

Using the Yale HPC Clusters

Robert Bjornson

Yale Center for Research Computing
Yale University

June 2017



What is the Yale Center for Research Computing?

- Independent center under the Provost's office
- Created to support your research computing needs
- Focus is on high performance computing and storage
- ~15 staff, including applications specialists and system engineers
- Available to consult with and educate users
- Manage compute clusters and support users
- Located at 160 St. Ronan st, at the corner of Edwards and St. Ronan
- <http://research.computing.yale.edu>

What is a cluster?

A **cluster** usually consists of a hundred to a thousand rack mounted computers, called **nodes**. It has one or two **login nodes** that are externally accessible, but most of the nodes are compute nodes and are only accessed from a login node via a **batch queueing system** (Slurm).

The CPU used in clusters may be similar to the CPU in your desktop computer, but in other respects they are rather different.

- Linux operating system
- Command line oriented
- Many cores (cpus) per node
- No monitors, no CD/DVD drives, no audio or video cards
- Lots of RAM
- Very large distributed file system(s)
- Connected internally by a fast network

Why use a cluster?

Clusters are very powerful and useful, but it may take some time to get used to them. Here are some reasons to go to that effort:

- Don't want to tie up your own machine for many hours or days
- Have many long running jobs to run
- Want to run in parallel to get results quicker
- Need more disk space
- Need more memory
- Want to use software installed on the cluster
- Want to access data stored on the cluster
- Want to use GPUs

Limitations of clusters

Clusters are not the answer to all large scale computing problems. Some of the limitations of clusters are:

- Cannot run Windows or Mac programs
- Not for persistent services (DBs or web servers)
- Not ideal for interactive, graphical tasks
- Jobs that run for weeks can be a problem (unless checkpointed)

Summary of Yale Clusters (June 2017)

	Omega	Grace	Farnam	Ruddle
Role	FAS	FAS	LS/Med	YCGA
Total nodes	1028	216	320+	156+
Total cores	8500	4700	5300	3000
Cores/node	8	20	8-20	20
Mem/node	36 GB/48 GB	128 GB	128-1500GB	
Network	QDR IB	FDR IB	10 Gb EN	
File system	Lustre	GPFS	GPFS	NAS + GPFS
Batch queueing	Torque	Slurm	Slurm	Slurm
Duo MFA?	No	No	No	Yes

Details on each cluster here:

<http://research.computing.yale.edu/hpc-clusters>

Migration to Slurm

We will be migrating all Yale clusters to Slurm this year.

Tentative schedule, coinciding with scheduled maintenance periods:

- Farnam: Done
- Ruddle: Done
- Grace: June 15
- Omega: ?

We are setting up small slurm clusters for testing:

- `grace-next.hpc.yale.edu`
- `omega-next.hpc.yale.edu` (coming soon)

Setting up a account

Accounts are free of charge to Yale researchers.

Request an account at:

<http://research.computing.yale.edu/account-request>.

After your account has been approved and created, you will receive an email describing how to access your account. This may involve setting up ssh keys or using a secure login application. Details vary by cluster.

If you need help setting up or using ssh, send an email to: hpc@yale.edu.

Ssh to a login node

To access any of the clusters, you must use **ssh**. From a Mac or Linux machine, you simply use the **ssh** command:

```
laptop$ ssh netid@omega.hpc.yale.edu  
laptop$ ssh netid@grace.hpc.yale.edu  
laptop$ ssh netid@farnam.hpc.yale.edu  
laptop$ ssh netid@ruddle.hpc.yale.edu
```

From a Windows machine we recommend using putty: <http://www.putty.org>.

For more information on using PuTTY see:

<http://research.computing.yale.edu/hpc-support/user-guide/secure-shell-for-microsoft-windows>

Understanding ssh key pairs

- We use key pairs instead of passwords to log into clusters
- Keys are generated as pair:
 - ssh-keygen: public (id_rsa.pub) and private (id_rsa)
 - putty/winscp: public (name.pub) and private (name.ppk)
- You should always use a pass phrase to protect private key!
- You can freely give the public key to anyone
- You can generate a new pair for each of your computers, or reuse the same pair.
- NEVER give the private key to anyone!

Outline of setting up keys

- 1 Generate key pair using ssh-keygen (linux/mac) or puttygen (windows)
- 2 Locate the public and private key files you generated
- 3 Use our tool to upload the public key (link on page shown below)
- 4 Wait 15 minutes for key to propagate
- 5 Connect using private key

This can be tricky. See

<http://research.computing.yale.edu/support/hpc/user-guide> for the exact steps or contact us if you have problems.

