



Introduction to Using the Terra Cluster



**HIGH PERFORMANCE
RESEARCH COMPUTING**
TEXAS A&M UNIVERSITY

HPRC Short Course



Outline

- Usage Policies
- Hardware Overview
- Accessing Terra
- File Transfers
- File Systems and User Directories
- Computing Environment
- Development Environment
- Batch Processing
- Common Problems
- Need Help?

Introduction

- Prerequisites:

- Basic knowledge of UNIX/Linux
- Slides from our UNIX/Linux short course are at:

https://hprc.tamu.edu/training/intro_unix.html

- Examples:

- Available in /scratch/training/Intro-to-terra directory
- Copy these files to your scratch directory!!!

```
cp -r /scratch/training/Intro-to-terra $SCRATCH/
```

Usage Policies

(Be a good compute citizen)

- It is illegal to share computer passwords and accounts by state law and university regulation
- It is prohibited to use Terra in any manner that violates the United States export control laws and regulations, EAR & ITAR
- Abide by the expressed or implied restrictions in using commercial software

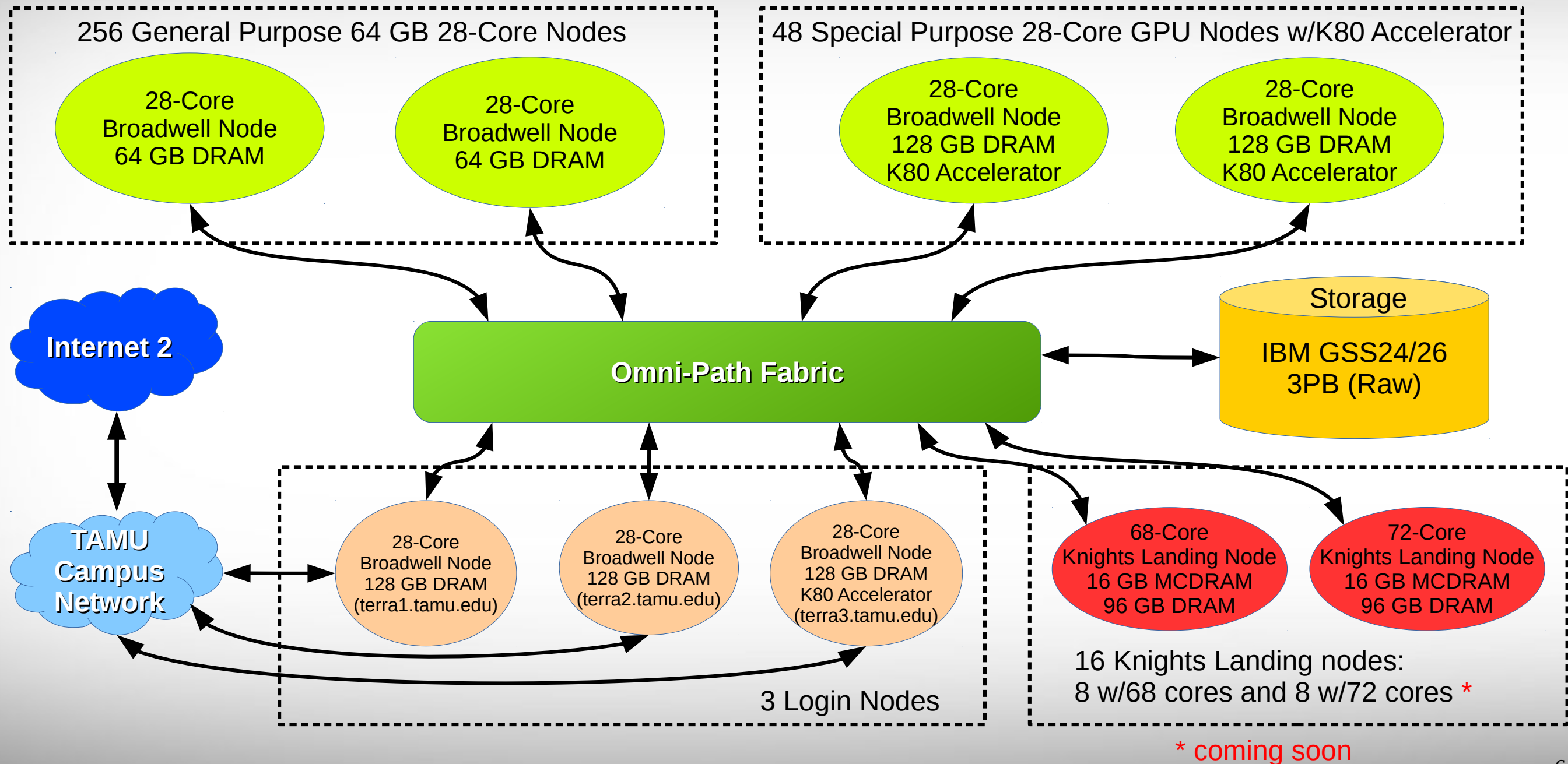
Terra – an x86 Cluster

A 9,716-core, 323-node cluster with:

