

FTS3: Quantitative Monitoring

This content has been downloaded from IOPscience. Please scroll down to see the full text.

2015 J. Phys.: Conf. Ser. 664 062051

(<http://iopscience.iop.org/1742-6596/664/6/062051>)

View [the table of contents for this issue](#), or go to the [journal homepage](#) for more

Download details:

IP Address: 188.184.3.56

This content was downloaded on 09/03/2016 at 08:03

Please note that [terms and conditions apply](#).

FTS3: Quantitative Monitoring

H Riahi¹, M Salichos¹, O Keeble¹, J Andreeva¹, A A Ayllon¹, A Di Girolamo¹, N Magini², S Roiser¹, M K Simon¹

¹European Organization for Nuclear Research, IT Department, CH-1211 Geneva 23, Switzerland

²Fermi National Laboratory, Batavia, IL 60510, USA

E-mail: Hassen.Riahi@cern.ch, Michail.Salichos@cern.ch, Oliver.Keeble@cern.ch, Julia.Andreeva@cern.ch, alejandro.alvarez.ayllon@cern.ch, Nicolo.Magini@cern.ch, Alessandro.Di.Girolamo@cern.ch, Stefan.Roiser@cern.ch, Michal.Simon@cern.ch

Abstract. The overall success of LHC data processing depends heavily on stable, reliable and fast data distribution. The Worldwide LHC Computing Grid (WLCG) relies on the File Transfer Service (FTS) as the data movement middleware for moving sets of files from one site to another. This paper describes the components of FTS3 monitoring infrastructure and how they are built to satisfy the common and particular requirements of the LHC experiments. We show how the system provides a complete and detailed cross-virtual organization (VO) picture of transfers for sites, operators and VOs. This information has proven critical due to the shared nature of the infrastructure, allowing a complete view of all transfers on shared network links between various workflows and VOs using the same FTS transfer manager. We also report on the performance of the FTS service itself, using data generated by the aforementioned monitoring infrastructure both during the commissioning and the first phase of production. We also explain how this monitoring information and network metrics produced can be used both as a starting point for troubleshooting data transfer issues, but also as a mechanism to collect information such as transfer efficiency between sites, achieved throughput and its evolution over time, most common errors, etc, and take decision upon them to further optimize transfer workflows. The service setup is subject to sites policies to control the network resource usage, as well as all the VOs making use of the Grid resources at the site to satisfy their requirements. FTS3 is the new version of FTS and has been deployed in production in August 2014.

1. Introduction

The Large Hadron Collider (LHC), located at CERN, became operational in spring 2010 and led to extraordinary physics results like the discovery of a Higgs particle. A significant increase of the LHC energy and luminosity is expected during the Run 2 leading to higher data rates. To scientifically exploit the data, the data processing requires the use of computing and storage resources from many additional centers apart from CERN. These are in fact coordinated by the Worldwide LHC Computing Grid (WLCG)[1] a worldwide distributed data Grid of over 150 compute and storage clusters varying in size. They are organized in a tiered structure accordingly to the MONARC model [2], where different tier levels correspond to different functions and where each sites serves one or several LHC



experiments (ATLAS, CMS, LHCb, and ALICE). The WLCG computing model is shown in figure 1. With the technological progress of wide area networks and consequently improved network connectivity, LHC experiments have already gradually moved away from the hierarchical association of Tier-2s to one Tier-1 during Run1.

