



# Data Handling and Analysis in ATLAS

*...and some other things*

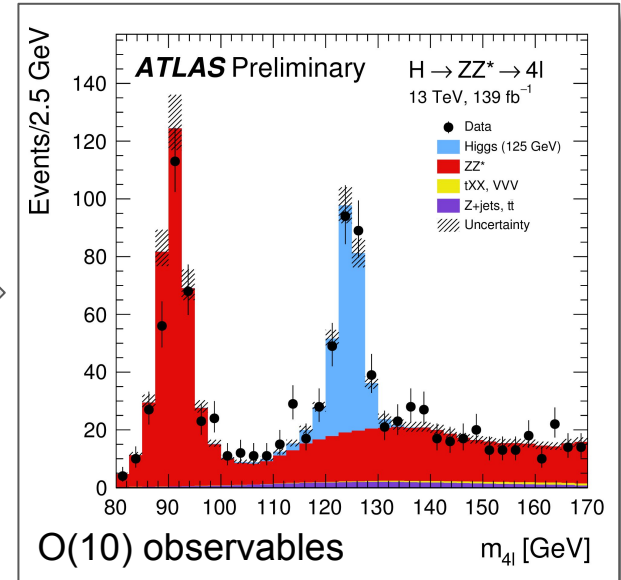
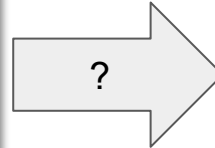
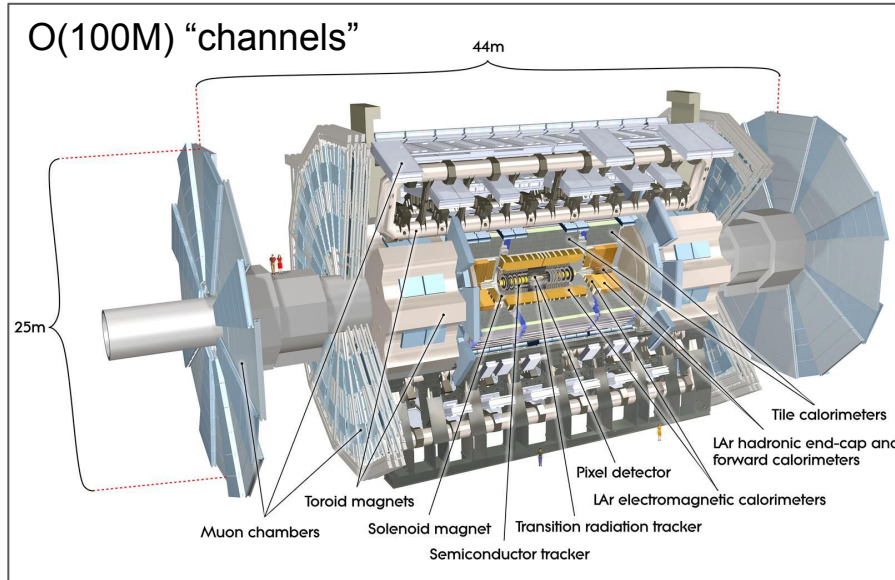
Attila Krasznahorkay



- Overview of HEP data analysis
- Status quo of HEP data analysis
- Future prospects for HEP data analysis
- Data analysis at “different” scales

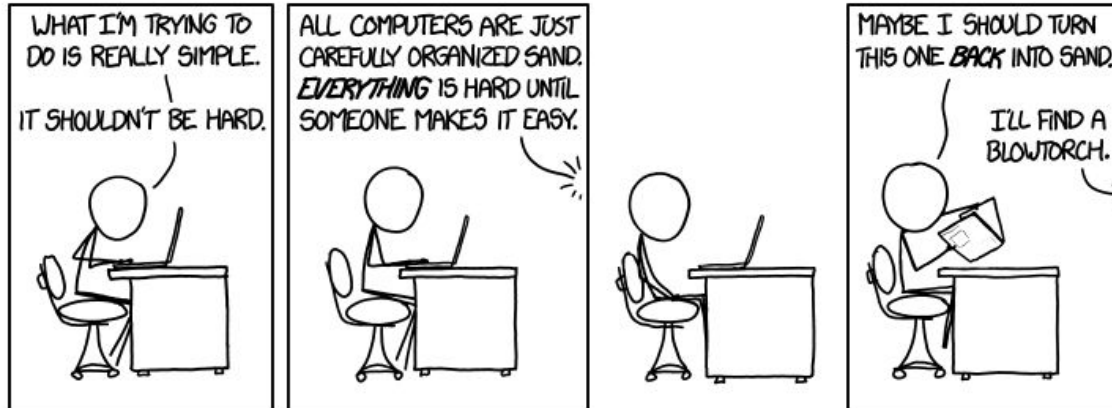
# HEP Data Analysis

- As a reminder: How we operate today was informed by decades of development. Most of the things that we do have a very good reason for being as they are...
- The goal is “simple”:



