

Data handling/ computing equipment – the heart of *the HPC* *cluster*

{ buy CPU cores, but also plan for data (space and performance-wise)

Maria Verina & Clement Onime

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

The Information and Communication Technology
Section (ICTS)



HPC Data Center/Server room



Data handling/computing equipment

- ✂ servers (nodes)
- ✂ storage,
- ✂ network switch (switches^s!)

ICTP HPC cluster *Argo*

- ◆ Servers

 - 149 nodes

 - 2736 cores (Intel x86_64 +GPU), heterogeneous (The story of growth.)

 - Total RAM: 7.5TB

- ◆ Storage: ~300TB NFS, 10Gbps, dedicated + common /home, /opt

- ◆ Switches:

 - 1 Gigabit Ethernet private cluster network

 - Infiniband: 40 Gbps QDR + Omnipath: 100 Gbps (*low latency* for MPI)

 - Management network at 100 Mbps

- ◆ racks: 5





Nodes, nodes, nodes

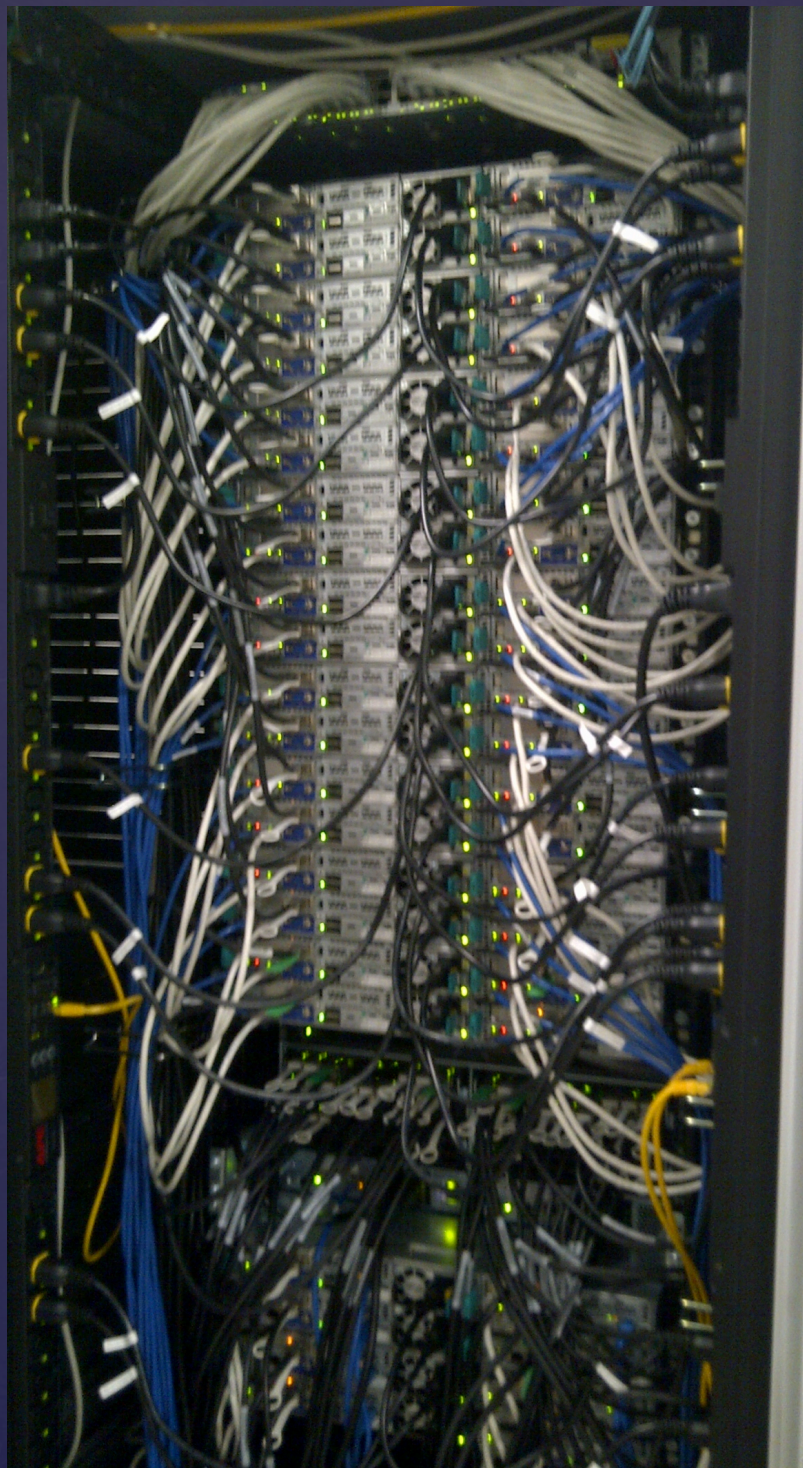


InfiniBand switch



Worker nodes, 1RU twin

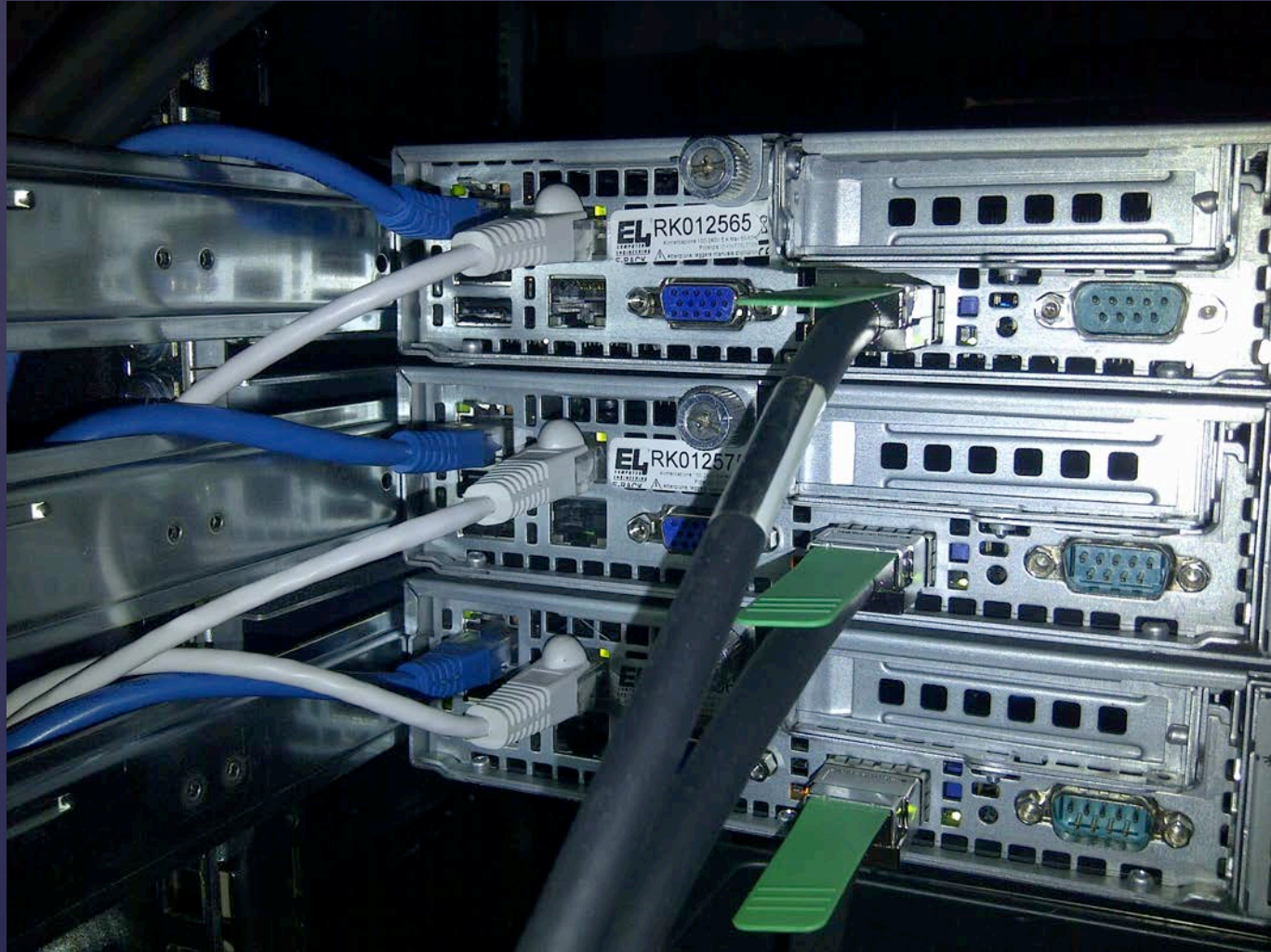




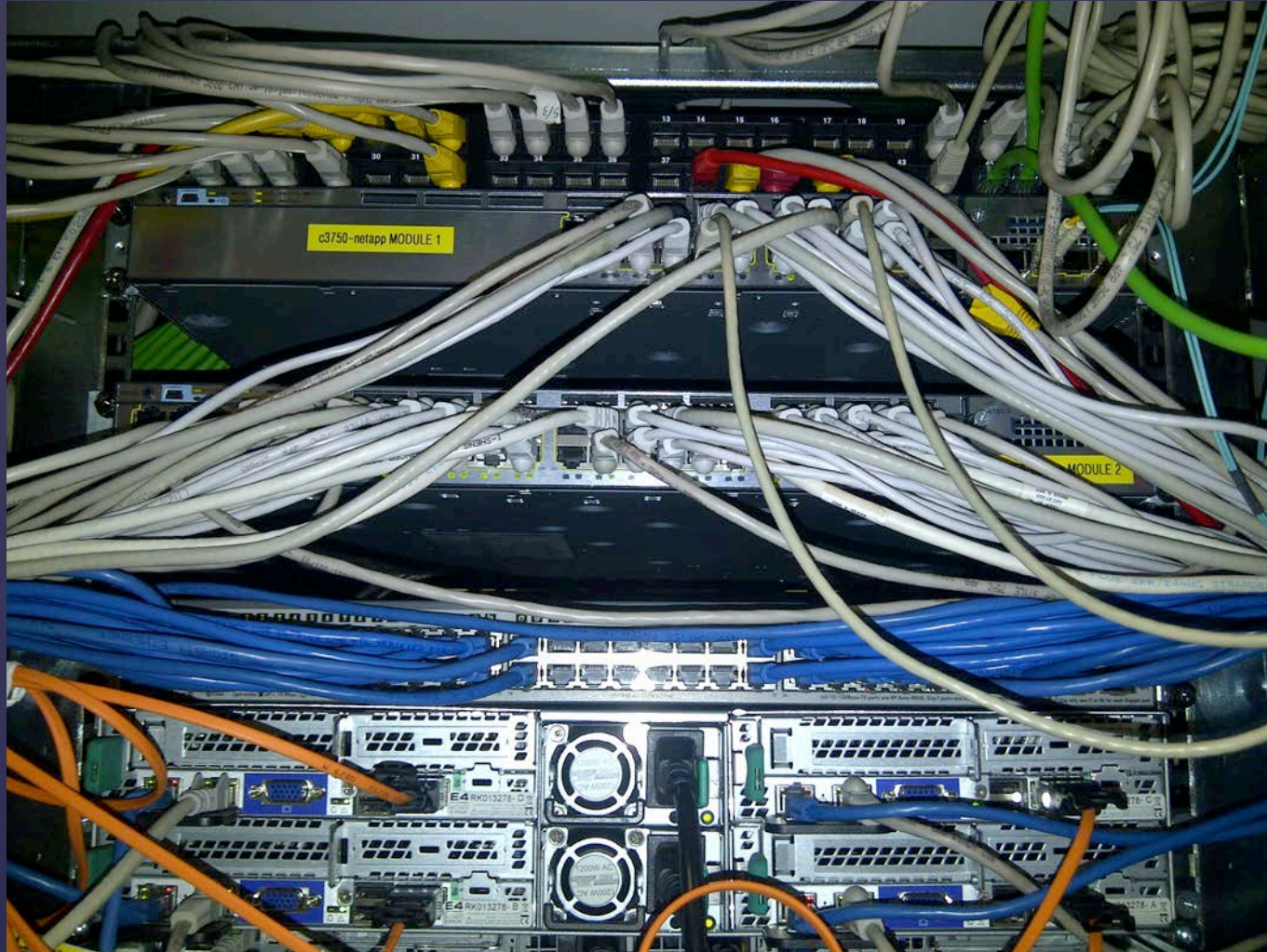
Master node network ports



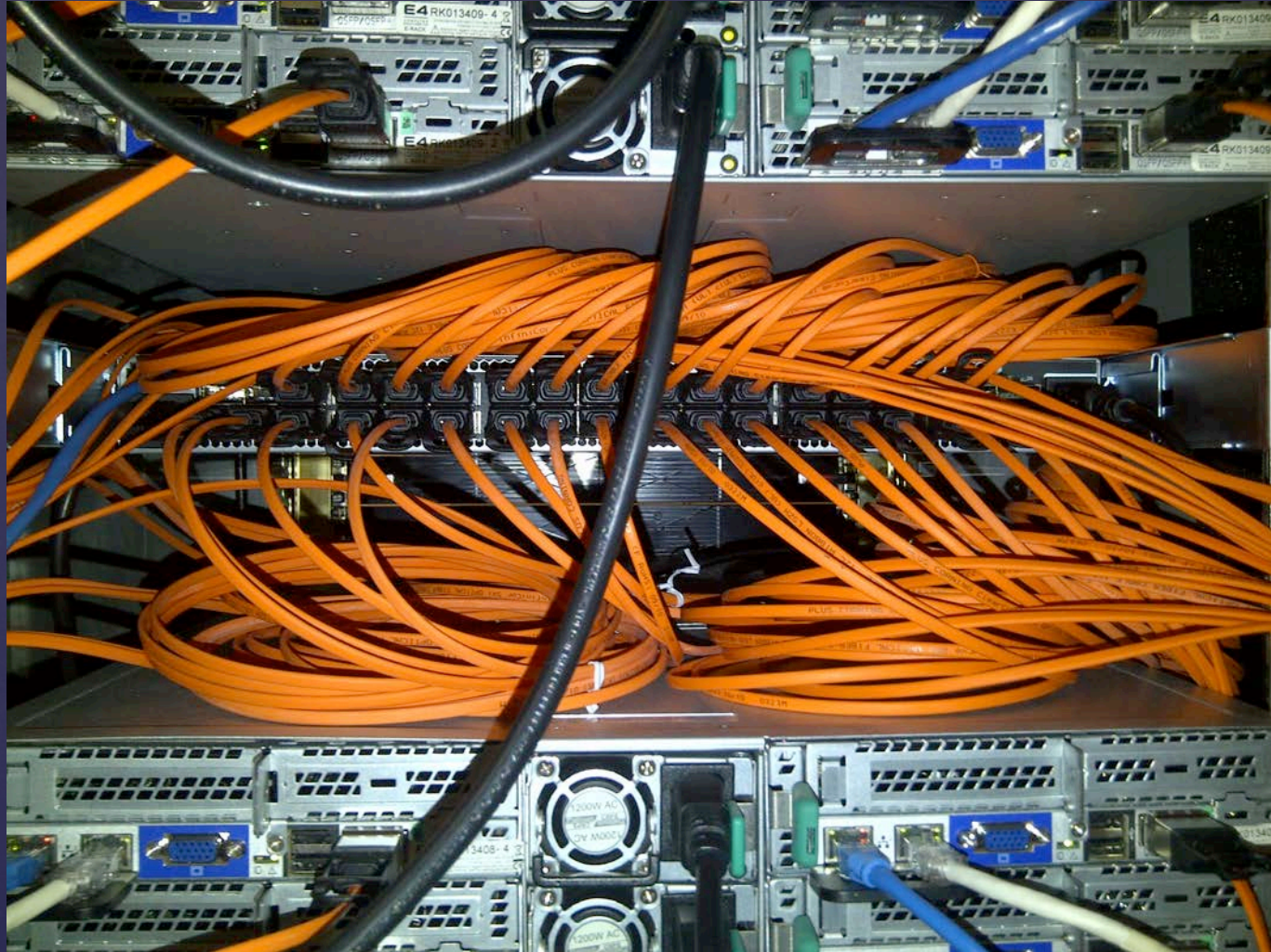
Compute node ports



Ethernet switch



Infiniband switch (back)



HPC equipment

- Servers (nodes)
- Network devices
- Storage
- ...

