



DSS

Data & Storage Services

CERN IT
Department

CASTOR and EOS status and plans

Giuseppe Lo Presti

on behalf of CERN IT-DSS group

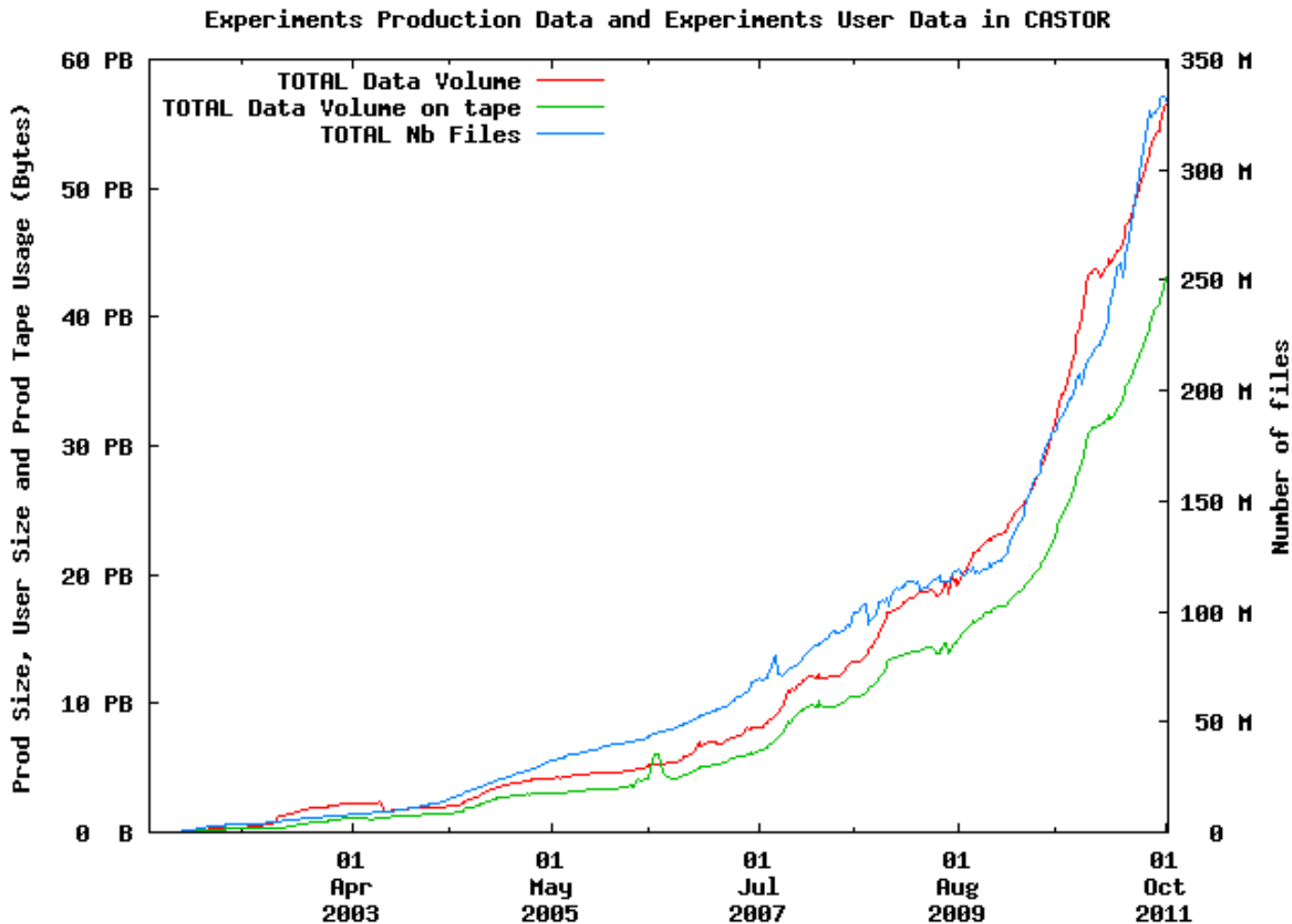


- CASTOR and EOS strategies
- CASTOR status and recent improvements
 - Disk scheduling system
 - Tape system performance
 - Roadmap
- EOS status and production experience
 - EOS Architecture
 - Operations at CERN
 - Roadmap/Outlook