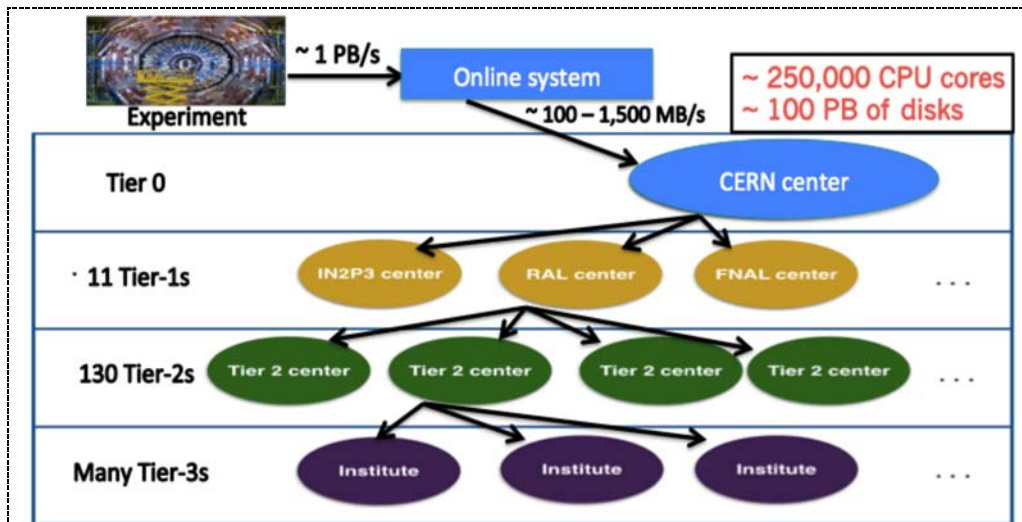


Figure 1. WLCG computing model

The data movement tools of LHC experiments rely on the WLCG File Transfer Service (FTS) for moving data from one site to another. FTS transfers around 15 PB of data each month and millions of files per day. FTS3 [3] is the current version of this service. In particular, with version 3, FTS allows moving from a structured topology to a full mesh. The WLCG data movement model based on FTS3 is shown in figure 2.

Given the complexity and the shared nature of the WLCG data transfer infrastructure as well as the new challenges that FTS3 will face in Run 2, it is important to measure its performance in a continuous way to identify and correct problems. The CERN IT-SDC (Support Distributed Computing) group has implemented a monitoring infrastructure for FTS3. This infrastructure intends to provide, with necessary measurements and statistics, means to analyze and improve the data transfer connection among the distributed computing sites. The following sections describe the main features of FTS3, its monitoring platform and how it has helped in the commissioning of the service as well as in operations and performance measurement of WLCG data transfers during the first phase of production.

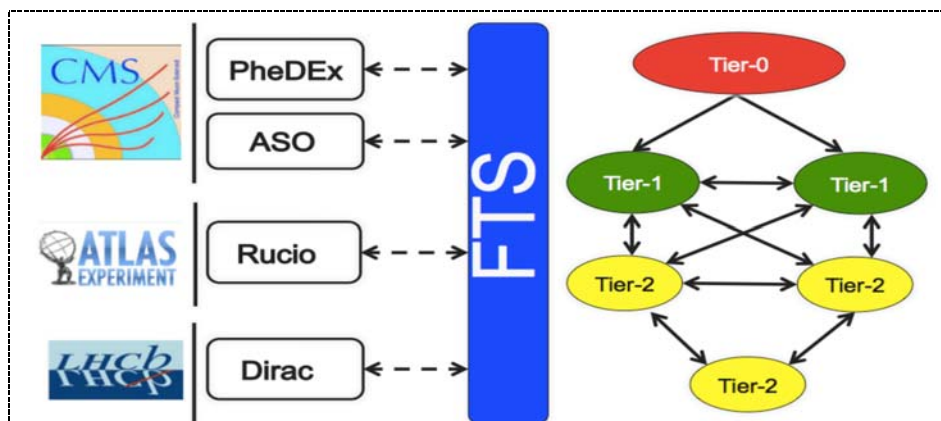


Figure 2. WLCG data movement model

2. Overview of FTS3 main features

As experiments' computing models and data access protocols are evolving, several new features have been requested and implemented in FTS3. The channel-less model has been introduced to overcome the static channel model of FTS2, where file transfers had to be performed on defined channels between sites. In FTS3, the number of parallel transfers of each link is optimized dynamically based on the throughput and efficiency of the link retrieved from its database. This strategy is known as adaptive optimization. For the support of new protocols, such as XRootD [4] and HTTP/WebDav [5], FTS3 relies on the GFAL2 [3] abstraction layer of Grid storage systems, which provides a plugin-based system.

In each Virtual Organization (VO), there are various transfer activities, namely users job output files, production data, testing activities, etc, with different latency and resources usage requirements. In FTS3, the resources sharing could be done between VOs as well as between the various activities in the VO.

FTS3 was deployed officially in production for WLCG on August 1st, 2014.

3. FTS3 monitoring

3.1. Architecture

The FTS3 monitoring infrastructure is composed of 3 different layers: Service monitoring, FTS Dashboard and WLCG Dashboard.

