

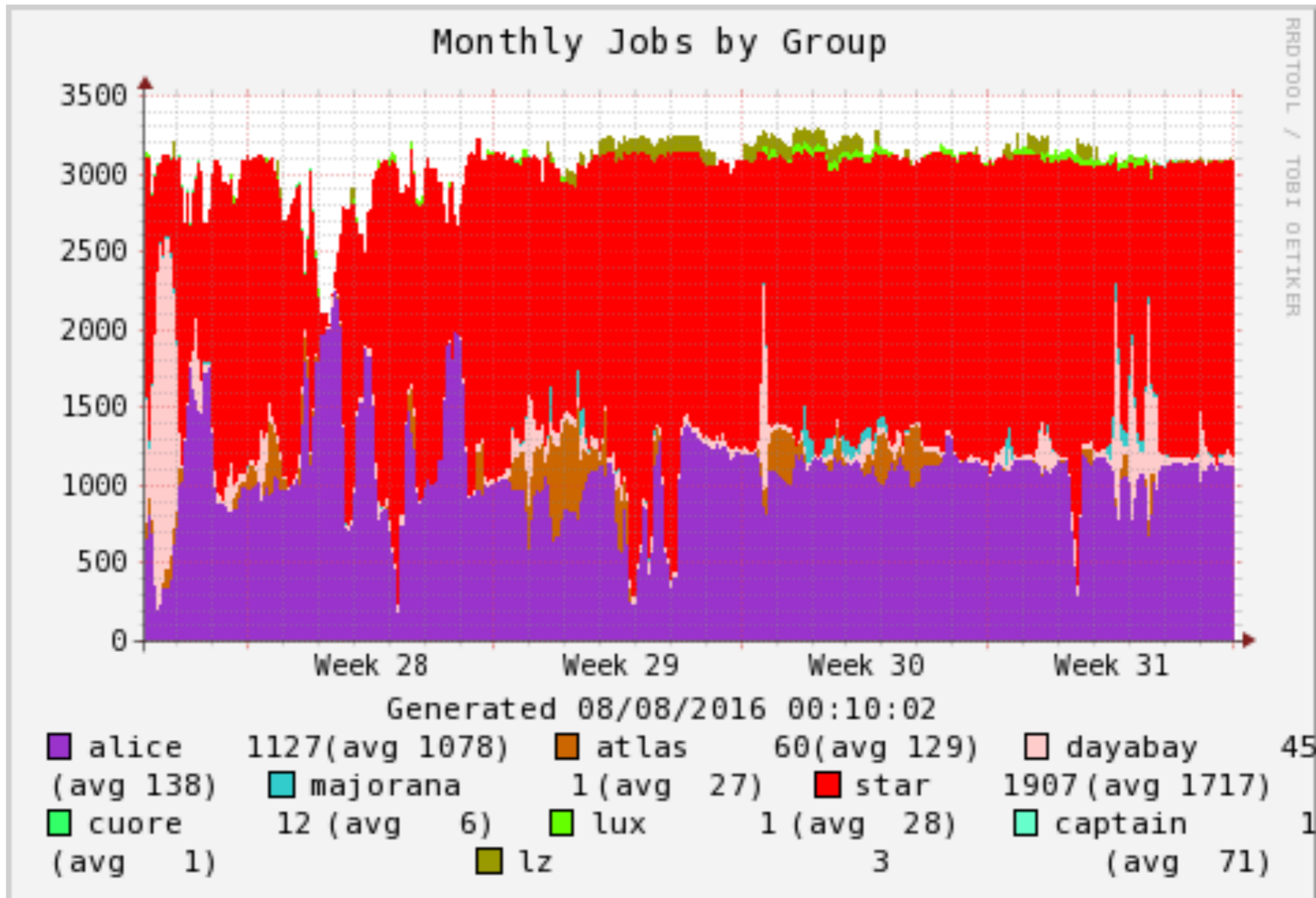
PDSF Users Meeting

- * PDSF performance
- * PDSF job queue oversubscribed
- * PDSF interactive nodes
- * past/future outages
- * announcements
- * toward common PDSF Shifter image
- * AOB

August 9, 2016

Jan Balewski

Utilization



PDSF queue capacity exceeded - INC0088407

Dates: 7-27, 8-2 , 8-9

Symptom: "I've been unable to submit more than a handful of jobs at a time since earlier today."

