

# Introduction to Commodity Linux clusters (part 2)

{ With examples from ICTP HPC cluster “Argo”

Maria Verina

Abdus Salam International Center for Theoretical  
Physics (ICTP)  
Trieste, Italy

The Information and Communication Technology  
Section (ICTS)

Member of HPC team



# Credits:

- ◆ Partially based on:

“Installation and configuration of Linux Cluster”

(smr2613, ICTP 2014)

Addisu Gezahagn  
University of Trieste  
ICTP, Trieste  
[asemie@ictp.it](mailto:asemie@ictp.it)

- ◆ With additional elements from:

Clement Onime  
Antonio Messina  
Moreno Baricevic



# From Function to Structure

- ◆ Introduction to Commodity Linux clusters, part 1
- ◆ Introduction to Commodity Linux clusters, part 2

The problem/Function:

“make parallel/HPC apps run well for your users”

Cluster Structure:

... fairly complex ...

- ◆ But it **CAN** be done!



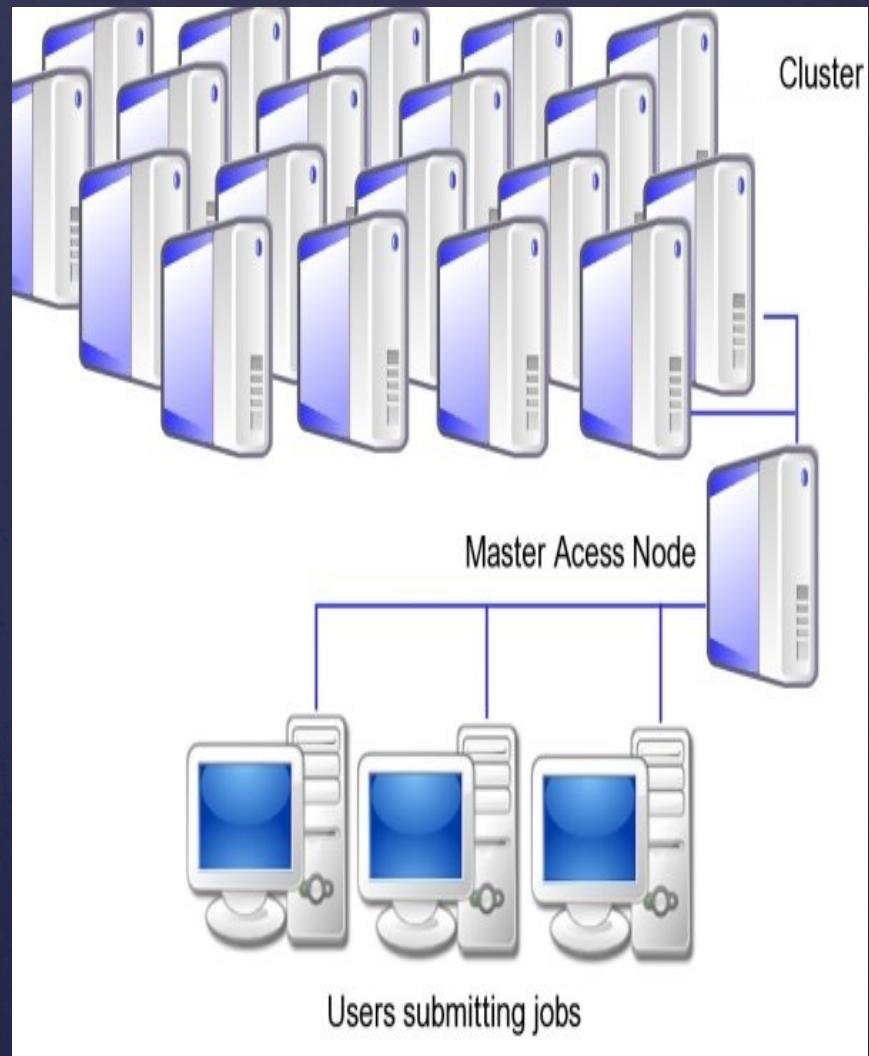
# Why Complex?

- ◆ Number of components and technologies
- ◆ New to IT team: think big
- ◆ Commodity parts, no Blue Print:
  - Each cluster is unique
- ◆ Life cycle:
  - Build/operate daily
  - Initial/extensions
    - (heterogeneous?)
- ◆ Parts “fit” together
- ◆ Best practices exist
- ◆ Standards/protocols exist
- ◆ Need broad range of expertise
- ◆ Team effort



# What is Cluster computer?

- It is a single **logical unit** consisting of multiple computers that are linked through a network



# ICTP HPC cluster experience

- ◆ ICTP HPC cluster “Argo”
- ◆ 198 nodes
  - 2640 cores
  - Intel x86\_64 +GPU
  - heterogeneous
- ◆ Total RAM: 6.8TB
- ◆ IB: 40 Gbps QDR
- ◆ Storage: ~300TB NFS
- ◆ racks: 5
- ◆ Users: ~200

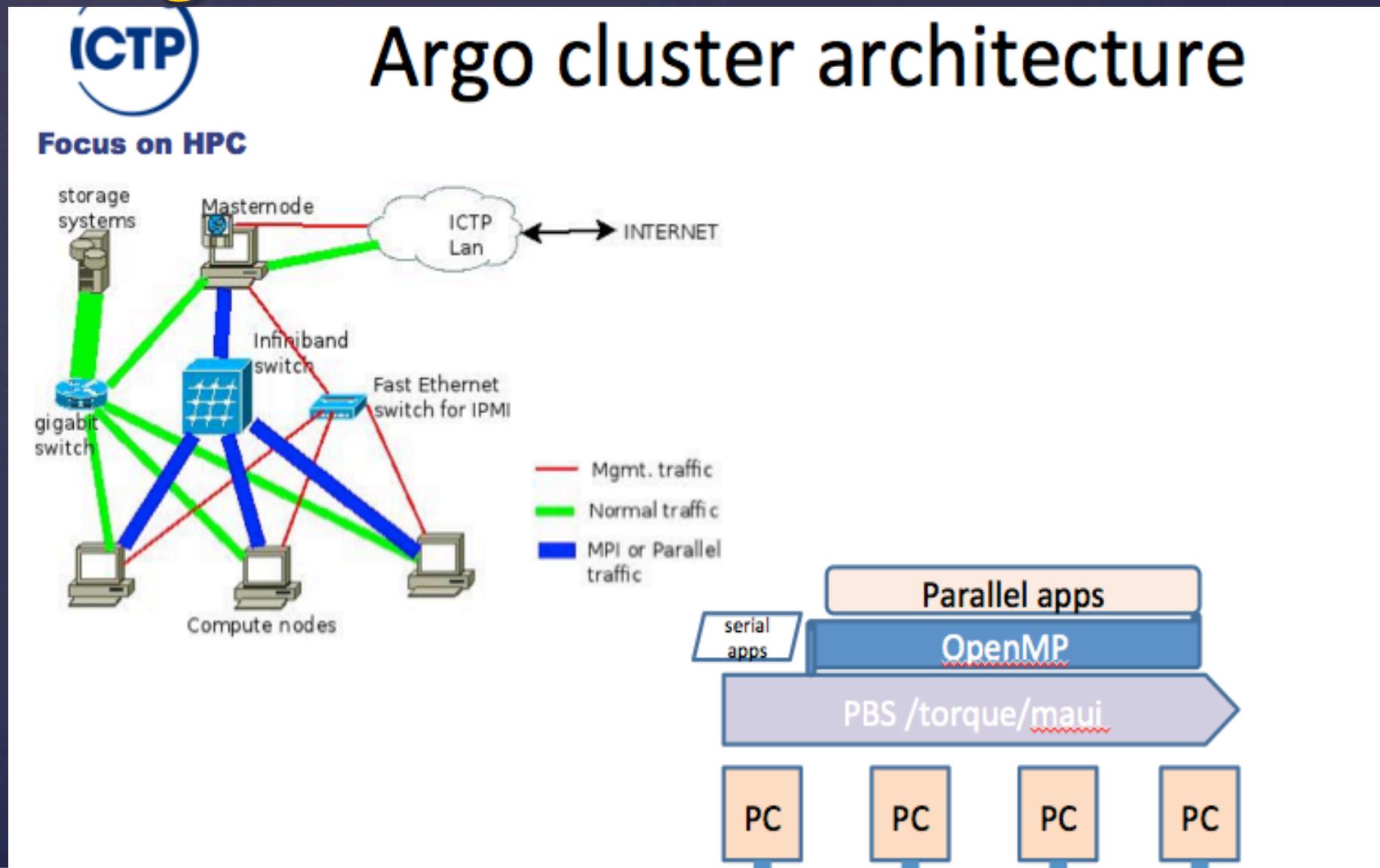
The story of growth.

