

Why we are here!

- You have a computational task of some sort
- Most people who come to CARC fall into these categories
 - ... for a publication you are working on
 - ... an upcoming class assignment
 - ... analysis of data for a government agency
- The common feature is that they require more compute/memory/storage resources than is commonly available.

Why we are here!

- To use the resources at CARC effectively there are things you have to learn! Some are commonplace but some are very specific to HPC and CARC.

By the end of the day you will know how to:

- Request time compute time on the clusters (PBS scripting)
 - Interact with contained environments (Anaconda, Modules, Singularity)
 - Get output from your compute jobs
 - Use storage appropriately (scratch vs user storage)
 - Write scripts to run embarrassingly parallel tasks
 - Write and run a simple MPI program for tightly coupled tasks
- Ultimately get a huge increase in the computing power you can apply.

Introduction to HPC at the UNM Center for Advanced Research Computing

Matthew Fricke (mfricke@unm.edu)

Research Assistant Professor

<http://unm.edu/~mfricke>

Please send corrections to me at mfricke@unm.edu.

Version 1.1.

1.1 Changelog

- Corrected various typos
- Changed font to make code clearer
- Reordered MAUI/Torque slides
- Corrected “conda list” to “conda env list”
- Added `–machinefile $PBS_NODEFILE` to `mpirun` example
- Added slide on torque queue status commands
- Modified pbs examples for wheeler instead of galles

What is High Performance Computing?

Scaling up

- NVIDIA DGX-2H (\$400,000 each, 81k CUDA cores, 10240 tensor cores)
https://www.nvidia.com/content/dam/en-zz/es_em/Solutions/Data-Center/dgx-2/nvidia-dgx-2h-datasheet.pdf



Scaling out

- Stampede 2
- <https://www.tacc.utexas.edu/systems/stampede2>
- \$30,000,000, 285,000 CPUs



Why it Matters to You

- Grant and publication reviewers know about these systems so there are no excuses for small sample sizes.
- Machine Learning is showing up everywhere from cosmology to firefighting. Machine Learning requires enormous resources to process huge datasets.

The Center for Advanced Research Computing's Mission

The UNM Center for Advanced Research Computing is the hub of computational research at UNM and one of the largest computing centers in the State of New Mexico. It is an interdisciplinary community that uses computational resources to create new research insights. The goal is to lead and grow the computational research community at UNM.

CARC provides not just the computing resources but also the expertise and support to help the university's researchers. This service is available to faculty, staff, and student researchers free of charge through support from the UNM Office of the Vice President for Research.

<http://carc.unm.edu>



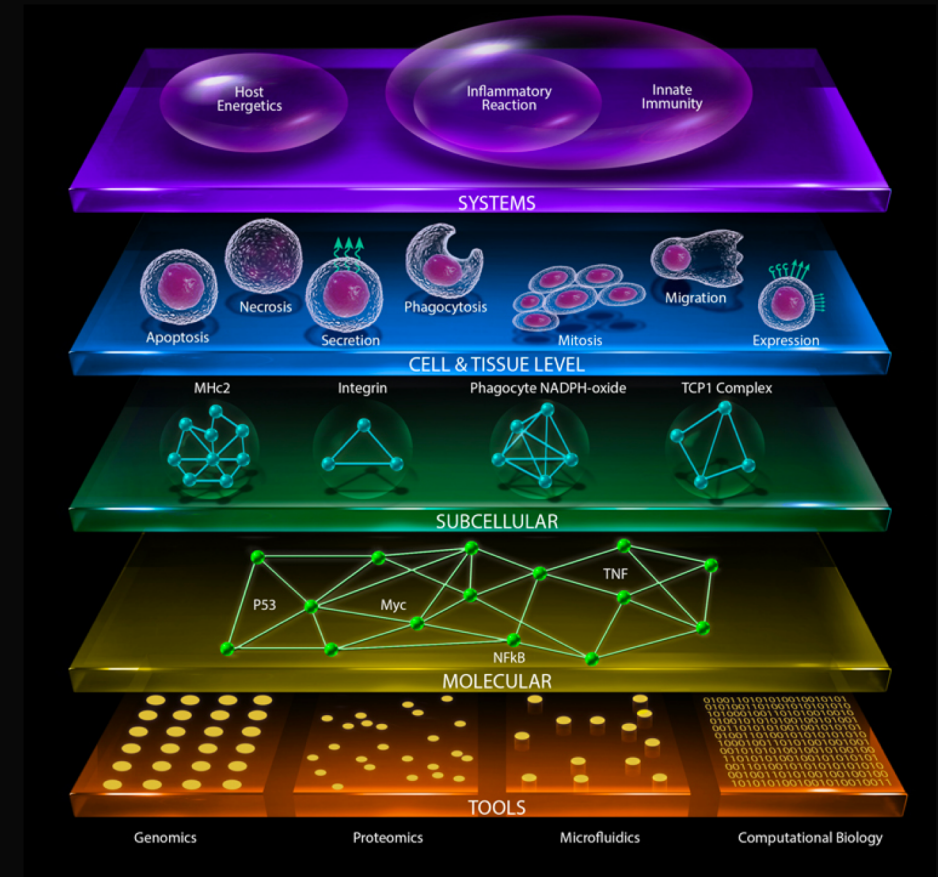
Big Data and Machine Learning

- Machine learning needs lots of everything.
- The current revolution in convolutional neural networks (Think Google, Self-driving cars, etc) is due to algorithms rooted in the 1960s being given huge training sets.
- Most of machine learning comes down to floating point matrix and vector operations. GPUs excel at those operations and are orders of magnitude faster at them than CPUs.
- Xena has dual Nvidia Tesla K40M GPUs for this purpose.



Biology

- Computational biology memory usage increases with input sizes. Rapid genotyping tools generate sequences faster and faster.
 - Hundreds of GB of RAM are becoming a normal requirement to complete these calculations.
 - The Taos cluster is dedicated to computational biology and has 440 CPUs and 300 GB per node.
 - Xena has 3 TB RAM nodes.
- Pandemic flu modelling
 - Tuberculosis antibiotic resistance
 - Pacific island bird genetics
 - Intra-species viral spread
 - NM Tree species mapping from LASER scans

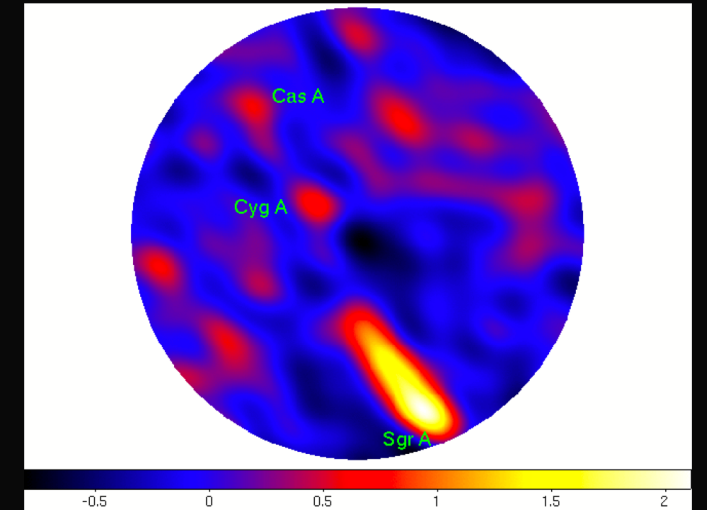


Physics

- Wheeler is a general purpose scale-out machine used by biophysicists, cosmologists, and many others.
- Gibbs is primarily used by computational chemists.



