

# Computing in High Energy Physics

*ISTC-CERN-JINR Summer School*

*CERN, 30 September 2011*

**Alexei Klimentov**

Brookhaven National Lab

# Acknowledgements

- Thanks to I.Bird, J.Boyd, S.Campana, T.Cass, B.Dahmes, D.Duelmann, M.Ernst, A.Farbin, I.Fisk, R.Heuer, M.Lamanna, N.Neufeld, S.Panitkin, L.Robertson, J.Shiers, M.Schulz, R.Walker and T.Wenaus whose slides were used in this presentation.
- Thanks to many colleagues from the ATLAS experiment at LHC

# *Act I*

## *A Bit of History*

# Outline

- Data Acquisition
- Bulk data processing
- Data Analysis
- IBM PC
- Internet & Web
- Linux

*TLA – Three Letters Acronym*

- *CM – Computing Model*
- *HEP – High Energy Physics*
- *LHC – Large Hadron Collider*

Related talk :

J.Andreeva : LHC Computing Grid

# *Data Acquisition*

# Tycho Brahe and the Orbit of Mars

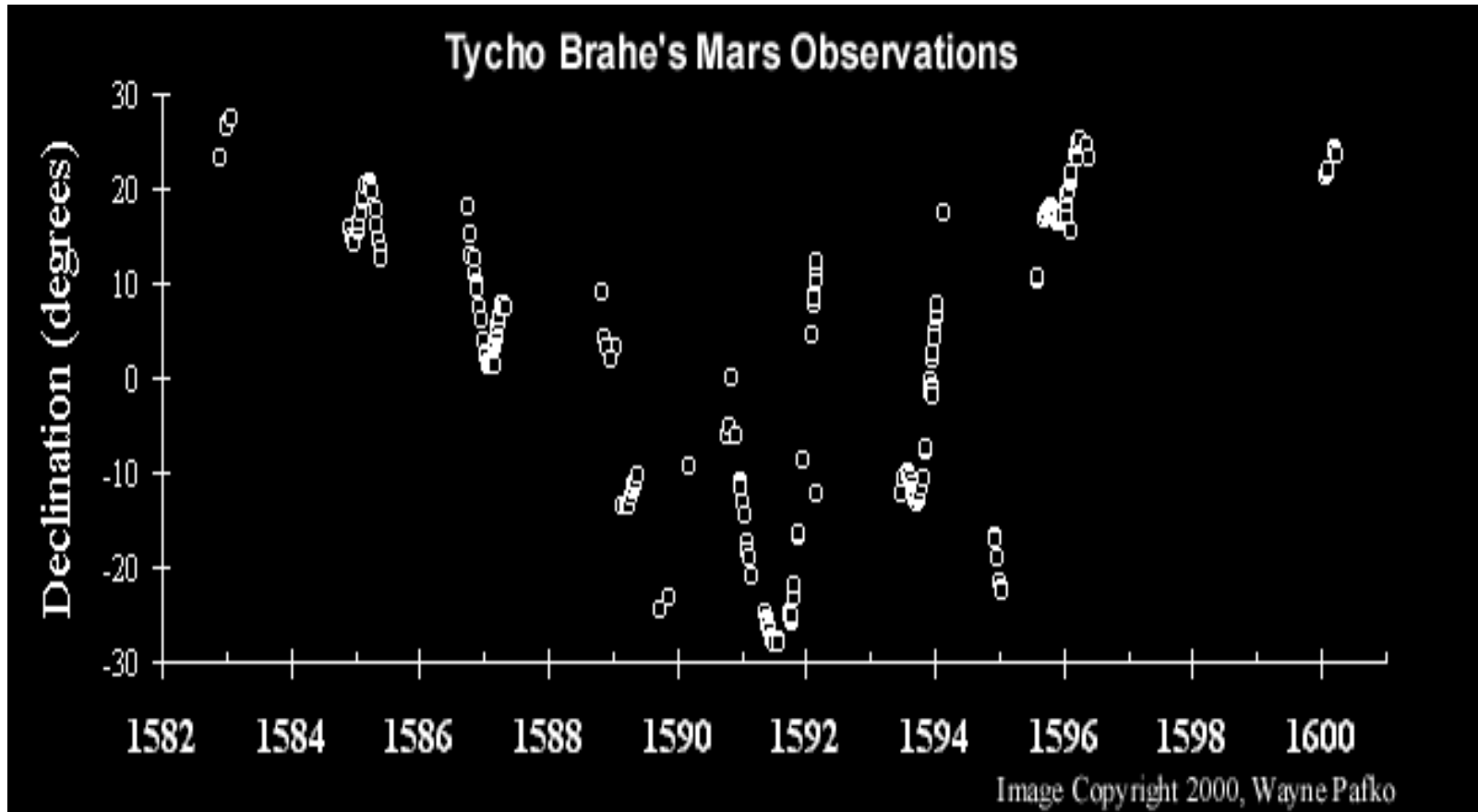
*I've studied all available charts of the planets and stars and none of them match the others. There are just as many measurements and methods as there are astronomers and all of them disagree. What's needed is a long term project with the aim of mapping the heavens conducted from a single location over a period of several years.*

Tycho Brahe, 1563 (age 17).



- First measurement campaign
- Systematic data acquisition
  - Controlled conditions (same time of the day and month)
  - Careful observation of boundary conditions (weather, light conditions etc...) - important for data quality / systematic uncertainties

# The First Systematic Data Acquisition



- Data acquired over 18 years, normally every month
- Each measurement lasted at least 1 hr with the naked eye
- Red line (only in the animated version) shows comparison with modern theory

# *Data Analysis*





- Tycho invited Johannes Kepler (1571-1630) to analyse the Mars data, he did it, eventually paving the way for Isaac Newton theory of universal gravitation

# *Data Processing*

A

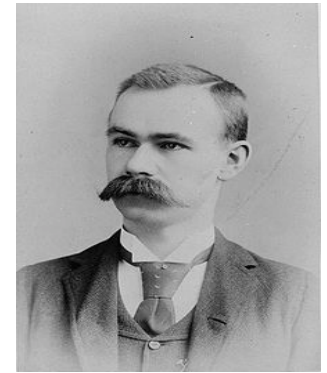
- According to the 1900s...



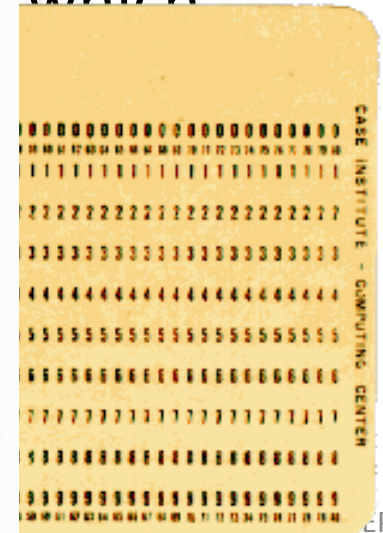
9/30/2011



[Jacquard-card Making.]



ed in late  
a  
s to  
ces of  
which



EP

# Hollerith's Successes

- In 1890 Hollerith founded a company called the Tabulating Machine Company.
- In 1911, his company merged with two other companies to create the Computing-Tabulating-Recording Company.
- Under the direction of **Thomas Watson, Sr**, CTR would change its name in 1924 to **International Business Machines**. Hollerith's machine would provide the basis for IBM's success and make him the father of information processing.



# Thomas Watson Jr

I think there is a world market for maybe five computers  
– 1943



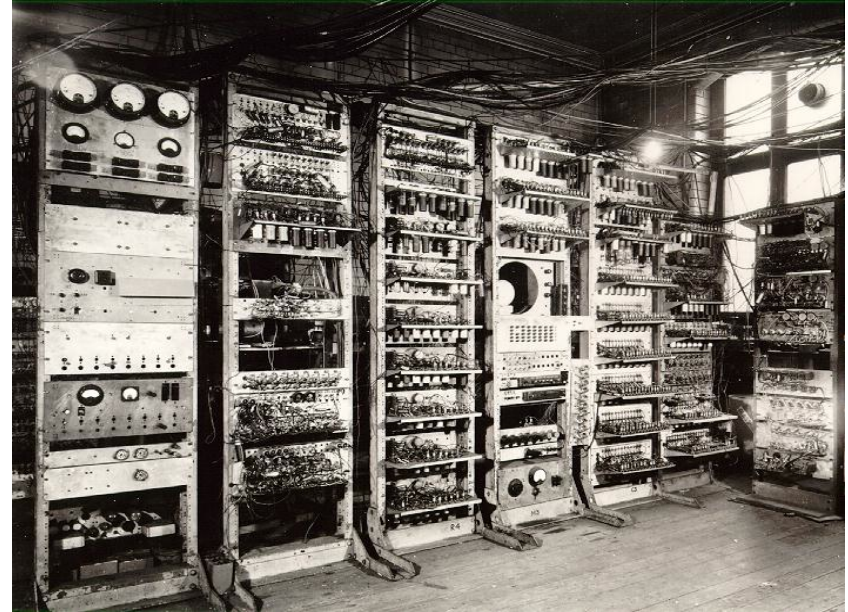
| Name                           | First operational | Numeral system        | Computing mechanism | Programming   |     |
|--------------------------------|-------------------|-----------------------|---------------------|---|-----|
| Zuse Z3 (Germany)              | May 1941          | Binary floating point | Electro-mechanical  | Program-controlled by punched 35 mm film stock (conditional branch)             |     |
| Atanasoff–Berry Computer (US)  | 1942              | Binary                | Electronic          | Not programmable—single purpose   |     |
| Colossus Mark 1 (UK)           | February 1944     | Binary                | Electronic          | Program-controlled by patch cables and switches                                 |     |
| Harvard Mark I – IBM ASCC (US) | May 1944          | Decimal               | Electro-mechanical  | Program-controlled by 24-channel punched paper tape (but no conditional branch) | No  |
| Colossus Mark 2 (UK)           | June 1944         | Binary                | Electronic          | Program-controlled by patch cables and switches                                 | No  |
| Zuse Z4 (Germany)              | March 1945        | Binary floating point | Electro-mechanical  | Program-controlled by punched 35 mm film stock                                  | Yes |

# Computers: Predictions

“Computers in the future may weight no more than 1 ton”

– Popular mechanics, 1949

The first electronic computer was named Colossus (~1944) and weighed approximately one ton



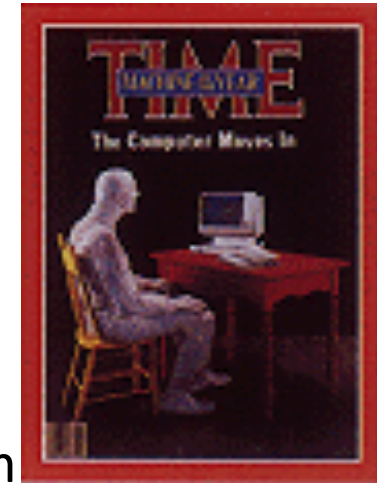
# Ken Olson on the PC



- “There is no reason for any individual to have a computer in his home.”
  - Ken Olson, president, digital equipment corporation, (DEC) 1977.
- Ironic that DEC was subsequently taken over by COMPAQ...
- ... and COMPAQ was taken over by Hewlett-Packard (HP) in Jan 2002.
- ... and HP announced that the company has intention to sell this part of the business after 2012



# The IBM PC



- On August 12, 1981, IBM released their new computer, the IBM PC.
- In July of 1980, IBM representatives met with Microsoft's Bill Gates to talk about an operating system for the PC.
- In 1983, Time Magazine named the PC “Man of the Year.”
  - i.e. just over 1 year from the launch.
- Jan , 1986. IBM PC used for HEP experiment at CERN (L3) to calibrate Hadron calorimeter utilizing natural Uranium radioactivity
- Now. A laptop per physicist
  - ATLAS collaboration has 3000+ physicists and engineers



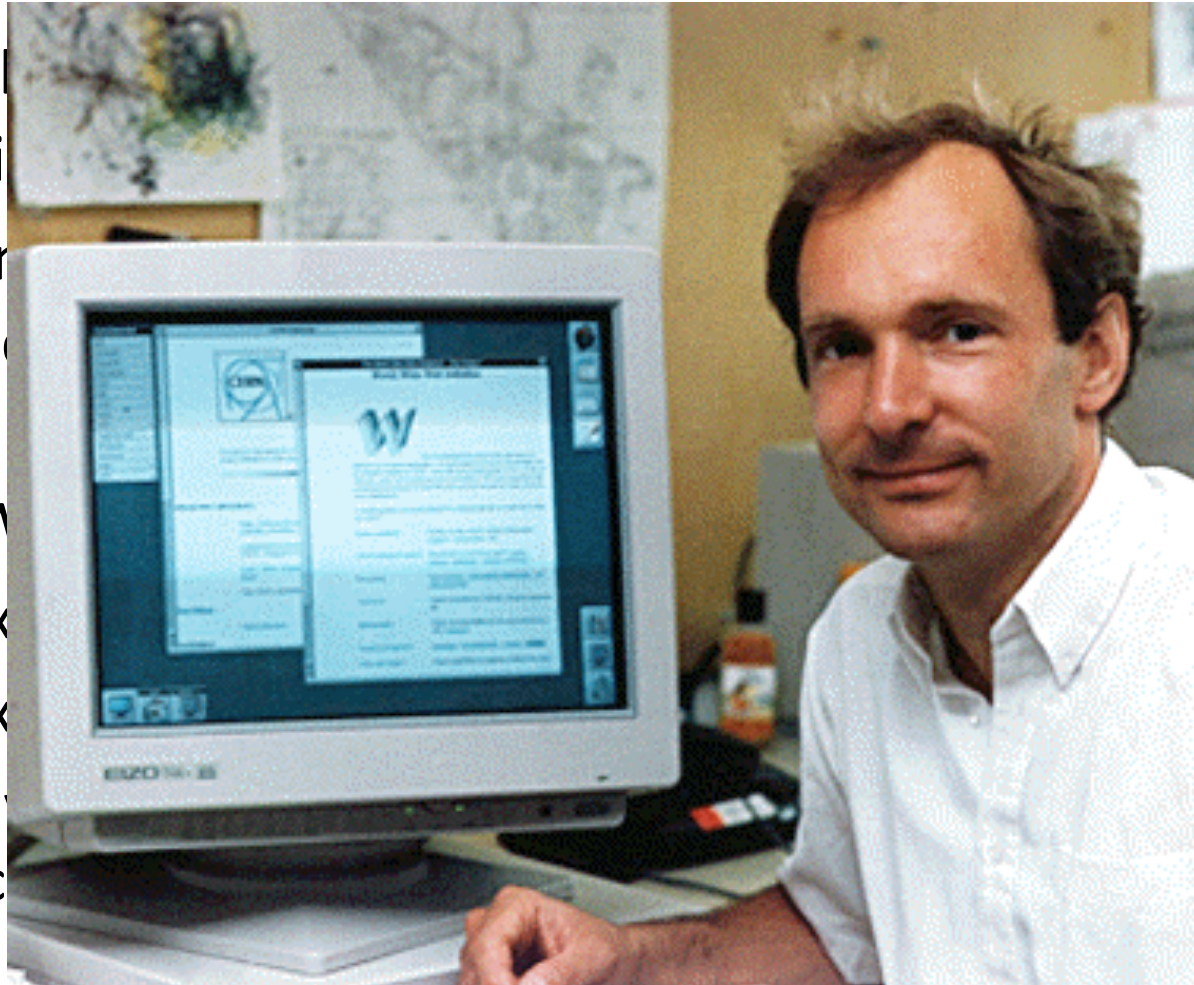
# *The Internet, Web, Linux*

# Internet Timeline

- 1957: sputnik, ARPA (Advanced Research Project Agency)
- Early 1960s: papers on packet switching, ideas for a “Intergalactic Computer Network”
  - Those ideas contained almost everything that composes the contemporary Internet
- Late 1960s: ARPANET
  - Original design speed: 2.4kbps
  - 4 sites (aka 4 IMP, Interface Message Processor)
- Early 1970s: network control protocol
- 1 January 1983: move to TCP/IP
  - Originally 32 bit addresses
- 1986: US NSF develops NSFNET
  - Originally 56Kbps links
  - Today leading backbone of internet
- 2011 : LHC Optical Private Network 10 Gbps link(s) between CERN and 11 Physics centers (aka Tier-1)

# Birth of the Web

- Original p
- Resubmi
- Various r
- (MOI) et
- 1991:
  - WWW
- 1992 – ex
- 1993 – ex
  - Large
  - Applic



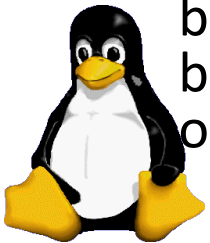
on

outing

# LINUX



- 1991 – the first version of Linux operating system
  - Posted announcement to the Minix group on USENET, and made the Linux source code available to other nerds free of charge. Programmers everywhere started adding their own improvements, and eventually companies like Red Hat, Corel, Caldera, and TurboLinux began selling their own versions of Linux.
  - The open-source nature of Linux is its greatest strength. Instead of having paid programmers devising improvements and looking for bugs from 9-to-5 with tight deadlines and budgets and memos from bosses, Linux is perpetually being tinkered with by the most obsessed and enthusiastic high-tech hobbyists and experts.
- 2011 – the primary Operating System in High Energy Physics Centers



