



STANFORD RESEARCH COMPUTING CENTER



slurm on Sherlock

workload manager

Slurm User Group 2019

17-18 September 2019 - Salt Lake City, Utah, USA

Kilian Cavalotti

Stanford Research Computing Center





Hi!

My name is **Kilian** Cavalotti

I'm HPC Technical Lead & Architect at Stanford University

I manage **Sherlock**, the Stanford shared HPC cluster

Stanford Research Computing Center

Our mission

Build & support a
comprehensive program and
capabilities to advance
computational and
data-intensive research at
Stanford







Sherlock is

- ▶ a **shared** HPC cluster, operated by the Stanford Research Computing Center
- ▶ available at **no cost** to all Stanford Faculty members and their research teams to **support sponsored research**
- ▶ a **condo** cluster, where PIs can become **owners**
 - ▷ they get their own Slurm partition
 - ▷ they get access to all the other owner's nodes when they're not in use (owner jobs preempt background jobs)

A little bit of **history**

2014 | Sherlock opens for production
initial seed of 120 nodes, funded by the University Provost

2016 | Sherlock reaches capacity
FDR Infiniband fabric maxed out, ~800 PI-owned nodes

2017 | Sherlock 2.0
complete hardware, platform, software and services refresh

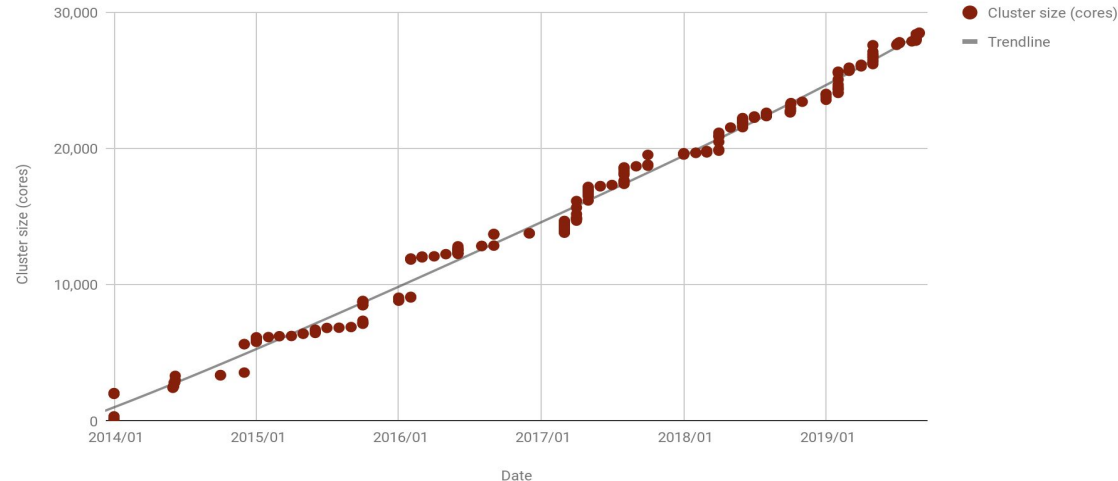
2019 | Sherlock 2.0 reaches capacity
EDR IB fabric maxed out

2020 | Sherlock 3.0
hardware refresh, HDR IB fabric rollout



Successful model

Cluster size (cores) vs. Date



- ▶ 15x growth in 5 years
- ▶ >90% of Sherlock's nodes are owners nodes

Sherlock

as of today, things change on a weekly basis

- ▶ **1,502** compute nodes
16 → 64 CPU cores, 64 → 3TB RAM
from Ivy Bridge to Skylake
- ▶ **780** GPUs
from K20s to V100 SXM2 & TITANs
4x or 8x GPUs per node
- ▶ **2** IB fabrics
2:1 FDR + 2:1 EDR
- ▶ **Slurm** 18.08.8
CentOS 7.6



Overall

29,192 CPU cores

6 PB scratch + 12 PB long-term storage

>4,500 users

from all the 7 schools, SLAC, Stanford institutes, etc.

>720 PI groups, 125 owner groups from astrophysics to the music department

over 3,500 support tickets per year

And more...