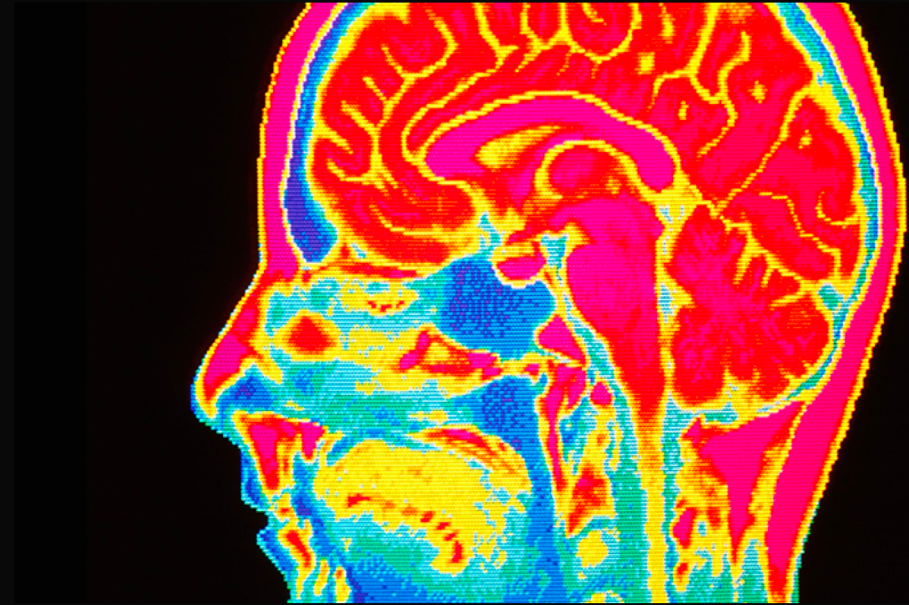
Molecular simulation of
new photovoltaic materials



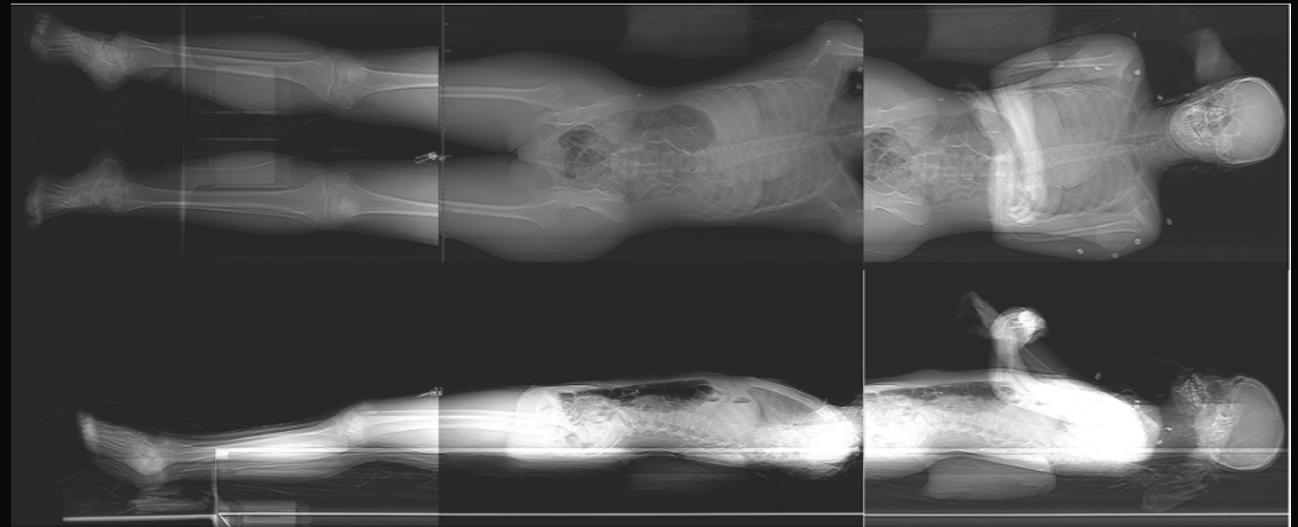
Long Wavelength Array Radio
Telescope
Data Processing on Wheeler

When users need a lot of network bandwidth or storage

- Lots of data comes in the form of large images
- Largest online pathology database (15k people, 15k image each)
- MRI and FMRI image databases



- Specialized Network Access
 - Mapping the Great Firewall





Albuquerque Integrated Reporting System

2016101	Finding approximate symmetries in graphs
2016100	Simulation of impact on armor plate using Velodyne
2016099	Cyber-infrastructure Performance Modeling
2016098	Effective refinement of protein structure
2016097	LANL HIGRAD
2016096	Coarse-grained modeling of biomolecules
2016095	Resilient Composites
2016093	Computational Investigation of Ru Photochromes
2016090	Reinforcement learning neuropathologies underlying psychiatric sequelae in Traumatic Brain Injury
2016087	Driving forces and dynamics between small molecules and amyloid- Aβ aggregates
2016086	Polyurethane Foam Modeling
2016084	An Interval Multi-level Monte Carlo Method for Reliability Analysis of Imprecise Probabilistic Systems
2016083	Modeling Immune System Cell Search Processes
2016081	RNA transcriptome sequencing of T-cells exposed to Uranium and Arsenic
2016080	Predicting Progression to Alzheimer's Disease
2016079	Relating plant traits to biomass dynamics in New Mexico aridlands
2016078	Statistical methods for investigating large scale gene environment interaction
2016077	TB Genomic Analysis
2016076	LiDAR-based tree identification in Northern New Mexico

Fricke, Matthew
Hogeveen, Jeremy P
Christodoulou, Christos
ANDEROGLU, Osman
Sorrentino, Francesco
Shen, Yu-Lin
Bridges, Patrick G
Nishima, Wataru
Poroseva, Svetlana
He, Yi
Taha, Mahmoud Reda
Rack, Jeffrey J
Hogeveen, Jeremy P
Chi, Eva Y
Tjiptowidjojo, Kristianto
Motamed, Mohammad
Fricke, Matthew
Schilz, Jodi R
Calhoun, Vince
Whitney, Kenneth
Luo, Li
Wearing, Helen

<u>2016050</u>	<u>Investigating the impact of metal contaminants in environmental microbial populations</u>	Cerrato, Jose M.
<u>2016049</u>	<u>COMPUTATIONAL INVESTIGATION OF PARAMETER SPACE APPLICABLE TO NUCLEAR FACILITY SAFEGUARDS USING THE UNM CENTER FOR ADVANCED RESEARCH COMPUTING</u>	Arthur, Edward D
<u>2016047</u>	<u>Multiscale Mechanistic Model to Study Nanotherapy Delivery in Tumors</u>	Bearer, Elaine L
<u>2016046</u>	<u>Designing of Fuel Performance Experiments to be Performed Using the Annular Research Reactor (ACRR)</u>	Lee, Youho
<u>2016045</u>	<u>Model building the HELP physical unclonable function</u>	Plusquellic, Jim F
<u>2016044</u>	<u>Genomic Comparisons of Multipartite Symbiosis: Understanding the metabolic basis of parasitism</u>	Kamel, Bishoy S
<u>2016042</u>	<u>QIIME analyses of microbial sequences acquired from saliva samples</u>	Carroll-Portillo, Amanda
<u>2016041</u>	<u>Deterministic and Bayesian Seismic source inversion involves</u>	Appelo, Daniel EA
<u>2016040</u>	<u>Consequence-based Impact Rating using ANSYS</u>	Moreu, Fernando
<u>2016039</u>	<u>Machine-Learning Design of Novel Photovoltaic Materials Based on Conceptual Understanding of Electronic Structure Calculations</u>	Talipov, Marat R
<u>2016036</u>	<u>Melt migration in continental interiors</u>	Roy, Mousumi
<u>2016035</u>	<u>Monte Carlo Simulation for Weak Neutron Sources</u>	Goss, Vanessa
<u>2016034</u>	<u>Programmable Nanowalkers: Models and Simulations</u>	Stefanovic, Darko
<u>2016033</u>	<u>Deep Learning and Differential Geometry</u>	Huang, Hongnian
<u>2016032</u>	<u>Analyzing Neuronal Coordination during a task of Behavioral Flexibility in a Model of Fetal Alcohol Spectrum Disorder</u>	Brigman, Jonathan L
<u>2016031</u>	<u>“Mountain Lions on the Edge: Integrating Conservation into Urban Planning through Predictive Modeling”</u>	Milne, Bruce T
<u>2016029</u>	<u>Discrete Element Modeling of Drilled Shafts in Granular Materials</u>	Ng, Tang-Tat
<u>2016028</u>	<u>Exploration of optical rogue wave phenomena in dielectrics as a function of the intrinsic randomness or disorder</u>	Mafi, Arash
<u>2016026</u>	<u>Differential Splicing by Sex in DNA Repair Genes</u>	Berwick, Marianne
<u>2016021</u>	<u>Atlantic salmon microbiome</u>	Salinas, Irene
<u>2016019</u>	<u>Small Area Population Estimates</u>	Rhatigan, Robert
<u>2016018</u>	<u>Differential Gene Expression in Cancer</u>	Trujillo, Kristina

