



<https://hps.vi4io.org>

Patrick Höhn, 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
- 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

































# Interactive Jobs

## s run: Interactive jobs

- `--x11` adds X11 (GUI) forwarding. This requires that you connect to the front-end with `ssh -Y` and your local machine supports X-Windows.
- `-p int` use the interactive partition. In `int` the nodes have no slot limit. They will take jobs until their load crosses a specified threshold, so jobs start immediately.
- `--pty` interactive mode

# Interactive X11 Job

## Running Matlab

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

- The job will be dispatched and as soon as an available node is found and the Matlab interface will start.

# Interactive Console Job

## Running python interactively

```
> ssh login-mdc.hpc.gwdg.de
> module load anaconda3
> srun --pty -p int python3
>>> import socket, os
>>> print(socket.gethostname())
>>> print(os.system("slurm_resources"))
```

## 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.





# 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

```
1 #!/bin/bash
2 #SBATCH -p medium
3 #SBATCH -t 10:00
4 #SBATCH -o job-%J.out
5
6 slurm_resources
```

Submit with:

```
sbatch <script name>
```

# 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!
- only lines starting exactly with #SBATCH are parsed
- you can copy files, load modules, and do any scripting you want
- for MPI, use `srun` or `mpirun` to start your program

## Useful options for batch job submission

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

**squeue** Show the job queue.  
`squeue -p <partition>, --me`

**scancel** Cancel Job  
`scancel <JobID>`  
`scancel -p <partition>|-u $USER`

## Batch Script for Task Distribution in medium partition

```
1 #SBATCH -p medium
2 #SBATCH -N 10
3 #SBATCH --ntasks-per-node 24
4 #SBATCH -o job-%J.out
5
6 module purge
7 module load intel-oneapi-compilers intel-oneapi-mkl intel-oneapi-mpi namd
8
9 srun namd2 +setcpuaaffinity 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). Very fast, no automatic file deletion, but also no backup! Files may have to be deleted manually when we run out of space.

**\$HOME** Available everywhere, permanent, with backup. Personal disk space can be increased. Comparably slow.



# Recipe: Combine shared memory and MPI

## Running hybrid jobs

```
1 #SBATCH -p medium
2 #SBATCH -N 5
3 #SBATCH --ntasks-per-node=4
4 #SBATCH --cpus-per-task=6
5 #SBATCH -o job-%j.out
6
7 module purge
8 module load openmpi/gcc
9
10 export OMP_NUM_THREADS=$SLURM_CPUS_PER_TASK
11
12 srun hybrid_job
```

## Longer or shorter jobs

### The --qos parameter

- Default maximum runtime: 2 days
- --qos= <qos> can select a QoS
- Two extra QoS available:
  - short for shorter jobs (max. 2h), has higher priority, limited job slots
  - long longer jobs (max. 5d), limited job slots.

### But my job is even longer

- try parallelizing more
- break it down into smaller steps
- check, if your software supports checkpoints
- check again!
- contact us

## 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`

# Using the gpu partition

## GPU parameters

`-G | --gpus=[type:]<n>` requests n GPUs of type

`--gpus-per-task=[type:]<n>` requests n GPUs of type per task

`--gpus-per-node=[type:]<n>` requests n GPUs of type per node

- CPUs are evenly distributed for every GPU
- Available types are:
  - ▶ rtx5000
  - ▶ v100
  - ▶ gtx1080
- See: `sinfo -p gpu --format=%N,%G`

# 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















