

FLTECH Cluster Documentation

Patrick Ford
August 2010

Recent Work	2
User Account List	4
Power System	6
Viz Node	6
New NAS Installation	7
Spares Kits	7
Regular software updates and maintenance	7
Important Contacts and Links	9

Recent Work

Kernel Patch

On July 28th 2010 I patched the kernels on the compute nodes. For stability reasons, the kernel stays at version 2.6.18 for our CentOS 5 version, but patches append a number to the end (in our case 2.6.18-194) and fix important stability and security issues. Generally, OSG will let the site security officer know if a patch is essential to grid security, otherwise it is not a high priority update. To update the kernels on all machines (note that cluster-fork is run on the CE):

CE/SE/NAS: # **yum update kernel kernel-headers kernel-devel**

Compute nodes: # **cluster-fork yum update kernel kernel-headers kernel-devel -y**

NOTE: I had to add some postrouting and packet forwarding stuff to iptables in order for the compute nodes and NAS to access the internet through the CE (as they don't have an external network interface themselves). See the following lines from /etc/sysconfig/iptables on the CE.

***nat**

**-A POSTROUTING -o eth1 -j MASQUERADE
COMMIT**

**-A FORWARD -i eth1 -o eth0 -m state --state NEW,RELATED,ESTABLISHED -j
ACCEPT**

-A FORWARD -i eth0 -j ACCEPT

Condor

On July 29th 2010 I upgraded condor on the CE using the following guide.

<https://twiki.grid.iu.edu/bin/view/Tier3/CondorSharedInstall>

I made this decision due to the fact that condor RPMs are now not relocatable, and hence I couldn't upgrade the 7.2.0 version. I installed 7.4.2 from the RHEL 5.0 source and followed the instructions exactly, except the condor installation is in /mnt/nas0/condor. This configuration installs condor on a shared directory on the NAS that is accessible by the compute nodes, so they run the condor binaries (and daemons) from the same source as the CE. The configuration files for the CE (central manager) and the nodes are also stored on the shared directory (/mnt/nas0/condor/condor-etc), making configuration much easier. We no longer have to install and configure condor using the Rocks kickstart distribution (i.e. extend-compute.xml) every time we make a change to the config, instead only a few commands are required to set up the nodes to use the shared condor version once after they kickstart-install. The current extend-compute.xml configuration is located in:

/home/install/site-profiles/5.0/nodes/extend-compute.xml

The file: **/mnt/nas0/condor/condor-etc/condor_config.cluster**

contains all the important central configuration options for the whole cluster, the condor_config.CE and condor_config.compute can specify options unique to the different systems.

NOTE: I believe that it was necessary to add the following lines to IPtables on the CE and all compute nodes.

```
-A INPUT -m state --state ESTABLISHED,NEW -p tcp -m tcp --dport 9000:10000 -j ACCEPT
```

```
-A INPUT -m state --state ESTABLISHED,NEW -p udp -m udp --dport 9000:10000 -j ACCEPT
```

Also, I tried to set up the password security system for the condor pool, but it was throwing a lot of authentication errors that I could not debug, so I disabled it. This is not essential, since condor only recognizes machines that have the same domain config and are part of the Rocks cluster database. It only provides an extra layer of security. Condor appears to be working very well now under the new version (7.4.2) and configuration, any issues we were having with held grid jobs seems to have disappeared.

Before performing the installation, I dumped the result of `condor_userprio -all -allusers` to a file in the root home directory called `Condor-7.2-usertime` so that we know what the total usage was before the new condor rollover.

The root home directory also contains a script called `NewCondorSetupScript` that lists the commands that need to be executed on the compute nodes so that they'll interface with the shared condor.

CRAB

On July 8th 2010 I upgraded our version of CRAB to 2.7.3 at the request of samir. This was a simple matter of getting the tarball from <https://twiki.cern.ch/twiki/bin/viewauth/CMS/CrabClientRelNotes273> and extracting it under `/mnt/nas0/OSG/APP/crab`

Note that I installed it under the `uscms01` user (`su to root then su uscms01`). After doing this, I ran `./configure` and everything worked. I attempted to update `gliteUI` because I thought that it was required for CRAB, but it turns out that this is not necessary unless a user has to run jobs on an european grid site. We won't get enough benefit from this yet to justify setting this up.

