

xCAT 2.0 Linux Beta Release Cookbook

04/12/2008

Table of Contents

1.0 Release Description	3
1.1 Function supported:	3
1.2 Function not supported:	4
1.3 Prerequisites:	4
1.4 Licensing	4
2.0 Installing xCAT 2.0 Software	5
2.1 Download OS and Create Repository	5
2.2 Downloading and Installing xCAT 2.0	6
2.2.1 If Your Management Node Has Internet Access:	6
2.2.1.1 Download Repo Files	6
2.2.2 If Your Management Node Does Not Have Internet Access:	6
2.2.2.1 Download xCAT2.0 and Its Dependencies:	6
2.2.2.2 SetupYUM repositories for xCAT and Dependencies	6
2.2.3 Install xCAT 2.0 software & Its Dependencies	7
2.2.4 Test xCAT installation	7
2.2.5 Update xCAT 2.0 software	7
2.2.6 Setup Yum for OS installs	7
3.0 xCAT 2.0 Commands	8
3.1 xCAT Tables	10
3.2 Table edit commands	11
3.3 Using the node* commands	12
3.3.1 Using xCAT object definition commands	13
3.4 Using xCAT hardware commands	15
3.4.1 Hardware discovery	15
3.4.2 Hardware Control	17
4.0 Setup the Management Node	17
4.1 Using Service nodes	17
4.1.1 Switching to PostgreSQL Database	18
4.2 Define the service nodes in the database	21
4.2.1 Define Service Nodes and bmc in nodelist table	21
4.2.2 Define Service Nodes in noderes table	21
4.2.3 Define Service Nodes in ipmi table	21
4.2.4 Define Service Nodes and bmc in nodehm table	21
4.2.5 Define Service Nodes and bmc in nodetype table	22
4.2.6 Define Service Nodes in site table	22
4.2.7 Define Service Node OS and Profile attributes	22
5.0 Setup Services	22

5.1 Setup networks Table	22
5.2 Setup DNS	22
5.3 Setup AMM	23
5.4 Setup DHCP	24
5.5 Startup TFTP	25
6.0 Define Compute Nodes	25
6.1 Setup the nodelist Table	26
6.2 Setup the noderes Table	26
6.2.1 Sample noderes table	27
6.2.2 Setting up which services run on the Service Nodes	27
6.3 Setup nodetype table	28
6.3.1 Sample nodetype table	28
6.4 Setup passwords in passwd table	28
6.5 Setup deps Table (need?)	28
6.5.1 Sample deps table	29
6.6 Get MAC addresses for all nodes	29
6.7 Setup Conserver	29
6.8 Makehosts	29
7.0 Installing xCAT nodes	30
7.1 Installing xCAT Service Nodes (diskfull)	30
7.2 Installing the Compute Nodes (diskfull)	31
7.2.1 Check the site table attribute	31
7.2.2 Discover Nodes	31
7.2.3 Install Nodes	32
7.3 Install the Service Nodes (stateless)	33
7.3.1 Build the service node stateless image	33
7.3.2 Install the Service Nodes	34
7.3.3 Test Service Node installation	34
7.4 Installing x86-64 Compute Nodes (stateless)	34
7.4.1 Building the stateless image	34
7.4.2 Install the stateless image	36
7.4.3 Update Stateless image	36
7.5 Installing PPC64 Compute Nodes (Stateless)	37
7.5.1 iSCSI install	37
7.5.2 Build PPC64 Stateless image	37
7.5.3 Install PPC64 Stateless image	38
7.5.4 Update PPC64 Stateless image	39
8.0 Using xCAT Notification Infrastructure	40
8.1 Using xCAT Monitoring Plug-in Infrastructure	42
9.0 xCAT Architecture	44
9.1 Client/Server	45
9.2 Flow	45
10.0 References	46

1.0 Release Description

xCAT 2.0 is a complete rewrite of xCAT 1.2/1.3 implementing a new architecture (see description at end of this document). All commands are client/server, authenticated, logged and policy driven. The clients can be run on any OS with Perl, including Windows. The code has been completely rewritten in Perl, and table data is now stored in a relational database. As the default database, we are including SQLite with the xCAT OSS rpm. To use the new Service Node feature, you must be using Redhat 5 or Fedora 8 and install and setup the PostgreSQL Database. See instructions below in Chapter “xCAT Hierarchy using Service Nodes chapter”.

The code is being released as RPMs and SRPMs . For the beta release, there is support for x86_64 hardware (IPMI and Blades) and ppc64 hardware (js and qs blades) . The OS must be RedHat 5, CentOS5, Fedora 8 or SLES 10.

The beta code should not be used for production work .

There is an alpha release for AIX available: See http://sourceforge.net/docman/display_doc.php?docid=86730&group_id=208749

1.1 Function supported:

- Tools to manipulate the database tables: tabdump,tabrestore,tabedit, chtab , nodefs, nodech, nodeadd, noderm, chdef, mkdef, lsdef, rmdef,
- Cluster setup commands: makehosts, makedhcp ,makeconservercf
- Notification commands (infrastructure allowing users to register for xCAT database table changes): regnotif, unregnotif
- Monitoring commands (monitoring plug-in infrastructure allowing plug-in third party monitoring software to the xCAT cluster): monstart, monstop, monls.
- Hardware control commands : lsslp, rscan, rpower, reventlog, rinvs, getmacs, rvitals
- Install commands : rnetboot
- Parallel remote and remote copy commands : xdsh, xdcp, xdshbak, psh. xdsh/ xdcp is now packaged with xCAT.
- Node discovery and diskfull and diskless deployment of CentOS5 and RHEL5, Fedora 8 on the supported hardware (see Prerequisites).

- Hierarchical install (diskfull/diskless) using Service Nodes on Redhat 5 or Fedora 8.
- For a list of all 2.0 xCAT commands run *rpm -ql xcat-client*.
- Manpages
- noderanges -ability to enter ranges of nodes on commands. See man noderanges for syntax
- nodegroups – ability to define groups of nodes to be used with the commands
- Diskless/Stateless install
- Data abstraction commands to make creating node and other database definitions easier. See Chapter 3.3.1, “Using xCAT Object Definition Commands”.

1.2 Function not supported

- No imaging
- No flash
- pSeries hardware control using HMC, IVM, FSP for Power5 and Power6 hardware
- Web GUI interface

1.3 Prerequisites:

- Hardware requirements:
 - x3455, x3550, x3650, x3455, LS21, HS21, LS41, x336, x346, ppc64
 - no SOL for x386 or x486
 - Must be IPMI based, rack mounted unit.
 - Blades
 - Ethernet switch must be SNMP enabled for node discovery.
- Software supported
 - RedHat5, CentOS5, Fedora 8, SLES 10

1.4 Licensing

xCAT 2.0 is OSS with a EPL license. For license information visit

<http://www.opensource.org/licenses/eclipse-1.0.php>

2.0 Installing xCAT 2.0 Software

Install your xCAT management node with RedHat5, CentOS5, Fedora 8 or SLES 10 making sure to install **all** packages available with the distribution to reduce the number of dependency RPMs you need to track down.

Note: If you installed the xCAT2.0 alpha code, you should remove it from the system. Ensure your networks are setup correctly.

2.1 Download OS and Create Repository

1. Get OS ISOs and place in a directory, for example /root/xcat2 and fedora8:

```
mkdir /root/xcat2
cd /root/xcat2
wget
ftp://download.fedora.redhat.com/pub/fedora/linux/releases/8/Fedora/x86\_64/iso/Fedora-8-x86\_64-DVD.iso
wget
ftp://download.fedora.redhat.com/pub/fedora/linux/releases/8/Fedora/ppc/iso/Fedora-8-ppc-DVD.iso
```

2. Create YUM repository for OS RPMs:

```
mkdir /root/xcat2/fedora8
mount -r -o loop /root/xcat2/Fedora-8-x86_64-DVD.iso /root/xcat2/fedora8

cd /etc/yum.repos.d
mkdir ORIG
mv fedora*.repo ORIG
```

Create fedora.repo with contents:

```
[fedora]
name=Fedora $releasever - $basearch
baseurl=file:///root/xcat2/fedora8
enabled=1
gpgcheck=0
```

3. Install createrepo:

```
yum install createrepo
```

2.2 Downloading and Installing xCAT 2.0

2.2.1 If Your Management Node Has Internet Access:

2.2.1.1 Download Repo Files

YUM can be pointed directly to the xCAT download site.

```
cd /etc/yum.repos.d
wget http://xcat.sf.net/yum/core-snap/xCAT-core-snap.repo
wget http://xcat.sf.net/yum/dep-snap/rh5/x86\_64/xCAT-dep-snap.repo
```

2.2.2 If Your Management Node Does Not Have Internet Access:

2.2.2.1 Download xCAT2.0 and Its Dependencies:

Note: do the wget's on a machine with internet access and copy the files to this machine.

```
cd /root/xcat2
wget http://xcat.sf.net/yum/core-rpms-snap.tar.bz2
wget http://xcat.sf.net/yum/dep-rpms-snap.tar.bz2
tar jxvf core-rpms-snap.tar.bz2
tar jxvf dep-rpms-snap.tar.bz2
```

