

WELCOME TO CLUSTERS 101 WORKSHOP

Presenters:

Gladys Andino

Dan Dietz

Jieyu Gao

Steve Kelley

Xiao Zhu

Workshop material:

<https://www.rcac.purdue.edu/tutorials/clusters101>



ACKNOWLEDGEMENTS

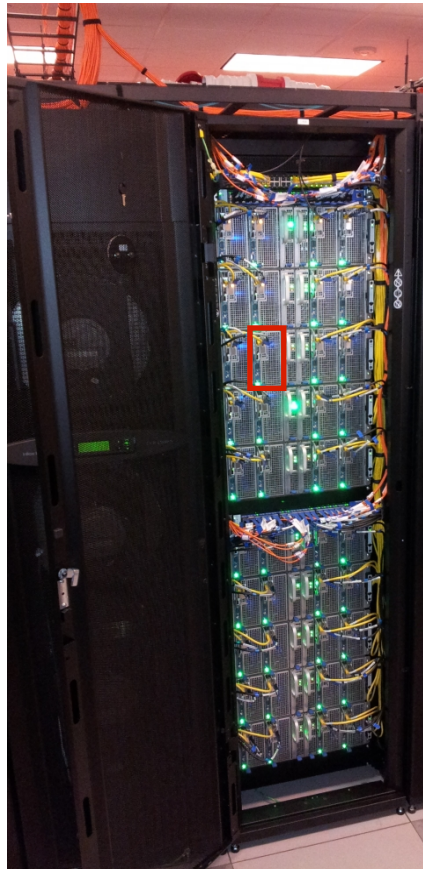
- The material in this workshop was prepared by the Purdue University ITaP Research Computing team.
- Special thanks to Dr. ***Boyu Zhang*** and ***Lev Gorenstein*** for their collaboration and slides preparation.

VOCABULARY

A CLUSTER IS A CLUSTER IS A CLUSTER?

- **Cluster:**

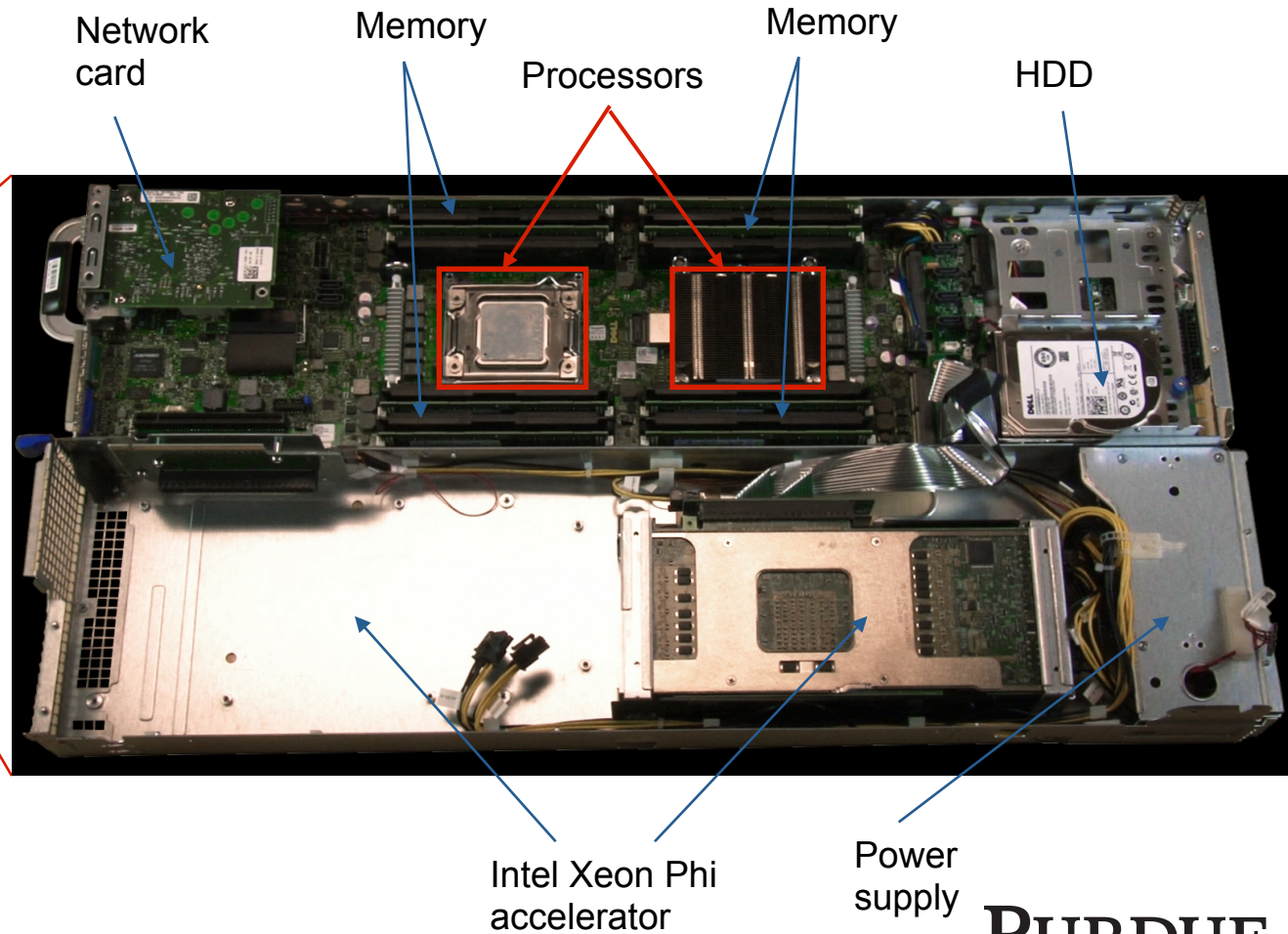
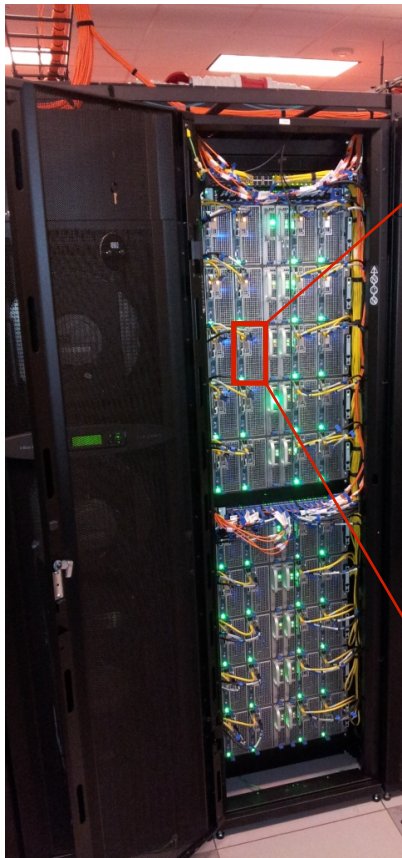
- Hardware (compute nodes + interconnect + storage)
- + Software (OS + compilers + libraries + apps + queue manager)
- + Infrastructure (front-ends + power + cooling + data center + staff)



VOCABULARY

NODES

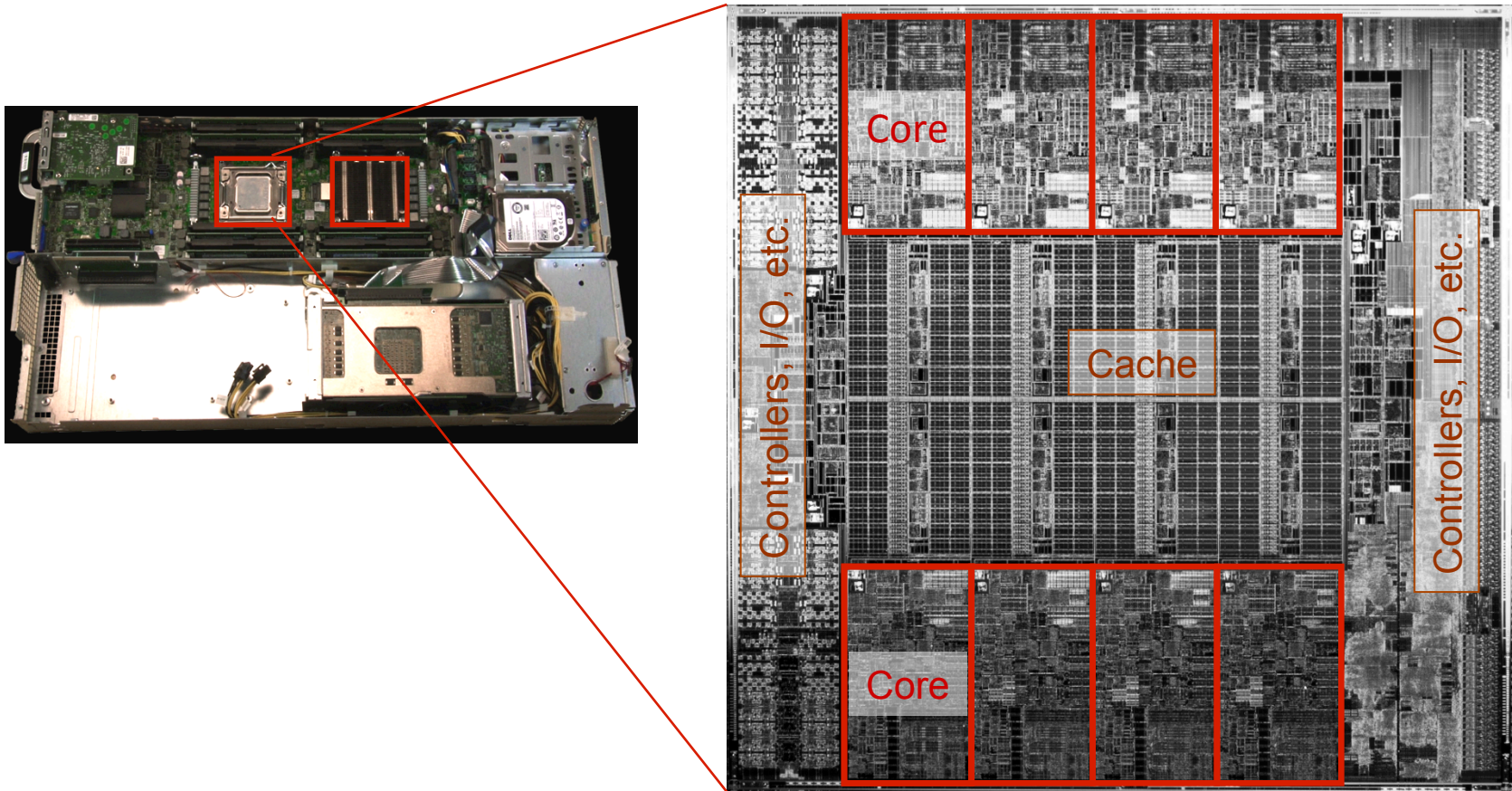
- **Node**: a computer



VOCABULARY

CPUS

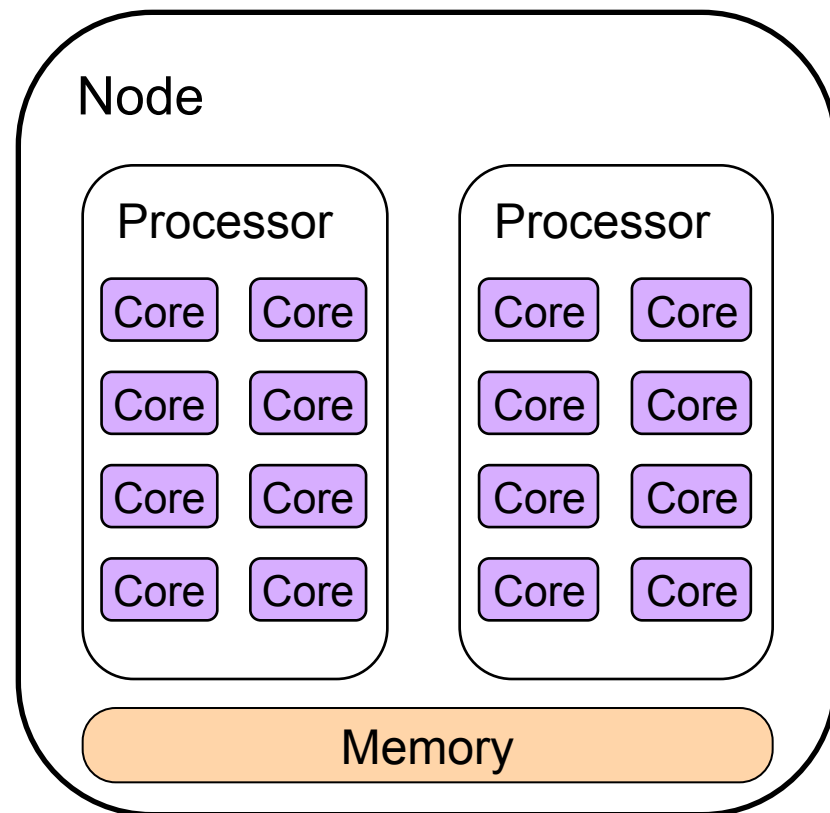
- **Processor**: a.k.a chip, a.k.a. physical CPU



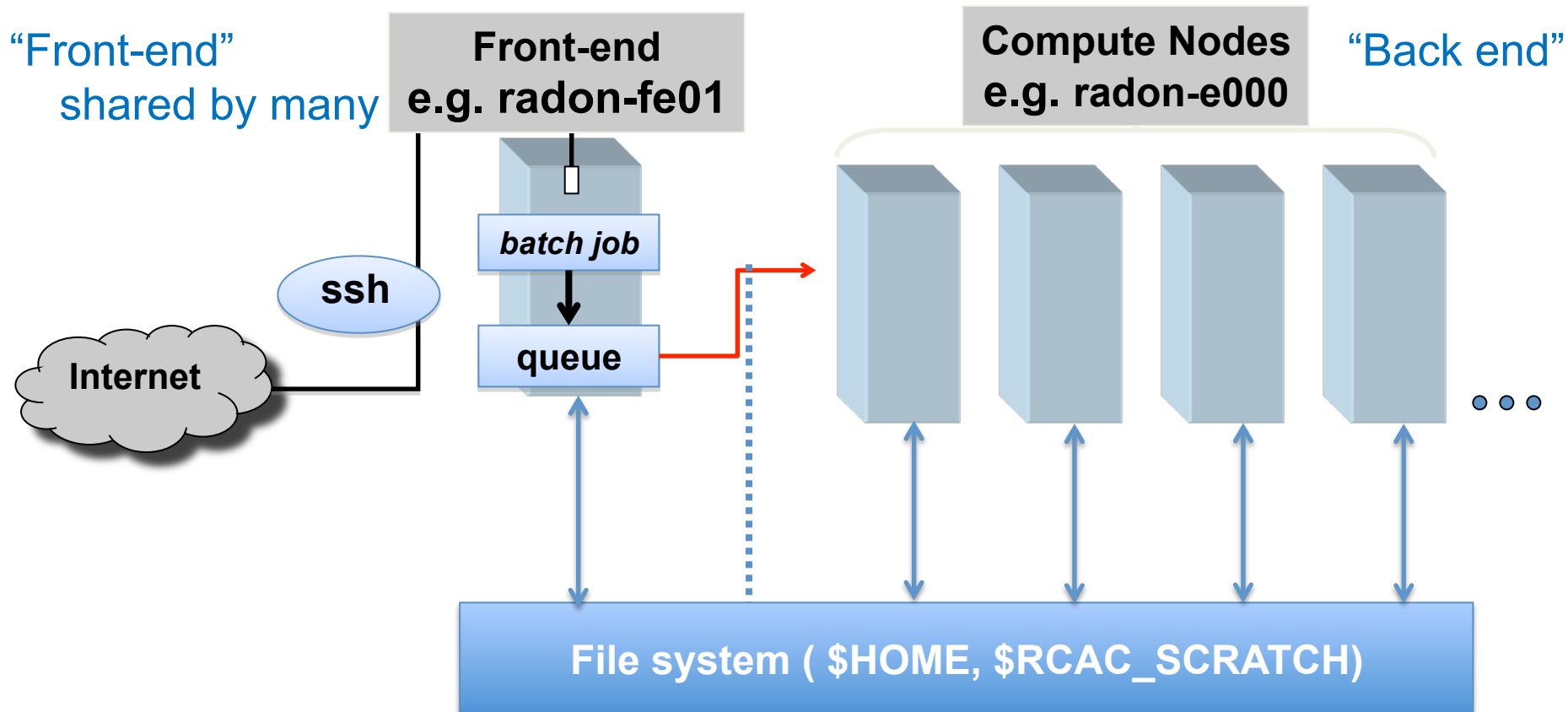
VOCABULARY

CORES

- **Processor Core**: individual compute unit (“slot”) on the chip
- You see:
 - 2 physical processors, 8 cores ea.
- Queuing system sees:
 - 16 logical processors
 - (hence 'nodes=X:ppn=16' effectively stands for “*processor cores per node*”).
- From now on, we will be mostly concerned with cores (logical processors), not physical chips
(“I ran on 5 CPUs” == “on 5 processors” == “on 5 cores”)



FRONT-END VS COMPUTE NODE



Running Jobs: The goal is getting to the compute nodes

VOCABULARY

IBM SYSTEM 360 (1 COMPUTER, 1 PROCESSOR)



TYPES OF NODES

COMMUNITY CLUSTERS

HPC (Rice): Multiple cores or nodes, probably MPI. Benefit from high-performance network and parallel filesystem. The vast majority of campus - 80% of all work!

