

Slurm Survival Guide

M. D. Jones, Ph.D.

Center for Computational Research
University at Buffalo
State University of New York

High Performance Computing I, Fall 2014

Links

Primary Slurm URLs:

- computing.llnl.gov/linux/slurm/slurm.html,
Slurm's original home at LLNL, now a bit outdated, historical reference
- www.schedmd.com,
Slurm's new primary site

These sites contain a lot of documentation and pointers to others, including various tutorials.

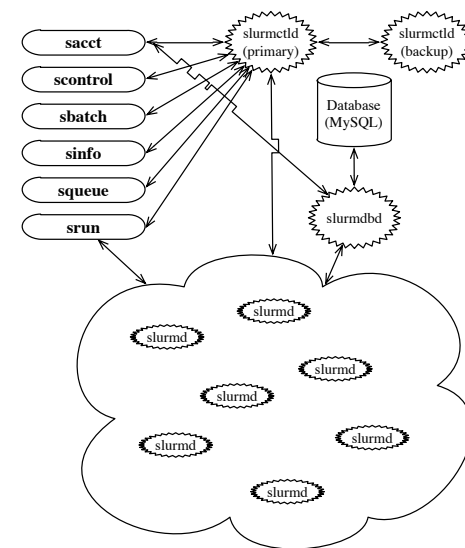
Slurm

Slurm (**S**imple **L**inux **U**tility for **R**esource **M**anagement) is a highly modular open-source software package for utilizing clusters in high-performance computing (HPC):

- Allocates exclusive and/or non-exclusive access to resources (computer nodes) to users for some duration of time for the user to utilize for their workload,
- Provides a framework for starting, executing, and monitoring work (typically a parallel job, but not necessarily) on a set of allocated nodes,
- Schedules resources by managing a queue of pending workloads submitted by users.

Slurm Architecture

Client commands, `sacct`, `scontrol` (mostly for administrative use), `salloc`, `sbatch`, `sinfo`, `squeue`, `srun` for user submissions and queries. Each compute node runs a **slurmd** daemon that is responsible for command and control on that node. **slurmctld** is the Slurm manager daemon, while **slurmdbd** handles recording the usage data (there is an alternative flat file option).



sinfo

`sinfo` is Slurm's client for querying clusters, nodes, and partitions. Note that partitions can overlap, so a node can belong to multiple partitions. If you run `sinfo` with just default options, on a large cluster you will see a large quantity of displayed information, so let us instead consider a simpler query, and focus on CCR's `debug` partition, which is a small set of nodes set aside for development and debugging work with rapid turnaround:

```
1 [rush:~]$ sinfo --partition=debug
2 PARTITION AVAIL  TIMELIMIT  NODES  STATE NODELIST
3 debug      up       1:00:00      7    idle d07n33s[01-02],d16n[02-03],k05n26,k08n41s[01-02]
```

From this we infer a wall time limit of 1 hour (hence the quicker turnaround for development jobs), and 7 nodes total - but what about more details on the nodes themselves?

squeue

`squeue` provides details on jobs in Slurm. On a very busy cluster, the default output can be overwhelming, much like `sinfo` it can be more helpful to carefully focus on points of interest. More often than not we are interested mainly in our own jobs:

```
1 [rush:~]$ squeue --users=$USER
2          JOBID PARTITION  NAME      USER ST      TIME  NODES NODELIST(REASON)
3 2482840  general-co  nb-mult  jonesm PD      0:00      2  (Resources)
```

The default output is rather terse, so if you need more detail you can resort to using the `--format` option:

```
1 [k07n14:~]$ squeue --format="%8i %8u %8a %10P %4D %4C %5m %9l %20S %6h %6Q %6b %12f %R" --users=$USER
2 JOBID  USER  ACCOUNT PARTITION  NODE CPUS MIN_M TIMELIMIT START_TIME  SHARED PRIORI GRES  FEATURES  NODELIST(REASON)
3 2482840 jonesm  ccrtech  general-co  2    32   3000  1:00:00  2014-07-17T11:31:57  no    50490  (null)  CPU-E5-2660  (Resources)
```

