# Challenges in high-energy physics computing

Andrei Gheata
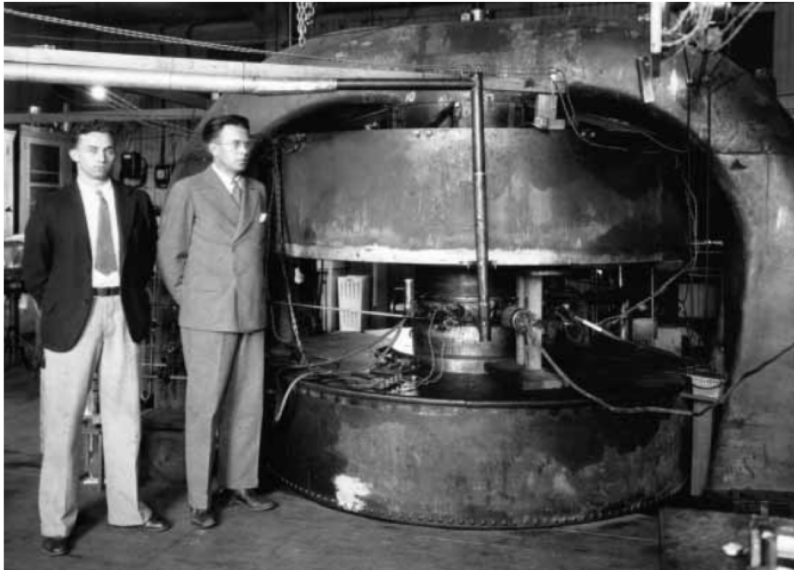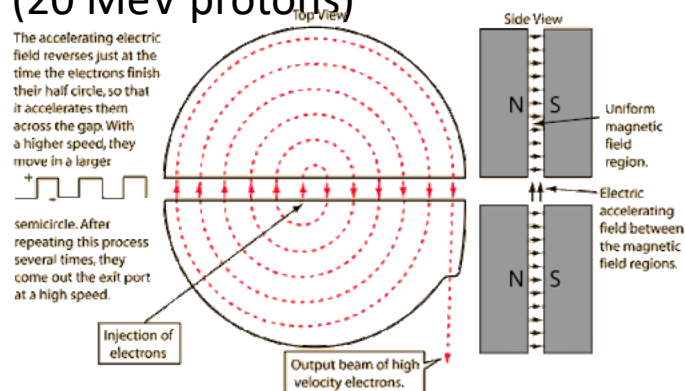
EP/SFT

CERN

# Outline

- Modern High Energy Physics (HEP) experiments
- Physics software
  - online processing
    - triggering, selection
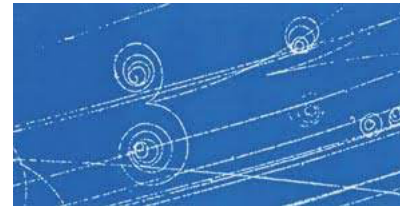  - offline processing
  - simulation
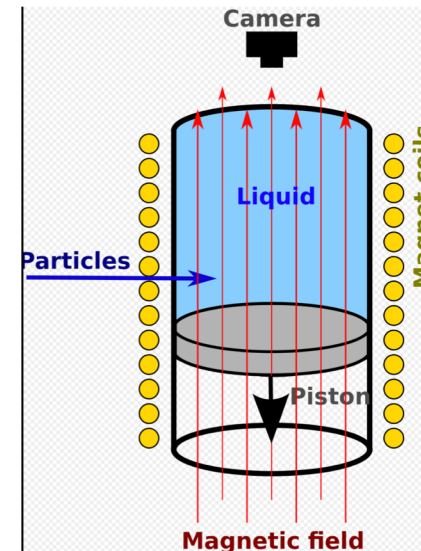- Conclusions

# Accelerators & particle physics



S. Livingstone and E. Lawrence
69cm cyclotron (20 MeV protons)



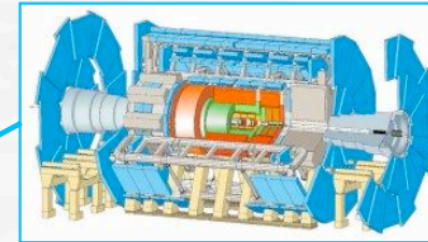Bubble chamber experiments
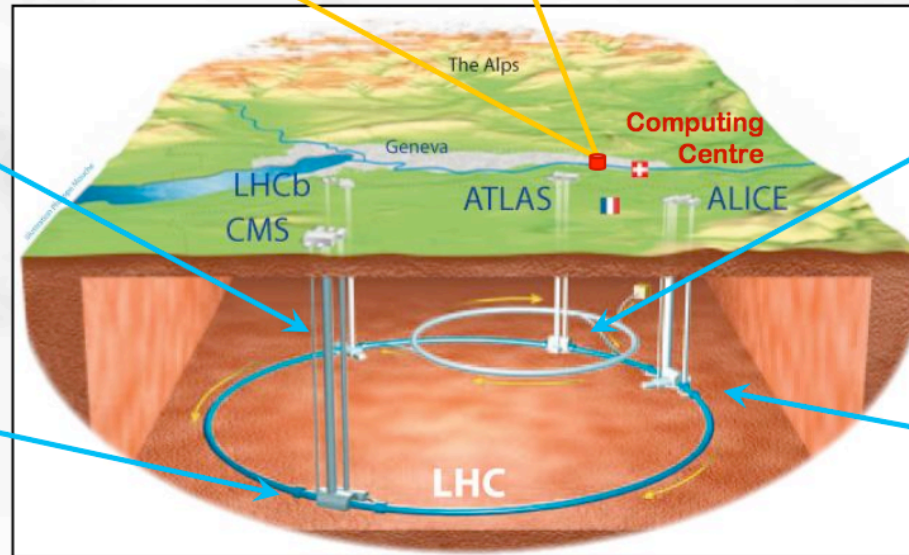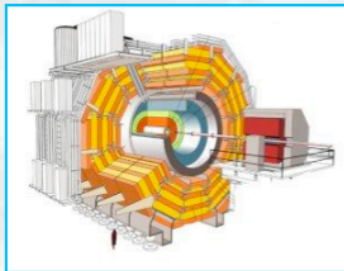




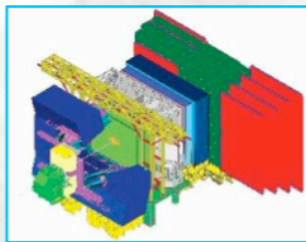BIg European Bubble Chamber (Microcosm)

# Big science –> big data

Exascale computing ($10^{18}$ bytes)

Grid/cloud distributed computing

500k processor cores

800 million pp collisions per second

# The data processing chain



Online processing

Trigger

1 ms          $10^6$ -$10^8$ sec

Collision     Detectors     Event fragments     Full event     Storage     Offline analysis

# Collisions at the LHC

**Bunch**

**Proton**

**Parton**
(quark, gluon)

**Particle**

jet      jet

Higgs       $e^+$

$e^+$                $e^-$

$Z^o$         $Z^o$

$e^-$

SUSY.....

| | |
|---|---|
| **Proton** - **Proton** | **2804 bunch/beam** |
| **Protons/bunch** | $10^{11}$ |
| **Beam energy** | **7 TeV ($7\times10^{12}$ eV)** |
| **Luminosity** | $10^{34}$cm$^{-2}$s$^{-1}$ |
| **Crossing rate** | **40 MHz (25 ns)** |
| **Collision rate ≈** | $10^{7}$-$10^{9}$ **Hz** |

# Searching the Higgs – needle in a haystack



'Flying garbage'

'Hard Scatter'

'Secondary scatter'

Proton-Proton collision at the LHC

# The needle in many, many haystacks



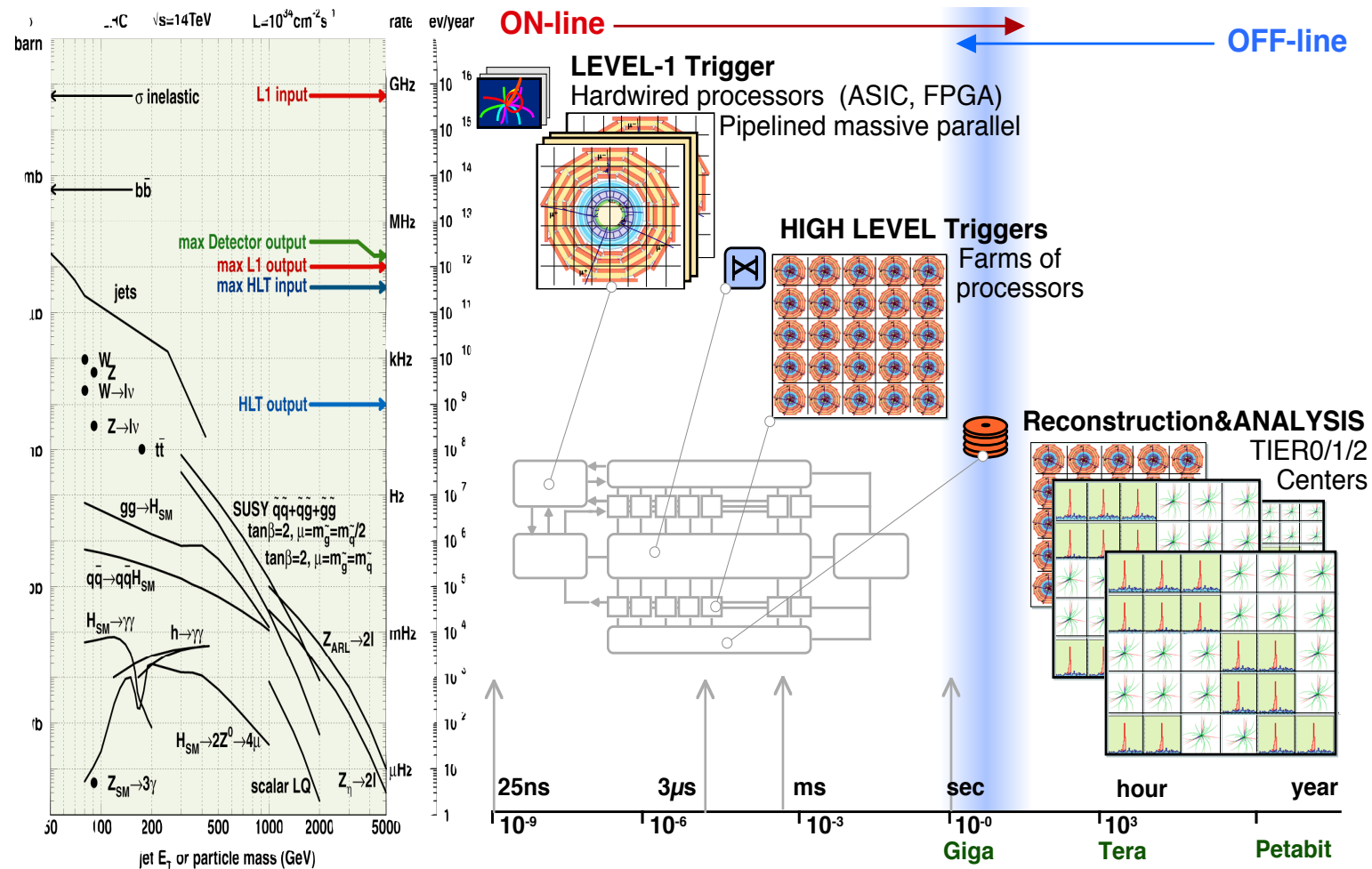- Cross sections (probabilities) of physics processes vary over many orders of magnitude
  - Inelastic: GHz
  - $W \to \ell \nu$: 100 Hz
  - t $t_{bar}$ production: 10 Hz
  - Higgs (125 GeV/c$^2$): 0.1 Hz
- Selection needed: 1:10$^{10-11}$

