# SLURM Version 2.3 and Beyond
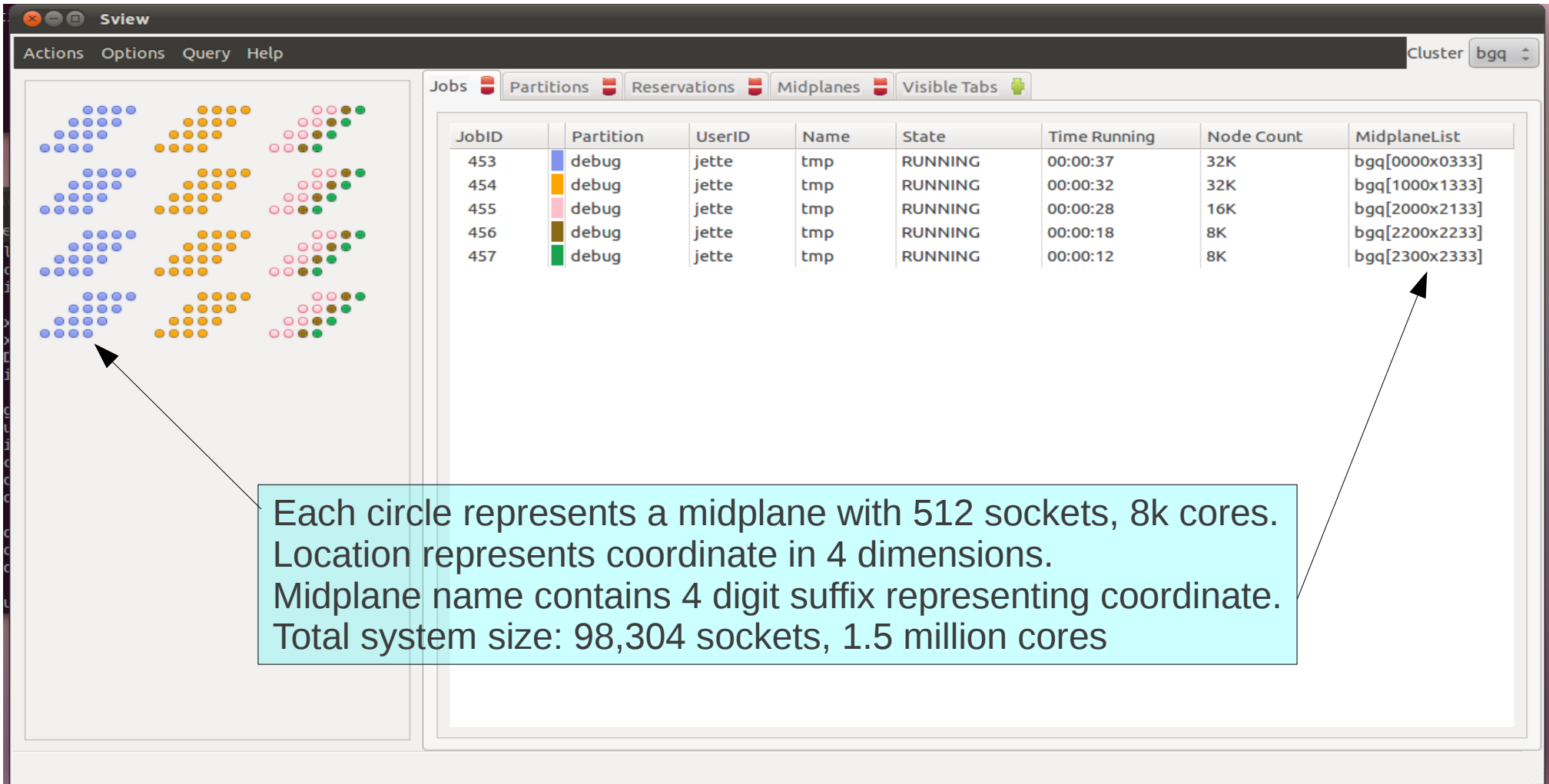
Morris Jette
jette@schedmd.com

SchedMD LLC

# SLURM Version 2.3

- Released September 9, 2011

- New systems supported:

  - Cray XE and XT systems

    - Runs over ALPS

    - Provides SLURM scheduling functions and *srun* wrapper for *aprun* for task launch

  - IBM BlueGene/Q systems (incomplete support)

    - Major changes from BlueGene/P

      - Completely new IBM API

      - 5-dimension torus topology

    - Work to be complete in version 2.4

# IBM BlueGene/Q and sview
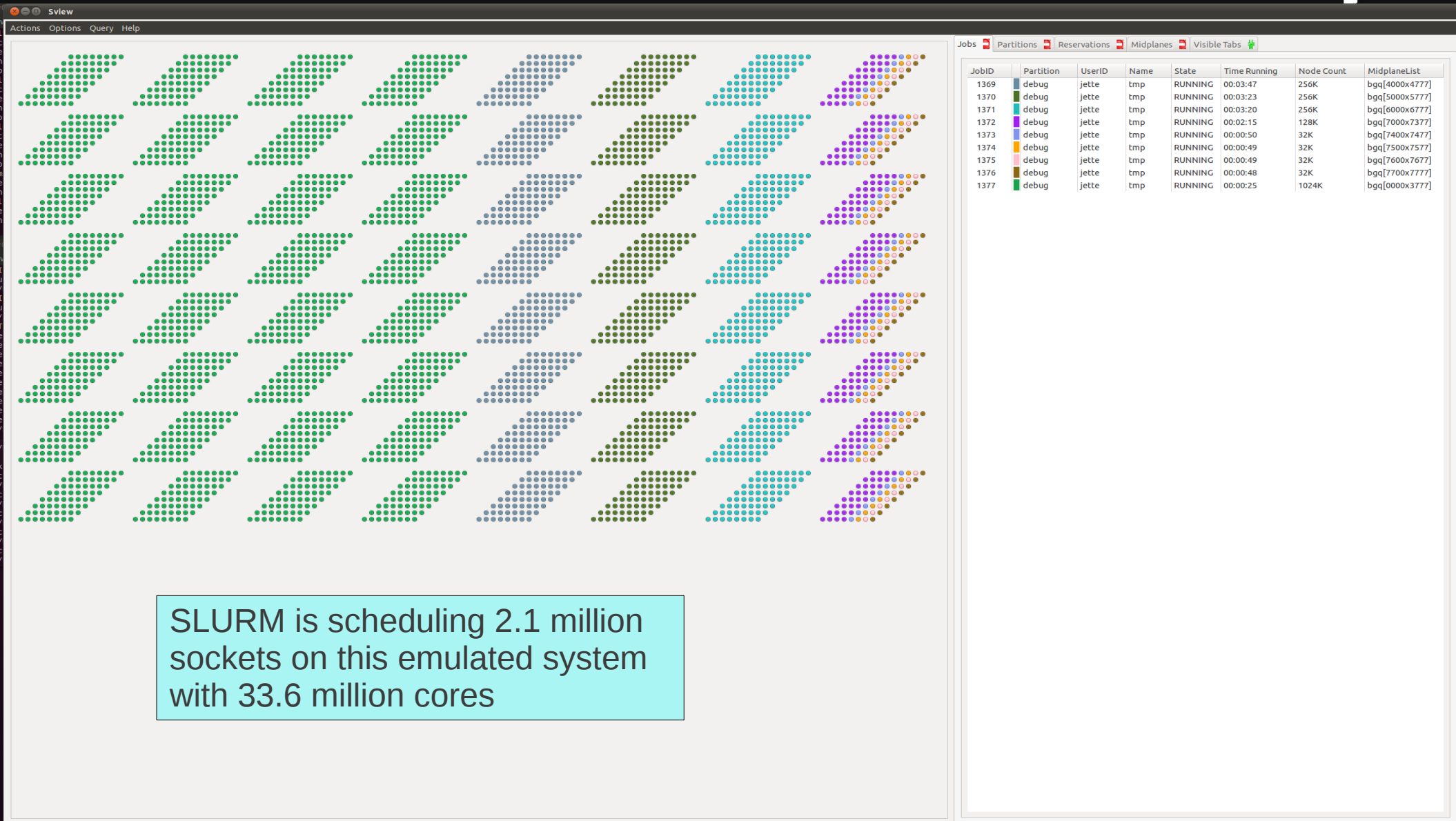## (LLNL's 96 rack Sequoia)



| JobID | Partition | UserID | Name | State | Time Running | Node Count | MidplaneList |
|-------|-----------|--------|------|-------|--------------|------------|--------------|
| 453 | debug | jette | tmp | RUNNING | 00:00:37 | 32K | bgq[0000x0333] |
| 454 | debug | jette | tmp | RUNNING | 00:00:32 | 32K | bgq[1000x1333] |
| 455 | debug | jette | tmp | RUNNING | 00:00:28 | 16K | bgq[2000x2133] |
| 456 | debug | jette | tmp | RUNNING | 00:00:18 | 8K | bgq[2200x2233] |
| 457 | debug | jette | tmp | RUNNING | 00:00:12 | 8K | bgq[2300x2333] |

Each circle represents a midplane with 512 sockets, 8k cores.
Location represents coordinate in 4 dimensions.
Midplane name contains 4 digit suffix representing coordinate.
Total system size: 98,304 sockets, 1.5 million cores

# Exascale BlueGene/Q
## (2,048 rack system)



SLURM is scheduling 2.1 million sockets on this emulated system with 33.6 million cores

| JobID | Partition | UserID | Name | State | Time Running | Node Count | MidplaneList |
|-------|-----------|--------|------|-------|--------------|------------|--------------|
| 1369 | debug | jette | tmp | RUNNING | 00:03:47 | 256K | bgq[4000x4777] |
| 1370 | debug | jette | tmp | RUNNING | 00:03:23 | 256K | bgq[5000x5777] |
| 1371 | debug | jette | tmp | RUNNING | 00:03:20 | 256K | bgq[6000x6777] |
