

Slurm 21.08 and Beyond

Tim Wickberg
SchedMD



Slurm User Group Meeting 2021

Agenda

All times are US Mountain Daylight (UTC-6)

Time	Speaker	Title
9:00 - 9:50	Jason Booth	Field Notes 5: From The Frontlines of Slurm Support
10:00 - 10:25	Nate Rini	REST API <i>and also</i> Containers
10:30 - 10:50	Marshall Garey	burst_buffer/lua and slurmscriptd
11:00 - 11:25	Nick Ihli	Slurm in the Clouds
11:30 - 11:50	Tim Wickberg	Slurm 21.08 and Beyond

Welcome




- Five separate presentations, five separate streams
- Presentations will remain available for at least two weeks after SLUG'21 concludes
- Presentations are available through the SchedMD Slurm YouTube channel
 - <https://youtube.com/c/schedmdslurm>
- Or through direct links from the agenda
 - https://slurm.schedmd.com/slurm_ug_agenda.html

Asking questions



- Feel free to ask questions throughout through YouTube's chat
- Chat is moderated by SchedMD staff
 - Tim McMullan and Ben Roberts
 - Also identified by the little wrench symbol next to their name
- Questions will be relayed to the presenter by the moderators
 - Some may be deferred to the end if they cannot be relayed in a timely fashion
 - Or some may be answered by the moderators in chat directly



Slurm 21.08 and Beyond

Tim Wickberg
SchedMD

Slurm Releases




- 20.11 - November 2020 - released during SC'20 BoF!
- 21.08 - August 2021
- 22.05 - May 2022
- 23.03 - February 2023

Slurm Release Schedule



- Slurm major releases come out every nine months
- Major release numbers are the two digit year, period, two digit month
 - 21.08 \Rightarrow 2021, August
- Maintenance releases, such as 21.08.1, come out roughly monthly for the most recent major release
- Two most recent major releases are still supported
 - These are 21.08 and 20.11 currently



Slurm 21.08 Release

Job submission command



- Un-modified command now captured by default
- Available through new '-o SubmitLine' output format in sacct

Job submission command

```
tim@blackhole:~$ sacct -o JobId,User,SubmitLine%50
```

JobID	User	SubmitLine
-----	-----	-----
1358	tim	sbatch --wrap sleep 1000 --exclusive
1358.batch		
1359	tim	sbatch --wrap sleep 1000 --exclusive -N 2

Store batch scripts in SlurmDBD



- New AccountingStoreFlags=job_script option in slurm.conf
 - As well as new AccountingStoreFlags=job_env
- 'sacct --batch-script' and 'sacct --env-vars' to fetch them

Store batch scripts in SlurmDBD

```
tim@blackhole:~$ sbatch --wrap "sleep 1000" --exclusive -N 2
```

```
Submitted batch job 1360
```

```
tim@blackhole:~$ sacct --batch-script -j 1360
```

```
Batch Script for 1360
```

```
-----
```

```
#!/bin/sh
```

```
# This script was created by sbatch --wrap.
```

```
sleep 1000
```

```
tim@blackhole:~$
```

New "PLANNED" node state



- PLANNED now shown instead of IDLE for nodes that are being held empty while waiting for a multi-node job to launch

New "PLANNED" node state

Slurm 20.11:

tim@blackhole:~\$ squeue

JOBID	PARTITION	NAME	USER	ST	TIME	NODES	NODELIST (REASON)
1359	general	wrap	tim	PD	0:00	2	(Resources)
1358	general	wrap	tim	R	1:49	1	node0004

tim@blackhole:~\$ sinfo

PARTITION	AVAIL	TIMELIMIT	NODES	STATE	NODELIST
general*	up	4:00:00	1	idle	node0005
general*	up	4:00:00	1	alloc	node0004

New "PLANNED" node state

Slurm 21.08:

tim@blackhole:~\$ squeue

JOBID	PARTITION	NAME	USER	ST	TIME	NODES	NODELIST (REASON)
1359	general	wrap	tim	PD	0:00	2	(Resources)
1358	general	wrap	tim	R	1:49	1	node0004

tim@blackhole:~\$ sinfo

PARTITION	AVAIL	TIMELIMIT	NODES	STATE	NODELIST
general*	up	4:00:00	1	plnd	node0005
general*	up	4:00:00	1	alloc	node0004

RS256 token support in auth/jwt



- Keys specified through a JWKS file
 - Such as those generated by AWK Cognito
- AuthAltParameters=jwks=/path/to/my.jwks
- May be used alongside existing HS256 support

