

IBG Research Computing Cheatsheet

Basic terminology

Slurm

This is the software that the RC clusters are built on. This is software that allows users to spread *jobs* across the individual computer *nodes* that make up the cluster. Slurm commands include `sbatch`, `sinfo`, `scontrol`, and `srun`.

Node

A single computer, which is a member of the cluster. Different nodes may have different duties, such as login nodes, file transfer nodes, and compile nodes. Different nodes may also have different characteristics, such as the number of processors and amount of memory.

Job

Something that is to be done on the cluster. For example, run an R script, or run a series of commands to extract certain subjects and SNPs, and then analyze the results.

Queue

Slurm puts *jobs* in a *queue* until they are ready to run.

Scheduler

The components of Slurm that decide which jobs from the queue to run, and when to run them.

Core

A core is the unit of computing used when requesting resources to be used by a *job*. In typical computing fashion, the term is *overloaded*, which means it has multiple definitions, and the correct one can only be determined by context. A computer contains one or more *CPU* (or processors) chips, which are distinct physical objects that sit in *sockets*. Each physical CPU, will contain multiple *cores*, each of which is also referred to as a CPU. So a single computer may have two CPU sockets, the CPU in each socket then each has 32 cores, for a total of 64 CPUs.

Partition

A set of nodes that have been grouped together in Slurm. Nodes within a partition usually share a common feature (each has a GPU), a purpose (nodes for testing), or an owner (on Blanca).

QOS (Quality of Service)

A set of QOS definitions is defined in Slurm. Every job will be associated with a particular QOS, which determines which partitions the job can run on, and also affects priority and resource limits.

Priority

Priority is a trait of a job that the scheduler considers when deciding what jobs from the queue to run next.

- Based on how long the job has been waiting
- The amount of resources the job needs, so the size of the job
- FairShare, which is based on how much total resources a user or project allocation has consumed
- Jobs do not strictly run in priority order—a “backfill” scheduler may run low priority small jobs ahead of higher priority large jobs, if the resources to run the large jobs are not yet available

Best practices

We all share resources—learning to effectively manage these resources not only increases the efficiency of your work but that of your colleagues. Here are a few principles to consider when allocating resources.

1. Memory and CPU usage should be proportional.

This principle is best demonstrated via example: the himem nodes `bnode0412- bnode0414` have 976gb of RAM and 64 logical cores each. For a job that requires 300gb of RAM, you'll want to request $300/976 * 64 \approx 20$ tasks (`--N 1 --mem 300gb --ntasks 20`; always request cpus in multiples of two, because of hyperthreading slurm will round up anyways). This helps maximize the amount of work we can accomplish simulatenously. E.g., if you were to submit 8 jobs requiring 10gb of RAM and 8 tasks each, you would be using all 64 cores (100% of cpu power) but only 80/976 gb of ram (8% of total memory), rendering the remaining RAM unusable by anyone. Total memory and cpus vary by node, so you should consider which nodes are best for a given task and use the appropriate constraints (see ‘View available nodes and their properties’).

2. Use the preemptable QOS when possible

Most of Blanca is idle at any given time—the preemptable QOS provides access to all nodes and has the dual benefits of increasing the amount of computing you can accomplish in a given time frame while freeing speciality hardware (himem nodes) for applications that can only node on those nodes. You should consider the preemptable QOS for any job that requires fewer than 170gb of ram and will complete in less than a day. Situations where you might NOT want/be able to use the preemptable QOS include: jobs that require more than day to finish, jobs that require more than 170gb of memory, debugging small instances of single jobs before submitting many preemptable jobs, small interactive sessions.

3. Prototype/debug before submitting large numbers of jobs

If you're unsure how much memory/how many cpus a particular task will require, make sure you figure it out before submitting many such jobs. Methods for checking memory use are provided in the next session. You can also directly ssh into any node where you have currently running jobs and use `top` or `htop` to examine cpu / memory use in real time. If you find that a job is taking longer than you expected, investigate and cancel if necessary rather than waiting for it to time out.

