




ATLAS Distributed Computing Evolution: Developments and Demonstrators Towards HL–LHC

David Cameron ¹, Mario Lassnig ², and David South ^{3,*}
on behalf of ATLAS Software and Computing

¹University of Oslo, P.b. 1048 Blindern, 0316 Oslo, Norway

²European Organisation for Nuclear Research (CERN), Geneva, Switzerland

³Deutsches Elektronen–Synchrotron (DESY), Hamburg, Germany

Abstract. The computing challenges at the HL–LHC require fundamental changes to the distributed computing models that have served experiments well throughout LHC. ATLAS planning for HL–LHC computing started back in 2020 with a Conceptual Design Report outlining various challenges to explore. This was followed in 2022 by a roadmap defining concrete milestones and associated effort required. Today, ATLAS is proceeding further with a set of "demonstrators" with focused R&D in specific topics described in the roadmap. The demonstrators cover areas such as optimised tape writing and access, data recreation on–demand and the use of commercial clouds.

1 ATLAS Distributed Computing today

The ATLAS [1] experiment at the LHC [2] uses a worldwide complex and distributed computing infrastructure comprising hundreds of thousands of computing cores and several hundred petabytes of storage, interconnected through high–speed networks. In addition to the WLCG [3] grid resources, which consistently surpass the pledged capacity, ATLAS effectively employs High–Performance Computing (HPC) and Cloud resources, alongside the High–Level Trigger (HLT) farm [4] in periods outside of data taking. Various workflows run in parallel, adapting to the changing focus of ATLAS activities, and the entire system is running non–stop, 24 hours a days, 7 days a week, 365 days a year. In February 2023, ATLAS achieved the significant milestone of reaching over 1 million concurrently running jobs for the first time.

Concerning the data themselves, over 1 billion files are managed by Rucio [5], with a total size of more than 750 PB and an interaction rate of over 400 Hz. These rates are dominated by pilot jobs running on the grid sites, which query the location of files in real time, with additional contributions from the various workflow and workload systems within the distributed computing infrastructure. Data movement is substantial, with an annual transfer of 1.2 EB and a combined upload and download of 2.7 EB. ATLAS accesses more than 5 PB of data daily for computation and daily data transfers between storage locations exceed 2 PB.

*e-mail: david.south@desy.de

© Copyright 2023 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

2 ATLAS HL–LHC milestones

The coming years are expected to bring a significant increase in scale, as well as new paradigms and approaches, reshaping the distributed computing landscape for the HL–LHC era. The LHC Experiments Committee (LHCC) performs a series of reviews of the Software and Computing plans of the LHC experiments, out of which ATLAS has now produced several publications. The ATLAS HL–LHC Computing Conceptual Design Report was published in May 2020 [6], outlining the challenges facing ATLAS Computing and Software in the HL–LHC Era and the general approaches to be taken to address them. This was followed in March 2022 by the ATLAS Software and Computing HL–LHC Roadmap [7], which provided more detailed information on the development work to be undertaken, as well as listing specific milestones and target dates. The full ATLAS HL–LHC Computing Technical Design Report is planned for 2025.

Several milestones related to distributed computing are defined within the ATLAS HL–LHC roadmap. These milestones are not static, are regularly reviewed and updated or expanded, and new milestones have been defined since publication. They can be broken down into essentially two types: Maintenance and Operations, and R&D.

Maintenance and Operations milestones are related to essential and necessary adaptations required due to distributed computing technology changes and evolution in areas such as storage technologies and access protocols, authentication (tokens), operating system changes, network capacity and so on. Further details on these milestones can be found in the roadmap publication [7].

The second category of milestones are associated to R&D projects, each attributed to either "conservative" or "aggressive" development scenarios. The success of these R&D projects is typically dependent on how much person power can be allocated to them and, as one would expect, those labelled as "aggressive" are expected to make the most impact but also require additional effort to see them fully to fruition. The remainder of these proceedings describes the main R&D milestones associated to ATLAS distributed computing and how these have been translated into so-called "demonstrators": a series of specific, short-term R&D projects and prototypes, which are used to appraise the progress and suitability of each set of milestones and deliverables.