- **256** 28-core compute nodes with two Intel 14-core 2.4GHz *Broadwell* processors and 64 GB of memory.
- **48** 28-core compute nodes with two Intel 14-core 2.4GHz *Broadwell* processors, 128 GB of memory, and one dual-GPU K80 accelerator.
- **8** compute nodes with one Intel Knights Landing 68-core 1.4 GHz processor, 16 GB of MCDRAM and 96 GB of memory. (coming soon)
- **8** compute nodes with one Intel Knights Landing 72-core 1.5 GHz processor, 16 GB of MCDRAM and 96 GB of memory. (coming soon)
- **3** 28-core login nodes with two Intel 14-core 2.4GHz *Broadwell* processors.
 - 1 login node has a dual-GPU K80 accelerator (terra3.tamu.edu).
- Nodes are interconnected with Omni-Path fabric in a two-level fat-tree topology.



Terra Schematic: 9,716-core, 323-node Cluster



Accessing Terra

- SSH is required for accessing Terra:
 - On campus: `ssh NetID@terra.tamu.edu`
 - Off campus:
 - Set up VPN: u.tamu.edu/VPnetwork
 - Then: `ssh NetID@terra.tamu.edu`
- SSH programs for Windows:
 - MobaXTerm (preferred, includes SSH and X11)
 - PuTTY SSH
- Terra has 3 login nodes. Check the bash prompt. `NetID@terra1 ~]$`
- Login sessions that are idle for 60 minutes will be closed automatically
- Processes run longer than 60 minutes on login nodes will be killed automatically.
- **Do not use more than 8 cores on the login nodes!**
- **Do not use the sudo command.** Contact us if you need help installing software.

<https://hprc.tamu.edu/wiki/index.php/HPRC:Access>

File Transfers with Terra

- Simple File Transfers:
 - scp: command line (Linux, MacOS)
 - rsync: command line (Linux, MacOS)
 - MobaXterm: GUI (Windows)
 - WinSCP: GUI (Windows)
 - FileZilla: GUI (Windows, MacOS, Linux)
- Bulk data transfers:
 - ***Will be available at later date via the login nodes.***
 - Recommended methods will likely be:
 - Globus Connect (<https://hprc.tamu.edu/wiki/index.php/SW:GlobusConnect>)
 - GridFTP

File Systems and User Directories

Directory	Environment Variable	Space Limit	File Limit	Intended Use
/home/\$USER	\$HOME	10 GB	10,000	Small to modest amounts of processing.
/scratch/user/\$USER	\$SCRATCH	1 TB	50,000	Temporary storage of large files for on-going computations. Not intended to be a long-term storage area.

- View usage and quota limits: the ***showquota*** command
- Also, only home directories are backed up daily.
- Quota and file limit increases will only be considered for scratch directories
- **Do not share your home/scratch directories.** Request a group directory for sharing files.

Computing Environment

- Paths:

Try "echo \$PATH"

- \$PATH: for commands (eg. /bin:/usr/bin:/usr/local/sbin:/usr/sbin:/home/netid/bin)
- \$LD_LIBRARY_PATH: for libraries
- Many applications, many versions, and many paths How do you manage all this software?!
- The solution: **module** (lmod)
 - Each version of an application, library, etc. is available as a module.
 - Module names have the format of package_name/version.

Application Modules

- Installed applications are available as modules which are available to all users (*except for restricted modules*)
- **module** commands
 - `module avail` #show all available modules
 - `module spider tool_name` #search all modules
 - `module key genomics` #search with keyword
 - `module load tool_name` #load a specific module
 - `module list` #list loaded modules
 - `module purge` #unload all loaded modules
 - `module load Python` #load the default version of a package
 - `module load Python/2.7.12-intel-2017A` #load a specific version (**recommended way**)
- It's a good habit to purge unused modules before loading new modules.
- **Avoid loading modules in your `.bashrc`**

Software

- Search module first:
 - *module avail*
 - *module spider software_name*
- Check Software wiki page (<https://hprc.tamu.edu/wiki/index.php/SW>) for instructions and examples
- License-restricted software: contact license owner for approval
- Contact us for software installation help/request

Development Environment - Toolchains

- Intel toolchain (eg. software stack) is recommended, which includes:
 - Intel C/C++/Fortran compilers
 - Intel Math Kernel Library
 - Intel MPI library
- Intel toolchain modules are named intel/version
- Recommended version intel/2017A
module load intel/2017A
- For applications that need gcc/g++, run *module spider GCC* to find available versions.