Data-Intensive Life Science (Snyder): Use entire node to get large amounts of memory. Less need for high-performance network. Needs large, fast storage.

USER ENVIRONMENT



USER ENVIRONMENT

HELP WITH THE ENVIRONMENT AT ITAP RESEARCH COMPUTING

- Read the User Guide

Radon: <https://www.rcac.purdue.edu/compute/radon/guide/>

Rice: <https://www.rcac.purdue.edu/compute/rice/guide/>

Conte: <https://www.rcac.purdue.edu/compute/conte/guide/>

Carter: <https://www.rcac.purdue.edu/compute/carter/guide/>

- Submit a ticket (rcac-help@purdue.edu)

USER ENVIRONMENT

HELP WITH THE ENVIRONMENT AT ITAP RESEARCH COMPUTING

[Home](#) [News & Events](#) [User Info](#) [Accounts](#) [Computation](#) [Storage](#) [Visualization](#) [Purchase](#) [Services](#) [About](#)



Rice Info

- › [About Rice](#)
- › [Quick Reference Card](#)
- › [User Guide](#)
- › [Frequently Asked Questions](#)
- › [John R. Rice Bio](#)

Rice – User Guide

[Expand All](#)

- › **1 Conventions Used in this Document**
- › **2 Overview of Rice**
 - › **2.1 Namesake**
 - › **2.2 Detailed Hardware Specification**
 - › **2.3 Node Interconnect Systems**
 - › **2.3.1 FDR InfiniBand**
- › **3 Accounts on Rice**
 - › **3.1 Purchasing Nodes**
 - › **3.2 Cluster Partner Services**
 - › **3.3 Obtaining an Account**
 - › **3.4 Login / SSH**
 - › **3.4.1 SSH Client Software**
 - › **3.4.2 SSH Keys**
 - › **3.4.3 SSH X11 Forwarding**
 - › **3.4.4 Thinlinc Remote Desktop**
 - › **3.5 Passwords**
 - › **3.6 Email**
 - › **3.7 Login Shell**

LOGGING IN

WINDOWS

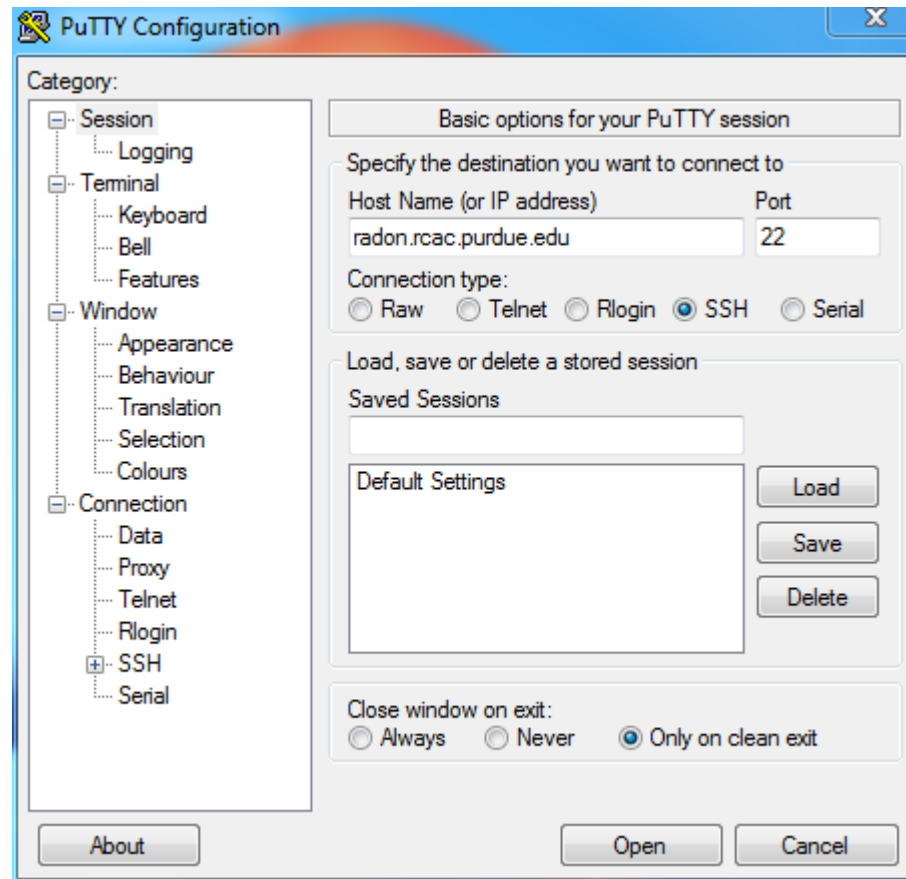
- Many clients are available for Windows
- We will use the PuTTY SSH client
- Download PuTTY, no install required
- <http://www.chiark.greenend.org.uk/~sgtatham/putty/download.html>
(or Google search “putty”)
- Download “putty.exe” for Intel x86 to your desktop



LOGGING IN

WINDOWS

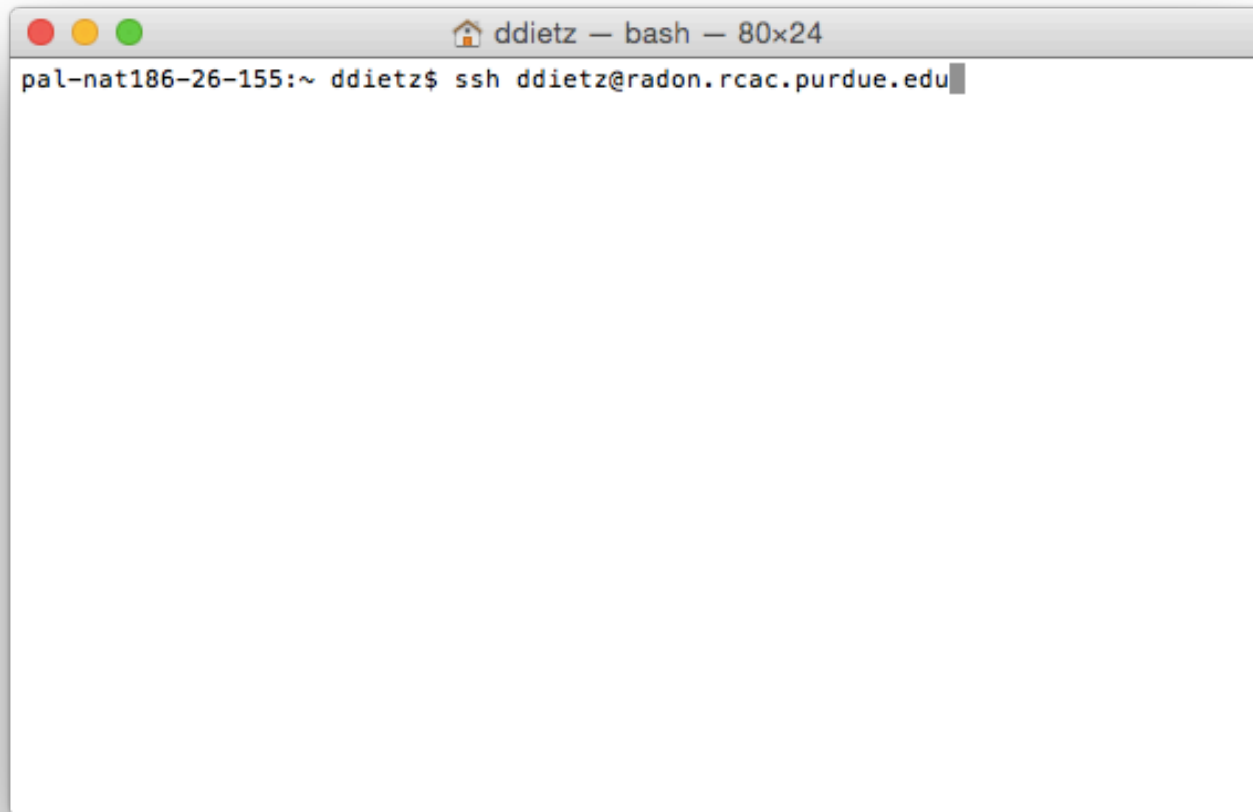
- Host Name for Radon is “radon.rcac.purdue.edu”



LOGGING IN

MAC

- Connect using `ssh username@radon.rcac.purdue.edu`



MODULES

WHAT IS A MODULE SYSTEM?

- All ITaP Research Computing clusters use a module system
- Module system provides for the dynamic modification of a user's environment
- Environment can be modified on a per-module basis using the module commands
- Set up the environment for running applications using one or two simple commands

```
$ module load bioinfo
```

```
$ module load bowtie
```


MODULES

RECOMMENDED ENVIRONMENT

- There are also modules for setting up one of several different programming environments. Most users will not need to change their programming environment
- Research Computing recommends a set of compilers, math libraries, and message-passing libraries for general optimal performance
- Recommended environment on each cluster contains different modules

```
$ module load devel
```

MODULES

KEY MODULE COMMANDS

\$ module avail

Running the command on the node will give you the most up to date list. The full list contains modules for compilers, MPI libraries, user applications and utilities.

\$ module avail python

“python” restricts the listing to the modules whose name starts with *python*.

```
$ module avail python
```

```
----- /opt/modules/modulefiles -----  
python/2.6.5                python/2.7.8_intel-15.0.2.164  
python/2.7.2                python/2.7.8_intel-15.0.3.187  
python/2.7.5_intel-13.1.1.163 python/3.4.1  
python/2.7.6_intel-14.0.2.144 python/anaconda(default)  
python/2.7.8_intel-15.0.1.133
```

MODULES

KEY MODULE COMMANDS

\$ module list	# list loaded modules
\$ module load <module>	# add a module
\$ module unload <module>	# remove a module
\$ module purge	# unload all modules

```
$ module list
No Modulefiles Currently Loaded.
```

```
$ module load devel
$ module load matlab/R2013a
```

```
$ module list
Currently Loaded Modulefiles:
1) intel/13.1.1.163    2) openmpi/1.6.3_intel-13.1.1.163    3) devel/20130708    4) matlab/R2013a
```

```
$ module unload matlab
$ module list
Currently Loaded Modulefiles:
1) intel/13.1.1.163    2) openmpi/1.6.3_intel-13.1.1.163    3) devel/20130708
```

```
$ module purge
$ module list
No Modulefiles Currently Loaded.
```

MODULES

WHAT DOES “MODULE LOAD” DO?

By loading a module, you will have just changed the environmental variables defined for that module. Typically this is \$PATH, \$MANPATH, and \$LD_LIBRARY_PATH. Also a few extra module-specific environmental variables could be set.

```
$ module load vmd
```

loads its module and adds to PATH and LD_LIBRARY_PATH

```
PATH=/apps/rhel6/vmd-1.9.1/bin:$PATH
```

```
LD_LIBRARY_PATH=/apps/rhel6/vmd-1.9.1/lib:$LD_LIBRARY_PATH
```

and sets

```
VMD_HOME=/apps/rhel6/vmd-1.9
```

```
LIBGL_ALWAYS_INDIRECT=y
```

```
$ ls -l $VMD_HOME/bin
total 182K
-rwxr-xr-x 1 vvergara 12K Jul 27 2012 vmd
-rw-r--r-- 1 vvergara 12K Feb 1 2012 vmd.csh
-rwxr-xr-x 1 vvergara 12K Feb 1 2012 vmd.sh
```

```
$ module show <module>    # show the commands in the module file
```

MODULES

QUICK SUMMARY OF MODULE COMMAND

\$ module help	#lists options
\$ module load <module>	#add a module
\$ module avail	#lists available modules
\$ module unload <module>	#remove a module
\$ module swap <mod1> <mod2>	#swap two modules
\$ module purge	#remove all modules
\$ module show <module>	#show module commands
\$ module help <module>	#module-specific help

MODULES

TRANSITION TO LMOD

- We recently switched to Lmod on our community clusters
- Lmod provides all of the functionality that the current module system does
- Supports software hierarchy
- Has nice features like “spider” and “save”

<code>\$ module spider</code>	<code># {lists all modules}</code>
<code>\$ module spider petsc</code>	<code>{list all versions of petsc}</code>
<code>\$ module save</code>	<code>{save personal environment}</code>
<code>\$ ml</code>	<code>{abbrev for module list}</code>
<code>\$ ml <module></code>	<code>{abbrev for module load}</code>
<code>\$ ml -<module></code>	<code>{return to system defaults}</code>
<code>\$ module reset</code>	<code>{return to system defaults}</code>

MODULES

TIPS FROM US

- Not load in `.bashrc/.profile/.login/.cshrc`
- Not load more than needed for current task
- Do `'module purge'` and reload when in doubt
- Check with `'module show ...'` what a given module does
- Check with `'module list'` what's actually loaded
- Script and automate your builds – even if it's just a 3-line snippet!