There are a host of options that you can specify for `sinfo`, this combination is handy for looking at the details on the nodes:

```
1 [rush:~]$ sinfo --partition=debug --long --Node
2 NODELIST      NODES PARTITION  STATE CPUS  S:C:T MEMORY TMP_DISK WEIGHT FEATURES REASON
3 d07n33s[01-02] 2    debug      idle   8    2:4:1  24000    0    23  IB,CPU-L none
4 d16n[02-03]    2    debug      idle   8    2:4:1  24000    0    23  IB,CPU-L none
5 k05n26          1    debug      idle  16    2:8:1 128000    0   102 CPU-E5-2 none
6 k08n41s[01-02] 2    debug      idle  12    2:6:1  48000    0    31  IB,CPU-E none
```

If you look at the man page for `sinfo`, you can see that the above command is just a handy short cut for using `sinfo --format`, with which we expand the feature set for the nodes:

```
1 [rush:~]$ sinfo --partition=debug --format="%14N %5D %9P %6T %3c %6z %6m %8d %6w %12f %6G %20E"
2 NODELIST      NODES PARTITION  STATE CPU  S:C:T MEMORY TMP_DISK WEIGHT FEATURES  GRES  REASON
3 d07n33s[01-02] 2    debug      idle   8  2:4:1  24000    0    23  IB,CPU-L5630 (null) none
4 d16n[02-03]    2    debug      idle   8  2:4:1  24000    0    23  IB,CPU-L5520 (null) none
5 k05n26          1    debug      idle  16  2:8:1 128000    0   102 CPU-E5-2660 gpu:2 none
6 k08n41s[01-02] 2    debug      idle  12  2:6:1  48000    0    31  IB,CPU-E5645 (null) none
```

sview

`sview` is a graphical user interface (GUI) that provides Slurm information on nodes, jobs, partitions, reservations, *etc.* - it provides much the same functionality as `squeue` and `sinfo`. Note that you have to have graphics forwarding enabled (X11 if you are accessing the cluster front end via `ssh`), or alternatively use a remote visualization system (although in this case the graphics acceleration is not really needed). Note that on a busy cluster, `sview` can be rather slow.

| Partition | Default | Part State | Time Limit | Node Count | Node State | NodeList |
|-----------------|---------|------------|------------|------------|------------|--|
| debug | no | up | 01:00:00 | 7 | idle | d07n33s02,d13s02,d13s02,k05n26,k05n41s01-02 |
| general-compute | yes | up | 3:00:00:00 | 703 | mixed | d07n04-11,13,20,24-31,34-40s01-02,d09n04-11,13,20,24-31,33-40s01-02,d13s03 |
| general-compute | | | | 118 | mixed | d07n06s02,d07n10s01,d07n11s01-02,d07n24s02,d07n25s02,d07n26s01,d07n30s01,d |
| general-compute | | | | 460 | allocated | d07n04s01-02,d07n05s02,d07n06s01,d07n07s01-02,d07n08s01-02,d07n09s01-02 |
| general-compute | | | | 119 | idle | d07n15s02,d07n16s01,d07n16s01,d07n20s01,d07n26s01,d07n36s01-02,d07n37s01,d |
| general-compute | | | | 3 | drainnet | d07n05s01,d09n04s01-02 |
| general-compute | | | | 1 | down* | k05n12s02 |
| general-compute | | | | 1 | drainnet | k05n20s02 |
| general-compute | | | | 1 | down | k14s02s01 |
| gpu | no | up | 3:00:00:00 | 32 | mixed | k05n11-12,20-21,30-31,39-40s01-02,k05n11-12,20-21,30-31,39-40s01-02 |
| gpu | | | | 11 | mixed | k05n11s01-02,k05n20s01-02,k05n21s01-02,k05n30s01-02,k05n31s01-02,k05n3 |
| gpu | | | | 20 | allocated | k05n11s01-02,k05n12s01-02,k05n20s01-02,k05n21s01-02,k05n30s01-02,k05n3 |
| gpu | | | | 1 | down* | k05n12s02 |
| largemem | no | up | 3:00:00:00 | 19 | mixed | k07n04-11,13,20,24-31,34-40s01-02,k05n26,k07n23-25,k05n25,k07n34 |
| largemem | | | | 9 | allocated | k07n34,k07n23-25,28-30s01 |
| largemem | | | | 11 | allocated | k07n05,09,13,17,22,26,30,k05n05,24,k07n26-27 |
| supporters | no | up | 3:00:00:00 | 703 | idle | d07n04-11,13,20,24-31,34-40s01-02,d09n04-11,13,20,24-31,33-40s01-02,d13s03 |
| test | no | up | 3:00:00:00 | 2 | idle | d13s01-02 |
| viz | no | up | 1:00:00:00 | 2 | mixed | f11s02,k05n28 |

