

STAR data management and XROOTD+SRM

(Data access and management at STAR)

Pavel Jakl¹ for STAR collaboration

¹Nuclear Physics Institute, Academy of Sciences of the Czech Republic

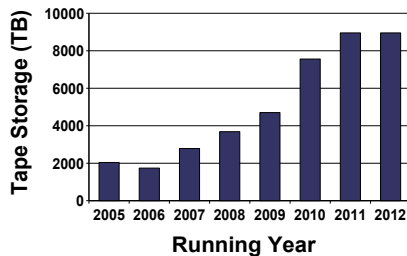
Lawrence Berkeley National Laboratory
California, USA

20th of November 2006

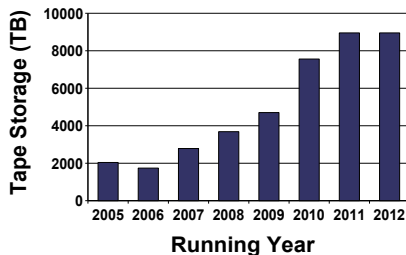
Outline

- 1 Storage challenges at STAR experiment
- 2 Past years experience and data model
- 3 Xrootd real production scenario
- 4 XROOTD+SRM integration
- 5 Summary

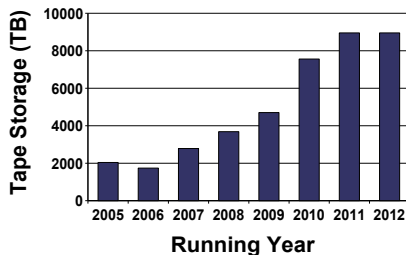
- over 1PB data per year at STAR



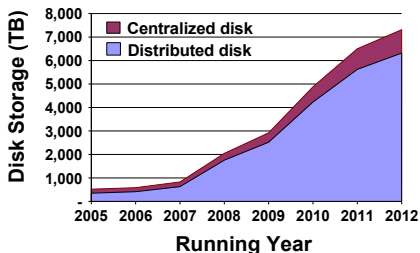
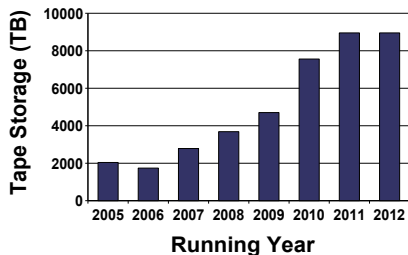
- over 1PB data per year at STAR
- *Permanent* location:
 - **tape** system (HPSS): offers several **PBs**
- *Temporary* locations:
 - **centralized** disk space(NFS area): **75 TB**
 - **distributed** disk space(spread on 500 nodes): **350 TB**



- over 1PB data per year at STAR
- *Permanent* location:
 - **tape** system (HPSS): offers several **PBs**
- *Temporary* locations:
 - **centralized** disk space(NFS area): **75 TB**
 - **distributed** disk space(spread on 500 nodes): **350 TB**
- **distributed** vs **centralized** disk:
 - ⊕ very low cost (factor of ~ 10)
 - ⊕ less human resources to maintain
 - ⊖ worse manageability (one has to build aggregation)
 - ⊖ none of current data management solutions allow to directly exploit distributed storage