2.12 PFlops

FP64, FWIW, YMMV

46 racks

17 server models

450 kW

4 CPU/GPU generations

96 IB switches

117 Slurm partitions

7,248 IB cables

100,000 jobs/day
peaks at 250 jobs/s

15 million jobs since Jan 01
job #50,000,000 submitted on Sep. 10



```
# scontrol show config
Configuration data as of 2019-09-04T16:46:51
AccountingStorageBackupHost = sh-sl02
AccountingStorageEnforce = associations,limits,nosafe
AccountingStorageHost = sh-sl01
AccountingStorageLoc = N/A
AccountingStoragePort = 6819
AccountingStorageTRES = cpu,mem,energy,node,billing,fs/disk,mem,pages,fs/swaps,strings,cores
AccountingStorageType = accounting_storage/slurddb
AccountingStorageUser = N/A
AccountingStoreJobComment = Yes
AcctGatherEnergyType = acct_gather_energy/rapl
AcctGatherFilesystemType = acct_gather_filesystem/lustre
AcctGatherInterconnectType = acct_gather_interconnect/ofed
AcctGatherNodeFreq = 300 sec
AcctGatherProfileType = acct_gather_profile/none
AllowSpecResourcesUsage = 0
AuthInfo = (null)
AuthType = auth/munge
BatchStartTimeout = 10 sec
JOB_TIMEOUT = 2019-08-29T09:00:46
CheckpointBufferType = (null)
CheckpointType = checkpoint/none
CheckpointType = sherlock
CheckpointType = (null)
```


Slurm **partitions**

- ▶ **Public partitions**

Anybody can use

- ▷ **normal:** regular CPU nodes
- ▷ **gpu:** GPU nodes
- ▷ **bigmem:** large memory nodes
- ▷ **dev:** interactive/debug nodes

- ▶ **Owner partitions**

Nodes purchased by PIs, one partition per owner

- ▶ **Owners partition**

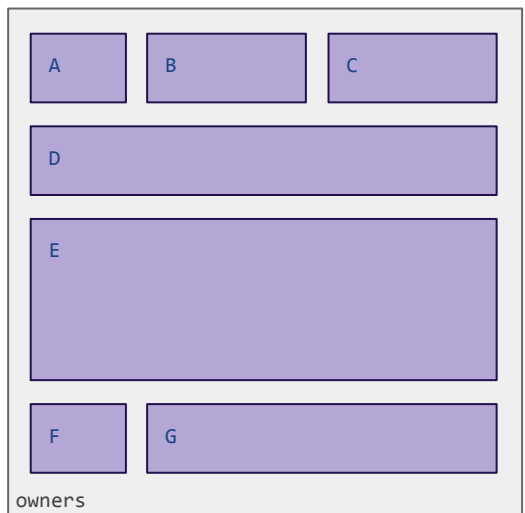
Preemptable partition for background owners' jobs

Partition structure



"Public" partitions

- ▶ All users can submit jobs there



Owner partitions

- ▶ PI A (only) can submit to partition "A"
- ▶ PI A can submit to "owners" and run on nodes from PI B
- ▶ When PI B submits to "B", PI A's job on B is preempted

Triggers

<https://slurm.schedmd.com/strigger.html>

triggers are generated when specific Slurm events occur, and can be used to run actions (scripts)

- ▶ events could be on the controller, the database daemon, jobs, nodes...
- ▶ we use triggers... a lot!
 - ▷ for notifications and monitoring
 - ▷ for automatic recovery actions (w/ NHC)

Triggers

RES_TYPE	RES_ID	TYPE	OFFSET	USER	FLAGS	PROGRAM
node	*	down	0	slurm	PERM	down.sh
node	*	drained	0	slurm	PERM	drained.sh
node	*	up	0	slurm	PERM	up.sh
slurmctld	*	primary_slurmctld_resumed_operation	0	slurm	PERM	primary_slurmctld_resumed_op.sh
slurmctld	*	primary_slurmctld_resumed_control	0	slurm	PERM	primary_slurmctld_resumed_control.sh
slurmctld	*	backup_slurmctld_resumed_operation	0	slurm	PERM	backup_slurmctld_resumed_operation.sh
database	*	primary_database_resumed_operation	0	slurm	PERM	primary_database_resumed_operation.sh
slurmdbd	*	primary_slurmdbd_resumed_operation	0	slurm	PERM	primary_slurmdbd_resumed_operation.sh
slurmdbd	*	primary_slurmdbd_failure	0	slurm	PERM	primary_slurmdbd_failure.sh
slurmctld	*	backup_slurmctld_failure	0	slurm	PERM	backup_slurmctld_failure.sh
slurmctld	*	primary_slurmctld_failure	0	slurm	PERM	primary_slurmctld_failure.sh
slurmctld	*	backup_slurmctld_assumed_control	0	slurm	PERM	backup_slurmctld_assumed_control.sh
database	*	primary_database_failure	0	slurm	PERM	primary_database_failure.sh
slurmctld	*	primary_slurmctld_acct_buffer_full	0	slurm	PERM	primary_slurmctld_acct_buffer_full.sh

When problems happen on nodes:

- ▶ NHC detects the issue, drains nodes
- ▶ script associated w/ events is triggered
- ▶ sysadmins are notified
- ▶ corrective actions are taken

The screenshot displays a series of log entries from the Slurm APP interface. A red arrow points from the 'drained.sh' program in the table to the first 'Node drained' event. The logs show a sequence of events: a node is drained, NHC reports memory errors, a corrective action (Scheduling reboot) is triggered, the node is drained again, NHC reports more errors, another corrective action is triggered, and finally, the node is brought back up.