Strategy:

- Keep Tier0/production activity in CASTOR
 - Not necessarily only tape-backed data
 - Typically larger files
 - Focus on tape performance
- Moving xroot-based end-user analysis to EOS
 - Disk-only storage
 - Focus on light(er) metadata processing



Generated Oct 11, 2011 CASTOR (c) CERN/IT

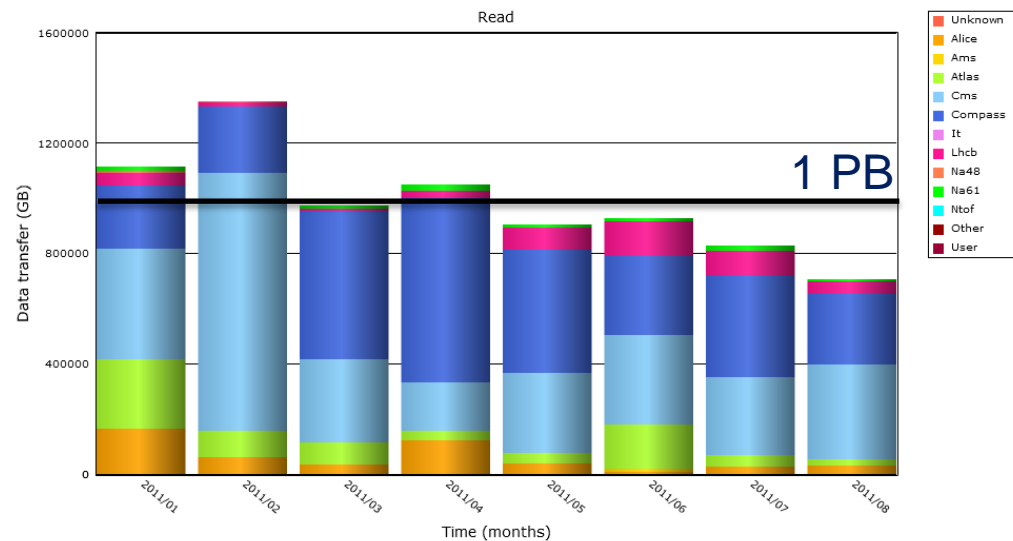
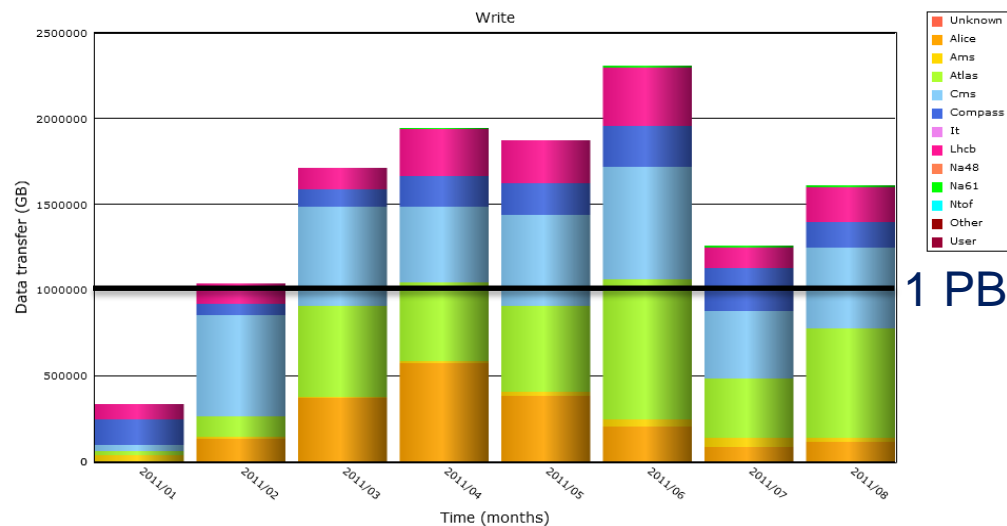
55 PB of data

