# ATLAS Operations: Experience and Evolution in the Data Taking Era

To cite this article: I Ueda and (forthe ATLAS collaboration) 2011 *J. Phys.: Conf. Ser.* **331** 072034

View the article online for updates and enhancements.

# ATLAS Operations: Experience and Evolution in the Data Taking Era

### I. Ueda for the ATLAS collaboration

The University of Tokyo, International Center for Elementary Particle Physics, 7-3-1 Hongo, Bunkyo-ku, Tokyo 113-0033, Japan

ueda@icepp.s.u-tokyo.ac.jp

**Abstract**. This paper summarises the operational experience and improvements of the ATLAS hierarchical multi-tier computing infrastructure in the past year leading to taking and processing of the first collisions in 2009 and 2010. Special focus will be given to the Tier-0 which is responsible, among other things, for a prompt processing of the raw data coming from the online DAQ system and is thus a critical part of the chain. We will give an overview of the Tier-0 architecture, and improvements based on the operational experience. Emphasis will be put on the new developments, namely the Task Management System opening Tier-0 to expert users and Web 2.0 monitoring and management suite. We then overview the achieved performances with the distributed computing system, discuss observed data access patterns over the grid and describe how we used this information to improve analysis rates.

## 1. Introduction

The Large Hadron Collider (LHC) at CERN has been delivering stable beams colliding at 7 TeV since the first collisions at the end of March in 2010. ATLAS [1], one of the general-purpose experiments at the LHC, has been taking data with a good efficiency, accumulating more than 10 pb$^{-1}$ of total integrated luminosity.

The ATLAS distributed computing system [2, 3] consists of three classes of "Regional Centres", namely, Tier-0, Tier-1 and Tier-2. The raw data acquired with the ATLAS detector (RAW) are recorded into tape at the Tier-0 at a nominal rate of 200 Hz with average event size of 1.6 MB. The RAW data are promptly sent to the Tier-1 centres and stored on tape so that each of the data has a copy on the Grid. After the calibration of the data is performed at the CERN Analysis Facility associated to the Tier-0, the first-pass processing of the RAW data with event reconstruction is carried out at the Tier-0 [4, 5]. The processing outputs Event Summary Data (ESD), Analysis Object Data (AOD) and other derived data (dESD, dAOD, NTUP, etc.) that are then distributed over the Grid. The average event size of ESD and AOD are approximately 1 MB and 100 kB, respectively.

The main roles of the Tier-1 centres are permanently storing the copy of RAW data on tape, storing the reconstruction outputs on disk for faster access, and perform the second and the later processing (i.e. reprocessing) of RAW data hosted at the site. Each Tier-1 centre has associated Tier-2 centres and forms a "cloud". The data distributed over the Grid on 10 Tier-1 and 37 Tier-2 centres are managed in an organised way with the clouds. The Tier-2 centres are the main facilities for user jobs and host input data for user analysis on disk. The data on Tier-2 centres are replicated from their associated

Tier-1 centres. The data to be distributed on the Grid are registered to the ATLAS distributed data management system (DDM) [6].
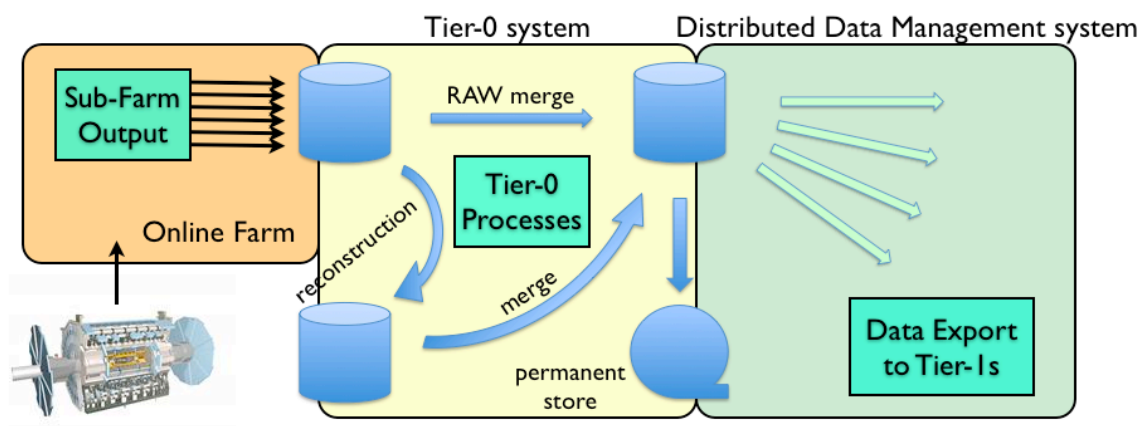
## 2. The ATLAS Tier-0
Within the ATLAS hierarchical multi-tier computing infrastructure, the most important role of the first level, the Tier-0, is accepting the raw data from the online system and ensuring it is archived to tape, merging small files to a size that is appropriate for the tape system. The other main responsibilities of the Tier-0 are the express processing of the raw data for support of calibration and alignment, the prompt first-pass processing of the raw data after calibration and alignment being done, archival of the derived data and registration of new data with the ATLAS Distributed Data Management system and other catalogues, and preparation of data for export to Tier-1 and calibration Tier-2 centres. The data flow of the Tier-0 system is summarised in the figure 1. The system has an overall I/O capacity of 6 GB/s including the internal accesses as well as the import from the online system and the export to the Tier-1 centres. The Tier-0 system is a critical part in the data processing chain and must operate in continuous 24/7 real-time mode with a high resilience to failures.