*'Just For Fun' is a humorous autobiography of Linus Torvalds*



Larry Ellison (ORACLE CEO) predicted the future of computing (~early 2000):

- *“There have been 3 generations of computing: mainframe, client-server and Internet computing*
- *There’ll be nothing new for one thousand (1000) years”*

Curiously enough, very soon Oracle declared Grid to be **“the next big thing”**

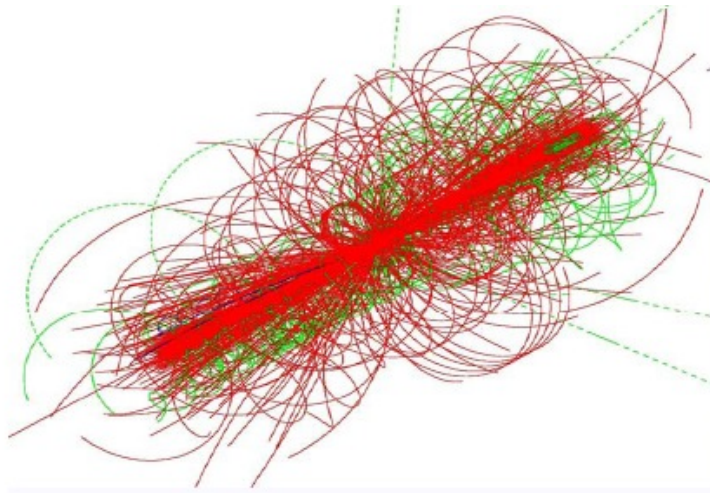
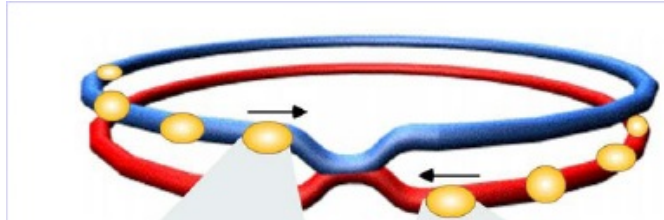


*Act II*

Large Hadron Collider

*Computing at LHC (looking at the data)*

# Proton-Proton Collisions at the LHC



- **2808 + 2808 proton bunches separated by 7.5 m**  
→ **collisions every 25 ns**  
**= 40 MHz crossing rate**
  - **$10^{11}$  protons per bunch**
  - **at  $10^{34}/\text{cm}^2/\text{s}$**   
 **$\approx 35$  pp interactions per crossing**  
**pile-up**
  - **$\approx 10^9$  pp interactions per second !!!**
  - **in each collision**  
 **$\approx 1600$  charged particles produced**
- enormous challenge for the detectors  
and for data collection/storage/analysis**

# Enter a New Era in Fundamental Science

Start-up of the Large Hadron Collider (LHC), one of the largest and truly global scientific projects ever, is the most exciting turning point in particle physics.

(R.Heuer. CERN Director General)



CMS

- 27 km ring of superconducting magnets operating at 1.9° Kelvin
- Colliding proton beams travel at 99.999999991% the speed of light
- Collisions at 3.5 TeV + 3.5 TeV generate temperatures a billion times hotter than the heart of the sun
- 40 MHz Beam interaction rate – every 25 ns



LHCb

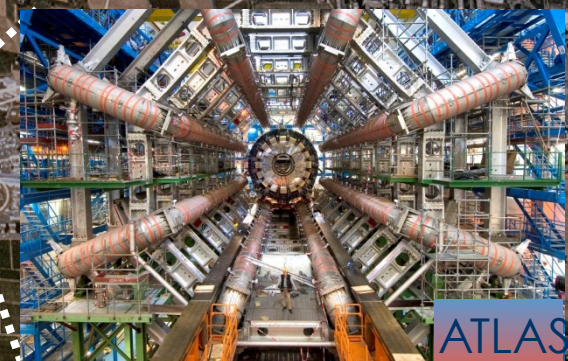


ALICE



LHC ring:  
27 km circumference

TOTEM  
LHCf  
MOEDAL



ATLAS

9/30/2011





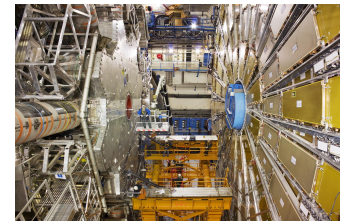
# A Thoroidal LHC ApparatuS

- ATLAS is one of the six particle detectors experiments at Large Hadron Collider (LHC) at CERN
- The project involves more than 3000 scientists and engineers in ~40 countries
- ATLAS has 44 meters long and 25 meters in diameter, weighs about 7,000 tons. It is about half as big as the Notre Dame Cathedral in Paris and weighs the same as the Eiffel Tower or a hundred 747 jets



9/30/2011

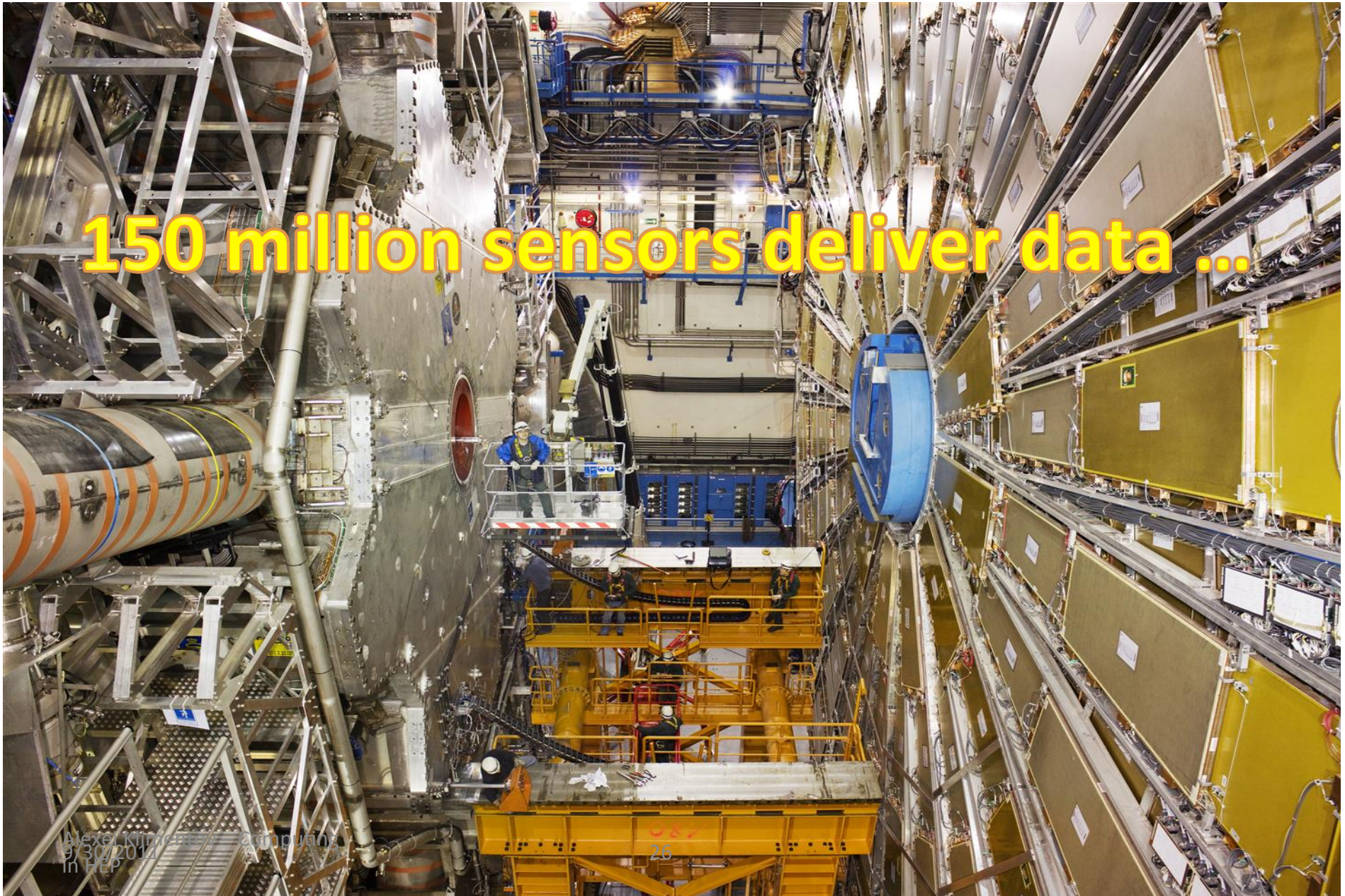
Alexei Klimentov – Computing in HEP



25

# ATLAS

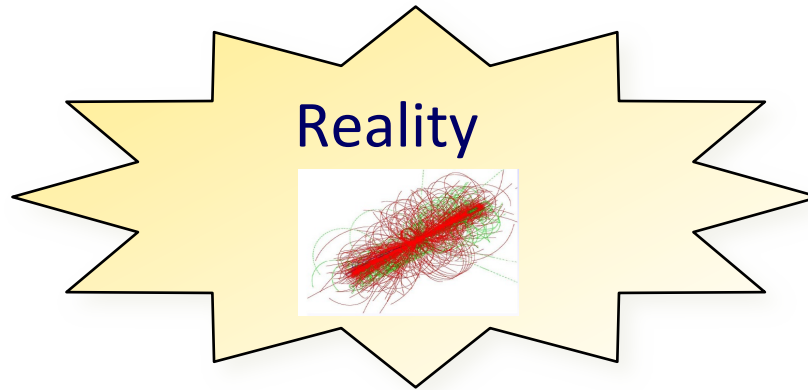
150 million sensors deliver data ...



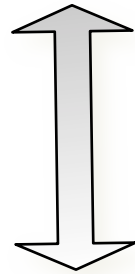
# Terminology

- Data is collected online (real-time)
  - Collision data recorded by the detectors
- Physicists analyze this data offline
  - Optimizing selection, estimating/modeling background, establishing limits, discovering New Physics, etc.
- The LHC delivers a lot of data, which we need to first select online for future analysis
  - Data filtering is done online and offline

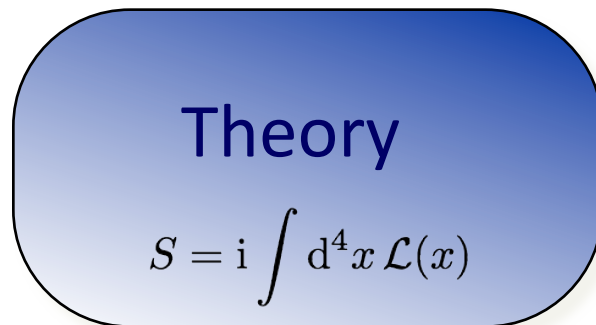
# Our Task



We use experiments to inquire about what “reality” (nature) does



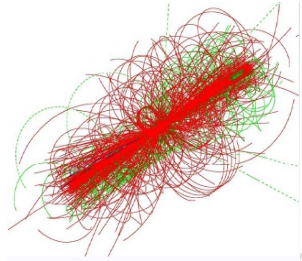
We intend to fill this gap



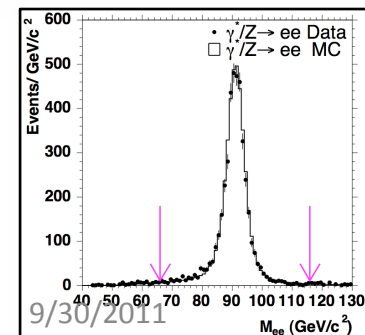
The goal is to understand in the most general; that’s usually also the simplest.

- A. Eddington

# Data Analysis Chain

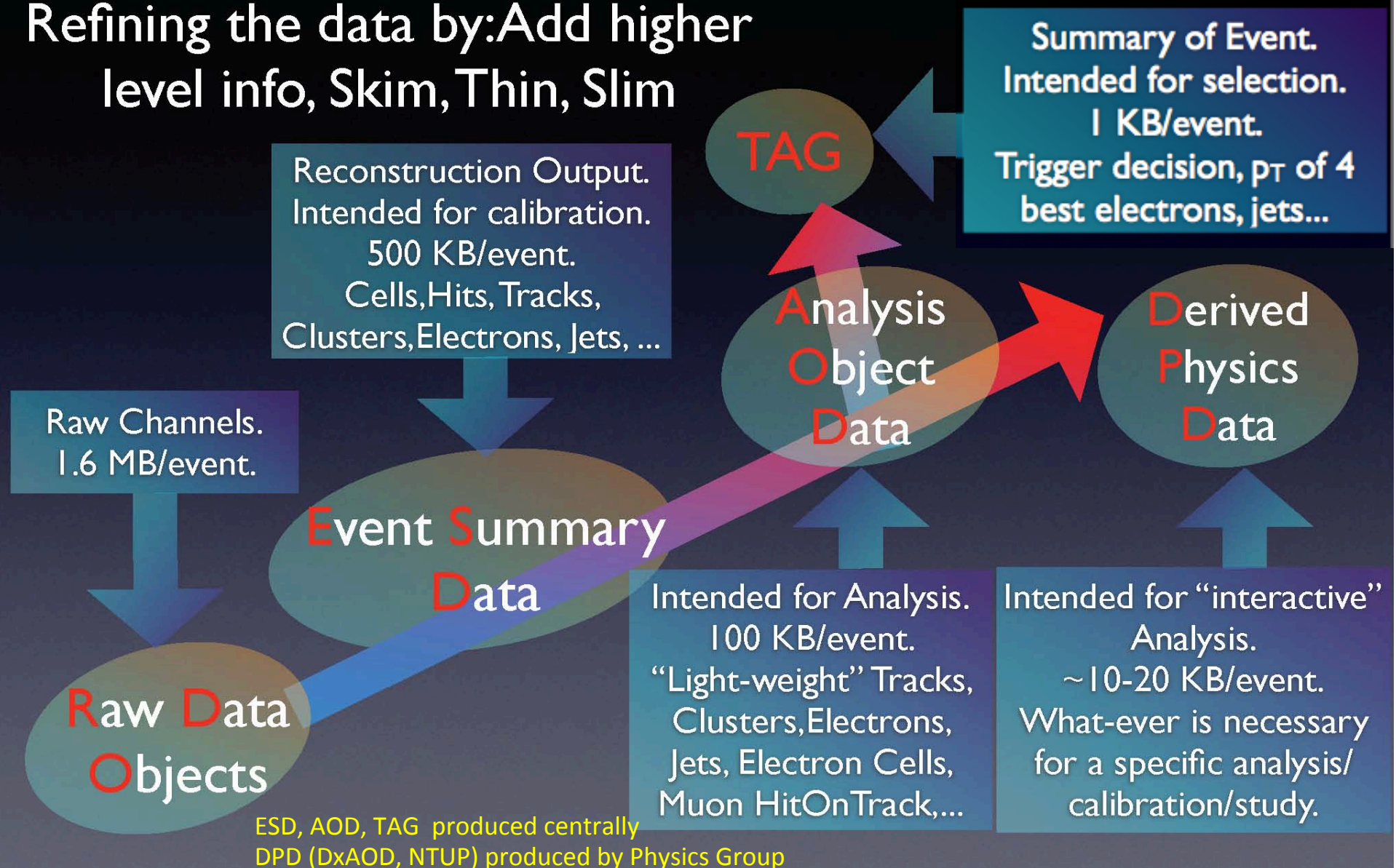


- Have to collect data from many channels on many sub-detectors (millions)
- Decide to read out everything or throw event away (Trigger)
- Build the event (put info together)
- Store the data
- Analyze them
  - reconstruction, user analysis algorithms, data volume reduction
- do the same with a simulation
  - correct data for detector effects
- Compare data and theory



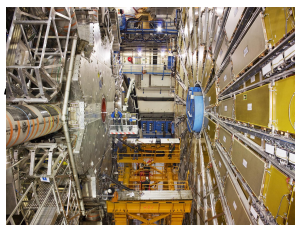
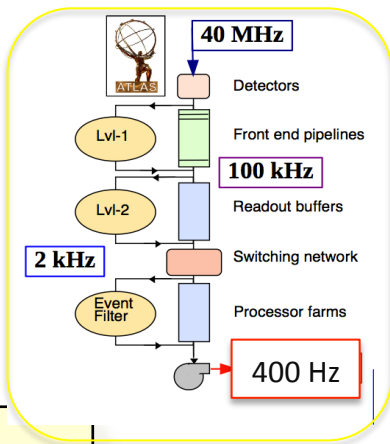
# The Event Data Model

Refining the data by: Add higher level info, Skim, Thin, Slim



# Data and Computation for Physics Analysis

Reduce data volume in stages  
 Select ONLY 'interesting' events  
 Initial data rate (25 ns) :  
 40 000 000 events/s  
 Selected and stored  
 400 events/s