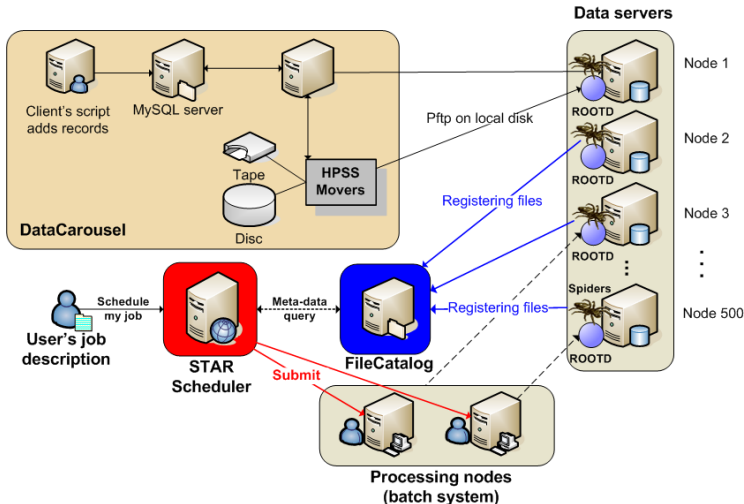


- over 1PB data per year at STAR
- *Permanent* location:
 - **tape** system (HPSS): offers several **PBs**
- *Temporary* locations:
 - **centralized** disk space(NFS area): **75 TB**
 - **distributed** disk space(spread on 500 nodes): **350 TB**
- **distributed** vs **centralized** disk:
 - ⊕ very low cost (factor of ~ 10)
 - ⊕ less human resources to maintain
 - ⊖ worse manageability (one has to build aggregation)
 - ⊖ none of current data management solutions allow to directly exploit distributed storage



ROOTD distributed data model

ROOTD - provides remote file access mechanism via TCP/IP-based data server daemon



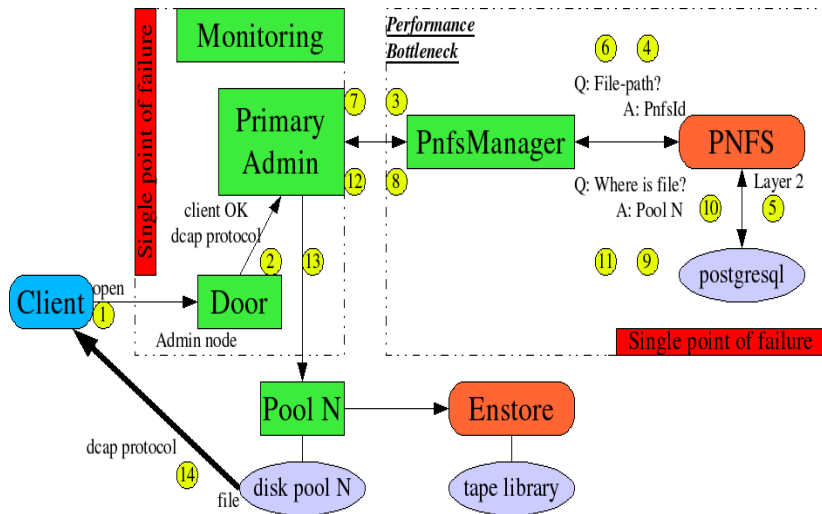
Is it rootd scalable and maintainable ?

- 1 **ROOTD knows only PFN**
 - rootd doesn't know where the data are located -> data needs to be cataloged and kept up-to-date
- 2 **Spidering scalability issues**
 - additional problems with "Spidering" of nodes when the storage grows (too many processes, db deadlocks etc.)
- 3 **Overloaded and not responding node**
 - rootd connection expires after defined time and job dies
- 4 **Job start time latency**
 - catalog is not updated accordingly when node is down for maintenance
 - job dies when requested files are deleted between the time "a" job is submitted and starts
- 5 **Static data population**
 - human interaction is needed to populate data from HPSS to distributed area
 - data-sets need to be watch (data-sets gets "smaller" in case of disk reset/format)

What are the requirements and goals ?

- distributed file systems providing high performance file-based access
- main goals:
 - Scalability** - can serve thousands of clients
 - Fault-tolerant** - adaptation to server crash or missing data
 - Flexible security** - allowing to run any security protocol
 - Load balancing** - sharing the load among multiple servers
 - MSS integration** - accessing files from permanent storage (such as HPSS)
 - Single global unique name-space** - span single name-space across multiple servers
 - Replica management** - determination of the location and multiplicity of data
 - Grid integration** - Consistent data management strategy: possibly talk to other DM tools, local or distributed on Grid
- two most popular solutions in HENP: **dCache** and **Xrootd**

