

## DDM Site Services: A solution for global replication of HEP data

---

**Fernando Harald Barreiro Megino<sup>1</sup>**

*E-mail: fernando.harald.barreiro.megino@cern.ch*

**Simone Campana**

*E-mail: simone.campana@cern.ch*

**Vincent Garonne**

*E-mail: Vincent.garonne@cern.ch*

**Alessandro di Girolamo**

*E-mail: Alessandro.di.girolamo@cern.ch*

**David Tuckett**

*E-mail: david.tuckett@cern.ch*

*CERN CH-1211*

*Genève 23, Switzerland*

*On behalf of the ATLAS Collaboration*

The ATLAS Distributed Data Management (DDM) is the system built on top of the World Wide LHC Computing Grid (WLCG) middleware and is responsible for the organization of the multi-Petabyte ATLAS data across more than 100 distributed grid sites. One particular component of the system - the DDM Site Services – is the set of agents responsible for the discovery and placement of ATLAS data between sites. DDM Site Services manage aggregated throughputs of over 6GB/s or one million file-transfers a day and have to work with extremely high reliability and availability. This contribution reports on the production experience acquired during the last 2 years of LHC data taking and show the changes, adaptations and improvements that we implemented on the system to guarantee a flawless service. In addition we will give an update on the service and activity monitoring frameworks that publish the information needed by shifters and experts. Since the implementation is based on common grid middleware, these proceedings can be interesting for any community that is planning their move to the grid or would like to benefit from the approach of one of the largest heavy user communities.

*EGI Community Forum 2012 / EMI Second Technical Conference,  
Munich, Germany  
26-30 March, 2012*

---

<sup>1</sup> Speaker

## 1. Introduction

While simple data replication tools exist – e.g. in the gLite or Globus software stacks –, heavy user communities need to use a more complex data management infrastructure given the volumes of data and number of files managed and transferred between a large number of grid sites. This allows them to make a balanced usage of their resources, while preventing the network and storage to overload.

In the case of the ATLAS experiment [1], the Distributed Data Management (DDM) system is built on top of the grid middleware and is responsible for the organization of the multi-Petabyte data across more than 100 distributed grid sites. One particular component of the system - the DDM Site Services – is the set of agents responsible for the discovery and placement of data between sites. This paper will give an overview of the workflow of these agents, motivate their full deployment model and explain the solutions adopted for service and activity monitoring in order to provide the best possible service experience for users and operators.

## 2. ATLAS Distributed Data Management

ATLAS Distributed Data Management (DDM) is the system that manages the experiment's detector, simulated and user data while enforcing the policies defined in the ATLAS Computing Model. It provides functionality for data placement, deletion and bookkeeping on a hierarchic grid model composed of around 100 sites with heterogeneous storage technologies.

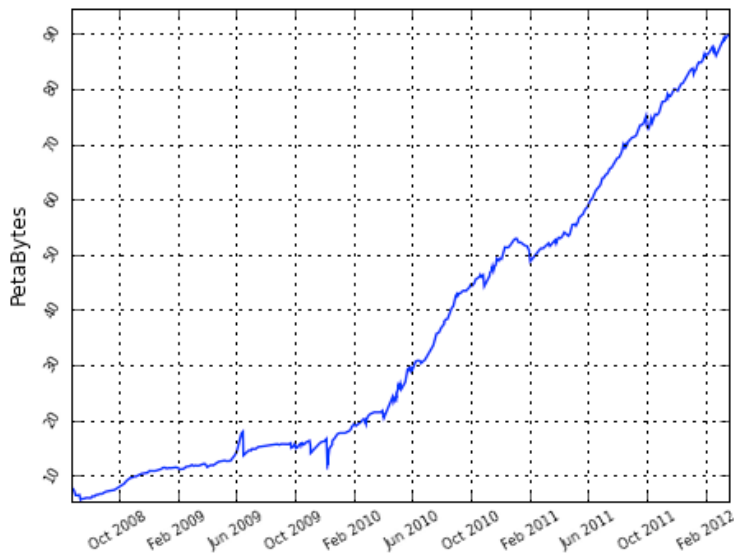


Figure 1 – Evolution of the replicated data volume ATLAS stores on the grid. Notice the increase of the slope at the beginning of 2010 when the LHC started taking data.

