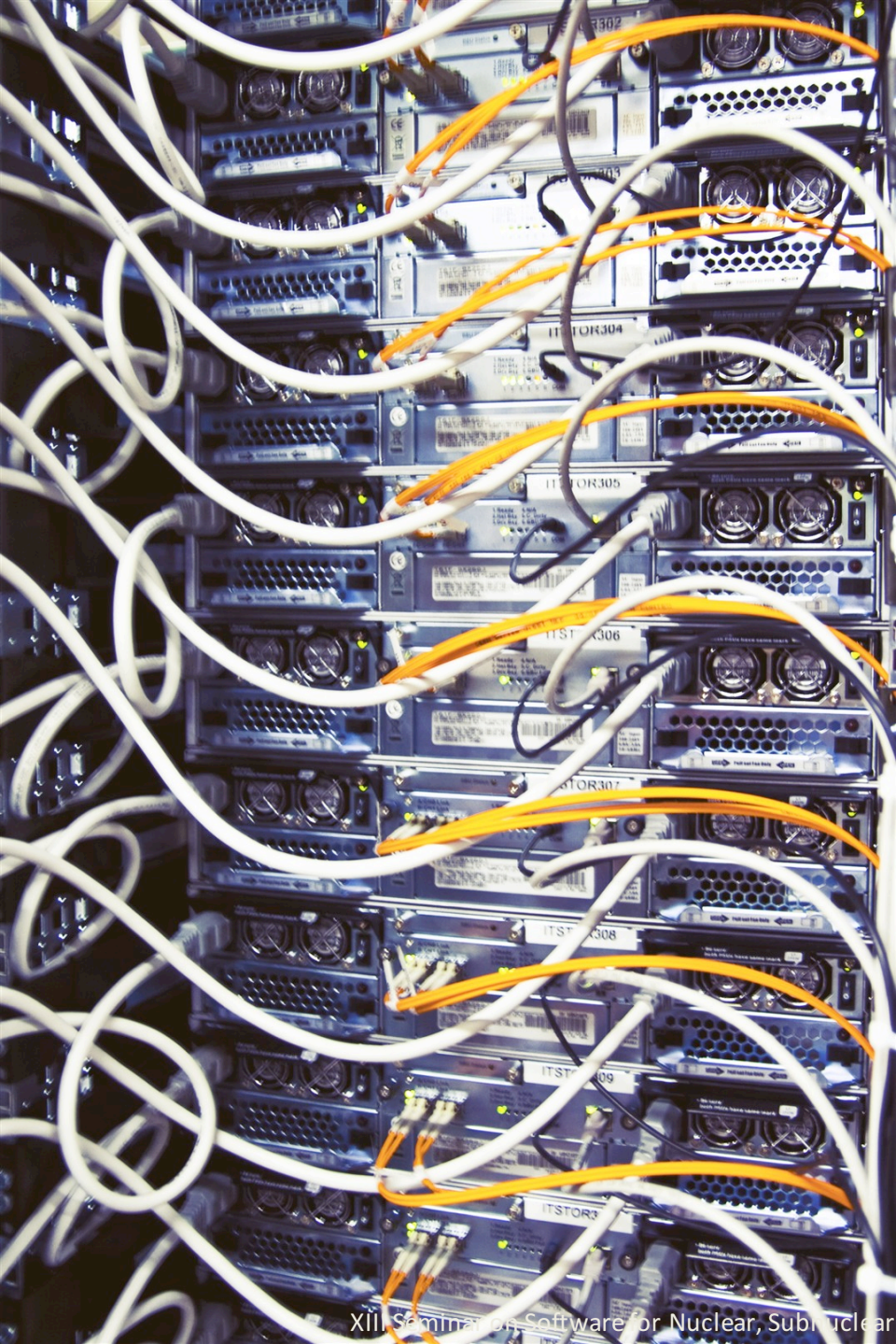


# *Data flow*









# *Numbers from the movie (2013)*

*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*







## *CERN Data Centre (Tier 0)*

*~ 85 000*

*~~73.000~~ 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, end-user analysis*

*> 2*

*~~1.5~~ million jobs running every day*

*50*

*~~10~~ GB/s global transfer rate*







# ...more numbers

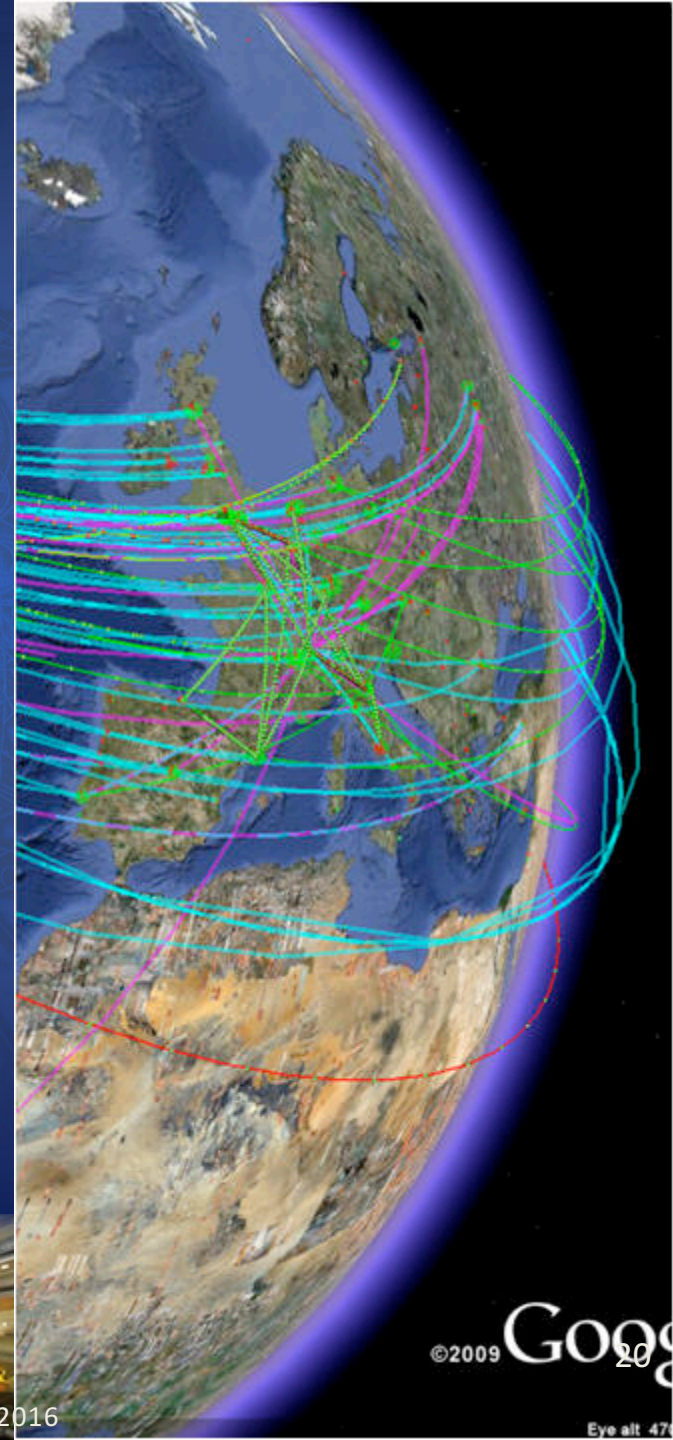
Global resources for 2016 are:

- 2,900,000 HS06 (~290.000 processor cores)
- 240.000 TB disk
- 250.000 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

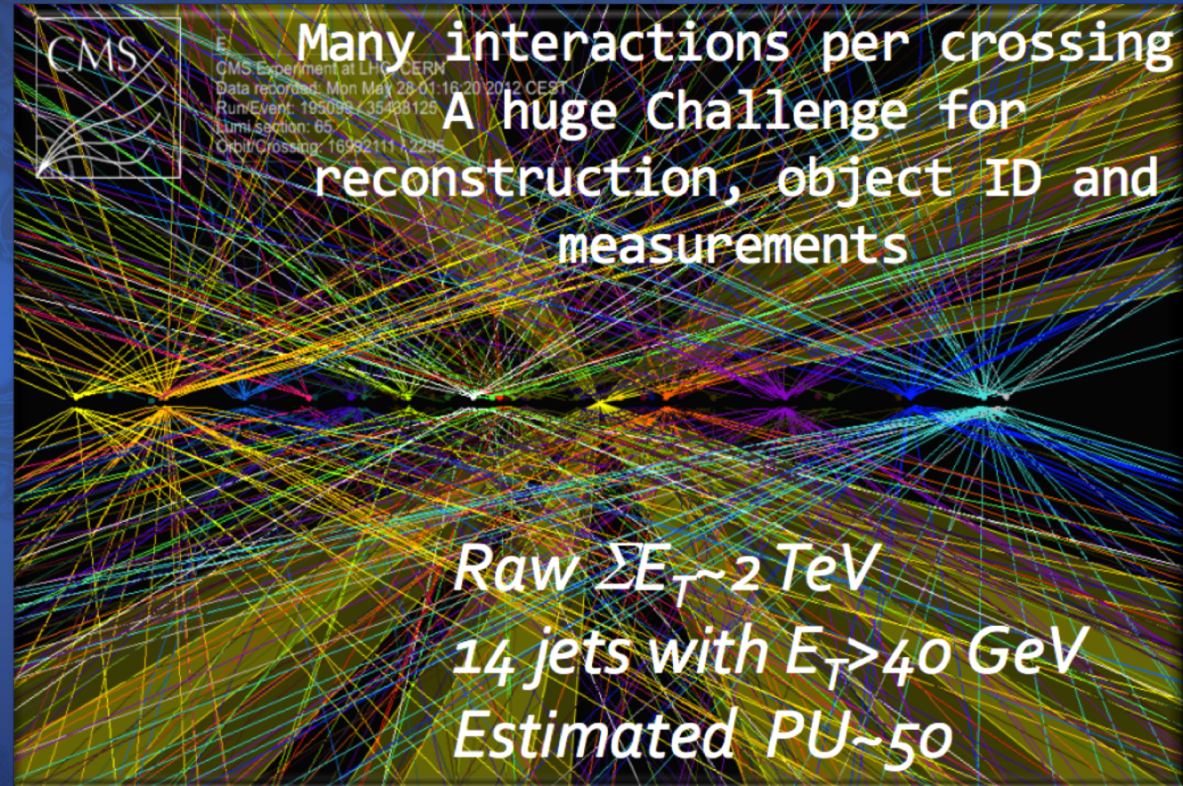




# Pile-up

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*)

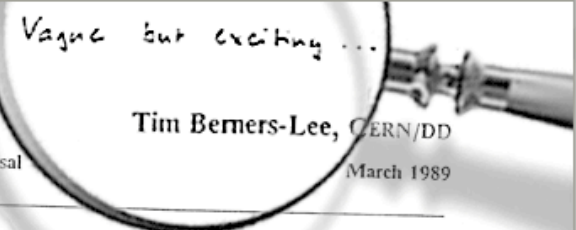


Pileup is increasing to **50** and eventually to more than 150 in HL-LHC



*How?*





CERN DD/OC

Information Management: A Proposal

Tim Berners-Lee, CERN/DD

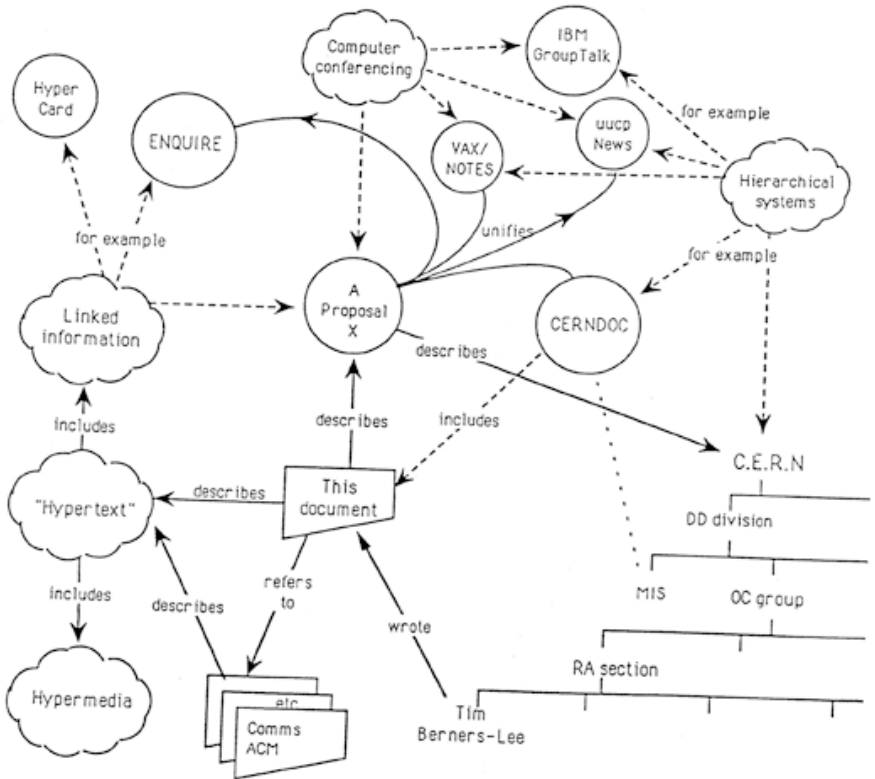
March 1989

# Information Management: A Proposal

## Abstract

This proposal concerns the management of general information about accelerators and experiments at CERN. It discusses the problems of loss of information about complex evolving systems and derives a solution based on a distributed hypertext system.

