

# ATLAS Distributed Computing Evolution

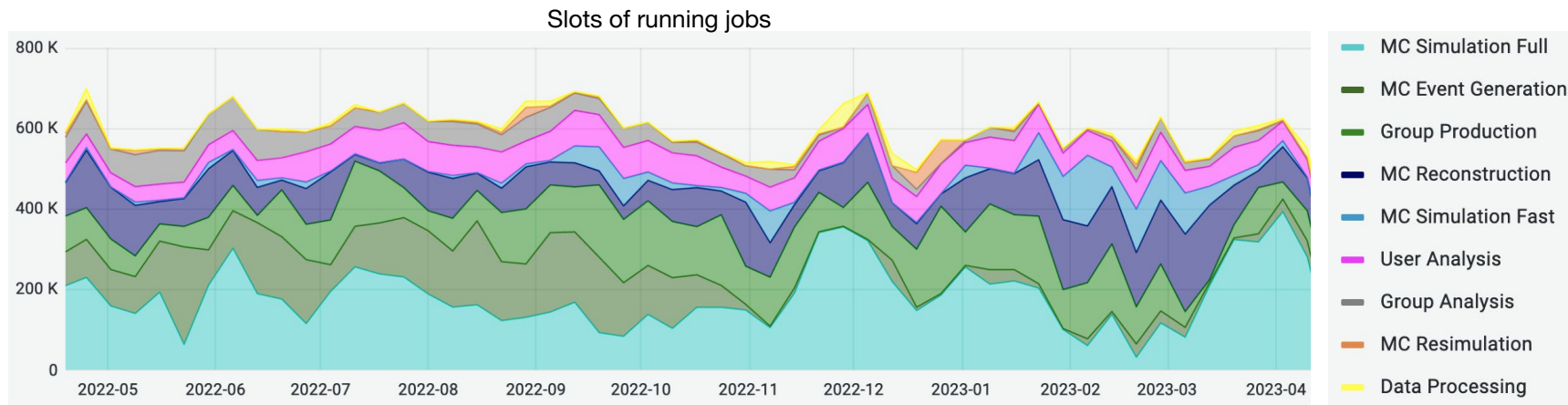
*Developments and demonstrators towards HL-LHC*

David Cameron, Mario Lassnig, David South  
on behalf of the ATLAS Computing Activity

