



Science & Technology
Facilities Council

RAL Site Report

Castor F2F, CERN

Matthew Viljoen



GridPP

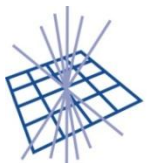
UK Computing for Particle Physics



Stuff to cover

(at some point)

- **General Operations report (Matthew)**
- **Database report (Rich)**
- **Tape report (Tim)**
- **Disk-only solution investigations (Shaun)**





Current group structure

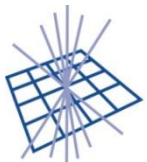
Scientific Computing Department (Adrian Wander)

Data Services group (Jens)

DB Team (Rich) and CASTOR Team (Matt)

Running
CASTOR

Matt 70% mgnt/ops
Chris 50% ops
Shaun 25% ops & problem solving
Rob 100% ops, CIP
Jens 10% CIP and Tim 50% Tape
... plus effort from DBAs
... and 1 FTE from Fabric Team
... and 0.5 FTE from Production Team





Current status

CASTOR

Prod Instances (all 2.1.12-10 + 2.1.11-9 NS):

- **Tier 1:** ATLAS, CMS, LHCb
and Gen (ALICE, H1, T2K, MICE, MINOS, ILC β)
- **Facilities:** DLS, CEDA

Test Instance (2.1.12/2.1.13):

- Preprod + Cert

CIP homegrown in LISP

DB Details from Rich...

Tape SL8500 with 2,453xA, 4,655xB, 1,998xC tapes for Tier 1





Setup

- **Headnodes**

10 SRMs (4 ATLAS, 2 CMS, 2 LHCb, 2 Gen)

“Stager/Scheduler/DLF” per instance and 2 NS

- **SRM** (4 ATLAS, 2 LHCb/CMS/Gen running 2.11)

- **Disk servers** (approx. 443)

10-40TB, RAID6 ext4,XFS

- **DB** Details later from Rich...

- **CIP** homegrown in LISP





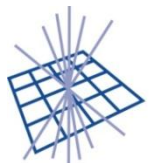
Stats (Nov '12)

VO	Disk (used/total)	Tape (used)
ATLAS	3.4/4.2PB	2.6PB
CMS	1.2/1.9PB	3.7PB
LHCb	1.3/2PB	1.2PB
Gen	0.3/0.4PB	0.85PB
Facilities	(no D1T0)	1.5PB (inc. D0T2)

Total used:

6.2PB

11.5PB



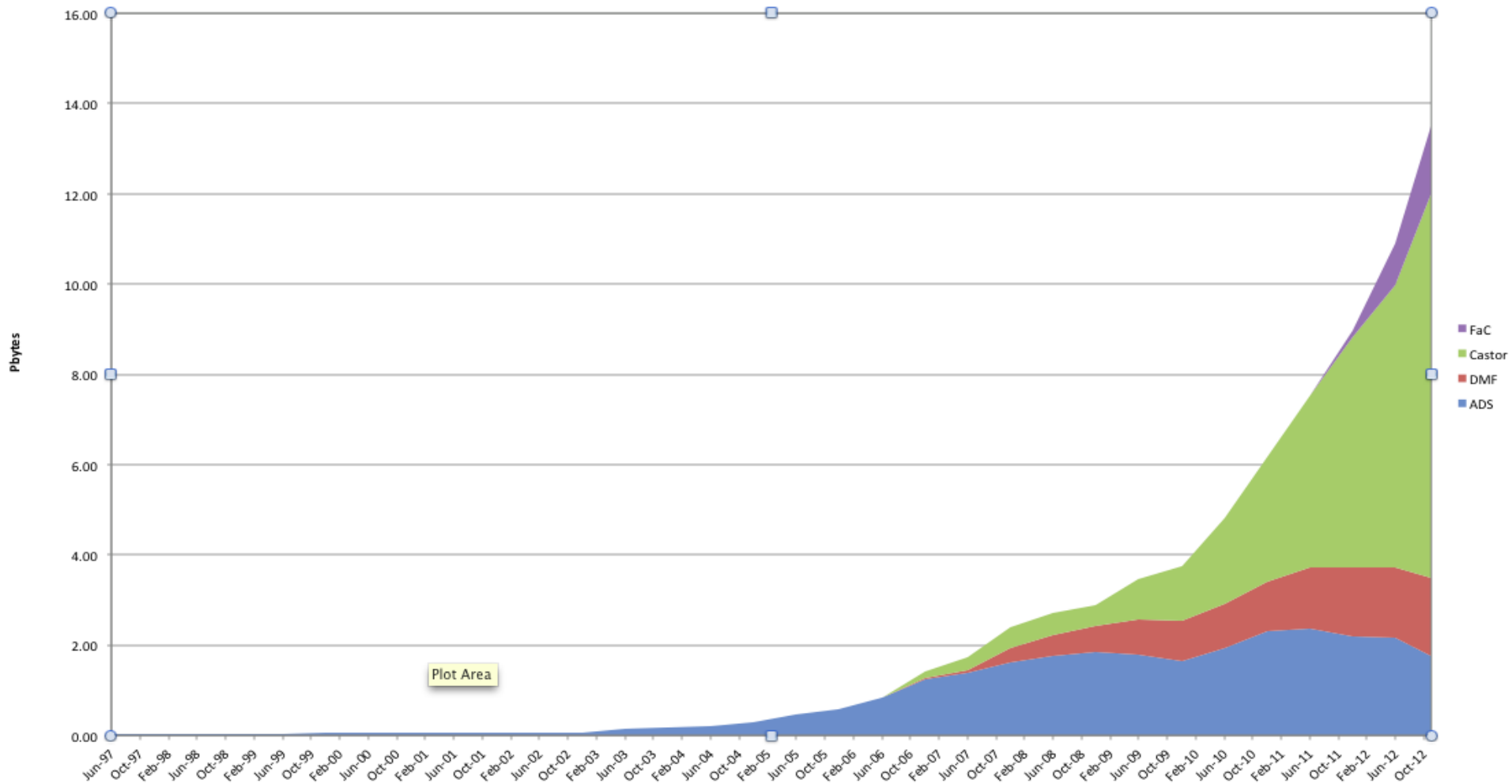
GridPP

UK Computing for Particle Physics



SL8500 usage (All RAL)

