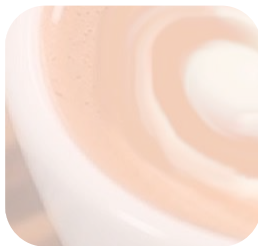
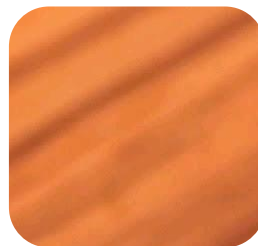


**CRAY**



**How-to write a xtpmd\_plugin for your  
Cray® XC™ system**

Steven J. Martin ([stevem@cray.com](mailto:stevem@cray.com))

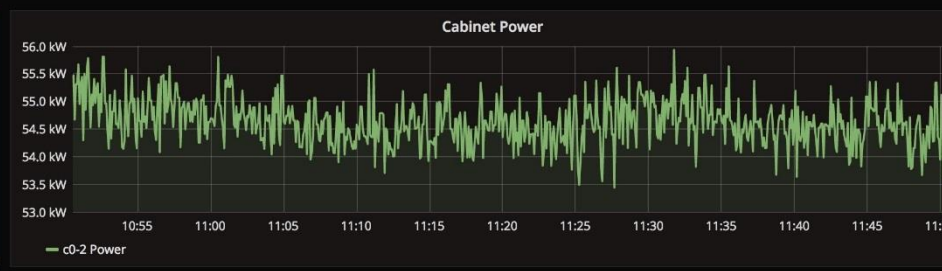
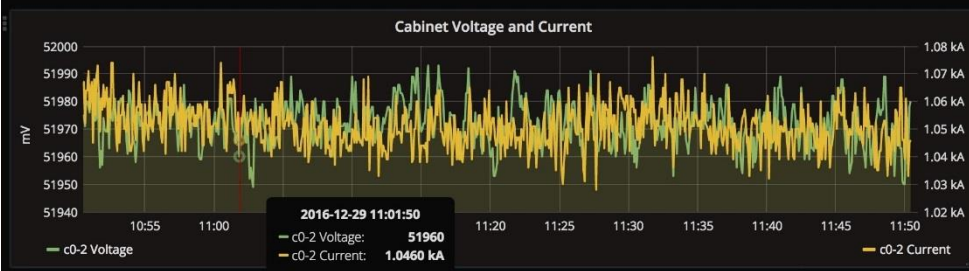
**CUG 2017. CAFFEINATED COMPUTING**

Redmond, Washington May 7-11, 2017

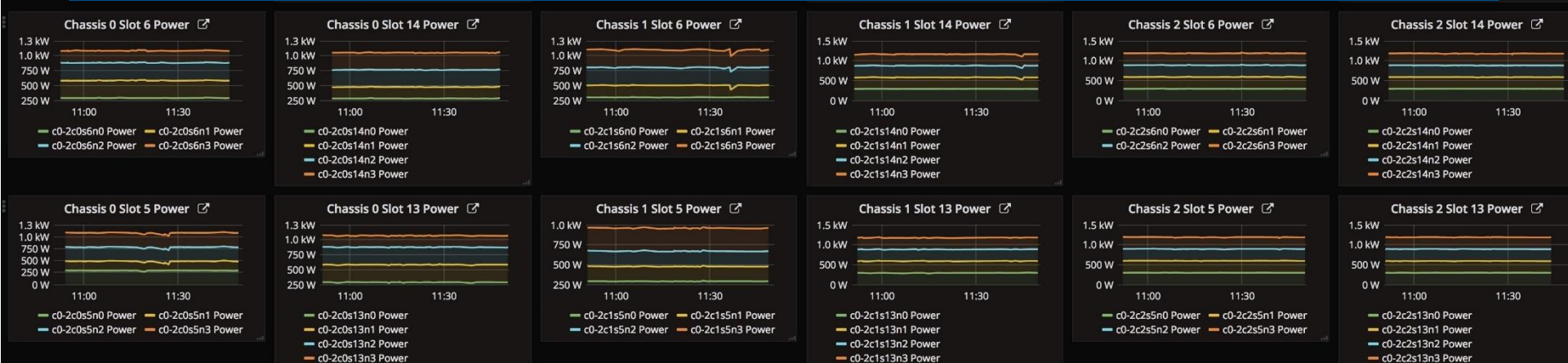
# Cray XC Telemetry Plugin Introduction

- **Enabling sites to get telemetry data off the Cray**
- **Plugin interface enables site specific customization**
- **This presentation and paper dive into details**

- I plan to move through the slides quickly
- And have time for question and discussion at the end



Why do a plugin, what kind of things might the data enable?



# Customer Driven Plugin History

- **Xtreme System Monitoring Collaboration meeting, Jan 28, 2016**
  - Requested access to SEDC and Cray high-speed power and energy data streams before data is injected into the Cray PMDB
- Cray previewed the plugin feature at CUG 2016 [Slides 59-74](#)
- Plugin using [Redis](#) Pub/Sub transport ([Hiredis](#) C-library) demonstrated at Trinity phase-2 factory trial
- Example plugin code ships with (SMW) 8.0 and newer software
- NERSC now using plugin in production on [Cori](#) system

# Streaming Data Available From Five Sources

- **Cabinet SEDC data**

- xtgetsedcvalues -l -t cc -c c0-0

- **Blade SEDC data**

- xtgetsedcvalues -l -t bc -c c1-0c1s3

- **Cabinet power and energy data**

- xtgetsedcvalues -l -t cc\_power -c c0-0

- **Blade power and energy data**

- xtgetsedcvalues -l -t bc\_power | grep c0-0c0s8

- **Job and application data**

Table II  
CC AND BC PMDB DATA

ID	Sensor Description	Unit
0	Cabinet Power	W
1	Cabinet Energy	J
2	Cabinet Voltage	mV
3	Cabinet Current	A
8	Cabinet Blower Power	W
16	HSS Power	W
17	HSS Energy	J
32	Node 0 Power	W
33	Node 0 Energy	J
36	Node 0 CPU Power	W
37	Node 0 CPU Energy	J
40	Node 1 Power	W
41	Node 1 Energy	J
44	Node 1 CPU Power	W
45	Node 1 CPU Energy	J
48	Node 2 Power	W
49	Node 2 Energy	J
52	Node 2 CPU Power	W
53	Node 2 CPU Energy	J
56	Node 3 Power	W
57	Node 3 Energy	J
60	Node 3 CPU Power	W
61	Node 3 CPU Energy	J
64	Node 0 Memory Power	W
68	Node 0 Memory Energy	J
76	Node 1 Memory Power	W
77	Node 1 Memory Energy	J
84	Node 2 Memory Power	W
85	Node 2 Memory Energy	J
92	Node 3 Memory Power	W
93	Node 3 Memory Energy	J

Table I  
BC SEDC DATA

ID	Sensor Description	Unit
1257	BC_I_ARIES_TEMP	degC
1300	BC_I_NODE0_CPU_TEMP	degC
1301	BC_I_NODE0_CPU_TEMP	degC
1308	BC_I_NODE0_CPU_CHI_DRM00	degC
1312	BC_I_NODE0_CPU_CHI_DRM00	degC
1636	BC_I_NODE50_VEM_TEMP	degC
1637	BC_I_NODE50_VEM_TEMP	degC
1796	BC_I_NODE0_PCH_THERMAL	degC
1200	BC_V_VDD0_0V	V
1201	BC_V_VDD0_1V	V
1202	BC_V_VDD0_1V_PLR_1V	V
1203	BC_V_VDD0_1V_GTP	V
1204	BC_V_VDD0_1V_BSS	V
1206	BC_V_VDD0_1V_BSS	V
1207	BC_V_VDD0_1V_BSS	V
1208	BC_V_VDD0_3V_PDC	V
1210	BC_V_VDD0_3V_BSS	V
1211	BC_V_VDD0_3V_MCRODA	V
1213	BC_V_VDD0_0V	V
1216	BC_V_VDD0_0V_STDBY	V
1218	BC_V_ARIES_VDD_VCORE	V
1219	BC_V_ARIES_VDD_1V0	V
1260	BC_V_ARIES_VDD_1V0	V
1261	BC_V_ARIES_VDD_1V0	V
1262	BC_V_ARIES_VDD_1V0	V
1263	BC_V_ARIES_VDD_1V0	V
1264	BC_V_ARIES_VDD_1V0	V
1265	BC_V_ARIES_VDD_1V0	V
1266	BC_V_ARIES_VDD_1V0	V
1267	BC_V_ARIES_VDD_1V0	V
1268	BC_V_ARIES_VDD_1V0	V
1269	BC_V_ARIES_VDD_1V0	V
1270	BC_V_ARIES_VDD_1V0	V
1271	BC_V_ARIES_VDD_1V0	V
1272	BC_V_ARIES_VDD_1V0	V
1273	BC_V_ARIES_VDD_1V0	V
1274	BC_V_ARIES_VDD_1V0	V
1275	BC_V_ARIES_VDD_1V0	V
1276	BC_V_ARIES_VDD_1V0	V
1277	BC_V_ARIES_VDD_1V0	V
1278	BC_V_ARIES_VDD_1V0	V
1468	BC_P_NODE0_CPU_CHI_DRAM_ACC	J
1469	BC_P_NODE0_CPU_CHI_DRAM_ACC	J
1470	BC_P_NODE0_CPU_CHI_DRAM_ACC	J
1471	BC_P_NODE0_CPU_CHI_DRAM_ACC	J
1500	BC_P_NODE0_CPU_VCC_ACC	J
1516	BC_P_NODE0_CPU_VCC_ACC	J
1716	BC_P_NODE0_GLOBAL_PDC_POWER	W
1266	BC_H_ARIES_VEM_FLT_BITS	status
1267	BC_H_ARIES_VEM_FLT_BITS	status
1268	BC_H_ARIES_VEM_FLT_BITS	status
1273	BC_H_NODE0_VOC_ECB_FAULT	status
1285	BC_H_BB_VERT_IVOCC_ECB_FAULT	status
1286	BC_H_BB_VERT_IVOCC_ECB_FAULT	status
1844	BC_I_NODE0_CPU_MEM_THROTTLE	%
1845	BC_I_NODE0_CPU_MEM_THROTTLE	%
1872	BC_I_NODE0_CPU_CPU_THROTTLE	N
1873	BC_I_NODE0_CPU_CPU_THROTTLE	N

# Job and Application Data

- Job and application information

- Job ID, app ID, user ID, timestamp
- Assigned nodes (Cray "nid")

```
enum {
```

```
APP_CMD_TYPE_START = 1,
```

```
APP_CMD_TYPE_END = 2,
```

```
APP_CMD_TYPE_SYNC = 3,
```

```
APP_CMD_TYPE_SUSPEND = 4,
```

```
APP_CMD_TYPE_RESUME = 5,
```

```
JOB_CMD_TYPE_START = 6,
```

```
JOB_CMD_TYPE_END = 7,
```

```
JOB_CMD_TYPE_SUSPEND = 8,
```

```
JOB_CMD_TYPE_RESUME = 9
```

```
};
```

Used by ALPS and Slurm

Used by ALPS at reservation time

```
{
  "ts": "2017-03-23T16:04:28.366738Z",
  "event": "APP_START",
  "userid": 27216,
  "job_id": "1723832.sdb",
  "apid": 3801382,
  "nid_count": 2,
  "nid_cname_array": [
    { "nid": 56, "cname": "c0-0c0s14n0" },
    { "nid": 57, "cname": "c0-0c0s14n1" }
  ]
}
```

# ALPS Example (Raw Format)

- **Job and application information**

- Job ID, app ID, user ID, timestamp
- Assigned nodes (or Cray “nid”)

```
enum {  
  APP_CMD_TYPE_START = 1,  
  APP_CMD_TYPE_END = 2,  
  APP_CMD_TYPE_SYNC = 3,  
  APP_CMD_TYPE_SUSPEND = 4,  
  APP_CMD_TYPE_RESUME = 5,  
  JOB_CMD_TYPE_START = 6,  
  JOB_CMD_TYPE_END = 7,  
  JOB_CMD_TYPE_SUSPEND = 8,  
  JOB_CMD_TYPE_RESUME = 9  
};
```

```
ts=1490734709583873,event=6,userid=1205,jobid=3880.sdb,apid=0,nids=',2'  
ts=1490734713023993,event=1,userid=1205,jobid=3880.sdb,apid=706691,nids=',2'  
ts=1490734718511193,event=2,userid=1205,jobid=3880.sdb,apid=706691  
ts=1490734719788394,event=7,userid=1205,jobid=3880.sdb,apid=0
```

# ALPS Example (JSON Format)

```
{
  "ts": "2017-03-31T20:12:32.601051Z",
  "event": "JOB_START",
  "userid": 7821,
  "job_id": "1745267.sdb",
  "apid": 0,
  "nid_count": 1,
  "nid_cname_array": [
    {
      "nid": 76,
      "cname": "c0-0c1s3n0"
    }
  ]
}

{
  "ts": "2017-03-31T20:12:39.043667Z",
  "event": "APP_START",
  "userid": 7821,
  "job_id": "1745267.sdb",
  "apid": 3860100,
  "nid_count": 1,
  "nid_cname_array": [
    {
      "nid": 76,
      "cname": "c0-0c1s3n0"
    }
  ]
}

{
  "ts": "2017-03-31T20:12:50.091905Z",
  "event": "APP_END",
  "userid": 7821,
  "job_id": "1745267.sdb",
  "apid": 3860100
}

{
  "ts": "2017-03-31T20:12:51.102143Z",
  "event": "JOB_END",
  "userid": 7821,
  "job_id": "1745267.sdb",
  "apid": 0
}
```



# Slurm Example



```
system:~> salloc -N 4  
salloc: Granted job allocation 369936  
salloc: Waiting for resource configuration  
salloc: Nodes nid00[122-125] are ready for job  
system:~> srun -N 2 hostname  
nid00123  
nid00122  
system:~> srun -N 4 hostname  
nid00124  
nid00125  
nid00123  
nid00122  
system:~> exit  
exit  
salloc: Relinquishing job allocation 3
```

1

```
{"ts":"2017-04-20T19:56:43.595141Z",  
  "event":"APP_START", "userid":26914, "job_id":"369936",  
  "apid":6056184802581266704, "nid_count":4,  
  "nid_cname_array":[{"nid":122,"cname":"c0-0c1s14n2"},  
                    {"nid":123,"cname":"c0-0c1s14n3"},  
                    {"nid":124,"cname":"c0-0c1s15n0"},  
                    {"nid":125,"cname":"c0-0c1s15n1"}]}
```

2

```
{"ts":"2017-04-20T19:57:10.623100Z", "event":"APP_START",  
  "userid":26914, "job_id":"369936",  
  "apid":369936, "nid_count":2,  
  "nid_cname_array":[{"nid":122,"cname":"c0-0c1s14n2"},  
                    {"nid":123,"cname":"c0-0c1s14n3"}]}
```

```
{"ts":"2017-04-20T19:57:11.543174Z", "event":"APP_END",  
  "userid":26914, "job_id":"369936", "apid":369936}
```

COMP

# Slurm Example

```
system:~> salloc -N 4
salloc: Granted job allocation 369936
salloc: Waiting for resource configuration
salloc: Nodes nid00[122-125] are ready for job
system:~> srun -N 2 hostname
nid00123
nid00122
system:~> srun -N 4 hostname
nid00124
nid00125
nid00123
nid00122
system:~> exit
exit
salloc: Relinquishing job allocation 3
```

3

```
{"ts":"2017-04-20T19:57:30.087104Z","event":"APP_START",
  "userid":26914, "job_id":"369936",
  "apid":10000369936, "nid_count":4,
  "nid_cname_array":[{"nid":122,"cname":"c0-0c1s14n2"},
                    {"nid":123,"cname":"c0-0c1s14n3"},
                    {"nid":124,"cname":"c0-0c1s15n0"},
                    {"nid":125,"cname":"c0-0c1s15n1"}]}

{"ts":"2017-04-20T19:57:31.007176Z","event":"APP_END",
  "userid":26914, "job_id":"369936",
  "apid":10000369936}
```

4

```
{"ts":"2017-04-20T19:57:35.783165Z","event":"APP_END",
  "userid":26914, "job_id":"369936",
  "apid":6056184802581266704}
```

# Slurm Example (Closer look at steps 1 and 4)

```
system:~> salloc -N 4  
salloc: Granted job allocation 369936  
salloc: Waiting for resource configuration  
salloc: Nodes nid00[122-125] are ready for job  
...  
system:~> exit  
exit  
salloc: Relinquishing job allocation 369936
```

1

```
{"ts":"2017-04-20T19:56:43.595141Z","event":"APP_START",  
  "userid":26914,"job_id":"369936","apid":6056184802581266704,  
  "nid_count":4,"nid_cname_array":[  
    {"nid":122,"cname":"c0-0c1s14n2"}, {"nid":123,"cname":"c0-0c1s14n3"},  
    {"nid":124,"cname":"c0-0c1s15n0"}, {"nid":125,"cname":"c0-0c1s15n1"}]  
  ...  
{"ts":"2017-04-20T19:57:35.783165Z","event":"APP_END",  
  "userid":26914,"job_id":"369936","apid":6056184802581266704}
```

4



# Plugin Configuration File

- `/opt/cray/hss/default/etc/xtpmd_plugins.ini`
  - Configuration file the in smw release supports the `plugin_csv` example
  - `smw:/opt/cray/hss/default/pm/xtpmd_api`

```
smw:/opt/cray/hss/default/etc> egrep "=|[" xtpmd_plugins.ini | grep -v _path
[plugins]
# shmsize=4194304
# instances=csv;other
# instances=csv
[plugin_csv]
object=/opt/cray/hss/default/lib64/xtpmd_plugin_csv.so
log_dir=/tmp
pmdb_cc_enabled=yes
...
```



# Makefile Example

```
# Makefile for xtpmd_xjson plugin
OBJ = xtpmd_plugin_xjson.o
LIB = xtpmd_plugin_xjson.so

CFLAGS += -O3 -fPIC
CFLAGS += $(shell pkg-config --cflags glib-2.0 gthread-2.0 libpq)
LDFLAGS += $(shell pkg-config --libs glib-2.0 gthread-2.0 libpq)
LDFLAGS += -lm -lz

ALL: xtpmd_plugin_xjson.so

xtpmd_plugin_xjson.o: xtpmd_plugin_xjson.c xtpmd_plugin.h
    $(CC) $(CPPFLAGS) $(CFLAGS) -c -o $@ $<

xtpmd_plugin_xjson.so: xtpmd_plugin_xjson.o
    $(CC) -shared $< -o $@ $(LDFLAGS)
```

# Starting a Test Plugin Unsupervised



```
smw:~> ipcs | grep -v post
----- Message Queues -----
key          msqid  owner    perms  used-bytes  messages
----- Shared Memory Segments -----
key          shmid  owner    perms  bytes       nattch     status
0x7a060809  32769  crayadm  600    4194304     1
----- Semaphore Arrays -----
key          semid  owner    perms  nsems

smw:~> xtpmd_plugd 0x7a060809 4194304 plugin_xjson ./xtpmd_plugins.ini
```

# Starting a Test Plugin Unsupervised



```
smw:~> ipcs | grep -v post
```

```
----- Message Queues -----
```

```
key          msqid  owner  perms  used-bytes
```

```
----- Shared Memory Segments -----
```

```
key          shmid  owner  perms  nattch  status
```

```
0x7a060809  32769  0      0666  4194304  1
```

```
----- Semaphore Arrays -----
```

```
key          owner  perms  nsems
```

Note: Use 'Ctrl-D' to kill the test plugin, or kill -9 from another window...  
• 'Ctrl-C' is blocked.

```
smw:~> xtpmd_plugd 0x7a060809 4194304 plugin_csv ./xtpmd_plugins.ini
```

# Design Considerations Covered in the Paper

- Limiting the plugin's impact
- Library usage
- Time stamp formatting
- Translating binary fields
- Getting data off node
- Other formatting considerations

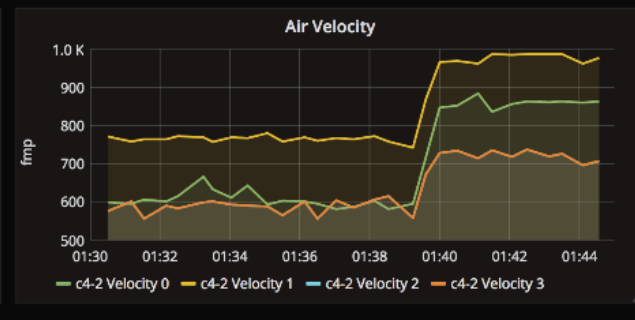
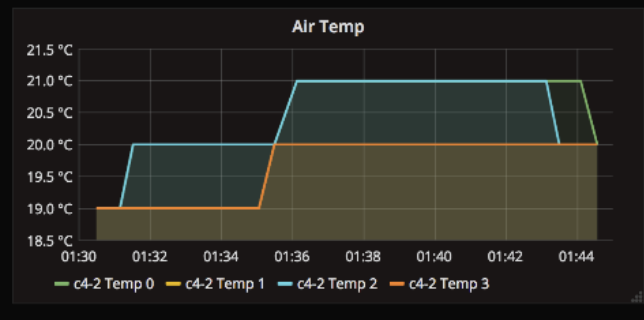
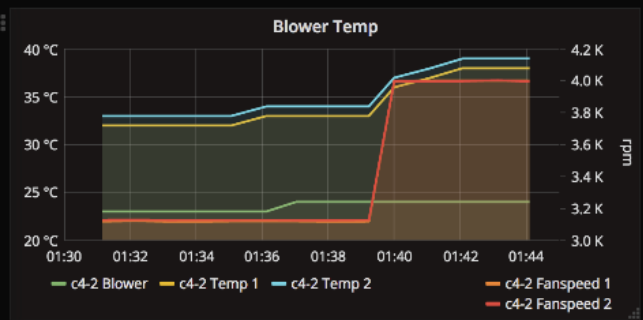


- **Plugin data sent into NERSC Center-wide Data Collect**
  - [https://cug.org/proceedings/cug2016\\_proceedings/includes/files/pap101.pdf](https://cug.org/proceedings/cug2016_proceedings/includes/files/pap101.pdf)
- **Implemented 5 plugins one for each data source**
  - instances = bp ; bc ; cp ; cc ; j o b
- **RabbitMQ is used to send data off the SMW**
  - Connection details in the configuration file
  - Changes can be made without recompiling code
- **Data written in JSON format**
  - Significantly increases [Elastic](#) ingest rate

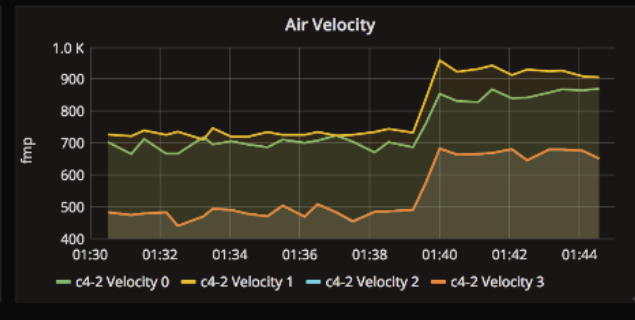
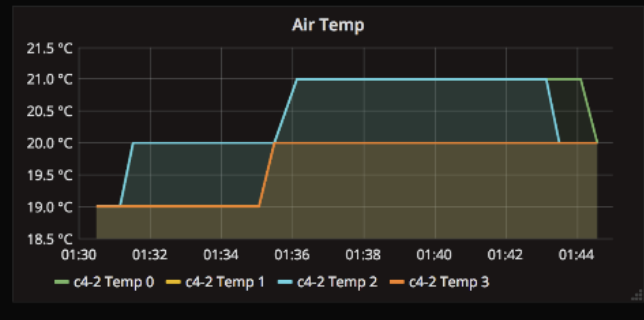
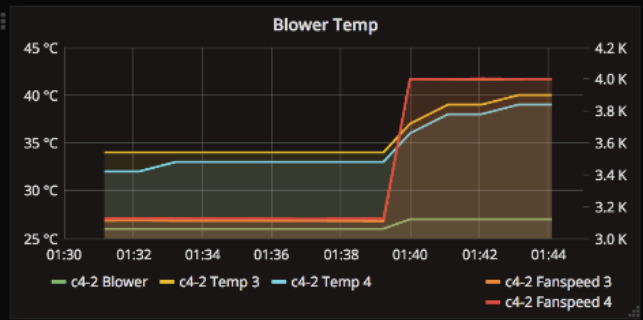
# Cori Problem Information Gathering

The next 5 slides are graphics from the paper  
**Collected** looking into conditions at the time of  
a reported thermal throttling event

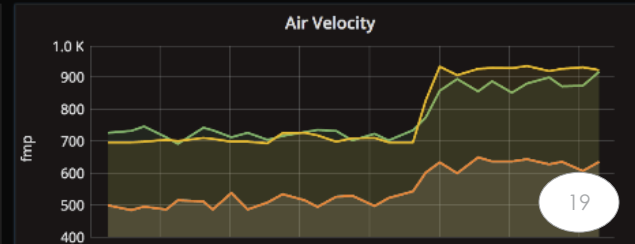
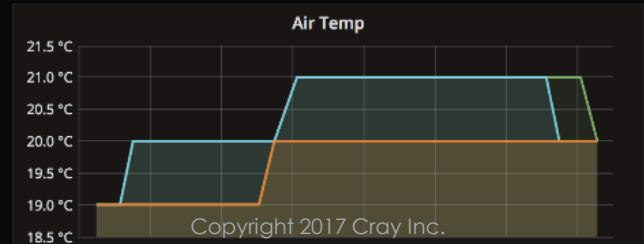
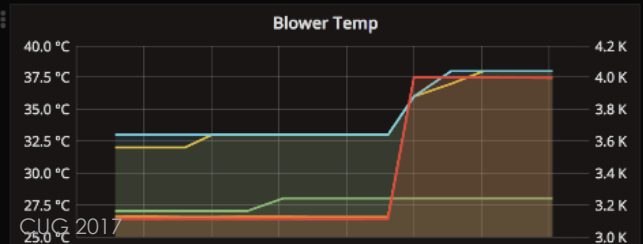
Chassis 0



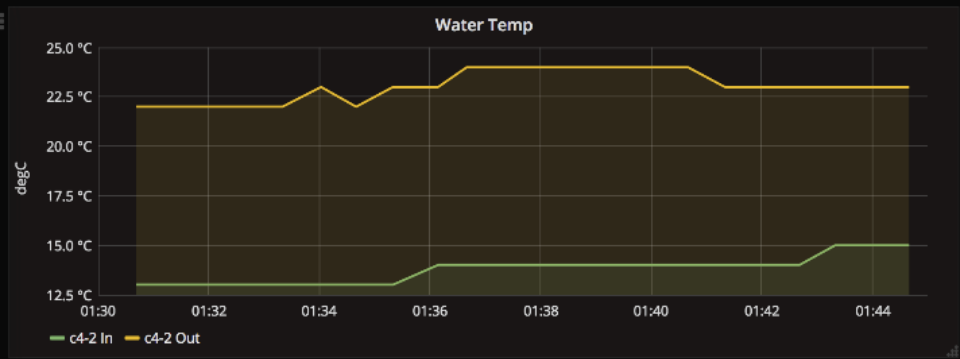
Chassis 1



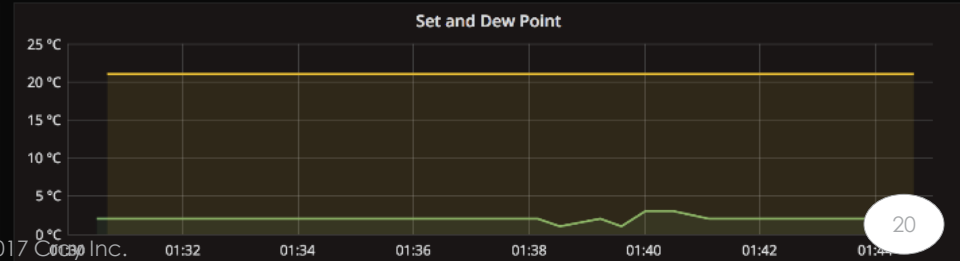
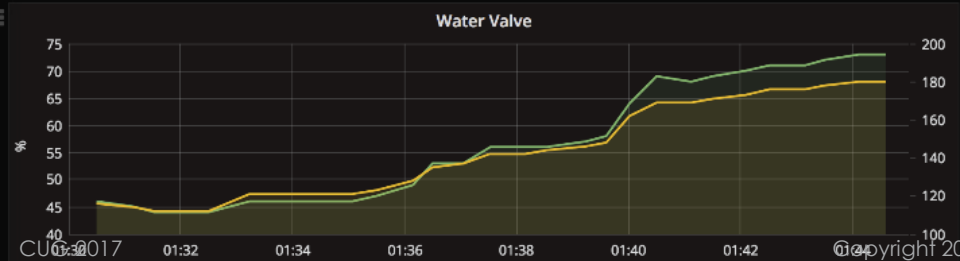
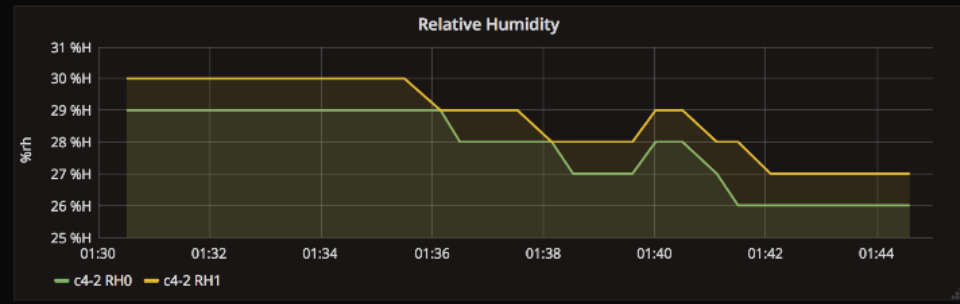
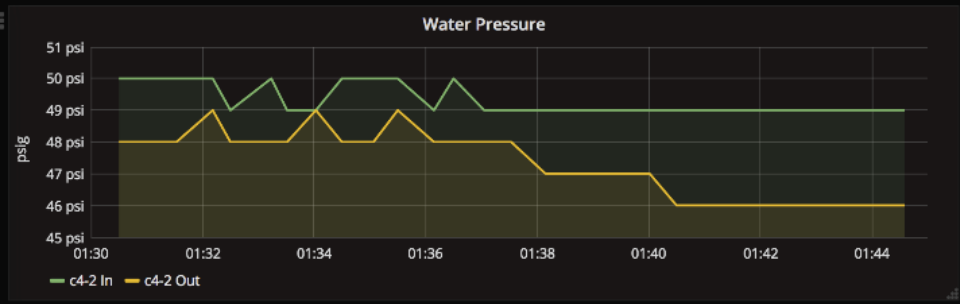
Chassis 2



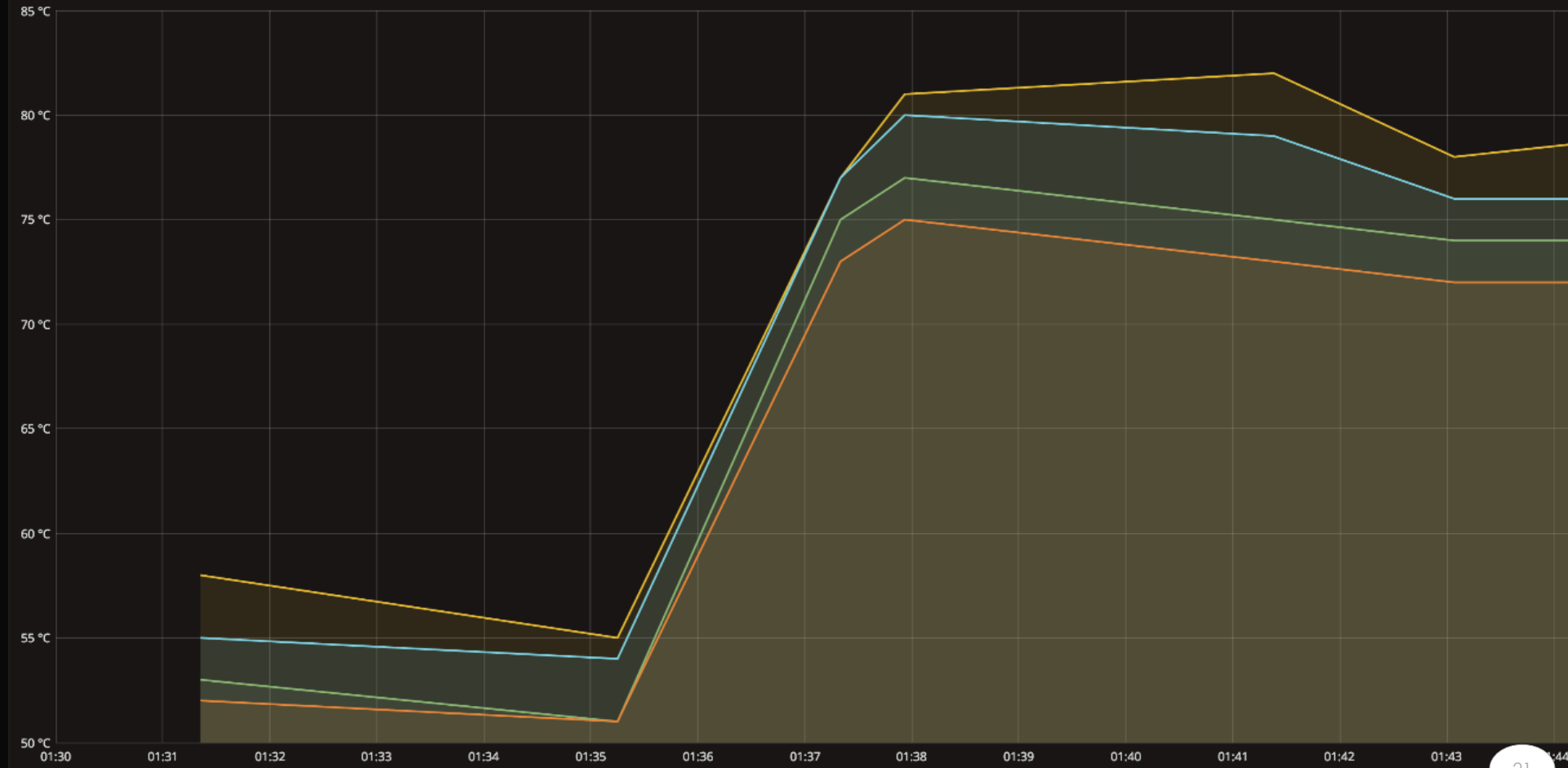
▼ Cabinet Water Temp

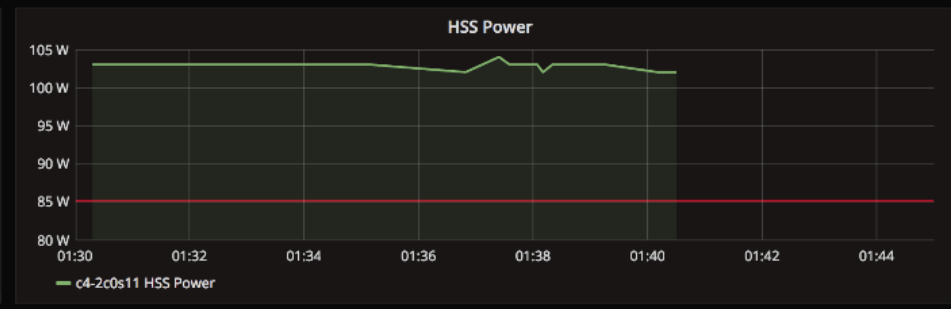
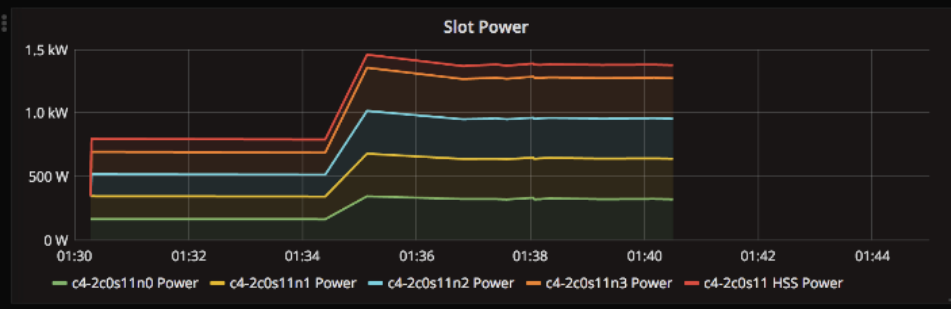


▼ Cabinet Water Pressure

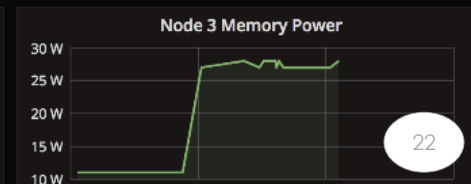
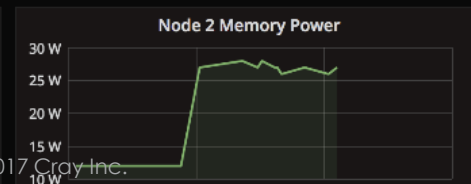
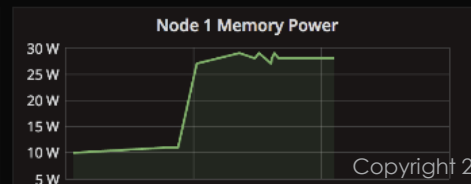
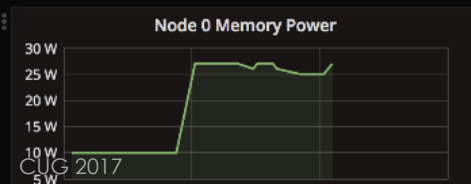
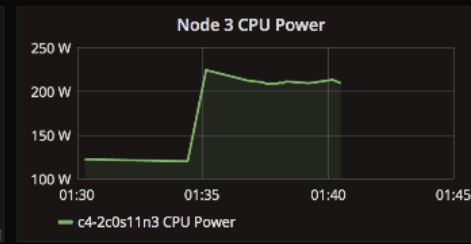
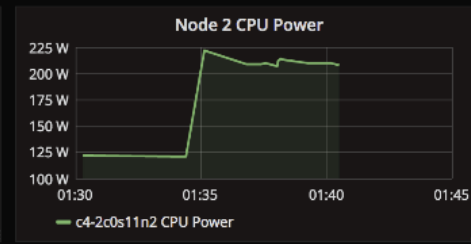
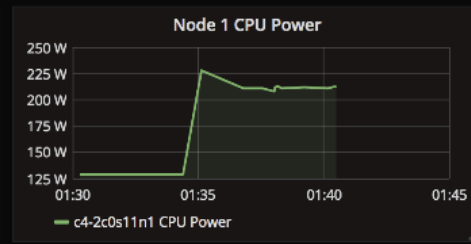
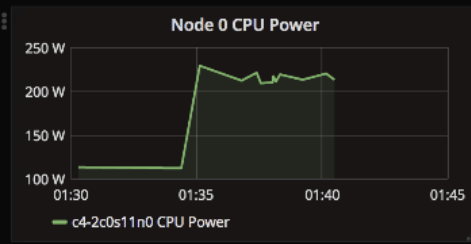
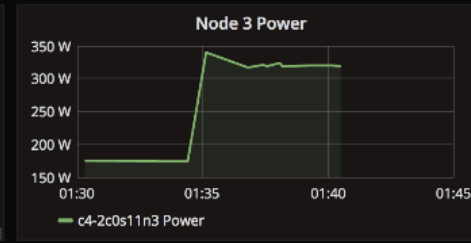
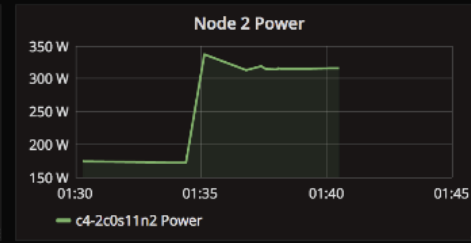
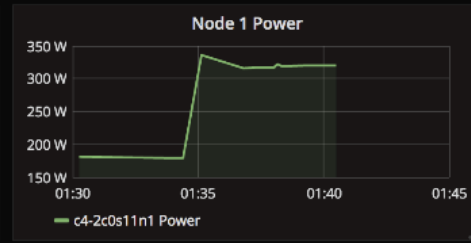
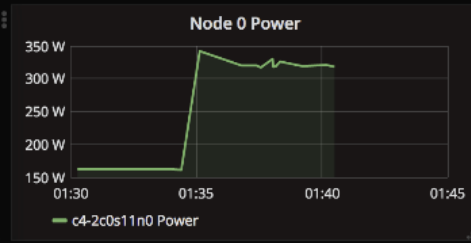


c4-2c0 - s11

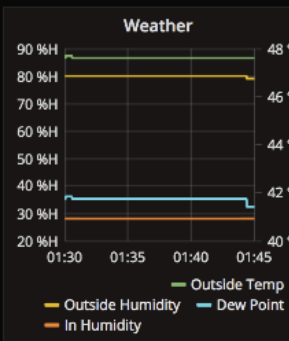
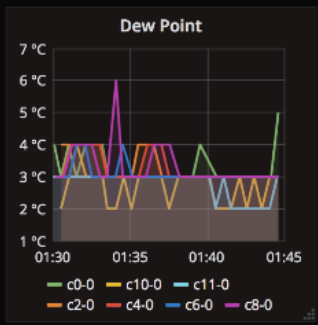
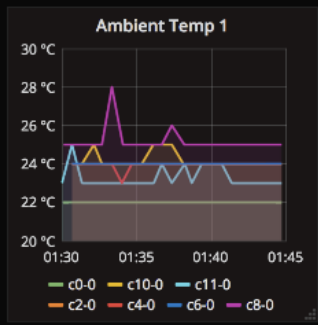
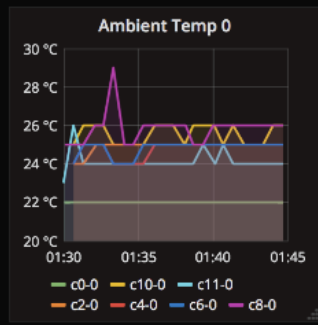
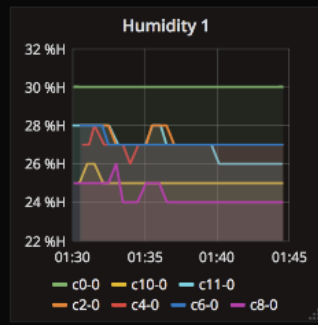
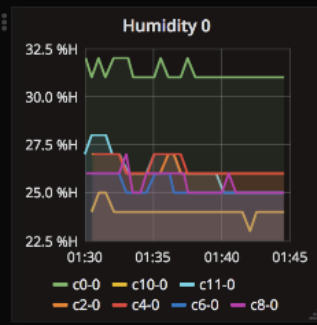




Slot

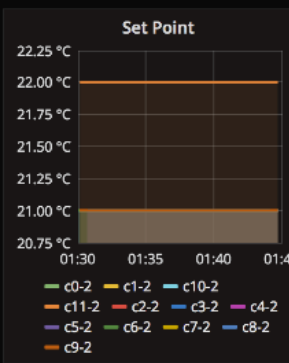
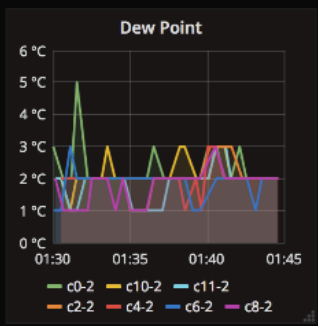
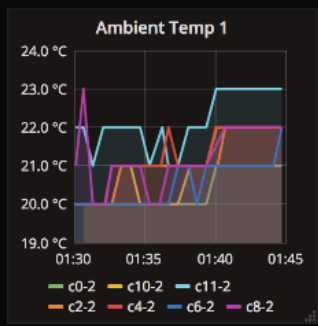
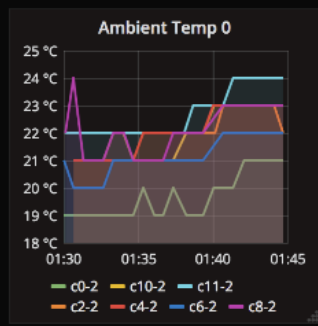
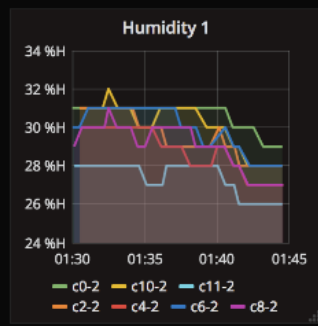
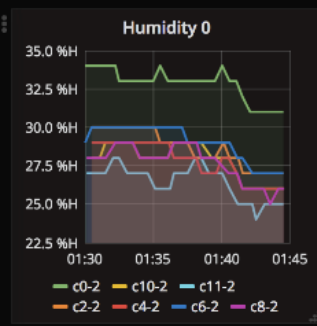


Row 0

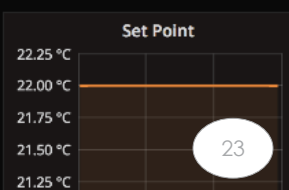
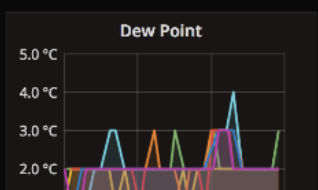
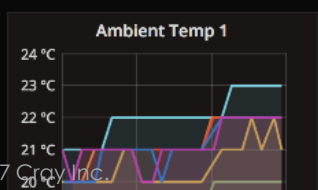
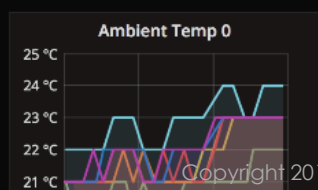
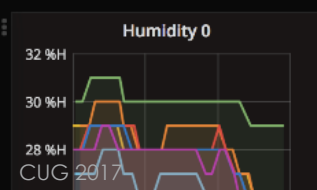


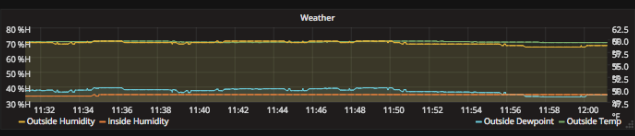
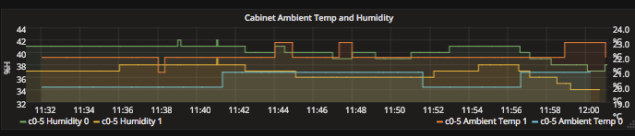
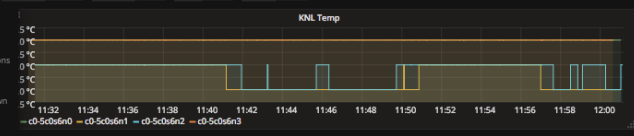
> Row 1

Row 2

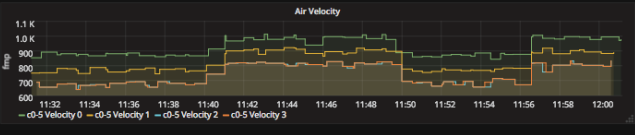
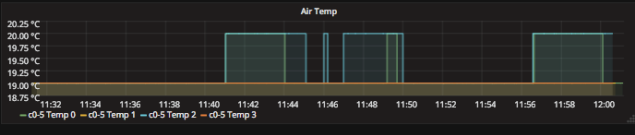
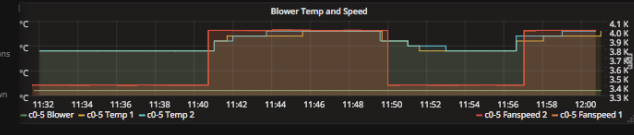


Row 3

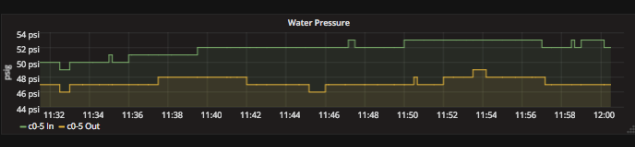
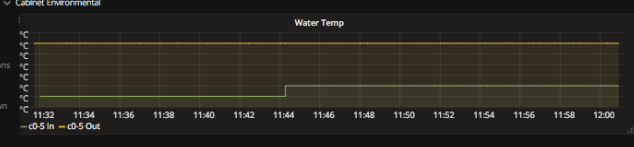




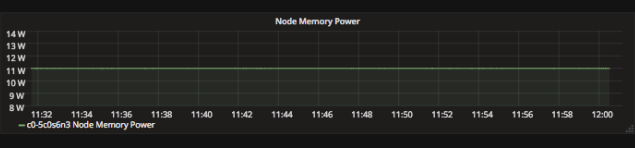
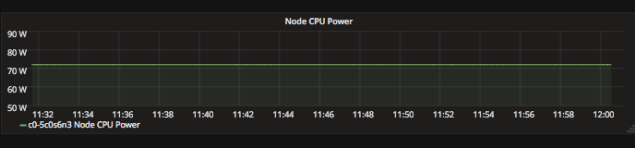
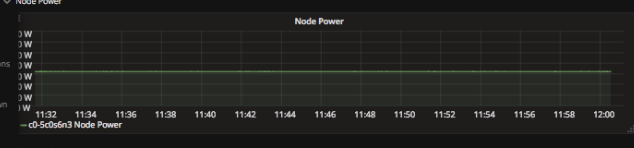
- Chassis 2 Environmental
- Chassis 1 Environmental
- Chassis 0 Environmental



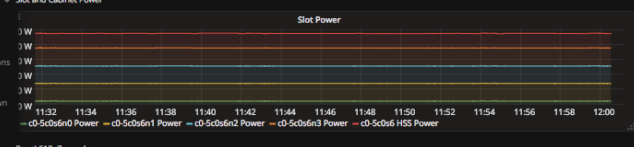
- Cabinet Environmental
- Node Power



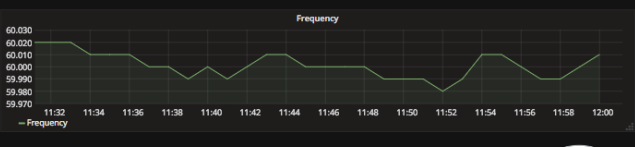
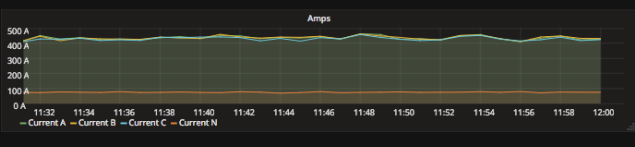
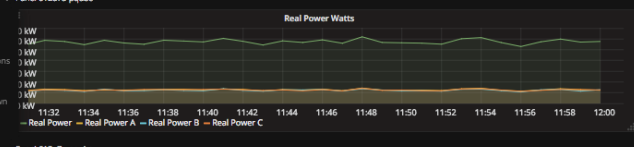
- Node CPU Power
- Node Memory Power
- Slot and Cabinet Power



- Panel 612a7a-pqube



- Panel 613a7a-pqube





# Wrap-up

- **Collaborative effort**
  - Cray, customers, and the user community
- **Allows streaming of**
  - Power, energy, thermal, and application meta-data
- **Making that data available**
  - System administrators
  - Application developers
  - HPC research community
- **As appropriate given site-level policy**

# Acknowledgment

- **NERSC**

- Specifically Cary Whitney for co-authoring this paper
- NERSC is leading the way in use of this new plugin feature

- **XTreme**

- Large system customer support is driving changes
- See also: [monitoring\\_wg@lists.cug.org](mailto:monitoring_wg@lists.cug.org)

- **HPC Community**

- Broad community push for improved monitoring capabilities
- See also: [EEHPC WG](#)

# Legal Disclaimer

*Information in this document is provided in connection with Cray Inc. products. No license, express or implied, to any intellectual property rights is granted by this document.*

*Cray Inc. may make changes to specifications and product descriptions at any time, without notice.*

*All products, dates and figures specified are preliminary based on current expectations, and are subject to change without notice.*

*Cray hardware and software products may contain design defects or errors known as errata, which may cause the product to deviate from published specifications. Current characterized errata are available on request.*

*Cray uses codenames internally to identify products that are in development and not yet publically announced for release. Customers and other third parties are not authorized by Cray Inc. to use codenames in advertising, promotion or marketing and any use of Cray Inc. internal codenames is at the sole risk of the user.*

*Performance tests and ratings are measured using specific systems and/or components and reflect the approximate performance of Cray Inc. products as measured by those tests. Any difference in system hardware or software design or configuration may affect actual performance.*

*The following are trademarks of Cray Inc. and are registered in the United States and other countries: CRAY and design, SONEXION, and URIKA. The following are trademarks of Cray Inc.: APPRENTICE2, CHAPEL, CLUSTER CONNECT, CRAYPAT, CRAYPORT, ECOPHLEX, LIBSCI, NODEKARE, REVEAL, THREADSTORM. The following system family marks, and associated model number marks, are trademarks of Cray Inc.: CS, CX, XC, XE, XK, XMT, and XT. The registered trademark LINUX is used pursuant to a sublicense from LMI, the exclusive licensee of Linus Torvalds, owner of the mark on a worldwide basis. Other trademarks used in this document are the property of their respective owners.*



# Q&A

Steven Martin  
stevem@cray.com

**CUG.2017.CAFFEINATED COMPUTING**

Redmond, Washington May 7-11, 2017