Queues: dedicated + common

Storage: dedicated + common



# Argo architecture



Taken from Clement's slide



# ICTP HPC Cluster Argo in pictures





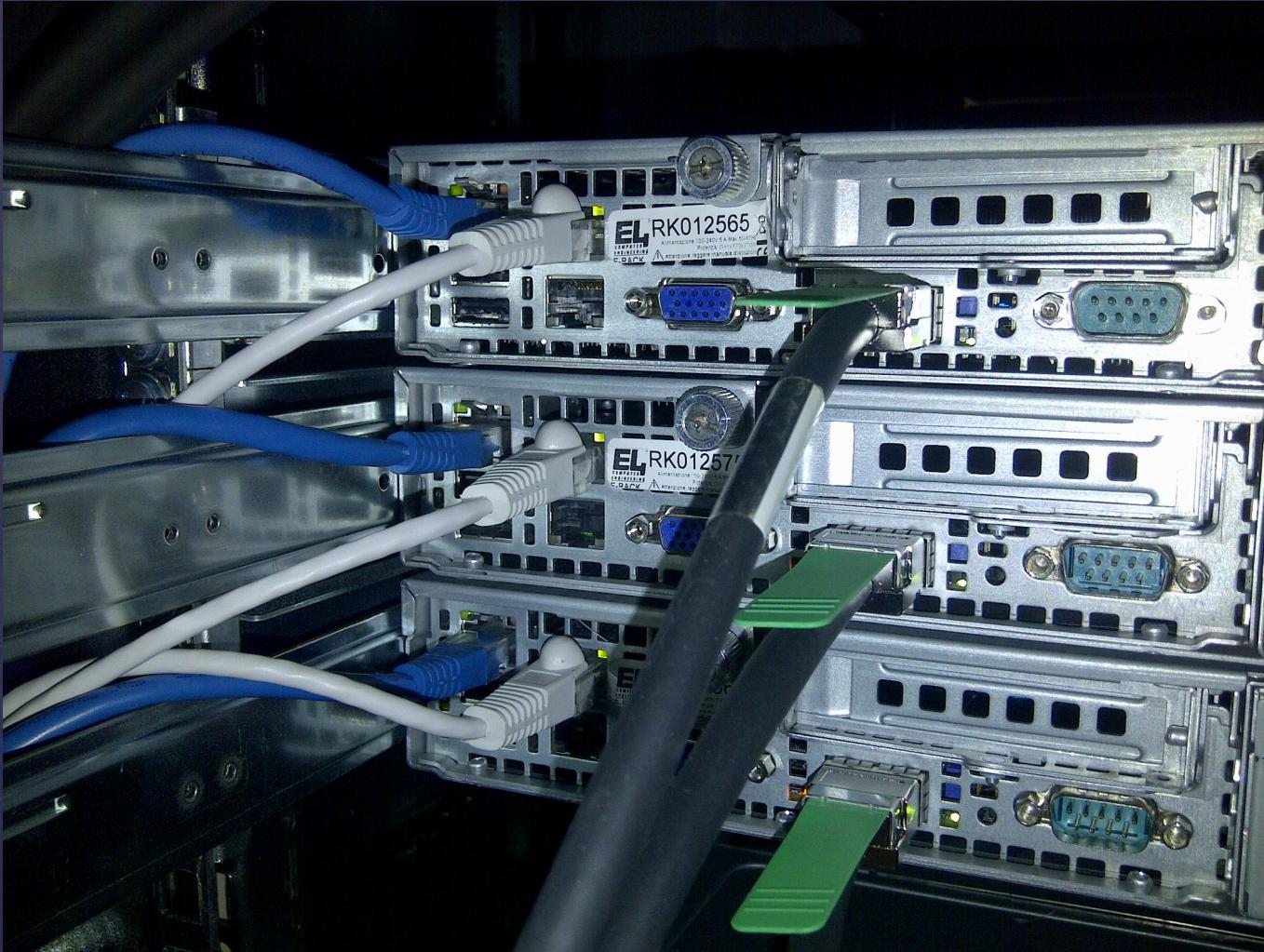


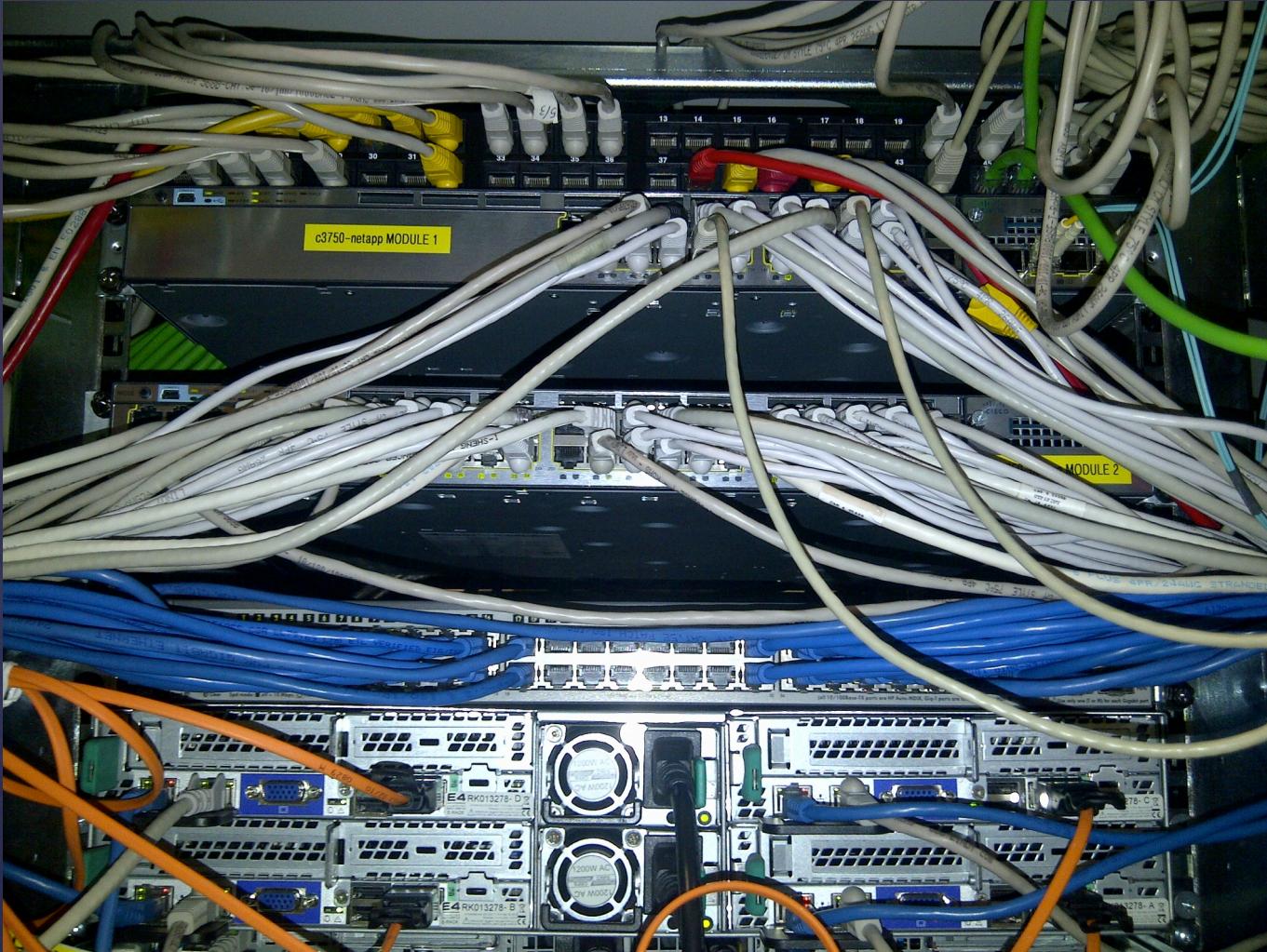
















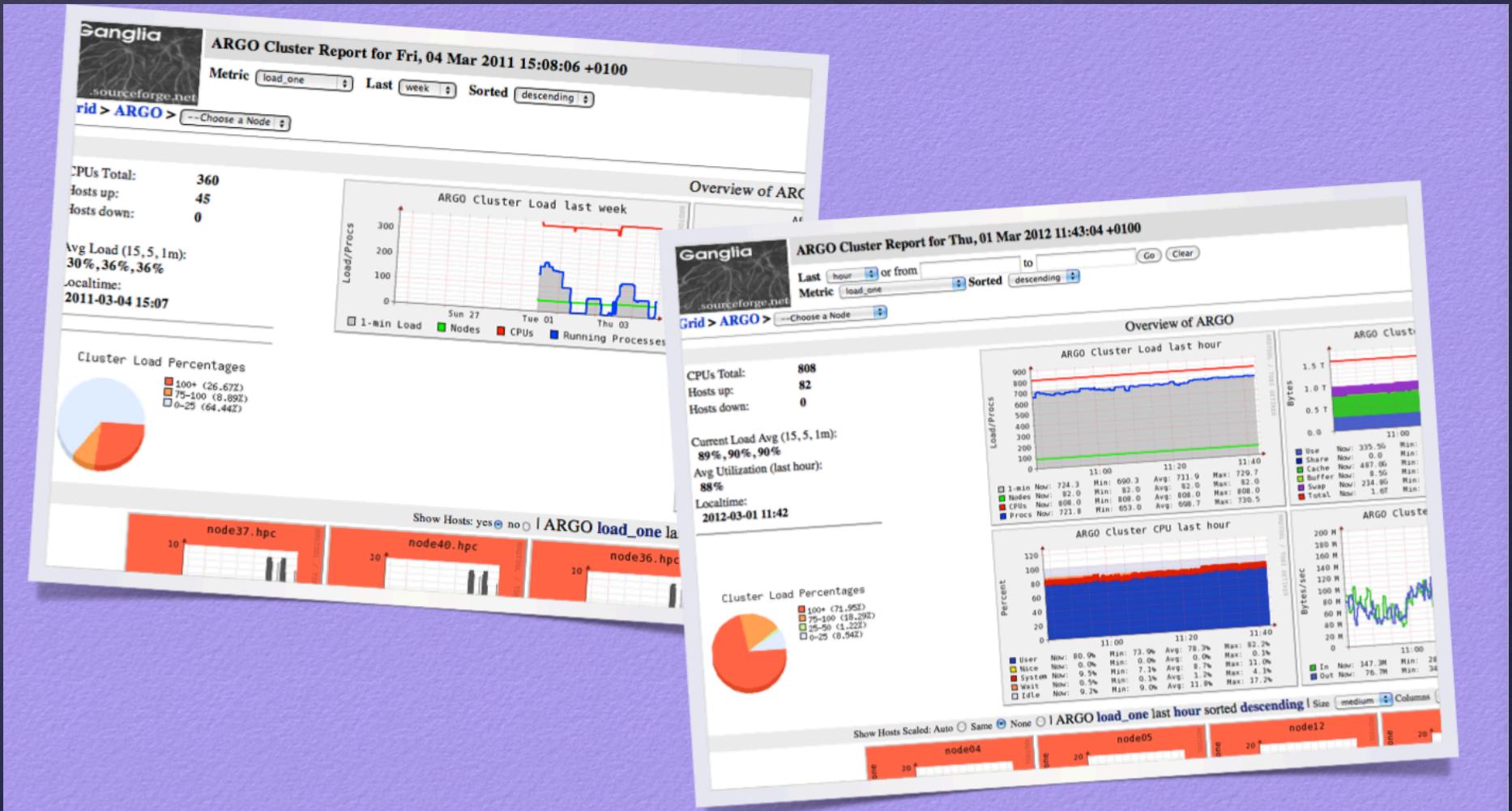


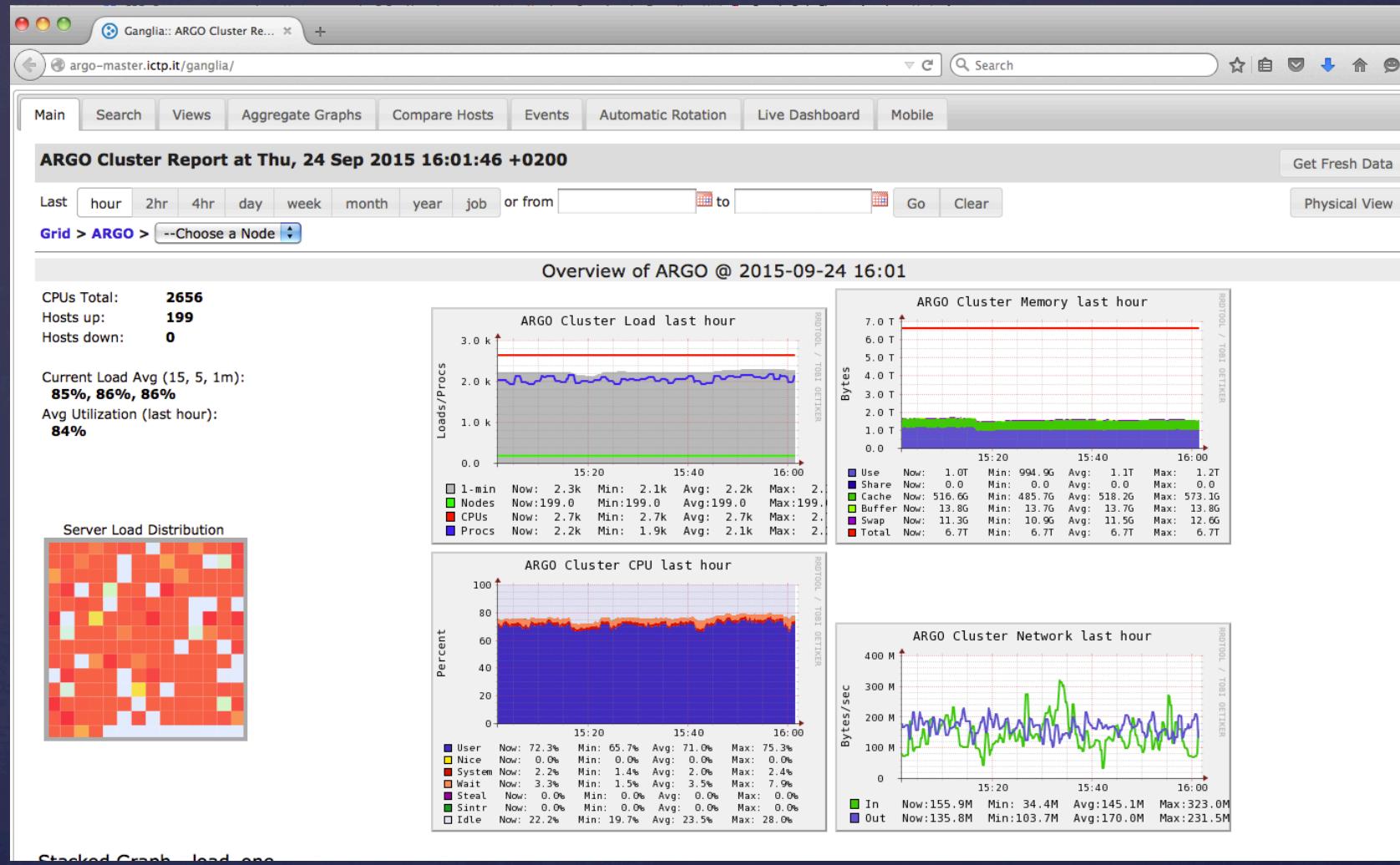




# snapshots from Argo's life

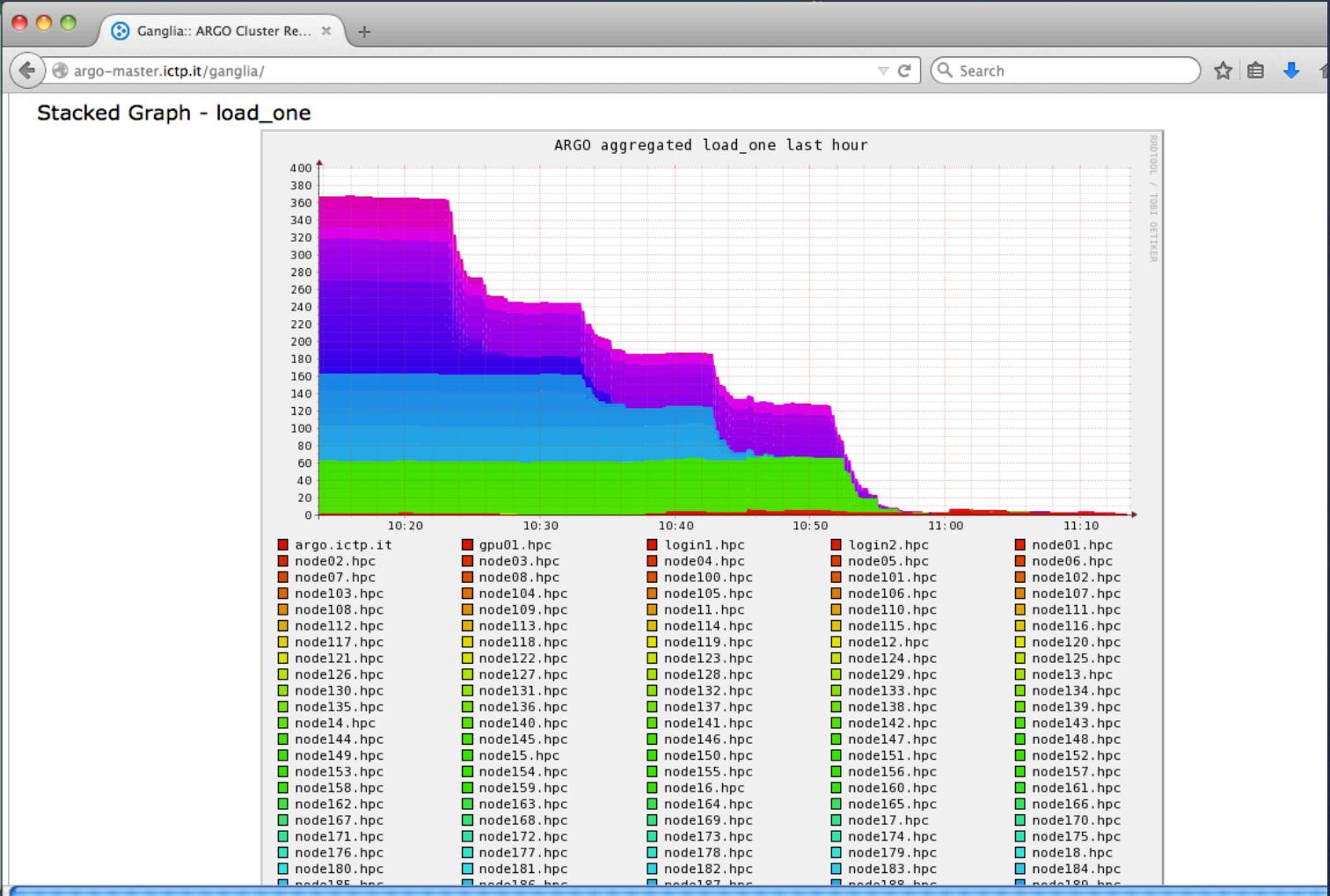


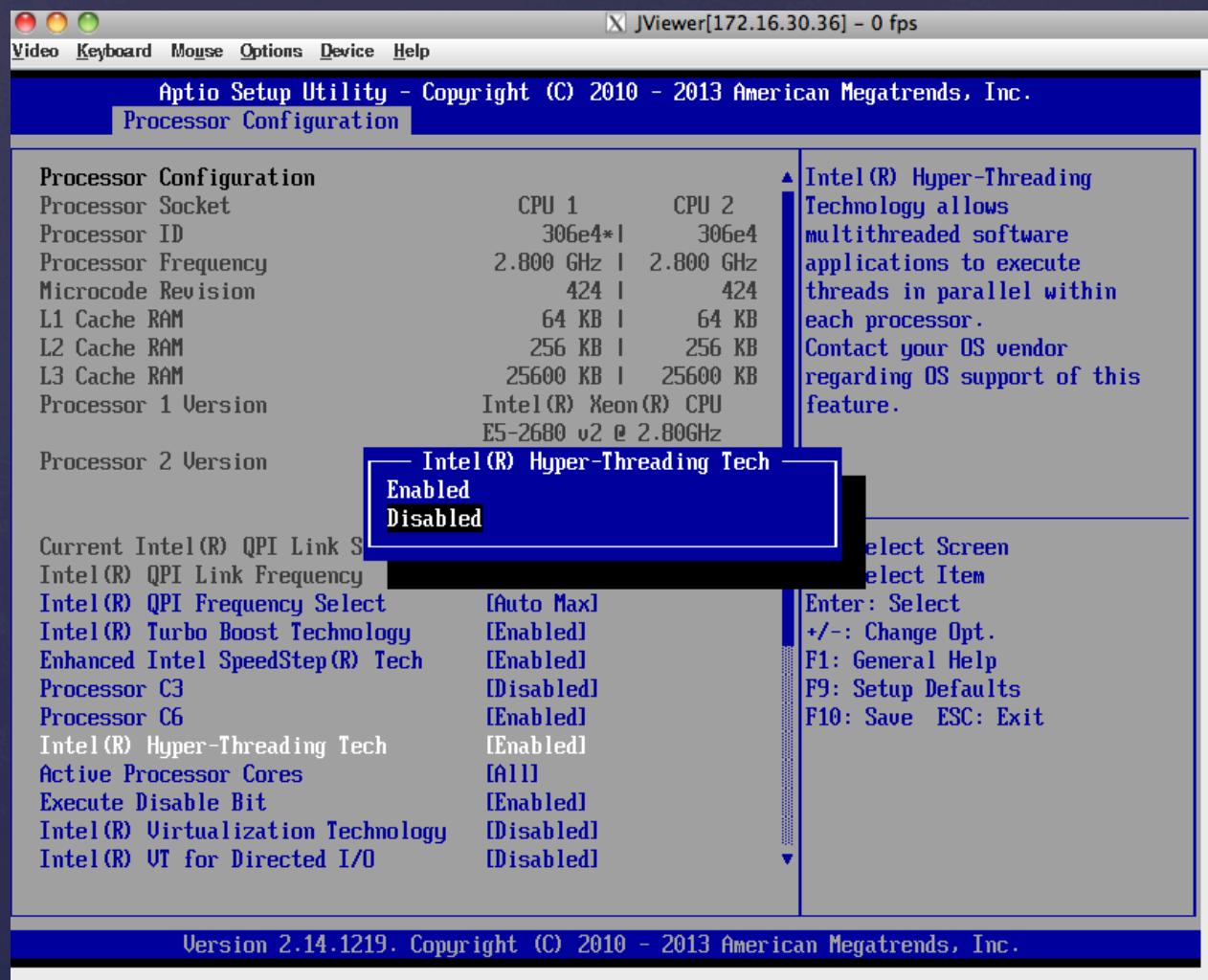




Stacked Graph - load\_avg







JViewer[172.16.30.37] - 10 fps

Video Keyboard Mouse Options Device VirtualKeyboard Help

