



# Tools and Techniques for Managing Clusters for SciDAC Lattice QCD at Fermilab

D Holmgren, R Rechenmacher, A Singh, S Epsteyn

Amitoj Singh      [amitoj@fnal.gov](mailto:amitoj@fnal.gov)

Fermi National Accelerator Laboratory, Batavia, IL





# FNAL SciDAC Lattice QCD clusters



80 node Pentium III cluster



176 node Xeon cluster



# Introduction

---

## # Tools

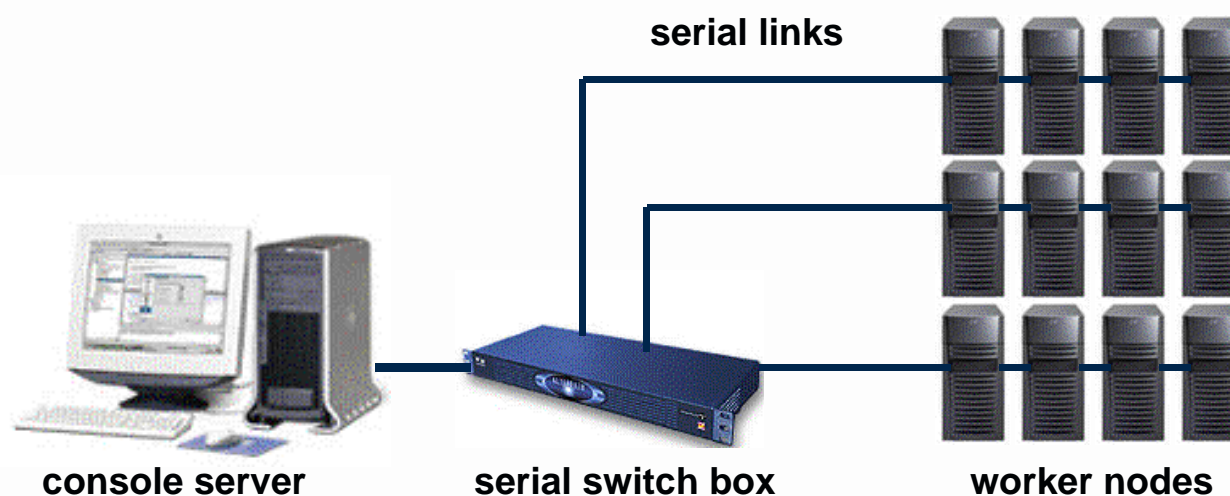
- n for hardware management tasks.
- n for OS installation/upgrade and reloading BIOS/firmware.
- n tools that work in conjunction with the PBS batch queue system.

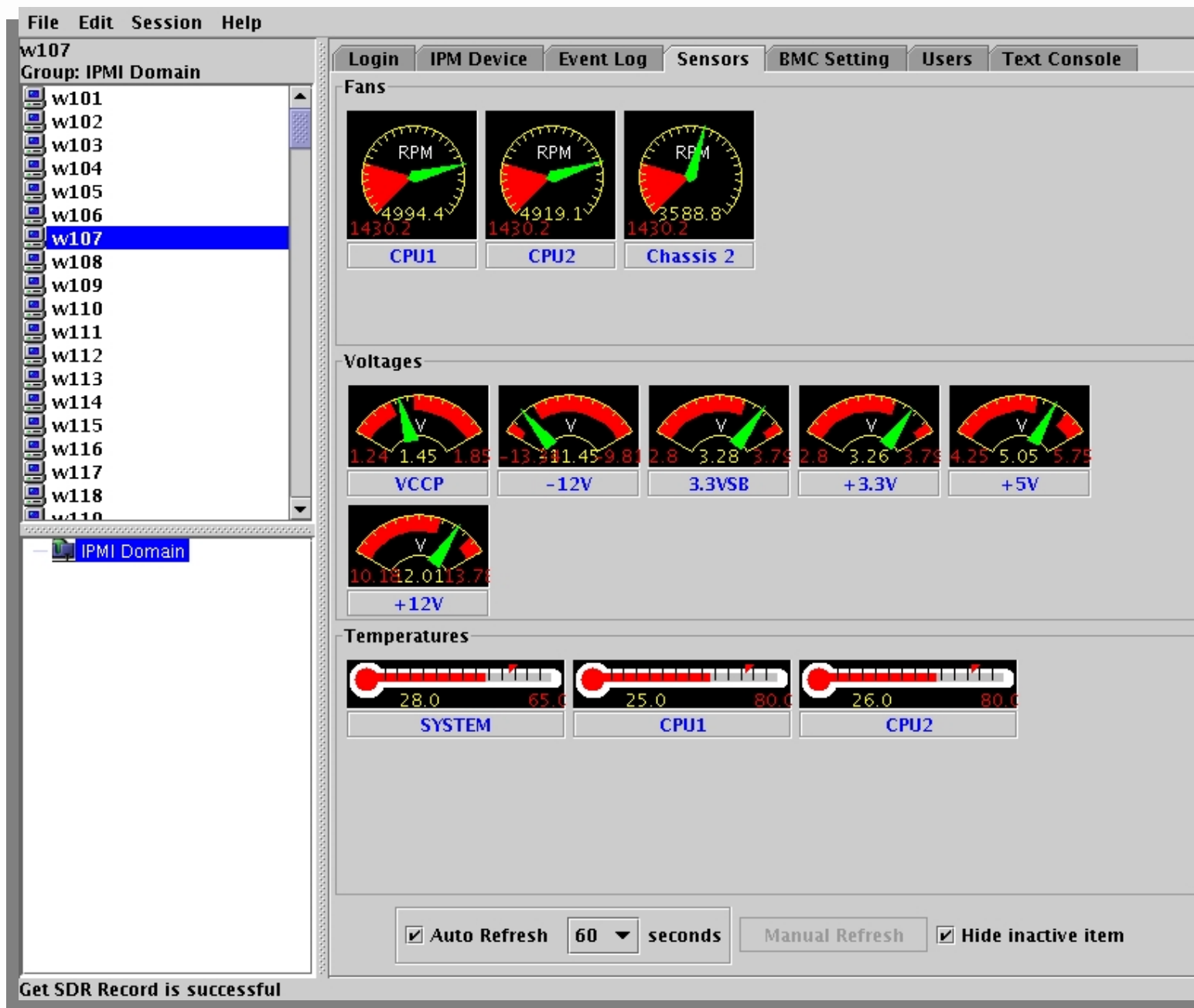
## # Integration of tools for remote administration.



# Remote node management

- # serial links to each worker node.
- # BIOS/console redirection.
- # support both IPMI 0.9 and 1.5 versions.
- # IPMI – remote power on, power off, reset.



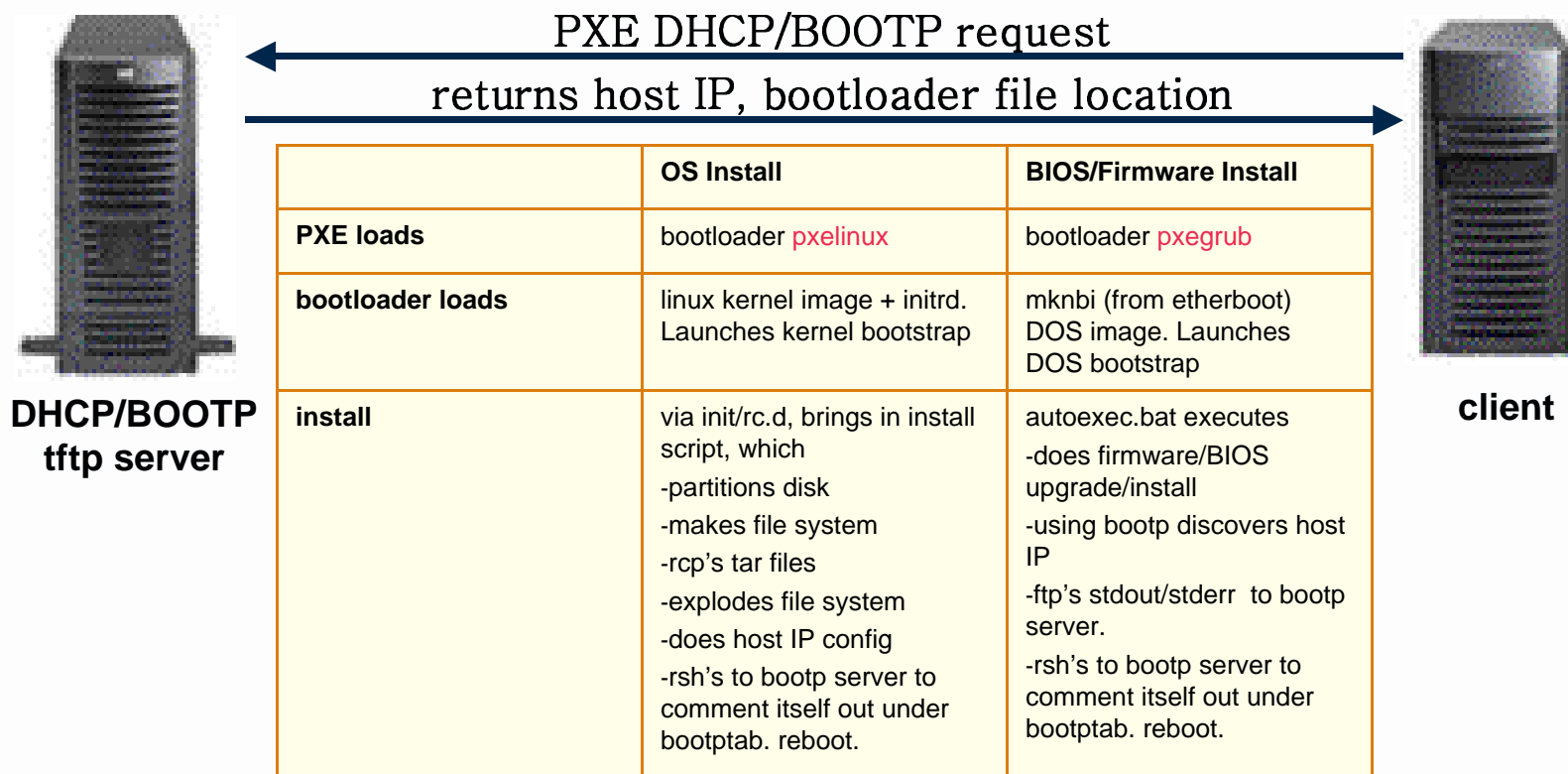


Super Micro's **IPMI View** – GUI based management Interface over LAN



# Network boot

## # PXE (Pre-boot Execution Environment)





# Fermi Tools

# **rgang** – (milc - minimum level of complexity) execute the same command on all of the nodes. Coded in python.

n two modes :

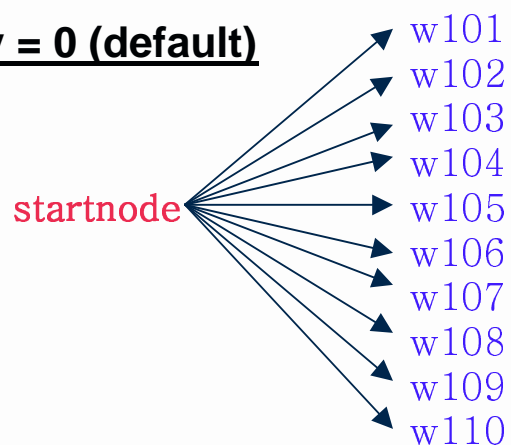
n command mode.

n copy mode.

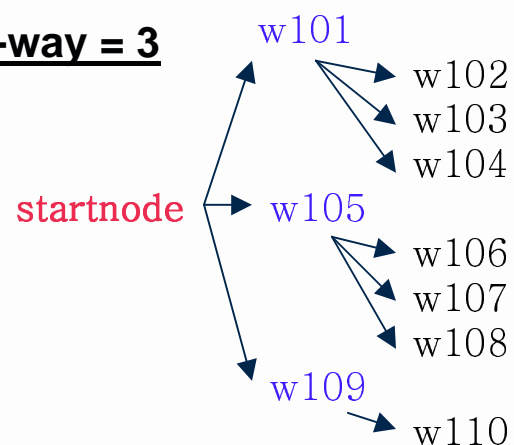
n n-way option.

n rgang can be used to install itself.

**n-way = 0 (default)**



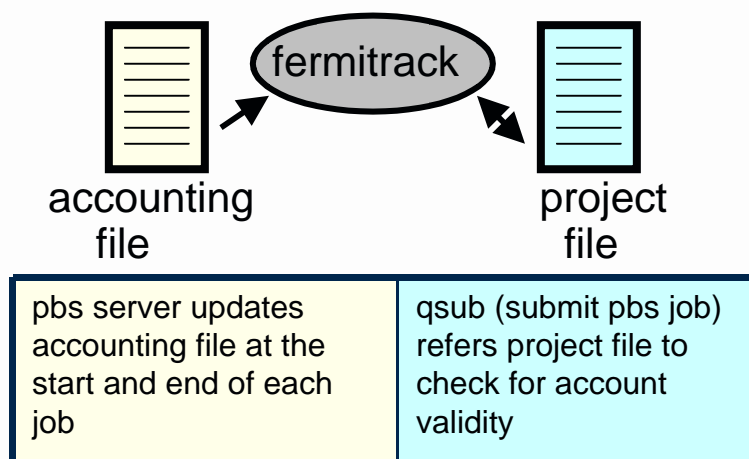
**n-way = 3**





# Fermi Tools

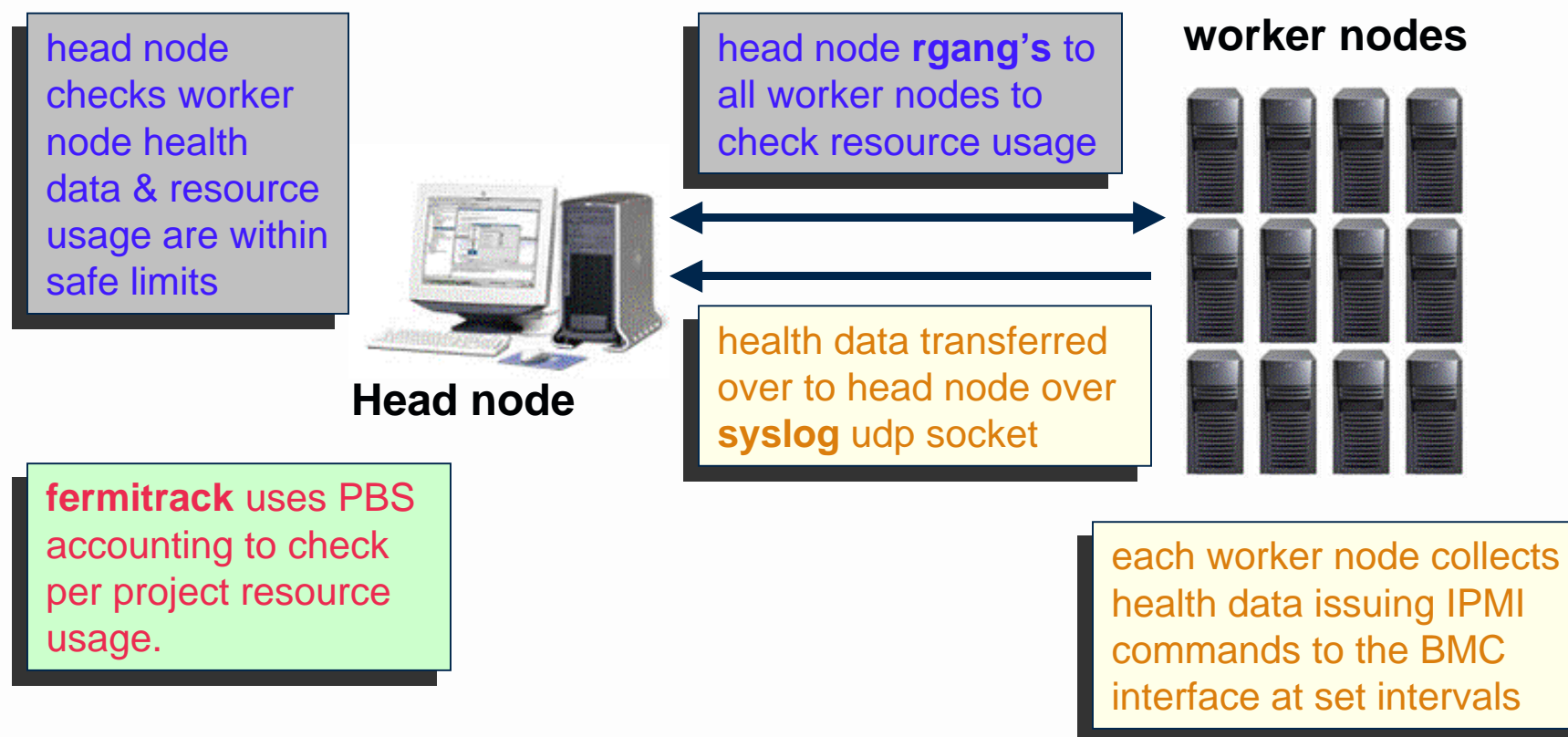
- # **fermistat** - list cluster resource usage by users and running jobs. Works in conjunction with PBS batch queue system.
- # **fermitrack** – poor man's project accounting. Works in conjunction with PBS batch queue accounting system.  
> qsub -A "myproject" -l nodes=n mypbsscript







# Integration-How it all comes together.



http://lqcd.fnal.gov



File Edit View Go Communicator Help

### lqcd.fnal.gov Node Map

w101	w102	w103	w104	w105	w106	w107	w108	w109	w110	w111
w112	w113	w114	w115	w116	w117	w118	w119	w120	w121	w122
w123	w124	w201	w202	w203	w204	w205	w206	w207	w208	w209
w210	w211	w212	w213	w214	w215	w216	w217	w218	w219	w220
w221	w222	w223	w224	w301	w302	w303	w304	w305	w306	w307
w308	w309	w310	w311	w312	w313	w314	w315	w316	w317	w318
w319	w320	w321	w322	w323	w324	w401	w402	w403	w404	w405
w406	w407	w408	w409	w410	w411	w412	w413	w414	w415	w416
w417	w418	w419	w420	w421	w422	w423	w424	w501	w502	w503
w504	w505	w506	w507	w508	w509	w510	w511	w512	w513	w514
w515	w516	w517	w518	w519	w520	w521	w522	w523	w524	w601
w602	w603	w604	w605	w606	w607	w608	nqcd0101	nqcd0102	nqcd0103	nqcd0104
nqcd0105	nqcd0106	nqcd0201	nqcd0202	nqcd0203	nqcd0204	nqcd0205	nqcd0206	nqcd0301	nqcd0302	nqcd0303
nqcd0304	nqcd0305	nqcd0306	nqcd0401	nqcd0402	nqcd0403	nqcd0404	nqcd0405	nqcd0406	nqcd0501	nqcd0502
nqcd0503	nqcd0504	nqcd0505	nqcd0506	nqcd0601	nqcd0602	nqcd0603	nqcd0604	nqcd0605	nqcd0606	nqcd0701
nqcd0702	nqcd0703	nqcd0704	nqcd0705	nqcd0706	nqcd0801	nqcd0802	nqcd0803	nqcd0804	nqcd0805	nqcd0806

free : 31 down : 0 offline : 1 reserve : 0 job-exclusive : 144 job-sharing : 0 Usage : 82%

### Job List

Ref Id	Job Id	Job Name	User	Time Use	S	Queue	Nodes
5	846799	falpha4.sh	trottier	13:35	R	workq	32
4	846807	par_4sc.sh	cdavies	01:02	R	workq	16
1	846801	falpha3.sh	trottier	13:32	R	workq	32
3	846800	falpha1.sh	trottier	13:34	R	workq	32
2	846797	falpha2.sh	trottier	13:37	R	workq	32

Status **BUSY**

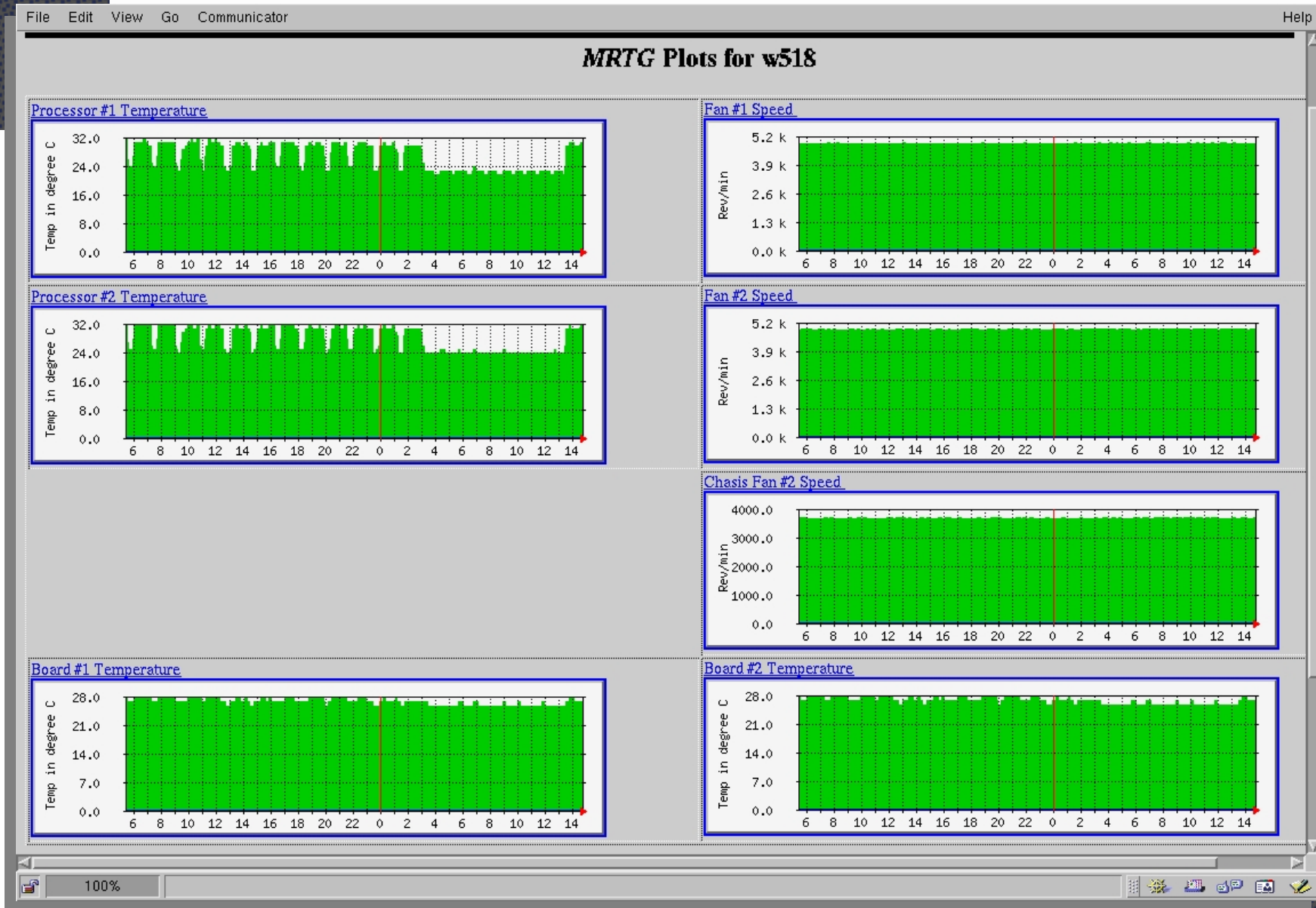
100%

http://lqcd.fnal.gov



all Health Statistics								
Hostname	Board Tmp1	Board Tmp2	Proc Tmp1	Proc Tmp2	Fan Speed1	Fan Speed2	Chasis Fan2	Warn Code
<a href="#">nqcd0101</a>	42	42	36	38	4819	4794	9999	0
<a href="#">nqcd0102</a>	41	41	37	38	4794	4794	9999	0
<a href="#">nqcd0103</a>	45	45	38	39	4794	4819	9999	0
<a href="#">nqcd0104</a>	45	45	39	40	4819	4844	9999	0
<a href="#">nqcd0105</a>	45	45	38	40	4869	4844	9999	0
<a href="#">nqcd0106</a>	44	44	38	38	4869	4844	9999	0
<a href="#">nqcd0201</a>	37	37	32	33	4819	4819	9999	0
<a href="#">nqcd0202</a>	39	39	34	35	4794	4844	9999	0
<a href="#">nqcd0203</a>	39	39	34	35	4794	4819	9999	0
<a href="#">nqcd0204</a>	40	40	35	35	4794	4794	9999	0
<a href="#">nqcd0205</a>	40	40	35	36	4794	4844	9999	0
<a href="#">nqcd0206</a>	39	39	34	35	4819	4794	9999	0
<a href="#">nqcd0301</a>	35	35	29	29	4718	4794	9999	0
<a href="#">nqcd0302</a>	35	35	30	31	4794	4768	9999	0
<a href="#">nqcd0303</a>	35	35	31	31	4743	4794	9999	0
<a href="#">nqcd0304</a>	37	37	31	32	4794	4794	9999	0
<a href="#">nqcd0305</a>	45	45	49	46	4819	4819	9999	0
<a href="#">nqcd0306</a>	44	44	48	47	4794	4794	9999	0
<a href="#">nqcd0401</a>	64	64	63	64	4869	4894	9999	0
<a href="#">nqcd0402</a>	67	67	68	67	4919	4844	9999	0
<a href="#">nqcd0403</a>	67	67	67	67	4869	4869	9999	0
<a href="#">nqcd0404</a>	64	64	66	66	4919	4844	9999	0
<a href="#">nqcd0405</a>	66	66	68	67	4869	4869	9999	0
<a href="#">nqcd0406</a>	64	64	64	65	4869	4894	9999	0
<a href="#">nqcd0501</a>	59	59	58	59	4844	4844	9999	0
<a href="#">nqcd0502</a>	61	61	60	61	4919	4819	9999	0
<a href="#">nqcd0503</a>	0	0	0	0	0	0	9999	
<a href="#">nqcd0504</a>	59	59	61	61	4844	4844	9999	0
<a href="#">nqcd0505</a>	62	62	63	63	4844	4894	9999	0
<a href="#">nqcd0506</a>	59	59	58	59	4844	4844	9999	0
<a href="#">nqcd0601</a>	54	54	54	53	4844	4819	9999	0

<http://lqcd.fnal.gov>





# Example

Ø **fermistat -l <jobid> | rgang - <command>**

Ø **fermistat -l 8345.job | rgang - uptime**

```
----- w605 -----  
rsh w605 'uptime'  
2:21pm up 9 days, 2:02, 0 users, load average: 1.99, 1.97, 1.91  
----- w606 -----  
rsh w606 'uptime'  
2:21pm up 9 days, 2:03, 0 users, load average: 1.99, 1.97, 1.91  
----- w607 -----  
rsh w607 'uptime'  
2:21pm up 9 days, 2:02, 0 users, load average: 1.99, 1.97, 1.91  
----- w608 -----  
rsh w608 'uptime'  
2:21pm up 9 days, 2:02, 0 users, load average: 2.00, 1.97, 1.91
```



# Example

Ø **fermistat -c <rgang style list-of-nodes>**

Ø **fermistat -c w1{01-04}**

pbsnodes -c w101

pbsnodes -c w102

pbsnodes -c w103

pbsnodes -c w104

Ø **fermistat -l <jobid> | fermistat -o -**

Ø **fermistat -l 8345.job | fermistat -o -**

pbsnodes -o w605

pbsnodes -o w606

pbsnodes -o w607

pbsnodes -o w608





## Conclusion

---

It works . . . .

<http://qcdhome.fnal.gov> - 80 node Pentium III cluster

<http://lqcd.fnal.gov> - 176 node Xeon cluster

<http://fermitools.fnal.gov> - rgang and other tools

Don Holmgren  
[djholm@fnal.gov](mailto:djholm@fnal.gov)

Ron  
Rechenmacher  
[ron@fnal.gov](mailto:ron@fnal.gov)

Amitoj Singh  
[amitoj@fnal.gov](mailto:amitoj@fnal.gov)

Simon Epsteyn  
[seva@fnal.gov](mailto:seva@fnal.gov)