Service monitoring is the lowest level of the monitoring infrastructure. Each individual FTS3 server provides a Web monitoring built on top of a RESTfull interface. Some experiments' transfer applications rely also on this interface to submit and track their transfers. Every single FTS server has been instrumented to send a monitoring message each time to report state transitions to the transport layer, ActiveMQ[6] message broker. There are 4 FTS3 production servers for WLCG and several pilot and development instances publishing together in average more than 5 M messages per day. Experiments' transfer applications could also subscribe to messages published by FTS. For example, the distributed data management system of ATLAS, Rucio[7], relies on these messages to get the state of their transfers avoiding heavy polling of the service.

The raw monitoring data sent by the various FTS3 instances, are consumed by a collector in FTS Dashboard and stored in a relational database. While allowing close to real-time consuming of the stored messages, ActiveMQ can also act as a buffer in front of the database and the messages can be consumed asynchronously dramatically improving the reliability of the monitoring chain. To improve the performance of the Web application, the raw data are aggregated into statistics with different time period granularities. All database access goes through well-defined Web APIs providing data in different human-readable formats: JSON and XML. End-users access the application through a Web interface from which they can parameterize their queries and get the results in different representations such as bar chart, matrix, etc. Experiments' applications could also access the FTS Dashboard through their Web APIs. For example the CMS Analysis Task Monitoring Dashboard [9] relies on these APIs to retrieve the transfers status of the outputs produced by the analysis jobs. The results are then integrated in the Web interface of the Task Monitoring Dashboard providing a user-friendly monitoring view of the different analysis job steps.

WLCG Transfers Dashboard is a cross-technology global view of WLCG data transfers. It aggregates statistics from FTS and XrootD Dashboards. It consists of a database-less application, which makes HTTP calls to other dashboards to retrieve the data and aggregate them together. The results are then displayed in the Web user interface.

Moreover, various experiments' analytics tools are also accessing data from the various layers of the architecture.

The FTS3 monitoring architecture is shown in figure 3.

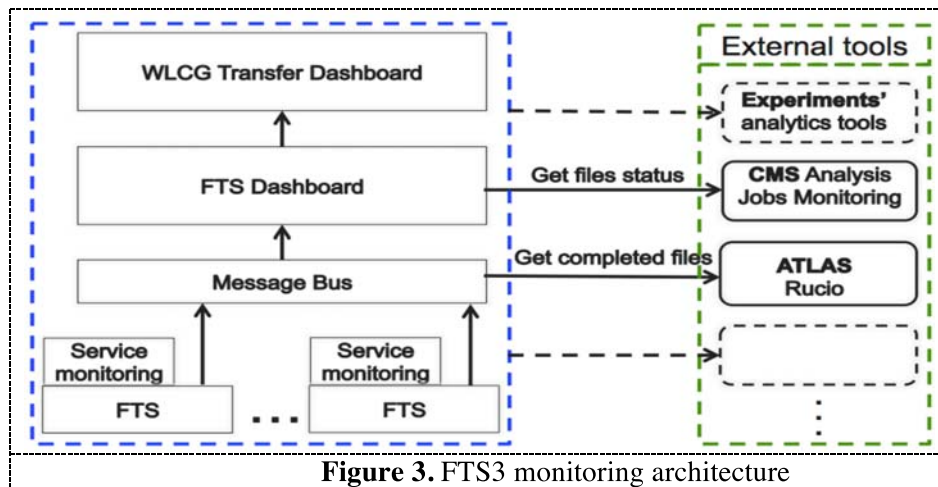


Figure 3. FTS3 monitoring architecture

3.2. Functionality

Each layer of the FTS3 monitoring offers functionality for different purposes.

FTS service monitoring provides:

- Per link/endpoint file transfers monitoring;
- Service configuration;
- Service performance monitoring. Figure 4 shows the performance of CERN pilot instance for 1 hour.

FTS Dashboard provides a global monitoring and statistics view of FTS transfers:

- Calculate transfer throughput and volume per VO/workflow, site, host and country. Figure 5 shows the aggregated transfers throughput per VO during the first month of FTS3 production for WLCG;
- Correlation of number of transfers and volume transferred;
- SRM overhead measurement;
- VO/workflow shares monitoring;
- Aggregate and report on common errors

WLCG Transfer Dashboard provides a global view of WLCG data transfers by aggregating XrootD/FTS transfers throughput/volume per VO, site, host and country.