Keywords: Hypertext, Computer conferencing, Document retrieval, Information management, Project control



# WWW

*In 1989 CERN had needs that were not addressed by existing tools*

*Tim Berners-Lee proposed a mechanism for information sharing in the scientific community: the **World Wide Web***





Today WWW is available to the  
entire society for free!

### Declaration

The following CERN software is hereby put into the public domain:

- W 3 basic ("line-mode") client
- W 3 basic server
- W 3 library of common code.

CERN's intention in this is to further compatibility, common practices, and standards in networking and computer supported collaboration. This does not constitute a precedent to be applied to any other CERN copyright software.

CERN relinquishes all intellectual property rights to this code, both source and binary form and permission is granted for anyone to use, duplicate, modify and redistribute it.

CERN provides absolutely NO WARRANTY OF ANY KIND with respect to this software. The entire risk as to the quality and performance of this software is with the user. IN NO EVENT WILL CERN BE LIABLE TO ANYONE FOR ANY DAMAGES ARISING OUT THE USE OF THIS SOFTWARE, INCLUDING, WITHOUT LIMITATION, DAMAGES RESULTING FROM LOST DATA OR LOST PROFITS, OR FOR ANY SPECIAL, INCIDENTAL OR CONSEQUENTIAL DAMAGES.

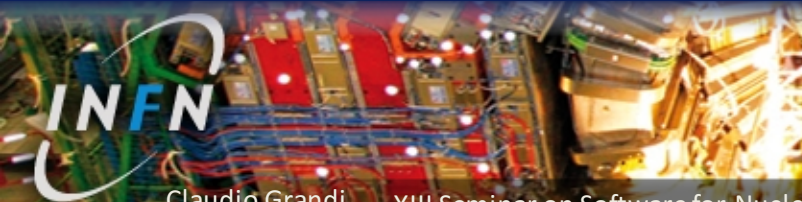
Geneva, 30 April 1993

W. Hoogland  
Director of Research

H. Weber  
Director of Administration

copie certifiée conforme

fait à Genève le 03-05-93





# The first picture on the web (1992)

## Collider

*I gave you a golden ring to show you my love  
You went to stick it in a printed circuit  
To fix a voltage leak in your collector  
You plug my feelings into your detector  
You never spend your nights with me  
You don't go out with other girls either  
You only love your collider  
Your collider.*

*(CERN Hardronic Festival – 1990)*



# *The first web-cam (1993)*



*Computer Laboratory,  
University of Cambridge*





# From Web to *Grid*

In the years 2000s the LHC community had to address the problem of how to manage the data that the experiments would produce

They started from an idea of a group of American computing scientists: the *Computing Grid*

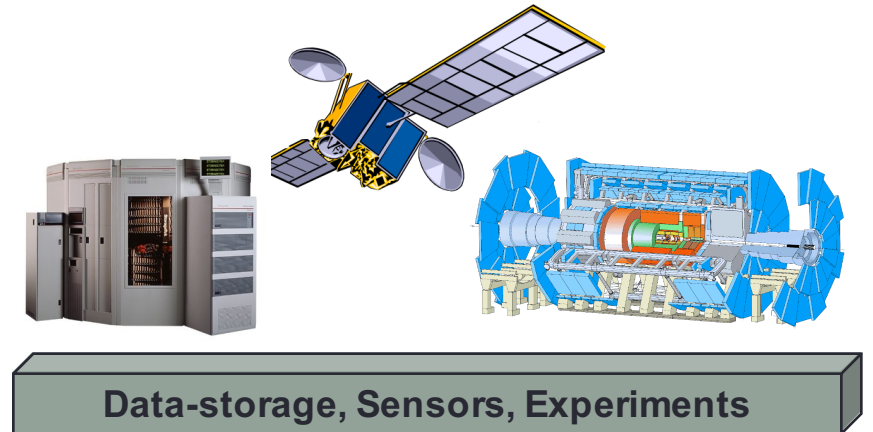
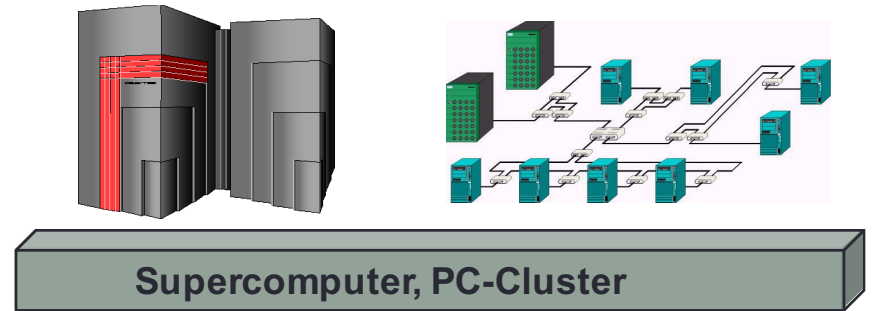
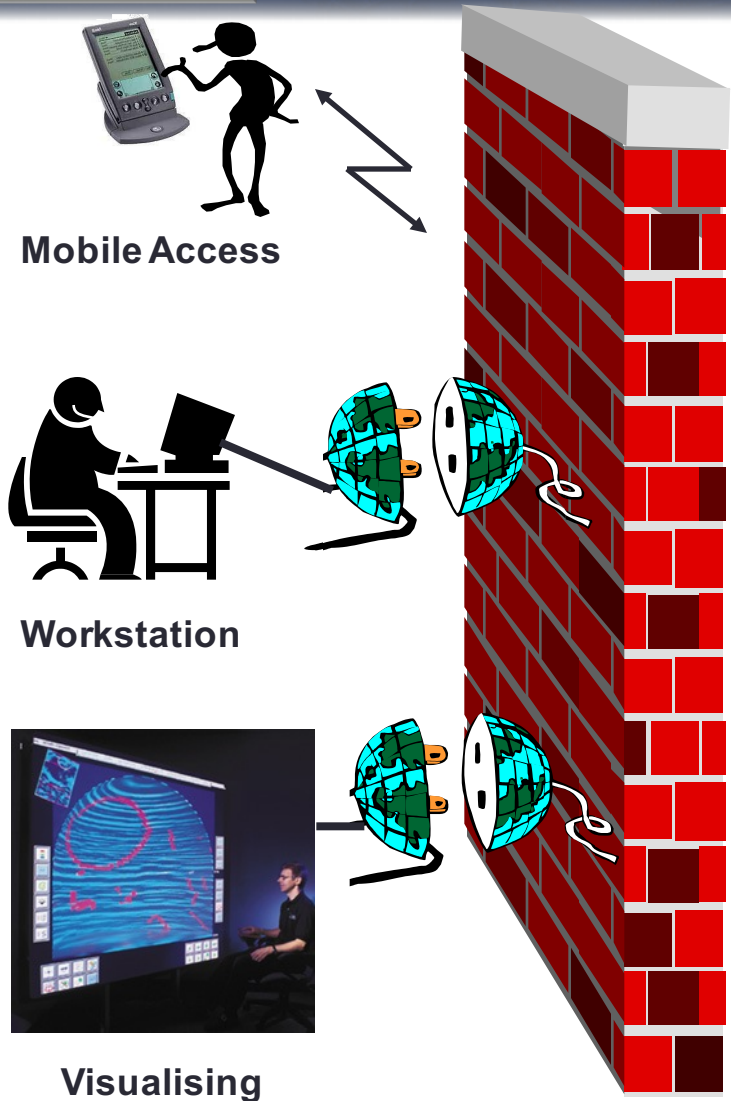
Computing resources are treated in the same way of the electrical power:

A computer is plugged to the network and gets what needed without knowing where it comes from

The *middleware* is a software layer between resources and users



# The Grid metaphore







# *A distributed system*

*Advantages of a **distributed system** (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 **manpower** available in different countries*

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



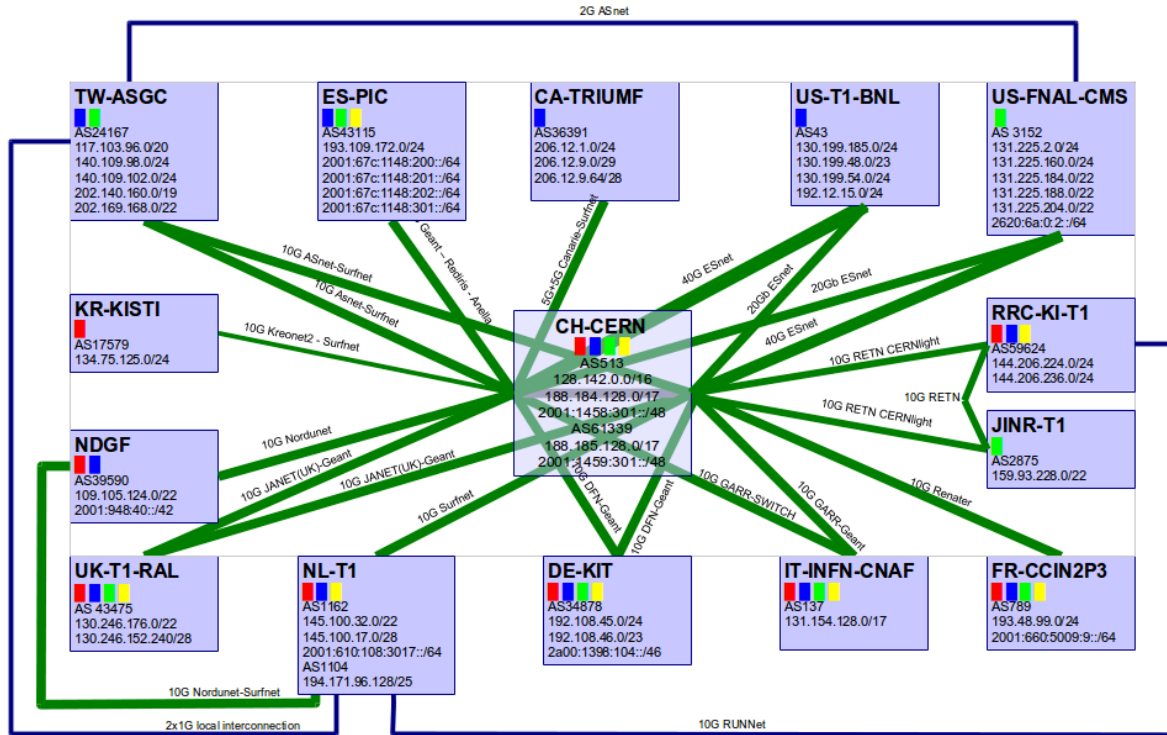
*Only a few technical details...*





# The network - LHCOPN

## LHCOPN



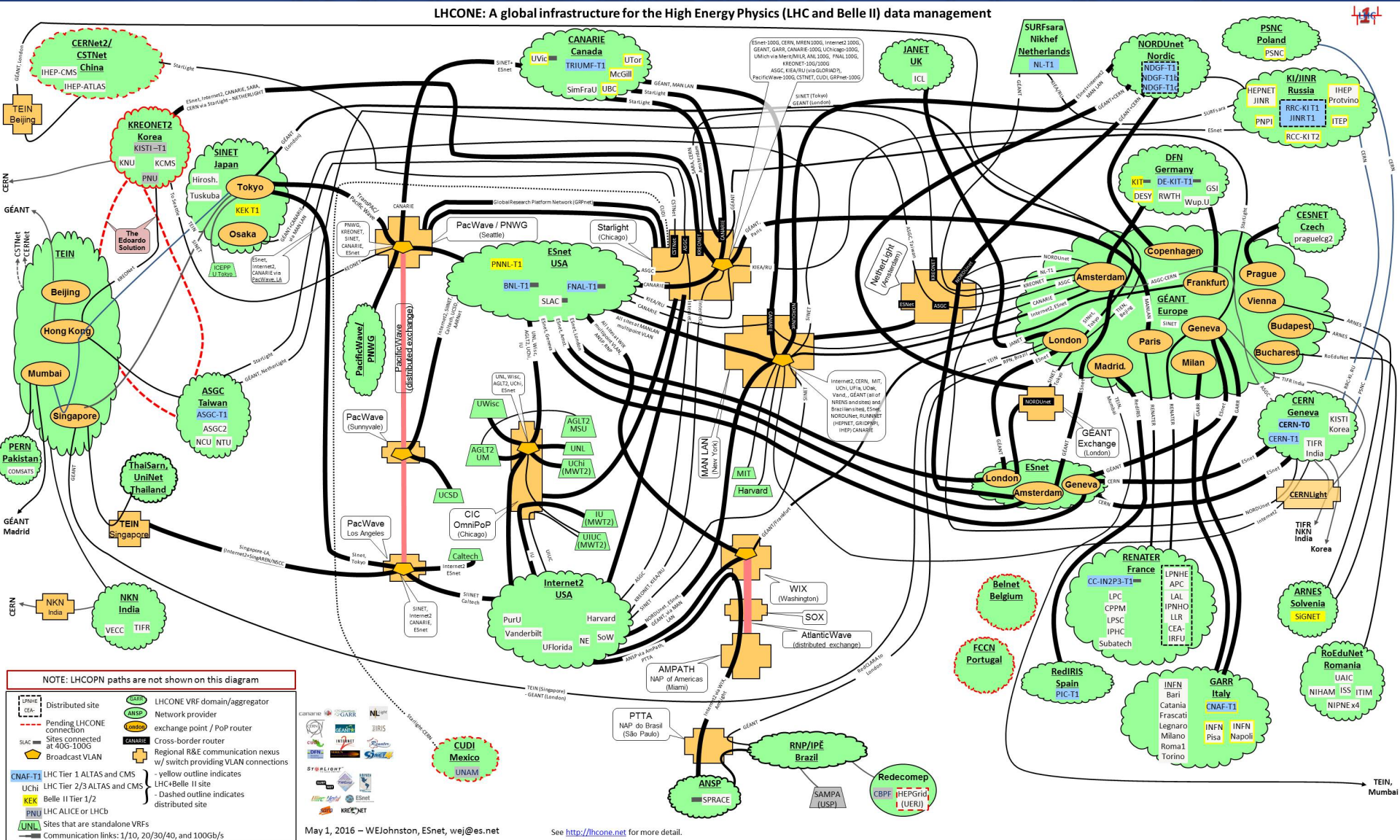
The network technology evolved significantly, offering adequate performance to support the distributed computing model