Memtest86+ v4.20 | Pass 27% #####

Core i7 (32nm) 2400 MHz | Test 46% #####

L1 Cache: 32K 80005 MB/s | Test #4 [Moving inversions, random pattern]

L2 Cache: 256K 31581 MB/s | Testing: 16G - 18G 24G

L3 Cache: 12288 22019 MB/s | Pattern: d0deded8

Memory : 24G 9796 MB/s |-----

Chipset : Core IMC 32nm (ECC : Detect / Correct) Scrub+ / BCLK : 133 MHz

Settings: RAM : 533 MHz (DDR3-1066) / CAS : 7-7-7-20 / Triple Channel

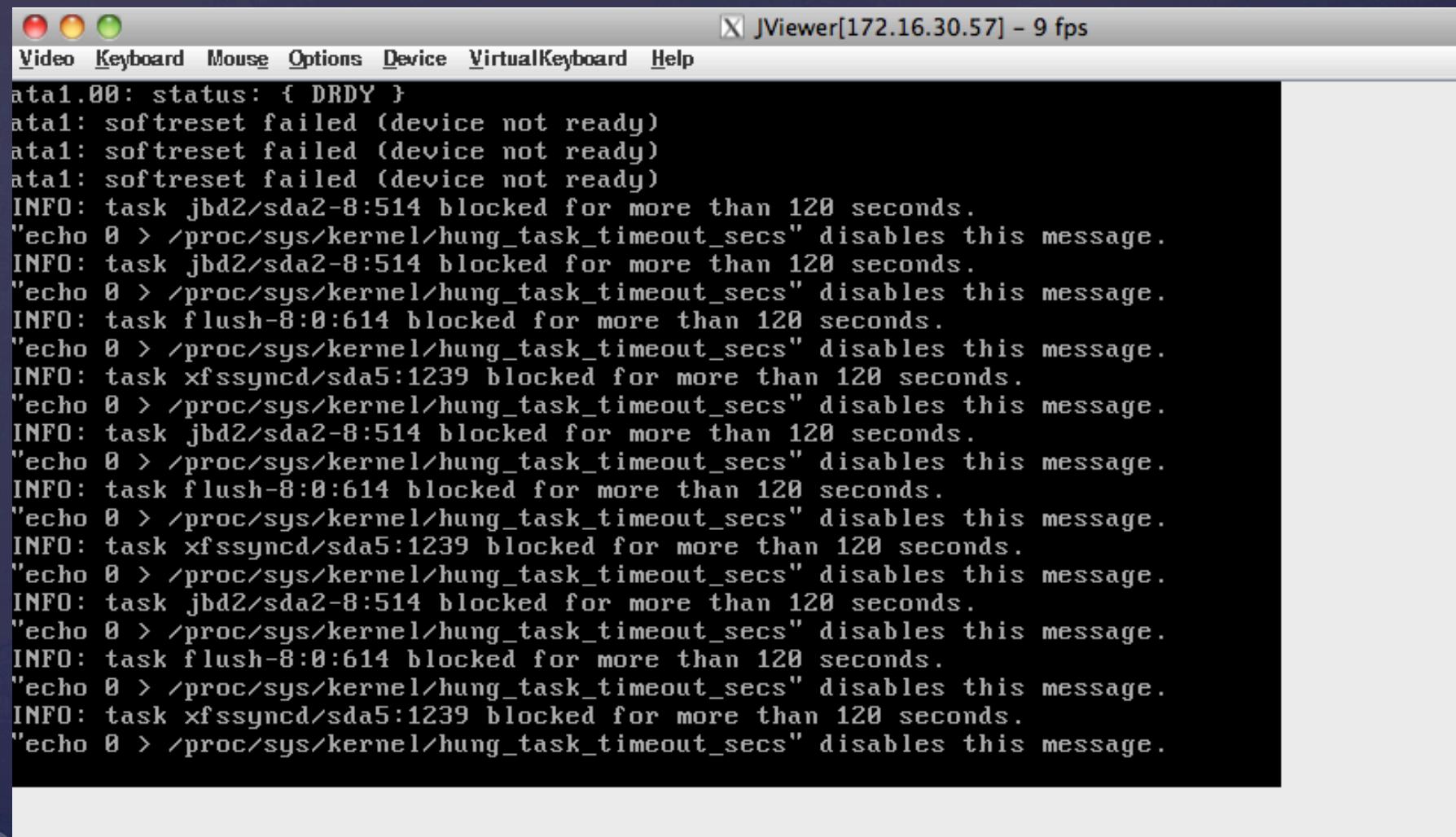
WallTime	Cached	RsvdMem	MemMap	Cache	ECC	Test	Pass	Errors	ECC Errs
19:02:48	24G	200K	e820	on	off	Std	4	0	