Sshing to Ruddle

- Ruddle has an additional level of ssh security, using Multi Factor Authentication (MFA)
- We use Duo, the same MFA as other secure Yale sites

Example:

```
bjornson@debian:~$ ssh rdb9@ruddle.hpc.yale.edu
Enter passphrase for key '/home/bjornson/.ssh/id_rsa':
Duo two-factor login for rdb9
```

Enter a passcode or select one of the following options:

1. Duo Push to XXX-XXX-9022
2. Phone call to XXX-XXX-9022
3. SMS passcodes to XXX-XXX-9022

Passcode or option (1-3): 1

Success. Logging you in...

More about MFA

- Don't have a smartphone? Get a hardware device from the ITS Helpdesk
- Register another phone: e.g your home or office phone, as a backup
- See <http://research.computing.yale.edu/support/hpc/user-guide/mfa>

Running jobs on a cluster

Two ways to run jobs on a cluster:

Interactive:

- you request an allocation
- system grants you one or more nodes
- you are logged onto one of those nodes
- you run commands
- you exit and system automatically releases nodes

Batch:

- you write a job script containing commands
- you submit the script
- system grants you one or more nodes
- your script is automatically run on one of the nodes
- your script terminates and system releases nodes
- system sends a notification via email

Interactive vs. Batch

Interactive jobs:

- like a remote session
- require an active connection
- for development, debugging, or interactive environments like R and Matlab

Batch jobs:

- non-interactive
- can run many jobs simultaneously
- your best choice for production computing

General overview of Slurm

Slurm manages all the details of cluster operation:

- Interactive node allocation
- Batch job submission
- Specifying and reserving the resources you need for your job
- Listing running and pending jobs
- Cancelling jobs
- Grouping node resources into partitions
- Prioritizing and scheduling jobs

For information specific to each cluster:

<http://research.computing.yale.edu/support/hpc/clusters>

Interactive allocations

```
srun -p interactive --pty bash
```

You'll be logged into a compute node and can run your commands. To exit, type `exit` or `ctrl-d`

Slurm: Example of an interactive job

```
farnam-0:~ $ srun --pty -p interactive --mem=8g bash
c01n01$ module load Bowtie2
c01n01$ bowtie2 -1 data/sample_R1.fastq -2 data/sample_R2.fastq -x h
c01n01$ exit
farnam-0:~ $
```

Batch jobs

- create a script wrapping your job
- declares resources required
- contains the command(s) to run

Slurm: Example of a batch script

```
#!/bin/bash

#SBATCH --mail-type all --mail-user robert.bjornson@yale.edu
#SBATCH -J bowtie
#SBATCH -c 20
#SBATCH -p general
#SBATCH --mem 8g

module load Bowtie2

time bowtie2 -p 20 -1 ../data/sample_R1.fastq \
  -2 ../data/sample_R2.fastq -x ../hg19/genome \
  > output.sam
```

Slurm: Example of a batch job

```
$ sbatch batch.sh
Submitted batch job 42
$ squeue -j 42
```

JOBID	PARTITION	NAME	USER	ST	TIME	NODES	NODELIST(REASON)
42	general	bowtie	rdb9	R	0:03	1	c13n10

The script runs in the current directory. Output goes to `slurm-jobid.out` by default, or use `-o`

Copy scripts and data to the cluster

- On Linux and Mac, use scp and rsync to copy scripts and data files from between local computer and cluster.

```
laptop$ scp -r ~/workdir netid@omega.hpc.yale.edu:
laptop$ rsync -av -e 'ssh -i ~/.ssh/id_rsa' ~/workdir \
    netid@grace.hpc.yale.edu:
```

- Cyberduck is a good graphical tool for Mac or Windows.
<https://cyberduck.io>
- For Windows, WinSCP is available from the Yale Software Library:
<http://software.yale.edu>.
- Graphical copy tools usually need configuration when used with Duo. See:
<http://research.computing.yale.edu/support/hpc/user-guide/transfer-files-or-cluster>

Summary of Slurm commands

Description	Command
Submit a batch job	<code>sbatch [opts] SCRIPT</code>
Submit an interactive job	<code>srun -p interactive --pty [opts] bash</code>
Status of a job	<code>squeue -j JOBID</code>
Status of a user's jobs	<code>squeue -u NETID</code>
Cancel a job	<code>scancel JOBID</code>
Get queue info	<code>squeue -p PART</code>
Get partition info	<code>sinfo -p PART</code>
Get node info	<code>sinfo -N</code>
Get job info	<code>sacct -j JOBID -l</code>