320M files

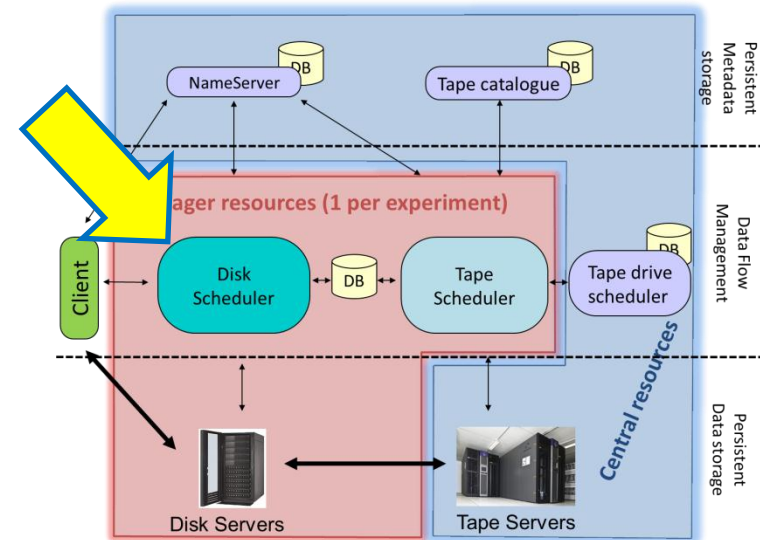
Peak writing speed: 6GiB/s
(Heavy Ion run, 2010)

Infrastructure:

- 5 CASTOR stager instances
- 7 libraries (IBM+STK), 46K 1TB tapes, ~5K 4TB or 5TB tapes
- 120 enterprise drives (T10000B, TS1130, **T10000C**)

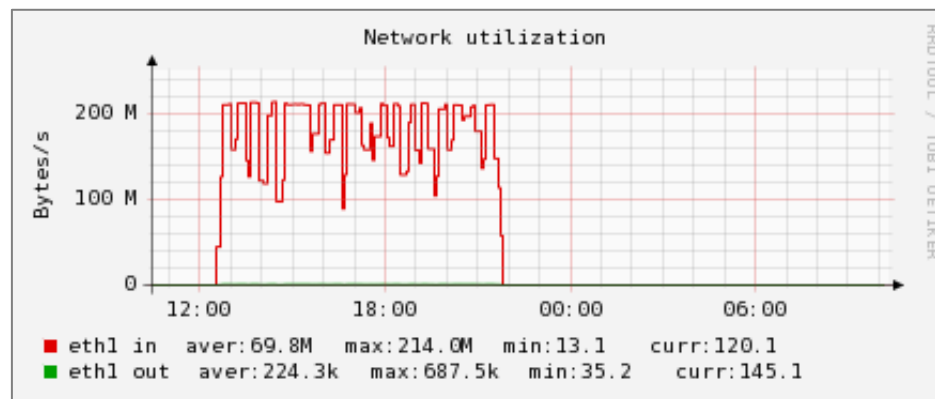


- Transfer Manager, replacing LSF in CASTOR
- Stress tested
 - Performances ~10x higher than peak production levels
 - Production throttled at 75 Hz (25 Hz per node)
- In production in all instances at CERN and at ASGC
 - Staged roll-in:
first ATLAS, then CMS,
then everybody else
 - Current release includes fixes
for all observed issues,
smooth operations since then



- Improving read performance
 - Recall policies already in production since ~1 year
- Improving write performance
 - Implemented **buffered Tape Marks** over multiple files
 - Theoretically approaching drive native speed regardless file size
 - Practically, different overheads limit this
 - Soon available for wide deployment
 - Currently being burned-in on a stager dedicated to Repack operations
 - Working on simplifying and optimizing the stager database, by using bulk interfaces
 - Expected timeframe for production deployment: **spring 2012**

- Measuring tape drive speed
 - Current data rate to tape: 60-80 MiB/s
 - Dominated by the time to flush the Tape Mark for each file
 - Average file size ~200 MB
 - Preliminary tests with an STK T10000C
 - Tape server with 10GigE interface
 - 195 MiB/s avg.
 - **214 MiB/s peak**