sbatch

`sbatch` is the primary way that most users submit jobs to Slurm. The syntax for usage is:

```
sbatch [options] script [args...]
```

where `script` is a conventional shell script.

- If `script` is omitted, then `sbatch` will attempt to read from standard input (very few users go this route)
- Various `options` can be supplied either on the command line, or in the `script` itself as `#SBATCH` directives
- See `man sbatch` for more details on the available `option`, some of which we will cover in the examples

salloc

`salloc` immediately requests an allocation of resources, generally interactively. Here is an example of a simple allocation request on the debug partition:

```
1 [rush:/d_mpi-samples]$ salloc --nodes=1 --ntasks-per-node=8 --partition=debug --time=00:10:00
2 salloc: Granted job allocation 2524554
3 [rush:/d_mpi-samples]$ srun hostname
4 d07n33s02
5 d07n33s02
6 d07n33s02
7 d07n33s02
8 d07n33s02
9 d07n33s02
10 d07n33s02
11 d07n33s02
12 [k07n14:/d_mpi-samples]$ export | grep SLURM
13 declare -x SLURM_JOBID="2524554"
14 declare -x SLURM_JOB_CPUS_PER_NODE="8"
15 declare -x SLURM_JOB_ID="2524554"
16 declare -x SLURM_JOB_NODELIST="d07n33s02"
17 declare -x SLURM_JOB_NUM_NODES="1"
18 declare -x SLURM_MEM_PER_CPU="3000"
19 declare -x SLURM_NNODES="1"
20 declare -x SLURM_NODELIST="d07n33s02"
21 declare -x SLURM_NODE_ALIASES="(null)"
22 declare -x SLURM_NPROCS="8"
23 declare -x SLURM_NTASKS="8"
24 declare -x SLURM_NTASKS_PER_NODE="8"
25 declare -x SLURM_SUBMIT_DIR="/ifs/user/jonesm/d_mpi-samples"
26 declare -x SLURM_SUBMIT_HOST="k07n14"
27 declare -x SLURM_TASKS_PER_NODE="8"
28 [rush:/d_mpi-samples]$ exit
29 exit
30 [rush:/d_mpi-samples]$ export | grep SLURM
31 [rush:/d_mpi-samples]$
```

Things to note when using `salloc`:

- `salloc` requests will be scheduled like any other Slurm job, which can take some amount of waiting time depending on the amount of resources requested,
- You get a new shell when the allocation succeeds, and you should exit that shell when finished, otherwise your shell environment will be somewhat confused,
- Most Slurm users find it easier to use `sbatch` with a dedicated script, rather than tying up a terminal with an `salloc`
- See `man salloc` for more details on the available options, which are essentially the same as for `sbatch`.