The DDM system represents a layer on top of the WLCG [1] middleware and provides the interface for data access to the production, analysis, users and other systems (see Figure 2). The heart of DDM is the Central Catalogues, which consist of an optimized Oracle database and a

cluster of web frontends [1]. The replication and bookkeeping unit is the *dataset* [4], a set of files, which allows reducing in one order of magnitude the number of entries in the database. The Central Catalogues store the following information:

- which datasets exist in the system: *Repository Catalogue*
- which files are contained in each dataset: *Content Catalogue*
- where the datasets are located: *Location Catalogue*
- which replication requests exist: *Subscription Catalogue*
- how the datasets and files are accessed: *Data Usage or Popularity Catalogue* [5]

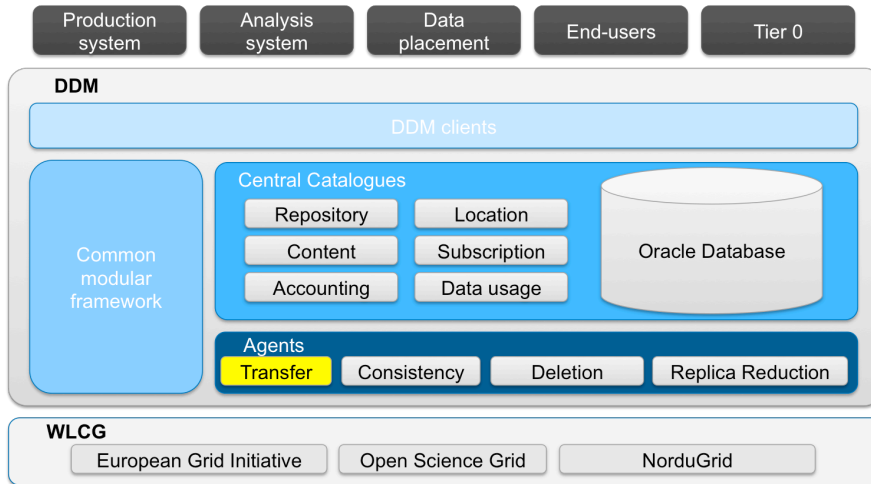


Figure 2 - ATLAS DDM architecture

Around the Central Catalogues there are a variety of services [6] that are in charge of different activities, such as the deletion agents, consistency agents or the data transfer agents. The latter are usually known as Site Services and the next sections will focus on how they work.

### 3. DDM Site Services

The DDM Site Services are the agents responsible for the data discovery and replication, being considered one of the heaviest clients that throttle data transfers on the grid. DDM Site Services are used to export beam data from CERN immediately to avoid its loss and to transfer real and simulated data across the grid to facilitate its analysis and reprocessing. As well, DDM Site Services takes care of transferring the analysis output to the final destination for a user community of a few thousand physicists. At the large scales of HEP data processing, it is important to provide a reliable system that is resilient to failures in the distributed environment so that the operations workload is minimized; the goal of the system is a transfer error rate of less than one fault per million transfers.

The DDM Site Services are in production since the beginning of the LHC data taking period and, based on the production experience, have been adapted to the evolution of the network infrastructure, new source selection policies, changes in the needs of particular storage elements or to include the requirements of emerging “off-grid” Tier3 sites.

The DDM Site Services build upon common grid middleware like FTS [8] (File Transfer Service), LFC [8] (LCG File Catalog) and SRM [9] (Storage Resource Manager). This

middleware is part of the gLite distribution and currently developed by the European Middleware Initiative [10].

### 3.1 Workflow

Each DDM Site Services instance consists of eight different types of independent agents that carry out a particular action and checkpoint their activity in an internal MySQL database, which serves as the central communication point.

0. Site Services query the dataset transfer requests from the Central Catalogues
1. For each dataset, Site Services resolves the file content
2. The locations of each dataset are obtained
3. Since the Central Catalogues only know the locations at dataset level and some files may already be available at the destination, the file lookup at source and destination is done by querying the LFC instances
4. Site Services selects the source sites for the files by evaluating past transfer statistics and submits the transfer requests to the FTS server.
5. FTS is used asynchronously and the transfer status for the files is polled
6. Copied files are registered to the LFC
7. Once the complete dataset has been registered, the new dataset replica is registered in the Central Catalogues.

In parallel, another agent is sending call-backs to the DDM Dashboard (see Section 5) updating the transfer status of the different files. At the same time service reports are written out, which are used to evaluate the health of the Site Services instances and displayed in the SLS monitoring (see Section 4).