<u>2016050</u>	<u>Investigating the impact of metal contaminants in environmental microbial populations</u>	Cerrato, Jose M.
<u>2016049</u>	<u>COMPUTATIONAL INVESTIGATION OF PARAMETER SPACE APPLICABLE TO NUCLEAR FACILITY SAFEGUARDS USING THE UNM CENTER FOR ADVANCED RESEARCH COMPUTING</u>	Arthur, Edward D
<u>2016047</u>	<u>Multiscale Mechanistic Model to Study Nanotherapy Delivery in Tumors</u>	Bearer, Elaine L
<u>2016046</u>	<u>Designing of Fuel Performance Experiments to be Performed Using the Annular Research Reactor (ACRR)</u>	Lee, Youho
<u>2016045</u>	<u>Model building the HELP physical unclonable function</u>	Plusquellic, Jim F
<u>2016044</u>	<u>Genomic Comparisons of Multipartite Symbiosis: Understanding the metabolic basis of parasitism</u>	Kamel, Bishoy S
<u>2016042</u>	<u>QIIME analyses of microbial sequences and</u>	roll-Portillo, Amanda
<u>2016041</u>	<u>Deterministic and Bayesian Seismic source</u>	pelo, Daniel EA
<u>2016040</u>	<u>Consequence-based Impact Rating using</u>	reu, Fernando
<u>2016039</u>	<u>Machine-Learning Design of Novel Photo</u> <u>Structure Calculations</u>	ipov, Marat R
<u>2016036</u>	<u>Melt migration in continental interiors</u>	y, Mousumi
<u>2016035</u>	<u>Monte Carlo Simulation for Weak Neutron</u>	ss, Vanessa
<u>2016034</u>	<u>Programmable Nanowalkers: Models and</u>	fanovic, Darko
<u>2016033</u>	<u>Deep Learning and Differential Geometry</u>	ang, Hongnian
<u>2016032</u>	<u>Analyzing Neuronal Coordination during a</u> <u>Disorder</u>	gman, Jonathan L
<u>2016031</u>	<u>“Mountain Lions on the Edge: Integrating Conservation into Urban Planning through Predictive Modeling”</u>	Milne, Bruce T
<u>2016029</u>	<u>Discrete Element Modeling of Drilled Shafts in Granular Materials</u>	Ng, Tang-Tat
<u>2016028</u>	<u>Exploration of optical rogue wave phenomena in dielectrics as a function of the intrinsic randomness or disorder</u>	Mafi, Arash
<u>2016026</u>	<u>Differential Splicing by Sex in DNA Repair Genes</u>	Berwick, Marianne
<u>2016021</u>	<u>Atlantic salmon microbiome</u>	Salinas, Irene
<u>2016019</u>	<u>Small Area Population Estimates</u>	Rhatigan, Robert
<u>2016018</u>	<u>Differential Gene Expression in Cancer</u>	Trujillo, Kristina

And a couple of
hundred more.
(TLDR)

Basics of HPC Systems

- Parallelism within CPUs
- Parallelism within GPUs
- Parallelism of CPUs and GPUs
- Parallelism of whole computers

Some useful Linux commands

- `ssh username@wheeler.alliance.unm.edu`
- `scp myfile.txt username@wheeler.unm.edu:~`
- `scp username@wheeler.unm.edu:~/myfile.txt .`
- CyberDuck, or WinScp
- CARC supports secure file transfer (SFTP), so choose that protocol if you use a graphical transfer program.
- Rsync
- `ls -lah`
- `find . -name "*.txt"`

Compartmentalization

- A challenge running large multi-user systems is supporting all the different software required for hundreds of projects.
- Compartmentalization keeps user software isolated.
- We use 3 general methods at CARC: Environment Modules, Anaconda, and Singularity (these are the current standards so what you learn here will translate to other HPC centers)

Environment Modules

- Open a secure shell on wheeler if you haven't already

```
$ ssh username@wheeler.alliance.unm.edu
```

- Display your environment variables

```
$ env
```

- All the modules do is set the environment variables for different software

Environment Modules

- Let's load the R module so we can use it.

```
$ module list
```

```
$ module load r-3.5.0-intel-18.0.2-python2-mkl-r6lx6yy
```

(use tab complete !!)

```
$ R
```


Environment Modules

- Let's load the R module so we can use it.

`$ module list` (Again)

Will take a while!

`$ module avail` (to show all available modules)

`$ module spider <software name>` (to find a software package)

https://lmod.readthedocs.io/en/latest/010_user.html

Conda

- CARC staff have to install the software and create the environment modules you just saw.
- Anaconda provides an environment manager called `conda` that allows you to install the software you need into your home directory.
- Conda works with python, perl, R, and theoretically any language

Conda – Hands On

- Let's setup a local install of numpy

```
$ module load anaconda
```

```
$ conda -V
```

```
$ conda create -n numpy_py3 python=3.4 numpy
```

Wait a while – introduce yourselves to your neighbor... believe there is a reason we are doing this...

Conda – Hands On

- Let's setup a local install of numpy

```
$ module load anaconda
```

```
$ conda -V
```

```
$ conda create -n numpy_py3 python=3.4 numpy
```

Now we have defined a conda environment called `numpy_py3` and installed numpy in our home directories. At CARC we pay for the commercial intel version of mkl (math kernel library). Installing numpy now means you will get a performance boost for later installs. Sometimes other packages will select the non-intel mkl so lets get the right version.

Conda – Hands On

- Let's setup a local install of numpy

```
$ module load anaconda
```

```
$ conda -V
```

```
$ conda create -n numpy_py3 python=3.4 numpy
```

- Now we can load the environment

```
$ source activate numpy_py3
```

Conda – Hands On

- Let's see if numpy is available
- Create a file called test_numpy.py with the following contents:

```
try:
```

```
    import numpy
```

```
    print("imported numpy")
```

```
except ImportError:
```

```
    print("numpy is not installed")
```

Conda – Hands On

- Let's see if numpy is available
- Create a file called test_numpy.py with the following contents:

```
try:
    import numpy
    print("imported numpy")
except ImportError:
    print("numpy is not installed")
```

```
$ python test_numpy.py
```

Conda – Hands On

Conda just installs the software under ~/.conda

```
$ conda env list
```

```
$ source deactivate numpy_py3
```


Docker and Singularity

- Singularity allows you to load converted Docker images on HPC systems.
- Docker is not secure so singularity locks down access to the host machine.
- Docker containers allow you to configure a whole virtual operating system environment (e.g. your software needs Ubuntu but Wheeler runs CentOS).
- Convert your docker image to singularity and you can run the container. There is a “QuickByte” (short tutorial on the CARC website) on how to do this.

15 mins Break

- Snacks and Drinks in the Seminar room

Submitting Jobs and Embarrassingly Parallel Problems

Hands On

- Write a very short program that takes as its first argument a file path
- Does something with the contents, can be as simple as printing the file's contents to standard out.
- If you are not a programmer you can use the python program we have provided (it reads matrix sizes from the input file, creates a random matrix, and outputs the inverted matrix.)

Multiuser Systems and Batch Scheduling

- TORQUE (PBS)
- MAUI

Interactive Mode

- ONLY FOR DEBUGGING!

```
$ ssh <username>@wheeler.alliance.unm.edu
```

```
$ qsub -I --l nodes=2:ppn=8
```

PBS Variables Provide Information

- In interactive mode try:

```
$ echo $PBS_O_WORKDIR
```

```
$ echo $PBS_NODEFILE
```

```
$ cat $PBS_NODEFILE
```

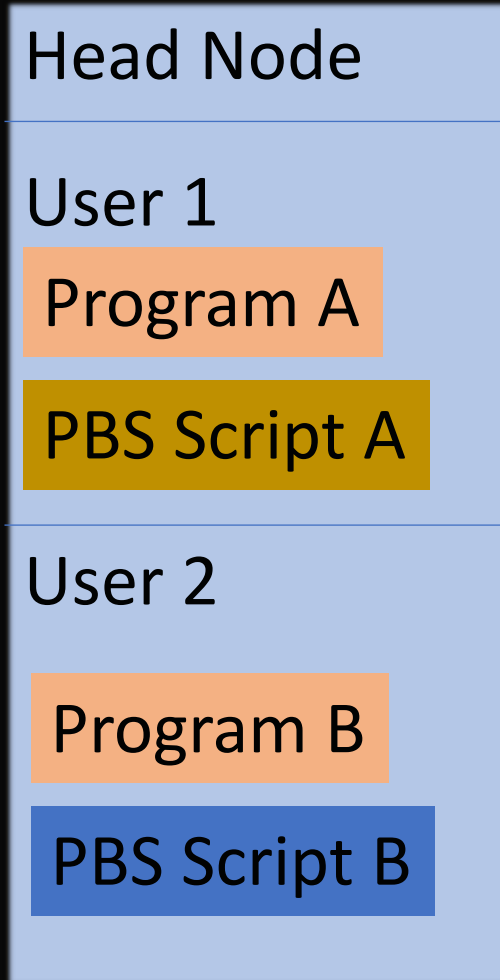
MAUI Scheduler

The scheduler looks at all the currently queued and running jobs and runs a backfill algorithm.

Smaller jobs in terms of number of CPUs and requested time are easier to schedule since there is more likely to be space for them.

