

Bright Cluster Manager

Advanced system management & monitoring made easy

... on Cray Systems

Dr Matthijs van Leeuwen
CEO

Martijn de Vries
CTO

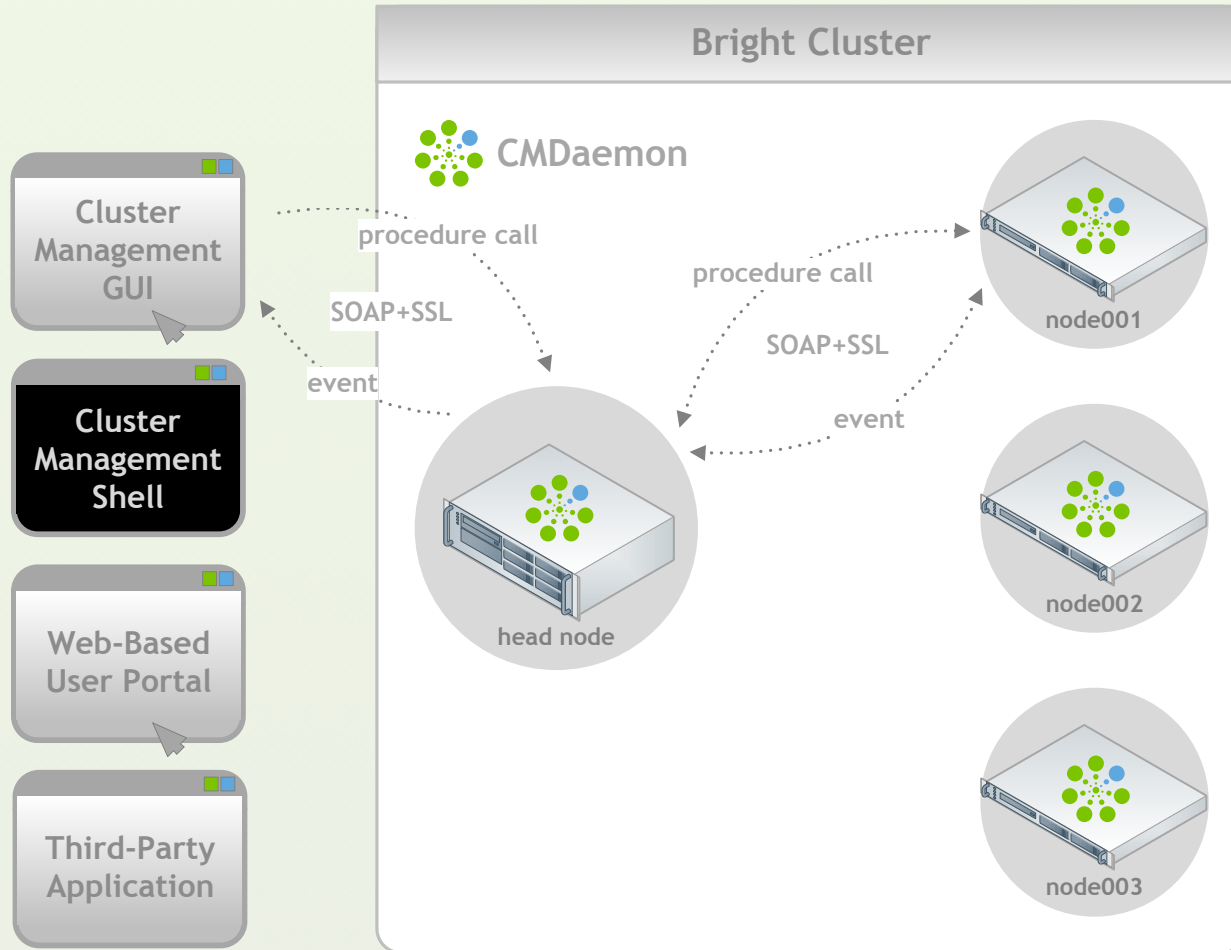


The Commonly Used “Toolkit” Approach

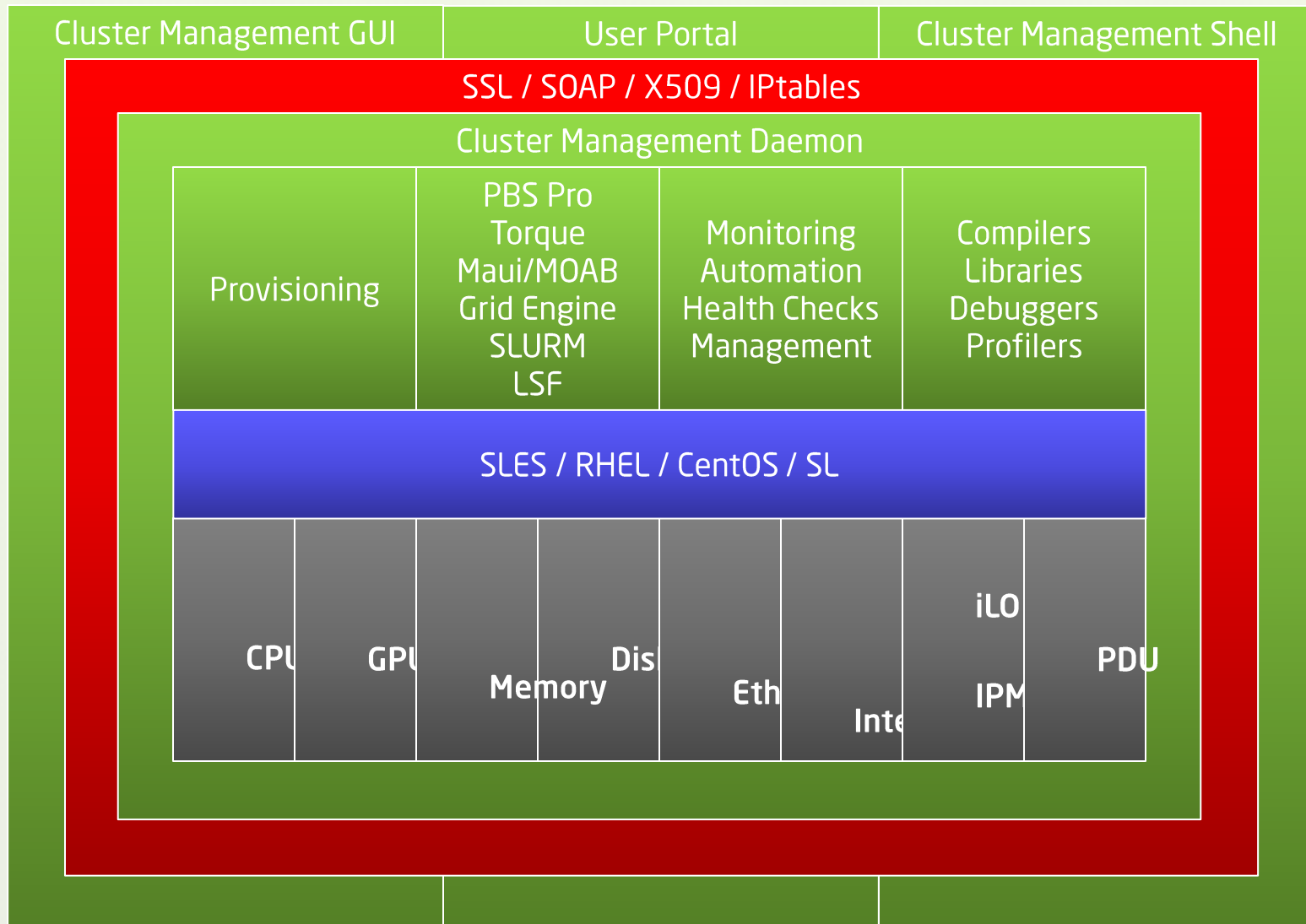
- Most HPC cluster management solutions use the “toolkit” approach (Linux distro + tools)
 - Examples: Rocks, PCM, OSCAR, UniCluster, CMU, etc.
 - Tools typically used: Ganglia, Cacti, Nagios, Cfengine, System Imager, xCAT, Puppet, Cobbler, Hobbit, Big Brother, Zabbix, Groundwork, etc.
- Issues with the “toolkit” approach:
 - Tools rarely designed to work together
 - Each tool has its own command line interface and GUI
 - Each tool has its own daemon and database
 - Tools rarely designed to scale
 - Tools rarely designed for HPC
- Making a collection of unrelated tools work together
 - Requires a lot of expertise and scripting
 - Rarely leads to a really easy-to-use and scalable solution

