



Science and  
Technology  
Facilities Council

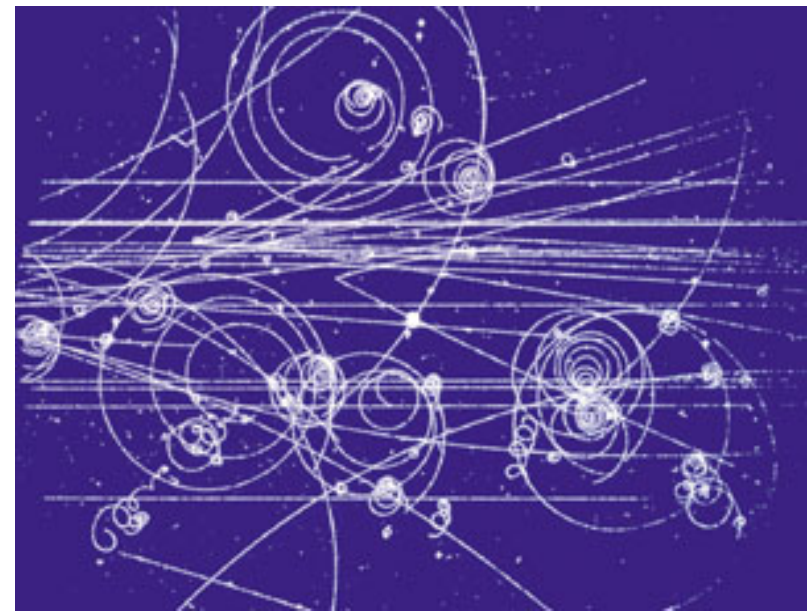
# The High Energy Physics Software Stack

Katy Ellis, PPD RAL

Efficient Computing for HEP, 17/02/2020

# Contents

- Personal introduction
- HEP terms and experiments
- Upcoming challenges
- HEP software stack
- Distributed computing
- Efficiencies



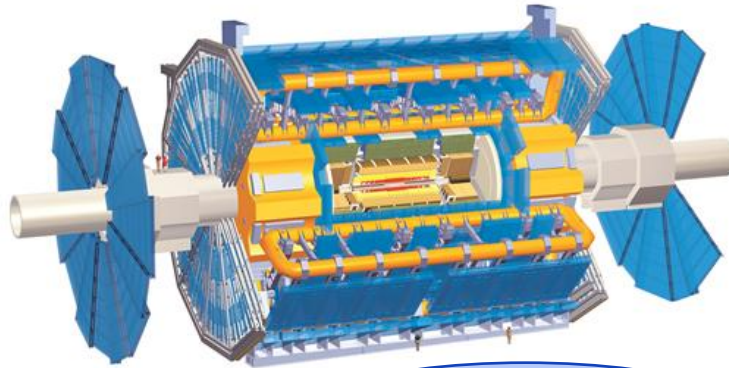
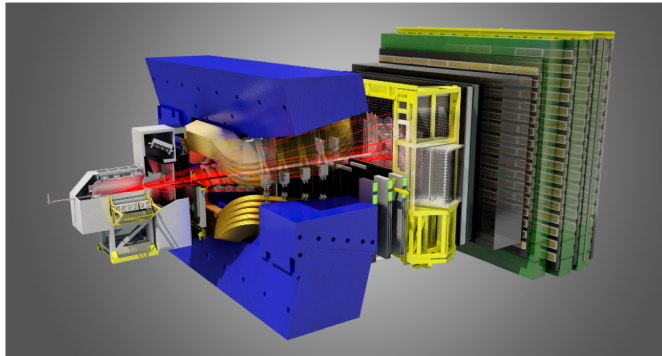


# Personal introduction

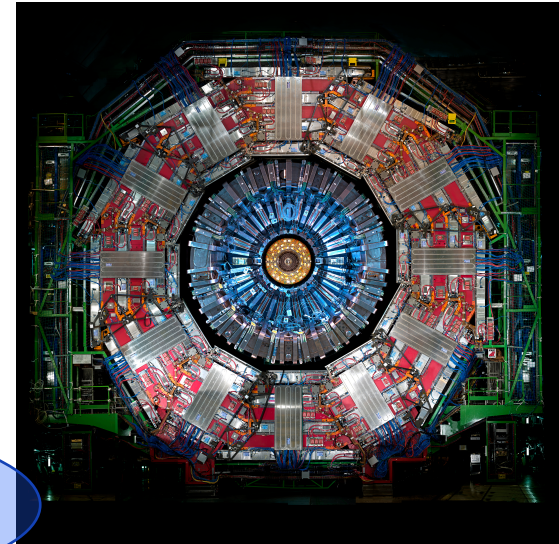
- CMS / Tier 1 Liaison at the Rutherford Appleton Laboratory (RAL)
  - Ensures that computing for the CMS experiment is performing well, particularly at RAL
  - Sits on the interface between CMS physicists and Tier 1 service administrators
- PhD in Experimental Particle Physics
- Various (computer) modelling jobs in industry and at UK labs
  - Stealth Scientist at Qinetiq
  - Tester of oil and gas reservoir simulation software at Schlumberger
  - Nuclear fusion power station modeler at Culham

# What is HEP?

- High Energy Physics, aka Particle Physics, is the study of subatomic matter particles and force-carrying particles.
- Several different types of experiment, many on a grand scale:
  - Accelerators (such as the Large Hadron Collider)
  - Neutrino