However the longer a job is in the queue the more priority it gets. This way every job runs eventually.

Workflow



Compute Node 01

Compute Node 02

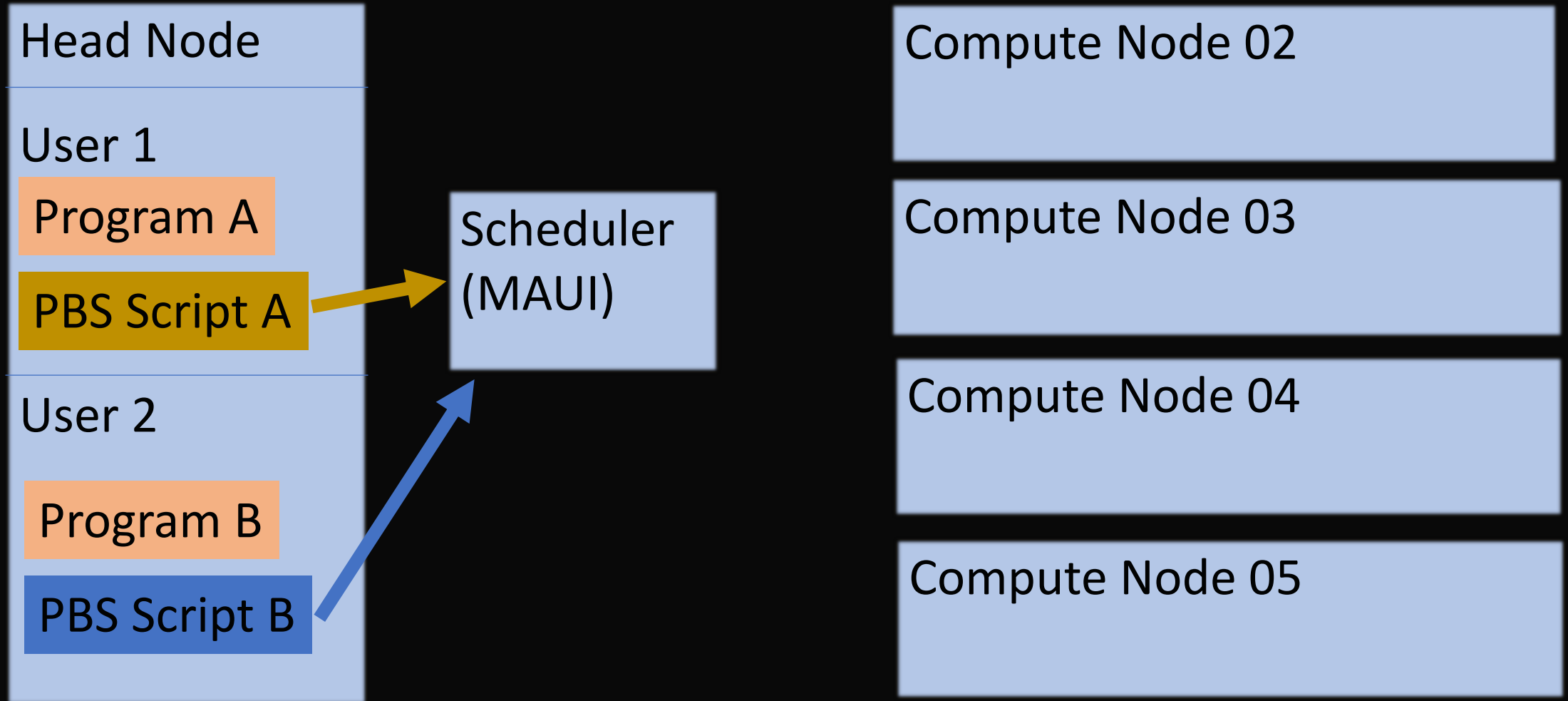
Compute Node 03

Compute Node 04

Compute Node 05

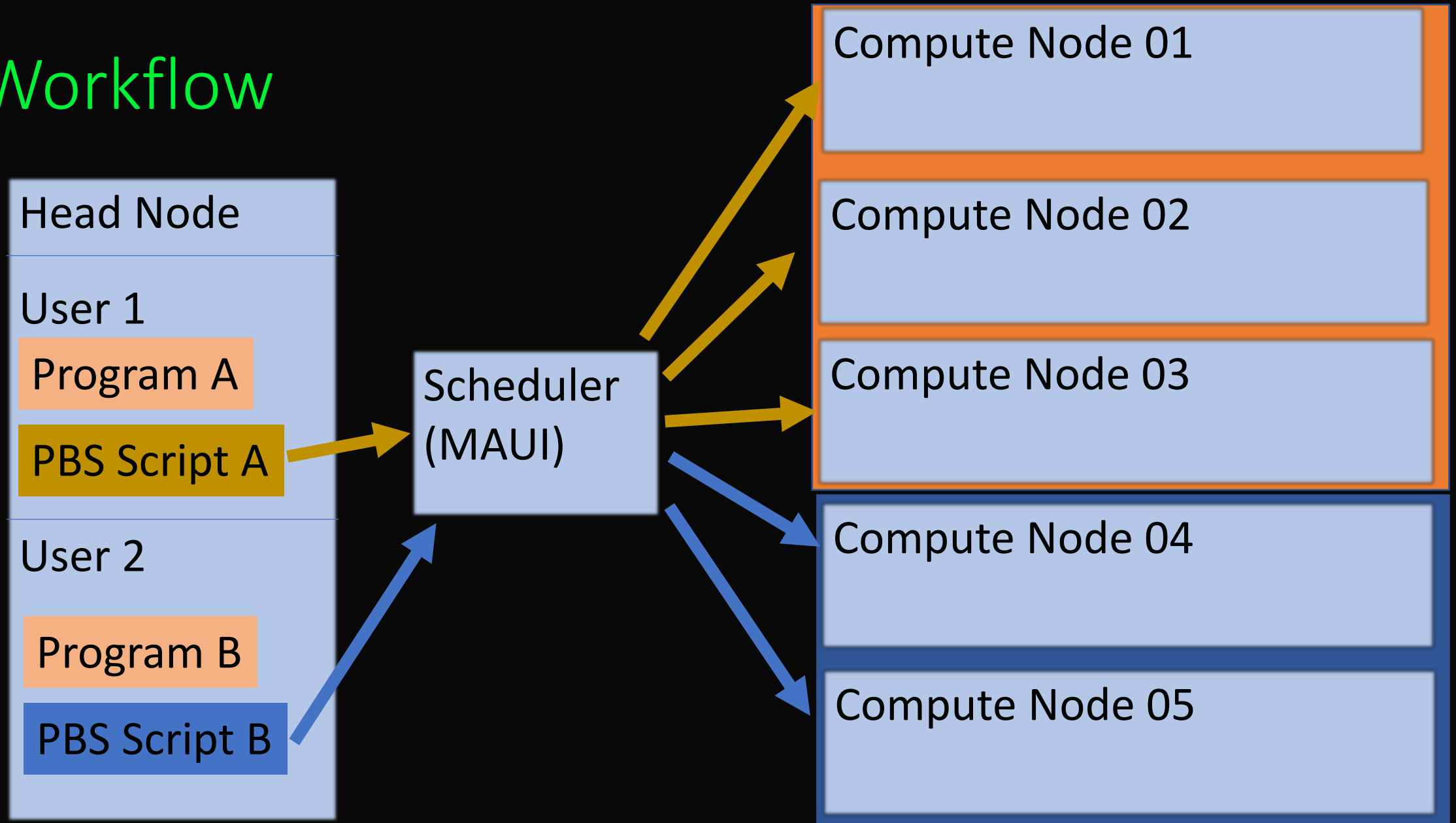
Shared filesystems – All nodes can access the same programs and write output

