# Data Analysis Working Group

## Analysis today

Paul Laycock

# Introduction to the Data Analysis Working Group

- **Aim - make publication of physics results more efficient, eliminate monotonous and laborious tasks from physics analysis**

- 1st priority - capture the requirements of analysis by direct consultation:
    - https://indico.cern.ch/event/782504/

- Second 1st priority - survey work of technology pioneers:
    - https://indico.cern.ch/event/789007/

- *18 excellent talks* which the three **DAWG** convenors will try to summarise
    - **What have we learned?  What can we improve? How bad is it?**

- Many thanks to the speakers, most of the material in this talk originated there
    - Credit goes to the original presentations (not always credited here, sorry!)
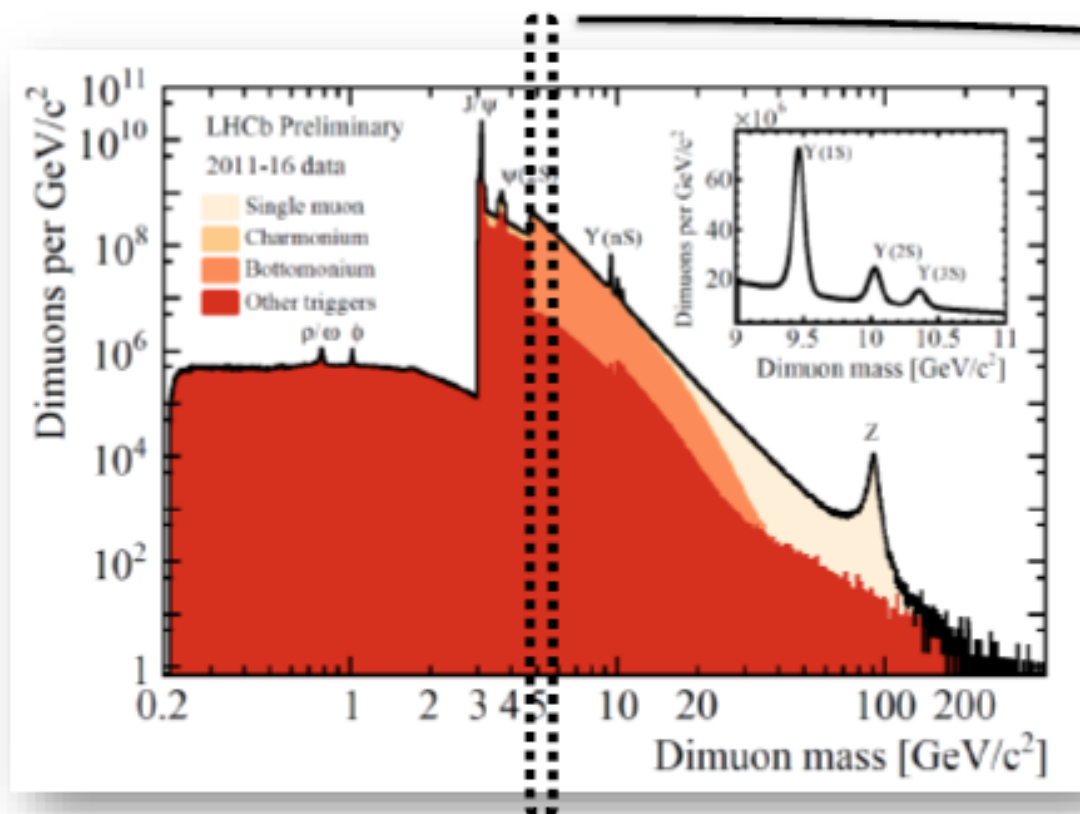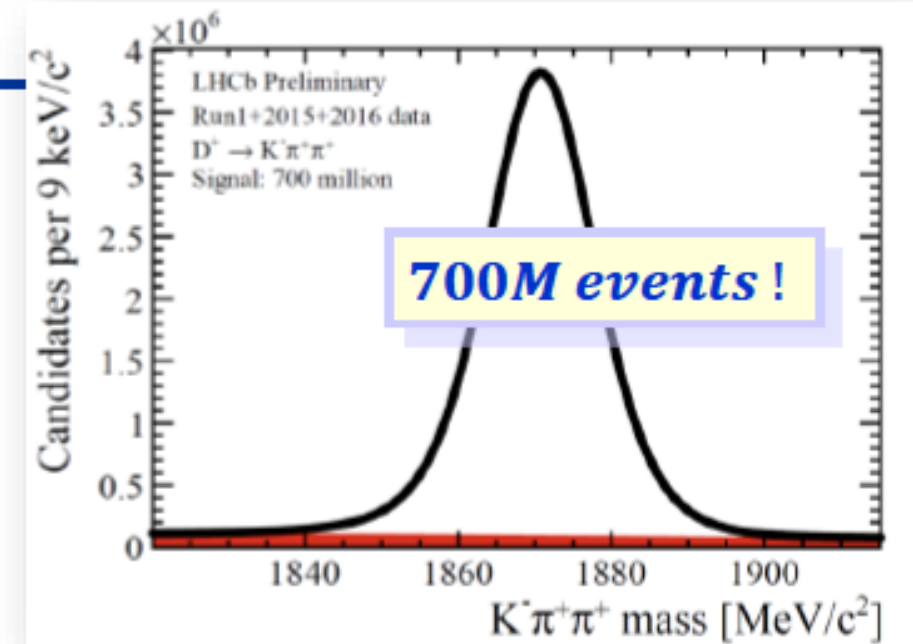
U.S. DEPARTMENT OF **ENERGY**
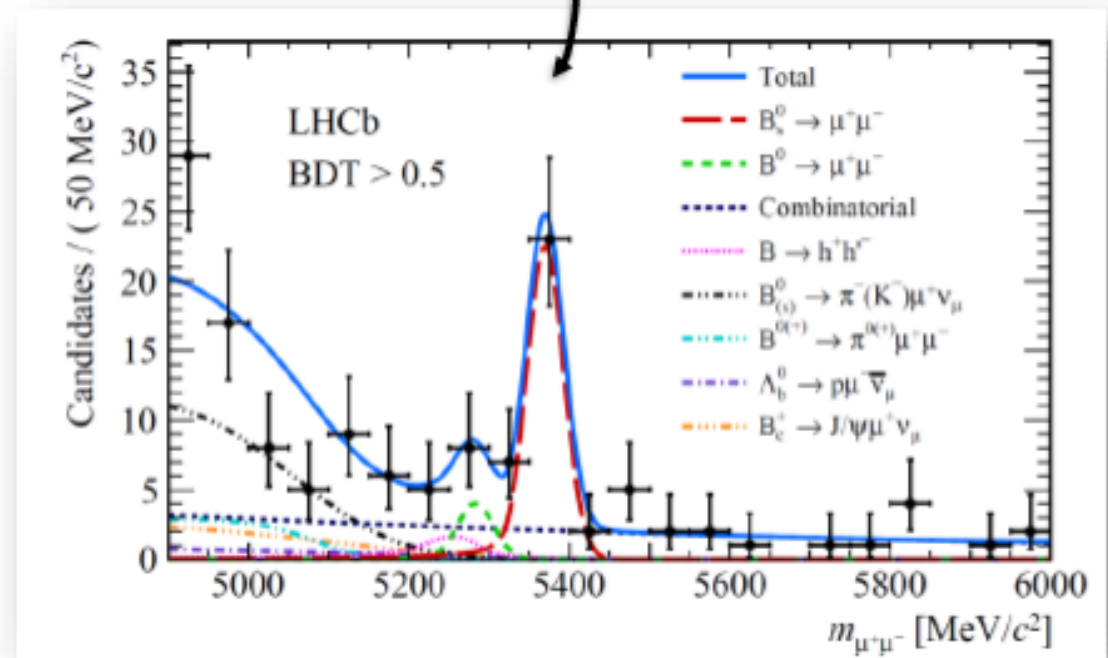
**BROOKHAVEN**
NATIONAL LABORATORY

