# The ATLAS distributed analysis system

You may also be interested in:

Distributed Data Analysis in the ATLAS Experiment: Challenges and Solutions
Johannes Elmsheuser and Daniel van der Ster

Examples of shared ATLAS Tier2 and Tier3 facilities
S González de la Hoz, M Villaplana, Y Kemp et al.

Lessons learned from the ATLAS performance studies of the Iberian Cloud for the first LHC running period
V Sánchez-Martínez, G Borges, C Borrego et al.

CMS computing operations during run 1
J Adelman, S Alderweireldt, J Artieda et al.

ATLAS Tier-3 within IFIC-Valencia analysis facility
M Villaplana, S González de la Hoz, A Fernández et al.

ATLAS Distributed Computing Operations: Experience and improvements after 2 full years of data-taking
S Jézéquel and G Stewart

Exploiting Virtualization and Cloud Computing in ATLAS
Fernando Harald Barreiro Megino, Doug Benjamin, Kaushik De et al.

Computing infrastructure for ATLAS data analysis in the Italian Grid cloud
A Andreazza, A Annovi, D Barberis et al.

Data federation strategies for ATLAS using XRootD
Robert Gardner, Simone Campana, Guenter Duckeck et al.

# The ATLAS distributed analysis system

**F Legger, on behalf of the ATLAS Collaboration**

Ludwig-Maximilians-Universität München, Munich, Germany
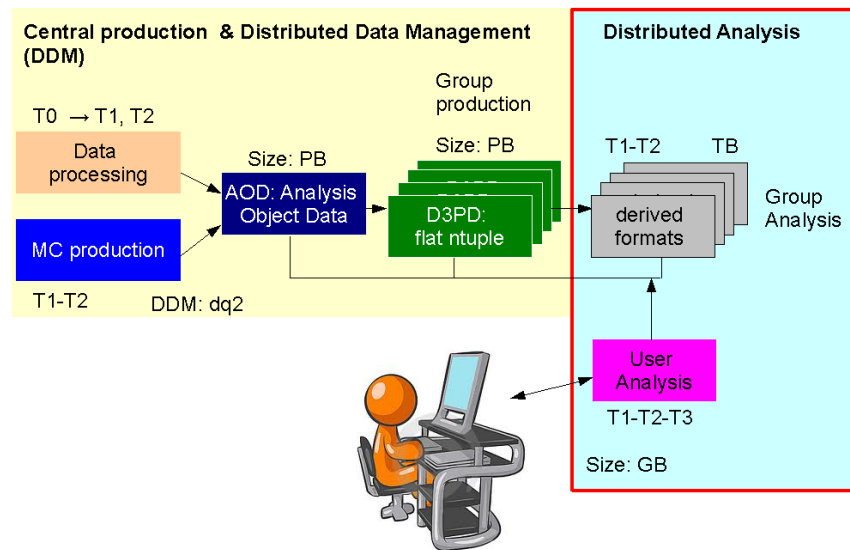
E-mail: `federica.legger@cern.ch`

**Abstract.** In the LHC operations era, analysis of the multi-petabyte ATLAS data sample by globally distributed physicists is a challenging task. To attain the required scale the ATLAS Computing Model was designed around the concept of Grid computing, realized in the Worldwide LHC Computing Grid (WLCG), the largest distributed computational resource existing in the sciences. The ATLAS experiment currently stores over 140 PB of data and runs about 140,000 concurrent jobs continuously at WLCG sites. During the first run of the LHC, the ATLAS Distributed Analysis (DA) service has operated stably and scaled as planned. More than 1600 users submitted jobs in 2012, with 2 million or more analysis jobs per week, peaking at about a million jobs per day. The system dynamically distributes popular data to expedite processing and maximally utilize resources. The reliability of the DA service is high and steadily improving; Grid sites are continually validated against a set of standard tests, and a dedicated team of expert shifters provides user support and communicates user problems to the sites. Both the user support techniques and the direct feedback of users have been effective in improving the success rate and user experience when utilizing the distributed computing environment. In this contribution a description of the main components, activities and achievements of ATLAS distributed analysis is given. Several future improvements being undertaken will be described.

