# **Computing in LHC**

#### José M. Hernández CIEMAT, Madrid



Taller Altas Energías 2014 Benasque, septiembre 2014



# Experimentación en Física de Partículas



### **Global effort, global success**



# Lots of collisions in LHC

- Low production probability for new physics (e.g. Higgs boson)
  - Higher collision energy, higher production probability: 7,8 → 13,14 TeV
- Highest possible collision rate
  - Proton flux (instantaneous luminosity)
  - ~3000 bunches of ~10<sup>11</sup> protons, transversal section ~10 µm, that cross every 25 nanoseconds (40 MHz)
  - ~25 collisions per bunch crossing
- Collision rate in LHC: ~1 GHz
- New Physics production rate: ~mHz

~1 out of 1.000.000.000.000 needs to be filtered



### Data volume

- Were all collisions stored, the data volume would be huge
  - Average event size: ~ 1 Mbyte
  - Number of events/s = 40 MHz
- ~ 1MB/event x 40x10<sup>6</sup>/s = 40 TB/s

(0.4 Zettabyte/year)

- The interesting events need to be filtered in real time down to a manageable rate
  - → "Trigger"

# **Trigger system**

- Nowadays we use a very complex electronic system
- Multi-level, data buffering, parallel processing
- First Level Trigger
  - Specialized hardware processors
  - Limited information, simple algorithms
  - 40 MHz  $\rightarrow$  100 kHz
  - 3.2 µs latency (128 pipeline),
     ~100 GB/s throughput
- High Level Trigger (software)
  - Linux PC farm
  - Flexible software algorithms
  - 100 kHz → 300 Hz
  - ~50 ms latency (5000 processors)



#### Lots of data

- Trigger accept rate: ~300 Hz
- Annual data volume:

   MB/event ×
   300 events/s ×
   10<sup>7</sup> s/year =
   3.000.000 Gbytes/year =

3 Pbyte/year



### That's a lot of data!

#### LEP, 1989-2000

- In more than a decade, LEP generated less than 1 per mille of 1 year-worth LHC data
- Tevatron, 1983-2011
  - It generated during its lifetime ~25% of 1 year-worth LHC data

### The challenge of LHC data

- Reconstruction and analysis add a new challenge to the data management
  - Raw data must be processed (reconstruction). This process generates a similar data volume
- Simulated data must be generated in order to understand the detector response, study predictions from theoretical models, compare with real data, etc
- In total, ~10-20 PB of data are produced annually per LHC experiment
- Reconstruction, simulation and analysis involve complex calculations that require a large amount of compute resources equivalent to hundreds of thousands of PCs

#### Storage on magnetic tape

Nowadays, in a magnetic tape of about 10 x 10 x 2 cm up to 5 TeraBytes of data can be stored





To manipulate tapes in an efficient way, robotic tape libraries are used to locate/allocate and read/write data

#### Storage on hard disks

#### Nowadays, hard disks in disk servers have a typical capacity of 4 TeraBytes





#### ~200 TB/fileserver 4 TB/disk, 10GE uplinks



#### **Compute servers**



#### quad/six core procs, blades $\rightarrow$ ~500-700 cores/rack



# **CERN Computing center**



# **CERN Computing center**



## **Divide and rule**

- The problem of HEP data processing and analysis can be divided and distributed
  - Each registered event can be processed independently
- A super-computer is not necessary
  - Very expensive facilities and with difficult access
  - Typically used for running a single complex application that can be parallelized and execute in many nodes simultaneously, requiring a high speed and low latency internal network between nodes and common access to memory → High Performance Computing
  - Stringent requirements of memory, processing, speed
  - Applications in meteorology, nuclear fusion, etc

# **Barcelona supercomputing center**



# **High Throughput Computing**

- The computing and memory requirements to execute a data processing application in LHC are modest
- Commodity hardware (e.g. ordinary PCs) can be used
  - Provision of computing resources becomes affordable
- What matters is the aggregate result of processing billions of events executing hundreds of thousands of jobs
- High Throughput Computing
- Proposal for the processing and analysis of LHC data:
  - Let's use computing resources available at the institutions participating in the LHC experiments
  - Let's develop a system to federate those heterogeneous and dynamic resources

