



Advanced HPC course

Agenda

1) Quick Overview of HPC (15 mins)

2) HPC Job Submission (60 min)

break (30 min)

3) Software Compile / Installs / Misc (60min)

Module 1: Quick Overview of HPC

What is HPC ?

- HPC, or high-performance computing, refers to the application of supercomputers or clusters of computers to computational problems that typically arise through scientific inquiry.
- HPC is useful when a computational problem:
 - **Is too large** to solve on a conventional laptop or workstation (because it requires too much memory or disk space) or ...
 - **Would take too long** (because the algorithm is complex, the dataset is large, or data access is slow) or ...
 - **Are too many** – High Throughput Computing

Reasons to use UCT HPC ?

- You have a program that can be recompiled or reconfigured to use optimized numerical libraries that are available on HPC systems but not on your own system.
- You have a "parallel" problem, e.g. you have a single application that needs to be rerun many times with different parameters.
- You have an application that has already been designed with parallelism.
- To make use of the large memory available.
- Our facilities are reliable and regularly backed up.

When not to use HPC ?

- Cannot host databases on HPC, flat file databases are allowed but not relational DBs
- Graphical User Interface (GUI) applications can be used but users are asked to be cautious.

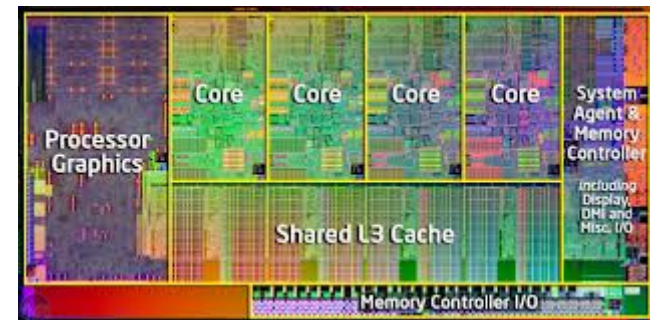
Parallelism on HPC

- Programs for HPC systems must be split up into many smaller “sub-programs” which can be executed in parallel on different processors
- Writing parallel software can be challenging, and many existing software packages do not support parallelism & may require development.

NOTE: Many tasks cannot be parallelised

What does HPC consist of ?

- HPC is the aggregation of computing resources.
 - Cores (cpus / sockets)
 - RAM
 - Disk
 - Interconnect