**event filter  
(selection & reconstruction)**

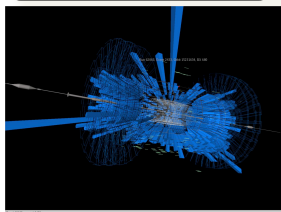
```

0x01e84c10: 0x01e8 0x8848 0x01e8 0x83d8 0x6c73 0x672 0x7400 0x0000
0x01e84c20: 0x0000 0x0019 0x0000 0x0000 0x01e8 0x4d08 0x01e8 0x5b7c
0x01e84c30: 0x01e8 0x87e8 0x01e8 0x8458 0x7061 0x656b 0x6167 0x6500
0x01e84c40: 0x0000 0x0019 0x0000 0x0000 0x0000 0x01e8 0x5b7c
0x01e84c50: 0x01e8 0x8788 0x01e8 0x8498 0x7072 0x6f63 0x0000 0x0000
0x01e84c60: 0x0000 0x0019 0x0000 0x0000 0x0000 0x01e8 0x5b7c
0x01e84c70: 0x01e8 0x8824 0x01e8 0x84d8 0x7265 0x6725 0x7370 0x0000
0x01e84c80: 0x0000 0x0019 0x0000 0x0000 0x0000 0x01e8 0x5b7c
0x01e84c90: 0x01e8 0x8838 0x01e8 0x84e8 0x7265 0x6725 0x7362 0x0000
0x01e84ca0: 0x0000 0x0019 0x0000 0x0000 0x0000 0x01e8 0x5b7c
0x01e84cb0: 0x01e8 0x8818 0x01e8 0x8558 0x7265 0x6725 0x6d65 0x0000
0x01e84cc0: 0x0000 0x0019 0x0000 0x0000 0x0000 0x01e8 0x5b7c
0x01e84cd0: 0x01e8 0x8798 0x01e8 0x8598 0x7265 0x6725 0x726e 0x0000
0x01e84ce0: 0x0000 0x0019 0x0000 0x0000 0x0000 0x01e8 0x5b7c
0x01e84cf0: 0x01e8 0x87c8 0x01e8 0x8568 0x7265 0x6725 0x726e 0x0000
0x01e84d00: 0x0000 0x0019 0x0000 0x0000 0x0000 0x01e8 0x5b7c
0x01e84d10: 0x01e8 0x87e8 0x01e8 0x8618 0x7265 0x6725 0x7400 0x0000
0x01e84d20: 0x0000 0x0019 0x0000 0x0000 0x0000 0x01e8 0x5b7c
0x01e84d30: 0x01e8 0x87a8 0x01e8 0x8658 0x7370 0x6974 0x7400 0x0000
0x01e84d40: 0x0000 0x0019 0x0000 0x0000 0x0000 0x01e8 0x5b7c
0x01e84d50: 0x01e8 0x8858 0x01e8 0x8698 0x7374 0x7269 0x6725 0x0000
0x01e84d60: 0x0000 0x0019 0x0000 0x0000 0x0000 0x01e8 0x5b7c
0x01e84d70: 0x01e8 0x87e8 0x01e8 0x86d8 0x7375 0x6273 0x7374 0x0000
0x01e84d80: 0x0000 0x0019 0x0000 0x0000 0x0000 0x01e8 0x5b7c
0x01e84d90: 0x01e8 0x7f0 0x01e8 0x8718 0x7377 0x6974 0x6368 0x0000
    
```

**event simulation**



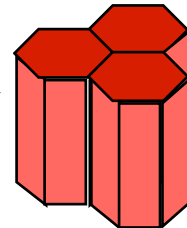
**event reconstruction**



**event summary data**



**batch physics analysis**

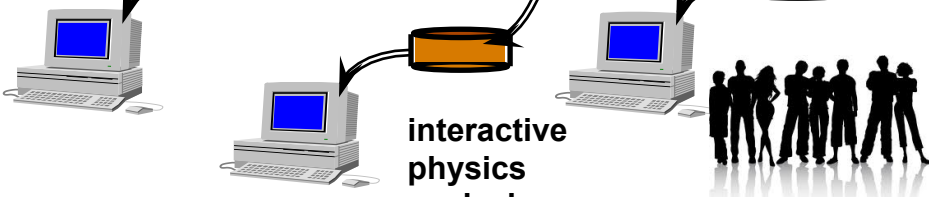


**processed data**

**analysis objects  
(extracted by physics topic)**



**interactive physics analysis**



# Experiment...

```
0x01e84c10: 0x01e8 0x8848 0x01e8 0x83d8 0x6c73 0x6f72 0x7400 0x0000
0x01e84c20: 0x0000 0x0019 0x0000 0x0000 0x01e8 0x4d08 0x01e8 0x5b7c
0x01e84c30: 0x01e8 0x87e8 0x01e8 0x8458 0x7061 0x636b 0x6167 0x6500
0x01e84c40: 0x0000 0x0019 0x0000 0x0000 0x0000 0x0000 0x01e8 0x5b7c
0x01e84c50: 0x01e8 0x8788 0x01e8 0x8498 0x7072 0x6f63 0x0000 0x0000
0x01e84c60: 0x0000 0x0019 0x0000 0x0000 0x0000 0x0000 0x01e8 0x5b7c
0x01e84c70: 0x01e8 0x8824 0x01e8 0x84d8 0x7265 0x6765 0x7870 0x0000
0x01e84c80: 0x0000 0x0019 0x0000 0x0000 0x0000 0x0000 0x01e8 0x5b7c
0x01e84c90: 0x01e8 0x8838 0x01e8 0x8518 0x7265 0x6773 0x7562 0x0000
0x01e84ca0: 0x0000 0x0019 0x0000 0x0000 0x0000 0x0000 0x01e8 0x5b7c
0x01e84cb0: 0x01e8 0x8818 0x01e8 0x8558 0x7265 0x6e61 0x6d65 0x0000
0x01e84cc0: 0x0000 0x0019 0x0000 0x0000 0x0000 0x0000 0x01e8 0x5b7c
0x01e84cd0: 0x01e8 0x8798 0x01e8 0x8598 0x7265 0x7475 0x726e 0x0000
0x01e84ce0: 0x0000 0x0019 0x0000 0x0000 0x0000 0x0000 0x01e8 0x5b7c
0x01e84cf0: 0x01e8 0x87ec 0x01e8 0x85d8 0x7363 0x616e 0x0000 0x0000
0x01e84d00: 0x0000 0x0019 0x0000 0x0000 0x0000 0x0000 0x01e8 0x5b7c
0x01e84d10: 0x01e8 0x87e8 0x01e8 0x8618 0x7365 0x7400 0x0000 0x0000
0x01e84d20: 0x0000 0x0019 0x0000 0x0000 0x0000 0x0000 0x01e8 0x5b7c
0x01e84d30: 0x01e8 0x87a8 0x01e8 0x8658 0x7370 0x6c69 0x7400 0x0000
0x01e84d40: 0x0000 0x0019 0x0000 0x0000 0x0000 0x0000 0x01e8 0x5b7c
0x01e84d50: 0x01e8 0x8854 0x01e8 0x8698 0x7374 0x7269 0x6e67 0x0000
0x01e84d60: 0x0000 0x0019 0x0000 0x0000 0x0000 0x0000 0x01e8 0x5b7c
0x01e84d70: 0x01e8 0x875c 0x01e8 0x86d8 0x7375 0x6273 0x7400 0x0000
0x01e84d80: 0x0000 0x0019 0x0000 0x0000 0x0000 0x0000 0x01e8 0x5b7c
0x01e84d90: 0x01e8 0x87c0 0x01e8 0x8718 0x7377 0x6974 0x6368 0x0000
```

fraction of RAW event

“Address” :

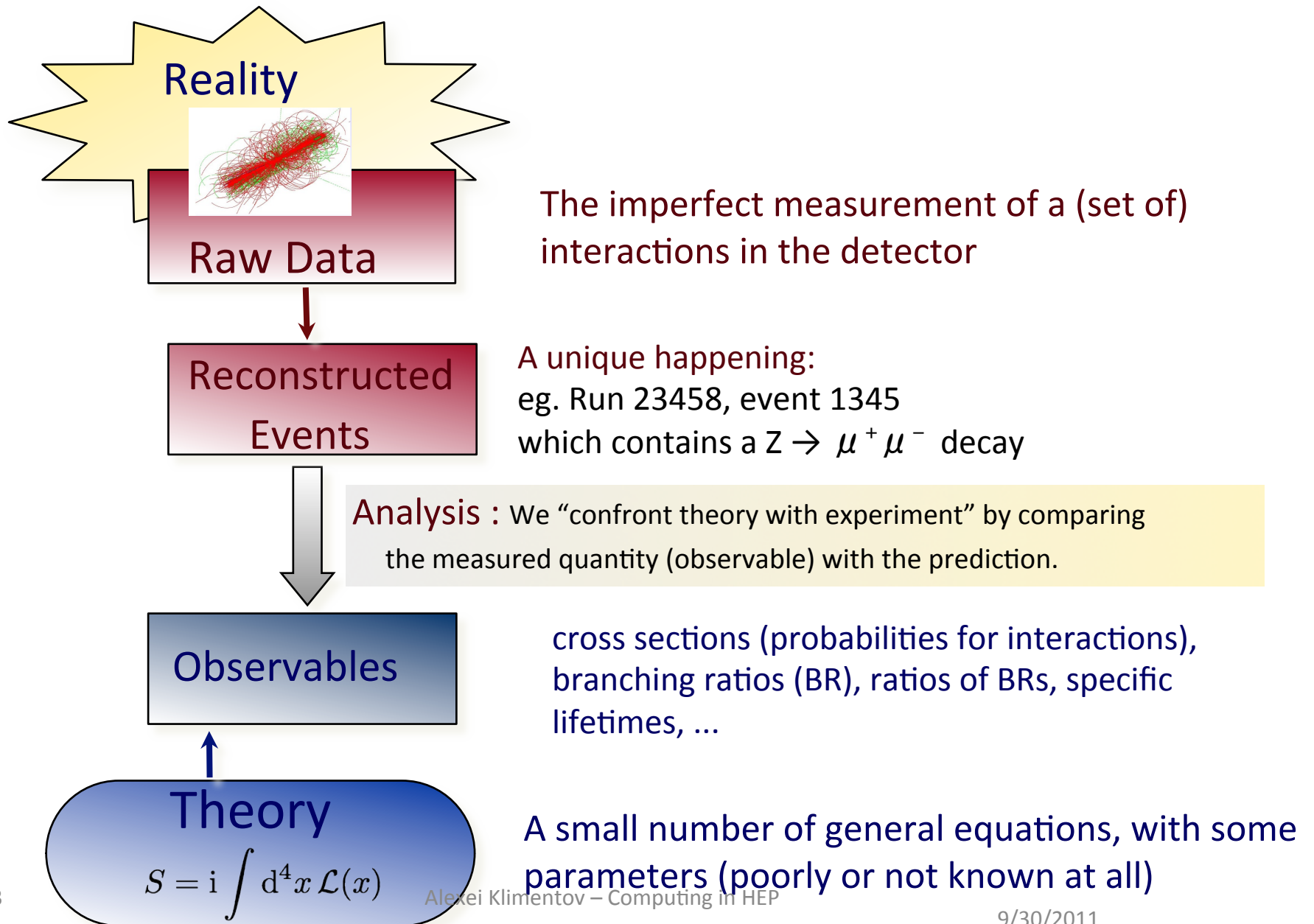
- which detector element took the reading

“Value(s)” :

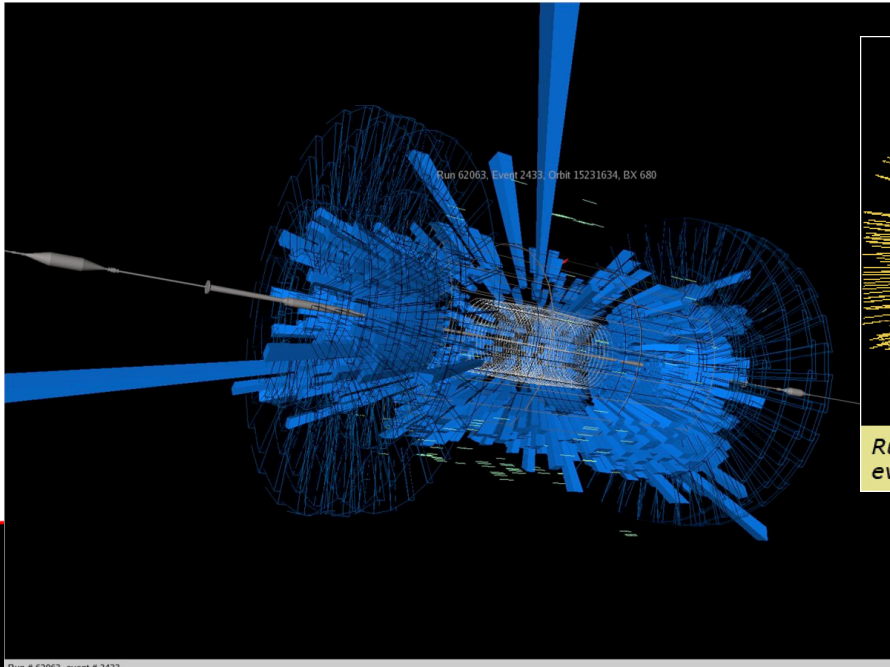
- what the electronics wrote out



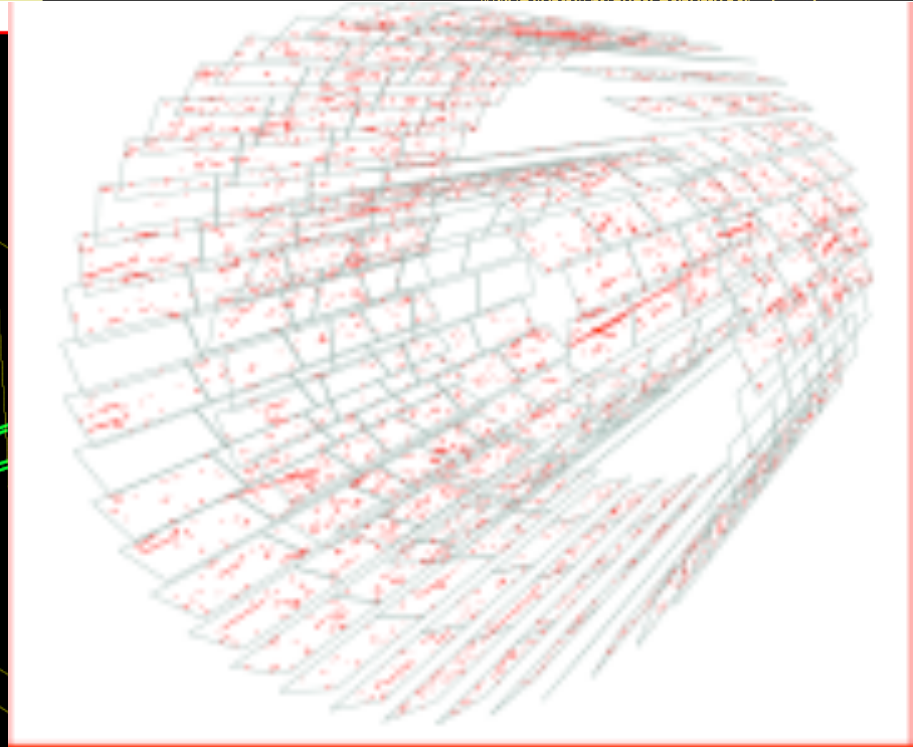
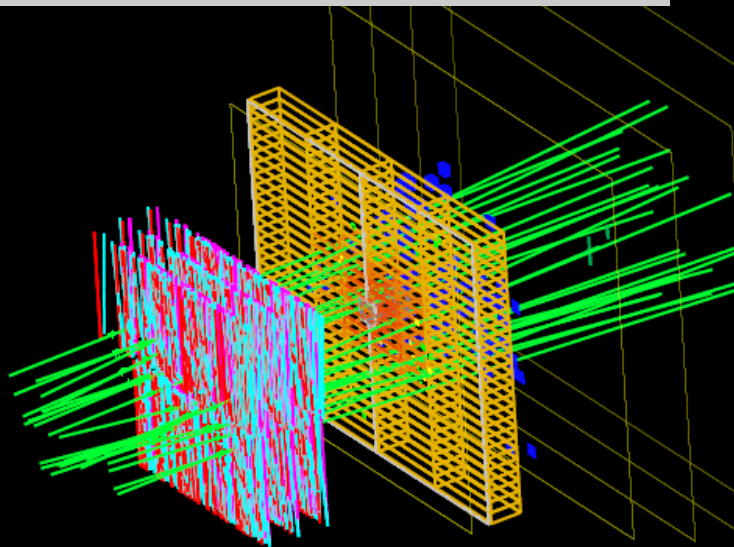
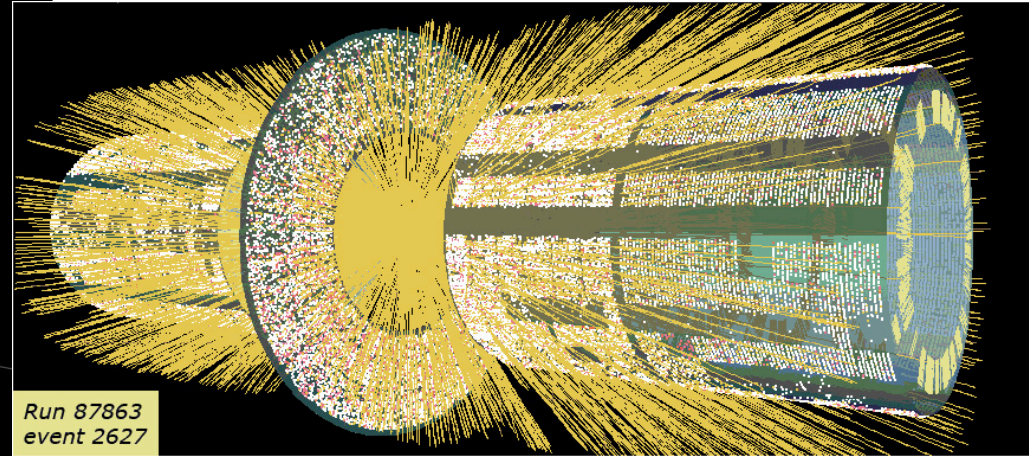
# Making the connection