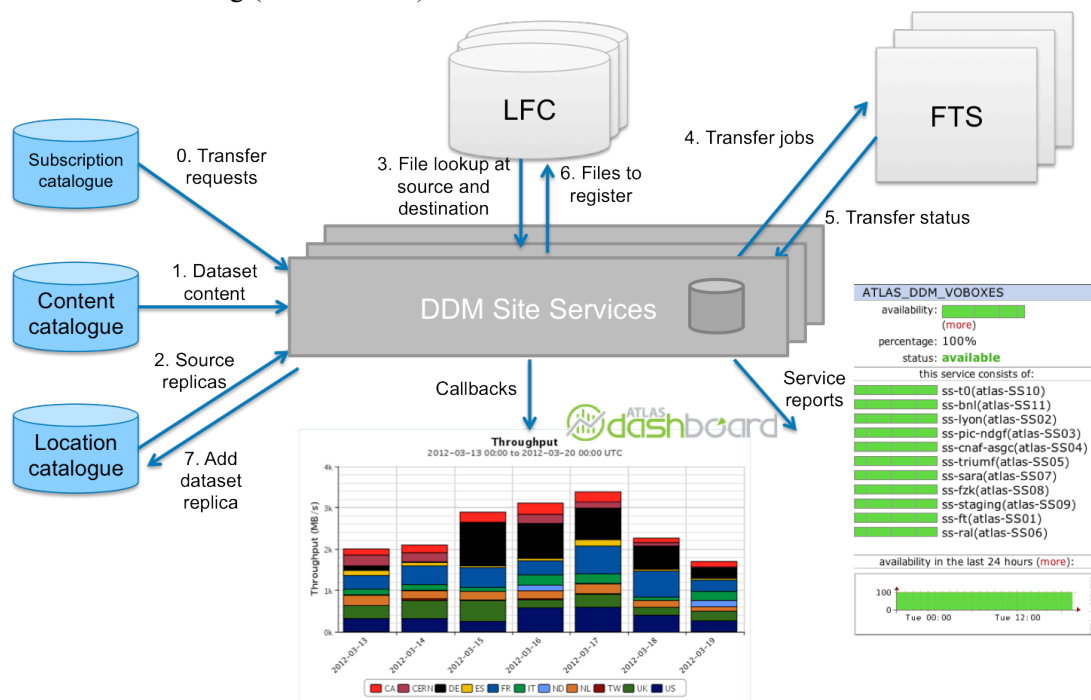


Figure 3 - DDM Site Services workflow

### 3.2 Deployment model

- **DDM Site Services:** DDM Site Services instances are hosted and operated centrally at CERN. There are 12 instances – each serving one *cloud* (organizational unit consisting of one Tier1 and O(10) Tier2s) or to carry out particular actions (e.g. stage from tape). There are also testbed and pre-production instances to evaluate code changes before putting them in production.

The Site Services instances use the standard hardware that is available in the CERN Computing Centre. Given the current virtualization process in the CERN Computing Centre, each instance of Site Services is being installed on 2 x 2-core Virtual Machines (VM) with 4GB of memory each. One VM is used for the MySQL database and the other VM for the agents. The main requirement is to have enough memory for the internal operations.

- **LFC:** In the past there used to be one Local File Catalogue per *cloud*, which was hosted in the Tier1s. If the LFC becomes unavailable, the complete cloud is unavailable not only for the data management, but also for workload management. The tendency is to consolidate all LFCs in a unique volume at CERN and create a mirror in a remote location [11].
- **FTS:** The FTS servers are located in each Tier1 and each one is responsible for the data transfer inside the cloud and into the cloud.

### 3.3 Protocols and interfaces

DDM and the grid middleware are currently coupled to the usage of SRM (Storage Resource Manager) interface and the gridFTP (FTP with authentication) transfer protocols. The first exploratory step in DDM Site Services was to evaluate the possibility of bypassing the SRM interface of certain storages and let FTS talk with the gridFTP endpoints directly. This can be the case for improving the performance in certain storages (e.g. EOS storage at CERN) or to integrate small Tier3 sites, which for simplicity would like to connect a gridFTP server to their storage directly. In the case of EOS, a significant increase of performance was observed during the data migration from the previous storage element. This was particularly evident when copying a large amount of small files, as SRM introduces a significant overhead during the initial preparation.

Two of the common storage elements in WLCG – DPM [12] and dCache [13]- have recently added an HTTP interface and it will be beneficial to evaluate the performance of these. In addition, it is of interest for the community to integrate grid storage elements with cloud computing object storages (e.g. S3).