4. Big jobs are fragile

Big here means consisting of many small tasks or consisting of one giant task that could be split up into smaller tasks. To the extent that you are able, you should discretize large tasks into many small tasks. This makes your workflow more fault tolerant (if some step goes wrong, you don't have to rerun everything) and easier to efficiently allocate resources (e.g., if a job that would take four days can be split into multiple jobs that take less than one day, you can use the preemptable QOS; this will most likely finish faster while simultaneously freeing speciality resources—i.e. himem nodes—so they'll be available to your colleagues). Simple strategies for reducing job fragility include: - avoid using loops within jobs (consider job-arrays instead) - avoid grouping multiple tasks (e.g., GWAS and summary statistic clumping) in a single job - read documentation to determine if jobs can be split into discrete chunks and then do so (e.g. using `--make-grm-part` in GCTA can allow you to run several fast, low memory, preemptable jobs, often simultaneously, instead of one multi-day job that requires a himem node)

5. Ask for help

If you're having difficulty figuring out how to get your jobs to run efficiently (or run at all) don't hesitate to ask for help from your colleagues.

Clusters

Alpine

- RC Documentation on Alpine
- Nodes owned by the university
- Partitions
 - `amilan` is the default, for general purpose use
 - `atesting` short test jobs
 - `ainteractive` interactive jobs
 - `acompile` building and installing software

Blanca

- RC Documentation on Blanca
- Nodes purchased by individual PIs
- Idle nodes are available to all, but owners get priority access
- IBG owns 7 nodes
- Partitions
 - `blanca-ibg` priority access to IBG owned nodes
 - `blanca` access all nodes in cluster
- QOSes
 - preemptable
 - * uses the `blanca` partition

- * ideal for large batches of jobs
- * jobs may be “preempted” off a node if the node’s owner starts a job
- * scheduling priority isn’t influenced by Fairshare
- * **this should be the QOS used for most tasks on Blanca**
- blanca-ibg
 - * the **default** QOS on Blanca, but not necessarily the best to use
 - * use for
 - jobs longer than 24 hours
 - interactive jobs
 - high priority small jobs, and preemptible is busy or broken
 - * a limited resource with only 7 nodes, so avoid if not necessary

Monitoring and Adjusting Slurm and Job Activity

Monitoring running jobs

View your running jobs

```
squeue -u $USER
```

- Additional fields displayed using the `-l` flag
- Custom output described in man page. E.g., I find `squeue -o "%.12i %.18j %.16q %.8T %.10M %.12l %.24R"` easier to read

View running jobs on a particular qos

```
squeue -q <QOS>
```

Measure performance of completed jobs

The following flags can be combined to suit your needs

Display jobs completed since a particular date

```
sacct -S <MMDD>
```

Display jobs’ timing and memory usage

```
sacct -o 'jobid%20,jobname%16,state,elapsed,maxrss'
```

Modifying running and submitted jobs

The `scontrol` command is used to change options on jobs already in the queue. Full documentation available online. The general format of `scontrol` is `scontrol [options] [command] [arguments]`. The exact arguments will depend on which commands are run * arguments * `<job_list>` a comma separated list of job IDs or a jobname, for example 2652,2653,2655 * `<job_id>` a single job ID, for example 2652 * `jobname=<name>` the name for one or more jobs, for example `jobname=my.jobs` * `<specification>` description of the state of a job or other slurm entity, see the full documentation under the `show` command to `scontrol` * `<entity>` a reference to a slurm job, node, etc., for example `job=2652`, `node=c3cpu-c11-u5-1`, or `node` * `scontrol show <entity>` show information about the specified entity, for example * `scontrol show job=2522` to show details about the job * `scontrol show node` to show information about all nodes * `scontrol hold <job_list>` prevents a job in the queue from starting, for example `scontrol hold`

2652,2654 or `scontrol hold jobname=my.job` * `scontrol release <job_list>` releases a job under hold, and allows it to run * `scontrol requeue <job_list>` puts a running or recently-failed job back in the queue * `scontrol update <specification>` make changes to a queued job that has not yet started by a specification similar to the output of the `scontrol show` command, for example `scontrol update job=2522 timelimit=3:00:00` to change the timelimit of 2522 to 3 hours * run time * QOS * node exclusion list * array task limits * dependencies

View available nodes and their properties

```
sinfo --Node -o "%.12N %.16P %.11T %.4c %.13C %.8e /%.8m %.30f"
```

- NODELIST is the name of the node, can be specified when submitting via `-nodelist`
- PARTITION is the name of the partition. You likely will never need to use this unless you want a particular partition of the preemptable queue
- STATE - mixed means that some of the CPUS are in use (see below)
- CPUS (A/I/O/T) - on a given node: allocated/idle/other/total nodes
- FREE_MEM / MEMORY free memory / total memory in MB (divide by 1024 to convert to GB)
- AVAIL_FEATURES - architectures/instruction sets requestable via the `--constraint` flag

Example jobs

Here are some example jobs as templates-the actual resources requested should be adjusted to fit your needs.