Memory SPD Informations

\*\*\*\*\*Pass complete, no errors, press Esc to exit\*\*\*\*\*

(ESC)Reboot (c)configuration (SP)scroll\_lock (CR)scroll\_unlock

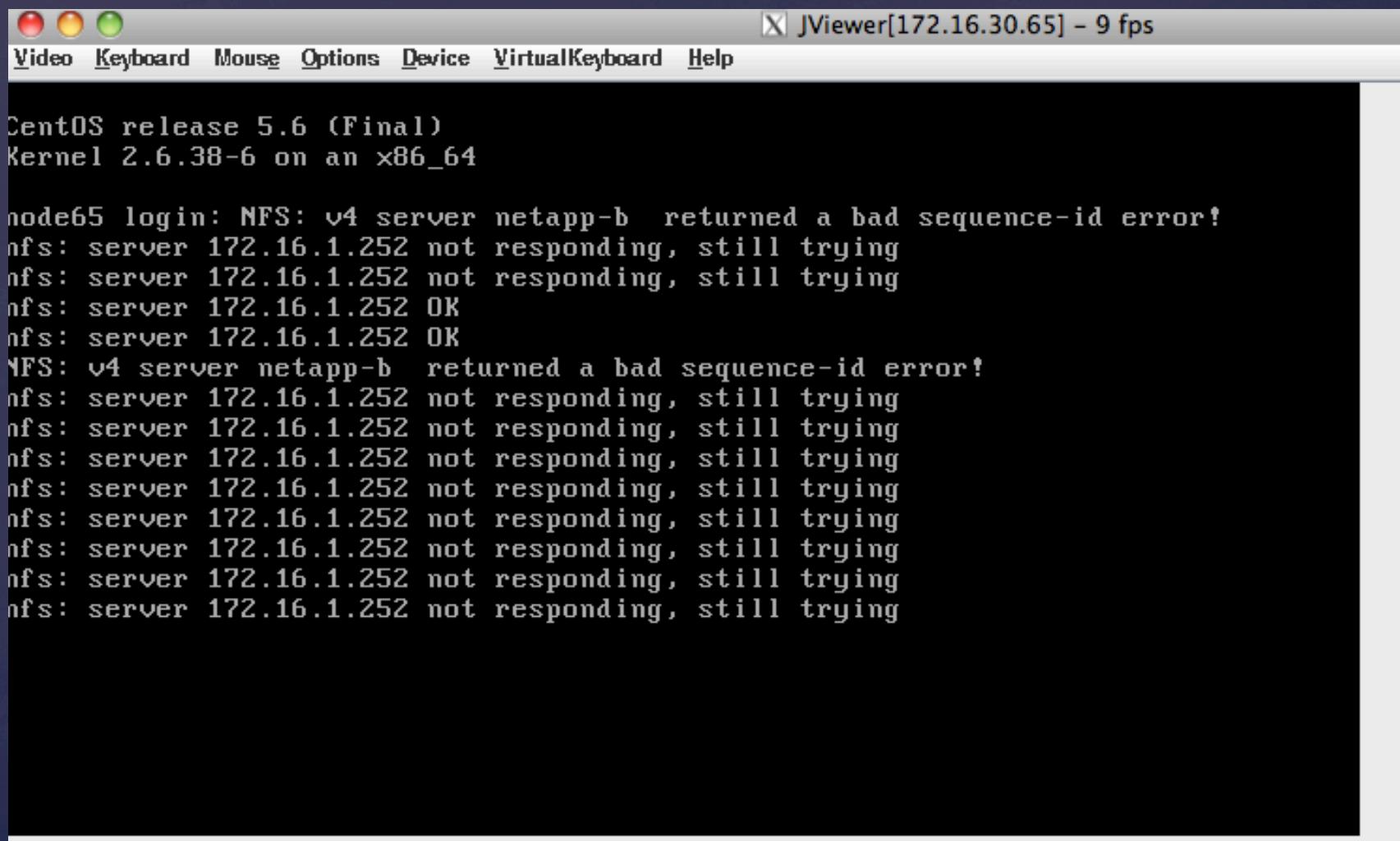




The screenshot shows a terminal window titled "JViewer[172.16.30.57] - 9 fps". The window has a menu bar with options: Video, Keyboard, Mouse, Options, Device, VirtualKeyboard, and Help. The main area of the window displays a continuous stream of kernel log messages. The messages are as follows:

```
ata1.00: status: { DRDY }
ata1: softreset failed (device not ready)
ata1: softreset failed (device not ready)
ata1: softreset failed (device not ready)
INFO: task jbd2/sda2-8:514 blocked for more than 120 seconds.
"echo 0 > /proc/sys/kernel/hung_task_timeout_secs" disables this message.
INFO: task jbd2/sda2-8:514 blocked for more than 120 seconds.
"echo 0 > /proc/sys/kernel/hung_task_timeout_secs" disables this message.
INFO: task flush-8:0:614 blocked for more than 120 seconds.
"echo 0 > /proc/sys/kernel/hung_task_timeout_secs" disables this message.
INFO: task xfssyncd/sda5:1239 blocked for more than 120 seconds.
"echo 0 > /proc/sys/kernel/hung_task_timeout_secs" disables this message.
INFO: task jbd2/sda2-8:514 blocked for more than 120 seconds.
"echo 0 > /proc/sys/kernel/hung_task_timeout_secs" disables this message.
INFO: task flush-8:0:614 blocked for more than 120 seconds.
"echo 0 > /proc/sys/kernel/hung_task_timeout_secs" disables this message.
INFO: task xfssyncd/sda5:1239 blocked for more than 120 seconds.
"echo 0 > /proc/sys/kernel/hung_task_timeout_secs" disables this message.
INFO: task jbd2/sda2-8:514 blocked for more than 120 seconds.
"echo 0 > /proc/sys/kernel/hung_task_timeout_secs" disables this message.
INFO: task flush-8:0:614 blocked for more than 120 seconds.
"echo 0 > /proc/sys/kernel/hung_task_timeout_secs" disables this message.
INFO: task xfssyncd/sda5:1239 blocked for more than 120 seconds.
"echo 0 > /proc/sys/kernel/hung_task_timeout_secs" disables this message.
```





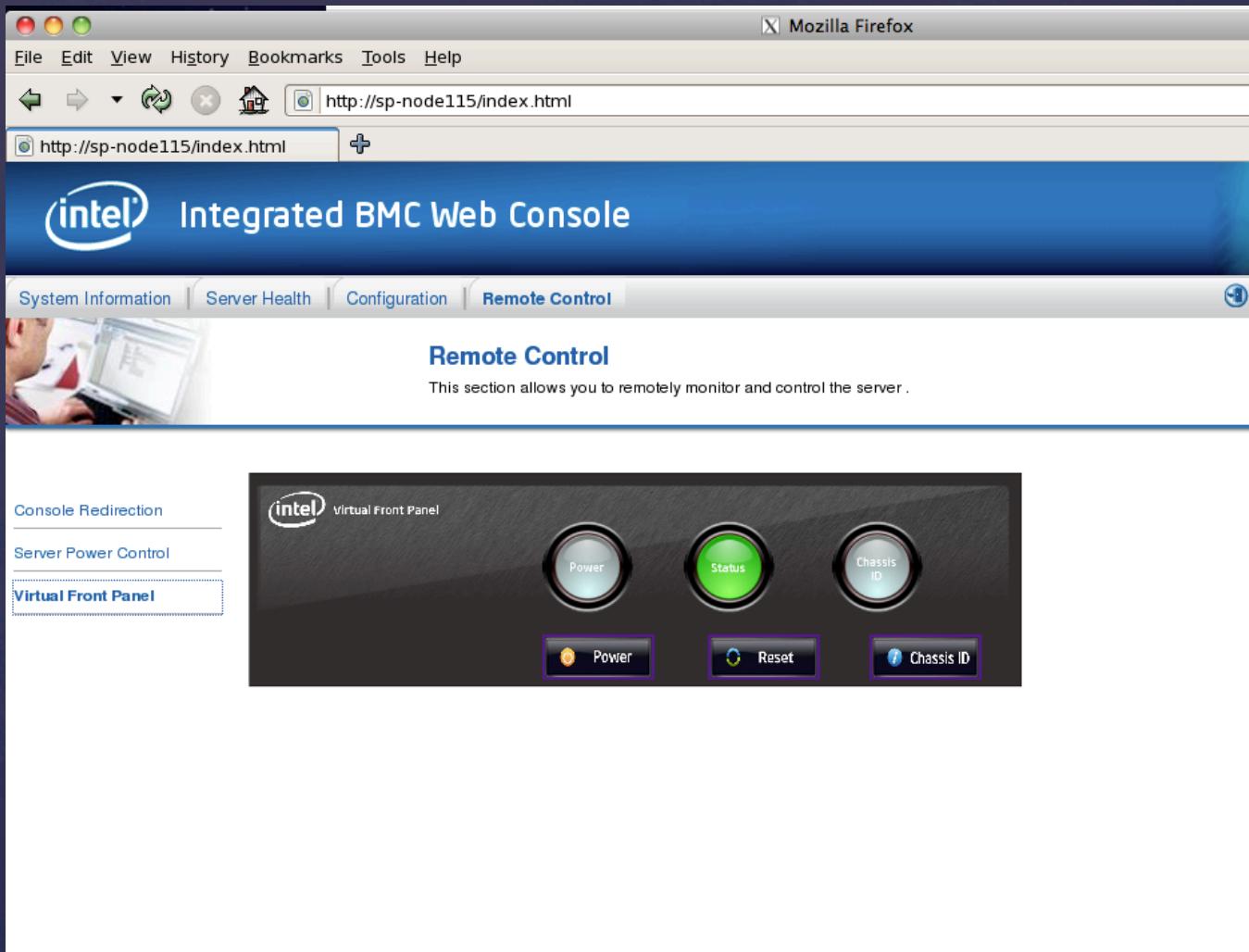
JViewer[172.16.30.65] - 9 fps

Video Keyboard Mouse Options Device VirtualKeyboard Help

```
CentOS release 5.6 (Final)
Kernel 2.6.38-6 on an x86_64

node65 login: NFS: v4 server netapp-b returned a bad sequence-id error!
nfs: server 172.16.1.252 not responding, still trying
nfs: server 172.16.1.252 not responding, still trying
nfs: server 172.16.1.252 OK
nfs: server 172.16.1.252 OK
NFS: v4 server netapp-b returned a bad sequence-id error!
nfs: server 172.16.1.252 not responding, still trying
```





The screenshot shows a Mozilla Firefox browser window displaying the Intel Integrated BMC Web Console. The title bar reads "Mozilla Firefox". The address bar shows the URL "http://sp-node169/index.html". The main content area is titled "Integrated BMC Web Console" with the sub-section "Server Health". A sub-header "Sensor Readings" is visible. On the left, there is a sidebar with links for "System Information", "Server Health" (which is active), "Configuration", and "Remote Control". There is also a "LOGOUT" link. The main content area displays a table of sensor readings with columns for Name, Status, Health, and Reading. The table includes rows for various sensors like Pwr Unit Status, Pwr Unit Redund, IPMI Watchdog, Physical Scty, SMI TimeOut, System Event Log, System Event, Button, VR Watchdog, SSB Therm Trip, IO Mod Presence, SAS Mod Presence, BMC Health, System Airflow, BB Inlet Temp, HSBP Temp, SSB Temp, and BB BMC Temp. The "Health" column uses color coding: red for Critical, green for OK, and yellow for Unknown. The "Reading" column shows numerical values such as 0x0020, 0x0001, 0x0000, and 20 degrees C. At the bottom of the table, there are "Refresh" and "Show Thresholds" buttons, and a text input field for setting auto-refresh intervals.

Name	Status	Health	Reading
Pwr Unit Status	reports there has been a soft power control failure	Critical	0x0020
Pwr Unit Redund	reports full redundancy has been regained	OK	0x0001
IPMI Watchdog	All deasserted	Unknown	Not Available
Physical Scty	All deasserted	OK	0x0000
SMI TimeOut	All deasserted	Unknown	Not Available
System Event Log	All deasserted	OK	0x0000
System Event	All deasserted	OK	0x0000
Button	All deasserted	OK	0x0000
VR Watchdog	reports it has been asserted	Critical	0x0002
SSB Therm Trip	All deasserted	OK	0x0000
IO Mod Presence	All deasserted	Unknown	Not Available
SAS Mod Presence	All deasserted	Unknown	Not Available
BMC Health	All deasserted	Unknown	Not Available
System Airflow	All deasserted	Unknown	Not Available
BB Inlet Temp	Normal	OK	20 degrees C
HSBP Temp	All deasserted	Unknown	Not Available
SSB Temp	All deasserted	Unknown	Not Available
BB BMC Temp	Normal	OK	23 degrees C



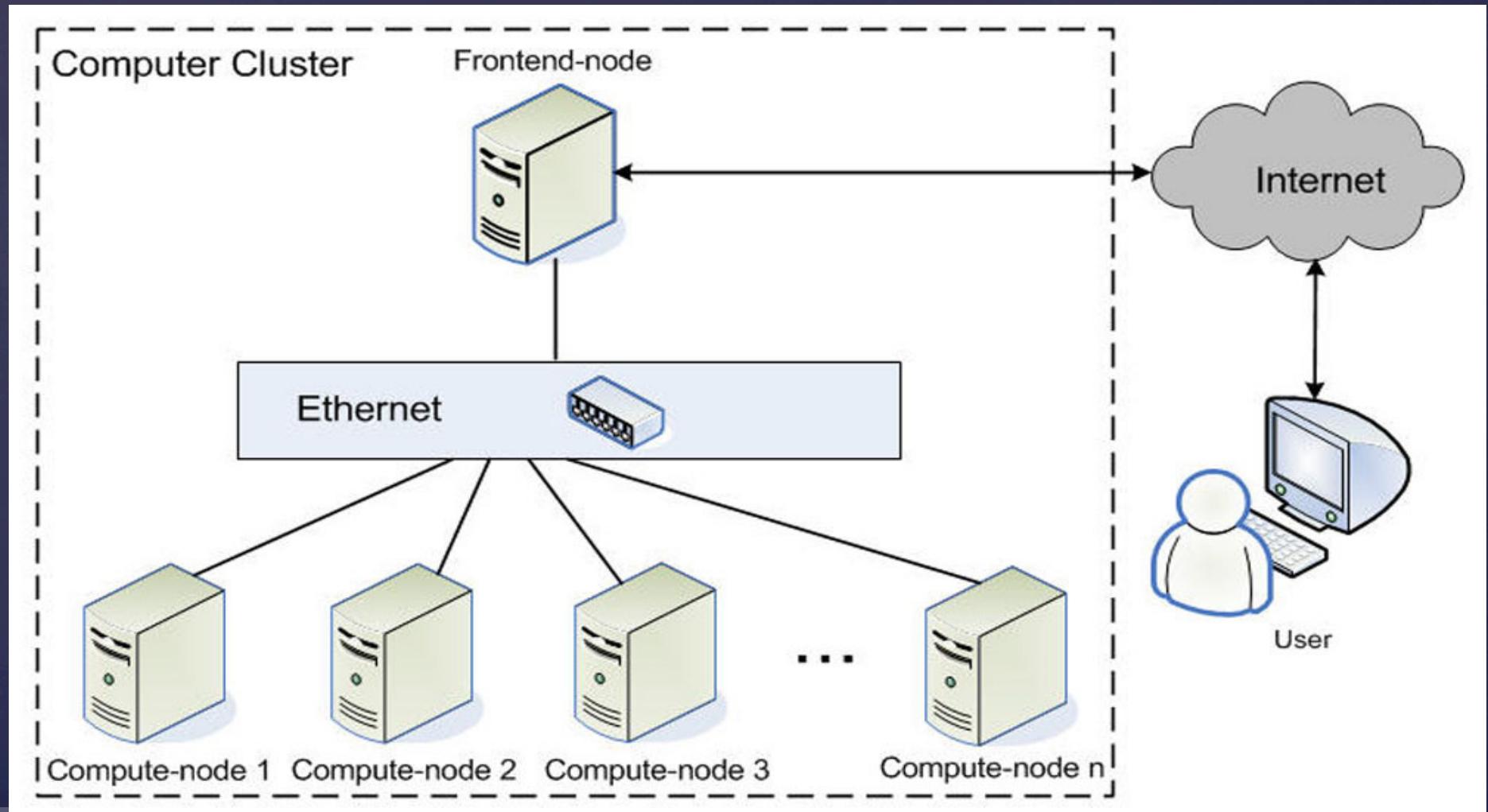
X JViewer[172.16.30.52] - 8 fps

Video Keyboard Mouse Options Device VirtualKeyboard Help

```
57 (xid=0x6273024d)
Jun 14 16:32:20 node52 dhclient[1635]: DHCPREQUEST on eth0 to 172.16.1.254 port
57 (xid=0x6273024d)
Jun 14 17:22:22 node52 kernel: e1000e 0000:0b:00.0: eth0: (PCI Express:2.5GT/s:Width x1) bc:ae:c5:28:21:7a
Jun 14 17:22:22 node52 kernel: e1000e 0000:0b:00.0: eth0: Intel(R) PRO/1000 Network Connection
Jun 14 17:22:22 node52 kernel: e1000e 0000:0b:00.0: eth0: MAC: 3, PHY: 8, PBA No: FFFFFF-0FF
Jun 14 17:22:25 node52 kernel: e1000e: eth0 NIC Link is Up 100 Mbps Full Duplex, Flow Control: None
Jun 14 17:22:25 node52 kernel: e1000e 0000:0b:00.0: eth0: 10/100 speed: disabling TSO
[root@node52 ~]# service network help
Usage: /etc/init.d/network {start|stop|status|restart|reload|force-reload}
[root@node52 ~]# service network reload
Shutting down interface eth0: [OK]
Shutting down interface ib0: [OK]
Shutting down loopback interface: [OK]
Bringing up loopback interface: [OK]
Bringing up interface eth0:
Determining IP information for eth0... failed; no link present. Check cable? [FAILED]
Bringing up interface ib0: [OK]
[root@node52 ~]# _
```



# Cluster structure



# Cluster Structure

- ◆ strictly speaking:
    - + hardware (nodes)
    - + network
    - + software
    - + applications running on it!
  - ◆ broadly speaking:

There is an “**eco-system**” around it:  
power, UPS  
cooling,  
rack space,  
documentation,  
health monitoring,  
backup  
user support
- 
- =cluster



# Cluster view

- ◆ User view (limited structure view)
  - ◆ “Power” user can obtain more performance

Common language:

job, jobID,  
node name,  
“why?”,  
running, producing output,  
job script, working directory

- ◆ HPC team view: all the architectural details, known issues



# Cluster Components, Hardware

- Computers(**Nodes**)
- Network devices
- Storage
- ...



