



<https://hps.vi4io.org>

Ruben Kellner, Markus Boden, Sebastian Krey,
Azat Khuziyakhmetov

High-Performance System Administration

Introduction to Slurm



Table of contents

1 Learning Objectives

Learning Objectives

After the course the students should be able to:

- Comprehend there are login and compute nodes with different functionality
- Utilize tools (e.g. FileZilla) to transfer data between the local and remote system
- Understand the basics of the HPC infrastructure
- Use a workload manager like SLURM to allocate HPC resources (e.g. CPUs) and to submit a batch job.
- Run parallel programs in an HPC environment.
- Compile programs on the HPC Cluster
- Run Interactive Sessions with a graphical user interface

Partitions

- Different compute nodes have different features
- Slurm differentiates using **Partitions**

Interactive X11 Job

Running Matlab

```
> ssh -Y login-mdc.gwdg.de
> module load matlab
> srun --x11 -p medium matlab
```

- The job will be dispatched and as soon as an available node is found and the Matlab interface will start.
- If you have your own license for Matlab then you need to place your `license.lic` file in `$HOME/.matlab/R2015a/licenses` directory (depending on the version you are using).

Interactive Console Job

Running R interactively

```
> ssh login-fas.hpc.gwdg.de
> module load r
> srun --pty -p medium R
```

Try it!

Assignment: Run a Program using Slurm

run R interactively and run the demo session

- log onto the cluster
- Load the R module
- Run R interactively on the medium partition
- call the list of demos and run one of them (hint you can quit the list with q)

run a Matlab session

- Start up X2go Client
- Load the Matlab module
- Run Matlab interactively on the medium partition with X11 forwarding

Basic Concepts 2

Serial job Job consisting of one task using one job slot.

SMP job Job with shared memory parallelization (often realized with OpenMP), meaning that all processes need access to the memory of the same node. Therefore uses several job slots **on the same node**.

MPI job Job with distributed memory parallelization, realized with MPI. Can use several job slots on several nodes and needs to be started with a helper program, e.g., `mpirun` or `srun`.

Resource selection: CPU

s run options for parallel (SMP or MPI) jobs.

- N <min>-<max>, Minimum and maximum node count. You can also
- nodes=<min>-<max> specify the exact number.
- n,--ntasks=<n> Number of tasks (not equally distributed!)
- tasks-per-node=<n> Tasks per node. If used with -n it denotes the maximum number of tasks per node.
- c,--cpus-per-task=<n> CPUs per tasks.

A note on -n vs. -c

Rule of thumb

- -c for single node jobs
- -n for MPI jobs

A note on -n vs. -c

Rule of thumb

- -c for single node jobs
- -n for MPI jobs

Rule of thumb 2

If you are unsure if your program uses MPI, then it does not.

How to avoid waiting

Reservation: hpcsa-course

Usage

Either: `--reservation=hpcsa-course` for each job

Or: `export SLURM_RESERVATION=hpcsa-course`

The latter has to be unset, if you want to submit to a partition besides medium.

Try it!

Exercises

Try these job configurations

- 1 10 tasks
- 2 10 tasks distributed over 3 nodes
- 3 3 nodes with 3 tasks each
- 4 1 task with 5 cores
- 5 2 tasks per node on 2 nodes with 4 cores per task

use `slurm_resources` script to get see the resources of your job

Resource Selection: Memory

srun options

- -mem <size[K|M|G|T] > Memory per node.

- -mem-per-cpu<size[K|M|G|T] > Memory per core.

■ without options:

- ▶ each partition has a DefMemPerCPU option
- ▶ can be retrieved via `scontrol show partition <name>`

Try it!

Exercise:

Play with the combination of number of cores or tasks, nodes and their effect on your available memory:

- 1 1 core and --mem 4G
- 2 3 tasks and 2 nodes, see effect of --mem and --mem-per-cpu
- 3 20 tasks, see distribution of memory over hosts.

Non interactive Jobs

Problem

- if you have big jobs, your queue time will be long
- srun needs you to stay logged in
- jobs can run for days

Non interactive Jobs

Problem

- if you have big jobs, your queue time will be long
- srun needs you to stay logged in
- jobs can run for days

Solution

Batch Jobs!

SBATCH: Using Job Scripts

A job script is a shell script with a special comment section.
The #SBATCH lines have to come first!

sbatch: Basic job script example

```
#!/bin/bash
#SBATCH -p medium
#SBATCH -t 10:00
#SBATCH -o job-%J.out
```

slurm_resources

Submit with:

Jobscripts

- A job script is essentially a normal script
- usually bash/shell, but can be any scripting language (R, python, perl)
- #SBATCH lines need to be at the top!
- you can copy files, load modules, and do any scripting you want
- for MPI, use `srun` or `mpirun` to start your program

More options

Solution

```
sbatch <slurm options> jobscript
```

- `--mail-type=<TYPE>` get mail notifications (type: BEGIN, END, etc.)
- `--mail-user=<address>` Default: `${USER}@gwdg.de`
- `-o/-e <file>` Store job output in file (slurm-<jobid>.out by default). `%J` in the filename stands for the jobid.

Slurm Commands

sinfo Info about the system and partitions.

`-p <partition>, -t <state>`

squeue Show the job queue.

