

## Distributed processing and analysis of ATLAS experimental data

---

**Dario Barberis<sup>1</sup>**

*Physics Department of the University of Genoa and INFN  
Via Dodecaneso 33, I-16146 Genova, Italy  
E-mail: Dario.Barberis@ge.infn.it*

**On behalf of the ATLAS Collaboration**

**Abstract:** The ATLAS experiment is taking data steadily since Autumn 2009, and collected so far over  $5 \text{ fb}^{-1}$  of data (several petabytes of raw and reconstructed data per year of data-taking). Data are calibrated, reconstructed, distributed and analysed at over 100 different sites using the World-wide LHC Computing Grid and the tools produced by the ATLAS Distributed Computing project. In addition to event data, ATLAS produces a wealth of information on detector status, luminosity, calibrations, alignments, and data processing conditions. This information is stored in relational databases, online and offline, and made transparently available to analysers of ATLAS data world-wide through an infrastructure consisting of distributed database replicas and web servers that exploit caching technologies. This paper reports on the experience of using this distributed computing infrastructure with real data and in real time, on the evolution of the computing model driven by this experience, and on the system performance during the first two years of operation.

*The 2011 Europhysics Conference on High Energy Physics-HEP2011  
Grenoble, Rhône-Alpes*

*July 21-27 2011*

---

<sup>1</sup> Speaker

## 1. Introduction

Distributed computing is essential for all LHC experiments and particularly for ATLAS, the largest multi-purpose experiment [1]. No institution within the collaboration can afford to fund and host the enormous computing infrastructure that is necessary to store and process all experimental data. The experience of the previous generation of HEP experiments, which used the LEP accelerator at CERN and the Tevatron accelerator at FNAL (Chicago) and were more than one order of magnitude smaller, led the LHC experiments to design upfront a distributed computing system, able to exploit optimally all available resources, independently of their geographical location. The Grid computing paradigm was adopted by the LHC community as the initial idea looked rather simple and elegant: each site provides common interfaces to its local batch system (the Computing Element, CE) and data storage system (the Storage Element, SE), and publishes its properties through a common Information System. A global authentication and authorisation framework guarantees the identity of the submitter of workload and his privileges (or lack thereof) on each site.

The Grid paradigm was modified for HEP from a compute-intensive system to a data-intensive one. Experimental data are precious, as the cost of building and operating each LHC experiment approaches one billion Swiss Francs, and data storage facilities (disks) are comparatively more expensive than data processing units (CPUs). Data storage also needs to be separated between archival storage (on tape) and online data (on disk). For reasons of robustness and also as a safeguard against data corruption, at least two copies of the same data must be kept on disk at different locations, and one on tape (two copies on tape for the "raw" data produced directly by the experiment).

Computing centres at large national HEP laboratories ("Tier-1" sites), ten in total for ATLAS, provide archival facilities on tape, several petabytes of disk space, and several thousands job processing slots. Smaller computing facilities, mostly placed at universities or their physics departments, provide disk space and processing facilities of differing size; there are over 70 such "Tier-2" sites for ATLAS. Local, batch and interactive, facilities ("Tier-3") are used for the final data analysis, usually consisting in preparing histograms and fitting functions from which final data to be published are extracted. The CERN laboratory is the source of all "raw" data and its computer centre ("Tier-0") holds a copy of all produced data and runs calibration, alignment and data reconstruction procedures in real time, before distributing the data to the Tier-1 sites. Selected data that are needed for specific physics analyses are then further distributed to Tier-2 sites.

## 2. Building blocks of the distributed computing infrastructure

Grid middleware includes all software components that are needed to provide remote and secure access to the computing resources. Several suites have been developed and deployed over the last ten years, all implementing server-client architectures [2]-[4]. The enormous amount of data generated by LHC (several petabytes per experiment per year) can be handled only by establishing hierarchical structures and cataloguing all the data [5]. Data files are

grouped into “datasets”, i.e. collections of all files containing statistically equivalent events in the same format and processing stage. As datasets can typically contain from 100 to 10000 files, in this way the cataloguing problem is reduced by 2 to 4 orders of magnitude. Each dataset is then created, replicated, moved or deleted as a single unit. Every operation on files or datasets must be registered in the central catalogues, in order to have at any point in time all information about data locations, access and popularity. Tools have been developed to analyse this information and automatically increase the number of replicas of the more popular datasets and decrease the number of replicas of least accessed datasets.

