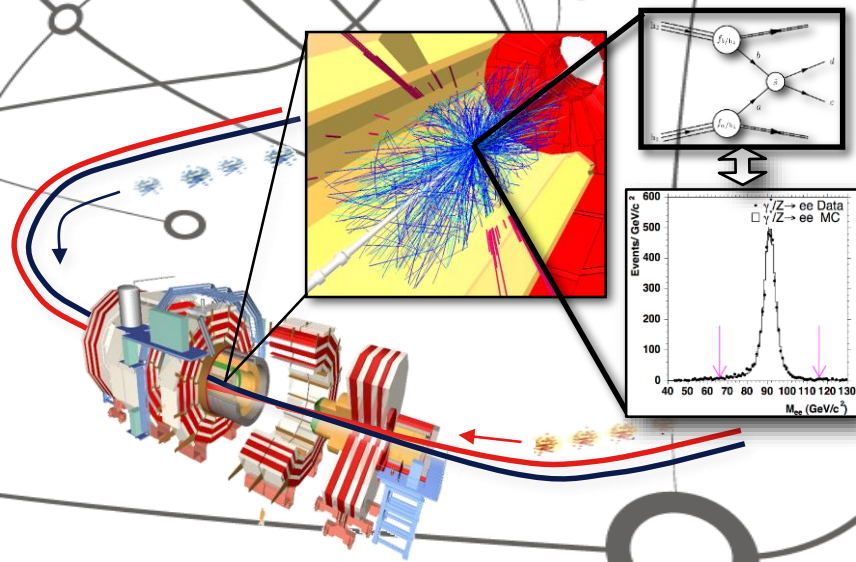




# Big Data Analytics and the LHC



Maria Girone  
CERN openlab CTO



CERNopenlab

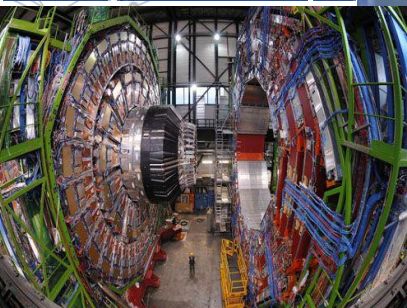


European Organization for Nuclear Research, founded in 1954 by 12 European countries

*"Science for Peace"*

~ 2300 staff  
~ 1600 other paid personnel  
~ 10500 scientific users

Budget (2014) ~1000 MCHF



**Member States:** Austria, Belgium, Bulgaria, the Czech Republic, Denmark, Finland, France, Germany, Greece, Hungary, Israel, Italy, the Netherlands, Norway, Poland, Portugal, Slovakia, Spain, Sweden, Switzerland and the United Kingdom

**Candidate for Accession:** Romania

**Associate Member in Pre-Stage to Membership:** Serbia

**Applicant States for Membership or Associate Membership:** Brazil, Croatia, Cyprus, Pakistan, Russia, Slovenia, Turkey, Ukraine

**Observers to Council:** India, Japan, Russia, Turkey, United States of America; European Commission and UNESCO



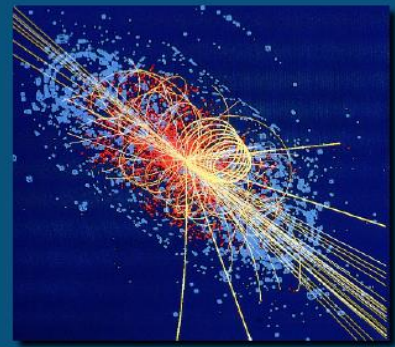
→ Interfacing between fundamental science and key technological developments



→ CERN Technologies and applications



Accelerating particle beams



Detecting particles

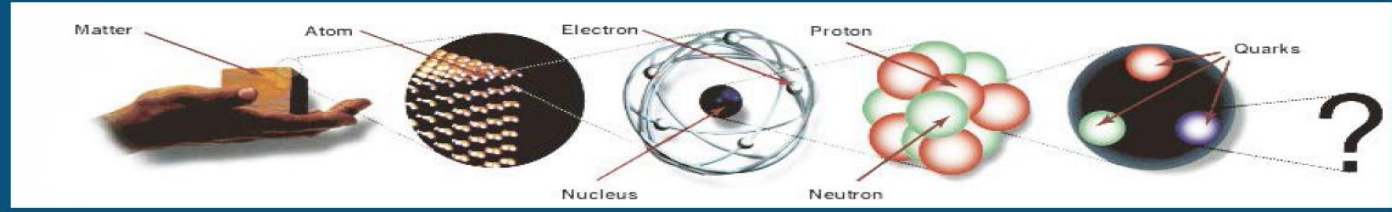


Large-scale computing (Grid)

# Number 1: "Particle physics"

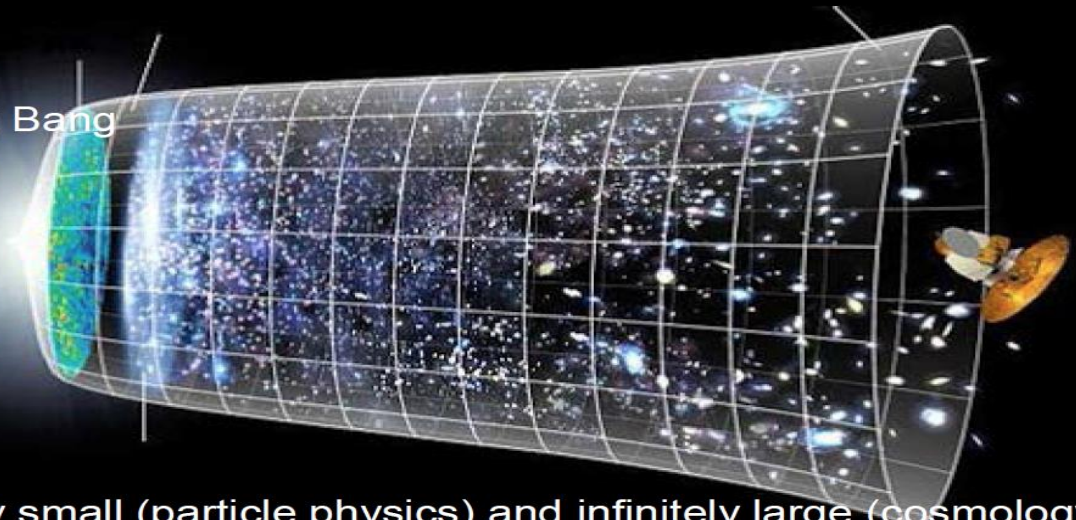
## ○ Quest to understand:

- Fundamental constituents of matter - Matter particles
- Interactions with which particles act on each other - Interactions
- Particles propagating the interactions - Messenger particles