Summary of Slurm sbatch/srun options

Description	Option
Partition	-p <i>PART</i>
Process count	-c <i>CPU</i>
Node count	-N <i>NODES</i>
Task count	-n <i>TASKS</i>
Wall clock limit	-t <i>HH:MM</i> or <i>DD-HH</i>
Memory limit	--mem 4g or --mem-per-cpu 4g)
Interactive job	srun --pty bash
Job name	-J <i>NAME</i>

Examples:

```
farnam1$ srun --pty -p general --mem-per-cpu 32g -t 24:00 bash
farnam1$ sbatch -p general --mem 8g -c 20 -t 3- job.sh
```

Partitions

- Compute nodes are grouped into partitions
- Each cluster has its own set of partitions
- Jobs are submitted to a specific partition
- Each partition has rules:
 - Who can submit jobs
 - How many cores and how much memory per user
 - Maximum walltime

`general/day/long` default, common usage

`gpu` nodes with gpus

`bigmem` nodes with large RAM

`PI` reserved for specific groups

`scavenge` uses idle nodes from other partitions (preempted)

- See cluster page for specifics

Nodes -N, Tasks -n, Cores -c

- The -N, -n, -c flags often cause confusion. They change the size of your allocation and its placement.
- By default, you get 1 task, having 1 core, running on 1 node
- -N *N* forces each task to be scheduled onto *N* nodes. Rarely useful.
- -n *n* specifies *n* tasks. This is normally only useful with MPI.
- -c *c* causes each task to have *c* cpus/cores. This is very useful for multithreaded programs

Controlling memory usage

- It's crucial to understand memory required by your program.
- Nodes (and thus memory) are often shared
- Jobs have (small) default memory limits that you can/should override
- Jobs exceeding their request are killed. Common errors: “bus error”, or “Exeeded step memory at some point”
- `/usr/bin/time -a cmd ...` will show memory usage (maxresident)
- `-mem=8g` requests 8GB per node
- `-mem-per-cpu=8g` requests 8GB per cpu. Multiplied by cores (`-c`)

To specify 8 cores on one node with 8 GB RAM for each core:

```
sbatch -c 8 --mem-per-cpu=8G t.sh
```

Controlling walltime

- Each job is assigned a maximum walltime
- If not specified, it will get the default walltime limit
- The job is killed if that is exceeded
- You can specify longer walltime to avoid the job being killed
- You can specify shorter walltime to get resources faster

To specify walltime limit of 2 days:

```
sbatch -t 2- t.sh
```

Modules

Software is setup with *module*

\$ module avail	<i>find modules</i>
\$ module avail name	<i>find particular module</i>
\$ module load name	<i>use a module</i>
\$ module list	<i>show loaded modules</i>
\$ module unload name	<i>unload a module</i>
\$ module purge	<i>unload all</i>
\$ module save collection	<i>save loaded modules as collection</i>
\$ module restore collection	<i>restore collection</i>
\$ module describe collection	<i>list modules in collection</i>

Example “module load” commands

```
module load Python/3.5.1-foss-2016b  
module load Perl  
module load Bowtie2/2.2.9-foss-2016a
```

You can ask for the default or specify a specific version:

```
module load Python  
module load Python/3.5.1-foss-2016b
```

Run your program/script

When you're finally ready to run your script, you may have some trouble determining the correct command line, especially if you want to pass arguments to the script. Here are some examples:

```
compute-20-1$ python compute.py input.dat
compute-20-1$ R --slave -f compute.R --args input.dat
compute-20-1$ matlab -nodisplay -nosplash -nojvm < compute.m
compute-20-1$ math -script compute.m
compute-20-1$ MathematicaScript -script compute.m input.dat
```

You often can get help from the command itself using:

```
compute-20-1$ matlab -help
compute-20-1$ python -h
compute-20-1$ R --help
```


Running graphical programs on compute nodes

Two different ways:

- X11 forwarding
 - easy setup
 - `ssh -Y` to cluster, then `qsub -Y/bsub -XF/srun -x11`
 - works fine for most applications
 - bogs down for very rich graphics
- Remote desktop (VNC)
 - more setup
 - allocate node, start VNC server there, connect via ssh tunnels
 - works very well for rich graphics

More information is here:

<http://research.computing.yale.edu/support/hpc/user-guide>

Cluster Filesystems

Each cluster has a number of different filesystems for you to use, with different rules and performance characteristics.

It is very important to understand the differences. Generally, each cluster will have:

Home : Backed up, small quota, for scripts, programs, documents, etc.