- Towards fully supporting small files
 - Buffered Tape Marks and bulk metadata handling
 - In preparation for the next repack exercise in 2012 (~40 PB archive to be moved)
- Further simplification of the database schema
 - Still keeping full consistency approach, No-SQL solutions deliberately left out
- Focus on operations

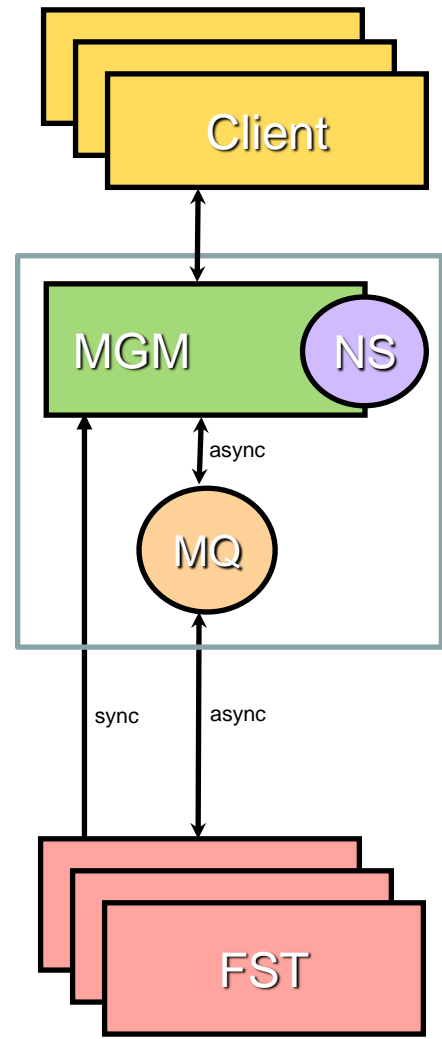
- CASTOR and EOS strategies
- CASTOR status and recent improvements
 - Disk scheduling system
 - Tape system performance
 - Roadmap
- **EOS status and production experience**
 - EOS Architecture
 - Operations at CERN
 - Roadmap/Outlook

- Easy to use standalone **disk-only** storage for user and group data with **in-memory** namespace
 - **Few ms** read/write open latency
 - Focusing on end-user analysis with chaotic access
 - Based on **XROOT** server plugin architecture
 - Adopting ideas implemented in *Hadoop*, *XROOT*, *Lustre* et al.
 - Running on low cost hardware
 - no high-end storage
 - Complementary to CASTOR

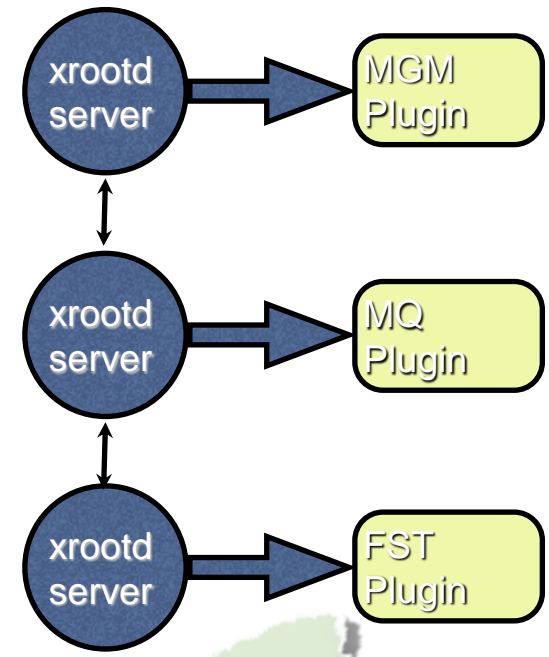
Management Server
Pluggable Namespace, Quota
Strong Authentication
Capability Engine
File Placement
File Location

Message Queue
Service State Messages
File Transaction Reports
Shared Objects (queue+hash)

File Storage
File & File Meta Data Store
Capability Authorization
Check-summing & Verification
Disk Error Detection (Scrubbing)