<https://hprc.tamu.edu/wiki/index.php/SW:Toolchains>

Terra

Modules and Toolchains

- Use modules based on the same toolchains in your job scripts

```
module load Python/2.7.12-intel-2017A  
module load OpenFOAM/2.4.0-intel-2017A  
module load Hypre/2.11.1-intel-2017A
```

- Avoid mixing modules from different tool chains in the same job script:

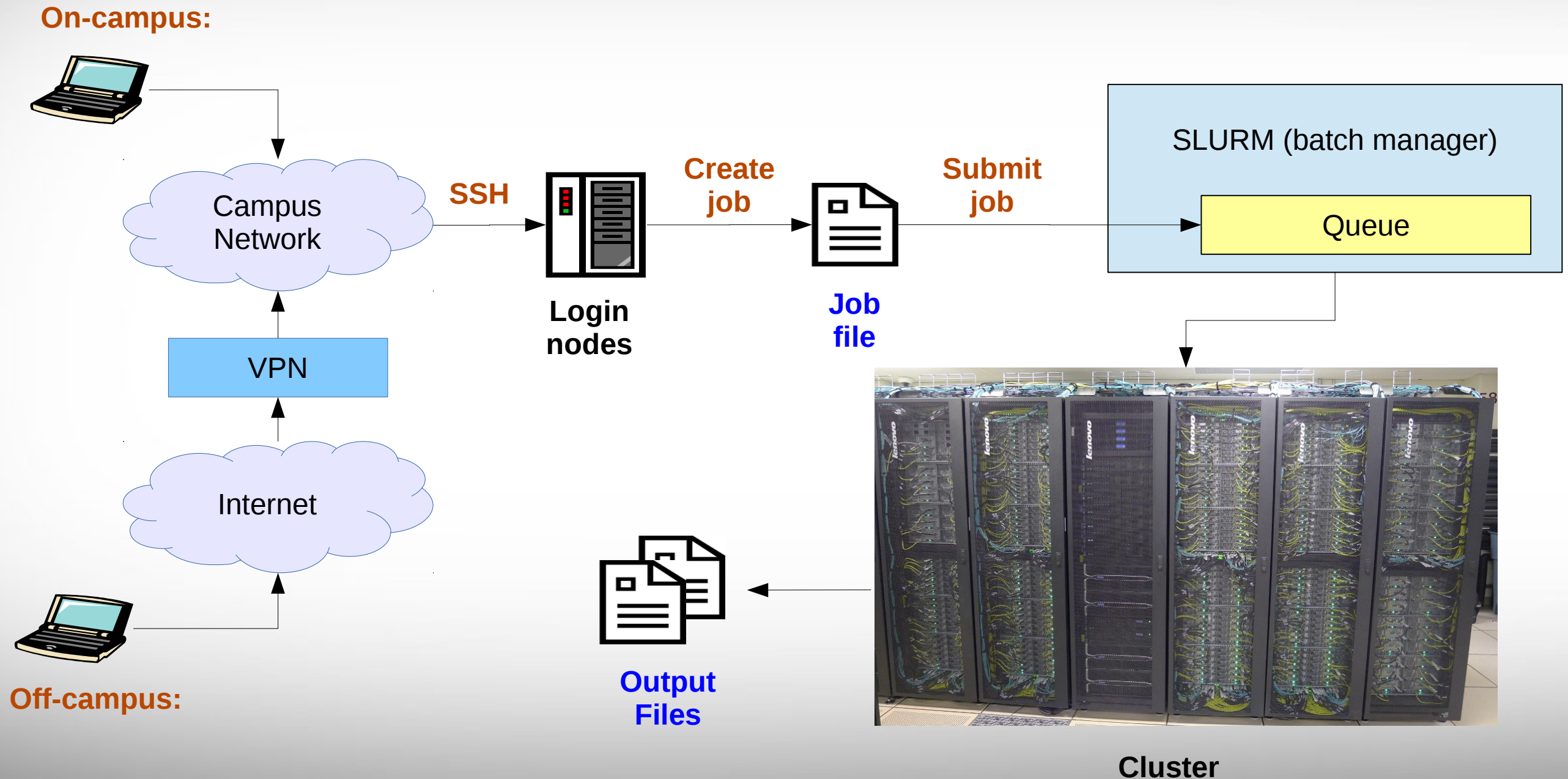
```
module load Python/2.7.12-intel-2017A  
module load OpenFOAM/2.4.0-intel-2016D  
module load Hypre/2.11.1-foss-2016D
```

- Same rule applies to compilers and libraries.

Development Environment: Compilers

- The commands to invoke each compiler are:
 - *icc* for C
 - *icpc* for C++
 - *ifort* for Fortran
- Man pages (documentation) are available for each compiler:
 - *man icc*
- Help for compiler options also available with *-help* option.
 - Also organized by categories (see *icc -help help* for more information).

Batch Computing on Terra



Batch Queues

- Job submissions are assigned to batch queues based on the resources requested (number of cores/nodes and wall-clock limit)
- Some jobs can be directly submitted to a queue:
 - If GPU nodes are needed, use the gpu queue
- Batch queue policies are used to manage the workload and may be adjusted periodically.