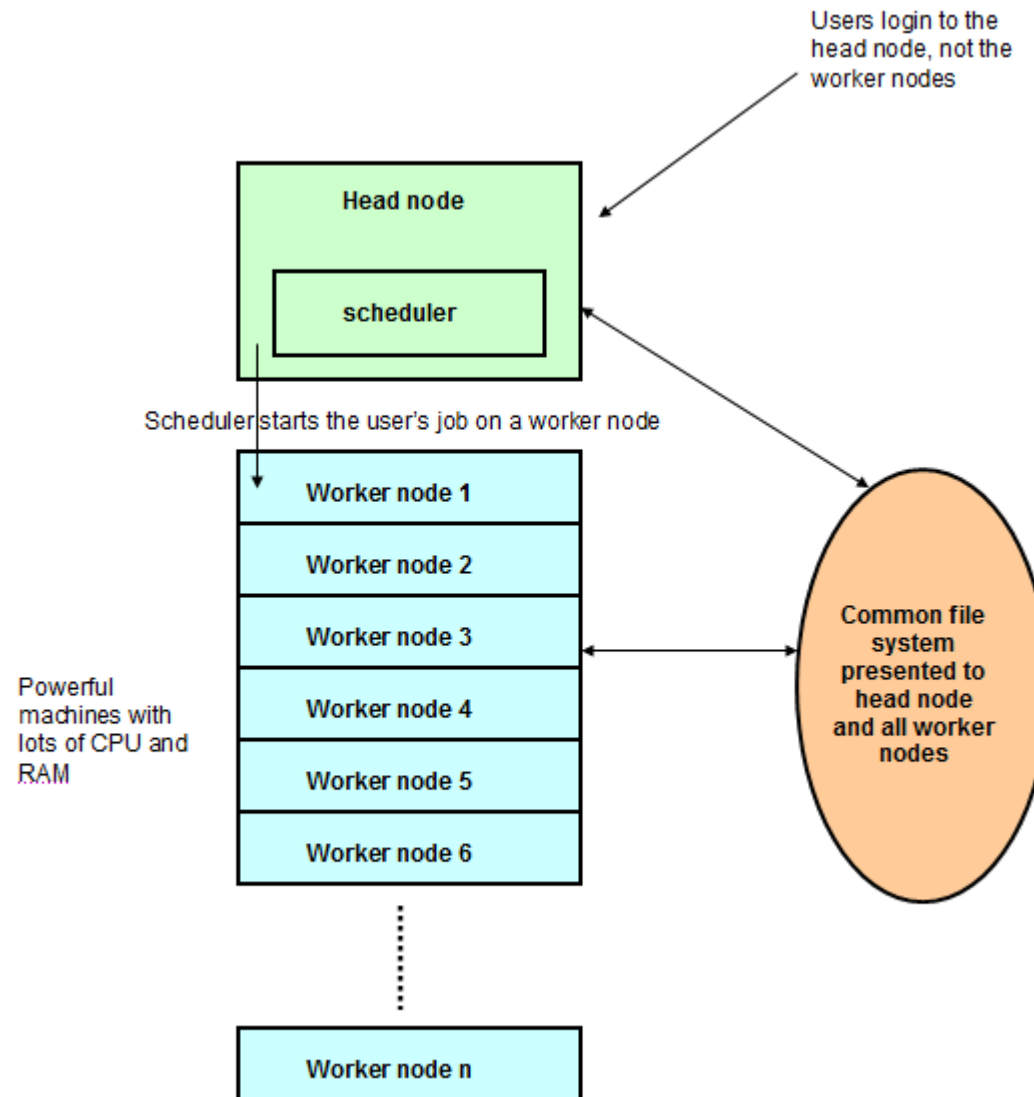
Cluster Architecture

- Operating system: Centos 7.4
- X86_64
- Scheduler: SLURM
- Worker nodes:
 - 27 Dell C6420 – Multi Core
 - 12 Dell C6145 – Many Core / dense array
 - 4 GPU servers (Tesla M2090 / K40)
 - 2 High Memory Machines – Dell R820 – 1TB RAM

BeeGFS Storage nodes:

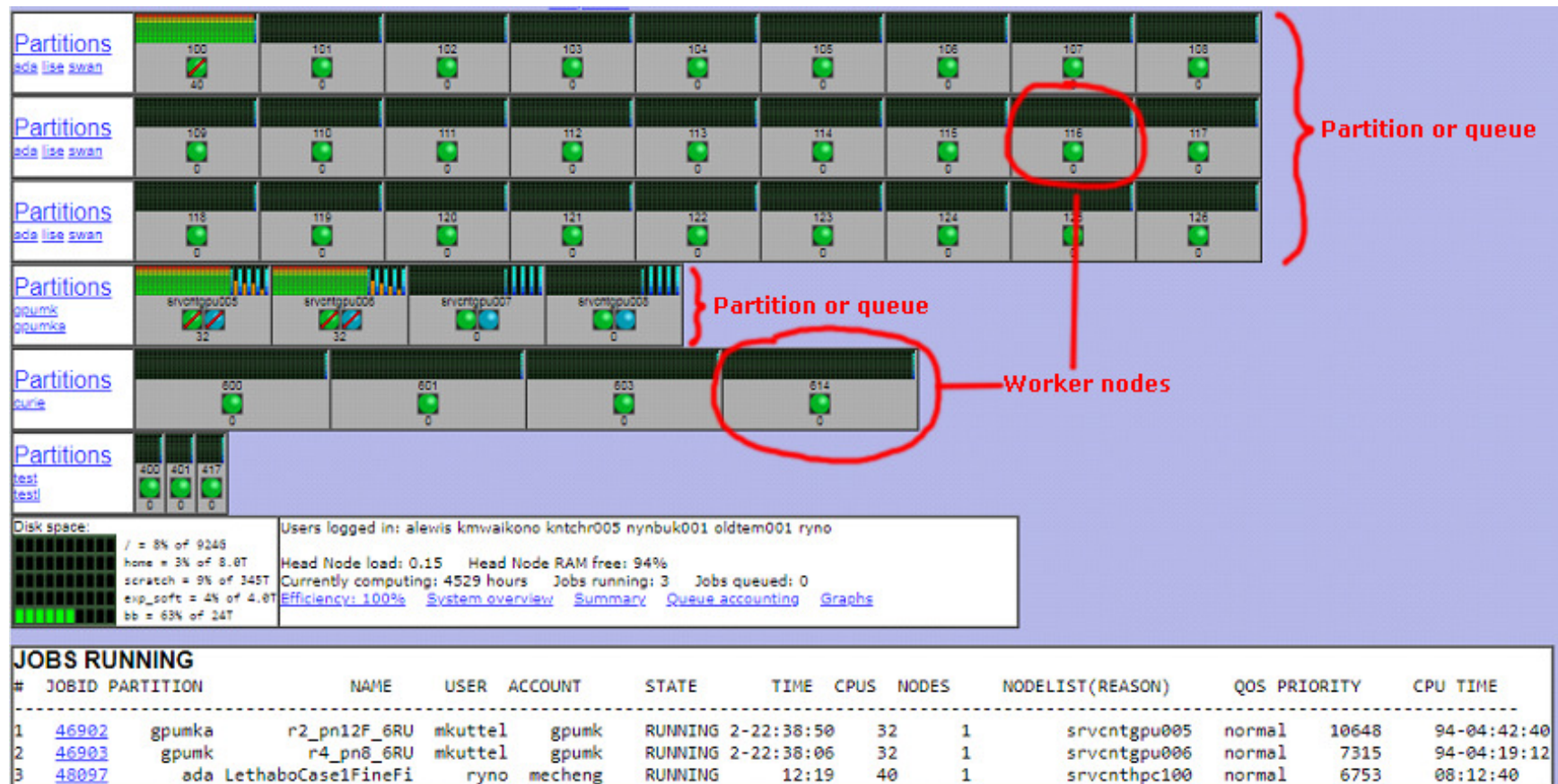
- 4 Dell R740s
- 360TB of scratch storage