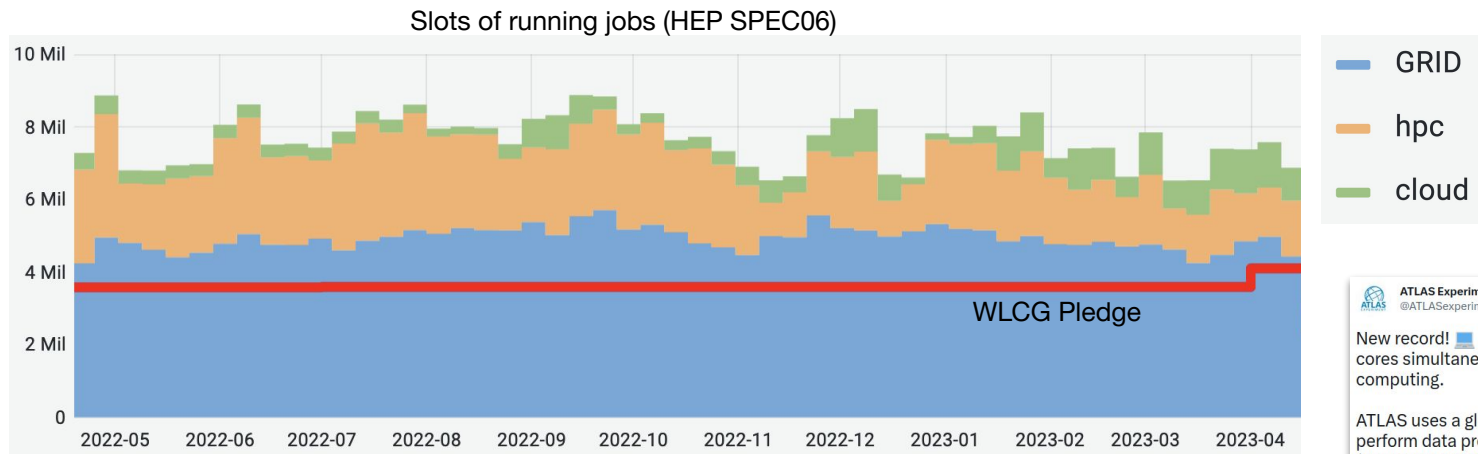


26TH INTERNATIONAL CONFERENCE ON COMPUTING IN HIGH ENERGY & NUCLEAR PHYSICS  
Norfolk, Virginia, USA      May 8-12 2023

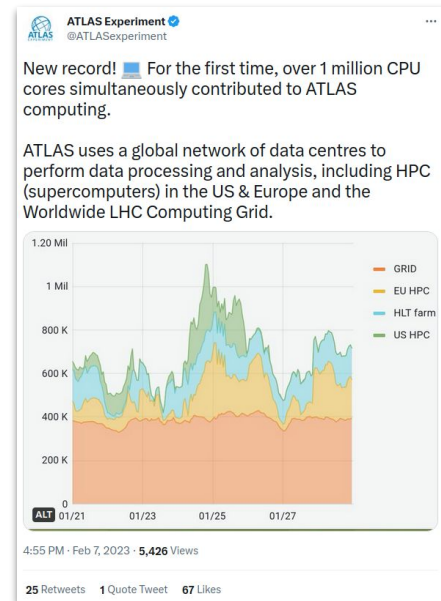




- Steady 600k+ running job slots, 24 hours a day, 365 days a year
  - Variety of job types, depending on the current focus of ATLAS activities
  - Mostly running as 8 cores per job, 2GB/core

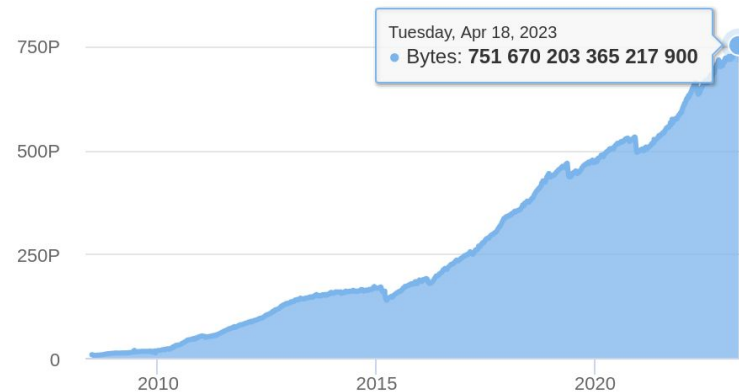


- Steady 600k+ running job slots, 24 hours a day, 365 days a year
  - Variety of job types, depending on the current focus of ATLAS activities
  - Mostly running as 8 cores per job, 2GB/core
- As well as the grid, which is consistently over pledge, ATLAS is also using HPC and Cloud resources as well as the HLT farm
  - In February we managed for the first time to reach over 1 million concurrently running jobs