# What Does a Physicist Want?

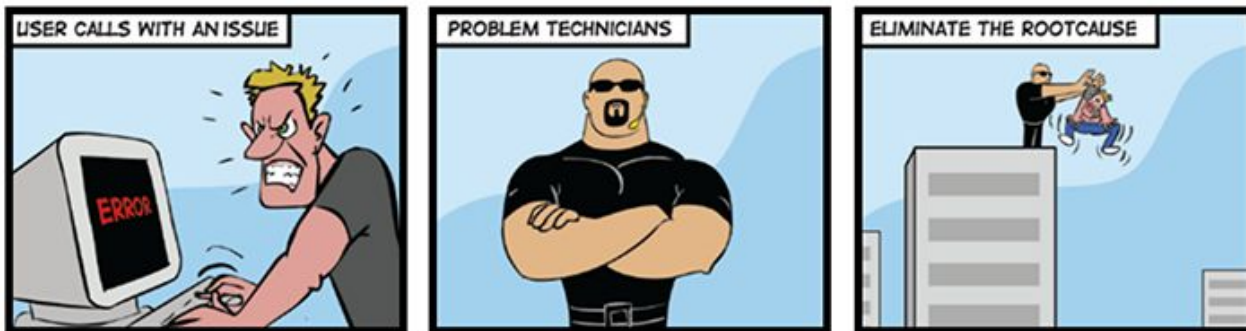
- To be able to go from ideas to results as quickly as possible
  - It should be possible to see the effects of “non-fundamental” changes to an analysis in  $O(15 \text{ min})$ 
    - If it takes hours (or even more) to get updated histograms after some minor change, the analysis will suffer
- To be able to express ideas in as easily understandable of a form as possible
  - Most physicists are not interested in computer science, it must be possible for them to reason about the collected/simulated data in a way that is as close to physics as possible



# What Do “Software/IT People” Want?

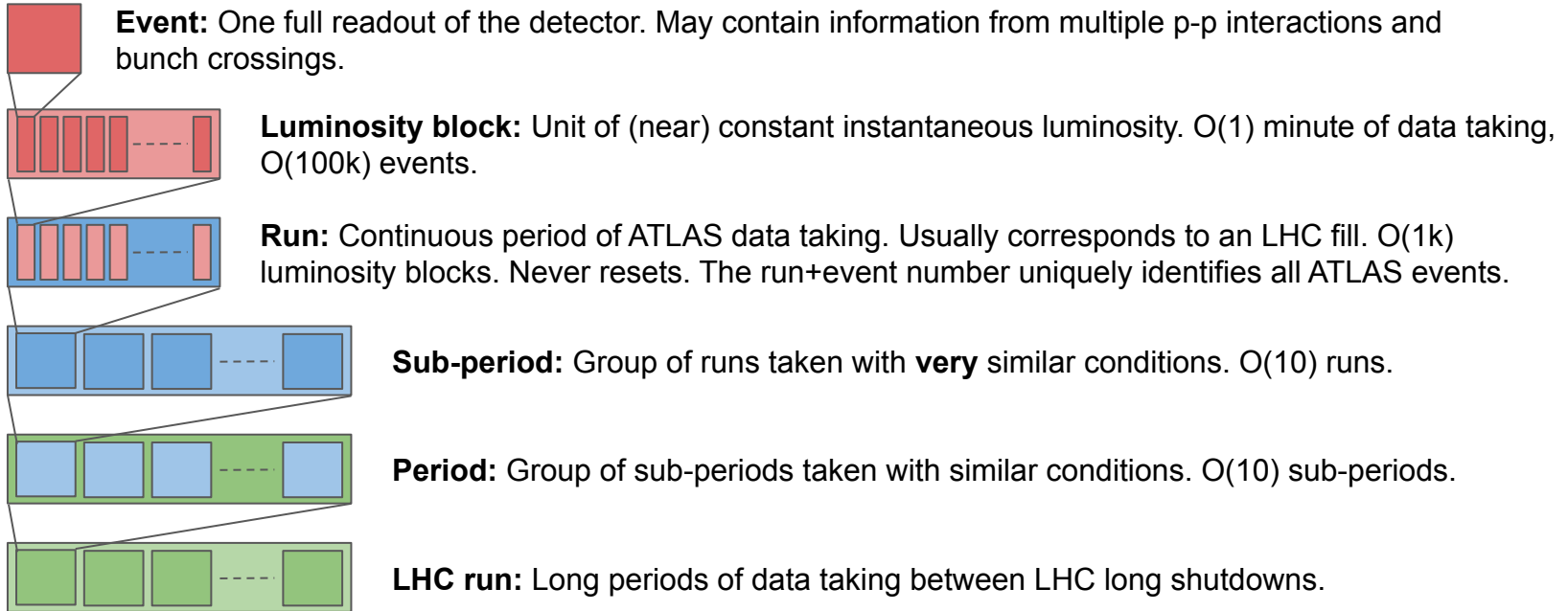


- To fit into their allocated budget... 😊
  - We have finite storage and computing resources available. The physics goals must be met while fitting into these.
- This is best served by:
  - Having a well defined data analysis model
  - Having the data processing code be as bug-free as possible
  - I.e. “reasonably” limiting what people can do...



