

HPC-101 Onboarding

Director: Gwen Jacobs
Sean Cleveland, Adriana Comerford,
Ron Merrill, David Schanzenbach

Information Technology Services
Cyberinfrastructure
University of Hawai'i

<https://www.hawaii.edu/its/ci/>
<http://datascience.hawaii.edu/>
uh-hpc-help@lists.hawaii.edu

March 8, 2021



UNIVERSITY OF HAWAII'

Part I – Introduction:

- 1 Terminology
- 2 University of Hawai'i High Performance Compute Cluster

Part II – Using Mana:

- 1 Connecting via SSH
- 2 Directories & Centralized Software
- 3 Job Scheduler

Part III – High Level View of Policies:

- 1 UH Policies
- 2 Mana Usage Policies



Part I

Introduction



UNIVERSITY OF HAWAII

1 Terminology

2 University of Hawai'i High Performance Compute Cluster

- Storage
- User Support
- Condo Program & Leasing



- **Node** – Another name for a server or computer
- **Login node** – A specialized node that users connect to in order to submit work to a computer cluster
- **Computer cluster** – A set of loosely or tightly connected nodes that work together so that, in many respects, they can be viewed as a single system¹
- **Data transfer node (DTN)** – Specialized nodes that minimize the impedance on the network to access the full capability of the network
- **Science DMZ (SciDMZ)** – A portion of the network configured with equipment and security policies in order to optimize for high-performance scientific applications rather than for general-purpose business systems or “enterprise” computing
- **Multi-factor Authentication (MFA)** – An authentication method in which a computer user is granted access only after successfully presenting two or more pieces of evidence or factors to an authentication mechanism, e.g., DUO

¹https://en.wikipedia.org/wiki/Computer_cluster

- **Symbolic Link (symlink)** – A file that contains a reference to another file or directory
- **Command-line interface/interpreter (CLI)** – A text-based user interface used to view and manage computer files
- **Message Passing Interface (MPI)** – A standard that is used by programs to pass messages between nodes
- **High Performance Compute (HPC)** – A computing paradigm in which applications are typically a tightly coupled parallel job that benefit from a low-latency interconnect
- **High Throughput Compute (HTC)** – A computing paradigm that focuses on the efficient execution of a large number of loosely-coupled tasks
- **Modified time** – The last time the file was modified (content has been modified)²
- **Shell script (script)** – A computer program designed to be run by a CLI³

²<https://unix.stackexchange.com/a/2465>

³https://en.wikipedia.org/wiki/Shell_script

1 Terminology

2 University of Hawai'i High Performance Compute Cluster

- Storage
- User Support
- Condo Program & Leasing



University of Hawai'i High Performance Compute Cluster

- The University of Hawai'i High Performance Compute Cluster - **Mana** is **free** to use for all active faculty, staff, and students affiliated with the University of Hawai'i
- Community acquired nodes are equally accessible to all users
- Nodes purchased through the **condo program** are shared with the community, but priority is given to the node owner and their agents
- Nodes may be **leased** from the community pool
- Additional permanent storage can be leased by faculty & staff

Mana Resource Summary

| | Nodes | CPU Cores | Memory | GPUs | Home/Group Space | Scratch Space | Storage for lease |
|-------|-------|-----------|---------|------|------------------|---------------|-------------------|
| Total | 345 | 8,452 | 62.9 TB | 120 | 80 TB | 211 TB | 1 PB |



① Home Storage

- **Purpose:** Personal storage for applications and active data that needs to persist on Mana

② Group/Lab Storage

- **Purpose:** Group/Lab storage allows users to share data & applications with a need for persistence on Mana

Permanent storage options on Mana have the following attributes:

- 50 GB default quota with a max quota of 300 GB
- Quota increases are re-evaluated annually
- Available on all nodes
- Freely available to users

- ❶ Lustre®
 - A parallel distributed file system
- ❷ NFS
 - A distributed file system

Scratch file systems on Mana have the following attributes:

- 5 TB quota
- **Purge policy** – 10 days based on file modify time
- Available on all nodes
- Freely available to users

Online documents & FAQ

Users are encouraged to look through the online documentation & FAQ prior to contacting ITS-CI directly. Many questions we receive are repeat questions and we try to capture them in our FAQ.

