

## **Procedures for restarting the MC farms.**

**Tomasz Wlodek**

*(UTA-HEP Computing #1)*

The following procedure applies both to rebooting the farm and to restarting it after power glitch.

### **Introduction**

There are three types of nodes in the farm: the server, file servers and execution nodes. The servers are cse000 and hepfm000. The file servers are hepfm001-3 and cse001. The remaining nodes are execution nodes.

All nodes mount their home directories on server. All execution nodes mount cache disks on file servers. Main server mounts all cache disks on all nodes.

When booting the farm machines will try to nfs mount external disks. When the nodes on which those disks sit are down, nfs will time out after few retries. However if two machines, which mount each other disks, boot in the same time, they can end up in an undefined state where they can neither mount their disks, nor time out. For that reason it is necessary to boot the nodes in a strict order.

### **Order of starting nodes**

1. When all nodes are powered off, connect a free terminal to one of file servers. Then power the file servers on, the one which has terminal connected as the last one. You will be able to follow his booting stage at the screen. Since you have powered it on as the last one, if this one is ready, the other one should be ready as well.

The fileservers will try to boot and nfs mount home directory. Since main server is down they will time out after few tries.

When they are up you can logon to fileservers as root. Logon and ping the other file servers to check whether they are alive. You cannot logon as mcfarm (or any other user) yet, since the yp server is not ready yet. (Because the main server, which hosts the yp server, is still down.) But you can logon to the node as root.

2. Connect the terminal to one of the production nodes. Power on the production nodes, the one with terminal attached as the last one, so that you know when it finishes other should be ready as well.

Each production node will try to mount fileservers - but since fileservers are ready the mount will be successful. They will try to mount /home on server, and will time out while doing this.

3. When all production nodes are up and running boot main server. It will mount everybody's cache disks.

### **After powering on the nodes**

Now all nodes are up, but file servers and production nodes did not mount /home directories yet. (Because when they were powered on /home was not yet available.) You can do now two things:

a) Logon to every node, become a root and issue *mount -a* command to mount /home from main server

Or

b) Logon to every node starting from 001 down to the last one, become root and do */sbin/shutdown -r 0 &* to reboot it again. This will accomplish the same thing: mount all remaining disks. My preference is for the last method, since it makes sure that all nodes are running.

### **Starting farm software**

When everything is done you should:

1. Logon to mcfarm on server and start farm software

*start\_farm*

2. Logon to root, go to /home/products/condor/sbin/ and do *./condor\_master* to start condor.

This should be done on every production node on the farm. (Not on file servers! – we do not want condor to run on file servers!) Logon to all production nodes (hepfm004-24 and cse002-11), become root and start condor on all of them.

3. Logon to sam account on server and do

```
source ~/products/etc/setup.sh
export $SAM_STATION=uta-hep
ups start sam_bootstrap
```

(This applies to hepfm farm only, as cse does not run sam)

4. Start mcfarm bookkeeping software: logon as mcfarm, goto /home/mcfarm/production\_logs/bookkeeper

and do

python start\_dolog

you will be prompted for mcfarm account password. After you give it program will go to the background.

### **Checking the file structure for disk errors.**

Occasionally Linux files can get corrupted. (Rarely, but we have seen this already.) As a precaution a file test should be done from time to time, when farm is not running

1. Logoff all users.
2. Logon as root from the terminal in swift center. (Remote logons as root are not allowed)
3. Unmount the disk you want to test from all machines which are using it.  
(*umount <disk name>*)
4. Stop the root daemons.
5. Check the disks using */sbin/fsck* or */sbin/e2fsck* commands. Those commands are described in man pages, read it before operating. On the server */sda1* is system disk / and */sda7* is the /home disk.
6. Reboot the server.

### **Warning about rebooting hepfm000**

Once upon a time a faulty piece of shareware software has damaged boot sector on hepfm000 disk. Ever since hepfm000 will not boot from its hard disk. For that reason we boot it using a floppy disk, inserted permanently into its floppy drive. Do not take this floppy out. The disk problem will be fixed in some distant future, when we change linux version to Linux 7.\*. For the time being booting hepfm000 must be done from floppy.