Basics

Interactive jobs

```
sinteractive --qos preemptable --mem 64gb --ntasks 12 --time 2:00:00
```

or

```
sinteractive --qos blanca-ibg --mem 64gb --ntasks 16 --time 2:00:00 --constraint
```

Preemptable job

```
#!/bin/bash
#SBATCH --qos=preemptable
#SBATCH --mem=100gb
#SBATCH --time=<D-HH:MM> or <HH:MM:SS> # must be at most one day
#SBATCH --ntasks=24
#SBATCH --nodes=1
#SBATCH --array=1-22
#SBATCH -J <jobname>
#SBATCH -o <output dir>/<jobname>_%a
```

Job arrays

Job arrays provide a clean method for submitting collections of jobs and are often preferable to either i. using a loop to submit multiple jobs or ii. using a loop within a single job to accomplish multiple tasks. Note that the resources you request apply to each job individually,

not to the collection of jobs. E.g., if each of 10 tasks requires 50gb of memory, you only need to request 50gb, not 500gb.

Job array with a single numeric index

```
#!/bin/bash
#SBATCH --qos=<QOS>
#SBATCH --mem=<MEM>gb
#SBATCH --time=<D-HH:MM> or <HH:MM:SS>
#SBATCH --ntasks=<Number of cpus>
#SBATCH --nodes=1
#SBATCH --array=1-22
#SBATCH -J <jobname>
#SBATCH -o <output dir>/<jobname>_%a
```

```
CHR=${SLURM_ARRAY_TASK_ID}
```

```
ml load <modules>
```

```
program <args> \
    --input <inputprefix>_"$CHR" \
    --output <outputprefix>_"$CHR"
```

Job array with a single non-numeric index

This script submit jobs in parallel for jobs with different inputs

```
#!/bin/bash
...
#SBATCH --array=0-2
#SBATCH -J <jobname>
#SBATCH -o <output dir>/<jobname>_%a
```

```
ii=${SLURM_ARRAY_TASK_ID}
```

```
inputArray=(phenoA phenoB phenoC)
```

```
input=${inputArray[$ii]}
```

```
ml load <modules>
```

```
program <args> \
    --input $input
```

Job array with a multiple simultaneous indices

This script submit jobs in parallel for jobs with different inputs/outputs

```
#!/bin/bash
...
#SBATCH --array=0-2
#SBATCH -J <jobname>
#SBATCH -o <output dir>/<jobname>_%a
```

```

ii=${SLURM_ARRAY_TASK_ID}

inputArray=(phenoA phenoB phenoC)
outputArray=(outputA outputB outputC)

input=${inputArray[$ii]}
output=${outputArray[$ii]}

ml load <modules>

program <args> \
    --input $input \
    --output $output

```

Job array with a multiple nested non-numeric indices

This script submit jobs in parallel for jobs with arbitrary nested lists of arguments (e.g., each model for each phenotype) using integer arithmetic. If this is unfamiliar, you can google “floor division bash” and “modulo bash”.

```

#!/bin/bash
...
#SBATCH --array=0-5
#SBATCH -J <jobname>
#SBATCH -o <output dir>/<jobname>_%a

ii=${SLURM_ARRAY_TASK_ID}

inputArray=(phenoA phenoB phenoC)
modelArray=(modelA modelB)

inputIndex=$((expr $ii % 3))
modelIndex=$((expr $ii / 3))

input=${inputArray[$inputIndex]}
model=${modelArray[$modelIndex]}

ml load <modules>

program <args> \
    --phenotype "$input" \
    --model $model

```

You can always double check that you didn’t screw something up by running a simple loop (in the shell):

```

inputArray=(modelA modelB)
modelArray=(phenoA phenoB phenoC)

for ii in {0..5}
do
    inputArray=(phenoA phenoB phenoC)
    modelArray=(modelA modelB)
    inputIndex=$((expr $ii % 3))
    modelIndex=$((expr $ii / 3))

```

```

input=${inputArray[$inputIndex]}
model=${modelArray[$modelIndex]}
echo input:"$input" model:"$model"
done