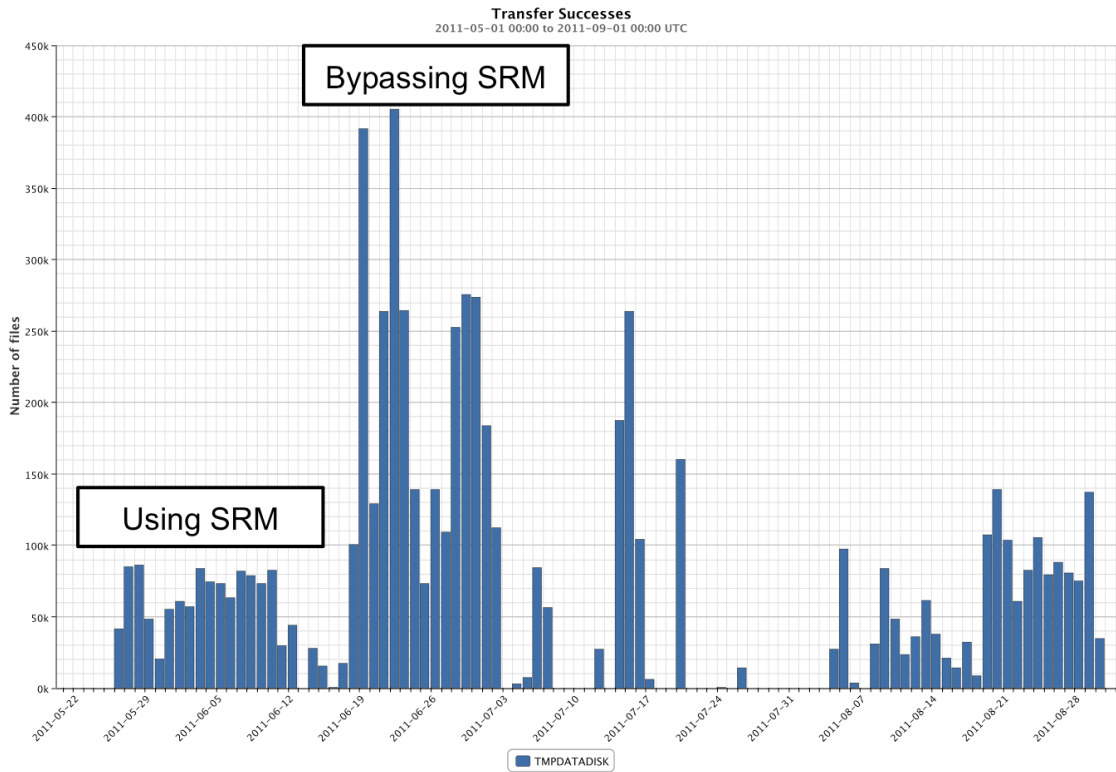


Figure 4 - Data migration at CERN from CASTOR to EOS

#### 4. Service monitoring

ATLAS Distributed Computing monitors several critical services through SLS [14]. These are not limited to the Distributed Data Management (DDM) and also include the Workload Management system PanDA, the ATLAS Metadata Interface (AMI) and the Frontier servers. Similar service monitoring based on SLS is used in other LHC experiments [15].

For the service monitoring of Site Services a simple agent is installed on each instance that counts particular events from the log files (ERROR and CRITICAL messages), parses service reports and publishes a synthetic report to SLS. Figure 5 shows the collected information that is periodically published and for which SLS keeps the history so that the service expert can view plots with the evolution of different metrics in order to pinpoint the cause of a service unavailability.

SLS periodically picks up the service reports from the specified URLs. Since Site Services and most DDM machines do not run web servers, we implemented a client that sends the service report through the MSG message queues (<https://tomtools.cern.ch/confluence/display/MIG/Home>) to a small web server (see Figure 6).