## ○ Ultimately describe:

- Birth of the Universe, the Big Bang
- Passed and future Evolution



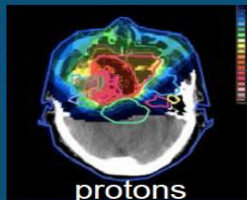
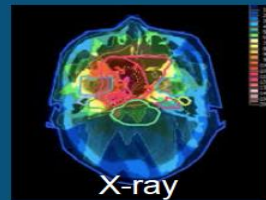
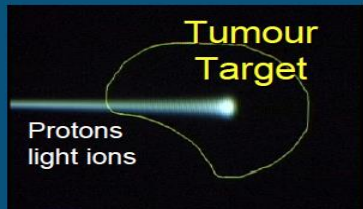
- Strong link between the infinitely small (particle physics) and infinitely large (cosmology)

# Number 2: Innovation

## Medical applications



## Hadron Therapy



Leadership in Ion Beam Therapy not only in Europe and Japan

Accelerating particle beams  
~30'000 accelerators worldwide  
~17'000 used for medicine

>100'000 patients treated worldwide (45 facilities)  
>50'000 patients treated in Europe (14 facilities)



Detecting particles

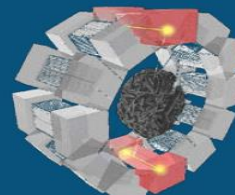


## Imaging

Clinical trial in Portugal, France and Italy for new breast imaging system (ClearPEM)



## PET Scanner



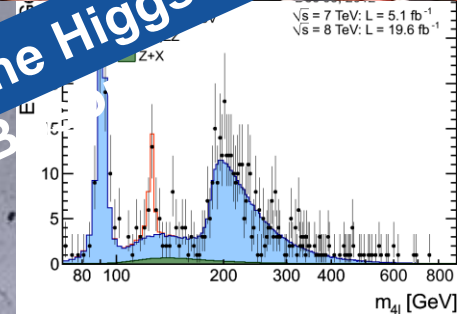
# CERN: A UNIQUE ENVIRONMENT TO PUSH TECHNOLOGIES TO THEIR LIMITS

In its 60 year life CERN has made some of the important discoveries in particle physics

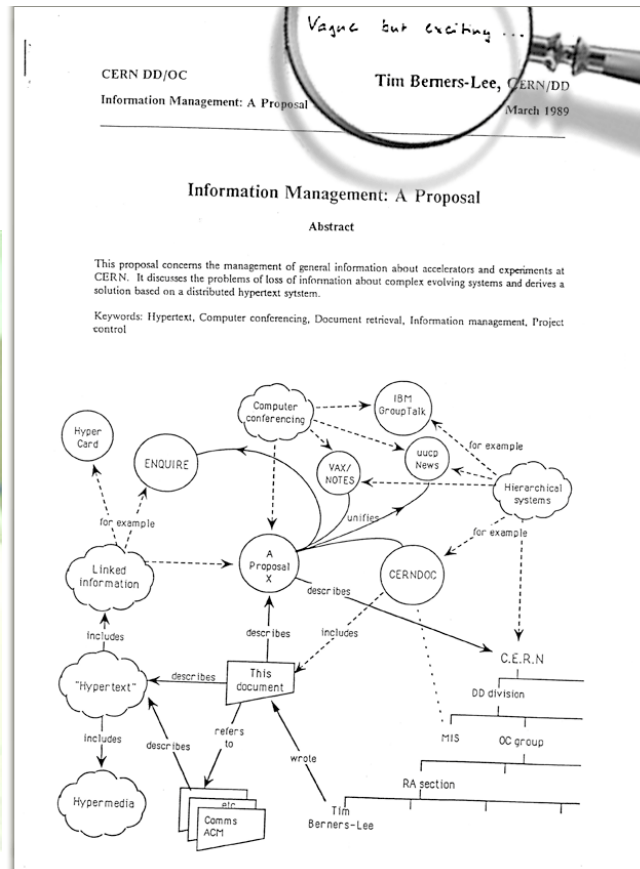
- Observation of the W and Z Bosons
- The number of neutrino families
- The Higgs Boson Discovery