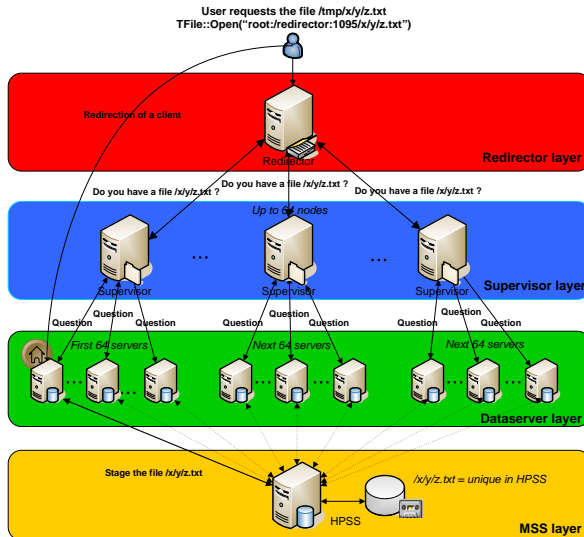
Why not dCache for STAR?



Solve rootd problems with xrootd features

- 1 **ROOTD knows only PFN** \Rightarrow **XROOTD knows "LFN"**
 - data are located within xrootd process and only LFN needs to be cataloged (reducing problem from 1:N to 1:1)
- 2 **Spidering scalability issues** \Rightarrow **XROOTD knows "LFN"**
 - no need to index available data on nodes, data are located within xrootd process
- 3 **Overloaded and not responding node** \Rightarrow **Load balancing**
 - xrootd determines which server is the best for client's request to open a file
- 4 **Job start time latency** \Rightarrow **Fault tolerance feature**
 - missing data can be again restored from MSS
- 5 **Static data population** \Rightarrow **Mass storage system plugin**
 - movement from **static** population of data to **dynamic**

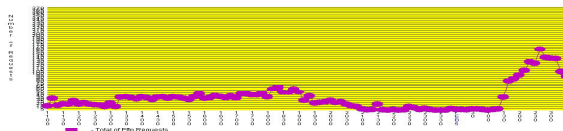
XROOTD architecture and request handling ?



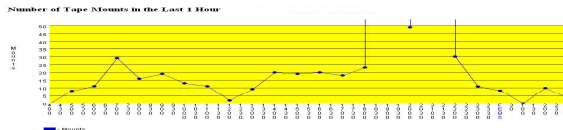
HPSS access pattern consequence

- requests to HPSS are not coordinated:

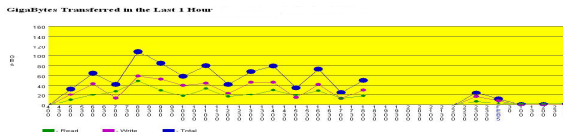
① increase number of requests



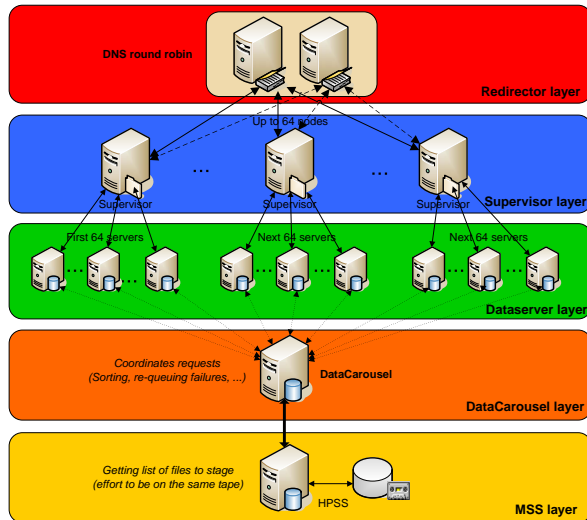
② increase tape mounts to maximum



③ decrease I/O Rate to zero

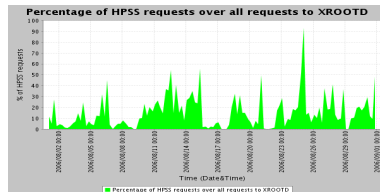
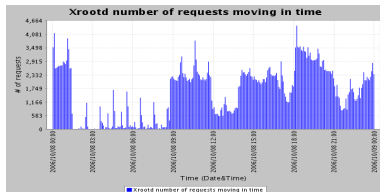