## 1. Introduction

During the first LHC run, the ATLAS experiment [1] collected more than 140 PB of data. The storage space and processing power required to analyze such a large amount of data led ATLAS to develop a computing model based on distributed resources [2]. The primary event reconstruction of ATLAS collider data takes place at the Tier-0, the computing facility at CERN, and at primary computing facilities worldwide, the so-called Tier-1s. The production of Monte Carlo (MC) simulated data is also done at Tier-1s, and at secondary facilities around the world, the Tier-2s. The output of the reconstruction of both experimental and simulated data is the Analysis Object Data (AOD). The total size of a single version of AODs is of the order of PBs. The AODs are used as primary format for analysis by about a quarter of ATLAS physicists. From AODs, the physics and Combined Performance (CP) groups derive reduced data formats, Derived Physics Data (DPD) [3], which are the result of dedicated slimming, skimming and thinning procedures. Skimming is defined as the reduction of events, whereas slimming and thinning involve the reduction of objects. Such activity is referred to as group production. The most common derived format is the D3PD, a ROOT flat ntuple. There exist several flavors of D3PDs, tailored to the needs of specific groups. The total size of a single version of D3PDs is also of the order of PBs. The above described activities are centrally managed, and will be generically referred to as production in the following. The output datasets from production activities are distributed to the Tier-1s, the Tier-2s, and additional analysis facilities, the Tier-3s. At Tier-1s,

Tier-2s and Tier-3s, both MC and LHC data can be analyzed by ATLAS users. These activities are referred to as Distributed Analysis (DA) in the following, and will be the main topic of this document. Usually, the data size is further reduced to a few TB, which can either be further analyzed on the Grid to produce the final ntuples (of the order of GB) or directly using local non-Grid resources (batch systems, PROOF farms, etc). The ATLAS computing model during Run 1 is schematically drawn in figure 1.
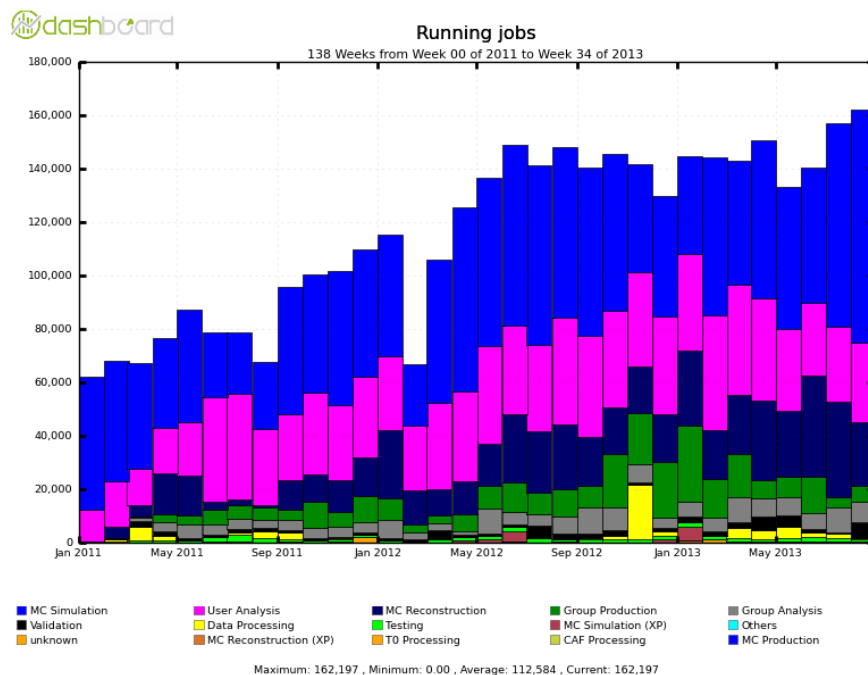