### 2.1. Tier-0 Architecture and Operation
The Tier-0 itself consists of two process entities and a database which preserves their states. The Tier-0 manager (TOM) is responsible for orchestrating the Tier-0 activity – it defines new tasks (a collection of jobs that transform one dataset into one or more new datasets) as new datasets arrive. The execution of defined tasks on Tier-0 resources is then handled by a workload management system called EOWYN, interfacing to the local CERN batch system through a plug-in called T-Zex [4]. Both systems are independent, exchanging information (e.g. about the list of task/jobs to be done and the datasets/files produced or imported into the system) in through the Tier-0 production database [4]. The Tier-0 interacts with the following external systems:
- The ATLAS distributed data management system [6],
- The event filter output processes (Sub-Farm Output) [3],
- The CERN mass storage system (CASTOR), and
- The CERN batch system (based on LSF).



**Figure 1.** Tier-0 processes in the ATLAS computing infrastructure.

The Tier-0 has been operating reliably since 2005, passed all the scaling tests and data challenges, and its functionality was continually extended [4]. The most significant additions in the period between 2009 and 2010 were:
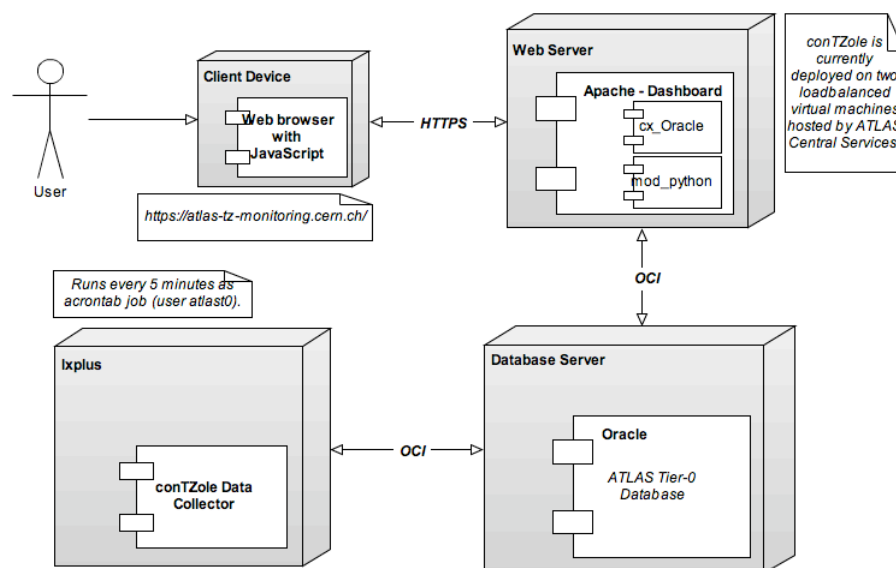
- the Task Management System, which makes the Tier-0 infrastructure available also for the detector, trigger, calibration and alignment groups,
- the new Tier-0 Monitoring system described below, and
- the inclusion of comprehensive TAG processing and uploading to distributed TAG databases hosted at several Tier-1 and Tier-2 centres.