Nodes, nodes, nodes

- Computing nodes/**Worker nodes**

homogeneous (initially)

heterogeneous: Not in one queue!

- Master node

- **Specialised** nodes (login, gpu, storage)

- ◆ compute node model

CPU's

Memory!

IB port

rack mountable

Service Processor (IPMI)

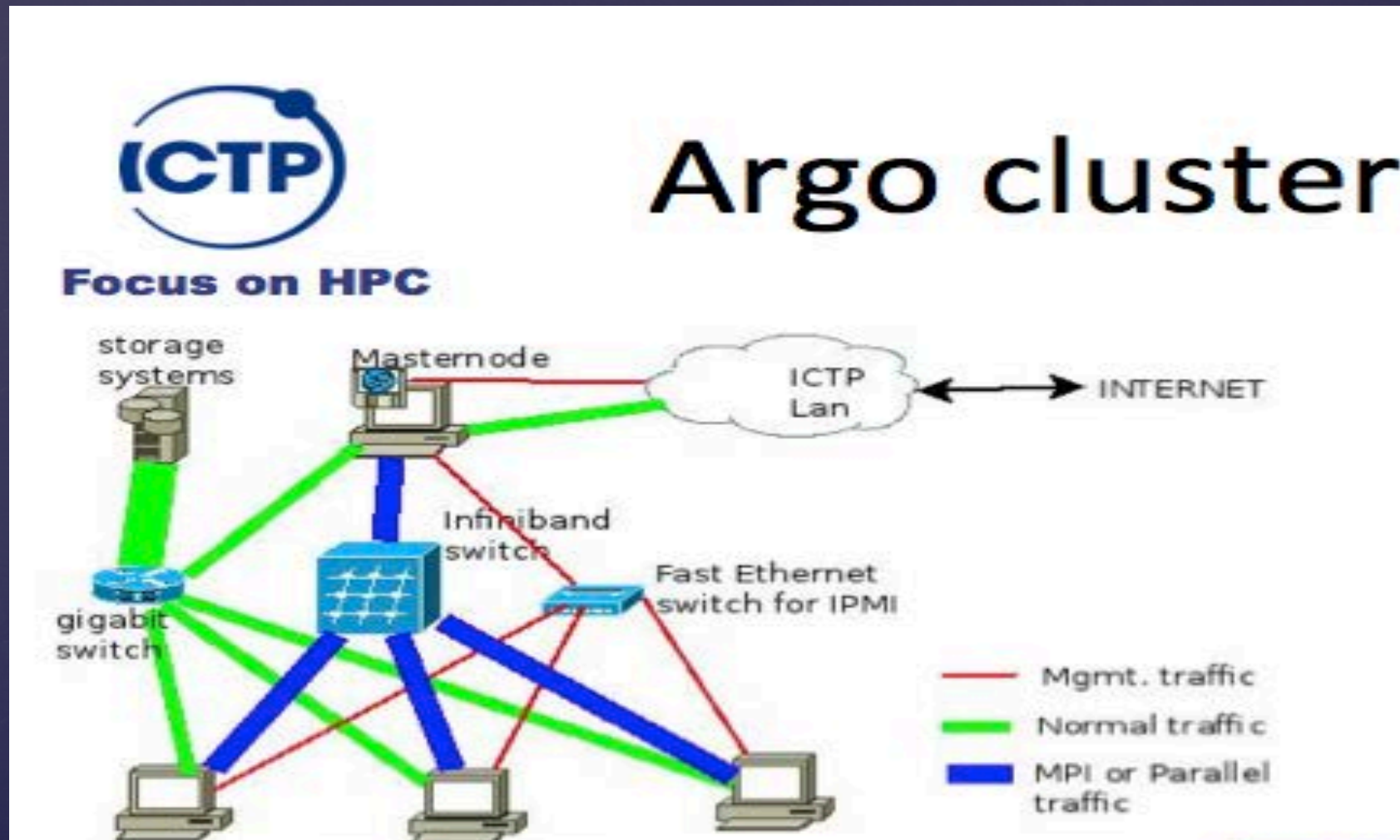
- ◆ **master** node model

differs

HPC equipment

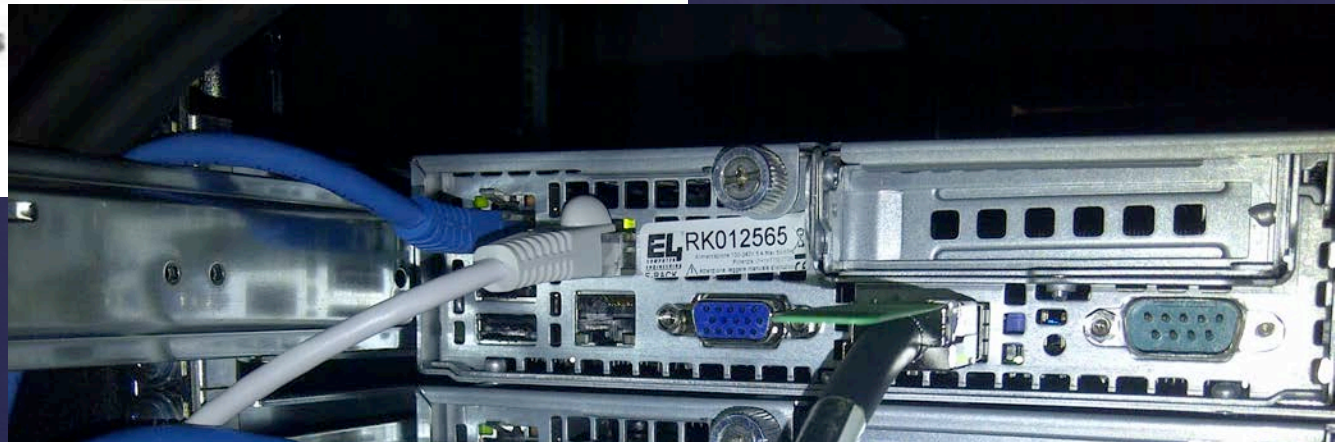
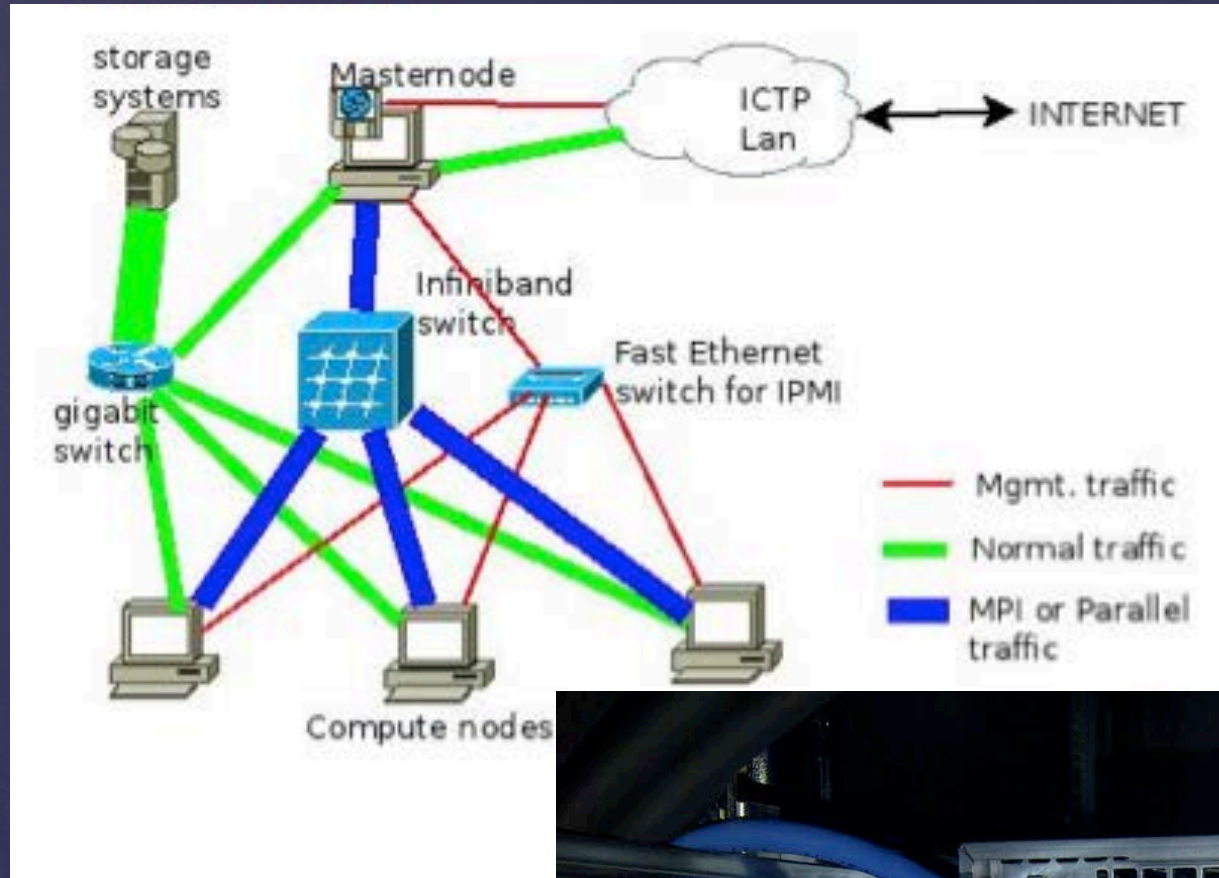
- Computers(Nodes)
- Network switches
- Storage
- ...

Traffic type vs Switch kind



credit: Clement O.