— T0-T1 and T1-T1 traffic  
— T1-T1 traffic only  
— Not deployed yet  
(thick) >= 10Gbps  
(thin) <10Gbps  
■ = Alice ■ = Atlas  
■ = CMS ■ = LHCB  
 p2p prefix: 192.16.166.0/24 - 2001:1458:302::/48  
 edoardo.martelli@cem.ch 20160322



# The network - LHCONE

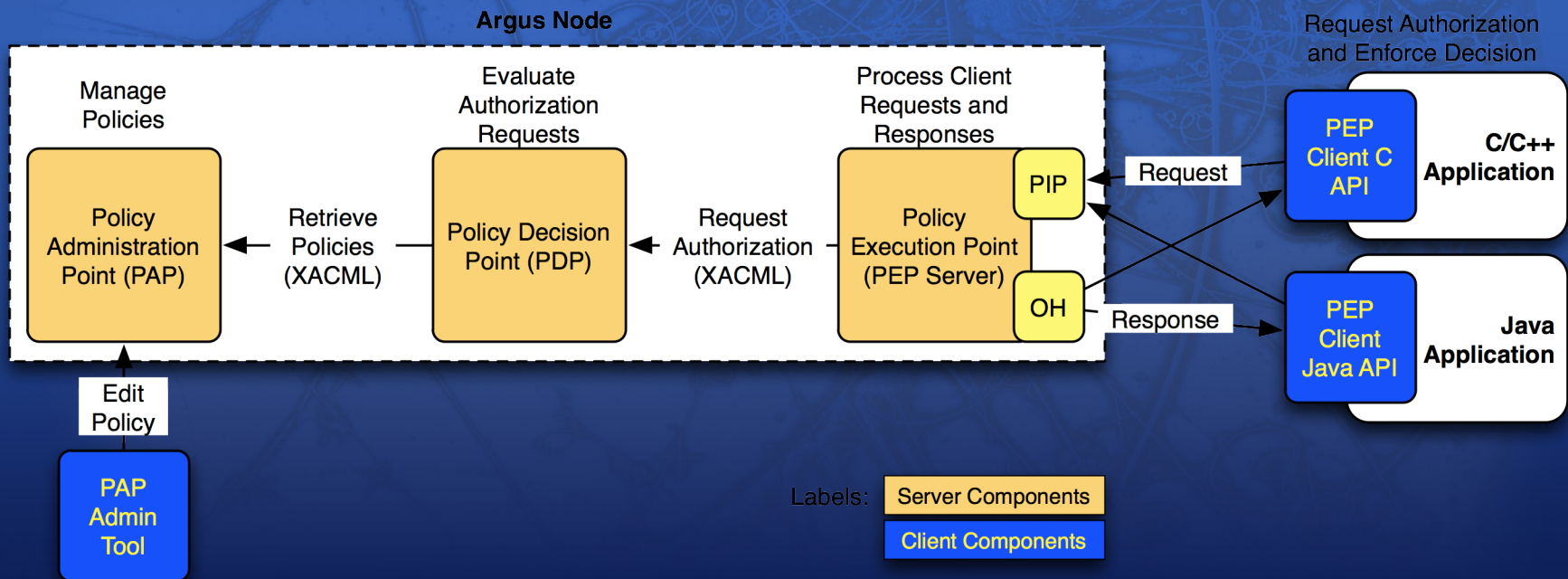
LHCONE: A global infrastructure for the High Energy Physics (LHC and Belle II) data management





# Grid Security management

- **Authentication** based on x.509 certificates
- **Authorization** based on *attribute certificates* (VOMS)
- **Policy management system** (ARGUS)



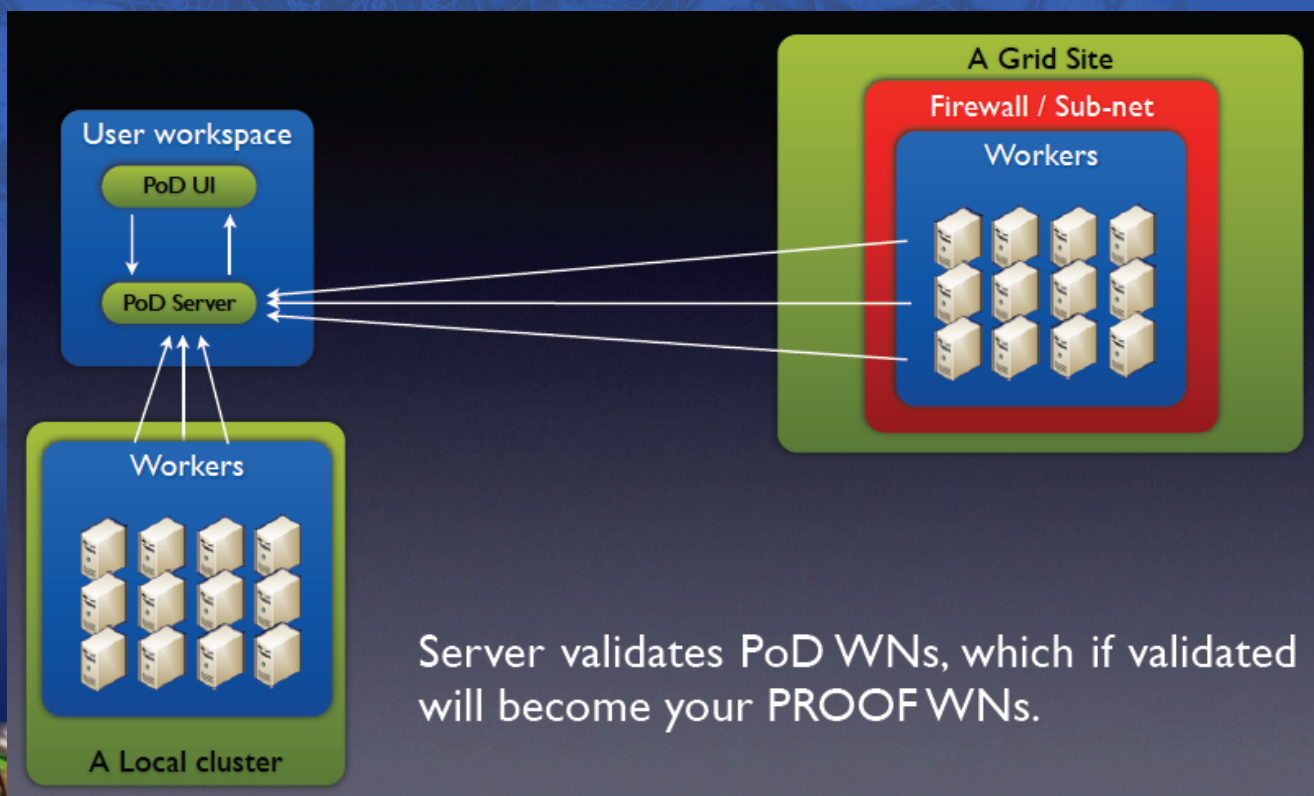
# Grid Computing management

Access is based on *batch jobs*: asynchronous execution

Dedicated interfaces allow to manage remote submissions as if local

*Interactive processing*

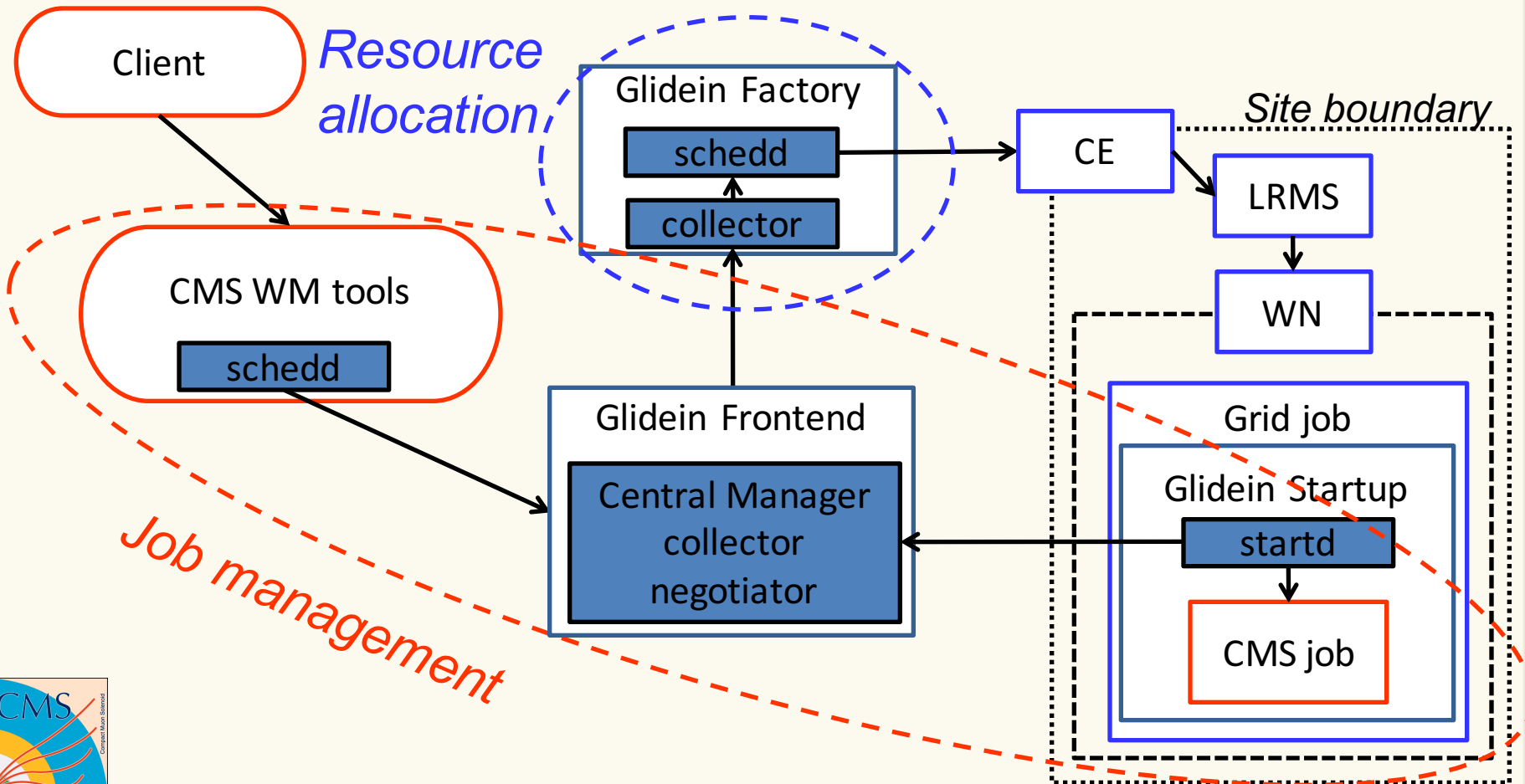
is limited and based on local resources or on systems able to manage part of the load in batch mode (e.g. PoD)





# The "pilot" model

*Separation of resource allocation and job management*





# Grid Data management

Heavily relying on *tape libraries* for persistent data storage

Accessible in a transparent way (nearline)

Dedicated interfaces to uniformly manage data on disk and on tape

Tools to manage the *transfer* of large amounts of data

*Local access* to data by jobs but today network performances allow transparent *remote access* on the Wide Area Network

