

# Data Management Software Development for ATLAS

C. Serfon, P. Calfayan, J. de Graat, G. Duckeck, J. Ebke, J. Elmsheuser, C. Galea, J. Kennedy, C. Kummer, C. Mitterer, D. Schaile, and R. Walker

Distributed Data Management (DDM) is one of the key component in ATLAS that is used by many other domains (Production system, user analysis...). In ATLAS, DDM is performed through the use of a software called DQ2 that interacts with different services and catalogs. DQ2 is composed of many different parts :

- Data movement service : Export and register data between sites.
- Deletion service : Perform deletion consistently on the Storage Elements and the various catalogs, based on deletion requests.
- Accounting system : Monitor the evolution of disk space versus time and different metadata (group, user, pattern).

An extension of the deletion service has been added by the LMU group into DQ2 to perform deletion using local protocols. This allows a faster and more efficient deletion than using the upper layer called SRM. Deletion with local method is now available for most of the Storage Element flavours (dCache, Castor...) and fully integrated into DQ2. It is now used to delete files at GridKa and CERN. Performance is much better than using remote access (in general >10 Hz vs ~1 Hz with SRM) as seen on the plot below.

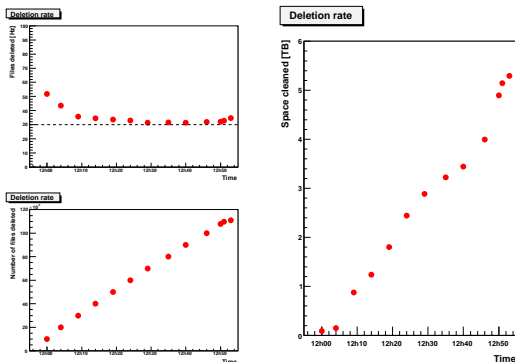


Fig. 1: Deletion of a batch of ~115 000 files at GridKa using local deletion. Deletion rate versus time (left) and space cleaned versus time (right).

In addition, a set of tools, described in [1], has been developed to help managing the DDM central operations at CERN :

- Identification and deletion of transient files/datasets from the Production System.
- Identification of corrupted files.
- Identification of files not existing anymore on the Grid.

Moreover development of consistency tools mentioned in [2] have been extended to include consistency check between Storage Elements and all catalogs (LFC, Central catalog) used in DQ2 (see figure 2).

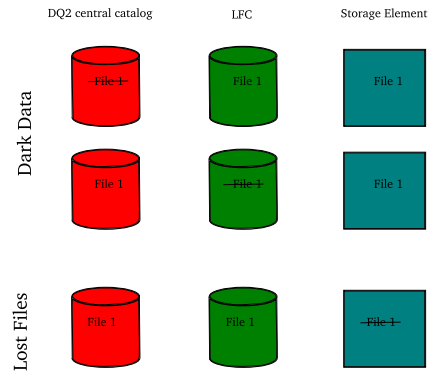


Fig. 2: Description of the various sources of inconsistencies. In the first 2 cases, the files are on the Storage Element but cannot be found on one of the catalogs (Dark Data). In the last case, the files are in the catalogs but cannot be found on the Storage Element (Lost Files).

New features including a remote check using Storage Element dumps is now available and allows to run regular consistency checks on all sites. Keeping the consistency is indeed particularly important both for users (Lost Files make jobs crashing) and for sites (Dark Data cannot be used and waste disk space). The tools are regularly used in particular in the German cloud (figure 3).

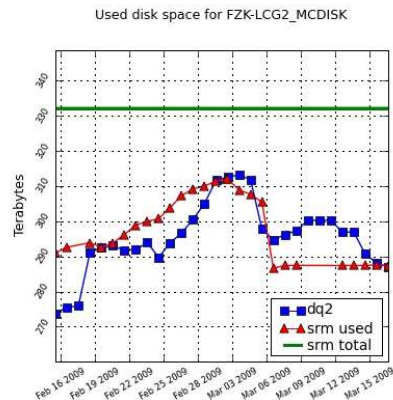


Fig. 3: Evolution of the Space seen on a site versus time : in red information from the Storage System and blue information from the catalogs. The effect of consistency checks (19th February, 3rd March, 15th March) can be seen.

## References

- [1] <https://twiki.cern.ch/twiki/bin/view/Atlas/DDMOperationsScripts>
- [2] Data Management software development for ATLAS, C. Serfon *et al.*, Annual report 2007