# Reconstructed events



Run # 62063, event # 2433

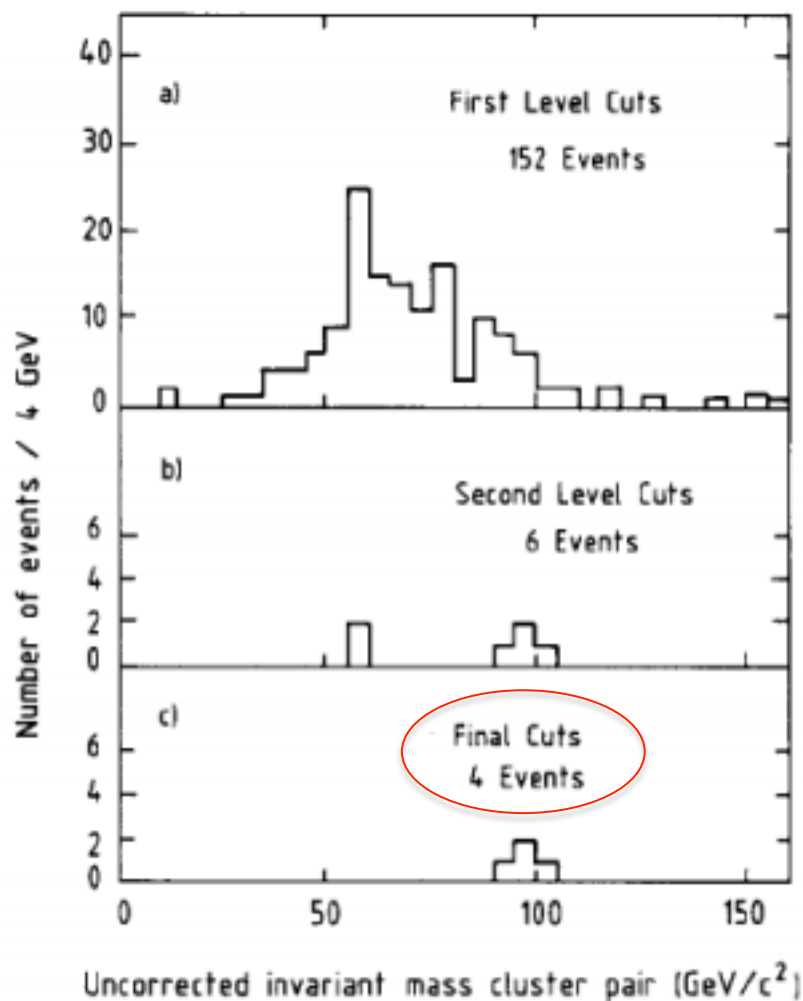


# UA1: observation of $Z \rightarrow e^+ e^-$

(May 1983)



The Nobel Prize in Physics 1984  
Carlo Rubbia, Simon van der Meer

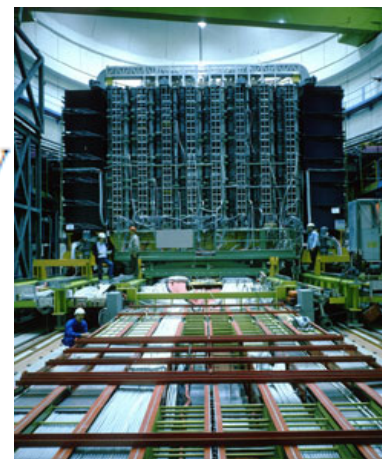


Two energy clusters ( $p_T > 25 \text{ GeV}$ )  
in electromagnetic calorimeters;  
energy leakage in hadronic calorimeters  
consistent with electrons

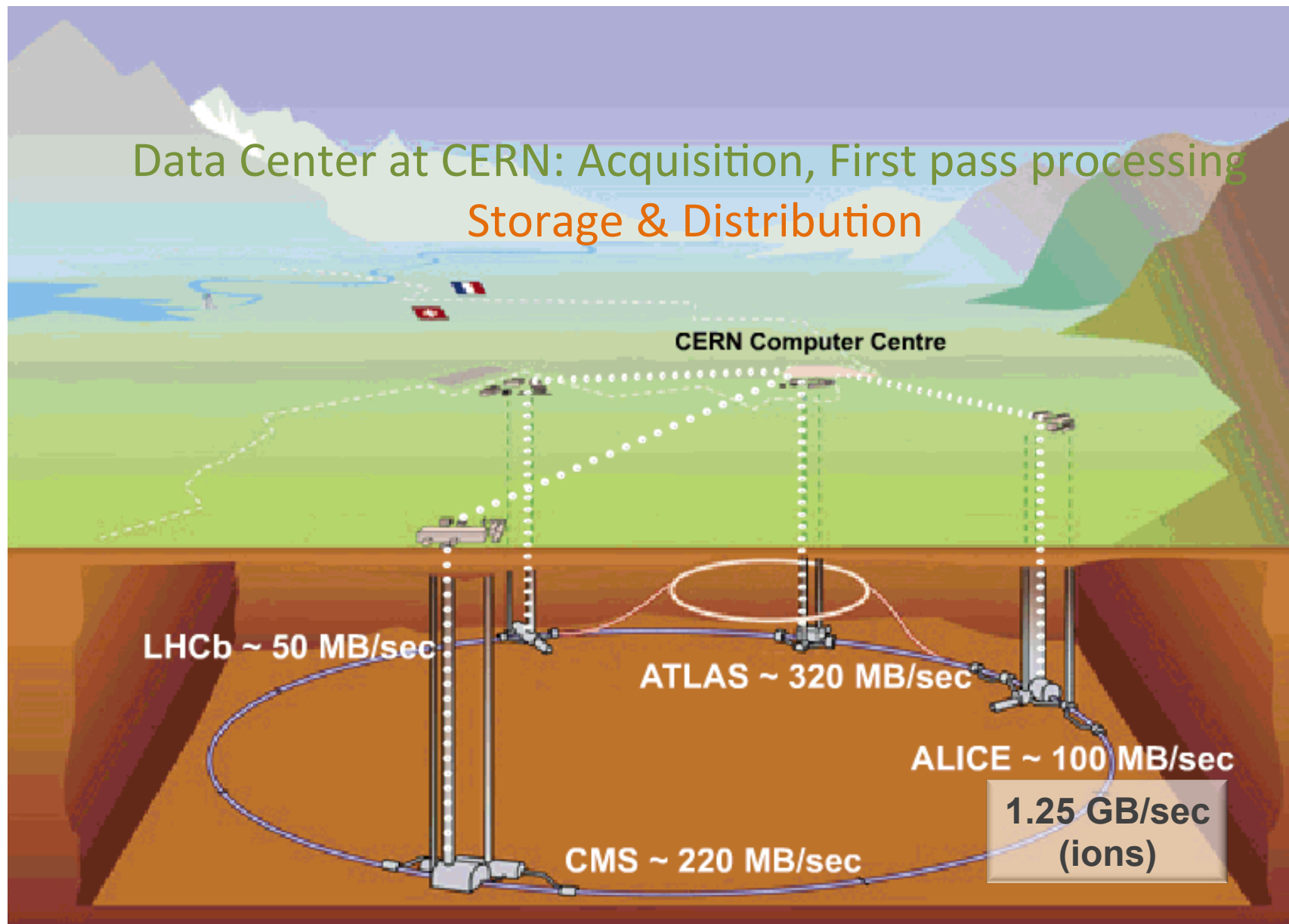
Isolated track with  $p_T > 7 \text{ GeV}$   
pointing to at least one cluster

