

Analysis tools

Pedro Ferreira da Silva, psilva@cern.ch - (CERN)

15th Course on Physics at the LHC

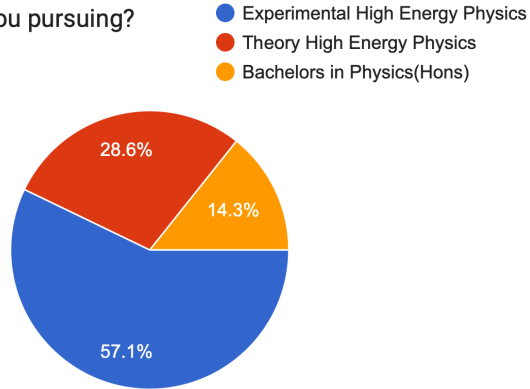
Monday 23rd March 2026



First things first: your interests!

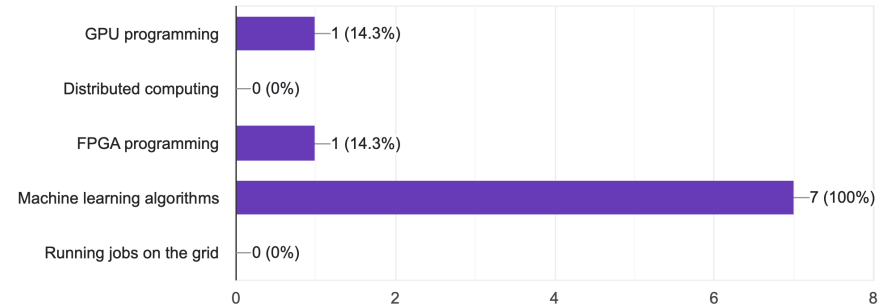
Which field are you pursuing?

7 responses



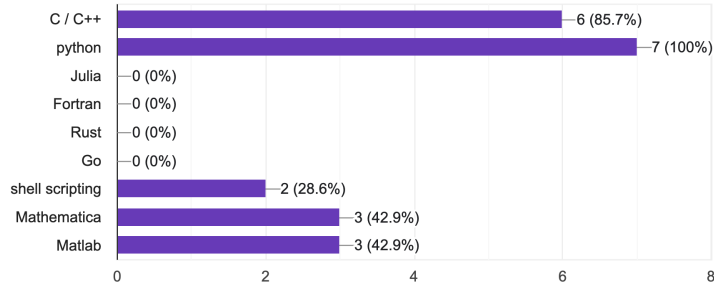
Are you familiar with any of these?

7 responses



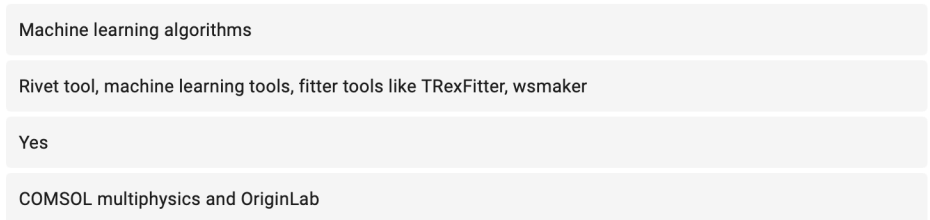
Which programming languages are you familiar with?

7 responses



What would you be interested in learning about in an Analysis tools session?

4 responses



How are data created at the LHC?

The high level trigger

Offline analysis

Hands-on

How are data created at the LHC?

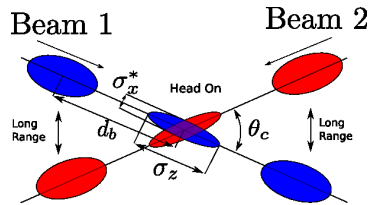
The high level trigger

Offline analysis

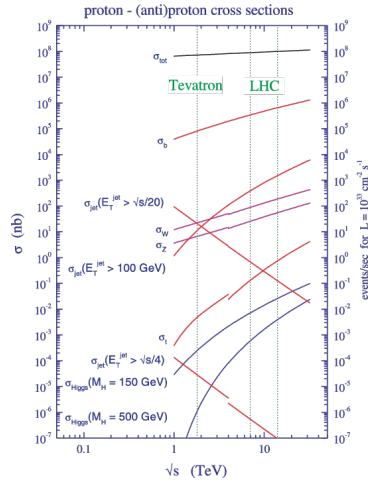
Hands-on

How are data created at the LHC

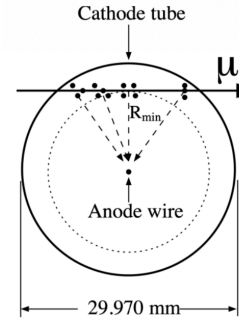
Bunches cross at the detector centre



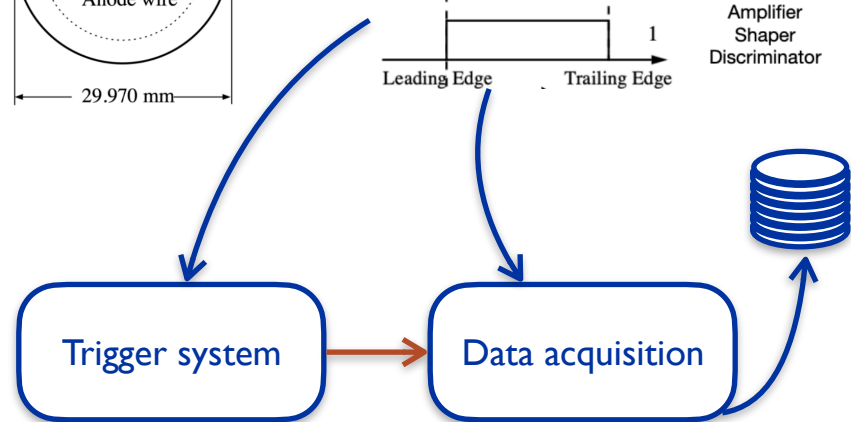
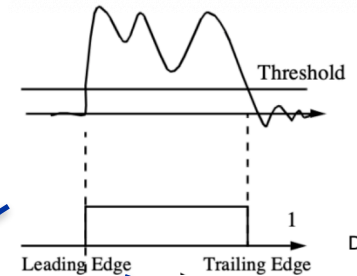
Beams interact at a rate dictated by the interaction cross sections