# ATLAS's (Current) Analysis Model

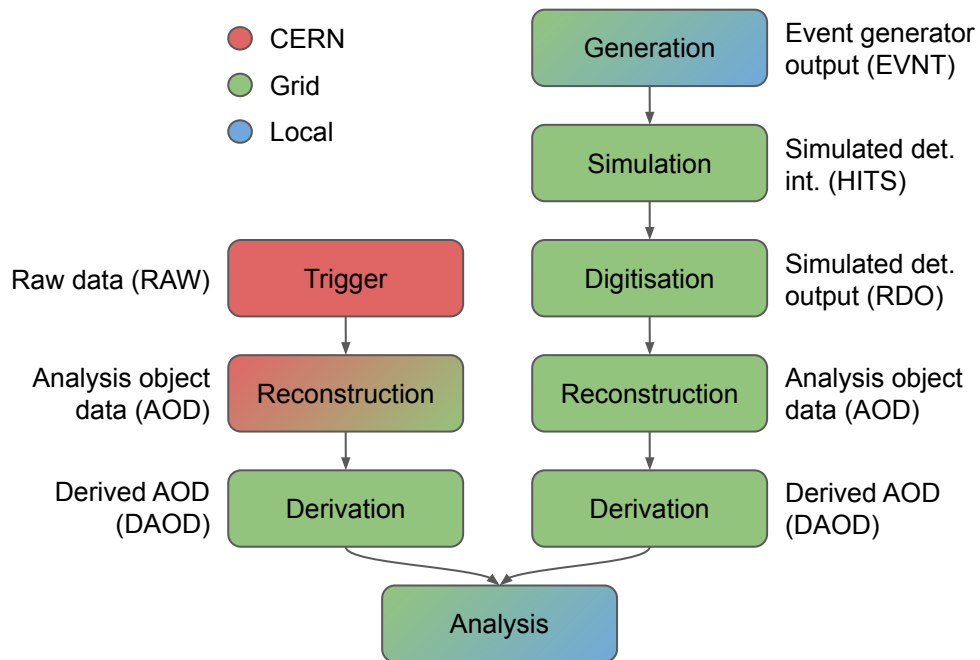
# Data Organisation in ATLAS



# The ATLAS Data Processing Model

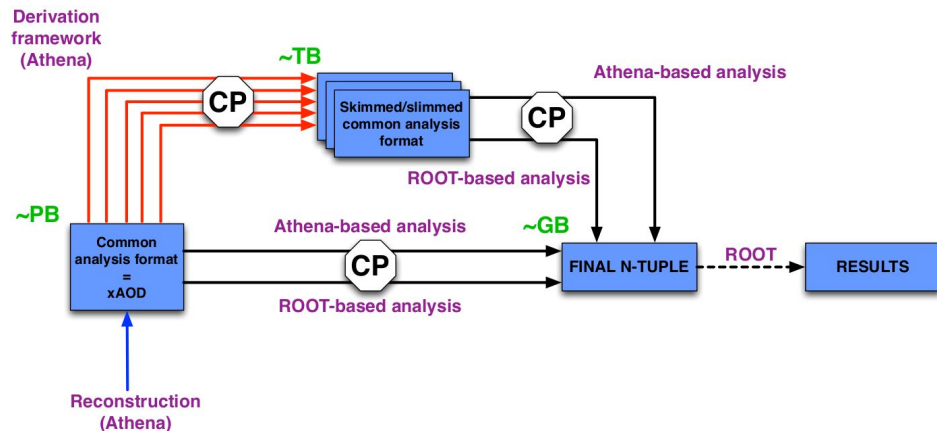


- We do what all other large HEP experiments do
  - Process the data in multiple steps, handling MC simulations in the “same way” as real data after the detector simulation
- With all steps producing specific file types, designed to efficiently hold that particular kind of information