Slurm APP 12:05
Node drained
1 node has been drained
sh-03-11

NHC: check_hw_mcelog: MCEs detected - 19 corrected memory errors detected (limit of 9/day).

Logs

Sep 5 04:21:43 sh-03-11 mcelog: MISC 90850002000228c ADDR 99744c000
Sep 5 04:21:43 sh-03-11 mcelog: TIME 1567682503 Thu Sep 5 04:21:43 2019
Sep 5 04:21:43 sh-03-11 mcelog: MCG status:
Sep 5 04:21:43 sh-03-11 mcelog: MCI status:
Sep 5 04:21:43 sh-03-11 mcelog: Error overflow
Sep 5 04:21:43 sh-03-11 mcelog: Corrected error
Sep 5 04:21:43 sh-03-11 mcelog: MCI_MISC register valid
Sep 5 04:21:43 sh-03-11 mcelog: MCI_ADDR register valid
Sep 5 04:21:43 sh-03-11 mcelog: MCA: MEMORY CONTROLLER MS_CHANNEL2_ERR
Sep 5 04:21:43 sh-03-11 mcelog: Transaction: Memory scrubbing error
Sep 5 04:21:43 sh-03-11 mcelog: MemCtrl: Corrected patrol scrub error
Sep 5 04:21:43 sh-03-11 mcelog: ror
Sep 5 04:21:43 sh-03-11 mcelog: STATUS cc000491000800c2 MCGSTATUS 0
Sep 5 04:21:43 sh-03-11 mcelog: MCGCAP 1000c19 APICID 20 SOCKETID 1
Sep 5 04:21:43 sh-03-11 mcelog: MICROCODE 42e
Sep 5 04:21:43 sh-03-11 mcelog: CPUID Vendor Intel Family 6 Model 62

Corrective action
Scheduling reboot
Sherbot : Today at 12:05

Node drained
1 node has been drained
sh-03-11

NHC: rebooting node (reboot_node)

Sherbot : Today at 12:05

12:07 Node drained
1 node has been drained
sh-09-02

NHC: rebooting node (reboot_node)

Sherbot : Today at 12:07

Node up
1 node is up again
sh-09-02

Sherbot : Today at 12:10

Slurm APP 12:18
Node drained
1 node has been drained
sh-09-02

Triggers

RES_TYPE	RES_ID	TYPE
node		* down
node		* drained
node		* up
slurmctld		* primary_slurmctld_resumed_operation
slurmctld		* primary_slurmctld_resumed_control
slurmctld		* backup_slurmctld_resumed_operation
database		* primary_database_resumed_operation
slurmdbd		* primary_slurmdbd_resumed_operation
slurmdbd		* primary_slurmdbd_failure
slurmctld		* backup_slurmctld_failure
slurmctld		* primary_slurmctld_failure
slurmctld		* backup_slurmctld_assumed_control
database		* primary_database_failure
slurmctld		* primary_slurmctld_acct_buffer_full

operations



Primary slurmctld resumed control



Backup slurmctld failure



On `slurmctld` / `slurmdbd` events:

- ▶ sysadmins are notified
- ▶ events are recorded

8/19

2019-08-23 09:45:48

Primary Slurm Controller resumed control

sherlock slurm primary_slurmctld

primary_slurmctld_resumed_control

last cycle time

Node weights

- ▶ get the best performance out of the box
 - ▷ provide the highest-end CPU when possible
- ▶ don't waste specialized resources
 - ▷ when not explicitly requested, allocate specialty nodes in last resort
- ▶ encode node characteristics in node weights
 - ▷ nodes with lower weight are selected first
 - ▷ more memory / GRES = higher weight, so nodes are selected last
 - ▷ more recent CPUs = lower weight, give the best performance by default

Weight mask

```
# Weight mask: 1 | #GRES | Memory | #Cores | CPUgen | 1
#   prefix is to avoid octal conversion
#   suffix is to avoid having null weights
#
# Values:
#   #GRES   none: 0   Memory  64 GB: 0   #Cores  16: 0   CPUgen  ???: 3
#           1 GPU: 1   96 GB: 1   20: 1   CSL: 4
#           2 GPU: 2   128 GB: 2   24: 2   SKX: 5
#           3 GPU: 3   256 GB: 3   28: 3   BDW: 6
# [...]
# BDW | 20c | 128GB
# nodeName=[...] \
#   Sockets=2 CoresPerSocket=10 \
#   RealMemory=128000 \
#   Weight=102161 \
#   Feature="CPU_MNF:INTEL,CPU_GEN:BDW,CPU_SKU:E5-2640v4,CPU_FRQ:2.40GHz"
```

- ▶ Example: 20-core Broadwell w/ 128GB RAM, no GPU

Weight=102161

Recurrent jobs

- user: “*can I do cron?*”
- sysadmin: “*nope, but you can Slurm it!*”