| 1372 | debug | jette | tmp | RUNNING | 00:02:15 | 128K | bgq[7000x7377] |
| 1373 | debug | jette | tmp | RUNNING | 00:00:50 | 32K | bgq[7400x7477] |
| 1374 | debug | jette | tmp | RUNNING | 00:00:49 | 32K | bgq[7500x7577] |
| 1375 | debug | jette | tmp | RUNNING | 00:00:49 | 32K | bgq[7600x7677] |
| 1376 | debug | jette | tmp | RUNNING | 00:00:48 | 32K | bgq[7700x7777] |
| 1377 | debug | jette | tmp | RUNNING | 00:00:25 | 1024K | bgq[0000x3777] |

# SLURM Version 2.3

- Support added for multiple front-end nodes

  - Improves fault-tolerance and performance for Cray and BlueGene systems

  - Jobs allocated to front-end nodes on a round-robin basis

  - New configuration file options

  - Scontrol modified to get/set front-end node state information

# SLURM Version 2.3

- Added ability to set default and maximum memory limits per partition instead of one value for the entire cluster

  - Provides better gang scheduling control (e.g. time-slice some partitions and not others)

- Added *GraceTime* to Partition and QOS data structures for job preemption

  - Gives job opportunity to gracefully terminate once preempted

- New plugins support Linux cgroup job container

  - Identifies and controls the processes in a job

  - Restrict use of CPUs, memory and device files

# SLURM Version 2.3

- Jobs can control network topology

  - Maximum number of leaf switches and maximum wait for that configuration

- Only current job dependencies are displayed

  - Satisfied dependencies are hidden for easier use

- Better estimates of pending job's start time

- Added ability to expand job sizes

  - Requires submission of new job that merges its resources into another job's resources