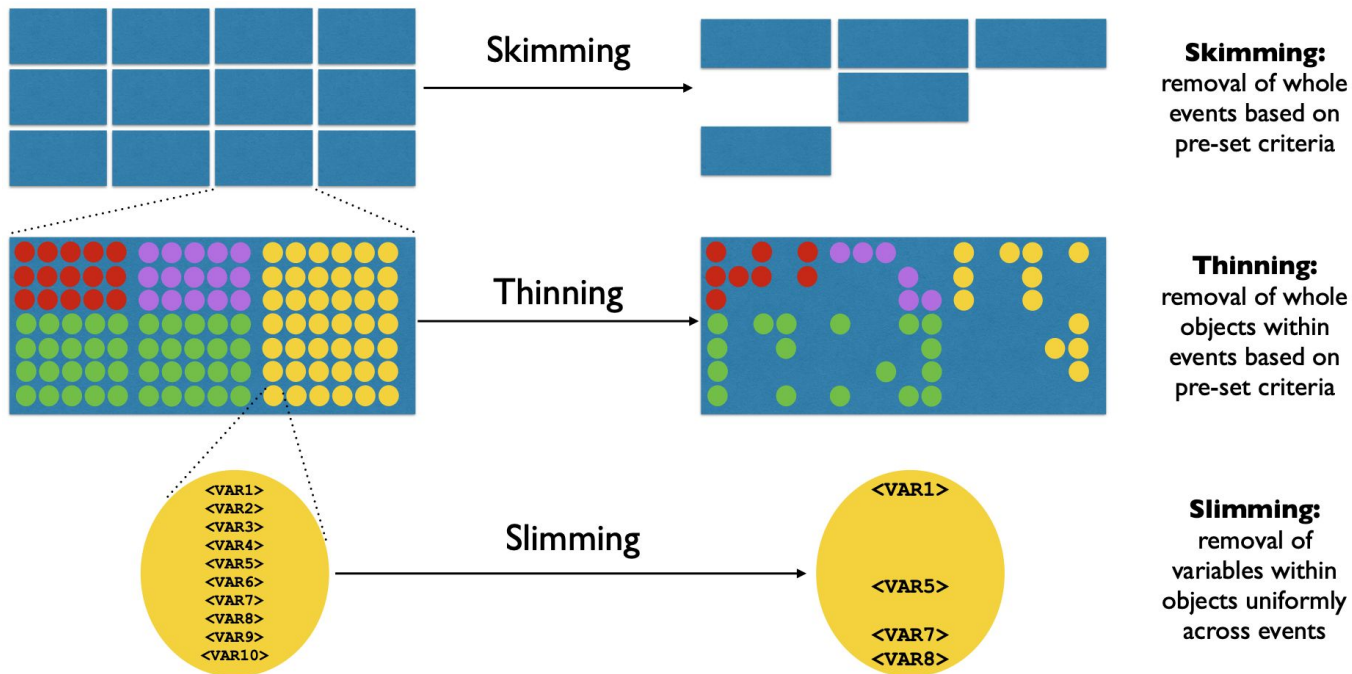


# The ATLAS Analysis Model



- Analysers are provided files/datasets with a single Event Data Model (EDM)
  - The same EDM can be used to create many different file formats with different amounts of information, object and event selections
- It is up to each analysis to apply the following during the processing of the (D)xAOD datasets
  - Final 4-momentum calibrations;
  - Efficiency corrections;
  - Systematic variations to the previous values.
- Most analyses write a custom small ntuple with the results of these operations

# ATLAS Derivations



- Performs all types of data reduction operations possible, while potentially applying fixes to the data not possible to do during analysis

# ATLAS Analysis Software Releases



TWiki > AtlasProtected Web > AtlasPhysics > AnalysisSoftwareGroup > AnalysisRelease (2019-11-26, JasonNielsen)

## Analysis Releases



Analysis releases are software releases specifically aimed at analysis use. Below is the key information about available analysis releases, as well as a glossary of useful commands.

### 21.2.101 (built on 2019-12-05)

- [28547](#) extending track smearing recommendation to 2018 data period (Analysis Recommendation)
- [28532](#) Do not run PRW for PHYSLITE (Analysis SUSYTools)
- [28523](#) fjvt variables for MVA added to EXOT5 (Derivation)
- [28488](#) 21.2 DAOD PHYSLITE MET Patches (Analysis Derivation [JetEtmis](#) Reconstruction)
- [28334](#) Add a refactored [MissingMassCalculator](#) (Analysis Tau)

## Native installers for Linux/MacOS

If you are running Linux or MacOS, you can use native installers to install one or more versions of [AnalysisBase](#):

- <http://cern.ch/akraszna/AtlasAnalysisReleaseInstaller> (Linux)
- <http://cern.ch/akraszna/AtlasAnalysisReleaseInstaller.dmg> (macOS)

These make it very simple, as they have a relatively up-to-date list of releases from which you can choose. The code

The screenshot shows the GitLab repository page for 'atlas/analysisbase'. It includes the repository name, a description 'ATLAS Standalone Analysis Release', and the last push time '4 hours ago'. On the right, there is a 'Docker commands' section with a 'Public View' button and instructions to push a new tag to the repository, with a code block showing 'docker push atlas/analysisbase:tagname'.