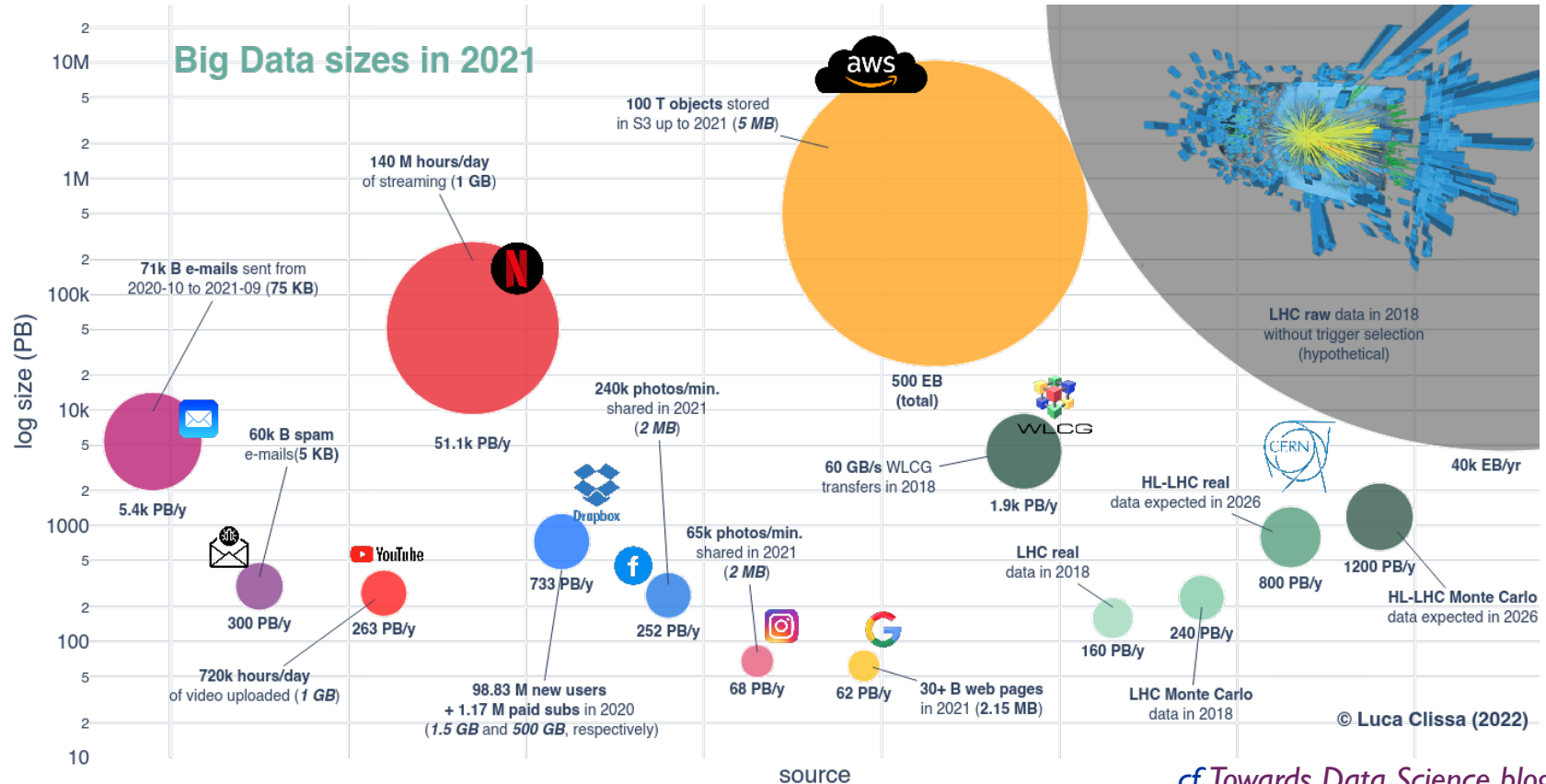
Particles cross the detector interacting with the material



Readout electronics identifies regions of interest and measures amplitudes, regional sums, times



How large are data?



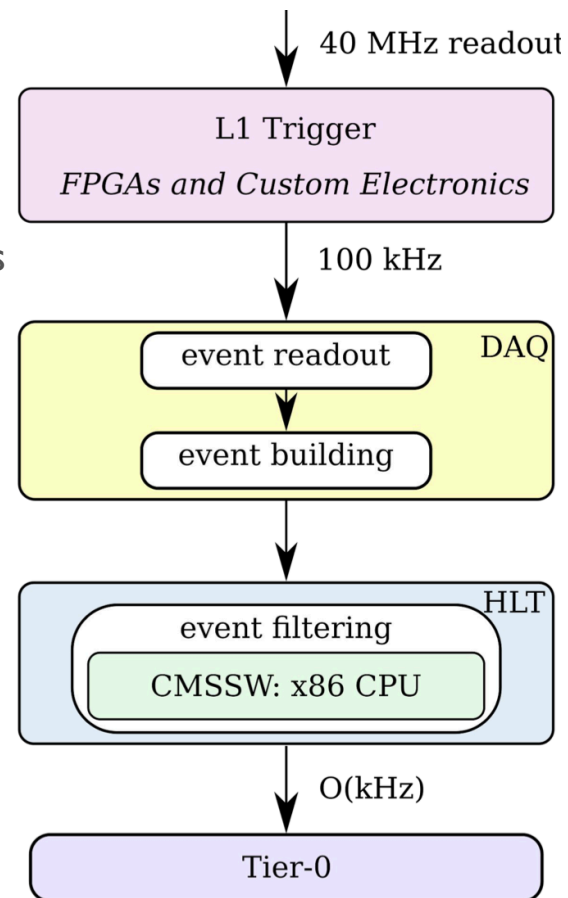
Several transformations until a physics object is formed

Data is produced at 40 MHz rate

- Undergoes stringent selections at Level-I and later HLT
- Increased complexity in reconstruction and calibration requirements
- Permanent storage to disk of RAW data for later reconstruction

With respect to Phase I of the LHC, Phase II foresees

- an increase of the LI rate to **750 kHz**
- more demanding reconstruction of more granular detectors in the presence of up to 200 pileup events
- quite a challenge from DAQ/Trigger point of view!



Several transformations until a physics object is formed

Local reconstruction

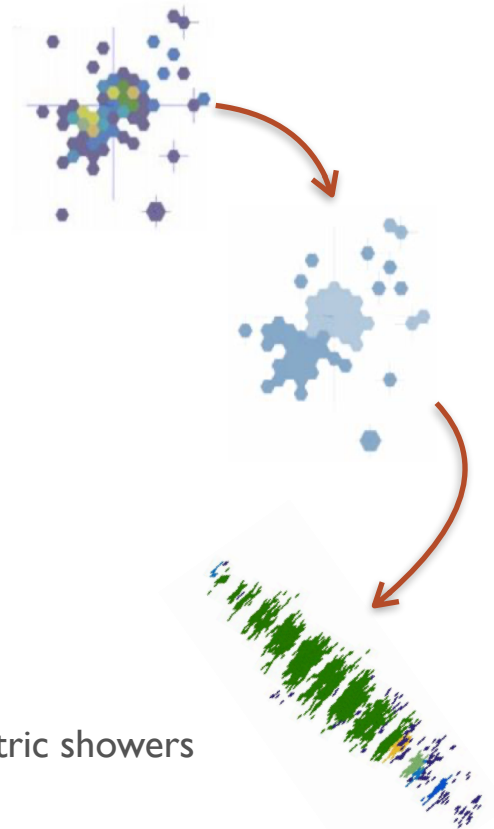
- Binary data (=DIGs) conversion to floating-point precision (=Hits)
- Usually a complicated function to convert amplitudes to energy or time
- Hits are typically a 5D object (E,t,x,y,z) with some associated quality flag

Clustering

- Reduction of complexity and noise
- Algorithms typically proceed by identifying seeds and aggregating neighbors
- Metric used to decide when to stop aggregation

Building blocks for Particle Flow

- Clusters (sometimes hits) are used as input for tracking and building calorimetric showers
- A preliminary identification can be assigned in some cases



How are data created at the LHC?