# Nodes

- Computing nodes/Worker nodes
    - homogeneous (initially)
    - heterogeneous: Not in one queue!
  - Master node
  - Specialised nodes (login, gpu,storage)
- ◆ compute node model
    - CPUs
    - Memory!
    - IB port
    - rack mountable
    - Service Processor (IPMI)
  - ◆ master node model
    - differs

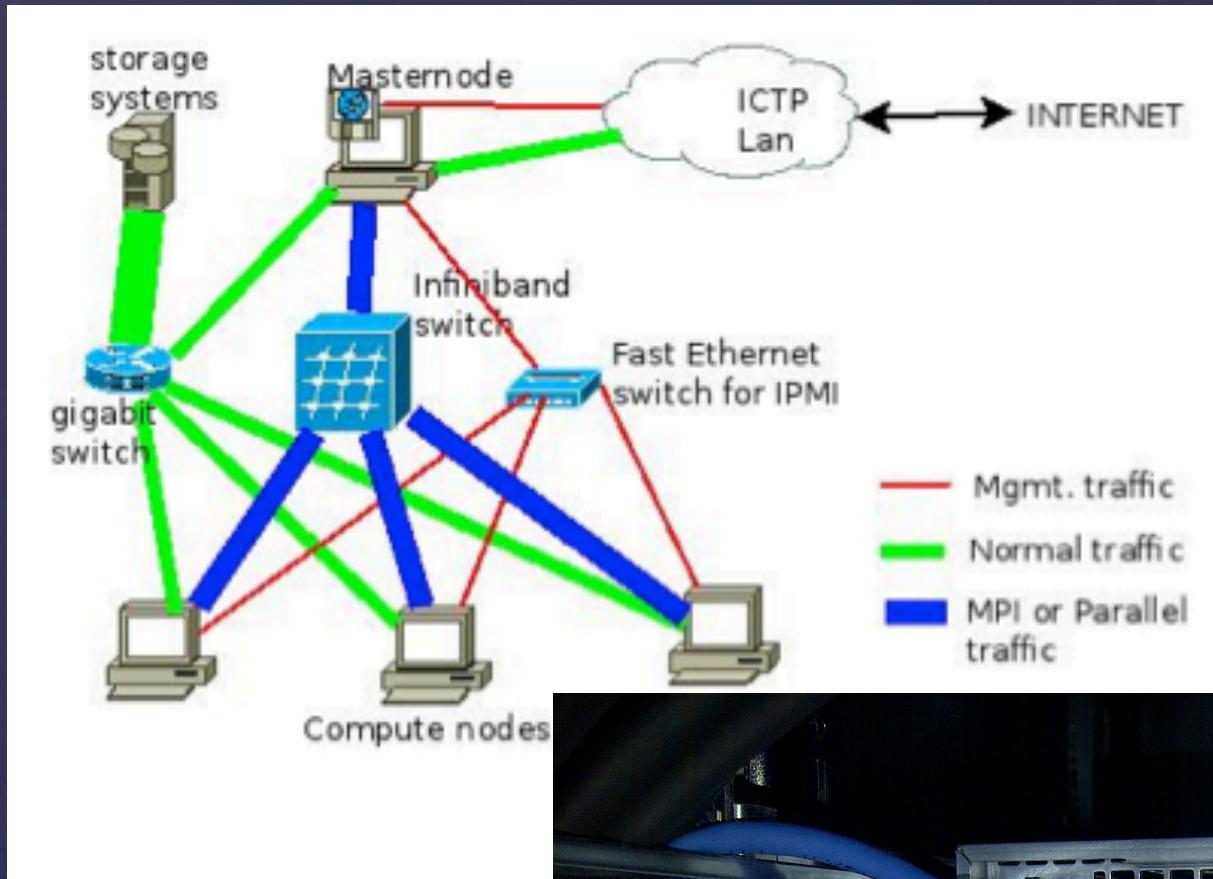


# Cluster Components, Hardware

- Computers(Nodes)
- Network devices
- Storage
- ...



# Network devices

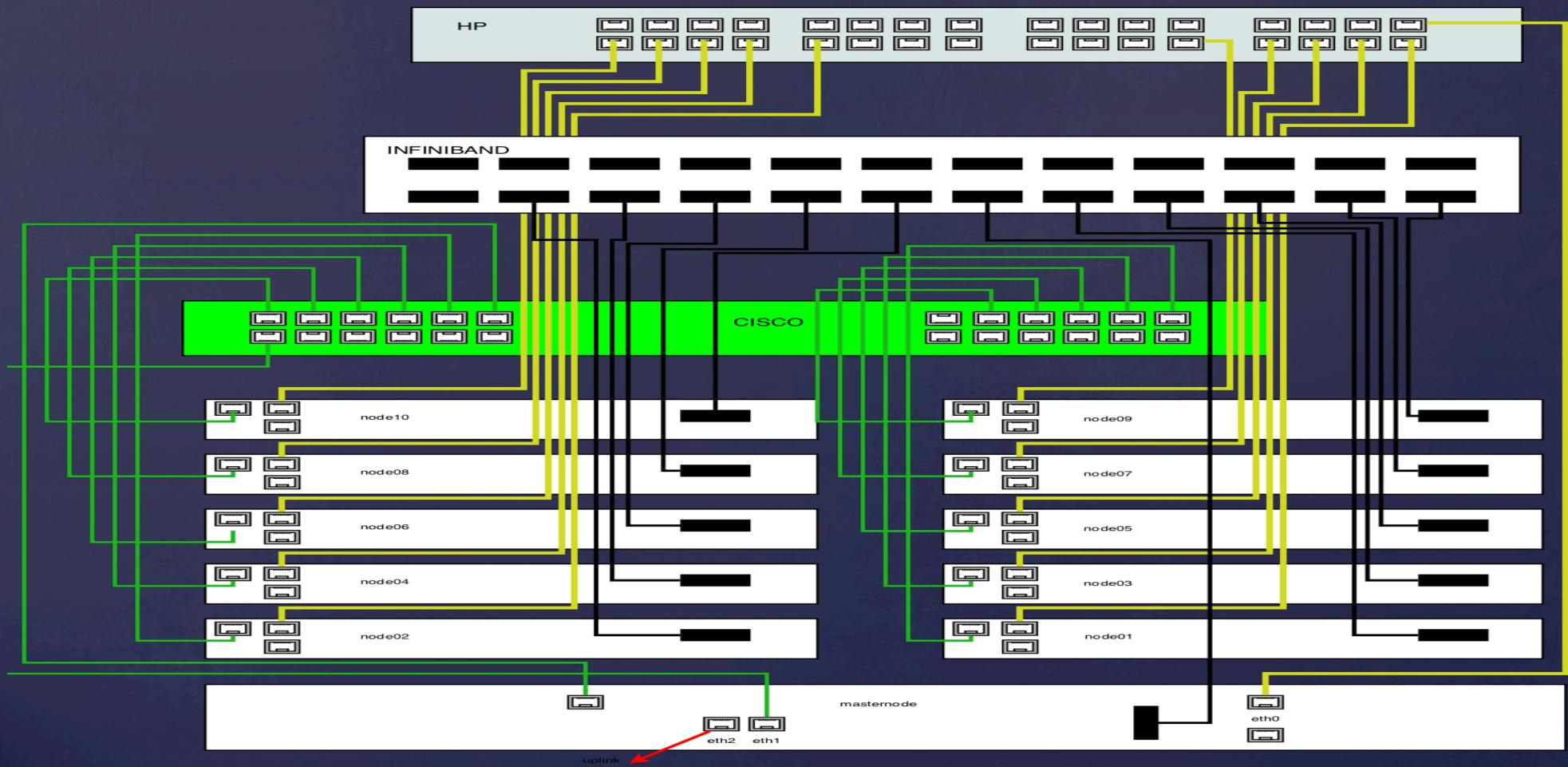


# Network Devices

- ◆ **GigaBit Ethernet** network (GETH)
  - Ssh/ job submission
  - Data I/O traffic (NFS)
  - good bandwidth, NOT low latency for HPC
- ◆ **Infiniband** (or Myrinet ) network
  - Message Passing Interface(MPI) for parallel computation
  - low latency and high bandwidth
- ◆ **Fast Ethernet** network
  - management traffic (IPMI to Service Processors)
- ◆ Switches and **cables**



# Network diagram of Addis HPC



Taken from Antonio's presentation



# Cluster Components, Hardware

- Computers(Nodes)
- Network devices
- Storage: from HDs  
to RAID group
- ...



# Redundant Array of Independent Disks( RAID)

Level	Useable capacity	Data protection
RAID0	$\text{Size}_{\min} * n$	None
RAID1	$\text{Size}_{\min}$	Failure of one single disk
RAID5	$\text{Size}_{\min} * (n - 1)$	Concurrent failure of one single disk
RAID6	$\text{Size}_{\min} * (n - 2)$	Concurrent failure of two disks
RAID1+0	$\text{Size}_{\min} * (n/2)$	Concurrent failure of more than two disks



Taken from Clement's slide

# Storage: Disk array

- a large group of hard disk drives (HDDs)
- RAID: to improve performance and data protection
- NFS: to mount these disks for the cluster.

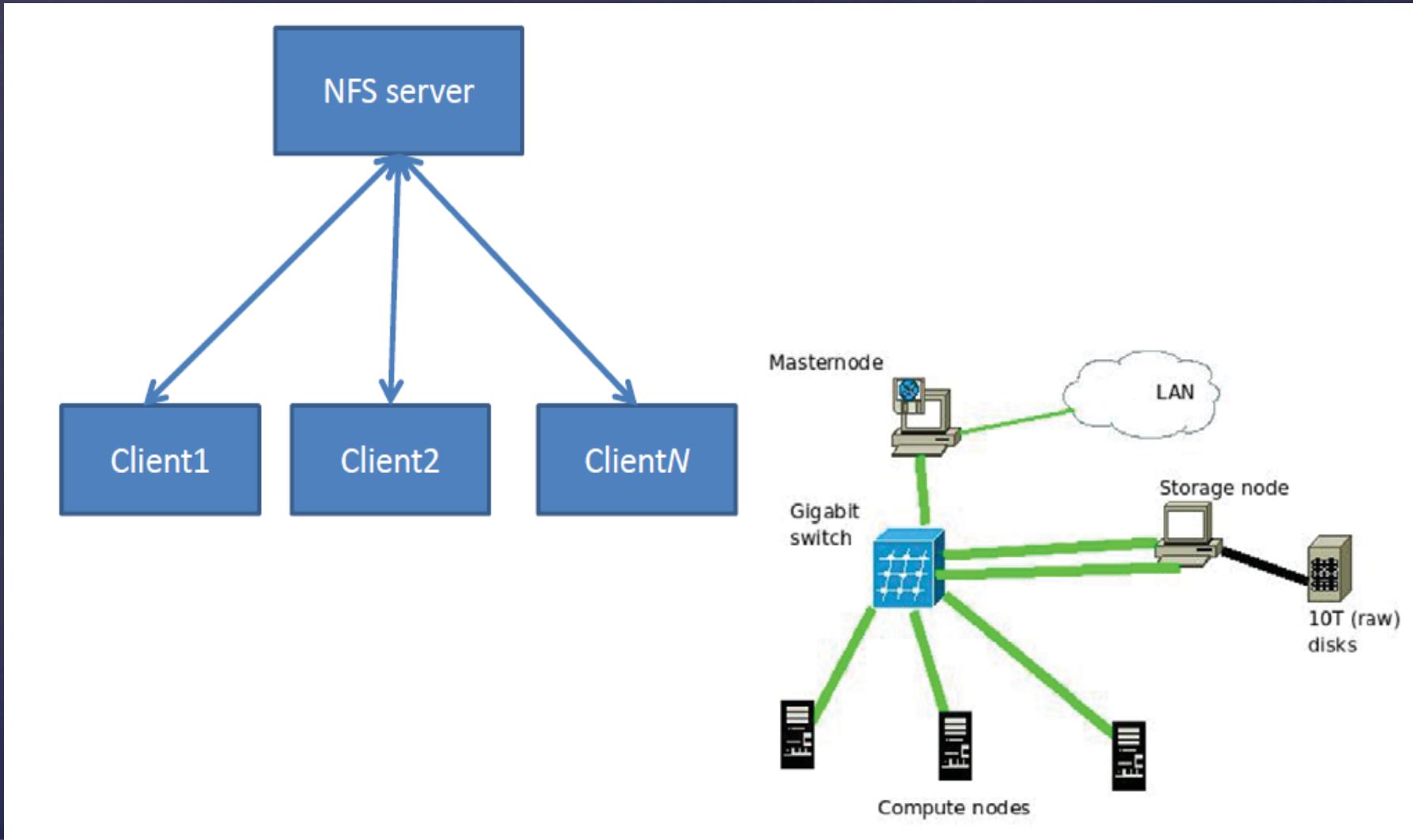


# Storage: Disk array

- HDDs work in parallel
- IOPS: the more the merrier!
- SSD in some cases
- place the data as near as possible to the computing power
- know your app data traffic(size and IOPS)



