



Birmingham Status and Plans

HEPSYSMAN Meeting, 14th June, 2017
Mark Slater, Birmingham University

As of right now, the Birmingham Grid Site consists of:

- ~1500 Cores providing ~17K HS06 ●
- 980TB Storage with another 200TB being prepared ●

Cluster management is done with a global (T2+T3) Puppet/Foreman instance.
Current running services include:



- Torque/CREAM batch system ●
- ALICE Storage on XRootD install ●
- All other storage on DPM ●
- Squid, BDII, APEL, ARGUS, VO Box ●
- 10 Gb/s link (soon to be 20 Gb/s) ●

The divide between experiments is ALICE 60%, ATLAS 30%, LHCb 5%, Other 5%

The current state of our local (Tier 3) systems is shown below:

- Batch Cluster Farm, 8 nodes, 32 (logical) cores, 48GB per node •
- 180TB 'New' Storage + 160TB 'Older' Storage + 160TB 'Oldest' Storage •
- ~80 (mostly) Fedora 24 Desktops •
- SL6 image (access through chroot) •
- Two F24 login nodes •
- Two Web servers •
- DHCP, Mail and LDAP servers •
- Share 1 10Gb/s link with the university •

The biggest change in recent months has been the loss of Matt Williams

There were a few teething problems during autumn due to this, mostly because of teaching commitments

However, thanks to improvements in monitoring and changes to some services, things are on a more even keel now!

With this in mind, I have been moving forward with plans for the future taking into account the reduced manpower:

- Making progress on server room rearrangement
- Integrating all monitoring into Grafana
- Switching to using VAC on Grid
- Switching to using ZFS for the storage
- Shifting all Tier3 storage to MooseFS

In recent weeks I have even managed to do some Ganga work!

After several discussions with people, I have decided to (gradually!) move all our Tier 2 storage from hardware RAID 6 to ZFS and not buy RAID6 cards for new storage

From my point of this has several benefits:

- Easy to monitor disk health across all systems
- Can use ~any disks in the RAID
- Cheaper to buy new hardware

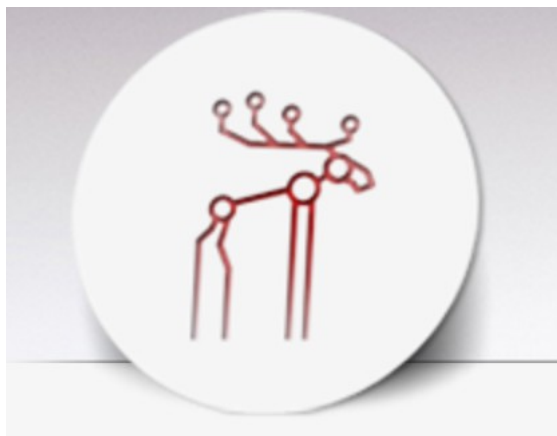
I have currently moved 100TB (prev. 40TB) of storage over to ZFS on Tier2 and have had no issues at present...



Since I took over 3 (4?) years ago, I have been attempting to find a good distributed solution for our Tier 3 storage to move away from basic NFS mounted RAIDs

I originally tried lustre but had no end of trouble trying to keep it working with the modern kernels that Fedora ships with

I eventually found 'MooseFS' which offered everything I needed:



- Very easy to setup and administer
- Very configurable
- Redundancy built in
- All done through fuse
- Can keep using disks until they die