2.2.2.2 Setup YUM repositories for xCAT and Dependencies

```
cd /root/xcat2/dep-snap/rh5/x86_64
./mklocalrepo.sh
cd /root/xcat2/core-snap
./mklocalrepo.sh
```

2.2.3 Install xCAT 2.0 software & Its Dependencies

```
yum clean metadata
yum install xCAT.x86_64
```

2.2.4 Test xCAT installation

```
source /etc/profile.d/xcat.sh
tabdump site
```

2.2.5 Update xCAT 2.0 software

If you need to update the xCAT 2.0 rpms later, download the new version of <http://xcat.sf.net/yum/core-rpms-snap.tar.bz2> (if the management node does not have access to the internet) and then run:

```
yum update xCAT.x86_64
```

If you have a service node stateless image, don't forget to update the image with the new xCAT rpms (see Install the Service Nodes (stateless)):

```
cp -pf /etc/yum.repos.d/*.repo /install/netboot/fedora9/x86_64/service/
  rootimg/etc/yum.repos.d
yum --installroot=/install/netboot/fedora8/x86_64/service/rootimg update
```

2.2.6 Setup Yum for OS installs

```
umount /root/xcat2/fedora8
chtab key=installdir site.value=/install
cd /root/xcat2
copycds Fedora-8-x86_64-DVD.iso
copycds Fedora-8-ppc-DVD.iso
```

Edit /etc/yum.repos.d/fedora.repo and change:

```
baseurl=file:///tmp/fedora8
to
baseurl=file:///install/fedora8/x86_64
```

3.0 xCAT 2.0 Commands

Note: use '`<xCAT command> -h`' for a usage message from each command. MAN pages are not available at this time.

xCAT Command	Description
chtab (Note: will be renamed to tabch in beta)	To add or update rows in a table. Allows you to add nodes, create groups, add attributes to the xCAT tables. <code>chtab node=devnode01 nodelist.group=all,compute</code> will add a new node devnode01 to the nodelist table and assign to the all and compute groups. <code>chtab key=rsh site.value=/usr/bin/ssh</code> will assign the site table rsh attribute to /usr/bin/ssh <code>chtab -d node=devnode01</code> will delete the previously create node from the nodelist table.
copycds	Copies Linux distributions and service levels to install directories.
chdef	Change xCAT data object definitions.
chvm	Changes HMC- and IVM-managed partition profiles. (not available for use)
dumpxCATdb	Dumps the xCAT database to the directory input
genimage	Generates the stateless image for the selected OS and node profile
getmacs	Collect node MAX addresses
lsdef	Use this command to list xCAT data object definitions.
lssl	Queries selected networked services configuration information.
lsvm	Lists partition profile information for HMC- and IVM-managed nodes. (not available for use)
makeconservercf	Make Conserver Configuration
mkdef	Use this command to create xCAT data object definitions.
makedhcp	Sets up the DHCP server.
makehosts	Creates entries in /etc/hosts for nodes. Node nodenames and ip addresses must be setup in the hosts table.
makenetworks	Builds the networks table
mkvm	Creates HMC- and IVM-managed partitions. (not available for use yet)
monls	Lists monitoring plug-in module names, status and description.
monstart	starts a monitoring plug-in module to monitor the xCAT cluster.

monstop	stops a monitoring plug-in module to monitoring the xCAT cluster
nodeadd	<p>Add a node to the cluster</p> <p>For example: <code>nodeadd <noderange> [table.column=value] [table.column=value]....</code></p> <p style="text-align: center;"><code>nodeadd blade1-blade7 nodelist.groups=all,compute</code></p> <p>nodeadd also supports some short cut tags:</p> <ul style="list-style-type: none"> groups is equivalent to <code>table.column = nodelist.groups</code> <ul style="list-style-type: none"> • <code>nodeadd blade1-blade8 groups=all,compute</code> mgt is equivalent to <code>table.column = nodehm.mgt</code> <ul style="list-style-type: none"> • <code>nodeadd blade7 mgt=blade</code> switch is equivalent to <code>table.colum= switch.switch</code> <ul style="list-style-type: none"> • <code>nodeadd blade8 switch=switch1</code>
nodech	Change node information
nodels	Display information about a node or range of nodes or all nodes
noderm	Remove Node
nodeset	Installs, boots the nodes uses pxe.
packimage	Packs the stateless image in the correct os and arch directory before the stateless install
psh	Runs a command across a list of nodes or nodegroups in parallel
rbeacon	Turns beacon on/off/blink or gives status of a node or a range of nodes.
rbootseq	For boot of Bladecenter node range. Change each node boot order.
regnotif	Register a Perl module or a command that will get called when changes occur in desired xCAT database tables. See Using xCAT Notification
restorexCATdb	Restores the xCAT database from the input directory
reventlog	Retrieves or clears remote hardware event logs
rinv	Retrieves hardware configuration information for a single or range of nodes
rmdef	Use this command to remove xCAT data object definitions.
rmvm	Removes HMC- and IVM-managed partitions. (not supported yet)
rnetboot	Will force an unattended network install for a range of nodes (diskless) .
rpower	Boots, resets, powers on and off and queries nodes
	Note: “boot” option not implemented yet. Use either “on” or “reset” options as appropriate.

rscan	Collects node information from hardware control point
rsetboot	rsetboot (IPMI) is a way to specify the singular device to try to boot only for the next power cycle
rsreset	Used to reset service processors out-of-band
rvitals	Retrieves hardware vital information from the on-board Service Processor for a range of nodes
tabdump	Display Database table information for table requested. tabdump with no input will display a list of all valid table names. tabdump -d <tablename> will list the fields of the table and their definitions
tabedit	Edit a table . Must export EDITOR to define your editor.
tabrestore	Restore a table from the table.csv template or from a tabdump output file.
unregnotif	Unregistered a Perl module or a command that was watching for changes of desired xCAT database tables.
updateSNimage	Add required xcat files to the Service Node stateless image before installing in a stateless environment
xdcp	Concurrently copies files to/from multiple nodes. See xdsh and xdcp man page for more information
xdsh	Concurrently runs remote commands on multiple nodes. All dsh code is now shipped with xCAT 2.0. If dsh rpms were obtained from the following website and installed for the alpha release, you should erase the csm.dsh* rpm. http://www14.software.ibm.com/webapp/set2/sas/f/csm/download/home.html . See xdsh man page for more information.
xdshbak	Presents formatted output from the xdsh command

3.1 xCAT Tables

Note: The Database Table Schema can be viewed in the /usr/lib/perl5/site_perl/5.8.3/xCAT/Schema.pm file or by running the tabdump command. tabdump -d <table> will give you definitions of the table attributes.

Table Name	Description
chain	Lists action that occur during node install, node boot . Used by nodeset.
deps	Dependency node table
hosts	List of hosts, alias hostname, ip addresses. Used to update /etc/hosts with makehosts

ipmi	Lists information on the nodes IPMI interface – bmc, username, password
iscsi	List information for setting up iSCSI
mac	Lists mac address for each node.
monitoring	Lists the monitoring plug-in module names that are monitoring the xCAT cluster.
monsetting	List the monitoring plug-in specific settings.
mp	This is the management processor network. Whereas the mpa.tab lists the adapter, this table lists devices that are networked off that adapter via daisy chained networks, or in the case of Blade Center, an internal network..
mpa	Lists the MPA, username and password for the nodes.
networks	Defines masks, gateways and DNS servers. Build my makenetworks command.
nodegroup	Lists information on all nodegroups defined
nodehm	Defines the hardware management method for each node.
odelist	Defines all nodes and groups.
nodepos	Node physical location
nodes	Installation resources for the node.
nodetype	Node install type (osversion, arch, type)
notification	Lists the Perl modules and commands that will get called for changes in certain xCAT database tables.
osimage	Contains information that describes a unique operating system image that may be deployed on a cluster node
passwd	user names and passwords used by xCAT scripts
policy	Table controls the policy for the execution of the xcat commands.
postscripts	Comma separated list of scripts that should be run on this node after installation or diskless boot. (TBD)
ppc	Store Series p hardware components – HMC, IVM, BPA, FSP, LPAR
ppcdirect	Contains direct-attached FSP hardware information
ppchcp	Contains HMC and IVM hardware information
site	Main xCat configuration file. Holds global information for the cluster.
switch	Lists switch interface(s) for the node.
vpd	Vital product data table. Holds machine serial number and model type.
xCATWorld	Sample xCat client – plugin program.

3.2 Table edit commands

To manage these tables directly, xCAT provides the **chtab**, **tabdump**, **tabrestore**, and **tabedit** commands.

The following are some basic examples of how to use the database table commands.