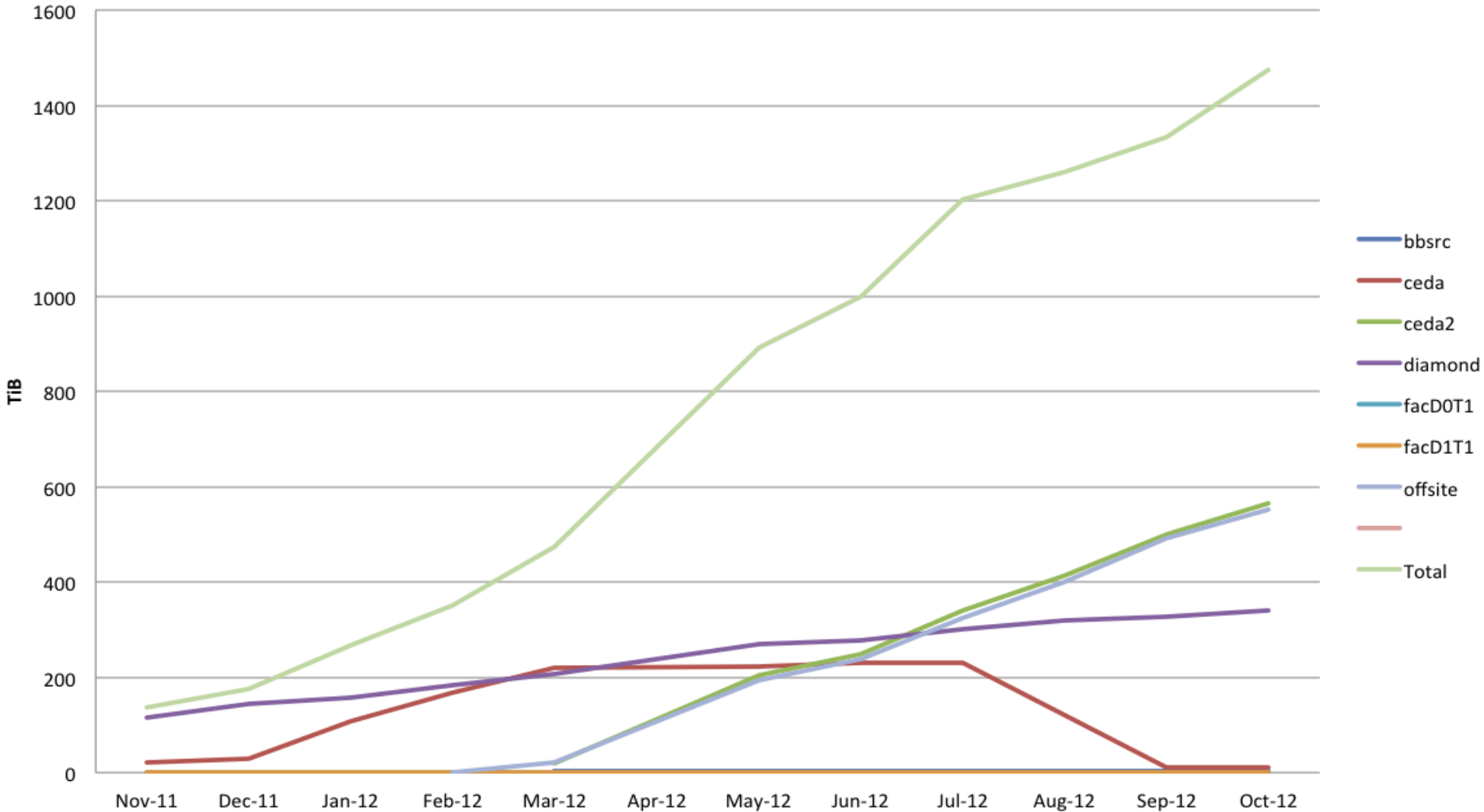
Total Holdings





Tape usage (Facilities)

FaC Totals





Recent changes

- **(2011/2012)Complete hardware refresh this year**
- **12Q2 Minor 2.1.11-* upgrades and Switch tape subsystem to Tape Gateway**
 - **Switch from LSF on all instances to Transfer Manager**
 - **No more licensing costs!**
 - **Better performance, and...SIMPLER!**
- **12Q3 Repack for LHCb 2300 A -> C**
- **12Q4 Major 2.1.12-10 stager upgrade**
 - **Introduction of global federated xroot for CMS, ATLAS**





Hardware Refresh

- **New SRMs, CASTOR + DB headnodes**
- **SL5 and Configuration Management System (CMS) - Quattor + Puppet - control throughout**

Leading to:

- **Improved overall performance**
 - **Switch over availability stats from SAM Ops to VO**
 - **No more ATLAS background noise in SAM tests (before, consistent <5% of miscellaneous ATLAS failures)**
- **Content Mgmt System (CMS) adoption (installation, DR, no more need to rely on system backups)**





CMS – a few words

- Before we relied on backups, now on re-installation
- A node can be reinstalled in <1hr
- Tier 1 solution is Quattor, supported by Fabric Team. CASTOR has always used Puppet (for config files).
- Now we use a Quattor/Puppet hybrid