The Higgs  
B



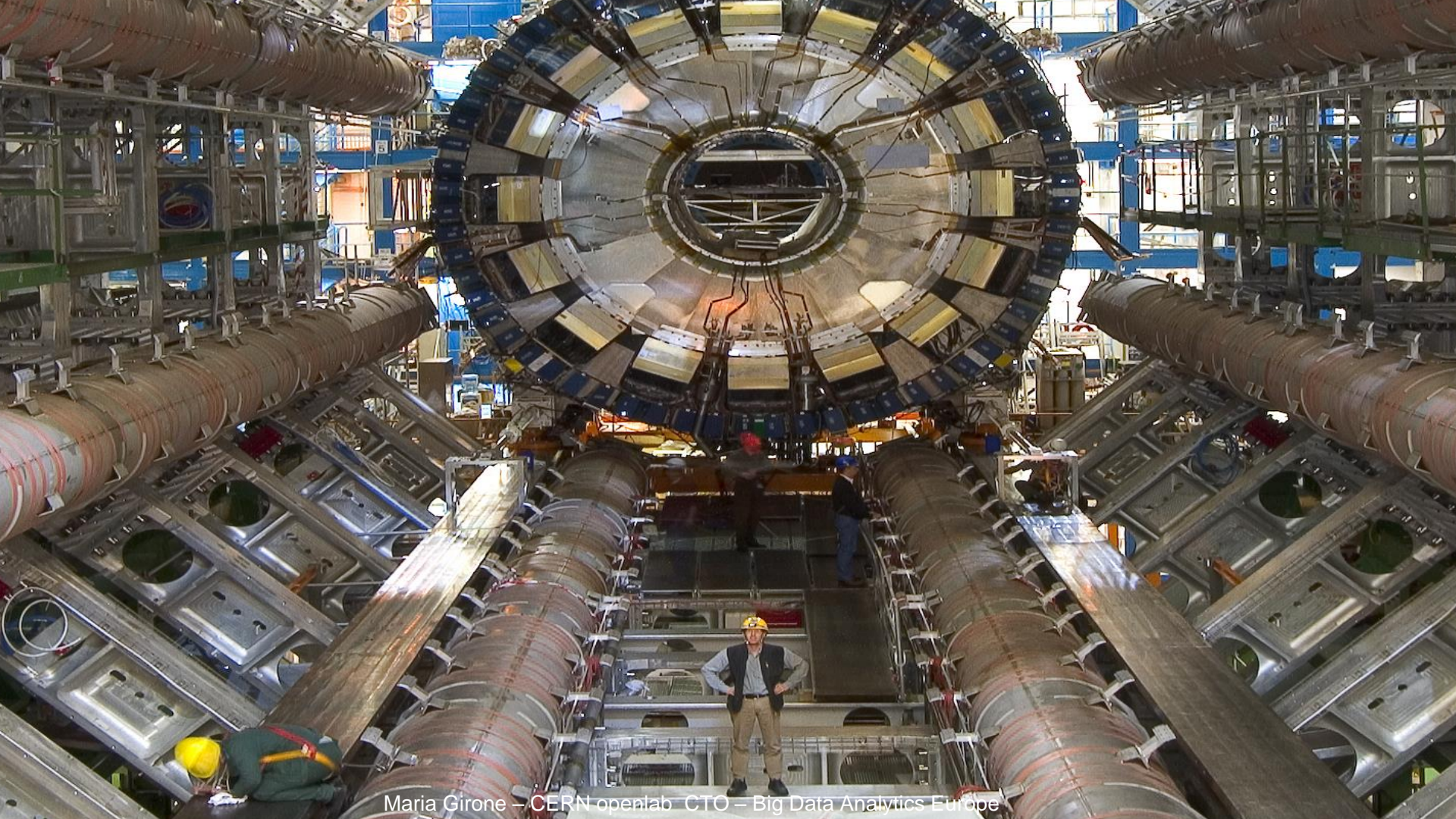
# CERN – Where the Web was Born



# The Large Hadron Collider (LHC)

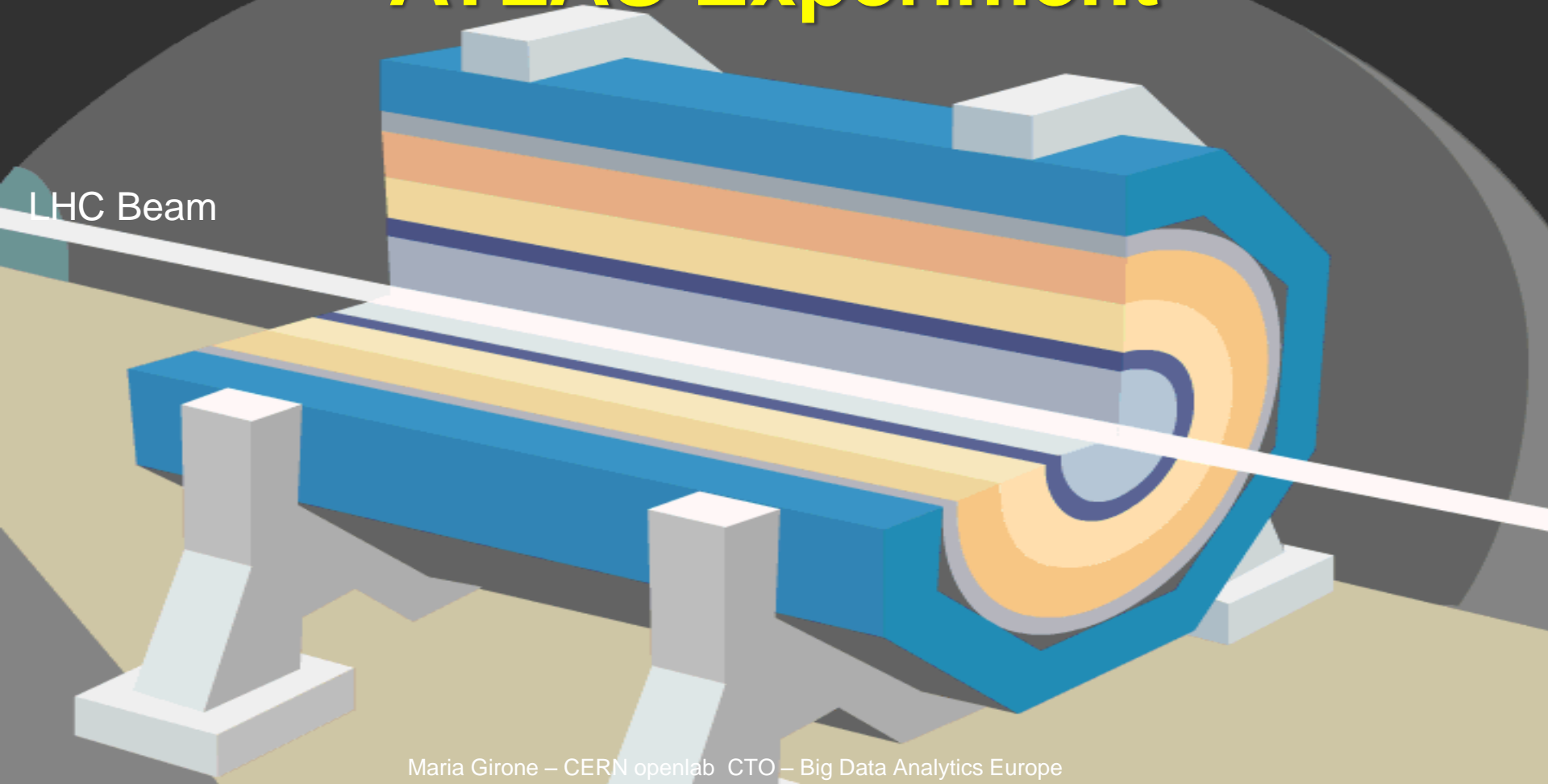






# ATLAS Experiment

LHC Beam



# Data from ATLAS

*Reduction factor of 1 million.*

**1 PB/sec from all sub-detectors**

**1 GB/sec raw data sent to Data Centre**

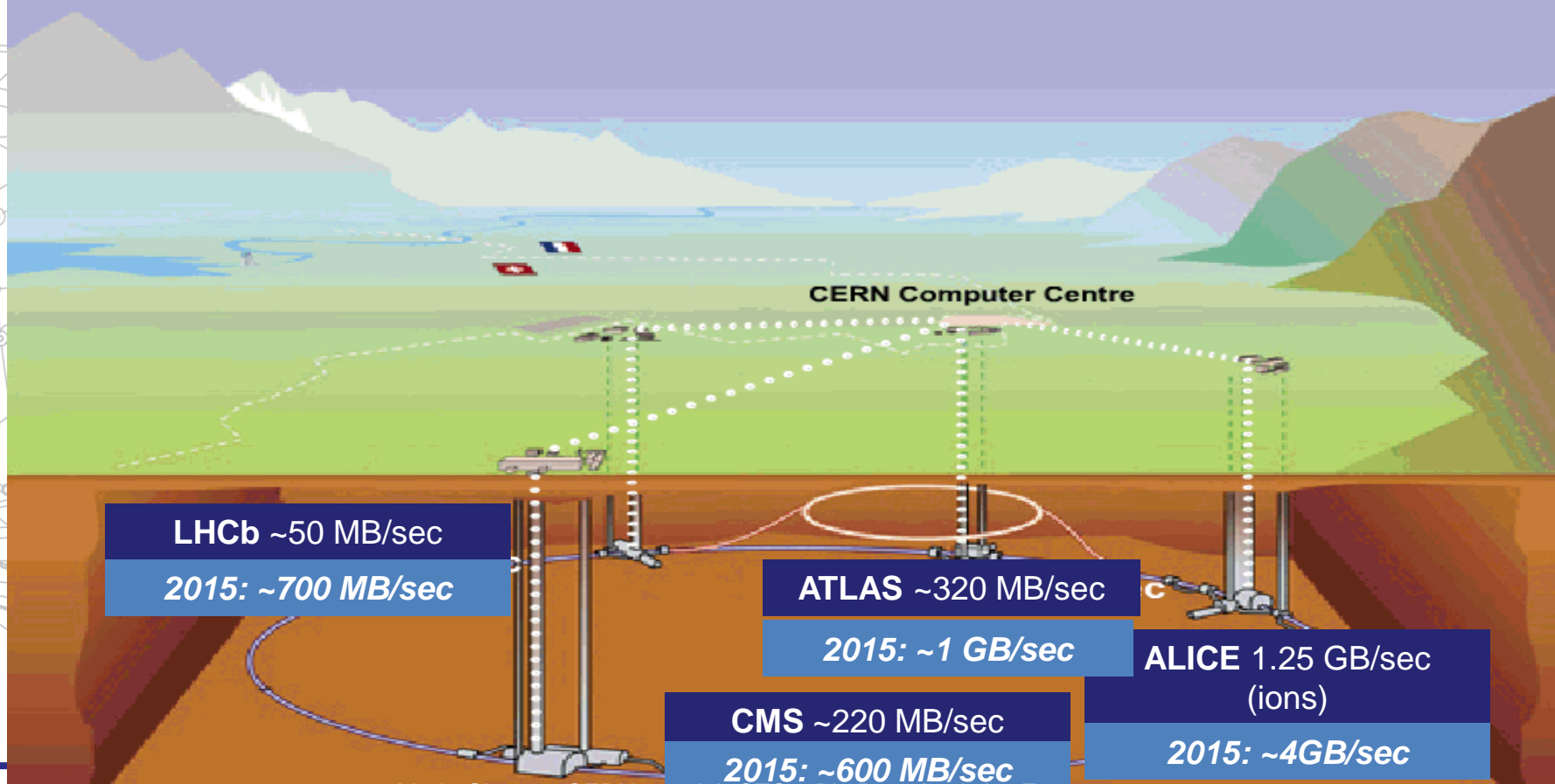
*Trigger and data acquisition*



*Event filter computer farm*