**Figure 1.** The ATLAS computing model during LHC Run 1.

## 2. Distributed Analysis activities during LHC Run 1

The recommended way for users to access distributed ATLAS resources is through the PanDA workload management system [4]. The system is based on pilot jobs with a central queue operating on WLCG [5], Open Science Grig (OSG) [6], and ARC (Advance Resource Connector) [7]. User jobs can be submitted with two different client tools, Ganga [8] and PanDA. The ATLAS software Athena [9] is installed at all Grid sites via a central installation system. The main analysis work flows include Athena analysis using AODs as input, and ROOT analysis on D3PDs. The input data (collider and MC) are distributed worldwide through the DQ2 [10] data management system. The output of the analysis jobs is temporarily stored on the scratch disks of the sites, from which it can be retrieved by the users using DQ2 client tools, or automatically transferred to group space disks.

The amount of resources for DA at the Tier-1s was initially set at 20% (in 2011) and currently ranges from 5% (minimum) to 50% (on site voluntary basis). The analysis share at Tier-2s and Tier-3s is 50% and 100% respectively. The usage of ATLAS Grid resources increased steadily in 2011, as shown in figure 2. From May 2012, the available resources were fully occupied by production and analysis jobs. The main production activities during Run 1 were MC simulation and reconstruction and group production. The ratio of production to analysis jobs is about 1/3.

More than 1600 users submitted jobs during the first run of the LHC. An average of 500000 jobs per day has been executed. The average duration of analysis jobs is less than an hour.
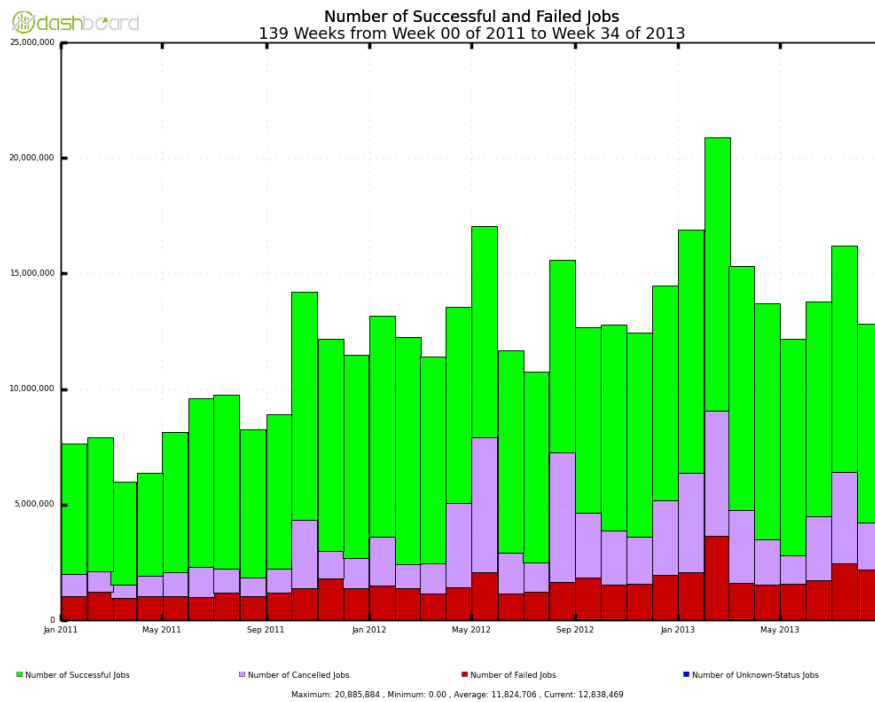
**Figure 2.** The number of concurrently running Grid jobs for the various production and DA activities since the beginning of 2011. Figure from [11].
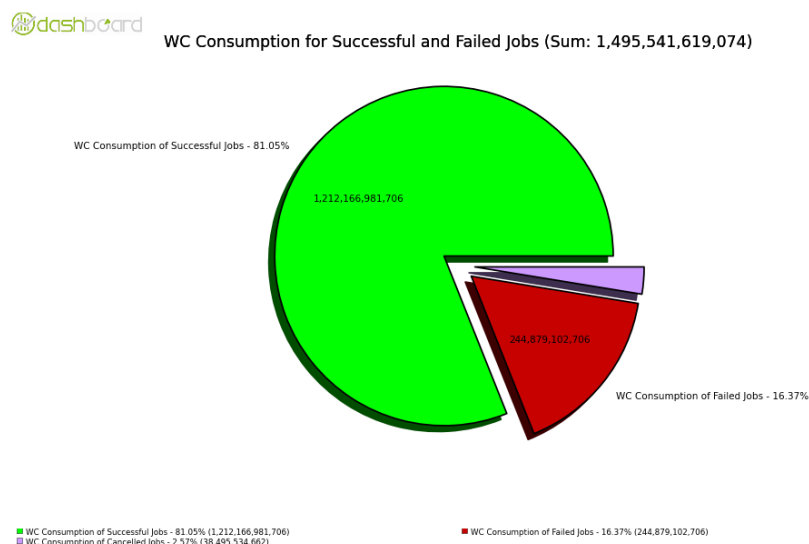
However there is a tail of jobs (less than 1% of the total) lasting longer (up to 12-15 hours). Athena-based jobs tend to have longer running times. In total, about 400 million analysis jobs were executed in the time period January 2011-August 2013, with an efficiency of $\sim 70\%$ (see figure 3). The wallclock consumption of failed and canceled user jobs is about 20% (see figure 4). Canceled jobs include jobs which are terminated by both users and the PanDA system. The main causes of failures for analysis jobs are user-related, however there is a 5% Grid-related error rate, which is also observed in production jobs, mostly due to failures of the storage systems. To mitigate the effect of Grid-related failures, automated monitoring and validation tools such as HammerCloud have been developed. The HammerCloud framework is currently used by ATLAS, CMS and LHCb, and is extensively described in [12–14]. Sites failing the HammerCloud tests are automatically excluded from job brokerage on a temporary basis, thus optimizing the usage of Grid resources. Moreover, failed jobs are categorized according to the exit code. Failures that are recognized as temporary are used to define jobs that can be automatically retried, thereby relieving the users from submitting new ones.