The high level trigger

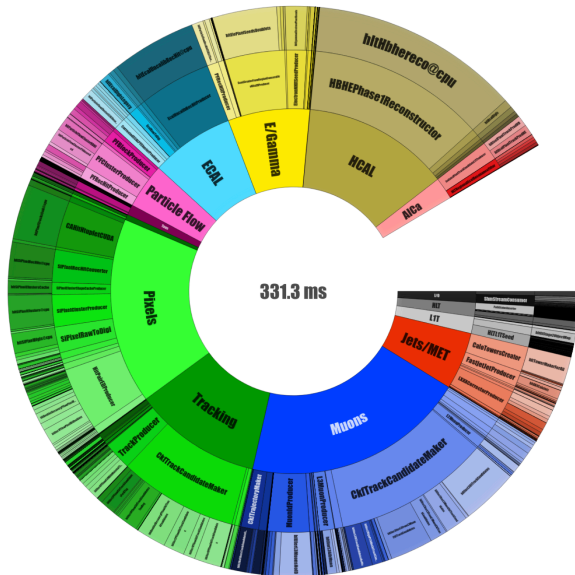
Offline analysis

Hands-on

The high level trigger

Full event reconstruction is performed in the HLT farms

- Several filter modules share the output reconstructed objects
- Used to decide whether an event should be kept or not
- The collection of trigger bits is called the “Trigger Menu”



Poster by Marco Montella

THE ATLAS RUN 3 TRIGGER MENU

MARCO MONTELLA [1]
ON BEHALF OF THE ATLAS TRIGGER GROUP

OUR SIGNATURES

- ELECTRONS** 270 HZ
- MUONS** 290 HZ
- TAUS** 160 HZ
- JETS & MET** 630 HZ

Specialty TRIGGERS

ABOUT US

ATLAS runs a two-level triggering strategy:

- **L1** → Hardware-based, coarse reconstruction, total accepted rate is 50-100 kHz
- **HLT** → Software-based, reconstruction precision approaching offline reconstruction

The Trigger Menu is limited by:

- **L1 RATE** → Constrained by dead time impacts the range of physics accessible
- **HLT CPU** → Limits the execution rate of high-precision reconstruction algorithms
- **TO CPU** → Limits the data volume that can be reconstructed promptly for endpoint analysis

L1/HLT limitations scale with luminosity:

- **END OF FILL** → Enhance signatures limited by L1 rate and/or HLT rate & CPU

OUR STREAMS

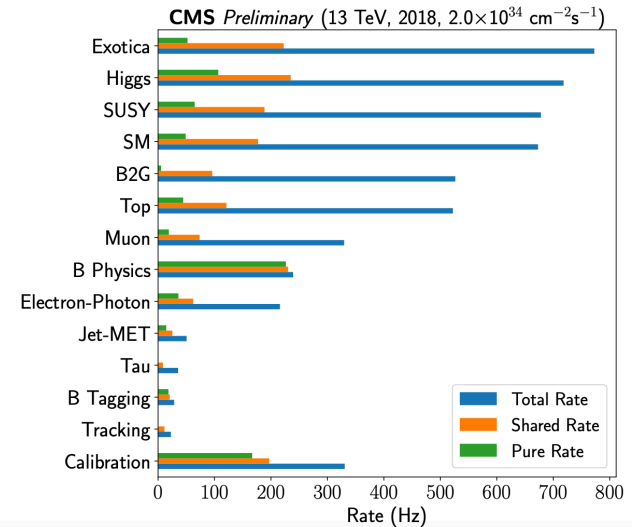
- MAIN** 1.7 KHZ
For prompt reconstruction Rate limited by L1, CPU & TO Resources
- DELAYED** 900 HZ
Historic B-Physics Delayed Processing when TO CPU available
- TRIGGER-LEVEL ANALYSIS (TLA)** 5+ KHZ
Historic B-Physics b-tag Reduced event content, HLT Objects only Minimal burden on bandwidth
- PARTIAL EVENT BUILDING (PEB)** 0.5 KHZ
Historic B-Physics b-tag Regional data around near physics objects identified by trigger

EMERGING JETS	10 HZ	DISAPPEARING TRACK	4 HZ	HEAVY IONS: RUN 2 THRESHOLDS PRESERVED!
HIGHLY IONIZING TRACK	5 HZ	DISPLACED OBJECTS	40 HZ	
ISOLATED TRACK	1 HZ	PARTIAL EVENT BUILDING	200 HZ	

The high level trigger - cont.

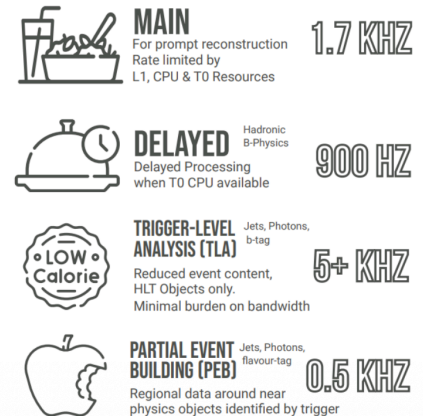
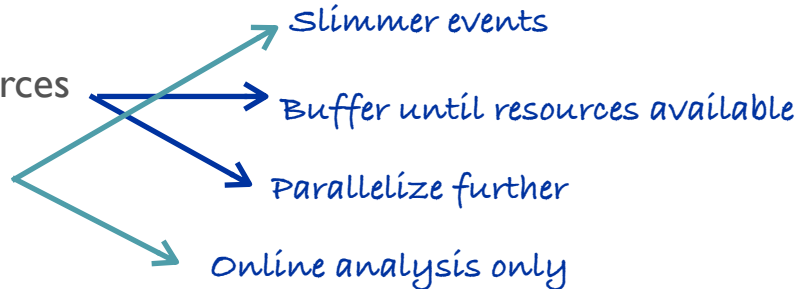
HLT also ends up limiting the capacity of storing all that data

- not more than 1.5 kHz can be output by the HLT farms
- Every extra selection (HLT bit) needs to be negotiated between different physics and physics objects groups
- Often a trigger bit is prescaled (allow only “one out of N”)