# What is analysis?



CSC 2016, Mol Belgium@CERN

# A question of many scales



Wildly different Challenges !

700M events !
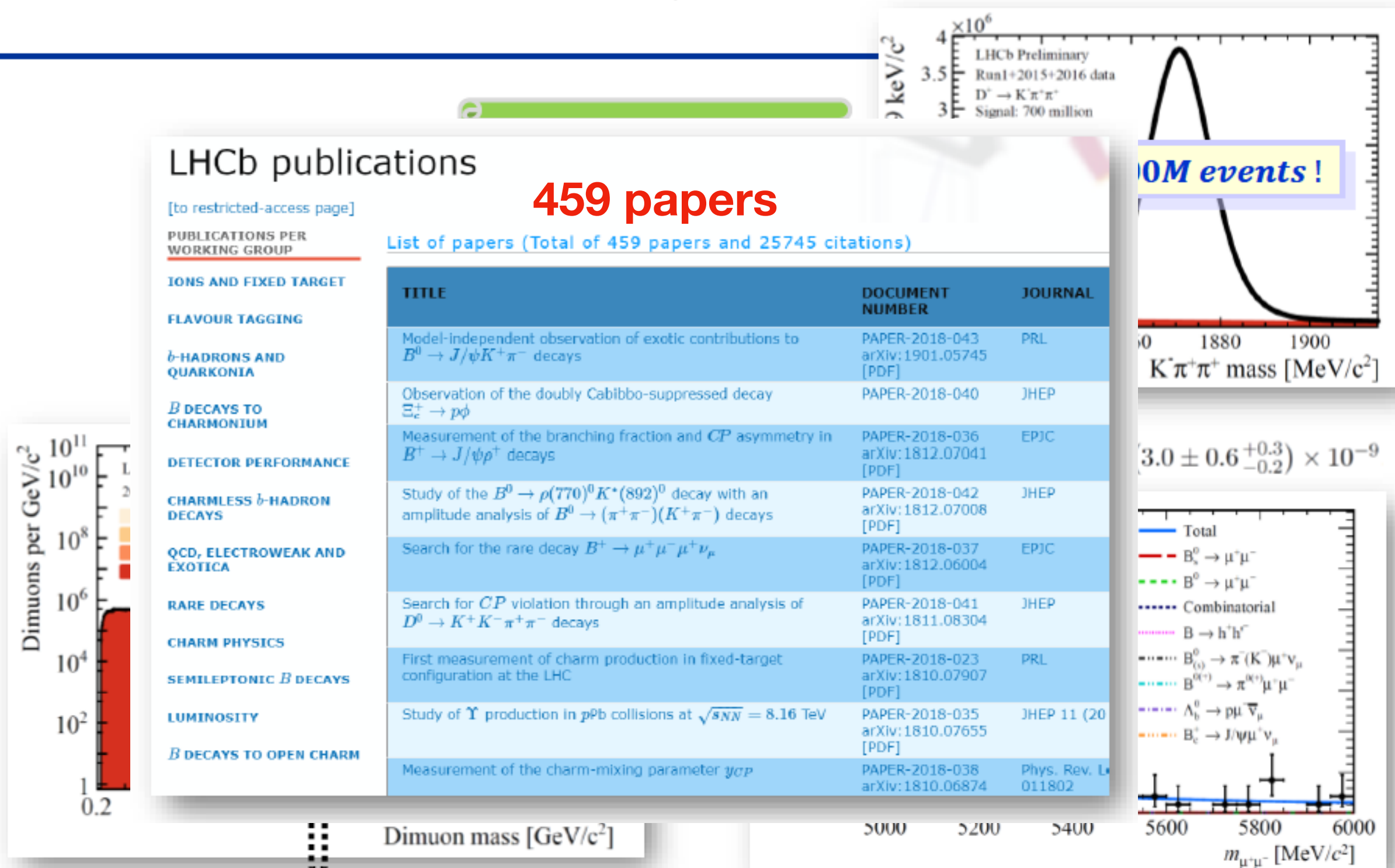
$$\mathcal{B}(B_s^0 \to \mu^+\mu^-) = (3.0 \pm 0.6^{+0.3}_{-0.2}) \times 10^{-9}$$