Workflow



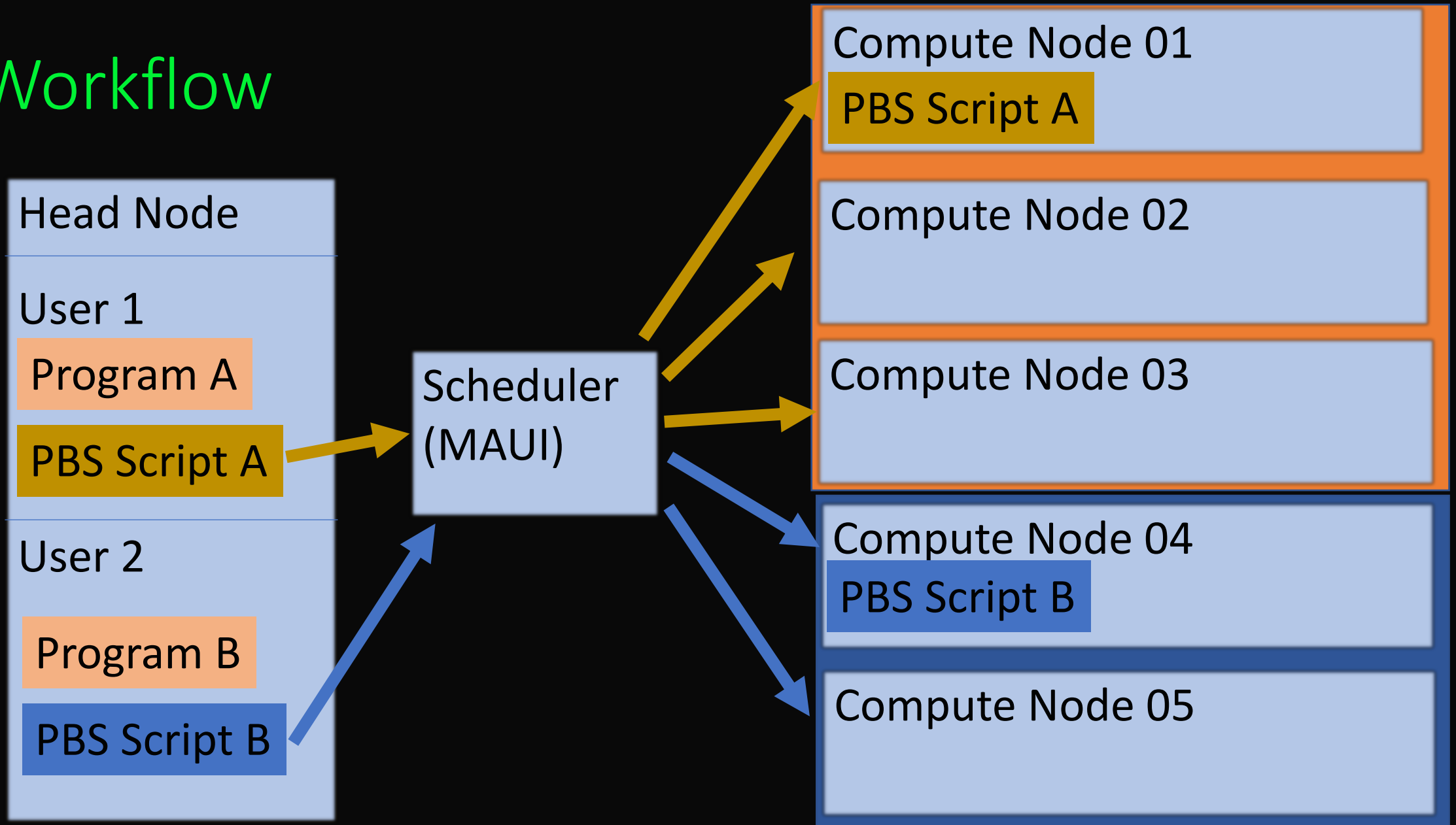
Shared filesystems – All nodes can access the same programs and write output

Workflow



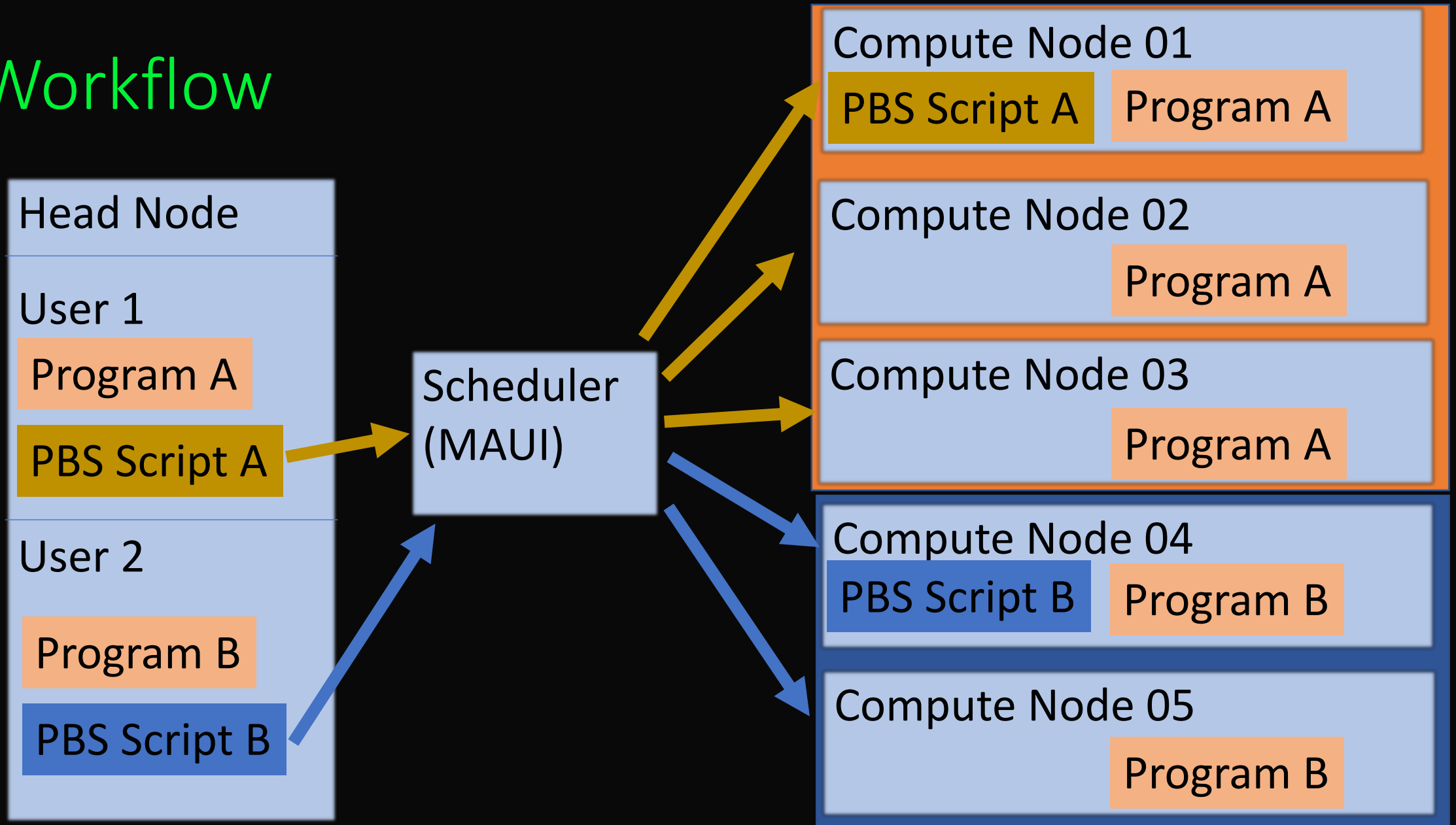
Shared filesystems – All nodes can access the same programs and write output

Workflow



Shared filesystems – All nodes can access the same programs and write output

Workflow




Shared filesystems – All nodes can access the same programs and write output

PBS Variables

- There are lots:

\$PBS_ENVIRONMENT	\$PBS_JOBID	\$PBS_MOMPORT	\$PBS_NP	\$PBS_O_HOME
\$PBS_O_LOGNAME	\$PBS_O_QUEUE	\$PBS_O_WORKDIR	\$PBS_VERSION	\$PBS_GPUFILE
\$PBS_JOBNAME	\$PBS_NODEFILE	\$PBS_NUM_NODES	\$PBS_O_HOST	\$PBS_O_MAIL
\$PBS_O_SERVER	\$PBS_QUEUE	\$PBS_VNODENUM	\$PBS_JOBCOOKIE	\$PBS_MICFILE
\$PBS_NODENUM	\$PBS_NUM_PPN	\$PBS_O_LANG	\$PBS_O_PATH.	\$PBS_O_SHELL
\$PBS_TASKNUM	\$PBS_WALLTIME			

Writing a Torque Batch Script

- Make a directory called ~/workshop
 - Edit a new file in the workshop directory, name it workshop_example.pbs
- 
- This REQUESTS time on the Core2 queue. We are asking for 2 nodes, and 2 cores on each node. We promise our job won't take more than 10 minutes. Email me when the begins, aborts, and ends (bae). Combine standard out and standard error into one file.

```
#!/bin/bash

#PBS -q Core2
#PBS -l nodes=2:ppn=2
#PBS -l walltime=00:10:00
#PBS -N ws_example
#PBS -S /bin/bash
#PBS -j oe
#PBS -m bae
#PBS -M my_email@unm.edu
#PBS -V
```

Writing a Torque Batch Script

```
#!/bin/bash

#PBS -q Core2
#PBS -l nodes=2:ppn=2
#PBS -l walltime=00:10:00
#PBS -N ws_example
#PBS -S /bin/bash
#PBS -j oe
#PBS -m bae
#PBS -M my_email@unm.edu
#PBS -V

echo $HOSTNAME
```

Everything that comes after the PBS preamble is executed on the first node you were allocated.

