



Configuring an IBM Cluster 1350 with xCAT Red Hat 4 x86_64

March 2006
Mark Weber

Contents

1. [Introduction](#)
2. [Tested configurations](#)
3. [Understanding xCAT's functions and features](#)
4. [Getting the xCAT software distribution](#)
5. [Understanding cluster components, connections and architecture](#)
6. [Configuring the Ethernet switch and VLANs](#)
7. [Installing the operating system on the management node](#)
8. [Configuring networking on the management node](#)
9. [Installing xCAT](#)
10. [Setup xCAT](#)
11. [Configuring the terminal servers](#)
12. [Configure xCAT](#)
13. [DHCP setup and configuration](#)
14. [Installing compute nodes](#)
15. [Serial Over LAN \(SOL\) setup](#)
16. [Clean up](#)
17. [Contributing to xCAT](#)
18. [Credits](#)
19. [Supporting documentation](#)

1. Introduction

The Extreme Cluster Administration Toolkit (xCAT) is a collection of mostly script based tools to build, configure, administer, and maintain Linux clusters. This document describes how to implement a Linux cluster on an IBM Cluster 1350 using xCAT v1.2.0-RC2 and other third party software. Software versions referenced in this document may be out of date, check for the latest versions before continuing.

xCAT is for use by IBM and IBM Linux cluster customers. xCAT is copyright © 2000, 2001, 2002 IBM corporation. All rights reserved. Use and modify all you like, but do not redistribute. No warranty is expressed or implied. IBM assumes no liability or responsibility.

2. Tested configurations

The following hardware components have been tested with xCAT v1.2.0 RC2:

x86 (i386, i486, i586, i686) supported distributions for IBM Cluster 1350

- Red Hat Enterprise Linux ES 3 U3*
- Red Hat Enterprise Linux WS 3 U3*

x86_64 (Opteron and EMT64) supported distributions for IBM Cluster 1350

- Red Hat Enterprise Linux ES 3 U3*
- Red Hat Enterprise Linux WS 3 U3*

x86 (i386, i486, i586, i686) distributions supported by xCAT:

- Red Hat 7.2
- Red Hat 7.3
- Red Hat 8.0
- Red Hat 9
- Red Hat Enterprise Linux AS 2.1
- Red Hat Enterprise Linux AS 2.1 U2
- Red Hat Enterprise Linux AS 2.1 U3
- Red Hat Enterprise Linux ES 2.1*
- Red Hat Enterprise Linux WS 2.1*
- Red Hat Enterprise Linux AS 3
- Red Hat Enterprise Linux ES 3*
- Red Hat Enterprise Linux WS 3*
- Red Hat Enterprise Linux AS 3 U1
- Red Hat Enterprise Linux ES 3 U1*
- Red Hat Enterprise Linux WS 3 U1*
- Red Hat Enterprise Linux AS 3 U2
- Red Hat Enterprise Linux ES 3 U2*
- Red Hat Enterprise Linux WS 3 U2*
- Red Hat Enterprise Linux AS 3 U3
- Red Hat Enterprise Linux ES 3 U3*
- Red Hat Enterprise Linux WS 3 U3*
- Red Hat Enterprise Linux AS 3 U4*
- Red Hat Enterprise Linux ES 3 U4*
- Red Hat Enterprise Linux WS 3 U4*
- Red Hat Enterprise Linux AS 4*
- Red Hat Enterprise Linux ES 4*
- Red Hat Enterprise Linux WS 4*
- Red Hat Fedora Core 1*
- Red Hat Fedora Core 2*

Red Hat Fedora Core 3*
CentOS 3.3 (Treat as RHAS3U3)
CentOS 3.4 (Treat as RHAS3U4) (CD and DVD)
SUSE Linux 8.1*
SUSE Linux 8.2*
SUSE Linux 9.0*
SUSE Linux 9.1*
SUSE Linux 9.2* (DVD Version only, non DVD missing KSH, 32-bit EM64T & Opteron Tested)
SUSE Linux SLES8
SUSE Linux SLES8 SP1
SUSE Linux SLES8 SP2a
SUSE Linux SLES8 SP3
SUSE Linux SLES9
SUSE Linux SLES9 SP1
SystemImager
Partimage

x86_64 (Opteron and EMT64) distributions supported by xCAT:

Red Hat Enterprise Linux AS 3*
Red Hat Enterprise Linux ES 3*
Red Hat Enterprise Linux WS 3*
Red Hat Enterprise Linux AS 3 U1*
Red Hat Enterprise Linux WS 3 U1*
Red Hat Enterprise Linux AS 3 U2*
Red Hat Enterprise Linux ES 3 U2*
Red Hat Enterprise Linux WS 3 U2*
Red Hat Enterprise Linux AS 3 U3* (64-bit EM64T & Opteron Tested)
Red Hat Enterprise Linux ES 3 U3* (64-bit EM64T & Opteron Tested)
Red Hat Enterprise Linux WS 3 U3* (64-bit EM64T & Opteron Tested)
Red Hat Enterprise Linux AS 3 U4* (64-bit EM64T & Opteron Tested)
Red Hat Enterprise Linux ES 3 U4* (64-bit EM64T & Opteron Tested)
Red Hat Enterprise Linux WS 3 U4* (64-bit EM64T & Opteron Tested)
Red Hat Enterprise Linux AS 4*
Red Hat Enterprise Linux ES 4*
Red Hat Enterprise Linux WS 4*
Red Hat Fedora Core 1*
Red Hat Fedora Core 2*
Red Hat Fedora Core 3* (64-bit EM64T & Opteron Tested)
CentOS 3.3 (Treat as RHAS3U3) (64-bit EM64T & Opteron Tested)
CentOS 3.4 (Treat as RHAS3U4) (64-bit EM64T & Opteron Tested) (CD and DVD)
SUSE Linux 9.0*
SUSE Linux 9.1*

SUSE Linux 9.2* (DVD Version only, 64-bit EM64T & Opteron Tested)
SUSE Linux SLES8
SUSE Linux SLES8 SP2
SUSE Linux SLES8 SP3
SUSE Linux SLES9 (64-bit EM64T & Opteron Tested)
SUSE Linux SLES9 SP1 (64-bit EM64T & Opteron Tested)
SystemImager
Partimage

IA64 (Itanium 1 and 2) distributions supported by xCAT

Red Hat 7.2
Red Hat Enterprise Linux AS 2.1 U2*
Red Hat Enterprise Linux AS 3*
Red Hat Enterprise Linux ES 3*
Red Hat Enterprise Linux WS 3*
Red Hat Enterprise Linux AS 3 U1*
Red Hat Enterprise Linux WS 3 U1*
Red Hat Enterprise Linux AS 3 U2*
Red Hat Enterprise Linux ES 3 U2*
Red Hat Enterprise Linux WS 3 U2*
Red Hat Enterprise Linux AS 3 U3*
Red Hat Enterprise Linux ES 3 U3*
Red Hat Enterprise Linux WS 3 U3*
Red Hat Enterprise Linux AS 3 U4*
Red Hat Enterprise Linux ES 3 U4*
Red Hat Enterprise Linux WS 3 U4*
Red Hat Enterprise Linux AS 4*
Red Hat Enterprise Linux ES 4*
Red Hat Enterprise Linux WS 4*
SUSE Linux SLES8
SUSE Linux SLES8 SP2
SUSE Linux SLES8 SP3
SUSE Linux SLES9*
SUSE Linux SLES9 SP1*

PPC64 (IBM JS20 only) distributions supported by xCAT:

Red Hat Enterprise Linux AS 3 U2*
Red Hat Enterprise Linux AS 3 U3*
Red Hat Enterprise Linux AS 3 U4*
Red Hat Enterprise Linux AS 4*
Red Hat Enterprise Linux ES 4*
Red Hat Enterprise Linux WS 4*
SUSE Linux SLES8 SP3aa*
SUSE Linux SLES9*

SUSE Linux SLES9 SP1*

PPC64 Node install tested only, however should work as management node.

The configuration examples shown in this document may need to be altered to suit any variances in the cluster and architecture, but the examples should give a good general idea of what needs to be done. Please do not use this document verbatim as an implementation guide. This document should be used as a reference to a custom implementation. Use the man pages, source and other documentation that is available to figure out why certain design or configuration choices are made and what different choices may be made. Because IBM Cluster 1350 clusters are preconfigured from manufacturing, this document covers very little of the hardware configuration that is required to implement a cluster. Additional documentation including hardware installation and configuration is available as a RedBook at [http://publib-b.boulder.ibm.com/Redbooks.nsf/9445fa5b416f6e32852569ae006bb65f/7b1ce6b3913caf b386256bdb007595e8?OpenDocument&Highlight=0,SG24-6623-00](http://publib.b.boulder.ibm.com/Redbooks.nsf/9445fa5b416f6e32852569ae006bb65f/7b1ce6b3913caf b386256bdb007595e8?OpenDocument&Highlight=0,SG24-6623-00). See <http://www.redbooks.ibm.com/redbooks.nsf/Redbooks?SearchView&Query=linux+cluster&SearchMax=4999> for additional information about implementing a cluster.

3. Understanding xCAT's functions and features

This section explains xCAT's uses and features.

Understanding what drives xCAT's design and architecture

xCAT's architecture and feature set have two major drivers:

1. **Real world requirements:** The features in xCAT are a result of the requirements met in hundreds of real cluster implementations. When users have had needs that xCAT or other cluster management solutions could not meet, xCAT has risen to the challenge. Over the last few years, this process has been repeatedly applied, resulting in a modular toolkit that represents best practices in cluster management and a flexibility that enables it to change rapidly in response to new requirements and work with many cluster topologies and architectures.
2. **Unmatched Linux clustering experience:** The people involved with xCAT's development have used xCAT to implement many of the world's largest Linux clusters and a huge variety of different cluster types. The challenges faced during this work has resulted in features that enable xCAT to power all types of Linux clusters from the very small to the largest ever built.

Understanding what type of Clusters xCAT benefits

xCAT works well with the following cluster types:

1. **High Performance (HPC):** such as computing physics, seismic, CFD, FEA, weather, bioinformatics and other simulations.
2. **Horizontal Scaling (HS):** such as Web farms.
3. **Administrative:** A very convenient platform, although non-traditional, to install and administer a number of Linux machines.
4. **Microsoft Windows and other Operating Systems:** With xCAT's cloning and imaging support, it can be used to rapidly deploy and conveniently manage clusters with compute nodes that run Windows or any other operating system.

xCAT's current features:

1. OS/Distribution support Any OS on compute nodes via OS agnostic imaging support.
2. Hardware Control Remote Power control (on/off state) through IBM Management Processor Network, BMC, and/or APC Master Switch.
3. Hardware Control Remote software reset (rpower).
4. Hardware Control Remote Network BIOS/firmware update and configuration on IBM hardware.
5. Hardware Control Remote OS console with pluggable support for a number terminal servers.