3 Integration of non–x86 resources

ATLAS today utilises a progressively diverse set of computing resources. The resource, workflow, and data management systems continuously evolve to effectively integrate such diverse and massively parallel, distributed, and heterogeneous systems, like the latest generation of HPC platforms or Cloud infrastructures. Furthermore, these resources have more recently begun to diversify from the more traditional x86–CPU architectures the experiment has relied upon since its inception.

Exploratory research and development by ATLAS related to GPU–based workflows continues to evolve. Although GPUs have not yet been integrated into ATLAS production, they are employed by specific dedicated analyses. An example of the role of GPUs in analysis using Cloud resources was presented at this conference [8].

At this time, the main emphasis in the large–scale deployment of non–x86 architectures is centred around ARM processors. The performance gains witnessed over the past decade have established ARM as a viable and compelling alternative. Notably, its commendable energy efficiency holds the potential for significant financial savings [9].

Substantial progress has been made by ATLAS with ARM in recent times, and further advancements are on the horizon. Both simulation and reconstruction Monte Carlo work-

flows have now been validated in small scale tests on ARM [10], involving 1M simulated events and a total of 1.3M reconstructed events using a variety of physics processes. The next step is the deployment of dedicated ARM resources at larger scale, which is underway at the time of writing via a collaboration between ATLAS and the University of Glasgow using around 1750 ARM cores to produce a substantially larger test MC sample of 50M events. Importantly, the integration of ARM resources into ATLAS distributed computing environment has encountered no insurmountable obstacles and once successfully validated at scale, it is expected that ARM resources can seamlessly coexist alongside Intel and AMD architectures.

4 Cloud resources and the ATLAS Google Project

ATLAS is presently engaged in a collaborative project with Google, investigating the seamless integration of a cloud site into the broader distributed computing activities and infrastructure. The necessary adaptations the ATLAS workflow management software ProdSys [11] and PanDA [12], as well as those in the data management software Rucio [5], are done with cloud-agnostic designs, rendering them adaptable and independent of specific cloud providers. User analysis, facilitated through platforms like Dask [13] and Jupyter [14], is also investigated

In addition to the integration of a Google cloud site into ATLAS distributed computing, other utilisation cases are investigated as part of the project. Dynamic and on-demand allocation of a significant number of computing slots, so called "cloud-bursting", is employed to produce MC samples on a short timescale. The potential for on-demand GPU hardware provisioning is also explored. Exploiting special resources available at Google is investigated in the context of specialised analysis workflows including machine learning, fitting procedures, special Monte Carlo simulations, among others.

A Total Cost of Ownership evaluation is also to be included as part of the Google project, to understand the consequences of the adoption of cloud resources by ATLAS. The current project with Google is a continuation of ATLAS distributed computing engaging with heterogeneous and opportunistic resources. Further details and insights into the ATLAS Google Project were presented at this conference, concerning the project as a whole [15] and the integration of Rucio with Cloud storage [16].

5 Smarter use of tape

The evolution of tape storage within ATLAS has shifted from its conventional role as pure archive storage towards a more dynamic and integrated approach. A prime example of this transformation is the successful implementation of the Data Carousel [17] mechanism in ATLAS production operations, a strategy that has been effectively in practice for several years. At the heart of the current ATLAS R&D tape activities is the concept of smart archiving, which serves as a central focus for research and development in our tape storage strategy. This involves a thorough exploration of optimising the placement of files on tape to enhance the efficiency of retrieval. Such optimisation holds the potential to positively impact the data throughput and latency of the Data Carousel system.