1. To see what tables exist in the xCAT database:

```
tabdump
```

2. To display the definition of the attributes of the nodelist table:

```
tabdump -d nodelist
```

3. To display the contents of the site table

```
tabdump site
```

4. To back up all the xCAT tables.

```
mkdir -p /tmp/xcatdb.backup
```

```
for i in `tabdump`;do echo "Dumping $i..."; tabdump $i  
>/tmp/xcatdb.backup/$i.csv; done
```

5. Add a new node “devnode01” to the “nodelist” table and assign it to the “all” and “compute” groups.

```
chtab node=devnode01 nodelist.group=all,compute
```

6. Assign the “site” table “rsh” attribute to “/usr/bin/ssh”.

```
chtab key=rsh site.value=/usr/bin/ssh
```

7. Delete the previously created node from the “nodelist” table.

i) `chtab -d node=devnode01 nodelist`

8. To restore database tables that were dumped with tabdump:

```
cd /tmp/xcatdb.backup
```

```
for i in *.csv;do echo "Restoring $i..."; tabrestore $i; done
```

3.3 Using the node* commands

These are a set of commands for adding (nodeadd), changing (nodech), listing (nodels) and removing (noderm) from the database.

The following are some basic examples of how to use the node* commands:

1. To add a node to the nodelist table with groups all:

```
nodeadd sn1 nodelist.groups=all
```

2. To change the node sn1 os definition in the nodetype table:

```
nodech sn1 nodetype.os=rhel5
```

3. To remove node sn1 from all database tables:

```
noderm sn1
```

4. To list all the nodes in the input noderange:
nodels sn1-sn10

3.3.1 Using xCAT object definition commands

In addition to managing the database tables directly xCAT also supports the concept of data object definitions. Data objects are abstractions of the data that is stored in the xCAT database. This support provides a conceptually simpler implementation for managing cluster data. It is also more consistent with other IBM systems management products such as Director, CSM, and AIX/NIM etc. The attributes and values defined in the data object definitions will still be stored in the database tables defined for xCAT 2.0. These data object definitions should not limit experienced xCAT customers from managing the specific tables directly, if they so desire. A new set of commands is provided to support the object definitions. These commands will automatically handle the storage in and retrieval from the correct tables.

The following data object types are currently supported.

- **site** - Cluster-wide information. All the data is stored in the *site* table.
- **node** - Information for a specific cluster node. The data for a node is stored in multiple tables in the database. The commands that are provided to manage these definitions automatically figure out which attributes are stored in which table. It is therefore not necessary to keep track of a large number of table names and attribute locations.
- **network** - A description of a unique network. This data is stored in the *networks* table.
- **monitoring** - A description of a monitoring plugin. This data is stored in the *monitoring* table.
- **notification** - Defines the Perl modules and commands that will get called for changes in certain xCAT database tables. The data is stored in the *notification* table.
- **group** - Defines a set of nodes. A group definition can be used as the target set of nodes for a specific xCAT operation. It can also be used to define node attributes that are applied to all group members. The group data is stored in multiple tables in the database.
- **policy** - Define the policies used when executing xCAT commands. The data is stored in the *policy* table.

There are four xCAT commands that may be used to manage any of the data object definitions.

- **mkdef** - Make data object definitions.
- **chdef** - Change data object definitions.
- **lsdef** - List data object definitions.

- **rmdef** - Remove data object definitions.

The following are some basic examples of how to use the database object definition commands. For more information on using these commands refer to the MAN pages.

- 1) To view the list of supported object definition types you can issue any of the commands with the “-h” option. Along with the usage you will also see a list of supported object types.

lsdef -h

- 2) To get a description of the attributes that can be defined for each object type you can issue the **lsdef** command with the “-t <object type>” option.

lsdef -h -t node

- 3) To get a list of all the objects currently defined.

lsdef -a

- 4) To get the details of a specific node definition.

lsdef -t node -l -o node01

- 5) To create a very simple node definition.

mkdef -t node -o node02 groups="all,aix"

- 6) To create a node group containing all nodes that have a “nodetype” attribute set to “compute”.

mkdef -t group -o computenodes -w nodetype= compute

- 7) To change the site definition.

chdef -t site -o clustersite rsh=/bin/rsh rcp=/bin/rcp installdir=/xcatinstall

- 8) To remove all node and group definitions.

rmdef -t node,group

- 9) To remove the group called hmcnodes.

rmdef -t group -o hmcnodes

In addition to the standard command line input and output the **mkdef**, **chdef**, and **lsdef** commands support the use of a stanza file format for the input and output of information. Input to a command can be read from a stanza file and the output of a command can be written to a stanza file. A stanza file contains one or more stanzas that provide information for individual object definitions. For example:

5. To create a set of definitions using information contained in a stanza file.

cat mystanzafile | mkdef -z

6. To write all node definitions to a stanza file.

```
lsdef -t node -l -z > nodestanzafile
```

The stanza file support also provides an easy way to backup and restore the cluster data.

For more information on the use of stanza files see the **xcatstanzafile** MAN page.

Note: In some cases the object definition commands may not be able to recognize changes that were made by updating the database tables directly by using the table commands. Generally speaking, the intermixing of the use of the two sets of commands is not recommended.

3.4 Using xCAT hardware commands

3.4.1 Hardware discovery

The following commands can be used to gather information about cluster hardware. See the MAN pages for additional details.

1. **rinv** - Retrieves hardware configuration information for a single or range of nodes and groups.

For example:

```
rinv node5 all
```

```
node5: Machine Type/Model 865431Z
node5: Serial Number 23C5030
node5: Asset Tag 00:06:29:1F:01:1A
node5: PCI Information
node5: Bus VendID DevID RevID Description Slot Pass/Fail
node5: 0 1166 0009 06 Host Bridge 0 PASS
node5: 0 1166 0009 06 Host Bridge 0 PASS
node5: 0 5333 8A22 04 VGA Compatible Controller 0 PASS
node5: 0 8086 1229 08 Ethernet Controller 0 PASS
node5: 0 8086 1229 08 Ethernet Controller 0 PASS
node5: 0 1166 0200 50 ISA Bridge 0 PASS
node5: 0 1166 0211 00 IDE Controller 0 PASS
node5: 0 1166 0220 04 Universal Serial Bus 0 PASS
node5: 1 9005 008F 02 SCSI Bus Controller 0 PASS
node5: 1 14C1 8043 03 Unknown Device Type 2 PASS
node5: Machine Configuration Info
node5: Number of Processors: 2
node5: Processor Speed: 866 MHz
node5: Total Memory: 512 MB
node5: Memory DIMM locations: Slot(s) 3 4
```

2. **rvitals** - Retrieves hardware vital information for a single or range of nodes and groups.

For example:

```
rvitals node5 all
```

```
node5: Frame Voltage (Vab): 201V
node5: Frame Voltage (Vbc): 203V
node5: Frame Voltage (Vca): 202V
node5: Frame Current (Ia): 19A
node5: Frame Current (Ib): 19A
node5: Frame Current (Ic): 20A
node5: System Temperature: 33 C (91.4 F)
node5: Running
```

3. **lsslp** - Queries selected networked services information within the same subnet. If the HMC/IVM that you are interested in discovering is on the same subnet as your Management Node, you can run the **lsslp** to discover and add his hardware to the xCAT database.

Note that the dependent programs **slp_query** and **libslp_client.so** are compiled modules required to perform SLP broadcasts. These modules can be obtained by posting a request to the xCAT mailing list (please specify the target O/S in the request).

For example:

```
lsslp -s HMC
```

```
device type-model serial-number          ip-addresses          hostname
HMC 7310CR2 103F55A 1.1.1.115 2.2.2.164 3.3.3.102 hmc01
HMC 7310CR2 105369A 3.3.3.103 2.2.2.103 1.1.1.163 hmc02
HMC 7310CR3 KPHHK24 3.3.3.154 2.2.2.110 1.1.1.154 hmc03
```

4. **rscan** - Collects node information from one or more hardware control points.

For example:

```
rscan hmc01
```

```
type name          id  type-model  serial- number  address
hmc  hmc01          7310-C05  10F426A  hmc01
```



```

    fsp Server-9117-MMA-SN10F6F3D 9117-MMA 10F6F3D
3.3.3.197
    lpar lpar3 4 9117-MMA 10F6F3D
    lpar lpar2 3 9117-MMA 10F6F3D
    lpar lpar1 2 9117-MMA 10F6F3D
    lpar p6vios 1 9117-MMA 10F6F3D

```

5. **getmacs** – Gathers adapter MAC information from cluster nodes.
For example:

```
getmacs node01
```

```
lpar4:
```

```

#Type Location Code MAC Address Full Path Name Ping Result

ent U9133.55A.10B7D1G-V12-C4-T1 8ee2245cf004 /vdevice/l-
lan@30000004 virtual

```

3.4.2 Hardware Control

The following commands can be used to control cluster hardware. See the MAN pages for additional details.

7. **rnetboot** – Initiate a network boot request on one or more cluster nodes.