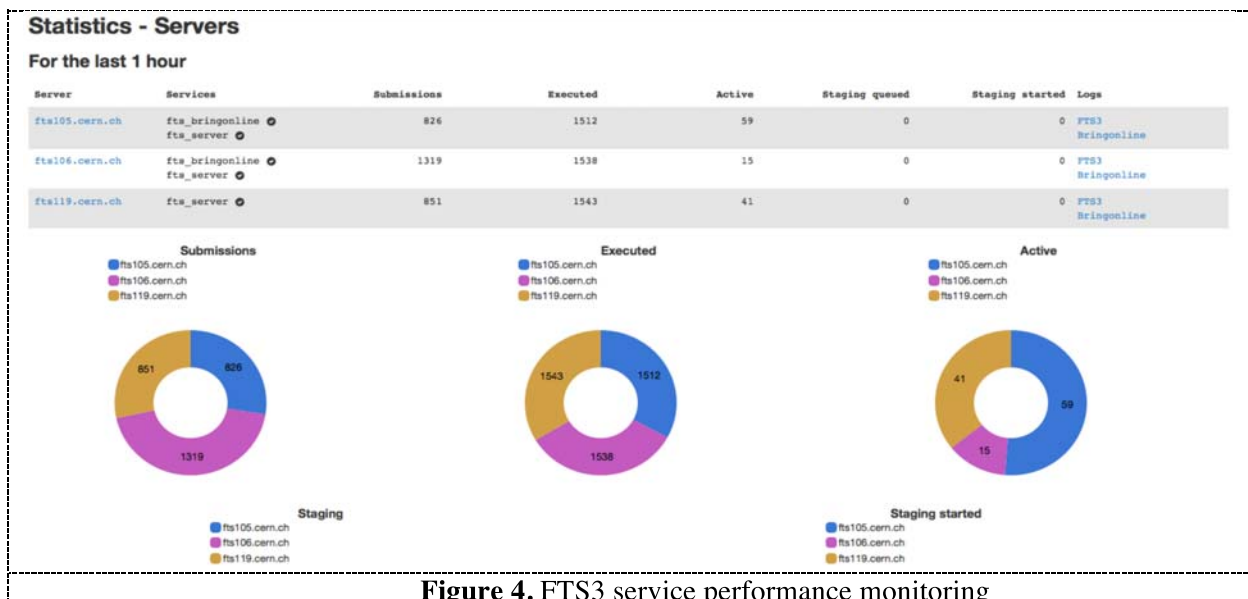
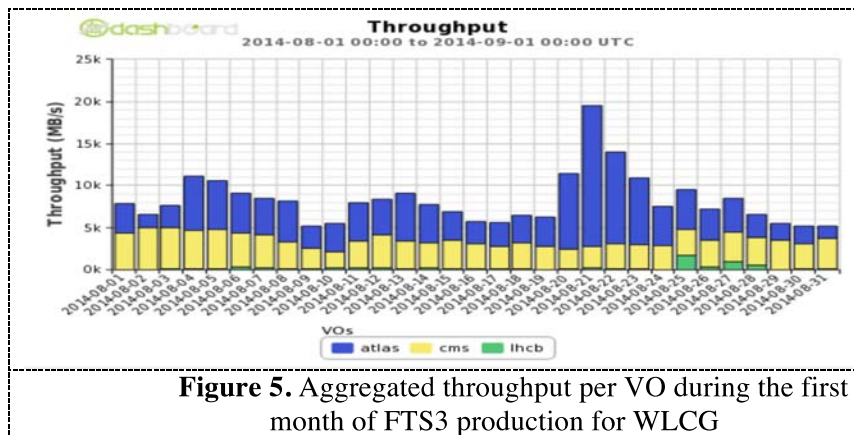