Current Queues

% sinfo

PARTITION	AVAIL	TIMELIMIT	JOB_SIZE	NODES(A/I/O/T)*	CPUS(A/I/O/T)*
short*	up	2:00:00	1-16	276/1/3/280	7714/42/84/7840
medium	up	1-00:00:00	1-64	276/1/3/280	7714/42/84/7840
long	up	7-00:00:00	1-32	253/0/3/256	7070/14/84/7168
gpu	up	2-00:00:00	1-48	39/9/0/48	1076/268/0/1344
vnc	up	12:00:00	1	39/9/0/48	1076/268/0/1344
staff	up	infinite	1-infinite	292/9/3/304	8146/282/84/8512
special	down	7-00:00:00	1-infinite	292/9/3/304	8146/282/84/8512

- For the NODES and CPUS columns:
 - A = Active (in use by running jobs)
 - I = Idle (available for jobs)
 - O = Offline (unavailable for jobs)
 - T = Total

Queue Limits

Queue	Job Max Cores / Nodes	Job Max Waltime	Compute Node Types	Per-User Limits Across Queues	Notes
short	448 cores / 16 nodes	30 min / 2 hrs	64 GB nodes (256) 128 GB nodes with GPUs (36)	1800 cores per user	
medium	1792 cores / 64 nodes	1 day			
long	896 cores / 32 nodes	7 days			
gpu	1344 cores / 48 nodes	2 days	128 GB nodes with GPUs (48)		For jobs requiring GPUs.
vnc	28 cores / 1 node	12 hours	128 GB nodes with GPUs (48)		For remote visualization jobs

Batch Queue Policies also at: <https://hprc.tamu.edu/wiki/index.php/Terra:Batch#Queues>

Consumable Computing Resources

- Resources specified in a job file:
 - Processor cores
 - Memory
 - Wall time
 - GPU
- Service Unit (SU) - Billing Account

Sample Job Script (structure)

```
#!/bin/bash ← This script will use the bash shell.

##ENVIRONMENT SETTINGS; CHANGE WITH CAUTION
#SBATCH --export=NONE
#SBATCH --get-user-env=L

##NECESSARY JOB SPECIFICATIONS
#SBATCH --job-name JobExample1
#SBATCH --time 01:30:00
#SBATCH --ntasks 2
#SBATCH --ntasks-per-node=2
#SBATCH --mem=2048M
#SBATCH --output Example1Out.%j

# this intel toolchain is just an example. recommended toolchain is TBD
module load intel/2016D ← Load the required module(s) first

# run program
./myprogram ← This is a command that is executed by the job
```

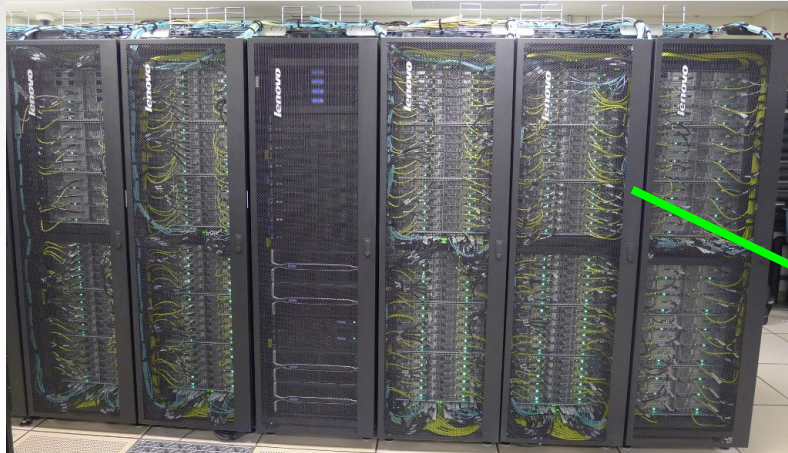
These lines (directives) describe your job to the job scheduler

