

# Lancaster Site Report

HEPSYSMAN 2011

Matt Doidge, Alex Finch & Peter Love

-more bullet points than Rambo's belt.

# The Team

- Rob Henderson, Alex Finch, Peter Love and myself, with our powers combined, provide support to the Lancaster HEP group.
  - The Griddy, Tier-2 side of things is governed by Peter, Rob and I.
  - Peter, Alex and Rob oversee the local cluster.
    - The loci of responsibilities are blurring.
- Roger Jones keeps us all in check.

# Local Cluster for Local People.

- ~20 clients, mostly running SL5.5 (although a few stragglers).
- Services include:
  - Web server, 20TB Fileserver, Backup services (Bacula + 4 AIT tapes in a 16-slot library).
  - Mail & NIS, local batch system (torque),
  - Print server (cups/samba) and printers.

# Tier Three-y stuff

- 18TB nfs server (same hardware as a poolnode).
- UIs (on VMs)
- “cvmfs” to give users access to atlas software areas.
- “Proof” cluster- 4x Dual X5650s in a Dell C6100 chassis, 48GB RAM and 12 core's per box, all running xproofd.
  - It frikkin flies.
  - Users love it.

# Tier-2: A Tale of Two Clusters

- The “old”:
  - 512 cores (~5000 HEPSPEC06), behind an lcg-CE, torque/maui batch system. NFS software area.
  - All ours.
- The “new”:
  - Shared HEC cluster, 1900 cores (~24k HEPSPEC06), behind a CREAM CE. LSF batch system.
  - Shared admin with ISS staff.
- Both use (separate) WN tarball installs and mounted software, nfs mounted on the old, panfs mounted on the new.

# Other Tier-2 Gumpf.

- DPM Storage Element, ~600TB storage online.
  - ~350TB of storage to come online.
    - Had one batch of storage with a lot of disk failures, Western Digital are getting involved.
  - All online pool nodes are now the 24-bay Supermicros, with a range of Raid Cards (Areca, Adaptec & 3ware).
  - Headnode is an overspecced beast.
    - 8-cores, 24GB RAM, /var on fast SAS's in a raid 10.

# Other Tier-2 Gumpf

- Like others we're moving to VMs for a bunch of services
  - As a lot of grid services are low-load but don't play nice with others.
  - Virtualised using KVM
    - Nothing exciting, the most complicated thing we did was bond each of the 2 virtual NICs to corresponding ones on the host.
- Very few SL4 services left, just the lcg-CE/torque server left.

# Recent CREAM Badness.

- Our CREAM CE recently kept dying with crazy high load coming from the mysql daemon.
  - Mysql tuning, Job Purging, cutting down on submissions from older condor clients and reducing the Blah Updater interval all helped stem the tide, but eventually things always overwhelmed whatever we took.
- And are job table was always bloated:
  - `select count(distinct jobId) from job_command;`  
25771  
1 row in set (9.45 sec)



# The Fix.

- We had found a (known) bug:  
<https://savannah.cern.ch/bugs/index.php?83749>
- 8 weeks ago we had moved our sandbox area.
  - This prevented all of the (20,000) jobs that had been run in the busy 10 days before the move from being deleted.
  - Worse the cream kept on trying to create leases for these jobs, expanding their DB footprint.
- The fix was to run the Job Purger using a hacked config that pointed to the old sand box area. Annoyingly simple.

# The Query that kicked us in the Cream:

```
select jstd.name, count(*) from job, job_status_type_description jstd,  
job_status AS status LEFT OUTER JOIN job_status AS latest ON  
latest.jobId=status.jobId AND status.id < latest.id WHERE latest.id IS  
null and job.id=status.jobId and jstd.type=status.type group by  
jstd.name;
```

# Zanshin

- We monitor things using central syslogging, logwatch mails, (local) nagios, ganglia, and cacti for the networking.
  - Plus all the external griddy monitoring: Panda, Steve's pages, UK Nagios, pilot factory logs, a host of other atlassy things.
- We use cfengine/pxe/kickstart/yaim to install and manage the “old” machines and kusu/yaim for the HEC machines.

# Networking.

- Both the local user and grid site are on separate, routable subnets
  - The Tier-2 is actually split across two, the “Grid” and the “HEC”.
  - The Tier-2 is further split between a private and public VLANs.
  - There is a 10G backbone to the Tier 2 network.
  - There is a 1-Gb dedicated light path from the Tier-2 to RAL, and we also share the University's link to Janet (although I believe we are capped at 1Gb/s).
  - All switches are managed by ISS.
  - We have access to the University Cacti pages.