*Storage Federations*





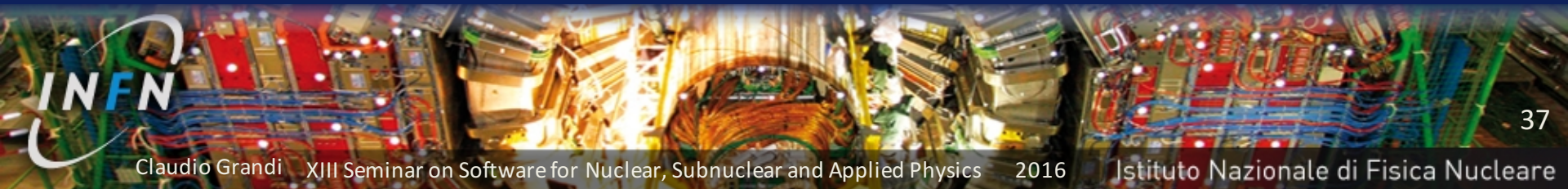
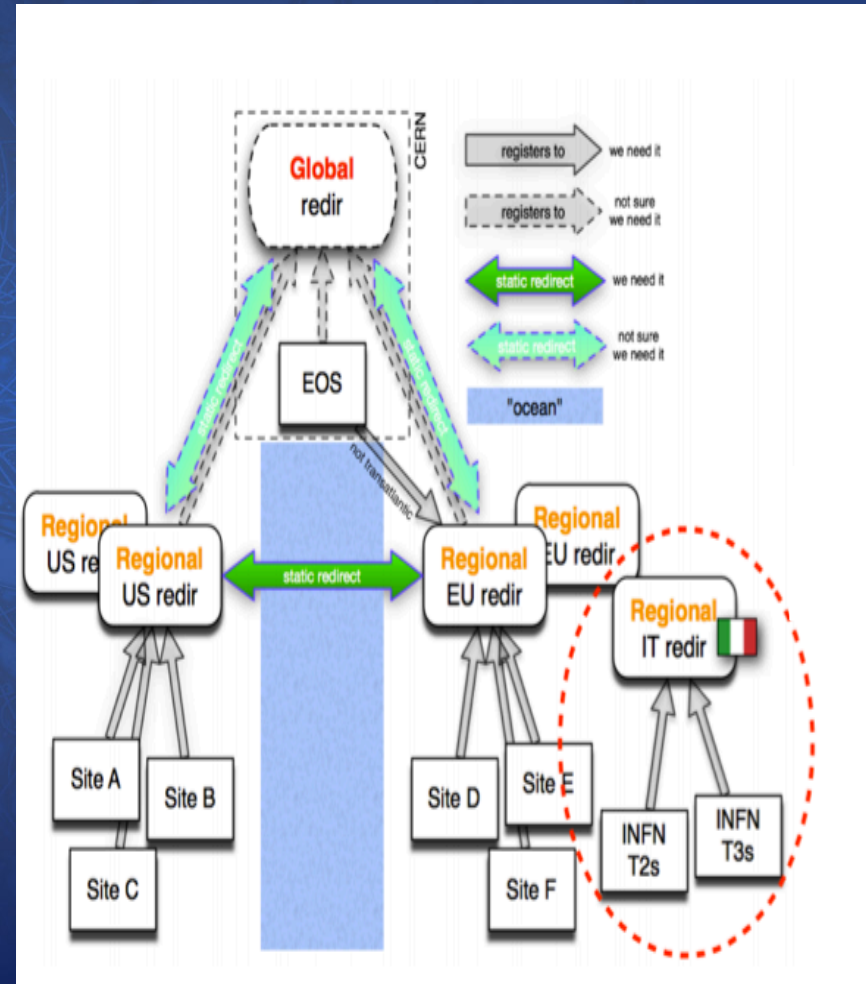
# Storage Federations

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)

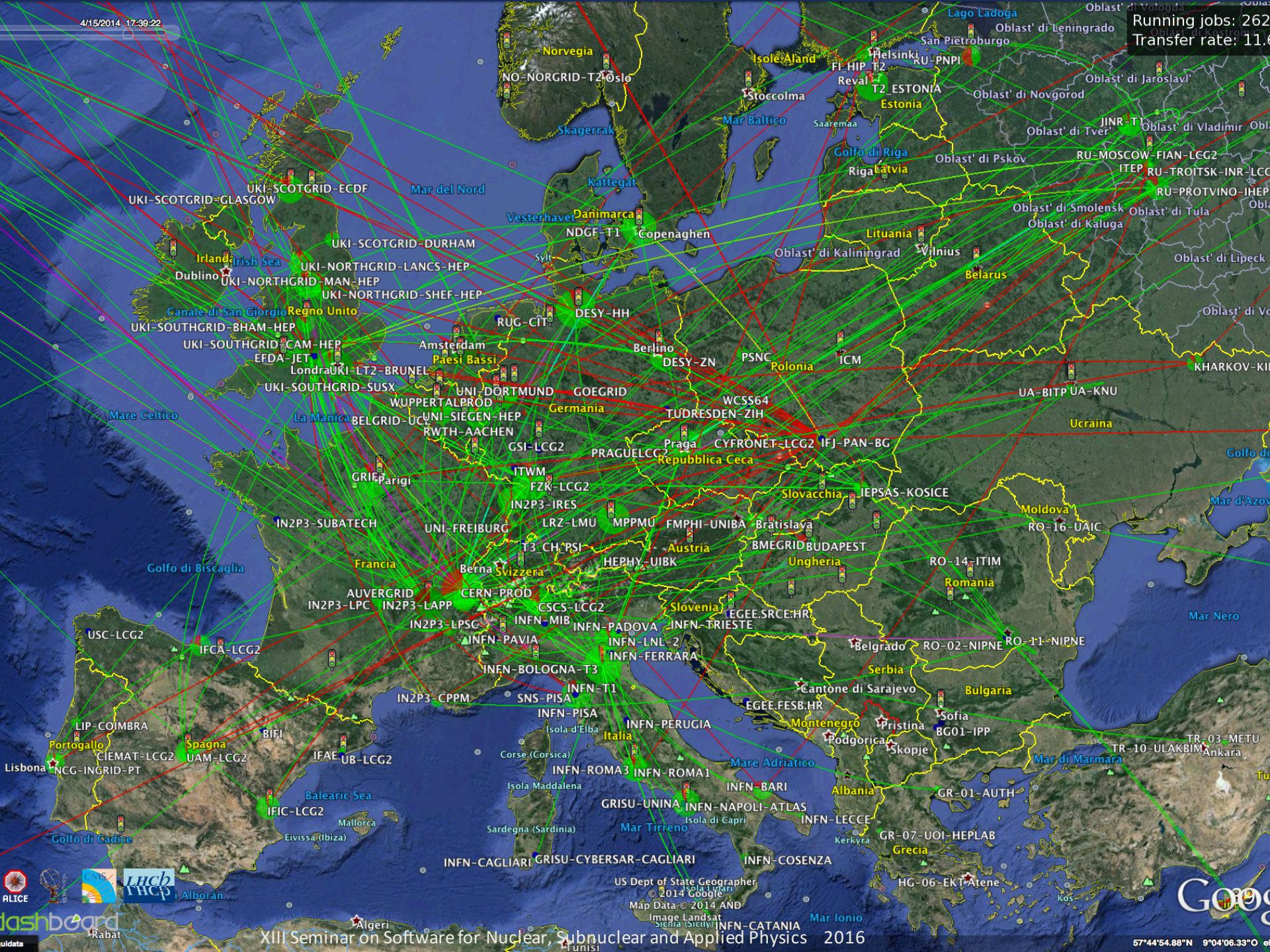
Currently in production for *xrootd*



*Let's see how it works...*











# Grid: an example of collaboration

Even though the HEP community has been dominant, the Grid has been thought and build for the whole **scientific community**

Projects as the European Grid Initiative (**EGI**), to which INFN participates, and the Open Science Grid (**OSG**) in the US provide computing resources to many scientific communities, and more.

Involvement also in the **industrial world**.



Long-tail of research

EGI Case Studies

- <http://www.egi.eu/case-studies>

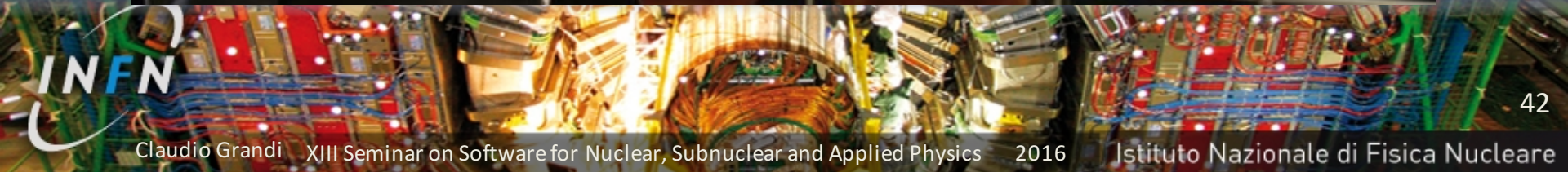
				
<b>Natural Sciences</b> Life Sciences, Earth Sciences, Mathematics, etc	<b>Physical Sciences</b> Physics, Astronomy, Chemistry etc	<b>Medical and Health Sciences</b> Medicine, Clinical sciences, etc	<b>Engineering &amp; technology</b> Material science, civil and mechanical engineering, etc	<b>Agricultural sciences</b> Veterinary sciences, food technology, etc





*Was that enough?*

# July 4<sup>th</sup> 2012





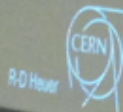
Global Effort → Global Success

Results today only possible due to  
extraordinary performance of  
accelerators - experiments - Grid computing

Observation of a new particle consistent with  
a Higgs Boson (but which one...?)

Historic Milestone but only the beginning

Global Implications for the future



[ credits: D.Bonacorsi ]



# *What about the years to come?*





# LHC roadmap

2015				2016				2017				2018				2019				2020				2021			
Q1	Q2	Q3	Q4	Q1	Q2	Q3	Q4	Q1	Q2	Q3	Q4	Q1	Q2	Q3	Q4	Q1	Q2	Q3	Q4	Q1	Q2	Q3	Q4	Q1	Q2	Q3	Q4



2022				2023				2024				2025				2026				2027				2028			
Q1	Q2	Q3	Q4	Q1	Q2	Q3	Q4	Q1	Q2	Q3	Q4	Q1	Q2	Q3	Q4	Q1	Q2	Q3	Q4	Q1	Q2	Q3	Q4	Q1	Q2	Q3	Q4

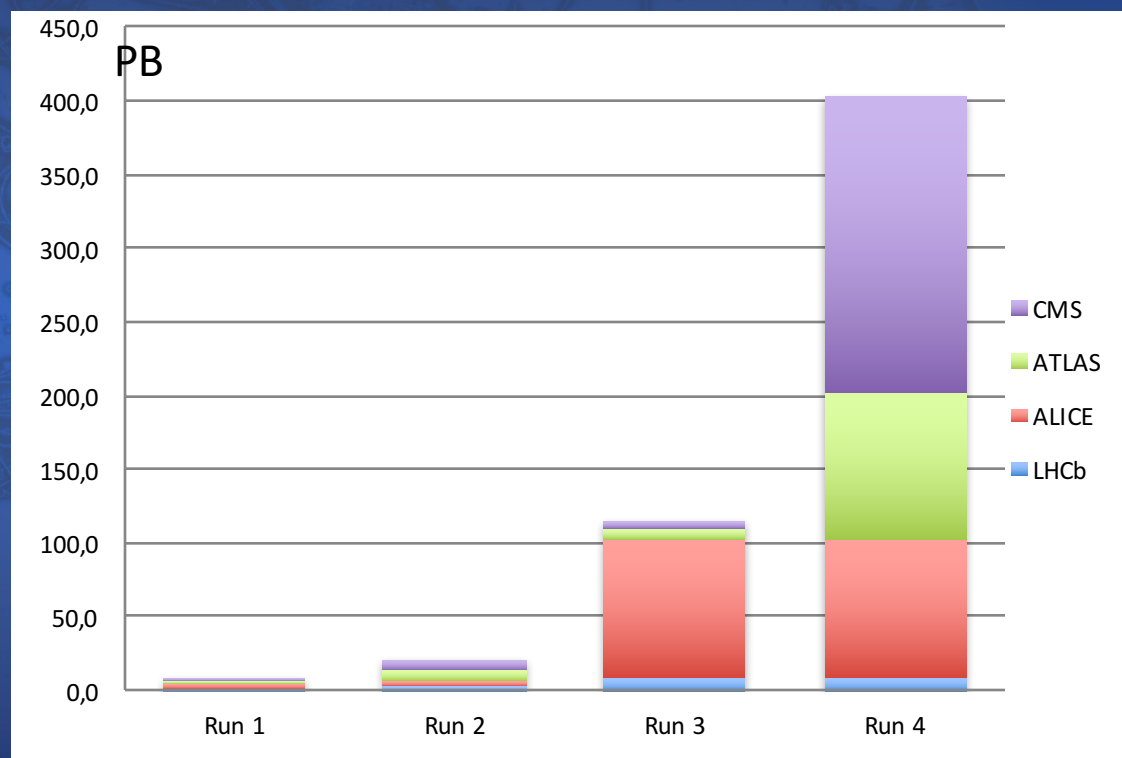


2029				2030				2031				2032				2033				2034				2035			
Q1	Q2	Q3	Q4	Q1	Q2	Q3	Q4	Q1	Q2	Q3	Q4	Q1	Q2	Q3	Q4	Q1	Q2	Q3	Q4	Q1	Q2	Q3	Q4	Q1	Q2	Q3	Q4



# Resource requests for the future

Significant increase  
in experiments'  
requests in the  
coming years: the  
“scary plot”



...but the buzz-word is “flat-budget”!







