# Federated Cluster Support

Brian Christiansen
SchedMD

SC 17

# Background

- Slurm has long had limited support for federated clusters
  - Most commands support a "--cluster (-M)" option to route requests to different clusters
- Submitted jobs are routed to one cluster
- Each cluster operates independently
  - Job IDs on each cluster are independent (two jobs can have same ID)
  - No cross-cluster job dependencies
  - No job migration between clusters
  - No unified view of system state, each cluster largely independent

# New Capabilities

- **Job Distribution**
  - Jobs distributed across federation
  - Unique job IDs
- **Unified Views**
  - Appear as one cluster
- **Easy Administration**
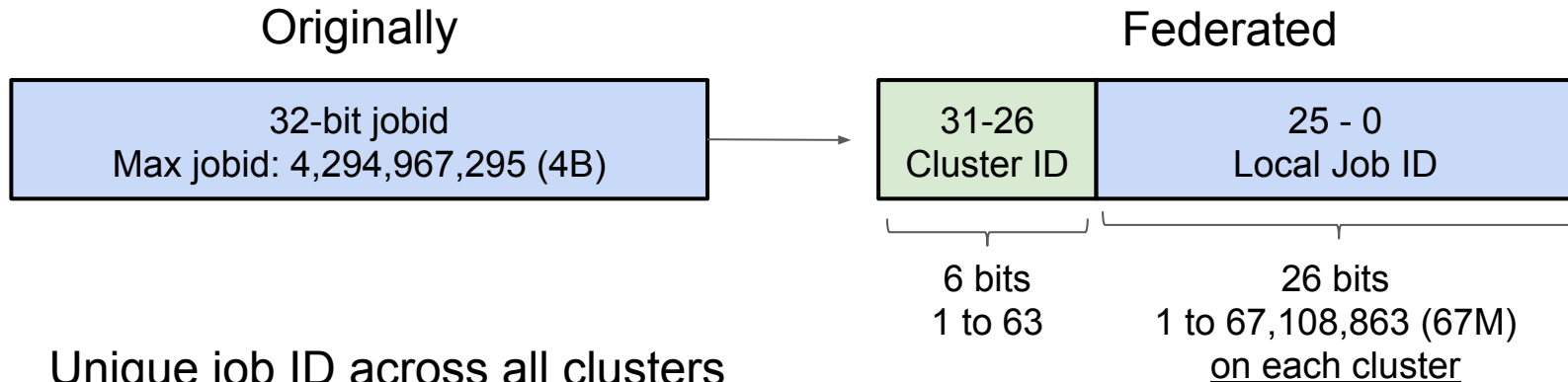  - Add/remove clusters to/from the federation with database commands

# Design Goals

- **Performance**: Little to no reduction in throughput of each cluster, performance scales with cluster count
- **Scalability**: No reduction in scalability of individual clusters, able to support many federated clusters
- **Ease of use**: Unified enterprise-wide view, minimize change in user interface
- **Stability**: No change in behavior for clusters not explicitly placed into a federation

# Eliminating the Bottlenecks

- Need mechanism to identify the cluster associated with a job ID without using slurmdbd lookup
- Make use of 32-bit job ID

| Originally | | Federated | |
|---|---|---|---|

| 32-bit jobid<br>Max jobid: 4,294,967,295 (4B) | → | 31-26<br>Cluster ID | 25 - 0<br>Local Job ID |
|---|---|---|---|

6 bits
1 to 63

26 bits
1 to 67,108,863 (67M)
on each cluster

- Unique job ID across all clusters
  - Large but unique: 134296327 (ClusterID:2 + JobID:78,599)

# Eliminating the Bottlenecks

- The max MaxJobID is now 67,108,863. Any pre-existing jobs will continue to run but new job ids will be within the new MaxJobID range. Adjust your configured MaxJobID value as needed to eliminate any confusion.
- Change was made in 17.02.