For example, to initiate a network boot of the node “node01”, enter:

```
rnetboot node01
```

8. **rpower** – Boots, resets, powers on and off, and queries node hardware, and devices.

For example, to power on a node, enter:

```
rpower -n clsn04 on
```

4.0 Setup the Management Node

4.1 Using Service nodes

If you do not plan on using service nodes, you can skip these sections on Using Service nodes and continue to use the SQLite Default database setup during the installation. Go to section 5 Setup Services.

In large clusters it is desirable to have more than one node (the Management Node) handle the installation of the compute nodes. We call these additional nodes service nodes. You can have one or more service nodes setup to install groups of compute nodes.

The service nodes need to communicate with the xCAT2.0 database on the Management Node and run xCAT command to install the nodes. The service node will be installed with the xCAT code and required the PostgreSQL Database be setup instead of SQLite Default database. PostgreSQL allows a client to be setup on the service node such that the service node can access (read/write) the database on the Management Node (Master Node) from the service node.

4.1.1 Switching to PostgreSQL Database

To setup the postgresql database on the Management Node follow these steps.

This example assumes:

192.168.0.1: ip of master

xcatdb: database name

xcatadmin: database role (aka user)

cluster: database password

192.168.0.10 & 192.168.0.11 (service nodes)

Substitute your address and desired userid , password and database name as appropriate.

The following rpms should be installed from the Fedora8 media on the Management Node (and service node when installed). These are required for postgresql.

1. yum install perl-DBD-Pg postgresql-server postgresql
2. Initialize the database :
service postgresql initdb
3. service postgresql start
4. su – postgres
5. -bash-3.1\$ createuser -P xcatadmin
Enter password for new role: cluster
Enter it again: cluster
Shall the new role be a superuser? (y/n) n

Shall the new role be allowed to create databases? (y/n) n
Shall the new role be allowed to create more new roles? (y/n) n

6. `$ createdb -O xcatadmin xcatdb`

7. `$ exit`

8. `cd /var/lib/pgsql/data/`

9. Edit the hba configuration file:

`vi pg_hba.conf`

#lines should look like:

`local all all ident sameuser`

IPv4 local connections:

`host all all 127.0.0.1/32 md5`

`host all all 192.168.0.1/32 md5`

`host all all 192.168.0.10/32 md5`

`host all all 192.168.0.11/32 md5` where 192.168.0.10 and 11 are service nodes.

10. `vi postgresql.conf`

11. set `listen_addresses` to `'*'`:

`listen_addresses = '*'` This allows remote access. **Note:Be sure and uncomment the line**

12. `service postgresql restart`

13. Backup your data to migrate to the new database

`#mkdir -p ~/xcat-dbback`

`dumpxCATdb -p ~/xcat-dbback`

14. `/etc/sysconfig/xcat` should contain these lines, substitute your cluster facing address for 192.168.0.1, and user and password are xcatadmin cluster in this instance

`XCATCFG='Pg:dbname=xcatdb;host=192.168.0.1|xcatadmin|cluster'`

`export XCATCFG`

`XCATROOT=/opt/xcat`

`export XCATROOT`

15. copy `/etc/sysconfig/xcat` to `/install/postscripts/sysconfig/xcat` for installation on the service nodes.

16. `/etc/xcat/cfgloc` should contain the following line, again substituting your info. This points the xCAT database access code to the new database.

`Pg:dbname=xcatdb;host=192.168.0.1|xcatadmin|cluster`

17. copy /etc/xcat/cfgloc to /install/postscripts/etc/xcat/cfgloc for installation on the service nodes.
18. chmod 700 /etc/sysconfig/xcat and /etc/xcat/cfgloc
19. . /etc/sysconfig/xcat #read the text into the current shell
20. You can add . /etc/sysconfig/xcat to a setup shell script in /etc/profile.d, so the XCATROOT and XCATCFG environment variables are setup when you login.

21. Start the xcatd daemon using the postgresql database
service xcatd restart
22. Restore your database: restorexCATdb -p ~/xcat-dbback to the postgresql database

23. Need to update the policy table:
Run this command to get correct Master node name known by ssl:
openssl x509 -text -in /etc/xcat/cert/server-cert.pem -noout|grep Subject
Subject: CN=mgt.cluster
Subject Public Key Info:
X509v3 Subject Key Identifier:

24. Update the policy table with mgt.cluster output from the command:
chtab priority=5 policy.name=<mgt.cluster> policy.rule=allow. Note this name must be an MN name that is known by the service nodes.

25. Check the database for the following settings:
[root@mn20 ~]# tabdump site
#key,value,comments,disable
"xcatiport","3002",,
"nameservers","11.16.0.1",,
"forwarders","9.114.8.1,9.114.8.2",,
"xcatdport","3001",,
"domain","foobar.com",,

"master","11.16.0.1",, where the Master node is the name or ip address known by the service nodes.

- [root@mn20 ~]# tabdump policy
#priority,name,host,commands,noderange,parameters,time,rule,comments,disable
"1","root",,,,,,"allow",,

"5","mn20" ,,,,,,"allow",, where mn20 is the output of step 26.

"2" ,,,,"getbmconfig" ,,,,"allow",,

"3" ,,,,"nextdestiny" ,,,,"allow",,

"4" ,,,,"getdestiny" ,,,,"allow",,

26. chkconfig postgresql on

27. service postgresql restart

4.2 Define the service nodes in the database

For this example, we have two service nodes rra000 and rrb000. To add the service nodes to the database run the following commands to add and update the service nodes' attributes in the site, nodelist and noderes tables. Note: service nodes are required to be defined with group "service". The commands below are using the group "service" to update all service nodes.

Note: For table attribute definitions run `tabdump -d <table name>`

4.2.1 Define Service Nodes and bmc in nodelist table

```
nodeadd rra000,rrb000 groups=service,ipmi, all
```

```
nodeadd rra000bmc,rrb000bmc groups=bmc,ipmi,all
```

4.2.2 Define Service Nodes in noderes table

```
chtab node=service noderes.netboot=pxe
```

```
chtab node=service noderes.servicenode="11.16.0.1"
```

```
chtab node=service noderes.tftpserver="11.16.0.1"
```

```
chtab node=service noderes.xcatmaster="11.16.0.1"
```

```
chtab node=service noderes.serialport=1
```

```
chtab node=service noderes.service="11.16.0.1"
```

4.2.3 Define Service Nodes in ipmi table

```
nodech rra000 ipmi.bmc=rra000bmc ipmi.userid=USERID
```

```
ipmi.password=PASSWORD
```

```
nodech rrb000 ipmi.bmc=rrb000bmc ipmi.userid=USERID
```

```
ipmi.password=PASSWORD
```

4.2.4 Define Service Nodes and bmc in nodehm table

```
chtab node=service nodehm.cons=ipmi
```

```
chtab node=service nodehm.mgt=ipmi nodehm.serialspeed=19200
nodehm.serialflow=hard
chtab node=bmc nodehm.mgt=ipmi
```

4.2.5 Define Service Nodes and bmc in nodetype table

```
chtab node=service nodetype.arch=x86_64 nodetype.os=fedora8
nodetype.nodetype=osi
chtab node=bmc nodetype.nodetype=rsa
```

4.2.6 Define Service Nodes in site table

```
chtab key=defserialport site.value=1
chtab key=defserialspeed site.value=19200
chtab key=xcatservers site.value=rra000,rrb000
```

4.2.7 Define Service Node OS and Profile attributes

```
chtab node=service nodetype.os=fedora8
chtab node=service noderes.primarynic=eth0 noderes.installnic=eth0
chtab node=service nodetype.profile=service
```

5.0 Setup Services

5.1 Setup networks Table

All networks in the cluster must be defined in the networks table. *makenetworks* runs during the xCAT install and updates the networks table. You should *tabdump networks* to ensure the setting are correct. If any need changing, *tabedit networks* table. Ensure the networks to be managed have the “dynamicrange” set to a hyphenated range of IP addresses to serve as staging for nodes being brought up. If any new networks are added, the *makenetworks* should be run again.

5.2 Setup DNS

Set nameserver, forwarders and domain in the site table

```
chtab key=nameservers site.value=192.168.100.1 (IP of mgmt node)
chtab key=forwarders site.value=172.16.0.1 (how to get to other DNS)
chtab key=domain site.value=foobar.com
```

Edit /etc/hosts:

```
127.0.0.1    localhost.localdomain localhost
::1         localhost6.localdomain6 localhost6
192.168.2.100 b7-eth0
192.168.100.1 b7
192.168.100.10 blade1
192.168.100.11 blade2
192.168.100.12 blade3
172.30.101.133 amm3
```

Run:

```
makenetworks
```

```
makedns
```

```
setup /etc/resolv.conf:
```

```
search foobar.com
```

```
nameserver 192.168.100.1
```

Start dns:

```
service named start
```

```
chkconfig --level 345 named on
```

5.3 Setup AMM

Note: xCAT will soon provide a script to replace this manual process.