6. Hardware Control Remote POST/BIOS console through the IBM Management Processor Network and with terminal servers.
7. Boot Control Ability to remotely change boot type (network or local disk) with syslinux.
8. Automated parallel install using scripted RedHat kickstart, SUSE Linux autoyast, on ia32, x86_64, ppc, and ia64.
9. Automated parallel install using imaging with other Linux distributions, Widows, and other operating systems.
10. Automated network installation with supported PXE NICs, with etherboot or BootP on supported NICs without PXE.
11. Monitoring hardware alerts and email notification with IBM's Management Processor Network and SNMP alerts.
12. Monitoring remote vitals such as fan speed, temperature, and more with IBM's Management Processor Network.
13. Monitoring remote hardware event logs with IBM's Management Processor Network/IPMI Interface.
14. Administration utilities such as parallel remote shell, ping, rsync, and copy.
15. Administration utilities such as remote hardware inventory with IBM's Management Processor Network.
16. Software Stack PBS and Maui schedulers to build scripts, documentation, automated setup, extra related utilities, and deep integration.
17. Software Stack Myrinet to automate setup and installation.
18. Software Stack MPI to build scripts, documentation, and automated setup for MPICH, MPICH-GM, and LAM.
19. Usability command line utilities for all cluster management functions.
20. Usability single operations can be applied in parallel to multiple nodes with very flexible and customizable group/range functionality.
21. Flexible support for various user defined node types.
22. Diskless support using warewulf.

4. Getting the xCAT software distribution

This section explains where and how to get the xCAT software distribution.

1. Download the latest version of xCAT. 3 of the 5 required packages can be downloaded from <http://www.alphaworks.ibm.com/tech/xCAT/> to the */opt* directory. The fourth file can be downloaded from <http://www-rcf.usc.edu/~garrick/> to the */opt* directory. The fifth file can be downloaded from <http://www.xcat.org/patch/> to your desktop. The patch file will be decompressed into */opt/xcat* after *./setupxcat* has been run.
2. Download the latest firmware and hardware configuration software from <http://publib.boulder.ibm.com/cluster/>.

Note: Additional documentation is accessible in the */opt/xcat/doc* folder once the *xCAT* *.tar* files are decompressed. For assistance with building, maintaining, and administering the xCAT cluster or an xCAT feature request, see the xCAT user mailing list, your IBM sales rep, or other IBM point of contact.

5. Understanding cluster components, connections, and architecture

This document is based on a basic 32 node cluster that uses serial terminal servers for out-of-band console access, an APC Master Switch, IBM's Service Processor Network for remote hardware management, Ethernet, and Myrinet as the basis of most of its examples. All network devices that must be statically set in the IBM Cluster 1350 are pre-configured using the manufacturing defaults listed in the "IBM Manufacturing defaults for all items in Cluster 1350 Clusters" table.

The following three examples describe some of the detail of this example cluster:

Components / Rack Layout

The following hardware is positioned in the rack, starting from the bottom and moving towards the top:

1. The Myrinet switch: Used for high-speed, low-latency inter-node communication. A cluster may not have Myrinet, if the cluster is not running parallel jobs that do heavy message passing.
2. Nodes 1-16: The first 16 compute nodes. Note that every 8th node has an MPA (Management Processor Adaptor) installed. The configuration may have RSA adapters, ASMA adapters, or BMCs. These cards enable the SPN (Service Processor Network) to function and remote hardware management to be performed. Newer machines do not require a RSA or MPA because they contain a built in BMC (Baseboard Management Controller) which uses the IPMI protocol for management. The BMC is internal hardware.
3. Monitor/Keyboard: This is for local input/output function.
4. Terminal servers: The terminals enable serial consoles from all of the compute nodes to be accessible from the management node. This feature is very useful during system setup and after setup administration. Serial Over LAN (SOL) can be used to emulate a terminal server setup if the cluster does not have a terminal server.
5. APC master switch: This enables remote power control of devices that are not part of the Service Processor Network, such as terminal servers, Myrinet switch, and ASMA adapters.

Ethernet Switch
node32
... nodes 27 - 31
node26
node25 MPA
node24
... nodes 19 - 24
node18
node17 MPA
Management Node
apc1 APC Master Switch
ts2 Terminal Servers
ts1
Monitor / Keyboard
node16
... nodes 11 - 17
node10
node09 MPA
node08
... nodes 03 - 07
node02
node01 MPA has MPA card
Myrinet Switch

6. The management node: The management node is where the rest of the nodes are installed from and the cluster is managed.
7. Nodes 17-32: The rest of the compute nodes with a Management Processor card every 8th node.
8. Ethernet switch

Networks

All IBM Cluster 1350 network devices that must be statically set are preconfigured with the manufacturing defaults listed in the “IBM Manufacturing defaults for all items in Cluster 1350 Clusters” table.

The following table lists the networks that are used in the rest of this document's examples.

Notes:

1. The listing of attached devices to the right.
2. The external network is the organization's main network. In this example, only the management node has connectivity to the external network.
3. The Ethernet switch hosts both the cluster and management network on separate VLANs.
4. The cluster network connects the management node to the compute nodes. A private class B network that has no connectivity to the external network is recommended during configuration. This is often the easiest way to configure the cluster and a good practice if the configuration might grow to more than 254 nodes.
5. The management network is a separate network used to connect all devices associated with cluster management such as terminal servers, BMC, and ASMA cards, to the management node.
6. Parallel jobs use the message passing network for interprocess communication. Our example uses a separate private class B network over Myrinet. If Myrinet is not being used, this network could be the same as the cluster network. For example, any required message passing could be done over the cluster network.

IBM Manufacturing defaults for all items in Cluster 1350 Linux Clusters.

A copy of this will ship with the cluster.

Hostnames and IP Addressing Scheme

This table shows the network addressing / hostnames used to identify the various IBM Cluster 1350 cluster components.

IP Address	Hostname	Component
172.20.0.1	mgt.cluster.com	management node eth0 (cluster VLAN)
172.30.0.1	mgt1.cluster.com	management node eth1(management VLAN)
172.29.0.1	bmc	mgt node alias eth0:1 (cluster VLAN)
172.29.101.1	bmc001	e325/eServer 326/xSeries 336/xSeries 346 node bmcs
172.20.1.1	storage001	storage nodes
172.30.2.1	triton001	FAST storage controllers
172.20.101.1	node001	compute nodes
172.40.101.1	hca001	HCA cards ib0 (cluster vlan)
172.20.4.1	user001	Usenodes
172.30.10.1	myri001	Myrinet or InfiniBand TopSpin TS120
	ib001	Voltaire 9024
172.30.20.1	ts001	Terminal Servers (MRV LX-32, LX-48)
172.30.30.1	rsa001	RSA cards
172.30.50.1	cisco3508-001 smc8624-001 smc8648-001 cisco3750-001	Cisco 3500, 3700 Series or SMC switches
172.30.60.1	apc001	APC
172.30.70.1	rcm001	RCM
172.30.80.1	cisco6503-001 Cisco6509-001 Force 10	Cisco 4000 series or 6500 series switches or Force 10
172.30.101.1	sm001/mm001	IBM BladeCenter switch/management modules

Compute Node IP Addressing: Example:.rack 1, node 1 = 172.20.101.1

	<u>Network</u>	<u>Rack Number</u>	<u>Node</u>
Rack1	172.20	101	1 to 84
Rack2	172.20	102	1 to 84

Note: Node numbering increases from bottom of rack upward and from left to right for IBM BladeCenters.

IBM eServer 325 BMC IP Addressing: Example:.rack 1, node 1 = 172.29.101.1

	<u>Network</u>	<u>Rack Number</u>	<u>Node</u>
Rack1	172.29	101	1 to 40
Rack2	172.29	102	1 to 40

IBM BladeCenter Switch/Management Module IP Addressing

<u>Network</u>	<u>BC Number</u>	<u>Bay Location Number</u>
172.30	101	1 to 4 (SM bay 1 - SM bay 4) or, 5 (MM ext) or, 6(MM int)
172.30	102	1 to 4 (SM bay 1 - SM bay 4) or, 5 (MM ext) or, 6(MM int)

Examples

172.30.104.3 is the switch module fitted to bay 3 in IBM BladeCenter 4

172.30.106.5 is the external port (eth0) for the management module in IBM BladeCenter 6

172.30.106.6 is the internal port (eth1) for the management module in IBM BladeCenter 6

InfiniBand HCA IP Addressing: Example:.rack 1, hca 1 = 172.40.101.1

	Network	Rack Number	HCA
Rack1	172.40	101	1 to 84
Rack2	172.40	102	1 to 84

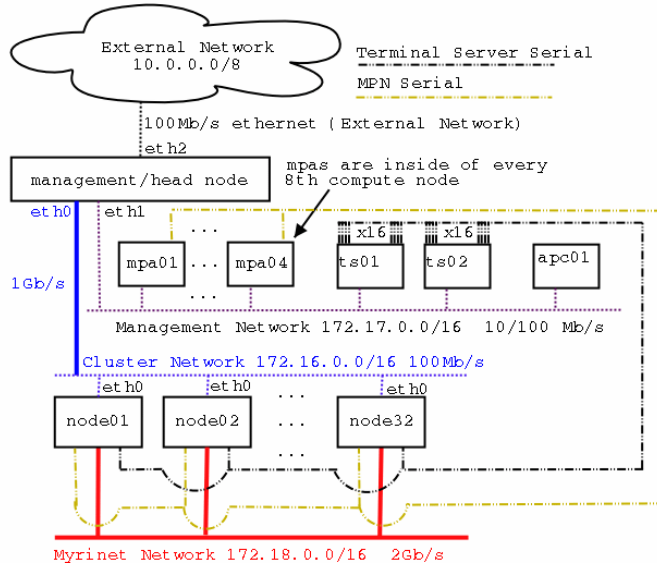
HCA numbering increases from bottom of rack upward and from left to right for IBM BladeCenters.

Manufacturing default userids and passwords.