# Foreseen evolution – LHC Run 3

## ATLAS and CMS

*Trigger rate is constant*

*50% increase in pile-up and luminosity → integrated luminosity doubles*

## ALICE

*DAQ rate in 50 kHz → 1 Tb/s...*

*...but data reduction of a factor of 20 on the O<sup>2</sup> farm*

## LHCb

*Software trigger only (30 MHz) → 2-5 GB/s to offline*

*In addition the Cherenkov Telescope Array experiment starts!*





# Italian resources in 2016

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

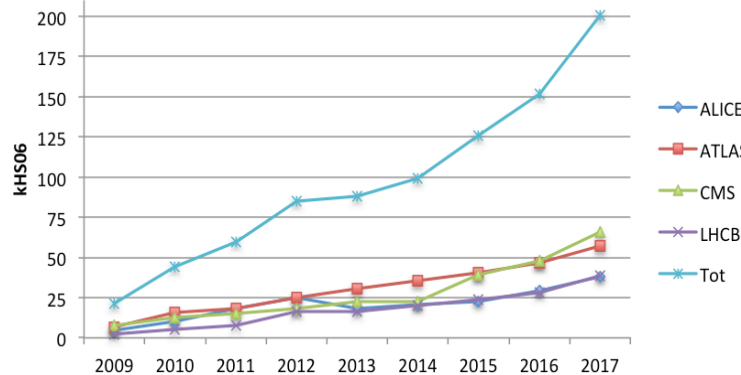
	<i>CPU (kHS06)</i>	<i>Disk (PB)</i>	<i>Tape (PB)</i>
<i>WLCG</i>	<i>2900</i>	<i>240</i>	<i>250</i>
<i>INFN</i>	<i>306</i>	<i>30</i>	<i>35</i>
<i>% INFN</i>	<i>11</i>	<i>12</i>	<i>14</i>





# CNAF evolution - LHC Run 1 & 2

### CPU - CNAF



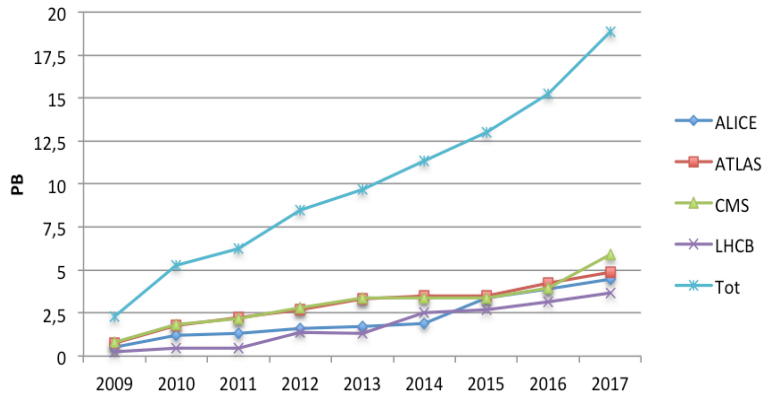
*Run2 is ok with the flat budget hypothesis:*

*CPU + 20 - 30%*

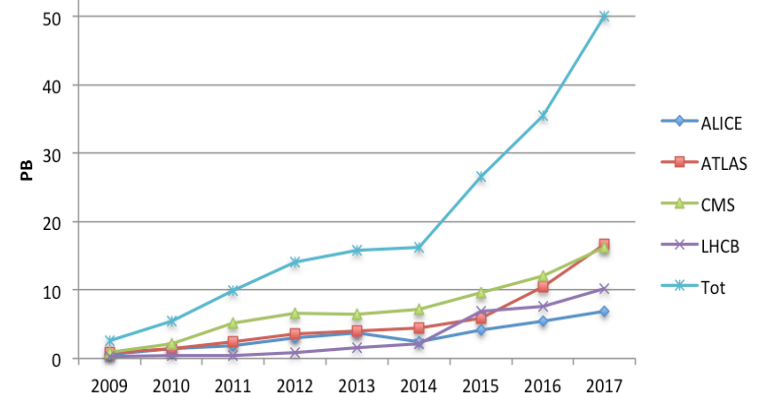
*Disk + 15 - 25%*

*Tape + 30% - 60%*

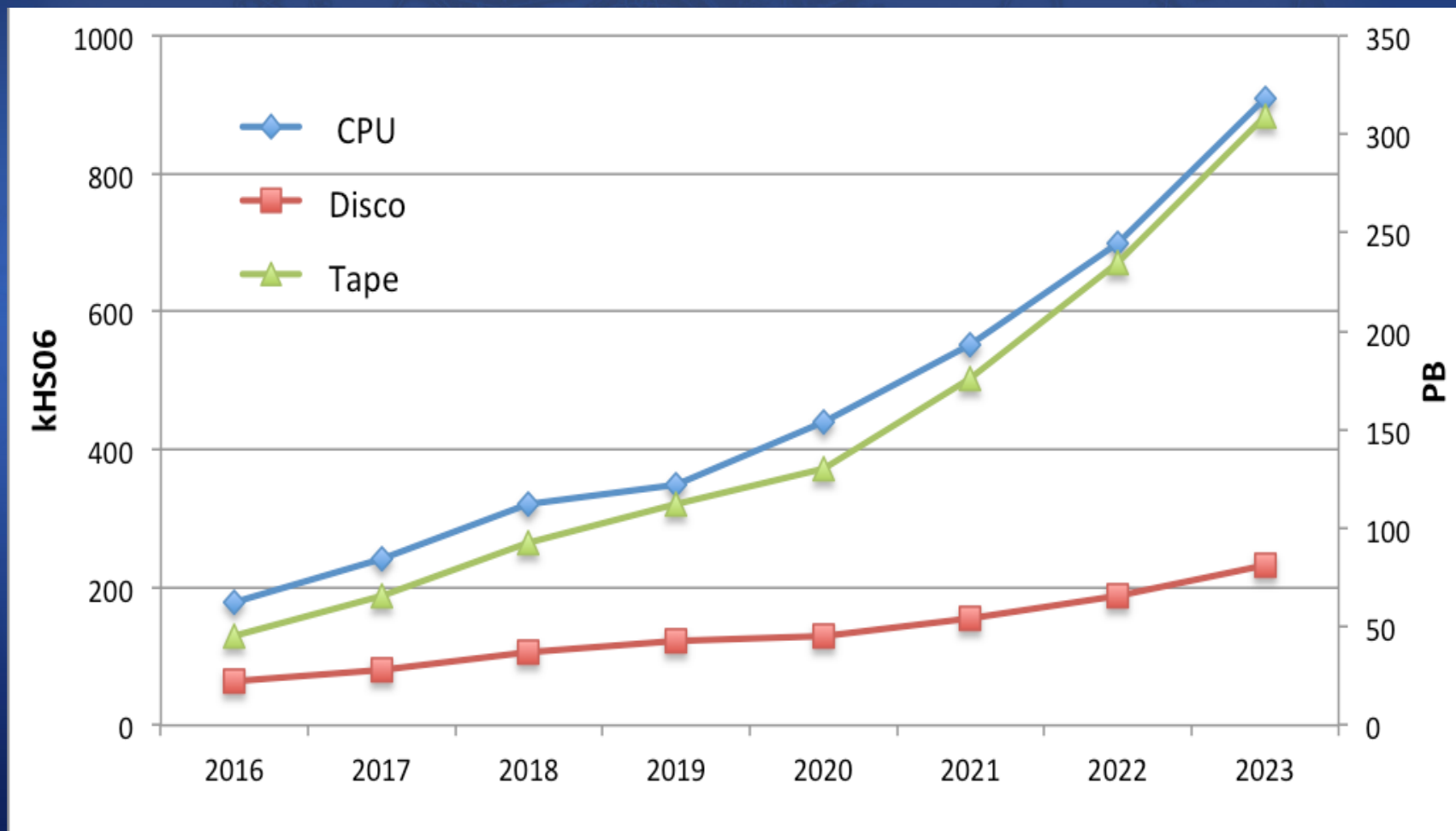
### Disco - CNAF



### Tape - CNAF

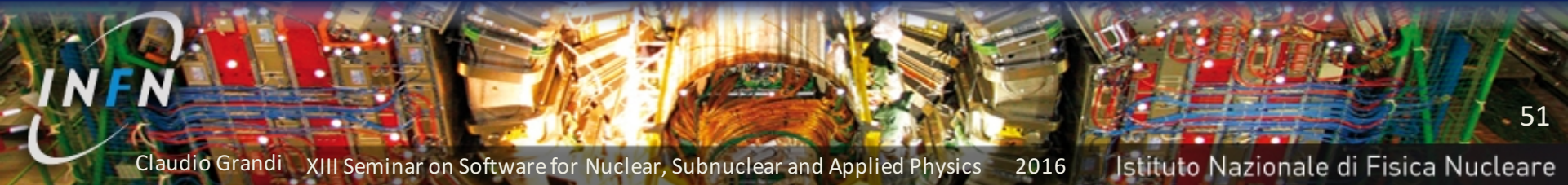


# *CNAF evolution up to LHC Run 3*





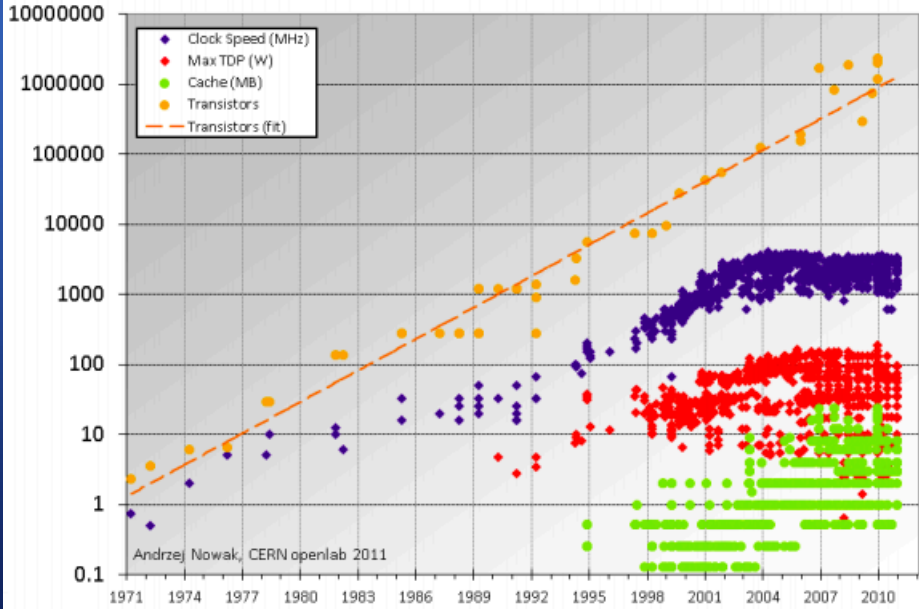
# *Does the technological evolution help?*



# CPU power

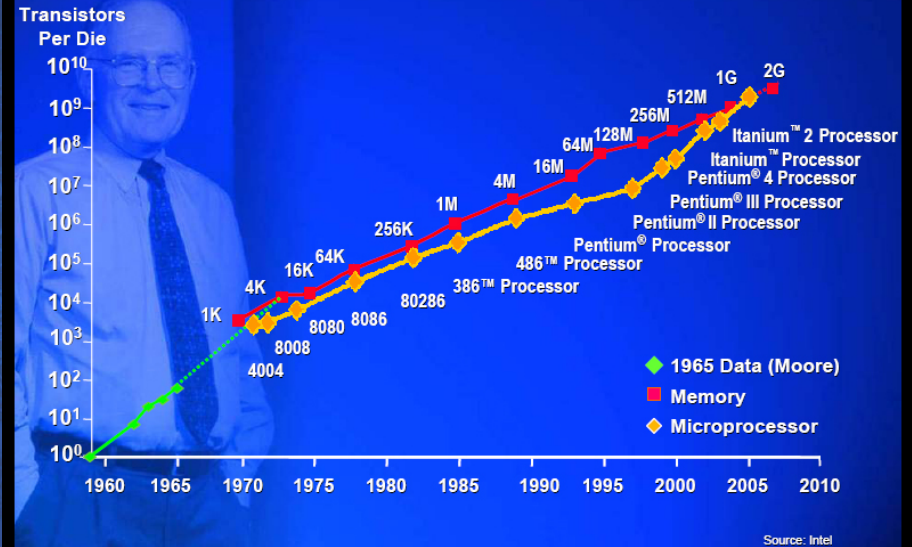
*Moore's law (CPU performance doubles every 18 months at the same cost) does not hold any more*