For blades, make sure our bladecenter management module is configured for the SNMP protocol.

```
telnet amm3
```

```
env -T mm[1]
```

```
users -1 -ap sha -pp des -ppw PASSWORD
```

```
users -1 -at set
```

As one line copy your id_rsa.pub key from the \$ROOTHOME/.ssh/id_rsa.pub file.

```
users -1 -pk -add ssh-rsa
AAAAB3NzaC1yc2EAAAABIwAAAQEA0u4zf9ULqp5jsZPiVlmcg8TWbPrrIyOK
+bMbHPmId0OEQvs6Opc12XqC4VF6POH8zEu6/YmpPphuDqhOmkou/TXxHgZJ
KQmZ/gFK7Fr9dFzbwA37eE0edeOK4WolwNZgH7t
+4Bm1fJ1sjELVIR1CjFSm59c6Fts83NKIeU6wuhEOzYG1UywyW1Aj/0rSLOk1pS
Fklhu9yXwt9RNVyQva7KKFhXFS51WaFRjyjEMU1Mc/AKaHYnNdehVSm3Bpks
dMIkOVC36/VCXdwqEZWkV0m1pgCIM4K8CPfQUyuP3iaBep2hLA6o8f4bwrXM
XAckrORWCKzFuiV3QoBCAJKxKPQ== root@mgt.cluster
```

NOTE: If you get error with -add type reset and try again. MM Bugs.

Test with:

```
ssh USERID@amm3 exit
```

Log off the management module and test the connection with a query command such as *rpower <noderange> stat* or *rinv <noderange> all*.

TIP to update firmware:

```
Put CNETCMUS.pkt in /tftpboot
```

```
telnet AMM
```

```
env -T mm[1]
```

```
update -v -i TFTP_SERVER_IP -l CNETCMUS.pkt
```

TIP for SOL to work best telnet to nortel switch and type:

```
/cfg/port int1/gig/auto off, for each port.
```

5.4 Setup DHCP

Setup dynamic range for your networks, for example:

```
chtab net=192.168.100.0 networks.dynamicrange=192.168.100.200-192.168.100.250
```


Define dhcp interfaces in site table:

```
chtab key=dhcpinterfaces site.value=eth1
```

Start dhcp:

```
service dhcpd restart
```

```
makedhcp -n
```

```
service dhcpd restart
```

Review the `/etc/dhcp.conf` file created to ensure all your network definitions are correct. Note that the node host definitions will no longer appear here, but rather will appear in the leases file (`/var/lib/dhcpd/dhcpd.leases`) after the initial DHCP request from the node. xCAT 2.0 sets up dhcp to use the OMAPI command shell to setup, query and change the dhcp configuration. See **man omshell**, and <http://linux.die.net/man/3/omapi> for more information.

5.5 Startup TFTP

```
mknb x86_64
```

```
service tftpd restart
```

6.0 Define Compute Nodes

Define the nodes in your cluster by using the `nodeadd` command. Ensure that all nodes, bmc's or management modules, and switches have hosts definitions, or the dhcp configuration will not update, and the `bmcsetup` will not receive meaningful data. (see `nodeadd` command in the xCAT Tables).

Check the 1350 default database template files in `/usr/share/xcat/template/e1350` directory to see if they apply to your environment. These templates, or templates you create from them, can be used to load the database xCAT tables using the ***tabrestore*** `<path to template>` command. The README, in the directory, explains how to use these files.

6.1 Setup the nodelist Table

The nodelist table contains a node definition for each node in the cluster.

You should devise a node naming convention and nodegroups that make it easy for you to use noderanges (see man noderange for the syntax) , especially for large cluster. Most of the xCAT commands take noderange and groups as input.

You can use nodeadd to add a range of redhat nodes:

```
nodeadd node1-node5 groups=all,rhel5
nodeadd node6-node10 groups=all,sles10
```

The nodelist table will look as follows:

```
tabdump nodelist:
#node,groups,status,comments,disable
"node1","all,rhel5",,,
"node10","all,sles10",,,
"node2","all,rhel5",,,
"node3","all,rhel5",,,
"node4","all,rhel5",,,
"node5","all,rhel5",,,
"node6","all,sles10",,,
"node7","all,sles10",,,
"node8","all,sles10",,,
"node9","all,sles10",,,
```

6.2 Setup the noderes Table

The noderes table will define for the node or nodegroup, the service node used to service the node or group, the type of network booting supported, the node which is the tftpserver, dhcpserver,etc as known by the node.

If you are using service nodes, (see sections Using Service nodes) for each node or nodegroup defined in the noderes table change the service node attribute in the noderes table to point to the name or ip address of it's service node.

So for nodes in group rhel5, assign rra000 service node to the node group and the xcatmaster will be the address that the node knows the service node by.

```
chtab node=rhel5 noderes.servicenode=rra000 noderes.xcatmaster=rra000
```

Define the services to run on the servicenode for the node group, for example to setup tftpserver and nfsserver

```
chtab node=rhel5 noderes.tftpserver=rra000 noderes.nfsserver=rra000
```

Whether or not you are using Service Nodes:

Define the type of network booting supported by this type of node (pxe,yaboot). If no service node, the xcatmaster is the Master Node.

```
chtab node=rhel5 noderes.netboot=pxe noderes.master="11.16.0.1"
```

Define the network adapters that will be used for deployment.

```
chtab node=rhel5 noderes.primarynic=eth0 noderes.installnic=eth0
```

6.2.1 Sample noderes table

Your noderes table will end up looking like this (if you use service nodes):

```
node,servicenode,netboot,tftpserver,nfsserver,monserver,kernel,initrd,kcmdline,nfsdir,serialport,installnic,primarynic,xcatmaster,current_osimage,next_osimage,comments,disable
```

```
"rhel5",
```

```
"11.17.0.1","pxe","11.17.0.1","11.17.0.1",,,,,,"1","eth0","eth0","11.18.0.1",,,,
```

6.2.2 Setting up which services run on the Service Nodes

Note: if in the noderes table you have an assigned servicenode for a node, and the field for the service (e.g nfsserver) is left blank, it is assumed that you want that service running on the defined service node. So you can either explicitly assign a service node to a node for any given service, or you can leave the fields blank and the service node assigned to the node will run all services for that node.

The settings for the services in the database will determine which services are setup on the service node. These services are setup when the xcatd daemon is started on the service node.

The services that are setup by xCAT on the service node are as follows:

- nfs (always setup)
- dns
- conserver
- tftp
- http (automatically installed)
- dhcp
- syslog (always setup)

6.3 Setup nodetype table

Define the OS and profile type for the node install.

```
chtab node=rhel5 nodetype.os=rhel5 nodetype.profile=compute
```

6.3.1 Sample nodetype table

Your nodetype table will look something like this:

```
#node,os,arch,profile,nodetype,comments,disable  
"service","fedora8","x86_64","service",,,  
"rhel5","rhel5","x86_64","compute",,,
```

6.4 Setup passwords in passwd table

Add needed passwords to the passwd table to support installs.

```
chtab key=system passwd.password=cluster passwd.username=root  
chtab key=blade passwd.password=PASSWORD passwd.username=USERID  
chtab key=ipmi passwd.password=PASSWORD passwd.username=USERID
```

Note: where key=system, this is the default root password when a node is installed (diskfull).

6.5 Setup deps Table (need?)

6.5.1 Sample deps table:

```
#node,nodedep,msdelay,cmd,comments,disable
```

6.6 Get MAC addresses for all nodes

```
rinv all macs |
```

```
perl -pi -e 's/([^\:]*):.*?ress (\d): (00(:[0-9A-F]{2}){5})/nodech \1 mac.mac=\3 #\2/' |  
grep \#1
```

```
tabdump mac to verify mac addresses in table.
```

6.7 Setup Conserver

Update the nodehm able (*tabedit nodehm*) to set fields for **cons**, **termport**, and **termserver** for your nodes. Currently, supported values for **cons** are “blade” and “ipmi”.

Run `makeconservercf` to generate a `conserver 8` configuration file. Review `/etc/conserver.cf`. Make sure you have valid “trusted” entries in the “`access{ }`” stanza for any host starting a console (most likely your management node).

```
service consver stop
```

```
service consver start
```

```
Test a few nodes with rpower and wcons
```

6.8 Makehosts

If you want *makehosts* to update the `/etc/hosts` file for the defined nodes, `bmcs/mms`, and switches, use `tabedit` to update the *hosts* table with the hostnames and ip addresses to be added to `/etc/hosts`. Then run *makehosts* .

7.0 Installing xCAT nodes

7.1 Installing xCAT Service Nodes (*diskfull*)

Before installing make sure that at least one node in the database has the service node you are going to installed defined as its service node.

Follow the normal steps for an OS install, see “ Adding and Installing Nodes”.