one Higgs on 10.000.000.000 collisions
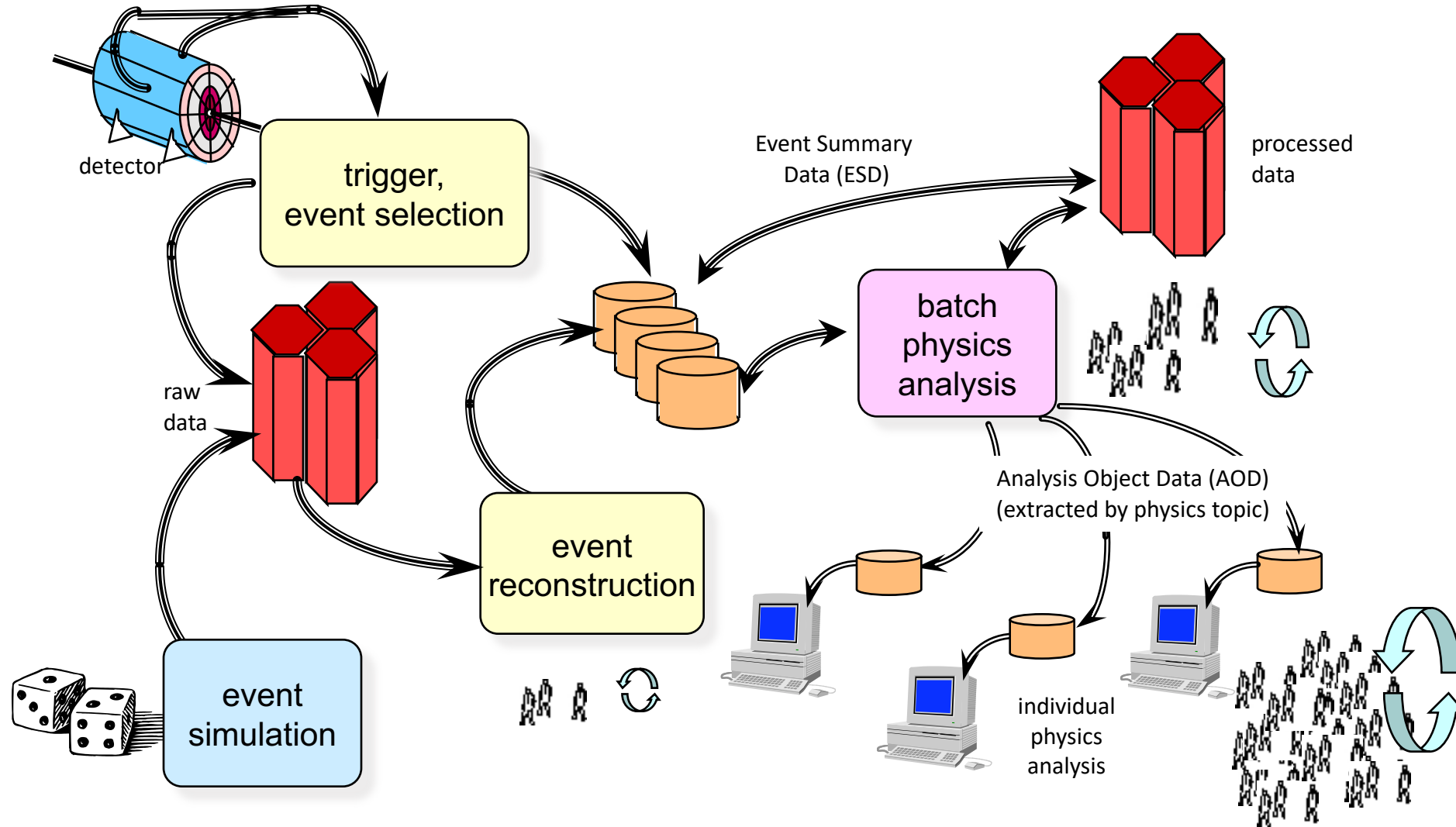
~3 million until 2017

# Physics Selection at LHC

# Physics software

- The scientific software needed to process this huge amount of data from the LHC detectors is developed by the LHC collaborations
  - Must cope with the unprecedented conditions and challenges (trigger rate, data volumes, etc.)
  - Each collaboration has written millions of lines of code
- Modern technologies and methods
  - Object-oriented programming languages and frameworks
  - Re-use of a number of generic and domain-specific 'open-source' packages
- The organization of this large software production activity is by itself a huge challenge
  - Large number of developers distributed worldwide
  - Integration and validation require large efforts

# Processing Stages



detector

trigger, event selection

Event Summary Data (ESD)

processed data

raw data

event reconstruction

event simulation

batch physics analysis

Analysis Object Data (AOD) (extracted by physics topic)

individual physics analysis

# Processing Stages - Trigger



detector

trigger,
event selection

Event Summary
Data (ESD)

processed
data

raw
data

batch
physics
analysis

event
reconstruction

Analysis Object Data (AOD)
(extracted by physics topic)

event
simulation

individual
physics
analysis

# Trigger Levels

o Level-1

  ~ 1:10$^4$

 o Hardwired processors (ASIC, FPGA, …)
 o Pipelined massive parallel
 o Partial information, quick and simple event characteristics (pt, total energy, etc.)
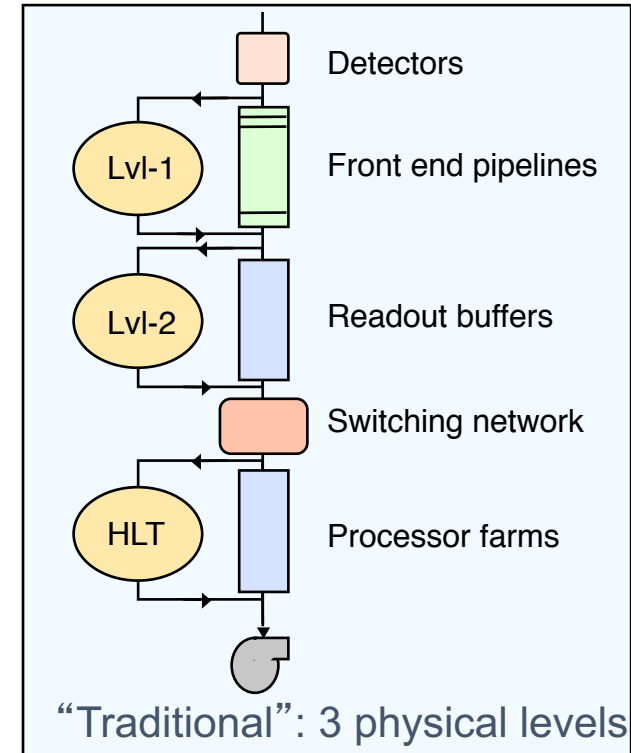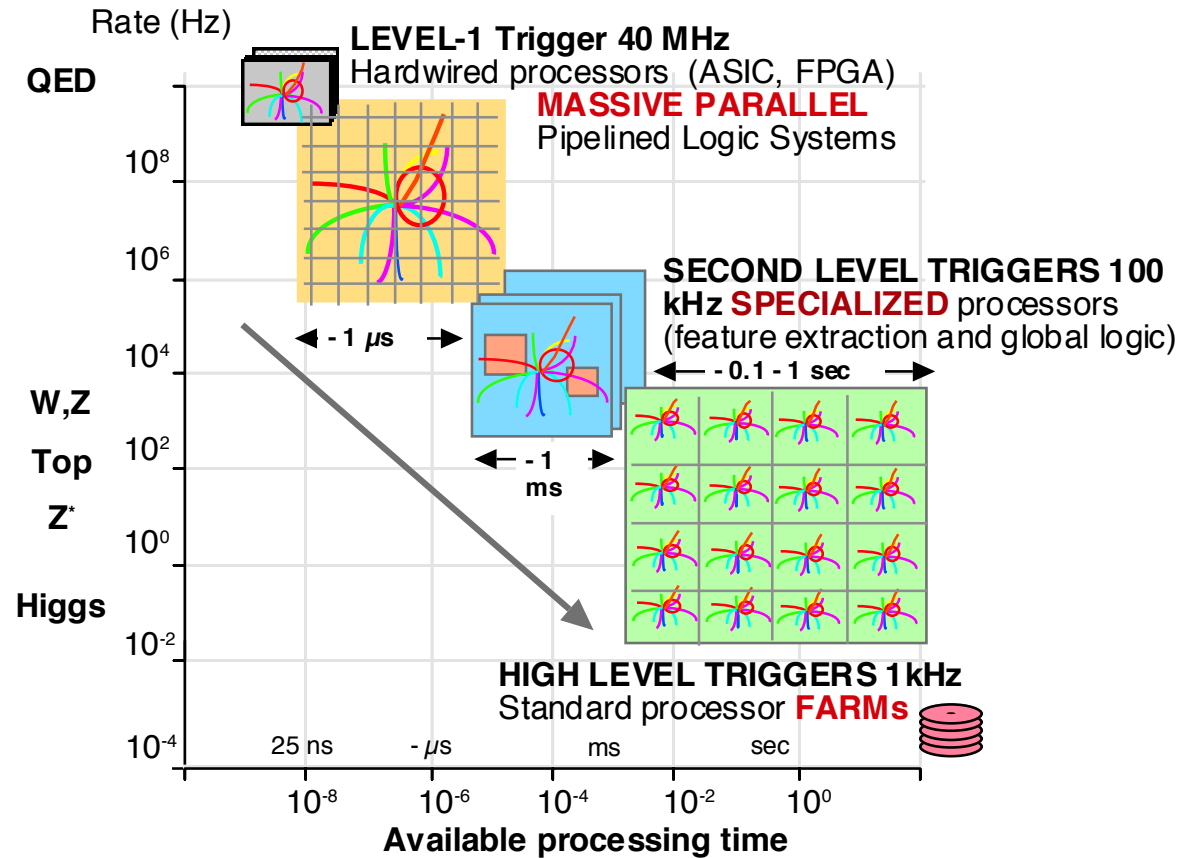 o 3-4 µs maximum latency

o Level-2 (optional)

  ~ 1:10$^1$

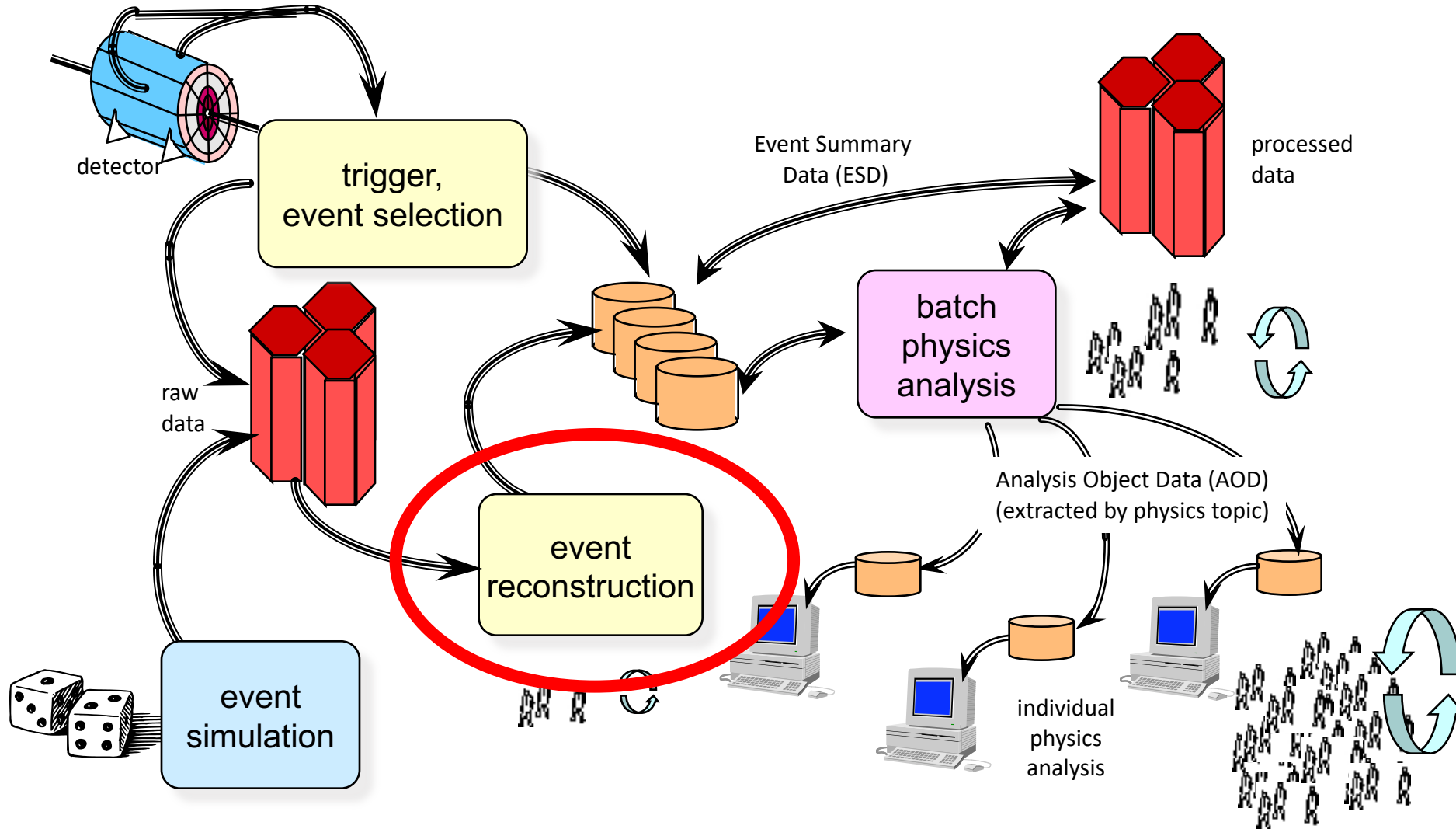 o Specialized processors using partial data