`salloc` can also be used to launch commands directly, either an application:

```
1 [rush:~/d_mpi-samples]$ salloc --nodes=1 --ntasks-per-node=8 --partition=debug \
2                               --time=00:10:00 srun hostname
3
4 salloc: Granted job allocation 2524894
5 d07n33s01
6 d07n33s01
7 d07n33s01
8 d07n33s01
9 d07n33s01
10 d07n33s01
11 d07n33s01
12 salloc: Relinquishing job allocation 2524894
13 salloc: Job allocation 2524894 has been revoked.
```

or say a handy extra terminal:

```
1 [rush:~/d_mpi-samples]$ salloc --nodes=1 --ntasks-per-node=8 --partition=debug \
2                               --time=00:10:00 xterm
3 salloc: Granted job allocation 2524915
```

fisbatch

`fisbatch` is a Slurm community script adapted to CCR - it is designed to automate the process of `salloc` followed by an `ssh` to the allocated node (with X forwarding enabled). To achieve a similar result to what we just did in the previous example:

```
1 [rush:~/d_mpi-samples]$ fisbatch --partition=debug --nodes=1 --ntasks-per-node=8 \
2                               --time=00:10:00
3 FISBATCH -- waiting for JOBID 2734481 to start on cluster=ub-hpc and partition=debug
4 !
5 FISBATCH -- Connecting to head node (d07n33s02)
6 [d07n33s02:/ifs/user/jonesm/d_mpi-samples]$
7 srun hostname
8 d07n33s02
9 d07n33s02
10 d07n33s02
11 d07n33s02
12 d07n33s02
13 d07n33s02
14 d07n33s02
15 d07n33s02
16 [d07n33s02:/ifs/user/jonesm/d_mpi-samples]$ exit
17 exit
18 [screen is terminating]
19 Connection to d07n33s02 closed.
20 FISBATCH -- exiting job
21 scancel: error: Kill job error on job id 2734481: Invalid job id specified
```

```
[rush:~]$ export | grep SLURM
declare -x SLURM_JOBID="2524956"
declare -x SLURM_JOB_CPU_PER_NODE="1"
declare -x SLURM_JOB_ID="2524956"
declare -x SLURM_JOB_NODELIST="d07n33s01"
declare -x SLURM_JOB_NUM_NODES="1"
declare -x SLURM_NODE_PER_CPU="3000"
declare -x SLURM_NODES="1"
declare -x SLURM_NODELIST="d07n33s01"
declare -x SLURM_NODE_ALIASES="(null)"
declare -x SLURM_SUBMIT_DIR="/ifs/user/jonesm"
declare -x SLURM_SUBMIT_HOST="k07n14"
declare -x SLURM_TASKS_PER_NODE="1"
[rush:~]$ srun --ntasks=1 --ntasks-per-node=1 hostname
d07n33s01
[rush:~]$
```

- The `xterm` still gets run on the cluster front-end, just like the previous example,
- You can again utilize Slurm's task launcher, `srun` to use the allocated resources,
- Once you exit the `xterm` the `salloc` will release the allocation.

srun

`srun` is Slurm's parallel task launcher. Note that `srun` can also create the allocation as well, so if you only need to run a single command interactively you can simply execute a single `srun`:

```
1 [rush:~/d_mpi-samples]$ srun --nodes=1 --ntasks-per-node=8 --partition=debug \
2                               --time=00:10:00 hostname
3 d07n33s01
4 d07n33s01
5 d07n33s01
6 d07n33s01
7 d07n33s01
8 d07n33s01
9 d07n33s01
10 d07n33s01
11 [rush:~/d_mpi-samples]$ module load intel-mpi
12 [rush:~/d_mpi-samples]$ mpicc -c cpi.impi.cpi.c
13 [rush:~/d_mpi-samples]$ export I_MPI_FMPI_LIBRARY=/usr/lib64/libmpi.so
14 [rush:~/d_mpi-samples]$ srun --nodes=1 --ntasks-per-node=8 --partition=debug \
15                               --time=00:10:00 cpi.impi
16 Process 0 on d07n33s01
17 Process 1 on d07n33s01
18 Process 2 on d07n33s01
19 Process 3 on d07n33s01
20 Process 4 on d07n33s01
21 Process 5 on d07n33s01
22 Process 6 on d07n33s01
23 Process 7 on d07n33s01
24 pi is approximately 3.1416009869231249, Error is 0.0000083333333318
25 wall clock time = 0.000337
```

or you can use `srun` in conjunction with `sbatch` and `salloc`.