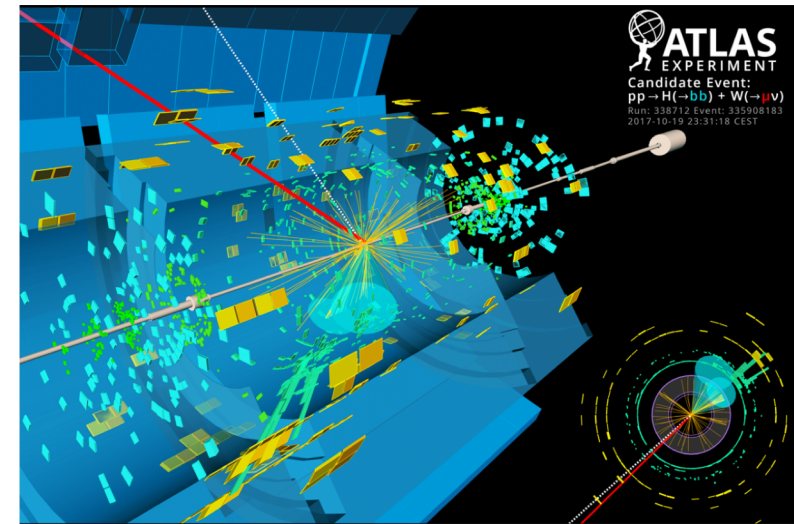


~100 million channels

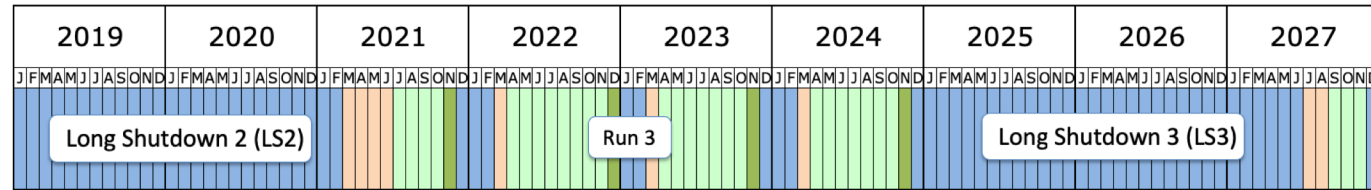


# Events

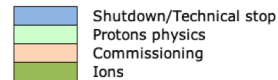
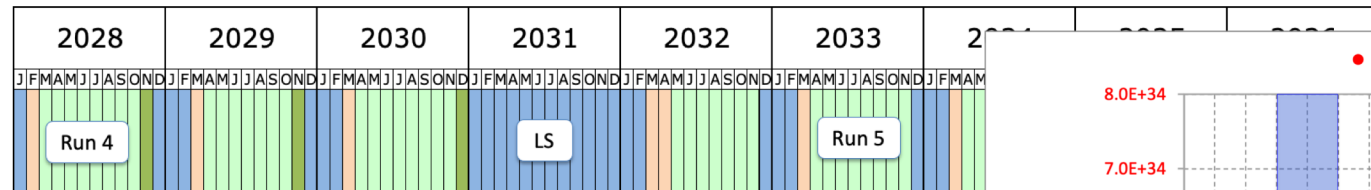
- ‘Event’ is used in HEP to represent a snapshot of the detector at a particular moment
  - Ideally this captures something interesting!
- The same word ‘event’ also describes a theoretical interaction created by Monte Carlo-based computer modelling
- As events are processed by different parts of the software stack they are stored in different formats.
  - RAW 1-2 MB/event, all the hits from the detector
  - AOD 100 kB/event, physics objects for analysis
  - Derived smaller formats



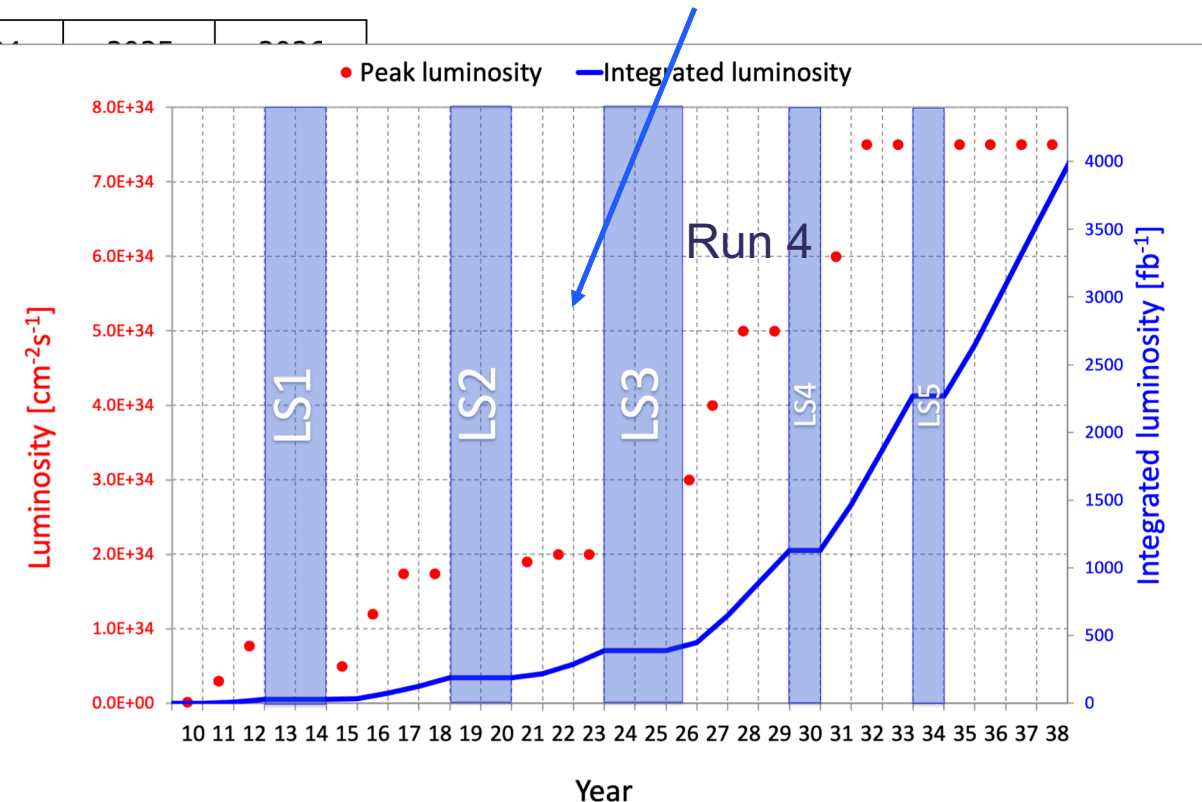
# Increase in LHC data rate



Plot has not been updated with delayed Run 3 schedule

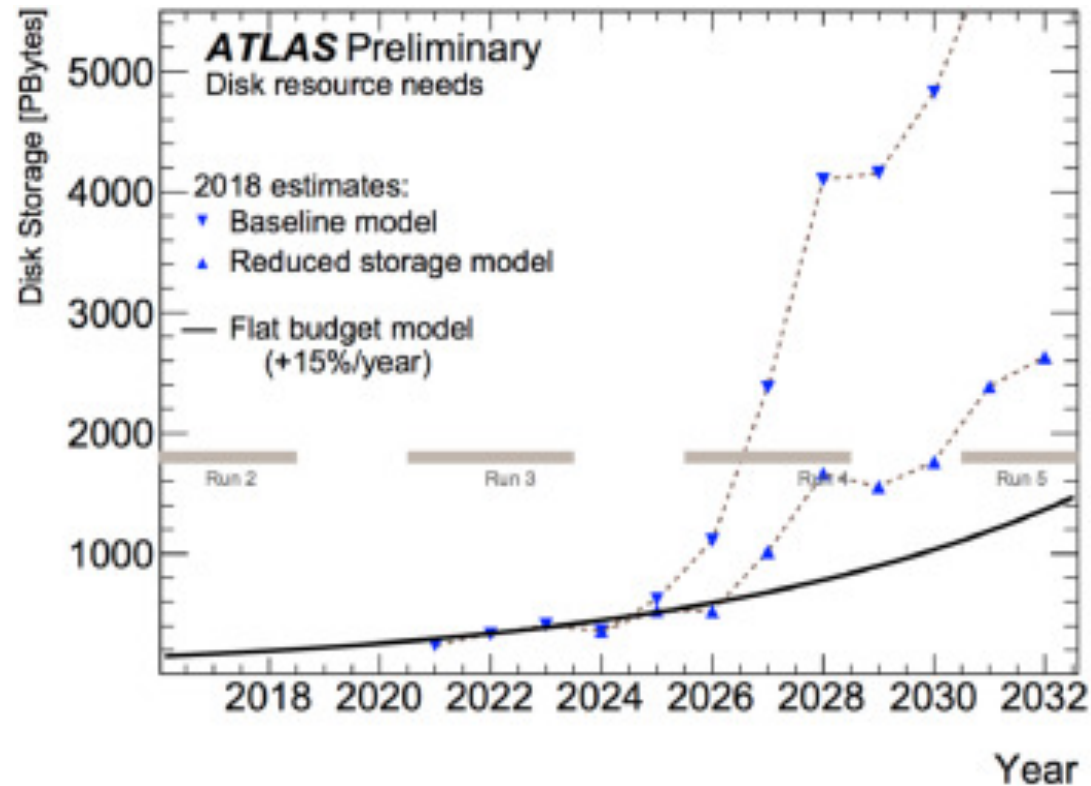


- Increased pile-up (number of collisions per bunch crossing)
- Increased detector resolution
- Increased recorded event rate