Managing your jobs

To submit your batch job:

```
$ qsub workshop_example.pbs
```

Writing a Torque Batch Script

Whatever is written to standard out and standard error is saved to <Job Name>.o<Job ID>
When the job ends.

If you want to see the output live, start your job with:

```
$ qsub -k oe <script_name.pbs>
```

Writing a Torque Batch Script

```
#!/bin/bash

#PBS -q Core2
#PBS -l nodes=2:ppn=2
#PBS -l walltime=00:10:00
#PBS -N ws_example
#PBS -S /bin/bash
#PBS -j oe
#PBS -m bae
#PBS -M my_email@unm.edu
#PBS -V

cat $PBS_NODEFILE
```

Everything that comes after the PBS preamble is executed on the first node you were allocated.

Writing a Torque Batch Script

```
#!/bin/bash

#PBS -q Core2
#PBS -l nodes=2:ppn=2
#PBS -l walltime=00:10:00
#PBS -N ws_example
#PBS -S /bin/bash
#PBS -j oe
#PBS -m bae
#PBS -M my_email@unm.edu
#PBS -V

$PBS_O_WORKDIR/workshop_example1
```

Everything that comes after the PBS preamble is executed on the first node you were allocated.

Job Arrays

```
#!/bin/bash

#PBS -q Core2
#PBS -l nodes=1:ppn=2
#PBS -l walltime=00:10:00
#PBS -N ws_example
#PBS -S /bin/bash
#PBS -j oe
#PBS -m bae
#PBS -M my_email@unm.edu
#PBS -V
#PBS -t 1-12%3

echo "$HOSTNAME - $PBS_ARRAYID"
```

-t allows you to schedule many jobs at once. -t 1-12%3 means run 12 jobs but only schedule 3 at a time.

Writing a Torque Batch Script

```
#!/bin/bash

#PBS -q Core2
#PBS -l nodes=1:ppn=2
#PBS -l walltime=00:10:00
#PBS -N ws_example
#PBS -S /bin/bash
#PBS -j oe
#PBS -m bae
#PBS -M my_email@unm.edu
#PBS -V
#PBS -t 1-12%3

$PBS_O_WORKDIR/ws_example1
```

Everything that comes after the PBS preamble is executed on the first node you were allocated.

Monitoring Jobs

Shows an overview of the queue status

```
$ qgrok
```

Lists all the jobs in the queue

```
$ qstat
```

Shows the resources requested by the jobs

```
$ qstat -a
```

Just the jobs submitted by a particular user

```
$ qstat -u <username>
```

-n shows the nodes being used by a running job

```
$ qstat -n -u <username>
```

-f gives detailed information about a particular job

```
$ qstat -f <job id>
```

Watch is a useful command that automatically updates the command that follows

```
$ watch qstat -n -u <username>
```

To see an estimate of how long it will be before a job starts and completes, enter:

```
$ /usr/local/maui/bin/showstart <job ID>
```

Writing a Torque Batch Script

```
#!/bin/bash

#PBS -q Core2
#PBS -l nodes=1:ppn=2
#PBS -l walltime=00:10:00
#PBS -N ws_example
#PBS -S /bin/bash
#PBS -j oe
#PBS -m bae
#PBS -M my_email@unm.edu
#PBS -V
#PBS -t 1-12%3

$PBS_O_WORKDIR/ws_example2 "my_inputfile$PBS_ARRAYID"
```

You can use this to execute your program with many different input files.
Here the program will run on files my_inputfile1, my_inputfile2, ... my_inputfile12

Hands on

Write some text files that will be input to the program you wrote earlier.
Give them all the same file extension: say .ws

If you want to use our canned matrix inverted program let us know.

Use the linux find command to send the paths to all the ws files to ws_example2

Play with the `--joblog` option to get some statistics on your job

Embarrassingly Parallel Problems

- Perfect case!
- All computation is independent so speedup is equal to the number of additional computers you throw at the problem.
- Especially good for generating lots of samples when you have a stochastic algorithm, simulation, or want to benchmark performance.

GNU Parallels – Input Driven

```
$ module load parallel-20170322-intel-18.0.2-4pa2ap6
```

```
$ find . -name "*.txt"
```

```
$ parallel echo ::: A B C ::: D E F
```

```
$ find . -name "*.txt" | parallel echo {}
```

```
$ find . -name "*.txt" | parallel echo {.}
```

```
$ find . -name "*.txt" | parallel echo {/.}
```

GNU Parallels – Monitoring Progress

- A logfile of the jobs completed so far can be generated with `--joblog`:

```
$ parallel --joblog $PBS_O_WORKDIR/job.log exit ::: 1 2 3 0  
$ cat $PBS_O_WORKDIR/job.log
```

- The log contains the job sequence, which host the job was run on, the start time and run time, how much data was transferred, the exit value, the signal that killed the job, and finally the command being run.

We are running “exit” so we can fake jobs that succeed and fail.

GNU Parallels – Resuming Jobs

- With a joblog GNU **parallel** can be stopped and later pickup where it left off. It is important that the input of the completed jobs is unchanged.
- Why would you want to do this...???

```
$ parallel --joblog $PBS_O_WORKDIR/job.log exit ::: 1 2 3 0
```

```
$ cat $PBS_O_WORKDIR/job.log
```

```
$ parallel --resume --joblog $PBS_O_WORKDIR/job.log exit ::: 1  
2 3 0 0 0
```

```
$ cat $PBS_O_WORKDIR/job.log
```

GNU Parallels – Resuming Jobs

- The previous command just ran the jobs that didn't finish
- This command reruns jobs that has a failing exit code

```
$ parallel --joblog $PBS_O_WORKDIR/job.log exit ::: 1 2 3 0
```

```
$ cat $PBS_O_WORKDIR/job.log
```

```
$ parallel -resume-failed --joblog $PBS_O_WORKDIR/job.log exit ::: 1 2  
3 0 0 0
```

```
$ cat $PBS_O_WORKDIR/job.log
```

We are running “exit” so we can fake jobs that succeed and fail.

GNU Parallels

- GNU Parallel is what you should be using to run many experiments, solve many independent instances of a problem.
- If you have 1000 input files and 100 CPUs allocated parallel will do all the scheduling for you to process those files.

GNU Parallels

- If you have 1000 input files and 100 CPUs allocated parallel will do all the scheduling for you to process those files.
- Remember though: Torque assigns you resources, parallels makes use of them. You have to use
`--sshloginfile $PBS_NODEFILE`
- To be sure parallels is using resources you were actually allocated

GNU Parallels – Environments

Recall that software may require particular environments. GNU Parallel by itself loses the environment in which it was called.

Use `env_parallel` to keep the environment. Need to tell it what shell you are using.

```
source `which env_parallel.bash`
```

Then you can use `env_parallel` exactly like `parallel`.

GNU Parallels

```
# Create a temporary unique directory in which to store the summary output for each job
TEMP_DIR=$(mktemp -d -p $PBS_0_WORKDIR)

# Setup Gurobi solver environment
module load parallel-20170322-intel-18.0.2-4pa2ap6
module load gurobi
module load anaconda
source activate gurobi

source `which env_parallel.bash`

# Use find to make a list of all the .graph files to pass to the integer program solver.
# Divide the work up among compute nodes using the GNU parallel tool. Use a local /tmp work directory.
# ":::: -" reads from stdin (find ... *.graph) to {1}, ":::: $EPSILON_VALUES_PATH" reads from the
epsilon parameter file to {2}, {1/.} fetches the input base filename only
find $GRAPH_INPUT_DIR -name '*.graph' | env_parallel --jobs 8 --joblog
$PBS_0_WORKDIR/ip_progress.joblog --resume-failed --sshloginfile $PBS_NODEFILE --workdir $(mktemp -d)
"python $PBS_0_WORKDIR/ip/solve_ip.py {1} $IP_METHOD {2} $TEMP_DIR/{1/.}.txt $TIME_LIMIT
$SOLUTION_OUTPUT_DIR" :::: - $EPSILON_VALUES_PATH
```

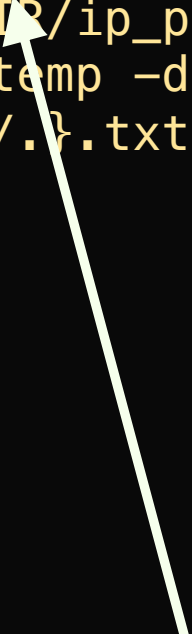
GNU Parallels

```
find $GRAPH_INPUT_DIR -name '*.graph' | env_parallel --jobs 8  
$PBS_0_WORKDIR/ip_progress.joblog --resume-failed --sshloginfile $PBS_NODEFILE --  
workdir $(mktemp -d) "python $PBS_0_WORKDIR/ip/solve_ip.py {1} $IP_METHOD {2}  
$TEMP_DIR/{1/.}.txt $TIME_LIMIT $SOLUTION_OUTPUT_DIR" :::: - $EPSILON_VALUES_PATH
```

Let's try to parse this command together...