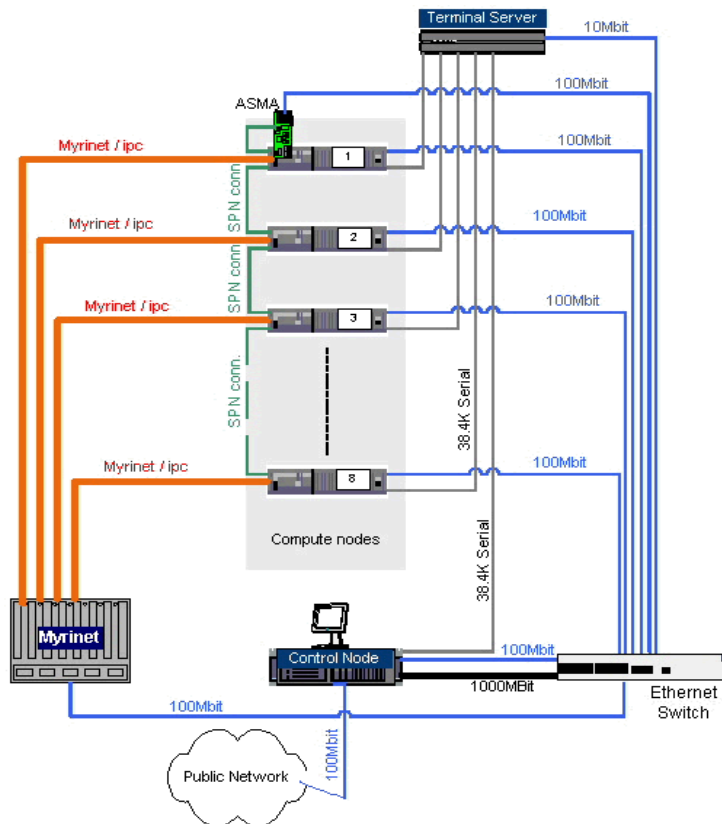
Hardware:	USERID	Password	Comments
APC Switch	apc	apc	Defaults
ITouch Terminal	access	system	Defaults
InReach Terminal	InReach	access	Defaults
Cisco Switch	(see comment)	ibm1350	just press Enter for userid
SMC Switch	admin	admin	Defaults
FAStT	(see comment)	infiniti	just press Enter for userid
RSAs	USERID	PASSWORD	(Note: the “0” in passw0rd is a zero)
TopSpin	super	super	defaults
Voltaire	(none)	123456	defaults
Force 10	(none)	(none)	defaults

Operating System:	USERID	Password
management server	root	ibm1350
nodes	admin	cluster
	root	cluster

Connections



Another Connections Diagram



Other architecture notes:

1. The compute nodes have no access to the external network.
2. The compute nodes get DNS, DHCP, and NIS services from the management node.
3. NIS is used to distribute username and password information to the compute nodes and the management node is the NIS master.
4. The management node is the only node with access to the management network.
5. PBS and Maui are used to schedule/run jobs on the cluster.
6. Users can only access compute nodes when the scheduler has allocated nodes to them and then only with SSH.
7. Jobs will use MPICH or LAM for message passing.

6. Configuring the Ethernet switch and VLANs

This section describes configuring the Ethernet switch. The examples are based on the Cisco 3750. The commands associated with the hardware actually configured in the Cluster 1350, may vary slightly from this documentation. Consult the documentation that came in the ship group for details on the specific switch used in the configuration.

Note: All switches associated with the IBM Cluster 1350 should be set up from manufacturing with default settings. Review the previous section for information.

Setup VLANs and Configure Ethernet Switches

If you have separate subnets for the management and compute networks, similar to the example in this documentation, configure the VLANs on the Ethernet switches. Use VLANs to separate the ports associated with the management and cluster subnets.

Connect the management node's COM1 to the switch's console port.

```
cisco> cu -l /dev/ttyS0 -s 9600
```

Login in and enable

Note: Refer to Section 5 for manufacturing defaults

Assign an IP to the default VLAN:

```
cisco> conf t
cisco> int vlan1
cisco> ip address 172.17.5.1 255.255.0.0
cisco> exit
```

Allow telnet access and set telnet password to “cisco”:

```
cisco> conf t
cisco> line tty 0 15
cisco> login password cisco
cisco> exit
```

Set enable and console password:

```
cisco> conf t
cisco> enable password cisco
cisco> exit
```

Set console password:

```
cisco> conf t
cisco> line vty 0 4
cisco> password cisco
cisco> exit
```

Extreme Networks:

```
> config vlan Default ipaddress 172.16.5.1 255.255.0.0
```

Changing Port Settings

Setup 'spanning-tree portfast', without this option, DHCP may fail because it takes too long for a port to come online after a machine powers up. Do not set spanning-tree portfast on ports that will connect to other switches. Do the following on each port on your switch:

```
cisco> conf t
```

```
cisco> int Fa/1
cisco> spanning-tree portfast
cisco> exit
cisco> int Fa/2
cisco> spanning-tree portfast
cisco> exit
```

Setting up Remote Logging

The switch sends all of its logging information to the management node's management interface. Remote syslog is not enabled at this point.

Cisco:

```
cisco> conf t
cisco> logging 172.17.5.1
cisco> exit
```

Extreme Networks:

```
> config syslog 172.17.5.1
```

Setting up VLANs

Cisco: For each interface with a VLAN (Fa0/1...Fa0/32) and gig ports.

```
cisco> interface Fa0/1
cisco> switchport mode access
cisco> switchport access vlan 2
cisco> exit
```

Extreme Networks:

```
> create vlan man
> config vlan man tag 2
> config vlan man ipaddress 172.17.5.1 255.255.0.0
> config Default delete port 1,2,3,4 # the ports you want in the managemnet VLAN
> config man add port 1,2,3,4
> show vlan
```

VLANs with Multiple Switches

Cisco:

```
cisco> configure terminal
cisco> interface Gi0/1
cisco> switchport mode trunk
cisco> switchport trunk encapsulation isl # you should prob use the standard encap
instead
cisco> exit
```

Extreme Networks:

```
unconfigure switch
config Default delete port (deletes unwanted ports in management VLAN)
create vlan cluster
config vlan cluster tag 2
config cluster add port (add ports cluster VLAN)
show vlan
save
```

Saving your changes

Make certain the switch configuration is saved in case the switch is rebooted.

Cisco:

```
cisco> write mem
```

Extreme Networks:

```
extreme> save
```

7. Installing the operating system on the management node

This section covers the steps necessary to install Linux on the management node.

Note: Before you install to prevent confusion disable all PCI adapters in BIOS.

The first step in building an xCAT cluster is installing Linux on the management node. See the following URLs for more information:

xSeries 346

http://www.ibm.com/pc/support/site.wss/search.do?free_text=Install&qtxbrand=IBM%2BPC%2BServer&qtxfamily=xSeries%2B346&qtxnav=es&qtxdoctype=Operating%20system%20installation

xSeries 336:

http://www-307.ibm.com/pc/support/site.wss/search.do?free_text=Install&qtxbrand=IBM+PC+Server&qtxfamily=xSeries+336&qtxnav=es&qtxdoctype=Operating%20system%20installation

eServer 326:

http://www-307.ibm.com/pc/support/site.wss/search.do?free_text=Install&qtxbrand=IBM%2BPC%2BServer&qtxfamily=eServer%2B326&qtxnav=es&qtxdoctype=Operating%20system%20installation

Notes:

1. Your management node may require specific drivers please consult the machine specific setup above for instructions.
2. For more detailed setup instructions or troubleshooting assistance on individual IBM eServer 326, xSeries 336, or xSeries 346, refer to the link above.

Create and Configure RAID Devices if Necessary

If you are using LSI, HostRaid, or ServeRAID devices in the management node, use the “LSI/HostRaid/ServeRAID flash/configuration” CD to update the LSI/HostRaid/ServeRAID firmware to version 4.84 and define the RAID volumes. If other nodes exist with hardware RAID, update and configure them now. The latest firmware can be downloaded from <http://www.pc.ibm.com/qtechinfo/MIGR-495PES.html>.

NIS Notes

Note: If you plan on interacting with an external NIS server, make sure the server supports MD5 passwords and shadow passwords. If it does not support these features, do not turn them on during the install of the management node.

Partition Notes

File System Type should be set up as *ext3*.

Recommended minimum drive partitioning scheme for the management node:

/boot (200 MB)

SWAP (1.5 x physical memory, not to exceed 2GB)

/var (2GB)

/ (the rest of the disk)

Select **No firewall**. This is for xCat installation purposes and can be changed after the configuration of the cluster has completed.

Select **Disable Selinux**.

Select **CUSTOMIZED SOFTWARE PACKAGES TO BE INSTALLED** from the software selection menu.

Scroll down to **Miscellaneous**.

Select the **Everything** option.

Note: If this is the first time installing Red Hat 4, everything is a check box at the end of the selection.

Install all 5, RHWS4 CDs.

You will be prompted during reboot to install and load the Extras CD.

It is recommended to create a normal user other than *root* during the install.

Start the newly installed system by rebooting and then login as root.

Open a terminal
>*updatedb*

Turn off unwanted services (general)

To turn off some of the network services turned on by default during the installation process, use the following commands:

To view installed services:
chkconfig --list | grep ':on'

To turn off a service:
chkconfig --level 0123456 <service> off

Turn off unwanted services (specific)

Use the following commands to turn off all unnecessary services:

```
chkconfig --level 0123456 autofs off  
chkconfig --level 0123456 isdn off  
chkconfig --level 0123456 iptables off  
chkconfig --level 0123456 ip6tables off  
chkconfig --level 0123456 rhnsd off  
chkconfig --level 0123456 rawdevices off  
chkconfig --level 0123456 kudzu off  
chkconfig --level 0123456 FreeWnn off  
chkconfig --level 0123456 arptables_jf off  
chkconfig --level 0123456 canna off  
chkconfig --level 0123456 cups off  
chkconfig --level 0123456 hpoj off  
chkconfig --level 0123456 alsasound off
```

8. Configuring networking on the management node

This section describes network setup on the management node.

Notes:

1. If you are using the onboard e1000/bcm5700 network interface card, download the latest driver from <http://publib.boulder.ibm.com/cluster/> and build it into the Linux kernel.
2. A USB storage device for use with the IBM eServer 326 and xSeries 336 servers will be required.

Remove *tg3* driver

```
>service network stop
>lsmod | grep tg3
>rmmod tg3
```

Download *<Driver>.src.rpm*

```
>rpm -ivh <Driver>.src.rpm
>cd /usr/src/redhat/SOURCES/
>tar -zxvf <Driver>.tgz
>make
>make install
>insmod <Driver.ko>
```

Configure the network devices using the Neat utility.

```
>neat
```

For example, in the Neat utility select:

```
ENT0
Statically Set IP
IP 172.20.0.1
Subnet Mask 255.255.0.0
OK
```

```
ENT1
Statically Set IP
IP 172.30.0.1
Subnet Mask 255.255.0.0
OK
```

```
New (will be BMC and alias for ENT0)
Broadcom ETH 0:1
IP 172.29.0.1
SM 255.255.0.0
OK
```

Reboot your machine and enable PCI devices in BIOS, the devices must now configure manually.