Important `srun` notes:

- You can use essentially the same Slurm options as with `salloc` and `sbatch`,
- `srun` is a **parallel** task launcher - note from the preceding example that it executed the same binary on all of the allocated cores,
- `srun` supports various MPI implementations, including Intel MPI as in the example above - in this case the extra library (`IMPI_PMI_LIBRARY`) overrides an existing one in Intel's MPI library to ensure proper task placement, see the `--mpi` option to `srun`,
- Note that options to `srun` can override existing options in a `salloc` or `sbatch` environment (.e.g., you can request all of the cores per node in an `sbatch` script, but then under-subscribe them).

The module Command

module command syntax:

```
-bash-2.05b$ module help

Modules Release 3.1.6 (Copyright GNU GPL v2 1991):
Available Commands and Usage:
+ add|load          modulefile [modulefile ...]
+ rm|unload         modulefile [modulefile ...]
+ switch|swap       modulefile1 modulefile2
+ display|show      modulefile [modulefile ...]
+ avail             [modulefile [modulefile ...]]
+ use [-a|--append] dir [dir ...]
+ unuse             dir [dir ...]
+ update
+ purge
+ list
+ clear
+ help              [modulefile [modulefile ...]]
+ whatis            [modulefile [modulefile ...]]
+ apropos|keyword   string
+ initadd           modulefile [modulefile ...]
+ initprepend       modulefile [modulefile ...]
+ initrm            modulefile [modulefile ...]
+ initswitch        modulefile1 modulefile2
+ initlist
+ initclear
```

Modules Software Management System

There are a large number of available software packages on the CCR systems, particularly the Linux clusters. To help maintain this often confusing environment, the **modules** package is used to add and remove these packages from your default environment (many of the packages conflict in terms of their names, libraries, etc., so the default is a minimally populated environment).

Using module in Batch

If you change shells in your batch script you may need to explicitly load the modules environment:

tcsh :

```
source $MODULESHOME/init/tcsh
```

bash :

```
. ${MODULESHOME}/init/bash
```

but generally you should not need to worry about this step (do a “module list” and if it works ok your environment should already be properly initialized).

Modules Highlights

| | |
|--------------------------------------|--|
| <code>module avail [string]</code> | list available modules (opt. <i>string</i>) |
| <code>module help name</code> | help text for <i>modulename</i> |
| <code>module list</code> | list modules in your environment |
| <code>module load name</code> | load selected <i>modulename</i> |
| <code>module unload name</code> | unload selected <i>modulename</i> |
| <code>module swap name1 name2</code> | swaps one module for another |

Build Your Own Modules

You can build your own modules, by using the `use.own` module:

```

1 [rush:~]$ module show use.own
2 -----
3 /util/Modules/modulefiles/use.own:
4
5 module-what is adds your own modulefiles directory to MODULEPATH
6 module use --append /user/jonesm/privatemodules
7 -----

```

- Create your `$HOME/privatemodules` directory
- Create your own module files (tc1-based) in this private repository, you can use the existing ones in `$MODULEPATH` as templates

