



Leibniz Supercomputing Centre
of the Bavarian Academy of Sciences and Humanities

Hands-on: Introduction to Multiuser Cluster Systems at LRZ

Theory & Practice

Florent Dufour · florent@lrz.de

Big Data & Artificial Intelligence
Leibniz Supercomputing Centre
October 10, 2022

Module description

It has been about 350 years that separate the original 1673 Leibniz mechanical calculator from today's Leibniz Supercomputing Centre's (LRZ) facilities located on the Campus Garching. And yet, the spirit has not changed. To quote the German mathematician:

"It is beneath the dignity of excellent men to waste their time in calculation when any peasant could do the work just as accurately with the aid of a machine."

This course module will allow participants to live up to their dignity by providing a comprehensive walkthrough and usage guide to such contemporary types of these machines that potentially fill whole buildings.

In a general overview, historical and current developments and trends in the space of scientific computing and cluster systems will be presented. This will, amongst others, address the following questions: How do modern cluster systems work and how are they architected? How did we come to High Performance Computing, High Performance Data Analytics and High Performance AI? What makes a system adequate to a specific workload? How are these systems operated and how are they made available to their users? In addition, typical interaction methods and usage patterns will be covered, including various possibilities of setting up user environments (e.g. environment variables and modules, user space package managers, containers) as well as tools for resource allocation (i.e. Slurm Workload Manager) and efficient parallelization (MPI, OpenMP, ...). Finally, an overview of different compute clusters as well as their background storage systems operated by LRZ will be provided. The requirements for acquiring access to these systems will be covered as well.

Participants will gain a good understanding of the characteristics of multiuser cluster systems in general and will practise basic methods of typical interaction. They will familiarize themselves with the landscape of cluster systems available at LRZ and this will allow them to choose the right system for their own compute projects.

This is document revision: 2022.2a

Latest revision available at: <https://doku.lrz.de/x/eQBvB>

Table of Contents

1	Get into the workbench	1
2	Interactive Jobs	3
3	Batch Jobs	5

List of Code snippets

1	We SSH into the workbench	1
2	We start and enter the SLURM cluster	1
3	We get familiar with our SLURM cluster	2
4	We run our first interactive job	3
5	We explore how jobs are scheduled on a SLURM cluster	3
6	A simple batch script	5
7	We run a simple batch script	5

Glossary

SLURM Simple Linux Utility for Resource Management.

VM Virtual Machine.

1 Get into the workbench

We have a workbench ready for you. It is a Virtual Machine (VM) in which we will run your own SLURM cluster. The first step consists in getting into the VM.

```

1 # ----- #
2 # On your local machine #
3 # ----- #
4
5 whoami
6 # florent
7
8 uname -a
9 # Darwin BADWLRZ-AB12345 20.6.0 Darwin Kernel Version 20.6.0: Mon Aug 30 06:12:21
   # PDT 2021; root:xnu-7195.141.6~3/RELEASE_X86_64 x86_64
10
11
12 ssh <user>@<IP>
13 # Replace <user> and <IP> with yours.
14
15 # The authenticity of host <IP> can't be established.
16 # ECDSA key fingerprint is SHA256:BZgws5BARCfwE6rDuN5i/aLgZMAKuC4si2D+ZuLN5gw.
17 # Are you sure you want to continue connecting (yes/no)? yes
18 # Warning: Permanently added <IP> (ECDSA) to the list of known hosts.
19 # <user>@<IP> password:
20
21 # Type your password when prompted, it's normal if you don't see what you type!
22
23 # ----- #
24 # On the compute cloud workbench #
25 # ----- #
26
27 # Welcome to Ubuntu 20.04.3 LTS (GNU/Linux 5.4.0-81-generic x86_64)
28
29 whoami
30 # <user>
31
32 uname -a
33 # Linux course-node 4.19.0-14-cloud-amd64 #1 SMP Debian 4.19.171-2 (2021-01-30)
   # x86_64 GNU/Linux
34
35 # Success! We can proceed from here!

```

Code snippet 1: We SSH into the workbench

Once into the VM, you'll be able to start and enter the toy SLURM cluster.

```

1 # We are in the VM
2 # The prompt looks like:
3 # hlrbkurs42@course-node:~$
4
5
6 # Let's start the cluster
7 slurm-start
8 # Starting mysql ... done
9 # Starting slurmdbd ... done
10 # Starting slurmctld ... done
11 # Starting c1 ... done
12 # Starting c2 ... done
13
14 # As you can see, we are starting a database,

```

```
15 # slurm services, as well as 2 compute nodes c1 and c2
16
17 # Let's enter the cluster
18 slurm-enter
19
20 # The prompt now looks like
21 # [root@slurmctld /]#
22
23 # Success, you're on the control node
24 # of the slurm cluster
```