LAB 1 – USING MODULES

- List current loaded modules
`$ module list`
- Loaded recommended modules on the cluster
`$ module load devel`
- List available modules
`$ module avail`
- Search available modules
`$ module avail xmgrace`
`$ module avail x`
- List the help information of one module
`$ module help vmd` #or any modules you are interested
- Show detailed of one module
`$ module show vmd`
- Load extra modules as you like
- Remove all loaded modules and reset environment
`$ module purge`

LAB 1 – USING MODULES

- Find if module `cufflinks` available on Radon.
- Search all modules containing "`mpi`" in their name.
- Search all modules that end with "`mpi`" in their name.
- Load module `trinity/2.2.0`
- Please finish the following sequences of module operations:
 1. clear all modules
 2. load default version `matlab`,
`python/2.7.8_intel-16.0.2.181`, `r/3.1.0` and `spark/2.0.0`
 3. save the current module setting in the name of `datajob`
 4. clear all modules and restore the all modules saved in `datajob`

TYPES OF WORK

SERIAL VS PARALLEL

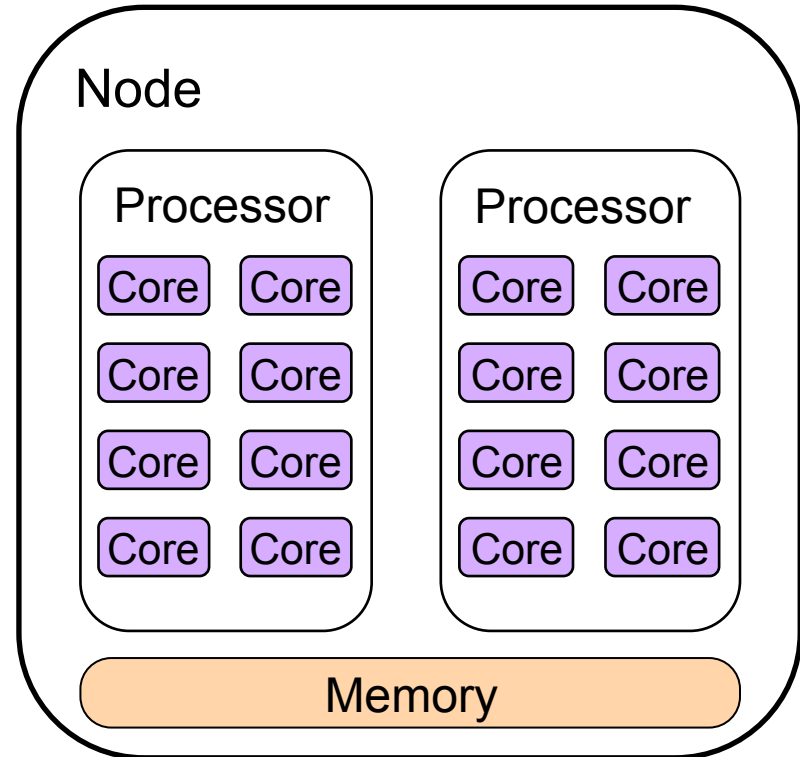


More than one mower!

Mow the lawn faster ...

... or mow a bigger lawn

Lawn mowers = cores



TYPES OF WORK

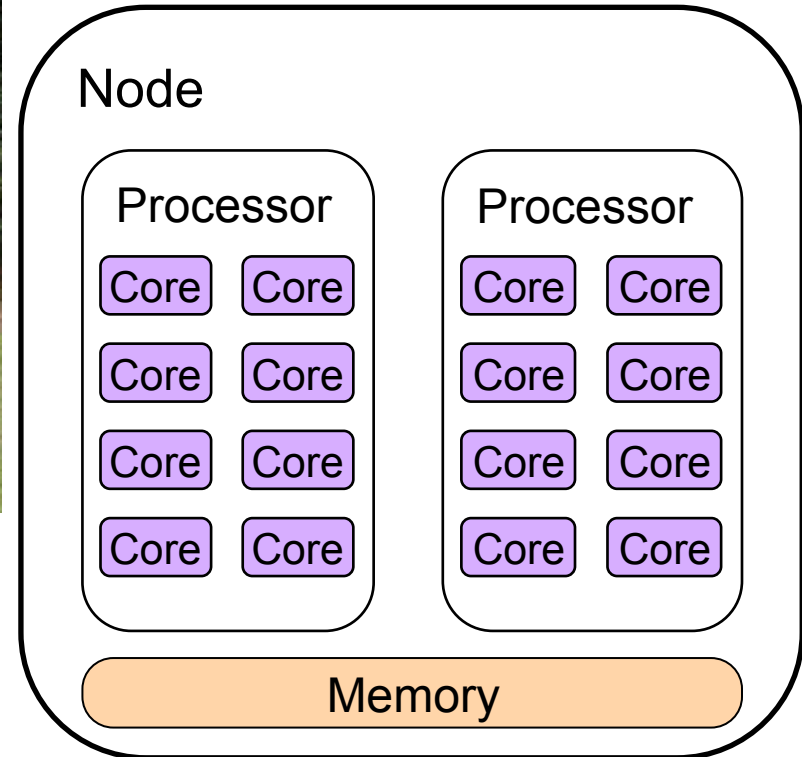
SERIAL VS PARALLEL

More than one mower!

Mow the lawn faster ...

... or mow a bigger lawn

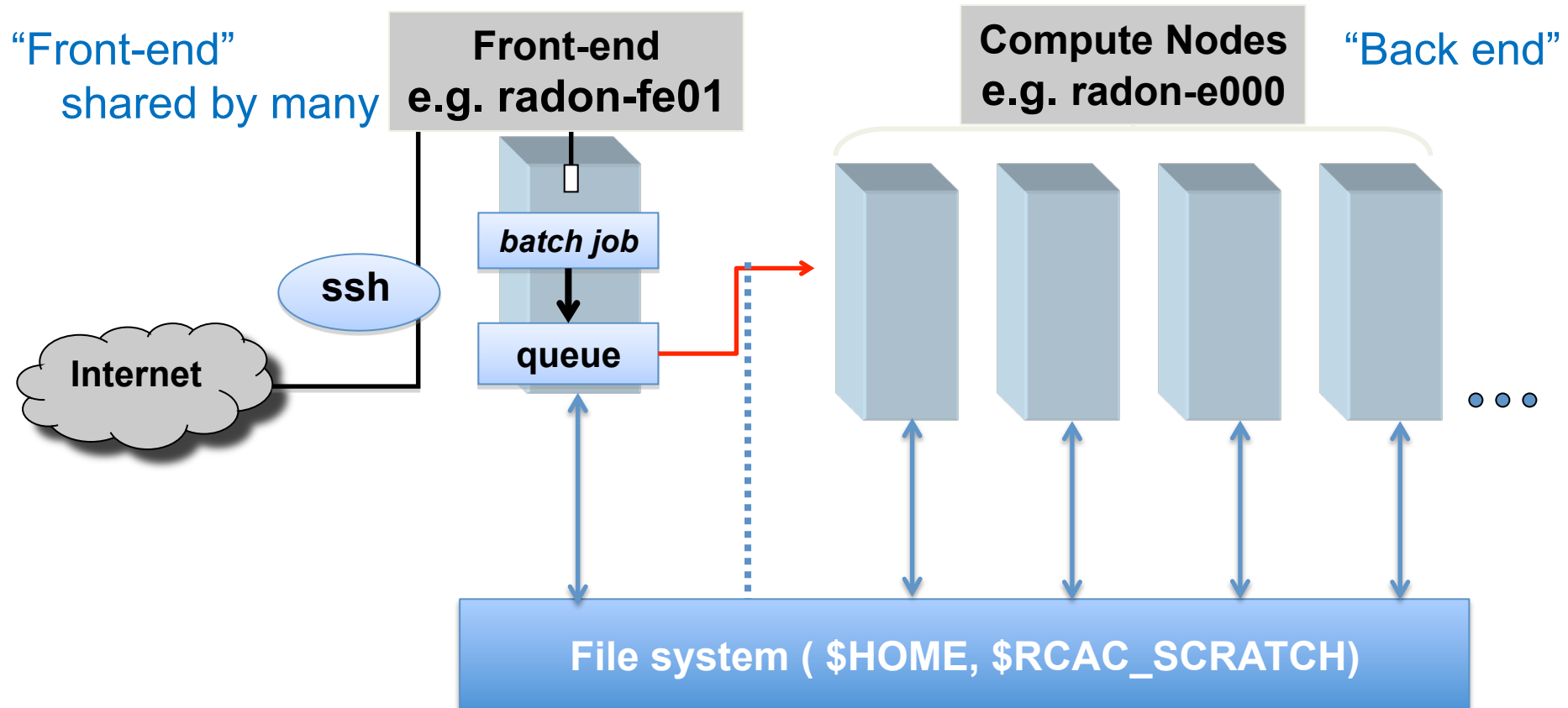
Lawn mowers = cores



PARALLEL COMPUTATIONS

- MPI – Message Passing Interface, multiple nodes
- OpenMP – confine yourself to one node
- Hybrid – mix MPI and OpenMP
- *Accelerators/coprocessors – GPU, Intel Xeon Phi*
- *Parallel R, Matlab, Hadoop & Spark*

FRONT-END VS COMPUTE NODE



Running Jobs: The goal is getting to the compute nodes

FRONT-ENDS: INAPPROPRIATE USE

- What does this mean?
- Do not do science on the front end
 - Either it's simply not possible
 - Or you'll annoy the system administrators and other users
- Instead: submit a batch job. We'll get to that.

FRONT-ENDS: APPROPRIATE USE

- Building software
 - are compilers also visible on compute nodes?
- Managing files
 - editing, transfers, `tar`, `gzip`, `hsi`
- Submitting, monitoring, managing batch jobs
- Launching interactive sessions
- Modest post-processing

THE KEYS TO GOOD CITIZENSHIP

- Remember you are sharing resources
 - login nodes, file systems, bandwidth
- Use components for intended purposes

JOBS, QUEUES, AND PBS



JOBS, QUEUES, AND PBS

WHAT IS A JOB?

- A job is simply a set of tasks to be performed by a cluster
- Instructs cluster precisely what to do to complete your work
- Self contained to be executed without any interaction
- Submit and forget!



JOBS, QUEUES, AND PBS

WHERE NOT TO RUN A JOB

- Remember, cluster front-end nodes are shared resources for
 - Creating, submitting, and monitoring jobs
 - File transfers
 - Preparing inputs
 - Editing and compiling code
 - Small-scale testing
- May be used by 50+ people simultaneously
 - Check out the **who** command
 - Jobs should have no interference by other people
 - Want jobs carefully arranged on compute nodes



[https://en.wikipedia.org/wiki/Rapid_transit#/media/File:Moscow_MetroCrowded_\(pixinn.net\).jpg](https://en.wikipedia.org/wiki/Rapid_transit#/media/File:Moscow_MetroCrowded_(pixinn.net).jpg)

JOBS, QUEUES, AND PBS

WHERE TO RUN A JOB

- Jobs are submitted to the cluster
- Cluster executes jobs on back-end compute nodes
- Jobs are carefully scheduled and arranged on the compute nodes



https://en.wikipedia.org/wiki/SNCF_TGV_Duplex#/media/File:TGV_Dupex_First_Class.jpg

JOBS, QUEUES, AND PBS

RESOURCE ALLOCATION

- Jobs need to specify the resources they require
 - Three basic units:
 - Number of nodes
 - Number of cores
 - Time
 - Memory
 - Other resources
- Cluster will allocate requested resources once they are available
- Job starts once resources are allocated

JOBS, QUEUES, AND PBS

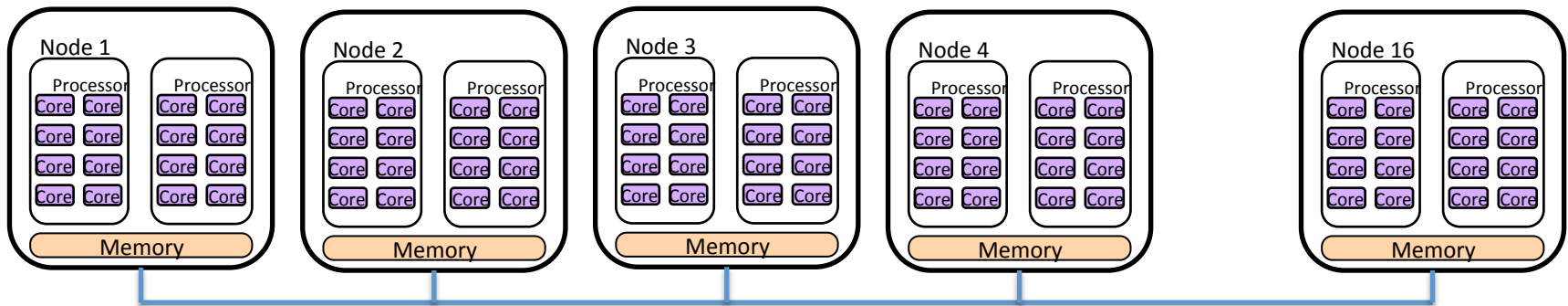
RESOURCE ALLOCATION – NODES & CORES

- The most basic unit you can request is nodes
 - Serial applications probably only need one node
 - Parallel applications may need 10s or 100s of nodes
- Your job should also request how many cores per node it needs
 - Cluster allows your job to use only the cores you request
 - Often your job should request all of the cores in a node
 - Remember, we want to avoid other people!

JOBS, QUEUES, AND PBS

RESOURCE ALLOCATION – NODES & CORES

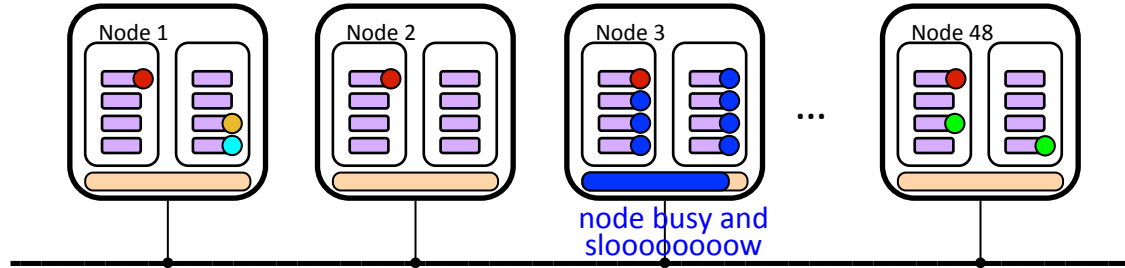
- Nodes are physically different pieces of hardware
- A job with 16 nodes and 16 cores each, is **not** 256 threads
 - 16 individual computers!
 - Program needs to be written with MPI to tie all the nodes together



JOB, QUEUES, AND PBS

RESOURCE ALLOCATION – NODES & CORES

- Nature of the job: how much inter-process communications?
- 1 core on 48 nodes, or for 48 cores on 1 node?
 - Communication inside a node is faster than between nodes
 - Collisions with other people's jobs ($48 \text{ nodes} \times 47 \text{ jobs} = 2256 \text{ jobs}$)
 - Most performant, most predictable, least vulnerable, least interfering



- Tighter packing is almost always preferred
- Newer ITaP clusters enforce exclusive nodes unless explicitly specified for single node jobs