In addition we need to install the xCAT rpms and dependencies on the Service Node:

- a) Create a directory /install/postscripts/xcat/RPMS/noarch
- b) Create a directory /install/postscripts/xcat/RPMS/x86_64
- c) The following rpms should be in /install/postscripts/xcat/RPMS/noarch
 - perl-Expect-1.20-1.noarch.rpm
 - perl-xCAT-2.0-*.rpm
 - xCAT-client-2.0-*.rpm
 - xCAT-nbkernel-x86_64-2.6.18_8-*.noarch.rpm
 - xCAT-nbroot-core-x86_64-2.0-*.noarch.rpm
 - xCAT-nbroot-oss-x86_64-2.0-*.noarch.rpm
 - xCAT-server-2.0-*.noarch.rpm
- a) The following rpms should be in /install/postscripts/xcat/RPMS/x86_64
 - atftp-0.7-1.x86_64.rpm
 - atftp-client-0.7-1.x86_64.rpm
 - atftp-debuginfo-0.7-1.x86_64.rpm
 - conserver-8.1.16-2.x86_64.rpm
 - conserver-debuginfo-8.1.16-2.x86_64.rpm
 - fping-2.4b2_to-2.x86_64.rpm
 - ipmitool-1.8.9-2.x86_64.rpm
 - ipmitool-debuginfo-1.8.9-2.x86_64.rpm
 - perl-IO-Tty-1.07-1.x86_64.rpm
 - xCATsn-2.0-*.x86_64.rpm

Run `nodeset <service nodename>`. Check /install/postscripts/<service nodename> to see if the servicenode is one of theb postinstall scripts that will be run. If not, make sure you have a node defined in the database with this service node designated as its service node in the noderes table.

```
rpower <service nodename> off
```

```
rpower <service nodename> stat - do this until you see status is off
```

```
rpower <service nodename> on
```

You can monitor progress with `nodestat < service nodename>` and tailing `/var/log/messages`

When the service node comes up, it should have services (tftp, dhcp, dns, nfs, conserver) started, you should be able to run database command like `nodels` or `tabdump` on the service node.

7.2 Installing the Compute Nodes (*diskfull*)

7.2.1 Check the site table attribute

```
[root@rh5 xCAT-core]# tabdump site
#key,value,comments,disable
"xcatdport","3001",,
"xcatiport","3002",,
"master","9.114.47.251",,
"domain","ppd.pok.ibm.com",,
"installdir","/install",,
"timezone","America/New_York",,
"nameservers","176.60.50.209",,
```

To change any of these values:

- a) `chtab key=domain site.value=<your domain name>`
For example: `chtab key=domain site.value=clusters.com`
- b) `chtab key=master site.value=<ip address on the cluster network of Master node>`
For example: `chtab key=master site.value=8.77.43.5`
- c) `chtab key=dhcpinterfaces site.value=<comma delimited list of nics to run dhcp>`
For example: `chtab key=dhcpinterfaces site.value=eth1`

7.2.2 Discover Nodes

- 1) Verify that correct chain table values are set:
nodels <noderange> chain.chain chain.ondiscover
If not, add for ipmi nodes:
chtab node=ipmi chain.chain="runcmd=bmcsetup,standby"
chain.ondiscover=nodediscover

Or for blades:

chtab node=blade chain.chain="standby" chain.ondiscover=nodediscover

- 2) Verify that switch table entries exist for the nodes:

nodels <noderange> switch.switch switch.port

If not, add correct values.

- 3) Create the initrd:

rm /tftpboot/pxelinux.cfg/* (mknb will not create new config files if they already exist)

mknb x86_64 (creates the netboot image and writes out the master parameter to the /tftpboot/pxelinux.cfg/default file).

- 4) Make sure your boot sequence is set to boot from network before harddrive:

rbootseq <noderange> list

If not, change it: ***rbootseq <noderange> f,c,n,h***

- 5) Power up the system using ***rpower <noderange> on.***

- 6) Within a few seconds of booting to the network, any BMCs should be configured and be setup to allow ssh. All nodes will be network booted (you can watch /var/log/messages for DHCP and TFTP traffic).

- 7) ***nodels <noderange> vpd.serial vpd.mtm mac.mac*** should show interesting data after discovery.

7.2.3 Install Nodes

- 1) Run ***copycds*** with full path to the ISO images

- 2) Run ***nodech*** (or ***tabedit***) to change nodetype OS and setup node profile :

nodech <noderange> nodetype.os=<os> nodetype.profile=compute

(for now only, the **compute** template file has been provided. See /usr/share/xcats/install/). Current possible values for **os**: rhels5, rhel5,centos5,fedora8

If using 64 bit distro, the **nodetype.arch** should have been populated with at discovery time. If not, set this value.

- 3) Run ***nodech*** (or ***tabedit***) to set noderes nfserver :

nodech <noderange> noderes.nfserver=<server>

(Note: may need to use your management server IP address instead of the hostname for the nfserver for now)

- 4) Also check the following fields to make sure they are set correctly and update as necessary:

a) **noderes.installnic** -- the Ethernet adapter on the node used for installation

b) **noderes.serialport** -- standard SOL for Blades "1", for IPMI nodes "0"

c) **nodehm.serialspeed** -- standard SOL for Blades "19200"

d) **nodehm.serialflow** -- standard SOL for Blades "hard"

- 5) Postscripts that will be run during node install are identified in **/etc/xcat/postscripts.rules** and located in **/install/postscripts**. Not all of the postscripts have been ported to xCAT 2.0 yet, so you may get some “script not found” messages during the postscript processing. Also, the **postage** and **postrules** commands have not been ported yet, so debug may take a little more effort.
 - 6) Run **makedhcp <noderange>** to setup for the installation.
 - 7) Run **nodeset <noderange> install**, to setup for installing the OS.
 - 8) Run **rpower <noderange> on** or **rpower <noderange> reset**, to boot the systems and start the network install process.
- ❖ The kexec to installers doesn't have the client scripts written yet, necessitating the reboot, if wanting to try kexec for now, you have to manually transfer the kernel, initrd, and run kexec -f with the right arguments to the xCAT nbfs environment)

7.3 Install the Service Nodes (stateless)

7.3.1 Build the service node stateless image

The service node stateless images must contain not only the OS, but also xCAT2.0. In addition a number of files are added to the image to support the postgresql database access from the service node to the Management node, and ssh access to the nodes that the service nodes services. Note: the following example assumes you are building the stateless image on the Management Node.

1. Check the service node packaging to see if it has all the rpms required.

```
cd /opt/xcat/share/xcat/netboot/fedora/
```

```
vi service.exlist and service.pklist
```

To add packages:

```
echo vi >>service.pkglis
```

```
echo dhcp >>service.pkglis
```

```
echo atftp >>service.pkglis
```

```
echo bind >>service.pkglis
```

```
echo nfs-utils >>service.pkglis
```

Include things you may need, for example by editing `service.exlist` and remove the following line:

```
./usr/lib/perl5*
```

Edit `compute.exlist`, if necessary, adding lines to remove unnecessary rpms.

2. Run image generation:

```
cd /opt/xcat/share/xcat/netboot/fedora/
```

```
./genimage -i eth0 -n tg3,bnx2 -o fedora8 -p service
```

3. Install xCAT code into the service node image:

```
rm -f /install/netboot/fedora8/x86_64/service/rootimg/etc/yum.repos.d/*
```

```
cp /etc/yum.repos.d/fedora.repo
```

```
/install/netboot/fedora8/x86_64/service/rootimg/etc/yum.repos.d
```

```
yum --installroot=/install/netboot/fedora8/x86_64/service/rootimg install xCATsn
```

4. Update the service node image with the additional files needed for setting up keys and postgresql db when installed:

```
updateSNimage -p /install/netboot/fedora8/x86_64/service/rootimg
```

5. Add automatic configuration of eth1 adapter on the service node

```
chroot /install/netboot/fedora8/x86_64/service/rootimg
```

```
bash-3.2# chkconfig --add xcatd-hack
```

6. Pack the image

```
packimage -o fedora8 -p service -a x86_64
```

7.3.2 Install the Service Nodes

```
nodeset service netboot
```

```
rpower service boot
```

7.3.3 Test Service Node installation

ssh to the service node.

Check to see that the xcat daemon `xcatd` is running.

Run some database command on the service node, e.g `tabdump site`, `nodels` and see that the database can be accessed from the service node.

Check that `/install` and `/tftpboot` are mounted on the service node from the Management Node.

7.4 Installing x86-64 Compute Nodes (stateless)

7.4.1 Building the stateless image

1. Check the compute node packaging to see if it has all the rpms required.

```
cd /opt/xcat/share/xcat/netboot/fedora/
```

```
vi compute.exlist and compute.pklist
```

To add packages:

```
echo vi >>compute.pkglis
```

Include things you may need, for example by editing `compute.exlist` and remove the following line:

```
./usr/lib/perl5*
```

Edit `compute.exlist`, if necessary, adding lines to remove unnecessary rpms.

