

# Fabric Infrastructure and Operations



# Experiment top 5 issues: CASTOR Status & Plans

#### GDB - May 2<sup>nd</sup> 2007

**Tony Cass** 



CERN - IT Department CH-1211 Genève 23 Switzerland WWW.cern.ch/it



# FIO Summary by Experiment



Experiment	Issue 1	Issue 2	Issue 3	Issue 4
CMS	Single request queue bottleneck	Priorities must work; need for capacity	Performance at CERN & Tier1s	Disk mover limitations
ATLAS	Tier0 shielded from users	SRM v2.2 availability	quotas	
LHCb		"Resource busy" from CASTOR	SRM v2.2 availability	
Alice	xrootd interface	Predictable latency for tape recall		

CERN - IT Department CH-1211 Genève 23 Switzerland www.cern.ch/it





## Summary by Issue



• Tier-0

- Stability & Performance, isolation from general users
- Tier-1
  - Stability & Performance
- SRM v 2.2 interface availability
  - including VOMS integration
- xrootd interface availability
- Quotas
- Predictable latency

CERN - IT Department CH-1211 Genève 23 Switzerland www.cern.ch/it





**CERN - IT Department** CH-1211 Genève 23 Switzerland www.cern.ch/it

### Tier-0 issues – I

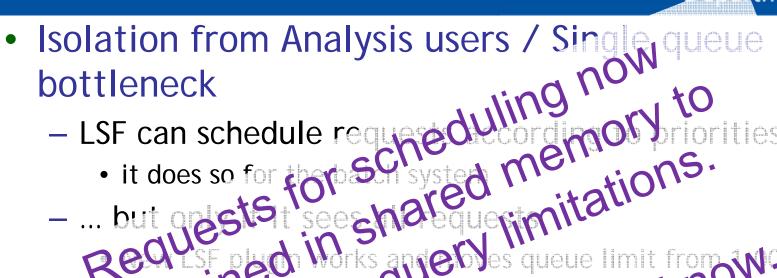


- Stability & Performance
  - Major problems occur if CASTOR cannot scherele new requests
    queue builds up and stager brack hown.
  - New LSF plugin improv le larce ce with queues.
    - - - 006 due to gs. All now fixed.
          - quests cannot be satisfied
            - Generally We to use of *disk1* pools; see later
          - Slow or unresponsive database server
            - main reason for ATLAS problems in March; improvement seen when memory added. New servers being introduced.





#### Tier-0 issues – II



CERN

- ... but onsoloti 2 content pluto to at post 30,0
- e separate nel stagers
  - efficiency (disk pools & tape drives). • but with reduced

**CERN - IT Department** CH-1211 Genève 23 Switzerland www.cern.ch/it



rtment



#### **CERN - IT Department** CH-1211 Genève 23 Switzerland www.cern.ch/it

### Tier-1 Issues



- Stability & Performance
  - As for Tier0, but with better version control for outside sites.
    - An improved build, test and deployment procedure is in place. However, CVS changes and increased modularity are not a priority compared to some other issues. Multiple version support not before end-2007.

#### Improved Disk1 support

- Highest priority after Tier-0 isola
  - Many issues already addresses load balancing, issue
- Recurs risk of a Recu wanted files searching through ify unwanted and files to delete.





CERN - IT Department CH-1211 Genève 23 Switzerland www.cern.ch/it

# Other



- SRM v2.2 Interface Availability
  - CASTOR interface passes all basic and use case tests.
  - Stress testing can start, but test stager instance at CERN is currently devoted to tests of the new LSF plugin in the context of the ATLAS taskforce.
  - Strong Authentication and VOMS integration next priority after disk1 issues
    - but no deployment in production before Q2 2008 (and probably later).

#### • Quotas

- Introducing support for quotas is not a priority
  - so nothing before end-2008 or 2009
  - and, even if done, quota checking performed only at moment of initial write request.

#### Predictable Latency

- Predictable latency for access to disk resources (a CMS requirement) believed achievable with scheduler interface.
- Predictable latency for access to files on tape is not achievable given likely resource levels
  - would require dedicated read-only drives per stager (experiment)
  - Production reload of data from tape should be managed
  - No guarantees for reload of user data from tape. Indeed, user requests (especially for small files) may have to be held until we have enough requests for a given volume.



## xrootd Interface Availability



- The xrootd interface to Castor is developed by SLAC, not CERN.
  - CERN made some slight changes to the underlying CASTOR software last year, but the key interface is developed by Andy Hanus
- We believe the intole? w working
  - but extensive testing is needed to be sure that ne likeľy
- overall Castor/xrootd combination in all use cases and can support (m load. ong term support for tols inter agreed agreed are portunded for this, and terface has yet to to provide support but
  - port currently depends on one person.



**CERN - IT Department** CH-1211 Genève 23 Switzerland www.cern.ch/it