

N.Ratnikova¹, C-H Huang¹, A.Sanchez-Hernandez², T.Wildish³, X. Zhang⁴
¹Fermi National Accelerator Laboratory, ²Centro Invest. Estudios Avanz, ³Princeton University, ⁴Institute of High Energy Physics, Beijing

During the first LHC run, CMS saturated one hundred petabytes of storage resources with data. Storage accounting and monitoring help to meet the challenges of storage management, such as efficient space utilization, fair share between users and groups, and further resource planning.

We present newly developed CMS space monitoring system based on the storage dumps produced at the sites.

Storage contents information is aggregated and uploaded to the central database. Web based data service is provided to retrieve the information for a given time interval and a range of sites, so it can be further aggregated and presented in the desired format. The system has been designed based on the analysis of CMS monitoring requirements and experiences of the other LHC experiments.

In this paper, we demonstrate how the existing software components of the CMS data placement system PhEDEx have been re-used, reducing dramatically the development effort.

1. Problem Overview

Efficient use of distributed resources would not be possible without knowing what data are stored at participating sites and how much space they occupy.

PhEDEx knows about centrally managed data at sites. However it does not know about temporary production files or data produced by users. Some sites have their own storage space monitoring - including users and group data. Still, there is no system for monitoring all CMS data across all sites.

CMS space monitoring system has been designed to provide a global view of the distributed storage based on the sites local storage information.

2. Space Monitoring Project

First prototype realized at the end of 2011 demonstrated a proof of concept for a global storage accounting and monitoring system based on storage dumps.

In the second prototype we kept the original design, but the system was fully re-implemented using PhEDEx components, which provided safe and efficient interfaces to the database and various types of storage, and common solutions to authentication, security, documentation, and system deployment.

Testing at pilot sites revealed some limitations due to several assumptions made in the prototype. The schema was enhanced and the APIs extended to resolve these limitations.

3. Components, Interfaces, and Information Flow

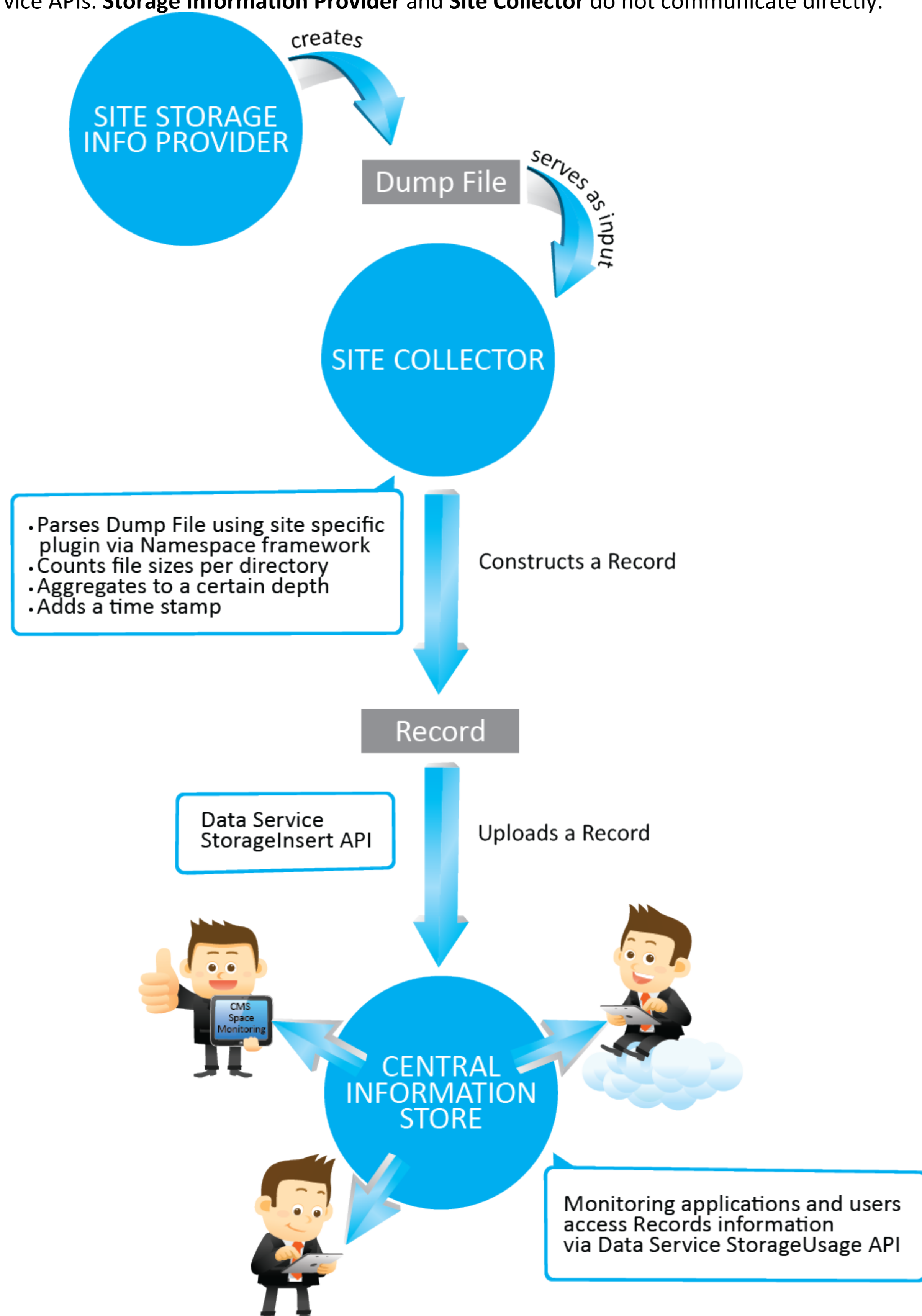
Storage Information Provider is a site and storage-specific service that produces storage dumps in the required format at regular intervals. For Dcache it can use the chimera-dump or pnfs-dump utilities.