Module Usage Examples

```

1 [rush:~]$ module list
2 Currently Loaded Modulefiles:
3 1) null      2) modules
4 [rush:~]$ module avail intel
5
6 ----- /util/Modules/modulefiles -----
7 intel/10.0      intel/13.0      intel-ipp/7.0.1      intel-mpi/4.1.0
8 intel/10.1      intel/13.1(default) intel-mpi/3.1        intel-mpi/4.1.1
9 intel/11.0      intel/14.0      intel-mpi/3.2        intel-mpi/4.1.3
10 intel/11.1      intel/9.0       intel-mpi/4.0        intel-tbb/3.0.3
11 intel/12.0      intel/9.1       intel-mpi/4.0.1      intel-tbb/4.1
12 intel/12.1      intel/9.1.040   intel-mpi/4.0.3
13 [rush:~]$ module avail intel/
14
15 ----- /util/Modules/modulefiles -----
16 intel/10.0      intel/11.1      intel/13.0      intel/9.0
17 intel/10.1      intel/12.0      intel/13.1      intel/9.1
18 intel/11.0      intel/12.1      intel/14.0      intel/9.1.040
19 [rush:~]$ module load intel/12.1
20 [rush:~]$ which icc
21 /util/intel/composer_xe_2011_sp1.11.339/bin/intel64/icc
22 [rush:~]$ module swap intel/12.1 intel/14.0
23 [rush:~]$ which icc
24 /util/intel/composer_xe_2013_sp1.2.144/bin/intel64/icc

```

More Module Information

- <http://modules.sourceforge.net>,
Home site for Modules source code
- <https://www.nersc.gov/users/software/nersc-user-environment/modules>,
NERSC Modules help page
- `man module`,
Module command man page
- `man modulefile`,
modulefile syntax man page

Simple OpenMP Example

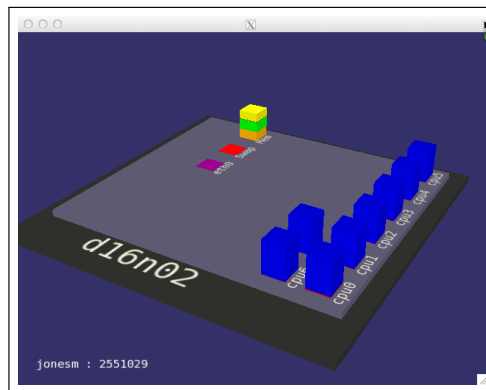
A simple example with an OpenMP based code:

```

1  #!/bin/bash
2  #SBATCH --nodes=1
3  #SBATCH --ntasks-per-node=1
4  #SBATCH --cpus-per-task=8
5  #SBATCH --exclusive
6  #SBATCH --constraint=CPU-L5520
7  #SBATCH --partition=debug
8  #SBATCH --time=00:10:00
9  #SBATCH --mail-type=END
10 #SBATCH --mail-user=jonesm@buffalo.edu
11 #SBATCH --output=slurmOMP.out
12 #SBATCH --job-name=omp
13 #
14 module load intel
15 module list
16 export | grep SLURM
17 echo "Running on $SLURM_NNODES node with $SLURM_CPUS_ON_NODE cores."
18 export OMP_NUM_THREADS=$SLURM_CPUS_ON_NODE # single node only for OpenMP
19 ./lap3-omp<<EOF
20 1024
21 1.e-4
22 EOF

```

- `slurmjobvis` view of preceding example, showing CPU/network activity and memory used (when you mouse over the columns it should display values)



Notes on the simple OpenMP Example:

- In this case the `--exclusive` flag is redundant, we are requesting all 8 cores in the L5520 nodes anyway,
- You can also request all 8 cores using `--ntasks-per-node=8`, and omitting `--cpus-per-task` since we run the application directly rather than through `srun` (which would otherwise interpret `--ntasks-per-node=8` to launch 8 versions of the application, each with 8 OpenMP threads)