Note : structure given the numerous "trigger lines" with different requirements
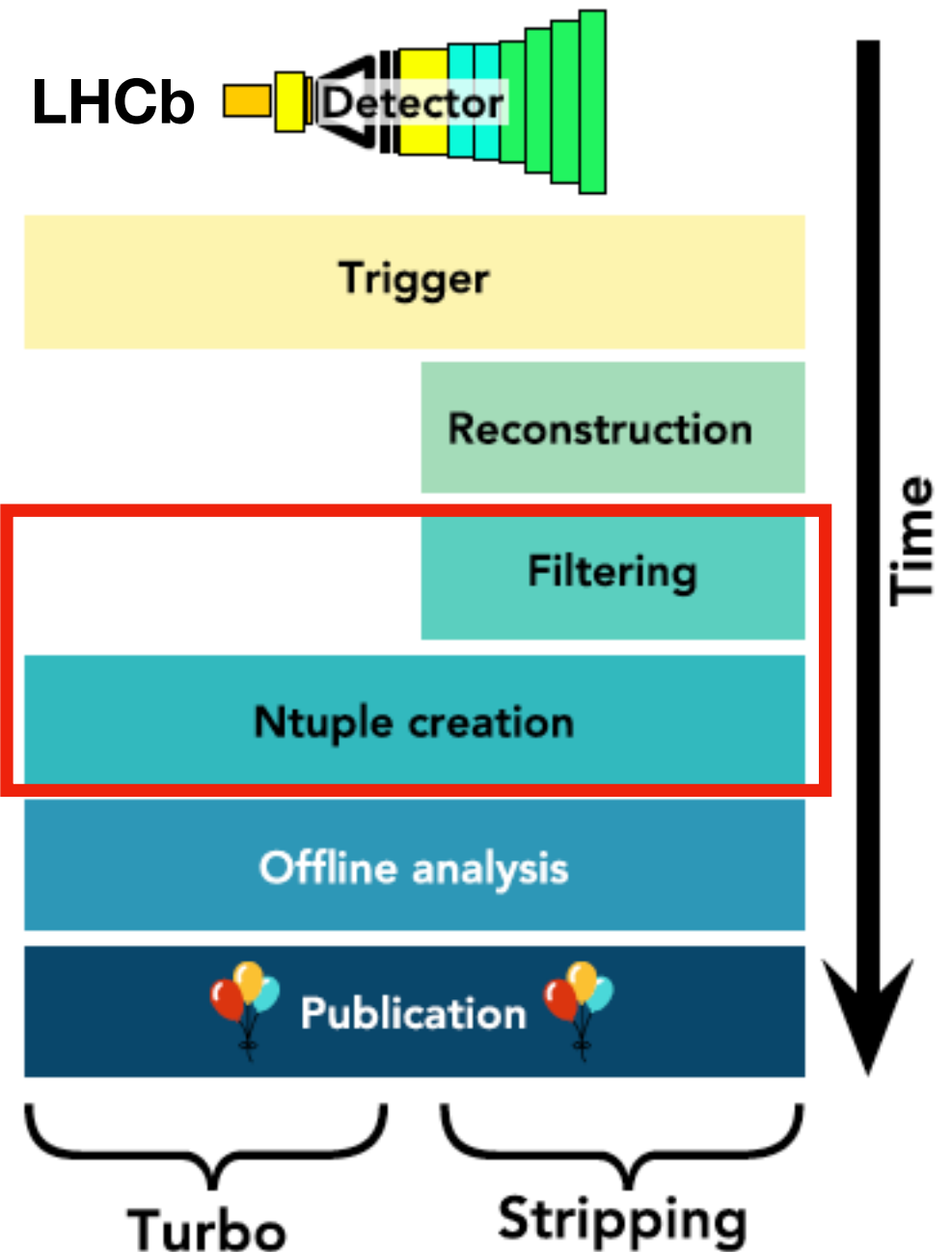
4