Main bottlenecks

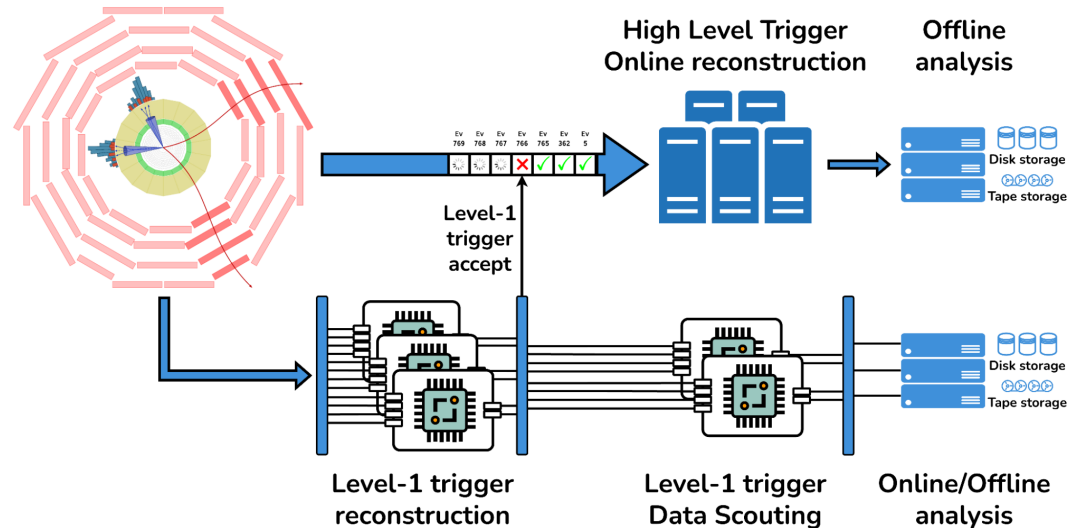
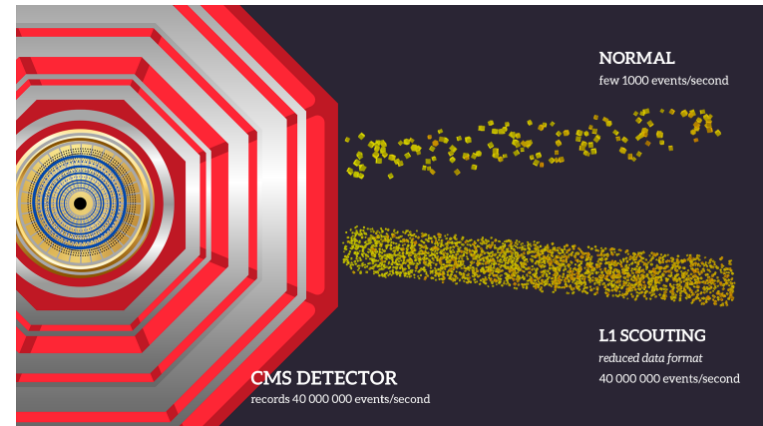
- Computing processing resources
- Limited storage



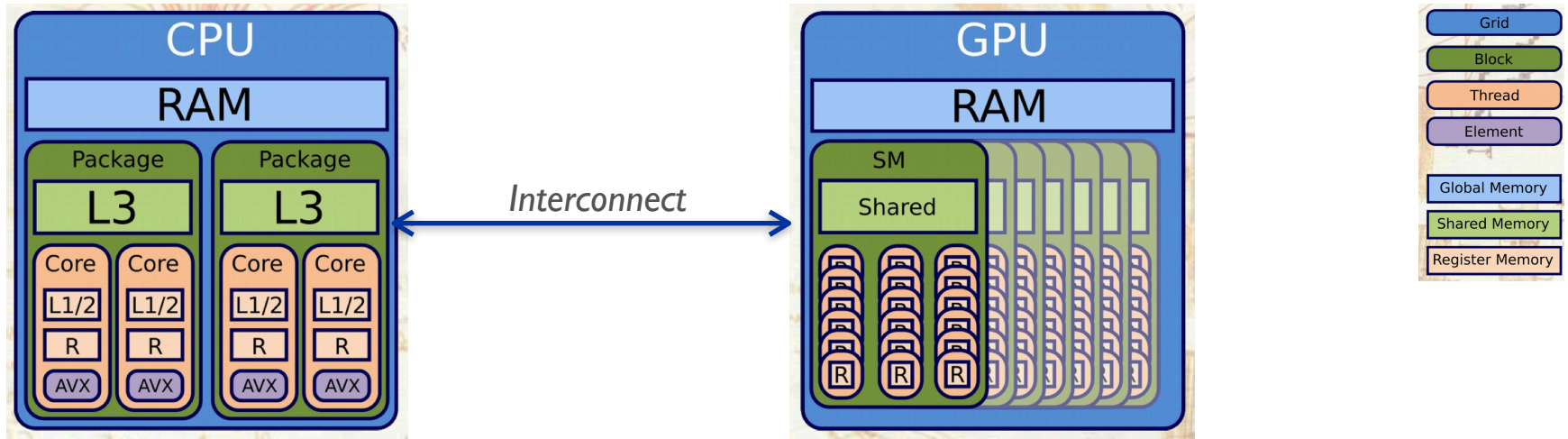
Data scouting at source

Direct scouting at L1 gives access to full 40 MHz rate

- Quite feasible for some analysis with less stringent calibration / resolution degradation (e.g. using muons)
- Opens the door to analyses which become ineffective with traditional triggers
- Potentially gives access to the analysis of multiple bunch crossings simultaneously
- Recent example in [arxiv:2601.20063](https://arxiv.org/abs/2601.20063)



Heterogenous computing approach



Latency

- Large cache
- Average bandwidth between CPU and host memory
- Complex control unit
- Low core count / Powerful Arithmetic Logic Unit (ALU)

Throughput

- Small cache
- High bandwidth between GPU cores and memory
- No complex control unit
- High core count

Heterogenous computing approach - cont.

Underlying objects need to be adapted for parallel computing

- Structure of arrays (SoA) divide a buffer of data into runtime sized columns
- Buffers host memory and are pinned to host or device
- Allocation and handling can be made using

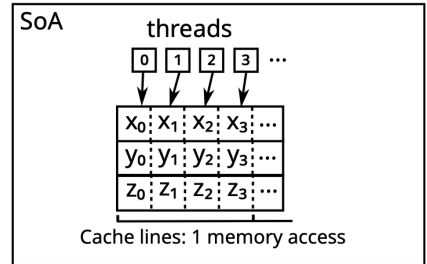
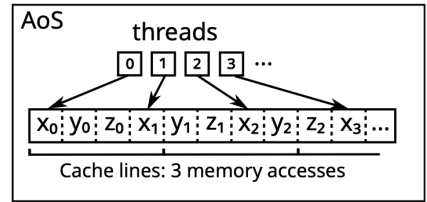
CUDA