# Configuration

- sacctmgr add federation <fedname> [clusters=<list>]
  - sacctmgr mod federation <fedname> clusters[+|-]=<list>
- sacctmgr mod cluster <cluster> ...
  - set federation=<federation>
  - set features=<features>
    - Features at a cluster level that can be requested by jobs
- Slurm.conf:
  - FederationParameters=fed_display to show federated view by default

# Configuration

- sacctmgr mod cluster <cluster> ...
  - set fedstate=<state>
    - ACTIVE: Cluster will actively accept and schedule federated jobs.
    - INACTIVE: Cluster will not schedule or accept any jobs.
    - DRAIN: Cluster will not accept any new jobs and will let existing federated jobs complete.
    - DRAIN+REMOVE: Cluster will not accept any new jobs and will remove itself from the federation once all federated jobs have completed. When removed from the federation, the cluster will accept jobs as a non-federated cluster.

# Configuration

- A cluster can only be part of one federation at a time
- Jobs can't span clusters

# Persistent Connections

- Clusters talk to each other over persistent connections
  - Reduces communication overhead -- only authenticate once
  - Broken connections detected immediately and established when needed
  - Controller and SlurmDBD use the same code

# Job Submission

- sbatch, salloc, srun supported
- Jobs submitted to local cluster
- Sibling jobs submitted to all "viable" clusters
  - viable == all clusters ||
  - --clusters=<clusters> & --cluster_constraint=<features>
- Job stays on the local cluster -- even if not viable -- to coordinate and route requests to/from sibling clusters
  - Job starts, updates, cancellations

# Job Submission

- Interactive Jobs (salloc/srun) example
    1. srun submits the job to a local cluster in a federation
    2. Local cluster will submit sibling jobs to sibling clusters
    3. A sibling cluster coordinates with the origin cluster and allocates nodes for the job
    4. The origin cluster removes/revokes any pending sibling jobs
    5. The sibling cluster directly notifies the srun, passing sibling and allocation information
    6. srun will talk directly to allocated nodes

NOTE: All compute nodes need to be accessible by each submission host!

NOTE: Interactive jobs can not be requeued

# Scheduling

- Federated jobs contain the locations of all "sibling" jobs
- Each cluster independently schedules each sibling job
- Coordinates with "origin" cluster to start job
  - The origin cluster is determined from the job id
  - Prevents multiple jobs from being started at the same time
  - Policies in place to handle if origin cluster fails
- Once sibling job is started, origin cluster revokes remaining siblings jobs
- Batch jobs can be requeued to federation

# Unified Views

- Unified views provided with --federation command line option
  - Made default with  FederationParameters=fed_display slurm.conf option
  - squeue, sinfo, sacct, sreport, sview etc.
  - --local, --clusters/-M options override federated view

```
$ export
SQUEUE_FORMAT2=jobarrayid:8,cluster:.8,statecompact:.4,origin:.8,siblingsviable:.16,siblingsactive:.16,timeu
sed:.8,numnodes:.6,nodelist:.12,reason:.15
$ squeue
JOBID     CLUSTER  ST  ORIGIN  VIABLE_SIBLINGS  ACTIVE_SIBLINGS   TIME NODES NODELIST       REASON
20132665    fed1   PD    fed3        fed1,fed3        fed1,fed3   0:00    5                  Priority
67109269    fed1   PD    fed1   fed1,fed2,fed3   fed1,fed2,fed3   0:00    5                  Resources
13421784    fed1   PD    fed2   fed1,fed2,fed3   fed1,fed2,fed3   0:00    5                  Priority
20132665    fed3   R     fed3        fed1,fed3             fed3   2:44    5 fed3_[6-10]      None
13421784    fed3   R     fed2   fed1,fed2,fed3             fed3   2:47    5 fed3_[1-5]       None
20132665    fed2   R     fed3             fed2             fed2   2:50    5 fed2_[1-5]       None
67109268    fed2   R     fed1   fed1,fed2,fed3             fed2   2:50    5 fed2_[6-10]      None
13421783    fed1   R     fed2             fed1             fed1   2:54    5 fed1_[6-10]      None
67109267    fed1   R     fed1   fed1,fed2,fed3             fed1   2:57    5 fed1_[1-5]       None
```

# Future

- Cross-cluster job dependencies
- Job-arrays span clusters in the federation (currently locked to local cluster)

# Questions?