# CERN Computer Centre (Tier-0): Acquisition, First pass reconstruction, Storage & Distribution



1 PB/s of data generated by the detectors  
Up to **30 PB/year** of stored data

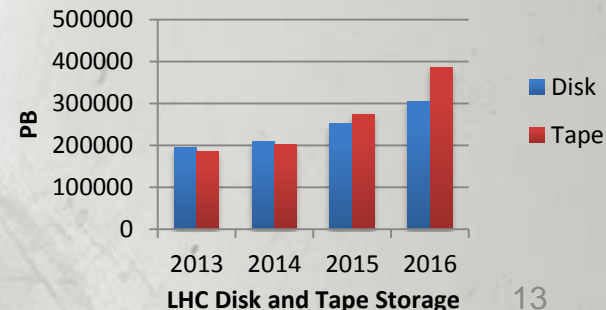
A distributed computing infrastructure  
of half a million cores working 24/7  
An average of 40M jobs/month

An continuous data transfer rate of 6 GB/s  
(**600TB/day**) across the Worldwide LHC Grid  
(WLCG)

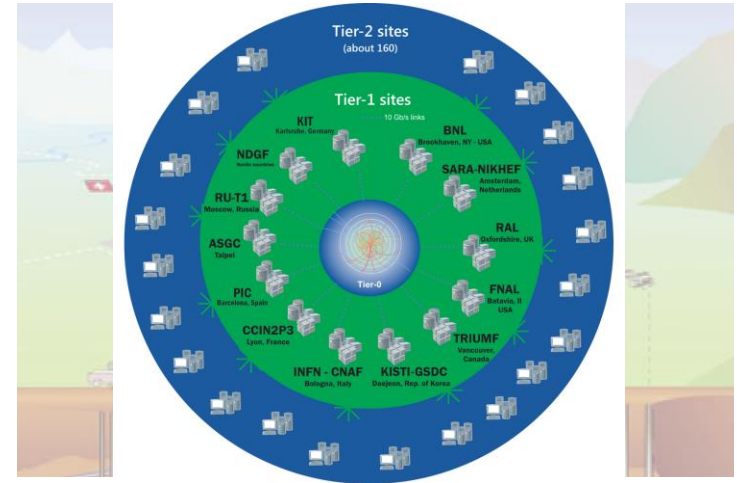
A sample equivalent to the accumulated data/simulation of the 10 year LEP program is produced 5 times a day

Would put us amongst the top Supercomputers if centrally placed

More than 100PB moved and accessed by 10k people



# Worldwide LHC Computing Grid



## Tier-0 (CERN):

- Data recording
- Initial data reconstruction
- Data distribution

## Tier-1 (12 centres):

- Permanent storage
- Re-processing
- Analysis

## Tier-2 (68 Federations, ~140 centres):

- Simulation
- End-user analysis
- 525,000 cores
- 450 PB

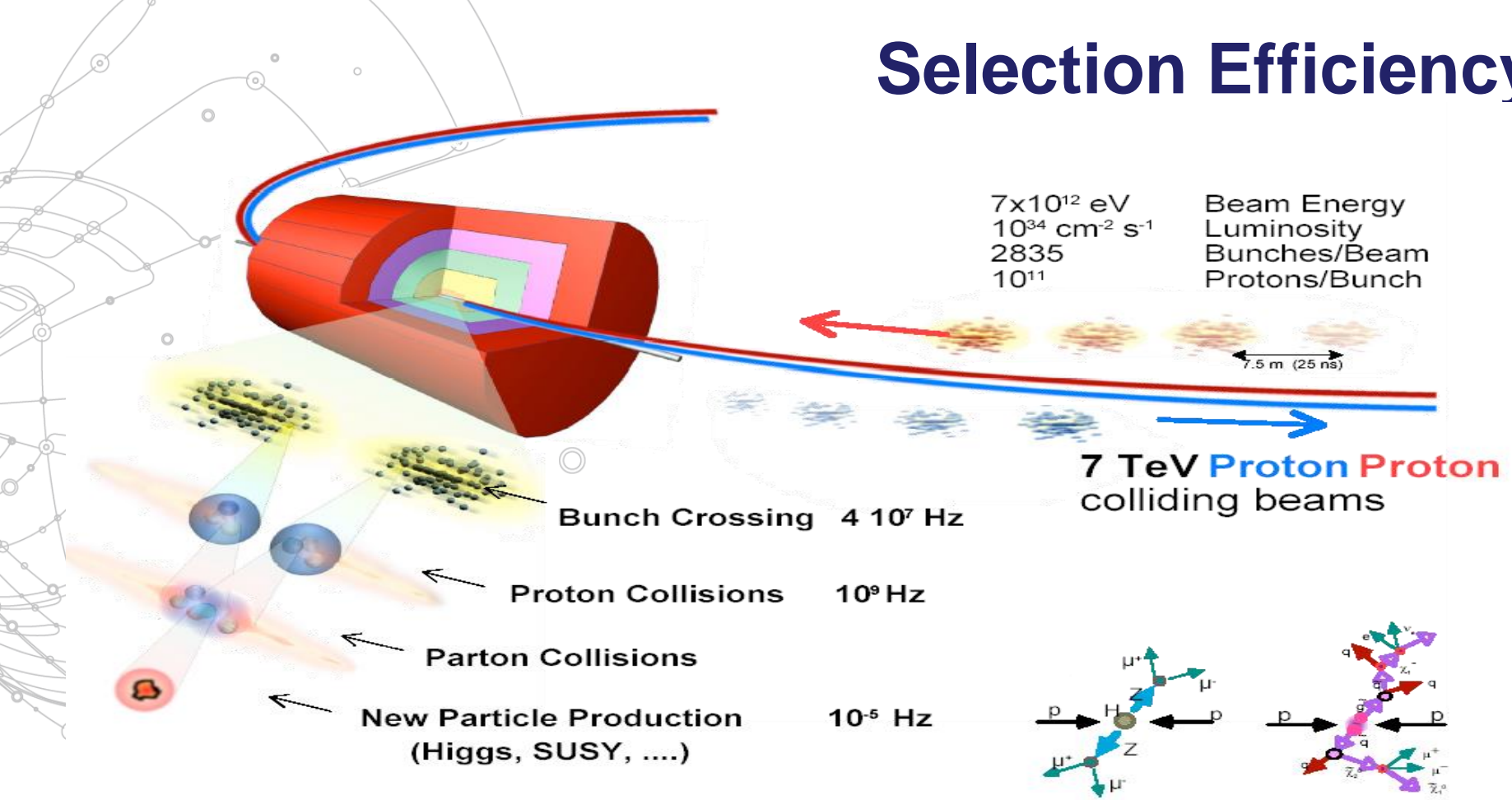
The background of the slide is a complex, abstract network diagram. It consists of numerous nodes, represented by small circles of varying sizes and colors (some white, some grey, some black), interconnected by a dense web of thin, grey lines. A prominent feature is a thick, dark grey line that forms a large, irregular loop on the left side of the image. The overall impression is one of a highly interconnected and complex system, likely representing data networks or computational processes.

# Big Data Analytics



CERNopenlab

# Selection Efficiency



**Selection of 1 event in 10,000,000,000,000**





# Data Reduction

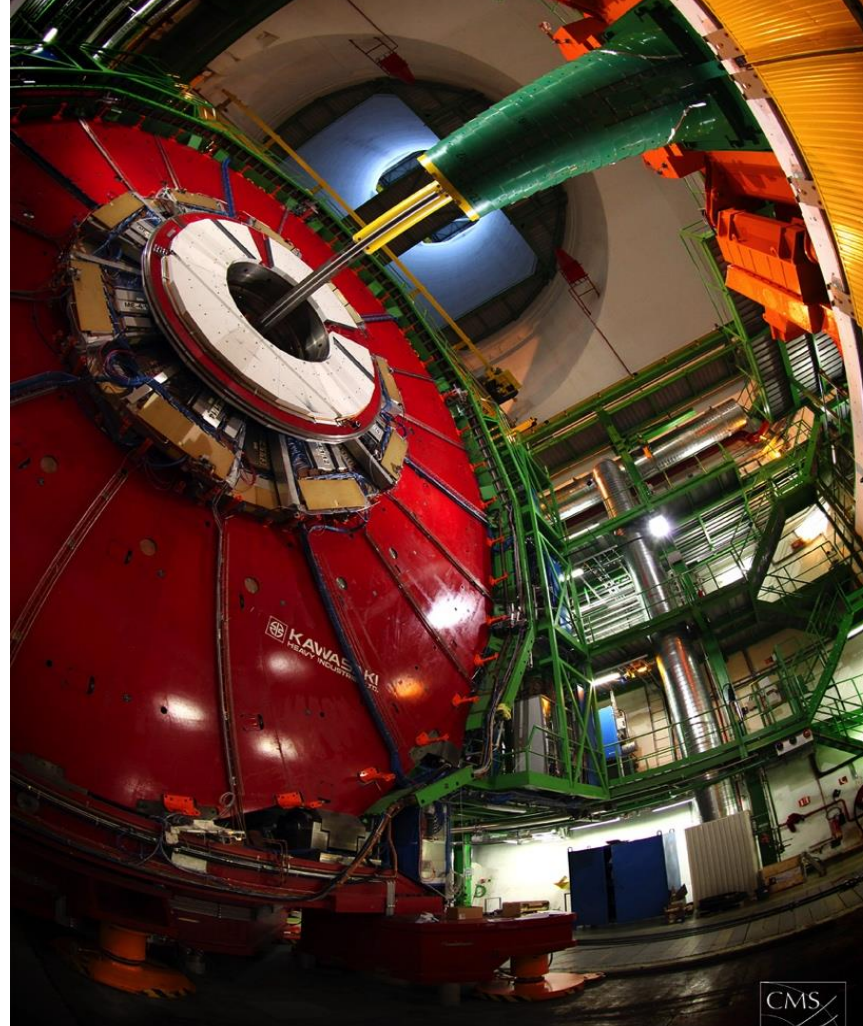
- 150 million active elements
- 20 (40) million bunch crossings per second
- O(1 PB/s) internal data rate



## Data reduction:

- Suppress electronic noise
- Decide to read out and save event, or throw it away (trigger)
- Build the event (assemble all data)

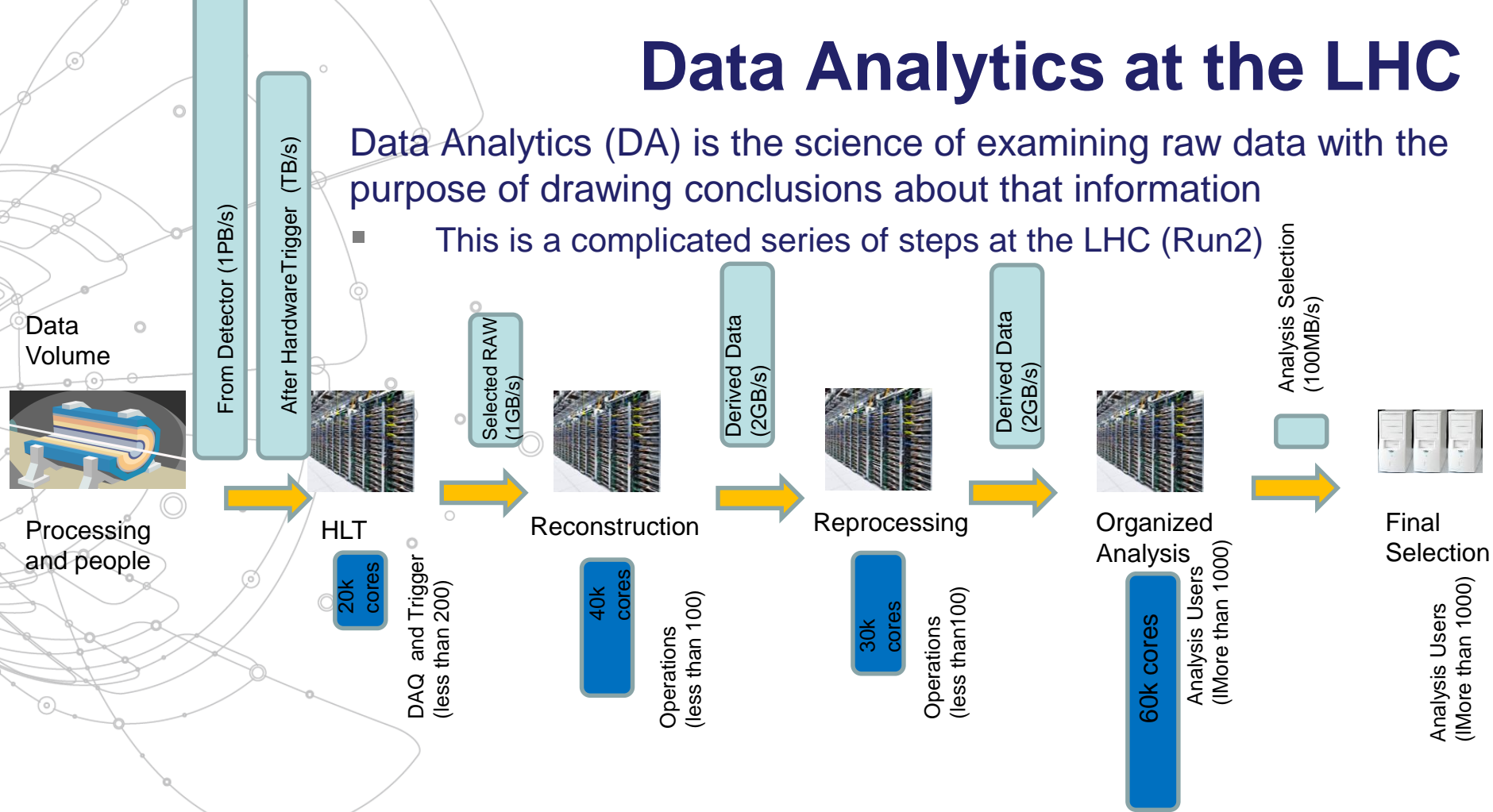
- O(1000 Hz) event rate
- O(1GB/s) data rate



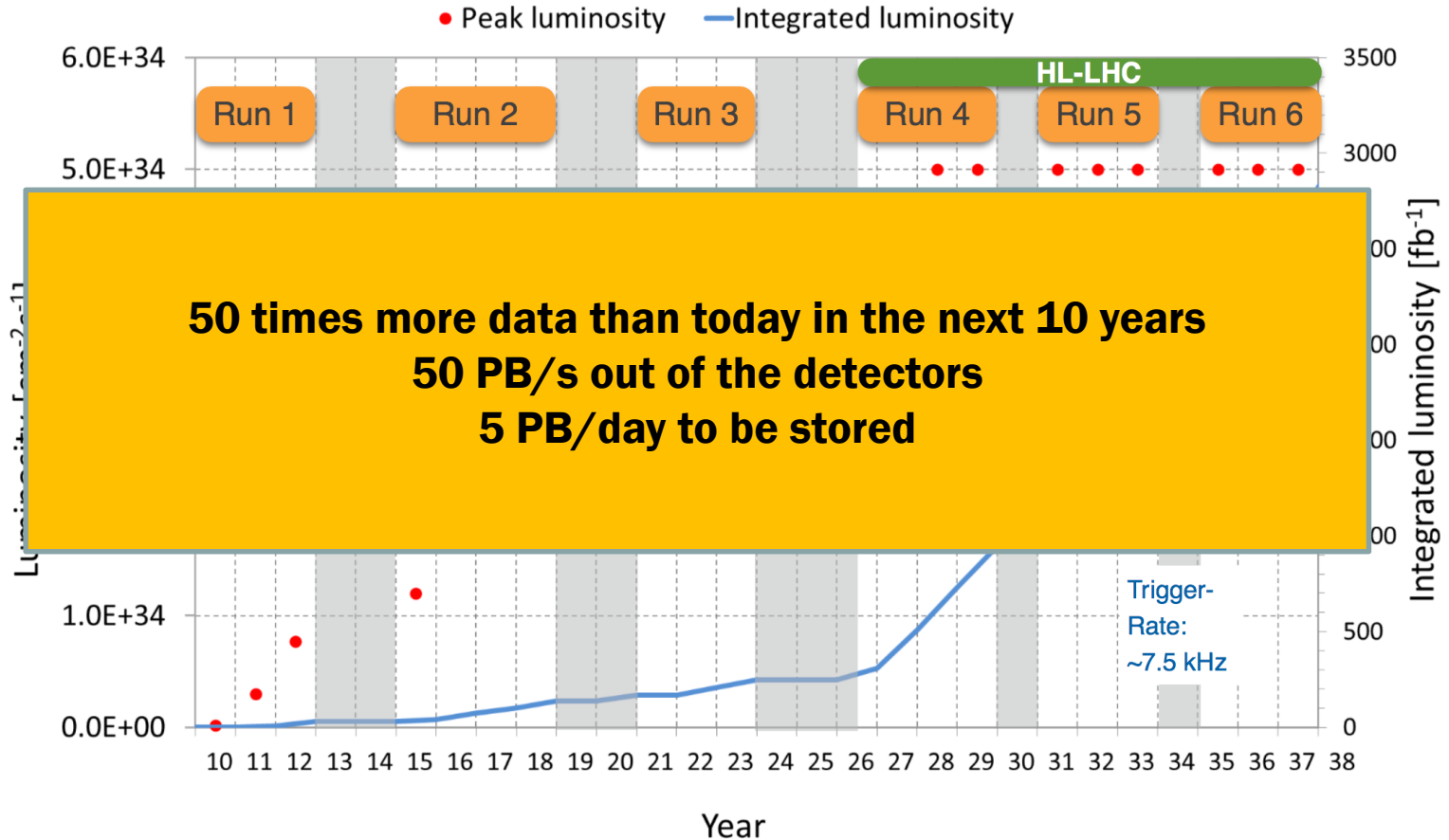
# Data Analytics at the LHC

Data Analytics (DA) is the science of examining raw data with the purpose of drawing conclusions about that information

- This is a complicated series of steps at the LHC (Run2)



# LHC Schedule

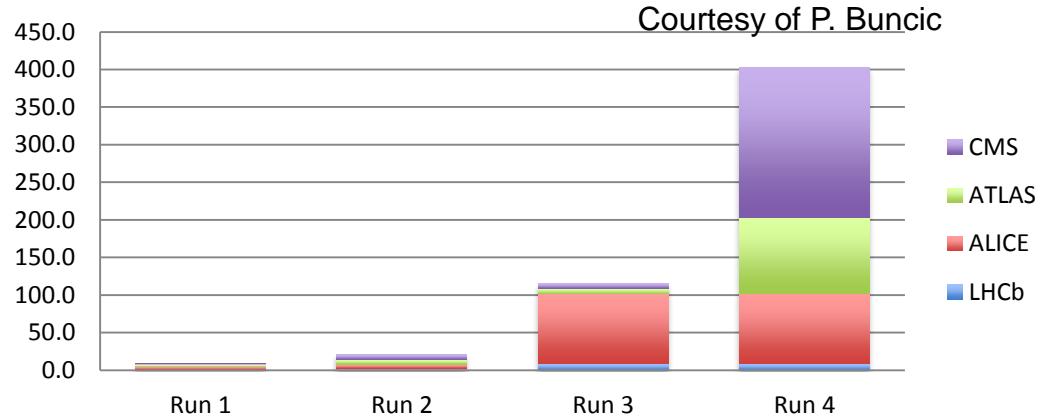


# LHC Run3 and Run4 Scale and Challenges



Raw data volume for LHC increases exponentially

- And with it processing and analysis load
- Current estimate by Run4 for technology improvements for flat budget is an **increase of a factor 8-10**



- › **LHCb and ALICE have big upgrades in Run3**
  - Event rate x 40-100 and factor 10 in volume
- › **ATLAS and CMS upgrade for Run4**
  - Event rate x 10 and big increase in volume

# Run3 and Run4 Scale and Challenges

- The increased data volume is combined with an increase of event complexity

- Resulting in a huge processing challenge

- › Example from CMS, but other experiments are similar

Detector	HLT output rate (kHz)	Data Reco.	Simulation			Total
			Detector sim.	Digi.	Reco.	
Phase-I	1	4	1	3.5	4	3
Phase-II (140)	5	100	5	47	100	65
Phase-II (200)	7.5	340	7.5	100	340	200

<https://cds.cern.ch/record/2020886>

- Total computing needs go up by a factor of 65-200 (wrt Run2)

- › Technology improvements only solve a factor of 10

- › **Code optimization** and **technology revolutions** are needed

# CERN openlab in a nutshell

A unique science – industry partnership to drive R&D and innovation with over a decade of success

- Evaluate state-of-the-art technologies in a challenging environment and improve them

- Test in a research environment today what will be used in many business sectors tomorrow

- Train next generation of engineers/employees

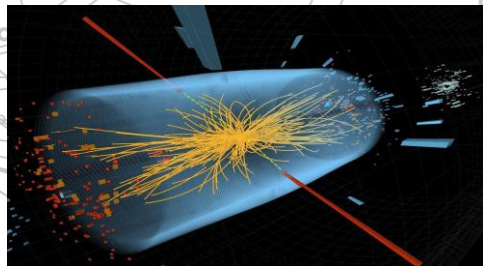
- Disseminate results and outreach to new audiences



# Data Analytics to the Rescue

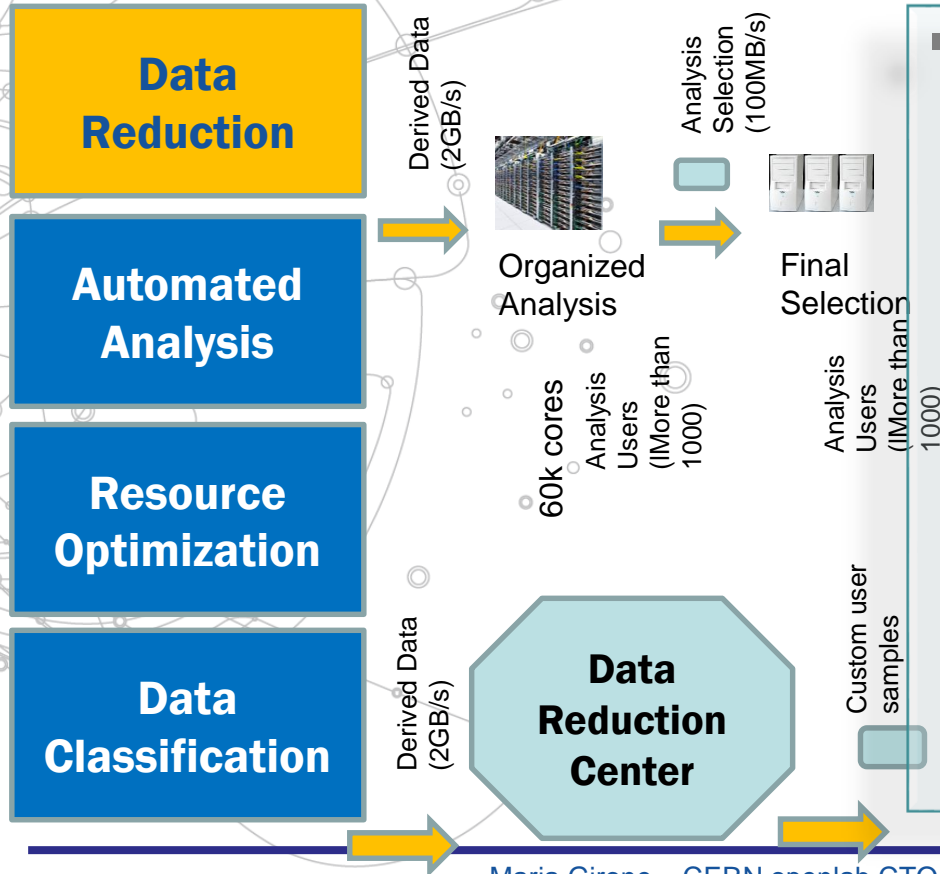
How to make more effective use of the data collected is critical to maximise scientific discovery and close the resource gap