- [Information & policies](#)
- [FAQ](#)
- [Video Tutorials](#)

Contact Information

If your question is not answered in our online documentation, please contact us at:
UH-HPC-Help@lists.hawaii.edu

- For batch jobs . . .
 - Job ID, path to submission script, submission command, error file location, output files
- For other problems . . .
 - State the problem, command issued, host, directory, remote host, error messages

What is the Condo Program & leasing?

Condo Program

- The condo program allows faculty & staff to purchase nodes and have them integrated with Mana
- Condo nodes can take advantage of networking and storage infrastructure that one may not typically have access to
- Condo nodes are managed and maintained by ITS staff
- Nodes are purchased with a five year warranty

Leasing

- Nodes leased from the community can have a contract period from one month, up to one year
- Node lessees are provided priority access to their leased hardware
- Leases are not considered an equipment purchase



Part II

Using Mana



UNIVERSITY OF HAWAII'

1 Connecting via SSH

2 Directories & Centralized Software

- User Directories
- Modules

3 Job Scheduler

- Terminology
- Slurm
- Commands
- Submitting Jobs
 - Interactive Jobs
 - Batch Jobs
- Partition Details
- Constraints & General Resources



Requirements

- Valid UH credentials
- Registered for [MFA/DUO](#)
- Familiarity with a SSH client & a file-transfer method
- Comfortable with the CLI

Connection Information

- Login node:
 - **uhhpc.its.hawaii.edu**
- DTNs:
 - **hpc-dtn1.its.hawaii.edu**
 - **hpc-dtn2.its.hawaii.edu**

Valid Credentials

- Your UH user name
- Accepted forms of authentication
 - UH Password + MFA
 - SSH key + MFA

Try and connect to the Mana Login node now using your SSH client



1 Connecting via SSH

2 Directories & Centralized Software

- User Directories
- Modules

3 Job Scheduler

- Terminology
- Slurm
- Commands
- Submitting Jobs
 - Interactive Jobs
 - Batch Jobs
- Partition Details
- Constraints & General Resources

Filesystem

```
[testuser@login001 ~]$ df -h
```

| Filesystem | Size | Used | Avail | Use% | Mounted on |
|--|------|------|-------|------|--------------------------------|
| fs01:/mnt/datastore/hpc/home/testuser | 50G | 128K | 50G | 1% | /home/testuser |
| fs02:/mnt/datastore/hpc/scratch/testuser | 5.0T | 256K | 5.0T | 1% | /mnt/scratch/nfs_fs02/testuser |
| 10.103.15.230@o2ib,10.103.15.231@o2ib:/scratch | 61T | 5.1T | 55T | 9% | /mnt/scratch/lustre_01 |

Home

```
[testuser@login001 ~]$ ls -l
```

total 1

| | | | | | | | | |
|------------|---|----------|----------|----|-----|----|-------|--|
| drwxr-xr-x | 3 | testuser | testuser | 23 | Jan | 15 | 20:38 | examples |
| lrwxrwxrwx | 1 | testuser | testuser | 19 | Jan | 15 | 20:38 | nfs_scratch -> /mnt/scratch/nfs_fs02/testuser |
| lrwxrwxrwx | 1 | testuser | testuser | 19 | Jan | 15 | 20:38 | lus_scratch -> /mnt/scratch/lustre_01/scratch/testuser |

- ~/examples contains example scripts to use as templates
- ~/nfs_scratch is a symlink to the NFS scratch file system
- ~/lus_scratch is a symlink to the lustre scratch file system

Modules

A tool to help users manage their Unix or Linux shell environment, by allowing groups of related environment-variable settings to be made or removed dynamically.⁴

Commands

- 'module avail' – list installed modules
 - 'module show <module name>' – Show what actions a module performs
 - 'module load <module name>' – Loads the named module
 - 'module spider <search string>' – search the modules list for a match
 - 'module list' – Show what modules are loaded
 - 'module purge' – Unload all loaded modules
-
- Hidden modules can be shown using the --show_hidden flag, e.g., 'module --show_hidden avail'
 - We create modules for frequently requested software packages for all users to access
 - Compilers, libraries, interpreters, applications are all added as modules
 - Users are encouraged to install software in their home/group directories
 - Modules can be listed on the login nodes, but loaded applications will only work on the compute nodes
 - Mana currently uses `lmod`

⁴ [https://en.wikipedia.org/wiki/Environment_Modules_\(software\)](https://en.wikipedia.org/wiki/Environment_Modules_(software))