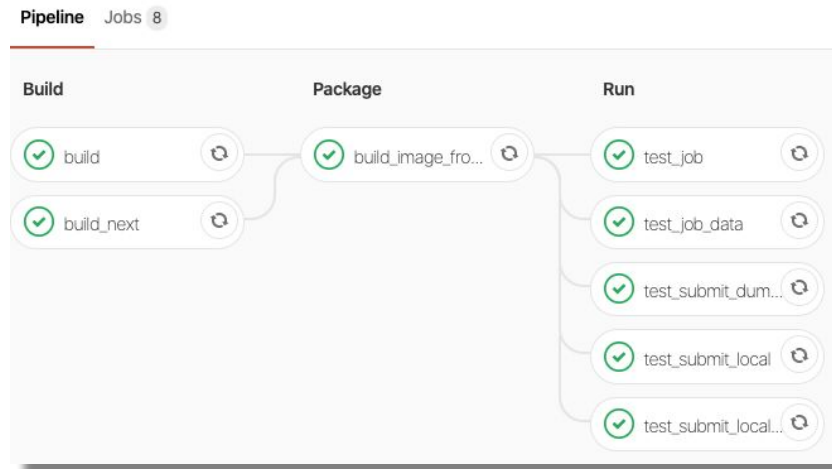
- As mentioned earlier, physicists need to perform non-trivial calibrations as part of their analysis
  - The “analysis releases” provide common code for doing so
- Built out of the same <https://gitlab.cern.ch/atlas/athena> repository as all offline/trigger code of ATLAS
  - But distributed in multiple additional ways on top of the “usual” installations on `/cvmfs/atlas.cern.ch` for SLC6/CentOS7

# Analysis Software (Repositories)

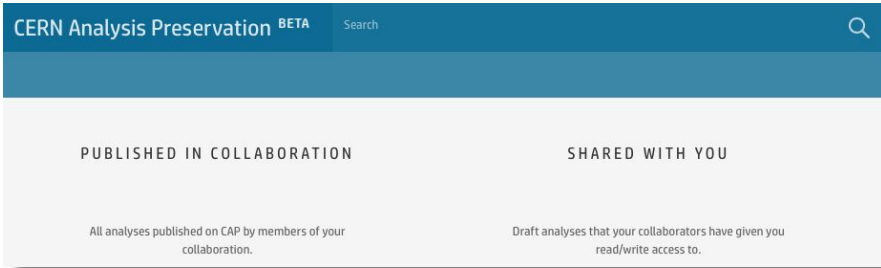


- As said earlier, the first step of most analyses produces a custom (ROOT ntuple) file format using the analysis release
  - These codes became a **lot** better since we started using Git, CMake, Docker, and a number of modern development tools
- Later steps of the analyses, operating on the custom ntuple, are usually a bit less organised
  - More on that later...

The screenshot shows the GitHub interface for the repository 'stop11-xaad' (Project ID: 57031). It includes a 'STOP code' logo, repository statistics (3,403 Commits, 9 Branches, 27 Tags, 1.6 GB Files), and a recent commit by Javier Montejo Berlingen titled 'Merge branch 'cherry-pick-DataCI' into 'master''.



# Analysis Preservation in ATLAS



reana

Home Examples Get Started Documentation News Contact

# reana

Reproducible research data analysis platform

Flexible

Run many computational workflow engines.



Scalable

Support for remote compute clouds.



Reusable

Containerise once, reuse elsewhere. Cloud-native.



Free

Free Software. MIT licence. Made with ❤️ at CERN.



- Analysis preservation is a big topic since a while
  - We need to do a better job with archiving how all our analyses were performed, so future generations would not have to re-learn everything from scratch
- The policy in ATLAS at the moment is to
  - Save a “runnable” version of the code operating directly on DxAOD files
    - And of all the code coming after
  - Save O(1 TB) of ntuples per analysis, on which the final plotting was done
- The goal is to make [REANA/RECAST](#) a success
  - Note that this is an ATLAS+CMS project by now...

# Machine Learning (in Analysis)



- Is a quite active field of course
- Used mainly for classification tasks currently
  - To separate different kinds of events/objects/etc.
- But is being used for regression tasks more and more
  - For instance to estimate energy correction values from a number of input variables
- Training is always happening outside of experimental frameworks
  - Even with TMVA this was always the case, and with Python based tools it didn't change
- Experimental-software-wise the challenge is implementing inference
  - Using the ML libraries directly is often not the best approach

## Keras: The Python Deep Learning library



### You have **Lightweight Trained Neural Network**

Keras is a high  
TensorFlow, C  
Being able to g

build passing coverage passed DOI 10.5281/zenodo.597221

### What is this?

The code comes in two parts:

1. A set of scripts to convert trained neural networks to a standard, JSON format



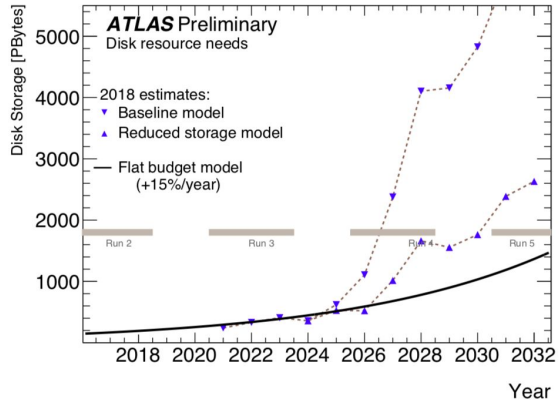
environment

Windows CPU succeeded Windows GPU succeeded Linux CPU succeeded Linux GPU succeeded MacOS CPU succeeded