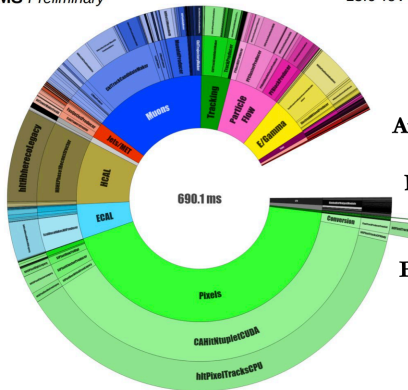
ALPAKA



AoS and SoA



CMS Preliminary CPU ONLY 13.6 TeV



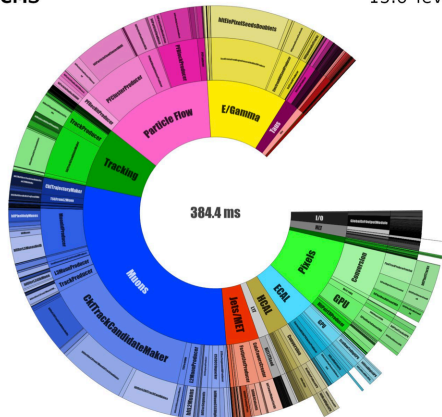
Event throughput +80%

Average time per event -40%

Power consumption: -30%

Experience gained towards heterogeneous offline reconstruction: Priceless

CMS Heterogeneous 13.6 TeV



Currently ~40% of the CMS HLT uses this approach

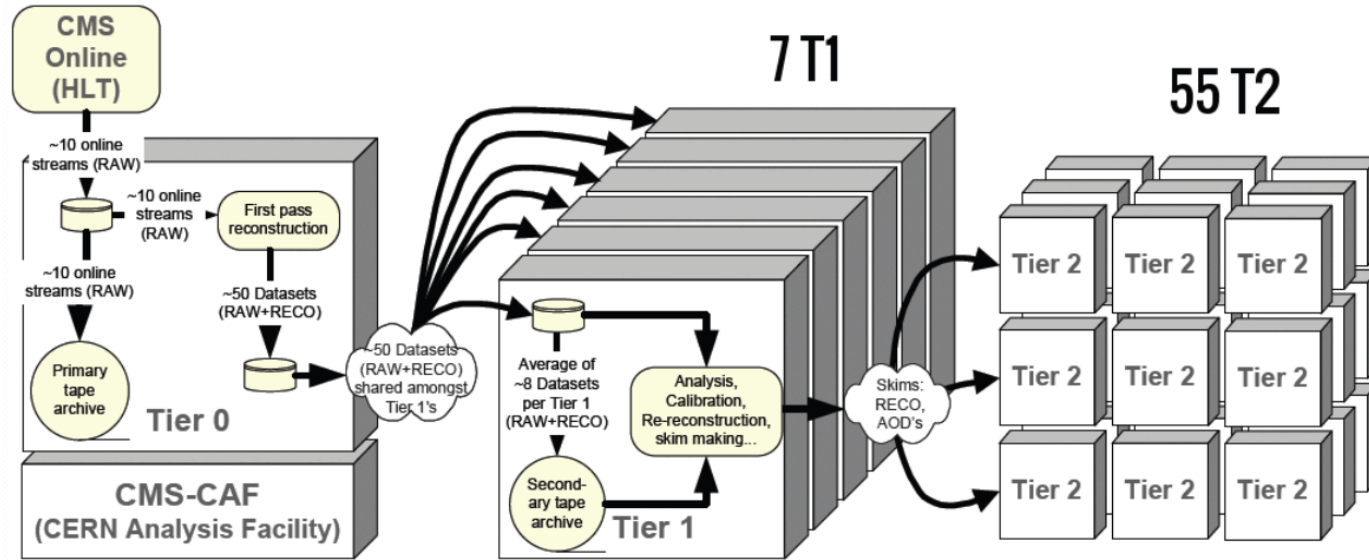
How are data created at the LHC?

The high level trigger

Offline analysis

Hands-on

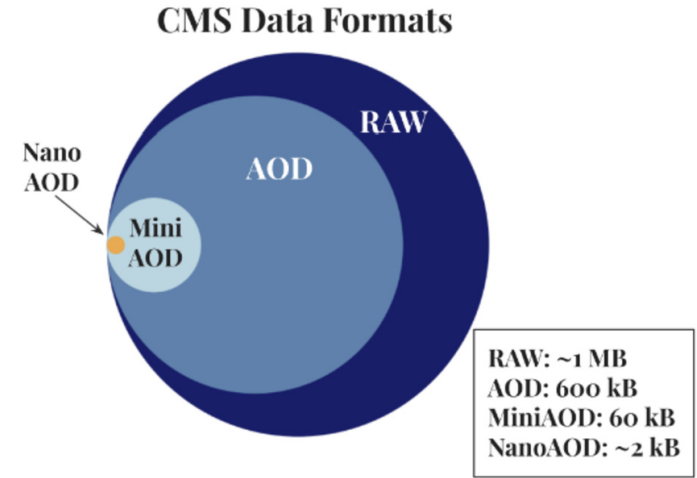
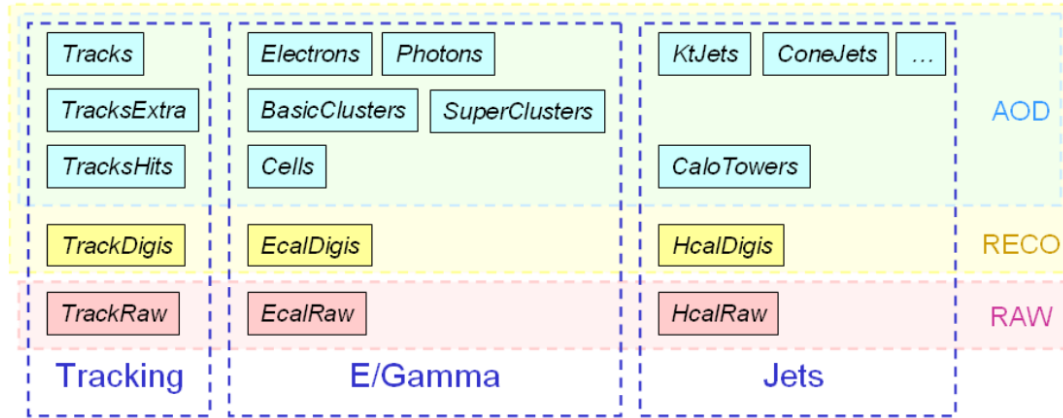
From HLT onwards: the computing ...



Grid computing is structured in Tiers

- CERN always preserves a copy of all events (RAW) on tape and first pass reconstruction on disk
- Tier 1 are used to run re-reconstruction but also analysis, calibration and further selection
- Tier 2 are used mostly for analysis
- Usage relies on World LHC Computing Grid ([WLCG](#))

