

LINUX WORKSHOP



05 RSB Computer Cluster

By Jiajia Li (ANU Biological Data Science Institute)

06/03/2025



**Australian
National
University**



Learning Objectives

- Learn the configuration of RSB computer cluster
- Learn the data storage policy of RSB cluster
- Learn the job scheduling system
- Write and run the variant calling pipeline





RSB Computer Cluster

The RSB computer cluster consists of 4 servers, including 2 CPU servers and 2 GPU servers. The servers work together and are controlled and scheduled by SLURM.

The specs of 4 servers:

1. `dayhoff.rsb.anu.edu.au`

- 1TB of RAM
- 196 Cores
- 100TB of data storage
- Ubuntu 20.04 Linux
- GPU, 2 Nvidia A30's
- DELL PowerEdge R750 + MD1400

2. `fisher.rsb.anu.edu.au`

- 1TB of RAM
- 56 Cores
- 50TB of data storage
- Ubuntu 20.04 Linux

3. `wright.rsb.anu.edu.au`

- 1TB of RAM
- 64 Cores
- 70TB of data storage
- Ubuntu 20.04 Linux

4. `thor.rsb.anu.edu.au`

- 170GB of RAM
- 80 Cores
- 9 NVIDIA GeForce RTX 2080
- Rocky Linux 8.9 (compat with RHEL8)

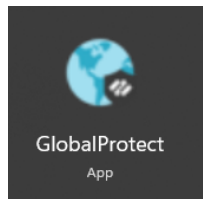




Assessing RSB servers

To access the RSB server, we need to:

1. Connect to GlobalProtect
2. Connect to server using `ssh` command
3. Using your UID as account name and ANU password



- `ssh (u number)@dayhoff.rsb.anu.edu.au`
- for example `ssh u2474733@dayhoff.rsb.anu.edu.au`



[illegible]



Data storage locations on RSB server

- Home directory: /mnt/data/(server)/home/UID, 100GB per user
- Groups directory: /mnt/data/(server)/home/groups, 500GB per group
- Projects directory: /mnt/data/(server)/home/projects, 250GB per project

Scratch Space: /mnt/data/(server)/home/scratch/...

- No limitation
- Data not backed up
- Files will be deleted after 130 days
- You can store temporary files here





Job Scheduling System - SLURM

A **job scheduling system**, also called **Workload Management System** or **Cluster Management System**, is a software designed to efficiently allocate and manage computing resources in a distributed computing environment.

These systems are commonly used by **high-performance computing clusters**, **data centres**, and other **large-scale computing infrastructures**.

Their primary purpose is to optimise the utilisation of available resources while ensuring **fair access** to those resources for multiple users.





Job Scheduling System - SLURM

The RSB cluster uses SLURM, which is an open-source project.

The NCI's supercomputer **Gadi** uses PBS Professional. It has a similar syntax to SLURM, and you can quickly learn PBS Pro after you learn SLURM.





SBATCH script

To submit a job to SLURM, you need to write a SBARCH script which includes a SBATCH header with several settings.

```
#!/bin/bash
#SBATCH --job-name=job_00
#SBATCH --output=/path/to/output/job_00.out
#SBATCH --error=/path/to/output/job_00.err
#SBATCH --partition=Standard
#SBATCH --time=1:00:00           # [hh:mm:ss]
#SBATCH --mem=5G
#SBATCH --nodes=1
#SBATCH --ntasks=1
#SBATCH --cpus-per-task=4
#SBATCH --mail-user=u_id@anu.edu.au
#SBATCH --mail-type=ALL
```

- `job-name` : the name of your job. You can name it anything.
- `output` : path and filename to store the output log file of this job. It contains information that should be printed on the screen.
- `error` : path and filea to store the error log file of this job. It contains all the error messages has in this job.
- `partition` : is the queue you want to submit your job to. It is used to separating jobs to different queues. We only have one partition on the cluster which is called `Standard` , so all the jobs will be submitted to the same queue.
- `time` : time limit for this job. Format in `days-hh:mm:ss` .
- `mem` : RAM size you want to allocate.
- `nodes` : each node is an independent computer. On our cluster, we have 3 nodes which is dayhoff, wright, and fisher. If you are using parallel processing, you may want to separate your sub-jobs to different nodes to spread the workload.
- `ntasks` : the number of tasks included in this job. It is used when doing parallel processing with `srun` command. If there is no `srun` command, your entire script would be 1 task.
- `cpus-per-taks` : number of CPUs you want to allocate for each task.
- `mail-user` : the email address to receive messages.
- `mail-type` : types of messages you want to receive. `ALL` for everything. `BEGIN` for job begins execution. `END` for job finishes. `FAIL` for job fails.





SBATCH script

On the cluster, you have to specify every directory and file from root /.

To use conda environment, write this in your SBATCH script:

```
source /opt/conda/bin/activate /mnt/data/wright/home/[u_id]/.conda/envs/[env-name]
```

The second path is where conda install packages in our environment.

You can use `cd ~/.conda/envs` to see what's inside.

Avoid using `cd` command in the SBATCH script, it sometimes doesn't work.





Submit a job

Let's save our SBATCH script to “job.sh”. The SBATCH script is also a shell script.

To submit a job, we run ``sbatch job.sh``.





Practise

Please set up the variant calling Conda environment on the cluster.

Download all needed packages.

Modify your previous shell script and submit it as a SLURM job.





Download files from Server to Local

You need to run this command on your local device, not the server.

Let's download the Final Variants files.

Copy from remote to Local:

```
scp u1234567@dayhoff.rsb.anu.edu.au:~/variant-calling/results/*_final_variants.vcf ~/variant-calling
```

Copy from Local to Remote:

```
scp ~/variant-calling/NexteraPE-PE.fa u1234567@dayhoff.rsb.anu.edu.au:~/variant-calling
```



THANK YOU

Contact Us

Jiajia Li
ANU Biological Data Science Institute

RN Robertson Building, 46 Sullivan's Creek Rd
Canberra ACT 2600

E jiajia.li1@anu.edu.au
W <https://bdsi.anu.edu.au/>



Australian
National
University