This is a single line comment and not run as part of the script

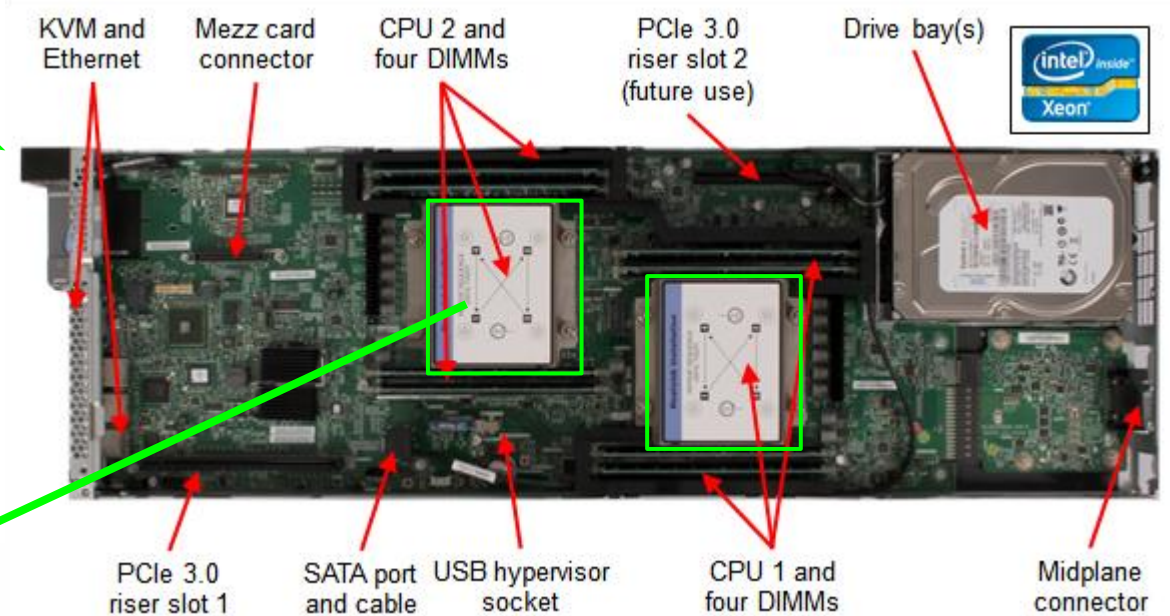
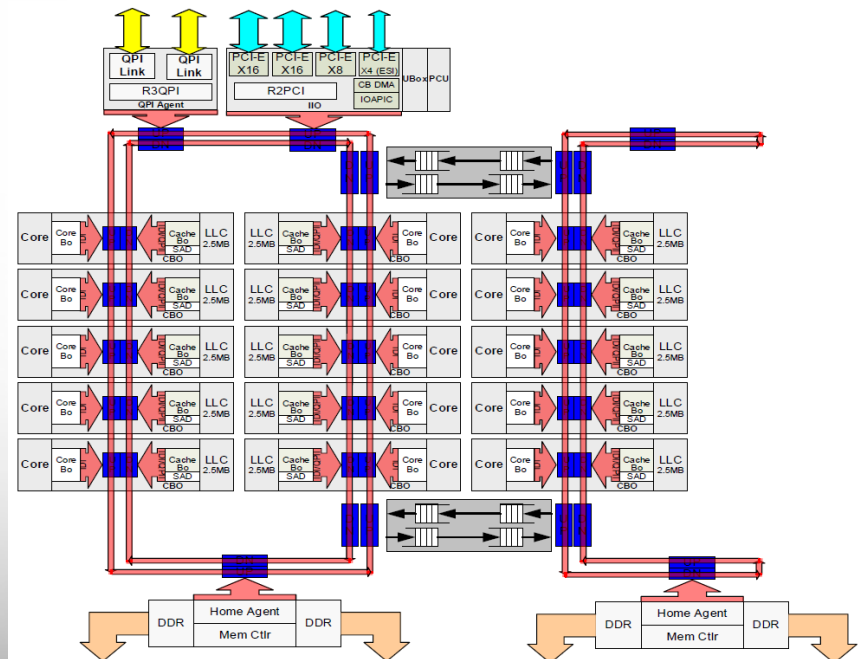
Important Job Parameters

#SBATCH --export=NONE #SBATCH --get-user-env=L	Initialize job environment.
#SBATCH --time HH:MM:SS	Specifies the time limit for the job.
#SBATCH --ntasks MM	Total number of tasks for the job.
#SBATCH --ntasks-per-node=NN	Specifies the maximum number of tasks to allocate per node
#SBATCH --mem=XXXXM	Sets the maximum amount of memory (MB) the job can use per node.

Compute Nodes

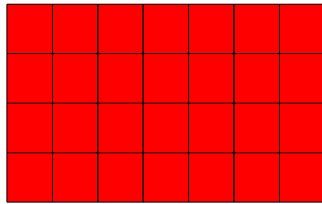


Part of Terra cluster.
Each green light is a node.



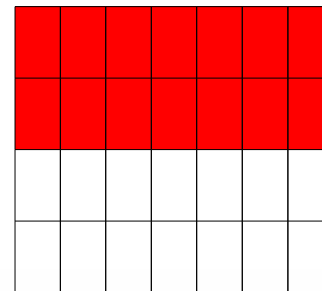
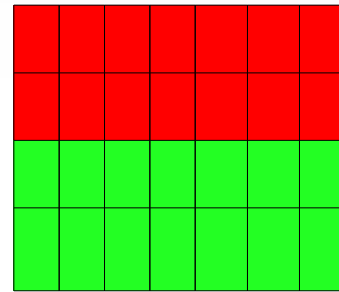
Each node has 28 processor cores.

Mapping Jobs to Nodes



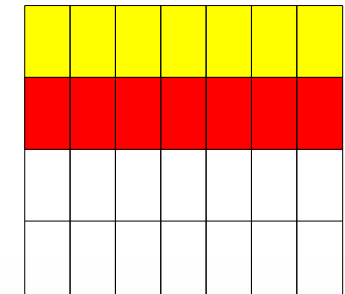
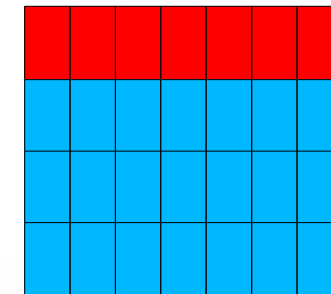
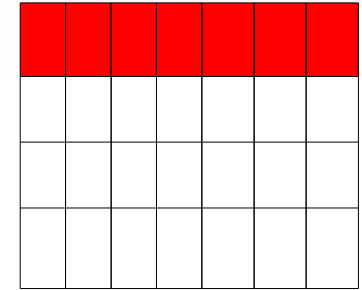
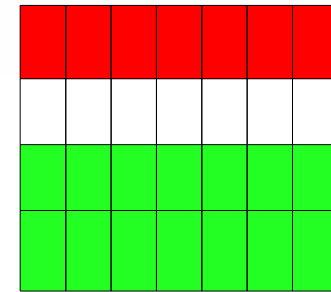
28 cores on
1 node