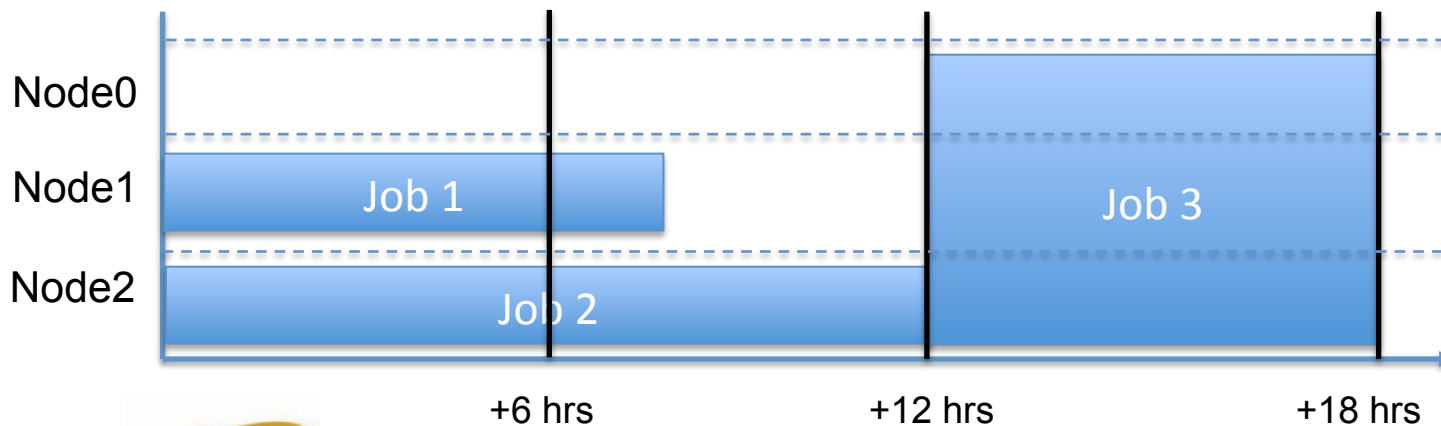
JOB, QUEUES, AND PBS

RESOURCE ALLOCATION – TIME

- The third basic unit jobs can request is walltime
- The cluster won't let your job run forever!
 - Cluster needs to have a time limit for your job
- Helps the cluster efficiently schedule you on the nodes
- Make request accurate with some safety buffer
 - Tune your walltime request through testing and experimenting



https://en.wikipedia.org/wiki/Hourglass#/media/File:Wooden_hourglass_3.jpg



JOB, QUEUES, AND PBS

RESOURCE ALLOCATION – MEMORY

- Jobs need memory, but...
- Clusters can **not** limit your job to an amount memory
 - All of a compute node's memory is available to your job
- All the memory is available to any other jobs on the node!
 - 1 core job needs 16 GB x 16 jobs = 256 GB
 - But wait, my node only has 32 GB!!
 - Jobs crash and burn

JOB, QUEUES, AND PBS

RESOURCE ALLOCATION – MEMORY

- Memory across multiple compute nodes is distributed
- Memory is discrete chunks of memory in different address spaces
- Can't just “add more nodes” to get more memory
- 16 nodes with 16GB memory each is **not** 256 GB of memory
 - 16 discrete chunks of 16GB memory
 - Can't malloc() 256GB of memory
 - Need MPI to tie memory together



JOB, QUEUES, AND PBS

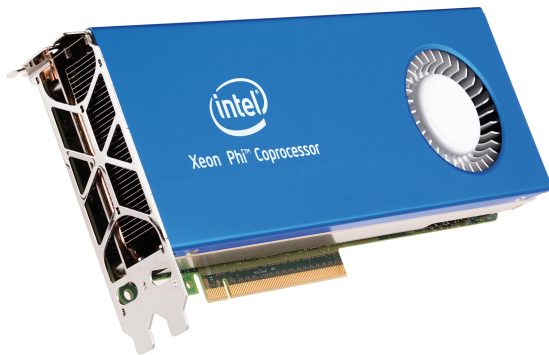
RESOURCE ALLOCATION – MEMORY

- Best bet is to always request all cores of a node
- Get full control of the node and its memory
- You can try requesting proportionate numbers of cores but **not** a guarantee!!
- You can always run multiple instances in a single job
 - Maximize utilization of the node
 - Only have to worry about yourself

JOBS, QUEUES, AND PBS

RESOURCE ALLOCATION – OTHER RESOURCES

- There are a number of other resources you can request in a job
 - Accelerators
 - Specific types of nodes
 - Certain licenses
 - MATLAB & Toolboxes



JOBS, QUEUES, AND PBS

WHAT IS A QUEUE

- Cluster can't run all jobs at once so sometimes you must wait
- Jobs are submitted into queues
- Jobs wait until cluster and queue has free resources
- Queues set constraints on jobs that can be submitted
 - Max nodes
 - Max walltime
 - Sets initial priorities
- Our community clusters have several



https://en.wikipedia.org/wiki/Queue_area#/media/File:Waiting_in_line_at_a_food_store.JPG

JOBS, QUEUES, AND PBS

TYPES OF QUEUES

- Owner queues
 - Named for your PI, lab, or group
 - Number of nodes set by what your PI buys into
 - Allows for longest jobs - 336 hours
 - Jobs should start in 4 hours or less
(assuming your queue isn't full of other jobs in queue)
- Not tied to specific nodes in the cluster
 - Allow you to use the number of nodes somewhere in the cluster
 - Whichever nodes are free, functioning, and type queue can access

JOBS, QUEUES, AND PBS

TYPES OF QUEUES

- Standby queue
 - Uses idle nodes from owner queues
 - Everyone on cluster gets access
 - Limited to 4 hours
 - Owners must be able get quick access from their queue
 - No limit on nodes
 - Probably won't get whole cluster though
 - Lowest priority jobs, no promises on turnaround time
 - Can be minutes or days
 - If you can run under 4 hours, go for it!

JOBS, QUEUES, AND PBS

TYPES OF QUEUES

- Debug queue
 - Allows for running small jobs for testing and debugging
 - Up to 2 nodes for 30 minutes at a time
 - Highest priority
 - Starting time <1 minute

9/9

0800 Antan started
1000 " stopped - antan ✓ { 1.2700 9.037 847 025
13'00 (032) MP-MC 1.982 1.000 9.037 846 845 connect
033 PRO 2 2.130476415 4.615925059(-2)
connect 2.130676415
Relays 6-2 in 033 failed special speed test
in relay 10,000 test.

1100 Relays changed
Started Cosine Tape (Sine check)
1525 Started Multi-Adder Test.

1545 Relay #70 Panel F
(moth) in relay.

1630 Antan started.
1700 closed down.

First actual case of bug being found.

Relay 2145
Relay 3370

JOBS, QUEUES, AND PBS

TYPES OF QUEUES

- In this workshop we are using Radon, a general access cluster
- Has none of the above queues!
- workq
 - Sort of like standby
 - Not competing with owner queues
 - Longer jobs
- For this workshop we have set up a special workshop queue
 - Jobs given priority and have several nodes reserved

JOBS, QUEUES, AND PBS

WHAT IS PBS?

- Portable Batch System
- Original developed for NASA in 1991
- Open source fork and commercial fork remain
- We use the open source fork, TORQUE

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

JOBS, QUEUES, AND PBS

WHAT DOES PBS DO?

- Resource manager
 - Jobs
 - Job submission
 - Job status
 - Executes jobs
 - Queues
 - Nodes

JOBS, QUEUES, AND PBS

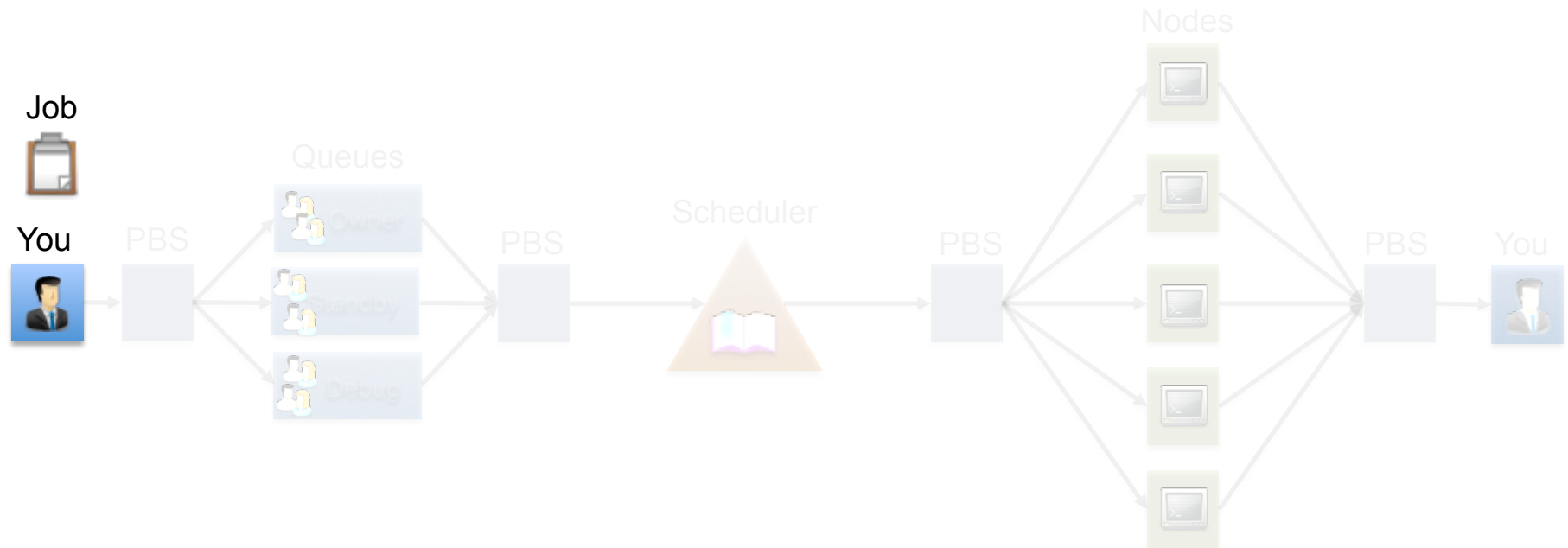
SCHEDULER

- Scheduler is actually separate software
- We use a scheduler called Moab
- Decides which jobs to run and where
- Typically “PBS” == “Scheduler” == “Batch system” == “Server”
 - We often interchange them in conversation
 - Referring to the whole system of servers, schedulers, managers

JOB, QUEUES, AND PBS

LIFE AS A JOB

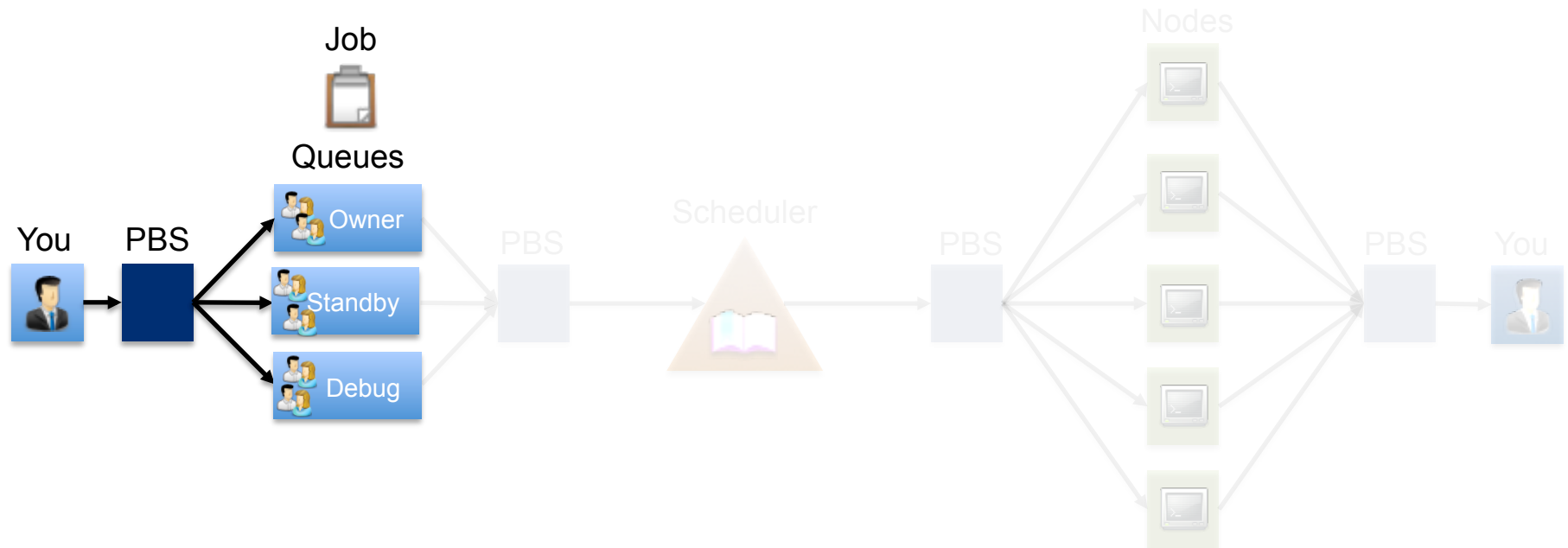
- You prepare a job submission script



JOBS, QUEUES, AND PBS

LIFE AS A JOB

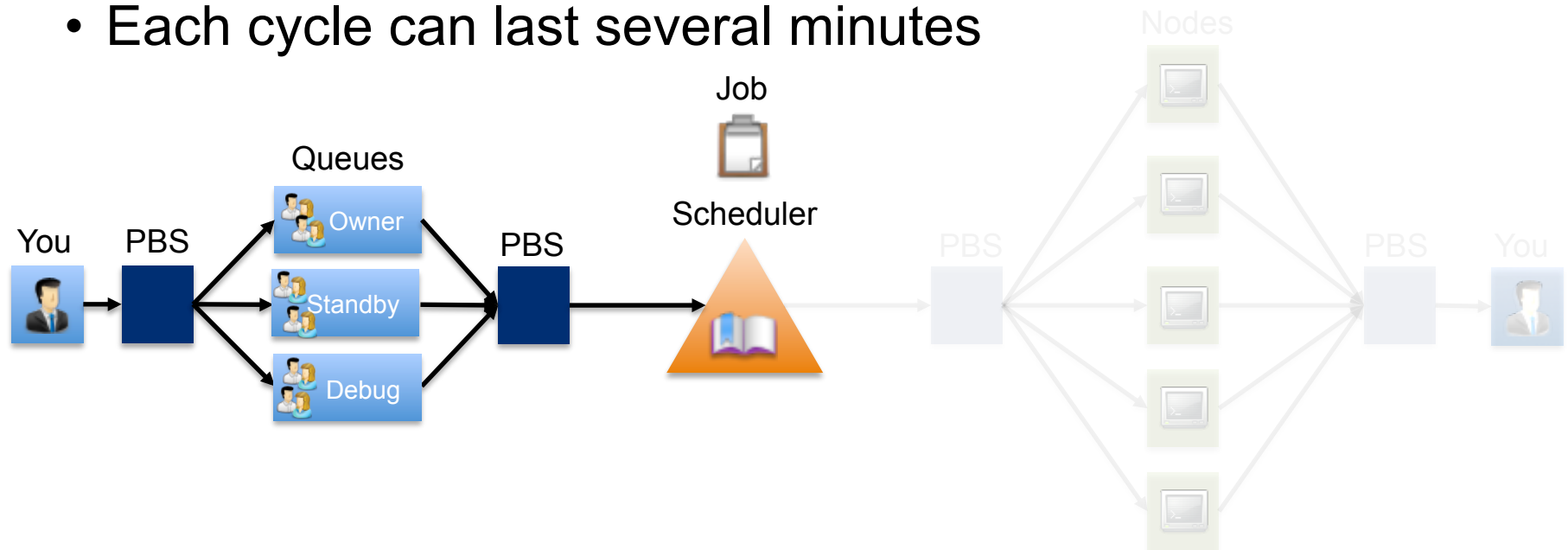
- You submit job script into PBS
- Job is placed in queue