PDSF load at the time:

```
pdsf9 $ sgeusers
---
r qw Eqw hr hqw Ehqw dr dt jobIDs Project
-----
1103 98 0 0 0 0 0 1201 alice
2 0 0 0 0 0 0 2 atlas
58 143 0 0 0 0 0 201 dayabay
40 4 0 0 0 0 0 44 lux
91 17 0 0 0 0 0 108 lz
83 969 0 0 0 0 0 1052 majorana
1875 19913 26 0 0 0 0 21814 star
-----
3252 21144 26 0 0 0 0 24422 Totals
---
```

Impact: members of other projects cut-off from submitting PDSF jobs for several hours

Reason: STAR changed use-pattern , lot of small jobs were submitted

Response: talk to STAR reps (Jeff, Jerome), ask STAR users to scale down, short term improvement

Solution: STAR will add & use job-array feature reducing the job count by factor of few

This week: PDSF load ~9k jobs (26k queue slots), not true any more (see next slide)

PDSF waiting queue in BUSY, again

pdsf11 \$ date
Tue Aug 9 06:36:17 PDT 2016

| r | qw | Eqw | hr | hqw | Ehqw | dr | dt | jobIDs | User-Project |
|-----|-------|-----|----|-----|------|----|----|--------|-----------------------|
| 19 | 100 | 1 | 0 | 0 | 0 | 0 | 0 | 120 | alicesgm alice |
| 0 | 154 | 0 | 0 | 0 | 0 | 0 | 0 | 154 | angrzej star |
| 0 | 675 | 0 | 0 | 0 | 0 | 0 | 0 | 675 | beizhen dayabay |
| 55 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 55 | ccuesta majorana |
| 0 | 500 | 0 | 0 | 0 | 0 | 0 | 0 | 500 | jinhui star |
| 0 | 5000 | 0 | 0 | 0 | 0 | 0 | 0 | 5000 | klandry star |
| 0 | 169 | 0 | 0 | 0 | 0 | 0 | 0 | 169 | kocmanek star |
| 0 | 84 | 0 | 0 | 0 | 0 | 0 | 0 | 84 | matonoli star |
| 0 | 44 | 0 | 0 | 0 | 0 | 0 | 0 | 44 | mjd majorana |
| 0 | 5000 | 0 | 0 | 0 | 0 | 0 | 0 | 5000 | mjsim majorana |
| 0 | 3384 | 0 | 0 | 0 | 0 | 0 | 0 | 3384 | nasim star |
| 130 | 920 | 0 | 0 | 0 | 0 | 0 | 0 | 1050 | qiuh star |
| 0 | 1501 | 0 | 0 | 0 | 0 | 0 | 0 | 216 | rlinehan lz |
| 0 | 1026 | 0 | 0 | 0 | 0 | 0 | 0 | 1026 | roliesha star |
| 121 | 9 | 0 | 0 | 0 | 0 | 0 | 0 | 130 | simkomir star |
| 0 | 1073 | 0 | 0 | 0 | 0 | 0 | 0 | 1073 | smizuno star |
| 0 | 100 | 0 | 0 | 0 | 0 | 0 | 0 | 100 | sss star |
| 0 | 4496 | 0 | 0 | 0 | 0 | 0 | 0 | 4496 | staremb star |
| 0 | 316 | 0 | 0 | 0 | 0 | 0 | 0 | 316 | tubiao star |
| 0 | 757 | 0 | 0 | 0 | 0 | 0 | 0 | 757 | vipul star |
| 0 | 568 | 0 | 0 | 0 | 0 | 0 | 0 | 568 | yezhenyu star |
| 0 | 1199 | 0 | 0 | 0 | 0 | 0 | 0 | 1199 | yshuai star |
| 327 | 31164 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 26239 Totals |

| r | qw | Eqw | hr | hqw | Ehqw | dr | dt | jobIDs | Project |
|-----|-------|-----|----|-----|------|----|----|--------|--------------|
| 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | admin |
| 19 | 100 | 1 | 0 | 0 | 0 | 0 | 0 | 120 | alice |
| 0 | 13 | 0 | 0 | 0 | 0 | 0 | 0 | 4 | atlas |
| 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | captain |
| 0 | 678 | 0 | 0 | 0 | 0 | 0 | 0 | 676 | dayabay |
| 0 | 3 | 0 | 0 | 0 | 0 | 0 | 0 | 3 | lux |
| 0 | 2497 | 0 | 0 | 0 | 0 | 0 | 0 | 217 | lz |
| 56 | 5044 | 0 | 0 | 0 | 0 | 0 | 0 | 5100 | majorana |
| 251 | 22828 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 20117 star |
| 327 | 31164 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 26239 Totals |

3 users can lock-out all other users from all experiments regardless of the job-slot shares per experiment

pdsf11 \$

Past/future outages

Today

PDSF down, batch system being moved to CRT

Future

week of August 29:

migration of the eliza3 and neweliza18 file systems to CRT

August 30:

NERSC quarterly maintenance

Announcements

Bi-weekly office hours in June

- Thursday, August 18 & September 1, 59-4016-CR

PDSF users meeting

- Tuesday, September 6, 11:00 - 12:01pm 59-3034-CR

PDSF interactive nodes

PDSF has 6 interactive nodes: pdsf6,...,11
each: 32 cores, 128GB RAM, 2.6 GHz clock
total : 192 cores

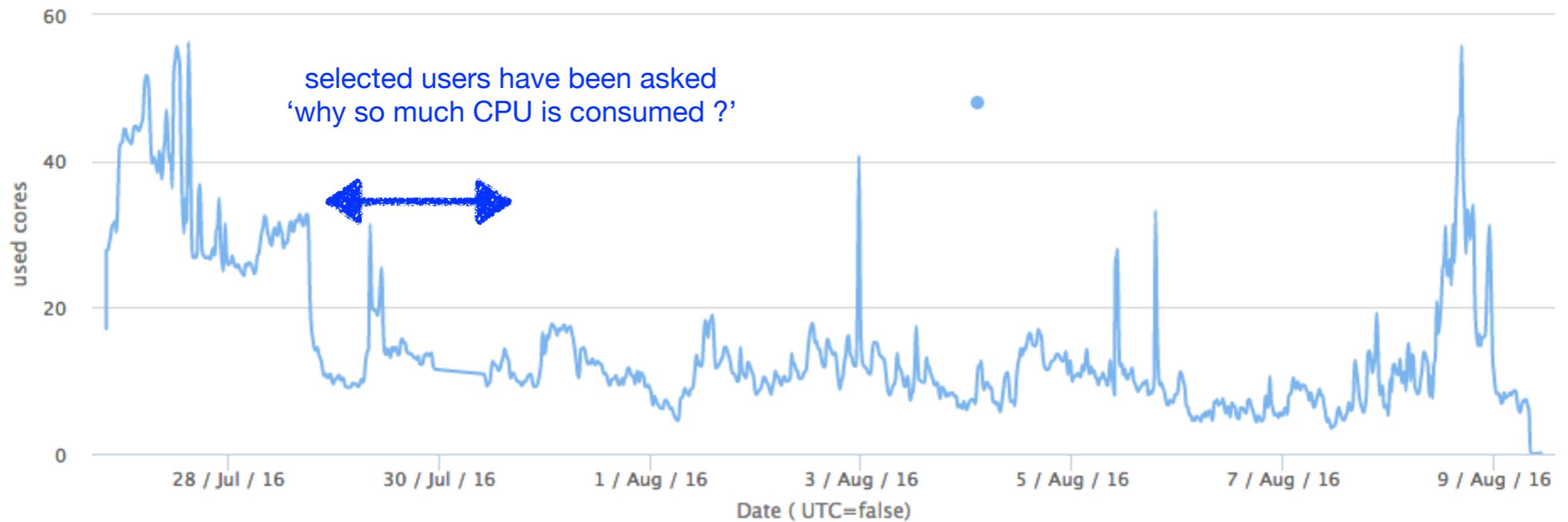
pds 9,10,11 were repurposed compute nodes when Mendel was split between OSF and CRT

Proposal :

- put pdsf 9,10,11 back to compute pool
- use only 3 interactive nodes, 96 cores total

'interactive' monitoring of load on PDSF interactive nodes

used cores, argegated over pdsf6...pdsf11



nominal capacity of 3 nodes x 32 cores=96

PDSF load recorder

work in progress - shall I make it permanent ?

Periodic snapshot of PDSF nodes load archived at MongoDB

<https://portal-auth.nersec.gov/pdsf-mon/m>

Table 1.: pdsf6-11

| | value |
|----------------|-----------|
| collection | janAbuseE |
| record cnt | 2 |
| heartbeat age | 652 (sec) |
| - - - | |
| num hosts | 6 |
| avr load/host | 4.6 of 32 |
| seen big users | 6 |
| avr big users | 4 +/- 2 |

when nodes were scanned last time
(every 20 minutes) →

sum of load /sum of hosts →

outstanding 'over limit#3' →

averaged over 90 minutes →

Proposed policy:

Definition of limit#3 : count in any of the following categories **exceeds 3**:

- * **number of processes** above 10%+ of CPU
- * **sum of CPU-load** for all your processes
- * **sum of consumed CPU-days** for all processes

The sum over all 6 PDSF interactive nodes is taken into account for each user.

List of big users Table 2.

| user | use core | num jobs | sum CPU hours | a job | one example of job name |
|----------|----------|----------|---------------|------------------------|-------------------------|
| carels | 1.0 | 1 | 372.3 | pdsf11:LUXSimExecutabl | |
| jennetd | 5.0 | 5 | 32.7 | pdsf7:athena.py | |
| massar | 0.9 | 1 | 3.4 | st.exe | |
| alfredso | 0.9 | 1 | | pdsf8:MaGe | |
| rlinehan | 2.0 | 2 | 208.8 | pdsf9:plot-hist | |
| zillay | 15.0 | 15 | 85.9 | pdsf9:root4star | |

1 or 2 jobs for 10+days →

5 jobs x 6 hours →

15 jobs x5 hours on the same node →

Example of monitoring page

Virtualization of software stack for experiments



docker



vmware®

I executed your ./dybinst trunk tests on my Docker image with nuwa software, compiled in SL6.4.

The setup: SL6.4, db credentials from Cheng-Ju as before, for those DBs:

```
[balewski@e1aaa35b1d2f nuwa]$ pwd
/home/balewski/nuwa
```

```
[balewski@e1aaa35b1d2f nuwa]$ grep host ~/.my.cnf
host = dybdb1.ihep.ac.cn
host = dayabaydb.lbl.gov
host = dybdb1.ihep.ac.cn
host = dybdb2.ihep.ac.cn
host = dayabaydb.lbl.gov
```

DAYABAY - The D



I can write to /home/balewski/nuwa and have internet connection:

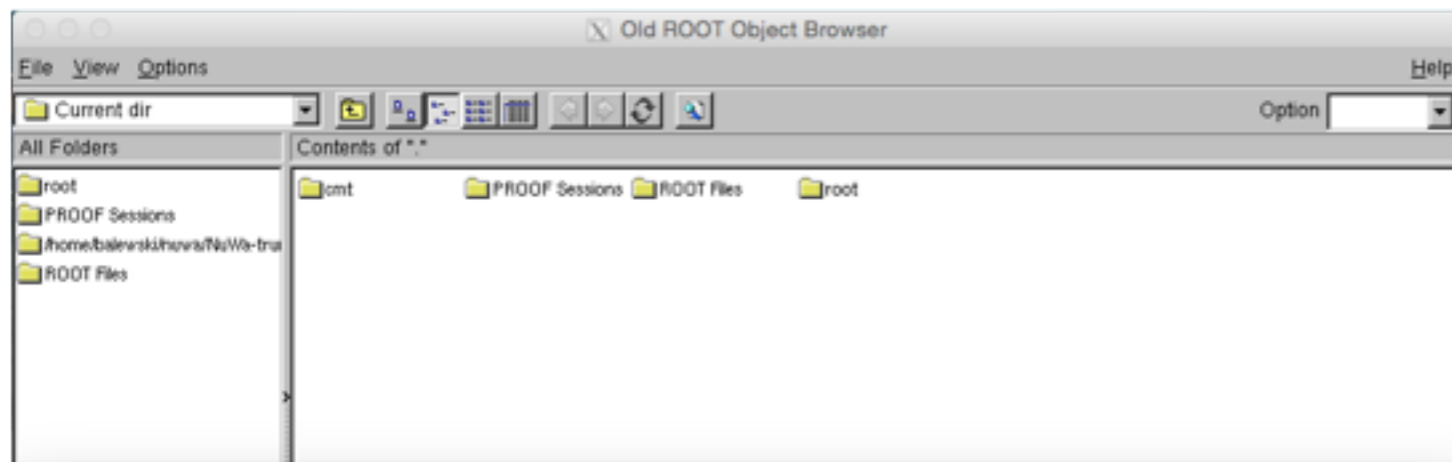
Again, most of the test have passed:

```
[balewski@e1aaa35b1d2f nuwa]$ grep ok Log-trunk-tests-Docker1 |nl |tail
550 Generate muon hit sim sample ... ok
551 Mix results of previous jobs ... ok
552 test_fmcp11a.test_213geom ... ok
553 test_fmcp11a.test_dummy ... ok
554 test_fmcp11a.test_nuwa ... ok
555 Running Elecsim, TrigSim, and ReadoutSim with low charge input ... ok
556 Running DigitizeAlg with prompt/delayed charge to test SimHeader splitting ... ok
557 test_KUP11a.test_KUP11a_EH1 ... ok
558 test_KUP11a.test_KUP11a_EH2 ... ok
559 test_KUP11a.test_KUP11a_EH3 ... ok
```

The 'ok' count is 559.

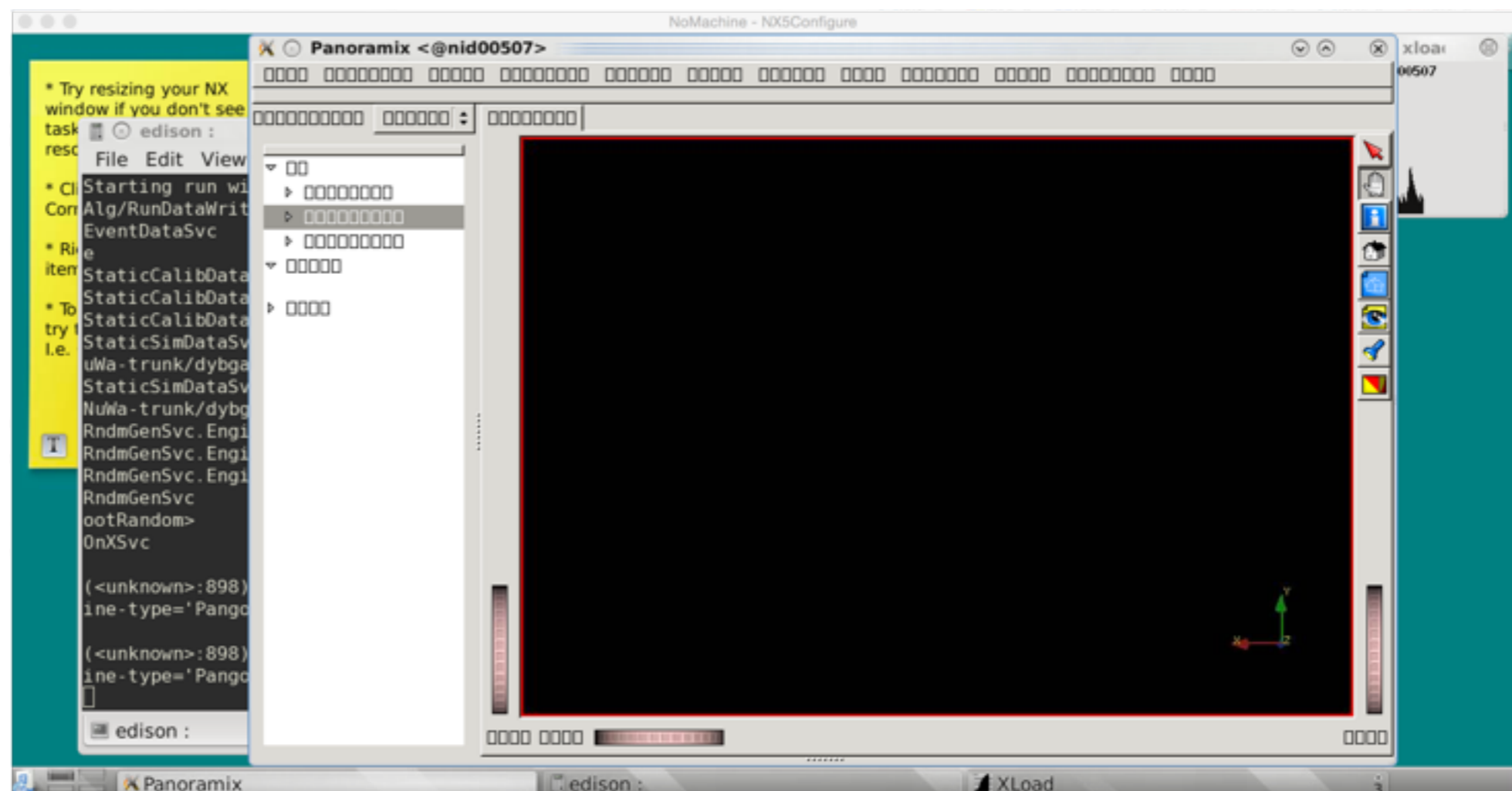
The word FAIL is counted 58 times, meaning there are about **30 failures**

DayaBay on Edison



```
RT_COMMAND=nosetests -v with subversion /usr/bin/svn using RT_PATH /usr/bin:/usr/local/bin:/bin:/usr/bin
pwd ; echo -n ; echo -n
/global/u2/b/balewski/nuwa/NuWa-trunk/dybgaudi/Production/KUP11a
/global/homes/b/balewski/nuwa/installation/trunk/dybinst/scripts/runtest.sh: line 161: nosetests: command not found
Command: popd
~/nuwa
Command: date
Thu Aug 4 20:00:20 EDT 2016
bash-4.1$ vi /global/homes/b/balewski/nuwa/install
bash-4.1$ eval
bash-4.1$ eval ls
dybinst dybinst-20160722-210245.log dybinst-rece
bash-4.1$ vi /global/homes/b/balewski/nuwa/install
bash-4.1$ pwd
/home/balewski/nuwa
bash-4.1$
bash-4.1$
bash-4.1$ uname -a
Linux nid00112 3.0.101-0.35.1_1.0502.8640-cray_ar
bash-4.1$ pwd
/home/balewski/nuwa
bash-4.1$ cd NuWa-trunk
bash-4.1$ source setup.sh
bash-4.1$ cd dybgaudi/DybRelease/cmt
bash-4.1$ source setup.sh
bash-4.1$ root -l
root [0] new TBrowser
```

```
You are using the old ROOT browser! A new version
Select the "New Browser" entry from the "File" me
"Browser.Name:" from "TRootBrowserLite" to "TRo
(class TBrowser*)0x1466b10
root [1]
```



Performance - DayaBay - simulations

time ./dybinst trunk tests **fmcp11a**

| | pdsf10 | Docker/Laptop | Shifter/Edison |
|------------------------------|---------------------|---------------|----------------|
| | SL5.3 | SL6.4 | SL6.4 |
| test_fmcp11a.test_sample_00 | fail | ok | ok |
| test_fmcp11a.test_sample_14 | ok | ok | ok |
| test_fmcp11a.test_sample_20 | ok | ok | ok |
| test_fmcp11a.test_sample_28 | ok | ok | ok |
| Generate muon hit sim sample | ok | ok | ok |
| Mix results of previous jobs | ok | ok | ok |
| test_fmcp11a.test_213geom | fail | fail | fail |
| test_fmcp11a.test_dummy | ok | ok | ok |
| test_fmcp11a.test_nuwa | ok | ok | ok |
| real time | similar performance | | |
| user time | | | |
| sys time | | | |

Optimization on Shifter in progress ...

Performance - DayaBay - data analysis

time ./dybinst trunk tests **kup11a**

| | mc0101 | Docker/Mac | Shifter/Edison |
|---|------------------------------|------------|----------------|
| | SL5.3 | SL6.4 | SL6.4 |
| test_KUP11a.test_KUP11a_EH1 ... ok | ok | ok | ok |
| test_KUP11a.test_KUP11a_EH2 ... ok | ok | ok | ok |
| test_KUP11a.test_KUP11a_EH3 ... ok | ok | ok | ok |
| test_KUP11a_cori.test_KUP11a_EH1 ... ok | ok | ok | ok |
| test_KUP11a_cori.test_KUP11a_EH2 ... ok | ok | ok | ok |
| test_KUP11a_cori.test_KUP11a_EH3 ... ok | ok | ok | ok |
| real time | different performance | | |
| user time | | | |
| sys time | | | |

Optimization on Shifter in progress ...

AOB