#SBATCH --ntasks 28
#SBATCH --tasks-per-node=28



28 cores on 2 nodes

#SBATCH --ntasks 28
#SBATCH --tasks-per-node=14



28 cores on 4 nodes

#SBATCH --ntasks 28
#SBATCH --tasks-per-node=7

Job Resource Examples (node vs memory)

Requests 8 tasks (2 per node). The job will span 4 nodes. The job can use up to 4 GB per node.

```
#SBATCH --ntasks=8
```

```
#SBATCH --tasks-per-node=2
```

```
#SBATCH --mem=4096M
```

Request 4 whole nodes (112 cores, 28 cores per node). The job can use up to 56 GB per node.

```
#SBATCH --ntasks=112
```

```
#SBATCH --tasks-per-node=28
```

```
#SBATCH --mem=57344M
```

Job Memory Requests

- Must use one of the following lines to request memory for your job:

#SBATCH --mem=XXXXM # memory per node in MB

#SBATCH --mem-per-cpu=XXXXM # memory per cpu in MB

- On 64GB nodes, usable memory is at most 56 GB. The per-process memory limit should not exceed 2048 MB for a 28-core job.
- On 128GB nodes, usable memory is at most 112 GB. The per-process memory limit should not exceed 4096 MB for a 28-core job.

Job File (Serial Example)

```
#!/bin/bash
##ENVIRONMENT SETTINGS; CHANGE WITH CAUTION
#SBATCH --export=NONE           #Do not propagate environment
#SBATCH --get-user-env=L       #Replicate login environment

##NECESSARY JOB SPECIFICATIONS
#SBATCH --job-name=JobExample1  #Set the job name to "JobExample1"
#SBATCH --time=01:30:00        #Set the wall clock limit to 1hr and 30min
#SBATCH --ntasks=1             #Request 1 task
#SBATCH --mem=2560M            #Request 2560MB (2.5GB) per node
#SBATCH --output=Example1Out.%j #Send stdout/err to "Example1Out.[jobID]"

##OPTIONAL JOB SPECIFICATIONS
#SBATCH --account=123456       #Set billing account to 123456
#SBATCH --mail-type=ALL        #Send email on all job events
#SBATCH --mail-user=email_address #Send all emails to email_address

# load Intel toolchain
module load intel/2017A

# run program
./myprogram
```

Job File (multi core, single node)

```
#!/bin/bash
##ENVIRONMENT SETTINGS; CHANGE WITH CAUTION
#SBATCH --export=NONE           #Do not propagate environment
#SBATCH --get-user-env=L       #Replicate login environment

##NECESSARY JOB SPECIFICATIONS
#SBATCH --job-name=JobExample2  #Set the job name to "JobExample2"
#SBATCH --time=6:30:00         #Set the wall clock limit to 6hr and 30min
#SBATCH --nodes=1              #Request 1 node
#SBATCH --ntasks-per-node=8    #Request 8 tasks/cores per node
#SBATCH --mem=8G               #Request 8GB per node
#SBATCH --output=Example2Out.%j #Send stdout/err to "Example2Out.[jobID]"

##OPTIONAL JOB SPECIFICATIONS
#SBATCH --account=123456       #Set billing account to 123456
#SBATCH --mail-type=ALL       #Send email on all job events
#SBATCH --mail-user=email_address #Send all emails to email_address

# load Intel toolchain
module load intel/2017A

# run program
./my_multicore_program
```

Job File (multi core, multi node)

```
#!/bin/bash
##ENVIRONMENT SETTINGS; CHANGE WITH CAUTION
#SBATCH --export=NONE           #Do not propagate environment
#SBATCH --get-user-env=L       #Replicate login environment

##NECESSARY JOB SPECIFICATIONS
#SBATCH --job-name=JobExample3  #Set the job name to "JobExample3"
#SBATCH --time=1-12:00:00      #Set the wall clock limit to 1 Day and 12hr
#SBATCH --ntasks=8             #Request 8 tasks
#SBATCH --ntasks-per-node=2    #Request 2 tasks/cores per node
#SBATCH --mem=4096M            #Request 4096MB (4GB) per node
#SBATCH --output=Example3Out.%j #Send stdout/err to "Example3Out.[jobID]"

##OPTIONAL JOB SPECIFICATIONS
#SBATCH --account=123456       #Set billing account to 123456
#SBATCH --mail-type=ALL        #Send email on all job events
#SBATCH --mail-user=email_address #Send all emails to email_address

# load Intel toolchain
module load intel/2017A

# run program with MPI
mpirun ./my_multicore_multinode_program
```