Denyhosts

On June 21st 2010 I set up a program to prevent ssh brute force attacks called `Denyhosts`. This program parses the `/var/log/authpriv` file to monitor login attempts, and blocks an ip address that makes 5 incorrect attempts on the root password, or 25 incorrect attempts on a valid user account, or 5 incorrect attempts on an invalid user account (i.e. one that doesn't exist). This was set up on both the CE and SE and appears to have stopped all brute-force hack attempts on the cluster by adding the blacklisted ip addresses to `/etc/hosts.deny`. The number of failed attempts is reset periodically to prevent the banning of valid users. The guide that I used to set this program up is located here:

http://www.howtoforge.com/preventing_ssh_dictionary_attacks_with_denyhosts

OpenGL

Will Bittner, who is creating our GUI for viewing muon tomography data, needs several up to date OpenGL libraries for graphics processing. I attempted to update these versions on the cluster only to find that the earlier versions are pre-requisites for over 300 critical packages on the CE CentOS install. As such, CentOS does not allow updating past a certain version for compatibility reasons. The ubuntu machine in the lab can be used for OpenGL processing (named vx4d by Will), and the plan was to set it up to tunnel the x-server (linux graphics engine) over ssh as Ubuntu has all the required libraries for his code. Still, this issue has forced us to investigate the addition of a Viz server to the cluster running a different OS than Rocks.

Crontab Fix and Phedex proxy update

The user crontab on the SE was fixed by editing /etc/exports to include no_root_squash option for the CE /export mount. Now the script **/sandbox/phedex/gridcert/proxyrenew.sh.forpatrick** can be run by the phedex user to renew the phedex proxy whenever necessary (via crontab).

Note that my membership with CMS has expired, so xenia needs to replace my certificate on the phedex user with hers and edit the proxyrenew script with her password.

User Account List

Condor wall hours since 1/08/09 until condor upgrade to 7.4.2. All other accounts are system or grid accounts.

Account Name	Full Name	Main Group(s)	Condor wall hours
pford	Patrick Ford	users, wheel	0.1
g4hep	Geant4 HEP group	mtusers	58346.46
idiaz	Ismael Diaz	users	707949.72
rpena	Rafael Pena	users	0
kgnanvo	Kondo Gnanvo	mtusers	0
rhoch	Richie Hoch	mtusers	91.17
sguragai	Samir Guragain	osgusers	46894.50
mzhang	Ming Zhang	users	7716.10
uscms01	CMS Grid User	osgusers	390514.56
glow	GLOW Grid User	osgusers	60197.20

Account Name	Full Name	Main Group(s)	Condor wall hours
osg	OSG Grid User	osgusers	1292.74
geant4	Geant4 Grid User	osgusers	3471.80
mwood	Matt Wood	users	46982.46
holuseyi	Hakeem Oluseyi	users	0
xfave	Xenia Fave	users, wheel	0
dthomas	David Thomas	users	1195.11
xluo	Xi Luo	users	70673.60
lgrasso	Lenny Grasso	users	0
fit	Test Account	fit	0
zsaleh	Ziad Saleh	users	0
ecramer	Eric Cramer	users	0
rsnihur	Rob Snihur	users	0
ywu	Yujun Wu	users	0
kgamayunov	Konstantin Gamayunov	users	0
bkosar	Burcu Kosar	users	8950.68
rromain	Randy St-Romain	users	0
jfischer	Johanna-Laina Fischer	users	0
djohnson	Doug Johnson	users	0
hohlmann	Marcus Hohlmann	users	0
aquintero	Amilkar Quintero	users	0
hkalakhety	Himali Kalakhety	users	0
gemhep	GEM HEP Group	mtusers	0
g4cs	Geant4 CS Group	mtusers	0

Account Name	Full Name	Main Group(s)	Condor wall hours
bdorney	Brian Dorney	users	0
crobinson	Curtis Robinson	users	0
jstevens	John Stevens	users	0
wbittner	Will Bittner	mtusers, svnuser	0
pzuo	Pingbing Zuo	users	21365.38
		TOTAL	1425641.58