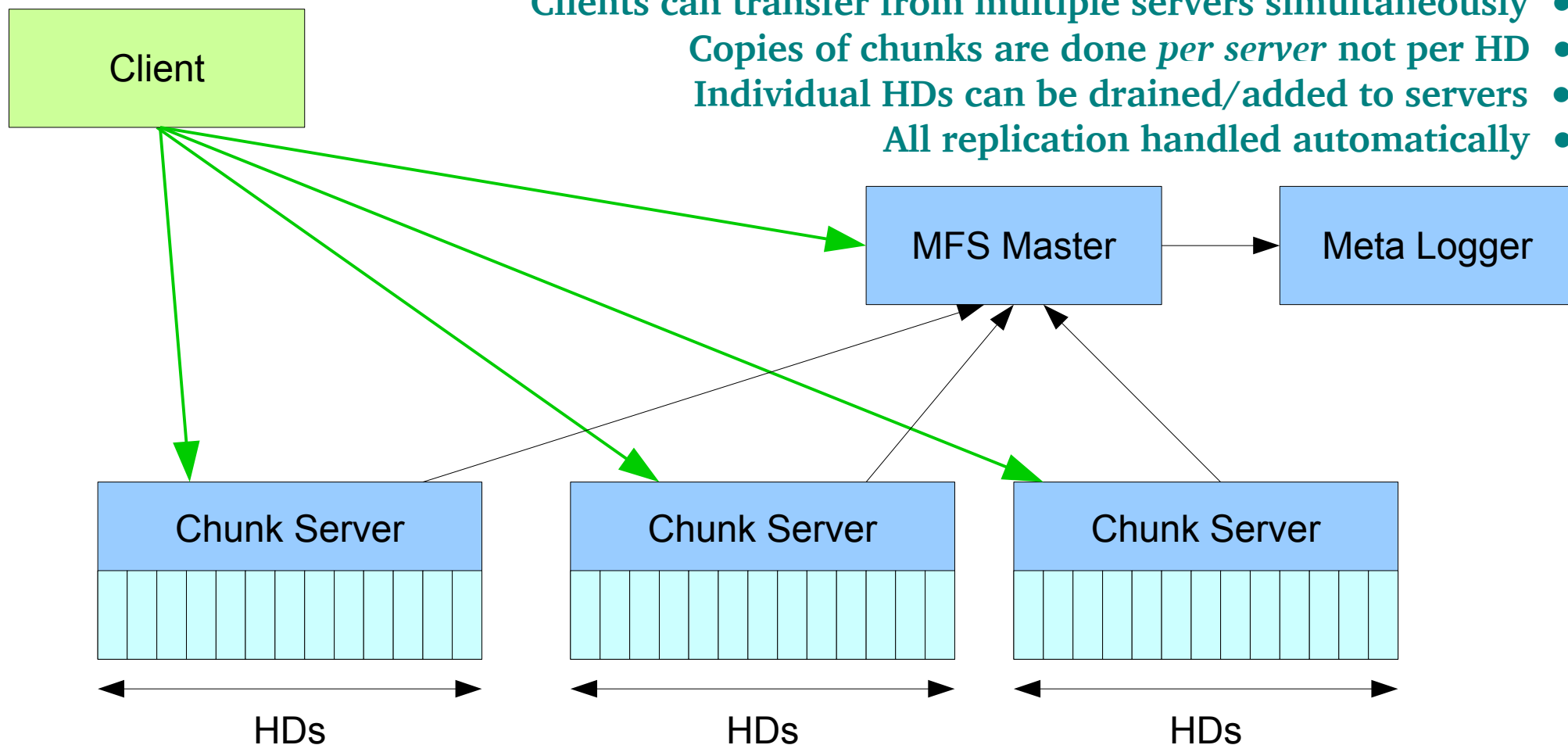
We currently have ~470TB giving (with my setup) 235TB usable space



MooseFS on Tier3 Storage

Moose works by dividing every file in to 'chunks' and then copying these chunks to the 'chunk servers' given the appropriate policy:

- Can have 1+ copies of each chunk
- Clients can transfer from multiple servers simultaneously
- Copies of chunks are done *per server* not per HD
- Individual HDs can be drained/added to servers
- All replication handled automatically