Figure 4. FTS3 service performance monitoring



4. Use of the FTS3 monitoring

4.1. Commissioning of FTS3 features

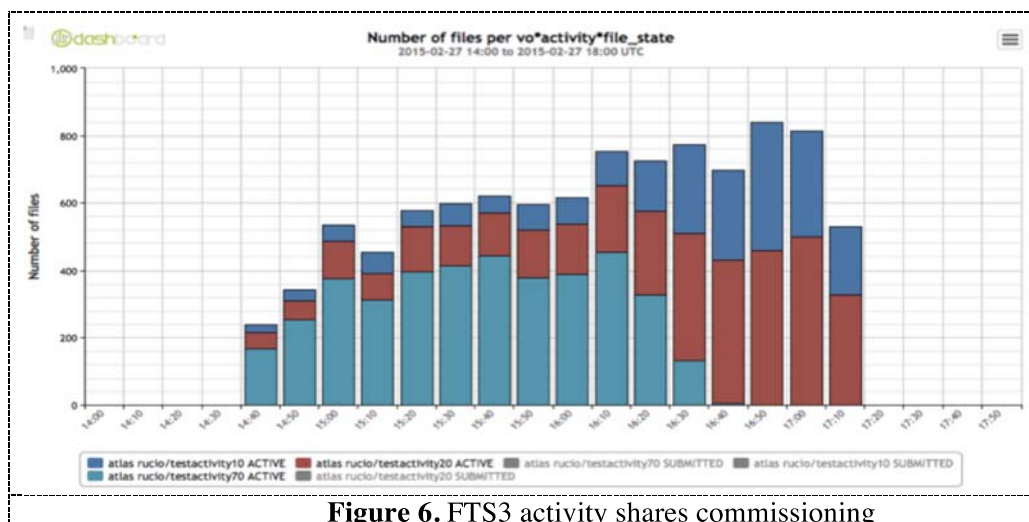
FTS3 comes with several outstanding new features for WLCG such as the VO activity shares and adaptive optimization previously described in section 2. Particular care has been required for the commissioning of these features given their major impact on resources sharing and data transfers performance.

The shares of resources for the VO activities are configured per FTS instance. Weights are assigned to experiments' transfers based on the VO strategy and latency requirements of their workflows. The monitoring of the VO activities shares has required first to include in FTS Dashboard the ability to monitor the VO transfers aggregated by activity. Then a new view has been implemented to aggregate the transfers by queued and active states.

For the commissioning of this feature, fake activities have been configured in the pilot instance of FTS3 at CERN setting different weights to each of them as shown below. The configuration is in JSON format.

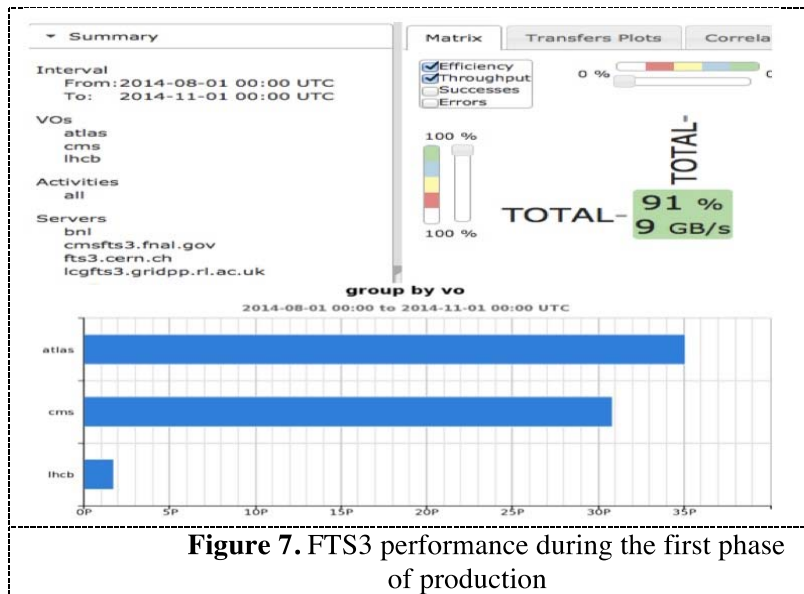
```
{ "vo": "atlas", "active": true, "share": [{"testactivity10": 0.1}, {"testactivity70": 0.7}, {"testactivity20": 0.2}] }
```

Then transfers are tagged to the various fake activities configured and submitted to FTS. As shown in figure 6, at each time during the test, the ratio of the weights of various activities set in the configuration matches the ratio of numbers of active transfers in each activity.