# A question of many scales



459 papers

Eduardo Rodrigues                    HSF Data Analysis WG Meeting, CERN, 23rd Jan. 2019

# Analysis workflows



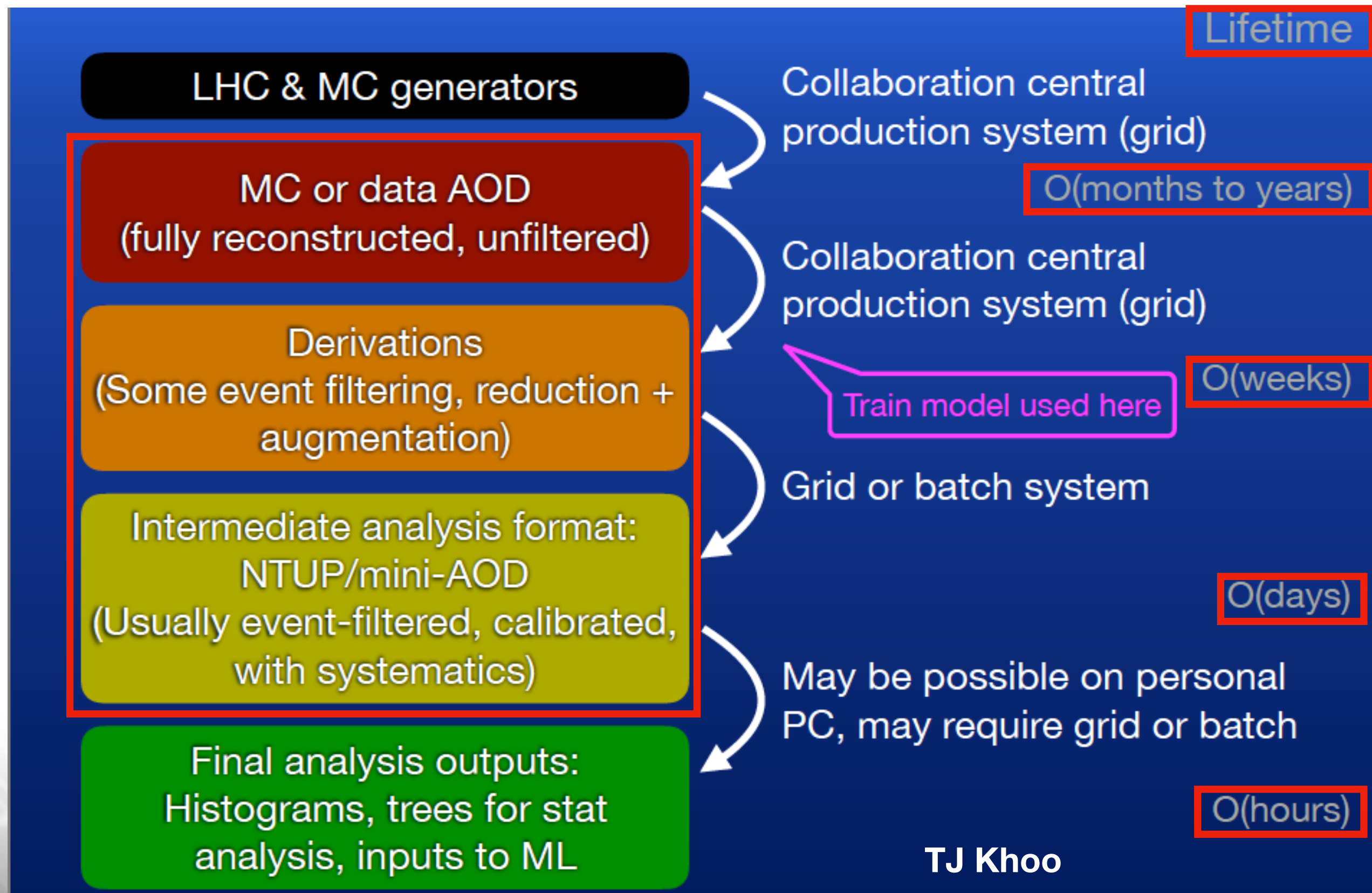- Tend to focus on **heavy lifting** here rather than the final stages
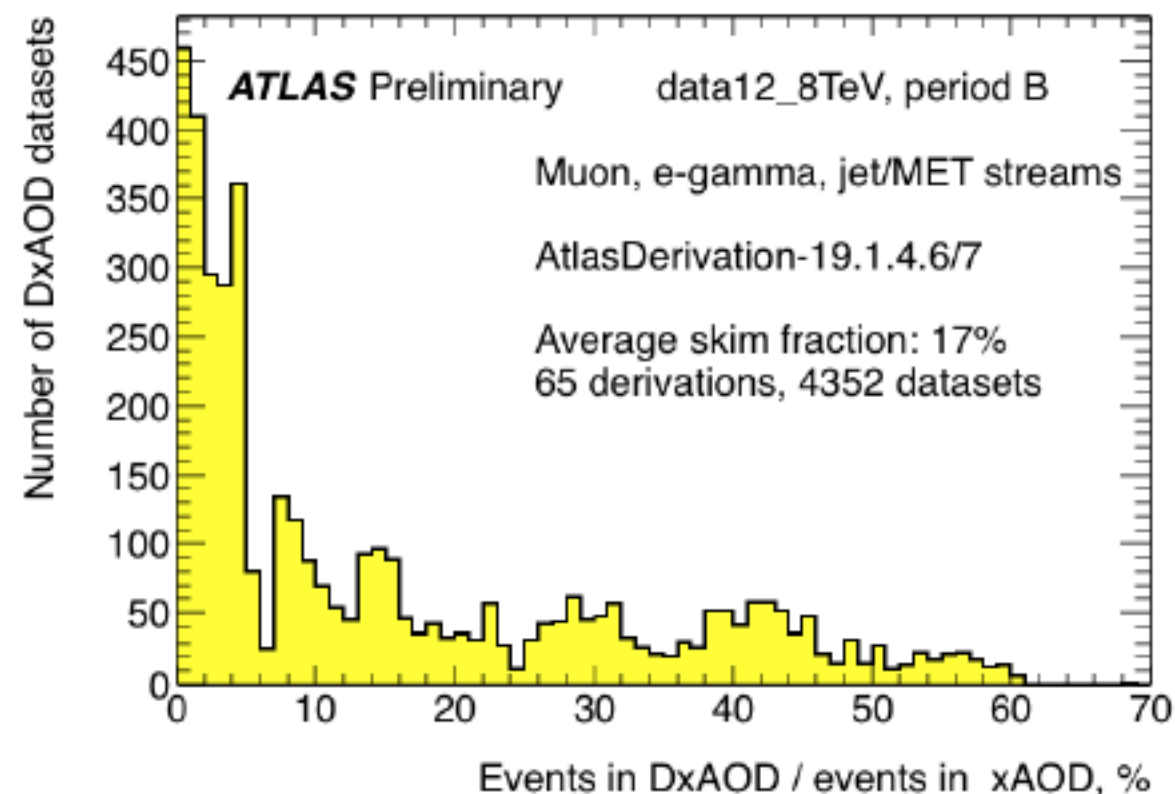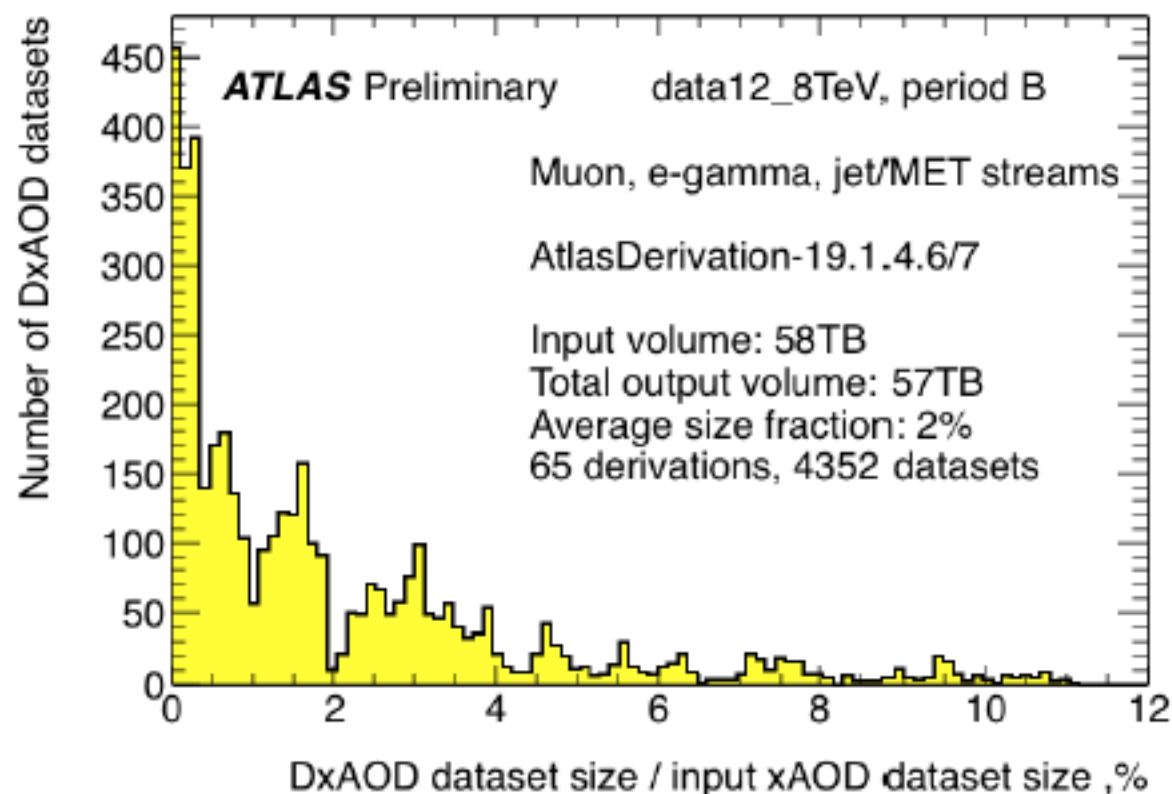  - More on statistical inference et al in **Andrea's** talk

# Repeated heavy lifting

# Data reduction trains - Fact and fiction



ATLAS Preliminary    data12_8TeV, period B

Muon, e-gamma, jet/MET streams

AtlasDerivation-19.1.4.6/7

Input volume: 58TB
Total output volume: 57TB
Average size fraction: 2%
65 derivations, 4352 datasets