	Cron job	Recurrent job
Dedicated resources for the task	✗	✓
Persistent across node reinstallations	✗	✓
Unique, controlled execution	✗	✓
Precise schedule	✓	✗

Recurrent jobs

```
#!/bin/bash
#SBATCH --job-name=cron
#SBATCH --begin=now+7days
#SBATCH --dependency=singleton
#SBATCH --time=00:10:00
#SBATCH --mail-type=FAIL

./my_cron_script

## Resubmit the job for the next execution
sbatch $0
```

- ▶ weekly execution
 - ▷ execution deferred by at least 7 days
- ▶ ensure unique running instance
- ▶ email notification on failure
- ▶ works great for backup, data sync jobs

Persistent jobs

- user: *“can I have a database server?”*
- sysadmin: *“nope, but you can Slurm it!”*
- ▶ run until explicitly `scancel` 'ed
- ▶ signal itself before end of allocated time
- ▶ resubmit itself on termination
- ▶ re-runs right away (maybe)
 - ▷ more adapted to partitions with lower wait times
- ▶ persistent `$JOBID`!

Persistent jobs

```
#!/bin/bash
#SBATCH --job-name=persistent
#SBATCH --dependency=singleton
#SBATCH --time=24:00:00
#SBATCH --signal=B:SIGUSR1@90

# catch the SIGUSR1 signal
_resubmit() {
    echo "$(date): job $SLURM_JOBID received SIGUSR1 at $(date), re-queueing"
    scontrol requeue $SLURM_JOBID
}
trap _resubmit SIGUSR1

./my_app &
wait
```

- ▶ signal `SIGUSR1` trapped by the `sbatch` step
 - ▷ main app needs to be running in the background for signal to be trapped
- ▶ `scontrol requeue` ensures persistent `$JOBID`
 - ▷ for easier job dependencies
- ▶ works great for database servers (MySQL, PostgreSQL...)

Wait, a database server in a job?

Sure, why not?

- ▶ resources are guaranteed and limited
- ▶ database service can run with user privileges
- ▶ db service jobs can be relocated to different compute nodes based on availability
- ▶ complete step-by-step instructions at
 - ▷ <https://www.sherlock.stanford.edu/docs/software/using/mariadb/>
 - ▷ <https://www.sherlock.stanford.edu/docs/software/using/postgresql/>

dir=\$DB_DIR

de, note
ule, and

rs/

Table of contents

Introduction

More documentation

MariaDB on Sherlock

Single-node access

Preparation

Start the server

Run queries

Multi-node access

Enable network access

Secure access

Start MariaDB in a job

Connect to the running instance

Persistent DB instances

Command **cache**ing

- user: “*are my jobs running? Let's* `watch -n.1 squeue`”
- `slurmctld`: “*aaargh!*”

```
# Pseudo shell
_cache_cmd() { # $1: cmd, $2: cache lifetime (sec)
  if cached_output is still valid; then
    echo "$cached_output"
  else
    execute $1 and cache output for $2 seconds
  fi
}

# define command alias
squeue() { _cache_cmd /usr/bin/squeue 10 "$@"; }
```

- ▶ goes in users' profile, and cache `squeue` results for 10s
- ▶ works for all status commands (`sstat`, `sacct`, `sinfo`)

Job script archiving

- user: “my job failed! why?”
- sysadmin: “what does your submission script look like?”
- user: “uuuh...”
- ▶ we archive job scripts
 - ▷ in PrologSlurmctld

```
SLURM_SPOOLDIR="/var/spool/slurm.state"  
SLURM_JOBSTORE="/share/admin/logs/slurm/jobs"  
  
jobid=${SLURM_ARRAY_JOB_ID:-$SLURM_JOBID}  
j_hsh=${jobid: -1}  
cp $SLURM_SPOOLDIR/hash.${jobid: -1}/job.$jobid/environment $JOB_STORE/env  
cp $SLURM_SPOOLDIR/hash.${jobid: -1}/job.$jobid/script $JOB_STORE/script
```

Job submit plugin

- ▶ for a long time: *we don't need no job submit plugin!*
- ▶ nowadays: *how did we ever worked without it?*
- ▶ we use it to:
 - ▷ automatically assign licenses to jobs
 - ▷ drop un-authorized options (`--reboot`, `--exclusive` in shared partitions)
 - ▷ provide helpful messages to users (when their job is rejected)

```
[kilian@sh-ln04 login ~]$ srun -p gpu --pty bash
srun: error: =====
srun: error: ERROR: missing GPU request, job not submitted
srun: error: =====
srun: error: Jobs submitted to the gpu partition must explicitly request GPUs, by using
srun: error: the --gres option.
srun: error: -----
srun: error: QOSMinGRES
srun: error: Unable to allocate resources: Job violates accounting/QOS policy (job submit
limit, user's size and/or time limits)
```