JOBS, QUEUES, AND PBS

LIFE AS A JOB

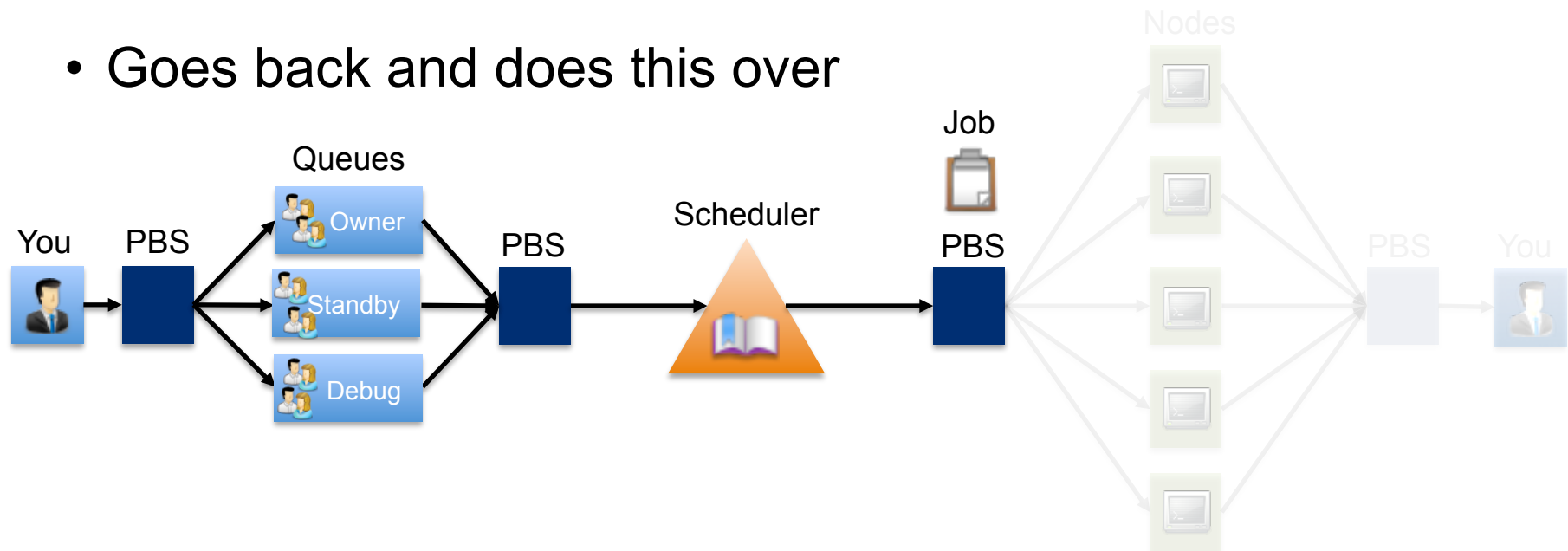
- Scheduler iteratively asks PBS for new jobs, status of old jobs, and status of nodes
- Each cycle can last several minutes



JOBS, QUEUES, AND PBS

LIFE AS A JOB

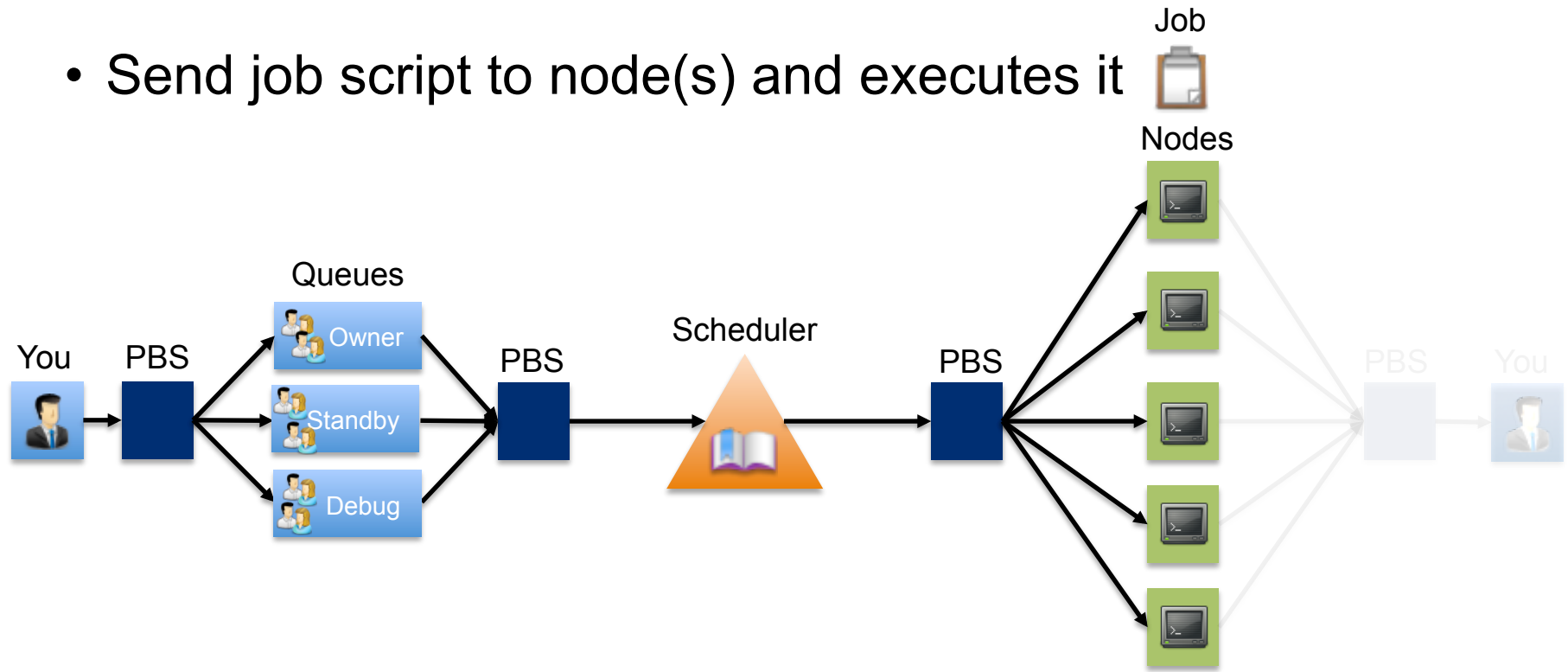
- Scheduler considers jobs, acts as a handler
- Tells PBS to start jobs one at a time
- Goes back and does this over



JOBS, QUEUES, AND PBS

LIFE AS A JOB

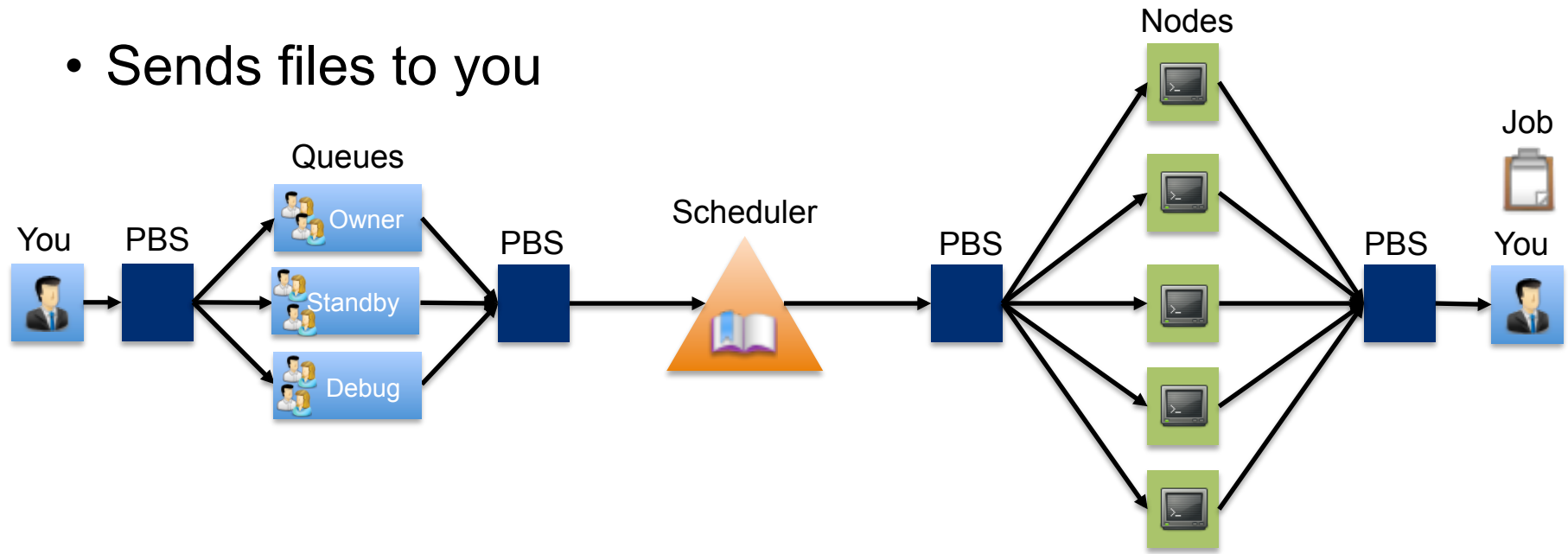
- PBS checks health of the node
- Send job script to node(s) and executes it



JOBS, QUEUES, AND PBS

LIFE AS A JOB

- PBS watches for job script to complete on nodes
- Collects job output files from nodes
- Sends files to you



HOW TO CREATE A SUBMISSION SCRIPT



HOW TO CREATE A SUBMISSION SCRIPT

COMPONENTS OF A SUBMISSION FILE

- PBS directives
 - Specify resources needed such as number of nodes, cores
- Module load
 - Set up paths, libraries
- PBS environment variables
 - Set by PBS, can be used in your submission script
- Customized commands
 - Your job to run

HOW TO CREATE A SUBMISSION SCRIPT

DIRECTIVES TO PBS

- A way to set PBS job attributes
- Appear at the top of your submission file
- Common PBS job attributes include:

PBS directives	Description
#PBS -l nodes=2:ppn=8	Number of nodes and cores required
#PBS -q workshop	The destination queue of the job
#PBS -N test	The name of the job
#PBS -l walltime=00:10:00	The estimated maximum walltime of the job, the job will be killed beyond this walltime

HOW TO CREATE A SUBMISSION SCRIPT

LOAD MODULES

- Unload all currently loaded modules that may be accidentally loaded from sourced files such as `~/ .bashrc`
- Load on the compute node(s) any relevant modules

```
module purge  
module load matlab/R2014b  
module load impi/5.1.1.109
```

HOW TO CREATE A SUBMISSION SCRIPT

CUSTOMIZED COMMANDS

- Specify everything needed for your job to run
- A simple job to print the hostname of the compute node on which the job is running:

```
/bin/hostname
```

- A slightly more complicated job to run a R script in your submission file:

```
# my_script.r is the R script that you want to run,  
it takes two arguments: arg1 and arg2
```

```
Rscript my_script.r arg1 arg2
```

HOW TO CREATE A SUBMISSION SCRIPT

PBS RELATED ENVIRONMENT VARIABLES

- Can be used within your job submission files
- Common PBS environment variables include:

Name	Description
PBS_O_WORKDIR	Absolute path of the current working directory when you submitted this job
PBS_JOBID	Job ID number assigned to this job by PBS
PBS_JOBNAME	Job name supplied by the user
PBS_NODEFILE	File containing the list of nodes assigned to this job

HOW TO CREATE A SUBMISSION SCRIPT

PBS_O_WORKDIR

```
#!/bin/sh -l
#  FILENAME:  subjob.sub

#  Change to the directory from which you originally
#  submitted this job
cd $PBS_O_WORKDIR

#  Print out the current working directory path
pwd
```

HOW TO CREATE A SUBMISSION SCRIPT

PBS_JOBID

```
#!/bin/sh -l
#  FILENAME:  subjob.sub

#  Change to the directory from which you originally
#  submitted this job
cd $PBS_O_WORKDIR

#  Make a directory for each job to store the job
#  output and error file
mkdir $PBS_JOBID
myprogram 1> $PBS_JOBID/myprogram.out 2> \
    $PBS_JOBID/myprogram.err
```

HOW TO CREATE A SUBMISSION SCRIPT

PBS_JOBNAME

```
#!/bin/sh -l
# FILENAME: subjob.sub
#PBS -N input_800

# Change to the directory from which you originally
# submitted this job
cd $PBS_O_WORKDIR

# I have an input file called input_800 (the
# PBS_JOBNAME)
./hello $PBS_JOBNAME
```

HOW TO CREATE A SUBMISSION SCRIPT

PBS_NODEFILE

```
#!/bin/sh -l
# FILENAME: subjob.sub

# Change to the directory from which you originally
# submitted this job
cd $PBS_O_WORKDIR

# Run a MPI program
mpixexec -n 40 -machinefile $PBS_NODEFILE ./mpi_hello
```

HOW TO CREATE A SUBMISSION SCRIPT

Tip: BACKGROUND TASKS

- Option 1: one process per PBS job
 - Drawbacks: overhead!
- Option 2: pack lots of processes into a loop inside a single PBS jobs
 - Drawbacks:
 - need to be careful not to overwhelm the node at once,
 - slightly more complicated submission script.
 - Huge gain in efficiency.
 - NB: if you start something with a ‘&’ in background, remember a ‘wait’ at the end!
- Scripted PBS submissions... use some `sleep` between `qsub`’s! (0.2-0.3s)

```
#PBS .....  
command1 &  
command2 &  
...  
command16 &  
wait
```

JOBS MONITORING COMMANDS



EXAMPLE JOB SUBMISSION SCRIPT

subjob4cluster101.sub

```
#!/bin/sh -l
#PBS -N blastx_nr
#PBS -q workshop
#PBS -l nodes=1:ppn=8
#PBS -l walltime=03:00:00
```

```
module purge
module load bioinfo
module load blast
module list
```

```
cd $PBS_O_WORKDIR
pwd
```

```
cat subjob4cluster101.sub
```

```
date +"%d %B %Y %H:%M:%S"
```

```
# this is a test example
echo " "
```

```
blastx \
-query sequences.fasta \
-out seq.blastx_nr.fmt6 \
-db nr -num_threads 8 \
-outfmt 6 \
-evalue 1E-06 \
-max_target_seqs 5
```

```
echo " "
```

```
date +"%d %B %Y %H:%M:%S"
```

HOW TO CHECK FOR AVAILABLE RESOURCES

qlist

List all queues I can use and their current status and limits.

```
$ qlist
```

Current Number of Cores

Queue	Total	Queue	Run	Free	Max Walltime
=====	=====	=====	=====	=====	=====
workq	352	900	334	18	336:00:00
workshop	8	0	0	8	4:00:00

HOW TO SUBMIT AND MONITOR YOUR JOB

qsub

Once you have a job submission script, you may submit this script to PBS using the qsub command. PBS will find an available compute node or set of compute nodes and run your job there, or leave your job in a queue until some become available.

```
$ qsub subjob4cluster101.sub
787233.radon-adm.rcac.purdue.edu

# here the number indicates the job id that PBS
# assigns to the job, this number can be used to track
# your job.
```