The failure rate of user jobs is flat over time. However peaks of user jobs are observed during busy times (usually before important conferences). As shown in figure 5, the number of submitted user jobs peaks at these times, whereas the number of pending jobs for production activities is constant. To avoid clogging the user analysis queues in the future, a more dynamic balance between production and analysis resources is needed. The use of opportunistic non-ATLAS resources such as commercial cloud services has been integrated in the PanDA production system. The use of High Performance Computing (HPC) is currently being studied and prototyped.

User support is provided by a team of expert shifters on a dedicated mailing list. Support is ensured through two daily 8-hour shifts in the European and American time zones. The support team provides the link between the user community, the Grid developers, and the sites.
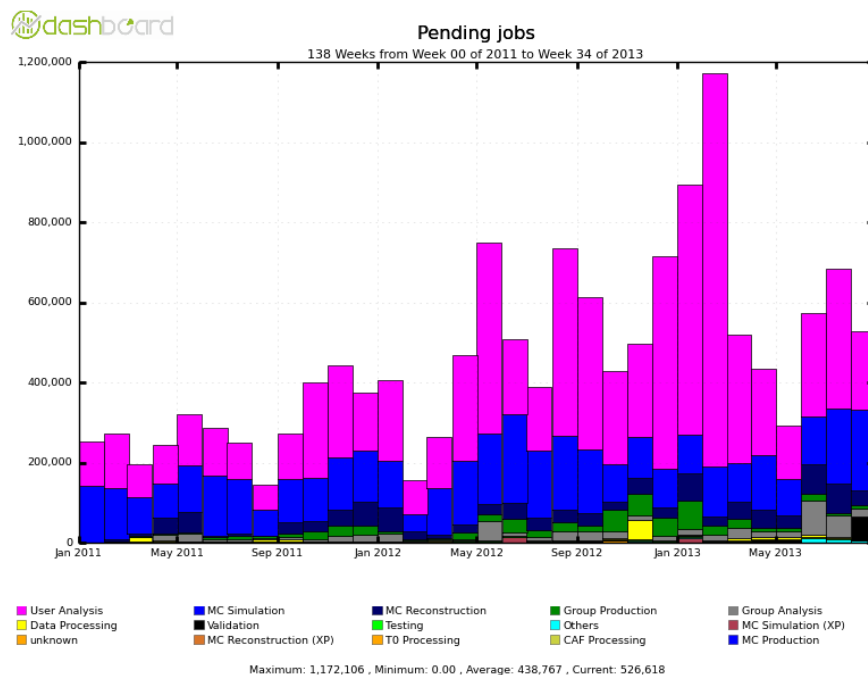
**Figure 3.** The number of completed, failed and canceled user Grid jobs since the beginning of 2011. Figure from [11].