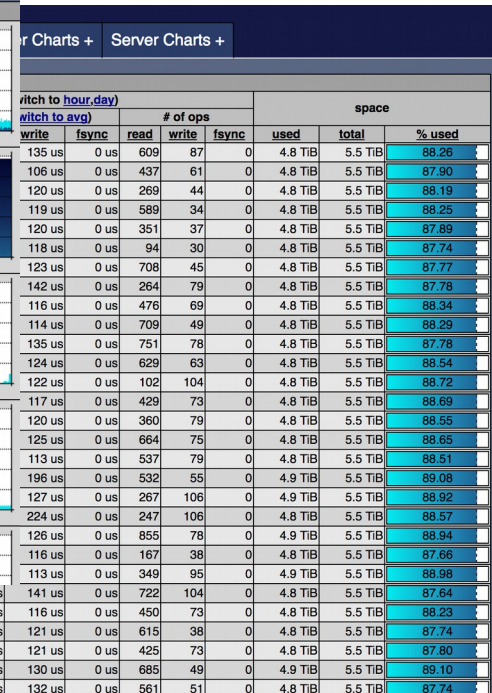
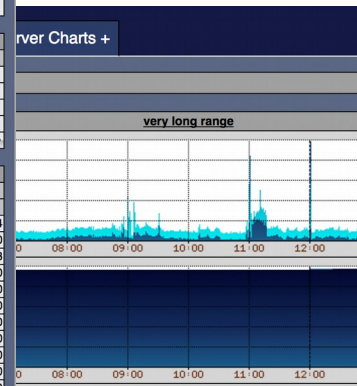
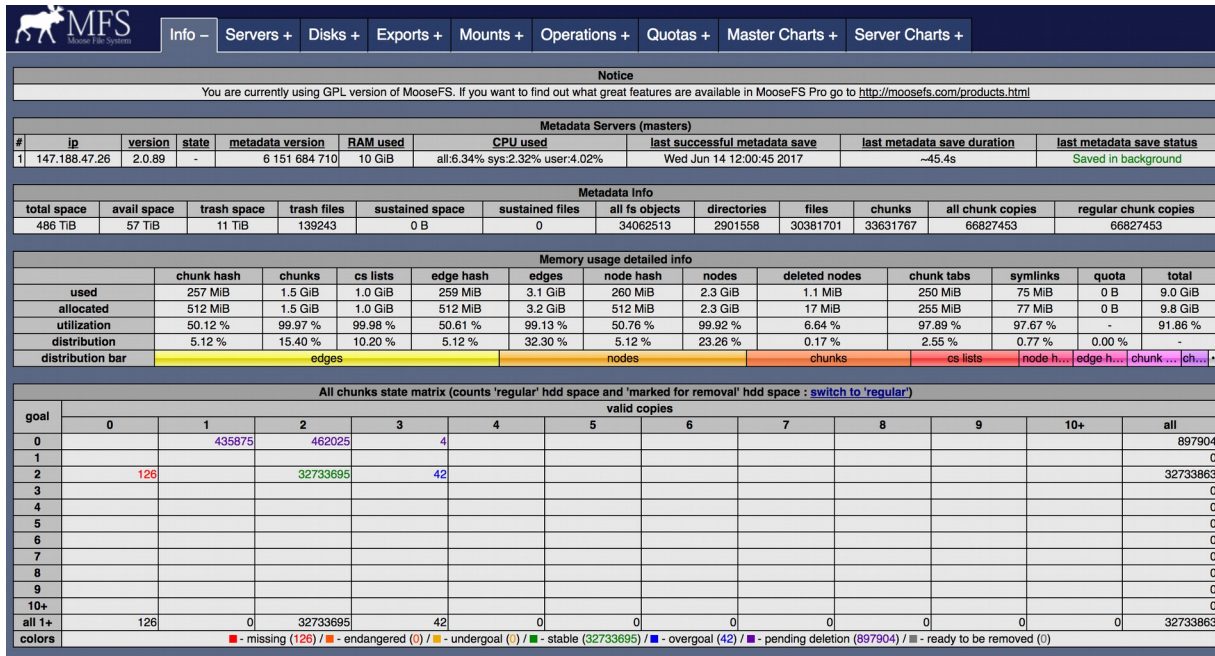




MooseFS on Tier3 Storage

Monitoring is handled through a very good web interface that gives:

- Overview of chunk replication progress (e.g. under/over goal)
- Individual disk space across all servers
- Load across all servers
- Client/Mount point information
- Number and type of commands issued by clients



Making Monitoring Easier

Before Matt left, he installed Grafana and started setting up monitoring pages. I've been continuing this work to cover both T2 and T3 machines:

- Graphite/Carbon system is incredibly easy to setup •
- Can monitor everything I want and easily add more •
- Grafana makes setting up dashboards trivial •



Moving Workers to VAC

The biggest ongoing change is that I'm switching all the workers from Torque/CREAM to VAC. Again, this has many benefits for us:



- Very easy to setup (after initial teething problems)
- Don't have to worry as much about OS updates, etc.
- Minimal ongoing administration required
- Don't have to run CREAM, Torque, APEL
- Reduces complexity of other services (Squid, BDII, Argus)
- Overall a significant reduction in manpower required

Drawbacks I've currently encountered:

- Initial setup did have problems (mostly because of me!)
- Much harder to overprovision due to HD and memory reqs being 'enforced'
- I found I needed a Squid per VM factory/Worker

Current status of VAC at Bham is that I have shifted ~50% of the site over

Many Thanks to Andrew McNab for helping me through the setup!

Generally, the install and setup of VAC was very easy. I just followed the instructions on the web page:

<https://www.gridpp.ac.uk/vac/admin-guide.html>

Fundamentally though, after installing appropriate libvirt tools, it's just a case of installing a single RPM

The configuration is managed through a small handful of easy-to-understand config files.