Implemented as plugins in **xrootd**



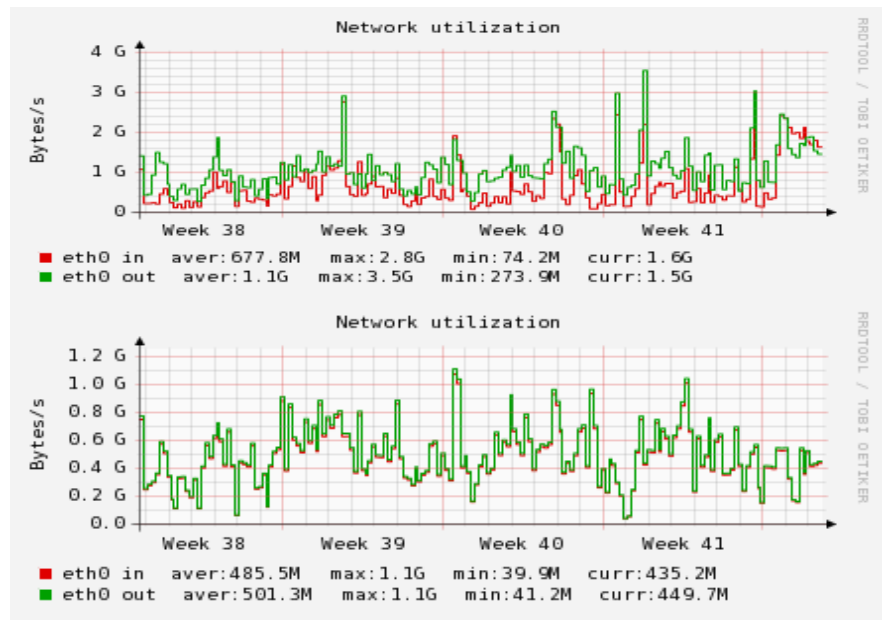
No DB Backend required!

- EOS uses XROOT as primary file access protocol
 - The **XROOT framework** allows flexibility for enhancements
- Protocol choice is not the key to performance as long as it implements the required operations
 - **Client caching matters most**
 - Actively developed, towards full integration in ROOT
- SRM and GridFTP provided as well
 - BeStMan, GridFTP-to-XROOT gateway

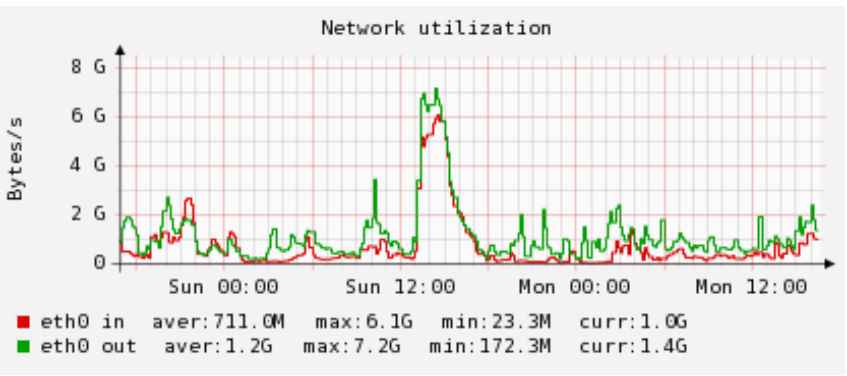
- Storage with single disks (JBODs, no RAID arrays)
 - redundancy by s/w using cheap and unreliable h/w
- Network RAID within disk groups
 - Currently file-level replication
- Online file re-replication
 - Aiming at reduced/automated operations
- Tunable quality of service
 - Via redundancy parameters
- **Optimized for reduced latency**
 - Limit on namespace size and number of disks to manage
 - Currently operating with **40M** files and **10K** disks
- Achieving additional scaling by partitioning the namespace
 - Implemented by deploying separated instances per experiment

- Failures don't require immediate human interventions
 - Metadata server (MGM) failover
 - Disks drain automatically triggered by I/O or pattern scrubbing errors after a configurable grace period
- Drain time on production instance < 1h for 2 TB disk (10-20 disks per scheduling group)
 - Sysadmin team replaces disks 'asynchronously', using admin tools to remove and re-add filesystems
 - Procedure & software support is still undergoing refinement/fixing
- **Goal:** run with best effort support