o High Level

 o Software running in processor farms
 o Complex algorithms using complete event information
 o Latency at the level of fractions of second  ~ 1:10$^2$
 o Output rate adjusted to what can be afforded
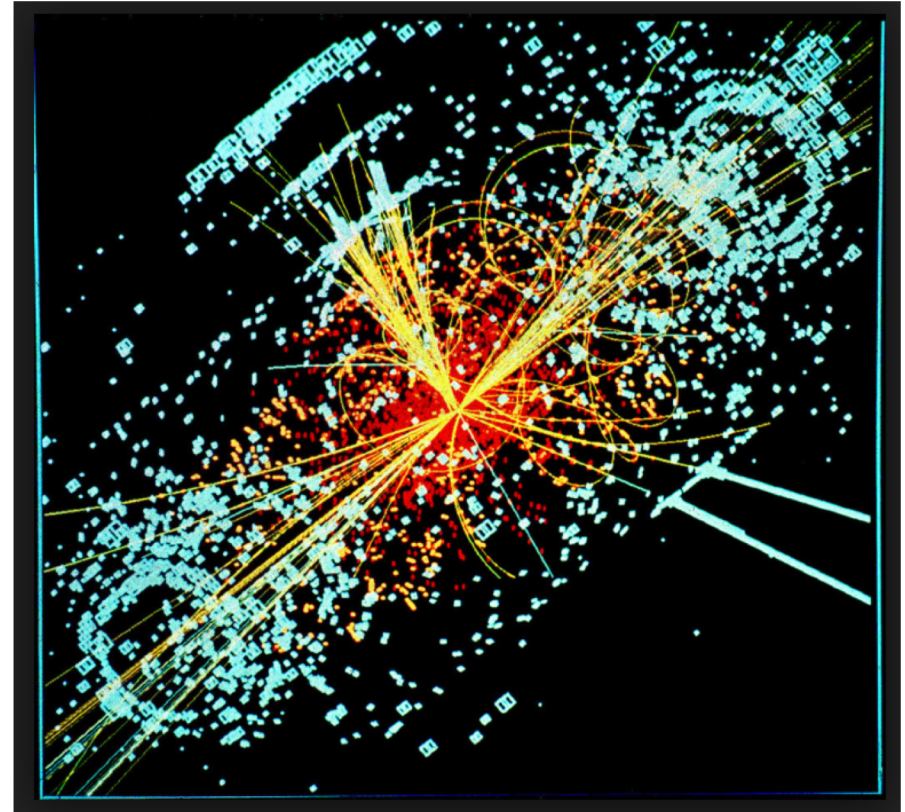
# Trigger Levels and Rates



Rate (Hz)

QED

**LEVEL-1 Trigger 40 MHz**
Hardwired processors (ASIC, FPGA)
**MASSIVE PARALLEL**
Pipelined Logic Systems

$10^8$

$10^6$

**SECOND LEVEL TRIGGERS 100
kHz SPECIALIZED** processors
(feature extraction and global logic)

- 1 $\mu$s

$10^4$

- 0.1 - 1 sec

W,Z

Top

$10^2$

- 1
ms

$Z^*$

$10^0$

Higgs

$10^{-2}$

**HIGH LEVEL TRIGGERS 1kHz**
Standard processor **FARMs**

$10^{-4}$

25 ns          - $\mu$s          ms          sec

$10^{-8}$   $10^{-6}$   $10^{-4}$   $10^{-2}$   $10^{0}$

**Available processing time**

Detectors

Lvl-1

Front end pipelines

Lvl-2

Readout buffers

Switching network

HLT

Processor farms

"Traditional": 3 physical levels

14

# Processing Stages - Reconstruction



detector

trigger, event selection

raw data

event simulation

event reconstruction

Event Summary Data (ESD)

processed data

batch physics analysis

Analysis Object Data (AOD) (extracted by physics topic)
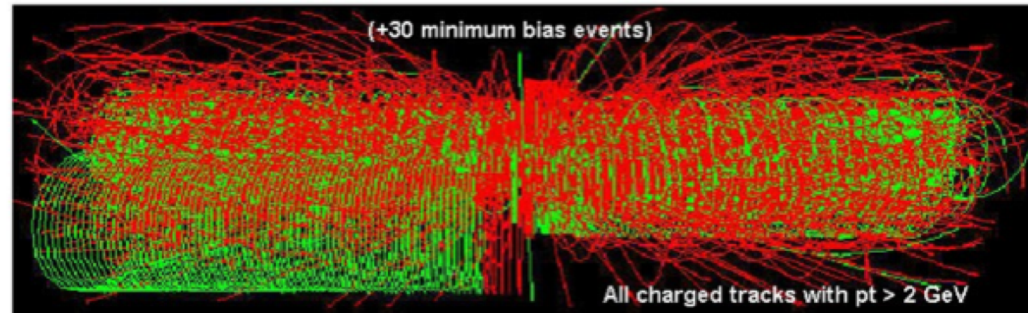
individual physics analysis

# What is reconstruction

- Tracker 'hits' form a puzzle
  - Which tracks created them?
- Each energy deposition is a clue
  - There are thousands of measurements in each snap-shot
- The experiment's reconstruction must obtain a solution!
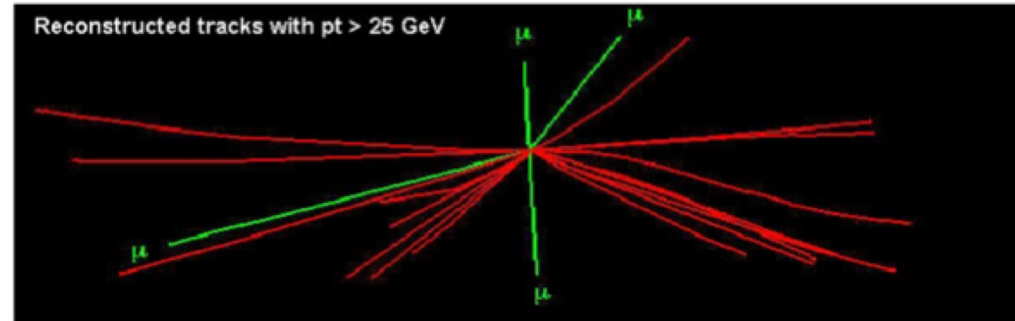  - In well measured magnetic field
  - Matches the traces to tracks

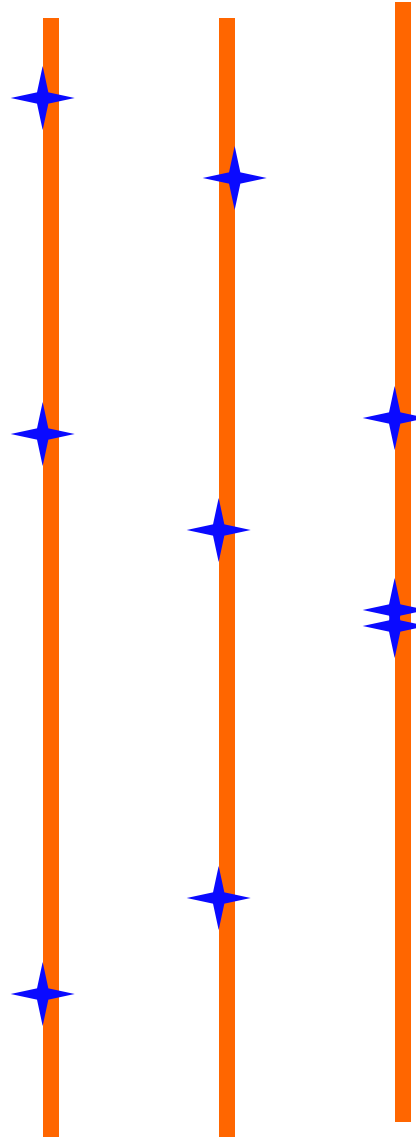# The Reconstruction challenge

**Starting from this event**



(+30 minimum bias events)

All charged tracks with pt > 2 GeV

**Looking for this "signature"**


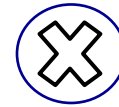
Reconstructed tracks with pt > 25 GeV

→ **Selectivity: 1 in $10^{13}$**
(Like looking for a needle in 20 million haystacks)

