# Building a Linux Cluster

## CUG Conference
### May 21–25, 2001

**by**
**Cary Whitney**
**Clwhitney@lbl.gov**

# Outline

- **What is PDSF and a little about its history.**

- **Growth problems and solutions.**

  - **Storage**

  - **Network**

  - **Hardware**

  - **Administration**

  - **Software**
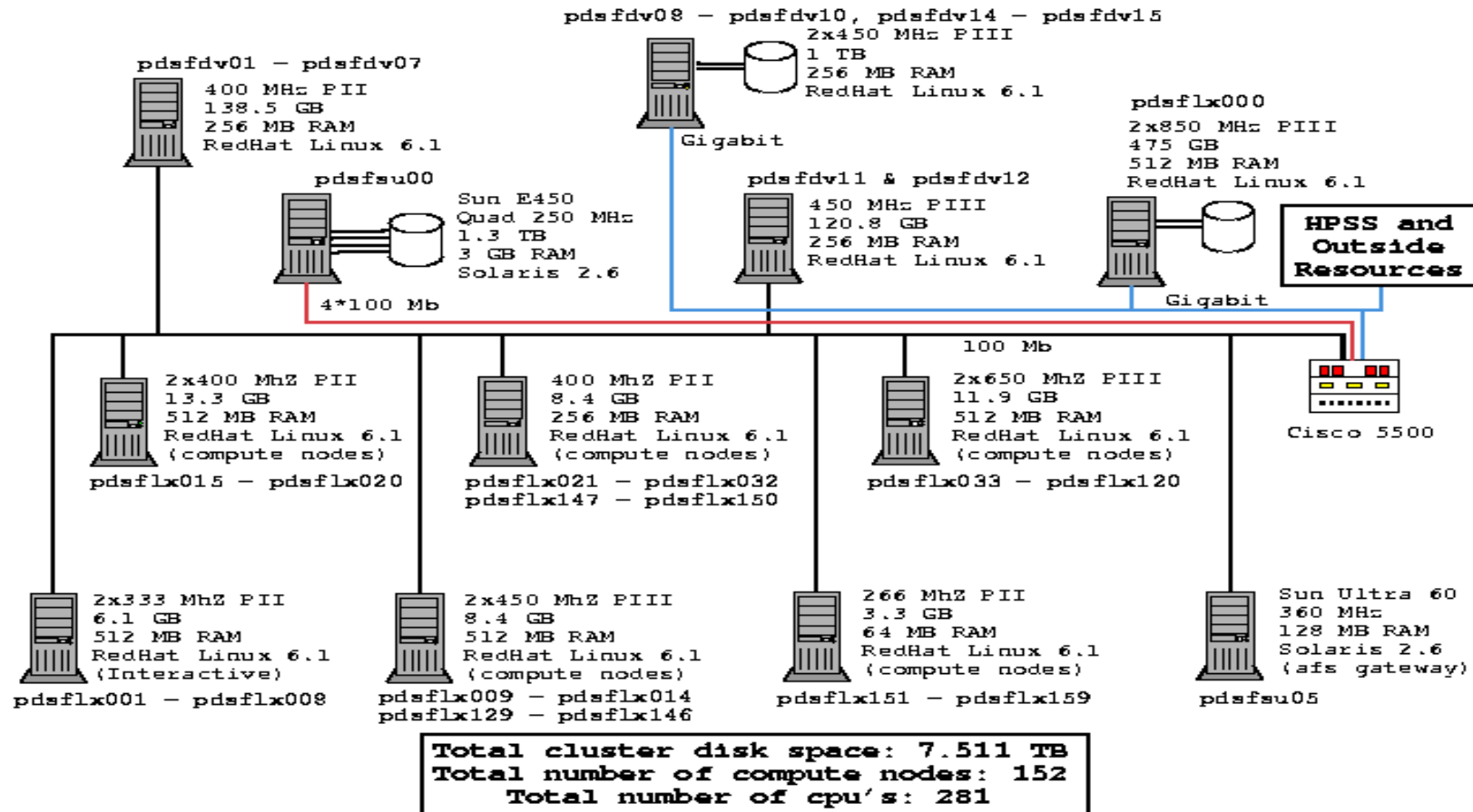
- **In the future**

- **Conclusion**

# PDSF

- **Linux cluster who's primary service is to the High Energy Physics (HEP) community**

  - **Parallel Distributed System Facility (PDSF)**

  - **Runs on commodity hardware**

  - **Takes advantage of open source software**

  - **Multiple user communities running on the same cluster**

  - **Application are Embarrassing Parallel (Seti@Home)**

  - **Fast ethernet and some gigabit interconnects**

# HEP

- HEP community
  - **Large datasets (up to 1 PB in size)**
  - **Embarrassing Parallel problem (no need for fast interconnects between machines or multiple processors per machine)**
  - **Can exploit commodity hardware market**
    - Dual processor instead of quad or larger
    - Not effected by limit of 32 bit architecture
    - Can run in under 4 GB of memory
    - Does not need checkpoint/restart capabilities
  - **Experiments span multiple labs and countries (100's to 1000's researchers)**