DxAOD dataset size / input xAOD dataset size ,%



ATLAS Preliminary    data12_8TeV, period B

Muon, e-gamma, jet/MET streams

AtlasDerivation-19.1.4.6/7

Average skim fraction: 17%
65 derivations, 4352 datasets

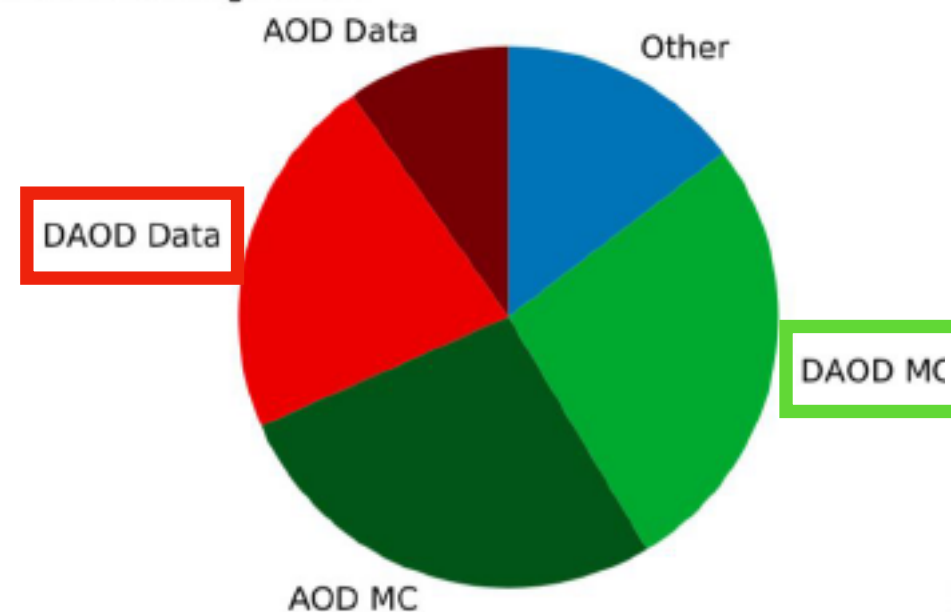Events in DxAOD / events in xAOD, %

**CMS nano-AOD ~1kB/event**

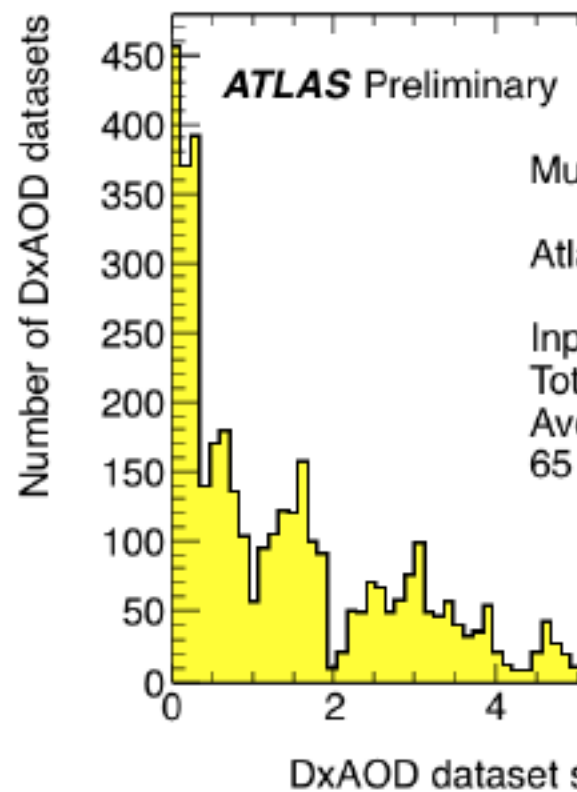- **expected to cover > 50% of analyses**

**ATLAS derivations**

- **2014** (*top*): small efficient data format

- **2028** (*right*): more than half the storage



ATLAS Preliminary. 2028 Disk resource needs
Reduced storage model

AOD Data    Other
DAOD Data
DAOD MC
AOD MC

# Data reduction trains - Fact and fiction



**ATLAS deriva**

- **2014** (*top*):
  - small efficie
- **2028** (*right*):
  - more than h

# ALICE analysis trains



Data taking + calibration → Reconstruction → Event Summary Data (ESD) → AOD "filtering" → Analysis Object Data (AOD)

Common tasks — Analysis 1 — Analysis 2 — Analysis 3 — Analysis 4...

**ALICE** successfully driving analysis trains over **AOD** input

# Going nano



Data taking + calibration → Reconstruction → Event Summary Data (ESD) → AOD "filtering" → Analysis Object Data (AOD) → Nano ROOT TTrees



Common tasks — Analysis 1 — Analysis 2 — Analysis 3 — Analysis 4...