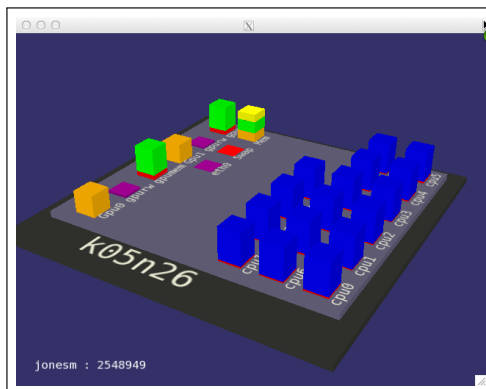
MPI Parallel Job With GPUs

```

1  #!/bin/bash
2  #SBATCH --partition=debug
3  #SBATCH --time=01:00:00
4  #SBATCH --nodes=1
5  #SBATCH --ntasks-per-node=16
6  #SBATCH --gres=gpu:2
7  #SBATCH --output=slurmNAMD.out
8  #SBATCH --mail-user=jonesm@buffalo.edu
9  #SBATCH --mail-type=END
10 #SBATCH --job-name=stmv
11 #
12 module load namd/2.9-MPI-CUDA
13 # CUDA version requires MPI, see /util/slurm-scripts/slurmNAMD for
14 # other options
15 echo "NAMDDHOME = $NAMDDHOME"
16 export I_MPI_DEBUG=4
17 export | grep I_MPI
18 # construct namd nodes file from $PBS_NODEFILE
19 #
20 export MY_NODEFILE=tmp.$SLURM_JOBID
21 srun -l hostname -s | sort -n | awk '{print $2}' > $MY_NODEFILE
22 NPROCS=$(cat $MY_NODEFILE | wc -l)
23 NNODES=$(cat $MY_NODEFILE | uniq | wc -l)
24 echo "Running $NPROCS processes."
25 export NAMDDNODEFILE=namd.$SLURM_JOBID
26 echo "group main" > $NAMDDNODEFILE
27 NODES=$(cat $MY_NODEFILE)
28 for node in $NODES; do
29     echo "host $node >> $NAMDDNODEFILE"
30 done
31 #
32 # run STMV benchmark
33 #
34 export I_MPI_PMI_LIBRARY=/usr/lib64/libpmi.so
35 srun namd2 stmv.namd
36 [ -e $MY_NODEFILE ] && \rm $MY_NODEFILE
37 [ -e $NAMDDNODEFILE ] && \rm $NAMDDNODEFILE

```


- `slurmjobvis` view of preceding example, showing CPU and GPU activity
- Note that the `--partition=debug` currently only has one such GPU node, for more nodes and longer running times, you can instead submit to `--partition=gpu` (32 12-core nodes)



Things to note from the ensemble example:

- You need the `wait` command to wait on the jobs sent to the background,
- This example can also be done with an ensemble of individual Slurm jobs, or with a **job array**,
- The individual runs need to take the same amount of time, otherwise some cores will be idled waiting on completion of the slowest run, leading to a poor load balance.

Advanced Example: Multiple Serial Runs

This example runs an ensemble of a serial application, allowing Slurm to map the individual tasks to whatever resources are available:

```
1 #!/bin/bash
2 #SBATCH --tasks=64
3 #SBATCH --time=00:10:00
4 #SBATCH --mail-type=ALL
5 #SBATCH --mail-user=jonesm@buffalo.edu
6 #SBATCH --output=slurmDIST.out
7 #SBATCH --job-name=DIST
8 #
9 module list
10 export | grep SLURM
11 NPROCS=`srun -l hostname -s | wc -l`
12 echo "Running $NPROCS processes."
13 for ((i=0;i<SLURM_NTASKS;i++)); do
14     echo "Task $i"
15     srun --cpus-per-task=1 --exclusive --nodes=1 --ntasks=1 ./rp p11_4.$i > out.$i &
16 done
17 wait
```

Advanced Example: Multiple Serial Runs With a Job Array

A job array version of running the same ensemble of a serial application:

```
1 #!/bin/bash
2 #SBATCH --nodes=1
3 #SBATCH --ntasks-per-node=1
4 #SBATCH --array=0-31
5 #SBATCH --time=00:10:00
6 #SBATCH --mail-type=ALL
7 #SBATCH --mail-user=jonesm@buffalo.edu
8 #SBATCH --output=slurmARRAY-%j.out
9 #SBATCH --job-name=ARRAY
10 #
11 module list
12 export | grep SLURM
13 NPROCS=`srun -l hostname -s | wc -l`
14 echo "Running $NPROCS processes."
15 srun --nodes=1 --ntasks=1 ./rp p11_4.$SLURM_ARRAY_TASK_ID > out.$SLURM_ARRAY_TASK_ID
```


Notes on Slurm array job example:

- Note that you get multiple jobs (the array) with a single `sbatch` submission, each has its own identifier and output file,
- Load balancing is not an issue in this case, as each job in the array will get its own dedicated resources,
- So when do you use the single job with an embedded ensemble? That depends partially on your preferences, the workload, and queueing policy. Sometimes parallel jobs get higher priority, so bundling the jobs together can actually get higher overall throughput.