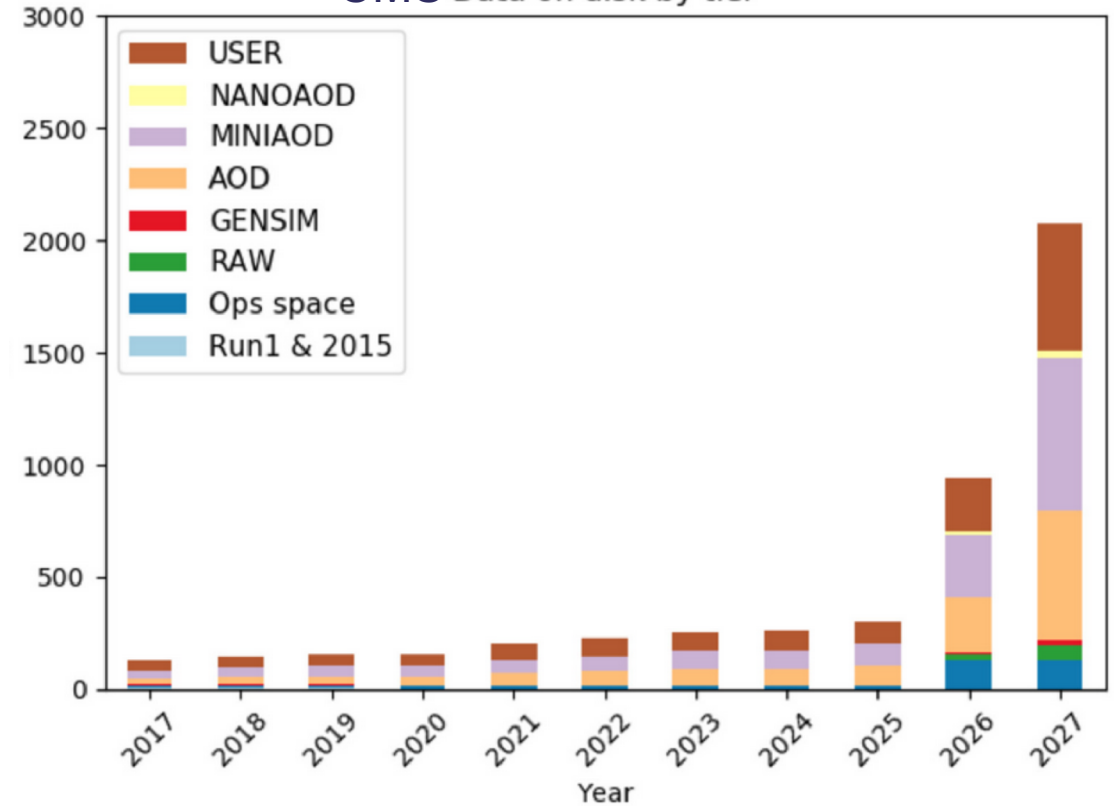




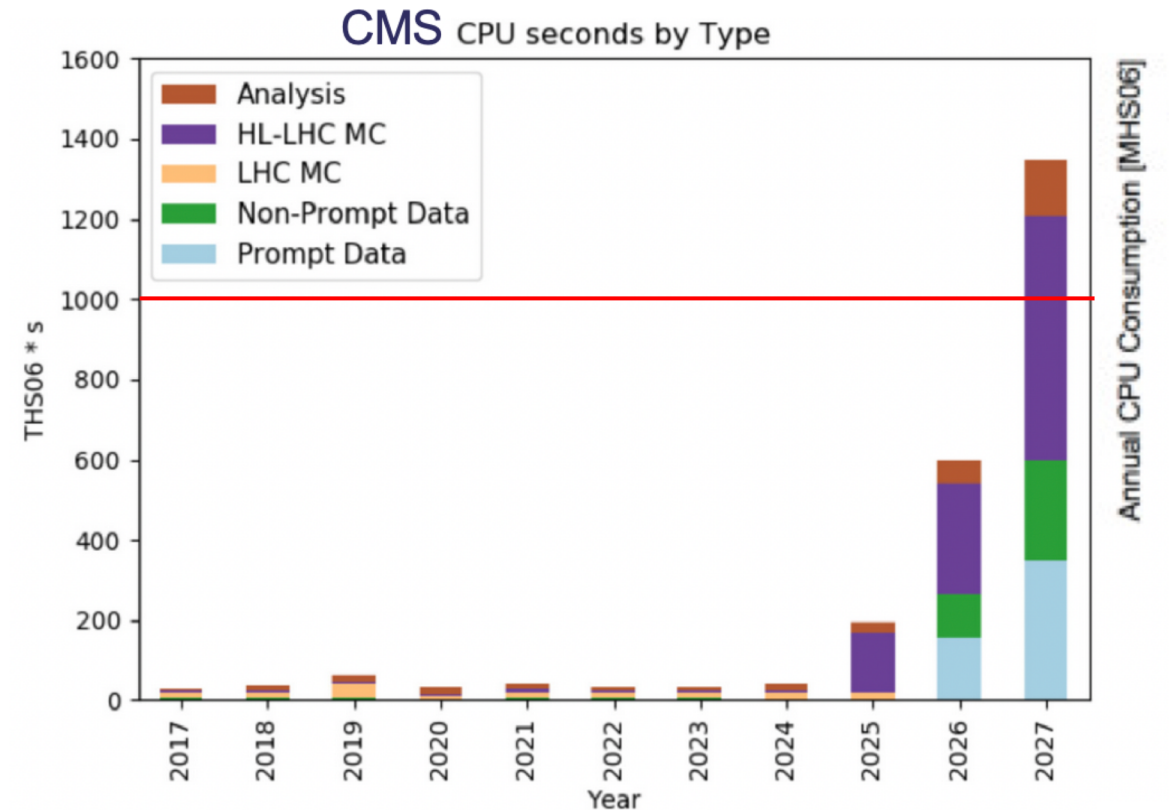
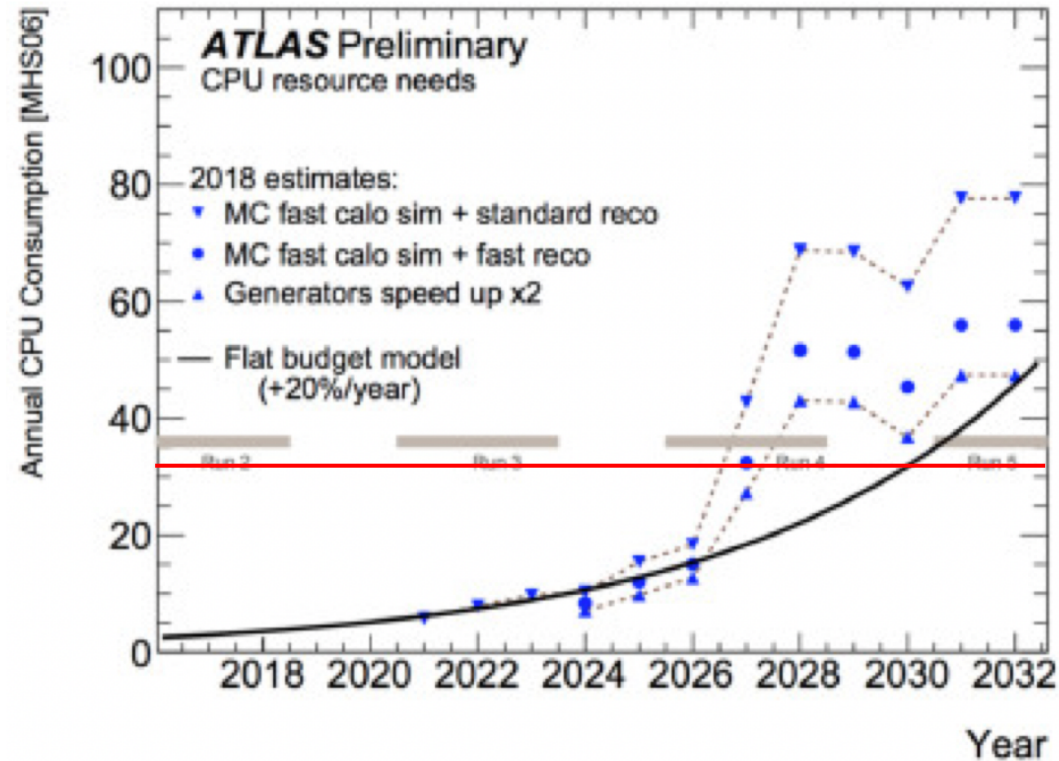
# Disk resource estimation