Intel Processor features



Source: Andrzej Nowak – CERN OpenLab

## Moore's Law - 2005



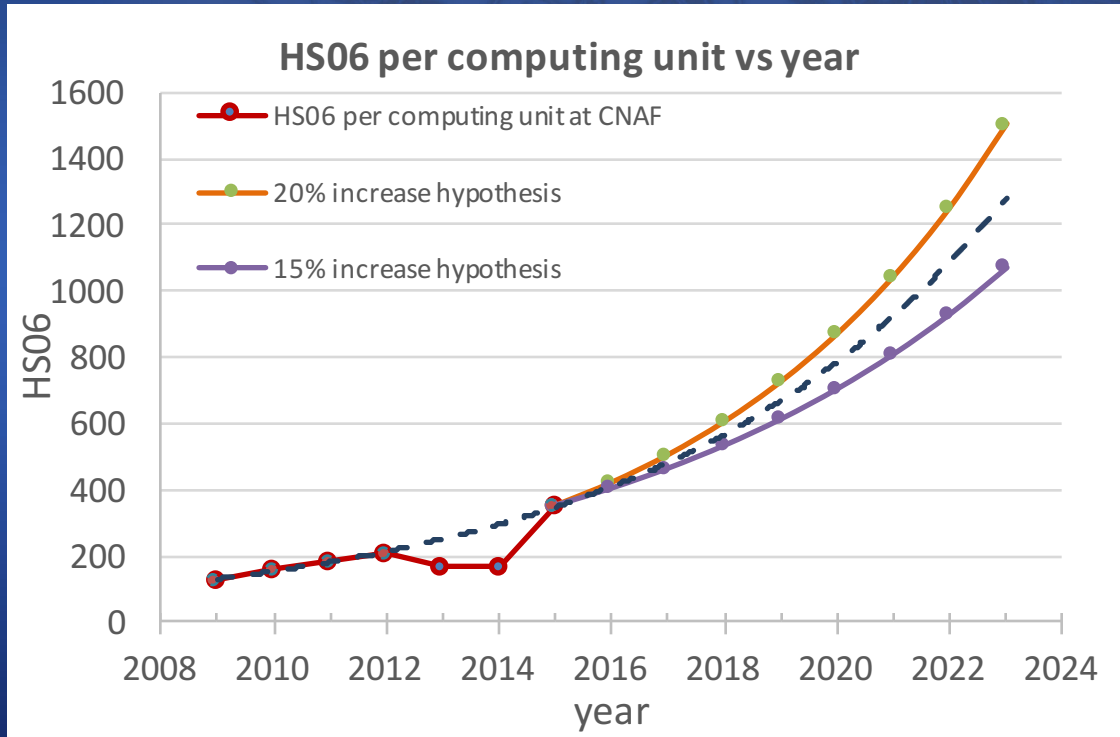
Source: Intel

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





# CPU power

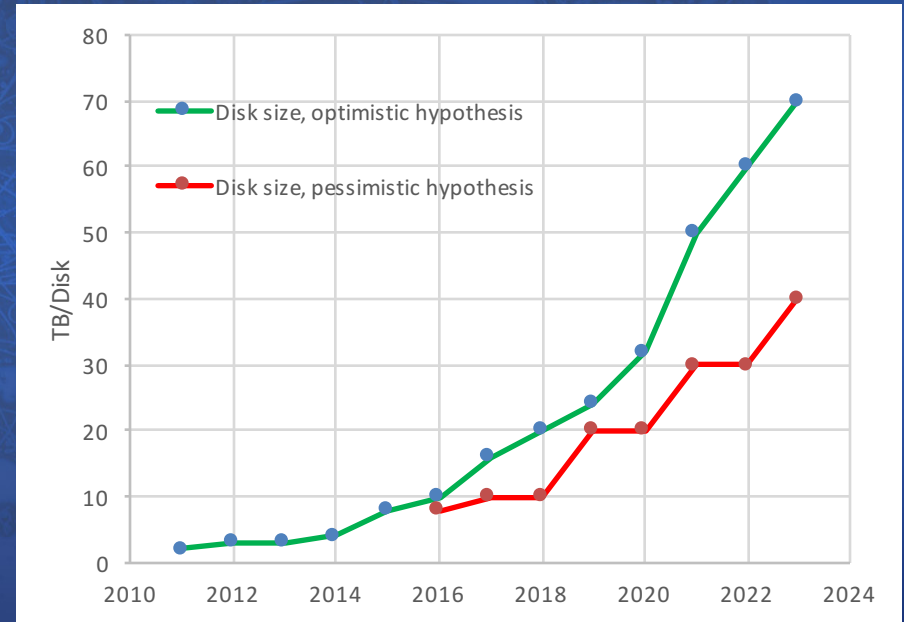
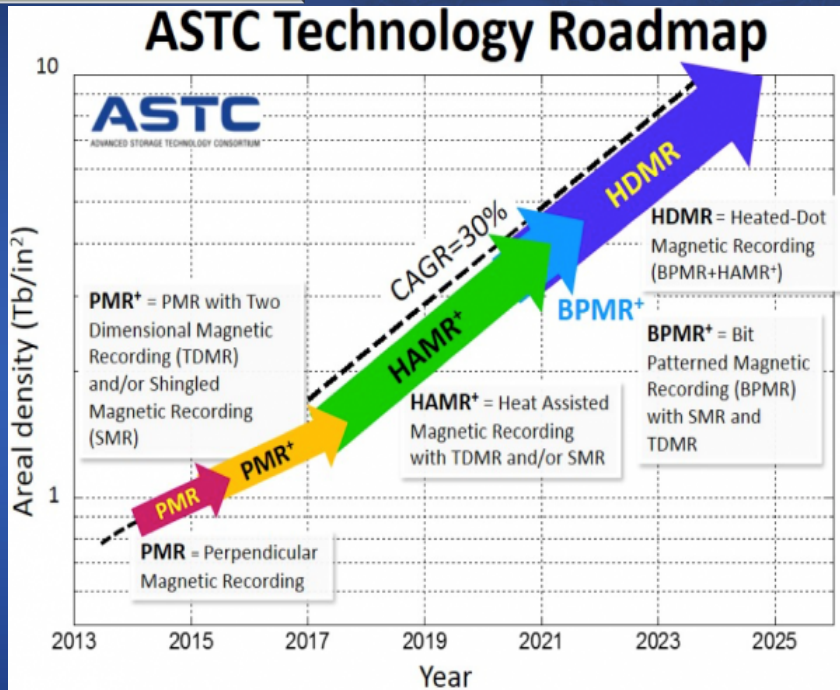


*Starting from the actual power of the nodes bought by CNAF in 2009-2015 we estimate an increase between 15 and 20%*



# Disk

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



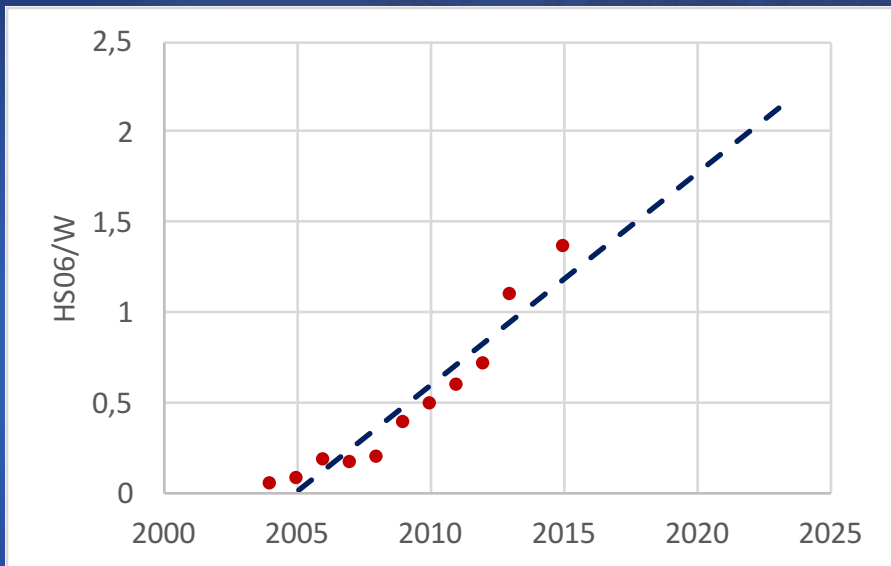
*Extrapolation is more difficult for disk because there are technology changes foreseen*

*The number of disks may not need to increase*



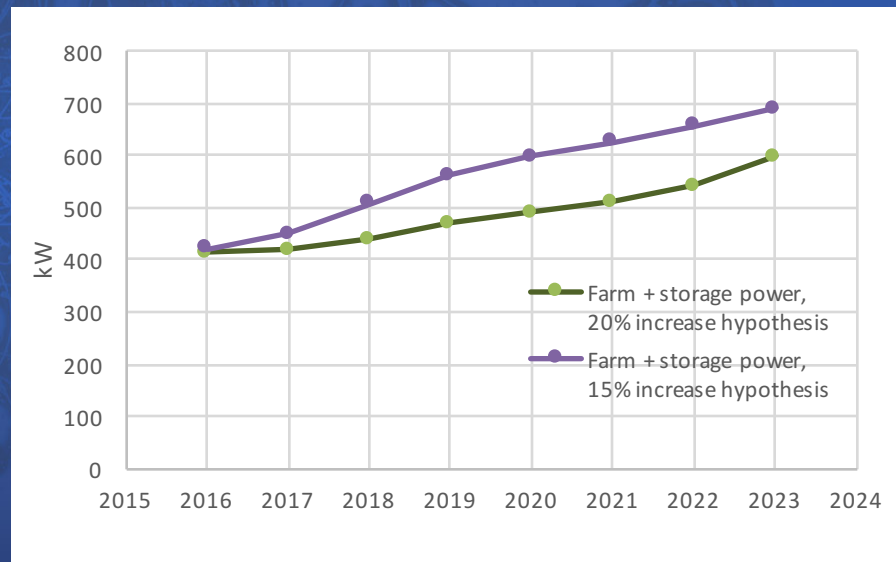


# Electrical power



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

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



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



# 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*





*Can't just follow the evolution of  
currently used technologies!*

# Are we different from the rest of the world?

SKA1-Low  
660 PB/yr

Business emails sent  
3000 PB/yr  
(Unstructured content)

YouTube  
15 PB/yr

Kaiser  
Permanente  
30 PB/yr

LHC Phase 1  
Raw data  
100 PB/yr

LHC Phase 2  
Raw data  
400 PB/yr

LHC  
data  
15 PB/yr

Google  
Search  
100 PB/yr

Facebook uploads  
180 PB/yr

Circle  
DK





*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 offers most of the functionalities needed by HEP computing*

*Commercial and industrial world*

*offers solutions that are being integrated*

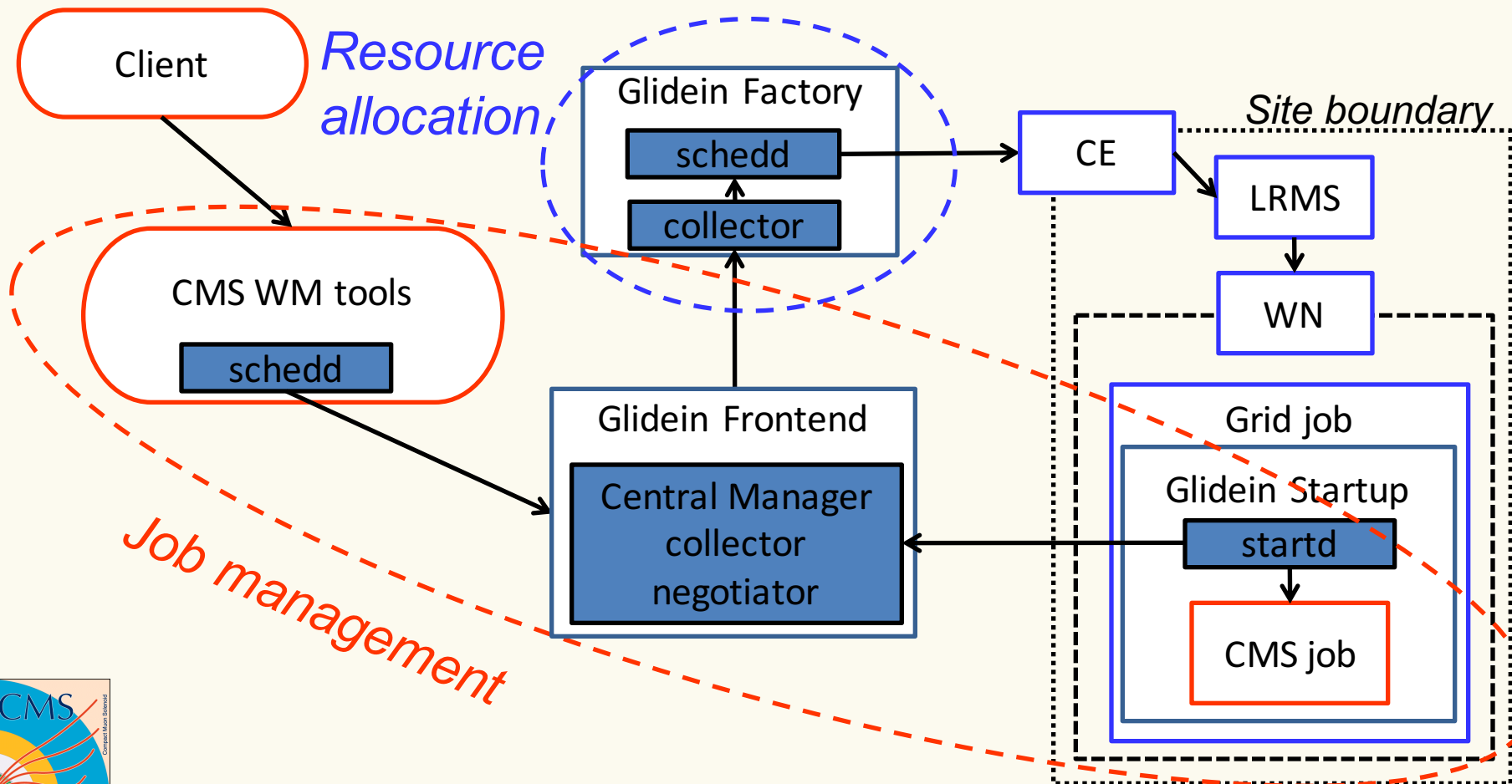
*Actually there is a lot of Grid in the Cloud!*