`-p <partition>, --me`

scancel Cancel Job

`scancel <JobID>`

`scancel -p <partition>|-u $USER`

Exercises

Exercise

Write your own Job script.

- Use `echo`, `hostname`, and `sleep X` (sleep for X seconds) to generate output or have it running for a longer time.
- Have the job send you an email. Advanced: Take a look at the different mail-type options. What do they do?
- Write the output to a different file. Redirect output and error into different files. Advanced: Take a look at the filename pattern options. Include node and job name in the output file.

Task Distribution

Distributing tasks in the medium partition

```
#SBATCH -p medium
#SBATCH -n 240
#SBATCH -o job-%J.out

module purge
module load intel/compiler intel/mkl intel/mpi namd

srun namd2 +setcpuaffinity apoal.namd
```

This will spread tasks among many nodes.

Task Distribution fixed

Distributing tasks in the medium partition

```
#SBATCH -p medium
#SBATCH -N 10
#SBATCH --ntasks-per-node 24
#SBATCH -o job-%J.out

module purge
module load intel/compiler intel/mkl intel/mpi namd

srun namd2 +setcpuaffinity apoal.namd
```

Memory is faster then network!

Try to spread your tasks to as little nodes as possible.

Job Disk Space Usage Options

- `/local` Local hard disk of the node. SSD based and therefore a very fast option for storing temporary data. Automatic file deletion. A temporary directory is created on all nodes at `$TMP_LOCAL`.
- `/scratch` Shared scratch space, available on most nodes, but there are two instances (use `-C scratch` or `-C scratch2`). Very fast, no automatic file deletion, but also no backup! Files may have to be deleted manually when we run out of space.
- `/scratch/ssd` special extra fast scratch file system only on `scratch1`. Ideal for temporary data in jobs spanning multiple nodes. Automatic file deletion. A per-job directory is created at `$TMP_SCRATCH`.
- `$HOME` Available everywhere, permanent, with backup. Personal disk space can be increased. Comparably slow.

Recipe: Using /scratch

```
#!/bin/bash
#SBATCH -p medium
#SBATCH -n 24
#SBATCH -N 1
#SBATCH -C scratch
#SBATCH -t 1-00:00:00

export g09root="/usr/product/gaussian/g09/d01"
source $g09root/g09/bsd/g09.profile

if [ ${TMP_SCRATCH} -a -d ${TMP_SCRATCH} ]; then
    export GAUSS_SCRDIR=${TMP_SCRATCH}
else
    export GAUSS_SCRDIR=$TMP_LOCAL
fi

g09 myjob.com myjob.log
```

Try it!

Exercise

Write a job script, where you

- create a scratch directory
- copy data from your home file system to the scratch directory
- run a job with the data
- copy the results back
- delete the scratch directory

If you do not have a program/data to try this on, there is a small python program in `/scratch1/projects/scc-course/` and a bit of input data.

The fat+ partition

The fat+ partition contains:

- 5 nodes with 1.5Tb Memory
- 1 node with 2Tb Memory

Usage recommendations:

- Work your way up. Start in fat and only use fat+ if your jobs runs out of memory.
- Use `reportseff` or `profit-hpc`, see if your job really is memory bound
- When unsure, ask us!
- `--mem` or `--mem-per-cpu` is mandatory
- You might get angry mails from me, if you waste resources here

Recipe: Combine shared memory and MPI

Running hybrid jobs

```
#SBATCH -p medium
#SBATCH -N 5
#SBATCH --ntasks-per-node=4
#SBATCH --cpus-per-task=6
#SBATCH -o job-%J.out

module purge
module load openmpi/gcc

export OMP_NUM_THREADS=$SLURM_CPUS_PER_TASK

srun hybrid_job
```


More Slurm Commands

`scontrol show [partition|node|job] <x>` where x should be a node name, JobID or partition name.

`sprio` Priority information about pending jobs

`sacct` Get information about a job after it finished

`-j <jobid>`

`--format=JobID,User,JobName,MaxRSS,Elapsed,TimeLimit`

Debugging

- take a look at your output files, while the job is running:
 - ▶ `tail -f /path/to/output`
- take a look at the jobs, while it is running
 - ▶ you can ssh into every node that currently calculates your job
 - ▶ use `htop` to see the processor and ram usage

Digression: Directory Structure 1

- Convention: Executables are stored in “bin”, shared libraries in “lib” directories
- Directories in “\$PATH” are searched for binaries, directories in “\$LD_LIBRARY_PATH” for libraries
- Two strategies:
 - 1 Put everything directly under \$HOME/bin, \$HOME/lib
 - Easy to setup search paths
 - Difficult to remove software packages
 - 2 Put each software into its own subdirectory
 - Easy to remove software (with “rm -rf <subdirectory>”)
 - Difficult to setup search paths

Digression: Directory Structure 2

- Or combine both strategies:
 - ▶ Put each software in its own subdirectory
 - ▶ Use “`ln -s`” to link everything to `$HOME/bin` and `$HOME/lib`, respectively
 - ▶ Use “`export LD_LIBRARY_PATH=$HOME/lib:$LD_LIBRARY_PATH; export PATH=$HOME/bin:$PATH`” in your shell and scripts
 - ▶ Use “`find $HOME/bin $HOME/lib -xtype l -delete`” after removing software