The PanDA jobs submission framework manages the large number of jobs (several hundred thousand jobs per day), interacting strongly with the data management tools, directing the jobs to the sites where they can run fastest, and collecting the outputs to the site indicated by the job owner [6]. "Production" jobs are submitted centrally to produce simulated events or reprocess real events with better calibrations or reconstruction code, when available. "Analysis" jobs are submitted by any ATLAS member who wishes to analyse the data; most importantly, all data and computing facilities are available to all members of the experiment, independently of the geographical locations and institute affiliations. All jobs submitted to the Grid are identical "pilot" jobs; once the pilot starts execution, it checks the environment (site, software availability, memory, disk space) and gets from a central database (the "task queue") the highest priority job that can be run on that machine. In this way there is no danger of submitting a job to a site that (at run time) is unable to run that job, or to queue a job in a given site when other sites are free to run it; also, the job priority as defined by the experiment is strictly and automatically enforced.

Not all experimental data are in event data files; in order to properly process and analyse the events, calibration, alignments and other time-dependent detector conditions data that are stored in relational databases are needed. ATLAS stores the conditions data in Oracle databases at CERN and replicates them to Oracle databases placed at the five largest Tier-1 sites; data are further exported to other sites using the web services technology, with a local cache on each site. The Frontier system [7] developed initially for the CDF experiment at Fermilab and then for CMS at CERN and finally adopted also by ATLAS and LHCb, consists of web servers in front of Oracle databases, and Squid caches placed at each site. In this way it is possible to run any job on any site without having to worry about database access or overloading central Oracle servers.

### 3. Operational experience

Over half million jobs per day are run over the Grid (without counting local batch and interactive usage); over 90k CPU cores are at any time used by ATLAS. “Production” consists of jobs simulating the physics of the experiment and the detector response, in order to compare experimental results with theoretical models, and of data reprocessing jobs, which are run when better software of calibration constants become available. “Analysis” jobs are submitted to the Grid by all Collaboration members who want to select and analyse processed data according to their own criteria. Data are available to be analysed within a day of data-taking.

A very large and distributed computing system can never be absolutely stable. Hardware failures and network interruptions happen continuously, leading to the permanent loss or temporary unavailability of data files and CPU resources. The only way to provide a robust service to the users is to implement as much redundancy and automatic failovers as financially possible, and shield the users from local site problems. First of all, all data should be replicated on disk to at least two different sites, so that if one site is off or that disk has a hardware failure, the data are still available. Data lost through hardware failures must be replicated again automatically in order to re-create the second copy. Checksums must be checked after each data transfer and compared to the original values in the data catalogues. Jobs that fail for Grid reasons must be retried automatically on a different site. The data provided by monitoring tools are used automatically by data transfer and job brokering systems in order to avoid problematic sites; for example the ATLAS “HammerCloud” [8] suite sends regular test analysis jobs to all sites and switches off job brokering to malfunctioning sites, keeps sending test jobs and when the problem is solved job brokering is re-activated automatically. Similarly, data transfer functional tests provide information for data brokering tools.

Several R&D projects have been recently started, tracking the evolution of computing technology, particularly in the fields of parallelisation of jobs for many-core processors, virtualisation, NoSQL databases as back-ends for Grid tools and interfacing to cloud computing infrastructures. The aim of this work is to continuously improve and optimise the present tools and develop and test new generations of these tools that will eventually replace the existing infrastructure, following the general trends of computing technology.

## References

- [1] The ATLAS Collaboration, *The ATLAS Experiment at the CERN Large Hadron Collider*, *JINST* **3** (2008) S08003, doi: 10.1088/1748-0221/3/08/S08003.
- [2] Roy, A. et al., *Building and testing a production quality grid software distribution for the Open Science Grid*, *Journal of Physics: Conference Series* **180** (2009) 012052
- [3] Eerola, P. et al., *The NorduGrid production Grid infrastructure, status and plans*, in *Proceedings of Fourth IEEE International Workshop on Grid Computing* (2003) 158–165. doi:10.1109/GRID.2003.1261711
- [4] Laure E. et al., *Middleware for the next generation Grid infrastructure*, in *Proceedings of the Conference on Computing in High-Energy Physics*, Interlaken (Switzerland) 2004.
- [5] Branco, M., *Managing ATLAS data on a petabytes scale with DQ2*, *Journal of Physics Conference Series* **119** (2007) 062017, doi: 10.1088/1742-6596/119/6/062017
- [6] Maeno, T., et al., *PanDA: distributed production and distributed analysis system for ATLAS*, *Journal of Physics Conference Series* **119** (2007) 062036, doi: 10.1088/1742-6596/119/6/062036
- [7] Dykstra, D. and Lueking, L., *Greatly improved cache update times for conditions data with Frontier/Squid*, *Journal of Physics Conference Series* **219** (2009) 072034, doi: 10.1088/1742-6596/219/7/072034
- [8] Van der Ster, D., et al., *Functional and large-scale testing of the ATLAS distributed analysis facilities with Ganga*, *Journal of Physics Conference Series* **219** (2009) 072021, doi: 10.1088/1742-6596/219/7/072021