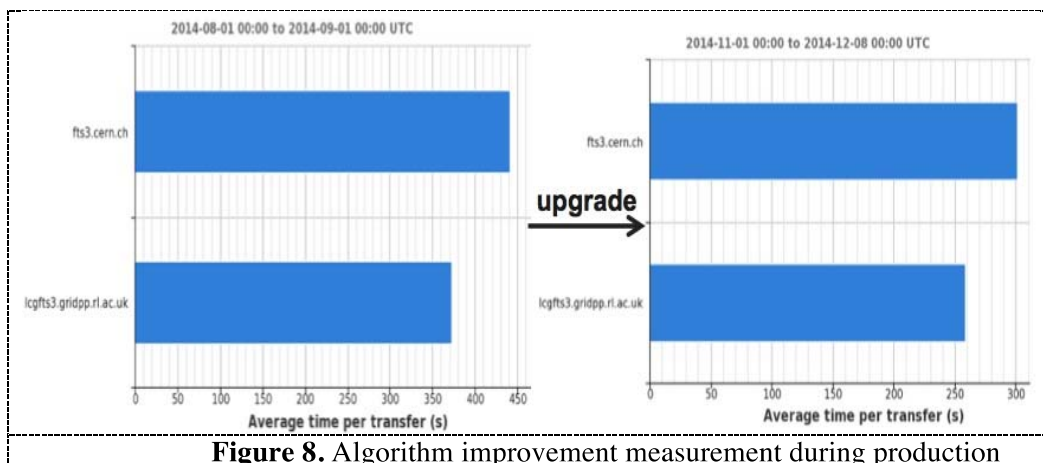
4.2. Measurement of the WLCG transfers performance

The FTS3 Dashboard provides views to measure the global performance of the WLCG transfers. As shown in figure 7, the service has shown very good performance during the first phase (3 months) of production. The aggregated throughput rate is ~ 9 GB/s with more than 90% of efficiency for the transfer of ~ 70 PB of data. Among the ~10% of failures, only 2% were caused by a service issue, e.g. because the transfer process was dying which has been fixed.



4.3. Measurement of service algorithm improvements

During the first months of production, a major improvement of the FTS algorithm has been included in the FTS 3.2.27 release. For previous releases of the service, the number of TCP streams per transfer was statically set, based on the file size. An analysis of the transfers performance has shown that it is not the best strategy. So the algorithm was changed to set the number of TCP streams per file transfer to the best number after experimenting all options over the link. The measurement of the impact of such a change was possible within the FTS Dashboard. The transfers have been categorized into various classes based on files sizes. Then for each class, the average transfer time for the same number of files and volume transferred has been measured before and after the upgrade. Figure 8 shows that for files > 2 GB the average transfer time per file has decreased by nearly 30% independently of the FTS instance after the upgrade to FTS 3.2.27.



4.4. Troubleshooting transfer issues

In 2014, CMS ran a computing exercise, Computing, Software, and Analysis (CSA14), to ensure its readiness for LHC Run 2. The goal of this exercise was to test the full chain of the end-user analysis workflow at a scale as close as possible to the scale necessary for LHC Run 2. AsyncStageOut (ASO)[10] is the distributed user data management system for CMS Analysis. During this challenge, the FTS Dashboard has been crucial to identify the infrastructural issues and for the troubleshooting of users' transfers through ASO. Figure 9 shows the transfer failure rate for 1 week grouped by destination host during this challenge. In particular, it identifies a high failure rate of transfers towards the CMS Tier-2 site, T2_ES_IFCA. The error classification functionality provided by FTS Dashboard allowed a quick discovery of the source of the issue. As shown in figure 10, the quota of one user was exhausted in the destination storage. The configuration of the scale tests has been fixed and the exercise has been recovered.