Scratch : Not backed up. Automatically purged. For temporary files.

Project : Not backed up. For longer term storage.

Local HD : /tmp For local scratch files.

RAMDISK : /dev/shm For local scratch files.

Storage@Yale : University-wide storage (active and archive).

Consider using local HD or ramdisk for intermediate files. Also consider avoiding files by using pipes.

For more info:

<http://research.computing.yale.edu/hpc/faq/io-tutorial>

Large datasets (e.g. Genomes)

- Please do not install your own copies of popular files (e.g. genome refs).
- We have a number of references installed, and can install others.
- For example, on Ruddle: `/home/bioinfo/genomes`
- If you don't find what you need, please ask us, and we will install them.

Wait, where is the Parallelism?

Sbatch can allocate multiple cores and nodes, but the script runs on one core on one node sequentially.

Simply allocating more nodes or cores DOES NOT make jobs faster.

How do we use multiple cores to increase speed?

Two classes of parallelism:

- Single job parallelized (somehow)
- Lots of independent sequential jobs

Some options:

- Submit many batch jobs simultaneously (not good)
- Use job arrays, or dSQ (much better)
- Submit a parallel version of your program (great if you have one)

dSQ (aka Dead Simple Queue)

- Useful when you have many similar, independent jobs to run
- Automatically schedules jobs onto a single PBS allocation

Advantages

- Handles startup, shutdown, errors
- Only one batch job to keep track of
- Keeps track of status of individual jobs
- More flexible than job arrays
- Automatically schedules jobs onto a single PBS allocation
- dSQ replaces a previous tool “SimpleQueue”

Using dSQ

- 1 Create file containing list of commands to run (jobs.txt)

```
prog arg1 arg2 -o job1.out  
prog arg1 arg2 -o job2.out  
...
```

- 2 Create launch script

```
module load dSQ  
dSQ --taskfile jobs.txt [slurm args] > run.sh
```

- 3 Submit launch script

```
sbatch run.sh
```

For more info, see <http://research.computing.yale.edu/support/hpc/user-guide/dead-simple-queue>

Parallel-enabled programs

- Many modern programs are able to use multiple cpus on one node.
- Typically specify something like `-p 20` or `-t 20` (see man page)
- You must allocate matching number of cores: `-c 20`

```
#!/bin/bash
#SBATCH -c 20

myapp -t 20 ...
```

MPI-enabled programs

- Some programs use MPI to run on many nodes
- Impressive speedups are possible
- MPI programs are compiled and run under MPI
- MPI must cooperate with Slurm

```
#!/bin/bash
#SBATCH -N 4 -n 20 --mem-per-cpu 4G

module load MPI/OpenMPI
mpirun ./mpiprogram ...
```


Best Practices

- Start Slowly
 - Run your program on a small dataset interactively
 - In another ssh, watch program with top. Track memory usage.
 - Check outputs
 - Only then, scale up dataset, convert to batch run
- Input/Output
 - Think about input and output files, number and size
 - Should you use local or ram filesystem?
- Memory
 - Use top or /usr/bin/time -a to monitor usage
 - Specify memory and walltime requirements
- Learn Linux! Almost all HPC is linux-based.
- Be considerate! Clusters are shared resources.
 - Don't run programs on the login nodes. Use a compute node.
 - Don't submit a huge number of jobs. Use simplequeue or job array.
 - Don't do heavy IO to /home.
 - Keep an eye on quotas. Don't fill filesystems.

Plug for scripting languages

- Learning basics of a scripting language is a great investment.
- Very useful for automating lots of day to day activities
 - Parsing data files
 - Converting file formats
 - Verifying collections of files
 - Creating processing pipelines
 - Summarizing, etc. etc.
- Python (strongly recommended)
- Bash
- Perl (if your lab uses it)
- R (if you do a lot of statistics or graphing)

To get help

- Send an email to: `hpc@yale.edu`
- Email me: `robert.bjornson@yale.edu`
- Read documentation at:
`http://research.computing.yale.edu/hpc-support`
- Come to office hours: `http://research.computing.yale.edu/support/office-hours-getting-help`

Resources

- This presentation: <https://github.com/ycrc/Intro-Bootcamp>
- Our other courses: Python, R, Git, Linux, etc:
<http://research.computing.yale.edu/training>
- Table of equivalent Batch Queueing commands:
<https://slurm.schedmd.com/rosetta.pdf>
- Linux: <http://www.ee.surrey.ac.uk/Teaching/Unix> or
<http://ryanstutorials.net>
- Slurm: <http://slurm.schedmd.com>