1 Connecting via SSH

2 Directories & Centralized Software

- User Directories
- Modules

3 Job Scheduler

- Terminology
- Slurm
- Commands
- Submitting Jobs
 - Interactive Jobs
 - Batch Jobs
- Partition Details
- Constraints & General Resources



- **Job Scheduler** – A tool/application to control and prioritize the execution order of unrelated jobs
- **Job** – Another name for a script or application that is to be executed
- **Job ID** – A number assigned to each job submitted to the job scheduler
- **CPU/Socket** – A processing unit in the node which may contain one or more cores
- **Core** – A processing element on a CPU (Multi-threading)
- **Task** – An instance of a running program or process (MPI)
- **Partition** – A group of nodes divided into possibly overlapping sets, which also contains constraints for the given set of nodes

<http://slurm.schedmd.com/quickstart.html>



Mana uses the Slurm Workload Manager to allocate nodes and assign jobs to them

How it works

Jobs are not executed in a **first in first out** manner. Instead, jobs are assigned a priority, which is continuously being re-evaluated for pending jobs.

Depending on load, some resources may go idle while waiting for sufficient free resources for a higher priority job. In these cases, the scheduler will use what is known as **backfilling** to fill in the idle machines with jobs that will not affect the start time of higher priority jobs.

https://en.wikipedia.org/wiki/Slurm_Workload_Manager

<http://slurm.schedmd.com/slurm.html>

Basic

- ***sbatch*** – Used to submit a job script for later execution
- ***srun*** – Used to submit a job for execution or initiate job steps in real time
- ***scancel*** – Used to cancel a pending or running job or job step

Informational

- ***squeue*** – Reports the state of jobs or job steps
 - ***sinfo*** – Reports the state of partitions and nodes managed by Slurm
 - ***sacct*** – Reports job accounting information about active or completed jobs
- Examples usage of the Slurm commands can be seen on schedmd's [quickstart](http://slurm.schedmd.com/quickstart.html)

<http://slurm.schedmd.com/quickstart.html>

Command

```
[login ~]$ srun -I30 -p sandbox -N 1 -c 1 --mem=6G -t 0-01:00:00 --pty /bin/bash
```

Options

- **-I30** – exit if resources are not available within the time period specified (30 seconds)
- **-p sandbox** – Submit my interactive job to the sandbox partition
- **-N 1** – Number of nodes requested (If omitted, default is 1)
- **-c 1** – Number of cores per task requested (If omitted, default is 1)
- **--mem=6G** –Memory allocated per node (See partition details for defaults)
- **-t 0-01:00:00** – How much time you are requesting (DD-HH:MM:SS)
- **--pty** – Execute initial task in pseudo terminal mode
- **/bin/bash** – Task to execute

Interactive jobs terminate when the specified time has elapsed or if you give the **exit** command.
Interactive jobs are good for testing, compiling and relatively short jobs.
Longer jobs should use a shell script and **sbatch**.



Batch job using sbatch

Command

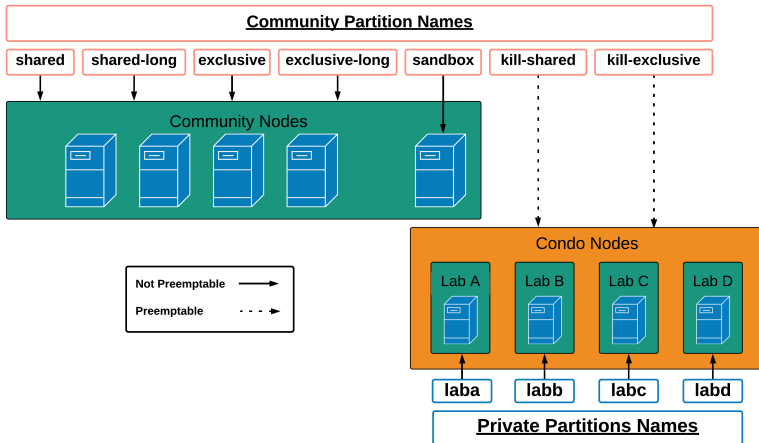
```
[login ~]$ sbatch <path to shell script>
```

Info

- Where sbatch is executed, becomes the jobs working directory
- Submission scripts are shell scripts that begin with special comments that are parameters for the scheduler
- Parameters are evaluated with the command-line taking precedent over what the shell script contains
- Jobs submitted with sbatch are assigned a job ID by Slurm

Please navigate to `~/examples/slurm/non_mpi` and try to submit the example batch submission script (`non_mpi.slurm`) using sbatch.