Code snippet 2: We start and enter the SLURM cluster

Let's explore some SLURM commands. You'll be using: `sinfo` and `squeue` quite a lot.

```
1 # Let's gather the information on the cluster
2 # Let's add the --cluster=all option
3
4 sinfo
5 # CLUSTER: linux
6 # PARTITION AVAIL  TIMELIMIT  NODES  STATE NODELIST
7 # normal*        up 5-00:00:00    2   idle c[1-2]
8
9 # We have access to 1 cluster named linux
10 # The cluster contains 1 slurm partition composed of 2 compute nodes
11 # Both nodes are idle at the moment, nothing is running
12
13 # Let's see if there is jobs queued
14 squeue
15 #JOBID PARTITION     NAME     USER ST       TIME  NODES NODELIST(REASON)
16
17 # It's empty!
18
19 # The cluster is yours
```

Code snippet 3: We get familiar with our SLURM cluster

Question

Can you SSH into the Linux cluster or datalab login nodes and compare the slurm partitions and nodes available there?

2 Interactive Jobs

Alright, it's about time we start using our cluster. Let's start interacting interactively with the cluster before submitting batch jobs.

```

1 # Let's run a job interactively
2
3 srun --pty bash # We keep the pty open and ask to run bash
4 hostname # The prompt changed, let's confirm we are on compute node c1
5
6 # From there I can do computation
7 # For example, python is available on the compute node
8
9 python
10 # Python 3.6.8 (default, Sep 13 2022, 09:02:49)
11 # [GCC 8.5.0 20210514 (Red Hat 8.5.0-10)] on linux
12 # Type "help", "copyright", "credits" or "license" for more information.
13 # >>> 2+2
14 # 4 # That's HPC (!)
15
16 exit()
17 exit
18
19 # We are back into the login node

```

Code snippet 4: We run our first interactive job

Cool, let's dive a bit more how jobs are scheduled now

```

1 # Let's keep c1 busy for a while
2 # We use "&" to run the job in the background
3 srun sleep infinity &
4
5 sinfo
6 # PARTITION AVAIL  TIMELIMIT  NODES  STATE NODELIST
7 # normal*      up 5-00:00:00      1  alloc c1
8 # normal*      up 5-00:00:00      1  idle  c2
9
10 # c1 is allocated, c2 is free
11
12 # Let's check the slurm queue
13 squeue
14 # JOBID PARTITION     NAME     USER ST       TIME  NODES NODELIST(REASON)
15 #      8   normal     sleep    root  R         2:11     1  c1
16
17 # What happens if we request more resources=
18 # Let's request a node and run bash interactively on it
19 srun --pty bash
20 hostname
21 # c2
22 # I am c2 because c1 is allocated
23
24 exit
25
26 # Let's now keep c2 busy a 30 secs
27 srun --pty sleep 30 & # make c2 busy 30 secs
28
29 # What happens if we submit a job while both nodes are allocated
30 srun --pty sleep 30 &
31 # srun: job 11 queued and waiting for resources
32 # We are queued

```



```
33
34 squeue
35 # JOBID PARTITION      NAME      USER ST      TIME  NODES NODELIST(REASON)
36 #    11    normal    sleep    root PD      0:00     1 (Resources)
37 #     8    normal    sleep    root R      14:37     1 c1
38 #    10    normal    sleep    root R       2:21     1 c2
39
40 sacancel 8 # Let's cancel job on c1
41 # srun: job 11 has been allocated resources
42 # Since c is freed, the next job is starting
43
44 sacct # Let's check the accounting of what we have been doing
45 # JobID      JobName  Partition  Account  AllocCPUS      State  ExitCode
46 # -----  -
47 # 7          bash     normal     root     1           FAILED  1:0
48 # 8          sleep   normal     root     1  CANCELLED+  0:0
49 # 9          bash     normal     root     1  COMPLETED  0:0
50 # 10         sleep   normal     root     1  RUNNING    0:0
51 # 11         sleep   normal     root     1  RUNNING    0:0
```

Code snippet 5: We explore how jobs are scheduled on a SLURM cluster

3 Batch Jobs

I hope SLURM starts to make sense to you! This might be the right time now to introduce you to batch jobs. This is how you would usually run your jobs. For this, you would write a batch script and submit it with the SLURM scheduler.

```
1 #!/bin/bash
2 #SBATCH --cluster=linux
3 #SBATCH --nodes=2
4 #SBATCH --output=batch.stdout
5 #SBATCH --error=batch.stderr
6
7 echo "Hello, I am a script run on node: `hostname`"
```

Code snippet 6: A simple batch script

And this is how you would run this job:

```
1 sbatch batch.sh
```

Code snippet 7: We run a simple batch script

Question

Can you submit 2 jobs and try submitting a third? Do you observe the same behaviour as when we did it interactively?