GPU mode SPANK plugin

- ▶ GPU compute modes

- ▷ EXCLUSIVE_PROCESS by default, good for most cases
- ▷ but some applications *need* multiple contexts on GPU
- ▷ changing GPU compute mode requires `root`

- ▶ https://github.com/stanford-rc/slurm-spank-gpu_cmode

- ▷ SPANK plugin to let users specify the GPU compute mode they need

```
$ srun --help
  --gpu_cmode=<shared|exclusive|prohibited>
          Set the GPU compute mode on the allocated GPUs to
          shared, exclusive or prohibited. Default is
          Exclusive

$ srun --gres gpu:1 --gpu_cmode=shared "nvidia-smi --query-gpu=compute_mode --format=csv,noheader"
Default
$ srun --gres gpu:1 --gpu_cmode=exclusive nvidia-smi --query-gpu=compute_mode --format=csv,noheader
Exclusive_Process
```

Expected **wait times**

- user: “*how long will my job wait?*”
- sysadmin: “*why don't we ask the scheduler?*”

```
$ ssh sherlock
[...]
```

Sherlock	status	OPERATIONAL	uptime : 99.989%
	usage	normal: 98.44%	use/tot: 1,764/ 1,792 cores
		global: 93.66%	use/tot: 26,732/28,540 cores

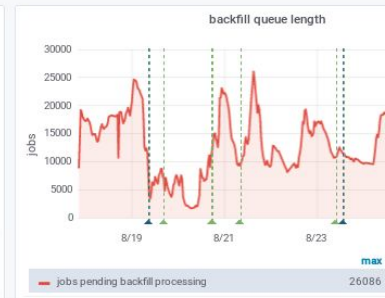
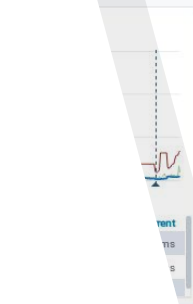
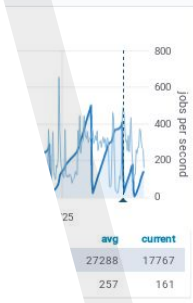
kilian	cur.jobs	0 RUNNING (0 core), 0 PENDING (0 core)	
	job wait	16 hours and 7 minutes in normal	

```
-----
```

- ▶ profile runs `srun --test-only` at login
 - ▷ gives a user-specific estimation of typical job wait time, including fairshare
 - ▷ consumes one jobid, but `_(ツ)_/`

Slurm dashboards

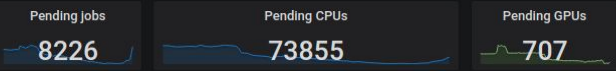
- ▶ overall usage
- ▶ scheduler internals
- ▶ queue info
- ▶ nodes states
- ▶ node utilization
- ▶ partition usage