UH HPC Partitions



Partitions

Details

| Partition | Max walltime | Jobs - total(running) | Max nodes per job | Default memory | Shared | Preemption |
|----------------|--------------|-----------------------|-------------------|----------------|--------|------------|
| gpu | 3-00:00:00 | ∞ | 1 | 512 MB | YES | NO |
| sandbox | 0-04:00:00 | $\infty(2)$ | 2 | 512 MB | YES | NO |
| shared | 3-00:00:00 | ∞ | 1 | 512 MB | YES | NO |
| shared-long | 7-00:00:00 | 5(2) | 1 | 512 MB | YES | NO |
| exclusive | 3-00:00:00 | ∞ | 20 | ∞ | NO | NO |
| exclusive-long | 7-00:00:00 | 5(2) | 20 | ∞ | NO | NO |
| kill-shared | 3-00:00:00 | ∞ | 1 | 512 MB | YES | YES |
| kill-exclusive | 3-00:00:00 | ∞ | 20 | ∞ | NO | YES |



Constraints, General Resources & Core Specialization

--constraint

Users can select nodes for a job by specifying a value or values for the constraint option. Only nodes having features matching the job constraints will be used to satisfy the request. Multiple constraints may be specified with "&" (AND), "|" (OR), etc.

--gres

Specifies a comma delimited list of generic consumable resources which a job should be granted access to.

--core-spec

On Mana, the scheduler by default withholds 1 core per node from use by users through a feature in Slurm known as **core specialization**. This is done to improve job performance. This can be overridden by users but is not recommended.^{1 2}

¹ https://slurm.schedmd.com/core_spec.html

² https://slurm.schedmd.com/SUG14/process_isolation.pdf

Information on Constraints, general resources and node specifications



Part III

High Level View of Policies



UNIVERSITY OF HAWAI'I

1 UH Policies

2 Mana Usage Policies

- Login nodes & DTN nodes
- Storage Policy
- User Account Life Cycle



Chapter 708, Hawaii Revised Statutes

- Access only by authorized people
- Should not be used in the act of committing a crime

University of Hawaii Executive Policy E2.210

- Protect your password (we will never ask or require your password)
- Computer resources should not be used to test or compromise systems without prior authorization
- University resources are intended to be used for institutional purposes and may not be used for private gain



1 UH Policies

2 Mana Usage Policies

- Login nodes & DTN nodes
- Storage Policy
- User Account Life Cycle



The following actions are acceptable on the shared systems in Mana:

- 1 File/Directory management [Login & DTN nodes]
- 2 Text editing with a text editor: vi/vim, emacs, nano [Login & DTN nodes]
- 3 Transferring files to and from the cluster (scp, rsync, SFTP, globus, aspera, lftp, etc.) [Login & DTN nodes]
- 4 Shall be used to submit batch and interactive jobs [Login nodes]
- 5 SSH shell access [Login & DTN nodes]

All other action shall take place on a compute node using an interactive session. If any actions outside of the sanctioned activities are detected, the following escalation will take place:

- 1 The process will be killed without notice
- 2 If the user continues this action, we will notify the user of violating the policy
- 3 If the user continues to ignore our warnings, the user may be banned from the cluster for a duration of time

- 1 ITS is not responsible for any data that is deleted or loss due to user error, hardware failure, administrator error. Users are responsible for their own data and are highly encouraged to backup data that is important to them at other storage locations
- 2 Files located on the scratch file systems are subject to a **10 day** purge policy that is based on the file/directory modified timestamp. **Files that are purge cannot be recovered.** Users are encouraged to copy their results off the scratch file system as soon as possible upon completion of their job
- 3 Home storage currently snapshots once a day per user in case of accidental file deletion



- 1 Accounts that are idle for 120 days will be locked
- 2 Accounts that have been locked can be unlocked upon request with no further action then emailing us a request to unlock their account
- 3 Accounts that are idle for a total of 180 days will be purged from the cluster
- 4 Users that have been purged may request a reactivation of their account, but may be required to take a quiz to confirm retention of how to use Mana
- 5 Users that will be leaving UH or are no longer active at UH, will need to have an active UH faculty/staff member in good standing, request a Mana account extension.



Where to find all Mana specific policies

- 1 Common systems & Storage
- 2 User Account life cycle
- 3 Security
- 4 For lease storage
- 5 Condo Program & Off Warranty condo nodes



Questions?



UNIVERSITY OF HAWAI'I