



LHC Computing Challenges Data and Workload Management

Alexei Klimentov, Sergei Kuleshov

ISP RAS Ivannikov Open Conference

December 2019

Outline

- ATLAS experiment at Large Hadron Collider
- Computing model in High Energy and Nuclear Physics
- Workflow and workload management
- Data Management
- High Throughput Computing (Grid) vs High Performance Computing
- HEP-Google R&D
- Conclusions

*HEP - High Energy Physics
R&D – Research And Development*

pp, B-Physics, CP Violation
(matter-antimatter symmetry)



LHCb



ATLAS

Lake of Geneva

Pt7: collimators

LHCb

ATLAS

SPS

Pt2: inj b1

ALICE

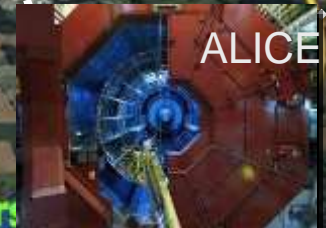
Pt6: dump

CMS

General Purpose,
proton-proton, heavy ions
Discovery of new physics:
Higgs, SuperSymmetry

Exploration of a new energy frontier
in p-p and Pb-Pb collisions
also a new frontier in data

Pt3: collimators



ALICE

Heavy ions, pp
(state of matter of early
universe)

Pt4: RF & BI

LHC 27 km






CMS

SUISSE
FRANCE

The Science Drivers for Particle Physics

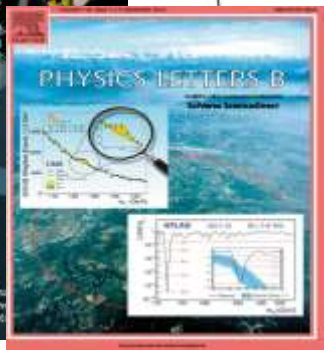
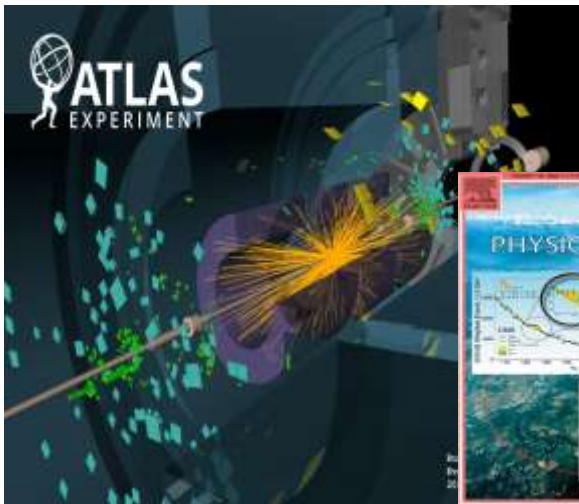
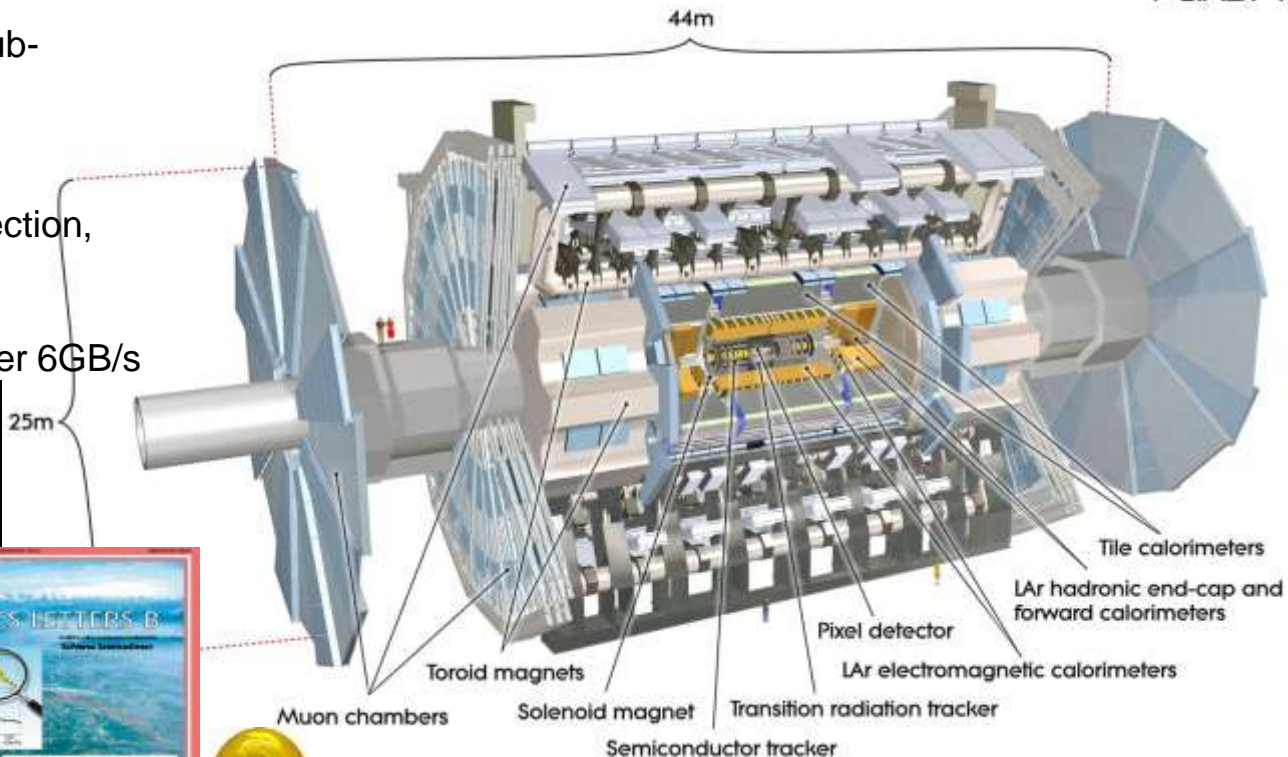
Five intertwined **science drivers**, compelling lines of inquiry that show great promise for discovery :

1. Use the *Higgs boson* as a new tool for discovery. 
2013
2. Pursue the physics associated with **neutrino** mass. 
2015
3. Identify the new physics of *dark matter*.
4. Understand **cosmic acceleration** : dark energy and inflation. 
2011
5. *Explore the unknown* : new particles, interactions, and physical principles.

ATLAS



- The largest detector
- Multiple components and sub-detectors
- 7000 tons
- 10 MW electric power
- 150M sensors measure direction, momentum and charge
- Collisions at 40MHz
 - Filtered to kHz or under 6GB/s



2013

What is this data?

Raw data:

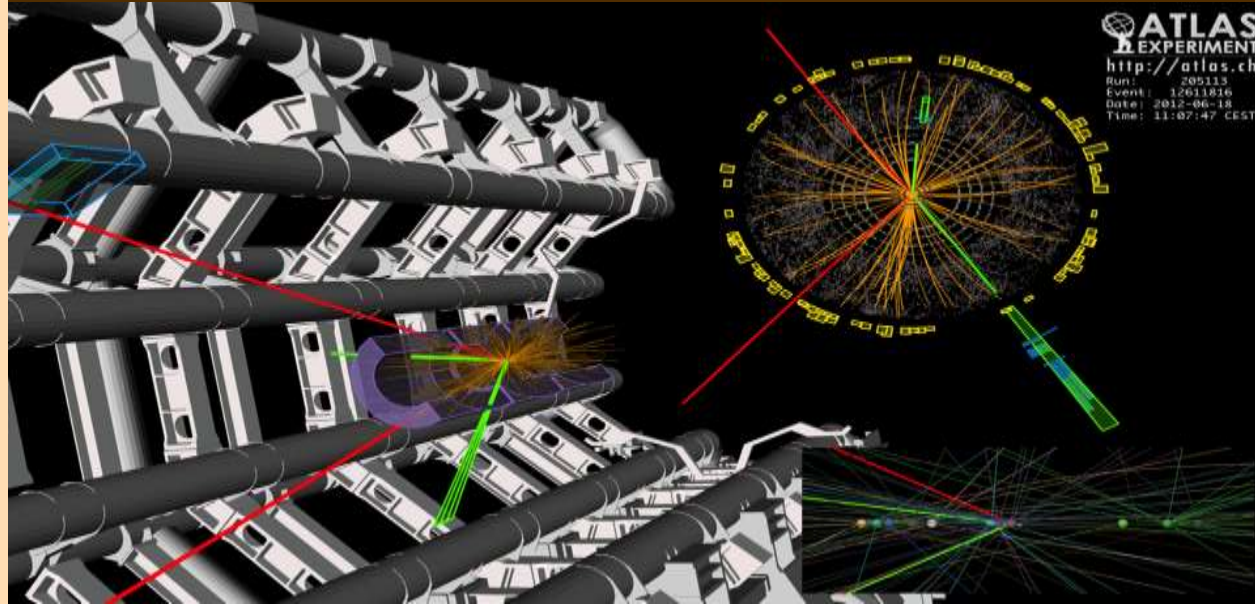
- ✓ Was a detector element hit?
- ✓ How much energy?
- ✓ What time?

Reconstructed data:


- ✓ Momentum of tracks (4-vectors)
- ✓ Origin
- ✓ Energy in clusters (jets)
- ✓ Particle type
- ✓ Calibration information

○ ...

- 150 Million sensors deliver data ... ~ 40 Million times per second
- Up to 6 GB/s to be stored and analysed after filtering




The data processing chain



$$S = i \int d^4x \mathcal{L}(x)$$

2 level, online system (HW+SW)



Reduce event rate from 40 MHz (60TB/s) to 1kHz (1.6GB/s) based on signatures
Event size ~1.6MB



Trigger

Raw data (RAW)

Reconstruction

Analysis Object Data (AOD)

Derivation

Organized production

Chaotic analysis

Detector data

Derived AOD (DAOD)

Analysis



Generation

Event generator output (EVNT)

Simulation

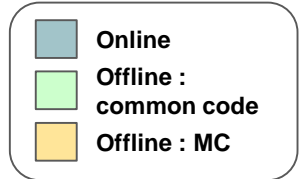
Simulated detector output (RDO)

Reconstruction

Analysis Object Data (AOD)