Isolated track with  $p_T > 7 \text{ GeV}$   
pointing to both clusters

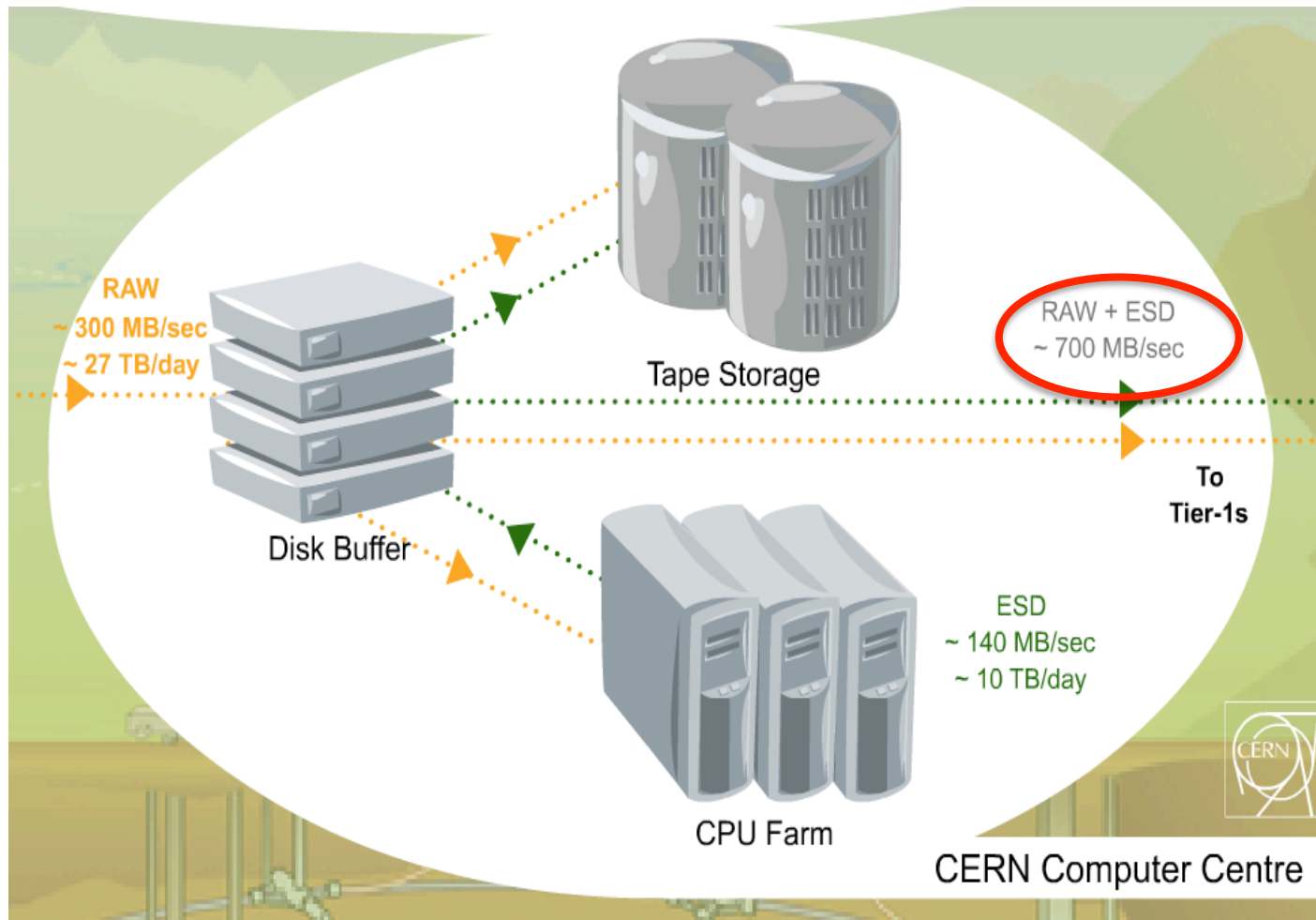
UA1 detector



## Data Center at CERN: Acquisition, First pass processing Storage & Distribution



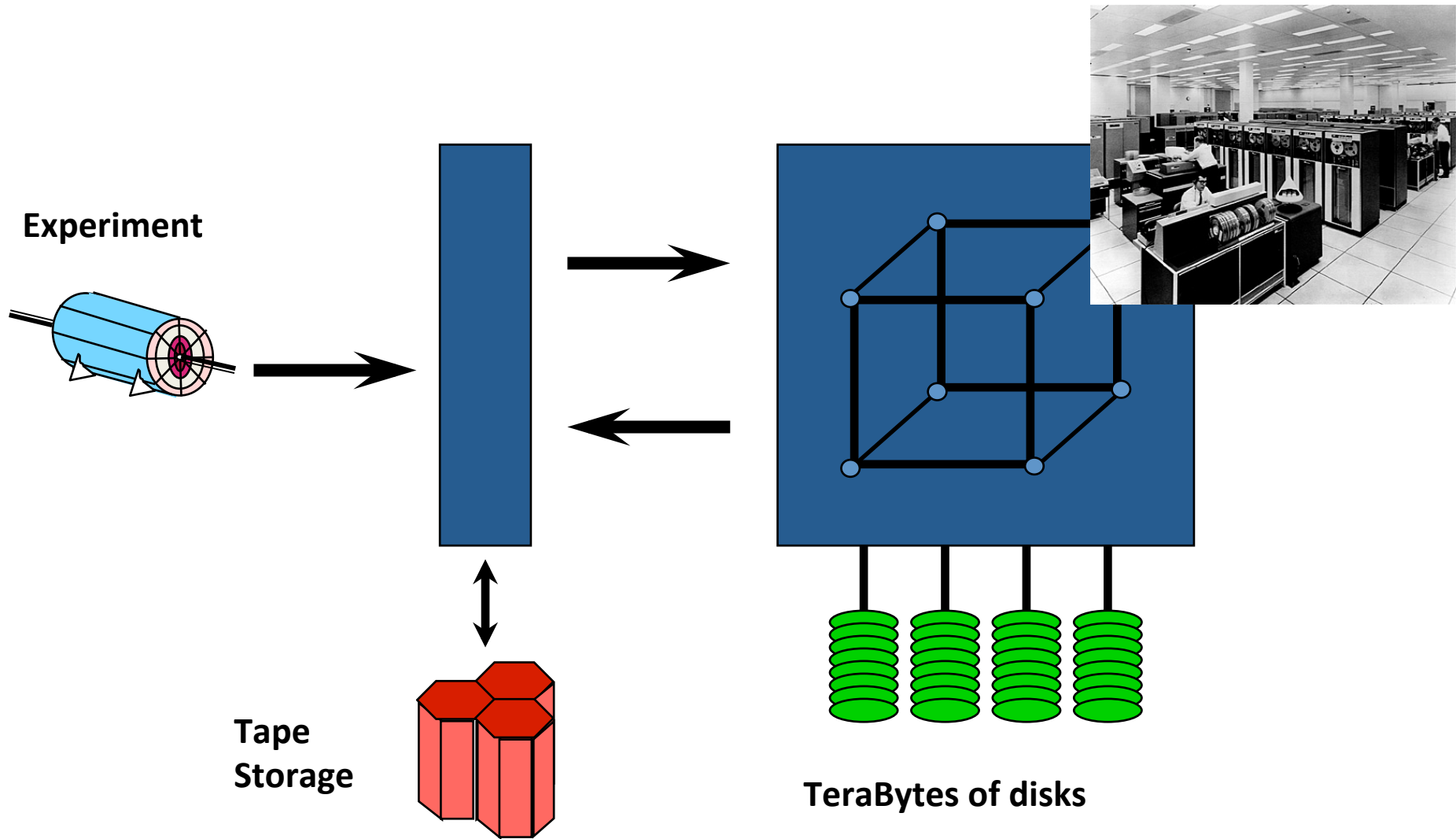
# Flow in and out of the center



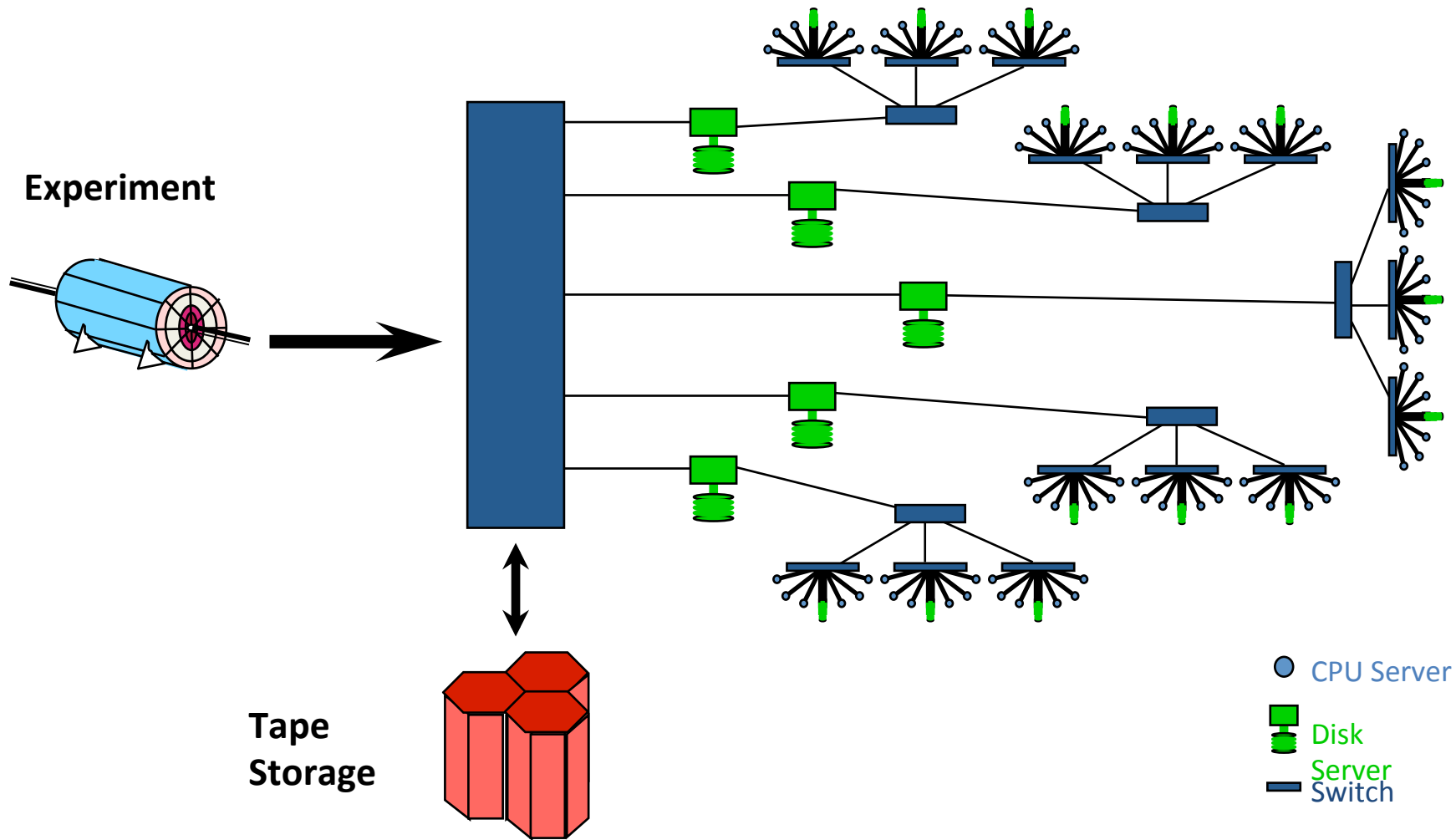
## *Act III*

# *Distributed Computing*

# Symmetric MultiProcessor Model



# Distributed Processing Model



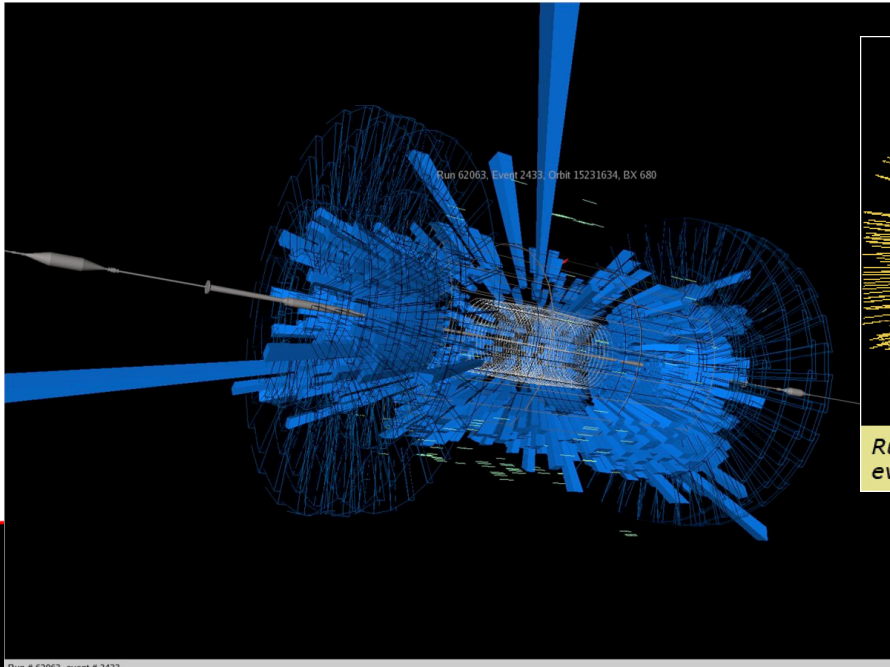


# Why Distributed Computing at LHC ?

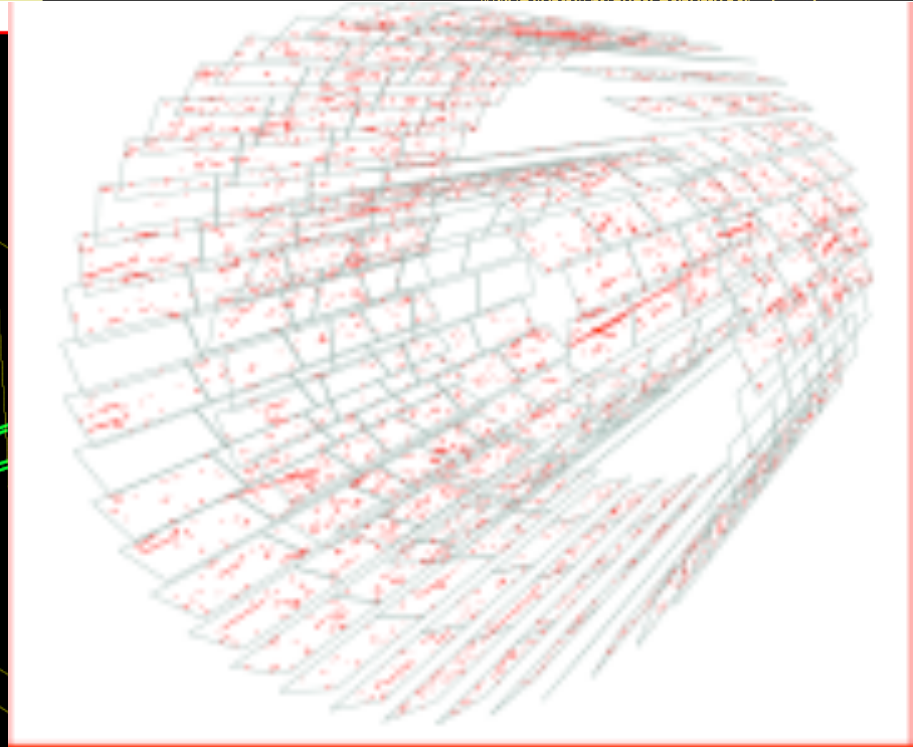
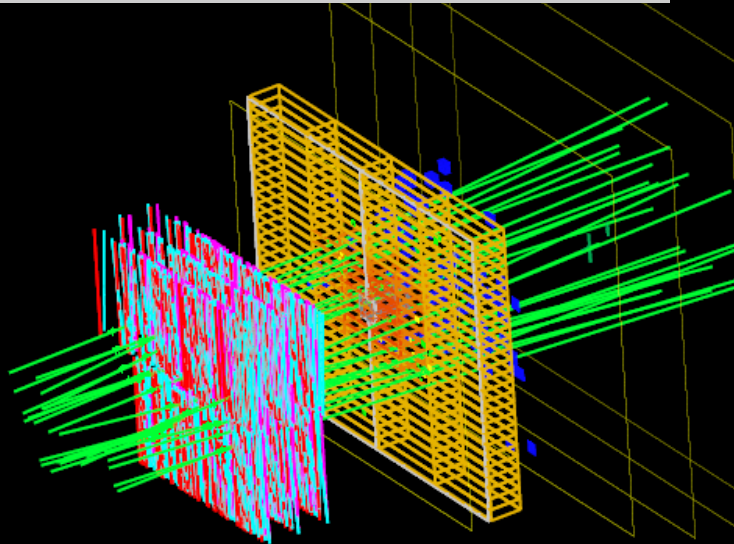
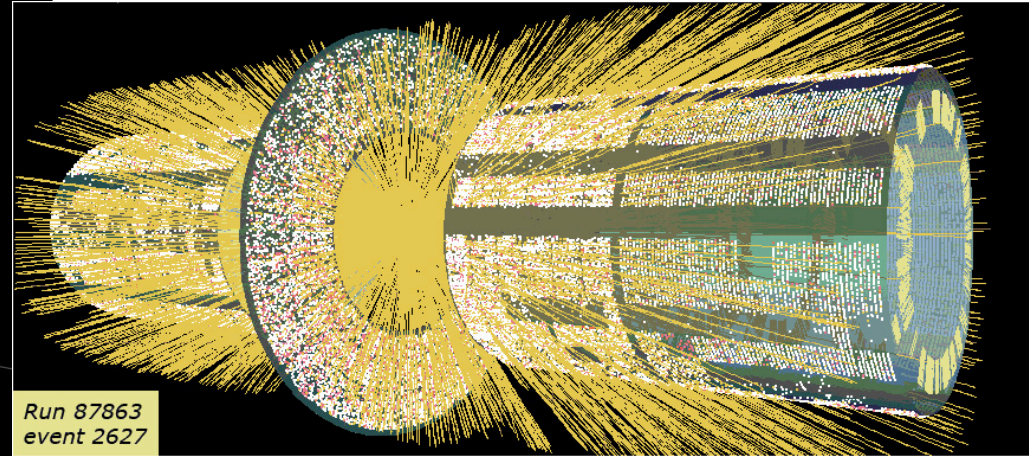
- From the start on it was clear that no center could provide ALL computing even for one LHC experiment
  - Buildings, Power, Cooling, Money .....
  - ATLAS Computing Requirements over time
    - 1995 : 100 TB disk space,  $10^7$  MIPS : Computing Technical Proposal
    - 2001: 1900 TB  $7*10^7$  MIPS : LHC Computing Review
    - 2007: 70000 TB  $55*10^7$  MIPS : Technical design report
    - **2010 LHC START**
    - 2011: 83000 TB  $61*10^7$  MIPS : recent request
- The High Energy Physics community is distributed and a most funding for computing is local
- Significant computing was available in many institutes
  - often shared with other research communities
- Both technical and political/financial reasons lead to the decision to build a distributed infrastructure for LHC computing

*MIPS – Million Instruction Per Second, used before ~2004/5, now it is replaced by SPECint benchmark*

# Complex events

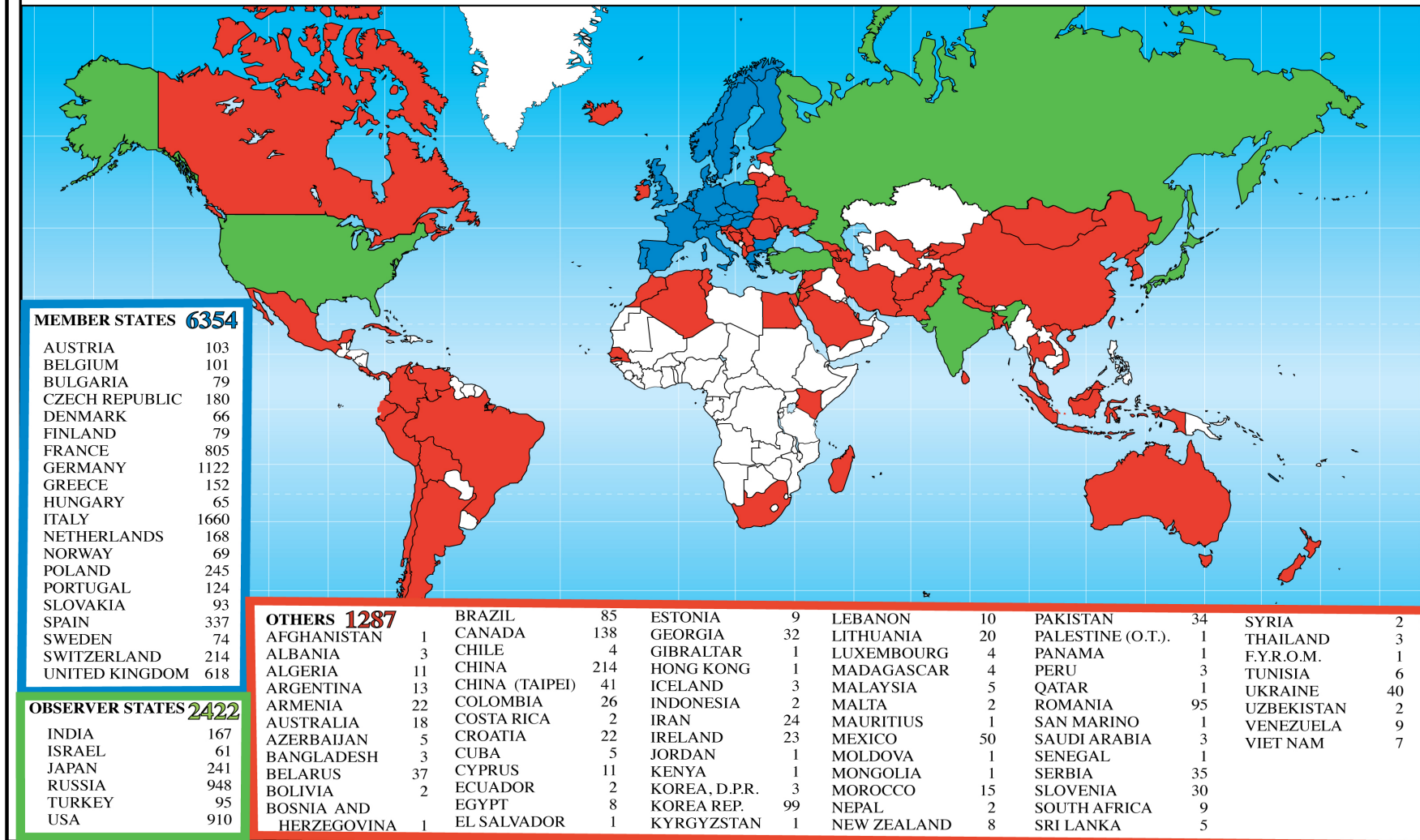


Run # 62063, event # 2433



# Complex Large Community

Distribution of All CERN Users by Nationality on 8 March 2011



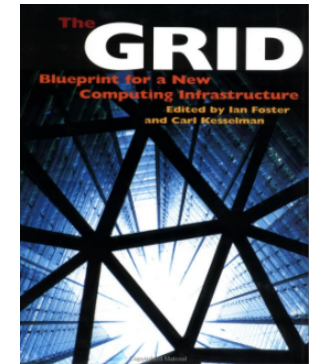
# Distributed Computing System for LHC.

- CERN's budget for physics computing was insufficient
  - Easy parallelism, use of simple PCs, availability of high bandwidth international networking .... make it **possible** to extend the distributed architecture to the wide area ....
- **AND**
  - The 6,000+ LHC collaborators are distributed across institutes all around the world with access to local computing facilities, ...
    - ... and funding agencies prefer to spend at home if they can
  - Mitigates the risks inherent in the computing being controlled at CERN, subject to the lab's funding priorities and with access and usage policies set by central groups within the experiments
- **ALSO**
  - Active participation in the LHC computing service gives the institute (not just the physicist) a continuing and key role in the data analysis
    - which is where the physics discovery happens
  - Encourages novel approaches to analysis .... ... and to the provision of computing resources

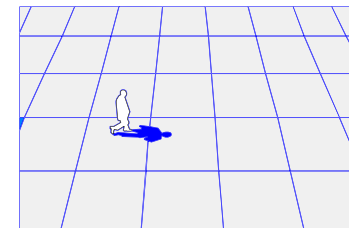
# *Act IV*

## *LHC Computing*

# What is a computing Grid ?



- There are many conflicting definitions.....
- **1998 The Grid by Ian Foster and Karl Kesselman**
  - Made the idea popular
- “coordinated resource **sharing** and problem solving in dynamic, **multi-institutional** virtual organizations. “
  - These are the people who started globus, the first grid middleware project
- From the user’s perspective:
  - I want to be able to use computing resources as I need them
    - I don’t care who owns resources, or where they are
    - Have to be secure
    - My programs have to run there
- The owners of computing resources (CPU cycles, storage, bandwidth)
  - My resources can be used by any authorized person (not for free)
    - Authorization is not tied to my administrative organization
- **NO centralized control of resources or users**



# LHC Computing Grid

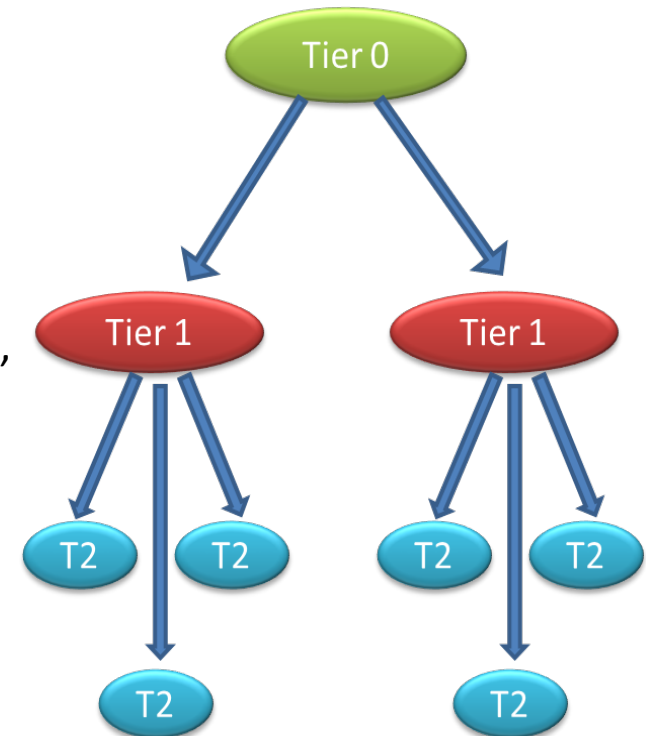


- Collaborating computing centres
- Interconnected with good networking
- Interfaces and protocols that enable the centres to advertise their resources and exchange data and work units
- Layers of software that hide all the complexity from the user
- So the end-user does not need to know where his data sits and where his jobs run
- The Grid does not itself impose a hierarchy or centralisation of services
- Application groups define **Virtual Organisations** that map users to subsets of the resources attached to the Grid