Notes on the `sbcast` example:

- The most common use for this is to run a data-intensive application from the local disk on each node (automatically set to `$SLURMTMPDIR`),
- Note the lack of an inverse operation to `sbcast`, hence the loop over the allocated nodes and the `ssh` to manually copy output files back to the submission directory (of course you could also do the same to replace the `sbcast`),
- Are you better off using the local disk or the shared file system? For a data-intensive application, the local disk i/o is more scalable.

Example Using `sbcast` File Staging

You can use Slurm's `sbcast` command to broadcast a file out to remote resources.

```

1  #!/bin/bash
2  #SBATCH --nodes=2
3  #SBATCH --ntasks-per-node=8
4  #SBATCH --constraint=CPU-L5520|CPU-L5630
5  #SBATCH --partition=debug
6  #SBATCH --time=00:10:00
7
8  #SBATCH --mail-type=END
9  #SBATCH --mail-user=jonesm@buffalo.edu
10 #SBATCH --output=slurmQ.out
11 #SBATCH --job-name=titan-test
12 #
13 module load titan
14 which titan
15 export I_MPI_DEBUG=4
16 export | grep I_MPI
17 export | grep SLURM
18 NNODES=$(srun hostname -s | uniq | wc -l)
19 echo "Number of nodes: $NNODES"
20 NPROCS=$(srun -l hostname -s | wc -l)
21 echo "Number of cores: $NPROCS"
22 NODELIST=$(srun --ntasks-per-node=1 --ntasks=$NNODES hostname -s)
23 echo "Allocated hosts: " $NODELIST
24 # copy all input files out to remote nodes' local scratch
25 for f in *.inp *.data ; do
26     sbcast $f $SLURMTMPDIR/$f
27 done
28 cd $SLURMTMPDIR
29 export I_MPI_PMI_LIBRARY=/usr/lib64/libpmi.so
30 srun titan
31 pwd
32 ls -lrt
33 # unfortunately slurm lacks an inverse to bcst for the output
34 for node in $NODELIST; do
35     ssh $node "cd $SLURMTMPDIR ; cp perf* pile* $SLURM_SUBMIT_DIR/"
36 done

```

Summary

Slurm is a very flexible resource manager, designed to scale to very large installations:

- Flexibility comes with a certain of complexity, there are multiple ways of accomplishing various tasks (e.g., running an interactive job, parallel task launching, etc.)
- Once you have some familiarity with Slurm, it becomes easier to automate your workflows

Options Not Exemplified

Things we have left out:

- There are many additional options to resource allocation requests, for different constraints, accounts, etc.,
 - Memory is one the most important of these resources - our examples utilized little memory, so the default allocated amount (currently 3GB/core) was enough. See the `--mem` option for setting the memory request per node if you need more,
- There are many ways to customize the outputs from the various Slurm client commands (see the various `man` pages),
- Priority and scheduling - a complicated study in its own right, and subject to many site policies and dependencies (see, for example, the Slurm documentation and `man sprio`).

CCR Slurm Specific Details/Notes

- `$SLURMTMPDIR` is a local scratch directory created on each compute node per Slurm job (generally `/scratch/$SLURM_JOBID`)
- Use the **debug** partition for development jobs (less than one hour of walltime, single or two nodes)
- Use the **gpu** partition for production jobs using GPUs
- Use the **largemem** partition when requesting the large memory 32-core nodes (a priority bump shortens the waiting time)