Job File (serial GPU)

```
#!/bin/bash
##ENVIRONMENT SETTINGS; CHANGE WITH CAUTION
#SBATCH --export=NONE           #Do not propagate environment
#SBATCH --get-user-env=L       #Replicate login environment

##NECESSARY JOB SPECIFICATIONS
#SBATCH --job-name=JobExample4  #Set the job name to "JobExample4"
#SBATCH --time=01:30:00        #Set the wall clock limit to 1hr and 30min
#SBATCH --ntasks=1             #Request 1 task
#SBATCH --mem=2560M            #Request 2560MB (2.5GB) per node
#SBATCH --output=Example4Out.%.#Send stdout/err to "Example4Out.[jobID]"
#SBATCH --gres=gpu:1           #Request 1 GPU
#SBATCH --partition=gpu        #Request the GPU partition/queue

##OPTIONAL JOB SPECIFICATIONS
#SBATCH --account=123456        #Set billing account to 123456
#SBATCH --mail-type=ALL        #Send email on all job events
#SBATCH --mail-user=email_address #Send all emails to email_address

# load Intel and CUDA toolchain
module load intel/2017A  CUDA

# run program
./my_gpu_program
```

OpenMP Jobs

- Must set ***OMP_NUM_THREADS*** to take advantage of the requested cores
- All processes run on the same node.

```
#!/bin/bash
##ENVIRONMENT SETTINGS; CHANGE WITH CAUTION
#SBATCH --export=NONE           #Do not propagate environment
#SBATCH --get-user-env=L        #Replicate login environment

##NECESSARY JOB SPECIFICATIONS
#SBATCH --job-name=JobExample5  #Set the job name to "JobExample2"
#SBATCH --time=6:30:00          #Set the wall clock limit to 6hr and 30min
#SBATCH --ntasks=1              #Request 1 task
#SBATCH --cpus-per-task=8       #Request 8 cpus/cores per task
#SBATCH --mem=8192M             #Request 8192MB (8GB) per node
#SBATCH --output=Example5Out.%j #Send stdout/err to "Example2Out.[jobID]"

# load Intel toolchain
module load intel/2017A

# set OpenMP number of threads to match job request
export OMP_NUM_THREADS=$SLURM_CPUS_PER_TASK

# run program
./my_multicore_program
```

MPI Jobs

- MPI programs may be run in batch jobs on multiple nodes
- Note, the mpirun command will know how many MPI tasks to launch from SLURM's node, task, and/or task per node directives.

```
#!/bin/bash
##ENVIRONMENT SETTINGS; CHANGE WITH CAUTION
#SBATCH --export=NONE           #Do not propagate environment
#SBATCH --get-user-env=L        #Replicate login environment

##NECESSARY JOB SPECIFICATIONS
#SBATCH --job-name=JobExample6   #Set the job name to "JobExample2"
#SBATCH --time=6:30:00          #Set the wall clock limit to 6hr and 30min
#SBATCH --ntasks=24             #Request 24 tasks
#SBATCH --ntasks-per-node=8     #Request 8 tasks/cores per node
#SBATCH --mem=8192M             #Request 8192MB (8GB) per node
#SBATCH --output=Example6Out.%j #Send stdout/err to "Example2Out.[jobID]"

# load Intel toolchain
module load intel/2017A

# run program with MPI
mpirun ./my_mpi_program
```


Submit the Job and Check Status

- Submit your job to the job scheduler

```
sbatch sample01.job
```

```
Submitted batch job 64152
(from job_submit) your job is charged as below
      Project Account: 122728110918
      Account Balance: 5000.000000
      Requested SUs:   10.38
```

- Summary of the status of your running/pending jobs

```
squeue -u $USER
```

```
% squeue -u $USER
JOBID  NAME      USER      PARTITION  NODES  CPUS  STATE  TIME  TIME_LEFT  START_TIME  REASON  NODELIST
64039  somejob   someuser   medium     4      112   PENDING  0:00  20:00      2017-01-30T21:00:4  Resources
64038  somejob   someuser   medium     4      112   RUNNING  2:49  17:11      2017-01-30T20:40:4  None     tnxt-[0401-0404]
```

Try yourself; copy examples: `cp -r /scratch/training/Intro-to-terra $SCRATCH/`

Job Submission and Tracking

Command	Description
<i>sbatch jobfile1</i>	Submit jobfile1 to batch system
<i>squeue [-u user_name] [-j job_id]</i>	List jobs
<i>scancel job_id</i>	Kill a job
<i>sacct -X -j job_id</i>	Show information for a job (can be running or finished)
<i>sacct -X -S YYYY-HH-MM</i>	Show information for all of your jobs since YYYY-HH-MM
<i>lnu job_id</i>	Show resource usage for a job

Node Utilization: *lnu*

lnu jobid # lists on stdout the CPU utilization and free memory across all nodes for an executing job.