HOW TO MONITOR THE STATUS OF YOUR JOB

`qstat -u myusername`

List all current jobs from the user *myusername*, where *myusername* is your Purdue login.

```
$ qstat -u gandino
```

```
radon-adm.rcac.purdue.edu:
```

Job ID	Username	Queue	Jobname	SessID	NDS	TSK	Req'd Memory	Req'd Time	S	Elap Time
-----	-----	-----	-----	-----	---	---	-----	-----	-	-----
787233.radon- adm.rcac.	gandino	workshop	blastx_nr	--	1	8	--	03:00:00	Q	--

HOW TO SUBMIT AND MONITOR YOUR JOB

checkjob

This is useful when your job is not running and you want to know why

```
$ checkjob 787233
```

```
job 787233
```

```
Aname: blastx_nr
```

```
State: Idle ...
```

```
BecameEligible: Mon Oct 19 17:15:04
```

```
SubmitTime: Mon Oct 19 17:13:34
```

```
...
```

```
NOTE: job req cannot run in partition radon-adm (available  
procs do not meet requirements: 0 of 8 procs found)
```

```
idle procs: 8   feasible procs: 0 ...
```

HOW TO SUBMIT AND MONITOR YOUR JOB

showstart

Will give you an estimate of when your job should start

```
$ showstart 787233
```

```
-----  
NOTE: The following information has been cached by the remote  
server and may be slightly out of date  
-----
```

```
job 787233 requires 8 procs for 3:00:00
```

```
*Estimated Rsv based start in 17:14:47 on Tue Oct 20 10:39:31
```

```
Estimated Rsv based completion in 20:14:47 on Tue Oct 20 13:39:31
```

* The estimated time will vary depending on which queue your are running the job

HOW TO SUBMIT AND MONITOR YOUR JOB

qdel

Stop and delete the job ID *myjobid*

```
$ qstat -u gandino
radon-adm.rcac.purdue.edu:
```

Job ID	Username	Queue	Jobname	SessID	NDS	TSK	Memory	Time	S	Time
787233.radon-adm.rcac.	gandino	workshop	blastx_nr	12208	1	8	--	3:00:00	R	01:00:11

```
$ qdel 787233
```

```
$ qstat -u gandino
radon-adm.rcac.purdue.edu:
```

Job ID	Username	Queue	Jobname	SessID	NDS	TSK	Memory	Req'd Time	Req'd S	Elap Time
787233.radon-adm.rcac.	gandino	workshop	blastx_nr	12208	1	8	--	3:00:00	C	--

GRAB THE LAB FILES

- Login to Radon
\$ `ssh <username>@radon.rcac.purdue.edu`
- Change to your working space \$RCAC_SCRATCH
\$ `cd $RCAC_SCRATCH`
- Copy the directory lab_cluster101 into your directory:
\$ `cp -r /depot/itap/cluster101-2016/sub_job ./`
- Move into the newly created lab_cluster101 directory:
\$ `cd sub_job`
\$ `pwd`
\$ `ls`

LAB 3 - RUN A BATCH JOB (QSUB)

- Open the batch script in an editor to see if you need to change it:
`$ nano subjob4cluster101.sub # or vi, or emacs, or just cat subjob4cluster101.sub`
- Launch the batch job
`$ qsub subjob4cluster101.sub`
- Monitor the job's status
`$ qstat -a -u <username>`
`$ checkjob -v <jobid>`
`$ qstat -a | more # hit space bar to advance`
- When job completes, take a look at results:
`$ ls # Note presence/names of output files`
`$ more subjob4cluster101.sub.oxxxxx # "xxxxxx" is your job's jobid`
`$ more subjob4cluster101.sub.exxxxx # "xxxxxx" is your job's jobid`

INTERACTIVE JOBS AND THINLINC



INTERACTIVE JOBS

INTERACTIVE JOB

Jobs are typically ran without user interaction, but some do need interaction.

```
$ qsub -I  
qsub: waiting for job 7867033.radon-adm.rcac.purdue.edu  
to start
```

The job will be queued like any other job, and may take some time to start.

```
qsub: job 7867033.radon-adm.rcac.purdue.edu ready  
$ hostname  
radon-a044
```

INTERACTIVE JOBS

WHY RUN INTERACTIVELY ON A COMPUTE NODE?

- Dedicated compute node (vs a shared frontend)
- Test code without impacting others
- Quicker develop / test / debug cycle
- Run GUI apps as a job
 - Matlab
 - IGV
 - Fluent

INTERACTIVE JOBS

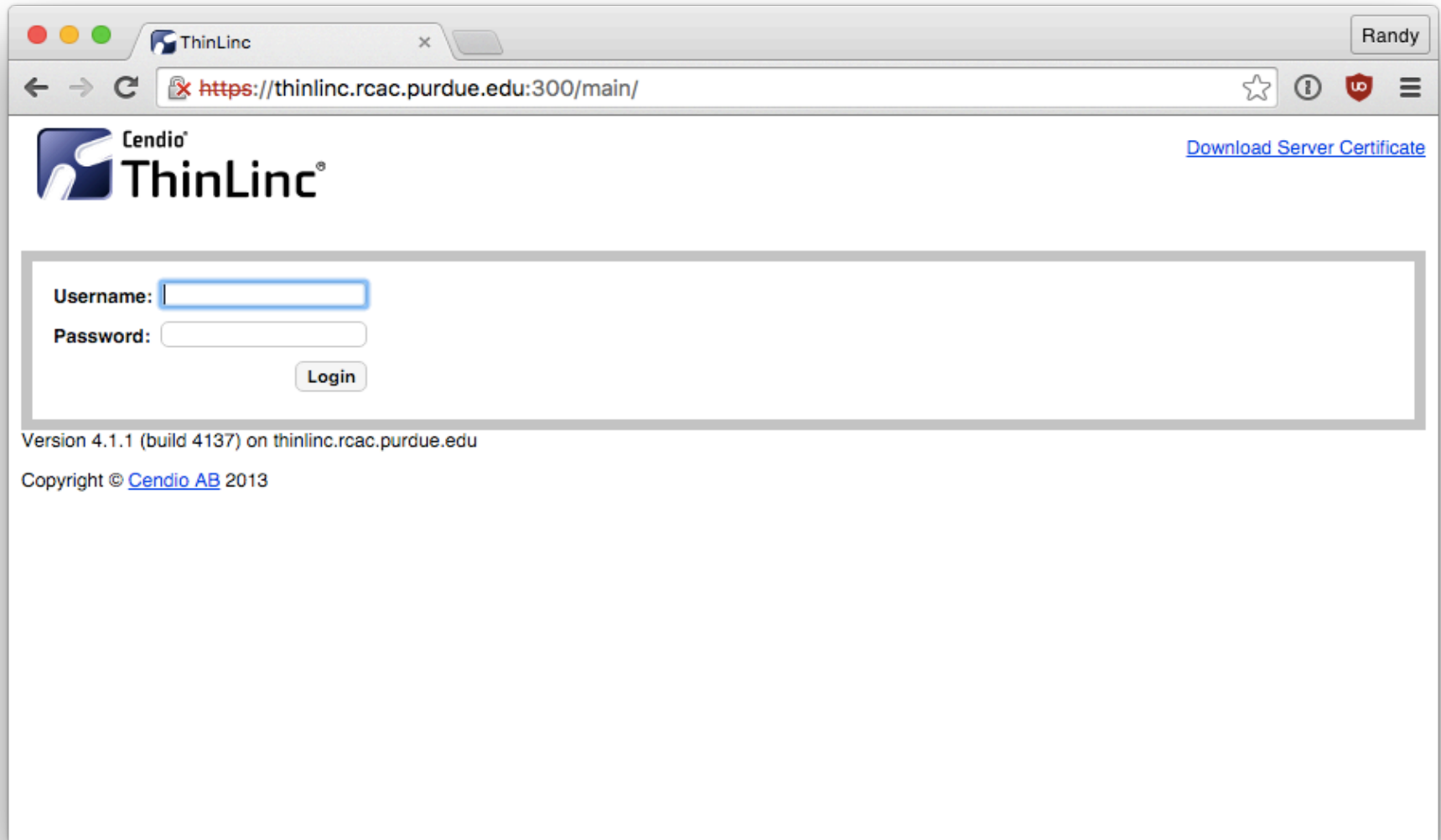
DISPLAY FORWARDING (X11)

To forward your display, use “-Y” option to `ssh`

```
$ ssh -Y radon.rcac.purdue.edu
```

We recommend using thinlinc.rcac.purdue.edu.

- Thinlinc keeps session alive after disconnect
- Faster than X11 forwarding
- Available off-campus, but runs as on-campus (useful for licenses)
- Standalone client is also available



A screenshot of a web browser window displaying the ThinLinc login page. The browser's address bar shows the URL <https://thinlinc.rcac.purdue.edu:300/main/>. The page features the Cendio ThinLinc logo at the top left and a link to "Download Server Certificate" at the top right. The main content area contains a login form with fields for "Username:" and "Password:", and a "Login" button. Below the form, it states "Version 4.1.1 (build 4137) on thinlinc.rcac.purdue.edu" and "Copyright © Cendio AB 2013".

ThinLinc

[Download Server Certificate](#)

Username:

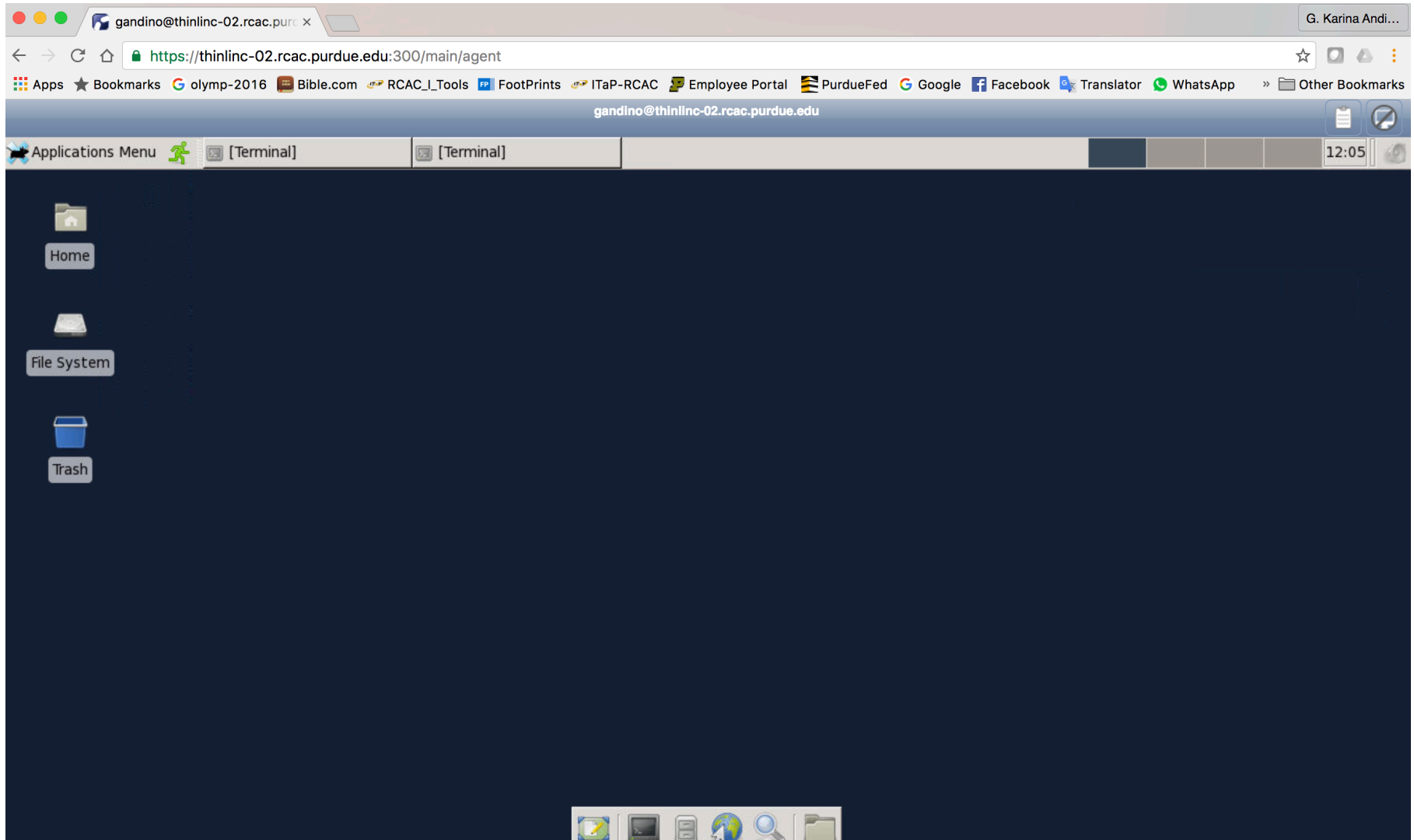
Password:

Login

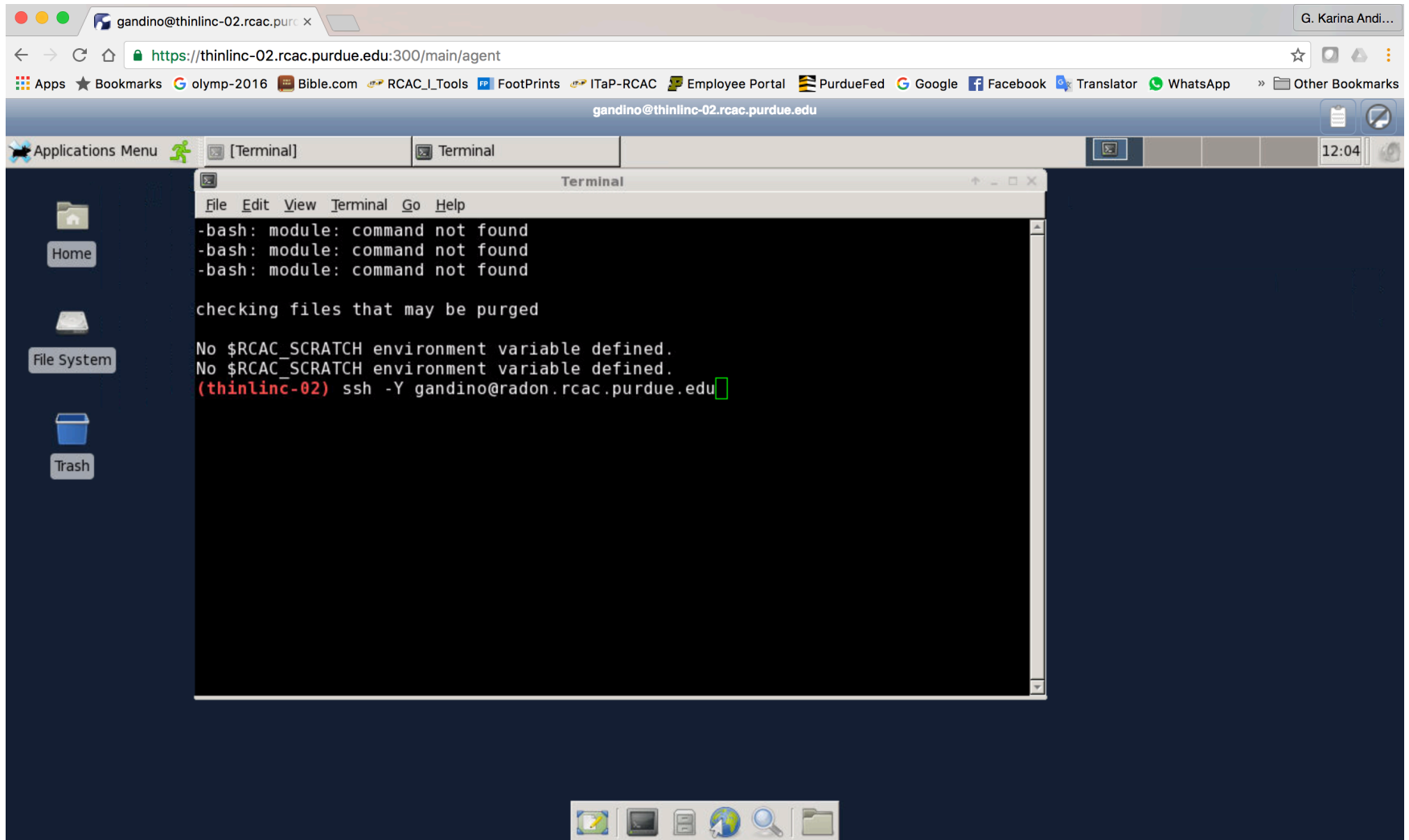
Version 4.1.1 (build 4137) on thinlinc.rcac.purdue.edu