# NFS architecture



Taken from Clement's slide

# NFS(Network File System)

- NFS Server/NAS exports.
- Clients "mount" into their file system.
- Under the mount point, file access is transparent.  
(for the user)
- /etc/export
  - /distro 10.1.0.0/16(ro,root\_squash)
  - /distro/centos 10.1.0.0/16(ro,root\_squash)
  - /home 10.1.0.0/16(rw,no\_root\_squash)

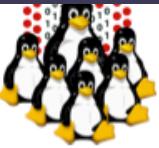


# NFS(Network File System)

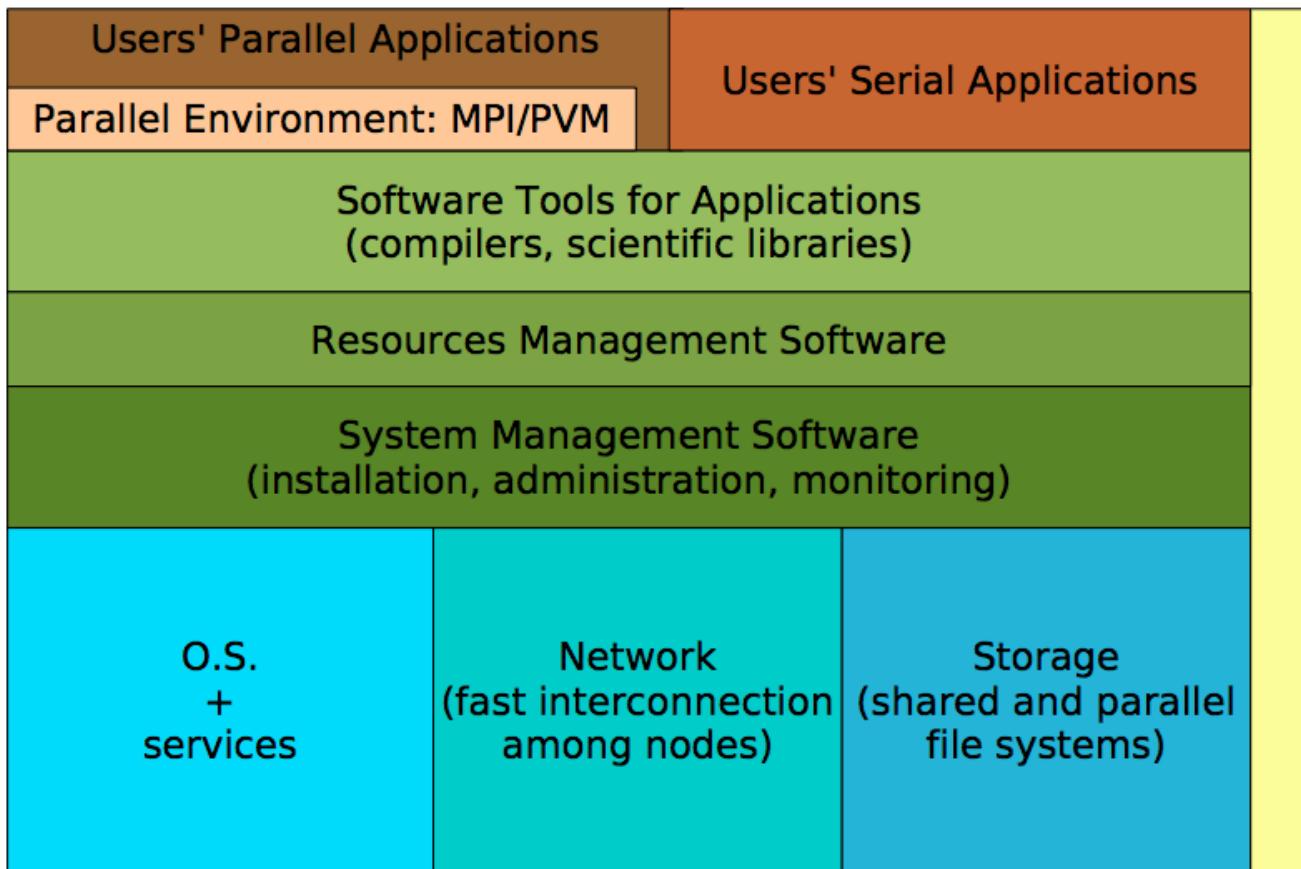
- ◆ Used for:
  - /home RW
  - /opt software RO
  - /distro install repository, RO
  - /scratch RW
  - /projects RW
- ◆ quotas
- ◆ auto-cleaning ?  
(agree!)
- ◆ data-plan
- ◆ /projects visible from Desktops! – out of cluster



# Cluster Components, Software



## HPC SOFTWARE INFRASTRUCTURE Overview



Taken from Moreno's slide

# Cluster Components, Software

- On Master node:
  - OS
  - MPI ?
  - Resource manager
  - Scheduler
  - Installation services
  - Compilations
  - ...
- ◆ what runs on Compute node?
  - OS
  - MPI
  - Resource manager client part
  - User apps
  - ganglia client



# Operating System

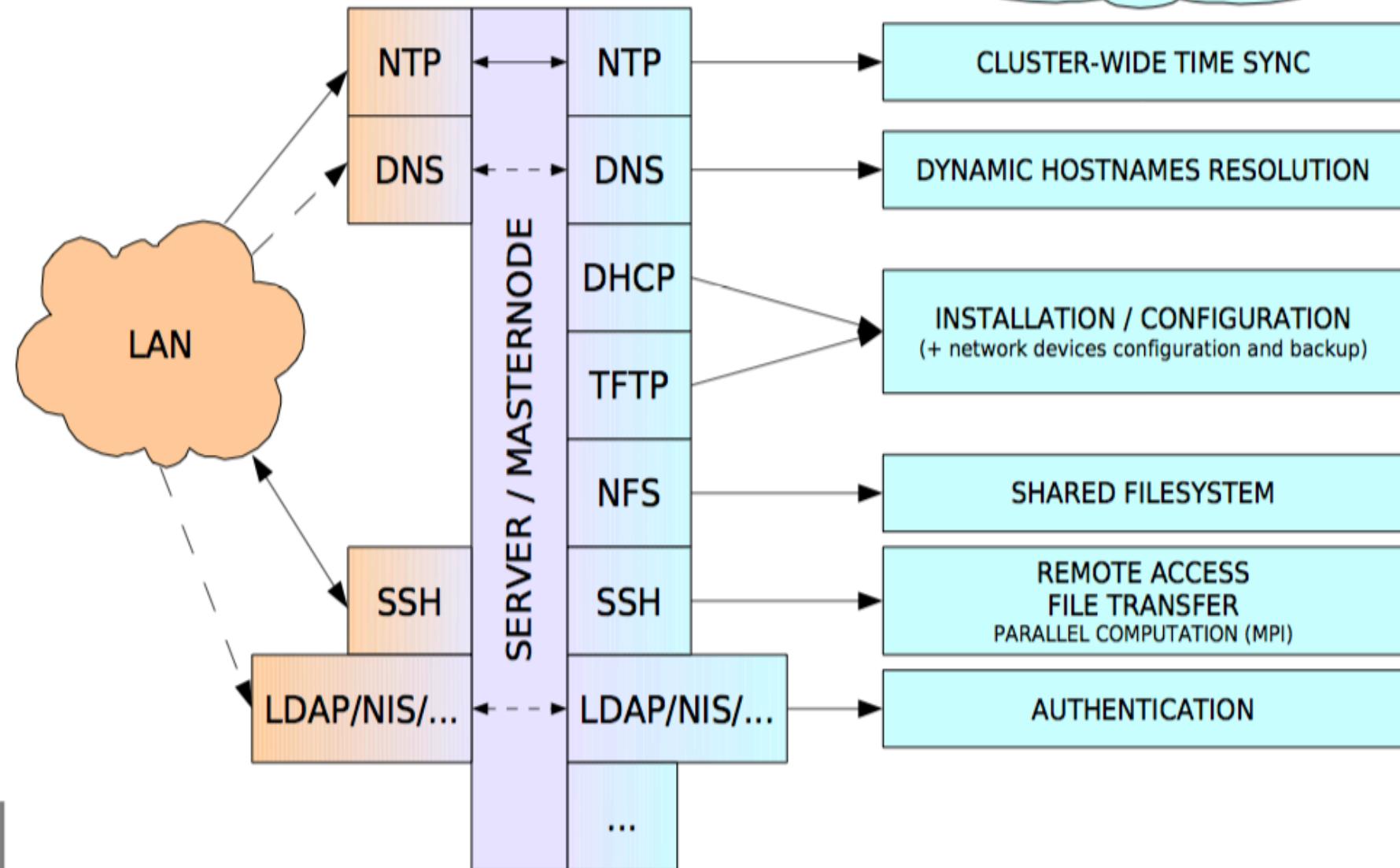
- Linux: Centos, Red Hat, Debian ...
- master node install:  
interactive (manual)
- Computing nodes:  
automated, hands-off mass  
install





# CLUSTER SERVICES

## Cluster services



Taken from Moreno's slide

# Services running on Master node

- DNS(/etc/hosts)
- DHCP
- TFTP
- NTP
- NFS (opt)
- ssh
- Torque
- Maui
- Ganglia
- Syslog
- http (opt)
- ...



# Common view cluster-wide

- ◆ user accounts (/etc/passwd)
- ◆ hostnames (/etc/hosts)
- ◆ mount points (/etc/auto.home)
- ◆ ...
- ◆ Pushed to nodes at regular intervals
- ◆ and at boot time



# /etc/hosts

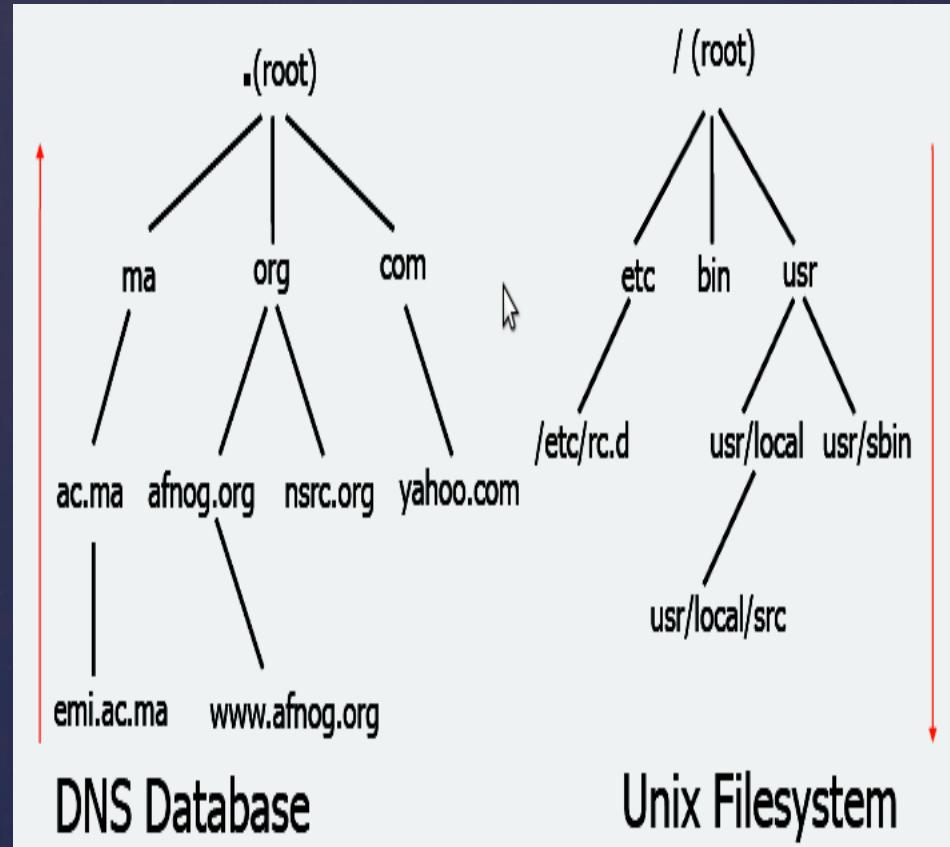
- plain file used to map host names to IP addresses
- In the absence of DNS name server, any network program on your system consults /etc/hosts
- can be rsynced to all nodes