Network devices



Swithes, switches, switches

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

Cables, cables, cables

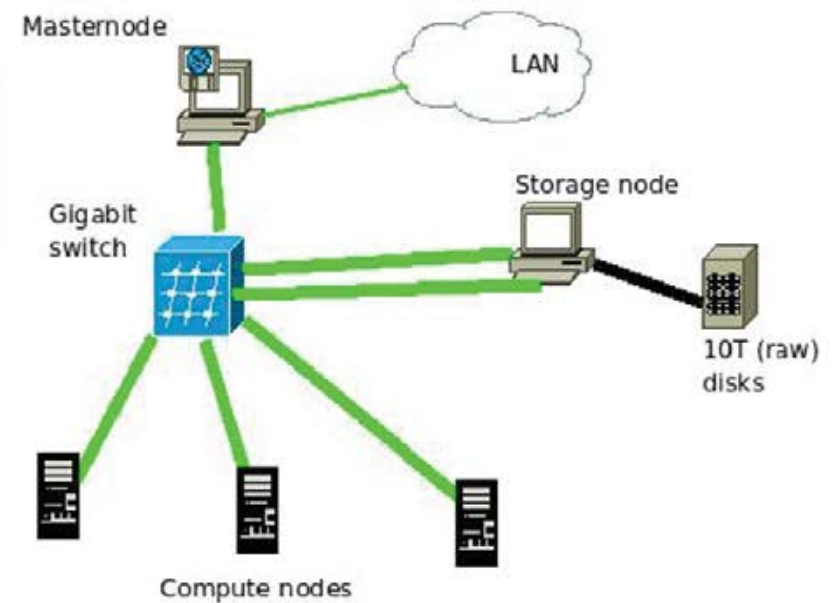
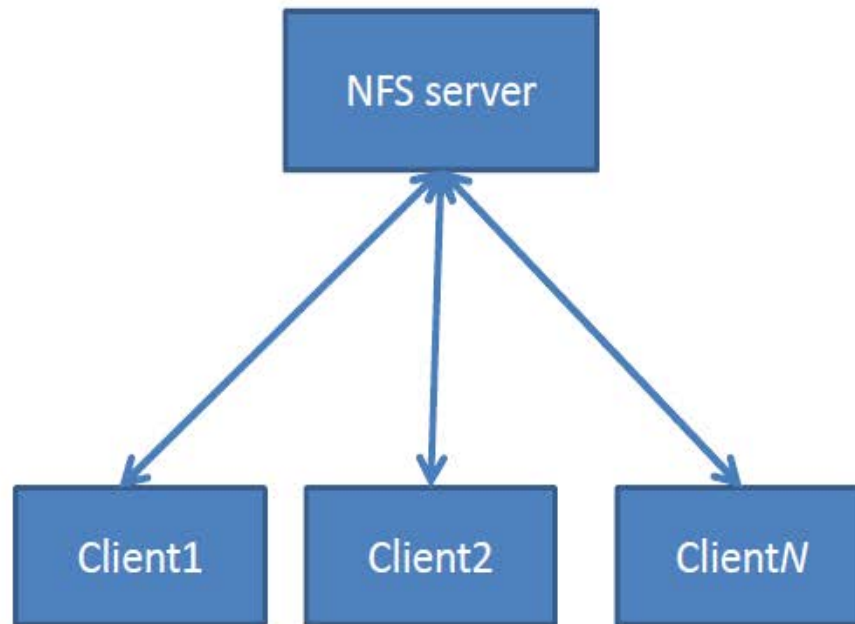
- ◆ use quality cables (heat on the back!)
- ◆ Use order, but no labels (heat!)
- ◆ Route cables away from (hot) air flow
- ◆ Use colour scheme (e.g. management cables are blu, uplinks are red)
- ◆ Use etherchannels where appropriate (redundancy!)
- ◆ use fat “pipes” to data (10G storage side)
- ◆ for power cables use firm cables (“locked”)

HPC equipment

- Computers(Nodes)
- Network devices
- Data Storage: HDs, RAID, NFS, Parallel File System

...data: centralization or distributed

NFS architecture



Credit: Clement O.

NFS-Network File System

- Centralised:
- NFS Server/NAS does “the export”.
- Clients do the “mount” of the export into a point inside their file system.
- Under the mount point, file access is transparent.
(for the user)