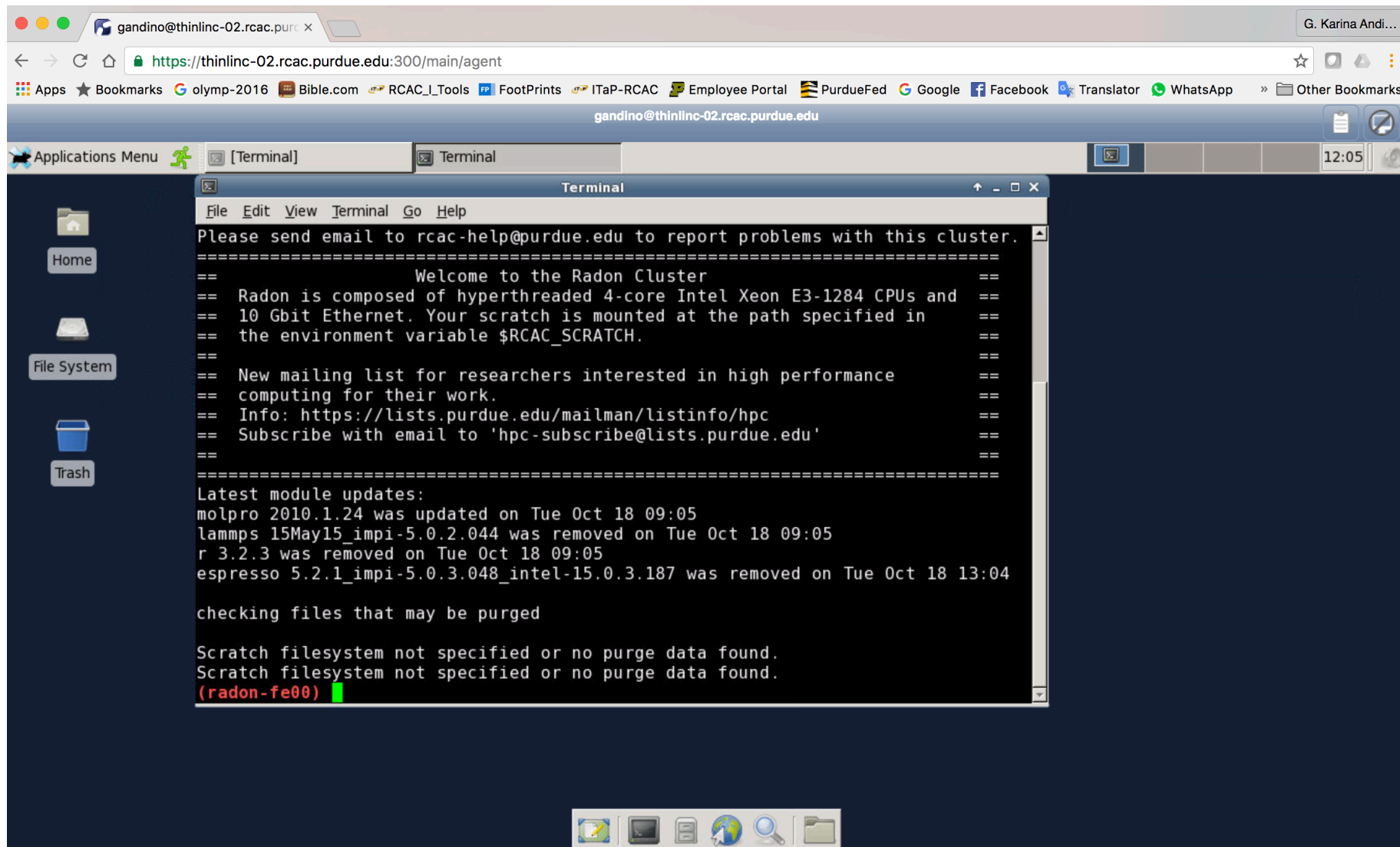
Copyright © [Cendio AB](#) 2013

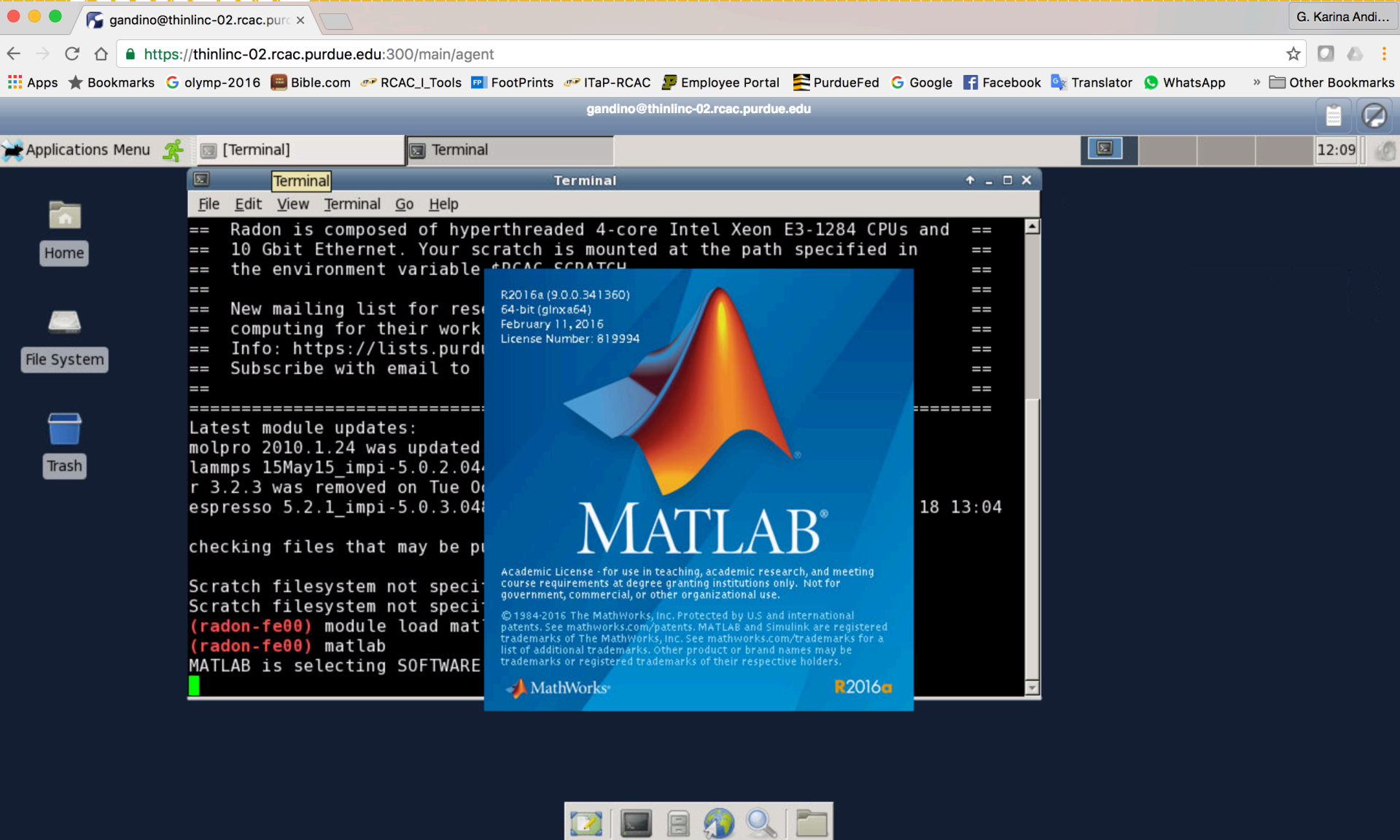
THINLINC



THINLINC







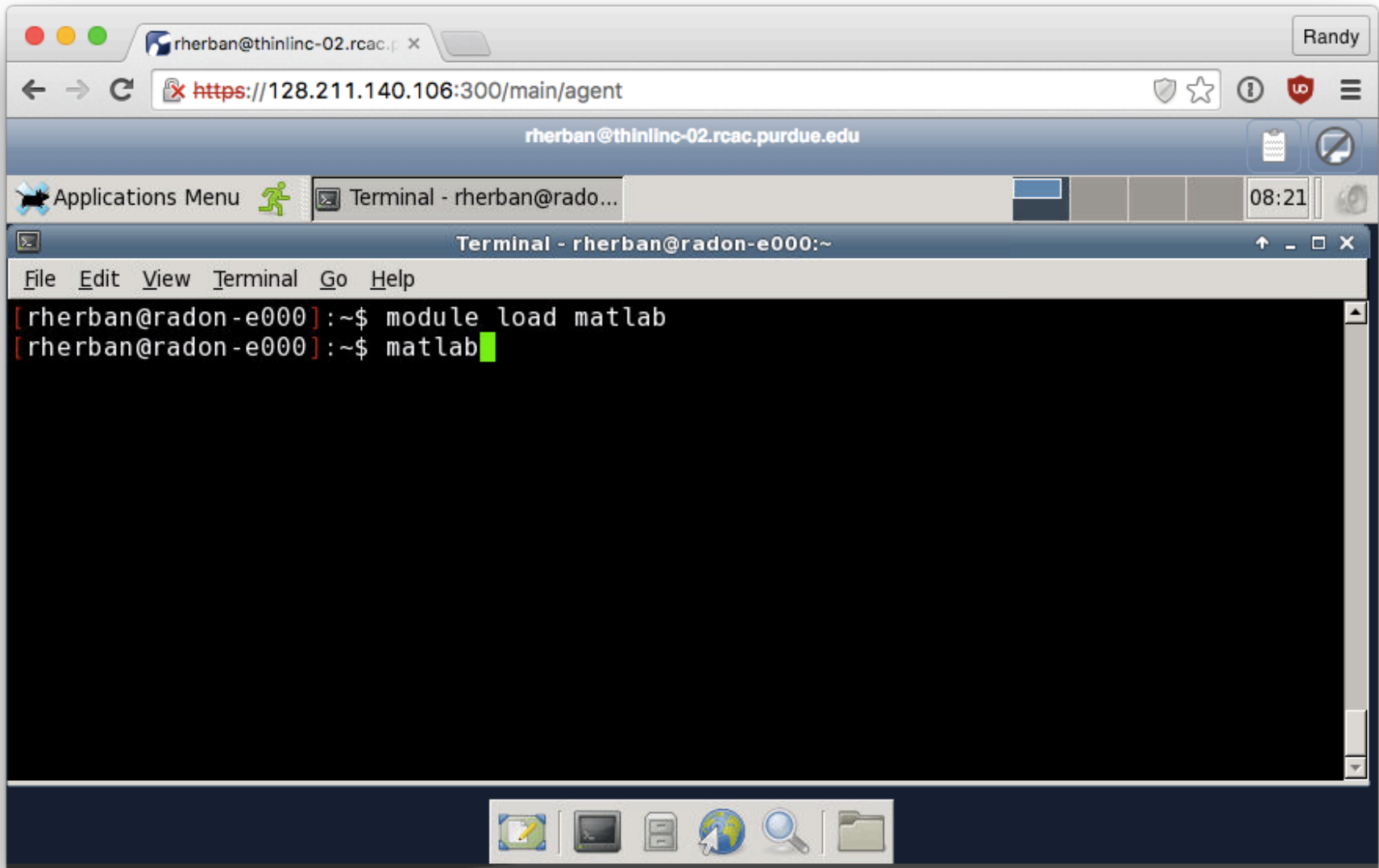
The screenshot displays a ThinLinc desktop environment. At the top, a web browser window is open to <https://thinlinc-02.rcac.purdue.edu:300/main/agent>. Below the browser, the MATLAB R2016a - academic use interface is visible. The MATLAB window includes a top toolbar with icons for file operations (New Script, Open, Save, etc.), a central menu bar with options like HOME, PLOTS, and APPS, and a bottom panel with several sub-windows. The 'Current Folder' window on the left shows a directory tree with folders like '90percen_virus', 'BioPerl-1.6.1', and 'class_project'. The 'Command Window' on the right displays a message about the 'New MATLAB Graphics System' in MATLAB R2014b, accompanied by a 3D surface plot. The 'Workspace' window at the bottom left is empty. The desktop background is dark blue, and the taskbar at the bottom shows various system icons and a clock displaying 12:10.