# How it works – a simple example

- Combine space points on first three planes (seeds)
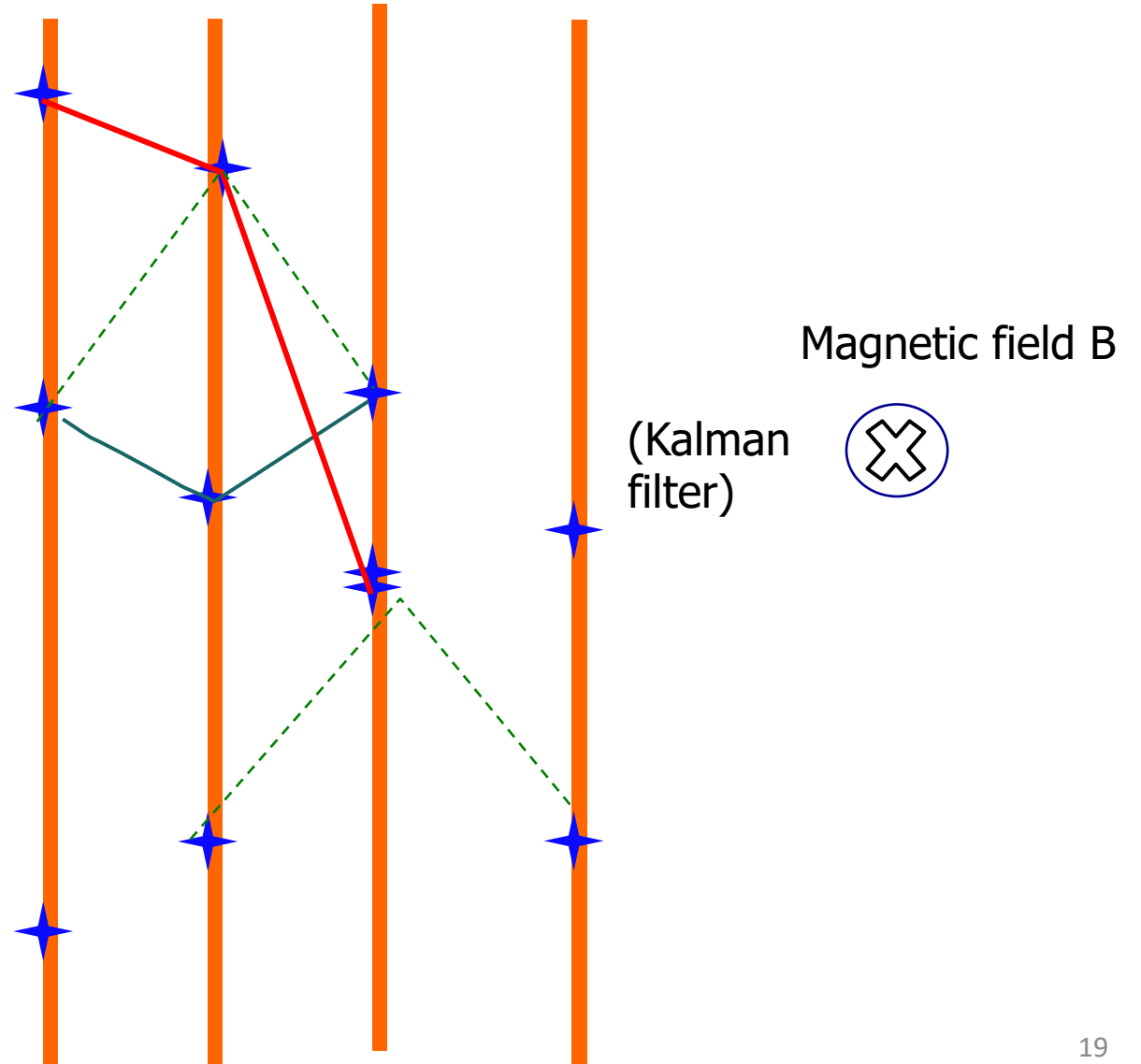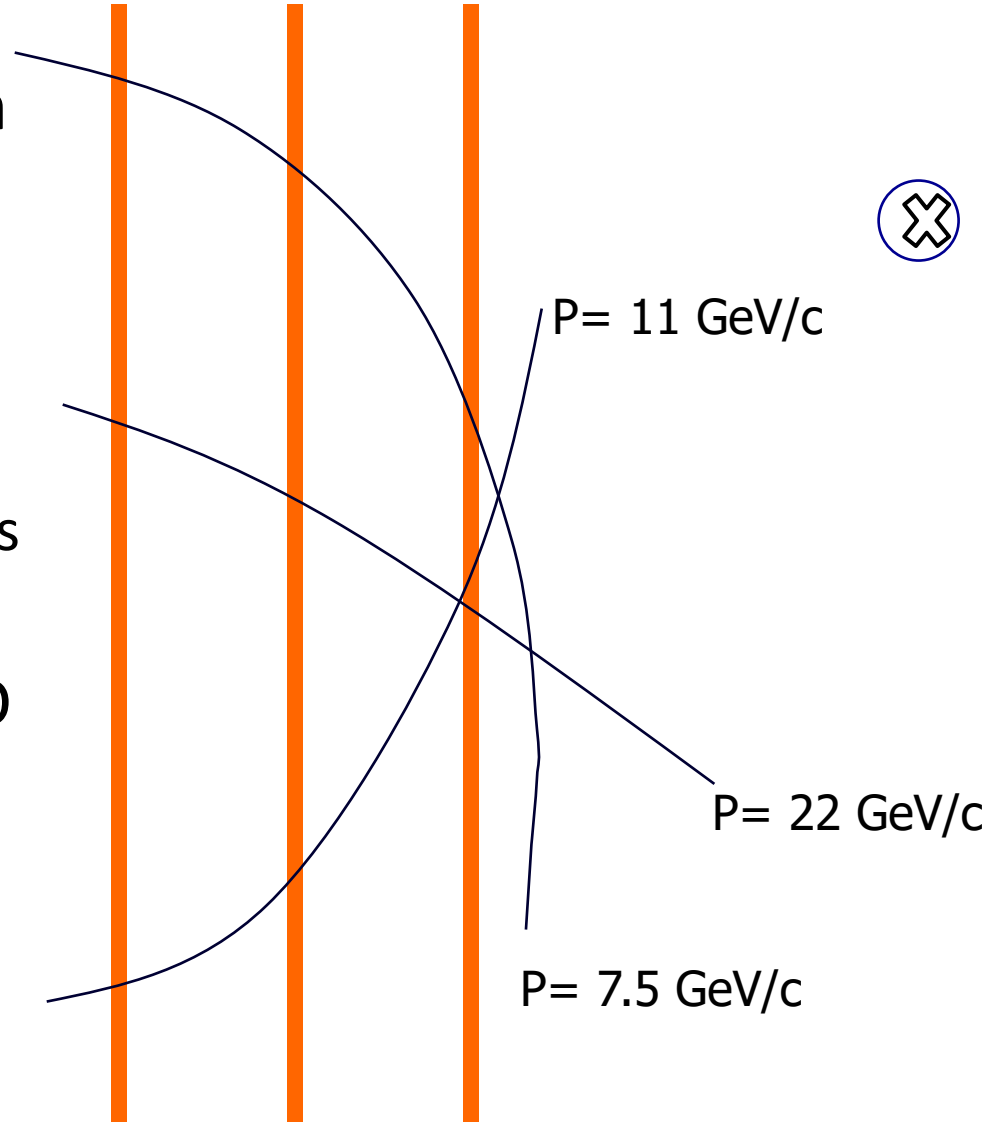
Magnetic field B

# How it works – a simple example

- Combine space points on first three planes (seeds)

- Prolongate to subsequent planes
  - Calculate differences between measured points and predictions

Magnetic field B
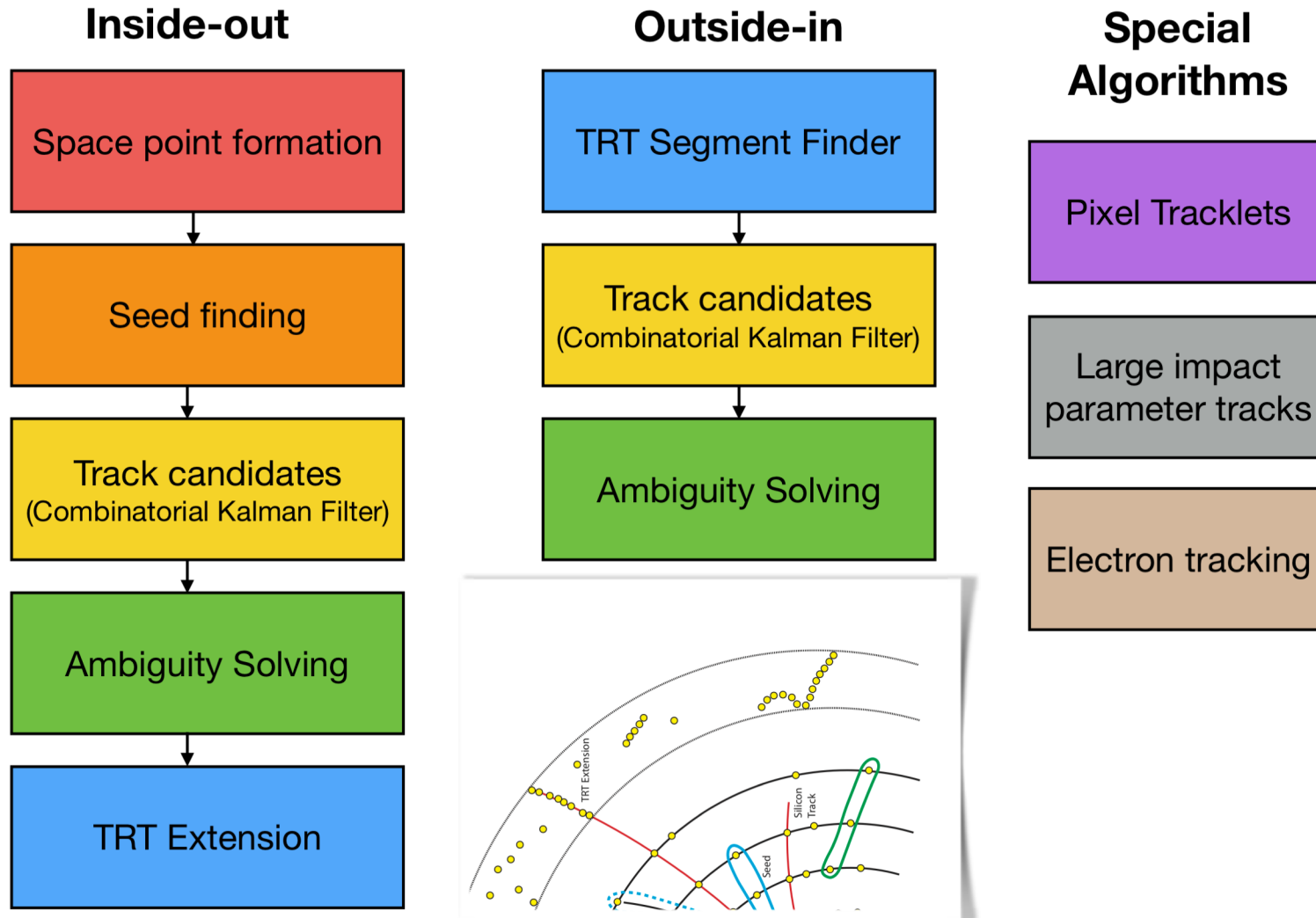
(Kalman filter)

# How it works – a simple example

- Combine space points on first three planes (seeds)

- Prolongate to subsequent planes
  - Calculate differences between measured points and predictions

- Do track fitting, using PID hypothesis
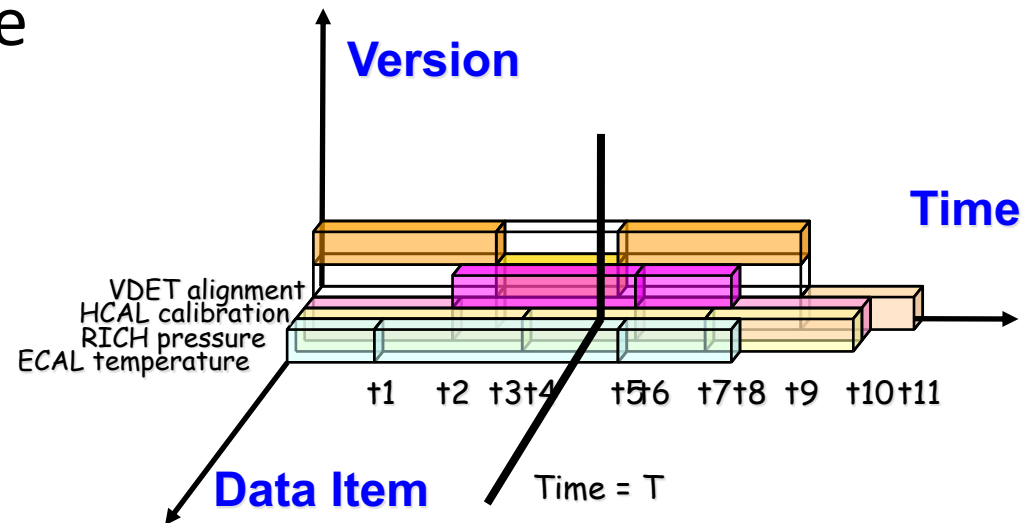  - Finally the track candidates are identified

P= 11 GeV/c

P= 22 GeV/c

P= 7.5 GeV/c

# ATLAS reconstruction procedure

**Inside-out**

Space point formation

↓

Seed finding

↓

Track candidates
(Combinatorial Kalman Filter)

↓

Ambiguity Solving

↓

TRT Extension

**Outside-in**

TRT Segment Finder

↓

Track candidates
(Combinatorial Kalman Filter)