NFS

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

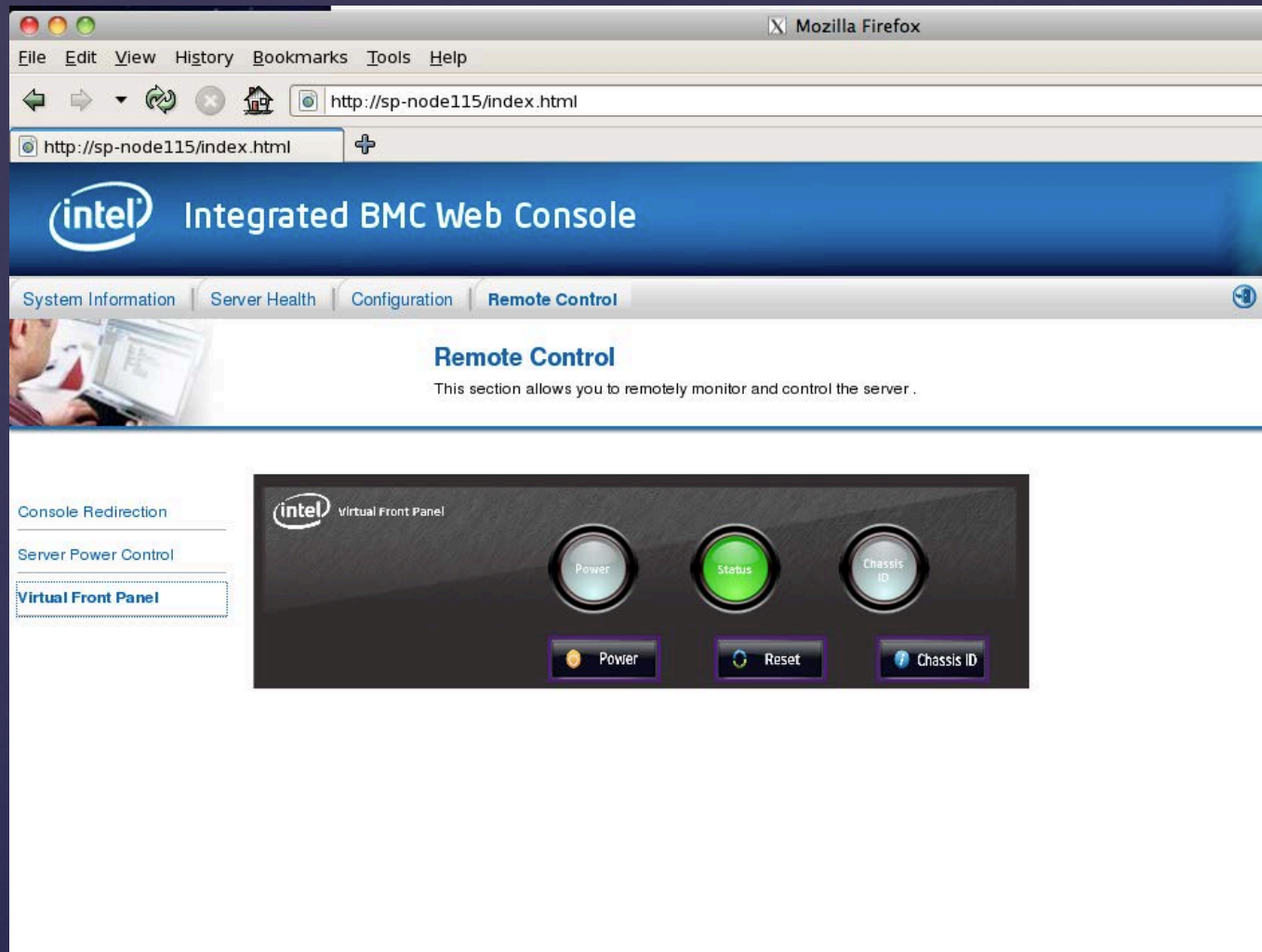
Parallel File systems

- ◆ Users would love them!
- ◆ they cost money and effort in setup and maintenance

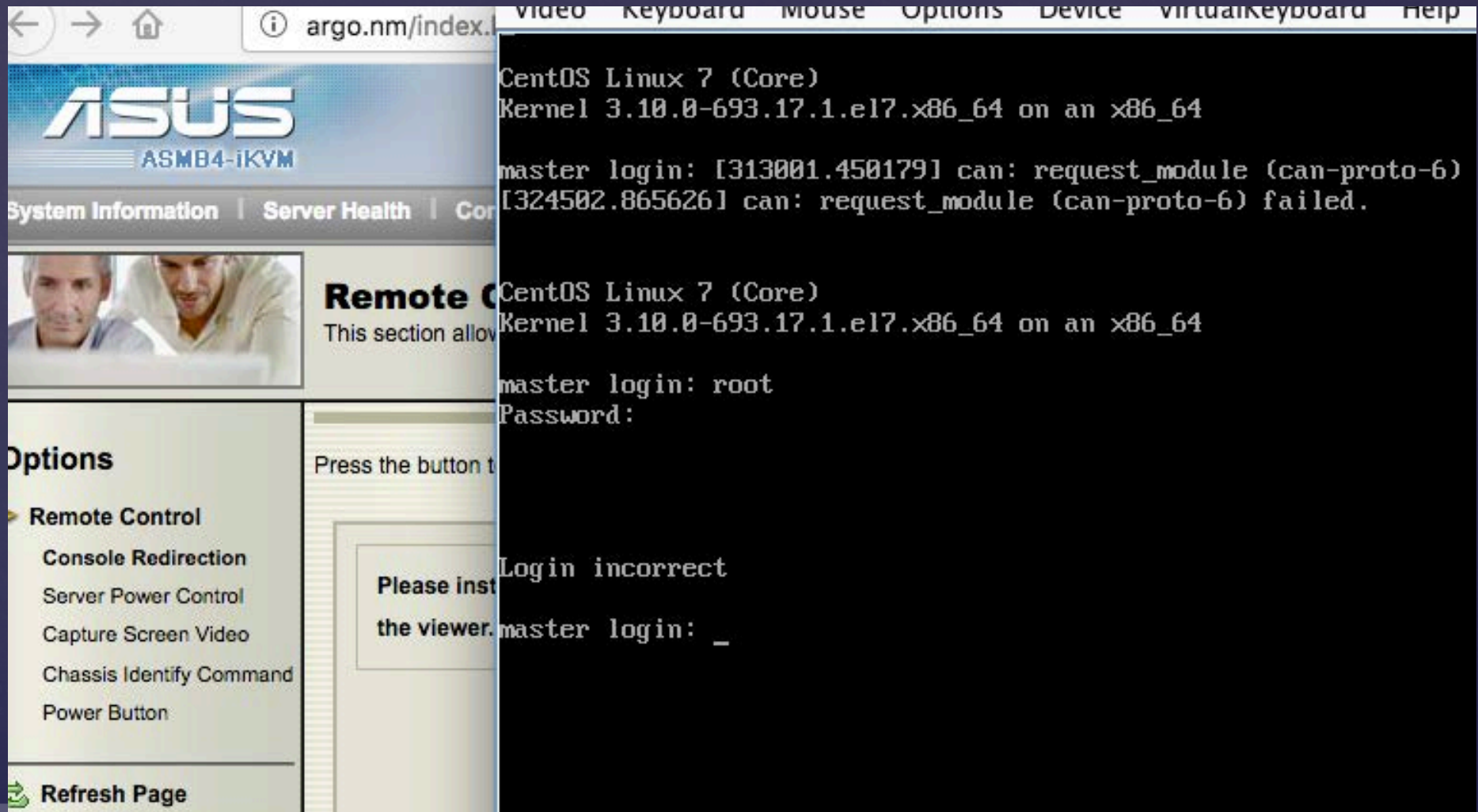
Equipment features

- & Choosing equipment features:
- & IPMI/ BMC,
- & dual power- supply,
- & warranty