127.0.0.1	localhost.localdomain	localhost
10.1.0.1	master.cluster	master
10.1.1.1	node1.cluster	node1
10.1.1.2	node2.cluster	node2



# DNS (Domain Name System)

- It is a network service that is used to convert domain names to IP addresses (and vice versa)
- DNS is hierarchical
- DNS administration is shared – no single central entity administers all DNS data



From Afnog Workshop

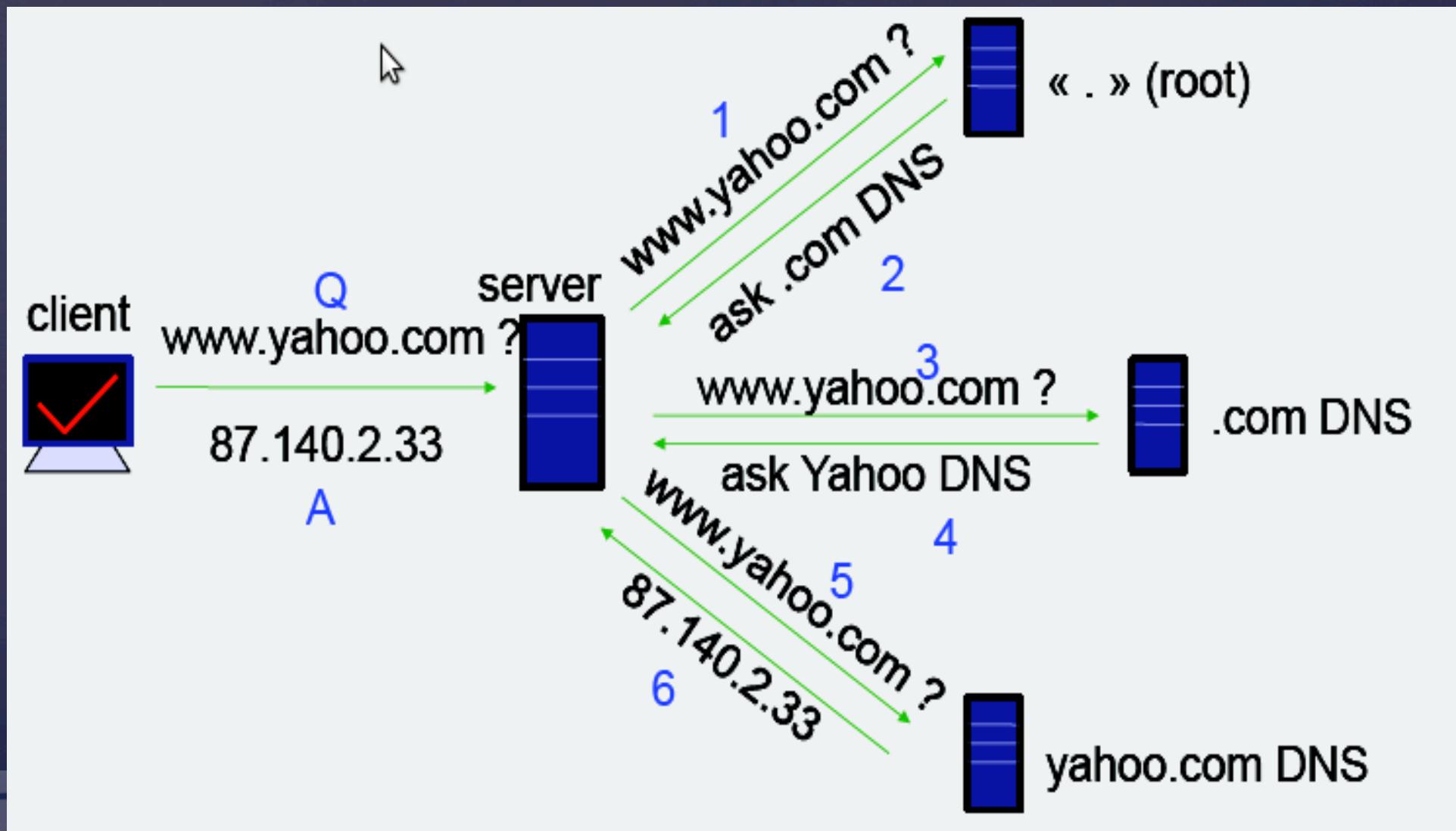


# How does DNS work?

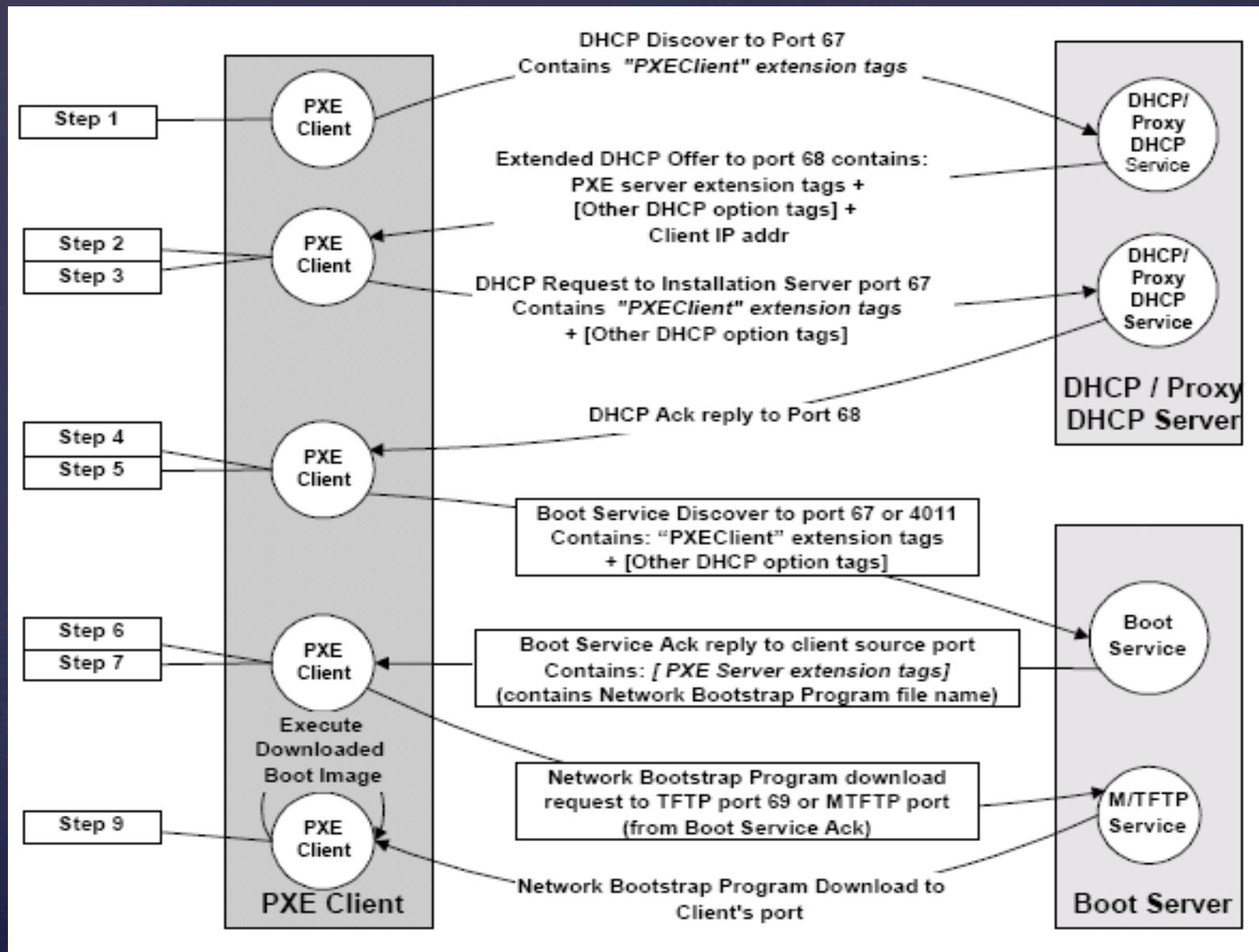
- The client (web browser, mail program, ...) use the OS's resolver to find the IP address - this is called a query
- The server being queried will try to find the answer on behalf of the client
- The server functions recursively, from top (the root) to bottom, until it finds the answer, asking other servers along the way - the server gets referred to other servers



# A DNS query



# Network booting timeline



# DHCP (Dynamic Host Configuration Protocol)

- DHCP allows hosts on a TCP/IP network to request and be assigned IP addresses, and also to discover information about the network to which they are attached
- DHCP provide a mechanism whereby the server can provide the client with information about how to configure its network interface (e.g., subnet mask), and also how the client can access various network services (e.g., DNS, IP routers, and so on).

```
subnet 239.252.197.0 netmask  
255.255.255.0 {  
    range 239.252.197.10  
        239.252.197.250;  
    default-lease-time 86400 max-  
        lease-time 172800;  
    option subnet-mask  
        255.255.255.0;  
    option broadcast-address  
        239.252.197.255;  
    option routers 239.252.197.1;  
    option domain-name-servers  
        239.252.197.2,  
        239.252.197.3;  
    option domain-name "isc.org";  
}
```



# DHCP configuration

- From master node we configure DHCP server which receives clients (computing nodes) requests and replies to them
- DHCP clients sends configuration requests to the server
- For Network booting, PXE acts as a DHCP client

```
& ddns-update-style none;
& ddns-updates off;
& authoritative;
& subnet 10.1.0.0 netmask 255.255.0.0 ^{
&   option domain-name "clusterXY";
&   option domain-name-servers 10.1.0.1;
&   option ntp-servers 10.1.0.1;
&   option subnet-mask 255.255.0.0;
&   option broadcast-address 10.1.255.255;
&   filename "/pxe/pxelinux.0";
&   next-server 10.1.0.1;
}
& host node1 { hardware ethernet ..... ;
fixed-address 10.1.1.1 ; option host-name
"node1" ; }
& host node2 { hardware ethernet ..... ;
fixed-address 10.1.1.2 ; option host-name
"node2" ; }
```



# PXE(Preboot eXecution environment)

- The specification describes a standardized client-server environment that boots a software assembly, retrieved from a network, on PXE-enabled clients.
- On the client side it requires only a **PXE-capable network interface controller (NIC)**, and uses network protocols such as DHCP and TFTP.
- TFTP server has to provide the following files for the clients to initialize the network booting:
  - ✓ pxelinux.0 - is used to load the operating system that is required to execute the assigned preboot work
  - ✓ vmlinuz – is a Linux kernel executable
  - ✓ Initrd.img – initial ramdisk contains various executables and drivers that permit the real root file system to be mounted



# TFTP(Trival File Transfer protocol)

- The protocol is extensively used to support remote booting of diskless devices
- The server is normally started by inetd, but can also run standalone
- TFTP services does not require an account or password on the server system.
- Due to the lack of authentication information, tftpd will allow only publicly readable files

- /etc/xinetd.d/tftp file

```
service tftp
{
    disable          = no
    socket_type     = dgram
    protocol        = udp
    wait            = yes
    user             = root
    server          = /usr/sbin/in.tftpd
    server_args     = -s /tftpboot -vvv
    per_source       = 11
    cps              = 100 2
    flags            = IPv4
}
```



# Network booting...

- One other thing that should be transferred through tftp is PXE configuration file that defines the menu displayed to the target host
- These configuration files can be stored /tftpboot/pxe/pxelinux.cfg directory in one of the following form, for a given ip and mac
  - ✓ /ftproot/pxe/pxelinux.cfg/01-88-99-aa-bb-cc-dd
    - ✗ If the mac address is 01-88-99-aa-bb-cc-dd
    - ✓ /tftpboot/pxe/pxelinux.cfg/C000025B
  - ✗ hexadecimal equivalent of ip address 192.0.2.91
  - ✓ /tftpboot/pxe/pxelinux.cfg/default
- Usually hexadecimal equivalent of pxe configuration file is used to install a node and the “default” is used when one wants to boot the node from the local hard disk

- /tftpboot/pxe/pxelinux.cfg/default

```
prompt 1
timeout 100
default local
label local
LOCALBOOT 0
label install
kernel vmlinuz
append initrd=initrd.img network
    ip=dhcp \
    ksdevice=eth0    ks=nfs:10.1.0.1:/
                    distro/ks/ks.cfg \
    load_ramdisk=1
    prompt_ramdisk=0 \
    ramdisk_size=16384 \
    vga=normal selinux=0
```