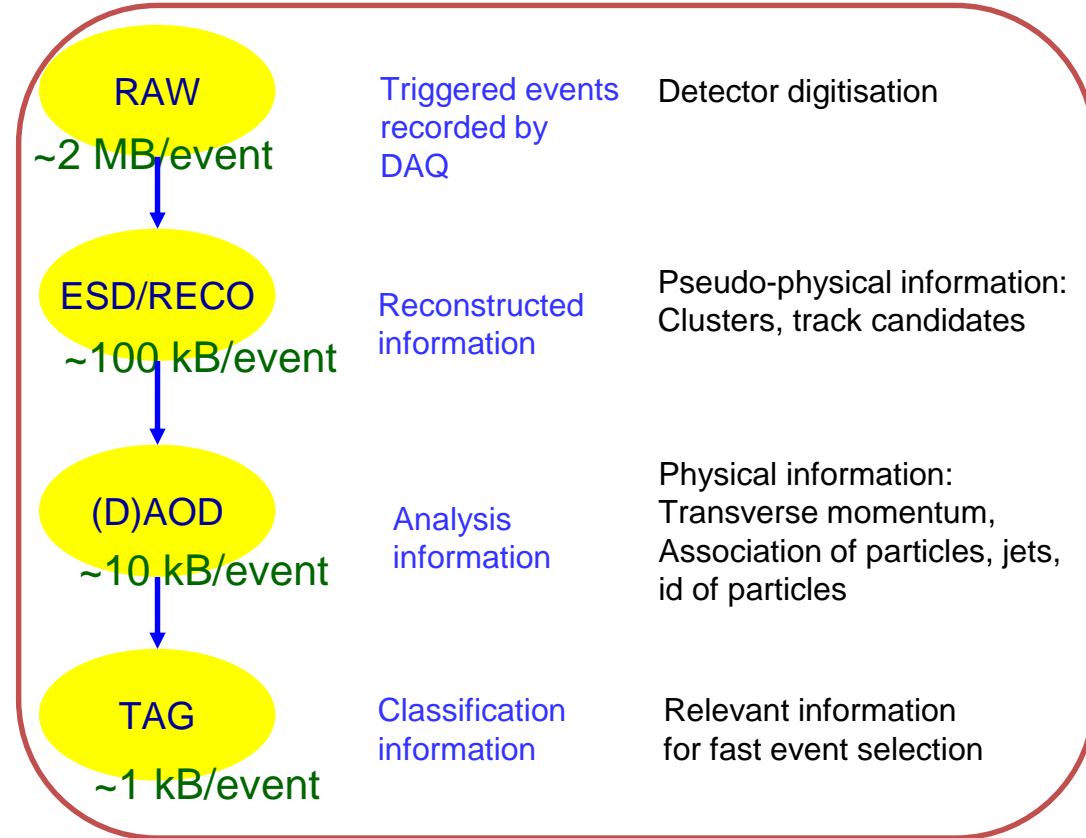
Derivation

Simulated data



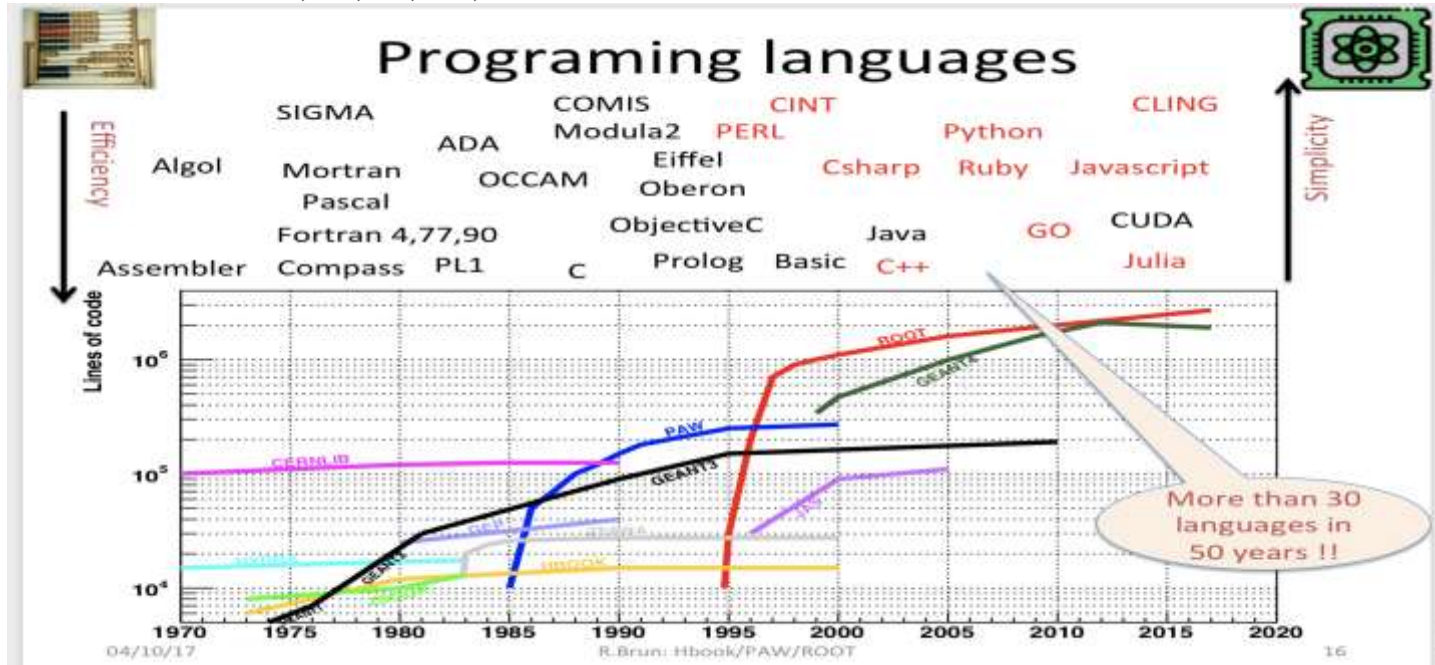
Data and Algorithms

- HEP data are organized as *Events* (particle collisions)
- Simulation, Reconstruction and Analysis programs process “one event at a time”
 - **Events are fairly independent**
→ **Trivial parallel processing**
- Event processing programs are composed of a number of algorithms selecting and transforming “raw” event data into “processed” (reconstructed) event data and statistics
- *ATLAS reconstruction and simulation code 5M LOC*
- *1000 software developers*

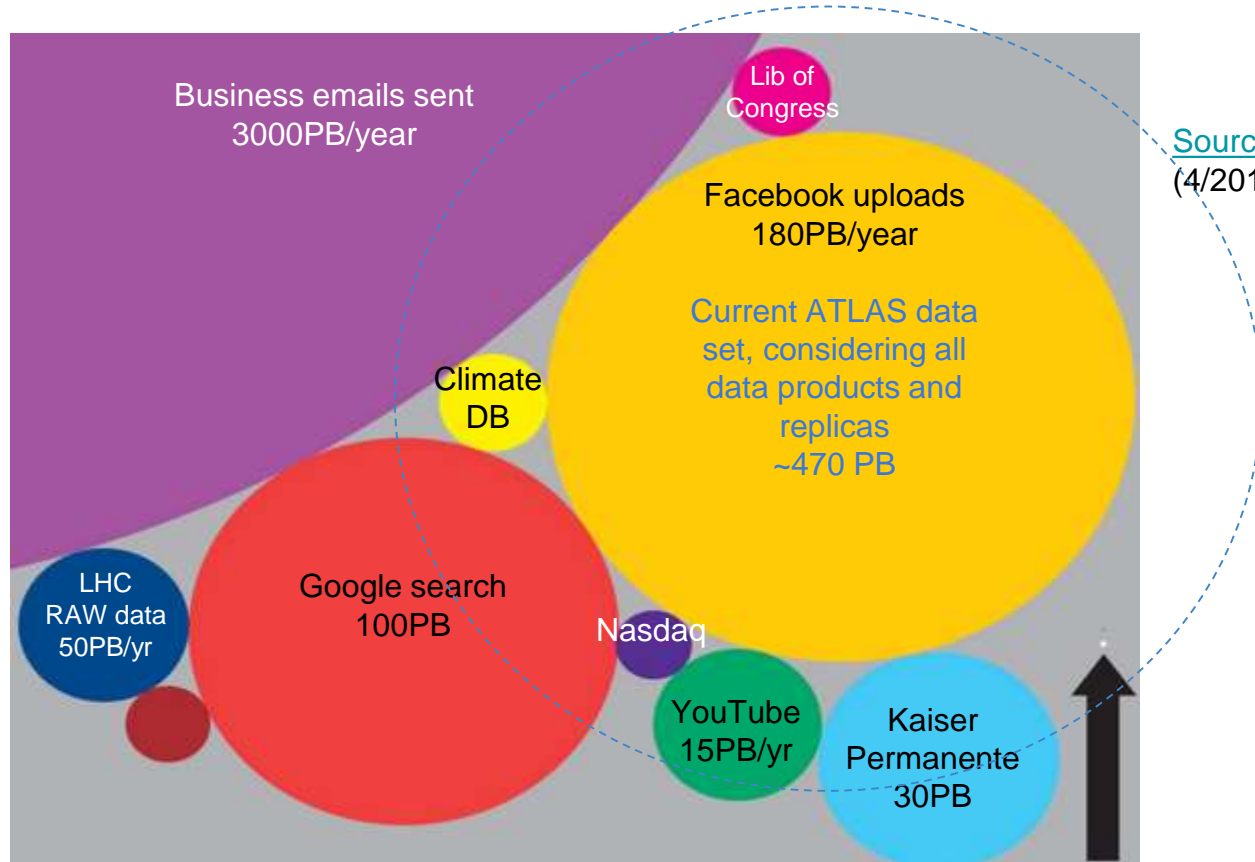


Software Tools, Libraries, Languages

- 1974 : Fortran IV, *HBOOK*, *GD3*, *HYDRA*, *CERNLIB*, *SIGMA*
- 1979 : Fortran 77, Pascal, *CERNLIB*, *GEANT2*
- 1984 : ADA, *ZEBRA*, *GEANT3*, *GHEISHA*, *PAW*, *X11*
- 1989 : *PHIGS*, *MOTIF*, FORTRAN90, C, Perl, ORACLE, Web
- 1994 : OO, Smalltalk, Eiffel, **C/C++**, Objectivity, *GEANT4*, *MOSAIC*, *ROOT*
- 1999 : JAS, Mathematica, Netscape, *PROOF*, *XROOTD*
- 2004 : Google, **Python**
- 2014 : GPUs/AI, ML, DL, GO, CLING

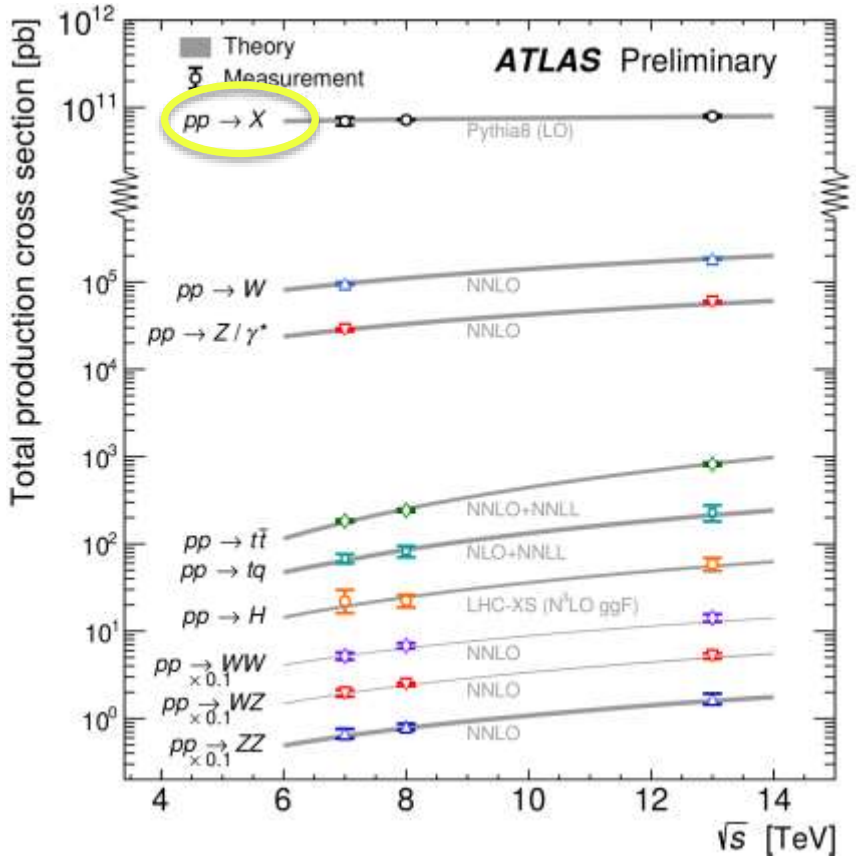


Magnitude of the problem. BigData



Source: [Wired magazine](#)
(4/2013 - a bit outdated)

Magnitude of the problem. How likely something interesting happen



- ◆ Total Production Cross Section (== probability) vs Energy in pp collisions
- ◆ Notice the logarithmic scale on the Y-axis: it spans 11 orders of magnitude
- ◆ E.g. you produce 10 Higgs bosons out of 10^{11} billions of collisions
- ◆ The probability increases logarithmically with energy
- ◆ Theory (lines) agrees very well with measurements (markers)

Paradigm shift in Particle Physics Computing in XXI century

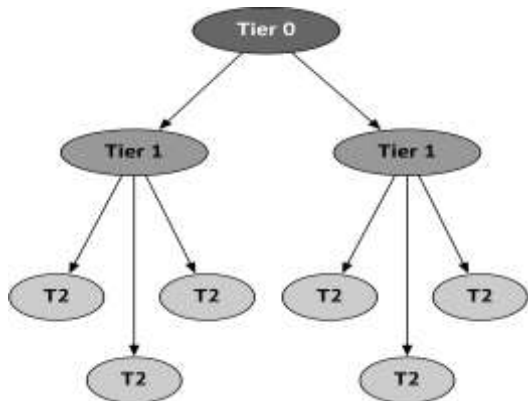
| Old paradigms | New ideas |
|---|---|
| <ul style="list-style-type: none">● Distributed resources are independent entities | <ul style="list-style-type: none">● Distributed resources are seamlessly integrated worldwide through a single submission system● Hide middleware while supporting diversity |
| <ul style="list-style-type: none">● Groups of users utilize specific resources (whether locally or remotely) | <ul style="list-style-type: none">● Access to all resources may be granted to all users |
| <ul style="list-style-type: none">● Fair shares, priorities and policies are managed locally, for each resource | <ul style="list-style-type: none">● Global fair share, priorities and policies allow efficient management of resources |
| <ul style="list-style-type: none">● Uneven user experience at different sites, based on local support and experience | <ul style="list-style-type: none">● Automation, error handling, and other features improve user experience● Central support coordination |
| <ul style="list-style-type: none">● Privileged users have access to special resources | <ul style="list-style-type: none">● All users have access to same resources |