>*Init 6*

After the machine reboots, log back in as *root* and configure the PCI Network Devices using the Neat utility. See example below.

```
New (will be external connection)
Select devxxxx
Automatically obtain IP address setting
Automatically obtain DNS info from provider
Forward
Apply
```

```
ACTIVATE each device
```

```
>Service Network Restart
```

```
>ifconfig | more
```

Example ifconfig scripts:

```
dev23991 Link encap:Ethernet HWaddr 00:10:18:0C:BB:93
inet addr:10.1.1.73 Bcast:10.1.1.255
Mask:255.255.255.0
inet6 addr: fe80::210:18ff:fe0c:bb93/64
Scope:Link
UP BROADCAST RUNNING MULTICAST MTU:1500
Metric:1
RX packets:22337 errors:0 dropped:0 overruns:0
frame:0
TX packets:7402 errors:0 dropped:0 overruns:0
carrier:0
collisions:0 txqueuelen:1000
RX bytes:2288336 (2.1 MiB) TX bytes:1179835 (1.1
MiB)
Interrupt:201 Memory:deff0000-df000000

eth0 Link encap:Ethernet HWaddr 00:0D:60:55:44:1A
inet addr:172.20.0.1 Bcast:172.20.255.255
Mask:255.255.0.0
inet6 addr: fe80::20d:60ff:fe55:441a/64
Scope:Link
UP BROADCAST RUNNING MULTICAST MTU:1500
Metric:1
RX packets:17495 errors:0 dropped:0 overruns:0
frame:0
TX packets:16985 errors:0 dropped:0 overruns:0
carrier:0
collisions:0 txqueuelen:1000
```

```

RX bytes:1360076 (1.2 MiB) TX bytes:12581705
(11.9 MiB)
Interrupt:169 Memory:dcff0000-dd000000

eth0:1 Link encap:Ethernet HWaddr 00:0D:60:55:44:1A
inet addr:172.29.0.1 Bcast:172.29.255.255
Mask:255.255.0.0
UP BROADCAST RUNNING MULTICAST MTU:1500
Metric:1
Interrupt:169 Memory:dcff0000-dd000000

eth1 Link encap:Ethernet HWaddr 00:0D:60:55:44:1B
inet addr:172.30.0.1 Bcast:172.30.255.255
Mask:255.255.0.0
inet6 addr: fe80::20d:60ff:fe55:441b/64
Scope:Link
UP BROADCAST RUNNING MULTICAST MTU:1500
Metric:1
RX packets:1153177 errors:2 dropped:0 overruns:0
frame:2
TX packets:1152679 errors:0 dropped:0 overruns:0
carrier:0
collisions:0 txqueuelen:1000
RX bytes:334595931 (319.0 MiB) TX bytes:80676258
(76.9 MiB)
Interrupt:169 Memory:daff0000-db000000

lo Link encap:Local Loopback
inet addr:127.0.0.1 Mask:255.0.0.0
inet6 addr: ::1/128 Scope:Host
UP LOOPBACK RUNNING MTU:16436 Metric:1
RX packets:2333 errors:0 dropped:0 overruns:0
frame:0
TX packets:2333 errors:0 dropped:0 overruns:0
carrier:0
collisions:0 txqueuelen:0
RX bytes:1820451 (1.7 MiB) TX bytes:1820451 (1.7
MiB)

```

Create your */etc/hosts* file.

```
>vi /etc/hosts
```

Note: This file is provided on the Cluster 1350 configuration disks from manufacturing that can be found in your ship group. Make sure all devices are entered, such as terminal servers, switches and hardware management devices.

The following is an example of the */etc/hosts* for the example cluster:

Note: It is recommended to insert the fully qualified domain name before the short name.

```
# Localhost
```



```

127.0.0.1      localhost.localdomain localhost
##### Management Node #####
# cluster interface (eth0) GigE
172.20.0.1  mgmt1.mydomain.com      mgmt1
# management interface (eth1)
172.30.0.1  mgmt2.mydomain.com      mgmt2
# external interface (eth2)
10.0.0.1    external.mydomain.com     external
##### Management Equipment #####
# RSA adapters. You might have ASMA cards instead
172.30.30.1  rsa001.mydomain.com      rsa001
172.30.30.2  rsa002.mydomain.com      rsa002
172.30.30.3  rsa003.mydomain.com      rsa003
172.30.30.4  rsa004.mydomain.com      rsa004
# Terminal Servers
172.17.2.1  ts01.mydomain.com        ts01
172.17.2.2  ts02.mydomain.com        ts02
# APC Master Switch
172.17.3.1  apc1.mydomain.com        apc01
# Myrinet Switch's Ethernet management port
172.17.4.1  myri01.mydomain.com     myri01
# Ethernet Switch
172.17.5.1  Ethernet01mydomain.com  Ethernet01c
172.16.5.1  Ethernet01.mydomain.com Ethernet01
##### Compute Nodes #####
172.20.101.1 node01.mydomain.com      node01
172.30.10.1  node01-myri0.mydomain.com node01-myri0
172.20.101.2 node02.mydomain.com      node02
172.30.10.2  node02-myri0.mydomain.com node02-myri0
172.20.101.3 node03.mydomain.com      node03
172.30.10.3  node03-myri0.mydomain.com node03-myri0
172.20.101.4 node04.mydomain.com      node04
172.30.10.4  node04-myri0.mydomain.com node04-myri0
172.20.101.5 node05.mydomain.com      node05
172.30.10.5  node05-myri0.mydomain.com node05-myri0
172.20.101.6 node06.mydomain.com      node06
172.30.10.6  node06-myri0.mydomain.com node06-myri0
172.20.101.7 node07.mydomain.com      node07
172.30.10.7  node07-myri0.mydomain.com node07-myri0
172.20.101.8 node08.mydomain.com      node08
172.30.10.8  node08-myri0.mydomain.com node08-myri0
172.20.101.9 node09.mydomain.com      node09
172.30.10.9  node09-myri0.mydomain.com node09-myri0
172.20.101.10 node10.mydomain.com     node10
172.30.10.10 node10-myri0.mydomain.com node10-myri0
172.20.101.11 node11.mydomain.com     node11

```

172.30.10.11 node11-myri0.mydomain.com node11-myri0
172.20.101.12 node12.mydomain.com node12
172.30.10.12 node12-myri0.mydomain.com node12-myri0
172.20.101.13 node13.mydomain.com node13
172.30.10.13 node13-myri0.mydomain.com node13-myri0
172.20.101.14 node14.mydomain.com node14
172.30.10.14 node14-myri0.mydomain.com node14-myri0
172.20.101.15 node15.mydomain.com node15
172.30.10.15 node15-myri0.mydomain.com node15-myri0
172.20.101.16 node16.mydomain.com node16
172.30.10.16 node16-myri0.mydomain.com node16-myri0
172.20.101.17 node17.mydomain.com node17
172.30.10.17 node17-myri0.mydomain.com node17-myri0
172.20.101.18 node18.mydomain.com node18
172.30.10.18 node18-myri0.mydomain.com node18-myri0
172.20.101.19 node19.mydomain.com node19
172.30.10.19 node19-myri0.mydomain.com node19-myri0
172.20.101.20 node20.mydomain.com node20
172.30.10.20 node20-myri0.mydomain.com node20-myri0
172.20.101.21 node21.mydomain.com node21
172.30.10.21 node21-myri0.mydomain.com node21-myri0
172.20.101.22 node22.mydomain.com node22
172.30.10..22 node22-myri0.mydomain.com node22-myri0
172.20.101.23 node23.mydomain.com node23
172.30.10..23 node23-myri0.mydomain.com node23-myri0
172.20.101.24 node24.mydomain.com node24
172.30.10.24 node24-myri0.mydomain.com node24-myri0
172.20.101.25 node25.mydomain.com node25
172.30.10.25 node25-myri0.mydomain.com node25-myri0
172.20.101.26 node26.mydomain.com node26
172.30.10.26 node26-myri0.mydomain.com node26-myri0
172.20.101.27 node27.mydomain.com node27
172.30.10.27 node27-myri0.mydomain.com node27-myri0
172.20.101.28 node28.mydomain.com node28
172.30.10.28 node28-myri0.mydomain.com node27-myri0

Verify the management node's network setup by:

1. Pinging all network interfaces, refer to the manufacturing defaults for verification.
2. Pinging other devices on all of the subnets, including the cluster, management, and any external devices.
3. Pinging and route through your gateway.

9. Installing xCAT

Follow the steps below to install xCAT on the management node.

1. Download the latest version of xCAT to the `/opt` directory if you have not already done so in section 4. Three of the five required packages can be downloaded from <http://www.alphaworks.ibm.com/tech/xCAT/> and the fourth can be downloaded from <http://www-rcf.usc.edu/~garrick/>. The fifth can be downloaded from <http://www.xcat.org/patch/> to your desktop. The patch file will be decompressed into `/opt/xcat` once `./setupxcat` has been run.
2. Unpack xCAT in to `/opt/`.

```
cd /opt
tar -xzyf xcat-dist-core-RCx.x.x.tgz
tar -xzyf xcat-dist-doc-RCx.x.x.tgz
tar -xzyf xcat-dist-ibm-RCx.x.x.tgz
tar -xzyf xcat-dist-oss-RCx.x.x.tgz
```

3. Use the following commands to setup xCAT.

```
>export XCATROOT=/opt/xcat
>cd $XCATROOT/sbin
>./setupxcat

>cd /opt/xcat
tar -xzyf xcat-1.2.0-RC2.2-patch.tgz
```

Update `Modules.conf` and replace **tg3** with **bcm5700**

```
>vi /etc/modules.conf
```

Enable time services (xntpd) on management node.

```
>mv -f /etc/ntp.conf /etc/ntp.conf.ORIG
```

Create a new `/etc/ntp.conf`:

```
>server 127.127.1.0
>fudge 127.127.1.0 stratum 10
>driftfile /etc/ntp/drift
```

Set time, date, and time zone with setup:

```
>chkconfig - - level 2345 ntpd on
>service ntpd restart
```

Note: It can take a few minutes before *ntpd* is working.

```
> ntpdate -q localhost
```

If working you should receive the following output:

```
server 127.0.0.1, stratum 2, offset -0.000002, delay 0.02570  
22 Jan 08:04:24 ntpdate[14540]: adjust time server 127.0.0.1 offset -0.000002 sec
```

If not working you will receive the following output:

```
no server suitable for synchronization found
```

Note: *setupxcat* must actually be run after xCAT *.tab* files are setup later on.

Add the xCAT Man Pages to *\$MANPATH* by adding the following line to */etc/man.config*:

```
> MANPATH /opt/xcat/man
```

Test out the man pages by entering:

```
> man site.tab
```

10. Setup xCAT

This section describes some of the xCAT configuration necessary for the 32 node example cluster. If configuring a cluster that differs from this example, you may have to change some settings. xCAT configuration files are located in */opt/xcat/etc*. You must setup these configuration files before proceeding.

Copy the configuration files to their required location.

Note: If this is an IBM Cluster 1350 the configuration files can be found in the ship group. Only copy the samples if the configuration files are not available.

```
> mkdir /install  
> cp /opt/xcat/samples/etc/* /opt/xcat/etc
```

Create a custom configuration by editing */opt/xcat/etc/** to work with the cluster. Read the man pages “*man site.tab*”, to learn more about the format of these configuration files. More detailed information on some of these files can be found in some of the later sections. The following are examples that work with the example 32 node cluster.

To find documented examples of the *.tab* files, go to the */opt/xcat/samples/etc* directory.