The **Site Collector** is a process running locally on the site, which checks whenever a new storage **Dump File** is available from the **Site Information Provider**, and passes it to the **Storage Insert** utility.

The **Storage Insert** parses the dump, counts file sizes per directory, aggregates their sizes to a certain level of depth defined in Configuration, and uploads it to an Oracle database at CERN.

The **Storage Insert** utility is provided centrally to the sites as part of the **Space Monitoring** package. It comes with a set of plugins for handling different formats of the **Dump File**.

Communication with the Oracle database, both to store and to retrieve the information, is realized via Data Service APIs. **Storage Information Provider** and **Site Collector** do not communicate directly.

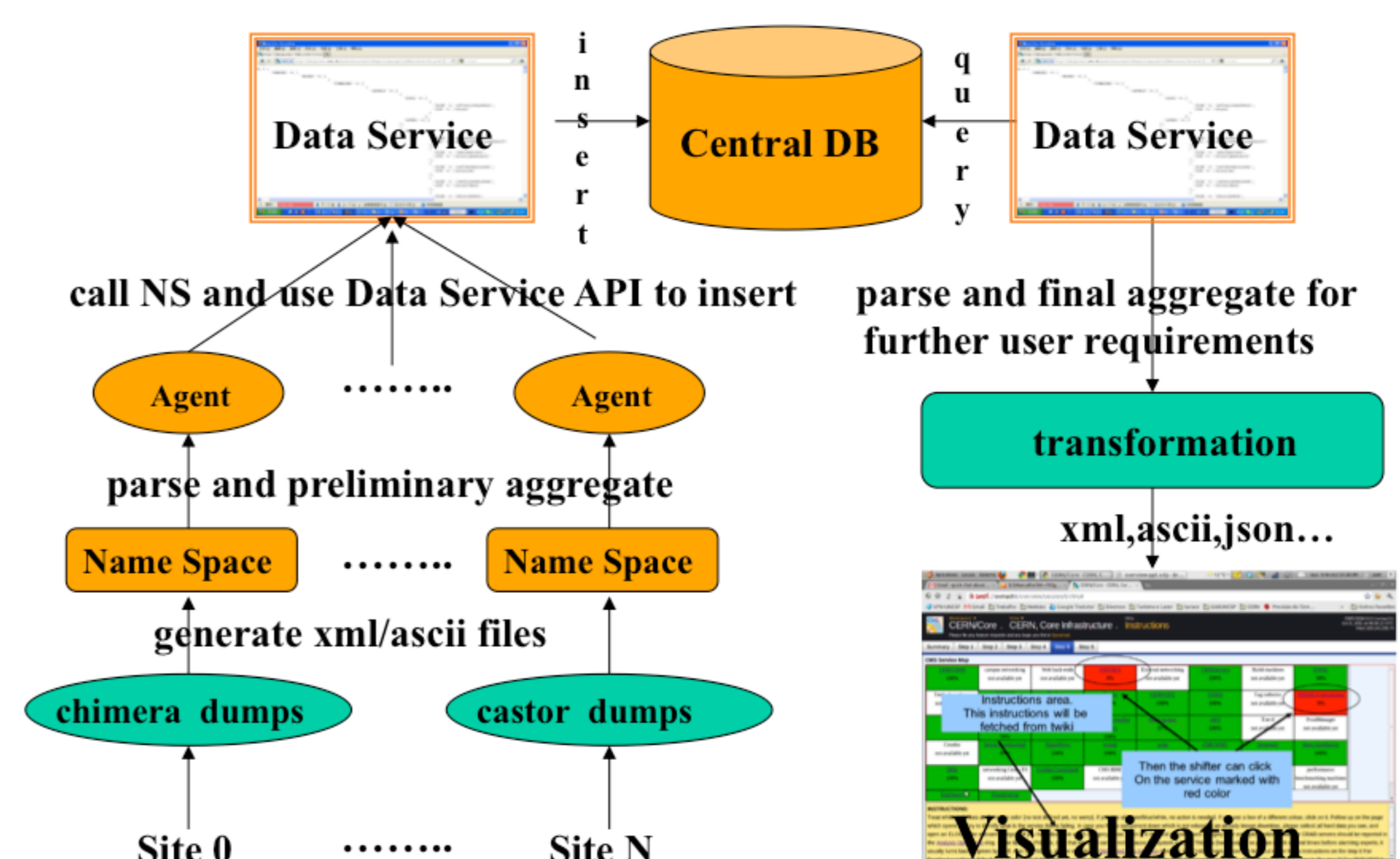


4. Infrastructure, deployment, support and operations

While Space Monitoring and PhEDEx use the same code base, they do not share the infrastructure: they use separate databases and Data Service instances.

The Space Monitoring code is maintained under the PhEDEx umbrella for practical reasons, but is packaged and distributed separately, it is not coupled with PhEDEx release cycles.

A dedicated Data service instance has been deployed by the CMS web services team.



5. Steps for the Site to deploy space monitoring

- Install Space Monitoring package on the system where storage dumps can be accessed.
- Make sure site is registered in the central information store
- Configure Site Collector to use one of the provided parsers, or write their own if needed
- Provide mapping between data types and storage locations for the configuration, part of this can be done automatically using information from the Trivial File Catalog used for the data transfers to the site