The Worldwide LHC Computing Grid

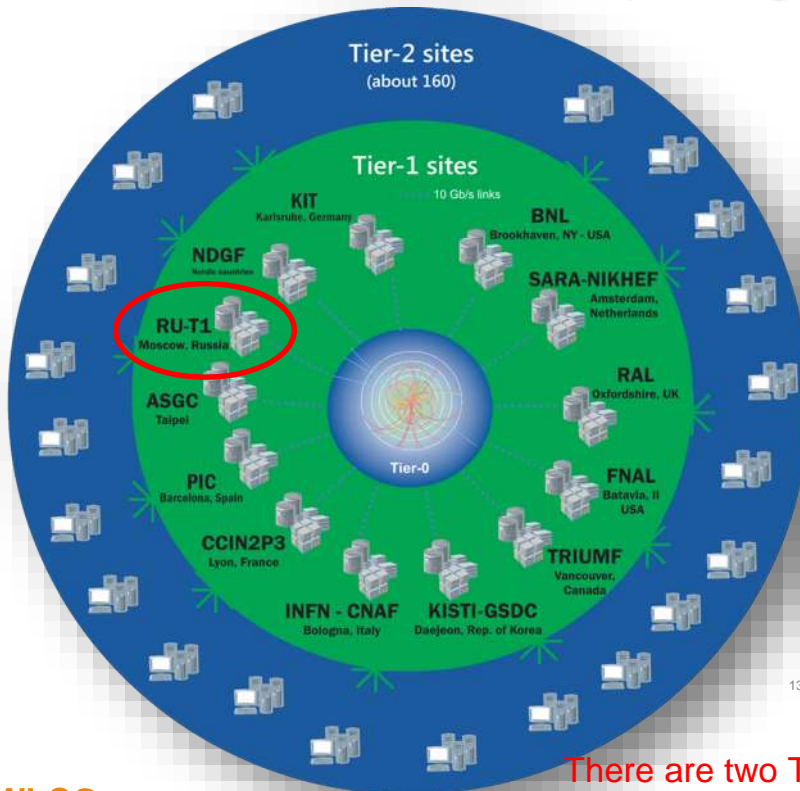
Tier-0 15%: (CERN and Hungary):
data recording,
reconstruction and
distribution

Tier-1 40%: permanent
storage, re-processing,
Analysis
T0 spill-over
HLT
MC Simulation
Derivation production

*MONARC - Models of
Networked
Analysis at Regional Centres for
LHC Experiments.*



12/24/2019



Tier-2 45%: Simulation,
end-user analysis
Re-processing
Derivation production

~170 sites,
42 countries

~750k CPU cores

~1 EB of storage

> 2 million jobs/day

10-100 Gb links

There are two Tier-1s in Russia : JINR and NRC KI

WLCG:

An International collaboration to distribute and analyse LHC data

Integrates computer centres worldwide that provide computing and storage resource into a single infrastructure accessible by all LHC physicists

Primary distributed computing software tools

Workflow Management:

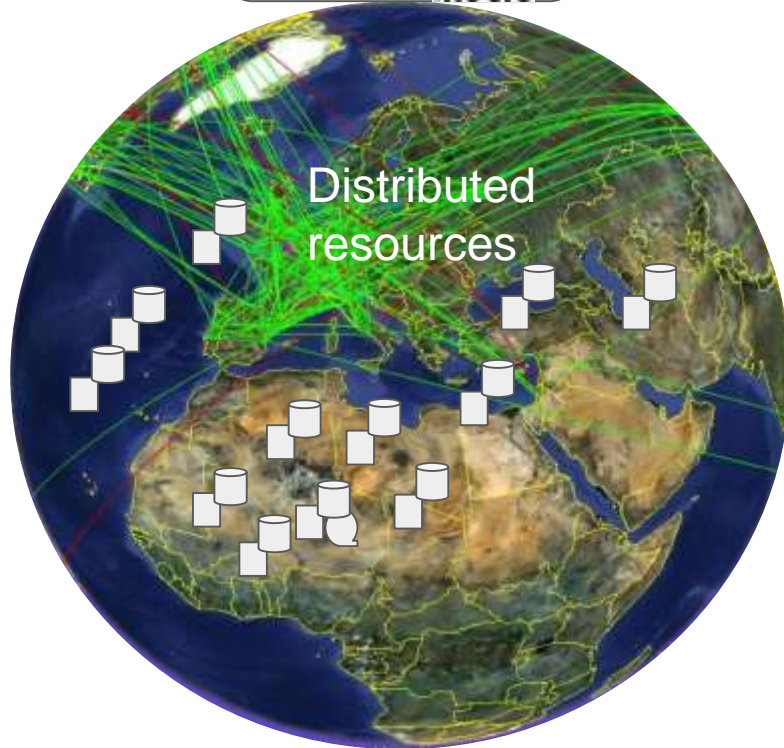
“translates” physicist requests into production tasks



Workload Management:

submission and scheduling of jobs & tasks

Monitoring production jobs & tasks, shares, users



Data Management:

bookkeeping and distribution of files & datasets

Information System

(ORACLE backend)

Queues and resources description

Databases: Conditions and data processing (ORACLE, mySQL, PostgreSQL)



Workflow and Workload Management. PanDA

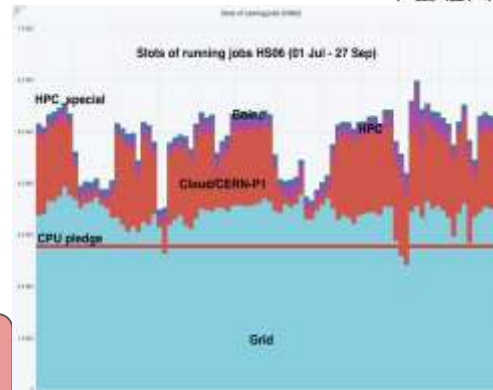
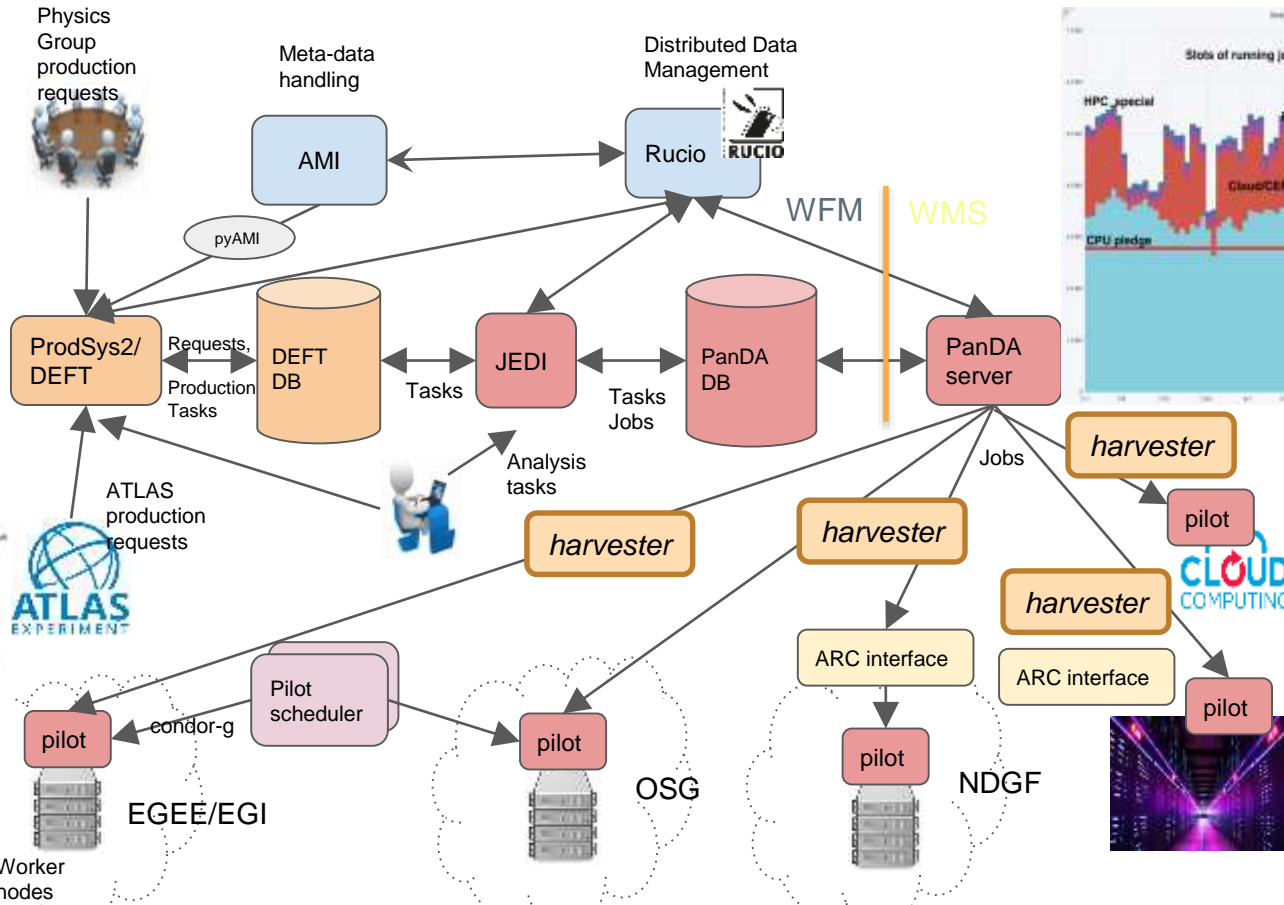
Global ATLAS operations

Up to ~1.2M concurrent jobs
25-30M jobs/month at >250 sites
~1400 ATLAS users

Orchestrate all ATLAS Workflows :

- MC Production
- Physics Groups WF
- Data reprocessing
- T0 spill-over
- HLT processing
- SW validation
- User's analysis
- ART

Shares/priorities



First exascale workload manager in HENP
1.4+ Exabytes processed yearly in 2014/18
Exascale scientific data processing today



HPCs

Support rich harvest of heterogeneous resources. Integrate WF and data flow

Workload Management System Summary and Lessons learned

- ***We designed and implemented a scalable, flexible, automated production that follows physics priorities***

- Steady state production 24x7x365 with ~300-350k cores across ~140 sites
- HPC peaks to >1M cores, demonstrating extreme scalability of PanDA
- The system orchestrates ~10 principal workflows and dozens of variants, with automated shares that follow ATLAS physics priorities and allocate work across global resources
- Also supporting over 1000 analysis users with fair sharing of resources

- ***Integrated workflow and dataflow***

- Moving >1 PB, >20 GB/s, 1.5-2M files per day
- 405PB disk+tape, 1+B files in total (and ~540PB in 2019)
- **PanDA processes over 1.5 Exabytes per year**

PanDA was adopted for many HEP and astro-particle physics experiments

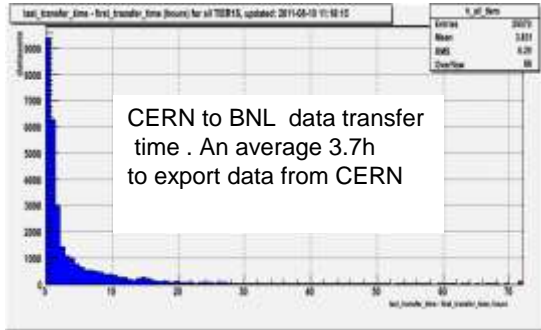
- WMS is designed by and serves the physics community
- WMS new features are driven by experiment operational needs
- WMS functionality is important as scalability
- Computing model and computing landscape in general has changed

*There are several systems with very well defined roles which are integrated for physics computing : Information system (AGIS), DDM (Rucio), WMS (ProdSys2/PanDA), meta-data (AMI), and **middleware (HTCondor, Globus...)**. We managed to have a good integration of all of them in PanDA.*

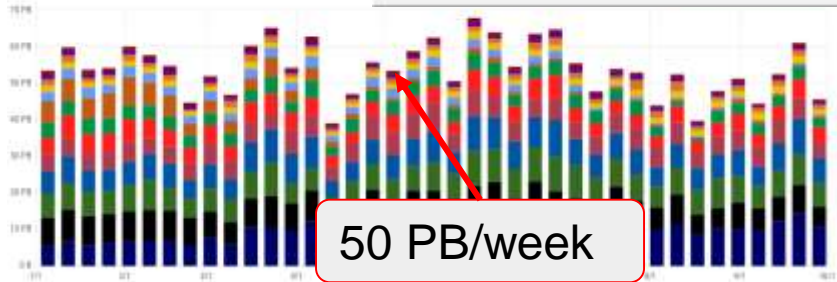


Data management. Rucio

- A few numbers to set the scale
 - 1B+ files, 460+ PB of data, 400+ Hz interaction
 - 120 data centres, 5 HPCs, 2 clouds, 1000 users
 - 500 Petabytes/year transferred & deleted
 - 2.5 Exabytes/year downloaded & uploaded

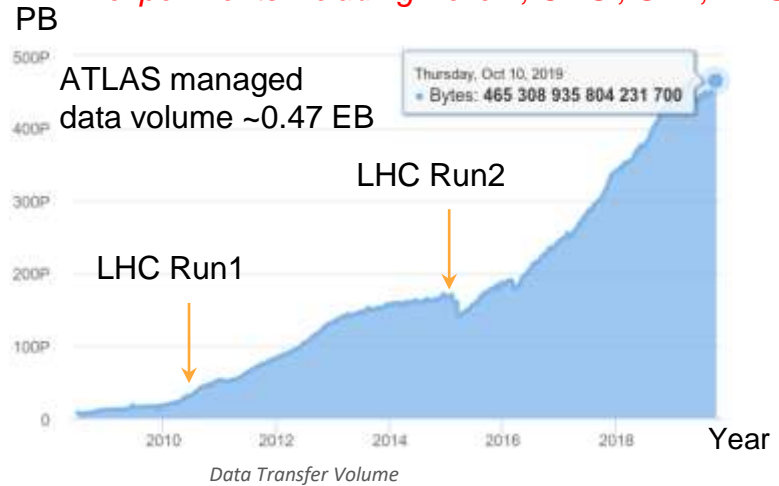


Data access volume

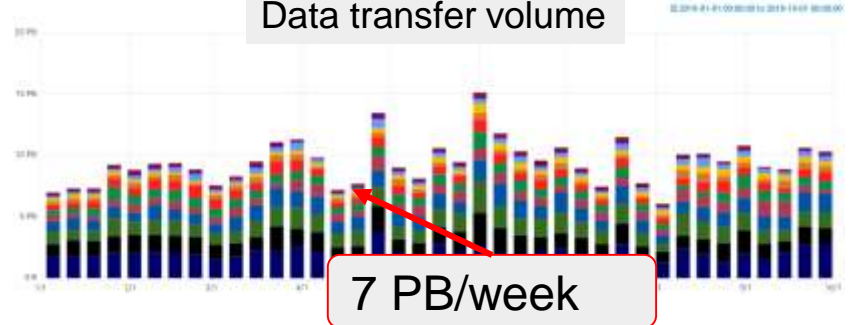


First exascale scientific data management system today

Rucio is evaluating or already in use for many experiments including Belle II, CMS , SKA, AMS



Data transfer volume





Monitoring and Analytics (bigpanda.cern.ch)

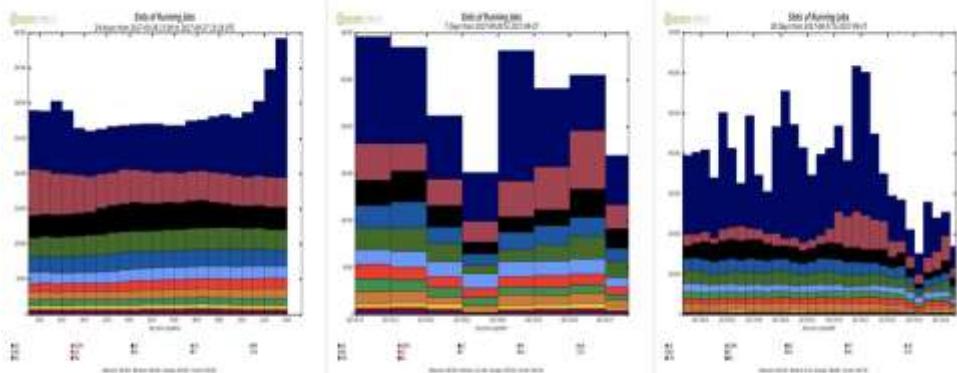
ATLAS PanDA [Dash](#) [Tasks](#) [Jobs](#) [Errors](#) [Users](#) [Sites](#) [Incidents](#) [Search](#) [Admin](#)

Friday, 15 Nov 2013 10:18:14 [Refresh](#) [Logout](#)

ATLAS PanDA monitor home

spanid:00 1511814 [Refresh](#) [Logout](#)

Global concurrent running job core counts, all sites, all job types, by cloud, last 1, 7, 30 days



Global concurrent running job core counts, all sites, all job types, by activity, last 1, 7, 30



- Upgrade
- Reprocessing default
- Data Derivations
- Event Index
- MC production
- MC Derivations
- Analysis
- HLT Reprocessing
- Heavy Ion
- Test
- Final production
- Validation
- MC Default
- MC 18



Task 11016615

Task ID: 11016615

Owner: [User]

Created Date: 2013-11-15 10:18:14

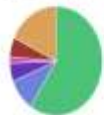
Modified: 2013-11-15 10:18:14

Production to complete: 2013-11-15 10:18:14

Current progress: 100%

Task completion: 100%

HL Analysis jobs per job type (count)



HL Analysis jobs per job type (count)



HL Analysis jobs per job type (event)

- DarkD-based analysis
- Event generation
- Reconstruction
- Non-T analysis
- Simulation
- uCCO based analysis



HL Analysis jobs per job type (event)



- DarkD-based analysis
- Event generation
- Reconstruction
- Non-T analysis
- Simulation
- uCCO based analysis

HL Analysis jobs per job type (event)

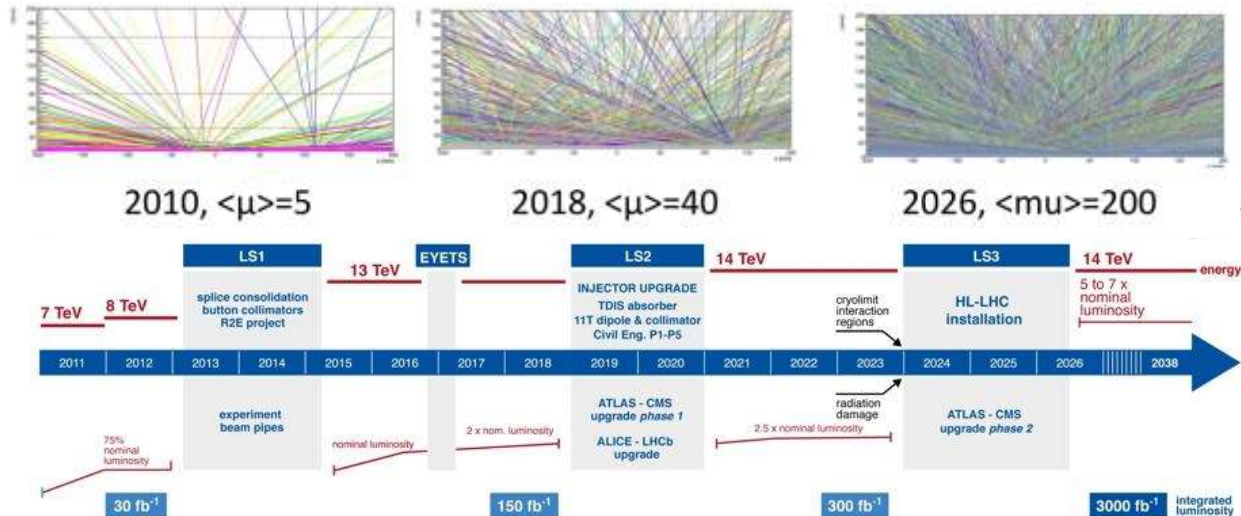


More in Tatiana's and Sasha's talks

High Luminosity LHC (HL-LHC). More Challenges Ahead

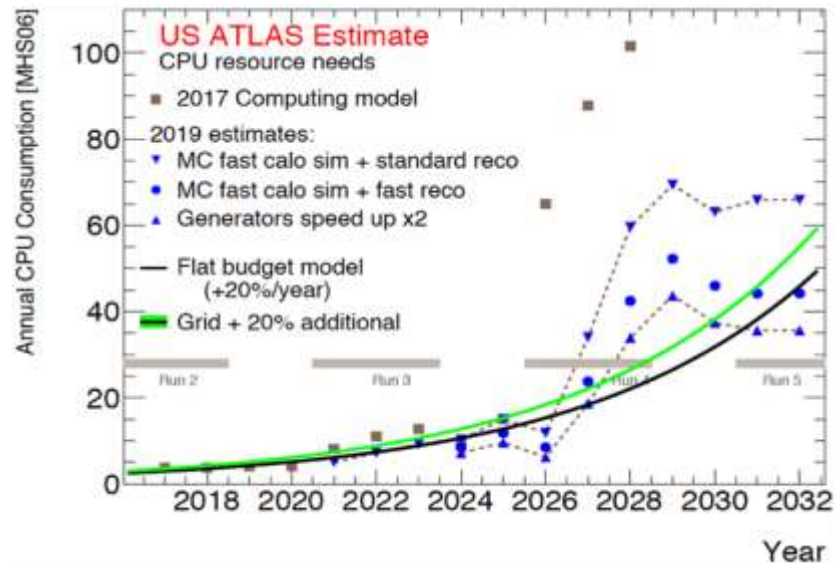
Starting in 2026, the LHC enters a new era

- 5-7x LHC original design luminosity, pileup up to 200 interactions/crossing
- Extensively upgraded, more complex detectors: 4-5x increase in event size
- Upgraded trigger system: up to 10x increase in event rate
- **A 10+ year program of precision and discovery physics with a ten-fold increase in integrated luminosity**

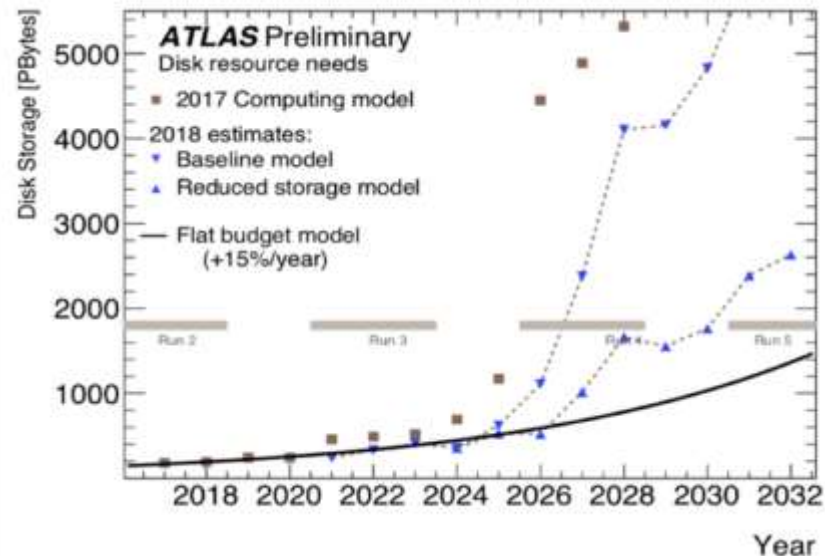


The High Luminosity LHC Challenge

Growth in CPU Needed



Growth in Disk Storage Needed



- High Luminosity LHC will be a multi-exabyte challenge where the envisaged Storage and Compute needs are a factor 10 to 100 above the expected technology evolution.
- LHC experiments have successfully integrated HPC facilities into its distributed computing system. “Opportunistic storage” basically does not exist for LHC experiments.
- The HEP community needs to evolve current computing and data organization models in order to introduce changes in the way it uses and manages the infrastructure, focused on optimizations to bring performance and efficiency not forgetting simplification of operations.

CPU Challenge. High Performance Computing vs High Throughput Computing

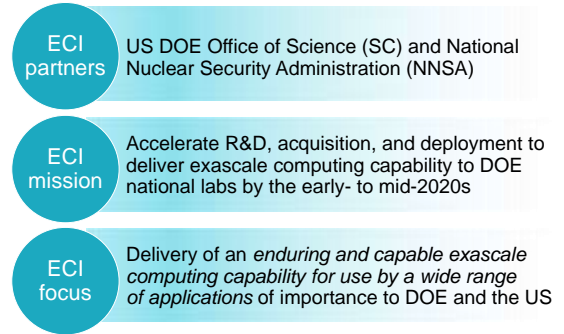
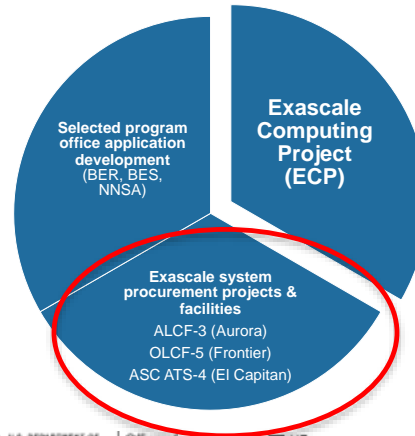
- High Performance Computing
 - Leadership Computing Facilities (Summit, Theta,...)
 - Massively Parallel Processors
 - Clustered Machines
- Grid – High Throughput Computing
 - LHC experiments are extremely successful using Grid
 - Variety of technologies and middleware
 - Processors homogeneous –x86

High Performance Computing Project in US

- *Get extra computing resources on demand*
- *The Worldwide LHC Computing Grid and a leadership computing facility (LCF) are of comparable compute capacity.*
 - *WLCG (ATLAS share): 300,000's x86 compute cores*
 - *Titan: 300,000 x86 compute cores and 18,000 GPUs*

DOE Exascale Program: The Exascale Computing Initiative (ECI)

Three Major Components of the ECI



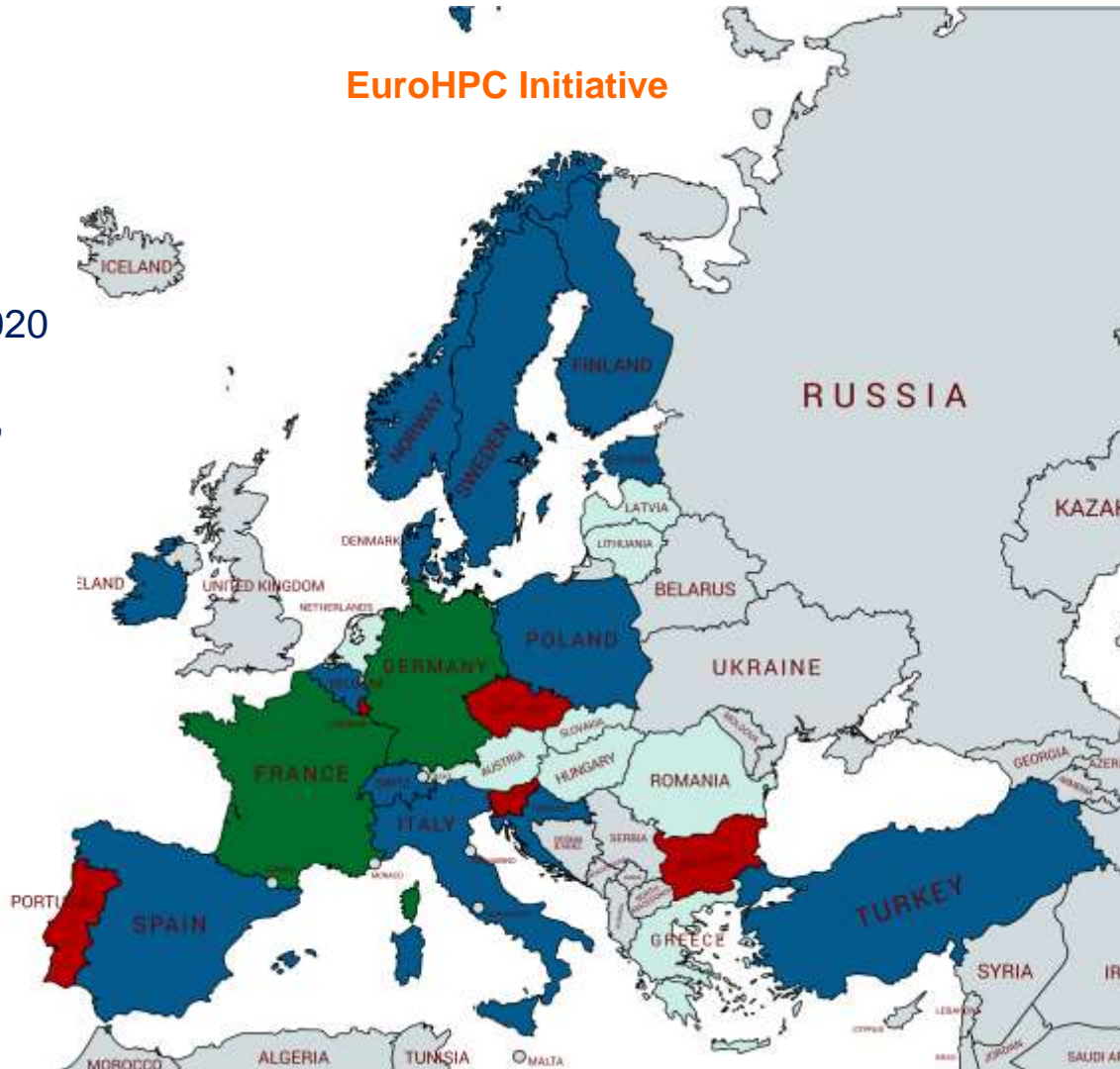
How do we efficiently integrate HPC resources and run canonical physics workflows on them ?

- ❑ **Five EuroHPC-JU petascale systems, 50+PFLOPS**, installed by 2020
- ❖ Meluxina, Deucalion, Euro-IT4I, Vega, PetaSC

- ❑ **Three EuroHPC-JU pre-exascale consortia 600+PFLOPS**, installed by 2020
- ❖ **Lumi, BSC, Leonardo**
- **Lumi (Kajaani Data Center, FI)** : Finland, Belgium, Czech Rep, Denmark, Estonia, Norway, Poland, Sweden, Switzerland

- ❑ **Two exascale sites**

EuroHPC Initiative



A blueprint for the new
Strategic Research Agenda
for High Performance Computing



April 2019

HPC-Grid integration. Research Challenges

- Motivation:
 - Demand for resources greatly outstrips supply
 - New payloads are being developed well suited for HPC architectures
- Challenge: Resource management approaches, policies and software are different for each HPC, and Grids, and Clouds
- Objectives:
 - How can we provide existing workload management capabilities (PanDA) uniformly across HPC, Grids and Clouds?
 - Use it to utilize otherwise unusable resources on HPC as well as traditional allocations queue-based resource management

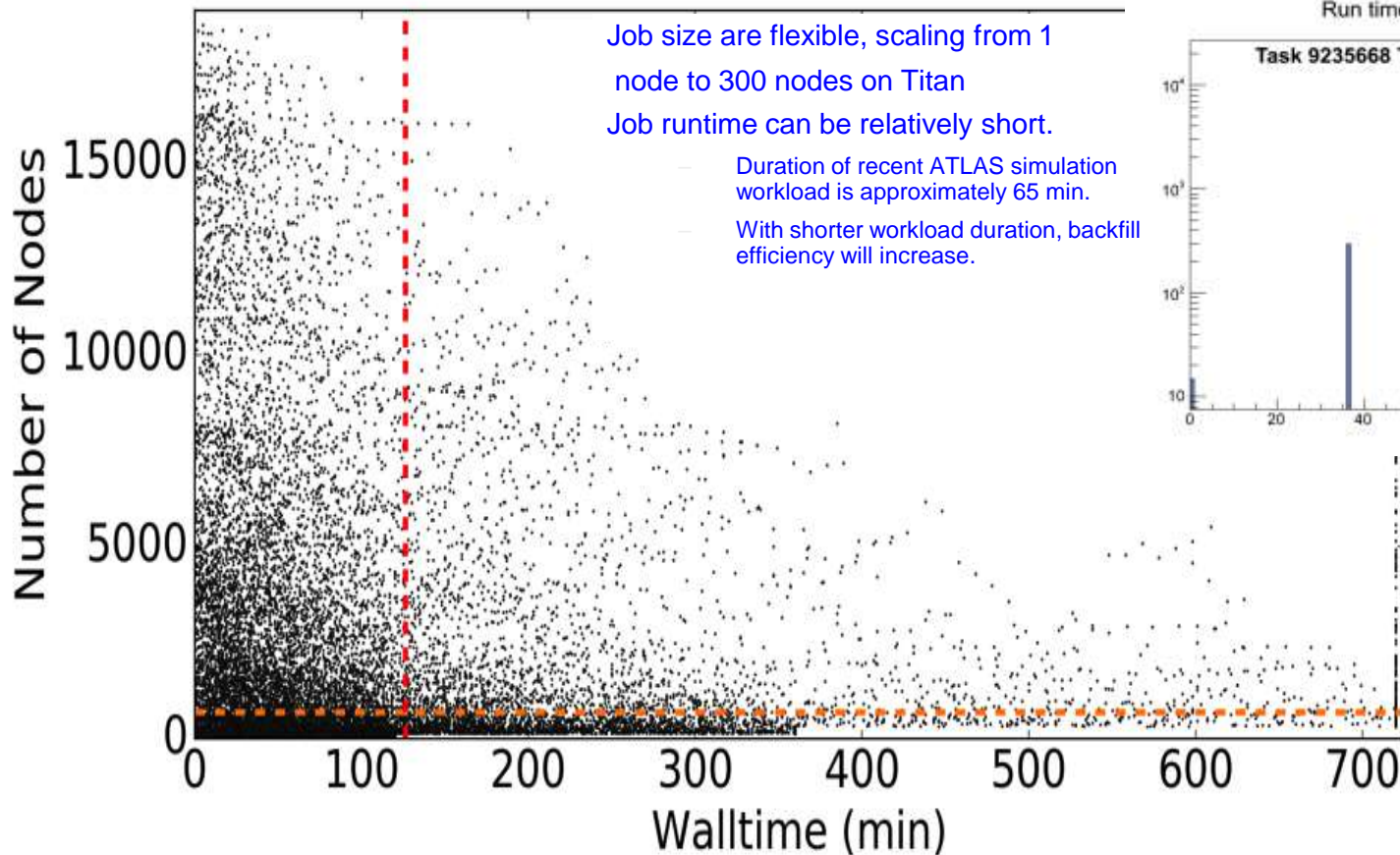
Deploying PanDA on Titan

- Titan's architecture, configuration and policies posed several challenges to the deployment of PanDA e.g., :
 - The default deployment model of PanDA Pilot is unfeasible on Titan:
 - Pilot requires communication with the PanDA Server to pull jobs to execute, but
 - This not possible on Titan because worker nodes do not offer outbound network connectivity
 - The specific characteristics of the execution environment require re-engineering of ATLAS application frameworks, e.g., absence of local storage on the worker nodes, and modules tailored to Compute Node Linux
- Given various constraints and challenges, the Monte Carlo detector simulation task is most suitable for execution on Titan
 - Accounts for \approx 60% to 70% of all the jobs on WLCG, making them a primary candidate for offloading
 - Task is mostly computational-intensive, requiring less than 2GB of RAM at runtime and small input data



**Titan: 27 PF
Accelerated Computing
World's Fastest**

Learning Backfill Capabilities on Titan



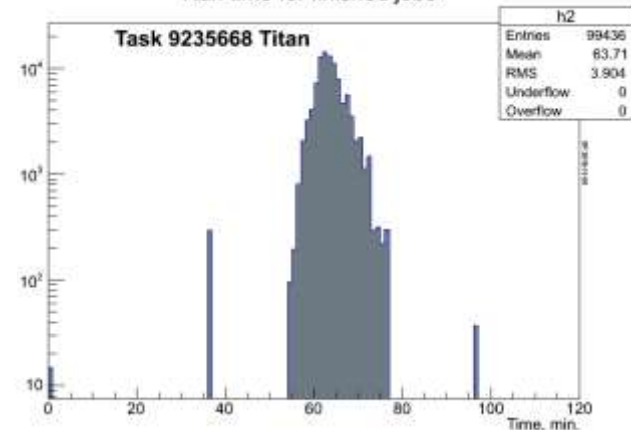
Job size are flexible, scaling from 1 node to 300 nodes on Titan

Job runtime can be relatively short.

Duration of recent ATLAS simulation workload is approximately 65 min.

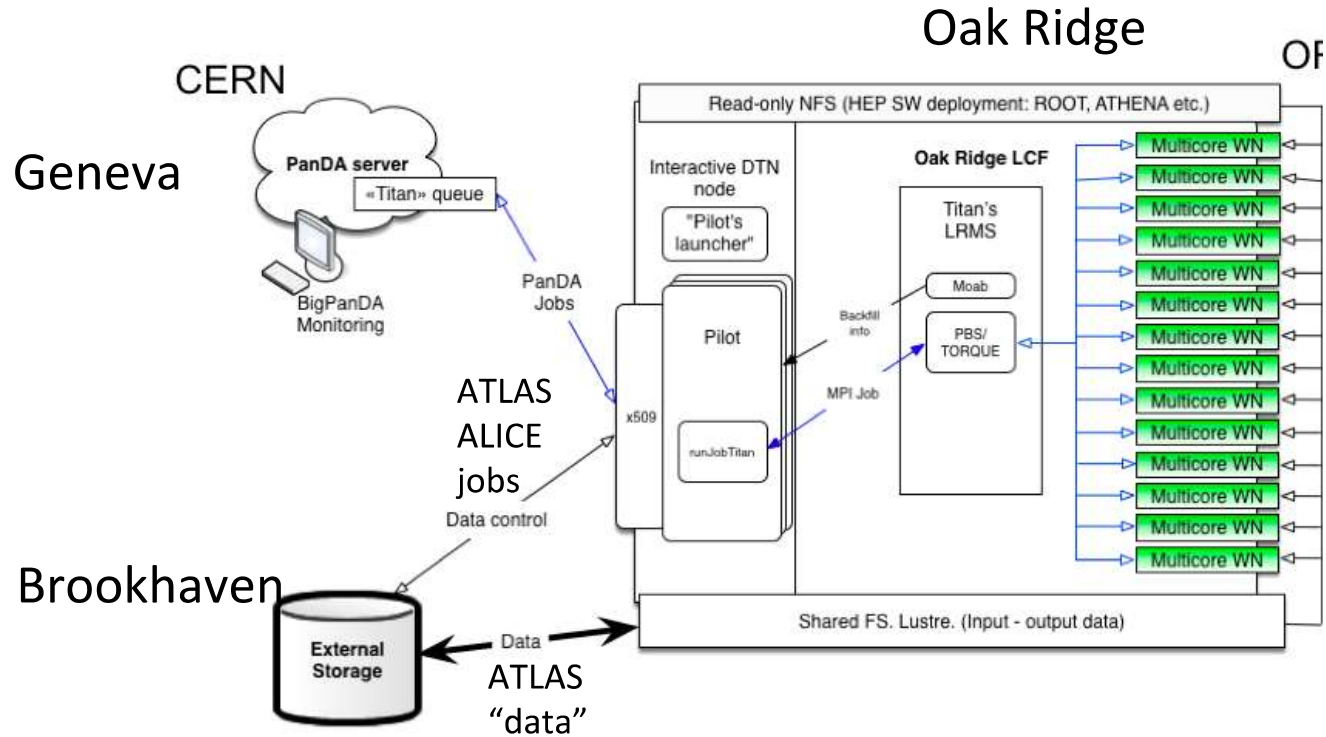
With shorter workload duration, backfill efficiency will increase.

Run time for finished jobs



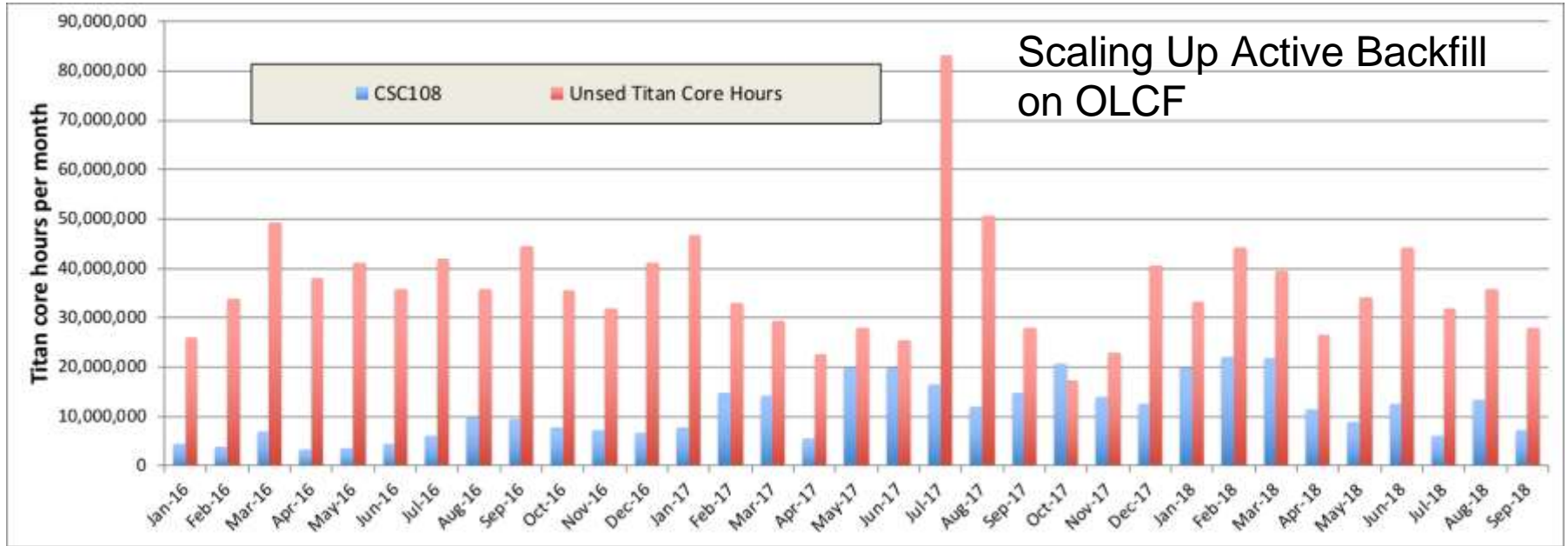
Task 9235668 Titan

OLCF Titan Integration with ATLAS Computing



D. Oleynik, S. Panitkin, M. Turilli, A. Angius, S. Oral, K. De, A. Klimentov, J. C. Wells and S. Jha, "High-Throughput Computing on High-Performance Platforms: A Case Study", IEEE e-Science (2017) available as: <https://arxiv.org/abs/1704.00978>

Concrete Results. New modes for use HPC for Data Intensive Sciences.



Consumed 370 Million Titan core hours from Jan. 2016 to Sep. 2018

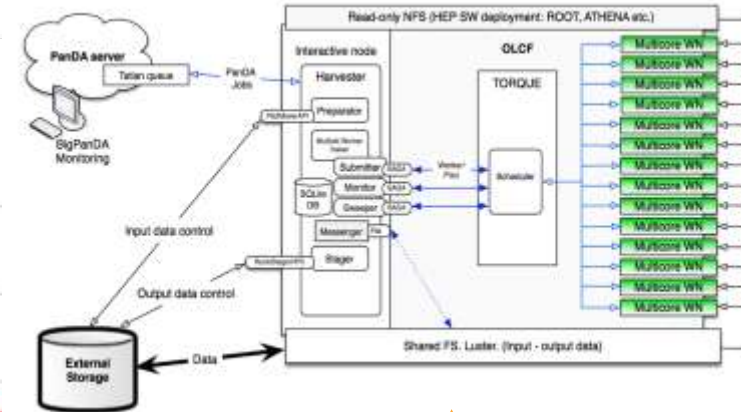
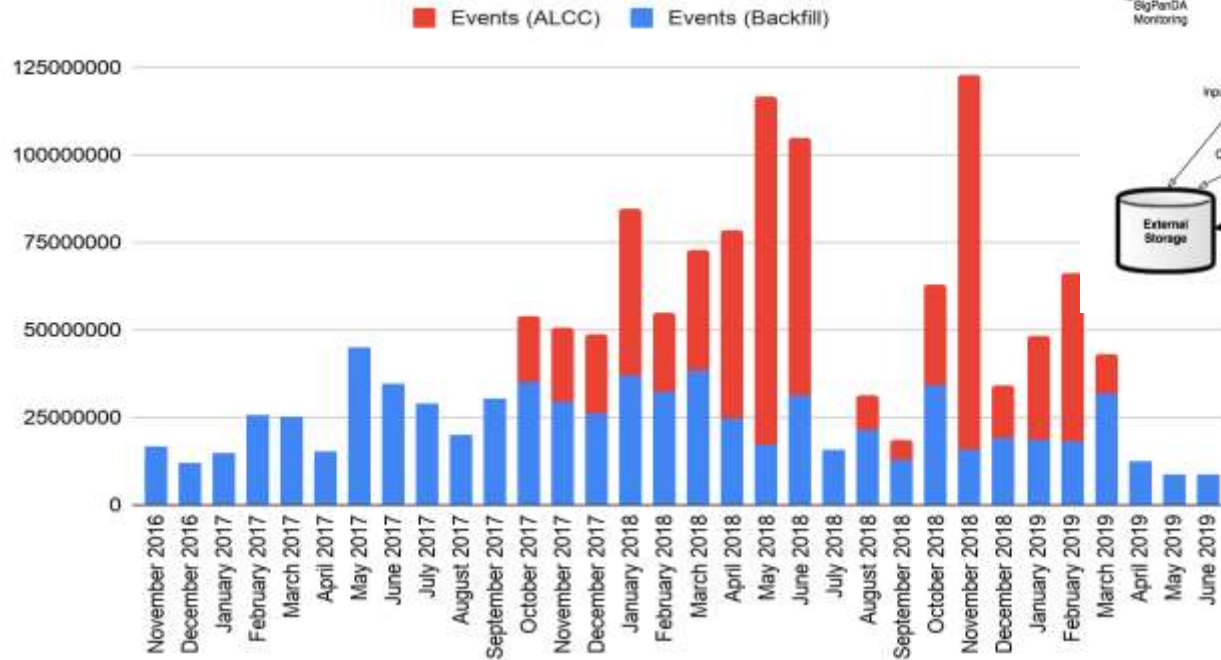
- This is 2.8 percent of total available time on Titan over this period

Remaining used backfill slots are often too short or too small for assigned ATLAS payloads

ATLAS production at OLCF (Nov. 2016 - Jun. 2019)

1,406,426,900 Events produced, 755,575,900 - "Backfill", 650,851,000 - "ALCC"

Events (Backfill) и Events (ALCC)



↑
New architecture for BigPanDA (Harvester) was developed to increase ALCC utilization

BigPanDA Workflow Management on OLCF for High Energy and Nuclear Physics and for Future Extreme Scale Scientific Application. BigPanDA project (2012-2019) a DOE ASCR and HEP funded project since 2012; a collaboration between BNL, UTA, ORNL and Rutgers University.



Quantum chromodynamics (QCD) is the component of the Standard Model of elementary particle physics that governs the strong interactions. It describes how quarks and gluons, the fundamental entities of strongly interacting matter, are bound together to form strongly interacting particles, such as protons and neutrons, and it determines how these particles in turn interact to form atomic nuclei.



The goal of the nEDM experiment at the Fundamental Neutron Physics ~~Beamline~~ at the Spallation Neutron Source (ORNL) is to further improve the precision measurement of neutron properties by a factor of 100 to search for violations of fundamental symmetries and to make critical tests of the validity of the Standard Model of electroweak interactions



The goal of the Large Synoptic Survey Telescope project is to conduct a 10-year survey of the sky that will address some of the most pressing questions about the structure and evolution of the universe and the objects in it:

- * Understanding Dark Matter and Dark Energy
- * Hazardous Asteroids and the Remote Solar System
- * The Transient Optical Sky
- * The Formation and Structure of the Milky Way



Molecular Dynamics: simulations of enzyme catalysis, conformational change, and ligand binding/release in collaboration with research group from University of Texas at [Arlington](#).



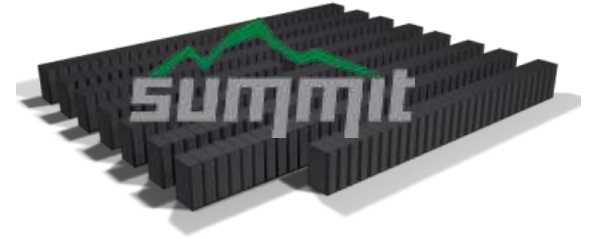
In collaboration with Center for Bioenergy Innovation at ORNL, the PanDA based workflow for epistasis research was established. Epistasis is the phenomenon where the effect of one gene is dependent on the presence of one or more modifier genes.



IceCube collaborators address several big questions in physics, like the nature of dark matter and the properties of the neutrino itself. IceCube also observes cosmic rays that interact with the Earth's atmosphere, which have revealed fascinating structures that are not presently understood.

Backfill on Summit

- First look at backfill opportunities on Summit
- Data collected by querying LCF scheduler on Summit for a number of days from May to July 2019
- Data show significant opportunities for backfill on Summit that may be utilized to run ATLAS workloads using approach similar to the one successfully used by BigPanDA project on Titan
 - Dynamic job size and duration shaping, based on real-time backfill information
- First statistical analysis of the data showed positive autocorrelations over large time intervals
- We are working on application of statistical (ARIMA, etc) and ML (LSTM, etc) models to the data, with the goal to develop smart, data driven algorithm for efficient backfill jobs submission



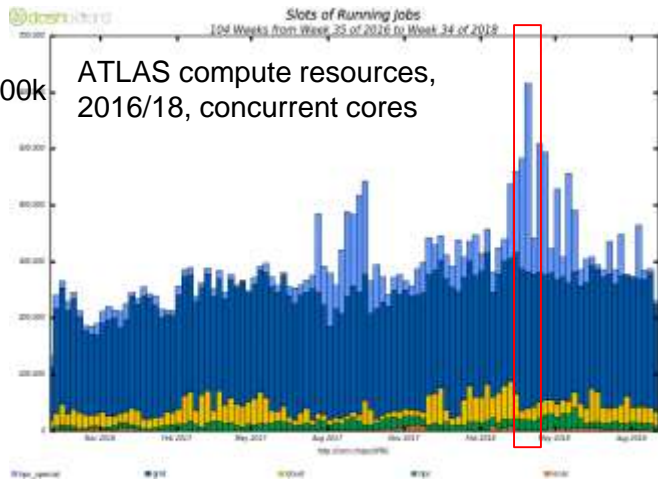
Summit: 200 PF
Accelerated Computing
5–10× Titan Performance

Grid-HPC integration. Concrete Results

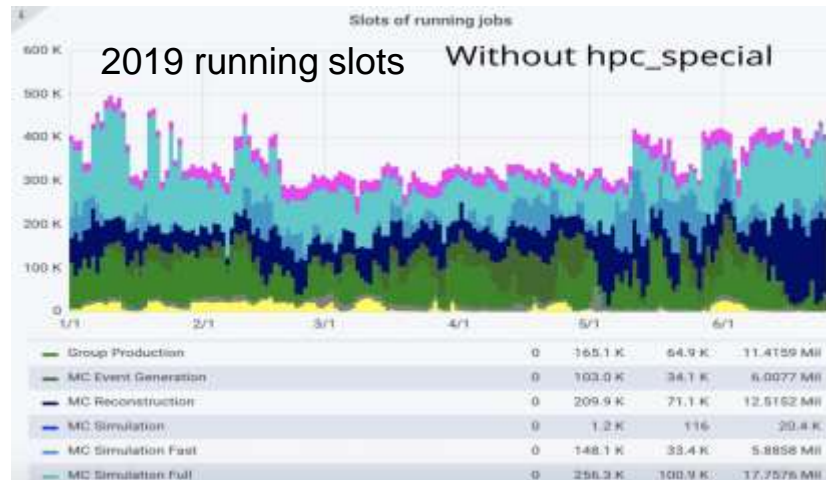
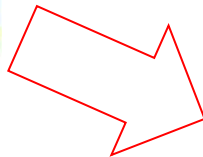
HPCs architecture, configuration and policies posed several challenges to the deployment of ATLAS distributed software components.

- The default model of PanDA pilot on the Grid was unfeasible for HPC/LCF
- **Harvester**, a new interface, common across resource types, between resource and workload manager was developed, as a result ALCC utilization was increased.
- **ARC software**, our backbone for HPC integration in Europe. Many EU HPCs are integrated via ARC software
- ATLAS Workload Management and Distributed Management Systems (**PanDA** and **Rucio**) have successfully coped with increased workload and traffic after HPC integration
- ATLAS **software releases** have been successfully built on HPCs (Titan and Summit)
- **Event Service**. High HPC utilization via fine grained workflows

HPC and Grid usage in ATLAS



A long history but a new era since 2018 : very large HPCs



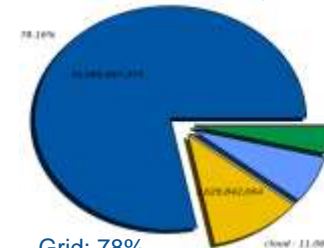
Light blue: “special” HPCs, where special means big, difficult to use, US DOE
 Dark blue: the grid
 Yellow: cloud resources including (dominantly) HLT
 Green: “regular” HPCs, meaning easier to use, European or US NSF

Zoom showing full size of scaling peak: 1.2M concurrent cores.

Our workload management system is highly scalable!



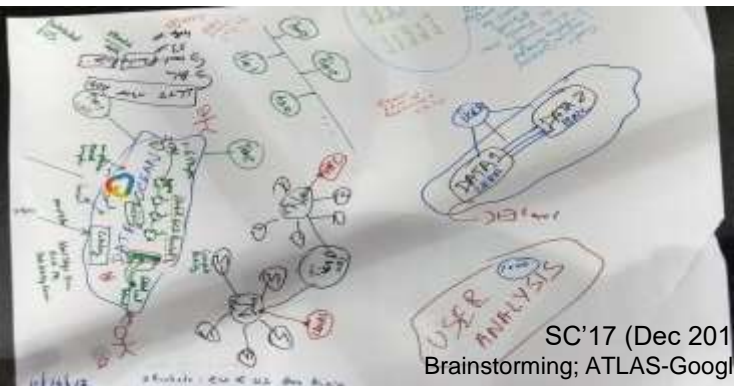
CPU HS06 shares, last year



Grid: 78%
 Cloud, HLT: 11%
 HPC special: 7%
 HPC regular: 4%

HEP-Google R&D. Motivation

- IT landscape has changed dramatically since end of XX century
- Technology sector is recognized as world leaders
 - Amazon, Google, Microsoft, Oracle, ... - already play significant role in worldwide scientific computing
- HENP data intensive computing challenges are (and have been) at the cutting edge of technology development
- **Foster partnerships with industries in research and development – and not just as late stage product adopters**
- The huge challenges at the HL-LHC have spurred new efforts in ATLAS to collaborate with technology partners
- **We are starting a new front in LHC R&D, with companies willing to invest in open source solutions**



SC'17 (Dec 2017)
Brainstorming; ATLAS-Google

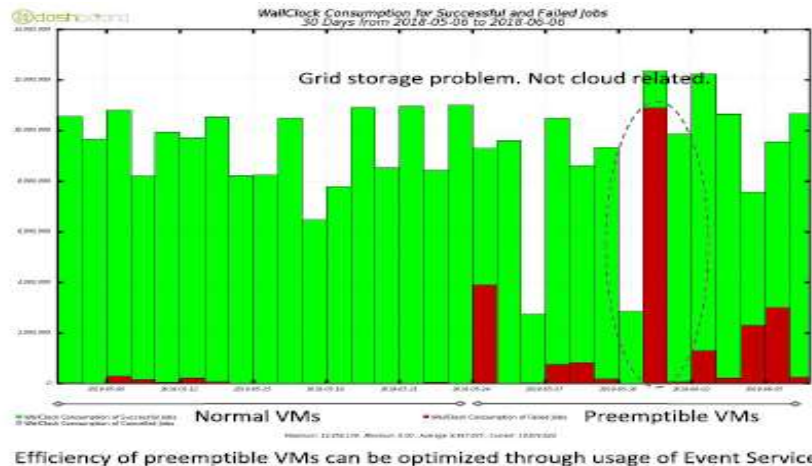
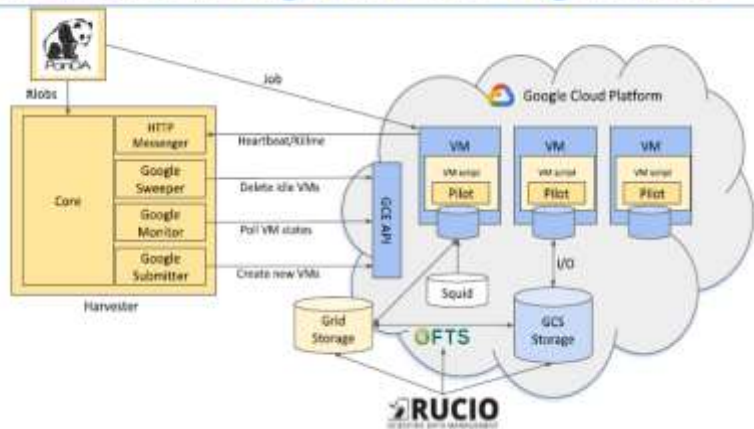
| | |
|----------------|---|
| Track 1 | Data Management across Hot/Cold storage |
| Track 2 | Machine learning and quantum computing |
| Track 3 | Optimized I/O and data formats |
| Track 4 | Worldwide distributed analysis |
| Track 5 | Elastic computing for WLCG facilities |

2018 PoC and the first realistic demo, 2019 5 Research tracks,

HEP-Google R&D. Accomplishments

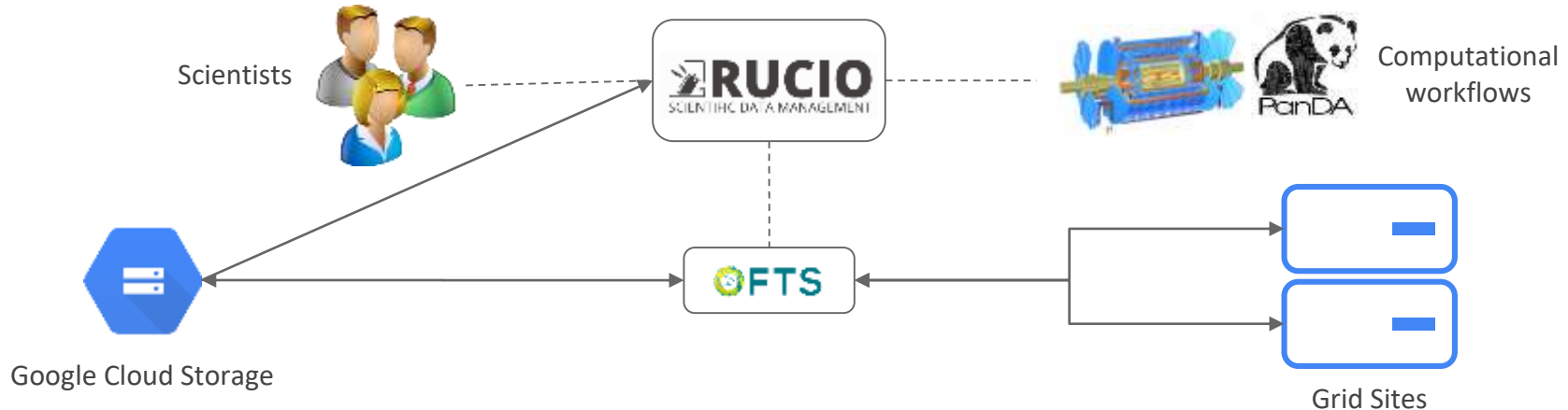
- ❖ ATLAS keeps multiple (expensive) copies of data for worldwide distributed analysis - R&D to use Google Storage
- ❖ Proof of concept project focused on analysis usage
 - ❖ Using Google storage transparently from ATLAS PanDA
 - ❖ Tested operating a 120 core Google cluster as PanDA resource
 - ❖ Successful with CPU at Google and data at CERN
 - ❖ Testing access of Google storage from ATLAS Tier 1 & Tier 2 sites

Job submission through Harvester edge service

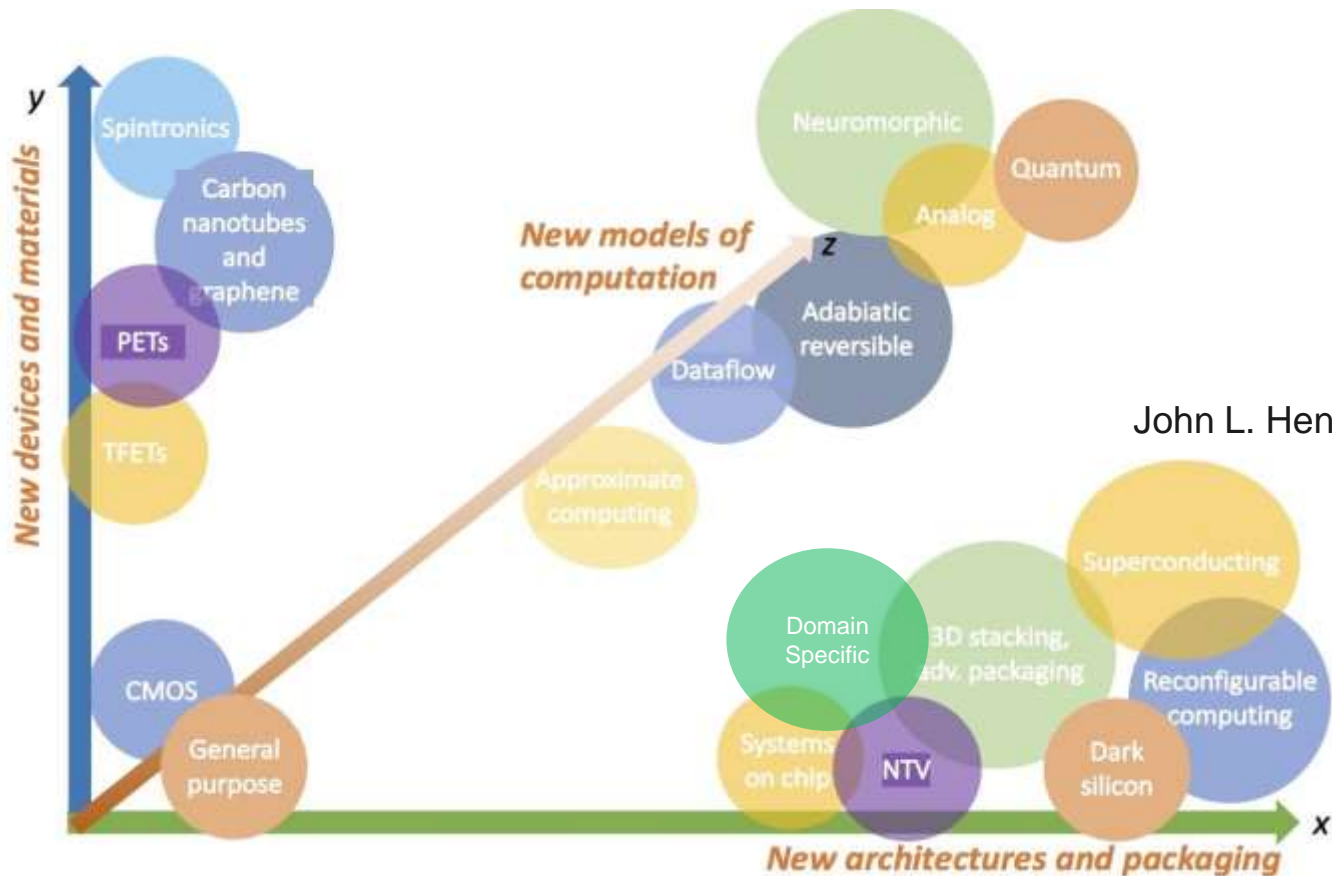


Data Ocean R&D. Getting data into Google Cloud Storage

- The ATLAS Data Management system *Rucio* orchestrates all experiment transfers
 - S3 used in the first iteration, since support is already available from both sides
 - Tests successful, however not usable for client-based access (key distribution, server-side signing)
 - Parallel third-party copy is rate-limited to 100MB/sec because we were not using the native GCS API
- Decision to move to GCP-native client-side signed URLs



A New Golden Age for Computing Architecture ?