↓

Ambiguity Solving

**Special Algorithms**

Pixel Tracklets

Large impact parameter tracks

Electron tracking

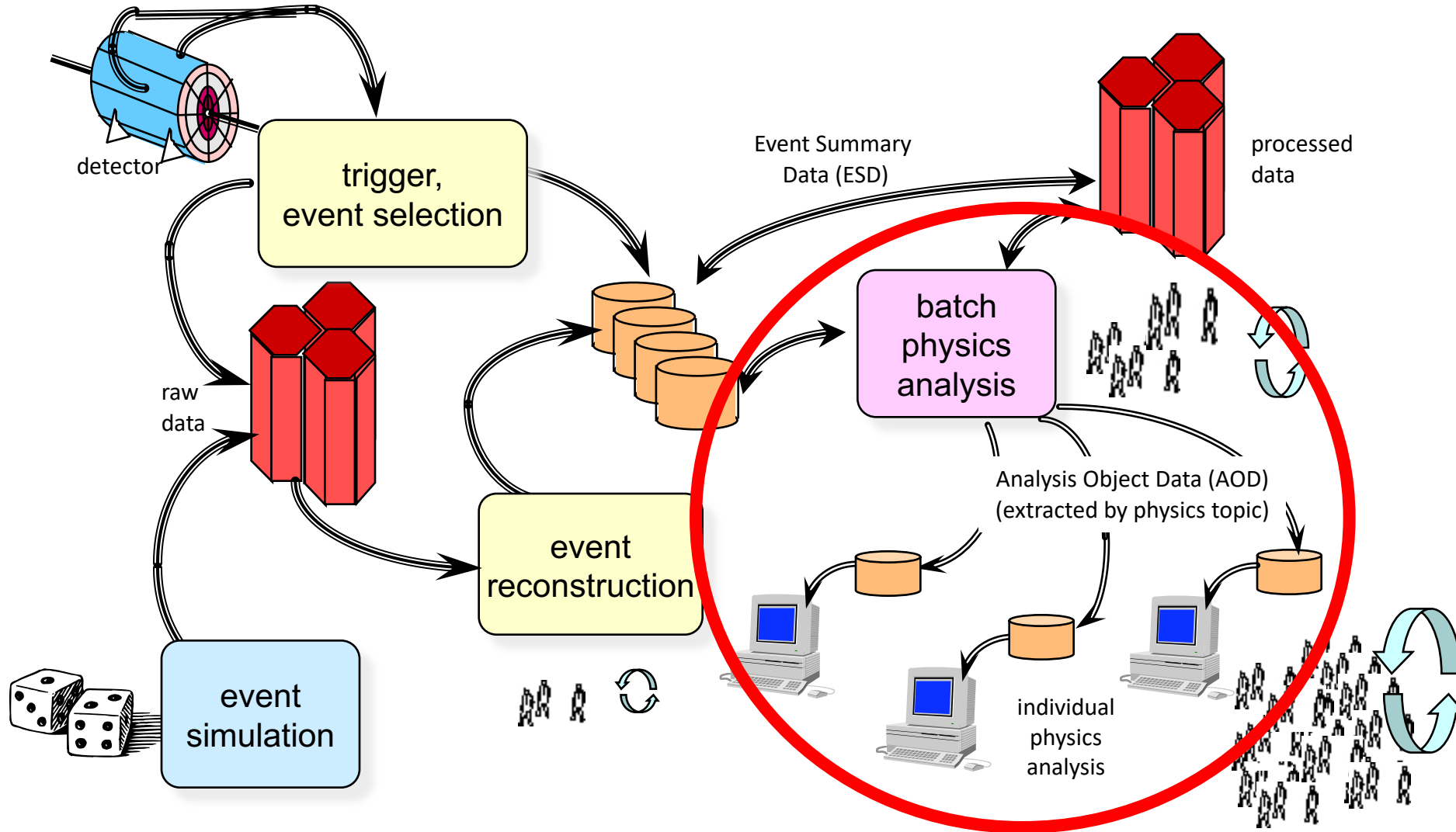# Detector conditions data

- Reflects changes in state of the detector with time

- Event Data cannot be reconstructed or analyzed without it

- Versioning

- Tagging

- Ability to extract slices of data required to run with job

- Long life-time

# Online and offline reconstruction

- Are collisions first-tagged really interesting enough to keep (given capacity constraints)?
  - Online reconstruction – seek to reconstruct 'as much as you can' quickly to enable decision
- Critical part of experiment – collisions which are not recorded are lost
- Later there is more time to reconstruct the contents of a collision – but this is also complex
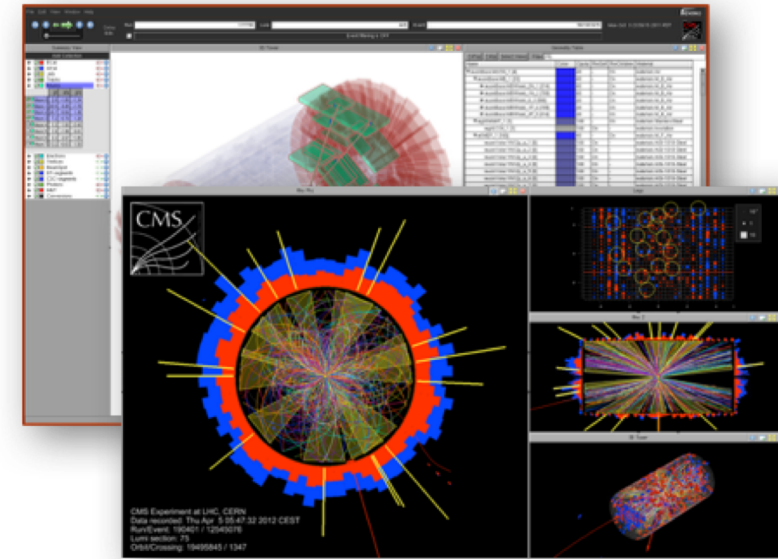
# Processing Stages - Analysis



detector

trigger,
event selection

Event Summary
Data (ESD)

processed
data

raw
data

event
reconstruction

event
simulation

batch
physics
analysis

Analysis Object Data (AOD)
(extracted by physics topic)

individual
physics
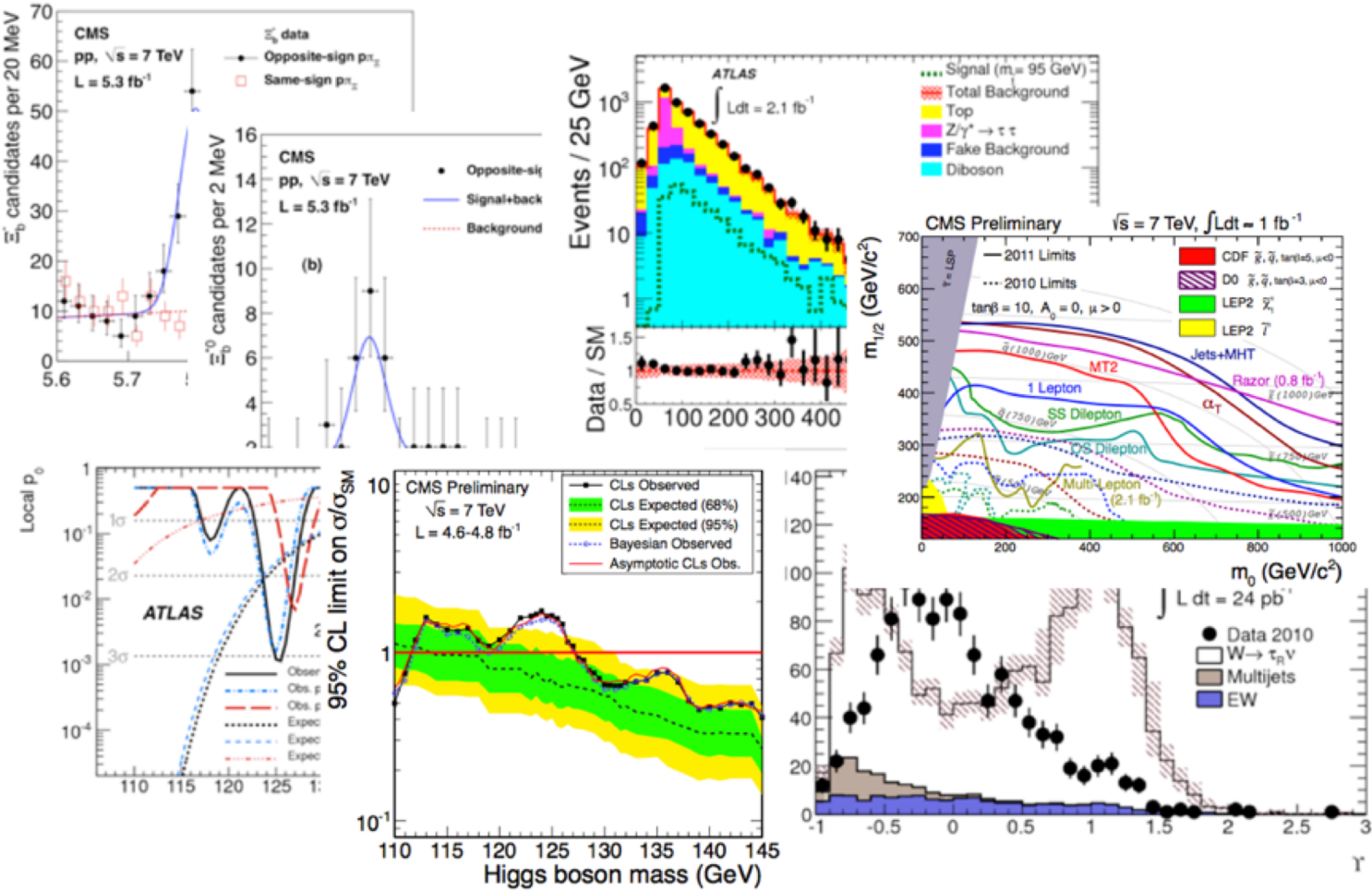analysis

# Data analysis

- <span style="color:red">Uses the results of Reconstruction</span>
  - The products are reconstructed tracks, energy deposits (calorimeters)
  - Hierarchy of data from original (RAW), to summary (AOD)
- Extract observables from data (e.g. invariant mass, particle correlations, …)
  - Understand errors and features, by comparing with simulation
  - Compare with physics hypothesis, theory predictions, explore new physics
  - Programmed mainly in C++ & Python
- An experiment's physics teams use the (large) <span style="color:red">pool of data</span>
  - No longer in one central location, but in multiple locations (cost, space of building, computers, disks, network) …. using the GRID

# ROOT



o "At the root of the experiments", project started in 1995

o Open Source project (LGPL3)
  o mainly written in C++; 4 MLOC

o ROOT provides (amongst other   things):
  o Interactive C++ interpreter (on top of LLVM and Clang)
  o Efficient data storage mechanism;   177 PB LHC data stored in ROOT (2015, now about 500 PB)
  o High-level interface for analysis in C++ and Python (RDataFrame)
  o Advanced statistical analysis algorithms
    o histogramming, fitting, minimization, statistical methods …
  o Scientific visualization: 2D/3D graphics, PDF, Latex
  o Geometrical modeler

# ROOT in plots

# Processing Stages - Simulation

# What is simulation?

- Simulation = doing 'virtual' experiment
- Take all the known physics
- Start from your 'initial condition' (two protons colliding)
- Calculate the 'final state' of your detector to get the 'experimental' results
  - Solve equations of motion, detector electronics response, etc
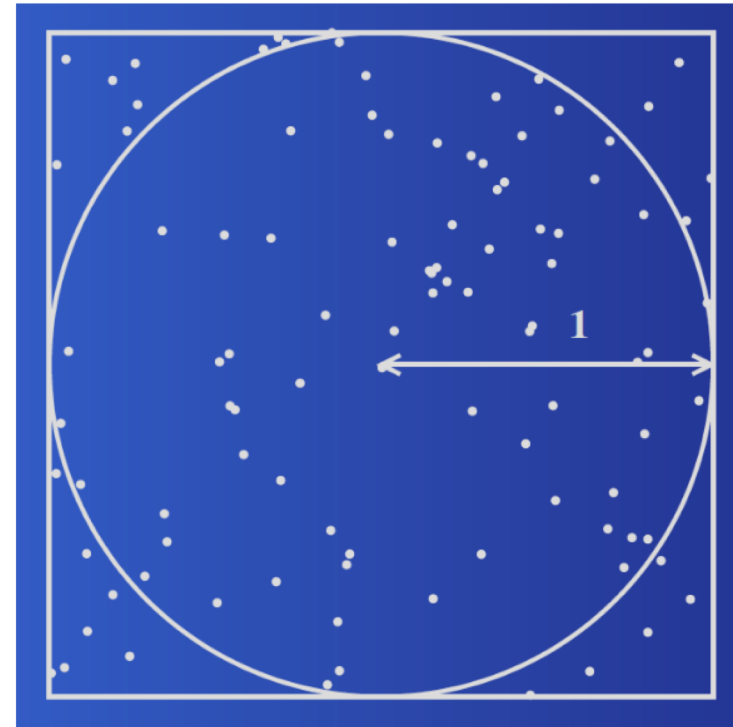- IMPOSSIBLE to be done analytically

# Monte Carlo simulation

- **What is Monte Carlo?**
  - Throwing random numbers
    - to calculate integrals
    - to pick among possible choices
- **Why Monte Carlo?**
  - complexity of the problem
  - lack of analytical description
  - need of randomness like in nature
    - Quantum mechanics: amplitudes => probabilities
      - Noting is certain, but anything that possibly can happen, will!
      - Want to generate events in as much detail as possible
      - get average and fluctuations right
      - make random choices, ~as in nature

# Laplace method of calculating π (1886)

- Area of the square = 4
- Area of the circle = π
- Probability of random points inside the circle = π / 4