Architecture









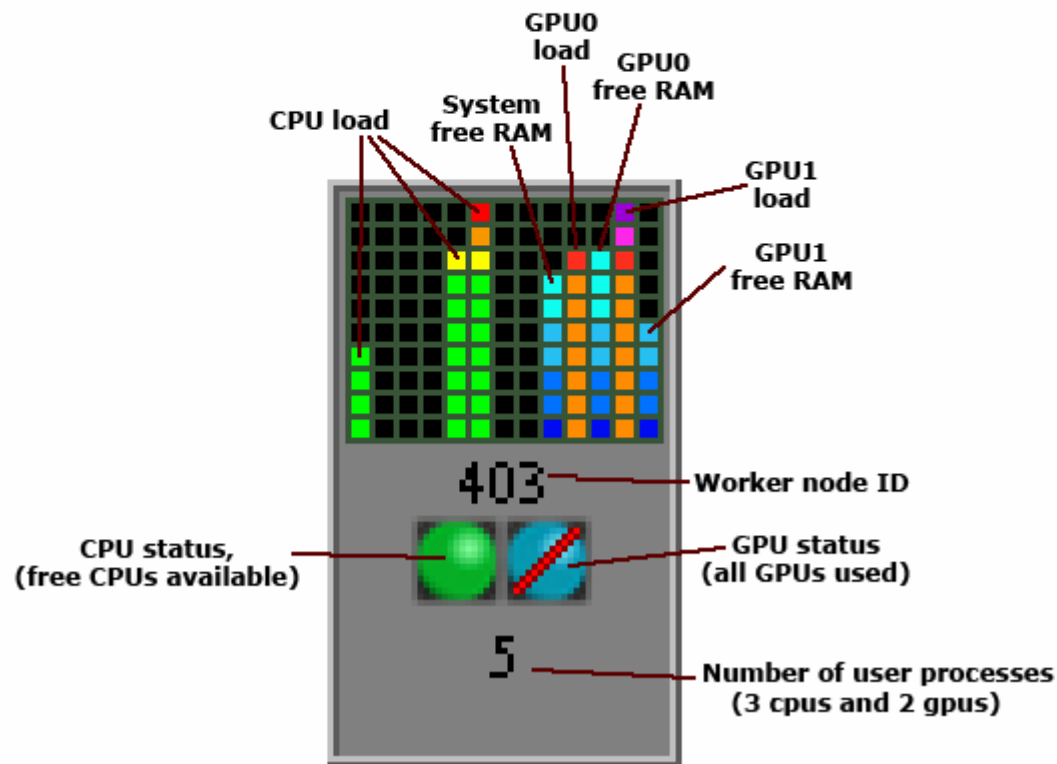
The dashboard

- To keep track of the cluster's status, workload and the jobs that are running go to: <http://hpc.uct.ac.za/db>



The dashboard

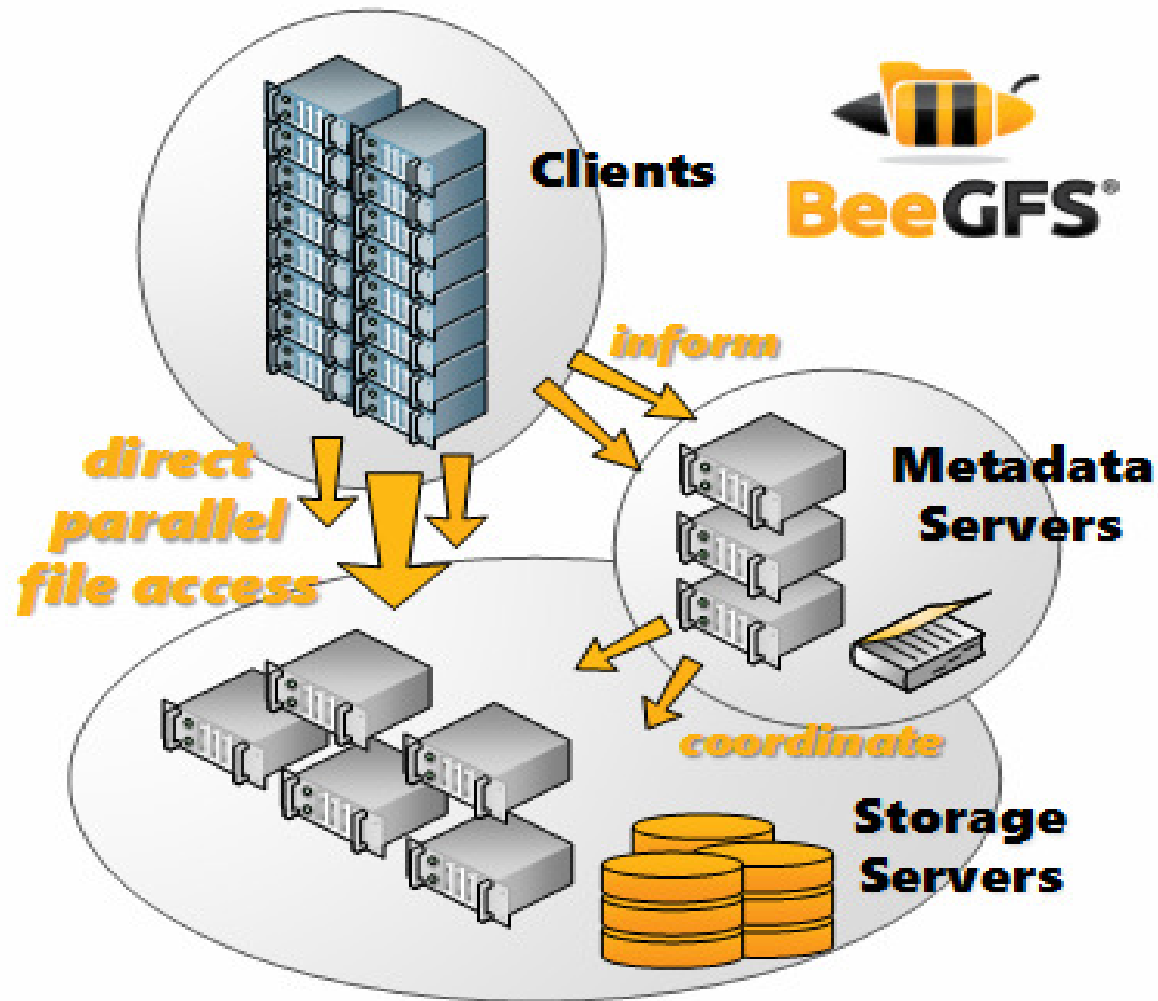
Icon	Value	Description
	Free CPUs	There are free CPUs, jobs may be submitted to this node.
	Job-exclusive	All CPUs are busy, the node is running but no further jobs may be submitted.
	Busy	Torque mom daemon or CPUs too busy to respond to further requests. Jobs are running but may be degraded.
	Down	Node down or PBS mom daemon offline or not responding, no jobs may be submitted.
	Free GPUs	There are free GPUs, jobs may be submitted.
	Busy	All GPUs are busy, the node is running but no further jobs may be submitted.