Note: If you have installed Java you may use the xTablePad or xTableWizard table generators in the */opt/bcat/lib* directory to generate your tab configuration files. You may

also use this application on a Windows machine but if the files are edited on the Microsoft Windows machine, the formatting may be wrong.

SUSE LINUX is packaged with Java installed.

Required tables:

site.tab

nodehm.tab

odelist.tab

nodepos.tab

noderes.tab

nodetype.tab

passwd.tab

postscripts.tab

postdeps.tab

snmptrapd.conf

networks.tab

mac.tab (loaded with non-collectable MACs, such as terminal servers, switches, and RSAs.)

mp.tab

mpa.tab

Required tables for clusters with terminal servers or SOL (Server Over LAN):

conserver.tab

conserver.cf

Required tables for clusters using Ethernet switches to collect MAC addresses (use the correct table for your switch):

cisco.tab

summit48i.tab

blackdiamond.tab

switch.tab

Required tables for clusters using IPMI management:

ipmi.tab

Required table for APC Master Switch:

apc.tab

Required table for APC Master Switch Plus:

apcp.tab

Required table for xCAT flash support:

nodemodel.tab

Required table for EMP support:

emp.tab

Required table for Baytech support:

baytech.tab

Required table for xCAT GPFS support:

gpfs.tab

Table for IPMI support. Required for systems having a different IPMI IP address than node address (for example, the IBM e325):

ipmi.tab

site.tab

/opt/xcat/etc/site.tab

site.tab control most of xCAT's global settings.

man site.tab for information on what each field means.

this example uses 'c' as a subdomain private to the cluster and

10.0.0.1 as the corp DNS server (forwarder).

rsh /usr/bin/ssh

rsh /usr/bin/ssh

rcp /usr/bin/scp

gkhfile /opt/xcat/etc/gkh

tftpdir /tftpboot

tftpserver xcat

modify domain to match your domain name

domain mydomain.com

dnssearch mydomain.com

**# nameserver - Comma delimited list of DNS name servers IP addresses, use your
#management node IP address. (172.16.n.100)**

nameservers 192.16.100.1

forwarders 10.0.0.1

**# nets - Comma delimited list of DNS network and netmask pairs colon delimited or
#NA. Required only if this cluster contains a primary DNS server. This list
determines what /etc/hosts entries are used to create the primary DNS server files.**

nets

172.16.0.0:255.255.0.0,172.17.0.0:255.255.0.0,172.18.0.0:255.255.0.0

dnsdir /var/named/chroot

dnschroot yes

**#dnsallow - Comma delimited list of DNS network and netmask pairs colon
#delimited or NA. Required only if this cluster contains a primary or secondary
DNS server. This list determines the access permissions for primary and secondary
DNS servers contained within this cluster.**

dnsallow 172.16.0.0:255.255.0.0,172.17.0.0:255.255.0.0,172.18.0.0:255.255.0.0

**#domainaliasip - IP address aliased to cluster DNS domain name or NA. Required
only if this cluster contains a primary DNS server. Use your management IP
address**

domainaliasip *172.16.100.1*

#mxhosts - Comma delimited list of FQDN mail exchange hosts for this cluster or #NA. Required only if this cluster contains a primary DNS server. Each node will be assigned *mxhosts* as the MX records for that host.

mxhosts *mydomain.com,man-mydomain.com*

#mailhosts - Comma delimited list of mail hosts aliases for this cluster or NA. Required only if this cluster contains a primary DNS server. Each host listed in #mailhosts will be aliased as *mailhost* for the purpose of providing the cluster with a single host name for all mail.

mailhosts *man-c*

#master - Master host/node name

master *man-c*

#homefs - Default global home file system.

homefs *man-c:/home*

#localfs - Default global local file system.

localfs *man-c:/usr/local*

pbshome */var/spool/pbs*

pbsprefix */usr/local/pbs*

#Pbserver - Name of the node which is running the PBS server.

pbserver *man-c*

scheduler *maui*

xcatprefix */opt/xcat*

keyboard *us*

#timezone – Current Linux time zone.

timezone *US/Eastern*

#offutc - UTC offset.

offutc *-5*

mapperhost *NA*

#serialmac - What serial port to use to collect MAC addresses.

serialmac *0*

serialbps *9600*

snmpc *public*

#snmpd - The IP address to collect SNMP traps.

snmpd 172.17.100.1
poweralerts Y

#timeservers - Comma delimited list of IP addresses for nodes to sync their clocks.

timeservers man-c
logdays 7
installdir /install
clustername Clever-cluster-name

#dhcpver - set this to 3 since we are using DHCP version 3

dhcpver 2
dhcpconf /etc/dhcpd.conf

**#dynamicr - This is the range of IP addresses assigned for node discovery.
Comment this out by placing “#” at the beginning of the line**

dynamicr eth0,ia32,172.30.0.1,255.255.0.0,172.30.1.1,172.30.254.254

#usernodes - A comma delimited list of nodes users are allow to login to.

usernodes man-c

#usermaster - The single node that users accounted are added to.

usermaster man-c

#nisdomain and nismaster. Set to NA, NIS is beyond the scope of this class.

nisdomain NA
nismaster NA

nisslaves NA
homelinks NA
chagemin 0
chagemax 60
chagewarn 10
chageinactive 0
mpcliroot /opt/xcat/lib/mpcli

#End of site.tab

nodelist.tab

/opt/xcat/etc/nodelist.tab

*nodelist.tab* contains a list of nodes and defines groups that can be used in commands.

Use *man nodelist.tab* for more information.

node01 all,rack1,compute,myri,mpn1
node02 all,rack1,compute,myri,mpn1

node03 all,rack1,compute,myri,mpn1
node04 all,rack1,compute,myri,mpn1
node05 all,rack1,compute,myri,mpn1
node06 all,rack1,compute,myri,mpn1
node07 all,rack1,compute,myri,mpn1
node08 all,rack1,compute,myri,mpn1
node09 all,rack1,compute,myri,mpn2
node10 all,rack1,compute,myri,mpn2
node11 all,rack1,compute,myri,mpn2
node12 all,rack1,compute,myri,mpn2
node13 all,rack1,compute,myri,mpn2
node14 all,rack1,compute,myri,mpn2
node15 all,rack1,compute,myri,mpn2
node16 all,rack1,compute,myri,mpn2
node17 all,rack1,compute,myri,mpn3
node18 all,rack1,compute,myri,mpn3
node19 all,rack1,compute,myri,mpn3
node20 all,rack1,compute,myri,mpn3
node21 all,rack1,compute,myri,mpn3
node22 all,rack1,compute,myri,mpn3
node23 all,rack1,compute,myri,mpn3
node24 all,rack1,compute,myri,mpn3
node25 all,rack1,compute,myri,mpn4
node26 all,rack1,compute,myri,mpn4
node27 all,rack1,compute,myri,mpn4
node28 all,rack1,compute,myri,mpn4
node29 all,rack1,compute,myri,mpn4
node30 all,rack1,compute,myri,mpn4
node31 all,rack1,compute,myri,mpn4
node32 all,rack1,compute,myri,mpn4
rsa01 nan,mpa
rsa02 nan,mpa
rsa03 nan,mpa
rsa04 nan,mpa
ts01 nan,ts
ts02 nan,ts
myri01 nan

mpa.tab

```
/opt/xcat/etc/mpa.tab  
#service processor adapter management  
#  
#type = asma,rsa  
#name = internal name (must be unique)  
# internal name should = node name
```

```

#      if rsa/asma is primary management
#      processor
#number = internal number (must be unique and > 10000)
#command = telnet,mpcli
#reset   = http(ASMA only),mpcli,NA
#dhcp    = Y/N(RSA only)
#gateway = default gateway or NA (for DHCP assigned)
#
rsa01  rsa,rsa01,10001,mpcli,mpcli,NA,N,NA
rsa02  rsa,rsa02,10002,mpcli,mpcli,NA,N,NA
rsa03  rsa,rsa03,10003,mpcli,mpcli,NA,N,NA
rsa04  rsa,rsa04,10004,mpcli,mpcli,NA,N,NA

```

mp.tab

```

/opt/xcat/etc/mp.tab
# mp.tab defines how the Service processor network is setup.
# node07 is accessed via the name 'node07' on the RSA 'rsa01', etc.
# man asma.tab for more information until the man page to mp.tab is ready
node01 rsa01,node01
node02 rsa01,node02
node03 rsa01,node03
node04 rsa01,node04
node05 rsa01,node05
node06 rsa01,node06
node07 rsa01,node07
node08 rsa01,node08
node09 rsa02,node09
node10 rsa02,node10
node11 rsa02,node11
node12 rsa02,node12
node13 rsa02,node13
node14 rsa02,node14
node15 rsa02,node15
node16 rsa02,node16
node17 rsa03,node17
node18 rsa03,node18
node19 rsa03,node19
node20 rsa03,node20
node21 rsa03,node21
node22 rsa03,node22
node23 rsa03,node23
node24 rsa03,node24
node25 rsa04,node25
node26 rsa04,node26
node27 rsa04,node27

```

node28 rsa04,node28
node29 rsa04,node29
node30 rsa04,node30
node31 rsa04,node31
node32 rsa04,node32

apc.tab

/opt/xcat/etc/apc.tab

apc.tab defines the relationship between nodes and APC
MasterSwitches and their assigned outlets. In our example,
the power for asma1 is plugged into the 1st outlet the
APC MasterSwitch, etc.

rsa01 apc1,1

rsa02 apc1,2

rsa03 apc1,3

rsa04 apc1,4

ts01 apc1,5

ts02 apc1,6

myri01 apc1,7

conserver.cf

/opt/xcat/etc/conserver.cf

conserver.cf defines how serial consoles are accessed. Our example
uses the ELS terminal servers and node01 is connected to port 1
on ts01, node02 is connected to port 2 on ts01, node17 is connected to
port 1 on ts02, etc.

man *conserver.cf* for more information

#

The character '&' in logfile names are substituted with the console
name. Any logfile name that does not begin with a '/' has LOGDIR
prepended to it. So, most consoles will just have a '&' as the logfile
name which causes */var/consoles/* to be used.

