# ilifu Online Training – Advanced #2 - Resource Allocation

Tinus Cloete

System Administrator & User Support, ilifu
University of Western Cape, March 2024

ilifu

IDIA
Inter-University Institute
for Data Intensive Astronomy

UNIVERSITY of the
WESTERN CAPE

# ilifu: a shared resource-limited cluster

1. Supports a diverse range of projects
   – Astronomy and Bioinformatics
   – Varying resource requirements

   *e.g.*
   *CPU's, Memory,*
   *Running Time, GPU's etc.*

2. Shared environment
3. Resource-limited

# Efficient Usage of Resources

– Resource Allocation: "Picking the right amount of resources for your jobs"

*e.g.*
*if a job uses 100 GB of RAM, don't want to request 232 GB*

| 100 GB | 132 GB |
|:------:|:------:|

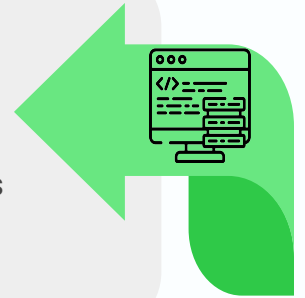– Best practices:  Resource Allocation Guide

# Services and Partitions

### Login node

Run SLURM & bash commands
cd, mkdir, ls, etc

### Jupyter/Dev. node

Development space
New code / workflows / routines
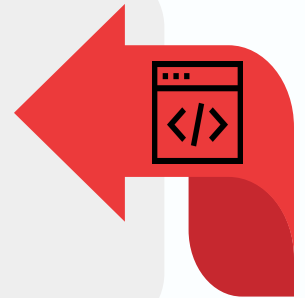Debugging / testing software

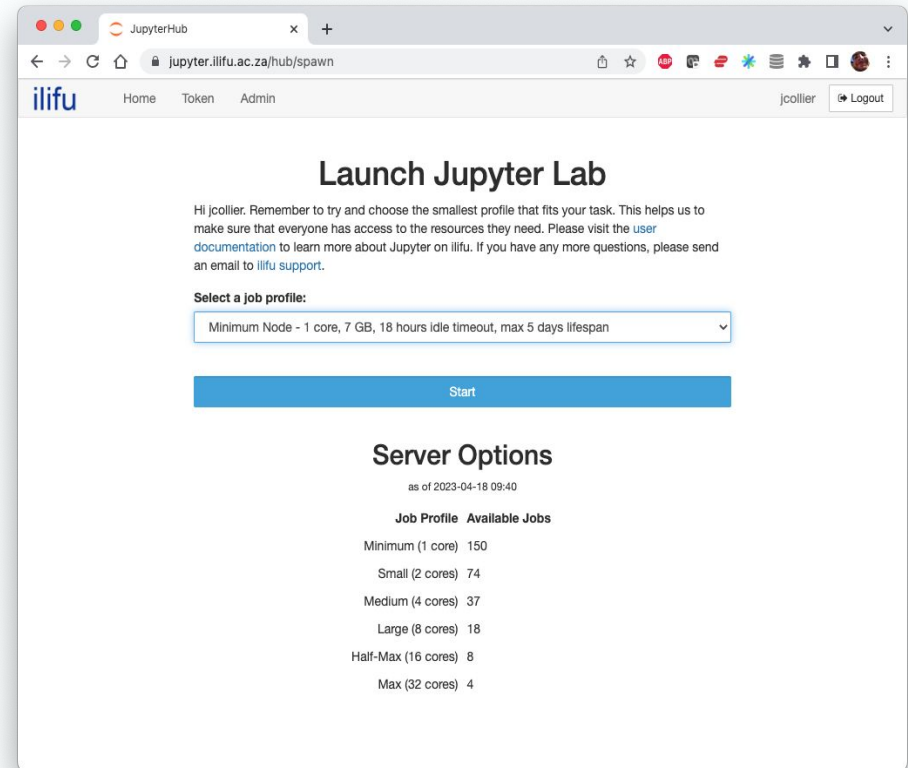### Main partition

Stable, computationally
heavy processing

### HighMem/GPU

For single-high memory jobs
or GPU resources

# Services and Partitions: Jupyter

- Jupyter (Jupyter.ilifu.ac.za)
  - Development space for writing, testing and debugging
  - New code, software, workflows or routines
  - Highly interactive Jupyter notebook environment
  - May be primary interface for stable workflows that shouldn't use Slurm
    - ➤ short analysis routines or other highly interactive workflows

# Jupyter: Resource Allocation

- Select job profile to match your requirements

- Memory often most important

- Jupyter shows current **memory usage** at the bottom

- Emailed about usage stats e.g. low memory usage

- Shut down your session

# Services and partitions: Devel

- [Devel](--partition=Devel)
  - Development of routines within shared resource environment
    - ➢ Submit jobs instantly / quickly
    - ➢ Resources shared, not solely allocated to your jobs
  - Interactivity via a shell
  - Generally for testing higher level workflows and pipelines
  - Access simply using the sinteractive command

```
tcloete@slurm-login:~$ sinteractive
Starting interactive Slurm session.
srun: job 9387238 queued and waiting for resources
srun: job 9387238 has been allocated resources
tcloete@compute-001:~$ |
```

# Services and partitions: Main

- Main partition
  - Default Slurm partition
  - Generally for stable, computationally-heavy workflows and pipelines
  - Can be used for:
    - Many small jobs OR
    - A few large jobs allocated many resources
  - For large workflows, better to first test on Devel or Jupyter

# Services and partitions: GPU and HighMem

- HighMem partition
  - Single high-memory jobs that can't be split into multiple jobs using MPI

- GPU partition
  - Jobs making use of GPUs
  - Not for jobs that only require CPUs (rather use Devel)



```
tcloete@compute-001: ~/demo/interactive_script
tcloete@slurm-login:~$ sinfo
PARTITION    AVAIL  TIMELIMIT  NODES  STATE NODELIST
Main*          up 14-00:00:0      1 drain* compute-002
Main*          up 14-00:00:0     20    mix compute-[012,021,101-105,201-203,205,216,220-221,225,229-230,233-234,238]
Main*          up 14-00:00:0     64  alloc compute-[003-011,013-020,204,206-215,217-219,222-224,226-228,231-232,235-237,239
-260]
Jupyter        up    infinite      5    mix jupyter-[002-006]
Jupyter        up    infinite      5  alloc jupyter-[001,007-010]
JupyterGPU     up 14-00:00:0      2  alloc gpu-[003-004]
HighMem        up 14-00:00:0      3    mix highmem-[001-003]
GPU            up 14-00:00:0      1    mix gpu-007
GPU            up 14-00:00:0      4  alloc gpu-[001-004]
GPU            up 14-00:00:0      2   idle gpu-[005-006]
GPUV100        up 14-00:00:0      1   idle gpu-005
Devel          up  5-00:00:00      1    mix compute-001
tcloete@slurm-login:~$
```

# Primary Resources

- [http://docs.ilifu.ac.za/#/tech_docs/resource_allocation](http://docs.ilifu.ac.za/#/tech_docs/resource_allocation)

- Primary resources
  1. CPU
  2. Memory
  3. Wall-time

- Notes
  – Nodes have 2 CPUs (sockets), each with 16 cores, all of which Slurm calls "CPUs"

# Allocating Resources

- How to allocate resources
  - Accurately determine your resource requirements
  - Use what you require

- Effect
  - Avoid wasting resources (allocated but not used)
  - Increase resource availability
  - Allow other users' jobs to run
  - Improves efficiency of Slurm scheduler
  - Decreased job wait times
  - Better fairshare priority for future job submissions.