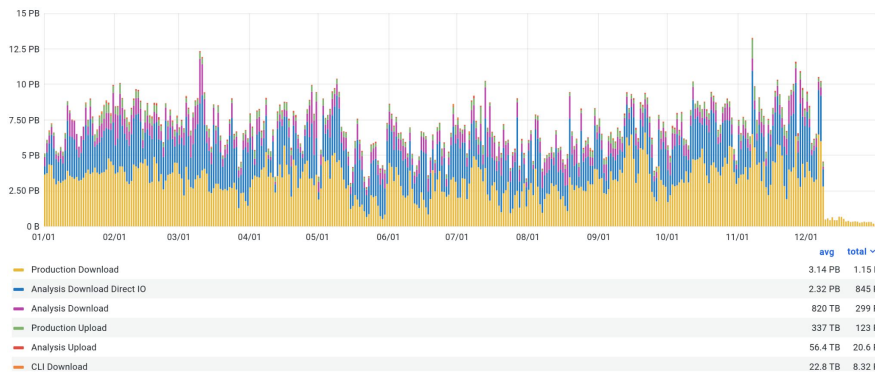


- A few numbers showing the scale of ATLAS data
  - 1B+ files, 750+ PB of data, 400+ Hz interaction
  - 120 data centres, 5 HPCs, 3 clouds, 1000+ users
  - 1.2 Exabytes/year transferred
  - 2.7 Exabytes/year uploaded & downloaded
- Expect an increase of at least one order of magnitude for the HL-LHC

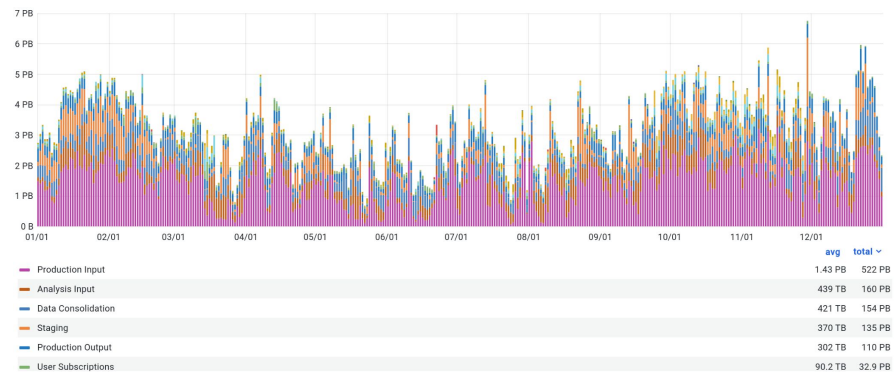
ATLAS data registered in Rucio



5+ PB/day data access for computation



2+ PB/day data transfers between storage



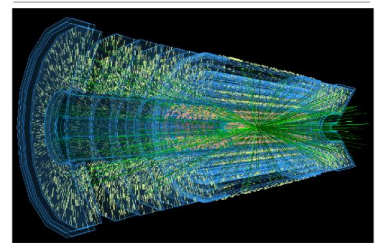
# Looking forward: The road to HL-LHC

- LHCC performs a series of reviews of the Software and Computing plans of the LHC experiments towards HL-LHC
  - The ATLAS HL-LHC Computing Conceptual Design Report was published in May 2020
- A follow up ATLAS Software and Computing HL-LHC Roadmap was published in March 2022 with clearly defined *milestones*

General Information			Key Software, Hardware Computing and Accelerator			Simulation			Data Access			Analysis			Development Concepts		
Item	Category	Task	Item	Category	Task	Item	Category	Task	Item	Category	Task	Item	Category	Task	Item	Category	Task
1.1	General	Finalize the HL-LHC Computing Conceptual Design Report (CDR) and its implementation plan	1.1.1	General	Finalize the HL-LHC Computing Conceptual Design Report (CDR) and its implementation plan	1.1.1	General	Finalize the HL-LHC Computing Conceptual Design Report (CDR) and its implementation plan	1.1.1	General	Finalize the HL-LHC Computing Conceptual Design Report (CDR) and its implementation plan	1.1.1	General	Finalize the HL-LHC Computing Conceptual Design Report (CDR) and its implementation plan	1.1.1	General	Finalize the HL-LHC Computing Conceptual Design Report (CDR) and its implementation plan



## ATLAS Software and Computing HL-LHC Roadmap



Reference:  
 Created: 1 October 2021  
 Last Modified: 22 February 2022  
 Prepared by: The ATLAS Collaboration  
 © 2022 CERN for the benefit of the ATLAS Collaboration.  
 Reproduction of this article or parts of it is allowed as specified in the CC-BY-4.0 license.