Power System

Currently we have two 6-outlet 15A 1U rackmount power strips per 10 compute nodes (4 power strips total). This was necessary due to the 120V outlet configuration on the Tripp-lite UPS. This setup is not very space or energy efficient. Server power supplies are capable of auto-switching between 120V and 240V line voltages, including 208V in the US. 208V only has to supply approximately half the current to a load than 120V, leading to lower temperatures and higher efficiency. My plan was to buy 2 power strips (either 1U or Vertical [zero-U]), one each for 10 compute nodes, they would plug directly into the L6-30 receptacles on the Tripp-lite UPS's, which will supply 30A at 208V, far more than enough to run 10 compute nodes.

This is not an essential upgrade, as the current power configuration works, it's just not ideal. The current cluster rack layout and power configuration is attached.

The quote for the new 36-drive NAS has an estimated power consumption of 858W/880VA. This will be powered from the 2100W APC UPS, the new NAS should reach, but not exceed, its rated capacity. Measurements of the current draw of the CE, NAS0, SE, and Cisco switch at full load confirm that there will be just enough power available for the new NAS. The NAS houses two redundant power supplies, so a splitter may be required to interface with the UPS. If the APC UPS is overloaded, the new NAS can be moved to one of the Tripp-lite UPS.

Viz Node

Preliminary tests have shown that the school (and even the cluster) network is a massive bottleneck for visualization applications, therefore customizing a machine to attach directly to the cluster will be a waste of money - as any performance gain from high performance hardware will be mitigated by the network latency. A workstation in the HEP lab that is accessed locally will be necessary to get a decent framerate on

OpenGL applications. It should be possible to export the NAS partition from the CE to the workstation.

A machine can be custom built, or pre-built, but custom will have a much lower price. The nVidia GTX470 (Fermi) is likely the best value for us, combined with a Core i7 quad core processor and at least 6GB RAM. This could possibly be built for under \$1000, and dual-monitors can be used.

New NAS Installation

The new NAS is configured to have two static drives in RAID1 configuration for the Rocks OS. The 36-drive array will be an independent extended volume. This is necessary in order to perform multiple tests and rebuilds of the array to optimize the configuration. The possible RAID configurations provided by the 3ware RAID card are RAID5, RAID6, and RAID50. Given the number of drives, either RAID6 or RAID50 should be considered. My plan was to create the array using RAID50 (3 striped RAID5 arrays with 1 hot spare each, this will lose 6 drives worth of free space) with GPT partitioning using the XFS filesystem, and NFS4 for networking. I would then perform dd read and write tests both on the system itself and over the network in order to determine the real-world copy and read speeds, this test would then be repeated using RAID6 and the results compared. I have a strong suspicion that RAID50 will work the best as it will give us far better write performance, and since only Grid data will be stored on the array - data redundancy is not as important as, for example, for user home directories.

Spares Kits

The spares kits from SiMech are adequate, and are customized for each chassis that they sell. One (or even two) spare 1TB drives would be a good idea to have on standby in case of drive failures, possibly a spare RAM stick and power supply. We also need to return the faulty 750GB drive to SiMech in exchange for another backup for NAS0.

Regular software updates and maintenance

When to update Rocks, Condor, OSG, BeStMan, PhEDEx, CRAB.
Cluster shutdown/restart

The cluster is relatively maintenance-free. The RSV and SAM tests provide a snapshot of the availability of the site, and indicate when any essential services are failing for some reason. This information can be found at the following links:

<http://myosg.grid.iu.edu>

(searching for Florida Institute of Technology)

http://dashb-cms-sam.cern.ch/dashboard/request.py/latestresultssmry?siteSelect3=T3&serviceTypeSelect3=vo&sites=T3_US_FIT&services=CE&services=SRMv2&tests=1301&tests=133&tests=111&tests=6&tests=1261&tests=76&tests=64&tests=20&tests=281&tests=882&exitStatus=all

Updating the Rocks OS requires a complete reinstall of the entire system, and I wouldn't recommend this unless there are serious security or stability issues found in Rocks 5.0.