**Figure 4.** The CPU wallclock consumption of completed, failed and canceled user Grid jobs since the beginning of 2011. Figure from [11].

The main data format used for analysis is the D3PD (existing in more than 100 flavors). AOD users are about 25% of the total ATLAS physicists. The size of one version of the most used D3PD formats is about three times the size of the corresponding AOD version. The CPU time needed to produce the most complex D3PDs ranges between 5-10 s/event, with a total of

**Figure 5.** The number of pending Grid jobs for the various production and DA activities since the beginning of 2011. Figure from [11].

30 s/event needed to produce all the official D3PD flavors. As a result, the same events exist in several formats and several copies. However, the full reconstruction of the event is only possible from the AODs. This model leads to a sub-optimal use of storage and computing resources. An optimization of the analysis data format is currently under study and will be discussed in section 3.

## 3. Future developments
During the first run of the LHC, a few limitations of the current computing model have become apparent. ATLAS' goal is to improve its overall production and analysis work flows during the long shutdown before LHC restarts operating in 2015. All the new developments have strong influence on DA activities.

The production system is being redesigned [15]. The new system will allow for more flexibility, and will be based on improved management tasks which will be more suited to handle the dynamic supervision of individual jobs. Scout jobs will be used to estimate the resources needed for each job, thus allowing a more efficient usage of Grid resources. Consequences for the users are a simplification of the client tools, resulting in shorter submission times. The job and task management will be implemented server side, which will make possible further improvements of the task monitoring (general overview, estimate of time-to-completion, automatic retrial and re-brokering of failed jobs).

To avoid the multiplication of different analysis formats, a new D3PD/AOD merged format, the so-called AODx, is under study [3]. Both Athena and ROOT-based analysis will be possible on the same input data. An estimate of disk space requirements using the new format is given in section 4. The data reduction will be centrally managed and included in the production activities, therefore freeing resources for analysis. A new analysis framework is under development to ease the application of CP recommendations and to provide a common event model. The new analysis

framework will become the recommended use case for the DA client tools.

Due to improved network bandwidth and reduced latency, remote file access through client software such as ROOT using the xrootd protocol is now being tested by ATLAS [16], and also other standard protocols based on https or WebDAV are under evaluation. The Federated ATLAS XrootD system (FAX) [17] is a storage federation aiming to treat Tier-1, Tier-2 and Tier-3 storage space as a single distributed storage system. FAX routes the client to the nearest site with the available requested data, the interaction is transparent for the user. Typical use cases include quickly scanning through large samples of data without copying them, running jobs at remote locations, fetching an unexpected missing file instead of failing the job, and processing data on non-ATLAS resources. Data-structure-aware caching mechanisms such as TTreeCache allows for good I/O performance. However, the use of FAX at large scale will have to be carefully evaluated before deployment, to avoid saturation of the site network bandwidth.

## 4. Disk space requirements

In this section, the disk space requirements for official data analysis formats at the end of 2015 are estimated. Two models are compared: the current model, where the official analysis data formats are the AODs and the D3PDs; and the new model using the merged AODx format described above. The following minimal assumptions are made:

- the AOD space requirements in 2012 are 0.24 (0.40) MB per data (MC) event;
- the same requirements are assumed for 2015 data and MC AODs;
- the size of the merged AODx is 1.25 times the size of the corresponding AOD;
- the size of the various D3PDs is 3 times the size of the corresponding AODs;
- the total size of Run 1 data and MC is twice the size of 2012 data and MC;
- one copy of Run 1 data and MC will be on disk;
- two copies of 2015 data and MC will be on disk, 2 versions are stored for each copy of the D3PDs or the AODx.