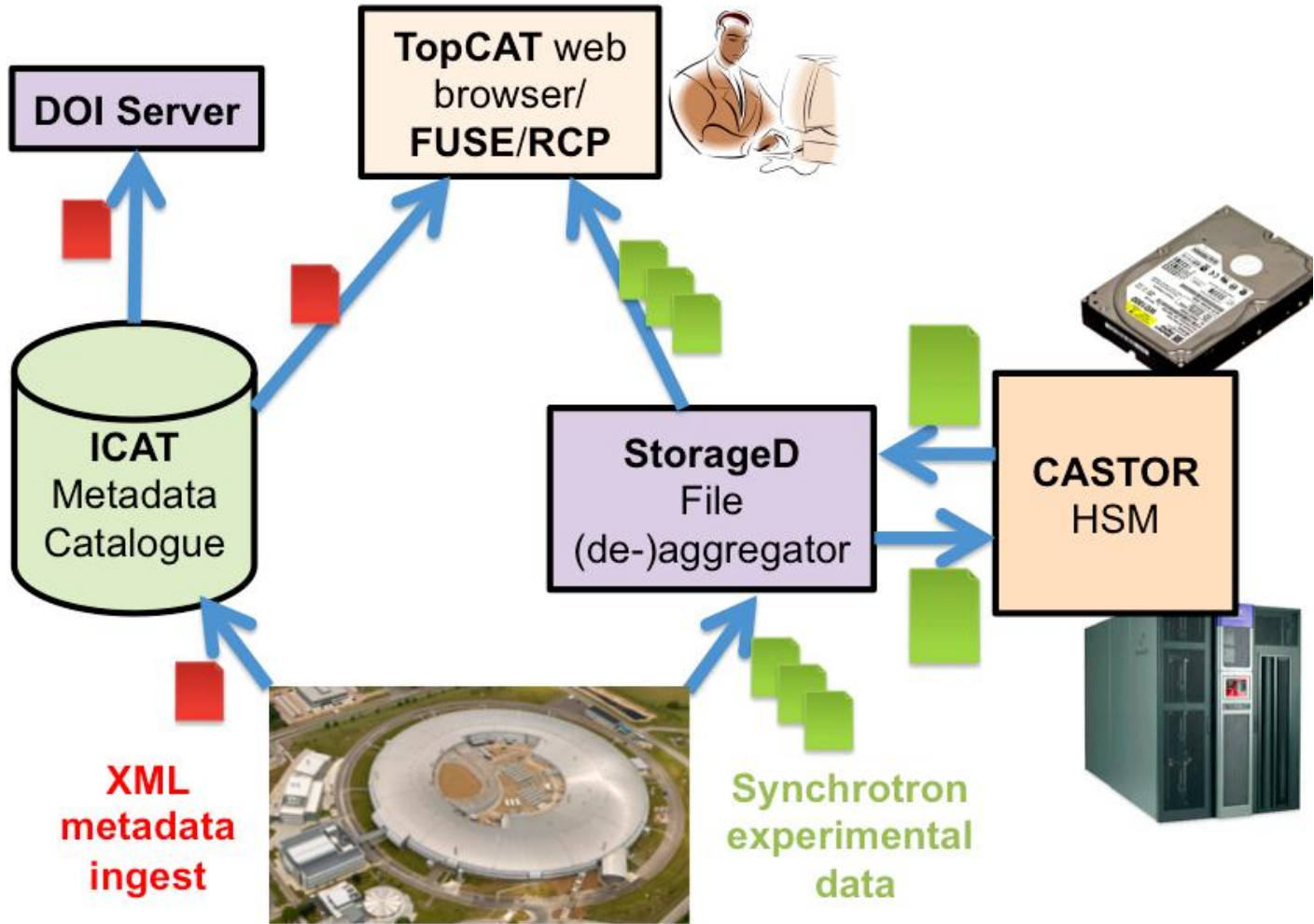
Content type	Examples
OS payload	glibc-2.5-81.el5_8.4, vim-enhanced-7.0.109-6.el5 RPMs
OS-level config	resolv.conf, crontab, admin accounts/keys
CASTOR payload	castor-stager-server-2.1.12-10, castor-vmgr-client-2.1.12-10 RPMs
CASTOR config	castor.conf, tnsnames.ora

Quattor
Puppet





Facilities Data Service Architecture for DLS





Remaining problem areas

- **Disk server deployment and decommissioning overheads**
Can extend automation with CMS
-> Shouldn't need to manually "bless" disk servers
- **Ongoing need for ORACLE database expertise**
Large number of different instances (4 prod, 3 test, Facilities...)
- **Lack of read-only mode with new scheduler**
- **Lack of disk server balancing**
- **Monitoring (based on Ganglia) currently difficult to use. CK looking at new solutions**
- **Quattor clunkiness and general reluctance to use it**
Better templates structure should help
Aqualon?

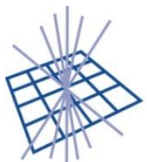




SIRs affecting CASTOR over last year

<https://www.gridpp.ac.uk/wiki/Category:Incidents>

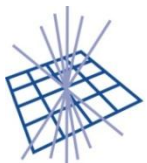
- 2011/10/31 Castor ATLAS Outage Bad ORA execution plan
- 2011/12/02 VO Software Server
- 2011/12/15 Network Break Atlas SRM DB ORA RM bug, now fixed
- 2012/03/16 Network Packet Storm
- 2012/06/13 Oracle11 Update Failure
- 2012/11/07 Site Wide Power Failure Power supply outage
- 2012/11/20 UPS Over Voltage Power intervention gone wrong





What next?

- **Full “off-site” database Dataguard backup**
- **Common headnode type, for improved:**
 - **Resiliency: easier to replace faulty node**
 - **Scalability: dynamically changing pool of headnodes**
 - **Doubling up daemons wherever possible**
 - > **Better Uptime – e.g. applying errata/rebooting etc.**
- **2.1.13 +SL6 upgrade in new year**





Further ahead...

- **Using virtualization more...**
Cert instance already virtualized
Virtualize by default (headnodes, tape servers, CIPs...)
VTL?
- **2013: New disk only solution alongside CASTOR**
Higher performance for analysis, easier to run
- **IPv6?**





To conclude...

CASTOR nice and stable nowadays

Rigorous change control at Tier 1 also helps!

Track record of good interventions

Comprehensive testing infrastructure paying dividends

Balance right between new functionality vs. stability

3-6 months trailing behind CERN head version

Good performance (esp. for tape). No plans to move away from CASTOR, alongside new “next-gen” disk storage solution