CMS Data on disk by tier



# CPU resource estimation



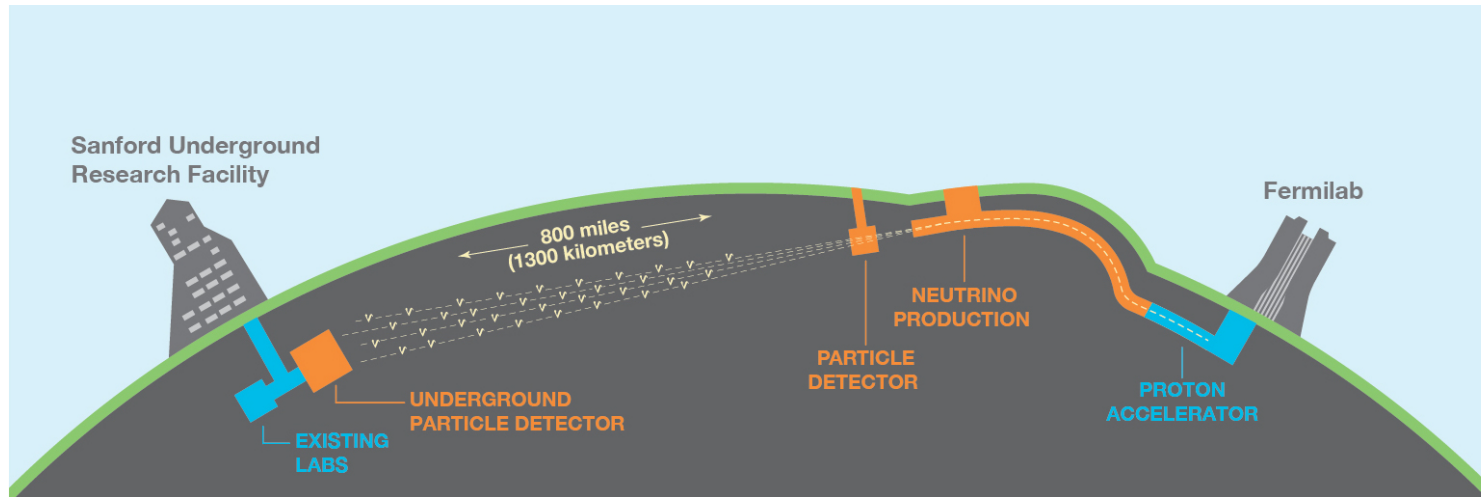
# Tape



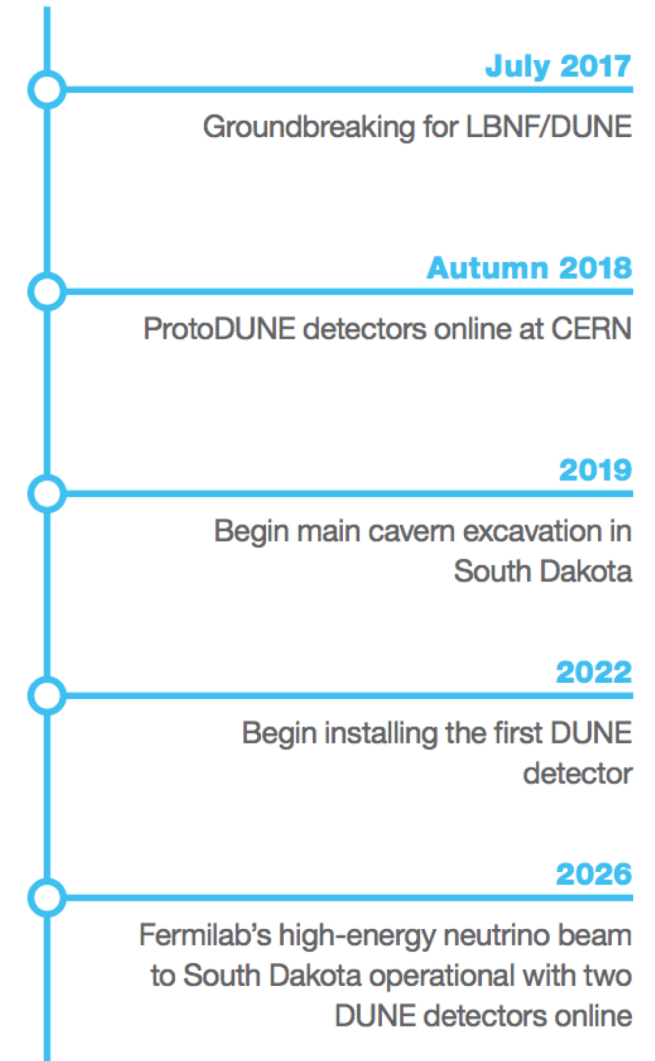
- 2 copies of the LHC experiments RAW data are archived on tape; one copy at CERN while the other is spread across the Tier 1 sites.
  - Other data (normally infrequently used) is also archived on tape.
- Cheap ( $\sim 1/4$  cost of disk) and reliable long term storage.
  - Experiments using tape to prevent them from filling up disk
  - Disadvantages if you want to read the data as it is slow to retrieve
- Tape will be critical for a cost effective computing model.
  - Projects like the ATLAS Tape Carousel aim to optimize recalling data from tape.



# DUNE – Deep Underground Neutrino Experiment



- Pre-trigger data rate is  $\sim 6$  TB/s.
- Limit on the output data rate is estimated to be about 30 PB/year or 8 Gbit/s steady state rate



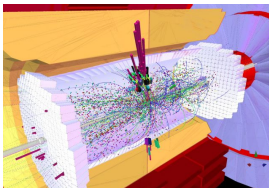


# HEP workflows

“Physics analysis”

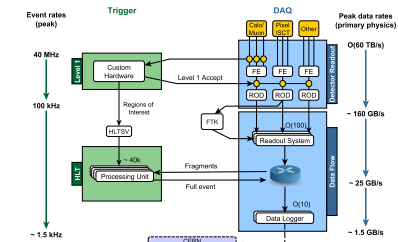
Data reduction

Reconstruction



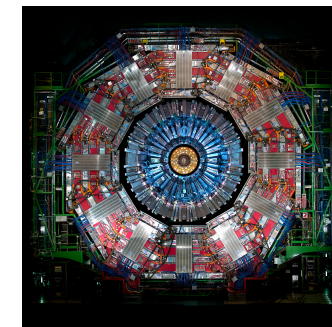
Simulation&DIGI

MC generators



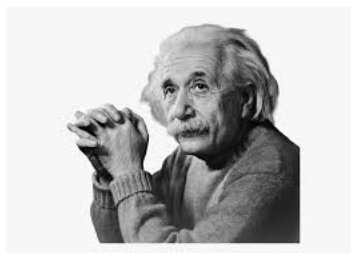
Trigger system

Experiment data



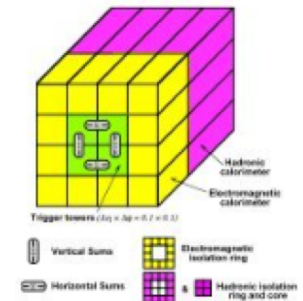
Theory

Experiment



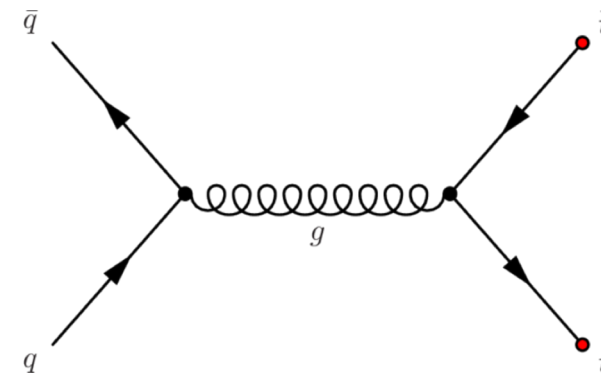
# Event rate and trigger

- The number of proton-proton **bunch-crossings** taking place in CMS and ATLAS is  $\sim 40$  million per second
- Trigger system determines which data events to keep
  - Decision must be made autonomously within an average time of a third of a second
- Experiments are currently limited by available computing to recording  $\sim 1000$  events/s.
  - First level on the detector
  - High Level Trigger run on  $\sim 60k$  cores (CMS)
  - Size of an event: 1-2 MB



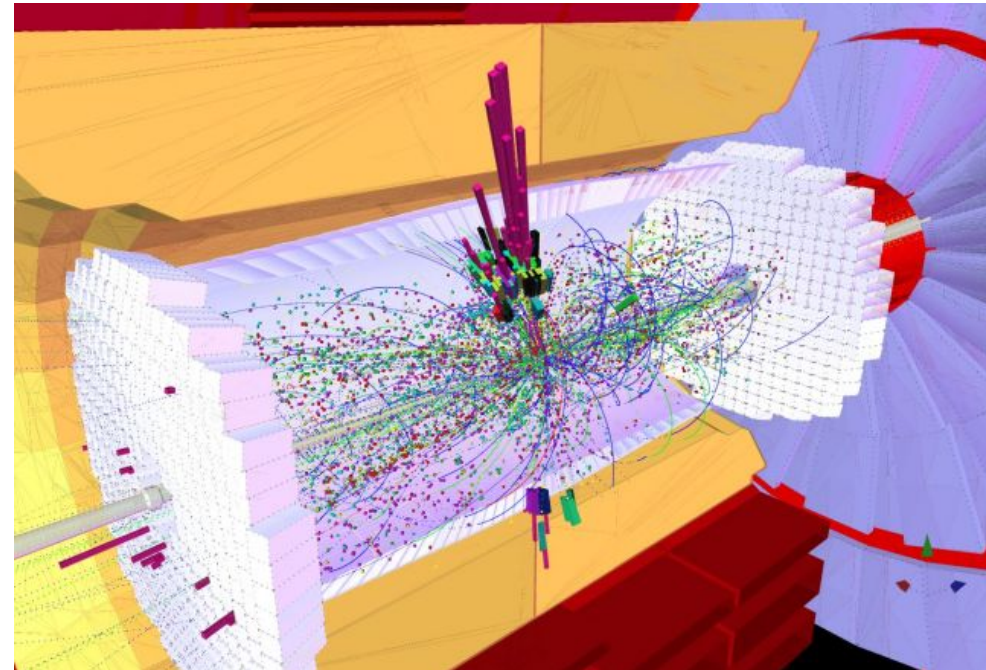
# Monte Carlo generation

- Computer program to model:
  - The “hard scatter”
  - “Final state” particles
  - Including “remnant” and background effects
- Outputs ‘4-vectors’ ( $p_x, p_y, p_z, E$ )
- More complex events require significantly more computing resource.
- Size of event: ~1 MB
- Require typically 1-10+ times as many events as measured experimentally.
- 10+ billion / year



# Simulation & Digitization

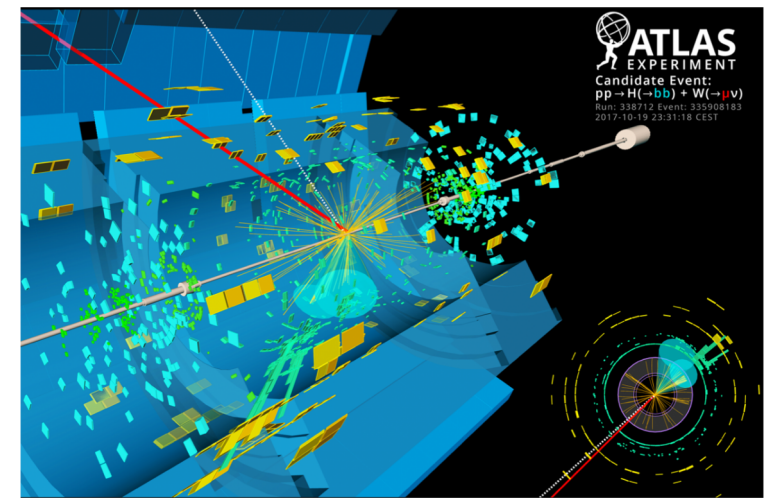
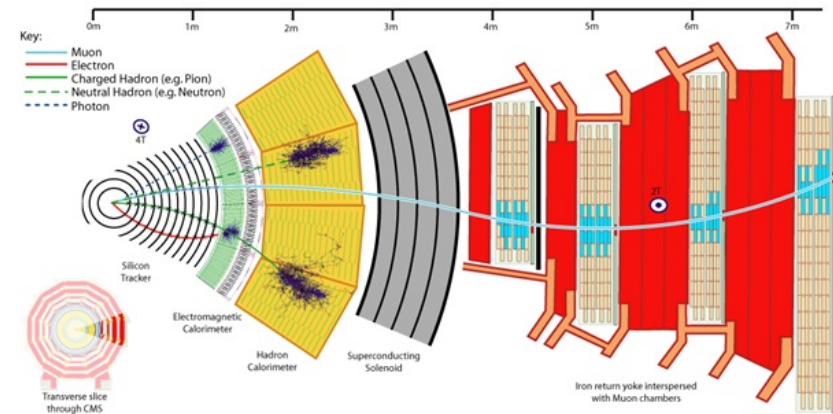
- The MC generator-level data is input to detector-modelling software to simulate e.g. magnetic fields, particle/detector interactions.
- Output is digitized so it resembles detector output, e.g. individual calorimeter cell hits, known detector inefficiencies.
- Including MC gen. step, uses ~40% of computing resources
- Size of event: 1.5 MB