# PDSF Layout

# History

- **Started at Superconducting Super–Collider**

  - **1991 1000 MIPS and 40 GB disk (HP and Sun)**

  - **1992 2000 MIPS and 80 GB disk**

  - **1993 8000 MIPS and 240 GB disk**

  - **1994 128 processors/12000 MIPS and 160 GB disk**

- **Moved to LBNL**

  - **1997 128 processors and 160 GB disk**

  - **1998 142 processors and 282 GB disk (Linux for disk vaults)**

  - **1999 66 processors and 658 GB disk (Move total to Intel)**

  - **2000 281 processors and 7.5 TB disk**

  - **2001 ~431 processors and 22.5 TB disk**

# Problems and Solutions

- **Networking**

- **Storage (Disk Space)**

- **Cluster Filesystem**

- **Administration**

  - **Configuration Management**

  - **Monitoring**

- **Hardware Density**

- **Users**

# Networking

- **Everything is based on fast ethernet**

- **Network bottlenecked for the NFS servers on the fast ethernet since a user can have up to 280 jobs running at a time**

  - **Tom Davis wrote the kernel bonding driver to bond two ethernet ports together when the max jobs was 120**

  - **Now we are running copper gigabit for added throughput**

# Storage

- Datasets outgrew 40 GB disk vaults
  - Used ide drives and linux to create 64GB disk vaults
  - Then upgraded to Raidzone's 15 drives plus 75 GB drives to create a 1 TB RAID 5 filesystem
- 1 TB + Linux NFS seems to be unstable under our work conditions
  - Limit system filesystem to under 1 TB currently 600 GB with 3ware
  - Looking into a SANS solution for greater needs

# Network Filesystem

- **Userland NFS proved too slow**
  - **Knfsd was introduced to help performance**
- **NFSv2**
  - **NFSv3 patches added but still running in v2 mode. Performance increased over standard kernel v2.**
  - **Checking into NFSv3**
- **But even NFSv3 has problems scaling**
  - **Looking into GFS or other network filesystem**

# Administration

- **Standalone configuration on each node had a problem with staying in sync**
  - **Installed with an NFS mounted /usr**
- **NFS mounted /usr has problems with RPM installs and local configuration files.**
  - **Planning on moving back to RPM since autorpm works better.**
  - **Cfengine to help maintain configuration files**

# Monitoring

- **SNMP polling was timing out because the number of nodes was increasing.**
  - **Implemented MON**
  - **MRTG was added to monitor the network**
- **MON worked but its interface was not friendly for operations staff and users**
  - **Checked out Big Brother/Sister**
  - **Now using Netsaint**

# NIS

- **Our central NIS server could not handle the load**
  - **Moving to our Sun box did not help**
  - **Setting up NIS broadcast between multiple servers only loaded the fastest responding server**
  - **Grouped several nodes to point to one server but this still has problems when a server goes down**
  - **Possible move to static files on compute nodes and NIS only for interactive nodes**

# Hardware Density

- Desktop mini-tower cases take up too much room in standard racks

    - Moved to 2U rack mount machines with dual cpus

    - Intel flip chips now allows 1U dual processor nodes

    - Care must be taken with cooling with high density

- Disk vaults where mini-towers with 4 ide drives

    - Moved to Raidzone hardware with 15 drives in 8U

    - New Raidzone are 15 drives 4U but restricted software at $22k

    - 3ware can provide 16 drives 6U at $15k

# Networking Hardware

- **Our Cisco Cat 5513 is full. 2 gigabit blades and 10 24 ports of 10/100 blades**
  - **Moving to a distributed network of small switches at the top of each rack with links back to a main switch**
  - **This creates less spaghetti wiring**
- **8 ports fiber gigabit blade density is not enough**
  - **Extreme Summit 7i is 28 copper + 4 fiber gigabit ports which can auto sense between 100/1000**

# Console Cabling

- **KVM switches are expensive and do not offer remote administration**
  - **Moved to Rocketport cards with linux console software**
  - **Enabled linux serial console. This allows access from the lilo prompt but not BIOS level stuff.**
  - **Serial consoles are rack based and not centralized**

# Power Cabling

- **Old machine room layout was 2 20 Amp circuits per rack**

    - **Since our density was increasing, when our new facility was being built we planned for 4 20 Amp circuits per rack. 8 nodes per circuit.**

    - **High bandwidth nodes are 5 per circuit**

    - **Power on flexible conduit under the floor so reconfiguration of power is easier**

    - **Key servers on UPS**

# User base

- **Most clusters have a single application/user base to deal with.**

- **PDSF is a general cluster supporting several groups.**

  - **Configuration is kept simple and not driven by any group**

  - **Local setups provide users with same look and feel of other clusters they may work on**

  - **Some customisations can be done only if it does not impact the general cluster (Adding light MPI work for one group's test run.)**

# Software

- LSF
  - **Main use is to provide fair sharing of compute nodes**
    - Groups buy hardware and thus shares in the cluster
    - We provide support and maintenance of the cluster
  - **Can provide some resource management (NFS servers)**
- **Group shared programs/code is placed into a 'group common' directory tree**
- **Extra system applications or applications shared by all group go into /usr/local**

# Future Projects

- Autorpm for installation
- Myrinet or other solution for increased MPI support
- SAN solutions involving fiber channel and/or gigabit (iSCSI)
- Journaling filesystem (xfs, ext3, reiser, jfs)
- Network filesystems (gfs, pvfs, afs, gpfs)
- LDAP for a replacement for NIS
- Remote power management
- Possible scheduling software replacement (GRD)

# Conclusion

- **Linux works**

  - **There are still problems but solutions are being worked upon**

  - **Large community of users and developers (Open source works)**

  - **Creating multiple clusters for different locations is cheaper for the our user base**

  - **Allows for better scaling of hardware since experiments will get bigger**