*More in the next lecture. J. Andreeva WLCG*

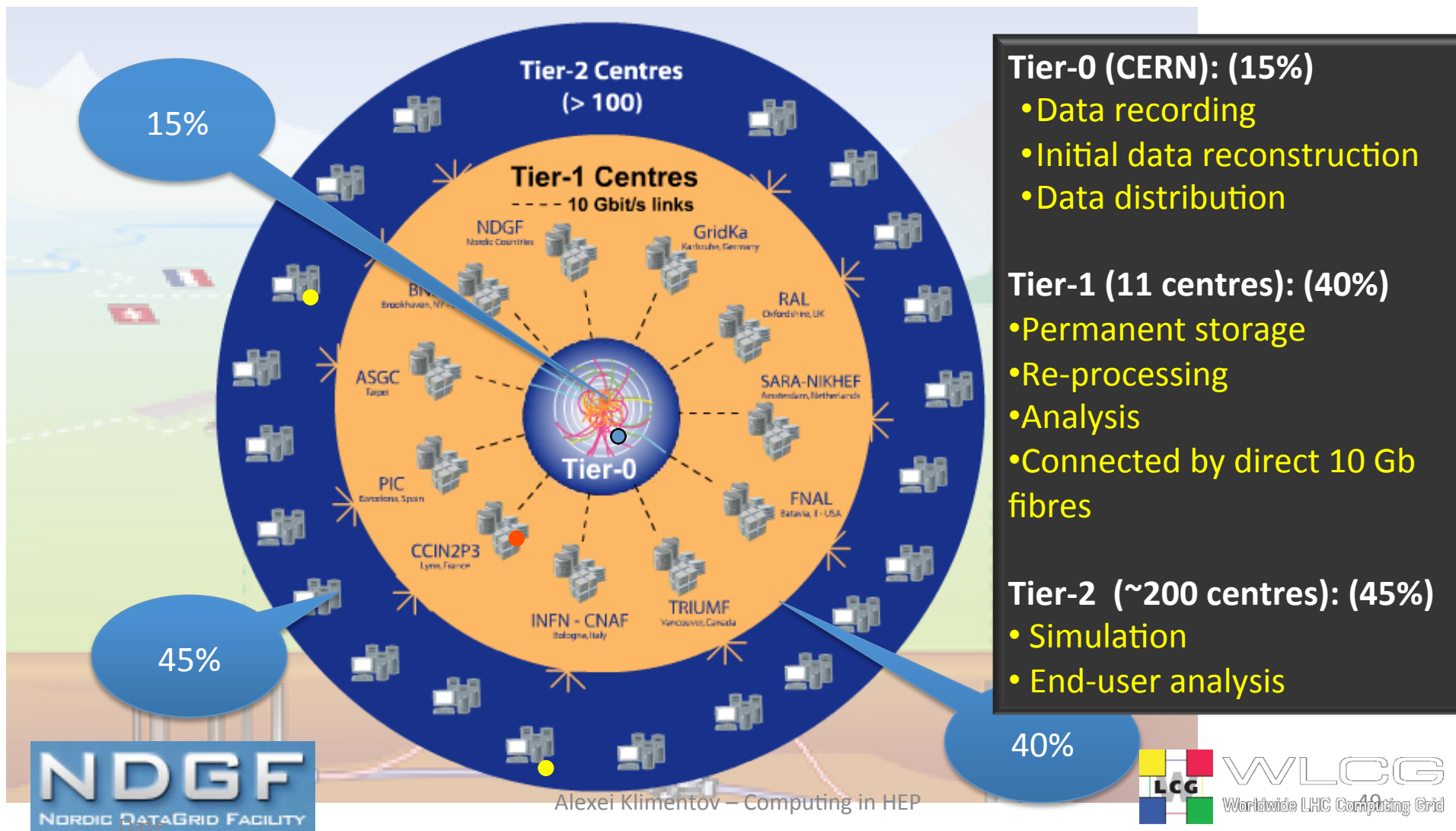
# MONARC Model

- 1998 – MONARC project
  - a distributed model
    - Integrate existing centres, department clusters, recognising that funding is easier if the equipment is installed at home
    - Devolution of control– local physics groups have more influence over how local resources are used, how the service evolves
  - a multi-Tier model
    - Enormous data volumes → looked after by a few (expensive) computing centres
    - Network costs favour regional data access
    - Simple model that HEP can develop and get into production ready for data in **2005**



Hierarchy in data placement







CERN



US-BNL



Amsterdam/NIKHEF-SARA



Taipei/ASGC



Bologna/CNAF



Ca-TRIUMF

### WLCG Collaboration Status

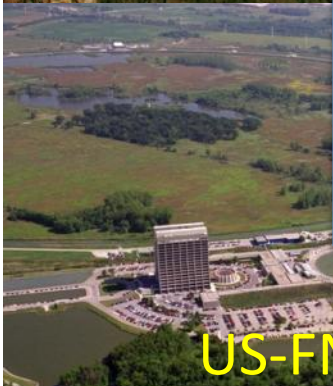
Tier 0; 11 Tier 1s; 68 Tier 2 federations (140 Tier 2 sites) + many T3 sites

Today we have 49 MoU signatories, representing 34 countries:

Australia, Austria, Belgium, Brazil, Canada, China, Czech Rep, Denmark, Estonia, Finland, France, Germany, Hungary, Italy, India, Israel, Japan, Rep. Korea, Netherlands, Norway, Pakistan, Poland, Portugal, Romania, Russia, Slovenia, Spain, Sweden, Switzerland, Taipei, Turkey, UK, Ukraine, USA.



NIDG



US-FNAL



De-FZK



Barcelona/PIC



Lyon/CCIN2P3

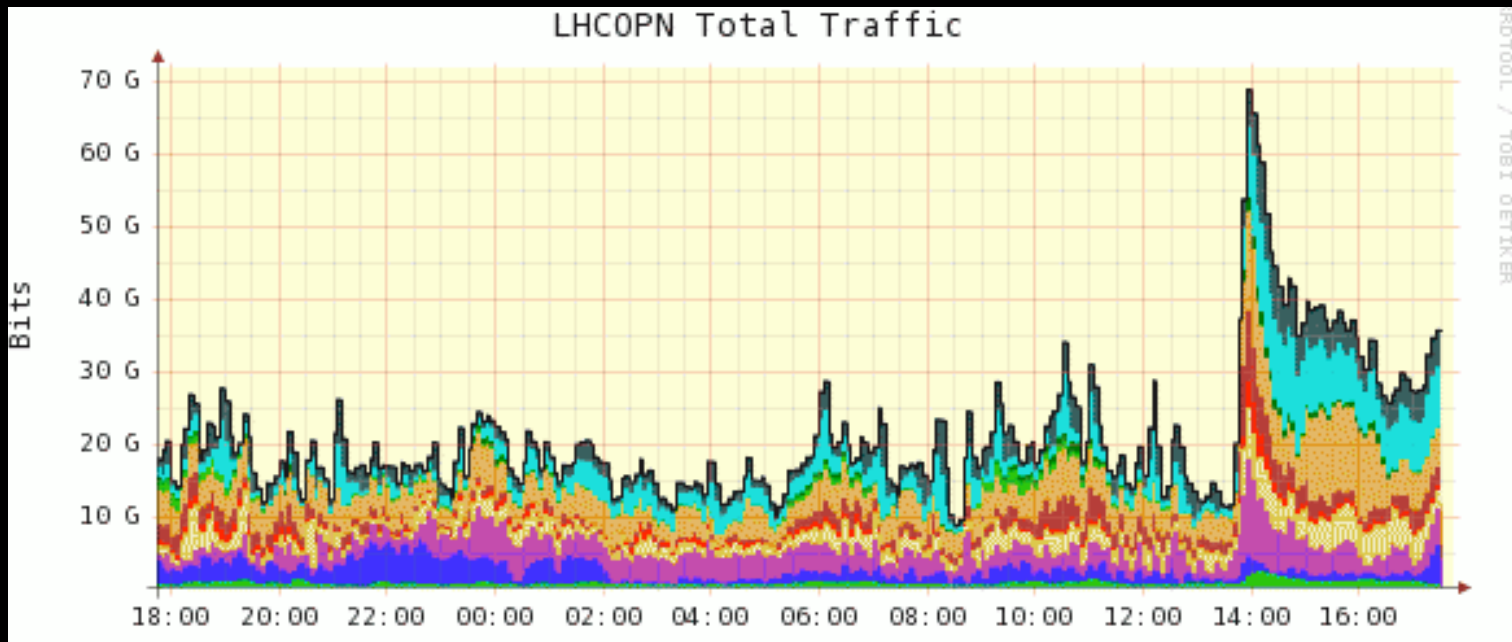


UK-RAL



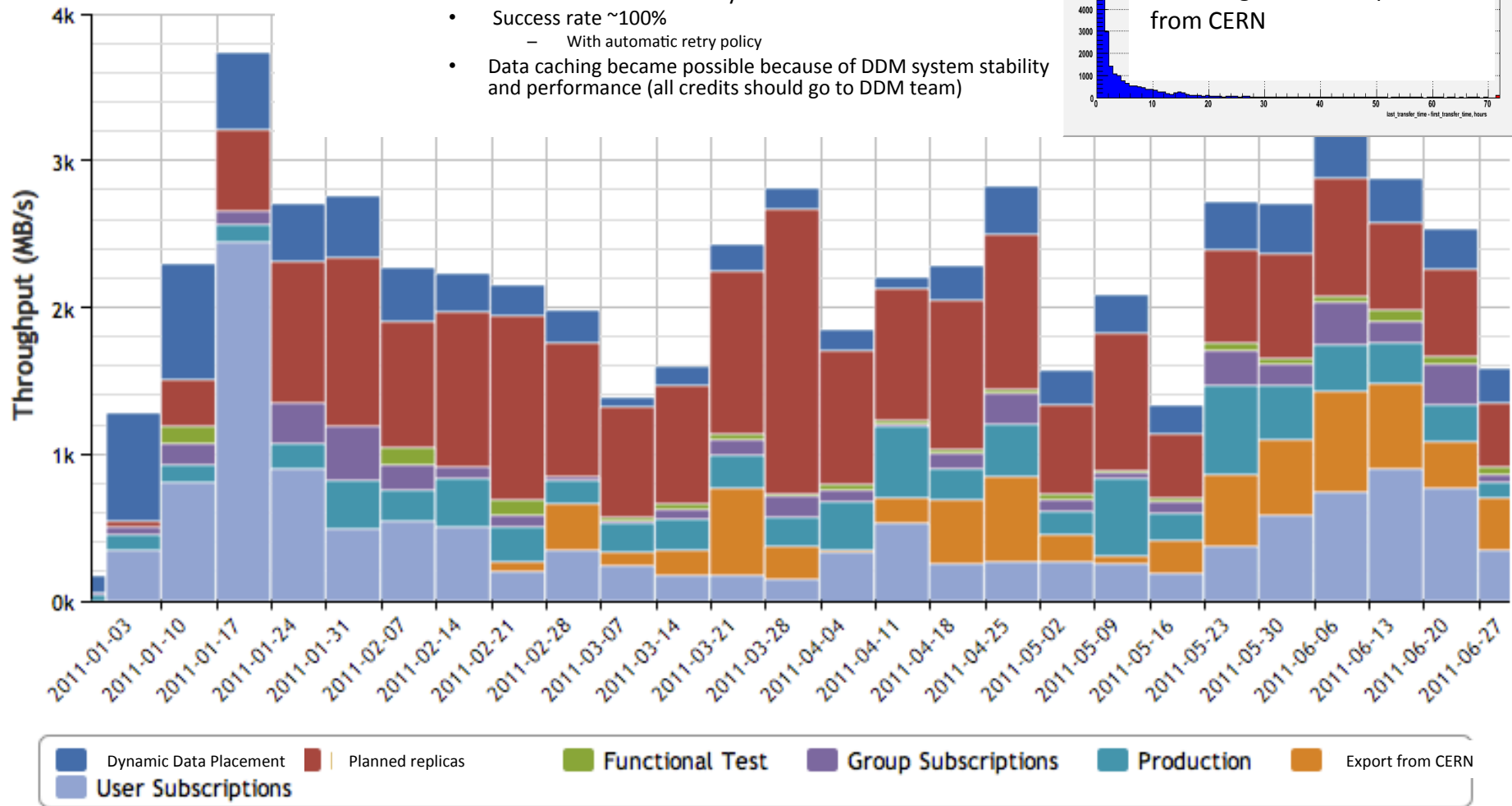
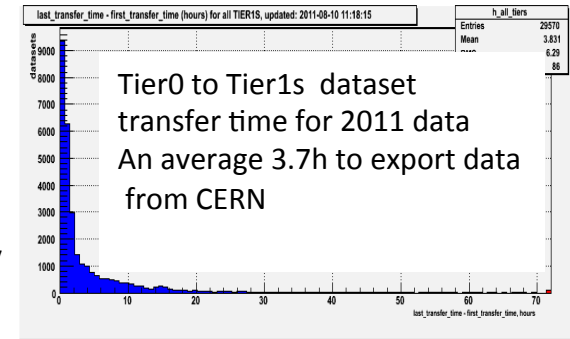
# 70 GB/s!

Traffic on OPN up to 70 Gb/s!  
- ATLAS reprocessing campaigns



# Data replication between ATLAS sites

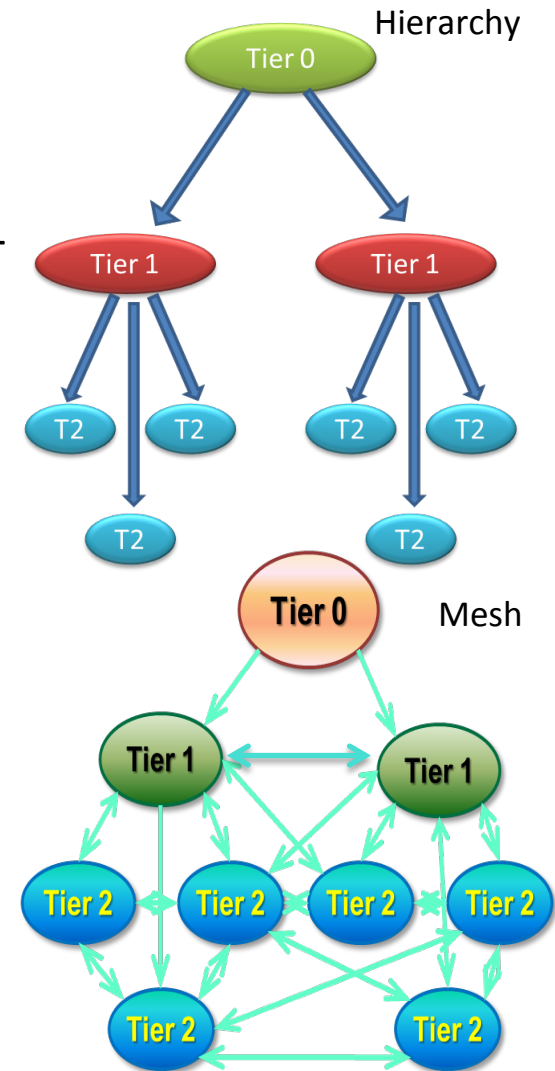
- 2011 Average Data Throughput per week in MB/s
  - Balance between planned replicas, dynamic placement and user requests
  - Excellent data transfer efficiency
    - Success rate ~100%
      - With automatic retry policy
    - Data caching became possible because of DDM system stability and performance (all credits should go to DDM team)