BeeGFS Parallel Storage

- Pure software solution for scale-out parallel network-storage.
- Each HPC node is connected with IB cables to the IB switch. The BeeGFS store is connected to the same switch.
 - /scratch
- Advantages : Very fast storage
- Disadvantages: No backups, “volatile” area, scrubbing policy available on the HPC website

BeeGFS Architecture



BeeGFS connected to HPC

- Parallel storage is connected via Infiniband (RDMA only). The only TCP connection which exists is for Admon / MGMT services.
- TCP is the backup protocol should RDMA (IB switch) fail.
- Headnode is identical to worker nodes and also has IB.
- Once your job executes on a worker node, traffic to the storage service is $\sim 100\text{Gb/sec}$.
- Parallel software can also communicate over IB if compiled with OpenMPI and PMIx.

Module 2:
Various Job Submission Methods
– Interactive

Software Required

Use your web browser to download Putty and PuttySCP from:
<http://www.putty.org>

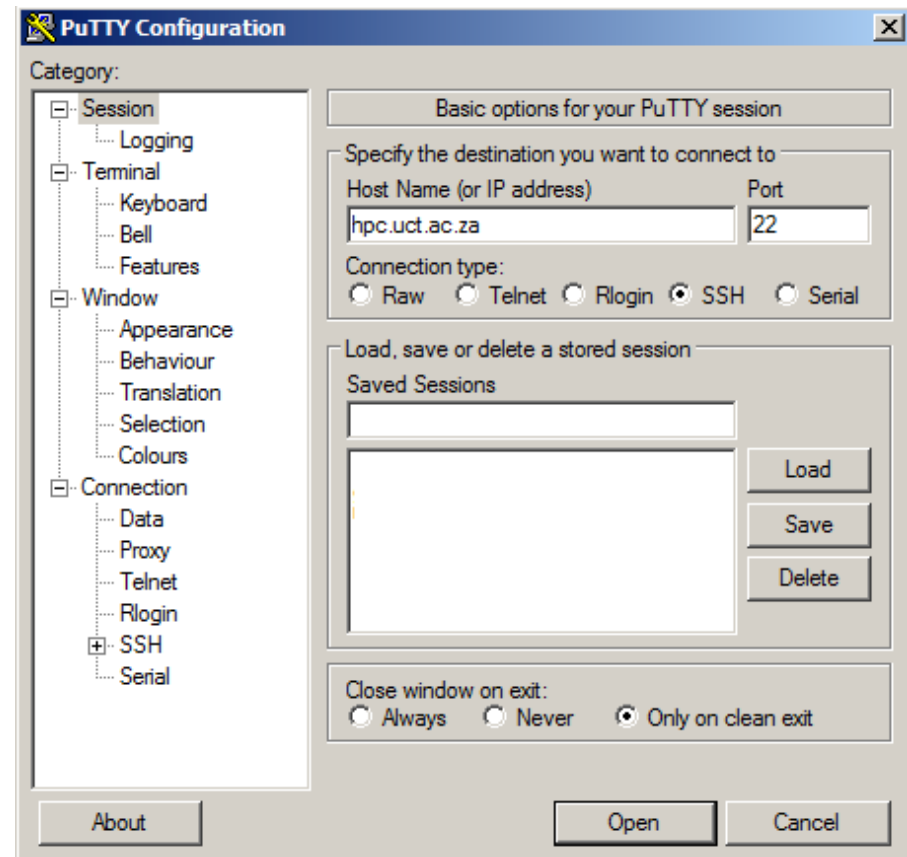
- Click on the “Download Putty” link and download:
 - putty.exe (a Telnet and SSH client)
 - WinSCP (GUI-BASED SCP)
- Double click to install on your PC.
- MacOS users may launch a terminal
- Xming for Windows (**Tick NoACL**) / MacOSX users may use Quartz

Course Credentials

- Start the putty telnet/ssh client by double clicking on putty.exe and connect to the HPC Machine

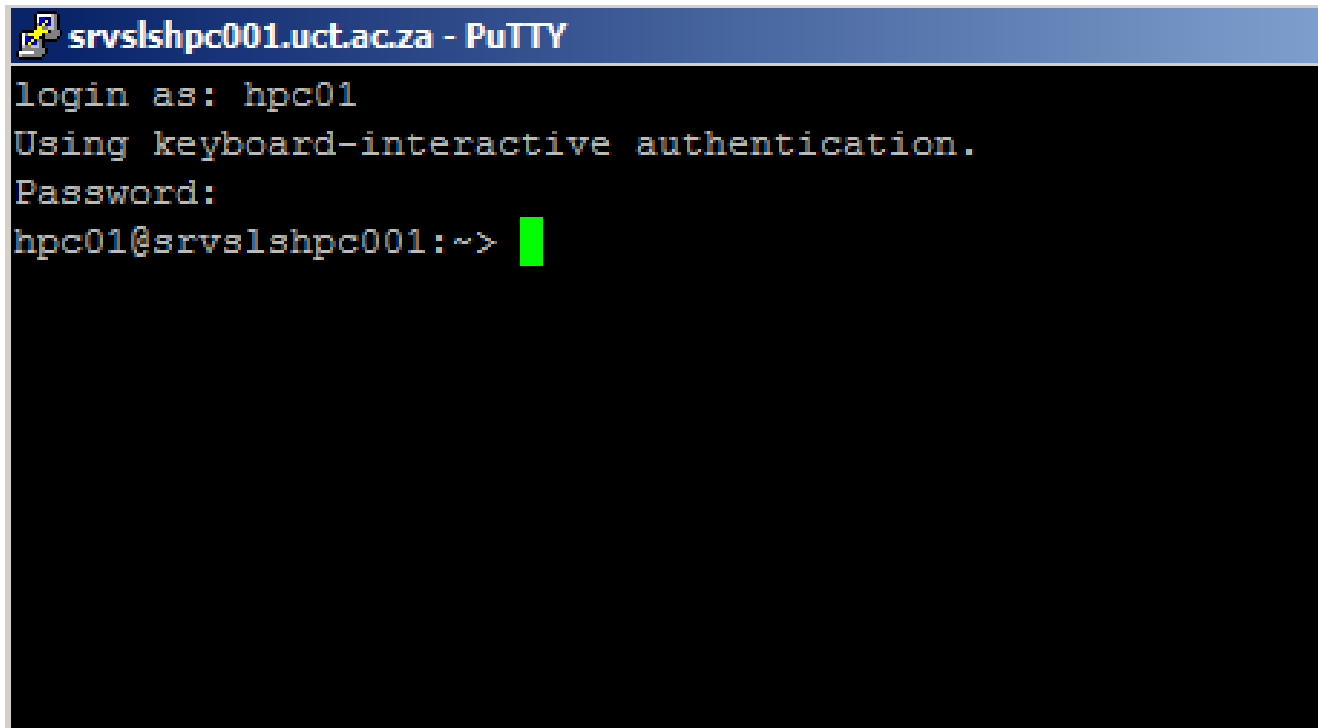
- Host: `hpc.uct.ac.za`
- Connection Type: `ssh`
- Port: `22`