### 2.2. New Tier-0 Monitoring System – conTZole

The new monitoring system "conTZole" has been developed primarily for the shifters and the Tier-0 operations team. The system is also useful to any ATLAS members to see how processing of a certain run goes. The conTZole project was initiated in 2009 with the aim of creating a new system which would replace all existing monitoring interfaces and merge their functionality into one interactive web accessible application, providing comprehensive real-time monitoring and management functionality. The first version was deployed for testing in summer 2009 and put into production in early 2010.

The system is currently based on the ARDA Dashboard web server (an extension of Apache server written in Python) [7], which connects directly to the Tier-0 database to retrieve real-time performance data, which are then served to the client usually in JSON format [8].
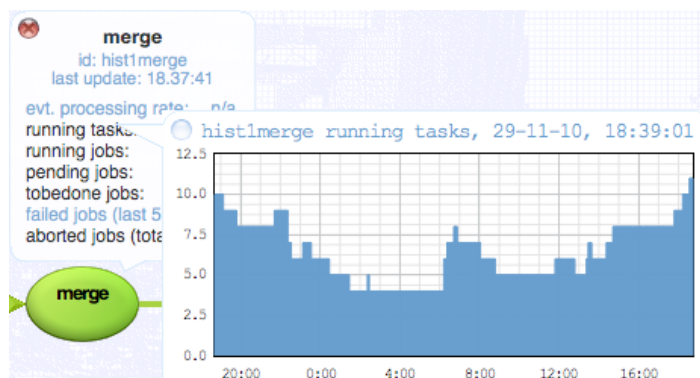
The thick client is responsible both for processing and visualization of received data. It is built using Web 2.0 technologies to achieve a high level of user interactivity. Data are periodically updated through asynchronous JavaScript calls (XMLHttpRequest) [9] and each element of the page is updated individually so that the whole page does not need to be refreshed and user settings are thus preserved. Additional (usually more detailed) data can be fetched on demand (figure 2).



**Figure 2.** conTzole deployment diagram.

For example, it may occur at the Tier-0 that up to 100 000 active jobs are defined at a given time (pending, running and finished jobs of active tasks). Displaying all of them at once would create an incomprehensible table and imply data transfers in orders of MBs. With the new interface, the user can see all active tasks (in the order of few hundreds) and only ask for jobs of those tasks behaving abnormally. Hence the number of fetched and presented data can be reduced by a factor of 1000. A listing of jobs appears directly under the relevant task and the user does not need to switch between pages – all data are conveniently on one place together with buttons to initiate corrective actions and links to view job definitions and log files for precise diagnostics of the problem (figure 3).

This approach allows users to access very detailed information for exact and precise diagnostics of any potential problem without being overwhelmed by unnecessary data and without a need to navigate

**Figure 3.** The conTZole Monitor – state of each processing step can be interactively inspected.

through dozens of disparate pages. The system alerts the user visually if the value of any of the monitored variables passes an alarm threshold. conTZole performs very well in production and received positive feedback from its users.
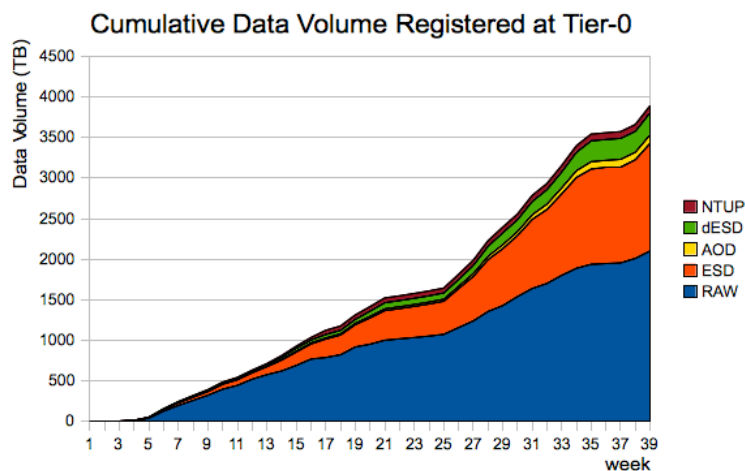
2.3. Tier-0 Performance

First collisions at the LHC started in late 2009. In 2010, the LHC ramped up to full operation, with a rapid increase in luminosity and thus the amount of recorded data. The 2010 run period was concluded with lead ion collisions in November and early December. The Tier-0 performed exceptionally well and experienced no problem which would have affected downstream clients. The system proved to be mature, reliable and resilient. The monitoring system allowed for early detection of problems with the external services the Tier-0 depends on and helped in resolving them before they could impact Tier-0 performance.