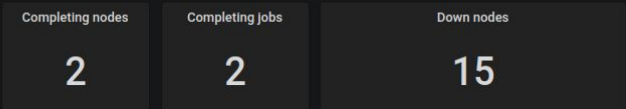
Running



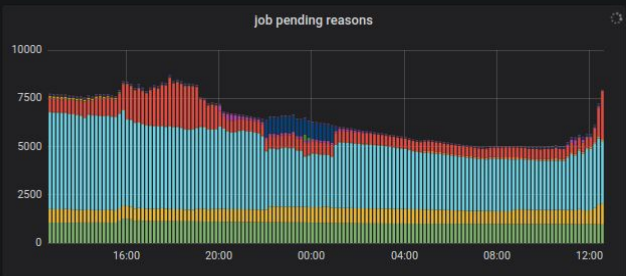
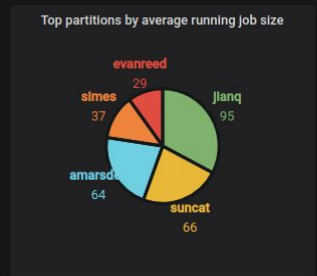
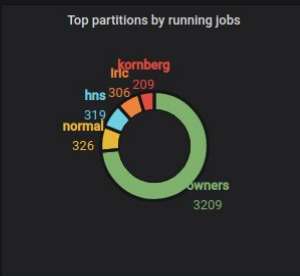
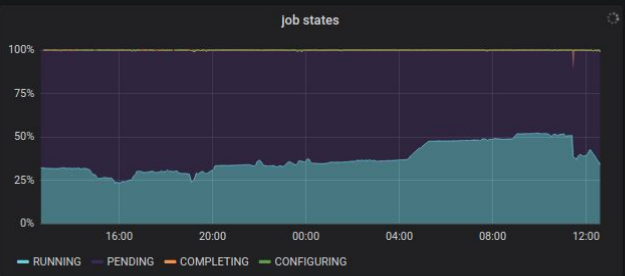
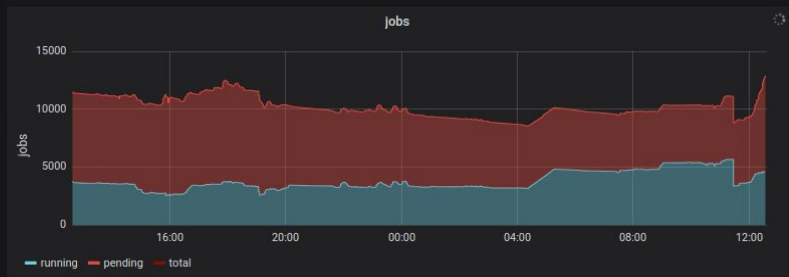
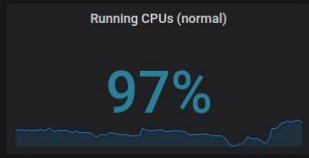
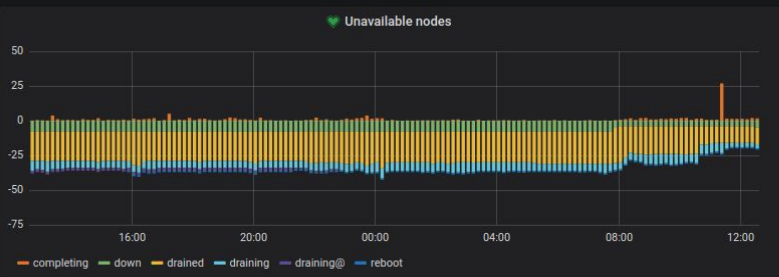
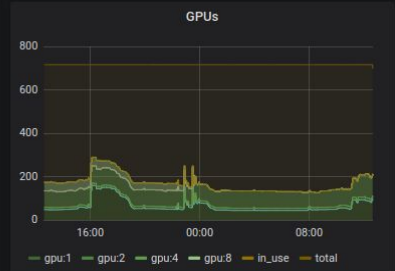
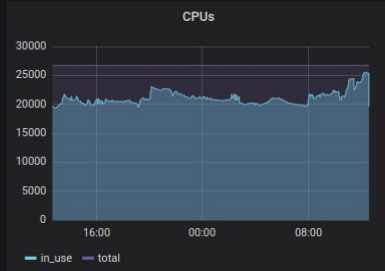
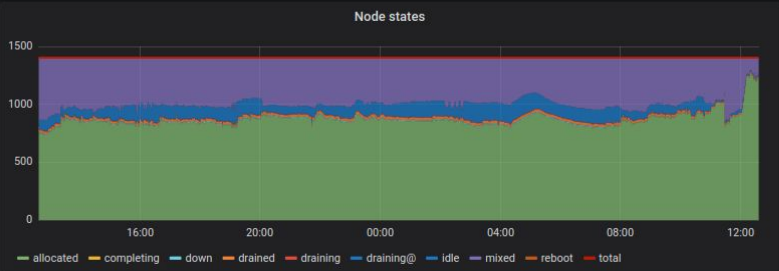
Waiting