2. Run image generation:

```
cd /opt/xcat/share/xcat/netboot/fedora/
```

```
./genimage -i eth0 -n tg3,bnx2 -o fedora8 -p compute
```

3. Edit `fstab` in the image

```
cd /install/netboot/fedora8/x86_64/compute/rootimg/etc
```

```
cp fstab fstab.ORIG
```

Edit `fstab`:

Change:

```
devpts /dev/pts devpts gid=5,mode=620 0 0
```

```
tmpfs /dev/shm tmpfs defaults 0 0
```

```
proc /proc proc defaults 0 0
```

```
sysfs /sys sysfs defaults 0 0
```

to:

```
proc /proc proc rw 0 0
sysfs /sys sysfs rw 0 0
devpts /dev/pts devpts rw,gid=5,mode=620 0 0
#tmpfs /dev/shm tmpfs rw 0 0
compute_x86_64 / tmpfs rw 0 1
none /tmp tmpfs defaults,size=10m 0 2
none /var/tmp tmpfs defaults,size=10m 0 2
```

4. Package the image

```
packimage -o fedora8 -p compute -a x86_64
```

7.4.2 Install the stateless image

Install the image on all the compute nodes

```
chtab node=rhel5 nodetype.profile=compute nodetype.os=rhel5
```

```
nodeset rhel5 netboot
```

```
rpower rhel5 boot
```

7.4.3 Update Stateless image

1. Update image using YUM

NOTE: before YUM/RPM commands type:

```
rm /install/netboot/fedora8/x86_64/compute/rootimg/var/lib/rpm/__db.00*
```

```
rm -f /install/netboot/fedora8/x86_64/compute/rootimg/etc/yum.repos.d/*
```

```
cp /etc/yum.repos.d/fedora.repo /install/netboot/fedora8/x86_64/compute/rootimg/etc/
yum.repos.d
```

Now install vi into the image:

```
yum --installroot=/install/netboot/fedora8/x86_64/compute/rootimg install vi
```

2. Update image using RPM

```
rpm --root /install/netboot/fedora8/x86_64/compute/rootimg -Uvh blah.rpm
```

3. Repackage

```
packimage -o fedora8 -p compute -a x86_64
```

4. Install:

```
nodeset rhel5 netboot
```

```
rpower rhel5 boot
```

7.5 Installing PPC64 Compute Nodes (Stateless)

7.5.1 iSCSI install

```
yum install yaboot-xcat scsi-target-utils
```

```
chtab key=iscsidir site.value=/install/iscsi
```

Pick one of the PPC64 nodes for the iSCSI install

Note: make sure the root userid and password are in the iscsi table

```
chtab node=ppc6401 iscsi.userid=root iscsi.password=cluster iscsi.server="11.16.0.1"
```

```
chtab node=ppc6401 noderes.nfsserver="11.16.0.1" (MasterNode)
```

```
chtab node=ppc6401 nodetype.os=fedora8 nodetype.profile=iscsi groups,=iscsi  
iscsi.server="11.16.0.1"
```

```
service tgtd restart
```

```
nodech ppc6401 iscsi.file=
```

```
setupiscsidev -s8192 ppc6401
```

```
nodeset ppc6401 install
```

NOTE: for reinstall:

```
chtab node=ppc6401 nodetype.profile=iscsi
```

7.5.2 Build PPC64 Stateless image:

1. Logon to the node

```
ssh ppc6401
mkdir /install
mount 11.16.0.1:/install /install
```

2. Create fedora.repo:

```
cd /etc/yum.repos.d
rm -f *.repo
```

Put the following lines in /etc/yum.repos.d/fedora.repo:

```
[fedora]
name=Fedora $releasever - $basearch
baseurl=file:///install/fedora8/ppc64
enabled=1
gpgcheck=0
```

3. Test with: yum search gcc
4. Copy the executables and files needed from the Management Node:

```
cd /root
scp 11.16.0.1:/opt/xcat/share/xcat/netboot/fedora/genimage .
scp 11.16.0.1:/opt/xcat/share/xcat/netboot/fedora/geninitrd .
scp 11.16.01./opt/xcat/share/xcat/netboot/fedora/compute.ppc64.pkglist .
```
5. Generate the image:

```
./genimage -i eth0 -n tg3 -o fedora8 -p compute
```

NOTE: iSCSI, QS22, tg3, all slow, take a nap

7.5.3 Install PPC64 Stateless image

On the Management Node:

1. Edit fstab in the image

```
cd /install/netboot/fedora8/ppc64/compute/rootimg/etc
cp fstab fstab.ORIG
```

Edit fstab:

Change:

```
devpts /dev/pts devpts gid=5,mode=620 0 0
tmpfs /dev/shm tmpfs defaults 0 0
proc /proc proc defaults 0 0
sysfs /sys sysfs defaults 0 0
```

to:

```
proc /proc proc rw 0 0
sysfs /sys sysfs rw 0 0
devpts /dev/pts devpts rw,gid=5,mode=620 0 0
#tmpfs /dev/shm tmpfs rw 0 0
compute_ppc64 / tmpfs rw 0 1
none /tmp tmpfs defaults,size=10m 0 2
none /var/tmp tmpfs defaults,size=10m 0 2
```

2. Pack the image

```
packimage -o fedora8 -p compute -a ppc64
```

3. Install the image on all the PPC64 nodes

```
chtab node=ppc64 nodetype.profile=compute nodetype.os=fedora8
nodeset ppc64 netboot
rpower ppc64 boot
```

7.5.4 Update PPC64 Stateless image

NOTE: before YUM/RPM commands type:

```
rm /install/netboot/fedora8/ppc64/compute/rootimg/var/lib/rpm/__db.00*
```

1. Update image using YUM

```
rm -f /install/netboot/fedora8/ppc64/compute/rootimg/etc/yum.repos.d/*
cp /etc/yum.repos.d/fedora.repo
/install/netboot/fedora8/ppc64/compute/rootimg/etc/yum.repos.d
```

Now install vi into the image:

```
yum --installroot=/install/netboot/fedora8/ppc64/compute/rootimg install vi
```

2. Update image using RPM

```
rpm --root /install/netboot/fedora8/ppc64/compute/rootimg -Uvh
/install/fedora8/ppc64/Packages/vim-minimal-7.1.135-1.fc8.ppc.rpm
```

3. Update the image by running genimage

Add packages to compute.ppc64.pkglist and rerun genimage

4. packimage -o fedora8 -p compute -a ppc64

8.0 Using xCAT Notification Infrastructure

With xCAT 2.0, you can monitor xCAT database for changes such as nodes entering/leaving the cluster, hardware updates, node liveness (to be added later) etc. In fact anything stored in the xCAT database tables can be monitored through the xCAT notification infrastructure. To start getting notified for changes, simply register your Perl module or command as the following:

regnotif *filename tablename -o actions*

where

filename is the full path name of your Perl module or command.

*tablename*s is a comma separated list of table names that you are interested in.

actions is a comma separated list of data table actions. 'a' for row addition, 'd' for row deletion and 'u' for row update.

Example:

regnotif /opt/xcat/lib/perl/xCAT_monitoring/mycode.pm nodelist,nodhm -o a,d

regnotif /usr/bin/mycmd switch,noderes -o u

Use the following command to view all the modules and commands registered.

tabdump notification

To un-register, just do the following:

unregnotif *filename*

Example:

unregnotif /opt/xcat/lib/perl/xCAT_monitoring/mycode.pm

unregnotif /usr/bin/mycmd

If the *filename* specifies a Perl module, the package name must be **xCAT_monitoring::xxx**. It must implement the following subroutine which will get called when database table change occurs:

processTableChanges(*tableop, table_name, old_data, new_data*)

where:

tableop Table operation. It can be 'a' for row addition, 'd' for row deletion and 'u' for row update.

tablename The name of the database table whose data has been changed.

old_data An array reference of the old row data that has been changed. The first element is an array reference that contains the column names. The rest of the elements are array references each contains attribute values of a row. It is set when the action is u or d.

new_data A hash reference of the new row data; only changed values are in the hash. It is keyed by column names. It is set when the action is u or a.