- An ATLAS HL-LHC Computing TDR is planned for 2024

Distributed Computing			
MID	DID	Description	Due
DC-1		Transition to tokens	Q4 2025
	1.1	Submission from Harvester to all HTCondor CEs with tokens	Q1 2022
	1.2	All users move from VOMS to IAM for X509	Q4 2022
	1.3	All job submission and data transfers use tokens	Q4 2025
DC-2		Storage evolution	Q4 2025
	2.1	No GridFTP transfers at any site	Q1 2022
	2.2	SRM-less access to tape	Q4 2025
	2.3	Recommended transition plan from DPM completed	Q4 2021
	2.4	Transition plan from all DPM sites	Q4 2022
DC-3	2.5	All sites moved away from DPM	Q2 2024
		Next operating system version	Q2 2024
	3.1	Ability to run on "future OS" on grid sites	Q4 2022
	3.2	Central services moved to "future OS"	Q4 2023
	3.3	(CentOS 7/8 EOL)	Q2 2024
DC-4		Network infrastructure ready for Run 4	Q4 2027
	4.1	Network challenge at 10% expected rate	Q4 2021
	4.2	Network challenge at 30% expected rate	Q4 2023
	4.3	Network challenge at 60% expected rate	Q4 2025
	4.4	Network challenge at 100% expected rate	Q4 2027
DC-5		Integrating next generation of HPCs	Q2 2023
	5.1	Integration of at least 2 EuroHPC sites	Q4 2022
	5.2	Integration of next generation US HPCs for production	Q2 2023
DC-6		Exploratory R&D on GPU-based workflows for next generation HPC	Q4 2023
DC-7		HL-LHC datasets replicas and versions management	Q2 2024
	7.1	Replicas and versions detailed accounting	Q4 2022
	7.2	DAOD replicas reduction	Q4 2023
	7.3	DAOD versions reduction	Q2 2024
DC-8		Data Carousel for storage optimization	Q4 2023
	8.1	Investigate with sites the cost of Tape infrastructure and the estimated cost in case of sensible increase of read/write throughput	Q4 2022
DC-9	8.2	Reduce the AOD on disk to 50% of the total AOD volume, using Data Carousel to orchestrate the stage from tape for DAOD production.	Q4 2023
		Disk management: secondary(cached) dataset	Q2 2023
	9.1	Evaluate the impact on job brokering and task duration if disk space for secondary data is reduced	Q2 2023
Maintenance & Operations		Conservative R&D	Aggressive R&D

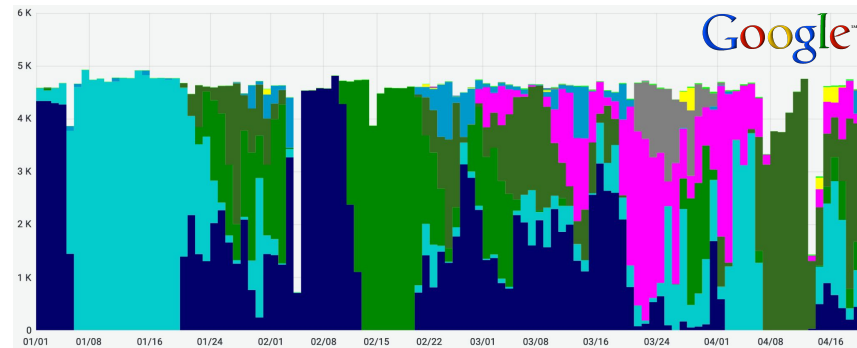
- Several HL-LHC milestones related to distributed computing
  - Not a static list! Regularly reviewed, updated and/or expanded
  - New milestones defined since the roadmap publication
  - Essentially two types: Maintenance and Operations, R&D
- ADC Maintenance and Operations Milestones
  - These are essential changes, needed “just to get by”
    - Tokens, evolution of storage technologies/access protocols, OS changes, network/data challenges, ..*
  - More details on these milestones in the [ATLAS report](#) at the WLCG Workshop in Lancaster last November and in other talks this week
- ADC R&D Milestones
  - Conservative R&D:** New developments achievable with current effort
  - Aggressive R&D:** New developments requiring extra effort
  - These are translated into “Demonstrators”, described in the following