The investigation into smart tape archiving involves three distinct steps. The initial phase involves defining relevant metrics, by forming a comprehensive understanding of our data access patterns, inclusive of both global and individual tape I/O metrics, as well as consolidation of the necessary metadata essential for streamlined archival processes. In a second step, functional tests are conducted to validate the process at several WLCG Tier-1 sites. This entails the propagation of the correct metadata so that the site can co-locate data through

the full operational stack, overseen by site experts of the underlying tape system. The final step entails real-world application testing within a production environment. This involves the creation of tasks with appropriately sized data samples, typically in the range of 100 TB, and the subsequent evaluation of the applied automatic co-location. Further details of the ongoing investigations of tape data recall performance [18] and smart data archiving [19] were presented at this conference.

6 Optimising use of disk space

A crucial part of optimising the utilisation of our finite disk space resources is to understand what types of data are stored on disk and under which conditions. This understanding categorises the data into three distinct types: Data that is pinned without a predefined lifespan, data temporarily occupying space as production inputs, and cached data containing additional replicas of frequently accessed information. Multiple mechanisms are at the disposal of the ATLAS data management team to regulate these categories, effectively steering resource allocation.

Recent progress in data management has led to a notably improved cached-to-persistent data ratio, and an investigation is planned into how the size variation of this cache affects critical aspects like job brokering and task duration. While the volume of both reconstructed (AOD) and detector simulation (HITS) data exhibit relative stability, the volume of analysis level derived data, so-called DAODs, experiences consistent growth due to constant ongoing production. Regular deletion campaigns attempt to keep this data volume in check, but substantial volumes of often unused data remain on disk due to exception requests. This situation necessitates labour-intensive procedures for both the distributed computing and physics groups, with extended publication timelines contributing to data retention on disk.

An effort to find the most effective solution to this issue has begun, which examines key elements such as the data lifetime, data popularity, data placement strategies, caching mechanisms, and the Data Carousel framework. The process is shaped by factors including data volume, access patterns, user requisites, resource availability, operational load, and so on. Three options emerge: Keep the data sets on disk or tape and schedule for deletion after a year; Archive the data onto tape, subsequently remove from disk, and recall as needed; Delete the data and reproduce on demand if and when it is needed in the future. The latter two options are actively investigated by ATLAS as part of a disk space optimisation demonstrator [19].

7 Towards a more sustainable distributed computing model

The comprehensive evaluation of our carbon footprint covers both ATLAS operational practices and the underlying infrastructure. The overarching objective here is transparency: The aim to offer a clear estimate of the global gCO₂ (carbon dioxide equivalent) consumption of the experiment. It is important to recognise that at the present time our capacity to directly influence our carbon footprint is limited. The intricacies of obtaining precise data are multifaceted, owing to the dynamic nature of gCO₂/kWh metrics which vary across nodes, regions, and dates [20]. It is recognised that whilst some level of carbon footprint may be unavoidable, carbon waste, arising from wasted CPU wall-time through failed or misconfigured jobs, is an area that can be tackled.

A multifaceted approach is applied towards reducing power consumption. One viable strategy involves load shedding during peak demand periods: Typically implemented twice daily for two hour intervals, this approach seeks to alleviate the load at peak. However, it is essential to implement this without draining nodes with substantial queues or killing jobs which have already run for a substantial amount of time.

A parallel effort focuses on reducing the CPU clock speed, potentially resulting in a significant 50-60% reduction in power utilisation. Another option is to exploit periods of reduced activity, such as holiday seasons, but reducing the overall capacity, aligning resource allocation with demand. One approach to this is to partition clusters, minimising the reliance on older nodes, particularly for tasks that are opportunistic or on-demand. The potential implications of all of these strategies on site commitments and pledges needs to be considered at some point.

The ATLAS sustainability drive also extends to the exploration of lower-energy platforms, with a specific emphasis on ARM, as detailed earlier. This direction carries substantial weight, as it may significantly affect resource planning in the future. In summary, ATLAS is actively taking steps towards environmentally-conscious computing by examining the carbon footprint of our infrastructure and taking proactive measures towards power reduction, in order to establish a robust foundation for a more sustainable and responsible distributed computing framework.