- There are currently ongoing projects in
  - › Accelerator system controls
  - › Data Storage and quality optimizations
- Organising projects on
  - › Data reduction
  - › Optimized formats
  - › Investigations for machine learning for analysis and event categorization
- CERN openlab is uniquely positioned to help in this area with connections to industry





# Analytics



After the upgrade LHC will collect large datasets. Investigating ways to more efficiently select events from the stream of data using “big data” techniques

- Through well established techniques like MapReduce can we reduce the computational load of analysis
- Need to reduce multi-petabyte datasets by a factor of 1000 based on physics selection criteria
  - › Performance, reproducibly, and completeness are all important

**Data  
Reduction**

**Automated  
Analysis**

**Resource  
Optimisation**

**Data  
Classification**



Reconstruction

20k  
cores

Operations  
(less than 50)

- The Data Quality Monitoring is a key to delivering high-quality data for physics. It is used both in the online and offline environments
  - Currently involves scrutinizing of a large number of histograms by detector experts comparing them with a reference
  - Aim at applying recent progress in Machine Learning techniques to the automation of the DQM scrutiny
- The LHC is the largest piece of scientific apparatus ever built
  - There is a tremendous amount of real time monitoring information to assess health and diagnose faults
  - The volume and diversity of information makes this an interesting application of big data analytics

**Data  
Reduction**

**Automated  
Analysis**

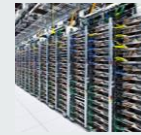
**Resource  
Optimization**

**Data  
Classification**

- Use machine learning techniques to predict how data should be placed and processing resources scheduled in order to achieve a dramatic reduction in latency for delivering data samples to analysts
- Design a system capable of using information about resource usage (disk access, CPU efficiency, job success rates, data transfer performances, and more) to make more automated decisions about resource allocation



Reconstruction



Reprocessing



Organized  
Analysis



# Analytics

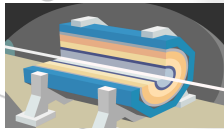
**Data  
Reduction**

**Automated  
Analysis**

**Resource  
Optimization**

**Data  
Classification**

Data  
Volume



All Data (TB/s)

- Investigate the possibility of performing real-time event classification in the high-level trigger system of LHC experiments
  - Extract information from events that would otherwise be rejected
- Uncategorized events might potentially be the most interesting, revealing the presence of new phenomena in the LHC data.
  - Event classification would allow both a more efficient trigger design and an extension of the physics program, beyond the boundaries of the traditional trigger strategies

# Summary and Outlook

- The LHC is planning to dramatically increase the volume and complexity of data collected by Run3 and Run4
  - This results in an unprecedented computing challenge in the field of High Energy Physics
- Meeting this challenge within a realistic budget requires rethinking how we work
  - Turning to industry and other sciences for improvements in data analytics
    - › Data reduction through MapReduce and automated analysis through machine learning techniques
- CERN openlab is in a unique position to engage with industry



#### EXECUTIVE CONTACT

Alberto Di Meglio, CERN openlab Head

[alberto.di.meglio@cern.ch](mailto:alberto.di.meglio@cern.ch)

#### TECHNICAL CONTACT

Maria Girone, CERN openlab CTO

[maria.girone@cern.ch](mailto:maria.girone@cern.ch)

Fons Rademakers, CERN openlab CRO

[Fons.rademakers@cern.ch](mailto:Fons.rademakers@cern.ch)

#### COMMUNICATION CONTACT

Andrew Purcell, CERN openlab Communication Officer

[andrew.purcell@cern.ch](mailto:andrew.purcell@cern.ch)

Mélissa Gaillard, CERN IT Communication Officer

[melissa.gaillard@cern.ch](mailto:melissa.gaillard@cern.ch)

#### ADMIN CONTACT

Kristina Gunne, CERN openlab Administration Officer

[kristina.gunne@cern.ch](mailto:kristina.gunne@cern.ch)

The background of the slide is a complex, abstract network diagram. It consists of numerous nodes, represented by small circles of varying sizes and colors (some white, some black, some grey), interconnected by thin, grey lines. Some lines are thicker and more prominent, forming a central structure that resembles a stylized letter 'S' or a similar shape. The overall appearance is that of a data network or a complex system architecture.

# Additional Slides



CERNopenlab

# Relative Size of Things

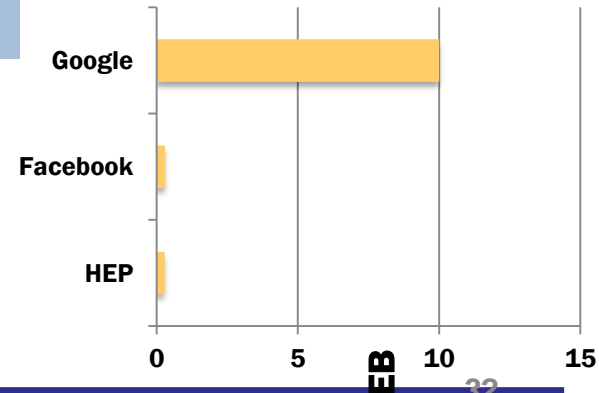
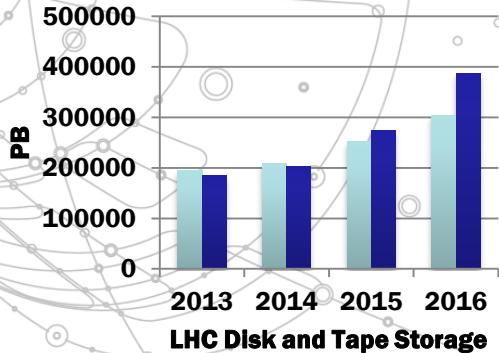
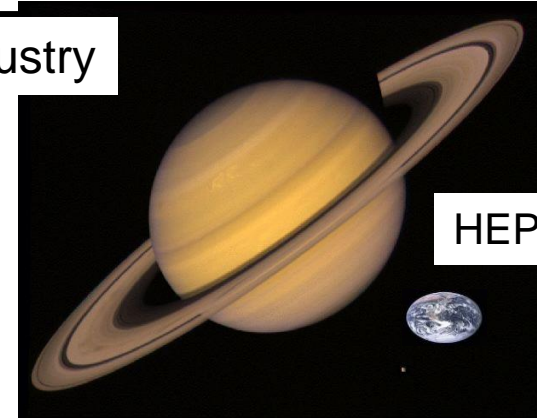
## Processing

- > Amazon has more than 40 million processor cores in EC2

## Storage

- > Amazon has  $2 \times 10^{12}$  unique user objects and supports 2M queries per second
- > Google has 10-15 exabytes under management
- > Facebook 300PB
- > eBay collected and accessed the same amount of data as LHC Run1

## Industry



**Our data and processing problems are ~1% the size of the largest industry problems**