# Evolution of the Computing Model.

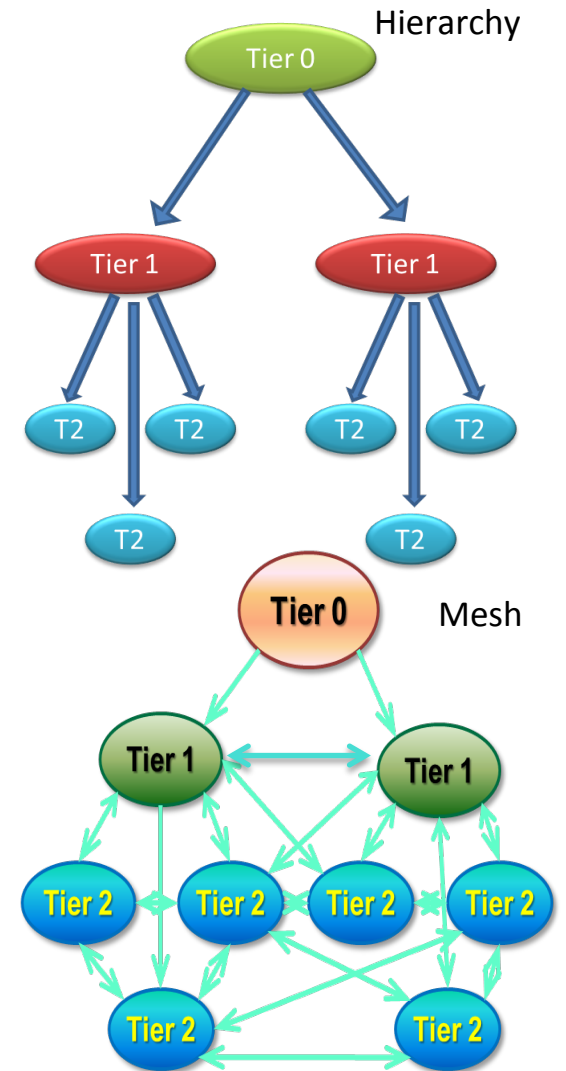
## Breaking Cloud Boundaries

- Network layout today differs from the Monarc model
  - Many T2s are connected very well with many T1s
  - Many T2s are not that well connected with their T1
- But breaking cloud boundaries is not THAT easy ...
  - Some links simply have limited bandwidth
    - In those cases, several hops will anyway be needed
  - Many more links to monitor (and fix)
    - A “one go” commissioning is not sufficient
  - Network between sites is not the only point here
    - Internal site network and storage configuration
    - Storage and Transfer Protocols overhead
- Start with “relaxing the model”
  - Rather than breaking



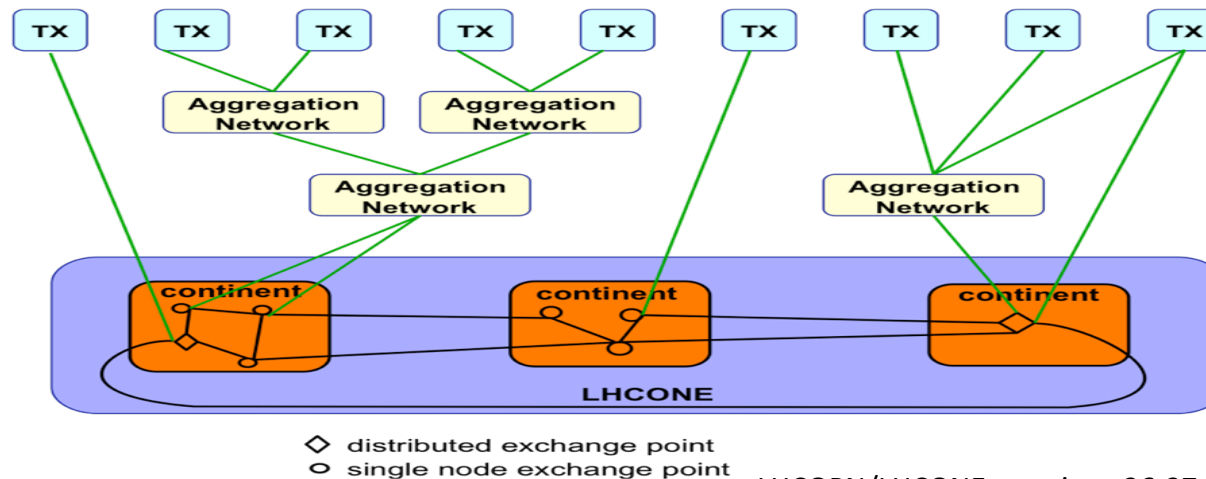
# LHC Open Network Environment. Data Models have changed

- Have moved/are moving away from the strict MONARC model. Evolution of computing models also require evolution of network infrastructure (Fisk-Bos document ~Jan 2010)
- 3 recurring themes:
  - Removing the hierarchy: Any site can replicate data from any other site
  - Dynamic data caching: Analysis sites receive datasets from any other site “on demand” based on usage pattern
    - Possibly in combination with pre-placement of data sets by centrally managed replication of whole datasets
  - Remote data access: local jobs accessing data stored at remote sites
    - Possibly in combination with local caching on a file or sub-file level
- Experiment-specific implementations of the above



# LHC Open Network Environment (LHCONE)

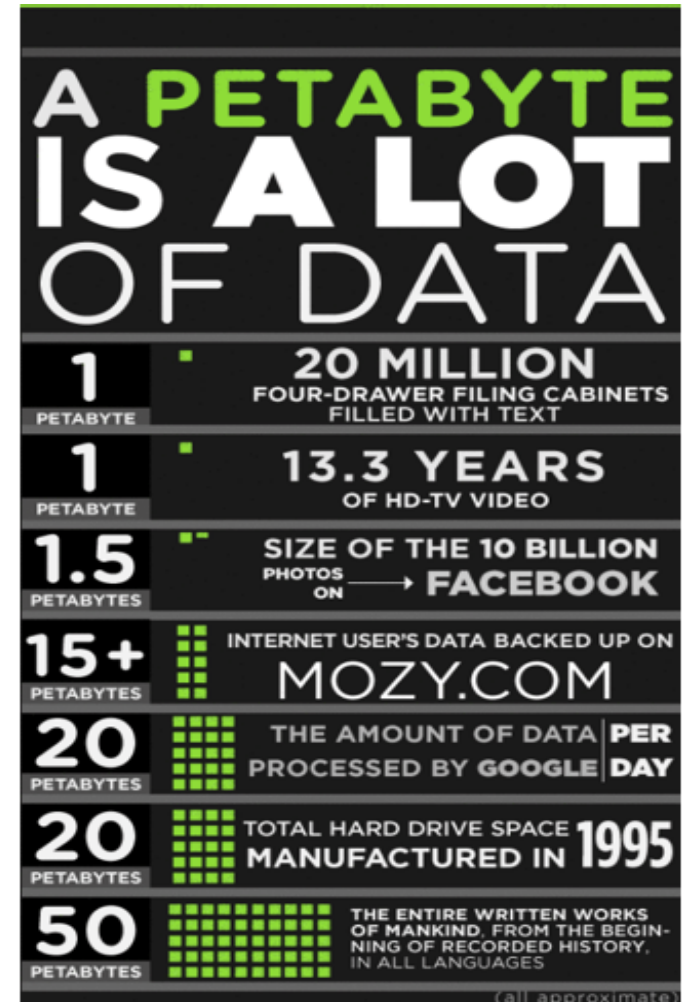
- Network providers, jointly working with the experiments, have proposed a new network model for supporting the LHC experiments, known as the LHC Open Network Environment (LHCONE).
- The goal of LHCONE is to provide a collection of access locations that are effectively entry points into a network that is reserved to the LHC T1/2/3 sites.
- LHCONE is not intended to supplant the LHCOPN but rather to complement it. It is not intended to let LHCONE carry Tier 0 (T0) – T1 traffic.
- Documentation available at
  - <http://lhcone.net>





# Data Volume and Data Storage

- 1 event size 1.6 MByte x Rate 400 Hz (Tycho's data rate ~100 Byte/h)
  - Taking into account 50% LHC duty cycle
    - Order of 4 PBytes of RAW data per year per experiment
    - Order of 10 PBytes of all data per year (RAW, reconstructed, meta-data)
    - 10 years of data taking : 100 PBYTES
- “New” physics is rare and interesting events are like single drop from the Jet d’Eau





1 in 10,000,000,000

Like looking for a single drop of water from the Jet d'Eau over 30 minutes



# Data Storage Challenges

- Technical Challenges
  - Lots of data and lots of capacity and lots of complexity
  - Combination of technologies : disk, tape, hierarchical storage (dcache, CASTOR,...)
  - Protecting data and making them accessible
- Data Analysis Challenges
  - Lots of people and lots of locations
  - Matching storage and processing resources
    - Mantra : “Jobs go to data” served us well for years

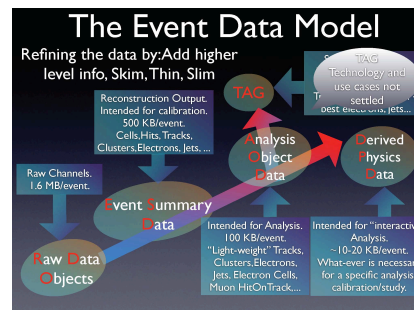
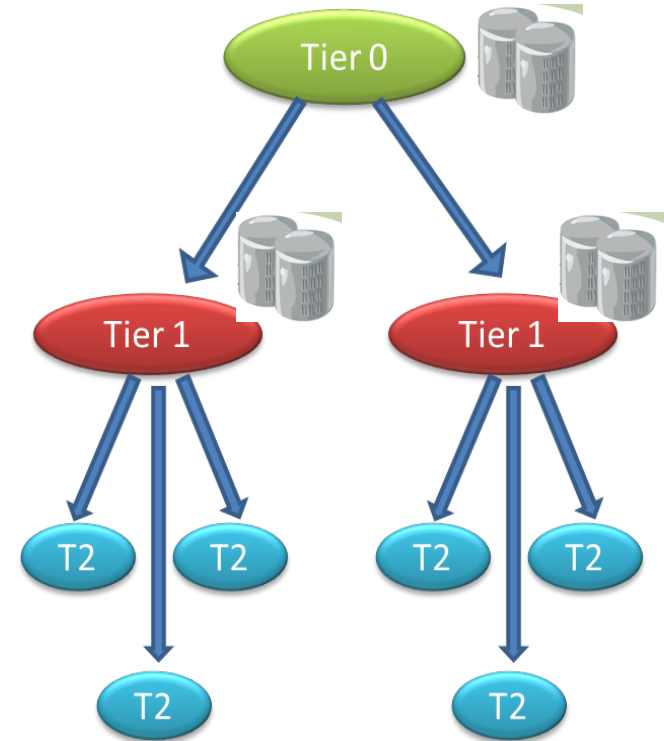
# Data Storage Evolution



Tapes originally were mounted by operators



Migrated to robots, but the latency for robots encourages careful planning of layout and placement what data is on disk and tape

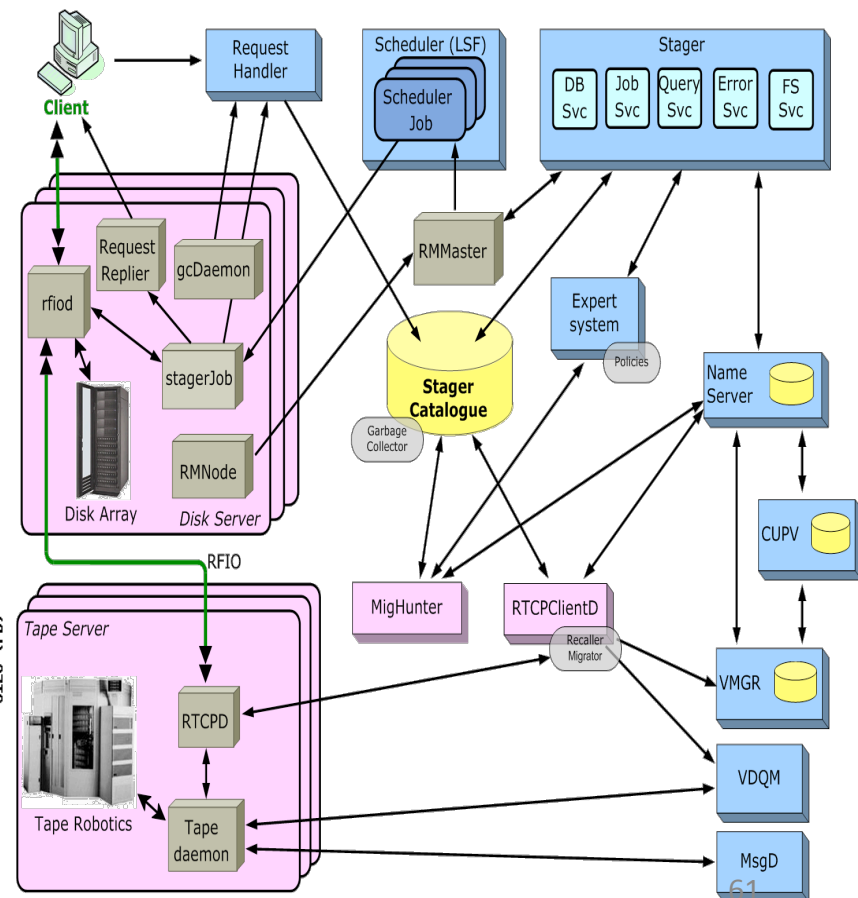
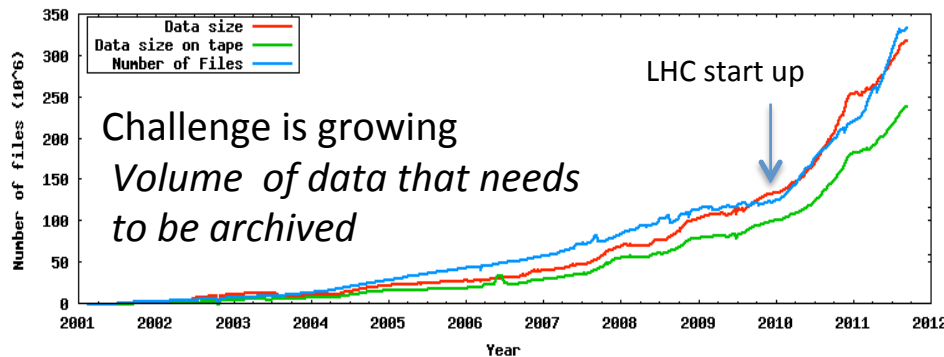


At the LHC most analysis work is conducted far away from the data archives

# CASTOR (cern.ch/castor)

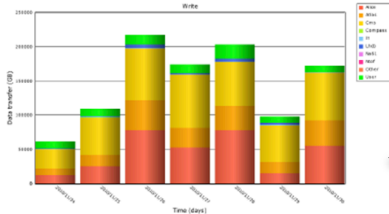
- The CERN Advanced STORAGE manager (CASTOR) is a hierarchical storage management system developed at CERN for physics data files.
- The design is based on a component architecture using a central database to save guard the state changes of CASTOR components. The 5 major functional modules are:
  - Stager: disk pool manager (allocating and reclaiming space; controlling client access; disk pool local catalogue)
  - Name Server: CASTOR name space (files and directories) including the corresponding file metadata (size, dates, checksum, ownership and ACLs, tape copy information). Command line tools modelled along Unix tools allows to manipulate the name space
  - Tape Infrastructure: under given conditions CASTOR saves files on tape. This is needed to provide data safety and to manage data storage larger than the available disks
  - Client: it allows you to upload, download, access and manage CASTOR data
  - Storage Resource Management: data access in LCG via the SRM protocol.

CASTOR at CERN statistics



# 1<sup>st</sup> year of LHC data

Tape recording: 220TB/day

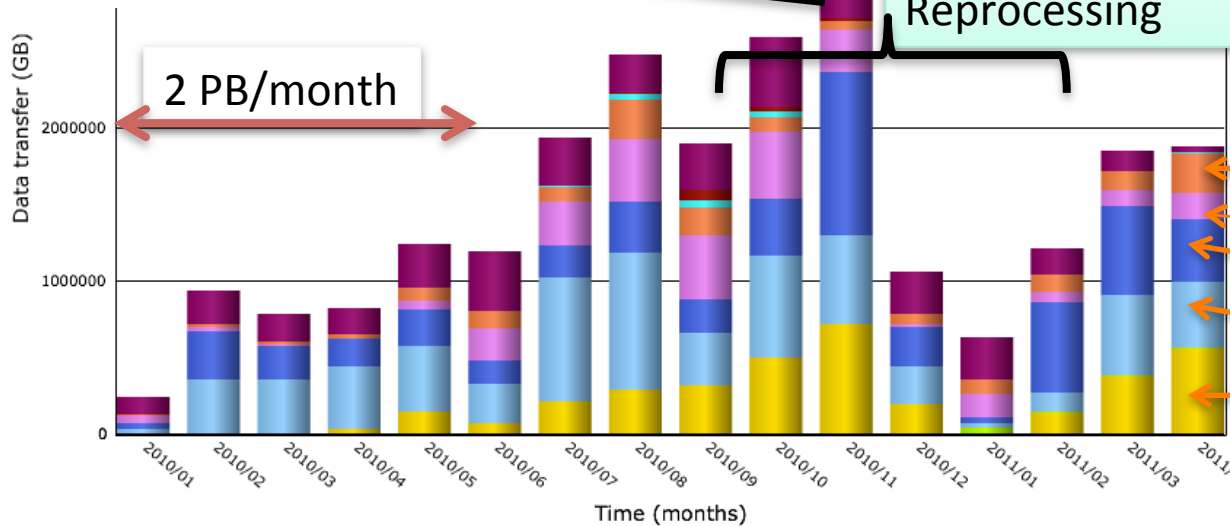


Data written to tape (GB/month): 2010-11

Stored ~ 15 PB in 2010

p-p data to tape at close to 2 PB/month

Peak rate: 225TB/day



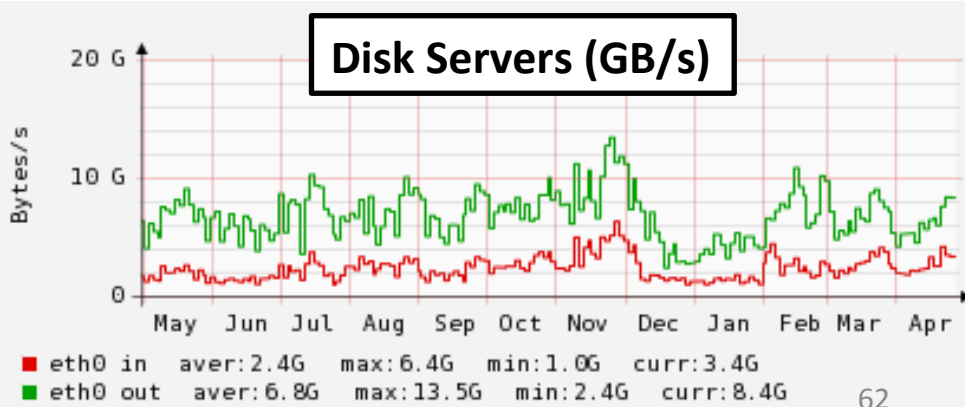
2 PB/month

2010 Reprocessing

HI

LHCb  
(COMPASS)  
CMS  
ATLAS  
ALICE

WLCG  
Worldwide LHC Computing Grid



Disk Servers (GB/s)