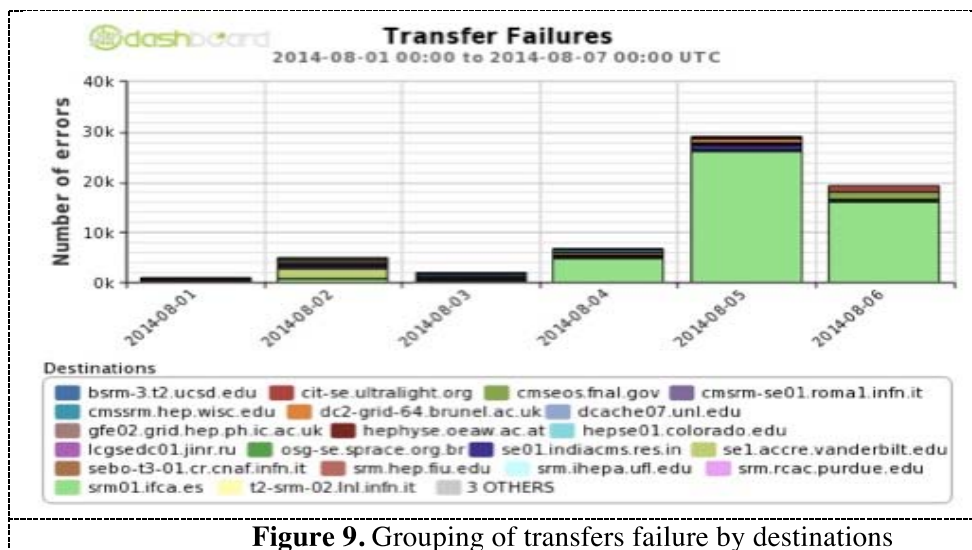


Figure 9. Grouping of transfers failure by destinations

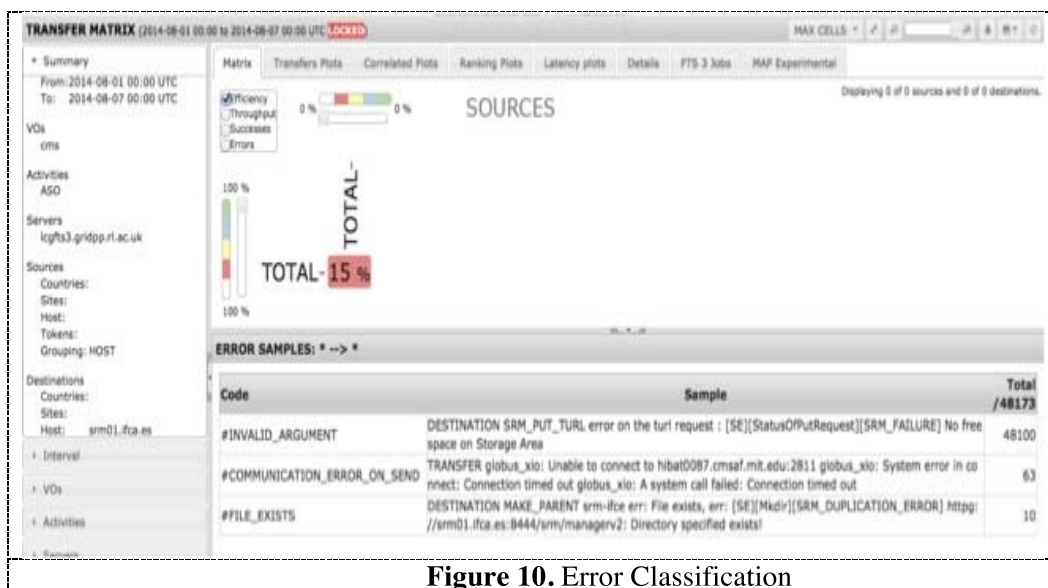


Figure 10. Error Classification

4.5. Transfer auto-tuning optimisation

To ensure readiness for LHC Run 2, the distributed data management team of ATLAS ran a data taking exercise. The goal was to expose problems that might prevent ATLAS from successfully and

speedily transfer the first data from LHC Run 2. RAW data is written first into the CERN EOS disk storage and then a backup copy is transferred to CERN CASTOR tape storage. The plot in figure 11 shows how the adaptive optimization algorithm is influenced by the achieved throughput (Figure 12) and failure rate (Figure 13) of the link, EOS --> CASTOR, and dynamically adjusts the load based on this information.