# Determining resource requirements

1. Determine parallelism of software
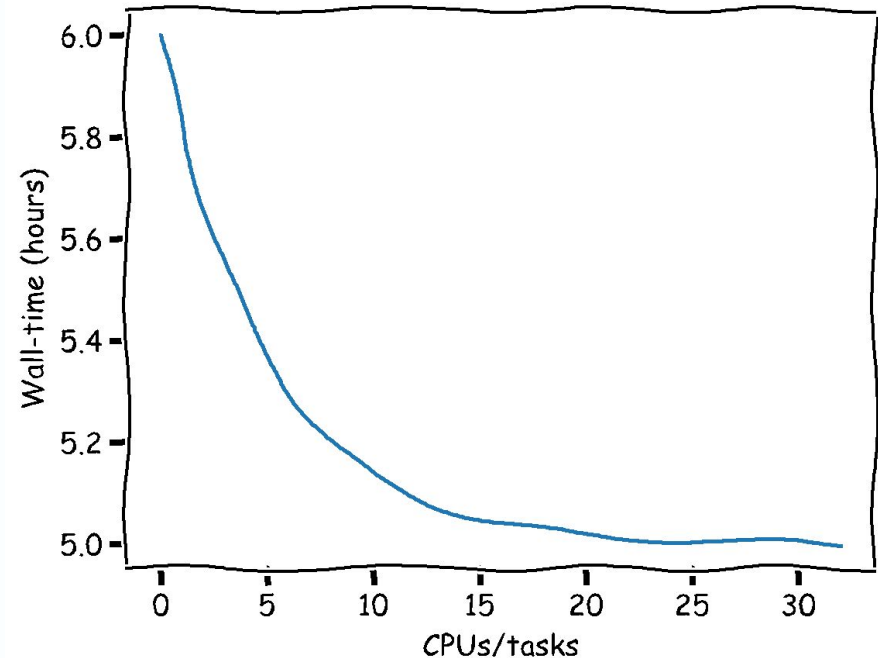2. Profiling previous similar jobs
3. Scaling up test jobs

# Determining resource requirements

- Determining parallelism of software
  - Many software packages only use 1 CPU
  - CPU-level parallelism: Max 1 Node of CPUs
  - Task-level parallelism: >= 1 Node

# Determining resource requirements

- Determining parallelism of software
  - Most parallel processing software doesn't scale linearly
  - Maximum performance often least efficient
    - i.e. shortest wall-time but large allocation necessary
  - Need to find middle ground
  - MPI jobs may perform worse for larger allocations (scatter/gather)
  - Most efficient generally to break into many small independent jobs
    - High-throughput approach

# Profiling previous similar jobs

- Find job ID
  - Job id is shown when you submitted your job
  - Can search for historical jobs
  - Display jobs named 'my-job' submitted during particular time range:
  - `sacct -X --name=my-job --starttime=YYYY-MM-DD --endtime=YYYY-MM-DD`
  - Omit job name (or end time) to show all jobs

- Once you have job ID, you can search for specific information about resource usage

# Profiling previous similar jobs

- Memory usage
  - Find MaxRSS statistic
    - Maximum memory usage of a job (sampled every 20 seconds)
    - Display MaxRSS for job ID 123456 compared to requested memory
    - `sacct -j 123456 --unit=G -o JobID,JobName,MaxRSS,ReqMem`
    - *Notes: 232 Gn = 232 GB per node; 7.25c = 7.25 GB per CPU*

  - Once memory requirement determined
    - Schedule future jobs with **~10-20% buffer**
      - Avoids out-of-memory (OOM) error
    - Avoid excessive usage of memory

# Profiling previous similar jobs

- CPU (and memory) usage
  - Determine used vs. allocated/requested
  - Show Slurm resource efficiency (seff) for job ID 123456
  - Shows % used vs. allocated (for memory, uses MaxRSS stat)
  - `seff 123456`
  - Can run this from Jupyter terminal (to determine resource selection)



```
tcloete@slurm-login:~$ seff 847197
Job ID: 847197
Cluster: ilifu-slurm2021
User/Group: jcollier/idia-group
State: COMPLETED (exit code 0)
Nodes: 1
Cores per node: 32
CPU Utilized: 1-15:22:40
CPU Efficiency: 71.93% of 2-06:44:48 core-walltime
Job Wall-clock time: 01:42:39
Memory Utilized: 213.77 GB
Memory Efficiency: 92.14% of 232.00 GB
```

```
tcloete@slurm-login:~$ seff 201280
Job ID: 201280
Cluster: ilifu-slurm2021
User/Group: jcollier/idia-group
State: COMPLETED (exit code 0)
Nodes: 1
Cores per node: 32
CPU Utilized: 00:00:09
CPU Efficiency: 1.17% of 00:12:48 core-walltime
Job Wall-clock time: 00:00:24
Memory Utilized: 519.09 MB
Memory Efficiency: 0.22% of 232.00 GB
```