**ALICE** plans to go nano like **CMS**, while **ATLAS** aiming for *10-50 kB/event*

U.S. DEPARTMENT OF **ENERGY**

**BROOKHAVEN**
NATIONAL LABORATORY

# Analysis workflows - best practice?



Event Rate (events/s)

High Level Trigger

Turbo | FULL | TurCal

Tape Storage

100% | 80%

Turbo | FULL | TurCal

Disk Storage

- *Turbo stream analysis (data scouting): reco and calibration done once in HLT no more reprocessing!*
  - in reality for LHCb, two HLT passes
    - TDAQ to record events to a big buffer
      - then prompt calibration
    - then second pass for data reduction
- *Q. How much physics bandwidth can go this route for other experiments?*

- *Centrally produced nano-format: no more reinventing the wheel for producing a data analysis format*
  - in reality, cannot accommodate all analyses, BUT important to use where it is possible (maybe even for ATLAS)
- *Q. How much physics bandwidth can go this route for all experiments?*

# Where: Power vs control

- Grid
  - Portability
  - Dataset sharing
  - Dataset access
  - Reliability

- Local cluster
  - Dataset access
  - Reliability
  - Portability
  - Dataset sharing

How easy for collaborator B to use collaborator A's submission scripts?

How easy for collaborator B to use collaborator A's job outputs?

How easy for collaborator A to use own job outputs?

How long before job outputs 100% available?

TJ Khoo

# Wherever: Power and control?

- Grid
  - Portability
  - Dataset sharing
  - Dataset access
  - Reliability

- Local cluster
  - Dataset access
  - Reliability
  - Portability
  - Dataset sharing

**Hiding the "how" is a common theme**

**see declarative analysis in Andrea's talk**

Do you use notebooks, whether standard Jupyter notebooks, or within JupyterLab?

146 responses



- Never
- Sporadically
- Every now and then
- Very often

26.7%
8.9%
7.5%
56.8%

**LHCb Data Analysis Survey, 2018**

# Analysis platforms



**Portal**

**Jupyter Notebooks**

**Web APIs**

- Data access via IVOA-standard protocols
- Same interfaces that support other aspects

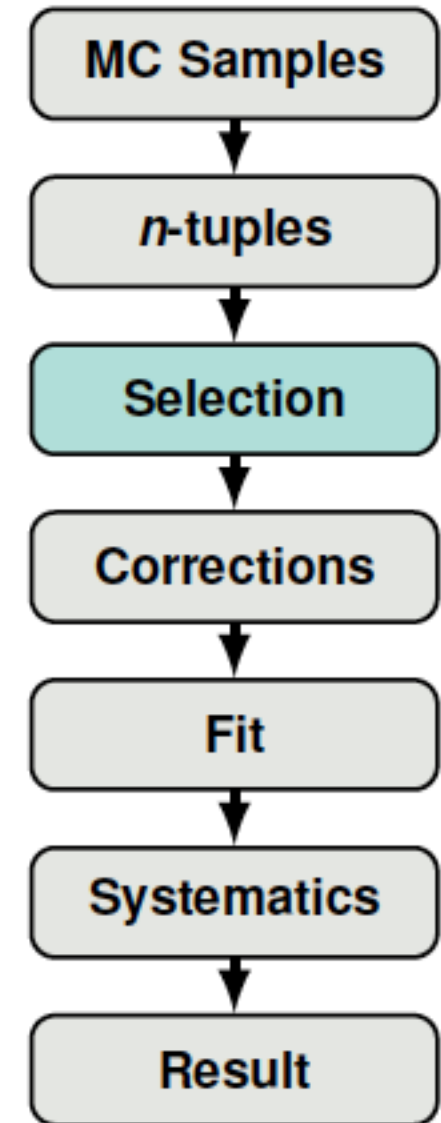- See also Lukas's <u>talk</u> from Tuesday

# Belle II analysis software stack

**Analysis of the *n*-tuples is done with Python:**

- *Pandas* and *numpy*
- *root_pandas* or *uproot* to load ROOT files
- *scikit-learn* or basf2 MVA package for MVA methods
- *matplotlib* for plots
- convert *n*-tuples to hdf5 files (these are loaded $\sim$10 times faster)
- data analysis in *jupyter notebooks*

**Why Python?**

- Well documented!
- Easy to integrate into the rest of the analysis
- Modern and nice interface…