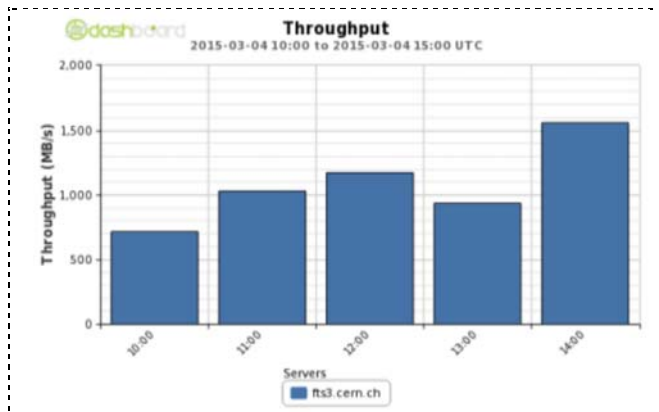
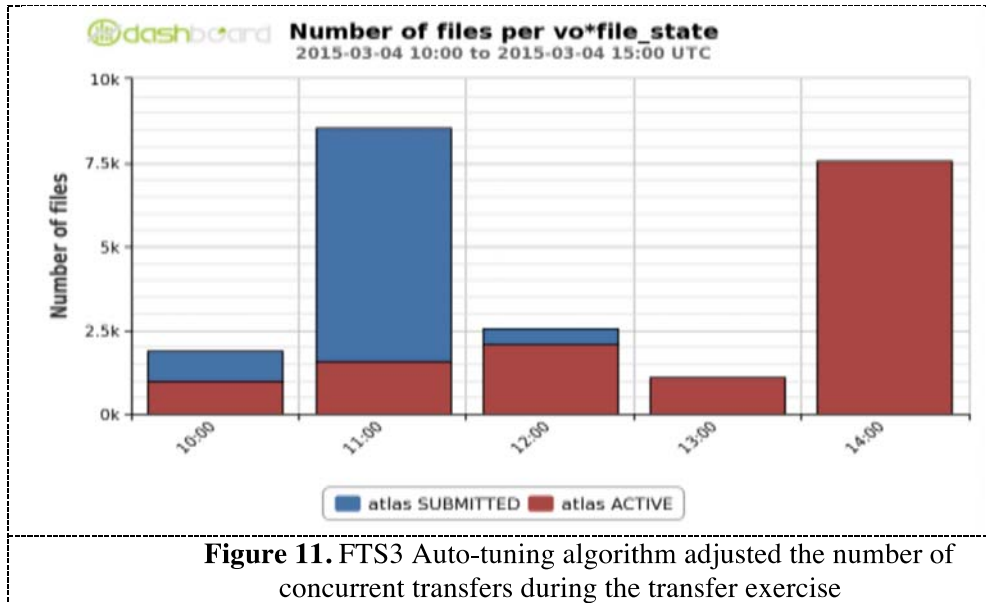


Figure 12. Achieved throughput during the transfer exercise



Figure 13. Failure rate during the transfer exercise

5. Conclusions

FTS3 was deployed in production for WLCG on August 1st, 2014 and has shown good performance so far. The FTS3 monitoring has been crucial for the service commissioning and during production for WLCG data transfers operations and performance measurement. Several experiments applications are relying on various layers of the modular architecture of FTS3 monitoring for near real-time monitoring of their transfers as well as for data analytics.

FTS3 monitoring will evolve to exploit the upcoming developments in network monitoring, perfSonar[11], for correlating network and FTS measurements for easier troubleshooting. The exploration of the historical data of FTS3 monitoring is also foreseen to improve the service efficiency and the performance of WLCG data movement.

References

- [1] Knobloch J *et al.* 2005 LHC computing Grid. Technical design report CERN-LHCC-2005-024
- [2] The MONARC Collaboration 2000 MONARC: Models of networked analysis at regional centers for LHC experiments. Phase 2 report CERN/LB-2000-001
- [3] Ayllon A A *et al.* 2014 FTS3: New Data Movement Service for WLCG. *J. Phys.: Conf. Ser.* **513** 032081
- [4] Bauerdick L A T *et al.* 2014 XRootd, disk-based, caching proxy for optimization of data access, data placement and data replication. *J. Phys.: Conf. Ser.* **513** 042044
- [5] Furano F *et al.* 2013 The Dynamic Federations: federate Storage on the fly using HTTP/WebDAV and DMLite. *PoS ISGC13* 015
- [6] ActiveMQ web page: <http://activemq.apache.org>
- [7] Garonne V *et al.* 2014 Rucio – The next generation of large scale distributed system for ATLAS Data Management. *J. Phys.: Conf. Ser.* **513** 042021
- [8] FTS Dashboard: <http://dashb-fts-transfers.cern.ch/ui/>
- [9] Karavakis E *et al.* 2010 CMS Dashboard Task Monitoring: A user-centric monitoring view. *J. Phys.: Conf. Ser.* **219** 072038
- [10] Riahi H *et al.* AsyncStageOut: Distributed user data management for CMS Analysis. Proceedings of CHEP'15 to be published by IOP *J. Phys. Conf. Ser.*
- [11] Grigoriev M *et al.* 2009 perfSONAR: Instantiating a Global Network Measurement Framework. *SOSP Wksp. Real Overlays and Distrib. Sys.* **ROADS'09**