Example:

```
% lnu 64033
JOBID  NAME      USER      PARTITION  NODES  CPUS  STATE  TIME  TIME_LEFT  START_TIME      REASON      NODELIST
64033  somejob  someuser  medium     4      112   RUNNING  4:42  15:18      2017-01-30T19:51:5  None      tnxt-[0401-0404]

HOSTNAMES  CPU_LOAD  FREE_MEM  MEMORY  CPUS(A/I/O/T)
tnxt-0401  24.17    36104    57344   28/0/0/28
tnxt-0402  25.78    33999    57344   28/0/0/28
tnxt-0403  26.29    36777    57344   28/0/0/28
tnxt-0404  25.36    36706    57344   28/0/0/28
```

Note: SLURM updates the node information every few minutes.

Job Environment Variables

- ***\$SLURM_JOBID*** = job id
- ***\$SLURM_SUBMIT_DIR*** = directory where job was submitted from
- ***\$SCRATCH*** = /scratch/user/NetID
- ***\$TMPDIR*** = /work/job.\$SLURM_JOBID
 - \$TMPDIR is local to each assigned compute node for the job
 - Local disk space is about 850GB
 - Use of \$TMPDIR is recommended for jobs that use many small temporary files

Check your Service Unit (SU) Balance

- Show the SU Balance of your Account(s)

```
myproject -l
```

```
=====
                        List of username's Project Accounts
-----
| Account      | Default | Allocation | Used & Pending SUs | Balance |
-----
|122728110918|        N| 50000.00 |          -10.38 | 49989.62 |
-----
```

- Use "**#SBATCH -A project_id**" to charge SUs to a specific project
- Run "**myproject -d accountNo**" to change default project account
- Run "**myproject -h**" to see more options

https://hprc.tamu.edu/wiki/index.php/HPRC:AMS:Service_Unit
<https://hprc.tamu.edu/wiki/index.php/HPRC:AMS:UI>

Job Submission Issues (SUs)

```
$ sbatch myjob
sbatch: error: (from job_submit) your account's balance is not sufficient to
submit your job
    Project Account: 123940134739
    Account Balance: 382.803877
    Requested SUs:   18218.666666667
```

- Insufficient SUs?
 - Ask PI to transfer SUs to you
 - Apply for more SUs (if you are eligible, as a PI or permanent researcher)

See Accounts section in FAQ: <https://hprc.tamu.edu/wiki/HPRC:CommonProblems#Accounts>

Debugging Job Failures

- Debug job failures using the stdout and stderr files

- `cat output.ex03.python_mem.2447336`

This job id was created by the parameter in your job script file

```
#SBATCH -o output.ex03.python_mem.%j
```

```
slurmstepd: error: Exceeded job memory limit at some point.
```

Make the necessary adjustments to SBATCH parameters in your job script and resubmit the job

Concurrent Program Execution in Jobs via Tamulauncher

- Useful for running many programs concurrently across multiple nodes within a job
- Can be used with serial or multi-threaded programs
- Distributes a set of commands from an input file to run on the cores assigned to a job
- Can only be used in batch jobs
- If a tamulauncher job gets killed, you can resubmit the same job to complete the unfinished commands in the input file
- Preferred over job arrays

<https://hprc.tamu.edu/wiki/index.php/Ada:Tamulauncher>

Common Job Problems

- Control characters (^M) in job files or data files edited with Windows editor
 - remove the ^M characters with: `dos2unix my_job_file`
- Did not load the required module(s)
- Insufficient walltime specified in #SBATCH -t parameter
- Insufficient memory specified in #SBATCH --mem or --mem-per-cpu parameters
- Memory specified is too large
- Running OpenMP jobs across nodes
- Insufficient SU: See your SU balance: `myproject -l`
- Insufficient disk or file quotas: check quota with `showquota`
- Using GUI-based software without setting up X11 forwarding
 - Enable X11 forwarding at login `ssh -X terra`
 - Or use VNC
- Software license availability: check license status with `license_status -s softwarename`

```
$ file jobfile.txt
jobfile.txt: ASCII text, with
CRLF line terminators
$ dos2unix jobfile.txt
dos2unix: converting file
jobfile.txt to UNIX format ...
$ file jobfile.txt
jobfile.txt: ASCII text
```

FAQ: <https://hprc.tamu.edu/wiki/index.php/HPRC:CommonProblems>

Need Help?

- Check the FAQ (<https://hprc.tamu.edu/wiki/index.php/HPRC:CommonProblems>) or the Terra User Guide (<https://hprc.tamu.edu/wiki/index.php/Terra>) for possible solutions first.
- Email your questions to help@hprc.tamu.edu. (Now managed by a ticketing system)
- Help us, help you -- we need more info
 - Which Cluster
 - UserID/NetID (*UIN is not needed!*)
 - Job id(s) if any
 - Location of your jobfile, input/output files
 - Application used if any
 - Module(s) loaded if any
 - Error messages
 - Steps you have taken, so we can reproduce the problem
- Or visit us @ 114A Henderson Hall
 - Making an appointment is recommended.



**HIGH PERFORMANCE
RESEARCH COMPUTING**
TEXAS A&M UNIVERSITY

Thank you.

Any questions?