GNU Parallels

```
find $GRAPH_INPUT_DIR -name '*.graph' | env_parallel --jobs 8  
$PBS_0_WORKDIR/ip_progress.joblog --resume-failed --sshloginfile $PBS_NODEFILE --  
workdir $(mktemp -d) "python $PBS_0_WORKDIR/ip/solve_ip.py {1} $IP_METHOD {2}  
$TEMP_DIR/{1/.}.txt $TIME_LIMIT $SOLUTION_OUTPUT_DIR" :::: - $EPSILON_VALUES_PATH
```



Pass the paths to all the files with a graph extension in the directory specified in the user shell variable \$GRAPH_INPUT_DIR to env_parallel.

GNU Parallels

```
find $GRAPH_INPUT_DIR -name '*.graph' | env_parallel --jobs 8 --joblog  
$PBS_0_WORKDIR/ip_progress.joblog --resume-failed --sshloginfile $PBS_NODEFILE --  
workdir $(mktemp -d) "python $PBS_0_WORKDIR/ip/solve_ip.py {1} $IP_METHOD {2}  
$TEMP_DIR/{1/.}.txt $TIME_LIMIT $SOLUTION_OUTPUT_DIR" ::: - $EPSILON_VALUES_PATH
```

The piped input is mapped to the first parameter to parallel referred to with "-".

Pipes the paths to all the files with a graph extension in the directory specified in the user shell variable \$GRAPH_INPUT_DIR to env_parallel.

GNU Parallels

... and referred to with {1}.

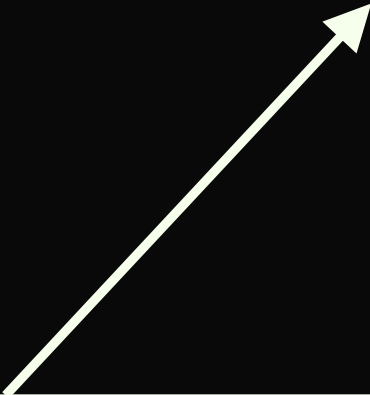
```
find $GRAPH_INPUT_DIR -name '*.graph' | env_parallel --jobs 8 --joblog  
$PBS_0_WORKDIR/ip_progress.joblog --resume-failed --sshloginfile $PBS_NODEFILE --  
workdir $(mktemp -d) "python $PBS_0_WORKDIR/ip/solve_ip.py {1} $IP_METHOD {2}  
$TEMP_DIR/{1/.}.txt $TIME_LIMIT $SOLUTION_OUTPUT_DIR" :::- $EPSILON_VALUES_PATH
```

The diagram consists of two white arrows. The first arrow originates from the box containing the text "... and referred to with {1}." and points to the "{1}" placeholder in the command "python \$PBS_0_WORKDIR/ip/solve_ip.py {1} \$IP_METHOD {2}". The second arrow originates from the box containing the text "The piped input is mapped to the first parameter to parallel referred to with \"-\"." and points to the "-" character in the command "python \$PBS_0_WORKDIR/ip/solve_ip.py {1} \$IP_METHOD {2} \$TEMP_DIR/{1/.}.txt \$TIME_LIMIT \$SOLUTION_OUTPUT_DIR" :::- \$EPSILON_VALUES_PATH".

The piped input is mapped to the first parameter to parallel referred to with "-".

GNU Parallels

```
find $GRAPH_INPUT_DIR -name '*.graph' | env_parallel --jobs 8 --joblog  
$PBS_0_WORKDIR/ip_progress.joblog --resume-failed --sshloginfile $PBS_NODEFILE --  
workdir $(mktemp -d) "python $PBS_0_WORKDIR/ip/solve_ip.py {1} $IP_METHOD {2}  
$TEMP_DIR/{1/.}.txt $TIME_LIMIT $SOLUTION_OUTPUT_DIR" :::: - $EPSILON_VALUES_PATH
```



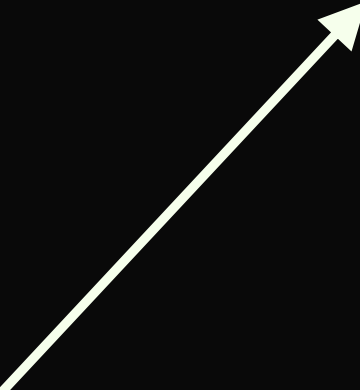
The second set of parameters is read from a file. In this example the path to the values is specified by the user variable `$EPSILON_VALUES_PATH`.

GNU Parallels

... and referred to with {2}.



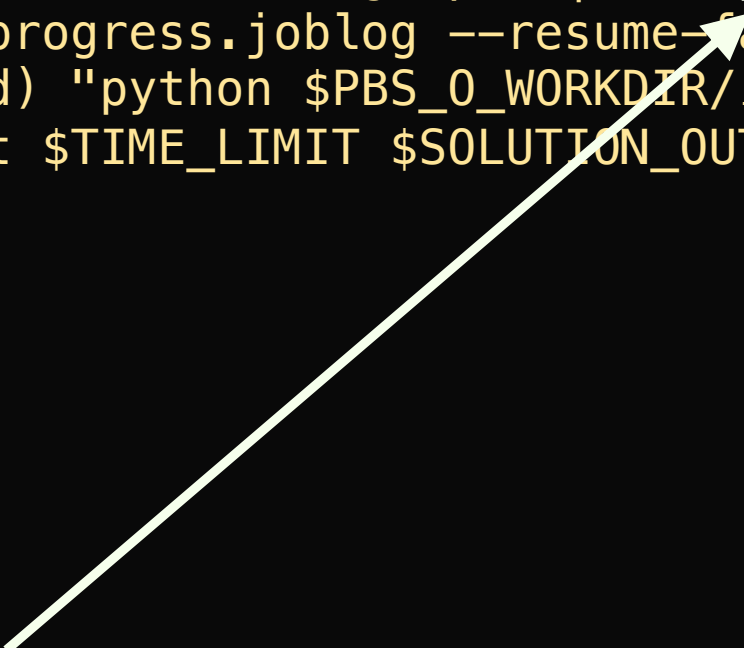
```
find $GRAPH_INPUT_DIR -name '*.graph' | env_parallel --jobs 8 --joblog  
$PBS_0_WORKDIR/ip_progress.joblog --resume-failed --sshloginfile $PBS_NODEFILE --  
workdir $(mktemp -d) "python $PBS_0_WORKDIR/ip/solve_ip.py {1} $IP_METHOD {2}  
$TEMP_DIR/{1/.}.txt $TIME_LIMIT $SOLUTION_OUTPUT_DIR" :::: - $EPSILON_VALUES_PATH
```



The second set of parameters is read from a file. In this example the path to the values is specified by the user variable `$EPSILON_VALUES_PATH`.

GNU Parallels

```
find $GRAPH_INPUT_DIR -name '*.graph' | env_parallel --jobs 8 --joblog  
$PBS_0_WORKDIR/ip_progress.joblog --resume-failed --sshloginfile $PBS_NODEFILE --  
workdir $(mktemp -d) "python $PBS_0_WORKDIR/ip/solve_ip.py {1} $IP_METHOD {2}  
$TEMP_DIR/{1/.}.txt $TIME_LIMIT $SOLUTION_OUTPUT_DIR" :::: - $EPSILON_VALUES_PATH
```



We are using `env_parallel` which passes the current shell environment to the jobs. In this example the user code uses shell variables.

GNU Parallels

```
find $GRAPH_INPUT_DIR -name '*.graph' | env_parallel --jobs 8 --joblog  
$PBS_0_WORKDIR/ip_progress.joblog --resume-failed --sshloginfile $PBS_NODEFILE --  
workdir $(mktemp -d) "python $PBS_0_WORKDIR/ip/solve_ip.py {1} $IP_METHOD {2}  
$TEMP_DIR/{1/.}.txt $TIME_LIMIT $SOLUTION_OUTPUT_DIR" ::: - $EPSILON_VALUES_PATH
```



Specifies the number of jobs to run on each node.