XROOTD with HPSS request coordination

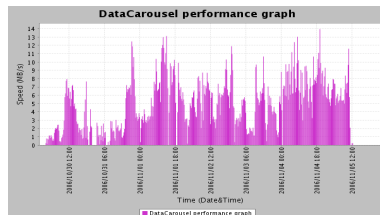
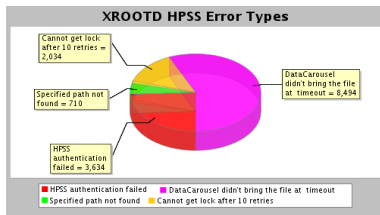


Xrootd in production and real analysis scenario

- possible to see up to 35 requests/sec to open a file, users use xrootd to access HPSS data-sets



- most of errors are caused by slow performance of HPSS



Analysis scenario

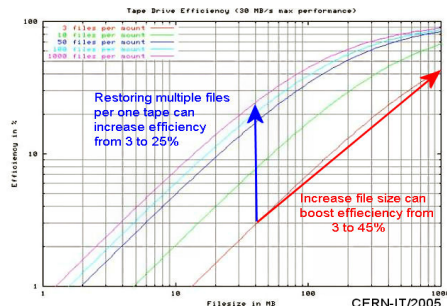
- users defines their jobs(=analysis) using job's description of **SUMS** (STAR Unified Meta Scheduler)
 - SUMS resolves their **meta-data query** (energy, collision etc.) into particular **physical data-sets** by handshaking with **STAR FileCatalog**
 - SUMS orders data-sets, splits them into specific sub-jobs and submits into batch system queue (~100-1000 files per one sub-job)
- STAR files are **compressed and structured** ROOT files (events sorted into tree structure)

	avg size of file	file's description
an event file	~284 MB	same size as daq files
an MuDST file	~88 MB	mostly used for analysis

- batch system controls the run time of the job by **Wall clock max-time**
- one file is restored from HPSS in average of **21 minutes**

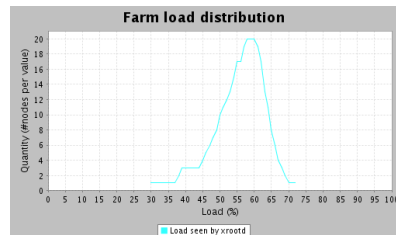
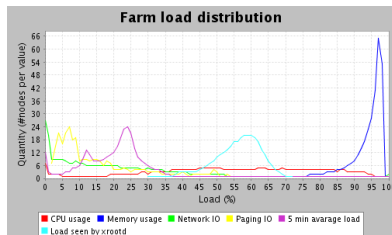
Reasons for slow HPSS performance

- drawback for user's side:
 - job hangs waiting for a file to be restored from HPSS and eventually is killed by batch wall clock (e.g $21 \times 1000 = 350$ hours)
- ① STAR files are **too small** (in avg. 90MB), should be 10 times bigger
- ② Xrootd **random** access and **sequential** processing causes excessive mounting of tapes
- ③ excessive mounting **destroys** tapes
 - DataCarousel already does a sorting of requests per tape, but not enough efficient
 - we need bigger list for sorting \Rightarrow **Pre-staging** of files
 - job publishes its whole intend for processing \Rightarrow usually files on the same tape



Understanding the load to increase the performance

- distributed system can have several choices to fulfill a incoming request (more replicas of a file etc.)
- the system needs to balance the load among many collaborating servers
- xrootd offers computation of the server workload as a flexible formula:
 - it is a combination of 5 main factors (cpu, memory etc.)
- how to setup the thresholds to represent STAR's environment ?
 - is it CPU-bound, Memory-bound environment?