1. Select SSH, X11, Enable X11
2. Click on session, top left.
3. Saved sessions: hpc
4. Click the **save** button.



Course Credentials

- Log into the training HPC system using the Test Account allocated to you, e.g.
 - **Account Name:** hpc0(n)
 - **Password:** will not be displayed as you type



```
srvslshpc001.uct.ac.za - PuTTY
login as: hpc01
Using keyboard-interactive authentication.
Password:
hpc01@srvslshpc001:~>
```

Download training material

```
wget http://hpc.uct.ac.za:/db/training-material.tgz
```

```
tar zxvf training-material.tgz
```

```
cd scripts
```

Modules

- Switching between multiple versions of the same application.
- Use Case: Single job requires functionality from one version of a application and functionality from another version of the same application.
- Sets up Library / Include / Bin / Custom Paths
- `module avail` - Lists all modules available
- `module load <module>` - Loads a specific module

Modules

- Can place in shell script or set on head node but then environment export required.
 - `--export=<environment variables [ALL] | NONE>`
- If using MPI and distributed memory then module command should be placed in your `.bashrc` file.

Standard Job Submission

```
#!/bin/bash
#SBATCH --account icts
#SBATCH --partition=ada
#SBATCH --nodes=1 --ntasks=1
#SBATCH --time=10:10:00

pwd
date
hostname
```

Rules:

No space before **#SBATCH**

No space between **#** and **SBATCH**

Exercise 1 (a)

Add **#SBATCH --mail-user=user@email**

Add **#SBATCH --mail-type=ALL**

sbatch standard-job.sh

Jobs Arrays

- Use Case: Lots of input files, not possible to submit manually.
- Common PBS Environment Variables
 - SLURM_ARRAY_TASK_ID=1
 - SLURM_ARRAY_TASK_COUNT=3
 - SLURM_ARRAY_TASK_MAX=3
 - SLURM_ARRAY_TASK_MIN=1

Exercise 1 (b)

```
sbatch --array=1-6 array-job.sh  
sbatch --array=2-6:2 array-job.sh  
sbatch --array=2,4,6 array-job.sh
```


Interactive Jobs

- `salloc`
- Dangerous as you're still on the head node.
- `srun` to run on cluster
- Type `exit` when done!!!

Exercise 1 (c)

```
salloc  
hostname  
srun hostname
```

Interactive Jobs

- `srun -- account=icts --time=10:00:00 --partition ada --ntasks=5 --pty bash -l`
Which we have shortened through a wrapper script to:
`sint --account=icts --time=10:00:00 --partition ada --ntasks=5`
- Use Case: Compiling, Debug Application / Testing,
- Advantage: Work directly on a worker node
- Disadvantage: CPU expensive. Get done and exit

Exercise 1 (d)

- `sint --account=icts --time=10:00:00 --partition ada --ntasks=5`

X support

- Must run putty with X forwarding
- Display variable must be set up and exported to job
- All variables exported in example below
- Hint, Start\Run\cmd\ipconfig

Exercise 1 (e)

- **export DISPLAY=137.158.X.Y:0.0**
- **sbatch X-job.sh**
- **qstat**
- **Kill eyes window**
- **qstat**

MPI Jobs

- Message Passing Interface (MPI) is used for communication among the nodes running a parallel program on a distributed memory system.
- Compile mpitest.c - `"mpicc -o mpitest mpitest.c"`
- `"sbatch mpi-job.sh"`
- Important to use srun from the same openmpi version.