#

LOGDIR=/var/log/consoles

#

list of consoles we serve

name : tty[@host] : baud[parity] : logfile : mark-interval[m|h|d]

name : !host : port : logfile : mark-interval[m|h|d]

name : |command : : logfile : mark-interval[m|h|d]

#

node01:!ts01:3001:&:

node02:!ts01:3002:&:

node03:!ts01:3003:&:

node04:!ts01:3004:&:

node05:!ts01:3005:&:

node06:!ts01:3006:&:

node07:!ts01:3007:&:

node08:!ts01:3008:&:

node09:!ts01:3009:&:

```
node10:!ts01:3010:&:
node11:!ts01:3011:&:
node12:!ts01:3012:&:
node13:!ts01:3013:&:
node14:!ts01:3014:&:
node15:!ts01:3015:&:
node16:!ts01:3016:&:
node17:!ts02:3001:&:
node18:!ts02:3002:&:
node19:!ts02:3003:&:
node20:!ts02:3004:&:
node21:!ts02:3005:&:
node22:!ts02:3006:&:
node23:!ts02:3007:&:
node24:!ts02:3008:&:
node25:!ts02:3009:&:
node26:!ts02:3010:&:
node27:!ts02:3011:&:
node28:!ts02:3012:&:
node29:!ts02:3013:&:
node30:!ts02:3014:&:
node31:!ts02:3015:&:
node32:!ts02:3016:&:
%%
#
# list of clients we allow
# {trusted|allowed|rejected} : machines
#
trusted: 127.0.0.1
```

conserver.tab

/opt/xcat/etc/conserver.tab

conserver.tab defines the relationship between nodes and conserver servers. Our example uses only one conserver on the localhost. Use *man conserver.tab* for more information.

```
node01localhost,node01
node02localhost,node02
node03localhost,node03
node04localhost,node04
node05localhost,node05
node06localhost,node06
node07localhost,node07
node08localhost,node08
node09localhost,node09
node10localhost,node10
```

node11localhost,node11
node12localhost,node12
node13localhost,node13
node14localhost,node14
node15localhost,node15
node16localhost,node16
node17localhost,node17
node18localhost,node18
node19localhost,node19
node20localhost,node20
node21localhost,node21
node22localhost,node22
node23localhost,node23
node24localhost,node24
node25localhost,node25
node26localhost,node26
node27localhost,node27
node28localhost,node28
node29localhost,node29
node30localhost,node30
node31localhost,node31
node32localhost,node32

nodehm.tab

/opt/xcat/etc/nodehm.tab

```
#  
#node hardware management  
#  
#power    = mp,baytech,emp,apc,apcp,NA  
#reset    = mp,apc,apcp,NA  
#cad      = mp,NA  
#vitals   = mp,NA  
#inv      = mp,NA  
#cons     = consver,tty,rtel,NA  
#bioscons = rcons,mp,NA  
#eventlogs = mp,NA  
#getmacs  = rcons,cisco3500  
#netboot  = pxe,eb,ks62,elilo,file:,NA  
#eth0     = eepr100,pcnet32,e100,bcm5700  
#gcons    = vnc,NA  
#serialbios = Y,N,NA  
#  
#node  
    power,reset,cad,vitals,inv,cons,bioscons,eventlogs,getmacs,netboot,eth0,gcons,ser  
ialbios
```



```

#TFTP      = Where is my TFTP server?
#          Used by makedhcp to setup /etc/dhcpd.conf
#          Used by mkks to setup update flag location
#NFS_INSTALL = Where do I get my files?
#INSTALL_DIR = From what directory?
#SERIAL      = Serial console port (0, 1, or NA)
#USENIS      = Use NIS to authenticate (Y or N)
#INSTALL_ROLL = Am I also an installation server? (Y or N)
#ACCT        = Turn on BSD accounting
#GM          = Load GM module (Y or N)
#PBS         = Enable PBS (Y or N)
#ACCESS      = access.conf support
#GPFS        = Install GPFS
#INSTALL_NIC = eth0, eth1, ... or NA
#
#node/group
      TFTP,NFS_INSTALL,INSTALL_DIR,SERIAL,USENIS,INSTALL_ROLL,ACCT,G
M,PBS,ACCESS,GPFS,INSTALL_NIC
#
compute man-c,man-c,/install,0,N,N,N,Y,Y,Y,N,eth0
nan      man-c,NA,NA,NA,NA,NA,NA,NA,NA,NA,NA,NA,NA,NA

```

nodetype.tab

nodetype.tab maps nodes to types of installs. The example below only uses one type. For more information, *man nodetype*.

Note: *nodetype.tab* can not contain comments.

/opt/xcat/etc/nodetype.tab

```

node01 compute73
node02 compute73
node03 compute73
node04 compute73
node05 compute73
node06 compute73
node07 compute73
node08 compute73
node09 compute73
node10 compute73
node11 compute73
node12 compute73
node13 compute73
node14 compute73
node15 compute73

```



```
node16 compute73
node17 compute73
node18 compute73
node19 compute73
node20 compute73
node21 compute73
node22 compute73
node23 compute73
node24 compute73
node25 compute73
```

Continue until the last node is entered.

passwd.tab

The file *passwd.tab* defines some passwords that will be used in the cluster, *man passwd.tab* for more information

```
/opt/xcat/etc/passwd.tab
```

```
cisco      cisco
rootpw     netfinity
asmauser   USERID
asmypass   PASSWORD
```

ipmi.tab

```
node001    bmc001,"",""
node002    bmc002,"",""
node003    bmc003,"",""
node004    bmc004,"",""
```

Continue until the last node is entered.
Check tabs by running *rpower* and *rbeacon* commands.

11. Configuring the terminal servers

For Cyclades AlterPath ACS see <http://www.cyclades.com> for configuration instructions.

This section describes setting up ELS and ESP terminal servers and conserver. If the cluster has either ELSes or ESPs, skip the instructions for the terminal server type. Terminal servers enable out-of-band administration and access to the compute nodes, for example watching a compute node's console remotely before the compute node can be assigned an IP address or after the network configuration is lost.

Notes:

1. The hardware associated with the Cluster 1350 serial network will be preconfigured from manufacturing.
2. Please see the last section in this document for Cluster 1350 orders not containing terminal servers.

Learn about conserver

Conserver's website: <http://conserver.com/>.

Shutdown conserver

Before setting up the terminal servers, make sure that the conserver service is stopped:

```
>service conserver stop
```

Setup terminal servers

This section describes how to configure the Equinox ELS terminal server.

conserver.cf setup

Modify */opt/xcat/etc/conserver.cf*

Each node has an entry similar to:

```
nodeXXX:!tsx:yyyy:&:
```

Where:

x = Terminal Server Unit number and

yyyy = Terminal Server port + 3000 e.g. node1:!ts1:3001:

& = access node1 via telnet to ts1 on port 3001. 'node1' should be connected to ts1's first serial port.

conserver.cf setup If needed.

Modify */opt/xcat/etc/conserver.cf*

Each node gets a line like:

```
node001:!ts001:7001:&:          (Cyclades)
node002:!ts001:7002:&:
```

Start Conserver

```
> service conserver start
```

Test if Conserver and terminal servers are working.

```
wcons -t <node range>
```

12. Configure xCAT

This section covers configuring xCAT on the cluster.

A restart of xCAT is required after the *.tab* files are installed

Use the following commands to setup xCAT:

```
>export XCATROOT=/opt/xcat
>cd $XCATROOT/sbin
>./setupxcat
```

If not done during the OS install, edit */etc/selinux/config*.
SELINUX=disabled

Build a DNS server:

```
>makedns master
```

Check DNS with:

```
>host mgt
```

The DNS should return the IP for mgt, *172.20.0.1*.

Enter non-collectable MACs, such as terminal servers, switches, and RSA adapters in *\$XCATROOT/etc/mac.tab*.

Notes:

1. Some network devices, such as the APC Master Switch, do not have the MAC address affixed to the unit. Some devices may have the MAC printed on a piece of paper in the manual. Before installing a device in to the rack, verify the MAC

address will be visible when racked. Some network devices have a serial port that may be used to obtain the MAC.

2. Manual non-collectable MAC entries in *mac.tab* do not require a *-eth0* appended, it is optional.

13. DHCP setup and configuration

This section covers installing and configuring DHCP on the cluster.

Collect the MAC addresses of the cluster equipment and place each MAC address that requires DHCP for an IP address into `/opt/xcat/etc/<MANAGEMENT_NET>.tab`. See the man page for `macnet.tab`.

Note: If using APC master switches, include their MAC addresses into this file.

Make the Initial `dhcpd.conf` configuration file.

```
> makedhcp -new
```

Edit `dhcpd.conf` and check for anything out of the ordinary.

```
> vi /etc/dhcpd.conf
```

Use the example below to verify the contents of `/etc/dhcp.conf`

```
#xCAT 1.2.0-RC2

authoritative;
ddns-update-style none;

option option-128 code 128 = string;
option option-150 code 150 = string;
option option-160 code 160 = string;
option option-192 code 192 = string;
option option-193 code 193 = string;
option option-194 code 194 = string;
option option-195 code 195 = string;

shared-network eth0 {

    filename                "/tftpboot/pxelinux.0";
    subnet 172.20.0.0 netmask 255.255.0.0 {
        max-lease-time      43200;
        default-lease-time  43200;
        option routers      172.20.0.1;
        option subnet-mask  255.255.0.0;
        option nis-domain   "cluster.com";
        option domain-name  "cluster.com";
        option domain-name-servers 172.20.0.1;
        option time-offset  -7;
        range                172.20.200.1 172.20.255.254;
    }
} #172.20.0.0/255.255.0.0 subnet_end#

subnet 172.29.0.0 netmask 255.255.0.0 {
    max-lease-time      43200;
    default-lease-time  43200;
    option routers      172.29.0.1;
```

```

        option subnet-mask          255.255.0.0;
        option nis-domain           "cluster.com";
        option domain-name         "cluster.com";
        option domain-name-servers 172.29.0.1;
        option time-offset         -7;
    } #172.29.0.0/255.255.0.0 subnet_end#
} #eth0 network_end#

shared-network eth1 {
    subnet 172.30.0.0 netmask 255.255.0.0 {
        max-lease-time             43200;
        default-lease-time         43200;
        option routers              172.30.0.1;
        option subnet-mask         255.255.0.0;
        option nis-domain          "cluster.com";
        option domain-name         "cluster.com";
        option domain-name-servers 172.30.0.1;
        option time-offset         -7;
    } #172.30.0.0/255.255.0.0 subnet_end#
} #eth1 network_end#

#shared-network all {
#} #all network_end#

```

Notes:

1. After using *getmacs* and then *makedhcp --allmacs*, an entry for each MAC address for each node will be listed in the *dhcp.conf*.
2. Usually DHCP should not run on the network interface that is connected to the rest of the network. In this case, remove the network section from *dhcpd.conf* that corresponds to the external network and then explicitly list the interfaces DHCPD should listen for in */etc/dhcpd.conf*.

Edit */etc/sysconfig/dhcpd*, with something similar to:

```
DHCPDARGS="eth0 eth1"
```

Notes:

1. The *dhcpver* field in *\$XCATROOT/etc/site.tab* must be set to match the version of *dhcpd* installed. Generally 2 for older Red Hat and 3 for SUSE Linux and newer Red Hat before you run *makedhcp*. If incorrect, correct and rerun *makedhcp --new --allmac*.
2. *\$XCATROOT/etc/networks.tab* must define each network that *dhcpd* is to support. Let *makedhcp* build it for you the first time, edit and rerun *makedhcp --new --allmac*.