- Expect an ever more heterogeneous environment as we approach the HL-LHC
  - Continue to evaluate how ATLAS resource, workflow and data management systems need to evolve to support massively parallel, distributed heterogeneous systems such as the next generation of HPCs
- R&D on GPU-based workflows remains so far at the exploratory level
  - GPUs are still not used in production by ATLAS, but are employed by a few dedicated analyses
    - *For more details on GPUs in analysis using Cloud resources, see the talk by [J. Sandesara](#)*
- Current ATLAS focus for large scale deployment of non-x86 is the use of **ARM processors**
  - Performance improvements seen over the last decade now make ARM a viable alternative
  - In particular, the energy efficiency can offer significant financial savings
- Significant progress in the last period - expect more to come this year
  - ATLAS Full Simulation now validated on ARM - other processes to follow - see talk by [J. Elmsheuser](#)
  - ATLAS is also looking at deployment of dedicated ARM resources at scale, outside of the pledge
  - Integration of ARM resources into distributed computing environment has no show stoppers
    - *If fully validated, can be fully integrated alongside Intel and AMD*

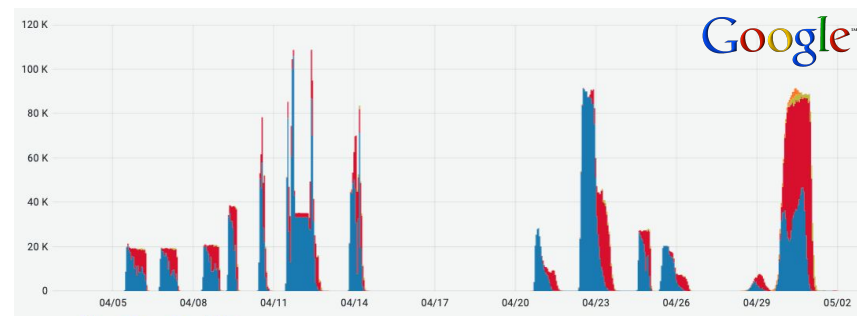


- ATLAS has a 1.5 year funded project with Google
  - Full integration of site into ADC activities and infrastructure
  - PanDA & Rucio developments are cloud-independent
  - Available for user analysis with Dask & Jupyter
- Total Cost of Ownership evaluation of Cloud resources
  - Understand the true cost breakdown for ATLAS
  - Limitations, optimal configuration, workflow adaptation
- Further investigation of interesting use cases
  - Cloud bursting: Dynamic/on-demand slot allocation
  - Network offloading: Use Google network for transfers
  - GPUs: On-demand GPU hardware
  - Special analysis workflows: ML, fitting, special MCs, ..
- The project continues ATLAS' long established use of opportunistic and heterogeneous resources
  - For more details on the ATLAS Google Project - see talk by [F. Barreiro Megino](#)
  - Rucio integration with Cloud storage - talk by [M. Lassnig](#)

Running various w/flows on 5k job slots at the Google site



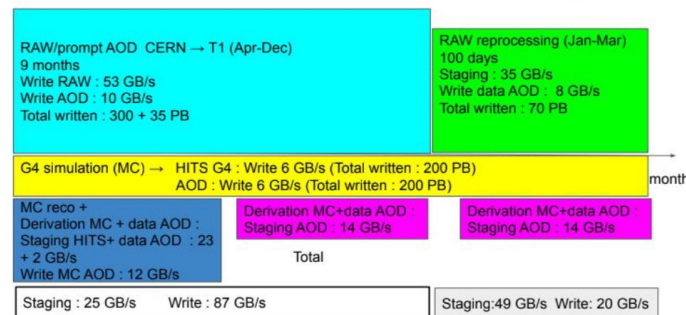
Bursting up to 100k slots for simulation and reconstruction jobs