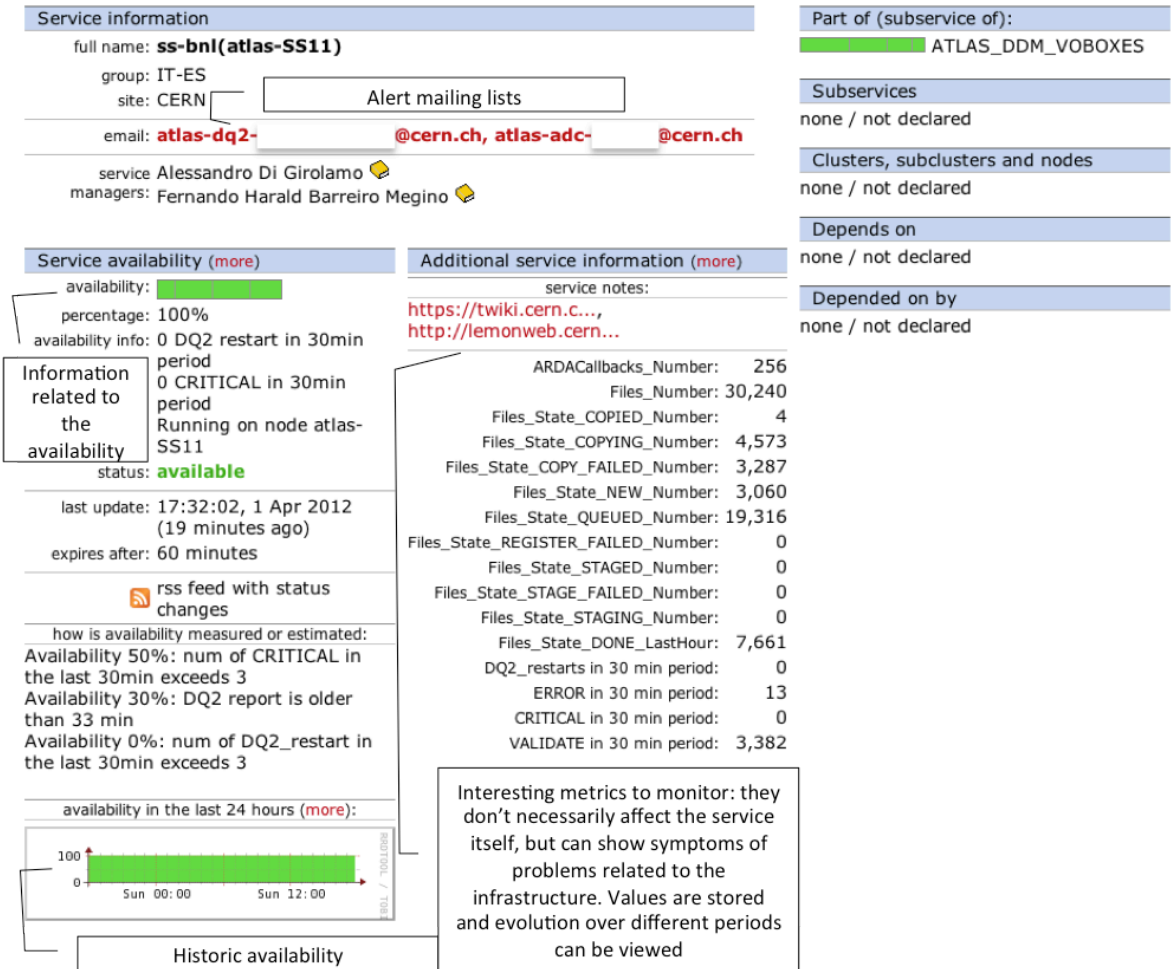


Figure 5 - Service monitoring information

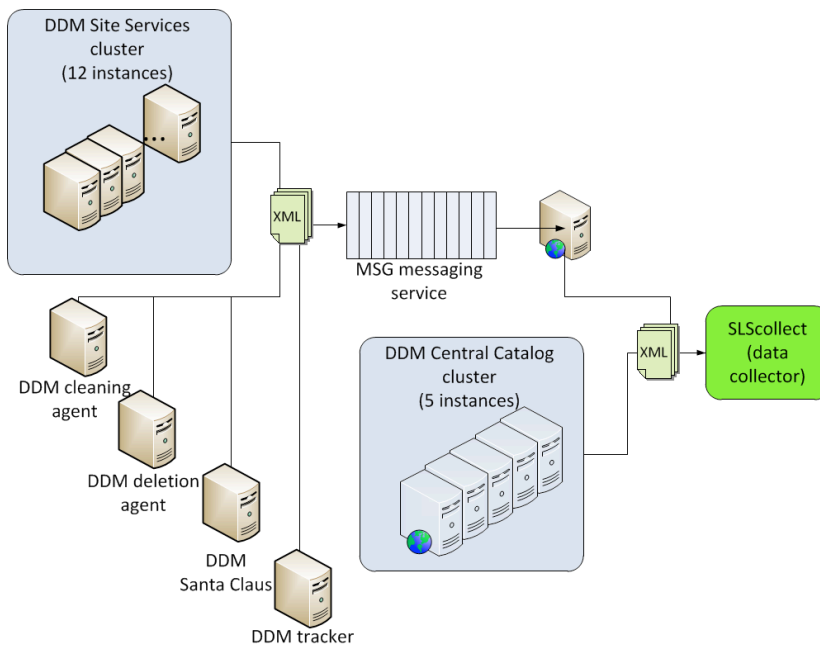


Figure 6 - DDM service monitoring infrastructure

## 5. Activity monitoring

The DDM Dashboard provides the monitoring for file transfers and registrations. The second version of the graphical user interface (<http://dashb-atlas-data.cern.ch/ddm2/>) is in an advanced state of development and is coexisting with the first version until all the functionalities are implemented and the first version can finally be phased out.

The implementation of the DDM Dashboard is based on the Dashboard team's own web framework and the *jQuery* and *highcharts* libraries. The graphical user interface has been reused for the foundations of the cross-experiment WLCG Transfer Monitoring.

The entry page (Figure 7) displays a simple matrix with the transfer statistics, where shifters can immediately spot malfunctioning links and debug the cause of the problem (e.g. a problem at a site) by drilling down in the matrix up to the needed level of detail. By clicking on the failures number, the error messages obtained from FTS will be displayed in a table.