```

Job array with a multiple nested indices, one numeric

You can frequently simplify things when one the lists you iterate over is numeric (e.g., each chromosome for each phenotype):

```

#!/bin/bash
...
#SBATCH --array=0-65
#SBATCH -J <jobname>
#SBATCH -o <output dir>/<jobname>_%a

ii=${SLURM_ARRAY_TASK_ID}

inputArray=(phenoA phenoB phenoC)

inputIndex=$(expr $ii / 22)
chrom=$(expr $(expr $ii % 22) + 1)

input=${inputArray[$inputIndex]}

ml load <modules>

program <args> \
    --phenotype "$input" \
    --genotypes <someprefix>_chr"$chrom"

```

Again, you can check that this works via a loop:

```

inputArray=(phenoA phenoB phenoC)

for ii in {0..65}
do
    inputIndex=$(expr $ii / 22)
    chrom=$(expr $(expr $ii % 22) + 1)
    input=${inputArray[$inputIndex]}
    echo chrom:"$chrom" pheno:"$input"
done

```

Defining custom functions

Commands that are lengthy to type in or that you frequently used can be turned into functions to save time. To do so, add something along the following lines to ~/.my.bashrc

```

function Sinfo() {
sinfo --Node -o "%.12N %.16P %.11T %.4c %.13C %.8e /%.8m %.30f"
}
export -f Sinfo

```

Then call `source ~/.my.bashrc` and you should be able to use `Sinfo` in place of the longer command above.

SSH connections to RC and X11 forwarding

SSH can be configured to use a single channel to connect to RC. This means that you only have to enter your password and use Duo once, and subsequent connections will reuse the existing channel.

Mac or Linux

On a Mac or Linux a persistent SSH is created by editing the `~/.ssh/config`, which is the config file in your home directory, and then in the `.ssh` directory. An entry needs to be added:

```
# These rules only apply for connections to login.rc.colorado.edu
Host login.rc.colorado.edu
# Setup a master ssh session the first time, and subsequent times
# use the existing master connection
ControlMaster auto
ControlPath ~/.ssh/%r@%h:%p
# Keep the ssh connection open, even when the last session closes
ControlPersist yes
# X forwarding. Remove this on a Mac if you don't want it to
# start X11 each time you connect to RC
ForwardX11 yes
# Compress the data stream.
Compression yes
# Send keep alive packets to prevent disconnects from
# CU wireless and behind NAT devices
ServerAliveInterval 60
```

Mac X11

XQuartz is an open-source X11 server for Macs. To install it, download and install the DMG file from the XQuartz project.

Windows

PuTTY

PuTTY can be configured to share a single SSH connection. The best way to do this is to create saved session for RC.

Open PuTTY, and put your username and RC login host into the Host Name field: `ralphie@login.rc.colorado.edu`. In the left menu, click on the SSH entry, and select the box at “Share SSH connections if possible”. As long as the original login window is open, additional logins can be created without re-entering your password.

If you would like to setup X11 forwarding, click the “+” at SSH, and then select the X11 entry. Select the box at “Enable X11 forwarding”. X11 forwarding requires a local X11 server to be installed on Windows.

Once the desired options have been selected, return to the “Session” entry in the left menu, and enter a name such as RC into the “Saved Sessions” field. Then click Save to save the session.

The saved session can be used by double clicking on it.

Bitvise

Bitvise has a freely available SSH client for Windows. Download and install the client software.

Once installed, run it and then click “New Profile” on the left, and give it a name, such as RC. By default it will save the profiles in your Documents folder.

Into the “Host” field enter `login.rc.colorado.edu`, and into the “Username” field put in your identikey username.

To setup X11 forwarding click on the “Terminal” tab at the top, and under “X11 Forwarding” select the box near “Enable”. This will require installing an X11 server.

Then click “Save profile” on the left.

Then click “Log in” at the bottom. That will bring up a prompt to enter your password, and then accept the Duo prompt on your phone.

That should bring up a terminal window and a file transfer window. To get more terminal windows, click on the “New terminal console” button on the left. Additional file transfer windows can be created by clicking the “New sftp window” button. It is not necessary to enter a password or use Duo for additional windows.

No additional settings are necessary to use a persistent connection, as Bitvise will automatically reuse the existing connection for future terminal and file transfer connections.

X11

Xming is available for free. On the Xming page go to the “Public Domain Releases” and then download and install the Xming and Xming-fonts packages.