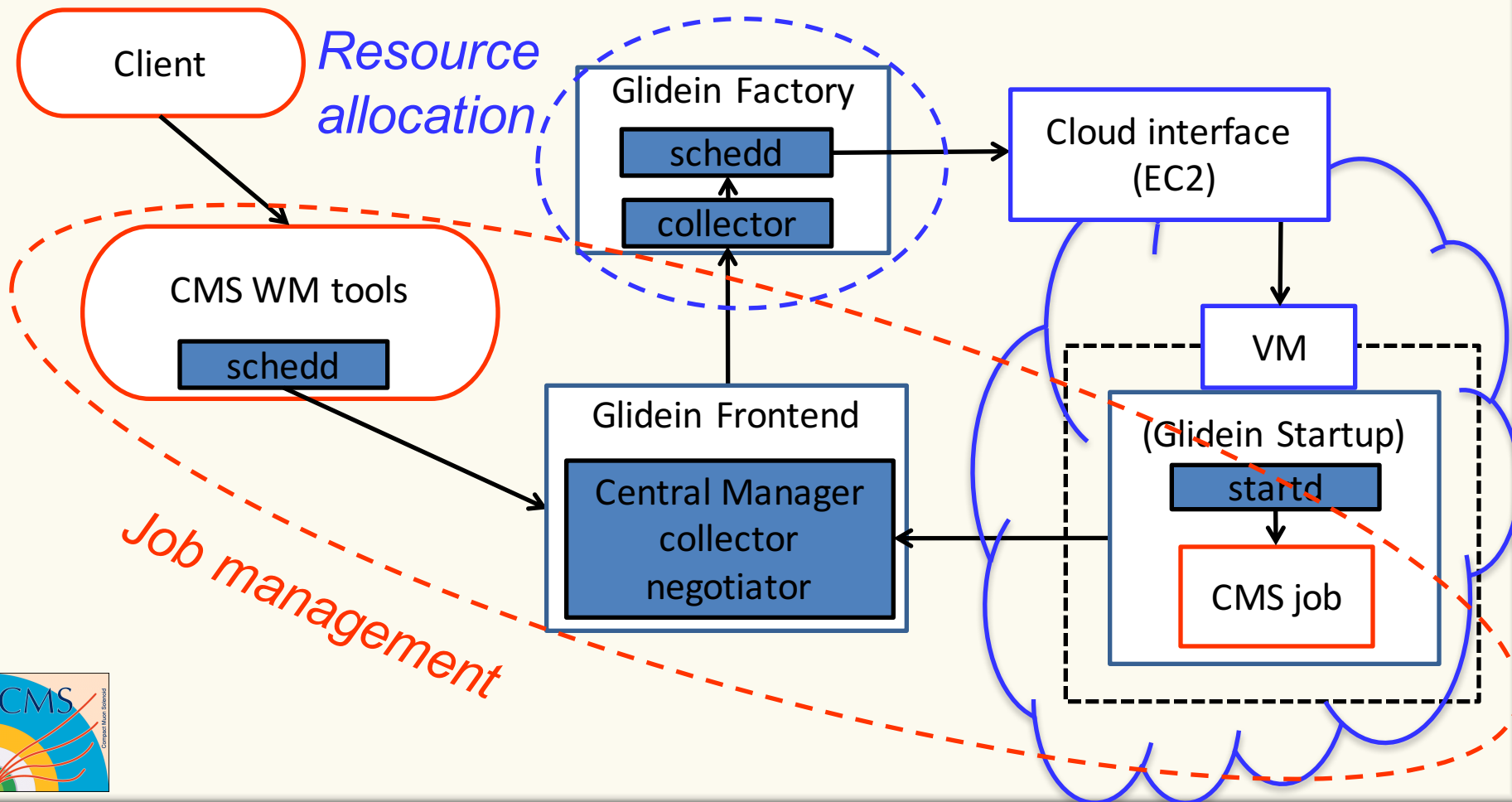
# From the Grid...

The “factory” harvests *job slots*



# ...to the Cloud

The “factory” harvests *machines* (or *containers*)







# Extension of a Computing Centre

*CNAF tested the ability to extend the computing center on external resources:*

Opportunistic: with transient *Aruba* resources (Arezzo)

Structural: *ReCaS/Bari*: extension and management of remote resources

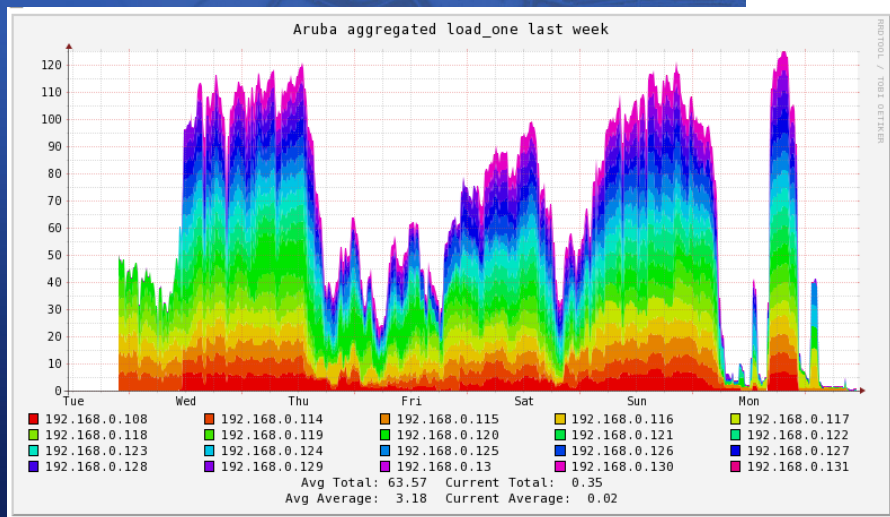
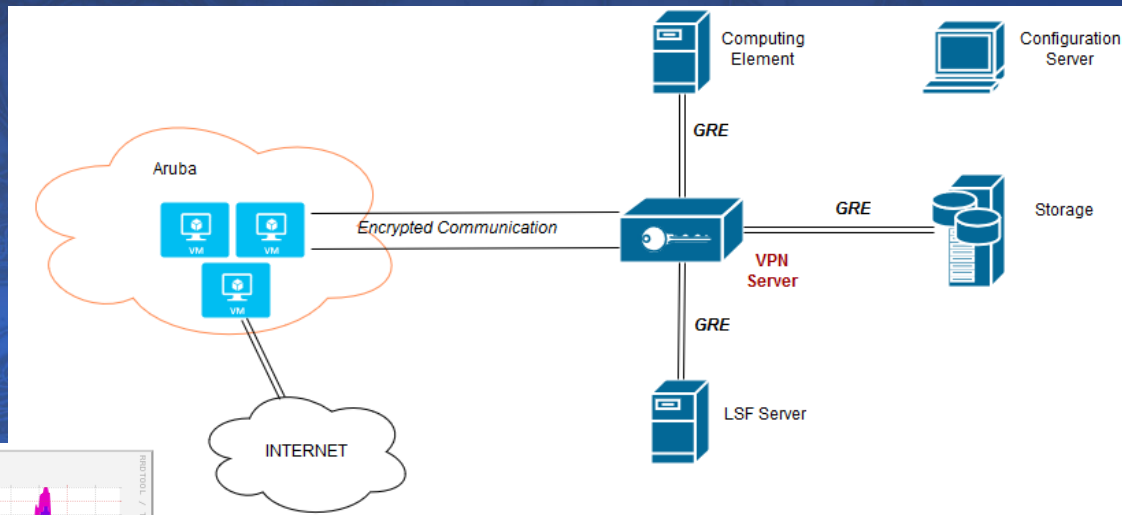
These will become *pledged* resources for CNAF



# Extension on external resources

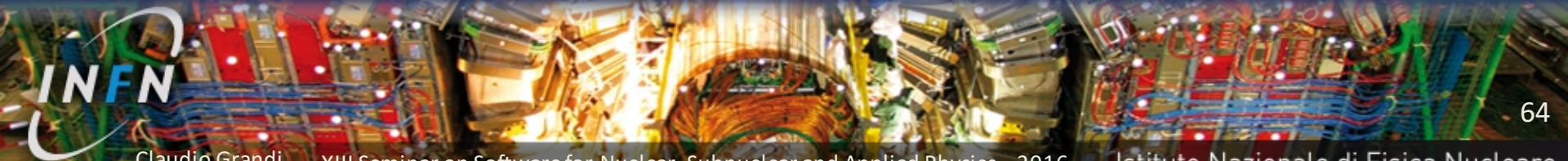
Use of *gateways* and *caches* over secure channels allowed to treat external resources as if internal

*Remote access to data*



First tests on *Aruba* with CMS in 2015

Very good efficiency for specific jobs (Monte Carlo)

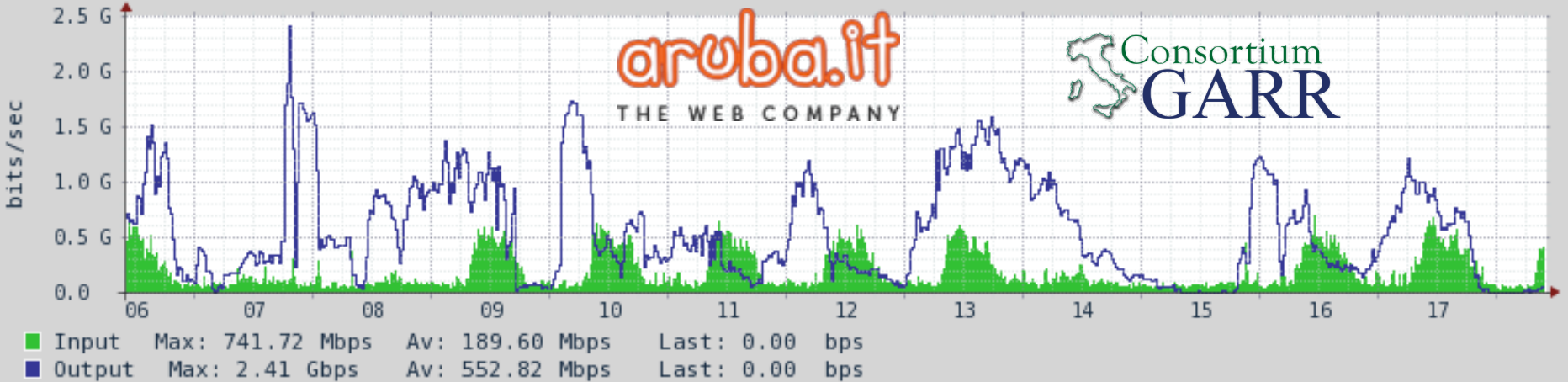






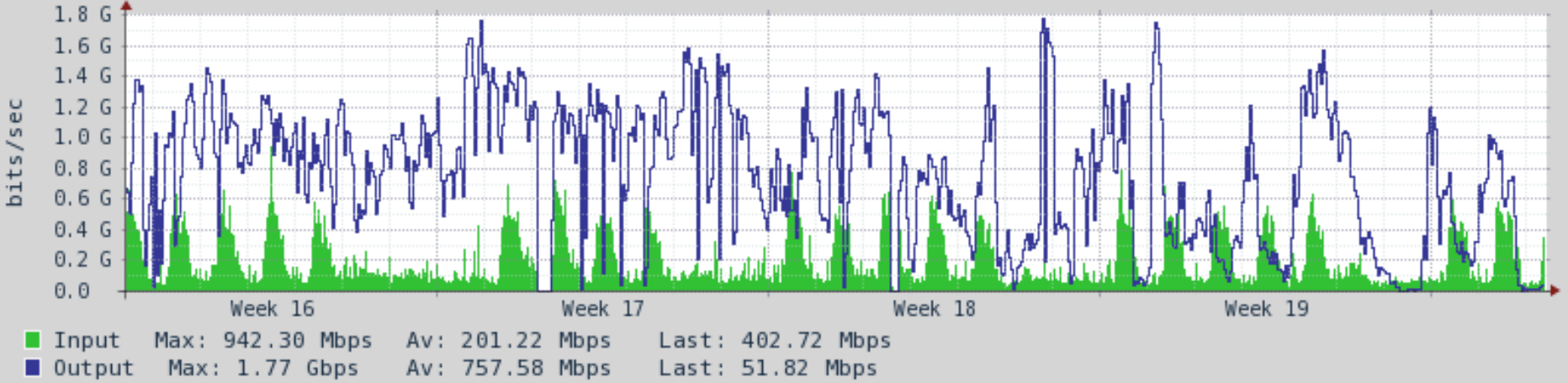
# Aruba - network traffic

AS Traffic: ARUBA-ASN (AS31034) on MIX-PRIMARIO [Week]



Updated on 09.45 am 18/05/16

AS Traffic: ARUBA-ASN (AS31034) on MIX-PRIMARIO [Month]