The data volume registered at Tier-0 this year is reaching nearly 4 PB (figure 4) and data export rate from Tier-0 surpassed 2 GB/s in daily averages with peaks of 4 GB/s in hourly averages. Since the overall I/O throughput for the Tier-0 storage system is limited to 6 GB/s, there have been times when it needed to throttle the export rate in order to avoid interference with the Tier-0 processes.
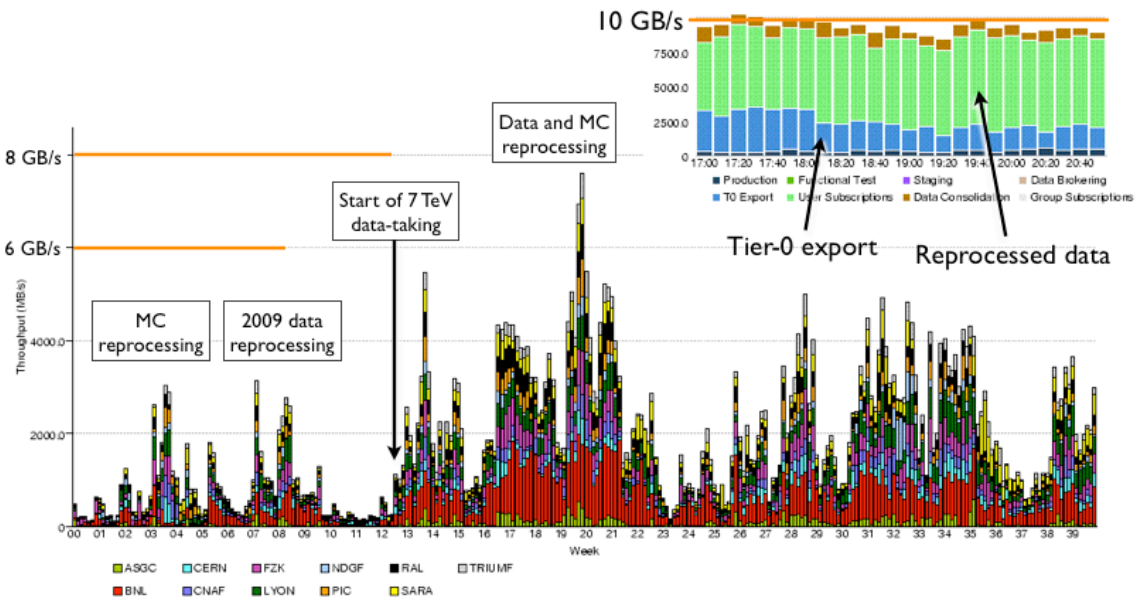
**3. Data distribution and processing over the Grid**

In addition to the Tier-0 export, there are other data transfer activities such as data movement for MC data production and distribution of reprocessed data. When all of those major activities happened all in



**Figure 4.** Cumulative data volume registered at Tier-0 over the year divided by the data type.

the same period, it was observed the over all throughput reaching nearly 8 GB/s in the daily averages and 10 GB/s sustaining several hours in the 10-minute averages (figure 5).

ATLAS has been able to sustain continued high rates of official production jobs exceeding 50,000 jobs at peaks. User analysis jobs on the grid have started rising since the start of 7 TeV collisions at the end of March. The system continues to scale up well and no problem has been observed in running those jobs [10].



**Figure 5.** ATLAS data distribution all over the Grid. The daily averages of the throughput over the year (the main figure) and the 10-minutes averages during the busiest period in the week 20.

### 3.1.  Data distribution revisited

The data distribution policies for the first year have been defined following the "ATLAS Computing Model" [2]. However, the model is not necessarily applicable to the "first year". For example, more studies on detector performance using ESD are carried out than on physics studies using AOD, and the original ESD are used more than expected while the 'derived' ESD (dESD) are not well tuned. The popularity accounting system was developed to record the number of accesses per dataset and per file from different activities. By monitoring the usage patterns of the data with the popularity accounting system, it was found that the usage patterns of the data types have been different from the presumption, and it was decided to change the distribution policies in such a way that the available disk spaces are filled constantly with the data that are possibly used, and then the data that are less used are removed later. For that purpose, a system for auto-cleaning that selects dataset replicas to be deleted based on the popularity accounting havs been developed. With this system the spaces on the grid are kept almost full, but not really full. The details on this system are described in [11].