# Reconstruction

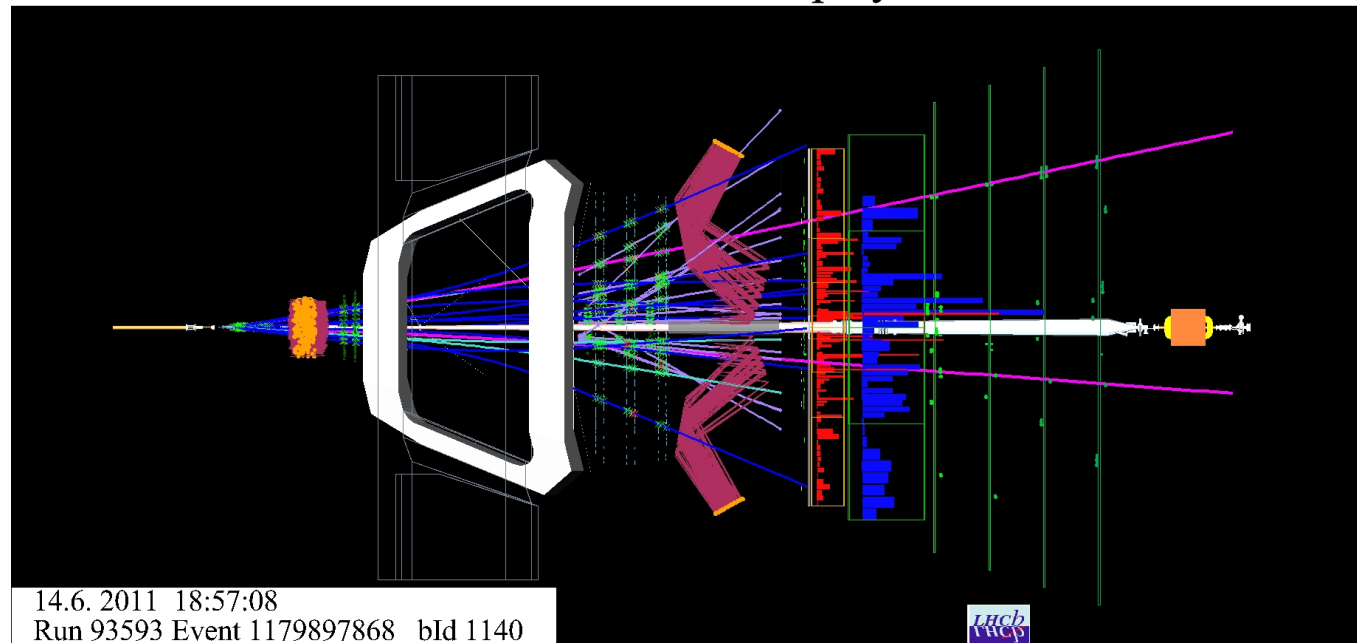
- Uses 'hits' from the detector or the digitized simulation.
  - Using information about where in the detector the hits occurred, energy deposited, track bending, etc.
  - Reconstruction algorithms define the 'physics objects'
    - Electrons, muons, taus, jets, photons, etc.
- Size of event: ~1.5 MB



# Reconstruction - LHCb

- LHCb do something different – calibrate and do reconstruction ‘live’, no hardware trigger
  - Smaller amount of data, not possible with CMS/ATLAS

LHCb Event Display



# Reprocessing?

- Better understanding of detector
  - Improved calibration
  - Ageing of detector
- Better software
- Produce updated event format

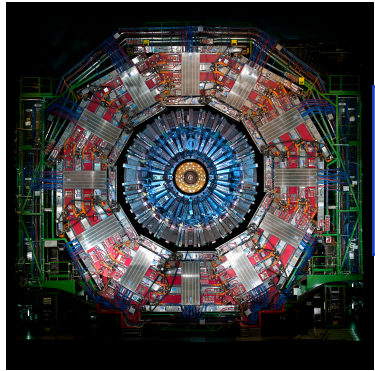
# Data reduction

- Lot of information is recorded, but not required for analysis
- An AOD (Analysis Object Data) file is a distilled data format, containing physics objects as a minimum.
- Data formats are shrinking
- Size of event:  $\sim 0.1$  MB
- Size of CMS 'nanoAOD' format:  $\sim 0.001$  MB

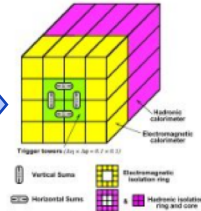




# Summary of data reduction



40 million/s



100,000/s \* 1.5 MB  
150GB/s

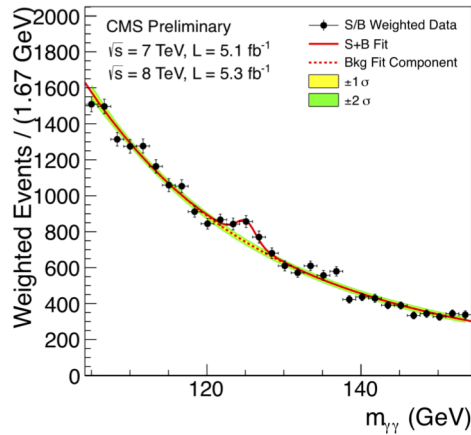


1000/s \* 0.75 MB  
750 MB/s

```
19 template<typename T>
20 unsigned int len1 = s1.size(), len2 = s2.size();
21 const size_t len1 = s1.size(), len2 = s2.size();
22 vector<unsigned int> col(len1+1, 0);
23
24 for (unsigned int i = 0; i < len1; i++) {
25     prevCol[i] = 1;
26     for (unsigned int j = 0; j < len2; j++)
27         col[j+1] = std::min(prevCol[i] + j + 1, len1);
28     col[j+1] = std::min(prevCol[i] + (s1[i]-s2[j]) * 2 + 1, len1);
29     col.swap(prevCol);
30 }
31 return prevCol[len2];
32
33 T* static void
34 ...
```

↑ “Online”

↓ “Offline”