- Random points : $N$
- Random points inside circle : $N_c$
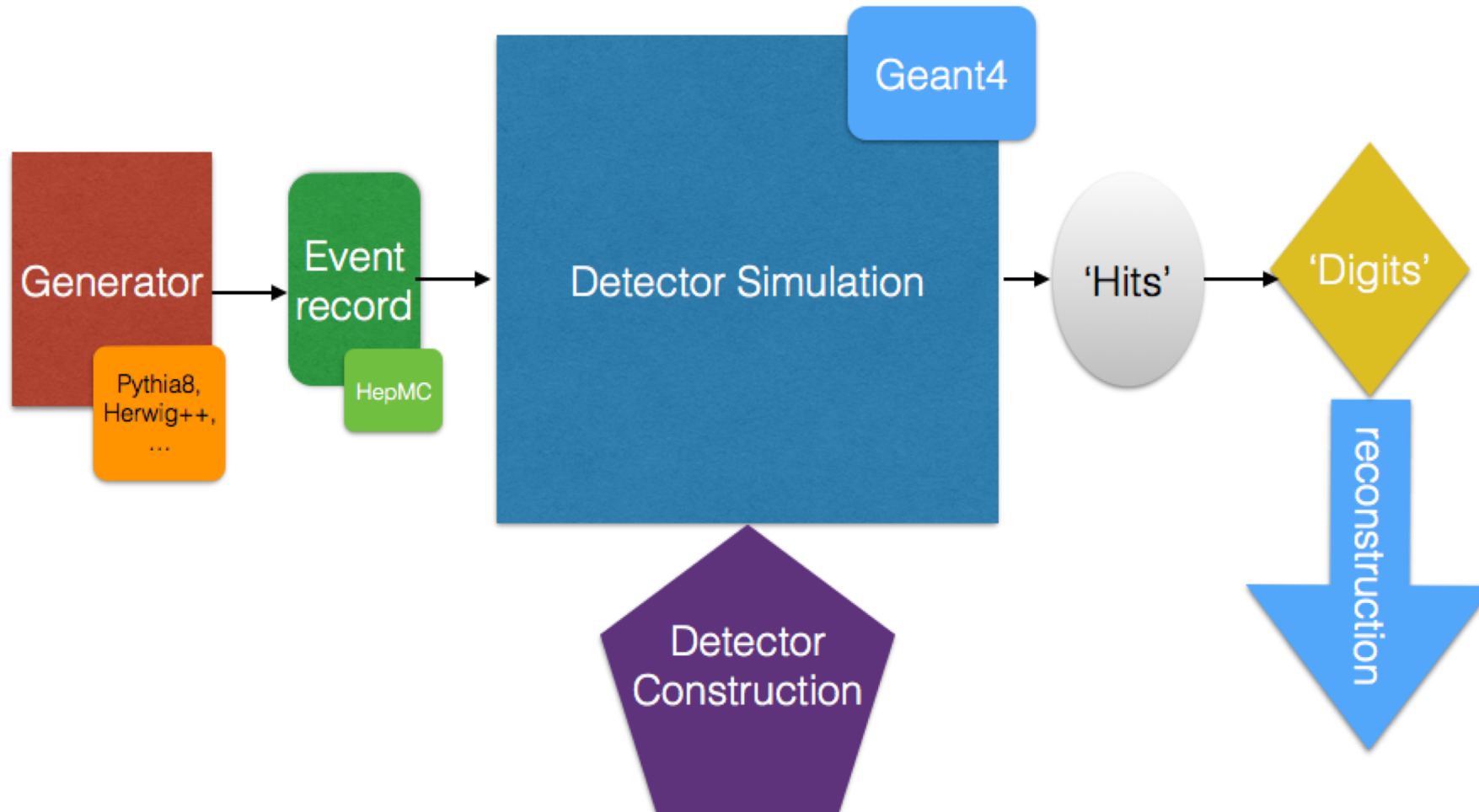
$π \sim 4 \, N_c \, / \, N$

# Why do we need simulation?

- To design the apparatus (detector) to fulfill its role
- To prepare the reconstruction and analysis of results
  - Training on 'known' (simulated) events (MC 'truth')
- To understand the results
  - We need to know what to expect to
    - Verify existing models
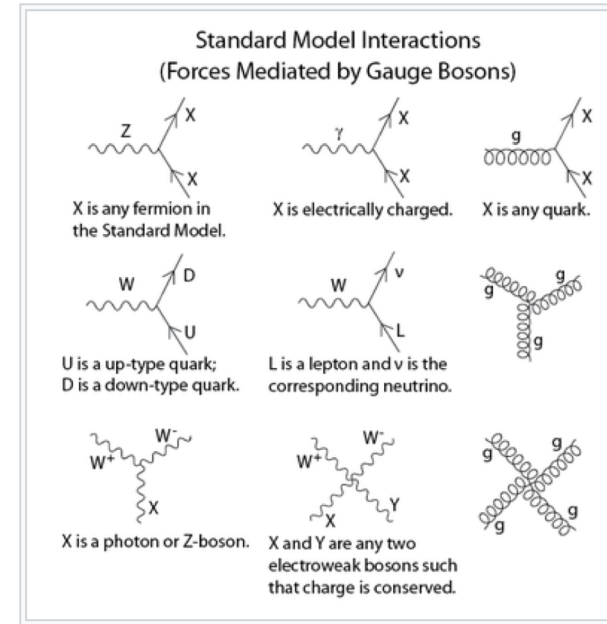    - Find new physics
  - Understand systematic errors

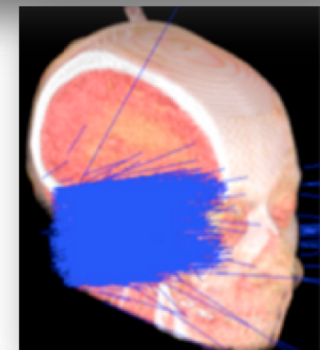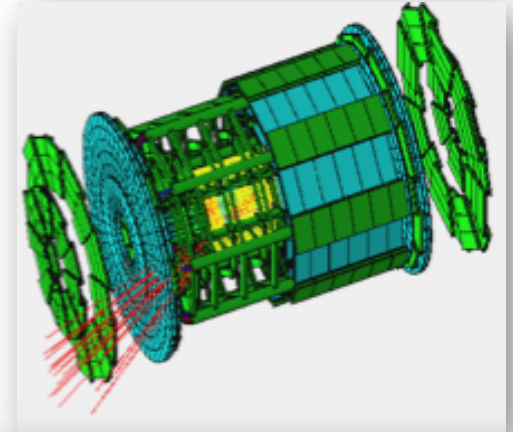# Simulation chain for HEP experiments

# Monte Carlo generators

- Simulate particles reaction in vacuum
  - knows nothing about the surrounding detector
- All Standard Model processes are included
- No propagation of particles, just generation of the products of the 'primary' collision
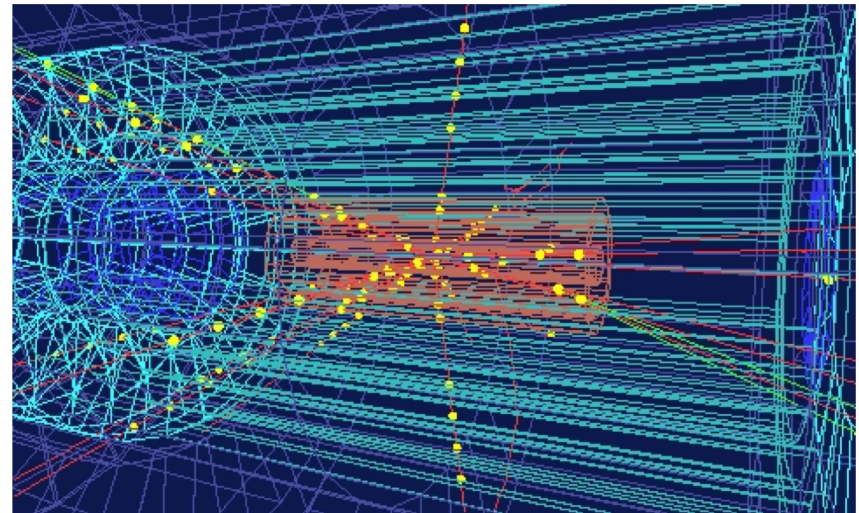- The output of the 'generators' is the input to the 'transport' code

# Transport Code: Geant4



- Geant4 is a toolkit (C++) for the simulation of the passage of particles through matter.

- Its areas of application include high energy, nuclear and accelerator physics, as well as studies in medical and space science

- In HEP has been successfully   employed for
  - Detector design
  - Calibration/alignment
  - Data analysis
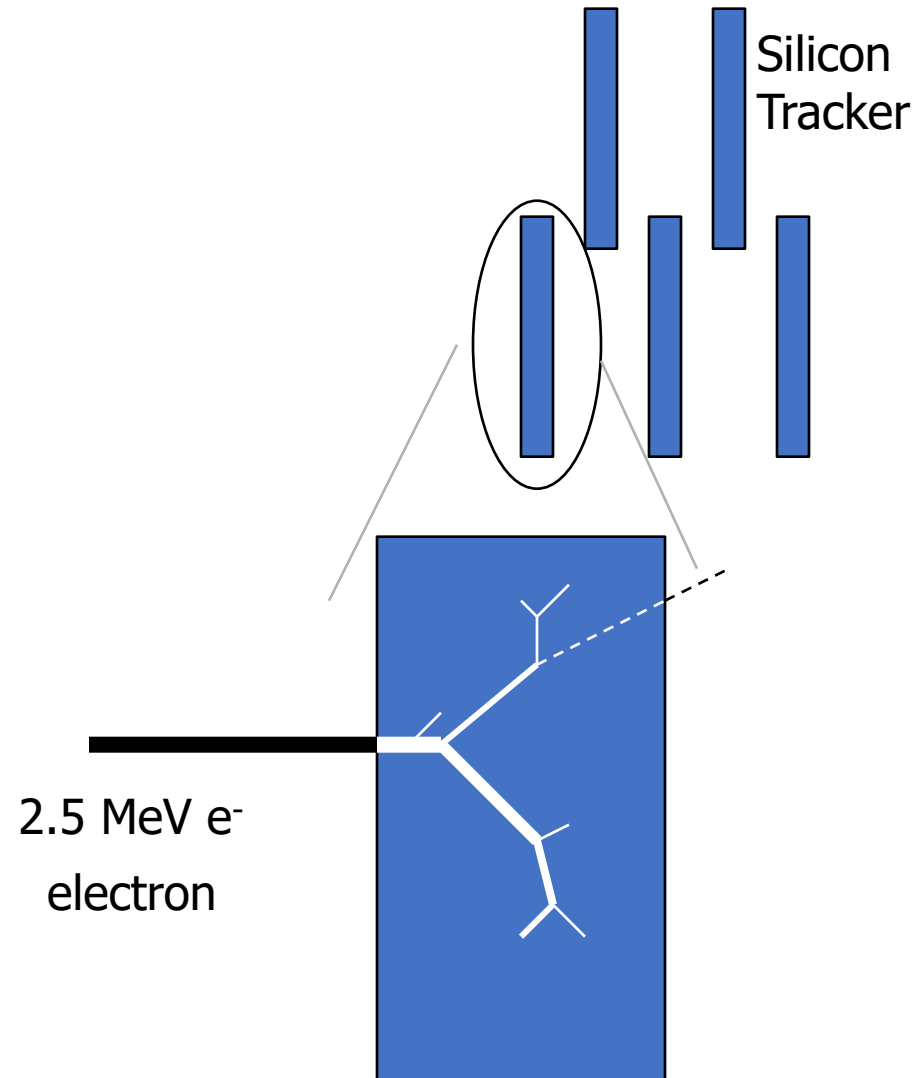


XMM-Newton

# What does Geant4 do?

- 'propagates' particles through geometrical structures of materials, including magnetic field

- simulates processes the particles undergo
  - creates secondary particles
  - decays particles