# **Grid Computing**

#### The Worldwide LHC Computing Grid (WLCG)

In LHC data are stored, processed and analyzed in a worldwide web of computing centers interconnected through Internet



#### Web → Computación Grid



# 20 years of history: from the WWW to the Grid

- 1989: Tim Berners-Lee proposes at CERN a project of global hypertext (http://)
- 1990: Windows 3.0 popularization of PC millions of people can create digital content
- 1991: First web site http://info.cern.ch
- 1993: 12 users of the 1<sup>st</sup> browser, 50 web servers
- 1994: barrabes.com
- 1995: Netscape on the stock exchange– Start of the .com bubble
- 2001: Collapse of the .com bubble
  - The world remains hyper-connected: the projects of Grid deployment at large scale become feasible
- 2009: First LHC data, analyzed by thousands of researchers around the globe using the **Grid**

#### "The Grid" (Ian Foster y Carl Kesselman, 1998)





- 1990's boom of accessible computing (PC, better communication networks, Internet, Linux, etc)
- State of computing similar to the development of electricity at the beginning of 1900
- The real revolution of electricity was the possibility to distribute it over a network
- The use of computational services should be as transparent as using a power plug
- Users don't need to know from where the computing power is coming from
- Computing revolution similar to the invention of the Web at CERN

# **Computing Grid**

- The Grid, integrating connectivity, computing power and information provides a virtual platform for computation and data management, in the same way the Web integrates resources to create a virtual platform for information
- The Grid makes it possible to dynamically link heterogeneous resources that support large scale processing and intensive use of resources and distributed applications
- The Grid must provided nontrivial quality of service (specified in service level agreements)



# **Grid architecture**

#### Resources

- Computers, storage, communication networks
- Heterogeneous, geographically distributed, dynamic

#### Middleware

- Software to connect and coordinate the resources
- Basic information services, security, data management and execution of tasks, monitoring

#### Applications

Interaction of the user with the Grid



#### **Grid Services at a glance**



# The Worldwide LHC Computing Grid (WLCG)

WLCG is the computing infrastructure that allows

- The connection of all LHC computing centers and its integration in a single "super-computer"
- Make accessible the computing resources to thousands of researchers who in turn are distributed around the globe



#### **Real globalization of LHC data**

# **Topology of WLCG**



WLCG: 150+ centers, 50+ countries, ~300k CPUs, ~ 200PB disk/tape, 10k users



#### Data distribution in real time

Data taken at CERN (**Tier-0**) are distributed to 11 main centers located around the world: the so-called **Tier-1s** 

CFRN

These centers keep a **backup copy of the raw data** and perform **organized mass data processing** (reconstruction)

CERN and the Tier-1 centers are connected through an optical private communications network with a bandwidth of **10.000 Megabits per second – more than 1000 ADSL lines** 

# LHC PN



(thin) <10Gbps edoardo.martelli@cern.ch 20 140620

# **PIC Tier-1 (Barcelona)**

The Spanish Tier-1 center: Puerto de Información Científica (PIC), Barcelona Managed by CIEMAT & IFAE 4k CPUs, 6 PB disk, 8 PB tape

> **PIC** port d'informació científica

#### **Tier-2 Centers**



### **RedIris: Spanish academic network**



http://www.rediris.es/conectividad/weathermap/



### Worldwide academic Internet network



### **Data processing and reduction**



#### WLCG 2010-2012 utilization



#### **Explosion of scientific data**

Particle physics is not the only scientific discipine in which there is an explosion of data

Scientific instruments get digitized, and its precision increases

Exponential increase of data in many fields

In few year we will have telescopes that will excrutinize the sky every night with an unprecendented precision 10 PetaBytes of images every year