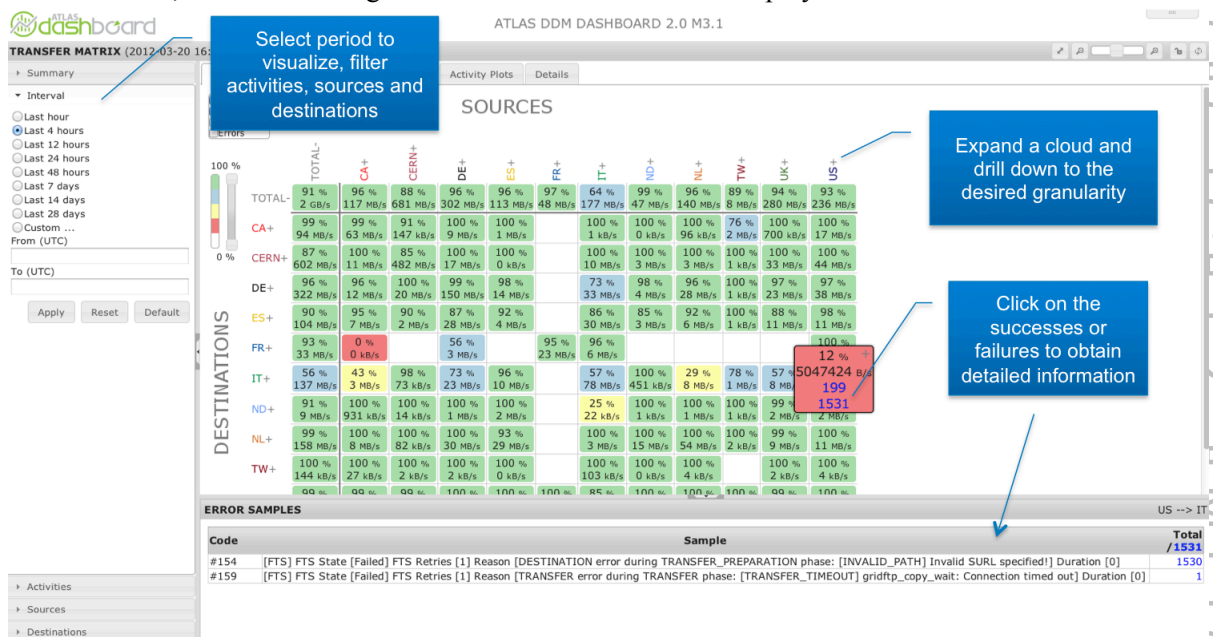


Figure 7 - Overview transfer matrix with drill down capabilities

Other DDM Dashboard views are shown in Figure 8 with the transfer statistics, efficiencies and number of failures are displayed. The user can apply filters on the displayed interval, activities, sources and destinations and view the information with the desired level of granularity (e.g. see by cloud or by site).



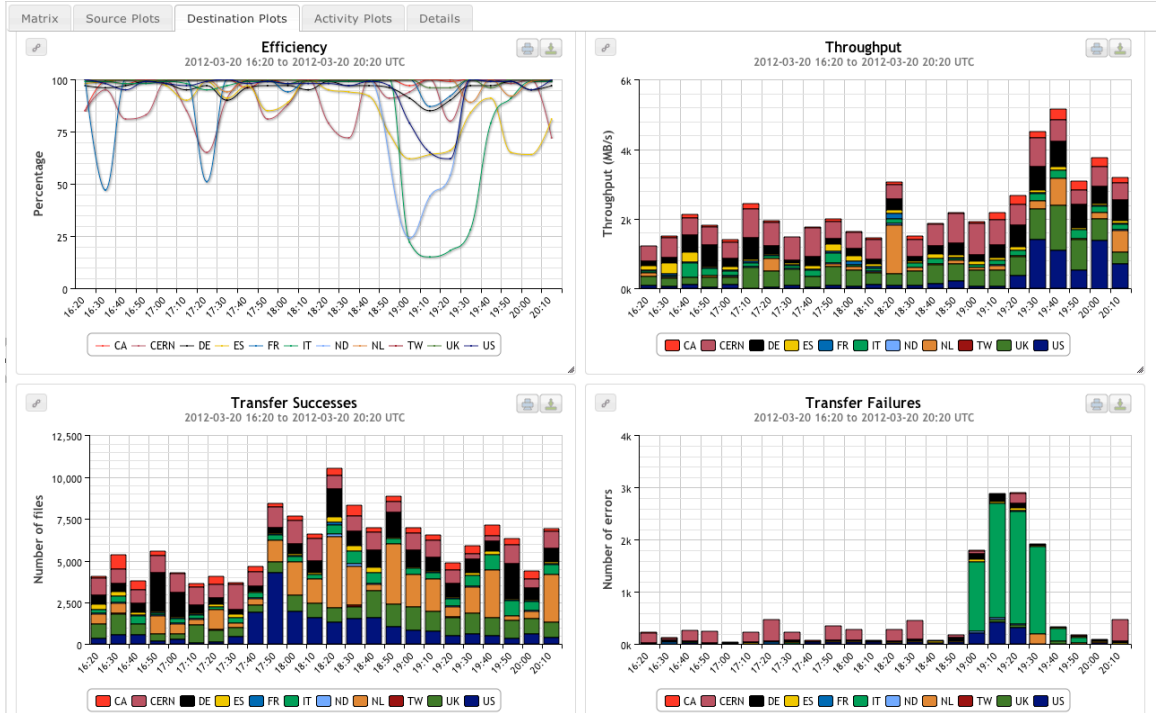


Figure 8 - Plots with transfer statistics

## 6. Related work

The challenges explained during these proceedings are common to other experiments and user communities. In particular, the other LHC experiments (ALICE, CMS and LHCb) also make use of the WLCG infrastructure and computing resources and have implemented related frameworks for workload and data management. ATLAS and CMS, are large multi-purpose detectors, while ALICE and LHCb have somewhat smaller collaborations and are more specialized. All the four experiments' frameworks have been developed and specialized for their particular workflows over more than a decade and several years of operations.