8 Summary

The ATLAS distributed computing infrastructure continues to deliver. It owes its success to a dedicated team that diligently manages day-to-day operations and simultaneously spearheads R&D for the HL-LHC phase. The HL-LHC era is poised to introduce operational changes, new features, and maintenance alterations, all evaluated through the lens of the HL-LHC Data Challenges. A series of R&D milestones pertinent to distributed computing have been outlined, and have been translated into a corresponding set of HL-LHC demonstrators.

These demonstrators include optimising resource utilisation, embracing non-x86 architectures, and the further integration HPC and cloud resources. The management of disk storage and the underlying data model is examined closely, to optimise how this finite resource is employed. ATLAS also continues to explore the potential of tape beyond its conventional archival role, now examining how data writing may be optimised. A more sustainable computing model is also investigated, including an assessment of the carbon footprint and additional efforts to curtail power consumption.

References

- [1] ATLAS Collaboration, The ATLAS Experiment at the CERN Large Hadron Collider, *JINST* **3** (2008) S08003, doi:10.1088/1748-0221/3/08/S08003
- [2] LHC – The Large Hadron Collider, <http://lhc.web.cern.ch/lhc/>, accessed 2023-08-17
- [3] J. Shiers, World LHC Computing Grid, *Computer Physics Communications*, **177** (2007) 219, doi:10.1016/j.cpc.2007.02.021
- [4] ATLAS Collaboration, Performance of the ATLAS trigger system in 2015, *Eur. Phys. J. C* **77** (2017) 5, 317, doi:10.1140/epjc/s10052-017-4852-3
- [5] M. Barisits et al, Rucio: Scientific Data Management, *Computing and Software for Big Science* (2019) 3:11, doi:10.1007/s41781-019-0026-3
- [6] ATLAS Collaboration, ATLAS HL-LHC Computing Conceptual Design Report, CERN-LHCC-2020-015; LHCC-G-178
- [7] ATLAS Collaboration, ATLAS Software and Computing HL-LHC Roadmap, CERN-LHCC-2022-005; LHCC-G-182
- [8] J. Sandesara et al., ATLAS Data Analysis using a Parallelised Workflow on Distributed Cloud-based Services with GPUs, *in these proceedings*
- [9] E. Simili et al., ARMing up for HEP, *in these proceedings*

- [10] J. Elmsheuser et al., The ATLAS Experiment Software on ARM, *in these proceedings*
- [11] F. H. Barreiro et al., The ATLAS Production System Evolution: New Data Processing and Analysis Paradigm for the LHC Run2 and High-Luminosity, *J. Phys. Conf. Ser.* **898** (2017) 052016, doi:10.1088/1742-6596/898/5/052016
- [12] F. H. Barreiro et al., PanDA for ATLAS Distributed Computing in the Next Decade, *J. Phys. Conf. Ser.* **898** (2017) 052002, doi:10.1088/1742-6596/898/5/052002
- [13] *Dask*, <https://www.dask.org>, accessed 2023-08-17
- [14] *JupyterHub*, <https://jupyter.org/hub>, accessed 2023-08-17
- [15] F. H. Barreiro et al., Accelerating Science: The usage of Commercial Clouds in ATLAS Distributed Computing, *in these proceedings*
- [16] M. Lassnig et al., Extending Rucio with Modern Cloud Storage Support: Experiences from ATLAS, SKA and ESCAPE, *in these proceedings*
- [17] M. Borodin et al., The ATLAS Data Carousel Project Status, *EPJ Web Conf.* **251** (2021) 02006, doi:10.1051/epjconf/202125102006
- [18] S. Misawa, Examining the Impact of Data Layout on Tape on Data Recall Performance for ATLAS, *in these proceedings*
- [19] X. Zhao et al., Updates to the ATLAS Data Carousel Project, *in these proceedings*
- [20] R. Walker, Sustainability in HEP Computing, *in these proceedings*