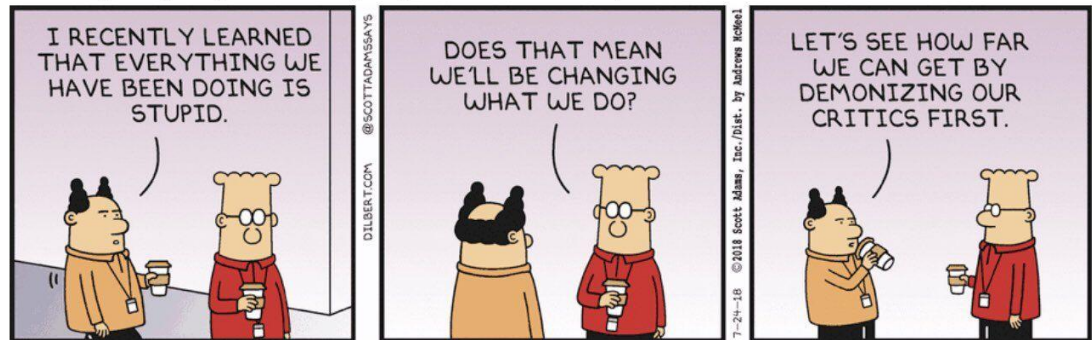
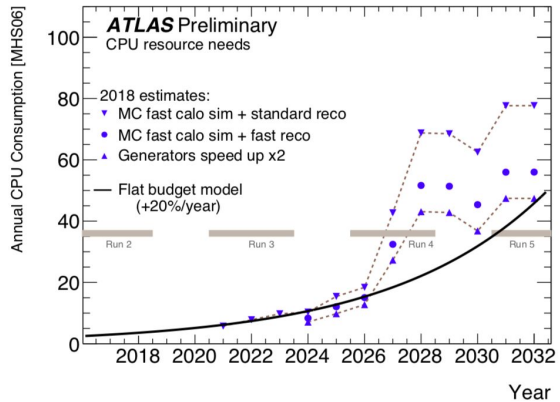
**ONNX Runtime** is a performance-focused complete scoring engine for Open Neural Network Exchange (ONNX) models, with an open extensible architecture to continually address the latest developments in AI and Deep Learning. ONNX Runtime stays up to date with the ONNX standard and supports all operators from the ONNX v1.2+ spec with both forwards and backwards compatibility. Please refer to [this page](#) for ONNX opset compatibility details.

# LHC Run-3 and Beyond

# Coping With More Data / Complexity



- Our current analysis model is serving us well at the moment
- But it will not scale to the HL-LHC
  - In Run-3 we could still survive without larger changes, but that would only delay the issue
- So we are updating our model once again...





- **Has a very similar goal as CMS's mini-AOD**
  - Provide a single data format that could cover the vast majority of analyses
    - With a <50 kB/event size
  - Will (hopefully) allow us to reduce the number of DAOD types from the current  $O(100)$  to  $O(40)$ 
    - A large number of technical formats will still need to remain
- **Fundamentally the same kind of file as all of our current DAODs**
  - But will only contain the variables absolutely necessary to perform all user level calibration tasks (with some very few additions)
    - This is where our EDM shines! We can keep using the same EDM while dropping variables. The EDM can also tell us which variables are used by a set of analysis tools.
- **Will not require a major change in analyses**
  - They will be expected to use DAOD\_PHYS datasets in the same way as the current physics DAODs
- **Will be our main workhorse for Run-3**

- Has a similar goal to CMS's nano-AOD
  - O(10 kB)/event
  - But still using the xAOD EDM!
    - Holding “nominally” calibrated objects instead of the “uncalibrated” ones shipped with DAOD\_PHYS
- Users will still need to use tools provided by the analysis releases to get the systematic variations on the pre-calibrated objects
  - So the analysis releases will still play a big role
  - But these tools are much lighter weight than the ones applying nominal calibrations, it will be possible to perform analyses in more ways than it is now
- We will provide DAOD\_PHYSLITE datasets already during Run-3
  - We have to figure out how to use them efficiently, as they will have to become our main way of physics analysis for Run-4

# Declarative Analysis



- You probably heard about Directional Acyclic Graphs (DAGs) recently...
  - This is one of the new buzzwords...
  - There is a growing trend in HEP to try to express all the analyses that we do in DAGs
- One of the easiest to use examples of this is [ROOT::RDataFrame](#)
  - Will probably make more and more use of such tools as we go forward
- But do exercise a bit of skepticism!
  - People are now trying to re-invent the exact same “analysis language” that we used with [PAW](#)...