Observing load distribution stability

- load distribution illustrated to be pretty stable over longer period of time

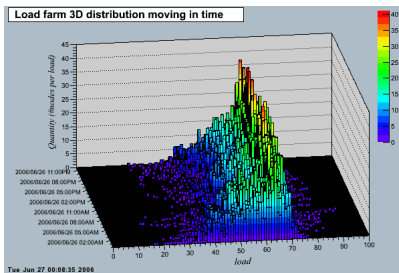


Figure: Week 26, Monday

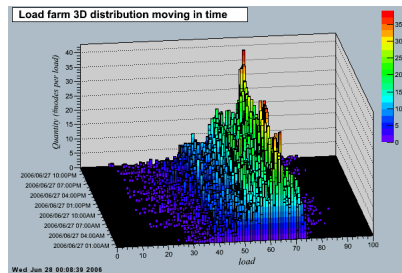
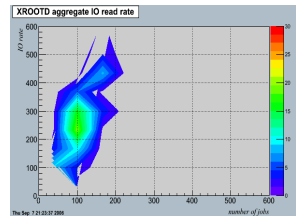
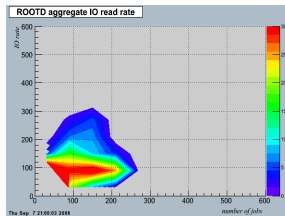
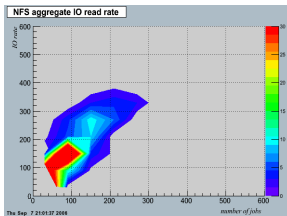


Figure: Week 26, Tuesday

Preliminary comparison of several solutions

- success of distributed file system relies on the ability to support increasing number of users with stable performance of individual file's operation
- it implies performance comparison with the following measurement's requirements:
 - relation of the aggregate IO throughput with the number of requests
 - identical measurement's conditions (same structure of files, compression etc.)



Motivation ...

XROOTD is not perfect and could be extended:

- does not bring files over from other space management systems (dCache, Castor etc.)
- always bring files from MSS, not from neighboring cache
- in large scale pools of nodes, clients could ALL ask for a file restore: lack of coordination or request "queue"
- no advanced reservation of space, no extended policies per users or role based
- no guarantee for stored files (no lifetime, no pinning of files)
- only access files (what about event-based access ?)
- other middle-ware are designed for space management.

Leveraging on other projects and targeted re-usable components ?

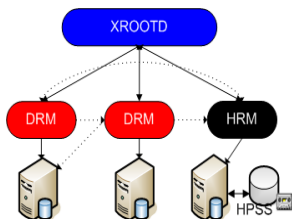
SRM functionality

- **SRM:** the grid middle-ware component whose function is to provide dynamic space allocation and file management on shared distributed storage systems
 - **Manage space**
 - Negotiate and assign space to users and manage *lifetime* of spaces
 - **Manage files on behalf of user**
 - Pin files in storage till they are released
 - Manage *lifetime* of files
 - **Manage file sharing**
 - Policies on what should reside on a storage or what to evict
 - **Bring the files from remote locations**
 - **Manage multi-file requests**
 - a brokering function: queue file requests, pre-stage

XROOTD+SRM integration overview

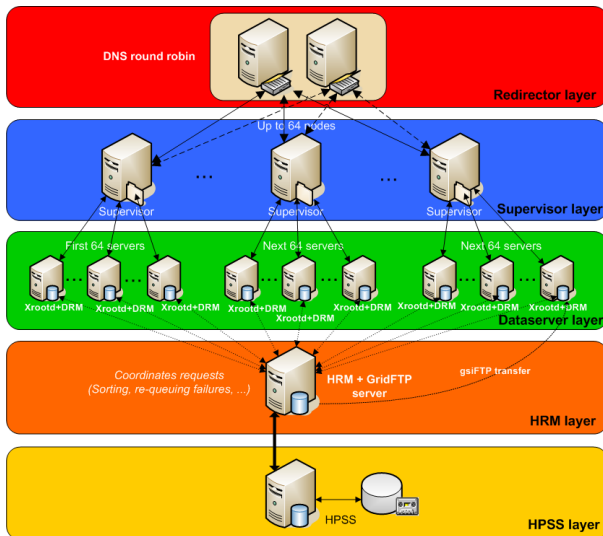
Types of storage resource managers:

- **Disk Resource Manager (DRM)**
 - Manages one or more disk resources
- **Tape Resource Manager (TRM)**
 - Manages the tertiary storage system (e.g. HPSS)
- **Hierarchical Resource Manager (HRM=TRM+HRM)**
 - An SRM that stages files from tertiary storage into its disk cache



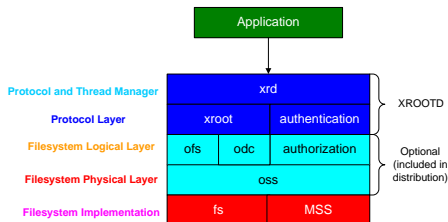
- **xrootd** is responsible for managing the disk **cluster** (aggregation, load balancing ...)
- **DRM** is responsible for managing the **disk cache**
- **HRM** is responsible for managing access to **HPSS**

XROOTD+SRM cluster overview



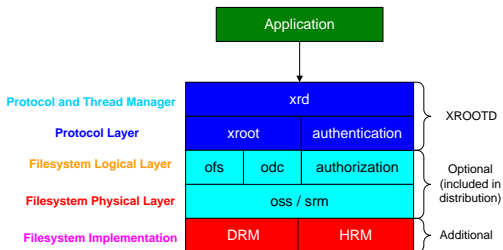
XROOTD components architecture

- xrootd architecture is very amenable to extensions (divided into several components and plug-ins)



- xrd** - provides networking support, thread management and protocol scheduling
- ofs** - provides enhanced first level access to file data (responsible for coordinating activities of oss, odc, auth)
- oss** - provides access to underlying storage system (controlled by ofs and invokes meta-data operations)

XROOTD+SRM components architecture



- oss component was externalized as a plugin
- several existing methods were virtualized (Create, Open, Close, Stage)
 - easy wrapping of concrete class implementation
- for example:
 - Create() uses DRM to create a file
 - Close() informs that the file is no longer in use

Current status ?

- The xrootd-SRM activity is a collaboration between:
 - **BNL/STAR** - proposed the idea, providing early demonstrations, performing measurement, provide enhancements
 - **SLAC** - providing changes to xrootd, providing distribution mechanism with xrootd
 - **LBNL/SDM** - providing plug-ins that interact with SRMs and changes to DRM for new functionality
- working version:
 - supports read, write mode into HPSS with all SRM functionalities as pinning, allocation etc.
 - thanks to **A. Romosan** and **A. Sim**
 - supports **one cache per a node** ⇒ multiple caches per one node has been recently developed (under testing)
 - planning to test it in **large scale** at BNL after multiple caches support and later production mode

Summary

- Xrootd is currently deployed on almost **500 nodes** (the biggest production deployment of xrootd) serving over **340TB** of distributed disk space
- load balancing and handshake with tape system make the system resilient to failures
- the system is used by users for **daily analysis**
- measurement of xrootd **aggregate IO** showed **competitive results** comparing to commercial Panasas(NFS)
- **large scale** testing of SRM+XROOTD is on the way
- planning a **regression tests** (performance comparison with previous version)
- planning **enhancements** such as **cache-to-cache** transfers and **interoperability** with other SRM-aware tools (such as DataMovers etc.)

References



P. Jakl, J. Lauret, A. Hanushevsky, A. Shoshani, A. Sim

From rootd to xrootd: From physical to logical file

Proc. of Computing in High Energy and Nuclear Physics (CHEP'06)



P. Jakl, J. Lauret, M. Šumbera

Managing widely distributed data—sets

Czech Technical University, FNSPE, Research report, 2006

<http://www.star.bnl.gov/~pjakl>



A. Romosan, A. Hanushevsky

XROOTD—SRM

SRM—Collaboration meeting, CERN, 2006