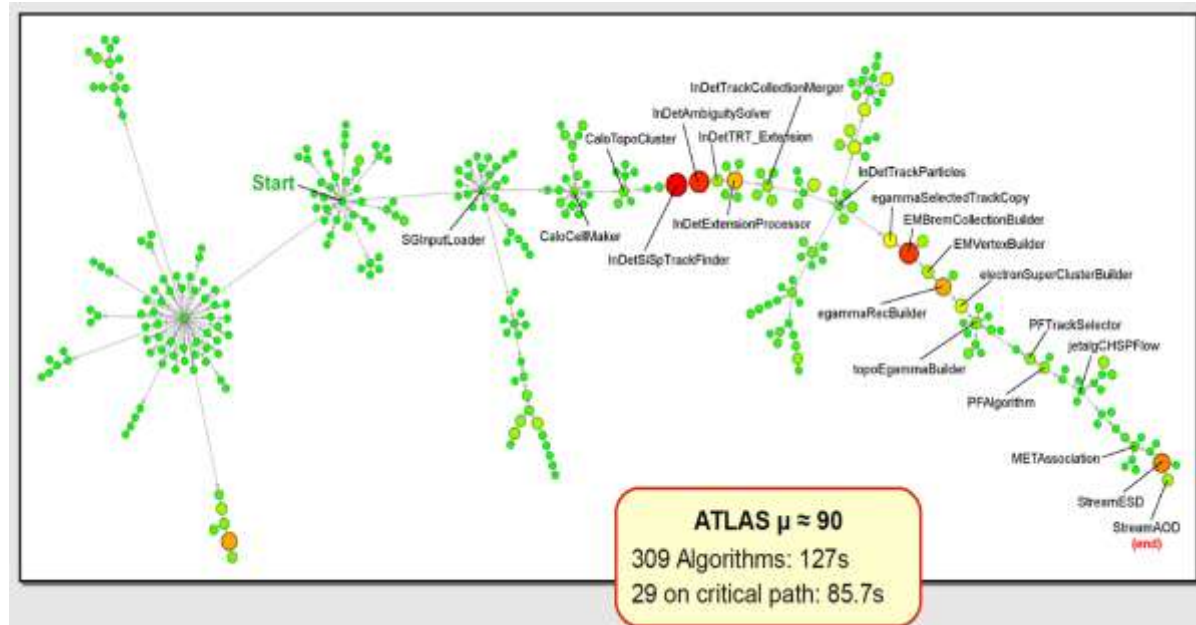
John L. Hennessy, David A. Patterson

Communications of the
ACM, February 2019, Vol.
62 No. 2, Pages 48-60

LHC Software. Role of Accelerators

Joint LHC experiments study of heterogeneous architectures

- Conclusion: accelerators help with throughput even when run at (reasonably) low efficiency
 - As long as one can run on O(10-100) events in parallel
 - Multi-threading needed for that



Side Evaluation of GPU in HLT farm completed in 2016:
GPU offloading service and test algorithms
Note: GPU for HLT cost-benefit analysis

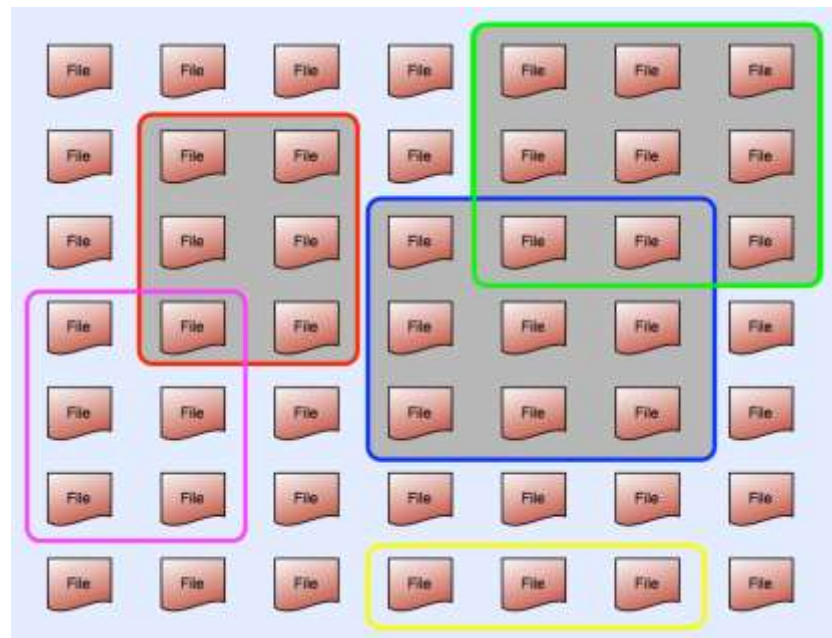
Conclusions

- LHC experiments are data and compute intensive
- Compute is done on distributed heterogeneous computing resources: the WLCG, HPC, ...
- Resources are managed centrally: data and workload management
 - We've learned many important concepts about both, many unknown unknowns problem retired/solved.
Known unknowns still remain
- LHC needs keep increasing and resources are spare! We need to be creative and extend the WLCG to any opportunistic resources we can get: Supercomputers, Cloud, Volunteer computing
- Commissioning prior to Run 4 in ~2026 is only **6 years away!**
- HEP's **compute-limited science** has driven us to (start to) develop early some of the capabilities we'll need for resource-constrained HL-LHC
- **CPU is probably on track** as long as we keep smart and keep working
 - PanDA deployment on supercomputers shows the potential of distributed high-throughput computing to be integrated with high-performance computing infrastructure. HPC themes in US and EU are very similar
- **Storage is the greater danger**
- **The Network** is our enabling foundation and will remain so
- **New architectures** is one of the most challenging topics for HEP
- Addressing the storage problem in the context of global computing on a high powered network foundation is in part **a big workload and data management challenge**

Backup slides

Rucio data concepts

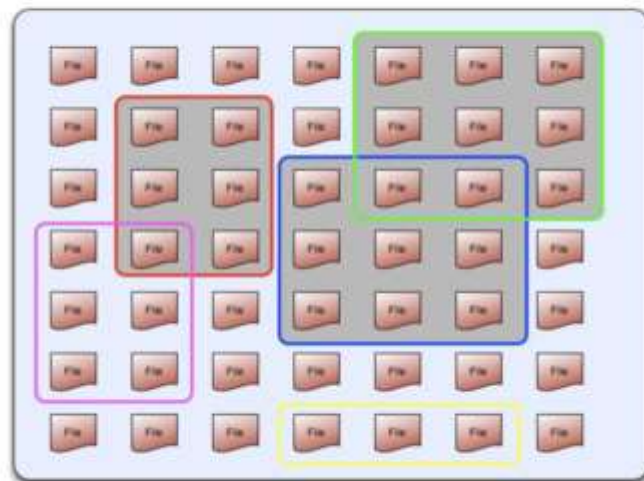
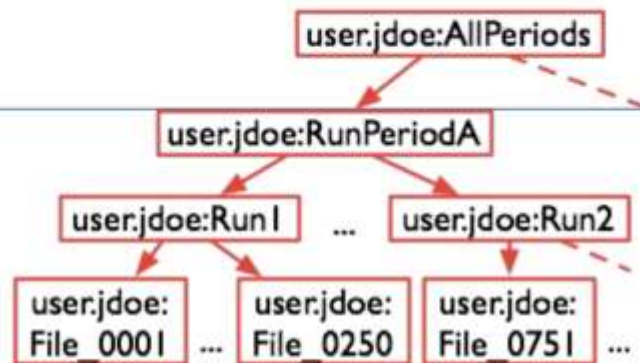
- Events: collisions
- Files: Collections of events (e.g. C++ objects)
- Datasets: logical grouping of files
 - Units of replication

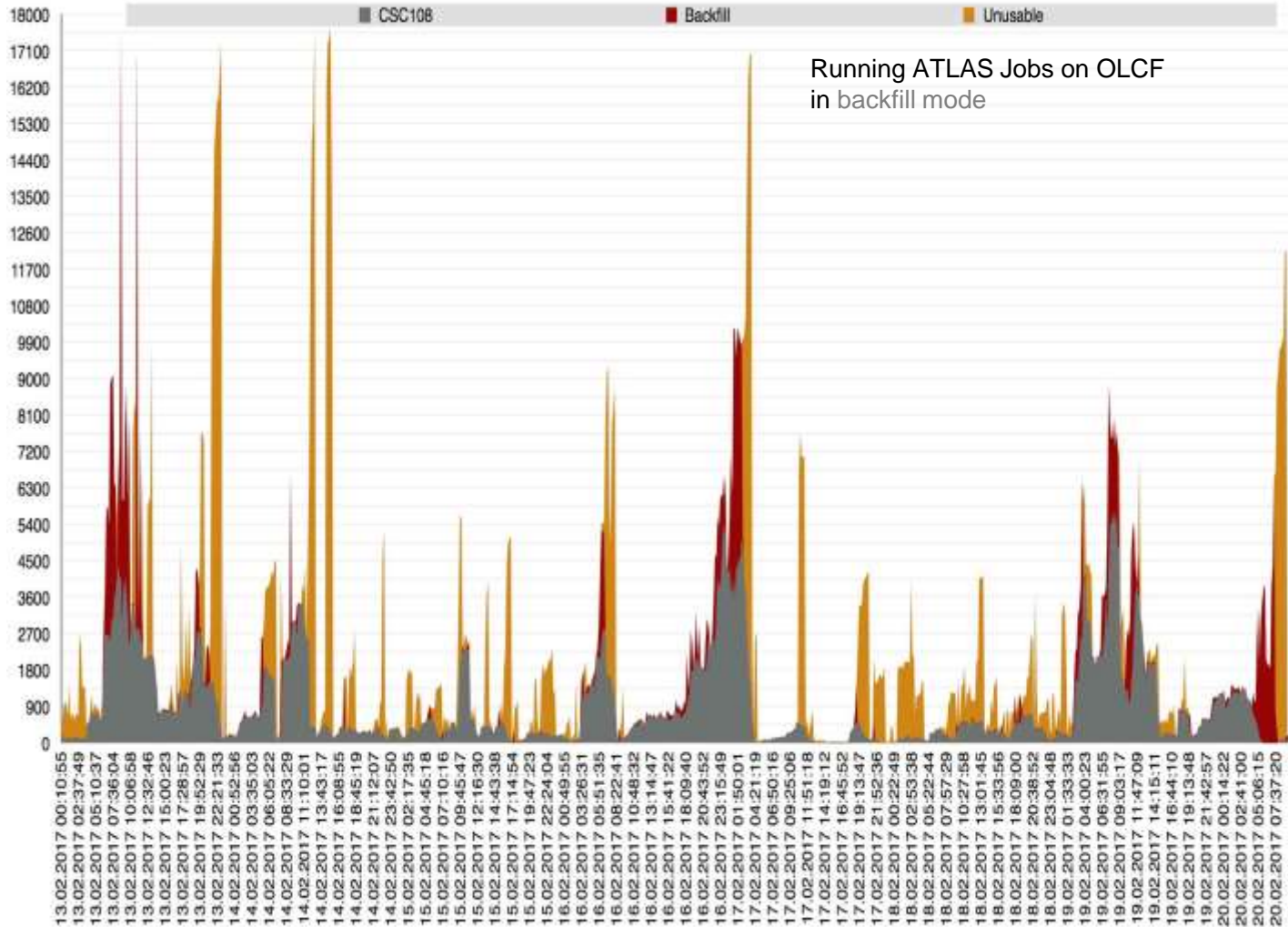


Data Hierarchy

- At the heart of everything is a file*
- Files are grouped into datasets
- Datasets are grouped into containers
 - Datasets only hold files
- Containers are grouped into containers
 - containers only hold datasets or containers
- Collections can be organised freely
 - Files can be in multiple datasets
 - datasets can be in multiple containers
 - containers can be in multiple containers

* sub-file support being explored





Acknowledgements

- Thanks to F.Barreiro Megino, D.Benjamin, I.Bird, K.Bhatia, R.Brun, P.Calafiura, S.Campana, K.De, A.Filipcic, M.Lassnig, T.Maeno, D.Oleynik, S.Panitkin, T.Schulthess, M.Schulz, P.Svirin, J.Wells, T.Wenaus for slides and materials. Thanks to members of the BigPanDA and PanDA projects, and colleagues from the ATLAS experiment at the LHC, CERN IT and WLCG colleagues