MC Samples
↓
*n*-tuples
↓
Selection
↓
Corrections
↓
Fit
↓
Systematics
↓
Result

# Belle II - best practice analysis code

## UI

### A simple example

```python
import basf2
from modularAnalysis import inputMdst, reconstructDecay, fitVertex, variablesToNtuple
from stdCharged import stdPi
from stdPhotons import stdPhotons

mypath = basf2.Path()

# configure modules
inputMdst("default", basf2.find_file('analysis/tests/mdst.root'), path=mypath)
stdPi("good", path=mypath)
stdPhotons("good", path=mypath)
reconstructDecay('rho0:myrhos -> pi+:good pi-:good', '0.5 < M < 1.0', path=mypath)
fitVertex('rho0:myrhos', path=mypath)
reconstructDecay('B0:myBs -> rho0:myrhos gamma:good', '5.0 < M < 6.0', path=mypath)

# output modules
momenta = ['px', 'py', 'pz']
variablesToNtuple('B0:myBs', momenta, path=mypath)

basf2.process(mypath)
```
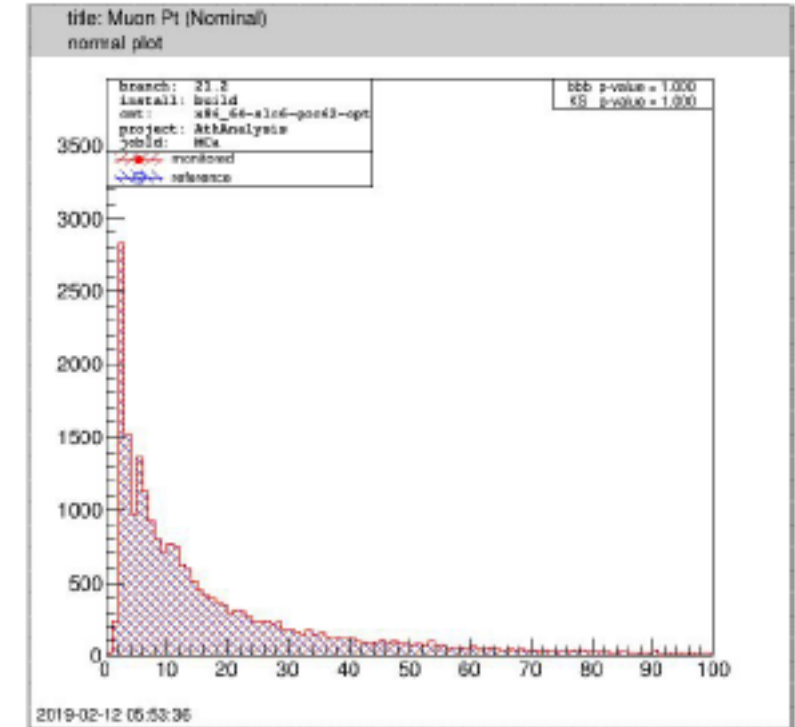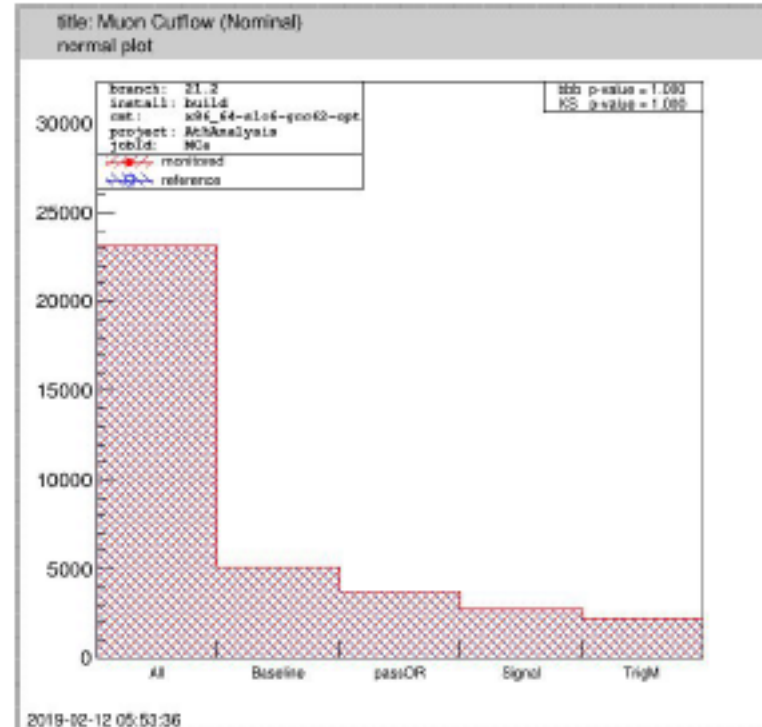
U.S. DEPARTMENT OF ENERGY

BROOKHAVEN
NATIONAL LABORATORY

# Beyond analysis functionality

**SUSYTools @ Hass AbouZeid**

**Configuration diff**
**High-level analysis comparison**

| | Default | Analysis 1 | Analysis 2 |
|---|---|---|---|
| EleBaseline.Pt | 10000. | 10000. | 7000. |
| EleBaseline.Eta | 2.47 | 2.47 | 2.47 |
| EleBaseline.Id | LooseAndBLayerLLH | LooseAndBLayerLLH | VeryLooseLLH |
| EleBaseline.CrackVeto | false | true | false |
| EleBaseline.z0 | 0.5 | 0.5 | |
| Ele.Et | 25000. | 10000. | 20000. |
| Ele.Eta | 2.47 | 2.47 | 2.47 |
| Ele.CrackVeto | false | true | false |
| Ele.Iso | Gradient | FCTight | Gradient |

**CI pipeline**
**Carefully control code**



Reference histogram updated automatically every night

# GUI Overkill

# Observations / Questions

- Analysis is diverse, but we see recurring themes and solutions

- **Reducing I/O for heavy lifting:**
- *Trains* an accepted solution, can more workflows use this concept?
- *Common nano-AOD* centrally produced, less reinventing the wheel on format
  - **Q**. How much bandwidth can go this route?
- *Turbo stream* calibrate once
  - **Q.** How much bandwidth can go this route?  How strong is the physics case to justify not doing that?
- Convergence on *Jupyter* notebooks as analysis platform, *hiding the how is good*

- Trend towards *declarative analysis*, especially for **LHCb/Belle II**
  - **Does anything prevent other experiments?**

- Addressing *systematics* is still a challenge, see **Andrea's** talk
  - Can we attack (some of) this as a community?  What is best practice?
  - Not covered - *Monte Carlo*

# But first, the next talk



**Winter is comin'**