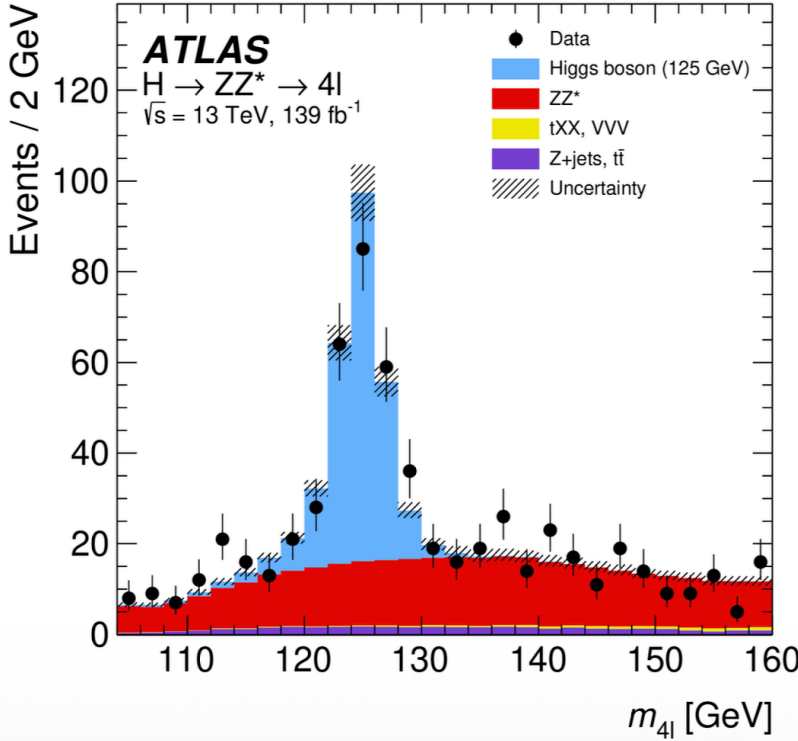
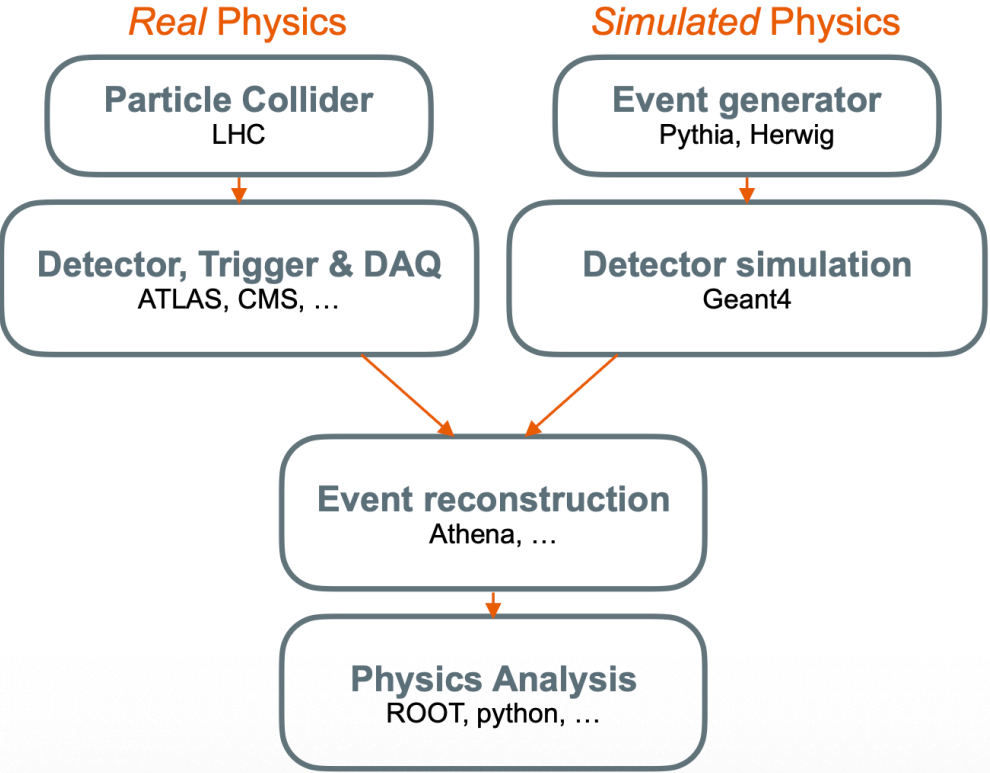
From HLT onwards: the computing ... and data tiers



LHC experiments use a variety of data tiers

- Depending on the analysis complexity it may or may not need some of the “core objects”
- It's fair to say that the vast majority only need the calibrated, final, physics objects

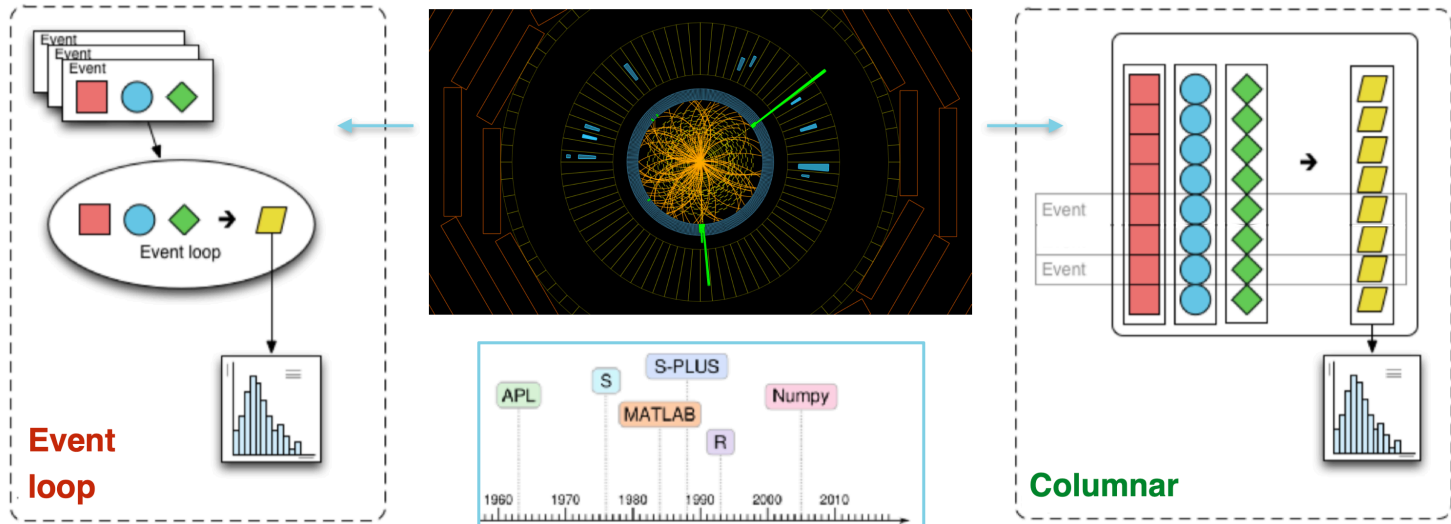
What's in a paper plot?



High energy physics analysis: some challenges

Data access is organized in an efficient way but it doesn't reflect their variety

- Events are highly structured and objects are complex with references to each other
- Complex objects can be inefficient on disk and memory (if not careful)



When doing analysis: understand where data is and how can I analyze it efficiently

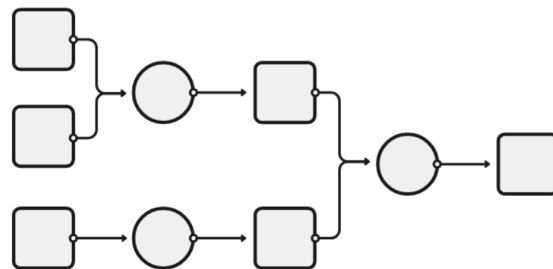
Evolutions in abstraction

Since some years HEP started moving to columnar analysis

- Parallelize as much as possible: multithreaded, GPU, distributed computing
- Actions are pre-defined and chained: execution is delayed

	Electrons	Muons	Jets
Event #1	■ ■		● ● ● ● ● ● ● ●
Event #2		■	● ● ●
Event #3	■ ■ ■	■ ■	● ● ● ● ● ● ● ●
		...	

```
for event in EVENTS:  
  for e in ELECTRONS: do_something  
  for m in MUONS: do_something  
  ...  
  with selected_objects: do_something_else  
  once finished, move to next event...
```



Scheduler:



Evolutions in abstraction

Since some years HEP started moving to columnar analysis

- Parallelize as much as possible: multithreaded, GPU, distributed computing
- Actions are pre-defined and chained: execution is delayed

```
void MyClass::Loop() {
    size_t nEvents;
    // load...

    for (Long64_t iEvent=0; iEvent<nEvents; iEvent++) {
        double MET_pt;
        int nElectron;
        double * Electron_pt;
        double * Electron_eta;
        // load...

        if ( MET_pt > 100. ) continue;

        for(size_t iEl=0; iEl<nElectron; ++iEl) {
            if ( Electron_pt[iEl] > 30. ) {
                hist->Fill(Electron_eta[iEl]);
            }
        }
    }
}
```

```
cut = (events.MET.pt < 100.) & (events.Electron.pt > 30.)
hist.fill(eta=events.Electron.eta[cut].flatten())
```

```
hist.compute()
```

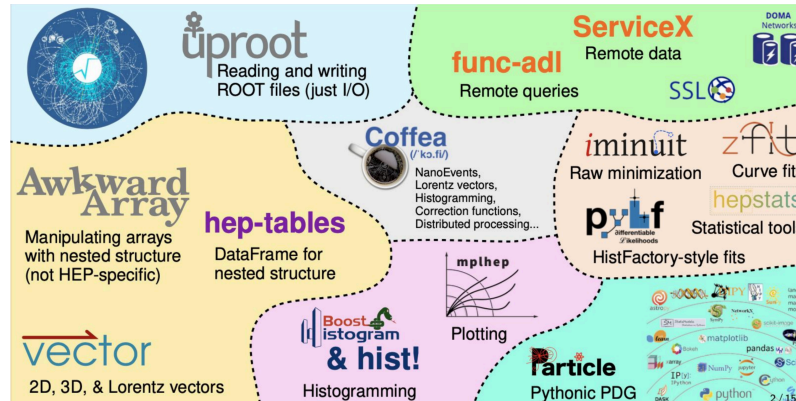
Current paradigms

Modern analysis software is factored and comprises a complex ecosystem

- each step should work well in isolation but needs to be chained in a system
- But require availability of many different tools and programming languages

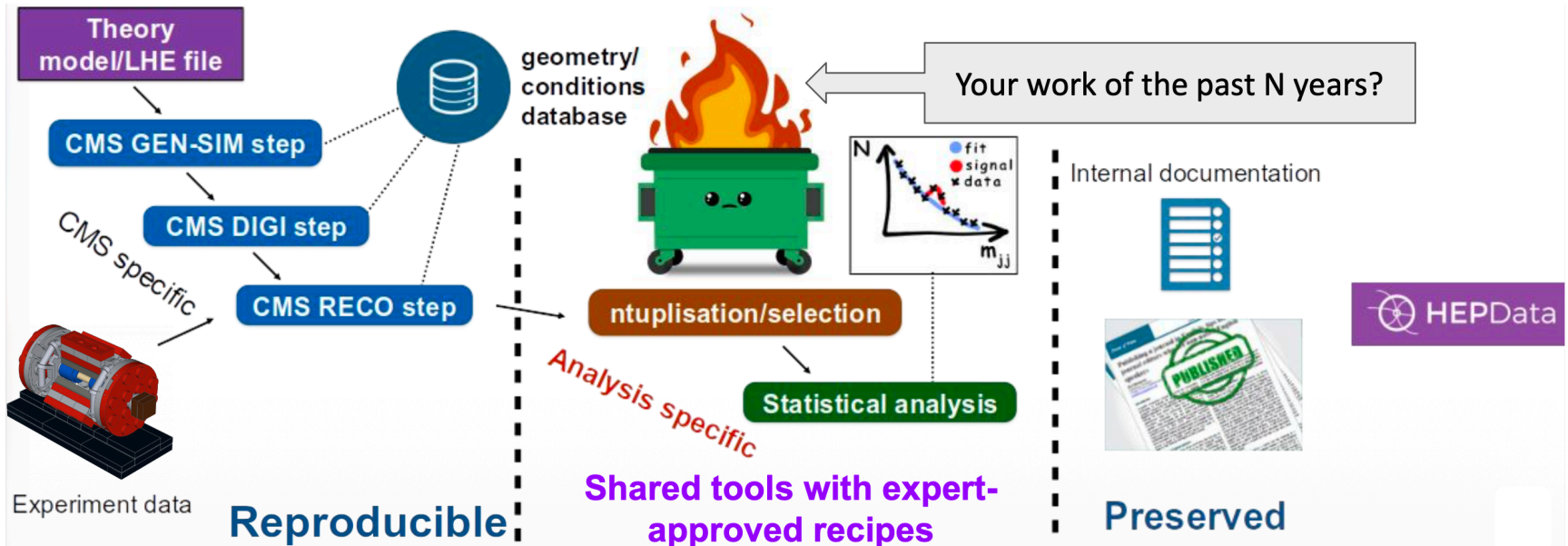
Managing all these tools as a user can be overwhelming

- Analysis facilities : CERN, FNAL, others can provide curation of software and hardware needed for the different experiments
- Web-based interface with access to jupyterlab but also terminal
- Users can complement baseline software with additional software packages as needed



Reproducibility rules

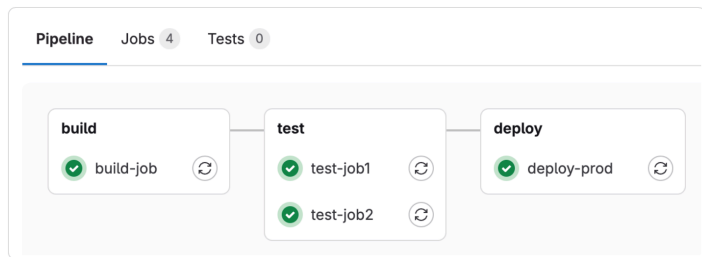
- 1) your closest collaborator is yourself 6 months ago (who doesn't reply to e-mails)
- 2) you shall re-run your analysis over and over again (until exhaustion)



Reproducibility handles

Continuous Integration/Deployment

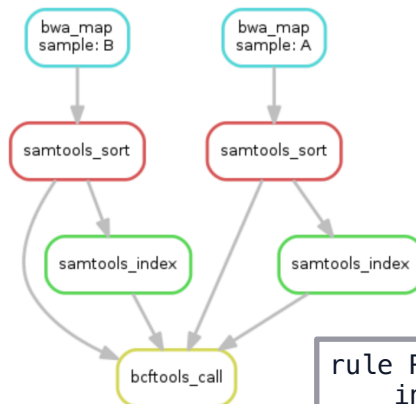
- All software scripts that you produce should be stored in git (or similar tool)
- New changes should not be merged blindly but rather tested
- Pipelines are there to help you out:
 - Describe the main steps that need to always work
 - Compare results with the previous
 - Track the changes
- A tutorial can be found in [here](#)



```
RULE_NAME:  
  extends: BASE_RULE  
  needs: [LIST_OF_RULES]  
  artifacts:  
    paths:  
      - Value  
  script:  
    - myscript arg1 arg2 ...
```

Workflow management tools

- Several on the market: [snakemake](#), [luigi](#), [law](#), [ploomber](#), etc.
- Describe each step separately
 - Factorize chained steps from those that can run in parallel
 - Pick-up flow from the last step
 - Repeat as many times as needed



```
rule RULE_NAME:  
  input:  
    Key = Value  
  output:  
    Value  
  shell:  
    myscript arg1 arg2 ...
```

How are data created at the LHC?