### **The Cerenkov Telescope Array**





- Array of cosmic gamma ray detectors
- 1 site in the norhern hemisphere (1 km<sup>2</sup>) y 1 site in the southern hemisphere (3 km<sup>2</sup>)
- 10-20 PB data annually

### Explosión de datos científicos

# In hospitals, diagnosis instruments are getting digitized as well

Databases with medical images are being created for research. Its size will increase at a rate of **dozens of PetaBytes per year** 



### Lots of data!



- 300 million of images uploaded to Facebook every day
- 400 millions of tweets posted every day
- 10<sup>3</sup> millions of visits to YouTube /
  - And 100 hours of video uploaded / min
- 10.000 millions of mobile phones connected in 2020
- Whatsapp processed 27.000 millions of messages in one single day in 2013



e

#### ... and from everywhere

- ADSL, wifi, 3G, 4G...
- Mobile phones, tablets, glasses (fridges, cars...)
- In 2012, the global volume of data reached 2,7 ZB



### **Big Data**



#### Gestión de los datos

• Where is LHC in Big Data Terms?



LHC Computing Upgrade and Evolution

# **From Grid to Cloud Computing**

- In industry, Grid computing has evolved into Cloud Computing
  - Computing Grids are still too specific and inflexible
  - Simplified management of installation and configuration of computing resources
  - Transparent environment for user applications
  - Pay per use business model
  - Interactive applications
  - More flexibility and elasticity in the allocation of resources



# Virtualization

- The technology that enables Cloud Computing
- A virtual computer inside a physical computer. Or several!
- User application is executed inside a process (virtual machine) that provides a virtual operating system



### Virtualization and flexibility

 Virtualization allows to start and move applications between servers according to needs



# **Cloud computing models**

#### Infrastructure as a service

- Basic model
- The provider supplies the physical infrastructure
- The user provides the image with the operating system and the application to be executed
- Platform as a service
  - The provider supplies the platform to develop and execute applications
- Software as a service
  - The provider supplies the applications



# **Everything in the Cloud**



### Services for industry and science



### The heart of the Cloud



# **Cloud Computing in WLCG?**

- Might be the natural evolution
- LHC experiment have added Cloud interfaces
  - Allows e.g. to access commercial clouds like Amazon
- Some sites have moved to Cloud-managed infrastructure
  - For the moment only a small fraction. Not really needed ...
- Cloud infrastructure has its own difficulties
  - No scheduling system
  - Management of VMs becomes experiment responsibility

Evolution of the WLCG Computing Grid

# **Hierarchical processing in LHC Run 1**





# **Evolution towards a more flexible model**

- More efficient use of the resources if centers could perform different types of workflows
  - Possible thanks to the improvements in reliability, performance and Internet connectivity
  - Tier-0
    - Re-use High Level Trigger farm (~10k CPUs) when no data taking
  - Tier-1
    - Execute analysis jobs, simulation, and even prompt reconstruction
  - Tier-2
    - Execute mass re-processing and skimming
- Distinguish center by availability level rather than by executed workflows
  - Tier-0: reaction time of minutes to a problem
  - Tier-1: available 24x7 and reaction time of hours
  - Tier-2: available 8x5 and reaction during working hours

## More efficient data processing

- Initial model of data access in WLCG
  - Data are statically distributed/replicated between sites
  - Processing jobs are sent to sites hosting the requested data which are accessed locally
- Advantages
  - Local access to data is a priori faster and more efficient
  - Data distribution is centrally controlled
- Disadvantages
  - Inefficient use of available CPUs globally
  - CPUs will be idle at a site if there are no jobs requesting data hosted at that site
  - Jobs can be in queue at a site because the requested data are not available at other sites
- Towards a hybrid model of dynamic data distribution and cache release together with remote access to data

#### **Remote access to data**

- Evolution towards a model of distributed storage where processing jobs can access data remotely through the WAN
- Latency and bandwidth limitations can impact on remote data access efficiency
- Lot of effort spent in optimizing remote data reads
  - Read ahead (prediction of data to be read next), vector reads (read parallelization)
  - Working on strategies à la *bit-torrent* to read files from multiple sources