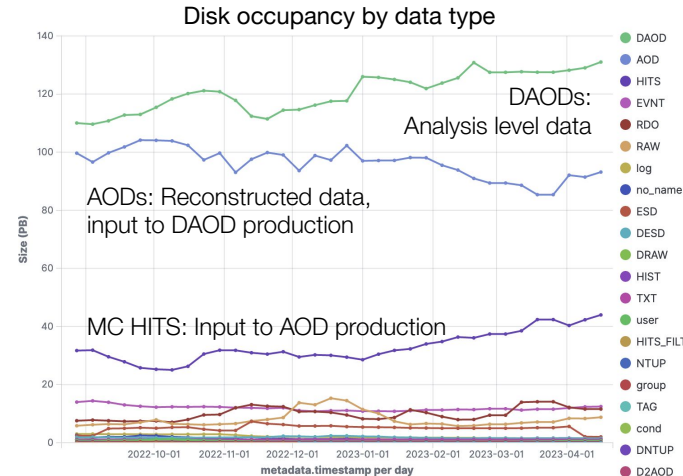
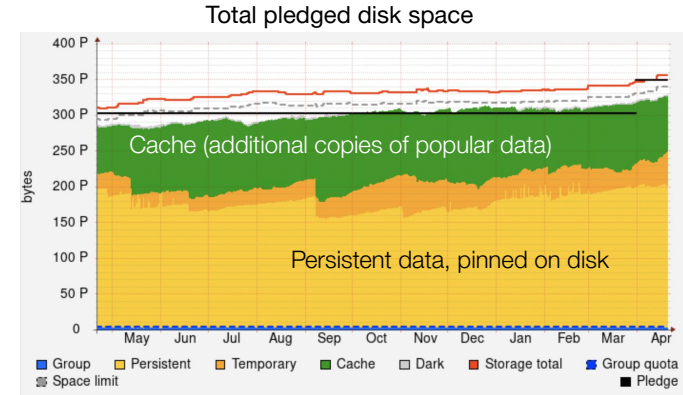
# Demonstrators: Smarter use of tape

- Tape has already moved from a pure archive storage towards more dynamic integration
  - Data Carousel mechanism successfully deployed in production by ATLAS for several years now
- Smart archiving: Core strategy R&D for our tape storage
  - Can we optimise file placement on tape for efficient retrieval?
  - Potential to improve Data Carousel throughput & latency
- First step: Definition of relevant metrics
  - Understand our data access patterns
  - Tape I/O metrics globally and individually
  - Consolidation of metadata required for efficient archival
- Second step: Functional tests to validate the process at a given T1 site
  - Propagation of metadata for site to colocate data through our stack (PanDA/Rucio/FTS/dCache/CTA)
  - Manual operations and monitoring by site experts of the underlying tape system
- Third step: Test real application in production
  - Creation of appropriately sized tasks with data samples in the 100 TB range
  - Assess effect of automatic colocation through tasks defined by production managers
  - *More on this subject in the talk by [X. Zhao](#) and the poster by [S. Misawa](#)*



Overall ATLAS T1 tape bandwidth estimates over the course of a typical year during Run 4

- Important to understand what exactly we have on disk at all times, to maximise use of limited space
  - Three categories: Pinned with no lifetime, temporarily occupying space (production input), cache containing additional copies of popular data
  - Several levers available to data mgmt team to control this
  - Recently have a much better cached-to-persistent ratio
    - *Evaluate the impact on job brokering and task duration if disk space for cached data is varied*
- AOD and HITS volume is rather stable, but DAOD grows from constant production, expect parallel data models during Run 3
  - Regular obsolescences, but large volumes of unused data kept on disk due to “lifetime model” exceptions
    - *Labour-intensive procedure for ADC and physics groups*
    - *Increase in time to publication leads to data being kept on disk*
- Ongoing R&D to find best solution, involving all facets
  - Lifetime model, data popularity, data placement, caching, data carousel
  - Volume, access patterns, user requests, available resources, ops load
  - Several options identified
    - *Delay: Keep datasets on disk/tape and delete after one year*
    - *Archive: Archive to tape, then delete from disk, recall as needed*
    - *Reproduce: Delete and reproduce when needed*