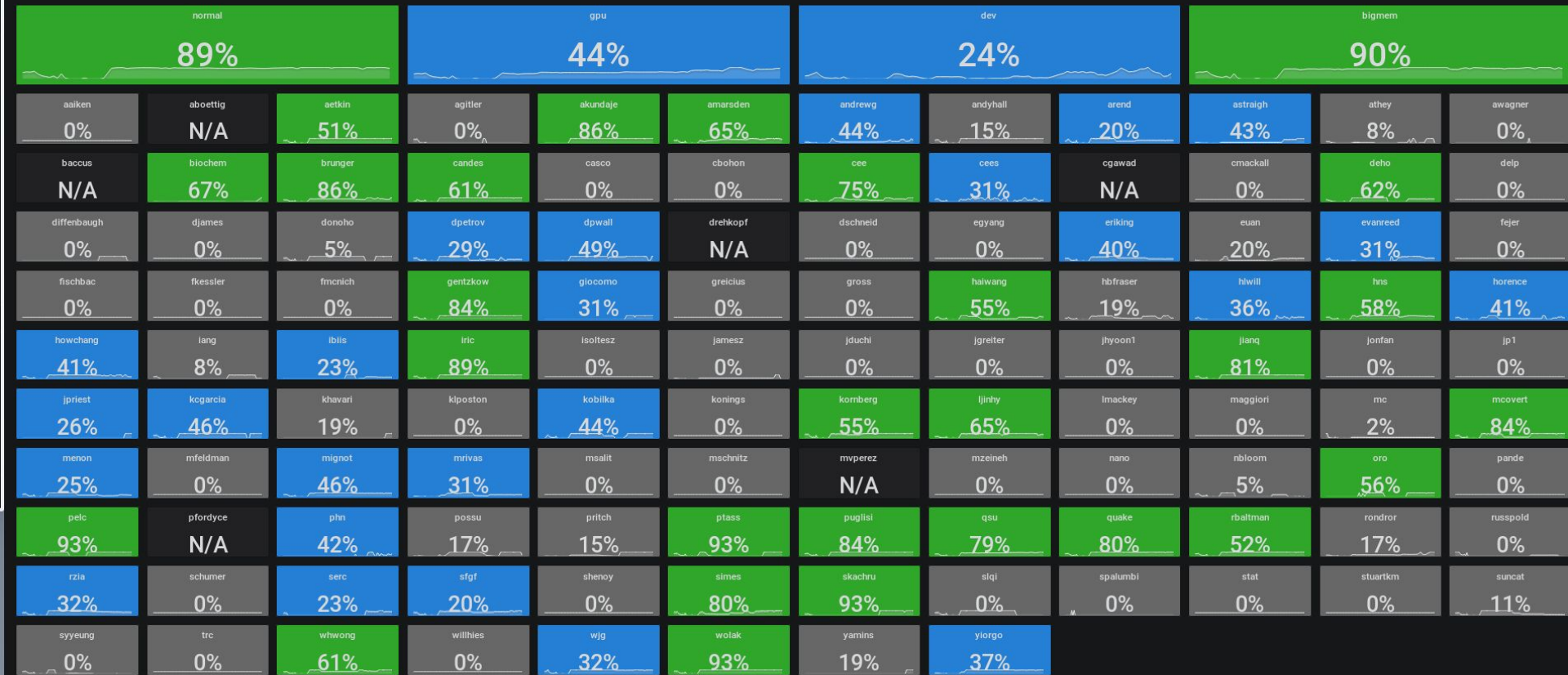


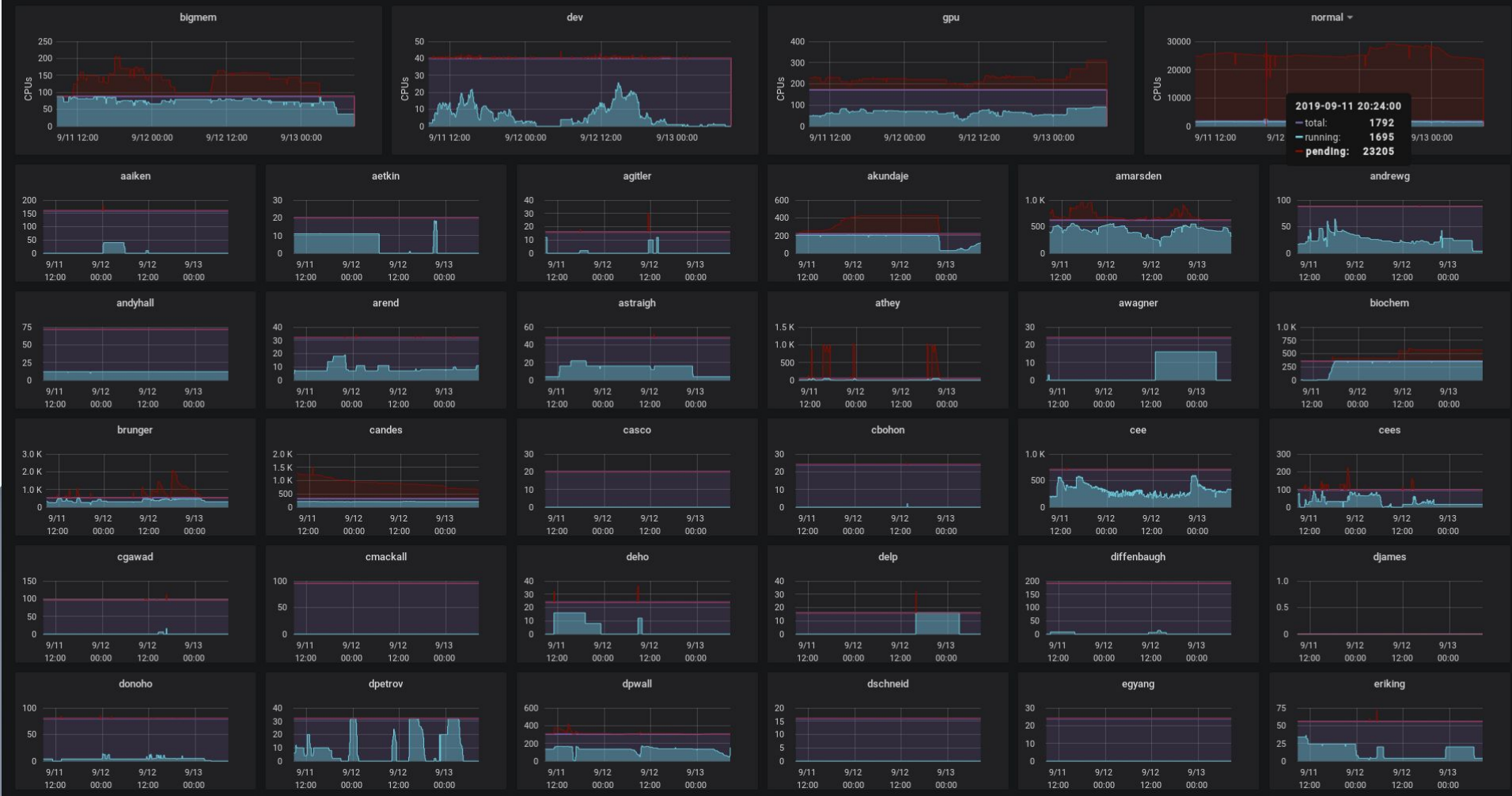
Down



Idle



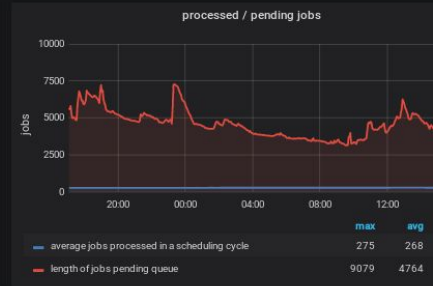
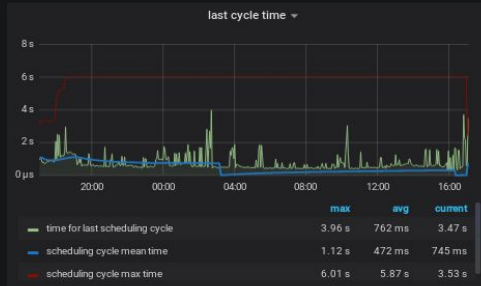
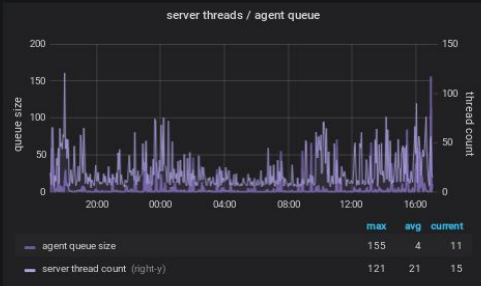




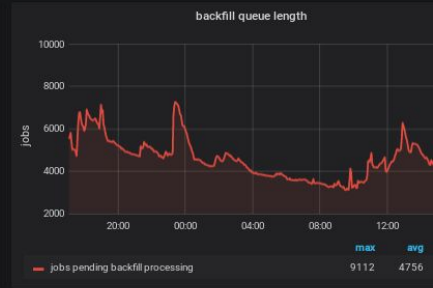
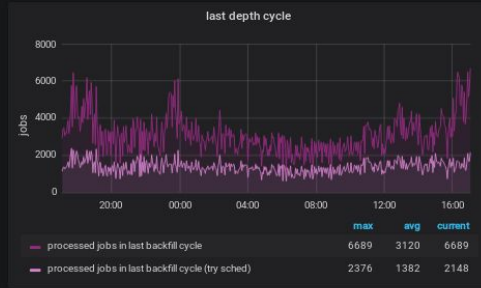
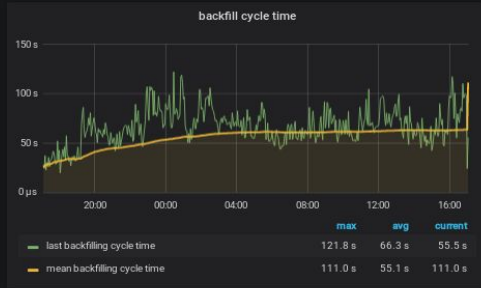
Jobs



Scheduler



Backfill Scheduler



Feedback

- ▶ Slurm is an amazing piece of software
 - ▷ 100,000 jobs/day on heterogeneous system with 1,000s of users and 100s of partitions → 95% utilization
- ▶ support is stellar
 - ▷ yes, you need support contract with SchedMD)
 - ▷ much better than much larger ISVs or HW vendors
 - ▷ a few *site down* situations resolved in a matter of hours

Feedback

- ▶ but... there **were** *site down* situations
 - ▷ a significant number of segfaults
 - ▷ major version updates are sometimes a challenge
 - ▷ DB conversions, unexpected config changes...
 - ▷ x.0 versions not for the faint of heart :)
- ▶ bugs.schedmd.com is a great source of information
 - ▷ can find a lot of details about Slurm internals there
 - ▷ we use it a lot (reported 122 bugs/requests so far)
 - ▷ did you know there's a bug report RSS feed?

THANKS!

Any question?

You can find me at kilian@stanford.edu

<https://github.com/stanford-rc>



Stanford
University