Access BMC over http



remote Console



The screenshot displays the ASUS iKVM remote console interface. The browser address bar shows 'argo.nm/index...'. The interface includes a navigation menu with 'System Information', 'Server Health', and 'Console'. The 'Console' section is active, showing a terminal window with the following text:

```
CentOS Linux 7 (Core)
Kernel 3.10.0-693.17.1.el7.x86_64 on an x86_64

master login: [313001.450179] can: request_module (can-proto-6)
[324502.865626] can: request_module (can-proto-6) failed.

CentOS Linux 7 (Core)
Kernel 3.10.0-693.17.1.el7.x86_64 on an x86_64

master login: root
Password:

Login incorrect
master login: _
```

The interface also features a 'Remote Control' section with options: 'Console Redirection', 'Server Power Control', 'Capture Screen Video', 'Chassis Identify Command', and 'Power Button'. A 'Refresh Page' button is located at the bottom left. A message box on the right side of the terminal area states: 'Please install the viewer.'

Mozilla Firefox

File Edit View History Bookmarks Tools Help

http://sp-node169/index.html

http://sp-node169/index.html

intel Integrated BMC Web Console

System Information **Server Health** Configuration Remote Control

LOGOUT

Server Health

This section shows you data related to the server's health, such as sensor readings and the event log.

Sensor Readings

This page displays system sensor information, including readings and status. You can toggle viewing the thresholds for the sensors by pressing the Show T...

Sensor Readings

Event Log

Power Statistics

Refreshing readings every 60 seconds

Select a sensor type category:

All Sensors

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 Sorty	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

Refresh Show Thresholds

Set auto-refresh in seconds (0 to disable).

Set



ipmitool (Command line)

```
[ictpadmin@nagios ~]$ ipmitool -H argo-master.nm -U admin -P  
xxx power status
```

Chassis Power is on

```
[root@master ~]# ipmiwrap sp-node133 power status
```

Chassis Power is on

```
[root@master ~]# ipmiwrap sp-node131 power off
```

```
[root@master ~]# ipmiwrap sp-node131 power on
```


Sensors reported

```
[root@master ~]# ipmiwrap sp-node133 sdr list
```

Pwr Unit Status	0x01	ok
Pwr Unit Redund	0x01	ok
Physical Scrty	0x00	ok
System Event Log	0x10	ok
System Event	0x00	ok
Button	0x00	ok
VR Watchdog	0x00	ok
SSB Therm Trip	0x00	ok
BB Inlet Temp	23 degrees C	ok
BB BMC Temp	30 degrees C	ok
P1 VR Temp	25 degrees C	ok
IB Temp	28 degrees C	ok
PS1 Status	0x01	ok
PS2 Status	0x01	ok
PS1 Input Power	420 Watts	ok
PS2 Input Power	410 Watts	ok
PS1 Temperature	31 degrees C	ok
PS2 Temperature	32 degrees C	ok
P1 Status	0x80	ok
CPU Missing	0x00	ok
Auto Shutdown	0x00	ok
Mem P1 Thrm Trip	0x00	ok
BB P5V STBY	4.96 Volts	ok
BB P1_8V AUX	1.79 Volts	ok
BB P3_3V STBY	3.17 Volts	ok



ipmi checks in Nagios

Service Status Details For Host Group 'IPMI-hosts'						
Host ▲▼	Service ▲▼	Status ▲▼	Last Check ▲▼	Duration ▲▼	Attempt ▲▼	Status Information
access-sv.nm	Connectivity	OK	04-12-2019 00:14:48	23d 19h 36m 14s	1/3	PING OK - Packet loss = 0%, RTA = 0.46 ms
	IPMI-TEMP	OK	04-12-2019 00:13:41	127d 6h 53m 51s	1/3	sensor type 'Temperature' Status: OK
adserver3.nm	Connectivity	OK	04-12-2019 00:14:42	65d 13h 52m 25s	1/3	PING OK - Packet loss = 0%, RTA = 0.51 ms
	IPMI-TEMP	OK	04-12-2019 00:13:42	65d 13h 55m 51s	1/3	sensor type 'Temperature' Status: OK
adserver4.nm	Connectivity	OK	04-12-2019 00:11:52	10d 12h 14m 57s	1/3	PING OK - Packet loss = 0%, RTA = 0.45 ms
	IPMI-TEMP	OK	04-12-2019 00:11:48	0d 7h 45m 1s	1/3	sensor type 'Temperature' Status: OK
argo.nm	Connectivity	OK	04-12-2019 00:13:02	28d 7h 33m 48s	1/3	PING OK - Packet loss = 0%, RTA = 0.58 ms

- ◆ Redundancy, High Availability
 - ◆ computes are many (redundant)
 - ◆ login is one (add more for HA)
 - ◆ master is one (opt HA)
 - ◆ NFS: choose robust, and performant
- ◆ Backups:
 - + master & /home
 - Nodes can be re-installed

Conclusions:

- ⌘ buy CPU cores, but also plan for data (space and performance-wise)
- ⌘ Choose equipment conformant to user needs (Ghz, TB, Gbps) and managment needs (IPMI!)
- ⌘ exercise order while connectiong the (numerous) elements

Thank You!

Maria Verina & Clement Onime

Questions ?