Configure all Ethernet switches, but block DHCP from in and out bound ports that are used to connect the cluster to a production environment. Please read the “[xCAT 1.1.0 Redbook](#)”, the “[cisco2950-HOWTO](#)”, and the “[force10-HOWTO](#)” found in */opt/xcat/doc* for more information.

Configure all terminal servers. Please read the “[terminalserver-HOWTO](#)”.

Restart *conserver* if you are using terminal servers or SOL. If you are using an IBM BladeCenter without SOL do not use *conserver*.

```
>service conserver restart
```

Setup stage boot image.

For x86 and x86_64 type:

```
>cd /opt/xcat/stage  
>./mkstage
```

For ia64 type:

```
>cd /opt/xcat/stage  
>./mkstage-ia64
```

Collect the MAC addresses of the compute nodes and create entries in *dhcpd.conf* for them.

Prepare to monitor stage2 progress

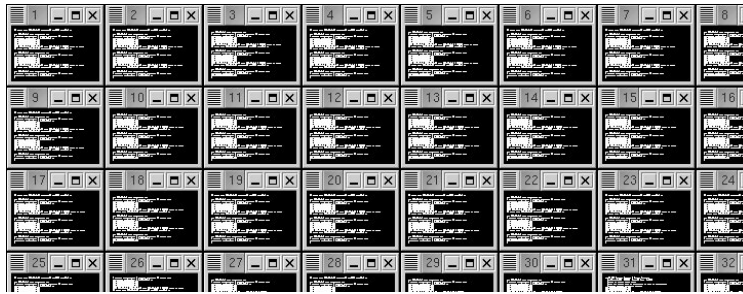
```
> wcons -t 8 compute (or a subset like rack01)  
> tail -f /var/log/messages
```

Be aware of system messages, it is a very good way to stay informed about the cluster.

Manually reboot the compute nodes.

During the boot process, the machines should PXE boot syslinux, obtain a dynamic IP address, and then load a Linux kernel and a special RAM disk that contains a script to print the machine's MAC address to the console.

Observe the output of the *wcons* windows. If the terminal servers are working correctly, the machines boot their kernels and display an image similar to the one below:



A closeup:

```
3
--- 172.30.0.1 ping statistics ---
1 packets transmitted, 1 packets received, 0% packet loss
round-trip min/avg/max = 0,2/0,2/0,2 ms
MAC-00:02:55:C6:C7:A8-MAC
MAC-00:02:55:C6:C7:A8-MAC
MAC-00:02:55:C6:C7:A8-MAC
PING 172.30.0.1 (172.30.0.1) from 172.30.1.3 : 56(84) bytes of data,
64 bytes from 172.30.0.1: icmp_seq=0 ttl=255 time=0,2 ms

--- 172.30.0.1 ping statistics ---
1 packets transmitted, 1 packets received, 0% packet loss
round-trip min/avg/max = 0,2/0,2/0,2 ms
MAC-00:02:55:C6:C7:A8-MAC
MAC-00:02:55:C6:C7:A8-MAC
MAC-00:02:55:C6:C7:A8-MAC
PING 172.30.0.1 (172.30.0.1) from 172.30.1.3 : 56(84) bytes of data,
64 bytes from 172.30.0.1: icmp_seq=0 ttl=255 time=0,2 ms

--- 172.30.0.1 ping statistics ---
1 packets transmitted, 1 packets received, 0% packet loss
round-trip min/avg/max = 0,2/0,2/0,2 ms
MAC-00:02:55:C6:C7:A8-MAC
```

The wcons windows are xterms. When viewing a large number of consoles on the screen at the same time, the xterms come up with the “unreadable” font size. Xterms have a feature that allows a user to change the size of the font very easily. This allows the user to enlarge a specific view when a screen of unreadable consoles is displayed. To do this, move the mouse over the text portion of the xterm in question, press and hold **Ctrl+Right click** the mouse. The following menu will be displayed:



Move the mouse down to select a larger font and then release the mouse button as shown:



Using this xterm feature, you can switch to a large font for detailed viewing and back to the smaller font to view all the consoles at once.

Press **Ctrl+E** to access additional terminal functionality.

Use the following command to collect the MACs once all the compute nodes are displaying their MAC addresses out of their serial consoles.

```
> getmacs compute
```

Note: This *.tab* file should be in the configuration files that came with the IBM 1350 cluster.

Use the following command to close the wcons windows.

```
> wkill
```

Manually reboot each node and use the following command to collect MAC addresses:

```
>getmacs <noderange>  
or  
>getmacs compute
```

```
node1-eth0 00:07:E9:93:F8:DD  
node1-eth1 00:00:5A:9A:DB:7C  
node2-eth0 00:07:E9:93:F8:DD  
node2-eth1 00:00:5A:9A:DB:7C
```

When the message “*Auto merge mac.lst with /opt/xcat/etc/mac.tab(y/n)?*” appears, type **Y**.

Each node will be suffixed with the interface of the collected MAC. Please do not alter.

Notes

1. Do not alter the *mac.tab* entries for collected MACs. It is critical that the stored node names remain untouched. If necessary, you may change the MAC.
2. Multiple *getmacs* commands will corrupt *mac.tab*. Only run one instance at a time.
3. Some operating systems report eth0 and eth1 differently than xCAT *getmac* reports. The settings may need to be reversed manually in *mac.tab*, however this may have a negative impact on other non-switched entries. Verify your settings making corrections.

```
perl -pi -e 's/(nodeprefix.*)-eth0/$1-ethfoo/' mac.tab  
perl -pi -e 's/(nodeprefix.*)-eth1/$1-eth0/' mac.tab  
perl -pi -e 's/(nodeprefix.*)-ethfoo/$1-eth1/' mac.tab
```

Note: Currently only the serial-based (rcons) method of connecting MACs will collect multiple MAC/node. A future version of xCAT will address this limitation.

Exception: IBM BladeCenter mpcli2 and bcmm getmacs methods can collect both MAC addresses.

Note: For IBM BladeCenter please use bcmm method in *nodehm.tab*.

MAC addresses may be collected without a terminal server.
Configure *cisco3500.tab* with an example of the following:

```
node01 Ethernet01,1
node02 Ethernet01,2
node03 Ethernet01,3
node04 Ethernet01,4
```

Make *nodehm.tab* have entries like:

```
nodexx mp,mp,mp,mp,mp,conserver,mp,mp,rcons,cisco3500,bcm5700,vnc
```

Make sure the switch has a hostname and DNS resolves.

Verify that the nodes plugged into the switch ports match those in *cisco3500.tab*.

Example: *node1 port1 node2 port2*

Make sure you can ping the switch, telnet to it and login. Make sure the password you set on the switch is the same in *passwd.tab*. Put the nodes in stage2. Power them on and *getmacs* as usual. The *getmacs* command issues the show mac-address-table on the switch and grabs the MACs from it.

For other switches, *switch.tab*, *getmacs.switch.snmp*, and *getmacs.switch* are required. Place *getmacs.switch.snmp* and *getmacs.switch* into the *opt/xcat/lib* directory and make sure they are executable by using *ls -l* to verify.

Place *switch.tab* into the *opt/xcat/etc* directory:

For example SMC alters the *switch.tab* as follows: (see examples in *switch.tab* for SMC and other switches. This will be the future way of setting up switches.)

```
nodexxx      smc8648-001,18,NA
|
|              |
|              |      smc port number
|              |
|              |      smc name-switch number (as named in your other tab files & hosts)
Node
```

Edit the *nodehm.tab*.

Here is an example of one that is set up for using RCONS as a method for *getmacs* (not necessarily the way yours will look but just an example of how the *nodehm.tab* file may appear):

```
node1 mp,mp,mp,mp,mp,conserver,mp,mp,rcons,pxe,eepr100,vnc,Y,NA,NA,def
```

Here is an example of using new *getmacs*:

Edit the appropriate entry to point to the switch scripts (this will be what tells *getmacs* to use *getmacs.switch* script).

```
node1 mp,mp,mp,mp,mp,conserver,mp,mp,switch,pxe,eepr100,vnc,Y,NA,NA,def
```

Build */etc/dhcpd.conf* with MAC entries:

```
makedhcp -allmac
```

For all IBM xSeries nodes with IBM management processors and the IBM eServer 325 and 326, excluding IBM BladeCenter, read “[managementprocessor-HOWTO](#)” found in */opt/xcat/doc* for more information. For IBM BladeCenter, use *mpname noderange*.

```
nodeset noderange stage3
```

Reboot each node manually after all MACs collected and DHCP server restarted.

Read the “[managementprocessor-HOWTO](#)” and “[IBM BladeCenter-NOTES](#)” for information on testing and troubleshooting all nodes management processors if applicable.

Test systems management:

```
rpower noderange stat
```

```
rbeacon noderange on
```

Note: Not all servers have a blinking light.

Copy the Red Hat Install CD(s) by inserting the CDs and then run *copycds* and follow the prompts.

Example: *copycds <namecd1>.iso, <namecd2>.iso, <namecd3>.iso*

Notes:

1. When the CDs are entered and a prompt for “auto run” appears, select **No**.
2. You may also use *copycds* to copy the contents of one or more *.iso* files.

Copy the “post” files for Red Hat.

Copy install files from the xCAT distribution to the post directory that is used during unattended installs:

```
>cp /opt/xcat/samples/etc/post* /opt/xcat/etc
>cp /opt/xcat/install/rhas4/x86_64/base/compute.tmpl ..
```

Enter the following commands to enable remote logging:

```
> cp /opt/xcat/samples/syslog.conf/etc
> touch /var/log/pipemessages
> service syslog restart
```

Setup snmptrapd

snmptrapd received messages from the SPN.

```
> chkconfig snmptrapd on
> service snmptrapd start
```

The following command creates a SSH keypair for root with an empty passphrase which sets up root's SSH configuration, copying *keypair* and *config* to */install/post/.ssh* so that all installed nodes will have the same root *keypair/config*. This allows you to install and log into nodes.

```
>gensshkeys root
```

Setup NFS and NFS Exports by making */etc/exports* look similar to the following:

```
/install node*(ro,sync,no_root_squash)
/tftpboot node*(ro,sync,no_root_squash)
/usr/local node*(ro,sync,no_root_squash)
/opt/xcat node*(ro,sync,no_root_squash)
/home node*(rw,sync,no_root_squash)
```

Turn on NFS by using the following commands:

```
> chkconfig nfs on
> service nfs start
> exportfs -ar # (to source)
> exportfs # (to verify)
>echo "/install *(ro,async,no_root_squash)" >>/etc/exports
>service nfs restart
```

Notes:

1. If you do not have a Myrinet read the “myrinet-how to” document in */opt/xcat/doc*. For more detailed information, read the “[nodeinstall-HOWTO](#)” and “[systemimager-HOWTO](#)” for details on node install and diskless installs.

2. If installing a node from disk, use *rinstall* or *winstall*. Only install 32 nodes at a time or use staging. For more information, read the “man” pages on *rinstall* and *winstall*.

14. Installing compute nodes

This section covers the installation of the compute nodes.