- calculates the deposited energy along the trajectories and allows to store the information for further processing ('hits')
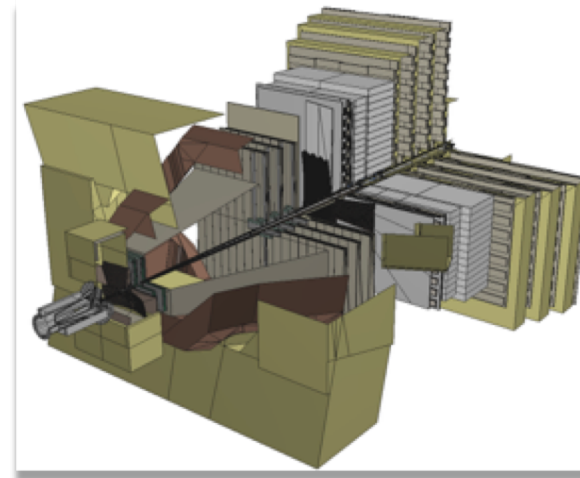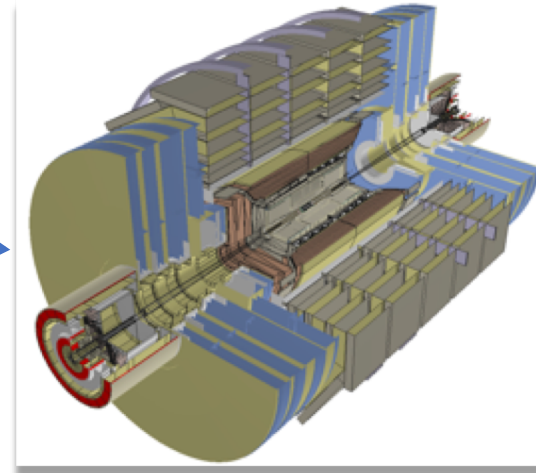
# Simulation ingredients

- ## We model
  - ### Detector's Geometry
    - Shape, Location, Material

  - ### Physics interactions
    - All known processes
      - Electromagnetic
      - Nuclear (strong)
      - Weak (decay
  - we 'shoot' particles and 'propagate' them through the modeled detector

$$\sigma_{total} = \Sigma\ \sigma_{per\text{-}interaction}$$
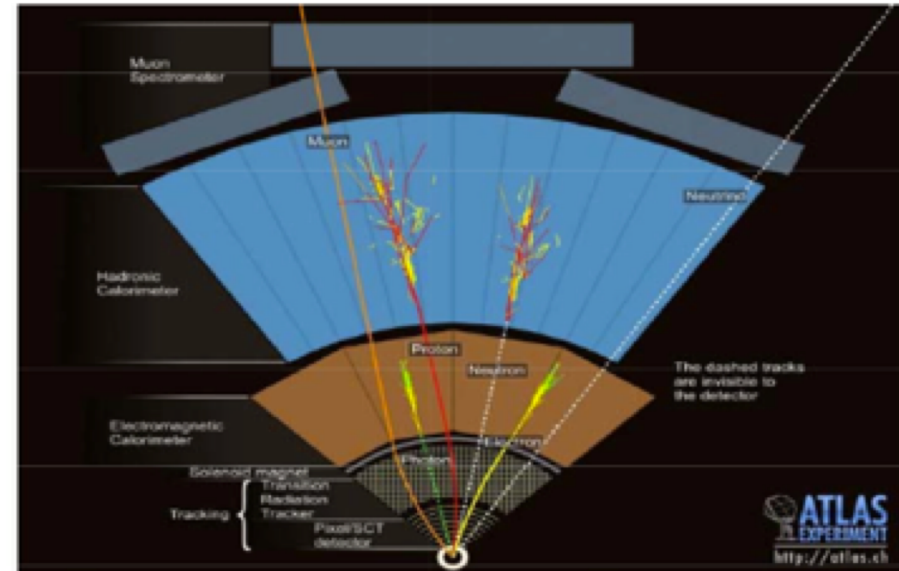
Silicon Tracker

2.5 MeV e⁻ electron

# Geometry and Materials

- How to implement (efficiently) this in your computer program?
- You need 'bricks'
  - 'solids', 'shapes'
  - you need to position them
  - you want to 'reuse' as much as possible the same 'templates'
- Database of Materials
  - National Institute of Standards (NIST)
- Magnetic Fields
  - numerical integration of the equation of motion (Runge-Kutta method)
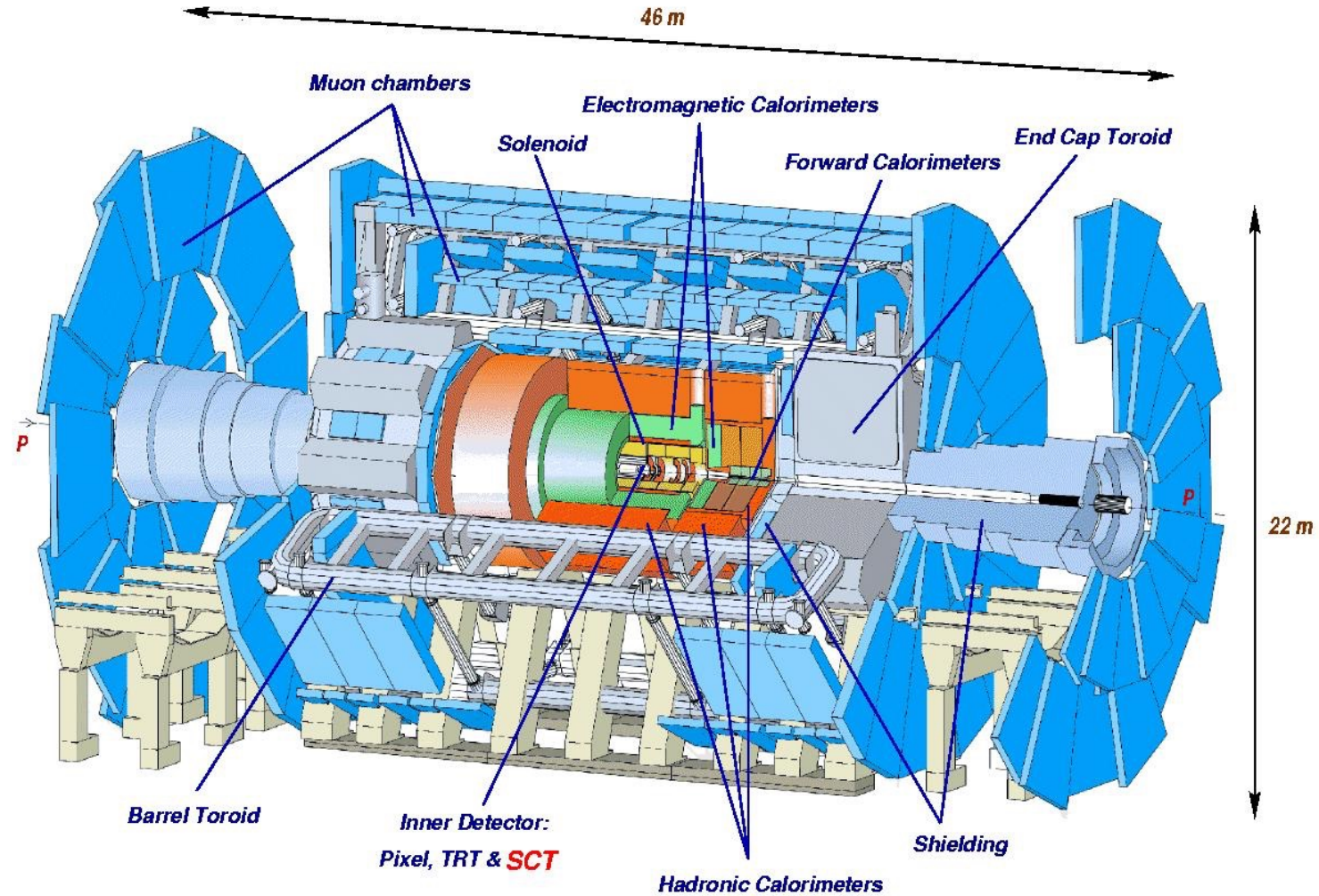
# Physics…

- What happens to particles in matter?
- We want to model the physics we know
  - each possible physics process provides the "interaction length" compared with distance to next geometrical boundary
    - the smallest wins
  - generating a "final state" and secondaries tracks
- Electromagnetic
  - gammas and charged particles
- Hadronic
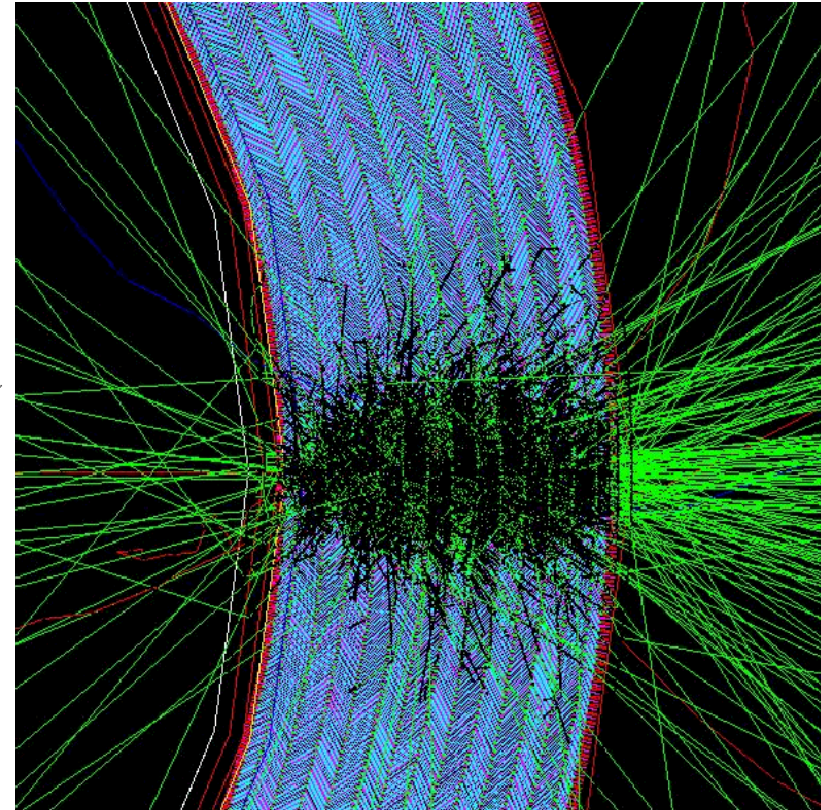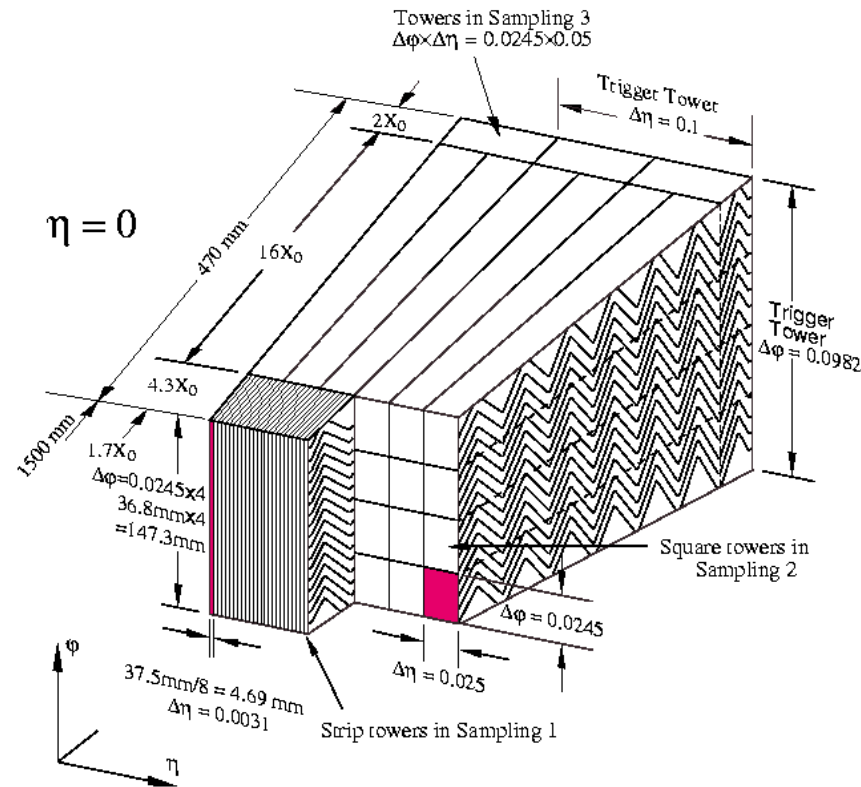  - neutrons, mesons (K,π), muons, …



Because of the detailed geometries, the detailed physics and the required precision the simulation is very CPU hungry