MPI Jobs

Exercise 1 (f)

- `module load mpi/openmpi-4.0.1`
- Compile `mpitest.c` - `"mpicc -o mpitest mpitest.c"`
- `"sbatch mpi-job.sh"`
- Here we are running 4 cores on 1 node
- To run on more than 1 node you may need to modify your `.bashrc` file and add the module load command or use the `--export` command

MPI Jobs

Exercise 1 (g)

- We could also use salloc:
salloc --ntasks=4 --nodes=2
module load mpi/openmpi-4.0.1
hostname (we're still on the head node)
srun mpitests (this runs on cluster nodes)
exit (please don't forget to relinquish resources!!!)
- Why don't we get an even spread of node resources, ie why are we given 3 cores on 1 node and 1 core on another?
- ...to preserve as many low usage nodes as possible.
- To insist on an even spread add --ntasks-per-node=2

Modules Exercise 1 (g)

module avail	Shows all modules available
module load python/anaconda-python-2.7	Environment modified for application
which python	Location for which binary
module unload python/anaconda-python-2.7	Unload the module

```
#SBATCH --job_name "Tea Time"  
#SBATCH --time=00:30:00
```


Module 3: Software Compile / Installs / Misc

The HPC software repository does not contain my software

- All software resides in /opt/exp_soft and shared between the HPC worker nodes. Please do not install on /scratch.
- **Problem:** I have a RPM file but cannot install because I do not have root privileges. **Solution:**

```
rpm --prefix=/home/username/install-dir -i app.rpm
```
- **Problem:** I have the source but its such a mission to compile. **Solution:** (1)Make a list of dependencies, (2)download install, (3) Compile and view logs
- Roadmap: UCT HPC Continuous Integration environment for keeping software up-to-date

Let us establish a HPC interactive session

- Please do not run software installs from the head node!
- Start an interactive job:

```
sint --account=icts --time=10:00:00 --partition ada --ntasks=5
```

PEAR - Paired-End reAd mergeR

- Software for merging raw illumina paired-end reads
- One of many Open Source tools in the Bio-Informatics software catalogue.
- It is one of the most popular tools currently being used on our HPC.
- Quick and simple to compile.
- .. however a lot of people are put off by how long it takes to compile applications, GCC being one of them.

Let's compile some software

- `mkdir ~/pear-install`
- Change directory into `~/training-material/software-src/`
- Uncompress with `" tar xfvz pear-0.9.6-src.tar.gz "`
- Change directory into `pear-0.9.6-src`
- `"./configure -- help "` for a list of features and tuning parameters
- `"./configure --prefix=/home/username/pear-install/"`
- `"make -j 5"` - Compile the application. `" -j 5 "` ??
- `"make install"` - Install the compiled binary / lib / include

Working remotely with screen

- Allows you create additional virtual terminals inside a single process called " Screen "
- Use Cases:
 - Works great for unreliable internet connections
 - Long running compilations / file copies
- Execute the command called " screen "
- "ctrl + a +c " - Create additional terminals
- "ctrl +a + n or p" - Move back / forward between tty
- "ctrl +a +d " - Detach from a screen session
- "screen -r -d " - detach and re-attach
- "screen -x " - reattach but keep my remote sys active

Being put off from screen because it doesn't scroll

- " termcapinfo xterm|xterms|xs|rxvt ti@:te@ "

Road Map for UCT HPC 2020

- New\expanded data center.
- Capability to allow departments to purchase their own hardware to intergrate into HPC.
- Provide long term, slower storage directly connected to HPC.
- Provision S3 block storage.

Thank You

Apply for a HPC account

<http://hpc.uct.ac.za>