??? \* 0.1 MB



1000/s \* 1.4 MB  
1.4 GB/s

Permanent storage

# Full HEP software stack

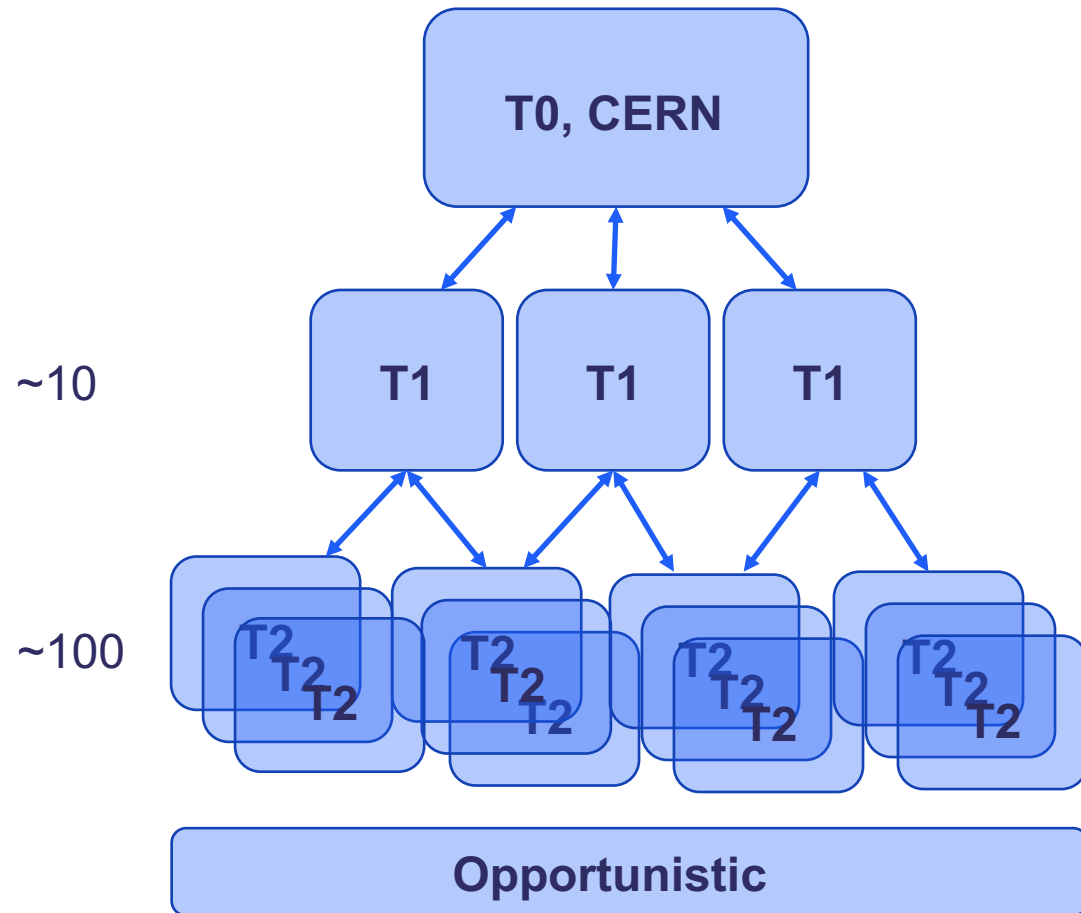
- Mostly talked about the top two layers... but these sit on infrastructure.
- Experiment computing teams coordinate jobs.
- The Worldwide LHC Computing Grid runs the rest at labs and universities.

Layer	Responsible	Experiment 1	Exp. 2	Exp.3
6	Physicist end users	Selecting data, writing analysis code	...	...
5	Physics coders	Analysis frameworks, reconstruction code, calibration code	...	...
4	Computing teams	Central job submission, transfer management	...	...
3	WLCG	Setting up and running middleware, experiment support		
2	WLCG	Software infrastructure for CPU and storage		
1	WLCG	Installing, maintaining and decommissioning physical hardware		

# 1 – Physical hardware

- Due to the way funding works we have computing resources spread across many sites.
  - Universities often provide a lot of additional resources/support.
- This is the opposite of the industry model which uses a small number of very large data centres.
  - The WLCG is consolidating storage at fewer sites but CPU will remain at a large number of sites.
- The result:
  - Resource provision is very heterogenous.
  - Hardware in production for a long time (very slow to move to new tech).

# Computing is distributed



- 150k CPUs HLT
- 100k CPUs, 100 PB disk, 270 PB tape
- 200k CPUs, 230 PB disk, 530 PB tape
- 240k CPUs, 230 PB disk

# Network and data access

- Major sites connected via 100 Gbit/s networks
  - But many smaller sites are 10-40 Gbit/s
- Move towards increased number of jobs accessing remote data
  - Can the network support such an increase in traffic?
  - Especially during data-taking?





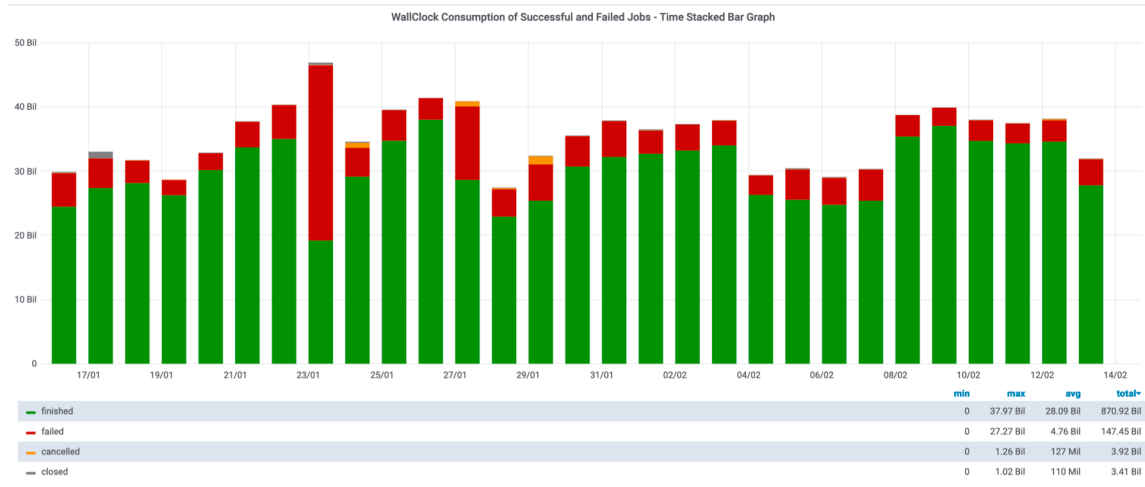
# 2, 3 – Software infrastructure and middleware

- Software infrastructure is the services run directly on top of the physical hardware
  - Storage infrastructure (e.g. dCache, Ceph)
  - Worker node infrastructure (e.g. HTCondor)
- Middleware provides access to the sites resources.
  - CVMFS for software access
  - Frontier/Squid for conditions data
  - CE for access to the batch systems.

# 4 – Experiment distributed computing

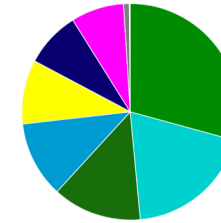
- Workflow management systems like Panda, WMAgent/CRAB and DIRAC/Gaudi run millions of jobs every day on the Grid.
  - Processing/production teams convert physics requests into jobs and check that they produce valid results.
- Data management systems like Rucio and FTS ensure data is correctly located.
- Operations team monitor systems to look for failures.