- Bright Cluster Manager takes a much more fundamental & integrated approach
 - Designed and written from the ground up
 - Single cluster management daemon provides all functionality
 - Single, central database for configuration and monitoring data
 - Single CLI and GUI for ALL cluster management functionality
- Which makes Bright Cluster Manager ...
 - Extremely easy to use
 - Extremely scalable
 - Secure & reliable
 - Complete
 - Flexible
 - Maintainable

Architecture



Bright Cluster Manager – Elements



Bright Cluster Manager *User Portal*

MESSAGE OF THE DAY

This is the message of the day. Feel free to edit this to your liking (in `/var/www/html/motd.php`).

On the right, you will see download and contact information. If there is no contact information available, you can set it in CMGUI/CMSH. Alternatively, you can modify `/var/www/html/contact.php`.

DOCUMENTATION

[Bright Computing website](#)

[Administrator manual](#)

[User manual](#)

CONTACT

James Smith
System Administrator
Tel: (408) 389-1922
james.smith@uni.edu

CLUSTER OVERVIEW

Uptime	9 days 8 hours 31 min	Memory	1.2 GiB out of 8.3 GiB total
Nodes	2 ↑ 6 ↓ 1 ⚡	Swap	0 B out of 32.7 GiB total
Devices	0 ↑ 1 ↓ 0 ⚡	Load	0.3% user
Cores	3 ↑ 3 total		0.2% system
Users	0 out of 2 total		99.4% idle
Phase Load	N/A ampere		0.1% other
Occupation Rate	3.3%		

WORKLOAD OVERVIEW

Queue	Scheduler	#Slots	#Nodes	#Running	#Queued	#Failed	#Completed	Avg. Duration	Est. Delay
short.q	Slurm	0	256	32	43	0	482	00:07:27	00:09:05
medium.q	Slurm	0	128	5	11	0	41	02:15:00	04:16:00
long.q	Slurm	0	128	8	13	0	91	08:09:00	15:13:00

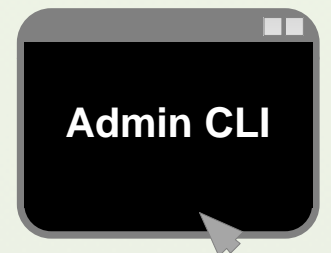
Graphical User Interface (GUI)

- Offers administrator full cluster control
- Standalone desktop application
- Manages multiple clusters simultaneously
- Runs natively on Linux & Windows



Cluster Management Shell (CMSH)

- All GUI functionality also available through Cluster Management Shell
- Interactive and scriptable in batch mode



RESOURCES

My Clusters

Seismic Houston

Switches

switch01

switch02

switch03

switch04

switch05

Networks

externalnet

ipminet

mpinet

slavenet

storagenet

Power Distribution Units

apc01

apc02

apc03

apc04

Software Images

default-image

Node Categories

slave

Head Nodes

demohead1

demohead2



Welcome to Bright Cluster Manager

Seismic Oslo



Modified: No

Host: oslo.seismic.com:8081

Connected: No

Certificate: /root/oslo.pfx



Seismic Abu Dhabi



Modified: No

Host: abudhabi.seismic.com:8081

Connected: No

Certificate: /root/.cm/cmgui/admin-abudhabi.pfx



Seismic Houston



Modified: No

Host: localhost:2581

Connected: Yes

Certificate: /root/.cm/cmgui/admin.pfx



Add a new cluster

EVENT VIEWER



All Events

	Ack	Time	Cluster	Source	Message
i		18/Sep/2009 17:05:53	Demo Cluster	demohead1	Service ntpd was restarted on demohead1
i		18/Sep/2009 17:05:47	Demo Cluster	demohead1	Service named was restarted on demohead1
i		18/Sep/2009 17:05:45	Demo Cluster	demohead1	Service postfix was restarted on demohead1
i		18/Sep/2009 17:05:45	Demo Cluster	demohead1	Service dhcpd was restarted on demohead1
i		18/Sep/2009 17:05:45	Demo Cluster	demohead1	Service maui was restarted on demohead1

RESOURCES

My Clusters

Demo Cluster

Switches

- switch01
- switch02
- switch03
- switch04
- switch05

Networks

- externalnet
- ipminet
- mpinet
- slavenet
- storagenet

Power Distribution Units

- apc01
- apc02
- apc03
- apc04

Software Images

- default-image

Node Categories

- slave

Head Nodes

- demohead1
- demohead2

Racks

Chassis

Virtual SMP Nodes

Slave Nodes

- node001
- node002
- node003
- node004
- node005
- node006
- node007
- node008
- node009

Demo Cluster

Overview

Settings

Failover

Rackview

Health

Parallel shell

License

Notes

Uptime: 45 days 3 hours 7 minutes

Nodes: 503 ↑ 7 ↓ 2 ⊖

GPU Units: 38 ↑ 0 ↓ 0 ⊖

Devices: 64 ↑ 0 ↓ 0 ⊖

Jobs: 45 running 67 waiting

Phase load: 783 A

CPU Cores:

3.93 K out of 4 K

GPUs:

13 out of 38

Memory:

7.32 TB out of 7.45 TB

Users:

13 out of 38

CPU Usage:

48% u 29% s 13% o 10% i

Occupation rate:

83.2 %

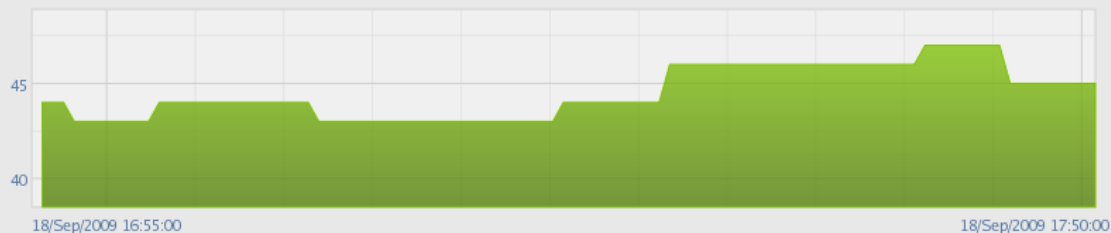
Disk Usage

Mountpoint	Used	Size	Use %
/	15.83 GB	37.25 GB	<div></div>
/boot	14.31 MB	99.18 MB	<div></div>
/home	832.6 GB	9.91 TB	<div></div>

Workload Management

Queue	Running	Queued	Error	Completed	Avg. Duration	Est. delay
short.q	32	43	0	482	7 hours, 27 minutes	9 hours, 5 minutes
medium.q	5	11	0	41	2 days, 15 hours	4 days, 16 hours
long.q	8	13	0	91	8 days, 9 hours	15 days, 13 hours

Metric: RunningJobs[all.q]

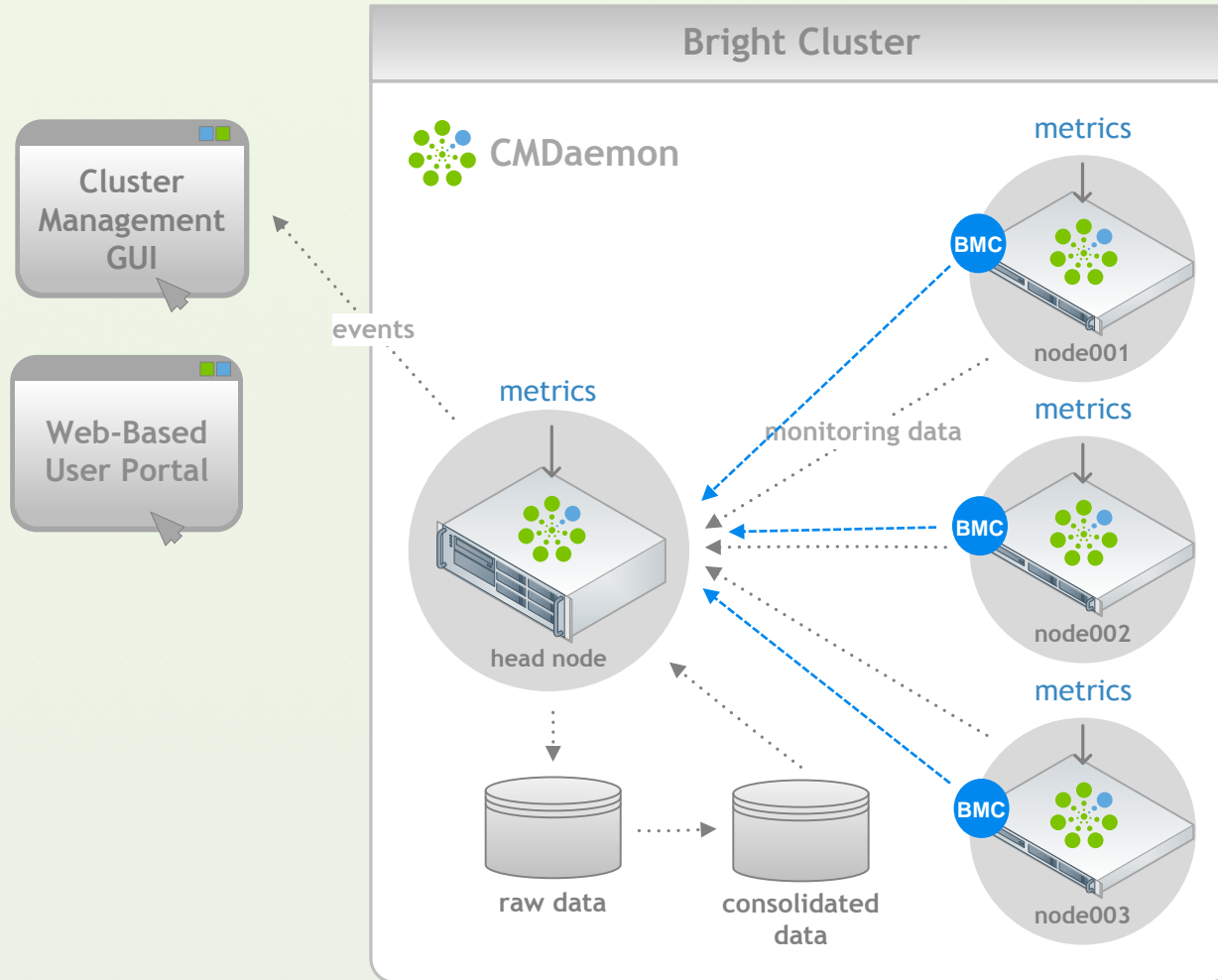


EVENT VIEWER

All Events

▼	Ack	Time	▲	Cluster	▼	Source	▼	Message	▼
ⓘ		18/Sep/2009 17:05:53		Demo Cluster		demohead1		Service ntpd was restarted on demohead1	
ⓘ		18/Sep/2009 17:05:47		Demo Cluster		demohead1		Service named was restarted on demohead1	
ⓘ		18/Sep/2009 17:05:45		Demo Cluster		demohead1		Service postfix was restarted on demohead1	
ⓘ		18/Sep/2009 17:05:45		Demo Cluster		demohead1		Service dhcpcd was restarted on demohead1	
ⓘ		18/Sep/2009 17:05:45		Demo Cluster		demohead1		Service maui was restarted on demohead1	

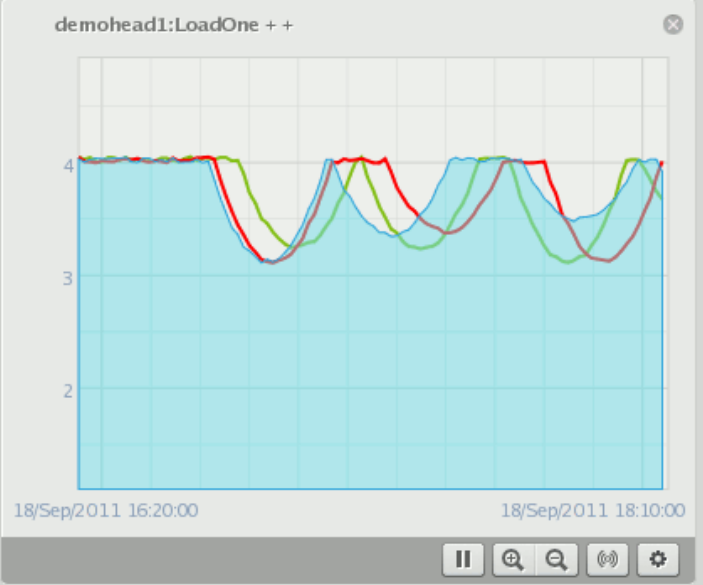
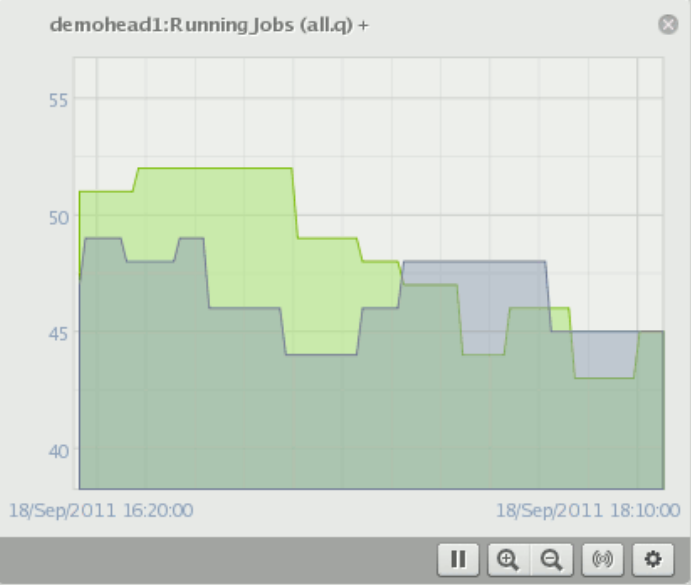
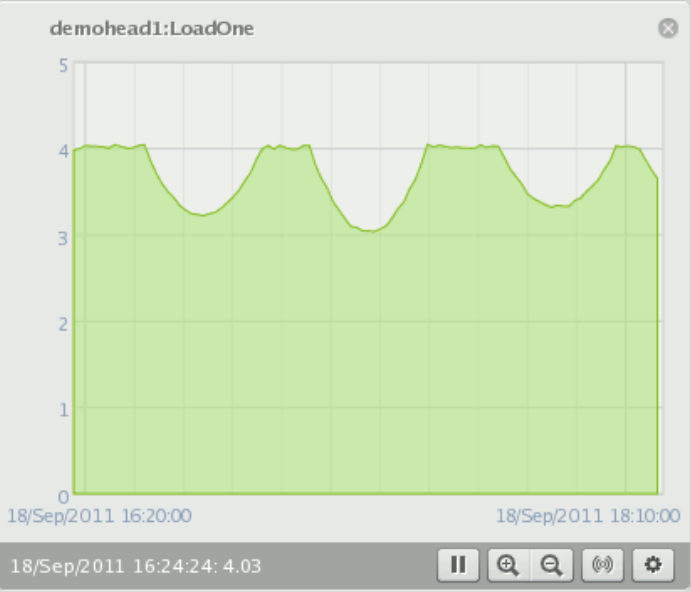
Architecture – Monitoring



RESOURCES

- demohead1
- CPU
- Disk
- Memory
 - BufferMemory (B)
 - CacheMemory (B)
 - MemoryFree (B)
 - MemoryUsed (B)
 - SwapFree (B)
 - SwapUsed (B)
- Network
- Operating System
 - CtxtSwitches (ctx_switch/s)
 - Forks (process/s)
 - LoadFifteen
 - LoadFive
 - LoadOne
 - ProcessCount
 - RunningProcesses
 - Uptime (s)
 - ldap
 - mysql
- Internal
- Workload
 - AvgExpFactor
 - AvgJobDuration[defq] (s)
 - CompletedJobs[defq]
 - EstimatedDelay[defq] (s)
 - FailedJobs[defq]
 - QueuedJobs[defq]
 - RunningJobs[defq]
 - failedprejob
 - schedulers
- Cluster
 - CPUCoresAvailable
 - DevicesUp
 - GPUAvailable
 - NetworkBytesRecv (B)
 - NetworkBytesSent (B)
 - NodesUp
 - OccupationRate (%)

Demo Cluster



Bright Cluster Manager

FileMonitoringViewHelp

RESOURCES

My Clusters

Seismic Houston

- Switches
 - switch01
 - switch02
 - switch03
 - switch04
 - switch05
- Networks
 - externalnet
 - ipminet
 - mpinet
 - slavenet
 - storagenet
- Power Distribution Units
 - apc01
 - apc02
 - apc03
 - apc04
- Software Images
 - default-image
- Node Categories
 - slave
- Head Nodes
 - demohead1
 - demohead2
- Racks
- Chassis
- Virtual SMP Nodes
- Slave Nodes
- Other Devices
- Node Groups
- Users & Groups
- Workload Management
- Monitoring Configuration
- Authorisation
- Authentication

Seismic Houston

OverviewSettingsFailoverRackviewHealthParallel shellLicenseNotes

U	Rack 1	Rack 2	Rack 3	Rack 4	Rack 5	Rack 6
01	demohead1	032	057	097098		231
02		033	058	099100		233234
03		034	059	101102		235236
04		035	060	103104		237238
05	demohead2	036	061	105106		239240
06		037	062	107108		241242
07		038	063	109110		243244
08		039	064	111112		245246
09		040	065	113114		247248
10		041		115116		
11	001	042	066	117118	169170	249250
12	002	043	067	119120	171172	251252
13	003	044	068	121122	173174	253254
14	004	045	069	123124	175176	255256
15	005	046	070	125126	177178	257258
16	006	047	071	127128	179180	259260
17	007	048	072	129130	181182	261262
18	008	049	073	131132	183184	263264
19	009		074		185186	265266
20	010		075		187188	267268
21	011		076		189190	269270
22	012		077		191192	271272
23	013		078		193194	273274
24	014		079		195196	275276
25	015	050	080	133134	197198	277278
26	016	051	081	135136	199200	279280
27	017	052	082	137138	201202	281282
28	018	053	083	139140	203204	283284
29	019	054	084	141142	205206	285286
30	020	055	085	143144	207208	287288
31		056		145146		

View:

Refresh

Setup

Temp CPU0 OC68.74 C

Temp CPU1 OC68.74 C

EVENT VIEWER

All Events

	Ack	Time	Cluster	Source	Message
i		18/Sep/2009 17:05:53	Demo Cluster	demohead1	Service ntpd was restarted on demohead1
i		18/Sep/2009 17:05:47	Demo Cluster	demohead1	Service named was restarted on demohead1
i		18/Sep/2009 17:05:45	Demo Cluster	demohead1	Service postfix was restarted on demohead1
i		18/Sep/2009 17:05:45	Demo Cluster	demohead1	Service dhcpd was restarted on demohead1
i		18/Sep/2009 17:05:45	Demo Cluster	demohead1	Service maui was restarted on demohead1

Ready

RESOURCES

My Clusters

Seismic Houston

Switches



switch01



switch02



switch03



switch04



switch05

Networks



externalnet



ipminet



mpinet



slavenet



storagenet

Power Distribution Units



apc01



apc02



apc03



apc04

Software Images



default-image

Node Categories



slave

Head Nodes



demohead1



demohead2

Slave Nodes

Other Devices

Node Groups



Users & Groups



Workload Management



Monitoring Configuration



Authorisation



Authentication



Seismic Houston

Overview

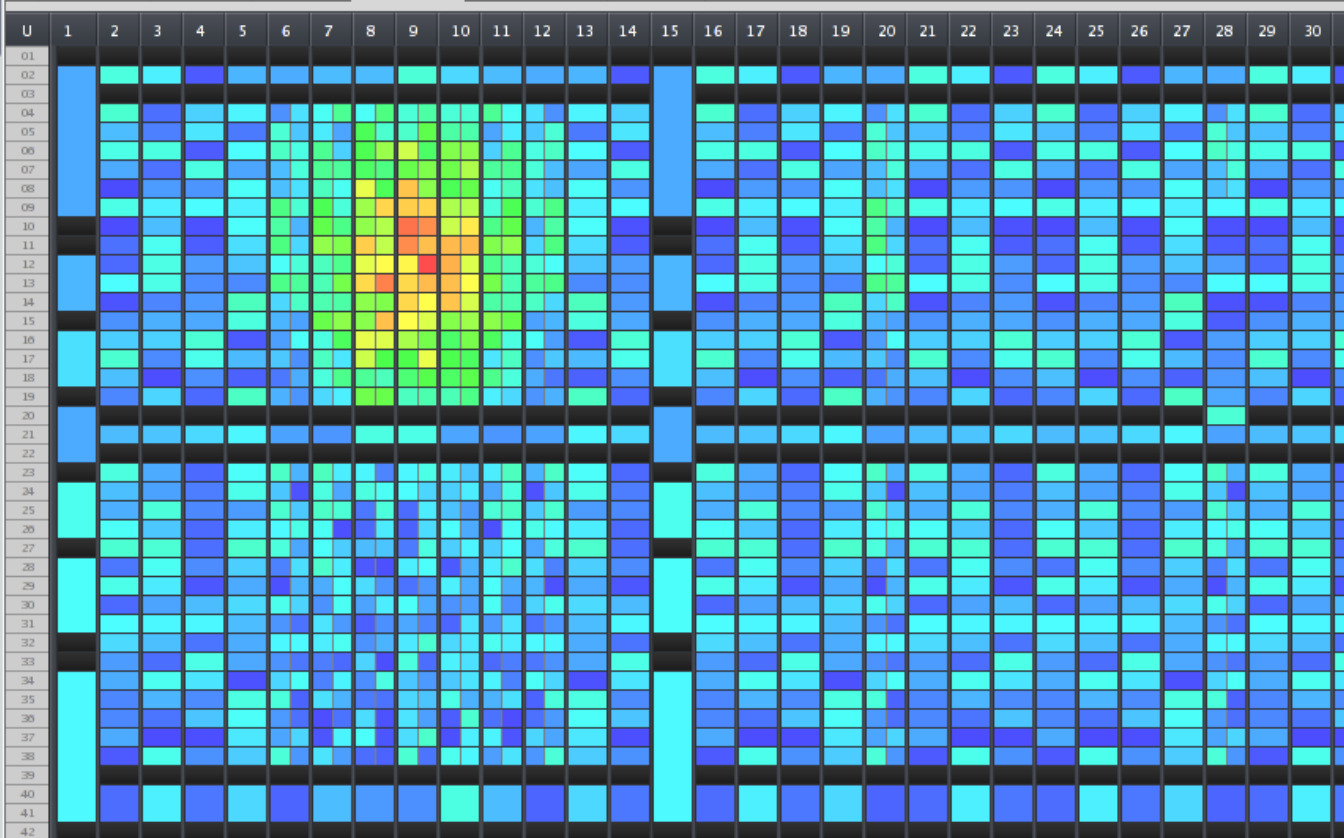
Settings

Failover

Rackview

Parallel shell

License



View:

☒ Live sampling[Refresh](#)

Metric 1: Temperature

30.01

69.34

EVENT VIEWER



All Events

	Ack	Time	Cluster	Source	Message
		18/Sep/2009 17:05:53	Demo Cluster	demohead1	Service ntpd was restarted on demohead1
		18/Sep/2009 17:05:47	Demo Cluster	demohead1	Service named was restarted on demohead1
		18/Sep/2009 17:05:45	Demo Cluster	demohead1	Service postfix was restarted on demohead1
		18/Sep/2009 17:05:45	Demo Cluster	demohead1	Service dhcpd was restarted on demohead1
		18/Sep/2009 17:05:45	Demo Cluster	demohead1	Service maui was restarted on demohead1

Bright Cluster Manager

FileMonitoringViewHelp

RESOURCES

My Clusters

Demo Cluster

- Switches
 - switch01
 - switch02
 - switch03
 - switch04
 - switch05
- Networks
 - externalnet
 - ipminet
 - mpinet
 - slavenet
 - storagenet
- Power Distribution Units
 - apc01
 - apc02
 - apc03
 - apc04
- Software Images
 - default-image
- Node Categories
 - slave
- Head Nodes
 - demohead1
 - demohead2
- Racks
- Chassis
- Virtual SMP Nodes
- Slave Nodes
- Other Devices
- Node Groups
 - Large Memory Nodes
- Users & Groups
- Workload Management
- Monitoring Configuration
- Authorisation
- Authentication

Monitoring Configuration

Demo Cluster

OverviewMetric ConfigurationHealth Check ConfigurationMetricsHealth ChecksActions

Category	Metric	Parameter	Threshold Bound	Action	Action Parameter
All Master Nodes	FreeSpace	/	< 10 GB	NotifyVendor	
All Master Nodes	FreeSpace	/	< 10 GB	SendEmail	administrator@localhost
All Master Nodes	FreeSpace	/home	< 10 GB	NotifyVendor	
All Master Nodes	FreeSpace	/home	< 10 GB	SendEmail	administrator@localhost
All Power Distributio...	PDULoad		> 32 A	SendEmail	datacenter_support@uni.edu
slave	Temperature		> 70	SendEmail	administrator@localhost
slave	Temperature		> 70	Shutdown	

Monitoring Rules Wizard

Select Category:

All Power Distribution Units

All Ethernet Switches

All Myrinet Switches

All IB Switches

All Master Nodes

All Rack Sensors

All Generic Devices

slave

Cancel

Previous

Next

Edit

Add rule

Remove

Refresh

Save

EVENT VIEWER

All Events

	Ack	Time	Cluster	Source	Message
		18/Sep/2009 18:30:06	Demo Cluster	demohead1	node003 Installing
		18/Sep/2009 18:29:39	Demo Cluster	demohead1	New certificate request with ID: 5
		18/Sep/2009 18:29:36	Demo Cluster	demohead1	node002 Installing
		18/Sep/2009 18:29:25	Demo Cluster	demohead1	New certificate request with ID: 4
		18/Sep/2009 17:05:53	Demo Cluster	demohead1	Service ntpd was restarted on demohead1
		18/Sep/2009 17:05:47	Demo Cluster	demohead1	Service named was restarted on demohead1

Ready

Bright Cluster Manager

FileMonitoringViewHelp

RESOURCES

My Clusters

Demo Cluster

Switches

switch01

switch02

switch03

switch04

switch05

Networks

externalnet

ipminet

mpinet

slavenet

storagenet

Power Distribution Units

apc01

apc02

apc03

apc04

Software Images

default-image

Node Categories

slave

Head Nodes

demohead1

demohead2

Slave Nodes

Other Devices

Node Groups

Users & Groups

Workload Management

Monitoring Configuration

Authorisation

Authentication

Users & Groups

Demo Cluster

Users

Groups

Modified	Name	Full Name	User ID	Home directory
	alex	Alex Bozarski	504	/home/alex
✓	james	James Watt	505	/home/james
✓	jodi	Jodi Johnson	503	/home/jodi
✓	kate	Kate Moss	506	/home/kate
✓	koen	Koen van den Bosch	502	/home/koen
✓	matthew	Matthew Ellis	507	/home/matthew

Edit User

Login Name:

matthew

User ID:

507

Full Name:

Matthew Ellis

Group ID:

matthew

Login shell:

/bin/bash

Home directory:

/home/matthew

Password:

.....

Retype password:

.....

Cancel

Ok

Edit

Add

Remove

Revert

Save

EVENT VIEWER

All Events

	Ack	Time	Cluster	Source	Message
i		18/Sep/2009 17:05:53	Demo Cluster	demohead1	Service ntpd was restarted on demohead1
i		18/Sep/2009 17:05:47	Demo Cluster	demohead1	Service named was restarted on demohead1
i		18/Sep/2009 17:05:45	Demo Cluster	demohead1	Service postfix was restarted on demohead1
i		18/Sep/2009 17:05:45	Demo Cluster	demohead1	Service dhcpcd was restarted on demohead1
i		18/Sep/2009 17:05:45	Demo Cluster	demohead1	Service maui was restarted on demohead1

Ready

RESOURCES

- My Clusters
 - Demo Cluster
 - Switches
 - switch01
 - switch02
 - switch03
 - switch04
 - switch05
 - Networks
 - externalnet
 - ipminet
 - mpinet
 - slavenet
 - storagenet
 - Power Distribution Units
 - apc01
 - apc02
 - apc03
 - apc04
 - Software Images
 - default-image
 - Node Categories
 - slave
 - Head Nodes
 - demohead1
 - demohead2
 - Slave Nodes
 - Other Devices
 - Node Groups
 - Users & Groups
 - Workload Management
 - Monitoring Configuration
 - Authorisation**
 - Authentication



Authorisation

Demo Cluster

Jadmin



CMMain

- ☐ Access
- ☐ GET_LICENSE_INFO_TOKEN
- ☐ GET_VERSION_TOKEN
- ☐ GET_SERVER_STATUS_TOKEN
- ☐ GET_CLUSTER_SETUP_TOKEN
- ☐ PING_TOKEN
- ☐ PCOPY_TOKEN
- ☐ UPDATE_SELF_TOKEN

CMSession

- ☒ Access
- ☐ CMSESSION_ADMIN
- ☐ GET_SESSION_TOKEN
- ☐ REGISTER_NODE_SESSION_TOKEN

admin (read-only)



CMMain

- ☒ Access
- ☒ GET_LICENSE_INFO_TOKEN
- ☒ GET_VERSION_TOKEN
- ☒ GET_SERVER_STATUS_TOKEN
- ☒ GET_CLUSTER_SETUP_TOKEN
- ☒ PING_TOKEN
- ☒ PCOPY_TOKEN
- ☒ UPDATE_SELF_TOKEN

CMSession

- ☒ Access
- ☒ CMSESSION_ADMIN
- ☒ GET_SESSION_TOKEN
- ☒ REGISTER_NODE_SESSION_TOKEN

node (read-only)



CMMain

- ☒ Access
- ☐ GET_LICENSE_INFO_TOKEN
- ☐ GET_VERSION_TOKEN
- ☐ GET_SERVER_STATUS_TOKEN
- ☐ GET_CLUSTER_SETUP_TOKEN
- ☐ PING_TOKEN
- ☐ PCOPY_TOKEN
- ☐ UPDATE_SELF_TOKEN

CMSession

- ☒ Access
- ☐ CMSESSION_ADMIN
- ☐ GET_SESSION_TOKEN
- ☒ REGISTER_NODE_SESSION_TOKEN



Add Profile

Revert

Save

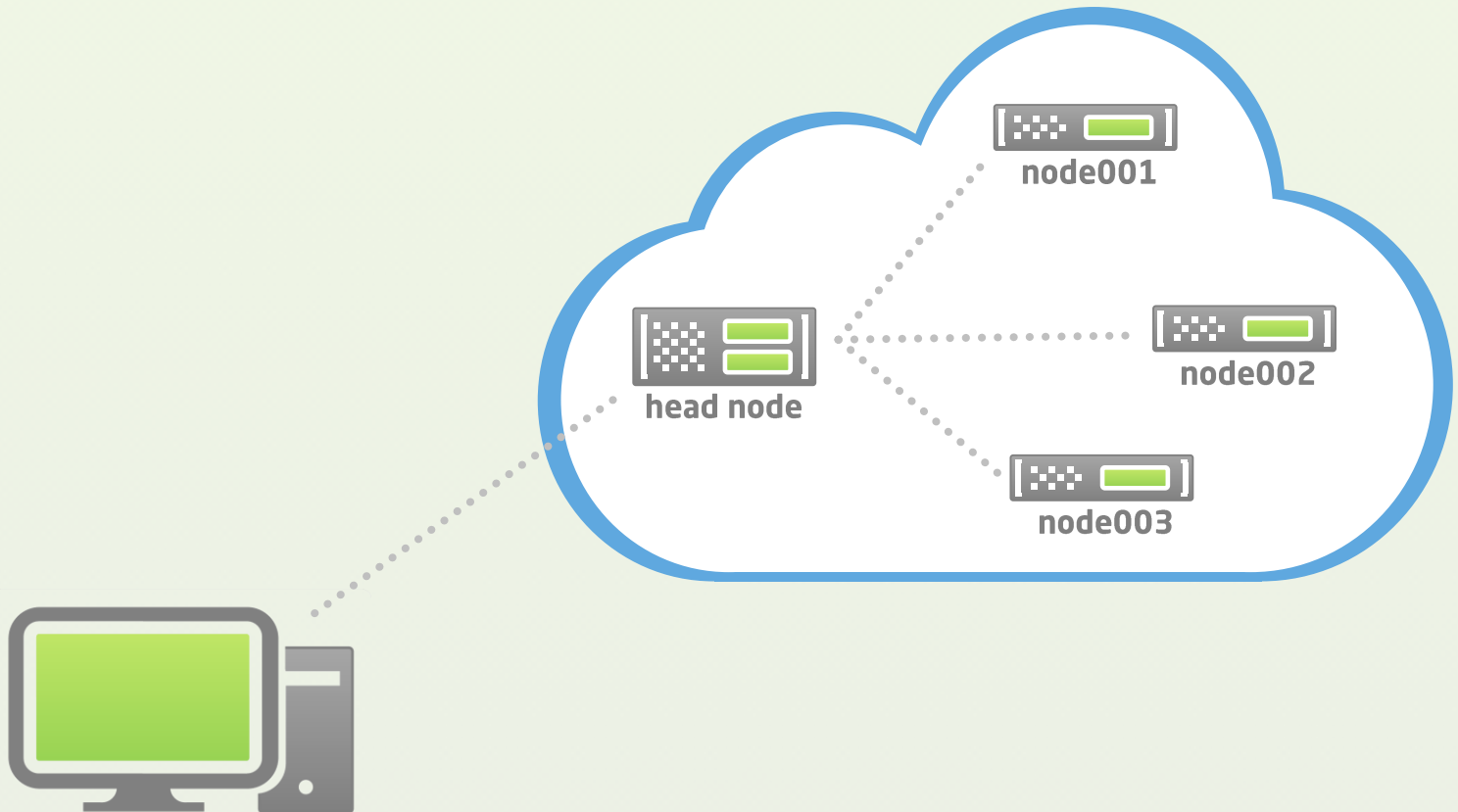
EVENT VIEWER



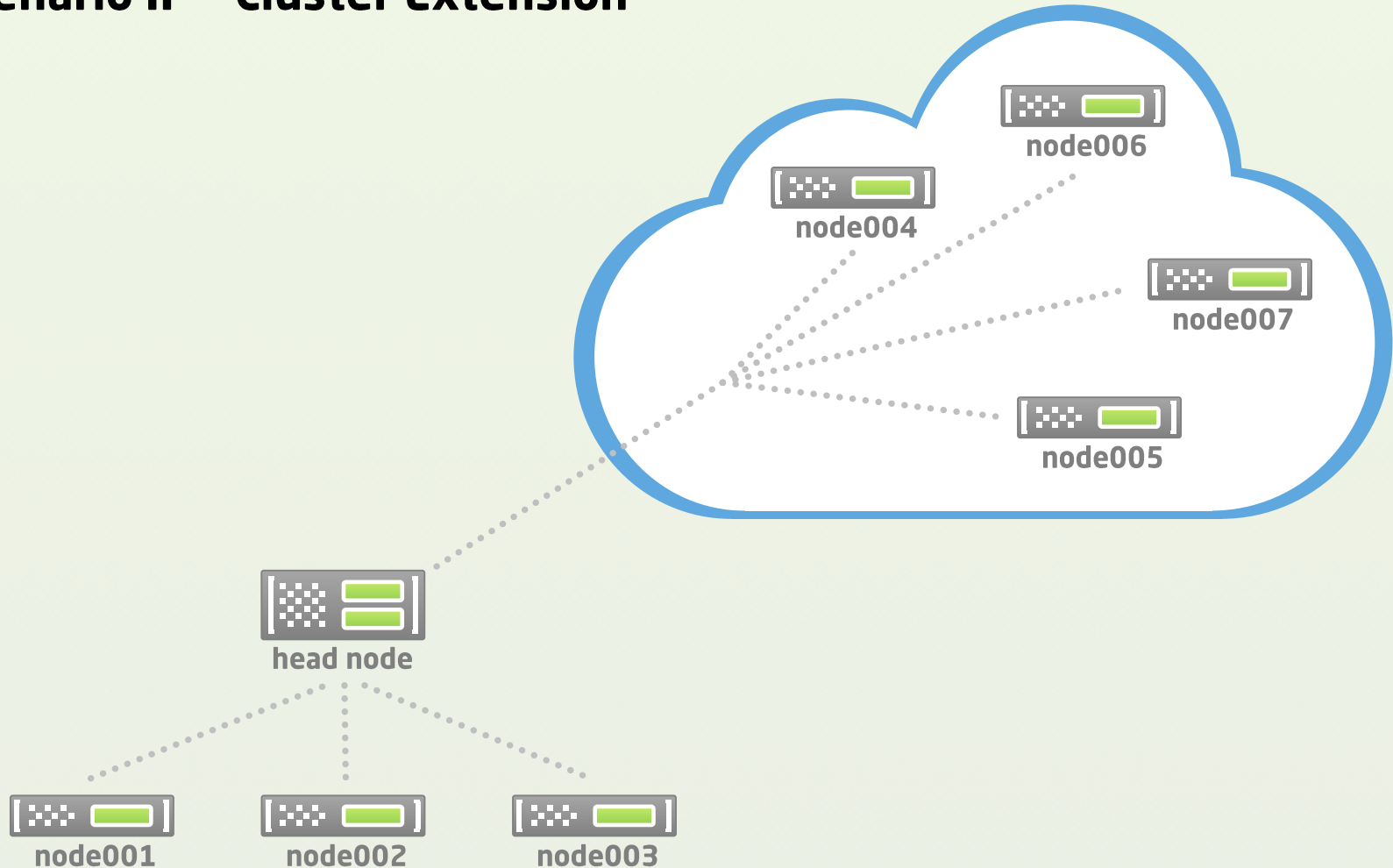
All Events

	▼	Ack	Time	▲	Cluster	▼	Source	▼	Message	▼
			18/Sep/2009 17:05:53		Demo Cluster		demohead1		Service ntpd was restarted on demohead1	
			18/Sep/2009 17:05:47		Demo Cluster		demohead1		Service named was restarted on demohead1	
			18/Sep/2009 17:05:45		Demo Cluster		demohead1		Service postfix was restarted on demohead1	
			18/Sep/2009 17:05:45		Demo Cluster		demohead1		Service dhcpd was restarted on demohead1	
			18/Sep/2009 17:05:45		Demo Cluster		demohead1		Service maui was restarted on demohead1	

Scenario I - "Cluster on Demand"



Scenario II - "Cluster Extension"



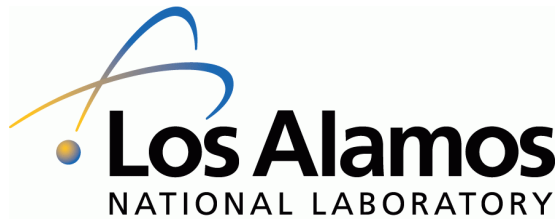
The Cray logo is rendered in a bold, blue, sans-serif font. The letters are thick and blocky, with a slight slant to the right.

&

Bright Cluster Manager

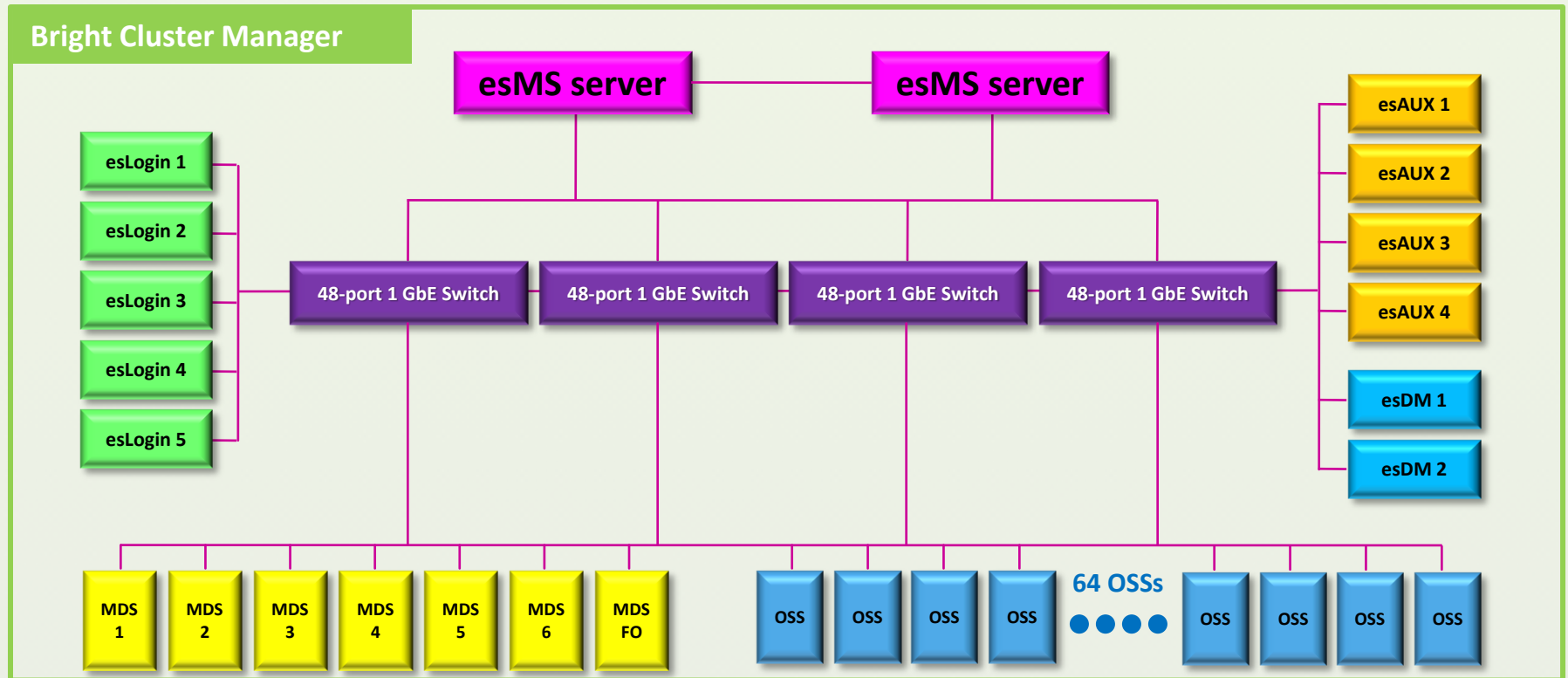
- Bright Cluster Manager default for Cray External Service nodes since 2010:
 - esMS
 - esLogin
 - esAUX
 - esDM
 - Lustre MDS
 - Lustre OSS
- Bright Cluster Manager considered for integration with Sonexion
- Bright Cluster Manager considered for Cray mainframe

Cray/Bright Customers



Bright on the Cray External Nodes

1. esMS servers are Bright head nodes in failover mode
2. All other servers are Bright slave nodes
3. Bright does provisioning, monitoring, alerting, automation, health checking, access control, Lustre failover, etc.



Bright on Cray XE6 / XK6

Project goal:

Investigate whether Bright Cluster Manager can be used to manage and monitor a Cray XE6

Effort made:

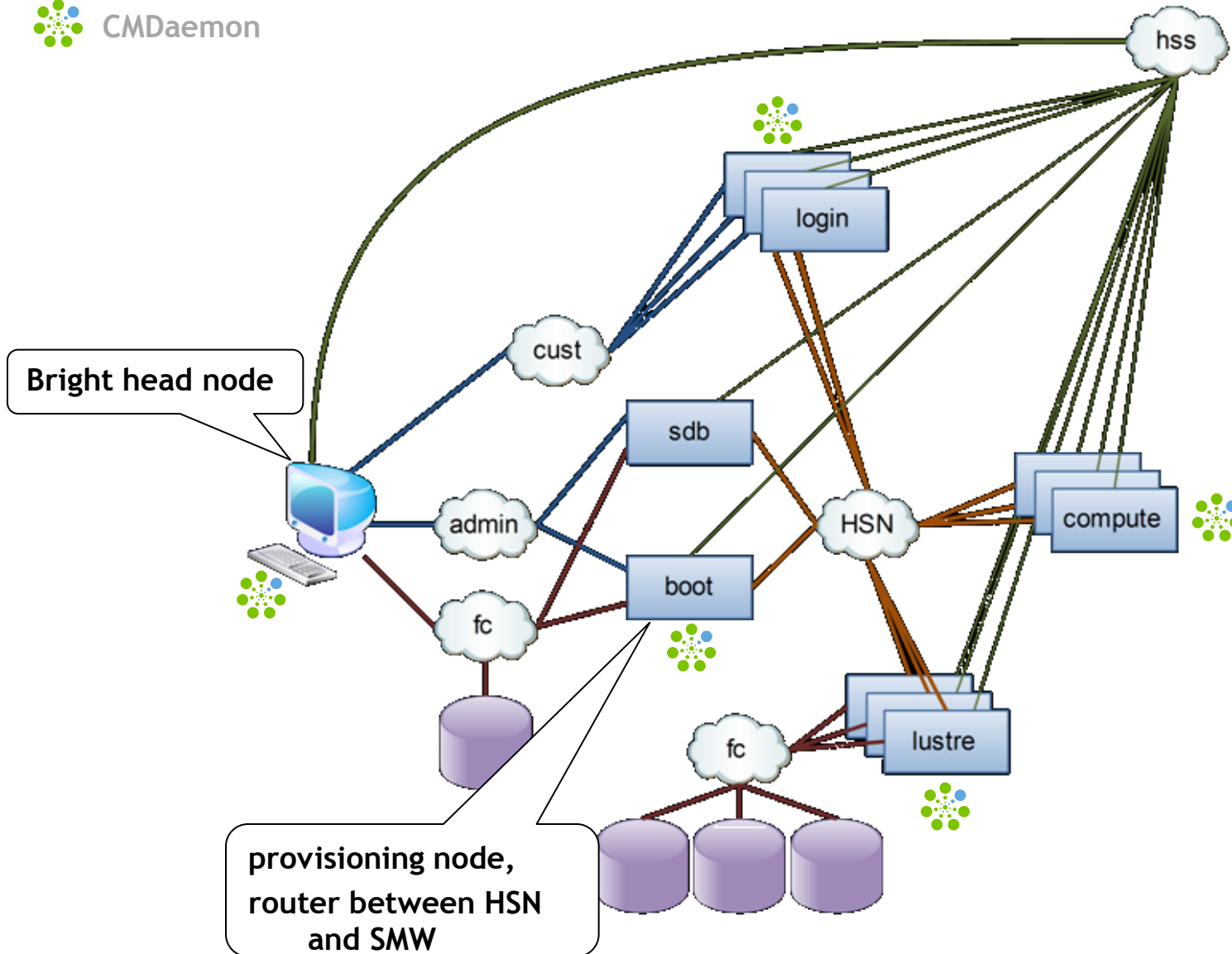
Worked almost 2 weeks with 2 developers

Result:

Bright 6.0b successfully manages XE6 making use of Cray Linux Environment infrastructure

High Level System Overview

 CMDaemon



- Default Cray kernel is used
- Bright Node Installer called from code in cpio boot image
- Node Installer:
 - Starts when node boots
 - Determines node identity based on Gemini MAC (could be NID)
 - Provisions software image into tmpfs filesystem
- Default software image about 2.7GB (can be reduced heavily)
- Parts of software image can be imported over NFS/DVS (minimal setup requires ~80MB)
- Root over NFS/DVS (now) also possible

- Power management:
 - allow components (e.g. nodes) to be reset
 - allow entire system to be reset powered on/off
 - using xtbootsys (probably need lower level utilities in future)
- Remote console:
 - allow console of nodes to be accessed
 - using xtcon
- Monitoring:
 - allow Cray hardware metrics to be monitored
 - using SEDC (would be good to get direct access to HSS instead)
- Health checking (*not done yet*):
 - Cray hardware health checks in Bright health checking framework
 - using xthealth

Why run Bright on Cray?

- Less steep learning curve for Administrators
- Single interface for managing mainframe and external service nodes (e.g. login, storage)
- Single head node (SMW) which manages everything (HA possible)
- Consistent software image on login and compute nodes
- Same solution across the data-centre which makes integration of Cray system easier
- Access to cutting edge features (e.g. cloud bursting, monitoring, health checking, GPU management, role based access control)

- Ability to scale cluster usually limited by head node providing vital services
- Bright philosophy: allow all services provided by head node to be off-loaded to multiple dedicated nodes
- Allow (re-)configuration on the fly by assigning *roles* to nodes
- Example: node can be turned into provisioning node by assigning it the *Provisioning role*
- Goal: Linear scaling in terms of node-count
- In large clusters head node is not responsible for anything

- CMDaemon resident memory size: 31MB
- 7.5 CPU core-seconds per day
- On 16 core node, less than 0.5s wall-clock time per day
- Just 15m wall-clock time lost over 5 years
- Metrics are sampled out-of-band where possible (e.g. through SEDC)
- Other metrics are sampled from within CMDaemon process (i.e. no fork())
- Monitoring configuration highly tunable
- Metric sampling synchronized as much as possible
- No measurable OS jitter at small scale, large scale remains to be tested

- Create clean installation procedure which integrates nicely into Cray installation procedure
- Migrate some services (e.g. named, LDAP) from SMW to boot node
- Support Cray component hierarchy natively:
 - Cabinet -> Cage -> Slot -> Node (*Cray*)
 - Rack -> Chassis -> ??? -> Node (*Bright*)
- Integrate Cray health checks into Bright health checking framework
- Let CLE tools such as “xtcli status” recognize nodes running Bright (currently reports nodes as down)
- Rack view which resembles physical layout of Cray system

- Allow nodes to be easily switched between classic Cray mode and Bright mode
- Extend range syntax in CMSH to support Cray-style hostnames
- Tighter integration with CLE (e.g. directly calling HSS)
- Integrate with Cray user environment (compilers, libraries, MPI)
- Improve integration of power management for individual components
- Test everything at scale

Questions?