# Profiling previous similar jobs

- Wall-time usage
  - Accurate estimation improves Slurm scheduler efficiency and may reduce your job wait time
  - Show used vs. requested wall-time for job ID 123456
  - `sacct -o jobID,jobName,Elapsed,TimeLimit`
  - Once wall-time requirement determined
    - Schedule future jobs with **~20-30% buffer** (avoids job timing out)
    - Avoid excessive wall-time
    - Contact support@ilifu.ac.za to see if we may increase your time limit

# Scaling tests

- Accurately estimating wall-time difficult to do

- Profile previous similar jobs, or

- Run test / scaling jobs
  - Start small test job (e.g. small subset of data)
  - Test the wall-time
    - Reasonable to over-allocate when running scaling test
    - Or if under-estimate, and test small enough, doesn't matter if crashes
  - Repeat process to see how resource usage scales
    - as a function of input (e.g. data volume)
    - as a function of CPUs / tasks (if doing parallel processing)
  - By the end, should have good idea of scaling and efficient choice
    - Allow for buffer for future jobs

# Scaling tests on running jobs

- Get MaxRSS for running job
  - `sstat -j 123456 -o MaxRSS`
  - Given in kB units. Divide by 1024² for GB


- Display real time stats on dashboard (`top` / `htop`)
  - <u>For sbatch</u>:
    First ssh into the login node using authentication forwarding.
    `ssh -A <username>@slurm.ilifu.ac.za`

    It's required to have a job running on a worker node.
    You can then ssh into that worker node (e.g. node 102)
    `ssh compute-102`
  - <u>For Jupyter</u>:  can simply open a new terminal.
  - Now Run:  `htop -u $USER`


- Can monitor real-time usage

# Maximum Resources

- If using **all** CPUs or memory, node becomes fully allocated
  - Any remaining CPUs / memory unavailable to other jobs (incl. your own)

E.g.
Typical worker node:   32 CPU and 232 GB RAM

Job Requesting:         2 CPU and 232 GB RAM  == Full Node
                        30 CPU **not accessible to other jobs**

If possible to split into two smaller jobs, if they ran on different nodes then:

 1 CPU and 116 GB                        1 CPU and 116 GB
 31 CPUs accessible                       31 CPUs accessible

# Account allocation

- Each ilifu project has a [Slurm account](#)

- Resource usage charged against account (affects [fairshare](#))

- View your accounts:
  - `shelp`

- View your default account
  - `sacctmgr show user $USER`

- Change default
  - `sacctmgr modify user name=${USER}`
    `set DefaultAccount=<account>`

- Set account (after #SBATCH for sbatch jobs)
  - `--account=b05-pipelines-ag`

# Resource Allocation Guide

DEMO TIME!

# Data Management Guidelines

- Hot off the press!

- https://docs.ilifu.ac.za/#/data/data_management

# Best practices

- Don't run software / heavy processes / scp on the login node
  - Only submit jobs and run SLURM commands (sbatch, srun, squeue, etc)
  - Use transfer.ilifu.ac.za to transfer data (external/internal), not login node

- Before running a large job, identify the available resources
  - Use sinfo. Don't hog the cluster. Reduce your allocation if possible
  - Increase likelihood of jobs running with less memory and less walltime

- Use sbatch (srun / screen / tmux / mosh are volatile)

- Cleanup files that aren't needed
  - Old raw data, temporary products, /scratch data, etc

- Don't place large files in your home directory (/users)

- Use Singularity (you cannot install software on the nodes)

# Thank you!

Thanks to Jordan Collier for letting me use his Slides

**Remember our support channels!**

support@ilifu.ac.za
https://docs.ilifu.ac.za