INTERACTIVE JOBS AND THINLINC

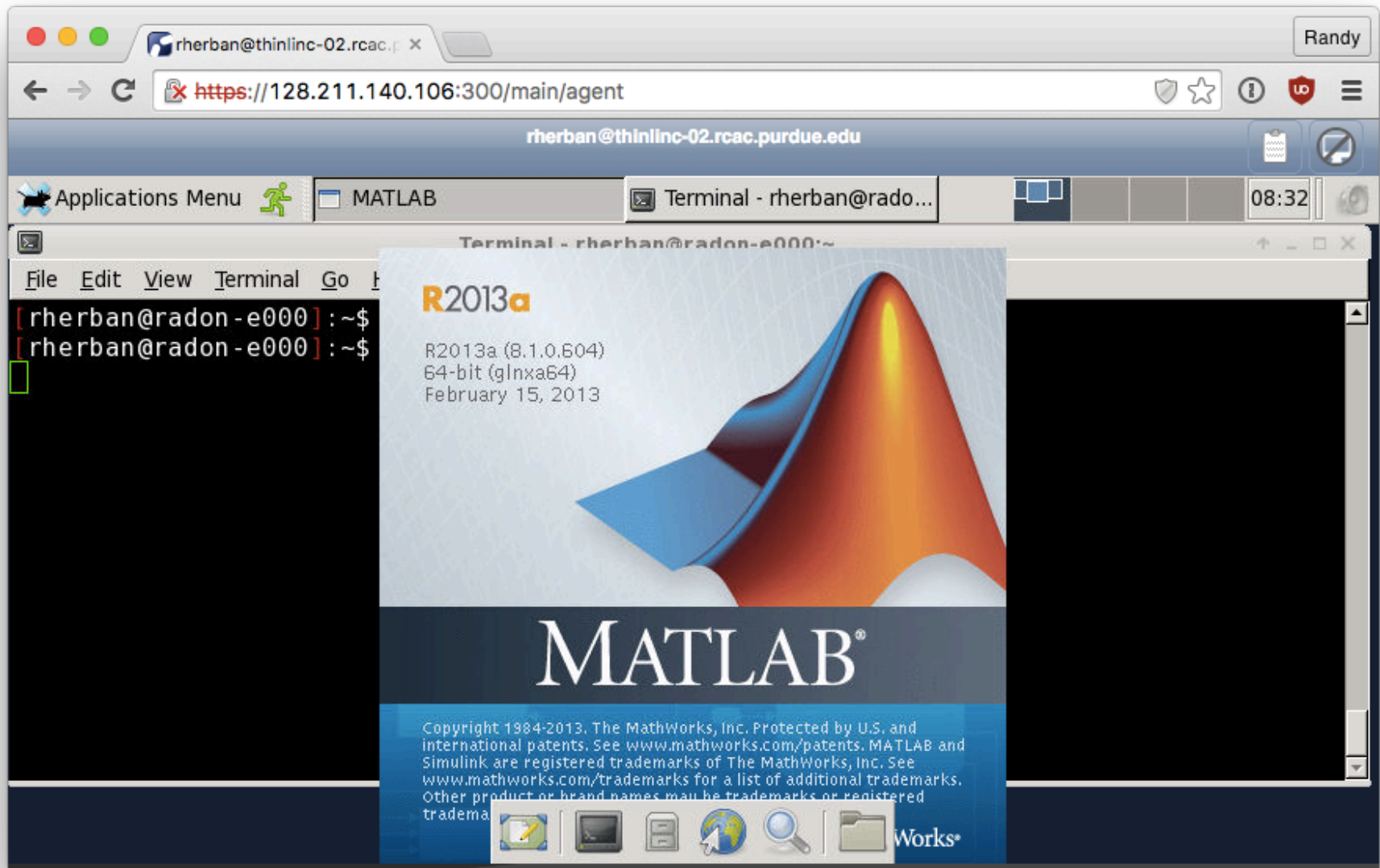
BUT WAIT, THIS JUST RAN ON THE FRONTEND!

```
$ ssh -Y radon.rcac.purdue.edu
$ echo $DISPLAY
radon-fe01.rcac.purdue.edu:27.0
$ qsub -I -q workshop -l nodes=1:ppn=8,walltime=03:00:00 -v
DISPLAY
qsub: waiting for job 903193.radon-adm.rcac.purdue.edu to start
qsub: job 903193.radon-adm.rcac.purdue.edu ready
$ echo $DISPLAY
radon-fe00.rcac.purdue.edu:27.0
```

INTERACTIVE JOBS AND THINLINC



INTERACTIVE JOBS AND THINLINC



WORKING SPACE AND STORAGE

HIERARCHY

- Types of storage and when to use them
- \$HOME vs \$RCAC_SCRATCH
- Local disk (/tmp) on each individual node
- Research Data Depot
- Fortress archive

WORKING SPACE AND STORAGE

SUMMARY

	\$HOME	\$RCAC_SCRATCH	/tmp	/depot/...	Fortress (HPSS)
Capacity	25 GB	Varies by cluster... 1 to 100 TB	150-400 GB	100 GB and up	unlimited
Resilience to hardware failures	yes	yes	no	yes	yes
Resilience to human errors	yes (snapshots)	no	no	yes (snapshots)	no
Subject to purging	no	yes	yes	no	no
Performance	medium	high	medium	medium	slow to very slow
Designed for HPC (running jobs off it)	no	yes	no	in moderation	-No (as main I/O) -Yes (for staging and archiving)
Common access within cluster	yes	yes	no	yes	yes
Common access across clusters	yes	no	no	yes	yes
Globus endpoint	yes	yes	no	yes	yes

WORKING SPACE AND STORAGE

DATA DEPOT, FORTRESS AND GLOBUS

- Data Depot - a high-capacity, fast, reliable and secure data storage service.
 - Basically, a “home directory for labs”
 - Redundant, fast, flexible, accessible from anywhere
 - Don’t have to own nodes
 - Data Depot space is mounted and accessible from all nodes and clusters
- To learn more: <https://www.rcac.purdue.edu/storage/depot/>

WORKING SPACE AND STORAGE

DATA DEPOT, FORTRESS AND GLOBUS

- Fortress - a large, long-term storage system
 - Does not support direct FTP, SFTP or SCP transfers
 - HSI - FTP-style interface without requiring any special user knowledge.
 - HTAR - a utility to aggregate a set of files into a single tar archive directly into Fortress.
- To learn more: <https://www.rcac.purdue.edu/storage/fortress/>

WORKING SPACE AND STORAGE

DATA DEPOT, FORTRESS AND GLOBUS

- Globus – a powerful and easy to use file transfer service
 - Transfer large amounts of research data to and from campus, or share data with collaborators around the country or around the globe.
- Purdue Globus Portal: <https://transfer.rcac.purdue.edu>

WORKING SPACE AND STORAGE

TIPS FROM US

- Types of storage and when to use them:
“capacity, speed, longevity – pick any two”
- ITaP Research Computing offerings are designed around these scenarios:
 - Personal codes, executables, scripts – develop and store in \$HOME
 - Lab-wide codes, executables, scripts, settings – in /depot
 - Data and results (input/output):
 - produce/analyze – in \$RCAC_SCRATCH
 - permanently store – in Fortress archive and/or in /depot
- myquota command
- Use variables instead of explicit paths
- **Snapshots are not a substitute for backups!**

FILE SYSTEM CITIZENSHIP

- Avoid running jobs from \$HOME
- Run demanding jobs from \$RCAC_SCRATCH
- Avoid frequent I/O when possible
- Minimize simultaneous I/O from many processes
- Learn to recognize/avoid other stressors
 - e.g. under-the-table stat (du, default ls) on big dirs
- Know when it's time to learn/use parallel I/O

GLOBUS



- Globus enables file transfers between two sites (“endpoints”)
 - Like hiring movers to pack your house and get it to new location
 - You don’t have to be involved
 - Get an update when it completes / fails
 - Goes directly between the two fast sites
 - Doesn’t use your bandwidth!
- Endpoints on all of our clusters, many other Universities/ labs and even on your own laptop!

GLOBUS



- Exercise:
- Log into Globus:
- transfer.rcac.purdue.edu
- globus.org

Log in to use Purdue Globus Web App

Use your existing organizational login

e.g. university, national lab, facility, project, Google or [Globus ID](#)

(Your Globus username and password used prior to February 13, 2016 is now Globus ID)

Purdue University Main Campus ▼

Continue



Globus uses CILogon to enable you to Log In from this organization. By clicking Continue, you agree to the [CILogon privacy policy](#) and you agree to share your username, email address, and affiliation with CILogon and Globus. You also agree for CILogon to issue a certificate that allows Globus to act on your behalf.

Didn't find your organization? Then use Globus ID to [sign up](#).

<https://transfer.rcac.purdue.edu>

PURDUE
UNIVERSITY

Information Technology at Purdue
Research Computing (RCAC)

Manage Data Groups Support Account

Transfer Files | Activity | Endpoints | Bookmarks | Publish | Console

RECENT ACTIVITY ○ 0 ▾ 0 ○ 0

Transfer Files

Endpoint ☆

Path Go

Endpoint ☆

Path Go

Start by selecting an endpoint.

Start by selecting an endpoint.

GLOBUS ENDPOINTS

Endpoint

purdue|

✕ 🔍 Cancel

Purdue Structural Biology - Saturn
owner: purduestbio@globusid.org
Purdue Cryo-EM Facility

Purdue Research Computing - Home Directories
owner: purdue@globusid.org
This endpoint provides access to home directories storage on all Purdue clusters - /home


Purdue Research Computing - Data Depot
owner: purdue@globusid.org
This endpoint provides access to the Purdue Research Data Depot storage - /depot

Purdue Radon Cluster
owner: purdue@globusid.org
This endpoint provides access to scratch storage on the Purdue Radon cluster - /scratch/radon

Purdue Peregrine1 cluster
owner: purdue@globusid.org
This endpoint provides access to scratch storage on the Purdue Peregrine1 cluster - /scratch/peregrine1

Purdue Hathi Cluster
owner: purdue@globusid.org
This endpoint provides access to scratch storage on the Purdue Hathi cluster - /hadoop/mnt/user

Purdue Hammer Cluster
owner: purdue@globusid.org
This endpoint provides access to scratch storage on the Purdue Hammer cluster - /scratch/hammer



GLOBUS TRANSFER

[Transfer Files](#) | [Activity](#) | [Endpoints](#) | [Bookmarks](#) | [Publish](#) | [Console](#)

Transfer Files

RECENT ACTIVITY ○ 0 ▾ 0 ○ 0

Endpoint



Path

Go



Endpoint

Path

Go



select all

↑ up one folder

↻ refresh list



unix101-2016

Folder


Start by selecting an endpoint.



GLOBUS TRANSFER

Endpoint

Purdue Structural Biology - Saturn owner: purduestbio@globusid.org Purdue Cryo-EM Facility
Purdue Research Computing - Home Directories owner: purdue@globusid.org This endpoint provides access to home directories storage on all Purdue clusters - /home
Purdue Research Computing - Data Depot owner: purdue@globusid.org This endpoint provides access to the Purdue Research Data Depot storage - /depot
Purdue Radon Cluster owner: purdue@globusid.org This endpoint provides access to scratch storage on the Purdue Radon cluster - /scratch/radon
Purdue Peregrine1 cluster owner: purdue@globusid.org This endpoint provides access to scratch storage on the Purdue Peregrine1 cluster - /scratch/peregrine1
Purdue Hathi Cluster owner: purdue@globusid.org This endpoint provides access to scratch storage on the Purdue Hathi cluster - /hadoop/mnt/user
Purdue Hammer Cluster owner: purdue@globusid.org This endpoint provides access to scratch storage on the Purdue Hammer cluster - /scratch/hammer



GLOBUS TRANSFER




[Transfer Files](#) | [Activity](#) | [Endpoints](#) | [Bookmarks](#) | [Publish](#) | [Console](#)

Transfer Files

RECENT ACTIVITY  0  0  0

Endpoint ☆

Path




select all  up one folder  refresh list 

unix101-2016

Folder

Endpoint ☆

Path

select all  up one folder  refresh list 

a

Folder

b

Folder

c

Folder

d

Folder

e

Folder

f

Folder

g

Folder

h

Folder

i

Folder

j

Folder

k

Folder

l

Folder

m

Folder

n

Folder

o

Folder

p

Folder

q

Folder

r

Folder

s

Folder

t

Folder

Navigate to your directory



GLOBUS TRANSFER

[Transfer Files](#) | [Activity](#) | [Endpoints](#) | [Bookmarks](#) | [Publish](#) | [Console](#)

Transfer Files

RECENT ACTIVITY ○ 0 ▾ 0 ⬡ 0

Endpoint ☆

Path Go

select none ⬅ up one folder ↻ refresh list ☰

📁 unix101-2016 Folder

Select file(s)

⏪ ⏩

▶

Endpoint ☆

Path Go

select all ⬅ up one folder ↻ refresh list ☰

Transfer

Refresh



GLOBUS TRANSFER

[Transfer Files](#) | [Activity](#) | [Endpoints](#) | [Bookmarks](#) | [Publish](#) | [Console](#)

Transfer Files

RECENT ACTIVITY  1  0  0

Transfer request submitted successfully. Task id: [2e22df00-961f-11e6-b0a4-22000b92c261](#)

Endpoint



Path

Go







Endpoint






Path

Go

select none  up one folder  refresh list 

 unix101-2016

Folder

select all  up one folder  refresh list 

- Can share select folders with outside collaborators
 - No need to have Purdue ID
- Can set up an endpoint on your personal laptop:
 - <https://www.globus.org/globus-connect-personal>

WORKFLOW

- Unless you are running a one-off simulation, you should stop and think about the workflow
- What's the optimal resource request?
 - Experimentation - scaling studies
- Do I break my workflow into multiple jobs?
- How does my data flow?
 - From original source
 - Working spaces
 - Post-processing/analysis
 - Archives

WORKFLOW

- Example – nightly weather forecasts
- Break into multiple jobs
 - Obtain raw data (current observations)
 - Pre-process raw data and prepare
 - Run forecast
 - Copy results somewhere or post-process into viewable images



GLOBUS & WORKFLOW

- Build your workflow

- Grab dataset before batch of jobs

- ```
rsync -rav /depot/group/data/experiment1 /scratch/radon/u/
user/
```

- Submit batch of jobs

- Save to /scratch/radon/u/user/results

- Copy results back to Depot or Fortress

- ```
ssh user@cli.globusonline.org transfer -- purdue#radon:/  
scratch/radon/u/user/results purdue#depot:/depot/group/  
data/experiment1-results
```

- Requires ssh key setup with Globus
 - Happens in the background; doesn't waste job time



GLOBUS & WORKFLOW

- Create GlobusID
- Go to globus.org and create Globus ID
- Log in with it and link your Purdue Credentials



GLOBUS SSH KEY

- Can use existing ssh key

```
$ cd ~  
$ ls -las .ssh  
$ cat .ssh/id_rsa.pub
```

- Or generate a new key

```
$ ssh-keygen -t rsa -f .ssh/id_globus -N ""  
$ cat .ssh/id_globus.pub
```

- Copy the key



GLOBUS SSH KEY

https://transfer.rcac.purdue.edu



Information Technology at Purdue
Research Computing (RCAC)



Manage Data

Groups

Support

Account



Account



Manage Your Consents



Logout

Identities



Manage These Identities



Add Linked Identity

Randy Herban via Globus ID

Name: Randy Herban

Username: rherban2@globusid.org

E-mail: rherban@purdue.edu

This is your **primary** identity.

This identity is managed by [Globus ID](#), where you can [change your password](#), [update your profile](#), and [manage SSH and X.509 keys](#).

Randy L Herban via Purdue University Main Campus

Name: Randy L Herban

Username: rherban@purdue.edu

E-mail: rherban@purdue.edu

This is a **linked** identity.



GLOBUS SSH KEY



Randy Herban (rherban2@globusid.org)

[Back to Account Information on the Purdue Globus Web App](#)

[Home](#) | [Log Out](#)

Manage SSH and X.509 Keys

[Add a New Key](#)



You have the following key attached to your Globus ID account:

SSH Public Key

[expand](#)

[delete](#)

[rename](#)



GLOBUS SSH KEY



Randy Herban (rherban2@globusid.org)

[Back to Account Information on the Purdue Globus Web App](#)

[Home](#) | [Log Out](#)

Add a New Key

[Back to Key Management](#)

Alias

Type ☐ SSH Public Key

Allows you to use the [Globus command line interface](#) via any standard `ssh` client. Some Globus-based applications may also use this SSH authorized key to provide access to their resources.

☐ X.509 Credential

Allows you to use the [Globus command line interface](#) via a `gsissh` client, and to activate endpoints with a delegated proxy of this X.509 credential via `gsissh` to the CLI.

Body

 Paste key here

Add Key

Cancel

ADDITIONAL WORKSHOPS

- Software Carpentry (Oct 27 & 28) – Sign Up Now!
 - Bash Automation
 - Git
 - R / Rstudio
 - <https://www.rcac.purdue.edu/news/896>
- XSEDE Workshops (Monthly)
 - Big Data (Nov 1)
 - GPU Programming (Dec 6)
 - MPI
 - Summer Boot Camp