The high level trigger

Offline analysis

Hands-on

Hands-on

We'll make use of the CERN analysis facility also known as



We'll go through three notebooks illustrating usage of different tools

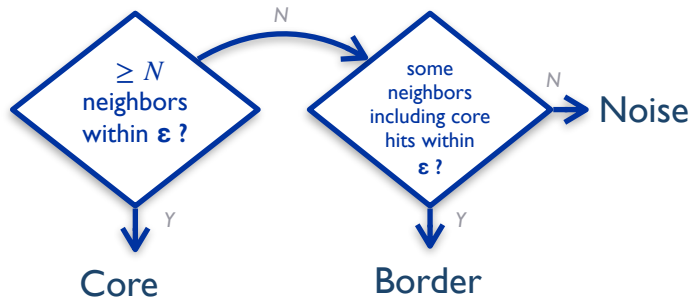
- Classic clustering examples
- A pseudo-analysis based on open data
- Running theory predictions and comparing to the result of your pseudo-analysis

First step: open <https://gitlab.cern.ch/psilva/courselhphysics> and follow the instructions therein

Example I: two clustering algorithms

DBSCAN

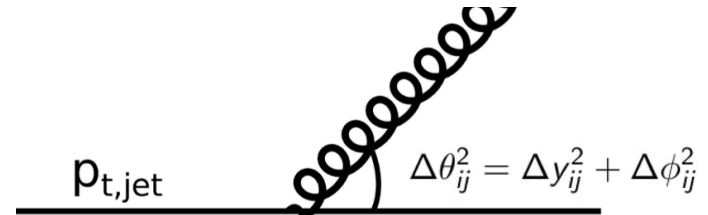
- Define reference radius (ϵ) and neighbors (N)
- Classify hits



- Connect core points into clusters if they are density connected
- Attach border points

FastJet

- Inspired by angular ordering of QCD emissions
- recursive aggregation of neighbors in the η - ϕ plane (recall its invariance properties)
- For each pair compute: $d_{ij} = \min(p_{T,i}^{2k}, p_{T,j}^{2k}) \frac{\Delta\theta_{ij}^2}{R^2}$
- Stop if $d_{ij} \geq p_{T,i}^{2k}$



Example I: two clustering algorithms

ClusteringAlgos.ipynb 9.83 KiB

Code

Preview



Crash course on LHC analysis tools - Part I : Clustering algos

In this notebook we inspect how two different clustering algorithms behave.

One (DBSCAN) is typically used in ML applications and uses a metric-based distance, while the other (FastJet) is typically used in HEP and adds an energy dimension.

You can learn more about these methods in the following references:

- <https://en.wikipedia.org/wiki/DBSCAN>
- <https://fastjet.fr>

The comparison is made using a toy MC which is used to demonstrate how tools such as numpy can be used for fast studies.

```
#you may need to install this plugin at start  
#!python -m pip install fastjet --user
```

```
from sklearn.cluster import DBSCAN  
import fastjet  
import numpy as np  
import matplotlib.pyplot as plt  
import mplhep as mh  
plt.style.use(mh.style.CMS)
```

Example II: corrections

ObjectCorrections.ipynb 5.29 KIB

Code

Preview



Exporting a correction to be used in an analysis

`correctionLib` provides a well-structured JSON data format for a wide variety of ad-hoc correction factors encountered in a typical HEP analysis.

It can be used in both python and C++. More details can be found in the [documentation page](#)

In generic terms a correction can be declared as a function of several variables.

One can use histograms, formulae, random number generators.

In this example we create a dummy corrector of the muon p_T scale.

In collaborations these are typically provided already by Physics Objects Groups.

However it's quite common that in an analysis additional corrections are needed.

Author: Pedro Ferreira da Silva (psilva@cern.ch)

```
#algebra tools
import numpy as np

#plotting tools
import matplotlib.pyplot as plt
import mplhep as mh
plt.style.use(mh.style.CMS)
```

Example III: measuring the Z transverse momentum

A2ZAnalysis.ipynb 26.54 KiB

Code

Preview



Crash course on LHC analysis tools - Part II : A-Z analysis

This is obviously not a real analysis, it's rather a collection of tools which you may find useful to include in a real analysis.

In this notebook we make use of a simulated sample of $Z \rightarrow \mu^+ \mu^-$ events and perform the following operations:

- select events of interest applying corrections and saving a summary locally
- measure the efficiency of a cut using a "tag-and-probe" approach
- correct detector-level effects by unfolding
- save the final result in a format that can be shared with theoreticians

Author: P. Silva (psilva@cern.ch)

Some useful tools

In the following we import tools which can be executed from python for simplicity. Most of them are widely used in high energy physics and it's good to get familiar with them. This is by no means an extensive list and it's mostly left as an appetizer

- [ROOT](#) - must-have in HEP: a library that spans efficient data storage and processing to math and physics tools
- [numpy](#) - offers a lot of useful tools to manipulate arrays efficiently
- [matplotlib](#) - a simple to use plotting tool
- [mplhep](#) - useful wrapper to beautify plots to publication standards
- [scinum](#) - scientific numbers with uncertainties and correlations, gaussian propagation and numpy support.
- [correctionlib](#) - provides a structured format and interface for correction factors
- [zfit](#) - one amongst several fitting tools, based on tensorflow
- [hepdatalib](#) - useful interface to export data measurements to [HepData](#)

Example III: theory predictions

Rivet8Comparisons.ipynb 8.06 KiB

Code

Preview



Crash course on LHC analysis tools - Part III - Theory predictions

Comparing our results to other models using RIVET and PYTHIA8

The Rivet toolkit (Robust Independent Validation of Experiment and Theory) is a system for validation of Monte Carlo event generators. It's written in C++ and has a large set of experimental analyses implemented as well as links to HepMC. It is used to preserve analysis for comparison and development of theory models. It is used by phenomenologists, MC generator developers, and experimentalists on the LHC and other facilities. You can find more details in [this page](#)

PYTHIA is a general-purpose Monte Carlo event generator widely used in high-energy physics collision events. It is able to describe collisions at high energies between different beams and contains several theory models for a number of physics aspects, including hard and soft interactions, parton distributions, initial- and final-state parton showers, multiparton interactions, fragmentation and decay. It can be used to interface with matrix element generators and provide the dressing of a full event. You can find more details in [this page](#)

The following notebook gives an example of running RIVET and PYTHIA. Both are available through `cvmfs` (see details [here](#)).

Nota bene: this notebook is an example: in practice we'll mostly call shell commands for compilation and running.

- We start by compiling our implementation of the analysis in RIVET : `include/Z_RIVET_ANALYSIS.cc`
- Then we run two different PYTHIA8 predictions with two different underlying event tunes. To that aim a PYTHIA8 helper is compiled in : `include/pythia_run.cc`
- The events are then analyzed with the RIVET plugin.
- The plots obtained are compared