# Job Expansion

```
$ salloc -N1 bash
salloc: Granted job allocation 65542
$ srun hostname
icrm1
```
⎱ Create original job allocation

```
$ salloc -N1 --dependency=expand:$SLURM_JOBID bash
salloc: Granted job allocation 65543
```
⎱ Create allocation for expanding original job

```
$ scontrol update jobid=$SLURM_JOBID NumNodes=0
To reset SLURM environment variables, execute
  For bash or sh shells:  . ./slurm_job_65543_resize.sh
  For csh shells:       source ./slurm_job_65543_resize.csh
$ exit
exit
salloc: Relinquishing job allocation 65543
```
⎱ Transfer additional resources to original job

```
$ scontrol update jobid=$SLURM_JOBID NumNodes=ALL
To reset SLURM environment variables, execute
  For bash or sh shells:  . ./slurm_job_65542_resize.sh
  For csh shells:       source ./slurm_job_65542_resize.csh
$ . ./slurm_job_$SLURM_JOBID_resize.sh
```
⎱ Update original job's environment variables (node count, node list, etc.)

```
$ srun hostname
icrm1
icrm2
$ exit
exit
salloc: Relinquishing job allocation 65542
```
⎱ Use expanded allocation

# SLURM Version 2.4 Plans

- Available 2$^{nd}$ quarter 2012

  - Pre-releases available monthly for development and test: *http://www.schedmd.com/#repos*

  - Latest code:  *https://github.com/SchedMD/slurm*

- Complete SLURM port to IBM BlueGene/Q

  - Work remaining for multiple jobs per block

    – Each c-node can run a different user's job

    – 5-dimensional torus supports very efficient job packing
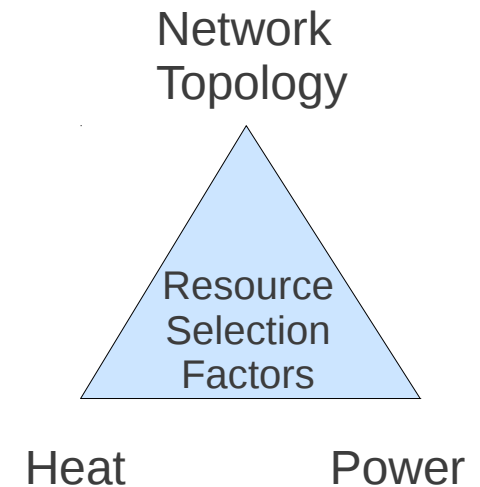
  - Work remaining for fault tolerance

# SLURM Version 2.4 Plans

- Enhanced job constraint support

  - Support multiple constraint counts:
    "--constraint=[rack1*2&rack2*2]"

- Cloud Bursting: Move overflow work to the cloud

  - Allocate, boot and start SLURM daemons in cloud

  - Add resources on demand, release idle resources

- Interface to IBM/Tivoli LoadLeveler

# Future Directions

- Power Management
  - Collect job power usage, optionally change for power
  - Estimate power needs of pending jobs (user input + historic data)
  - Manage workflow within available/dynamic power envelope

- Heat Management
  - Collect temperature data
  - Distribute high-power jobs to minimize hot-spots

- Failure Management
  - Proactive and Interactivet

Network
Topology

Resource
Selection
Factors

Heat            Power

SLURM:      Nodes tux10123 and tux10125 are failing
*Application: Can you give me two replacement nodes now?*
SLURM:      I can give you one node now and one more in 5 minutes
*Application: Can you extend my time limit by 5 minutes?*
SLURM:      Yes                               Credit: William Kramer, NCSA

SchedMD LLC
http://www.schedmd.com