The screenshot shows a Jupyter Notebook titled "CMS Meeting Notebook Example" with the following content:

CMS Meeting Notebook Example  
This is a simple example of using an ATLAS (D)xAOD file using ROOT::RDataFrame.  
First off, import ROOT to get going.

```
In [1]: import ROOT
Welcome to Jupyter 6.16/00
```

Now set up

```
In [2]: ROOT.xAOD
xAOD::In
Using the h
```

Since xAOD objects are unfortunately not directly "plottable" (at least in ROOT 6.16/00), we need to extract the variable that we want to plot, into a simple vector.

```
In [4]: elPT = df.Define( "ElectronsPt",
    """
    std::vector< float > result;
    result.reserve( Electrons.size() );
    for( const xAOD::IParticle & p : Electrons ) {
        result.push_back( p->pt() );
    }
    return result;
    """)
```

Now make a histogram, and draw it into a canvas. It has to be done in a few separate steps unfortunately...

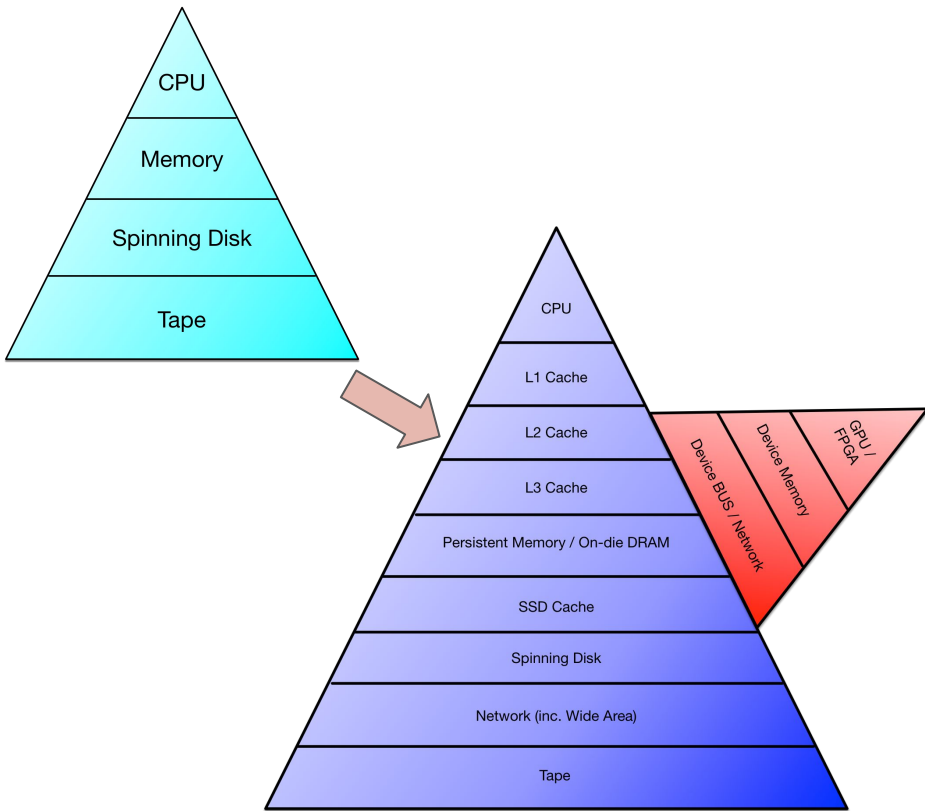
```
In [5]: hist = elPT.Histo1D( "ElectronsPt" )
In [6]: drawCanvas = ROOT.TCanvas()
        hist.Draw()
In [7]: %jsroot on
        drawCanvas.Draw()
```

The notebook displays a histogram titled "ElectronsPt" showing the distribution of electron transverse momentum. The x-axis is labeled "ElectronsPt" and ranges from 0 to 700, with a multiplier of  $\times 10^3$ . The y-axis ranges from 0 to 10000. A statistics box in the top right corner of the plot area shows the following values:

ElectronsPt	
Entries	15903
Mean	10736
Std Dev	20380

- ALICE is currently evaluating the usage of dedicated “analysis centres” for Run-3
  - All physics analyses would be done in a small number of dedicated sites
  - One idea is to use these with a UI similar to SWAN
    - Though I can’t imagine how they would survive without providing a classical batch system on these centres as well 😞
- This is certainly gaining some traction in HSF
  - “Closely controlled” analysis centres would make it easier to make use of heterogeneous hardware for physics analysis
- ALICE claims an amazing  $O(\text{PB})/\text{hour}$  analysis speed in their R&D on specialised setups
  - Will have to see how they can build an analysis model around this concept that will be compatible with how physicists develop code, and with analysis preservation in general

# Heterogeneous Computing



- Computing is quite a bit more complicated these days than it used to be... 😞
- All the LHC experiments are trying to figure out how to best make use of new types of hardware (GPUs, FPGAs, etc.) for data processing
- I do believe that how well we manage to figure this out, will make a big impact on how we analyse data on the Run-4 timescale

# Analysis at a Smaller Scale

# The “X17 Experiment”

- Don't want to go into any real detail about the experiment itself, just want to highlight a few computational details
- The size of this is many-many orders of magnitude smaller than those of the LHC experiments
  - Even much smaller than the planned [FASER Experiment](#)
  - $O(50)$  parameters read out for one “event”
- Its data processing uses many tools developed for much larger experiments though