In particular, CMS is the experiment with data volumes closest to ATLAS and PhEDEx is the service responsible for the CMS data replication. Please see Figure 8 and Figure 9 to get the overview of required transfer rates for the CMS and ATLAS experiments respectively. PhEDEx is based on FTS as well, but does not depend on the LFC. Another difference is that PhEDEx does not manage user data, while DDM does. PhEDEx has served the CMS experiment successfully during the past years of data taking.

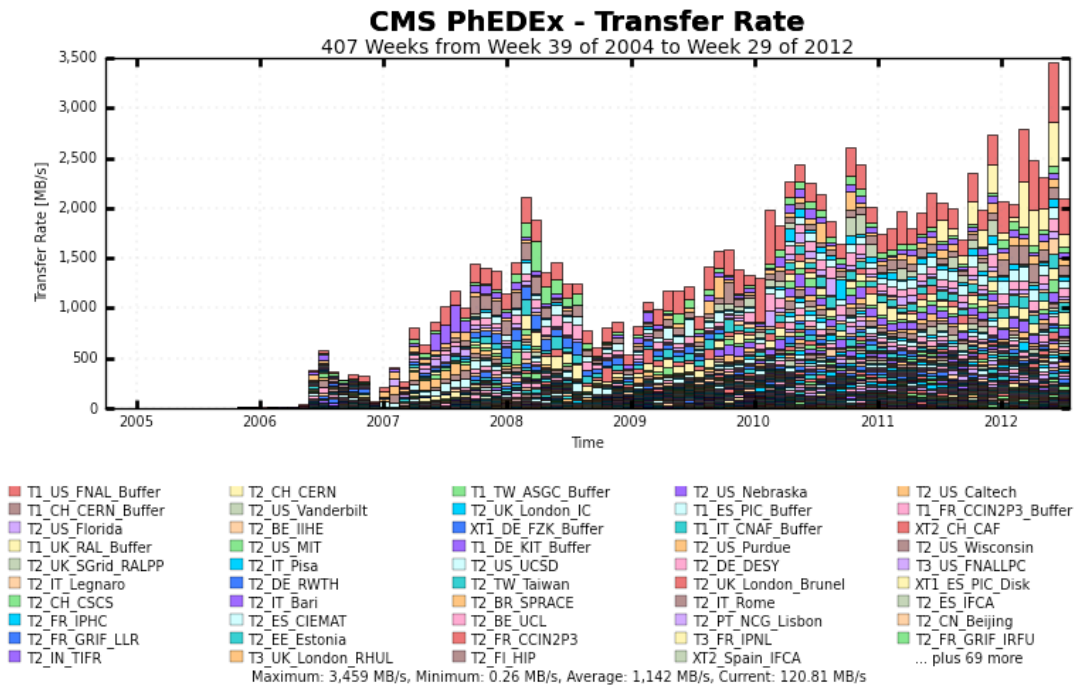


Figure 9 - Aggregated transfer rates in CMS between 2005 and July of 2012

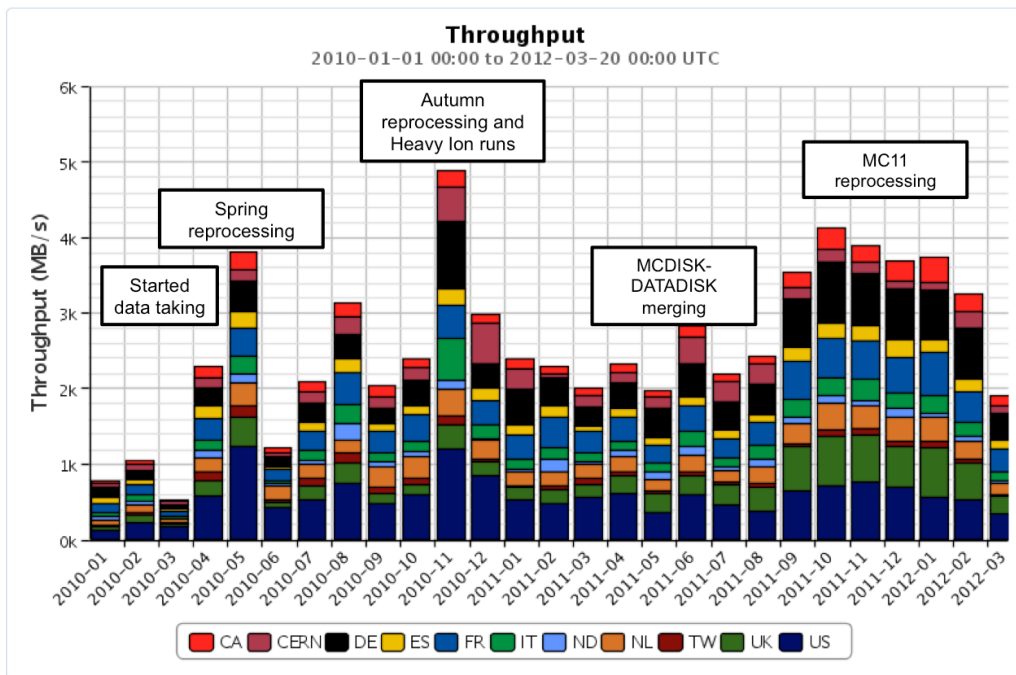


Figure 10 – Aggregated transfer rates in ATLAS between 2010 and beginning of 2012

### 7. Conclusions

To ensure further scalability, the core of the DDM Site Services has been designed as a set of independent agents, which work around an internal database to store the state. In this way indefinite instances of DDM Site Services can work in parallel as long as the central

bookkeeping system is able to sustain the load. Site Services have demonstrated to be able to cope with throughputs that exceed by far the initial design requirement of 2GB/s (see Figure 10). DDM Site Services have successfully coped with the needs during the first years of data taking and have proven high service reliability. However it is important to keep up with the evolution of new computing paradigms and technologies. One such evolution is the trending Cloud Computing model, which brings attractive features to improve the operations and elasticity of distributed computing. Although it is not clear yet how the adoption of the Cloud will impact the scientific community in the area of storage and data management, it is in the interest of DDM to evaluate how the ATLAS data organization model and the Cloud storage organization model can be integrated.

In the topic of monitoring we have worked in improving the service monitoring by publishing health reports to CERN IT's Service Level Status infrastructure, which provides an aggregated view of several services to the shifters.