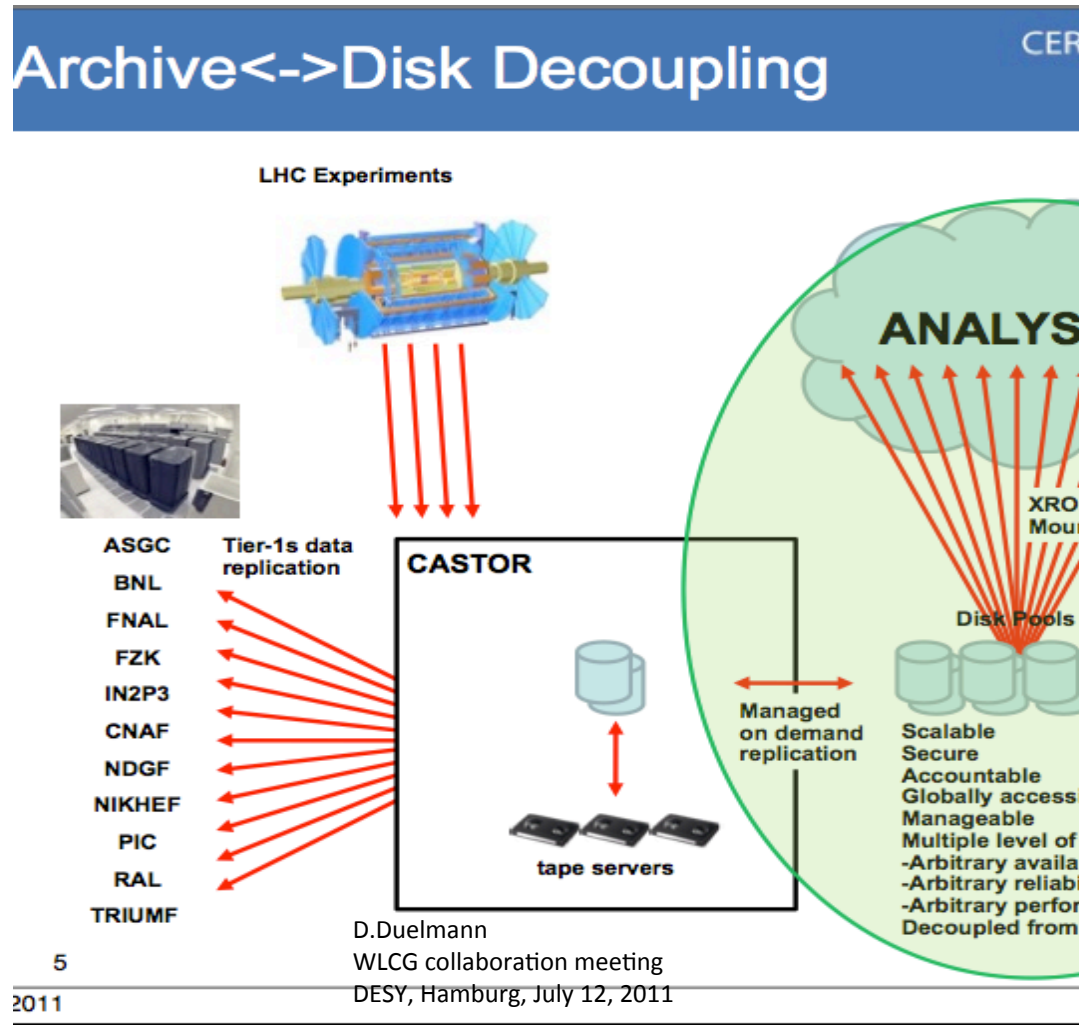
eth0 in aver:2.4G max:6.4G min:1.0G curr:3.4G  
eth0 out aver:6.8G max:13.5G min:2.4G curr:8.4G

## Tier 0 storage:

- Accepts data at average of 2.6 GB/s; peaks > 11 GB/s
- Serves data at average of 7 GB/s; peaks > 25 GB/s
- **CERN Tier 0 moves > 1 PB data per day**

# Data Storage Evolution.

## Archive<->Disk Decoupling



- Role separation and decoupling
  - Separating archive and disk activities
  - Separating Tier0 and analysis storage. Analysis storage can evolve more rapidly w/o posting risks for high priority T0 tasks
- CERN implementation – EOS
  - Mostly to eliminate CASTOR constraints
  - Intensive tests by ATLAS in 2010
  - ATLAS pools migration is done
    - Migrate ATLAS users, Grid storage to EOS

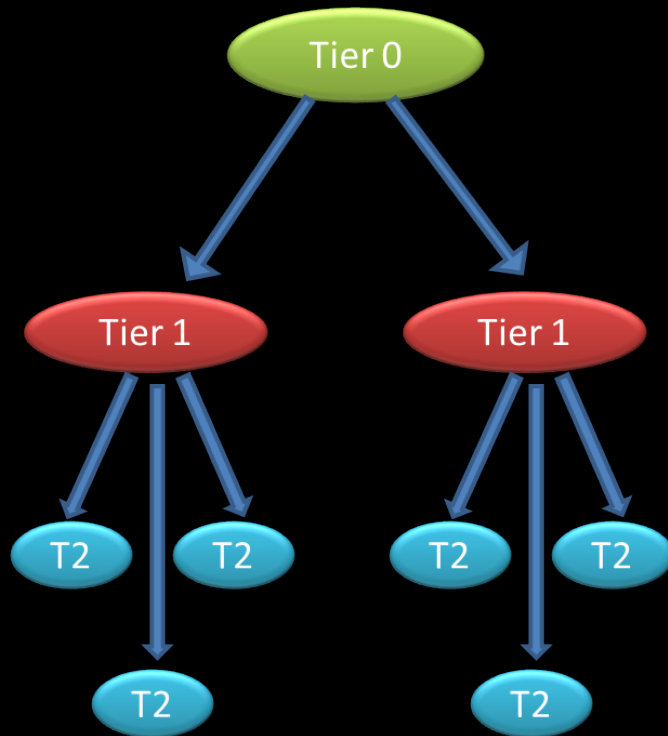
## *Final Act. What's Next ?*

Can we relax hierarchical model ?

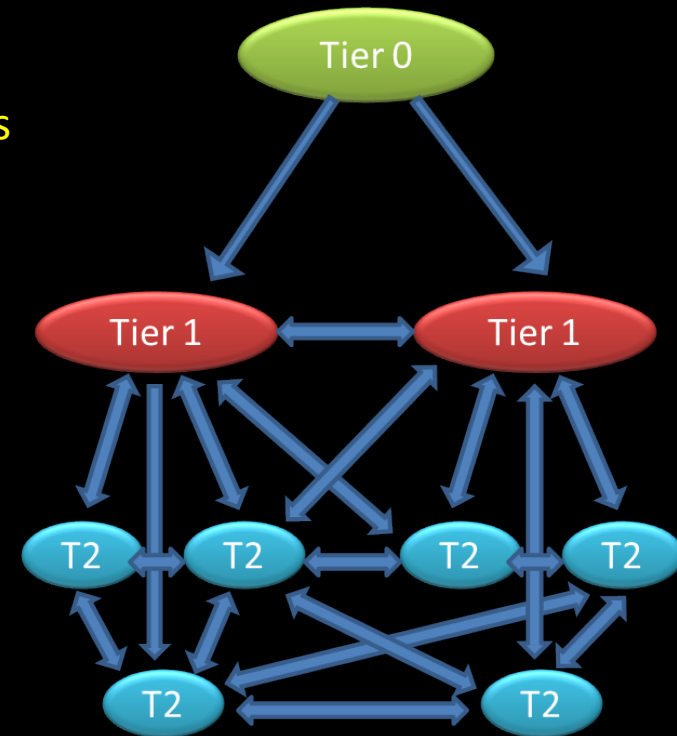
What is Cloud Computing ?



# Computing model evolution



Evolution of  
computing models



Wide area networks are very stable now

Hierarchy

Mesh

# Grids and Clouds.

- **At the end of the 90s High Energy Physics Community was a major computer user**
  - Having broken out of the cost/performance constraints of mainframes and minis at the beginning of the nineties
  - Having led the way in international high speed networking
  - Having exploited the power of the Web it had conceivedsurpassed only by a few other large sciences, the military and the spooks  
.... and the 2005 needs of LHC looked gigantic
- But the past decade has seen business and industry also exploit PC clusters, mass market disks, the Internet, the Web,..
- **the demand for computing power and storage has shot up**  
driven by **new applications**: search engines, web advertising, online commerce, digital libraries, interactive games, photo archives, social networking, ....
- **HEP is now a *relatively* small player**

- Google had 36 data centers in 2008\*\*
- LHC Grid has 200+



\*\*source: royal.pingdom.com

BUT

- One Google Data Center is estimated to cost ~\$600M
- An order of magnitude more than the new centre being planned at CERN

- *How many data centers does Google have? Nobody knows for sure, and the company isn't saying. The conventional wisdom is that Google has dozens of data centers. We're aware of at least 12 significant Google data center installations in the United States, with another three under construction. In Europe, Google is known to have equipment in at least five locations, with new data centers being built in two other venues.*



Google's data center at the Dalles on the Columbia river

- Microsoft's data center in Quincy, WA
  - 44K m<sup>2</sup> - 10 rooms each with up to 30K servers and storage for 6 trillion photos
- Yahoo, Amazon, IBM -- also building giant data centers



A look at the server area inside Microsoft's Quincy data center, courtesy of the BBC



- These major companies are expecting to build new markets for utility computing (clouds), software as a service (Google Apps), as well as absorbing the expansion of traditional computing services...
  - ... by offering **very cost-effective computing** –
    - Economies of scale (hardware, management, operation)
    - Efficient resource scheduling for high utilisation
    - Tax-efficient locations
    - Cheap and green energy

# Looking to the Future.



## Grids and Clouds

# Looking to the Future



Grids and  
Clouds

# Cloud Computing

- **IaaS (infrastructure as a service) paradigm, aka “Cloud computing” is a combination of and improvements in several key technologies in the enterprise level computing:**
  - *Advent of cheap multicore CPUs*
  - *Increase in power and cooling efficiencies (at components and facilities level)*
  - *Cheap high density storage*
  - *Strong security (computer security and facilities)*
  - *Improved reliability and fault tolerance of hardware and infrastructure*
  - *Virtualization*
  - *Convenient user and application interfaces*
  - *Cloud management and monitoring software*

# Cloud Technologies and ATLAS

- **Commercial clouds (Amazon, EC2, etc) as an additional resource for ATLAS**
- **Academic clouds (Magellan, etc)**
  - <http://magellan.nersc.gov/?p=878>
- **Adoption of virtualization by ATLAS computing facilities and possible conversion to cloud model of providing computing resources**

Main question - how useful are current cloud computing technologies for ATLAS Computing?



# EC2 Costs

Part of actual EC2 Bill for Panda related activities Apr. 2011

|  |                 | Totals   |
|--|-----------------|----------|
| Amazon Elastic Compute Cloud                                       |                 |          |
| US East (Northern Virginia) Region                                 |                 |          |
| Elastic IP Addresses   |                 |          |
| \$0.01 per non-attached Elastic IP address                         | 744 Hrs         | 7.44     |
| per complete hour  | »               | 7.44     |
| Amazon Simple Storage Service                                      |                 |          |
| US Standard Region   |                 |          |
| \$0.140 per GB - first 1 TB / month of storage used                | 387.825 GB-Mo   | 54.30    |
| \$0.01 per 1,000 PUT, COPY, POST, or LIST requests                 | 4,759 Requests  | 0.05     |
| \$0.01 per 10,000 GET and all other requests                       | 62,627 Requests | 0.06     |
|  | »               | 54.41    |
| Amazon Simple Notification Service                                 |                 |          |
|  | »               | 0.00     |
| Amazon Virtual Private Cloud                                       |                 |          |
|  | »               | 0.00     |
| AWS Data Transfer (excluding Amazon CloudFront)                    |                 |          |
| \$0.100 per GB - data transfer in per month                        | 129.790 GB      | 12.98    |
| \$0.000 per GB - first 1 GB of data transferred out per month      | 1.000 GB        | 0.00     |
| \$0.150 per GB - up to 10 TB / month data transfer out             | 2,533.798 GB    | 380.07   |
|  |                 | 393.05   |
| <b>Bill Summary</b>  |                 |          |
| Usage charges and monthly recurring fees during this billing cycle |                 | \$454.90 |

Panda server IP address

Panda monitor data

Panda monitor data transfers

# Cloud Resources

- **Commercial – pricey but big**
  - **Amazon EC2, Rackspace**
    - *CMS experience : 8 times more to do simulation on the cloud vs Tier-2*
    - *EC2 is costly so far, however with EC2 spot instances, the game might change*
- **Community clouds – free but small**
  - **Magellan LBNL, ANL**
  - **StratusLab reference cloud – 3M Euro EU project**
- **LXCLOUD@CERN**
  - **CERN will be a cloud provider ?**
- **National clouds or/and science cloud**
  - **large community resource with pledges to VOs**
  - **VOs must prove they can use it**
  - **Extra-Tier-1s vs national cloud**

# Looking to the Future. Grid vs Cloud

Les Robertson :

Clouds aim at efficient sharing of the hardware

- ✓ low-level execution environment, Isolation between users
- ✓ Operated as a homogeneous, single-management domain
- ✓ Straight-forward i/o and storage
- ✓ Expose only a high-level view of the environment - scheduling, data placement, performance issues are hidden from the application and the user



*After more than a decade of distributed computing (called Grid), an alternative approach - the centralized computing - is promoted by the industry (called Cloud) Major Companies (Microsoft, Google, IBM) are expecting to build new markets for cloud computing by offering very cost effective computing*

# Summary

- Grids are all about sharing.
  - groups distributed around the world can pool their computing resources
  - large centres and small centres can all contribute
  - users everywhere can get equal access to data and resources
- Grids are also flexible
  - place the computing facilities in the most effective and efficient places
  - exploiting funding wherever it is provided
- HEP and others have shown that
  - grids can support computational and storage resources on a massive scale
  - that can be operated around the clock
  - running hundreds of thousands of jobs every day
- The grid model has stimulated high energy physics to organise its computing
  - in a widely distributed way
  - building a collaboration involving directly a large fraction of the LHC members and their institutes
- This is the workhorse for production data handling for many years and as such must be maintained and developed through the first waves of data taking

# Summary

...BUT

- the landscape has changed dramatically over the past decade
  - The Web, the Internet, powerful PCs, broadband to the home, ...
    - have stimulated the development of new applications that generate a massive demand for computing remote from the user
    - .... that is being met by giant, efficient facilities deployed around the world
    - .... and creates a market for new technologies capable of operating on a scale equivalent to that of HEP
- Whether or not commercial clouds become cost-effective for HEP data handling is only a financial and funding-agency issue

BUT

- Exploiting the associated technologies is an obligation

## Could there be a revolution here for physics computing

# Computing in High Energy Physics (last page)

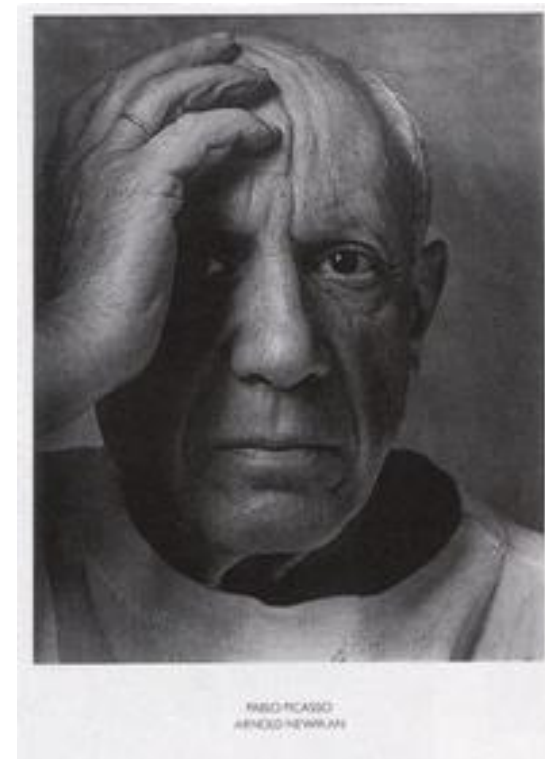
Computers are useless – they can only give you answers

– Pablo Picasso



9/30/2011

78



uting  
HEP