- Adjust levels of aggregation as necessary
- Start an agent or your own Site Collector scheduler to collect and feed the information to the Central Information Store

6. Summary

Recent development work includes enhancing database schema and interfaces to the space monitoring system. Data types have been introduced to allow for a targeted pre-aggregation of the information from the storage dump.

The CMS Space Monitoring project re-uses the generalized solutions and the code base from PhEDEx project: Data Service, Namespace framework, Agents framework, authentication and security model, packaging, deployment, and corresponding documentation.

It also re-uses storage-dump information which is currently used at many sites for different purposes, such consistency checking by central data operations, for local storage monitoring and troubleshooting, backup.

The re-use of the existing solutions helped to reduce dramatically the development efforts. Most work was required for the following tasks: understanding the problem and the requirements, creating the database schema, defining the interfaces, provide parsing and aggregation code.

Next steps will be to provide an easy way for the sites to deploy and configure the application and to join the global Space Monitoring system.

This will help to provide the necessary information for efficient storage resource management.

References

- N. Magini, CMS data operations, CHEP 2012
- N. Magini, The CMS data management system, CHEP 2013
- O. Gutsche, CMS Computing Operations During Run1, CHEP 2013
- N.Ratnikova, Data storage accounting and verification in LHC experiments, CHEP 2012
- T.Wildish, From toolkit to framework, the past and future evolution of PhEDEx. CHEP 2012