The only issues/gotchas I encountered were:

Firewall:

As I use puppet to manage iptables, it took a few tries to get every rule put in correctly

HS06, GOCDB entry:

VAC is able to send accounting records directly, however you must remember to add an appropriate GOCDB entry and the HS06 values for each worker node

Squid:

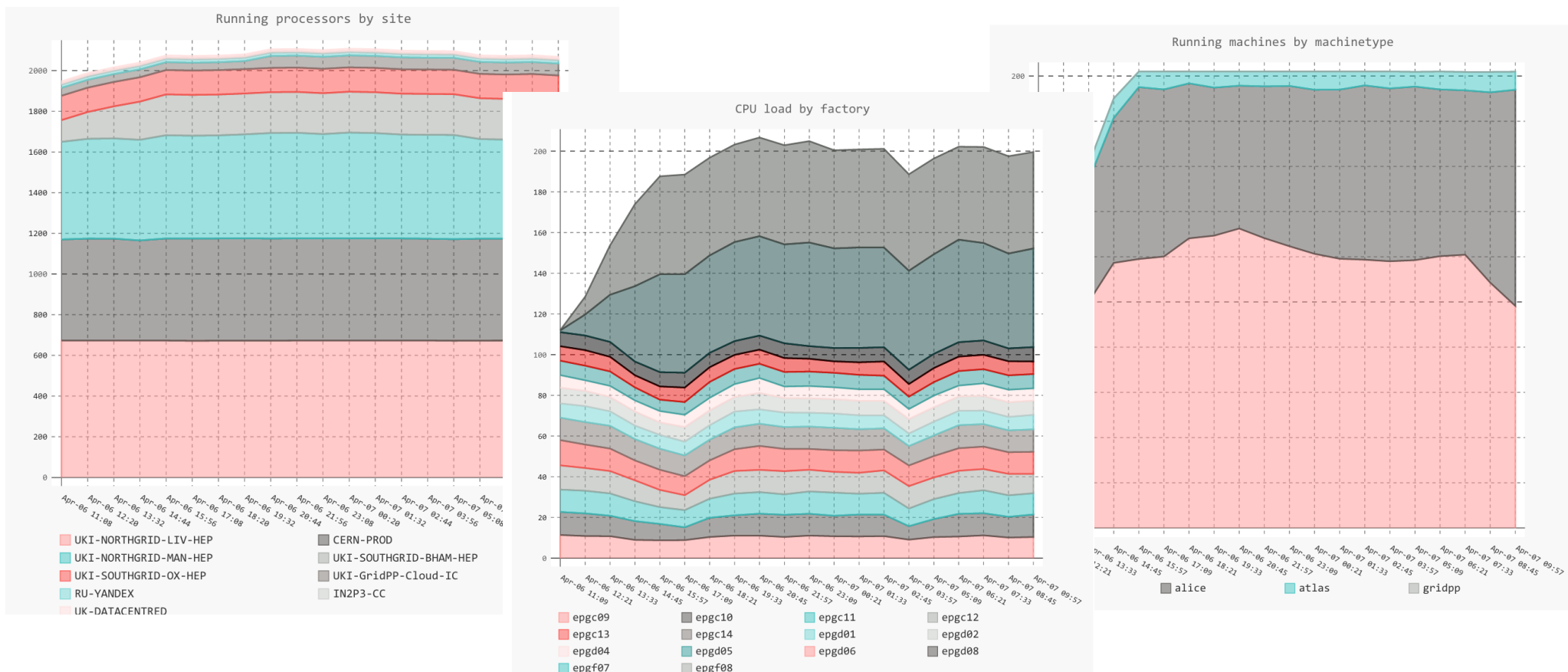
You will probably need multiple squids to cover the additional load because, as far as the squid is concerned, you will have a worker node per core. Our squid couldn't handle this and so I went to a squid per factory. Hopefully I can reduce this in the future.



There is very good overall VAC monitoring available here:

<http://vacmon.gridpp.ac.uk/1f4:15180::/>

From this you can drill down to your site and individual workers



In recent months I've been concentrating on putting things in place to make sys-admin tasks as easy as possible:

- Reorganised server room allowing for ease of installation and expansion •
- Easy monitoring of all aspects of the site via Grafana •
- Switching to ZFS over HW RAID 6 •
- Moving all workers to VAC •
- Switching to MooseFS from Lustre •

I hope to have completed all these tasks by the end of the year and will be in a much better position to keep on top of both Tier3 and Tier2 machines with the reduced manpower