# ATLAS Jobs

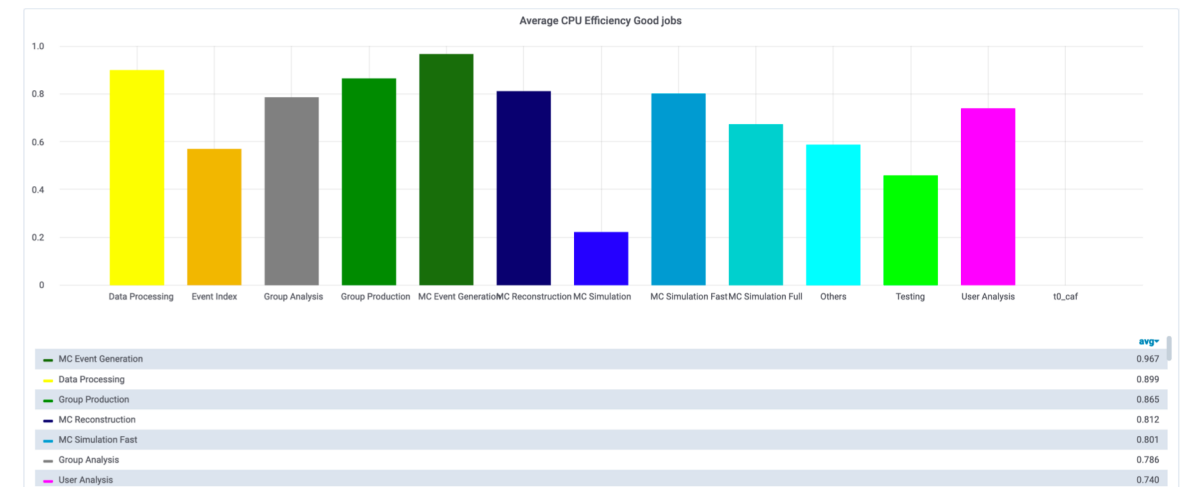


- Plots show ATLAS jobs run in the last 30 days.
- ~13% CPU Wall time lost due to failing jobs.
- Further CPU time lost due to job inefficiencies.

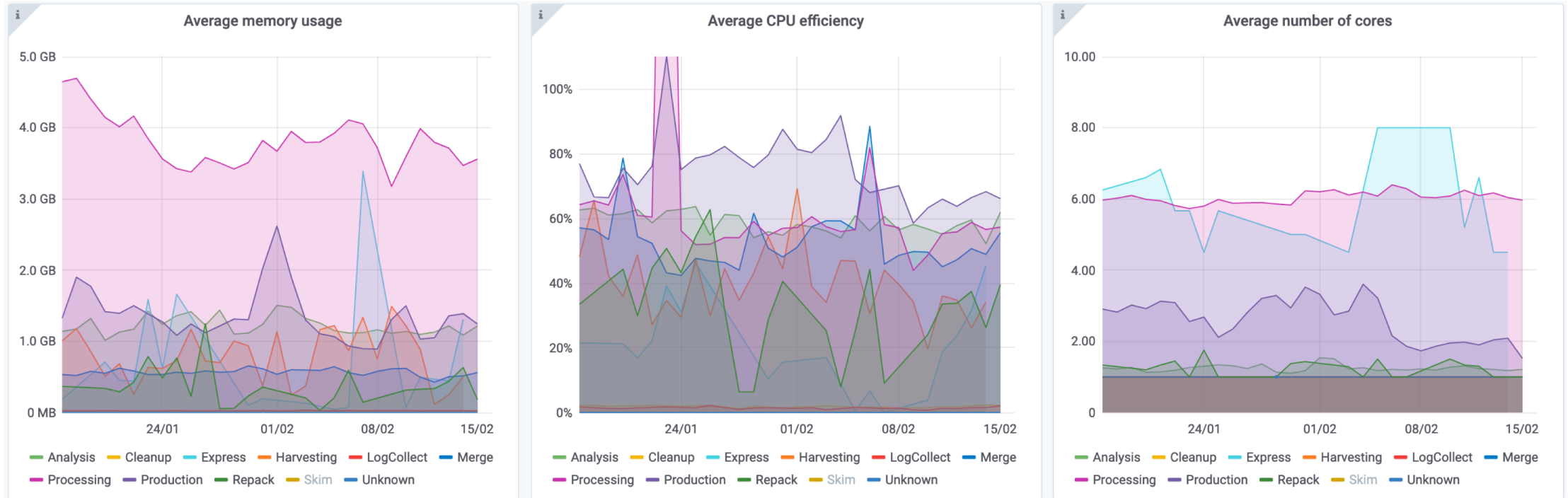
Wallclock Consumption: Successful Jobs in Seconds



	current	percentage
Group Production	258 Bil	29%
MC Simulation Full	170 Bil	19%
MC Event Generation	117.5 Bil	13%
MC Simulation Fast	100.6 Bil	11%
Data Processing	84.7 Bil	10%
MC Reconstruction	73.6 Bil	8%
User Analysis	69.7 Bil	8%
Group Analysis	7.65 Bil	1%
Testing	665 Mil	0%
MC Simulation	424 Mil	0%
Event Index	97.5 Mil	0%
Others	61.0 Mil	0%



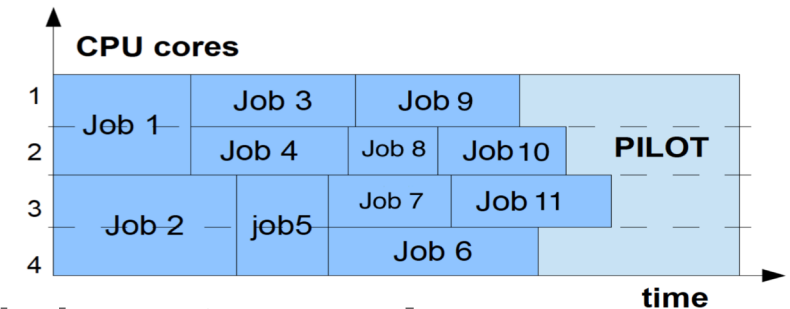
# CMS Jobs



- Plots show the range of different requirements of CMS jobs.

# Further causes of inefficiency

- Jobs run in containers (sometimes multiple layers)
  - Lots of advantages on heterogeneous hardware and environments
  - Overhead – what is the impact on efficiency, can it be improved?
- Jobs run in ‘pilots’
  - Easy allocation of single/multi core jobs
  - Is time wasted?
- Job mix, don’t run all the data intensive jobs at once!
- Are jobs accessing data efficiently?
  - What is lost by accessing data offsite?





# Final words

- HEP has a complex software stack, with a number of well-defined roles.
- High-throughput regime, not HPC
- Amount of data will increase slightly in the next few years, and then significantly after 2027.
- Both processing and volume need huge reductions.
  - Lots of places where efficiencies must be made.



Science and  
Technology  
Facilities Council

# Thank you

**Facebook:** Science and  
Technology Facilities Council

**Twitter:** @STFC\_matters

**YouTube:** Science and  
Technology Facilities Council

# References

- **Computing models in high energy physics:**
  - <https://www.sciencedirect.com/science/article/pii/S2405428319300449>
- **ATLAS Computing and software public results:**
  - <https://twiki.cern.ch/twiki/bin/view/AtlasPublic/ComputingandSoftwarePublicResults>
- **LHCb Tier-1 Resource review:**
  - <https://indico.cern.ch/event/862097/contributions/3631848/attachments/1973409/3284824/20200122-mcnab-lhcb-ral-review.pdf>