Updated on 09.45 am 18/05/16

# Remote management of resources

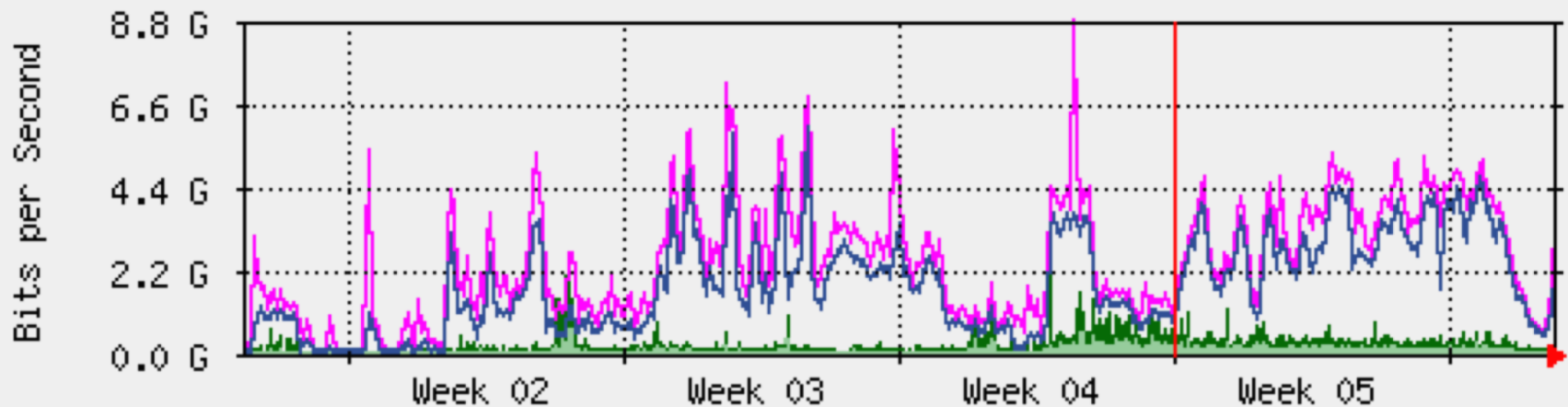
In principle require “*transparent use*” and performances comparable with a local execution

Special network configuration between CNAF and *ReCaS Bari*:

Level 3 Virtual Private Network and 20 Gb/s dedicated bandwidth

Use of caches for data apparently not performing enough

Better results with remote access to data





# Comparative Results

Queue	Nodetype	Njobs	Avg_eff	Max_eff	Avg_wct	Avg_cpt
Cms_mc	AR	2984	0.602	0.912	199.805	130.482
Alice	T1	98451	0.848	0.953	16.433	13.942
Atlas_sc	T1	1211890	0.922	0.972	1.247	1.153
Cms_mc	T1	41412	0.707	0.926	117.296	93.203
Lhcb	T1	102008	0.960	0.985	23.593	22.631
Atlas_mc	T1	38157	0.803	0.988	19.289	18.239
Alice	BA	25492	0.725	0.966	14.446	10.592
Atlas	BA	15263	0.738	0.979	1.439	1.077
Cms_mcore	BA	2261	0.444	0.805	146.952	69.735
Lhcb	BA	13873	0.916	0.967	12.998	11.013
Atlas_sc	BA	20268	0.685	0.878	24.378	15.658

Aruba

CNAF (local)

ReCaS Bari

# New architectures

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

→ Follow the market mainstream

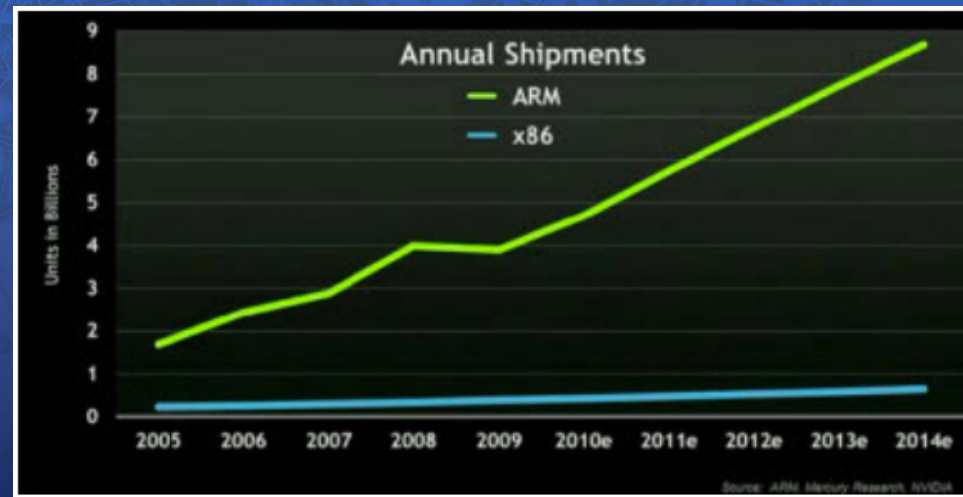
→ Use highly available architectures

ARM, ...

→ Exploit *parallelization*

Multi/many-core, GPGPU, ...

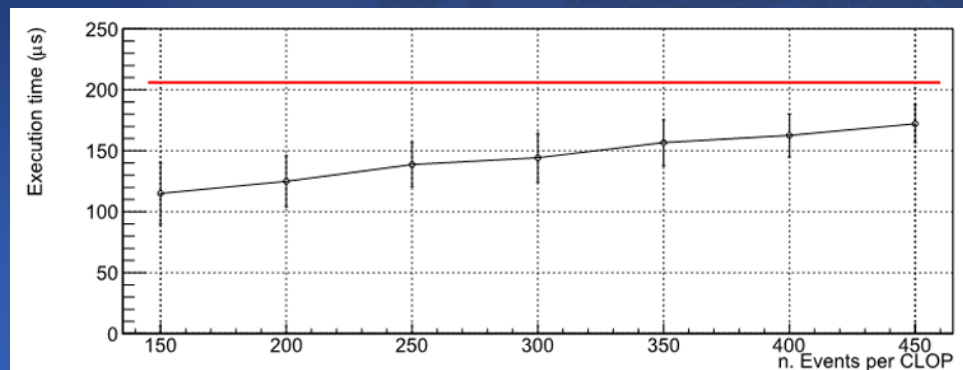
→ Use *low-power* architectures







# NA64 RICH pattern recognition

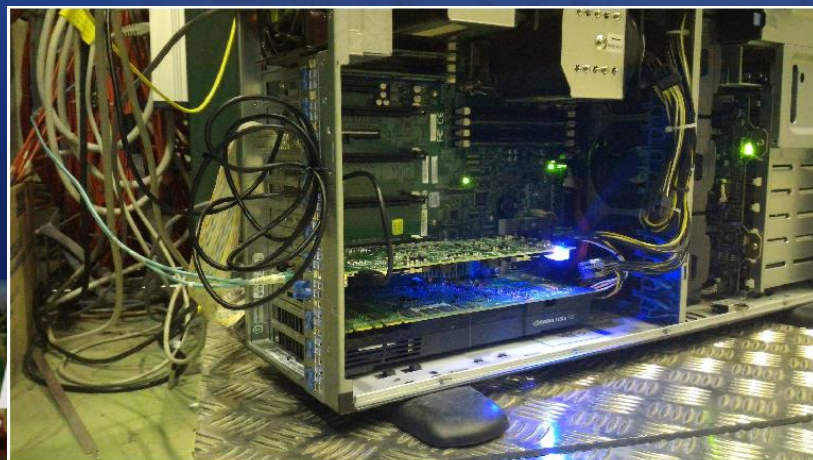
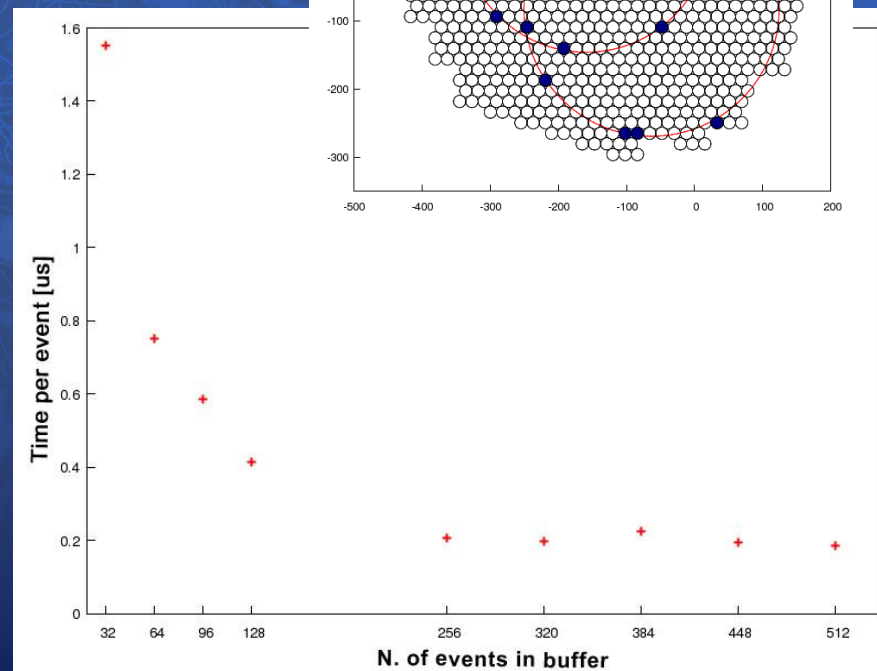
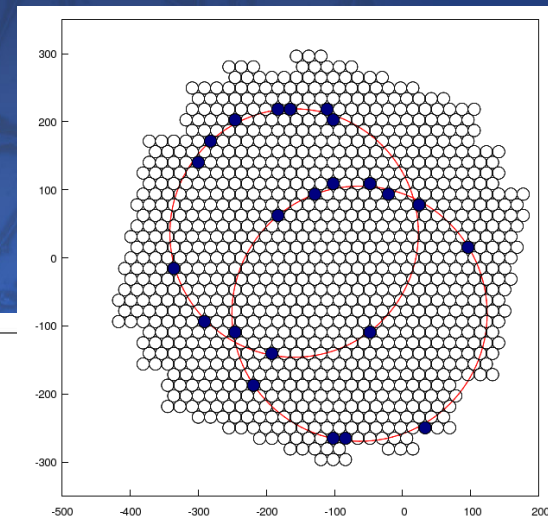


Algorithm parallelization

Almagest

Execution on NaNet-10

Based on GPU the Tesla K20c GPU



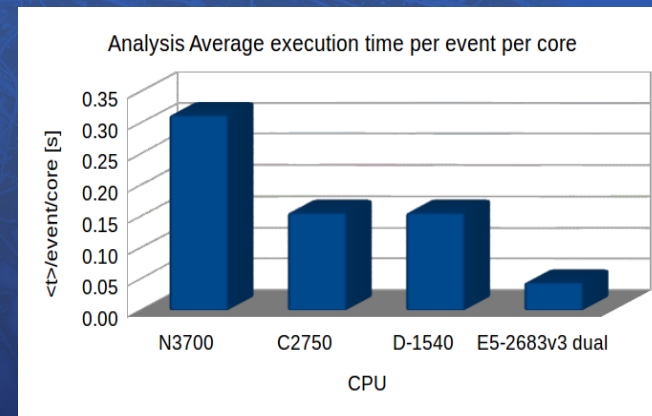
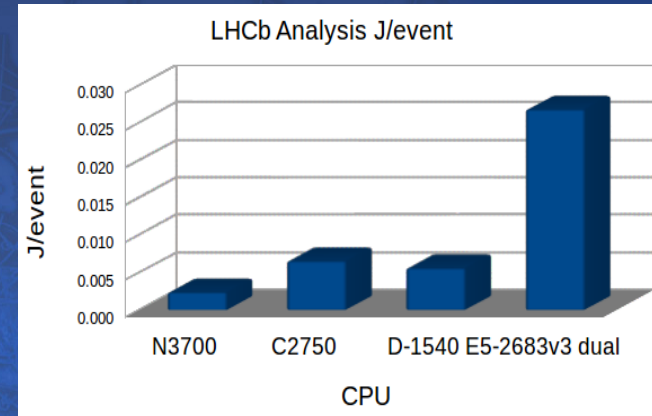
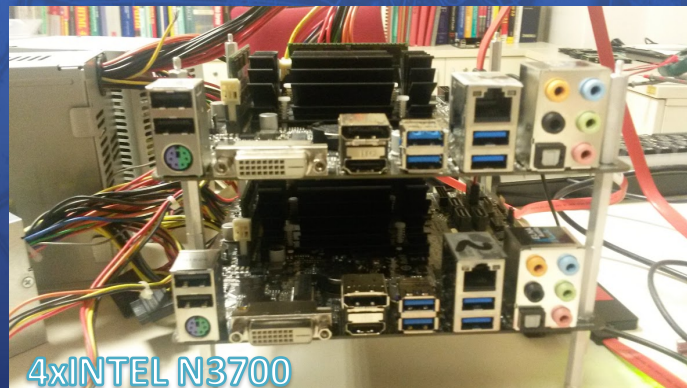


# Test on LHCb analysis

4xINTEL AVOTON C-2750  
4xINTEL XEOND-1540



Example of tests done by the *INFN-COSA* project on a small testbed at CNAF on Intel low power systems



CPU	Brand	Microarchitecture	Family	#	CORES	RAM (GB)	POWER (W)	HS06	HS/W
E5-2683v3	XEON	Haswell	(Reference)	2	56 (HT)	128	370	573	1.55
D-1540	XEON	Broadwell		1	16 (HT)	16	80	151	1.89
C2750	ATOM	Silvermont	Avoton	1	8	16	20	55	2.50
N3700	PENTIUM	Aimont	Braswell	1	4	16	7	28	4.00



*Concluding...*



*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*

*The world-wide and the Italian communities are very active*

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