If the file name specifies a command (written by any programming languages or scripts), when the interested database table changes, the info will be fed to the command through the standard input. The format of the data in the STDIN is as following:

action(a, u or d)

tablename

```

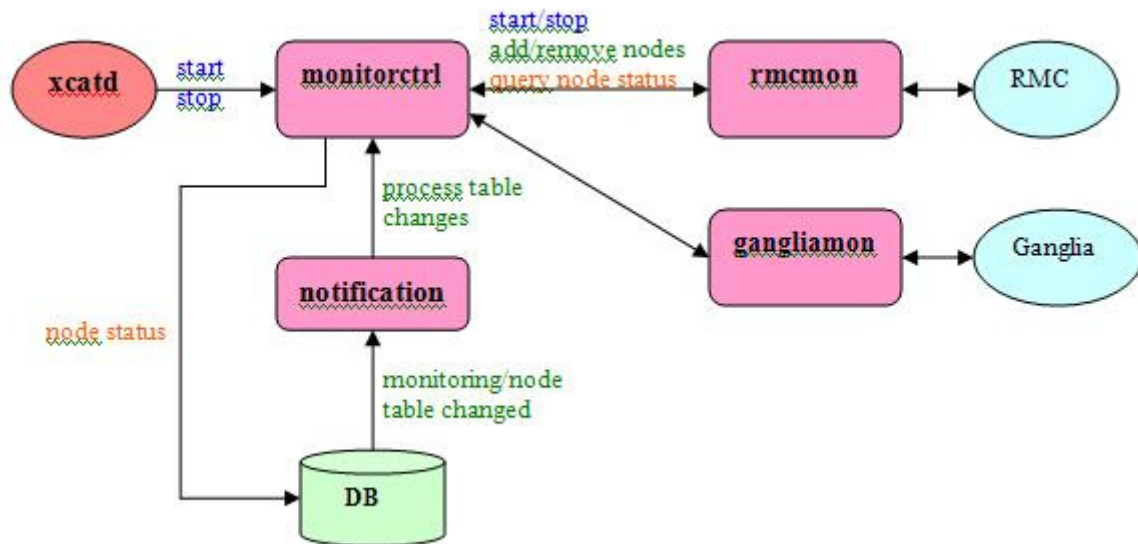
[old value]
col1_name,col2_name...
col1_val,col2_val,...
col1_val,col2_val,....
...
[new value]
col1_name,col2_name,...
col1_value,col2_value,...
...

```

The sample code can be found under `/opt/xcat/lib/perl/xCAT_monitoring/samples/mycode.pm` on a installed system.

8.1 Using xCAT Monitoring Plug-in Infrastructure:

With xCAT 2.0, you can integrate 3rd party monitoring software into your xCAT cluster. The idea is to use monitoring plug-in modules that act as bridges to connect xCAT and the 3rd party software. The functions of a monitoring plug-in module include initializing the 3rd party software, informing it with the changes of the xCAT node list, setting it up to feed node status back to xCAT etc. The following figure depicts the relationship and data flow among xcatd, plug-in modules and 3rd party monitoring software.



To use this infrastructure, first create a monitoring plug-in module and put it under `/opt/xcat/lib/perl/xCAT_monitoring/` directory. If the file name is xxx.pm then the

package name will be **xCAT_monitoring::xxx**. The following is a list of subroutines that a plug-in module must implement:

start
stop
supportNodeStatusMon
startNodeStatusMon
stopNodeStatusMon
addNodes
removeNodes
processSettingChanges
getDiscription

Please refer to `/opt/xcat/lib/perl/xCAT_monitoring/samples/tmplatemon.pm` for the detailed description of the functions.

Second, register the module in xCAT **monitoring** table using the following command:

monstart *name* [-n|--nodestatmon] [-s|--settings *settings*]

where

name is the monitoring plug-in module short file name without the extension. In this case xxx. Use ***monls -a*** command to list all the possible plug-in module names.

-n or --nodestatmon indicates it can help feeding the node status to xCAT. The node status is stored in the ***status*** column of the ***nodelist*** table.

-s or --settings specifies the plug-in specific settings. These setting will be used by the plug-in to customize certain entities for the plug-in or the third party monitoring software. The format of the setting string is: `[key=value],[key=value]...` Please note that the square brackets are needed. e.g. `[mon_interval=10],[toggle=1]`

Example:

monstart xxx -n (with feeding the node status to xCAT table)

or

monstart xxx (not feeding the node status to xCAT table)

Once it is registered, xCAT will automatically, through the plug-in module, start the 3rd party software for monitoring. To unregister the monitoring plug-in and stop the monitoring use this command:

monstop *name*

Example:

monstop xxx

Though you can write your own monitoring plug-in modules, over the time, xCAT will supply a list of built-in plug-in modules for the most common monitoring software. They are:

- xCAT (xcatmon.pm) (released in this beta)
- RMC (rmcmon.pm)
- Ganglia (gangliamon.pm)
- Nagios (nagiosmon.pm)
- SNMP (snmpmon.pm)

xcatmon.pm is included in this release. It provides node liveness monitoring using fping. This can be used if no other 3rd party software is used for node status monitoring. The *status* column of the *nodelist* table will be updated periodically with the latest node liveness status by this plug-in. To activate, use the monstart command:

```
monstart xcatmon -n -s [ping-interval=2]
```

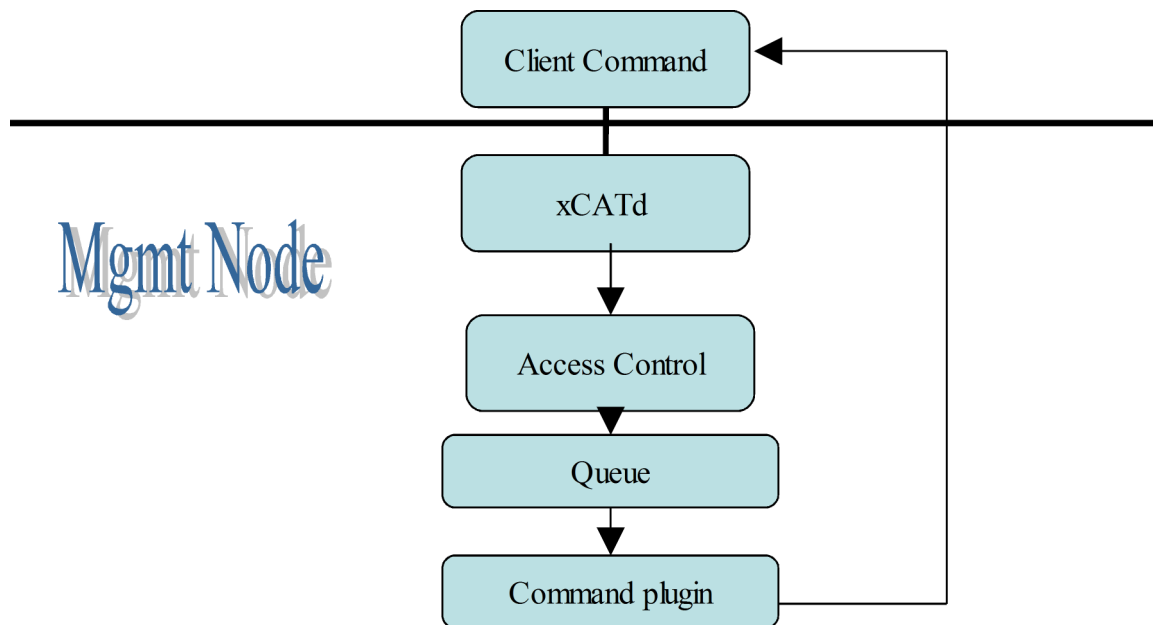
where 2 means that the nodes are pinged for status every 2 minutes.

9.0 xCAT Architecture

General/Overall Concepts

The heart of the xCAT 2.0 architecture is the xCAT daemon (xcatd) on the management node. This receives requests from the client, validates the requests, and then invokes the operation. The xcatd daemon also receives status and inventory information from the nodes

Client



9.1 Client/Server

9.2 Flow

- User invokes an xcat cmd on the client
- The cmds can either be a sym link to xcatclient or a thin wrapper that calls xcatclient.
- Some cmds will implement their own xcatclient function, if they have more processing than the generic xcatclient function supports. (e.g. xdsh/xdcp).
- The xcatclient function packages the info into xml and passes it to xcatd
- xcatd receives the request and forks to process the request
- The ACL/Role Policy Engine determines whether this person is allowed to execute this request. It evaluates the following info:
 - The cmd name and args
 - Who executed the cmd on the client machine
 - The hostname/IP address of the client machine
 - The node range passed to the cmd
- If the ACL check is approved, the cmd is passed to the Queue:
 - The queue can run the action in either of 2 modes. The client cmd wrapper decides which mode to use (although it can give the user a flag to specify):

- Keep the socket connection with the client open for the life of the action and continue to send back the output of the action as it is produced.
 - Initiate the action, pass the action ID back to the client, and close the connection. At any subsequent time, the client can use the action ID to request the status and output of the action. This is intended long running cmds.
- The Queue logs every action performed, including date/time, cmd name, arguments, who, etc.
- In phase 2, the Queue will support locking (semaphores) to serialize actions that should not be run simultaneously.
- To invoke the action, the xml is passed to the process_request() function of the appropriate plugin pm which contains the code for the function being run.
 - With the request examined per policy table, and noderange expanded to nodes, the request is passed in its entirety (including tags otherwise ignored) to a plugin's process_request function, which will receive two arguments, the first the aforementioned hash reference, the second a reference to a callback function to call per response message to send back.
 - The appropriate pm is chosen by loading all of the plugins from /usr/lib/xcat/plugins and invoking handled_commands to see which cmds each pm handles.
 - Data is returned from the command plugin back to the client command handle_response routine.

10.0 References

- a) XCAT 2.0 AIX Cookbook:
http://sourceforge.net/docman/display_doc.php?docid=86730&group_id=208749
- b) xCAT 2.0 RR Cookbook: