

# Bright Cluster Manager

Advanced HPC cluster management made easy

株式会社ベストシステムズ  
代表取締役 西 克也





## 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
  - Tools rarely designed for HPC
  - Tools rarely designed to scale
  - Each tool has its own command line interface and GUI
  - Each tool has its own daemon and database
  - Roadmap dependent on developers of the tools
- 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

**BAD IDEA**

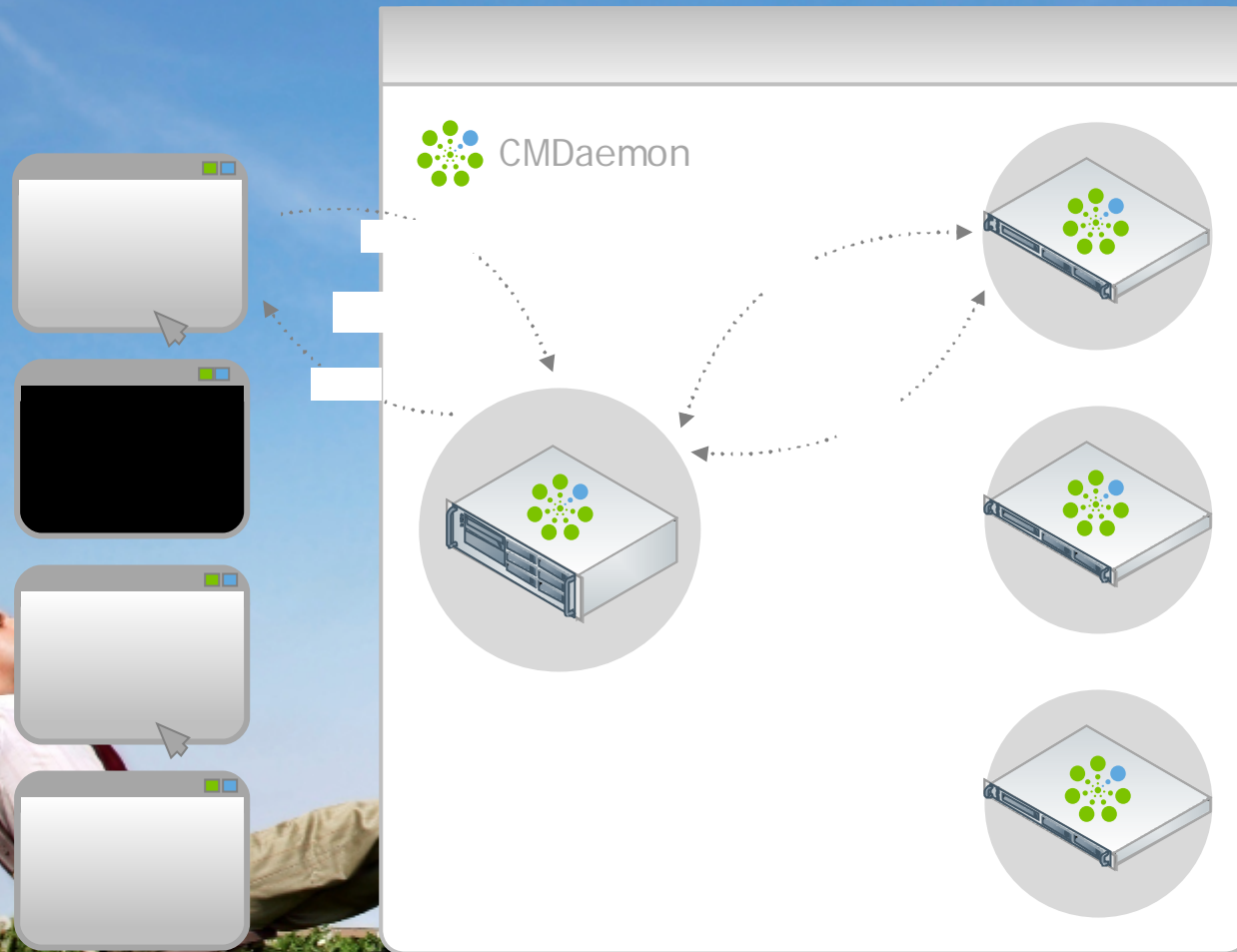


## About Bright Cluster Manager

- 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

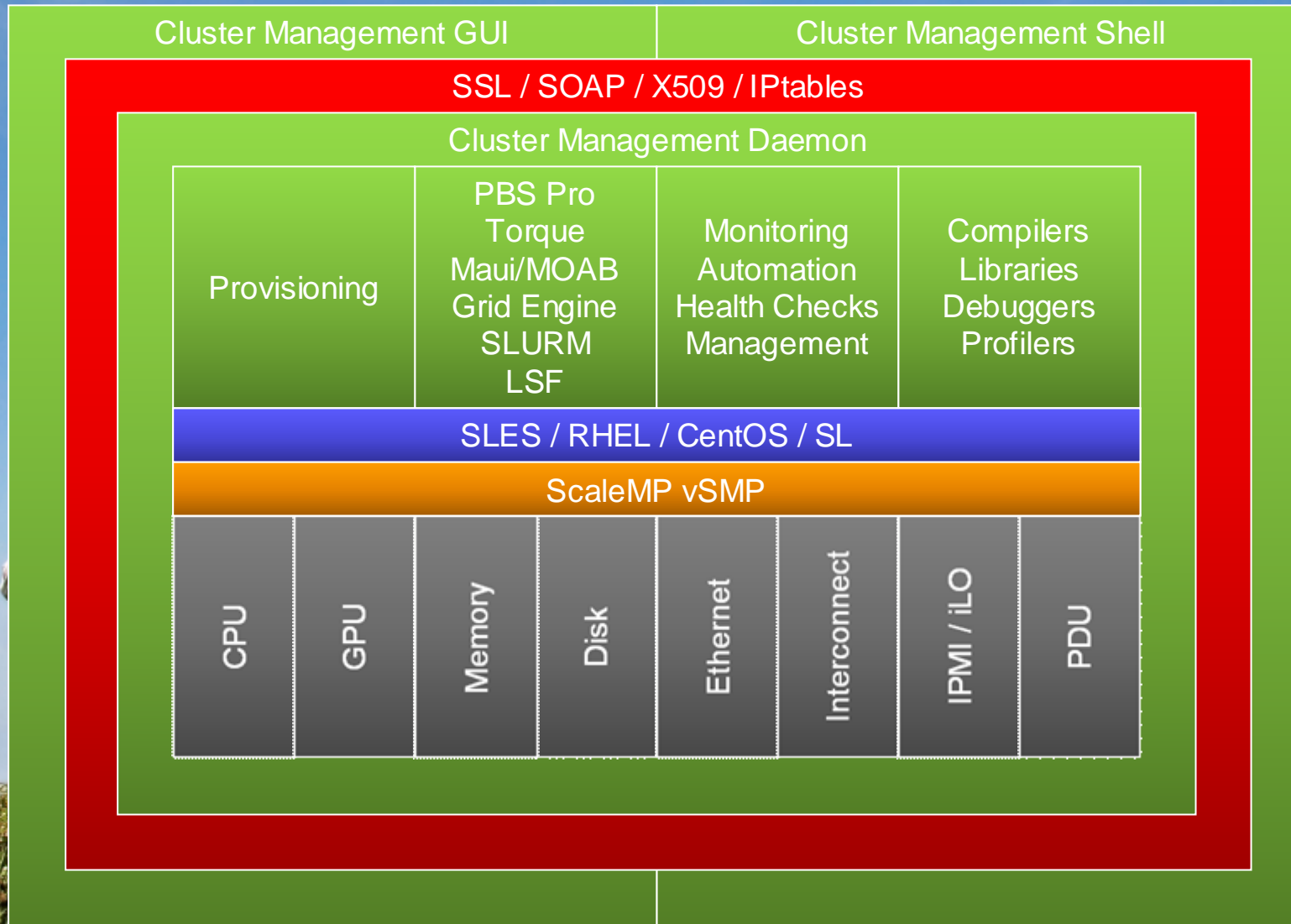


# Architecture





# Bright Cluster Manager — Elements



Bright Cluster Manager - J... x

https://demo.brightcomputing.com

**Bright Computing** Logged in as: **mal001** | [Logout](#)

[HOME](#) [WORKLOAD](#) [NODES](#) [GRAPHS](#)

## 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 [/usr/www/html/modd.php](#)).

On the right, you will see download and contact information. If there is no contact information available, you can set it in [CM3J/CM3H](#). Alternatively, you can modify [/usr/www/html/contact.php](#).

**DOCUMENTATION**

[Bright Computing website](#)

[Administrator manual](#)

[User manual](#)

**CONTACT**

James Smith  
System Administrator  
Tel: (400) 003-1922  
[james.smith@uni.edu](mailto:james.smith@uni.edu)

**CLUSTER OVERVIEW**

<b>Uptime</b>	9 days 8 hours 31 min	<b>Memory</b>	1.2 GiB out of 8.3 GiB total
<b>Nodes</b>	2 ↑ 6 + 1 ⊖	<b>Swap</b>	0 B out of 32.7 GiB total
<b>Devices</b>	0 ↑ 1 + 0 ⊖	<b>Load</b>	0.3% user
<b>Cores</b>	3 ↑ 3 total		0.2% system
<b>Users</b>	0 out of 2 total		99.4% idle
<b>Phase Load</b>	N/A (empty)		0.1% idle
<b>Occupation Rate</b>	3.3%		

**WORKLOAD OVERVIEW**

Queue	Scheduler	#Slots	#Nodes	#Running	#Queued	#Failed	#Completed	Avg. Duration	Est. Delay
short.q	Slurm	0	256	32	45	0	482	10:07:27	00:05:16
medium.q	Slurm	0	120	5	11	0	41	32:15:00	04:16:00
long.q	Slurm	1	128	0	12	1	81	18:08:11	1d:19:11



# Management Interface

## Graphical User Interface (GUI)

- Offers administrator full cluster control
- Standalone desktop application
- Manages multiple clusters simultaneously
- Runs on Linux, Windows, *MacOS X*\*
- Built on top of Mozilla XUL engine



## Cluster Management Shell (CMSH)

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





- **Welcome**
- License
- Kernel Modules
- Hardware Info
- Nodes
- Network Architecture
- Additional Networks
- Networks
- Nameservers
- Network Interfaces
- Subnet Managers
- Installation Source
- WorkLoad Management
- Disk Layout
- Time Configuration
- Authentication
- Console
- Summary



### License Information

Version	5.1
Edition	Advanced
Name	Bright 5.1 Cluster
Organization	Bright Computing
Unit	Development
Locality	San Jose
State	California
Country	US
Serial	2158
Valid from	15 Aug 2010
Valid until	16 Nov 2010
MAC address	?:?:?:?:?:?:?:?
Licensed nodes	512

### Installation mode

- Normal (recommended)
- Express

Remote Installation

Cancel

Go Back

Continue





## Overview of installation

---

- ✓ Mounting CD/DVD-ROM
- ✓ Partitioning harddrives
- ✓ Installing Cent OS 5
- ✓ Installing distribution packages
- ✓ Installing Bright Cluster Manager packages
- ✓ Configuring kernel and setting up bootloader
- ✓ Installing Cent OS 5 software image
- ✓ Installing distribution packages to software image
- ✓ Installing Bright Cluster Manager packages to software image
- ✓ Finalizing installation
- ✓ Initializing management daemon
- ✓ Installation Complete



Automatically reboot after installation is complete

[Install Log](#)

[Reboot](#)

**Bright Cluster Manager**

File Monitoring View Help

---

**RESOURCES**

- My Clusters
  - Seismic Houston
    - Switches
      - switch01
      - switch02
      - switch03
      - switch04
      - switch05
    - Networks
      - externalnet
      - ipm net
      - mpine:
      - slavenet
      - storagenet
    - Power Distributor 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/.cm/cmgui/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
		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 demchead1
		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

Ready



**Bright Cluster Manager**

File Monitoring View Help

---

**RESOURCES** Demo Cluster

Overview Settings Failover Rackview Health Parallel shell License Notes

**My Clusters**

- ▼ Demo Cluster
  - Switches
    - switch01
    - switch02
    - switch03
    - switch04
    - switch05
  - Networks
    - externalnet
    - lpinet
    - 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

**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:** 3,02 K out of 4 K

**GPUs:** 13 out of 38

**Memory:** 732 TB out of 7.45 TB

**Users:** 13 out of 38

**CPU Usage:** 48% u 29% s 14% r 10% i

**Occupation rate:** 13.2 %

**Disk Usage**

Mountpoint	Used	Size	Use%
/	15.53 GB	37.25 GB	<div style="width: 42%;"></div>
/boot	14.31 MB	39.18 MB	<div style="width: 37%;"></div>
/home	62.6 GB	5.91 TB	<div style="width: 1.1%;"></div>

**Workload Management**

Queue	Running	Queued	Error	Completed	Avg. Duration	Est. delay
shortq	32	43	0	482	7 hours, 27 minutes	9 hours, 5 minutes
mediumq	5	11	0	4	7 days, 11 hours	4 days, 15 hours
longq	8	13	0	91	8 days, 9 hours	15 days, 13 hours

**Metric:** RunningJobs[all]

---

**EVENT VIEWER** [Icons]

All Events

Acc	Time	Cluster	Source	Message
!	18/Sep/2009 17:05:44	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 17:05:48	Demo Cluster	demohead1	Service postfix was restarted on demohead1
!	18/Sep/2009 17:05:48	Demo Cluster	demohead1	Service dhond was restarted on demohead1
!	18/Sep/2009 17:05:48	Demo Cluster	demohead1	Service mail was restarted on demohead1

Ready

# Node Provisioning

## Image based

- Regular node image is a directory on the head node
- Unlimited number of images can be created
- Software changes for the regular nodes are made inside the image(s) on the head node
- Provisioning system ensures that changes are propagated to the regular nodes

## Nodes always boot over the network

- Regular nodes PXE boot into Node Installer, which
- Identifies node (switch port or MAC based)
- Configures BMC
- Partition disks (if any) and creates file systems (if needed)
- Installs or updates software image (if needed)
- Pivot the root from NFS to the local file system

**Bright Cluster Manager**

File Monitoring View Tools Help

**RESOURCES** **node001** Demo Cluster

Overview Tasks Settings System Information Services Process Management Network Setup FS Mounts FS Exports Roles

Power: On Off Reset

Operating Systems: Shutdown Restart

Add to node group: <new> Add Remove

Software image: Update node Synchronize image  
Install node Grab to different image

Workload: Drain Undrain

Access: Root Shell Remote Console

Watch: Open Close

Misc: Locate in rack Identify node  
Provisioning Log

Health: Check <all>

**EVENT VIEWER**

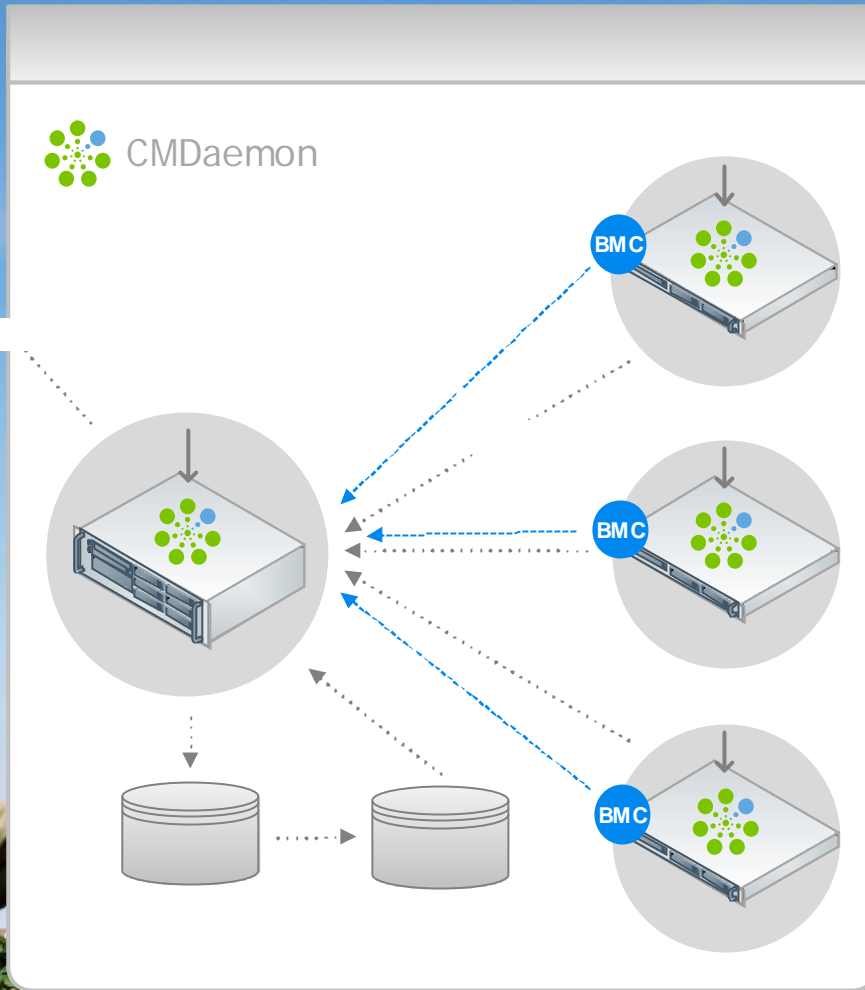
All Events

Time	Class	Name	Message
18/Sep/2008 17:01:53	Demo Cluster	demohead1	Service ntpd was restarted on demohead1
18/Sep/2008 17:01:47	Demo Cluster	demohead1	Service named was restarted on demohead1
18/Sep/2008 17:01:46	Demo Cluster	demohead1	Service postfix was restarted on demohead1
18/Sep/2008 17:01:46	Demo Cluster	demohead1	Service dhcpd was restarted on demohead1
18/Sep/2008 17:01:45	Demo Cluster	demohead1	Service mail was restarted on demohead1

Ready



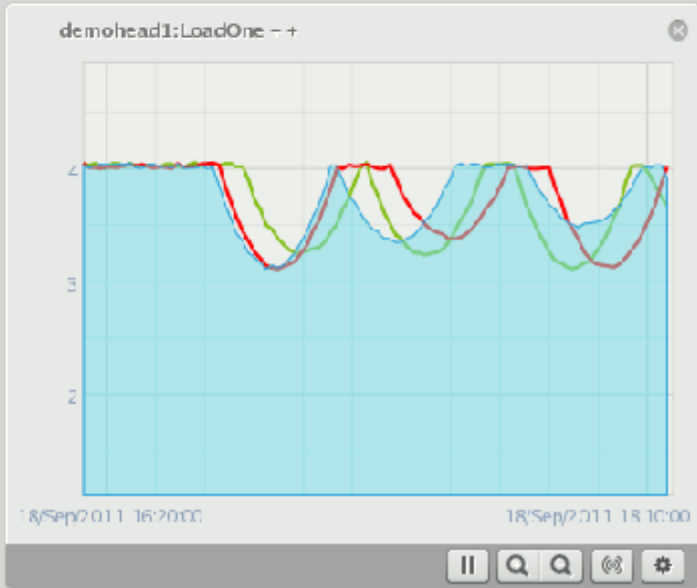
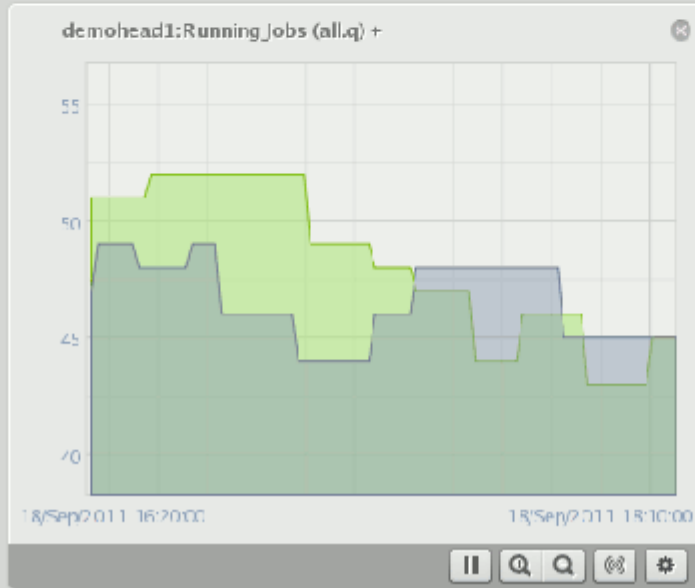
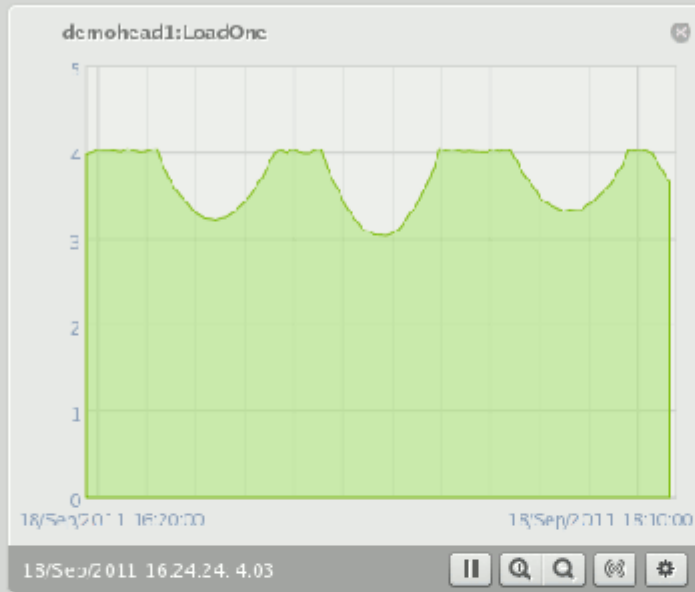
# Architecture — Monitoring



RESOURCES

- demohed1
  - CPU
  - Disk
  - Memory
    - BufferMemory (B)
    - CacheMemory (F)
    - MemoryFree (R)
    - MemoryUsed (S)
    - SwapFree (D)
    - SwapUsed (3)
  - Network
  - Operating System
    - CbtSwitches (ctx\_switch/s)
    - Forks (process/s)
    - LoadFifteen
    - LoadFive
    - LoadOne
    - ProcessCount
    - RunningProcesses
    - Uptime (s)
    - ldap
    - mysql
  - Internal
  - Workload
    - AvgExpFactor
    - AvgJobDuration[cefq] (s)
    - CompletedJobs[cefq]
    - EstimatedDelay[defq] (s)
    - FailedJobs[defq]
    - QueuedJobs[cefq]
    - RunningJobs[defq]
    - totalJob
    - schedulers
  - Cluster
    - CPU Cores Available
    - DevicesUp
    - GPU Available
    - NetworkBytesRecv (B)
    - NetworkBytesSent (D)
    - NodesUp
    - OccupationRate (%)

Demo Cluster



RESOURCES

Seismic Houston

- 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
      - Authorization
      - Authentication

UI	Rack 1	Rack 2	Rack 3	Rack 4	Rack 5	Rack 6
21	demohead1	087	087	087		211
22		088	088	089		212
23		089	089	090		213
24		085	080	085		214
25	demohead2	086	081	085		215
26		087	082	087		216
27		088	083	089		217
28		089	084	091		218
29		090	085	093		219
30		091	086	095		220
31	001	090	086	097	160	240
32	002	091	087	099	171	241
33	003	094	080	101	173	242
34	004	095	086	103	175	243
35	005	096	070	105	177	244
36	006	097	071	107	179	245
37	007	098	072	109	181	246
38	008	099	073	111	183	247
39	009		074		185	248
40	010		075		187	249
41	011		077		189	250
42	012		079		191	251
43	013		081		193	252
44	014		083	133	195	253
45	015	090	081	135	197	254
46	016	091	082	137	199	255
47	017	092	083	139	201	256
48	018	093	084	141	203	257
49	019	094	084	143	205	258
50	020	095	085	145	207	259
51	096	096	085	145	208	260

View: [Grid Icon] [List Icon] [Refresh] [Setup]



EVENT VIEWER

All Events	Act	Time	Cluster	Source	Message
1		18/Sep/2009 17:05:53	Demo Cluster	demohead1	Service mpd was restarted on demohead1
2		18/Sep/2009 17:05:47	Demo Cluster	demohead1	Service named was restarted on demohead1
3		18/Sep/2009 17:05:45	Demo Cluster	demohead1	Service postix was restarted on demohead1
4		18/Sep/2009 17:05:45	Demo Cluster	demohead1	Service dhcpd was restarted on demohead1
5		18/Sep/2009 17:05:45	Demo Cluster	demohead1	Service mail was restarted on demohead1

Ready

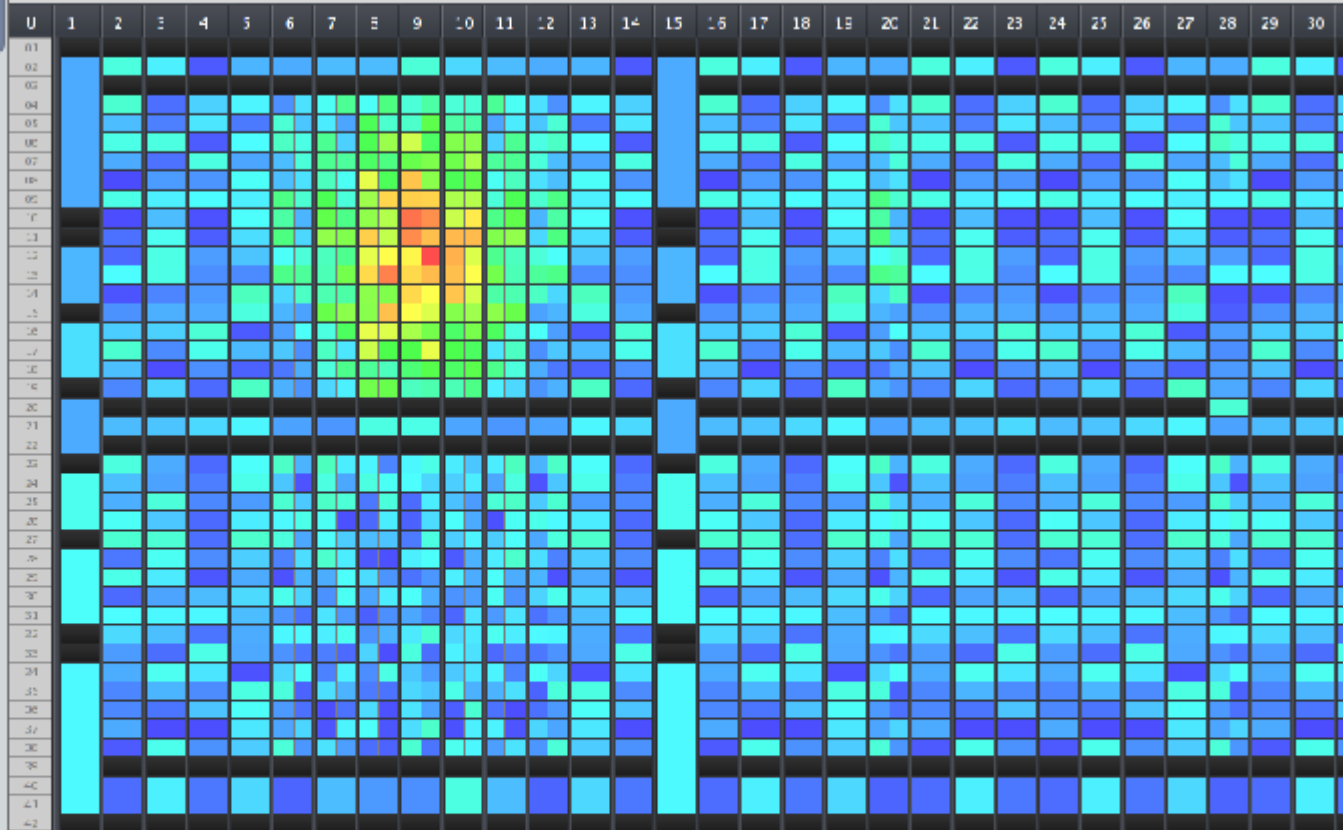


RESOURCES

- My Clusters
  - Seismic Houston
    - Switches
      - switch01
      - switch02
      - switch03
      - switch04
      - switch05
    - Networks
      - externalnet
      - pminet
      - mlinet
      - slaveret
      - 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:  Use sampling  Metric 1: Temperature 30.00 █ █ █ █ █ 69.74

EVENT VIEWER

Icon	Ack	Time	Cluster	Source	Message
	<input type="checkbox"/>	18/Sep/2008 17:05:53	Demo Cluster	demohead1	Service ntpd was restarted on demohead1
	<input type="checkbox"/>	18/Sep/2008 17:05:47	Demo Cluster	demohead1	Service named was restarted on demohead1
	<input type="checkbox"/>	18/Sep/2008 17:05:45	Demo Cluster	demohead1	Service postfix was restarted on demohead1
	<input type="checkbox"/>	18/Sep/2008 17:05:45	Demo Cluster	demohead1	Service dhcpd was restarted on demohead1
	<input type="checkbox"/>	18/Sep/2008 17:05:45	Demo Cluster	demohead1	Service mail was restarted on demohead1

# Workload Manager Integration

- Automatic installation
- Automatic configuration
- Sampling, analysis and visualization of workload manager statistics
- Consistent GUI, User Portal and CLI front-end to workload manager
- Bright cluster SOAP API provides consistent access to whole cluster, including workload manager
- Failover of workload manager
- Health checking

**RESOURCES**

- My Clusters
  - Demo Cluster
    - Switches
      - switch01
      - switch02
      - switch03
      - switch04
      - switch05
    - Networks
      - externalnet
      - ipm1net
      - mpl1net
      - 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
    - Authentication
    - Authentication

**Workload Management**

Jobs		QUEUES				
Modified	Name	Sched Job	User	Queue	Status	
	fluent	sge	jodi	med.um.q	queued	
	fluent	sge	jodi	med.um.q	queued	
	fluent	sge	jodi	med.um.q	queued	
	fluent	sge	jodi	med.um.q	running	
	gromacs	sge	alex	org.q	queued	
	gromacs	sge	alex	org.q	running	
	gromacs	sge	alex	org.q	running	
	gromacs	sge	alex	org.q	running	
	gromacs	sge	alex	org.q	running	
	gromacs	sge	alex	med.um.q	queued	
	hpc:	sge	kate	org.q	queued	
	hpc:	sge	kate	org.q	running	
	hpc:	sge	kate	org.q	running	
	magmasteel	sge	james	med.um.q	queued	
	magmasteel	sge	james	med.um.q	queued	
	magmasteel	sge	james	med.um.q	queued	
	magmasteel	sge	james	med.um.q	queued	
	magmasteel	sge	james	med.um.q	running	
	xhp	sge	mattew	short.c	running	
	xhp	sge	mattew	short.c	running	
	xhp	sge	mattew	short.c	running	

**EVENT VIEWER**

All Events	Alt	Time	Cluster	Source	Message
1		18/Sep/2009 17:05:53	Demo Cluster	demohead1	Service nmod was restarted on demohead1
1		18/Sep/2009 17:05:47	Demo Cluster	demohead1	Service nmed was restarted on demohead1
1		18/Sep/2009 17:05:45	Demo Cluster	demohead1	Service postfic was restarted on demohead1
1		18/Sep/2009 17:05:45	Demo Cluster	demohead1	Service d-rcpd was restarted on demohead1
1		18/Sep/2009 17:05:45	Demo Cluster	demohead1	Service mail was restarted on demohead1



# Cluster Health Management

- Goal: provide problem free environment for running jobs
- Four elements
  1. Cluster management automation
  2. Regular health checks
    - Actions that return PASS, FAIL or UNKNOWN
    - Can be associated with a settable severity and a message
    - Can launch an action based on any response value
  3. Prejob health checks
    - Let the workload manager hold the job very briefly
    - Check the health of each reserved node
    - If unhealthy, take the node offline, inform the system administrator
    - Let the workload manager reschedule the job to a different set of nodes
  4. Hardware stability & performance tests
    - Very wide range of tests
    - May include disk overwrites and reboot(s)
- All elements above are configurable and extensible



# Bright Cluster Manager for GPGPU

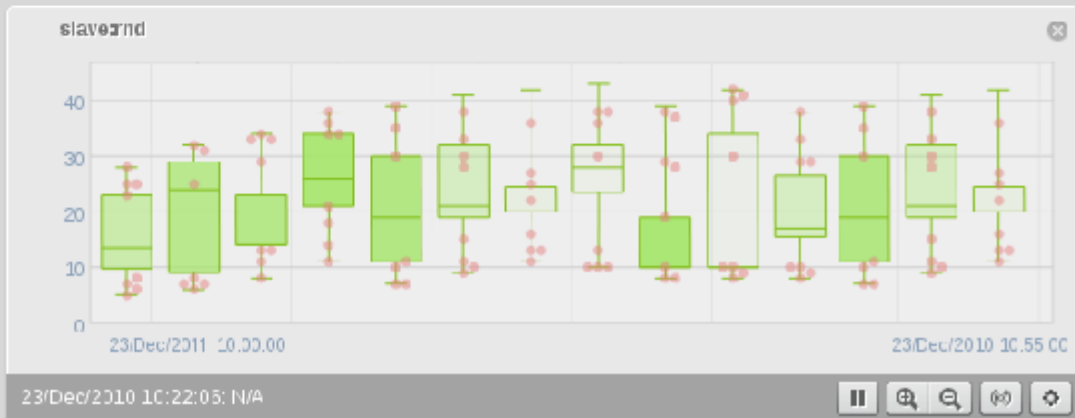
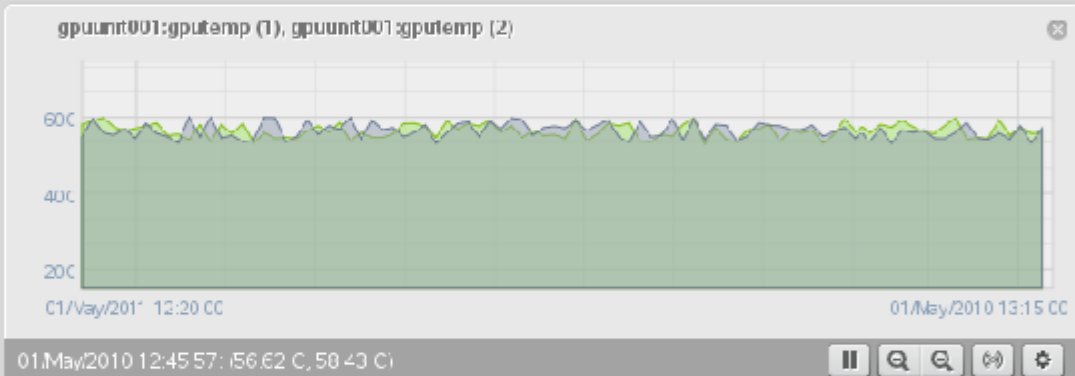
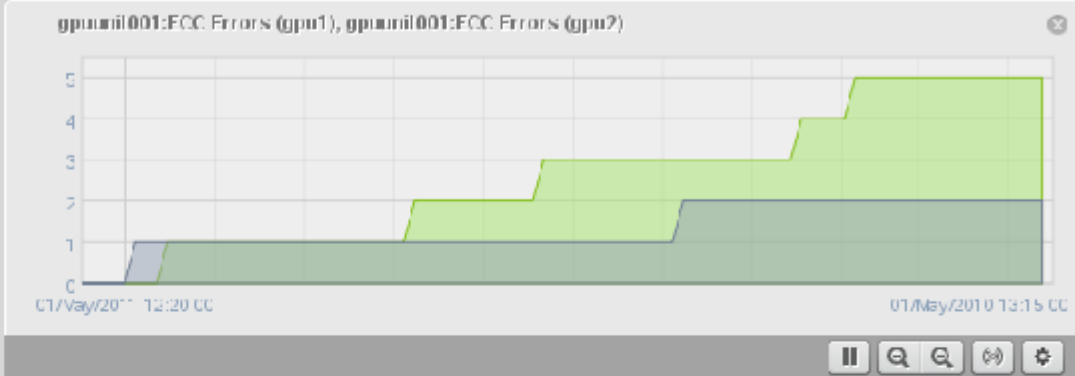
- CUDA & OpenCL redistribution rights
- Current and previous versions of CUDA & OpenCL
- Easy switching between CUDA & OpenCL versions
- CUDA driver automatically compiled at boot time
- Support for all NVIDIA GPUs



RESOURCES

- switch-04
- switch-05
- ibsw-tn-03
- ibsw-tn-04
- ibsw-tn-05
- ibsw-tn-06
- ibsw-tn-07
- ibsw-tn-08
- ibsw-tn-09
- ibsw-tn-10
- ibsw-tn-11
- ibsw-tn-12
- acc-01
- acc-02
- acc-03
- acc-04
- acc-05
- val-01
- val-02
- gpus-ni001
  - GPU
    - ECC Errors(gpu1)
    - ECC Errors(gpu2)
  - Environmental
    - gpu-temp[1] (C)
    - gpu-temp[2] (C)
    - gpu-temp[3] (C)
    - gpu-temp[4] (C)
- Operating System
- Internal
- Misc
- gpus-ni002
- gpus-ni003
- gpus-ni004
- gpus-ni005
- gpus-ni005
- gpus-ni007
- gpus-ni008
- gpus-ni009
- gpus-ni010
- gpus-ni011
- gpus-ni002
- gpus-ni003

GPU Demo Cluster





# The Future

## Cloud bursting I





# The Future

## Cloud bursting II





# The Bright Advantage

## Productivity & Efficiency

1. Easy to learn and use
2. Installation in less than 30 minutes
3. Full insight in and control over the cluster
4. All elements of the cluster are managed (servers, switches, networks, etc.)
5. Flexible provisioning (incremental, live, diskfull, diskless, IB-only, node discovery)
6. Comprehensive monitoring (graphs & rackview)
7. Powerful automation (thresholds, alerts, actions)
8. Vendor-independent workload manager integration
9. Integrated application development environment
10. Multi-cluster functionality
11. Easy, automatic updating from Linux & Bright repositories
12. Comprehensive GPU support
13. Rapid SMP deployment



# The Bright Advantage

## Uptime

1. Built-in support for unattended, reliable head node failover
2. Comprehensive cluster health checking framework
3. Powerful burn-in environment

## Performance

1. Single light-weight daemon
2. Daemons are optimized and synchronized

## Compliance & compatibility

1. Intel Cluster Ready
2. Audited by DICE and several customer (e.g. DoD, Pharma's)
3. Based on standard Linux distributions and kernels
4. Drivers included for most major hardware brands
5. Tried and tested for full compatibility with many ISV applications



# The Bright Advantage

## Scalability

1. Off-loadable provisioning
2. Efficient collection and processing of monitoring metrics
3. Tried & tested on largest clusters in the world

## Security

1. Automated security and other updates from PGP signed repositories
2. All internal + external communication encrypted using public/private key encryption through SSH/SSL
3. Authentication based on X509 certificates
4. Role-based access control
5. Auditing of all administrator write actions
6. Firewalls
7. Secure LDAP