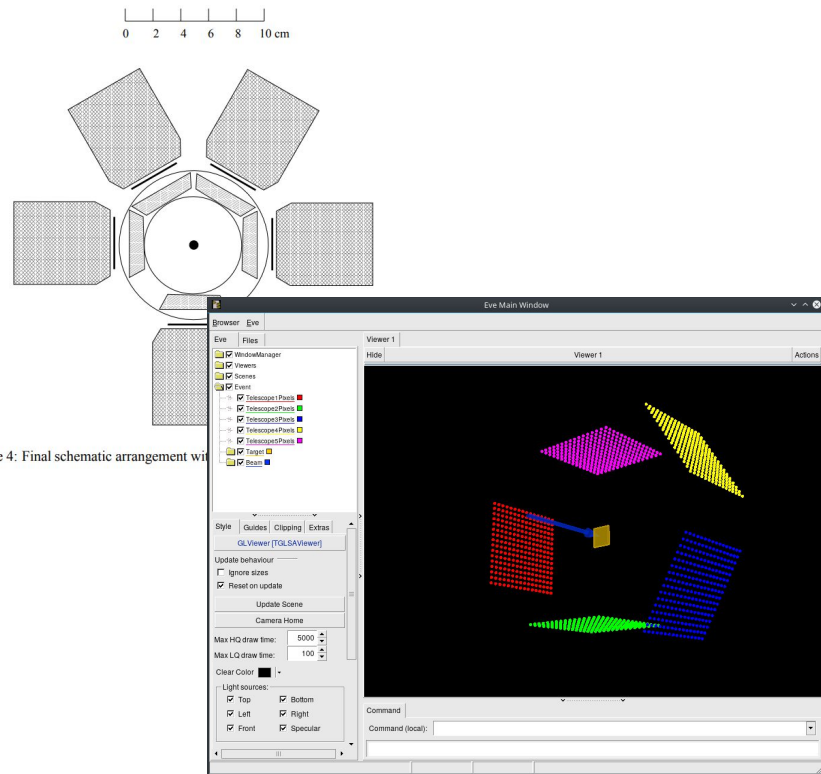
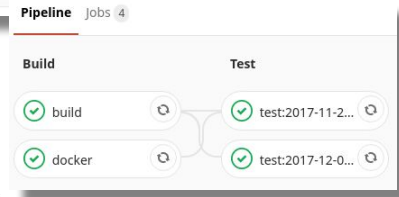
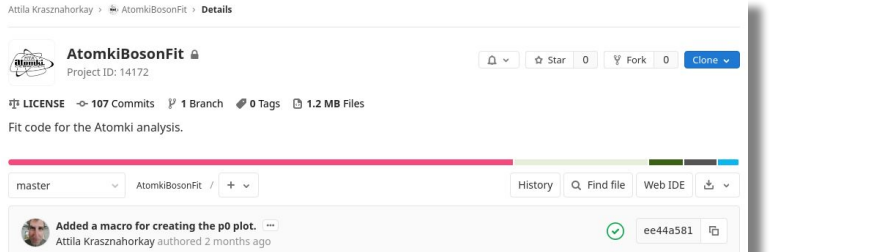
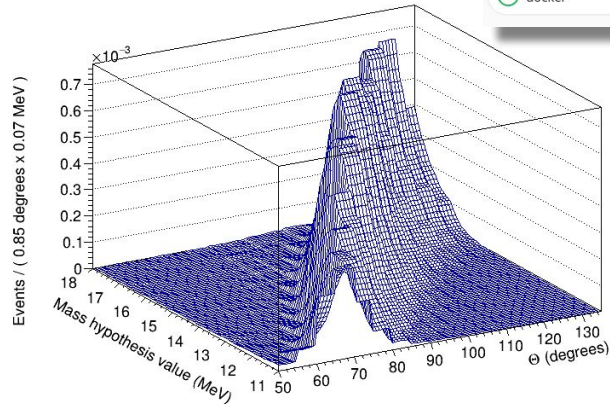


Figure 4: Final schematic arrangement with

# Statistical Analysis (for X17)



RoobosonPdf shape



- Apart from writing the DAQ software for the experiment, I wrote the code performing the statistical interpretation for it
  - Using many of the same coding skills that I learned / developed for ATLAS
- The C++/Python code used for the fitting / significance estimation lives in [gitlab.cern.ch](https://gitlab.cern.ch)
  - The repository's CI tests both the build of the code, and that it can successfully run on older distributions already analysed by it (histogram inputs are kept on EOS)



- **The way we analyse data at the LHC is once again changing**
  - What we did in Run-1 and 2 **were** the right thing to do. But now that we understand our detectors better, we need to do something else.
- **The computing landscape around us is changing**
  - Many smaller / simpler computing cores are becoming the norm. Our code has to adapt, as it currently performs very poorly on these systems.
- **Machine learning can definitely help us along this road, but is not a silver bullet**
  - For now just continue using it for the “classical” classification and regression applications, until somebody comes up with something better.



<http://home.cern>