# ATLAS

# ATLAS Calorimeter (a very, very small part of it)





. Kordas "Geant4 for the ATLAS EM calo" — CALOR2000, Annecy, 12 October 2000
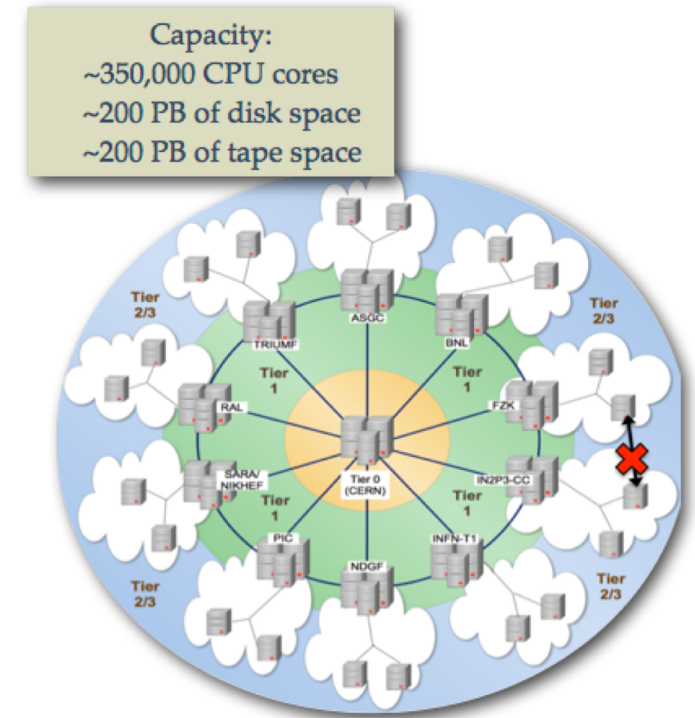
# Summary: data rates

o Particle beams cross every 25 ns (40 MHz)

    o Up to 25 particle collisions   per beam crossing (for Run2, higher for Run3)

    o Up to $10^9$ collisions   per second

o Basically 2 event filter/trigger levels

    o Hardware trigger (e.g. FPGA)

    o Software trigger (PC farm)

    o Data processing starts at readout

    o Reducing $10^9$ p-p collisions per second to O(1000)

o Raw data to be stored permanently: >15 PB/year

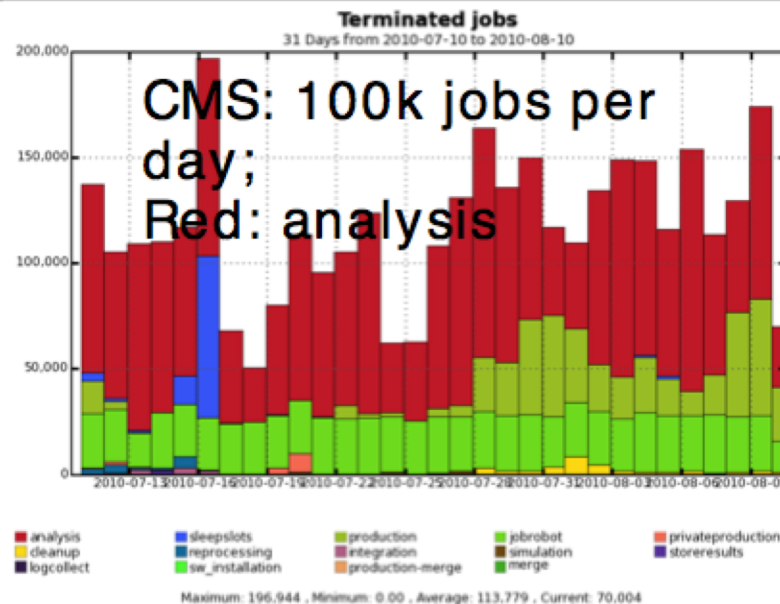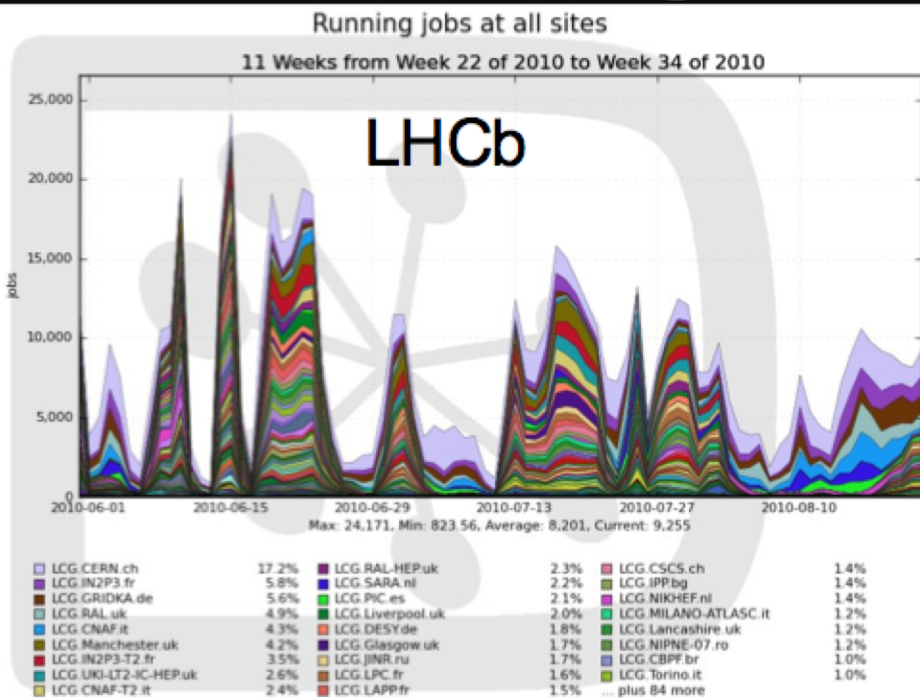| Physics Process | Events/s |
|---|---|
| Inelastic p-p scattering | $10^8$ |
| $b$ | $10^6$ |
| $W \rightarrow e\nu$ ; $W \rightarrow \mu\nu$ ; $W \rightarrow \tau\nu$ | 20 |
| $Z \rightarrow ee$ ; $Z \rightarrow \mu\mu$ ; $Z \rightarrow \tau\tau$ | 2 |
| $t$ | 1 |
| Higgs boson (all; $m_H = 120$ GeV) | 0.04 |
| Higgs boson (simple signatures) | 0.0003 |

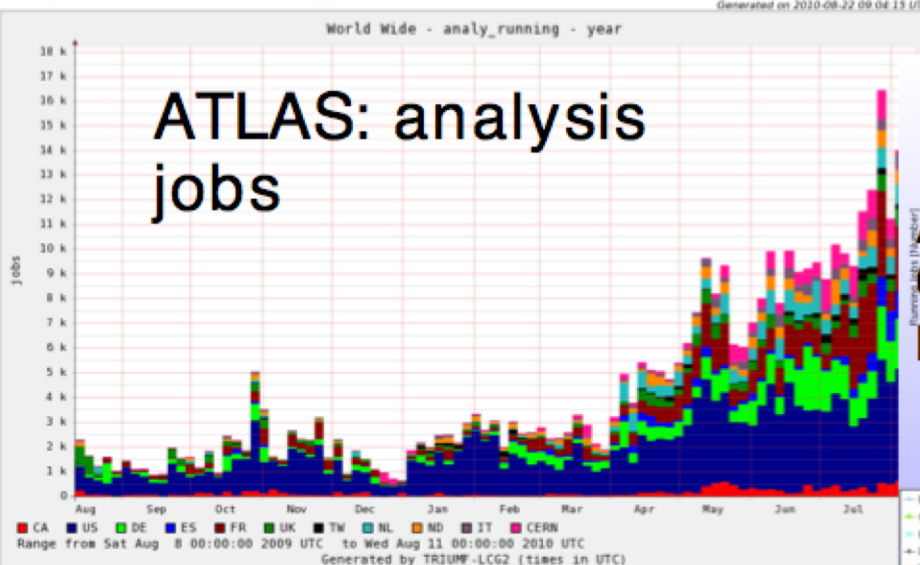This is our Big Data problem!!

# Big Data requires Big Computing

o The LHC experiments rely on distributed computing resources:
  o WLCG - a global solution, based on the Grid technologies/middleware.
    o distributing the data for processing, user access, local analysis facilities etc.
    o at time of inception envisaged as the seed for   global adoption of the technologies

o Tiered structure
  o Tier-0 at CERN: the central facility for   data processing and archival
  o 11 Tier-1s: big computing centers with   high quality of service used for most   complex/intensive processing operations   and archival
  o ~140 Tier-2s: computing centers across the   world used primarily for data analysis and   simulation.

o So far computing was not a limiting factor for the   Physics program of the LHC experiments
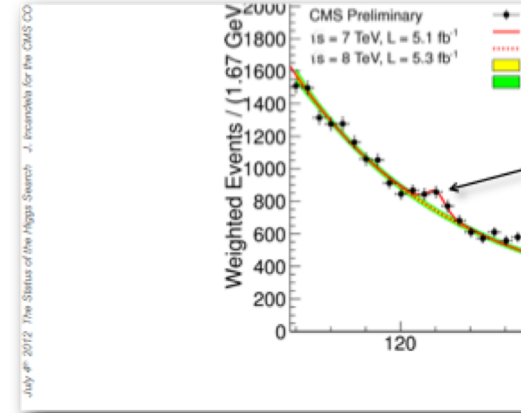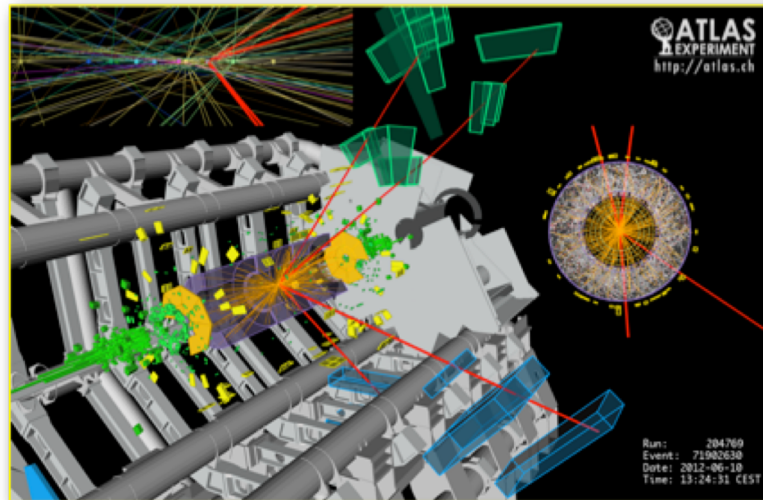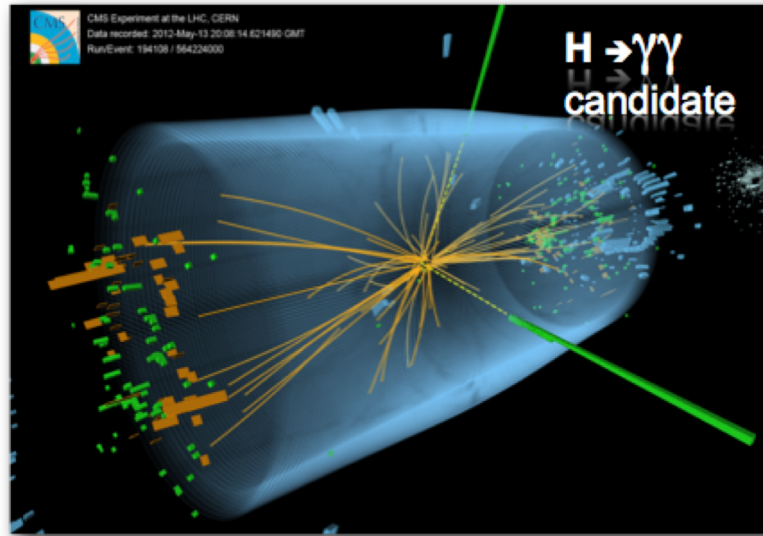


Capacity:
~350,000 CPU cores
~200 PB of disk space
~200 PB of tape space

# Running jobs on LCG

# A Success Story!



45

# Challenges for HEP Software

- High-luminosity LHC will produce 7x-10x today's event rate
  - More precise Higgs physics (5x), rare signals, new physics
  - Timescale: 2017-2018
  - Constant computing budget
  - Technology evolves, but we need to be able to make use of it
    - Massive parallelism, AI, hybrid computing, …
- Huge pressure for both experiment software systems and common software
  - Important R&D ongoing for experiment upgrades
    - Hardware and software
  - R&D for the common simulation tools

# Conclusion

- Modern HEP experiments would be impossible without computing
  - Online triggering and selection
  - Offline reconstruction, analysis and simulation
- Huge data volumes
- Distributed processing