Dynamic replicas of popular datasets are not taken into account, so the estimates shown in table 1 are to be taken as highly conservative. The minimal disk space requirements at the end of 2015 are: 22 PB (data) + 50 PB (MC) = 72 PB for the current AOD+D3PD model, and 7 PB (data) + 18 PB (MC) = 25 PB for the new AODx model. With the new model, ATLAS could save up to 60% in disk space.

**Table 1.** The parameters used for the estimation of the disk space requirements at the end of 2015, in the current AOD+D3PD model and in the new AODx model.

| Datasets | Events ($\times 10^9$) | Size (MB/event) | AOD (PB) | D3PD (PB) | AODx (PB) |
|---|---|---|---|---|---|
| 2012, data | 3.9 | 0.24 | 0.9 | 2.8 | 1.2 |
| 2012, MC | 4.5 | 0.40 | 1.8 | 5.4 | 2.2 |
| 2015, data | 5 | 0.24 | 2.4 | 14.4 | 6 |
| 2015, MC | 6 | 0.40 | 4.8 | 28.8 | 12 |

## 5. Conclusions

More than 1600 physicists processed ATLAS data on distributed resources during the first run of the LHC. More than 400M analysis jobs were executed in the time period January 2011-August 2013, for a daily average of half a million jobs. The efficiency of analysis jobs is about 70%. The wallclock consumption of successful jobs is 80%. Most failures are related to problems in the user code. Based on Run 1 experiences, ATLAS is planning several improvements in its computing model to cope with the increased amount of data expected during LHC Run 2. A new production system, event model and analysis data format are in the works. The new data analysis format in particular allows for a reduction of a factor 3 of the required disk space. Moreover, the concept of easy access to remote data will be explored. The distributed analysis infrastructure will adapt to the changes to further improve the system.

## Acknowledgements

## References

[1] The ATLAS Collaboration 2008 *JINST* **3** S08003
[2] 2005 *ATLAS computing: Technical Design Report* Technical Design Report ATLAS (Geneva: CERN)
[3] Laycock P *et al* 2013 *Int. Conf. on Computing in High Energy and Nuclear Physics 2013 Amsterdam*
[4] Nilsson P 2008 *Proceedings of XII Advanced Computing and Analysis Techniques in Physics Research* vol 1 p 27
[5] Enabling grids for e-science [Online] http://www.eu-egee.org
[6] Open science grid [Online] http://www.opensciencegrid.org
[7] Ellert M *et al* 2007 *Future Generation Computer Systems* **23** 219–240
[8] Moscicki J T *et al* 2009 *Computer Physics Communications* **180** 2303 – 2316 ISSN 0010-4655
[9] Athena - the ATLAS common framework, version 8 [Online] http://atlas-computing.web.cern.ch/atlas-computing/documentation/swDoc/AthenaDeveloperGuide-8.0.0-draft.pdf
[10] Branco M *et al* 2008 *Journal of Physics: Conference Series* **119** 062017
[11] The ATLAS dashboard [Online] http://dashb-atlas-job.cern.ch/dashboard/request.py/dailysummary
[12] van der Ster D C *et al* 2011 *Journal of Physics: Conference Series* vol 331 (IOP Publishing) p 072036
[13] Elmsheuser J *et al* 2012 *Journal of Physics: Conference Series* vol 396 (IOP Publishing) p 032111
[14] Legger F *et al* 2011 *Journal of Physics: Conference Series* vol 331 (IOP Publishing) p 072050
[15] Campana S 2013 *Int. Conf. on Computing in High Energy and Nuclear Physics 2013 Amsterdam*
[16] Elmsheuser J *et al* 2013 *Int. Conf. on Computing in High Energy and Nuclear Physics 2013 Amsterdam*
[17] Gardner R *et al* 2013 *Int. Conf. on Computing in High Energy and Nuclear Physics 2013 Amsterdam*