Modify the “kickstart” template file if needed and verify the correct version of RedHat is listed.

```
>cd /opt/xcat/install/rhws4/<architecture>/base/  
>cp /opt/xcat/install/rhws4/<architecture>/base/compute.tmpl ..
```

The following command makes the nodes PXE boot the RedHat “kickstart” image by altering the files in */tftpboot/pxelinux.cfg/*.

```
> nodeset compute install
```

Prepare to Monitor the Installation Progress

```
> wcons -t 8 compute
```

Note: A subset like rack01 may be substituted.

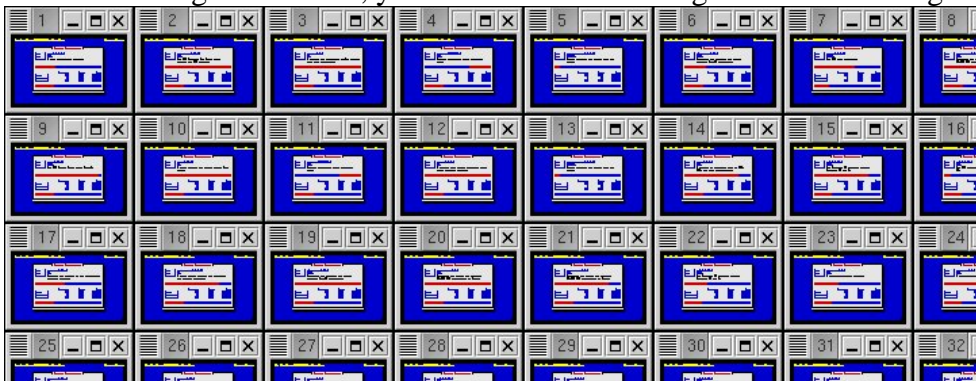
```
> tail -f /var/log/messages
```

Note: Be aware of any warning messages that may appear.

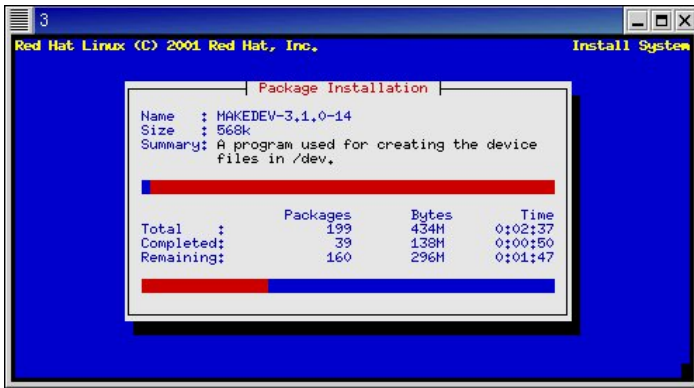
Reboot the Compute Nodes.

```
>rpower compute boot
```

When installing with *wcons*, you should see something like the following:



Close up:



Note: To install without terminal servers, Serial Over LAN must be configured.

15. Serial Over LAN (SOL) setup

Download and install “SMbridge RPM” from <http://www.ibm.com/support/docview.wss?uid=psg1MIGR-57729>.

This RPM needs to be installed on a management node (or the “crash cart” that has xCAT on it) since it is client software for the BMCs.

For IBM xSeries 336, 346, and eServer 326

Flash the Management Processor (BMC) to the latest version.

Flash BIOS to the latest version.

Remove the power cord for 10 seconds.

Restore the power cord.

Reboot and press **F1** to enter the BIOS configuration.

Configure the BIOS settings for optimum performance and then edit the Devices and I/O Ports as shown below:

Devices and I/O Ports

Serial port A: **Port 3F8, IRQ 4**

Serial port B: **Disabled**

Remote Console **Redirection**

- Remote Console Active: **Enabled**
- Remote Console COM Port: **COM 1**
- Remote Console Baud Rate: **19200**

To use Remote Console Text Emulation: VT100/VT220

Configure the text emulation settings as listed below:

Remote Console Keyboard Emulation: **VT100/VT220**

Remote Console After Boot: **Enabled**

Remote Console Flow Control: **Hardware**

Configure the “Startup” settings as listed below:

From the main menu, select **Startup Sequence** and configure the settings as follows:

First Startup Device: **CD ROM**

Second Startup Device: **Diskette Drive 0**

Third Startup Device: **Network**

Fourth Startup Device: **Hard Disk**

Wake On LAN: **Disabled**

Planer Ethernet PXE/DHCP: **Planer Ethernet 1**

Boot Fail Count: **Disabled**

Go to **Advanced Setup** then **CPU Options** and configure the settings as follows:

[Hyper-Threading Technology](#): Disabled

For IBM xSeries 326

Configure the optimum BIOS settings for the IBM xSeries 326 and include the following settings.

From the main menu, select **Console Redirection** and configure the settings as follows:

Console Redirection: **COM A**
Baud rate: **19.2 K**
FIFO Level: **14**
Console Type: **vt100**
Flow Control: **CTS/RTS**
Console Connection: **Direct**
Continue CR After Post: **On**

From the main menu, select **BMC** and configure the settings as follows:

IPMI Spec Version: **1.5**
BMC Firmware Version: **1.11**
Com port on BMC: **CLI**
Change Com port Setting: **No**
Clean System Eventlog: **Disabled**
System Firmware Progress: **Enable**
BIOS Post Watchdog: **Enable**

Additional settings for xSeries 336, 346, and eServer 326

From the main menu, select **Advanced Setup** then **Baseboard Management Controller (BMC) Settings** and configure the settings as follows:

System BMC Serial Port Sharing: **Enabled**
BMC Serial Port Access Mode: **Dedicated**

Save the settings.

Note: If switching to SOL, you must remove the power cord for 5 sec.

Tabs

conserver.cf

Conserver.cf will have to be altered to point to the SOL script for the specific node.

Example:

```
node001:/sol.eServer 326 node001::&:
node002:/sol.xSeries 336 node002::&:
node003:/sol.xSeries 346 node003::&:
ipmi.tab
```

There are a few different ways of approaching this

```
node123      bmc123, "", ""
```

```
node123      bmc123,
node123      bmc123,USERID,PASSWORD
```

Note: There is a zero in “PASSWORD”.

If you use quotations (“”) then you will have to enter “” for the userid and password when you start your *wcons* session.

If you leave the field blank then the userid and password should default to the definitions in the *passwd.tab* file.

We have also used the third example and placed the default userid and password (*USERID,PASSWORD*). Whatever you put in there will over ride the defaults and that is what you will have to enter on your *wcons* window for the node you intend to view.

nodehm.tab

Set up the *nodehm.tab* file to point to the *ipmi* tool (uses BMC).

Example:

```
node001
ipmi,ipmi,ipmi,ipmi,ipmi,conserver,NA,ipmi,switch,pxe,bcm5700,vnc,Y,ipmi,NA,19200
```

Note: The *ipmi* parameter in several of the fields. In this example we have also setup the baud rate for 19200. It has to match what is set in the **BIOS Setup** under **Remote Console settings**.

site.tab

RHEL 3.0 and below may cause a problem when using *wcons* to view the node console. The problem is that the console title will not show the node name and therefore confusion as to which node you are viewing may occur. To correct this you must turn off *bufferedcons* in the *site.tab* file, then the node name will display correctly in the title bar.

Example:

```
Bufferedcons no
```

WCONS

When you run *wcons <nodename>* , the screen will display:

connected....

login :

password :

Entry for login and password has to be the same as what is configured in the *ipmi.tab* file.

Example:

login : "" and *password : ""* is used as the user ID and password in the previous example of the *ipmi.tab*.

Note: You must have “smbridge RPM” installed.

Verify that the compute nodes installed correctly

>pping all

Update the SSH global known hosts file

> makesshghk compute

16. Clean up

This section covers the final installation steps and test information..

Copy the xCAT initialization files. This will enable some services to start at boot time and change the behavior of some existing services.

```
> cd /opt/xcat/rc.d
> cp atftpd portmap snmptrapd syslog /etc/rc.d/init.d/
```

There are other initialization files in */opt/xcat/rc.d* that may also be used, depending on the installation.

Move unneeded *.tab* files from */opt/xcat/etc/* to a temporary directory.

Test the cluster. Read the man pages for *rvitals*, *rinv*, and *rpower* and then try some of these commands on the cluster.

```
> psh compute date | sort
```

The output here will be a good way to see if SSH/gkh is setup correctly on all of the compute nodes, which is a requirement for most cluster tasks. If a node does not appear here correctly, you must go back and troubleshoot the individual node. Make certain the install process completed correctly by using *makesshgkh* and then test again with *psh*. The *psh* test should pass before continuing.

Additional test commands:

```
> rvitals compute ambtemp
> mpncheck compute
> pping all
> rbeacon ccompute on
```

17. Contributing to xCAT

Join the “xCAT-dev” mailing list and post your suggestions, bug-fixes and code by visiting <http://xcat.org/mailman/listinfo/xcat-user>.

18. Credits

This document was most recently modified:

03/01/2006

Original author Matt Bohnsack

Send additions and corrections to the editor jmweb@us.ibm.com, so this document can continue to be improved.

Thanks go out to the following people, who helped this document become what it is today:

Egan Ford for writing xCAT, Jarrod B Johnson, Mike Galicki, Andrew Wray, Chris DeYoung, Mark Atkinson, Greg Kettmann, Jay Urbanski, The people from POSDATA, Kevin Rudd, Tom Alandt, and Tonko L De Rooy for there continuing support and dedication to the development of xCAT,

19. Supporting documentation

Additional documentation may be located in */opt/xcat/doc*.

License

xCAT Support

xCAT Redbooks

xCAT Man Pages

OSS Licenses (Incomplete, WIP)

HOWTOs:

 xCAT Mini HOWTO (1.2.0) (Start Here)

 xCAT HOWTO (1.1.0) (Reference Only)

Hardware HOWTOs:

 Blade Center NOTES (1.1.7.2 and 1.2.0)

 Management Processor HOWTO (1.2.0)

 Stage1 HOWTO (1.2.0)

Switch/Terminal Server HOWTOs:

 Cisco 2950 HOWTO (1.2.0)

 Force 10 HOWTO (1.2.0)

 Myrinet-HOWTO (1.2.0)

 Terminal Server HOWTO (1.2.0)

Management Node HOWTOs:

 SUSE Linux Management Node HOWTO (1.2.0)

Node Install HOWTOs:

 Node Installation HOWTO (1.2.0)

 Imaging HOWTO (1.2.0)

 SystemImager HOWTO (1.2.0)

 Remote Flash HOWTO (1.2.0)

 Windows HOWTO (1.1.0)

 Diskless HOWTO (1.2.0)

 Warewulf HOWTO (1.2.0)

Software HOWTOs:

 HPC Benchmark HOWTO (1.2.0)

 GPFS HOWTO (1.1.0)

For more information, visit <http://www.xcat.org>.