#### Data Federation

- Federation of storage systems where a central service connects the client with the site hosting the requested data
- AAA in CMS, FAX in ATLAS

# **Data distribution**

- Initial model of static collocation of data
  - One or several copies of datasets (real data and simulations) are distributed among the Tier-2 sites for its analysis
- Simple model but with clear disadvantages
  - Manual process
  - It does not take into account which datasets are "hot" (accessed frequently) or "cold" (not accessed for a long time)
- Evolution towards a dynamic model
  - Processing jobs report the data read to a central service
  - This "data popularity service" is queried by the dynamic data collocation service to replicate hot datasets
  - Another service to delete cold datasets from the caches takes care of deleting replicas de datasets that have not been accessed for long time

# **Content Delivery Network**

- Evolution towards a model similar to the one used by Internet content providers
  - Video/audio streaming
- Servers of contents geographically distributed replicate/delete data on demand
- Bring data closer to the application
- Optimization of data access



A Content Delivery Network (CDN) tries to position as many caching server(s) as possible closest to the target audience (to whom the content needs to be served).

> Optimized connection (over the internet) between the origin server(s) and the cache server(s).

Un-Optimized connection (over the internet) between the cache server(s) and the consumers.

### **HEP vs CDN providers**

	NETFLIX	HEP
Bandwidth per client	I.5Mbit	IMB
Clients	IM*	100k cores
Serving	I.5Tbits	0.8Tbits
Total Data Distributed	I2TB	20PB

Similar serving requirements

LHC has a smaller number of clients, less distribution, a higher bandwidth per client and a larger total data volume

# **Using the CDN model**

- Sofware distribution of experiments to processing nodes
  - CERN Virtual Machine File system CVMFS
  - Distributed file system
- Access to calibration and alignment constants
  - FronTier, scalable system to access databases
- Hierarchical systems, scalable, based on data caches
  - The first copy into the cache takes some time, but subsequent accesses are very fast
  - Automatic distribution of new software. Single point of software installation
- Technology used by web servers, cache web proxies
  - Standard http protocol, security



#### **Processor evolution**



#### **Processor evolution**

- Speed of Intel/AMD processor stuck in ~3 GHz since about a decade
  - Too much dissipated heat
- Computing power increases by incrementing the number of processing cores
- Multi-core now → many-core soon → finer grained parallelism needed
- To take advantage of this architecture the application needs to be parallelizable and use
   Utilization of Cores the available cores simultaneously
  - Introduces complexity in the application
  - Non-parallelizable code introduces inefficiency
- Many or most of our codes require extensive overhauls
  - Being adapted: geant4, root, reconstruction code, exp. frameworks



## **Towards multi-core processing**

- Processing model during LHC Run 1
  - Each processing jobs uses a single core
  - In odes with N cores N processing jobs are executed in parallel
- LHC experiments' software is being re-factored to allow for parallelization
  - A single application will use several cores
- Many advantages
  - Better use of RAM and other resources (access to hard disk, network usage)
  - Decrease of the number of jobs the experiment's workload management system has to execute
  - Smaller number of files created



# LHC roadmap

Physics
Shutdown
Beam commissioning
Technical stop



Increasing amount of data and complexity

#### Upgrading LHC Computing in LS1 (2013-2015)

- The shutdown period is a valuable opportunity to asses
  - Lessons and operational experiences of Run 1
  - Computing demands of Run 2
  - The technical and cost evolution of computing
- Undertake intensive planning and development to prepare LHC Computing for 2015 and beyond
  - While sustaining steady state full scale operations
  - With an assumption of flat (at best) funding
- This has been happening internally to the experiments and collaboratively with CERN IT, WLCG, common software and computing projects
  - Upgrade in parallel to accelerator and detector upgrades to push the frontiers of HEP