We have also given a brief overview of the second version of the DDM Dashboard web frontend, which provides a very powerful and highly customizable user interface for the data transfer activity monitoring. Future work is to extend the DDM Dashboard for the other LHC experiments in order to reduce redundancy of monitoring tasks performed by the LHC experiments: this monitoring framework is known as the WLCG Transfer Dashboard [16].

### Acknowledgements

Miguel Branco as the original architect of DDM Site Services.

### References

- [1] I. Bird et al., LHC computing Grid. Technical Design Report, CERN-LHCC-2005-024
- [2] The ATLAS Collaboration, The ATLAS Experiment at the CERN Large Hadron Collider, 2008 JINST 3 S08003 doi:10.1088/1748-0221/3/08/S08003
- [3] Vincent Garonne et al, *Status, news and update of the ATLAS Distributed Data Management software project: DQ2*, CHEP, Taipei, Taiwan, October 18-22, 2010
- [4] Miguel Branco et al, *Managing ATLAS data on a petabyte-scale with DQ2*, CHEP, Victoria, Canada, September 02-07, 2007
- [5] Angelos Molfetas et al, *Popularity framework to process dataset tracers and its application on dynamic replica reduction in the ATLAS experiment*, CHEP, Taipei, Taiwan, October 18-22, 2010
- [6] Graeme Stewart et al, *Advances in Service and Operations for ATLAS Data Management*, ACAT, London, United Kingdom, September 05-09, 2011
- [7] Zsolt Molnár et al., *Next generation WLCG File Transfer Service (FTS)*, CHEP, New York, United States, May 21-25, 2012
- [8] A. Frohner, *Data Management in EGEE*, CHEP, Prague, Czech Republic, March 21-27, 2009
- [9] <http://www.ggf.org/documents/GFD.129.pdf>
- [10] <http://www.eu-emi.eu/>
- [11] Fabrizio Furano et al, *The ATLAS LFC consolidation*, CHEP, New York, United States, May 21-25, 2012
- [12] A. Alvarez et al, *Web enabled data management with DPM & LFC*, CHEP, New York, United States, May 21-25, 2012
- [13] P. Miller, dCache, agile adoption of storage technology, CHEP, New York, United States, May 21-25, 2012
- [14] Sebastian Lopienski, *Service Level Status Overview project*, HEPiX meeting, Jefferson Lab, Virginia, United States, October 10<sup>th</sup>, 2006
- [15] F. Barreiro Megino et al, *Service Monitoring in the LHC Experiments*, CHEP, New York, United States, May 21-25, 2012
- [16] J. Andreeva et al, *Providing WLCG Global Transfer Monitoring*, CHEP, New York, United States, May 21-25, 2012