RS256 token support in auth/jwt

```
tim@blackhole:~$ grep jwks /etc/slurm.conf
AuthAltParameters=jwks=/etc/slurm/jwks.json,jwt_key=/root/jwt_hs256.key
tim@blackhole:~$ cat jwks.json
{"keys":[{"alg":"RS256","e":"AQAB","kid":"fZqKj+4Zw9OhMC4XNtWWGQC8n8iDxVoy6HLMLkONNuY=",
"key":"7Lm5UDivRbAXNQ9-F15vVty1fA1jTTRrN9RJTlXoiFMJPGfgWqDHOWAIO2OtQur3bsGMckUQ_
7ZbRwZnbtMeDZ-QGAb-gWJ5mjxCegRD0xPC9QoulZzNDm3oB_56jsMDRuYUI6Q0qvC3QiXzurmNtUJwmRhE1mlTwQ
wc5b-b8mJBYHjIW3ROAAe3Onr9T7NPenQ1BzOi8DKYo35RwJEQYCz0hRsX2cpztOhBTDU5nvgkY1I6f1bQtgpmT6j
Z1HFjjX7IQGVCijU0W3F_rj-0JAccmFlskog3Vynos0cA7WRvQdJc2iMulznBAoeLsNRJ0rp0A361APDQQdcnoeI7C
9w","use":"sig"}],{"alg":"RS256","e":"AQAB","kid":"/zFkNPInxOO+4p7u2ccOSLQnMMxaulgPRr+3/0j
1YMs=", "key":"vMo6Ad50H8wOEvwIYyRXVXH7wB-aob9Um1GG2W-XCY4Eb7bSoqMDBTZZZgCb1IAzG
megs7QXuA50699Jfs0LrupC9TVB_zWkiU4DAIdB9RUEsBubmPCDJMobSK3L4UWVnqGdSf_c078CyyoumNSFhwRddo
tdzAKglRxMiCzvy3Zgldx3l3iNpeQRUTWJ_x8Du5eiirjqB4zdof9vwQ_DFVP0c9zRWZSheV7XD3lnqv1sBMVYZs
DxX_FBGU5flG8ExIZV2pV0jbHva7N1V6k3J69rwYfG5E9-d-JZKEXyIFMHPAl8zZQmUgEvXVusIJe6STJLKHgSZAw
a-eFKiQV6w","use":"sig"}]}
```

Improved cgroup subsystems



- Significant refactoring work for the task/cgroup, proctrack/cgroup, and jobacct_gather/cgroup plugins
- Still only supports cgroup v1
- All cgroup interactions now handled centrally
 - Preparation for future cgroup v2 support

burst_buffer/lua

- "Generic" "Burst Buffer" support
- Really a means of handling pre- and post- job setup
 - Asynchronously
 - Compute nodes not yet assigned
- Avoids wasting compute node time for large job starts by handling setup and teardown tasks while they are still running other jobs

burst_buffer/lua

- See separate presentation by Marshall for further details
- ... any suggestions for a better name than "burst_buffer"?

New 'slurmscriptd' process



- Developed alongside burst_buffer/lua
- fork()+exec() in slurmctld is **very** expensive for systems with high-throughput and high job counts
- Instead, the slurmscriptd process launches scripts on behalf of slurmctld
 - Limited to burst_buffer and PrologSlurmctld/EpilogSlurmctld
 - Expect to expand this in the future

Fixes to job_container/tmpfs

- job_container/tmpfs was snuck into the 20.11.5 maintenance release early
- Strong early adoption exposed a number of design issues with how the slurmd/slurmstepd shared responsibility for the namespaces
- Fixed with further refactoring in 21.08

json and yaml output



- sacct, sinfo, and squeue now have --json/--yaml options for output
- Uses same underlying serialization/translation code as slurmrestd, but in the standalone command
- Output only

json and yaml output

```
tim@blackhole:~$ sacct --json|jq .jobs|head -n 14
```

```
[
  {
    "account": "root",
    "comment": {
      "administrator": null,
      "job": null,
      "system": null
    },
    "allocation_nodes": 1,
    "array": {
      "job_id": 0,
      "limits": {
        "max": {
          "running": {
```

Shared libraries and 'srun --bcast'



- Added new feature to 'srun --bcast' to allow it to automatically identify and broadcast required shared libraries as part of job launch
- Creates a directory alongside the broadcasted executable, and prepends that into LD_LIBRARY_PATH as part of step launch
- Avoids "thundering herd" issues on parallel filesystems on massively parallel job launches
 - Single srun process reads the executable and libs

Shared libraries and 'srun --bcast'

- Enabled through BcastParameters=send_libs
 - Disabled by default
 - Or through 'srun --bcast --send-libcs ./my_program'
- New BcastExclude option can set system library directories to ignore
 - Defaults to "/lib,/usr/lib,/lib64,/usr/lib64"
 - No point in sending ld-linux-x86-64.so.2 or libc.so.6

OCI Container Support



- Initial support for launching processes in OCI containers
- See Nate's presentation for further details

Improved Job Step Throughput



- Significant performance improvements to job step launch
- Nicely complements past performance work with
`SlurmctldParameters=enable_rpc_queue`



Slurm 22.05 Roadmap

One Note



- Our published roadmap **only** includes committed development work
- SchedMD has several exciting projects in the works for 22.05, but unfortunately we can't share them yet

One Note



- We prefer this "no vaporware ever" approach
- Even though it means the roadmap is a bit sparse when contracts are in progress

"Preferred" node constraints

- A list of optional ("soft") constraints to be considered when selecting nodes for a job
 - Likely using "--prefer" as the option to salloc/sbatch/srun
 - Job launch will prefer those nodes, but fall back to any nodes if that cannot be satisfied immediately
 - Traditional "hard" constraints (--constraint) will always be respected


GPU Sharding

- Allow for cooperative GPU sharing between separate jobs
- Allows administrators to define a number of "Slices" for a GPU
 - Jobs can request between zero and all slices
 - All slices allocated to the job from a single GPU, cannot span between cards
- Caveat: no hardware enforcement
 - Jobs must cooperate effectively

AcctGatherInterconnect plugins



- Add support for gathering network statistics from OmniPath and Slingshot interconnects



Slurm 23.02 Roadmap

Truly Dynamic Nodes



- Move away from current FUTURE node handling
 - Support truly dynamic node addition and removal from the cluster
- Some underlying work will be in 22.05, but will not be ready until 23.02



Upcoming



- Slurm booth on the SC21 show floor - #3215
- Birds of a Feather Session ("the BoF") has been accepted
 - Fully virtual session has been requested



Questions?

End Of Stream



- Thanks for watching!