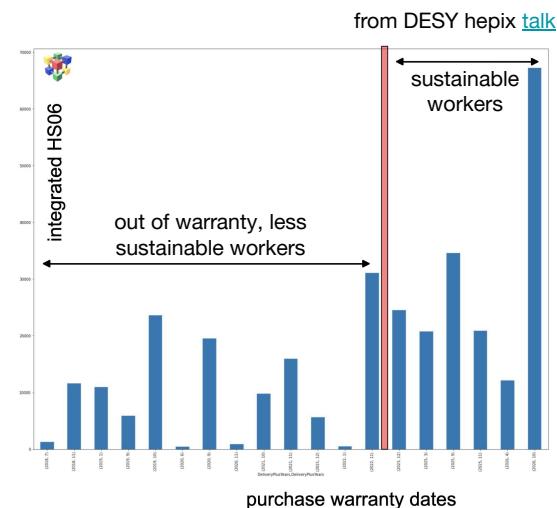
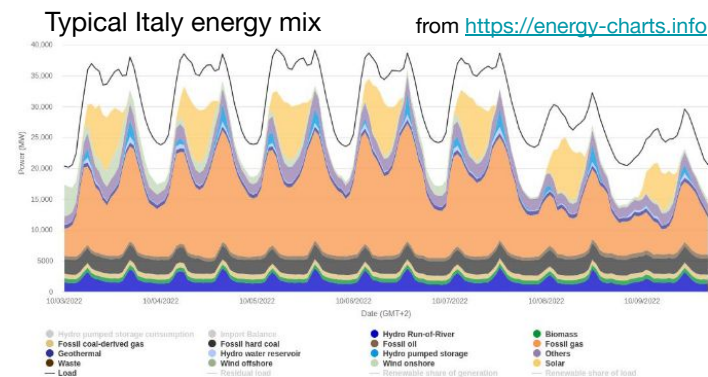


- Assess carbon footprint of our work and infrastructure

- Goal is transparency
- Show *estimate* of global gCO<sub>2</sub> usage
- We have no way to influence footprint currently
- Getting hard numbers is non-trivial
  - Variable gCO<sub>2</sub>/kWh per node, region and/or date
- Attacking carbon waste is another necessary angle

- Efforts to reduce power consumption

- Load shedding during peaks? Typically 2x day for 2 hours
  - We can't drain our nodes with long queues
  - Instead reduce CPU clock speed slightly? This could net 50-60% power reduction
- Reduce overall capacity during quieter periods, like Xmas
- Further ideas to partition clusters, reduce use of older nodes to opportunistic / on demand
  - Impact on pledges to be clarified
- As previously mentioned, investigate lower-energy platforms such as ARM
  - Expect this to influence resource planning going forwards



- The ATLAS computing infrastructure is in good shape
  - Made possible by a small but dedicated team following day-to-day operations
  - R&D for the HL-LHC is typically done by the same people
- Several HL-LHC milestones related to distributed computing have been defined and are regularly followed and updated
  - Maintenance and operation changes and new features will be imposed upon us
  - Many already during Run 3, assessed via the HL-LHC Data Challenges
- We are examining further ways to evolve our computing model, in terms of which resources we use and how we use them
  - Enabling the use of non-x86 architectures, cloud resources and the latest HPC
  - Expanding our investigations into the use of tape beyond archive
  - Improving our use and understanding of disk, and the data model employed
  - Evolution to a more sustainable model: Assess the carbon footprint, reduce the power consumption
- Look out for much more detail on the initiatives in other talks this week