Condor was recently updated, and generally our version should not fall more than two major releases behind (e.g. our version is 7.4, we should upgrade at or after version 7.8, odd numbered releases are development releases). With the new shared condor installation, upgrading the base code should be hassle-free since the configuration files are stored elsewhere.

OSG should be updated regularly, and the instructions for updating the 1.2 release can be found at:

<https://twiki.grid.iu.edu/bin/view/ReleaseDocumentation/OSG12UpdateInstructions>

Thankfully, updating is quite simple, but sometimes they change the skeleton config.ini files so copying over the old config.ini file is sometimes not possible, however the new file can be heavily based on the old one.

BeStMan can be upgraded using the same technique as the OSG CE upgrade, and the configuration is simpler.

Updating PhEDEx involves a complete reinstallation, and copying over of the old configuration files (contained in /sandbox/phedex/SITECONF/T3_US_FIT/PhEDEx). The installation of PhEDEx can be a little tricky, so I have not updated in a while, this doesn't seem to be a problem however. The installation instructions for both BeStMan and PhEDEx can be found in the SE installation guide. Obviously, we should not fall too far behind in releases, as the CMS software team are constantly improving the security, speed, and efficiency of PhEDEx.

The latest version of CRAB has been installed (2.7.3) and updating it is very simple. One must simply download the new tarball (such as from <https://twiki.cern.ch/twiki/bin/viewauth/CMS/CrabClientRelNotes273>) and extract it in its own versioned subdirectory in /mnt/nas0/OSG/APP/crab. Then enter the directory and run ./configure. Since updating is so simple, it can be done very regularly, but generally it's not necessary unless local users require new features.

Lastly, patching the linux kernel generally is only necessary if there is a serious security issue, and OSG or the T3 folks will be very vocal about it if this happens.

System restart

Shutting down the system is only really necessary for major system updates (such as kernel patches) or when there is an extended power outage (over 15 minutes). The

system may be taken down in any order as long as the CE is last (since it is the login machine). Run on the CE:

```
# cluster-fork /sbin/init 0
```

```
# ssh dev-0-0; /sbin/init 0
```

```
# ssh nas-0-0; /sbin/init 0
```

Confirm that all the servers are down either by checking ganglia or attempting to ssh.

```
# /sbin/init 0
```

The cluster can then only be restarted manually, following the given order:

1. Turn on the NAS and ensure it starts up, this is important to make sure that the other servers will mount the exported directories.
2. Turn on the CE and SE together, ensure they start up correctly and the /mnt/nas0 directories are mounted.
3. Turn on the compute nodes (all at once is fine).
4. Check that condor is running by running **\$ condor_q** or **\$ ps -ef | grep condor**
5. ssh into the SE and run (as user phedex):
\$ cd /sandbox/phedex; ./PHEDEX/Utilities/Master -config SITECONF/T3_US_FIT/PhEDEx/Config.Prod start
6. Everything else such as the OSG CE and SE services should have come up automatically. You can check by running and looking for globus processes:
\$ ps aux | grep globus
7. The system should now be fully operational, check the RSV and SAM tests in about an hour to make sure essential services are working.

Important Contacts and Links

Bockjoo Kim - Tier 2 center at UF, installs CMSSW releases on T3 sites.

Yujun Wu - System Admin at the Tier 2 center at UF, useful for system and network level info.

T3 hypernews - Requires a CMS account to access, all Tier 3 admins use this to support and help each other.

Rocks list - For any issues with Rocks, <https://lists.sdsc.edu/mailman/listinfo/npaci-rocks-discussion>

condor list - For any issues with Condor, <https://lists.cs.wisc.edu/mailman/listinfo/condor-users>

Doug Johnson - Part time Tier 3 support person, very knowledgeable about all the Grid middleware from OSG to CMS.

Malina Kirn - Admin at UMD Tier 3 site, graduate student, very helpful.

Rob Snihur - Tier 3 support person at Fermilab. Organizes the biweekly Tier 3 meetings and helps support Tier 3 sites.