



GridPP

UK Computing for Particle Physics

RAL Site Report

HEPiX Spring 2011, GSI

2-6 May

Martin Bly, STFC-RAL



Science & Technology Facilities Council

e-Science

- STFC Stuff
- RAL Stuff
- Building stuff
- Tier1 Stuff



- UK Govt Comprehensive Spending Review (CSR)
 - Was 2010-10-20
 - General: Level funding for Core Science, i.e., no increase with inflation
 - The CSR settlement allows STFC to pursue the high priority Science programme it outlined as a result of the 2009 prioritisation. In particle physics and astronomy, this was confirmed by PPAN following the CSR
- GridPP is an STFC project
 - Within the STFC programme, GridPP was rated as alpha-5, the highest priority, along with CMS and ATLAS (and other stuff)
 - The T1 is a high priority within GridPP
 - GridPP funding for the next 4 years (to 2015) has now been confirmed

- Email Addressing:

- Removal of old-style f.blogs@rl.ac.uk email addresses in favour of the cross-site standard fred.bloggs@stfc.ac.uk
 - (Significant resistance to this)
- No change in aim to remove old-style addresses but...
- mostly via natural wastage as staff leave or retire
- Staff can ask to have their old-style address terminated

- Exchange:

- Migration from Exchange 2003 to 2010 went successfully
 - Much more robust with automatic failover in several places
 - Mac users happy as Exchange 2010 works directly with Mac Mail so no need for Outlook clones
- Issue for exchange servers with MNLB and switch infrastructure
 - Providing load-balancing
 - Needed very precise instructions for set up to avoid significant network problems

- UPS problems
 - Leading power factor due to switch-mode PSUs in hardware
 - Causes 3KHz ringing on current, all phases (61st harmonic)
 - Load is small (80kW) compared to capacity of UPS (480kVA)
 - Most kit stable but EMC AX4-5 FC arrays unpredictably detect supply failure and shut down arrays
 - Previous possible solutions abandoned in favour of:
 - Local isolation transformers in feed from room distribution to in-rack distribution: Works! 😊

- Site

- Sporadic packet loss in site core networking (few %)
 - Began in December, got steadily worse
 - Impact on connections to FTS control channels, LFC, other services
 - Data via LHCOPN not affected other than by control failures
- Traced to traffic shaping rules used to limit bandwidth in firewall for site commercial tenants. These were being inherited by other network segments (unintentionally!)
- Fixed by removing shaping rules and using a hardware bandwidth limiter

- LAN

- Issue with a stack causing some ports to block access to some ip addresses: one of the stacking ports on the base switch faulty
- Several failed 10GbE XFP transceivers

- Summary of previous report:
 - 36 SuperMicro 4U 24-bay chassis with 2TB SATA HDD (10GbE)
 - 13 x SuperMicro Twin²: 2 x X5650, 4GB/core, 2 x 1T HDD
 - 13 x Dell C6100: 2 x X5650, 4GB/core, 2 x 1T HDD
 - Castor (Oracle) databases server refresh: 13 x Dell R610
 - Castor head nodes: 16 x Dell R410
 - Virtualisation: 6 x Dell R510, 12 x 300GB SAS, 24GB RAM, 2 x E5640
- New since November
 - 13 x Dell R610 tape servers (10GbE) for T10KC drives
 - 14 x T10KC tape drives
 - Arista 7124S 24-port 10GbE switch + twinax copper interconnects
 - 5 x Avaya 5650 switches + various 10/100/1000 switches

- One of two batches of the FY09/10 capacity storage failed acceptance testing: 60/98 servers (~2.2PB) ☹
 - Cards swapped (LSI -> Adaptec)
 - Acceptance testing completed
 - Released for production use
- After problems with one of the two batches of the FY08/09 capacity during commissioning (now resolved), the other batch has had issues in production resulting in data loss:
 - Single-drive throws cause array lock up and crash (array loss)
 - Whole batch (50/110) rotated out of production (data migrated)
 - Updated array firmware
 - Recreate arrays from scratch, new file systems
 - Undergoing hammer test
 - Eight drive throws in 3 months successfully handled

- Castor manages disk and tape storage
 - 12 million files (at March 2011)
 - Used/Total capacities: 3.2PB/5.2PB on tape and 3.2PB/7.5PB on disk
- Recent news:
 - Major upgrade during late 2010 introducing:
 - Checksums for all files, xrootd support, proper integrated disk server draining
 - Disk servers migrated to SL5/64bit with XFS capability
 - New (non-Tier1) production instance for Diamond synchrotron
- Coming up:
 - T10KC drive and tape media support
 - Need to update Library microcode
 - Need latest version of Castor to use these drives
 - Castor v2.1.10-1 for tape servers and some backend services
 - Move to new database hardware and better resilient architecture (using Oracle DataGuard) later this year
 - New service 'head nodes'

- Evaluating MS Hyper-V (inspired by CERN's successes) for services virtualization platform
 - Offers sophisticated management/failover etc without punitive cost of VMWare
- However as Linux admins, sometimes hard to know if problems are due to ignorance of the MS world 😊
- Struggled for a long time with iSCSI storage arrays (and poor support)
 - abandoned them recently and problems seem resolved
- Have learnt a lot about administering Windows servers....
- Ready to implement production platform

- Quattor
 - Batch and Storage systems under Quattor management
 - ~6200 cores, 700+ systems (batch), 500+ system (storage)
 - Significant time saving
 - Significant rollout on Grid services node types
- CernVM-FS
 - Major deployment at RAL to cope with software distribution issues
 - Details in talk by Ian Collier later this week
- Network future
 - Beginning to look at provision for Tier1 core network to continue to meet increasing data bandwidth requirements and resilience
 - Various mesh structures using lower cost components are attractive

Questions?



02/05/2011

RAL Site Report - HEPiX Spring 2011



Science & Technology Facilities Council
e-Science



- .
- Rollout of new hardware for services nodes
- Database architecture