### 3.2.  Extra Data Replication

The changes in the data distribution policies described in the previous section helped a lot in giving possibilities for user analysis to run at various sites. However, the demand from user analysis cannot always be fulfilled with pre-defined distribution patterns, and two extra data replication systems – on-demand replication and dynamic data placement – have been developed.

The on-demand replication system is a set of scheduled tools with a web interface for user requests. It has been working well with increasing number of users covering such cases that are not in pre-defined distribution patterns [12]. The dynamic data placement is a function in the distributed analysis system where user analysis jobs trigger replication of the input data to another site [10]. The system has just been introduced and is still to be tuned. In order to monitor those extra data transfers, an adjustment has been made to the data transfer monitoring tool so that transfers from different activities are distinguished with different colors. With this monitoring, it has been observed that the extra replication activities are not negligible comparing to the pre-defined data distribution, and it became possible to control the situation in case they are interfering the others.

## 4. Conclusions

ATLAS distributed computing system has been running stably with the large amount of data. The Tier-0 system has recorded and processed smoothly 2 PB of detector data and produced 4 PB of data. The data export from Tier-0 sustained at rates higher than 2 GB/s, reaching 4 GB/s at peaks, and the overall data transfer over the grid including the distribution of MC simulation and reprocessed data reached 10 GB/s without a problem. The ATLAS production system has sustained high rates of production jobs and growing number of user analysis jobs, and the system has scaled up well. The system has been continuing to evolve and improve with data access pattern measurement and auto-cleaning, enabling data distribution that matches the needs and aiming for a better environment for user analysis and faster achievement to physics results.

## References

[1]    The ATLAS Collaboration, Aad G *et al* 2008 The ATLAS Experiment at the CERN Large Hadron Collider *JINST* **3** S08003

[2]    Adams D *et al* on behalf of the ATLAS Collaboration 2005 THE ATLAS COMPUTING MODEL *CERN* ATL-SOFT-2004-007, CERN-LHCC-2004-037/G-085

       Jones R W L and Barberis D 2010 The Evolution of the ATLAS Computing Model *J. Phys.: Conf. Ser.* **219** 072037

[3]    The ATLAS Collaboration 2005 ATLAS computing : Technical Design Report *CERN* ATLAS-TDR-017, CERN-LHCC-2005-022

[4]    Elsing M, Goossens L, Nairz A and Negri G 2010 The ATLAS Tier-0: Overview and operational experience *J. Phys.: Conf. Ser.* **219** 072011

[5]    Barlow N *et al* Prompt Processing of LHC Collision Data with the ATLAS Reconstruction Software *J. Phys.: Conf. Ser.* This volume (Proc. of CHEP 2010 PS43-1-068)

[6]    Branco M *et al* 2008 Managing ATLAS data on a petabyte-scale with DQ2 *J. Phys.: Conf. Ser.* **119** 062017

       Garonne V *et al* Status, News and Update of the ATLAS Distributed Data Management Software Project: DQ2 *J. Phys.: Conf. Ser.* This volume (Proc. of CHEP 2010 PS41-5-296)

[7]    J Andreeva *et al* 2008 Dashboard for the LHC experiments *J. Phys.: Conf. Ser.* **119** 062008 (*Preprint* CERN-IT-NOTE-2007-048) and http://dashboard.cern.ch/

[8]    http://www.json.org/

[9]    http://www.w3.org/TR/XMLHttpRequest/

[10]   Maeno T *et al* Overview of ATLAS PanDA Workload Management *J. Phys.: Conf. Ser.* This volume (Proc. of CHEP 2010 PS25-3-076))

[11]   Molfetas A *et al* Popularity Framework to Process Dataset Tracers and Its Application on Dynamic Replica Reduction in the ATLAS Experiment *J. Phys.: Conf. Ser.* This volume (Proc. of CHEP 2010 PS48-1-300)

[12]   Titov M *et al* ATLAS Data Transfer Request Package (DaTRI) *J. Phys.: Conf. Ser.* This volume