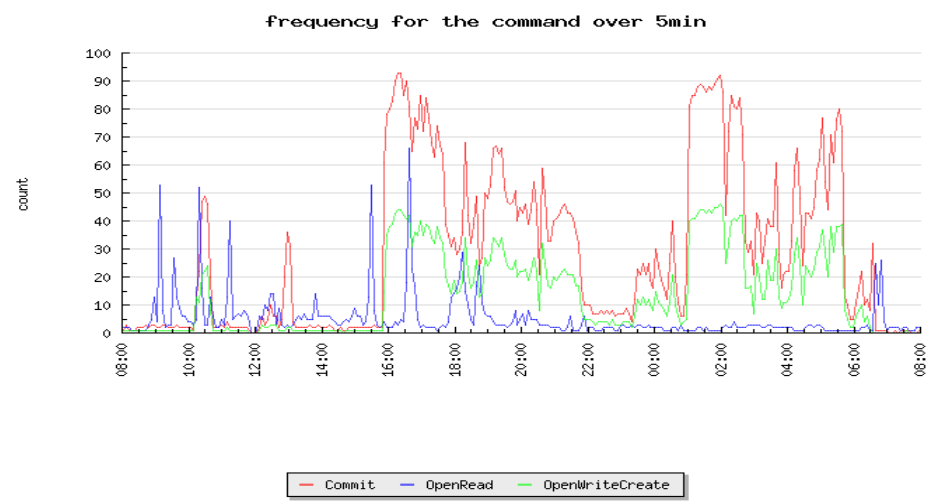
- Field tests done (Oct 2010 – May 2011) with ATLAS and CMS, production since summer
- EOS 0.1.0 currently used in EOSCMS/EOSATLAS
 - Software in bug-fixing mode, frequent releases though
- Pools migration from CASTOR to EOS ongoing
 - Currently at **2.3 PB** usable in CMS, **2.0 PB** in ATLAS
 - Required changes in the experiment frameworks
 - User + quota management, user mapping
 - Job wrappers
 - Etc.
 - Several pools already decommissioned in CASTOR
 - E.g. CMSCAF



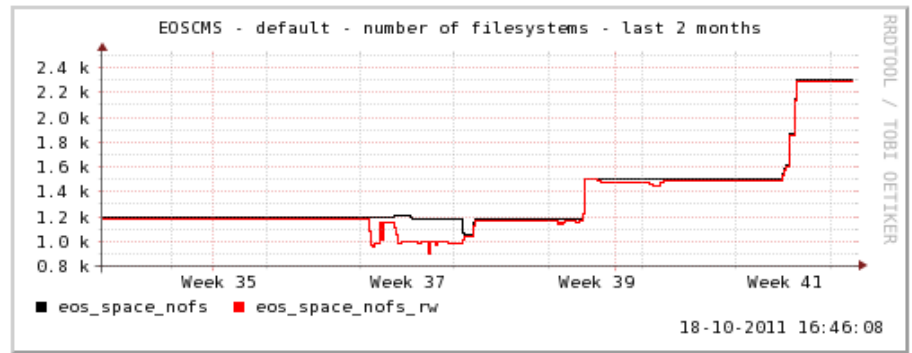
ATLAS instance: throughput over 1 month (entire traffic & GridFTP gw)



Pool throughput during a node drain



ATLAS instance: file ops per second



CMS instance: hardware evolution

- EOS 0.2.0 expected by end of the year
- Main Features
 - File-based redundancy over hosts
 - Dual Parity Raid Layout Driver (4+2)
 - ZFEC Driver (Reed-Solomon, N+M, user defined)
 - Integrity & recovery tools
 - Client bundle for User EOS mounting (krb5 or GSI)
 - MacOSX
 - Linux 64bit

- CASTOR is in production for the Tier0
 - New disk scheduler component in production
 - New buffered Tape Marks soon to be deployed
- EOS is in production for analysis
 - Two production instances running
 - result of very good cooperation with experiments
 - Expand usage and gain more experience
 - Move from fast development and release cycles to reliable production mode