#### Computing challenges for LHC Run2 (2015-2019)

- Computing in LHC Run1 was very successful but Run 2 from 2015 poses new challenges
- Increased energy and luminosity delivered by LHC in Run 2
  - More complex events to process
    - Event reconstruction time (CMS ~2x)
  - Higher output rate to record
    - Maintain similar trigger thresholds and sensitivity to Higgs physics and to potential new physics
    - ATLAS, CMS event rate to storage 2.5x
- Need a substantial increase of computing resources that we probably cannot afford





#### **Computing resources increase**



- Requests ~conform to expected flat budget
- ~25% yearly growth requested
- Benefit from technology evolution



66

#### **Computing strategy for Run2**

- Increase resources in WLCG as much as possible
  - Try to conform to constrained budget situation
  - Request ~25% yearly growth. Profit from technology evolution
- Make a more efficient and flexible use of the available resources
  - Reduce CPU and storage needs
    - Less reprocessing passes, less simulated events, more compact data format, reduce data replication factor
  - Intelligent dynamic data placement and remote data access
    - Automatic replication of hot data and deletion of cold data, remote I/O
  - Break down the boundaries between the computing tiers
    - Run reconstruction, simulation and analysis at Tier-1/Tier-2 indistinctly
  - Centralized production of group analysis datasets
    - Shrink 'chaotic analysis' to only what really is user specific
    - Remove redundancies in processing and storage, reducing operational workloads while improving turnaround for users
  - Evolve the data processing frameworks towards parallel processing
    - More efficient use of multicore processors

#### Access to new resources for Run 2

- Access to opportunistic resources
  - Unused capacities at Grid sites that allow opportunistic usage
  - Capacities provided to the experiments for a defined period of time at High Performance Computing Centres, etc
  - Significant increase in capacity with low cost (satisfy capacity peaks)
  - HPC clusters, academic or commercial clouds, volunteer computing
- Use HLT farm for offline data processing
  - ~10k cores
  - During extended periods with no data taking and even inter-fill periods
- Adopt advanced architectures
  - Processing in Run1 done under Enterprise Linux on x86 CPUs
  - Many-core processors, low-energy CPUs (e.g. ARM processors of mobile phones), accelerator cards (GPU)
  - Challenging heterogeneous environment
  - Parallelization of processing application will be key

# **Graphics Processing Unit (GPU)**

- Specialized processors initially used in graphics cards
- Used now as general purpose processors
- Hundreds of cores in a single card
- Good relation power consumption / computing power
- NVIDIA is the largest commercial provider. Inventor of GPU in 1999
- CUDA is the platform and the parallel programming model created by NVIDIA for the GPUs





# Looking ahead: LHC Run 4 (2024+)

- Run 4 w.r.t Run 2
  - Increase output rate 10x
    - 1 →10 kHz
  - Increase event processing time 2.5x
    - 40 → 140 pileup
  - Increase event size 2x
  - 25x CPU, 20x storage needs
- Expected increase of resources with flat budget
  - CPU doubling ~every 3 years (25%/year): 8x till 2024
  - Disk doubling ~every 4 years (20%/year): 5x till 2024
- About factor 3 (CPU) and factor 4 (disk) missing
- Need long term I+D+I to achieve a computing revolution needed to meet these huge requirements
  - In 1997, the Run 1 challenge was equally daunting. It took 10 years to develop the "WLCG computing revolution" to meet Run 1-2-3 requirements

#### Conclusions

- LHC distributed computing, based on Grid technologies, performed extremely well at all levels in Run 1
  - It has allowed experiments to produce great scientific results with unprecedented speed
  - We know how to deliver, adapting where necessary
  - Excellent networks, flexible and adaptable computing models and software systems paid off in exploiting resources
- Massive distributed computing has extended to industry and society thanks to the extraordinary development of Internet and the availability of commodity hardware
  - Big data, Cloud computing, ...
- LHC computing will have to face new challenges in the coming years
  - Difficult funding scenario
  - Optimization, flexibility, evolution