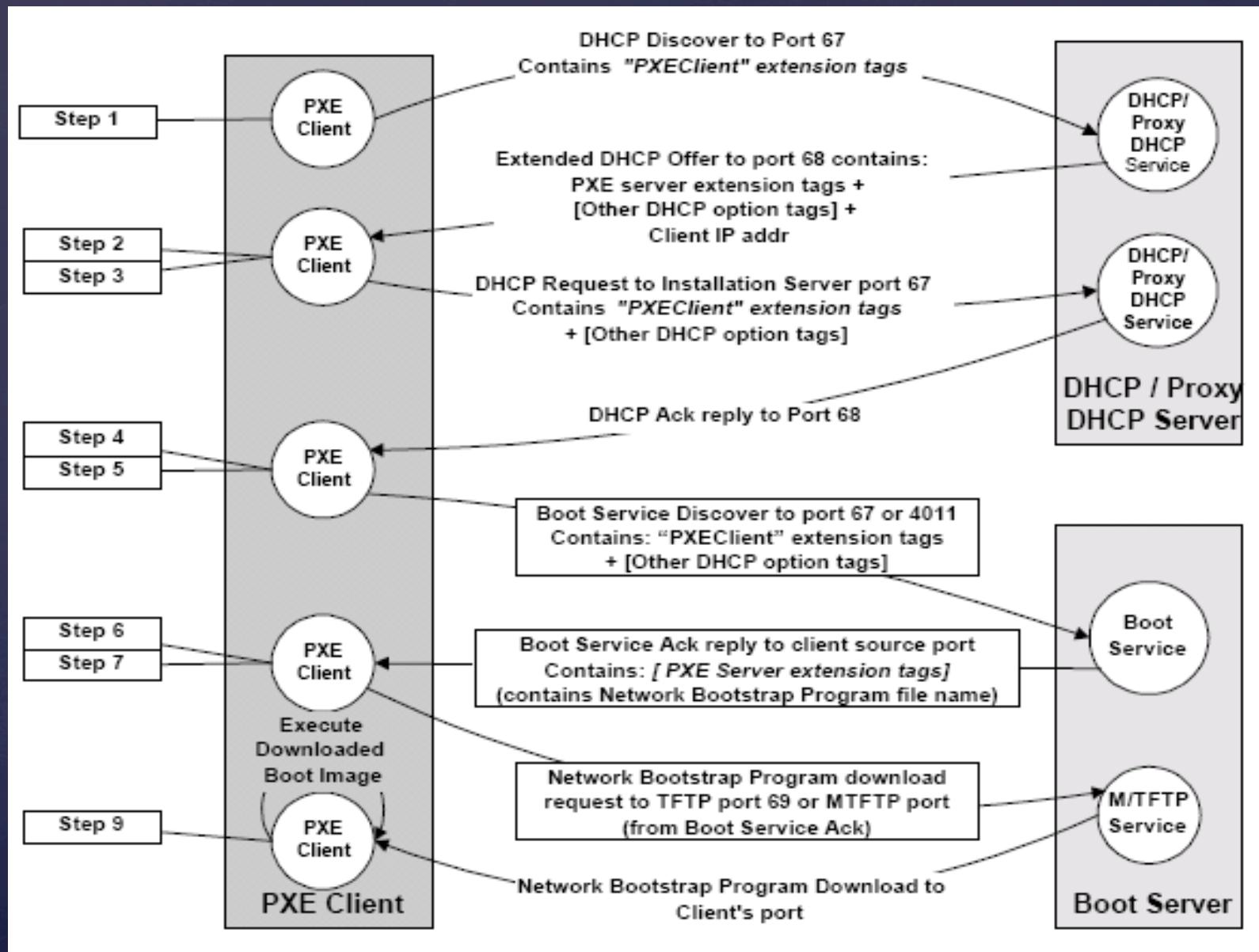
# Network booting “plan”

```
/  
`-- tftpboot/  
    '-- pxe/  
        '-- vmlinuz  
        '-- initrd.img  
        '-- memtest  
        '-- pxelinux.0  
        '-- pxelinux.cfg/  
            '-- 0A0A0101  
            '-- bootmsg.txt  
            '-- default -> default.local  
            '-- default.install  
            '-- default.local
```

Taken from M. Baricevic, 2013 slide



# Network booting timeline



# Kickstart Installations

- automated installation (all answers provided, no keyboard input needed)
- Install media: local CD-ROM, or via **NFS**, **FTP**, or **HTTP**
- The kickstart file is a simple text file, containing a set of instruction on how to install the OS + post install steps
- Kickstart Configurator application or by writing it from scratch
- <http://fedoraproject.org/wiki/Anaconda/Kickstart>



# Sample kickstart file

- Kickstart files allowing us to **configure** network, system configurations, HD partitioning and package selections
- In the pre-installations section one can choose **hardware setup** and configurations
- In the **post-installation** section one can include customizations and additional configurations.
- It is also possible to add lines of instructions that can stop the automated installation

```
#platform=x86, AMD64, or Intel EM64T
# System authorization information
auth --useshadow --enablemd5
# System bootloader configuration
bootloader --location=mbr
# Clear the Master Boot Record
zerombr
# Partition clearing information
clearpart --all --initlabel
# Use text mode install
text
# Firewall configuration
firewall --disabled
# Run the Setup Agent on first boot
firstboot --disable
# System keyboard
keyboard us
# System language
lang en_US
```



# Cluster Management tools

- Cluster wide commands
- Password-less environment
- Cluster wide file distribution and gathering
- Appropriate access privilege



# C3 tools

- Cluster Command Control (C3) tools are a suite of cluster tools that are useful for both administration and application support
- It includes tools for cluster-wide command execution, file distribution and gathering, process termination, remote shutdown and restart, and system image updates
- Example: cexec, cget, ckill , cpush and others



# Code execution Environment, modules

- It provides dynamic modification of a user's environment
- It allows a group of related environment variables to be made or removed dynamically
- It enables us to avoid errors that can be caused from changing \$PATH environment
- Some of frequently used module commands:
  - module avail
  - module list
  - module load/unload <package>
  - module purge



# Other Important Software

- Open MPI is a Message Passing Interface (MPI) library project combining technologies and resources from several other projects (FT-MPI, LA-MPI, LAM/MPI, and PACX-MPI).
- Compilers(GNU, Portland Group and Intel)
- Scientific libraries
- Resource manager and Scheduler  
PBS/Torque and Maui scheduler
- Ganglia for metrics collection ad presentation



# Health monitoring

- ◆ Internal monitoring:
  - ◆ services on master node up and running
  - ◆ node health (before the job)
- ◆ External monitoring:  
nagios

Service Status Details For Host 'argo-master'						
Host ▾	Service ▾	Status ▾	Last Check ▾	Duration ▾	Attempt ▾	Status Information
argo-master	2-DISKS-HEALTH	OK	09-29-2015 15:50:50	17d 6h 6m 52s	1/3	OK sda=PASSED sdb=PASSED
	HTTP	OK	09-29-2015 15:56:02	17d 6h 5m 13s	1/3	HTTP OK: HTTP/1.1 200 OK - 4104 bytes in 0.149 second response time
	HTTP-GANGLIA	OK	09-29-2015 15:55:02	9d 4h 5m 33s	1/3	HTTP OK: HTTP/1.1 200 OK - 133368 bytes in 2.986 second response time
	PBSNODES-DOWN	OK	09-29-2015 15:55:00	0d 1h 32m 29s	1/3	OK: pbsnodes down:0
	SSH-Check-load-8	OK	09-29-2015 15:54:16	10d 21h 53m 19s	1/3	OK: load (0.30) is below threshold (8/12) - load=0.30
	SSH_Disk_Free	OK	09-29-2015 15:53:35	44d 21h 4m 0s	1/3	OK: All Filesystems are below threshold (85/90%) [/=/32% /var=57% /boot=32% ]

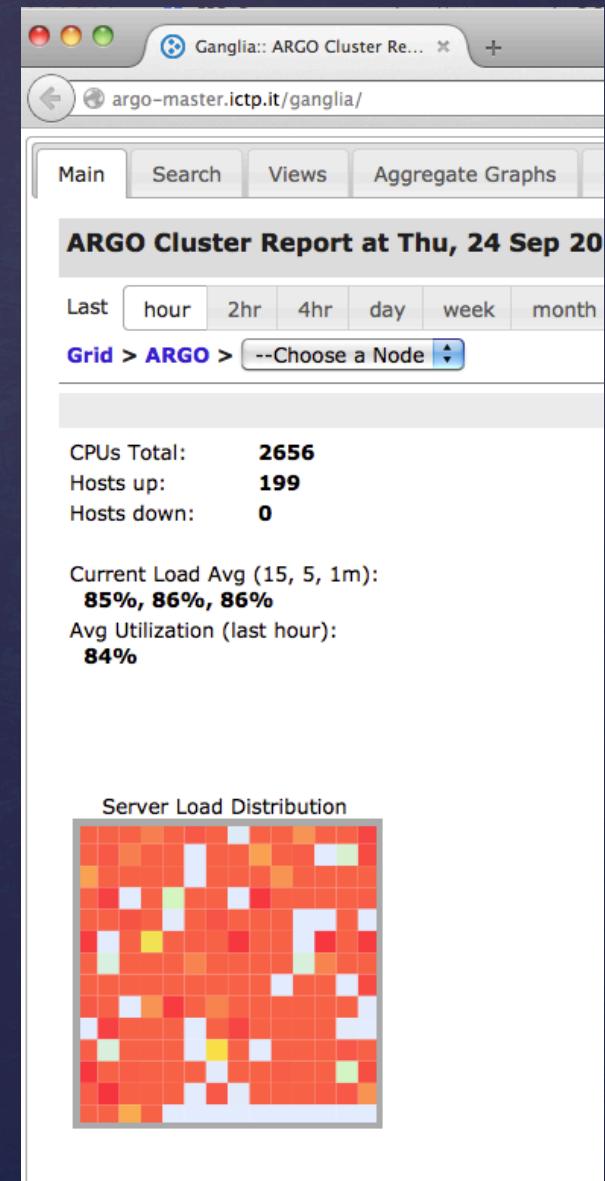
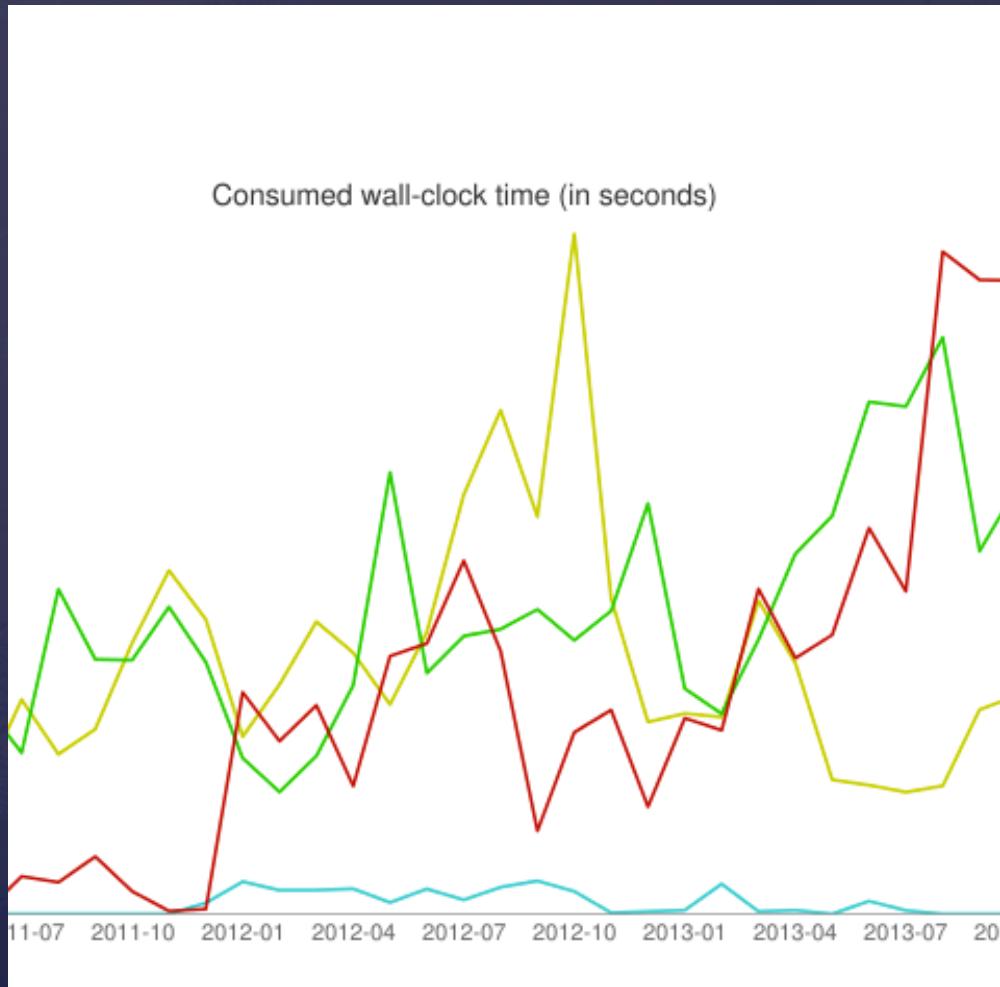


# Usage monitoring

- ◆ User view: cluster free ?
- ◆ Admin view: clusterutilised?



# Usage monitoring



# User support

- ◆ documentation
- ◆ presentations
- ◆ mailing list
- ◆ helpdesk/trouble ticket
- ◆ HPC team of people
  - ◆ general Linux,
  - ◆ shell scripting,
  - ◆ network management,
  - ◆ storage management,
  - ◆ batch system configuration,
  - ◆ compilation of scientific software



# Conclusion:

commodity cluster can be built and operated  
with success

The problem/Function:

“make parallel/HPC **apps run**  
well for your users”

Cluster Structure:

... fairly complex ...

- ◆ But it CAN be done!

The sum: hw + sw + motivated team

will make your **cluster run**



# Thank You!

# Questions ?