GNU Parallels

```
find $GRAPH_INPUT_DIR -name '*.graph' | env_parallel --jobs 8 --joblog  
$PBS_0_WORKDIR/ip_progress.joblog --resume-failed --sshloginfile $PBS_NODEFILE --  
workdir $(mktemp -d) "python $PBS_0_WORKDIR/ip/solve_ip.py {1} $IP_METHOD {2}  
$TEMP_DIR/{1/.}.txt $TIME_LIMIT $SOLUTION_OUTPUT_DIR" :::: - $EPSILON_VALUES_PATH
```



Record the progress to a file.

GNU Parallels

```
find $GRAPH_INPUT_DIR -name '*.graph' | env_parallel --jobs 8 --joblog  
$PBS_0_WORKDIR/ip_progress.joblog --resume-failed --sshloginfile $PBS_NODEFILE --  
workdir $(mktemp -d) "python $PBS_0_WORKDIR/ip_solve_ip.py {1} $IP_METHOD {2}  
$TEMP_DIR/{1/.}.txt $TIME_LIMIT $SOLUTION_OUTPUT_DIR" :::: - $EPSILON_VALUES_PATH
```



Record the progress to a file.

...and tell parallels to rerun any failed jobs listed in the joblog (those where the exit value not equal to 0).

GNU Parallels


```
find $GRAPH_INPUT_DIR -name '*.graph' | env_parallel --jobs 8 --joblog  
$PBS_0_WORKDIR/ip_progress.joblog --resume-failed --sshloginfile $PBS_NODEFILE --  
workdir $(mktemp -d) "python $PBS_0_WORKDIR/ip/solve_ip.py {1} $IP_METHOD {2}  
$TEMP_DIR/{1/.}.txt $TIME_LIMIT $SOLUTION_OUTPUT_DIR" :::: - $EPSILON_VALUES_PATH
```



Tell parallels which nodes we were
allocated to run our jobs.

GNU Parallels

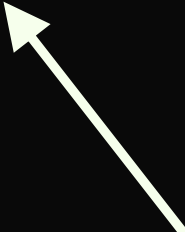
```
find $GRAPH_INPUT_DIR -name '*.graph' | env_parallel --jobs 8 --joblog  
$PBS_0_WORKDIR/ip_progress.joblog --resume-failed --sshloginfile $PBS_NODEFILE --  
workdir $(mktemp -d) "python $PBS_0_WORKDIR/ip/solve_ip.py {1} $IP_METHOD {2}  
$TEMP_DIR/{1/.}.txt $TIME_LIMIT $SOLUTION_OUTPUT_DIR" :::: - $EPSILON_VALUES_PATH
```



--workdir is set to a temporary directory.

GNU Parallels

```
find $GRAPH_INPUT_DIR -name '*.graph' | env_parallel --jobs 8 --joblog  
$PBS_0_WORKDIR/ip_progress.joblog --resume-failed --sshloginfile $PBS_NODEFILE --  
workdir $(mktemp -d) "python $PBS_0_WORKDIR/ip/solve_ip.py {1} $IP_METHOD {2}  
$TEMP_DIR/{1/.}.txt $TIME_LIMIT $SOLUTION_OUTPUT_DIR" ::: - $EPSILON_VALUES_PATH
```



The command that will be run in each parallel job. This program takes 6 arguments. Parallel will generate a job for all combinations of input parameter {1} and {2}. Argument 4 specifies an output path based on the input file name {1}. {1/.} gets just the input file basename.

GNU Parallels


```
find $GRAPH_INPUT_DIR -name '*.graph' | env_parallel --jobs 8 --joblog  
$PBS_0_WORKDIR/ip_progress.joblog --resume-failed --sshloginfile $PBS_NODEFILE --  
workdir $(mktemp -d) "python $PBS_0_WORKDIR/ip/solve_ip.py {1} $IP_METHOD {2}  
$TEMP_DIR/{1/.}.txt $TIME_LIMIT $SOLUTION_OUTPUT_DIR" :::: - $EPSILON_VALUES_PATH
```



Notice “::::” rather than “:::”. Four colons tells parallels that a list of parameters comes next.

GNU Parallels

```
find $GRAPH_INPUT_DIR -name '*.graph' | env_parallel --jobs 8 --joblog  
$PBS_0_WORKDIR/ip_progress.joblog --resume-failed --sshloginfile $PBS_NODEFILE --  
workdir $(mktemp -d) "python $PBS_0_WORKDIR/ip/solve_ip.py {1} $IP_METHOD {2}  
$TEMP_DIR/{1/.}.txt $TIME_LIMIT $SOLUTION_OUTPUT_DIR" :::: - $EPSILON_VALUES_PATH
```



So in this example, parallels will spawn a job for every combination of input file and value in the `$EPSILON_VALUES_PATH` file.

Parallels records its progress in a joblog file so it can pick up where it left off if the torque job runs out of time before all the jobs are complete, and the torque job needs to be resubmitted.

15 min break

MPI

Some parallel computations need the programs running on each core to communicate with each other.

The Message Passing Interface (MPI) supports this kind of communication.

We will look at a simple example using MPI.

Example Problem: Calculate π

Serial Calculation of π

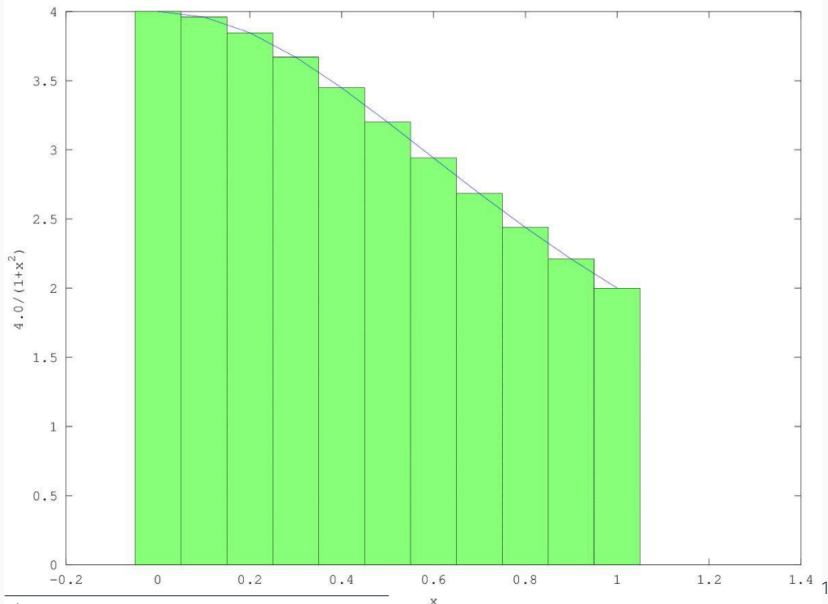
$$\int_0^1 \frac{4}{1+x^2} dx = \pi$$

And can be numerically approximated with:

$$\sum_{i=0}^N \frac{4}{1+x_i^2} \Delta x \approx \pi$$

As Δx gets smaller and N gets larger the approximation converges on π

Serial Program to Calculate π



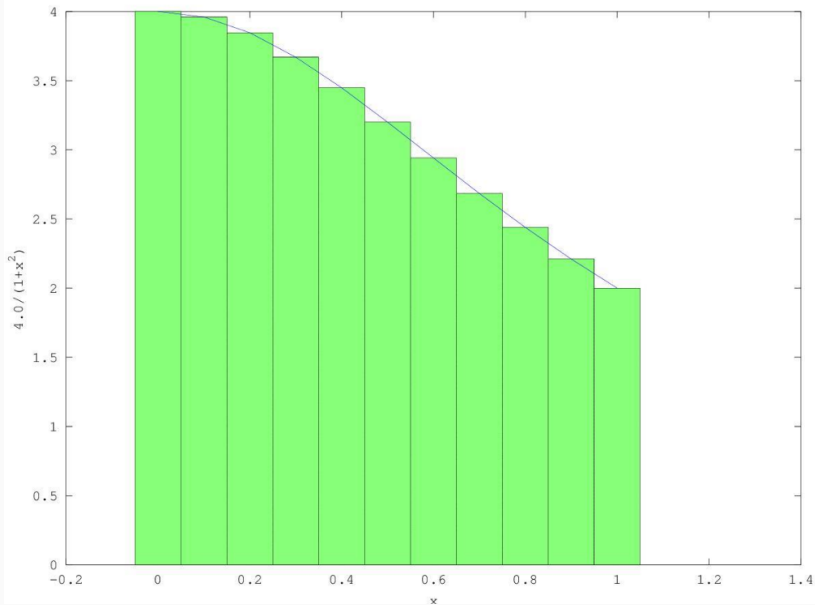
Serial Program to Calculate π

calcPiSerial.py

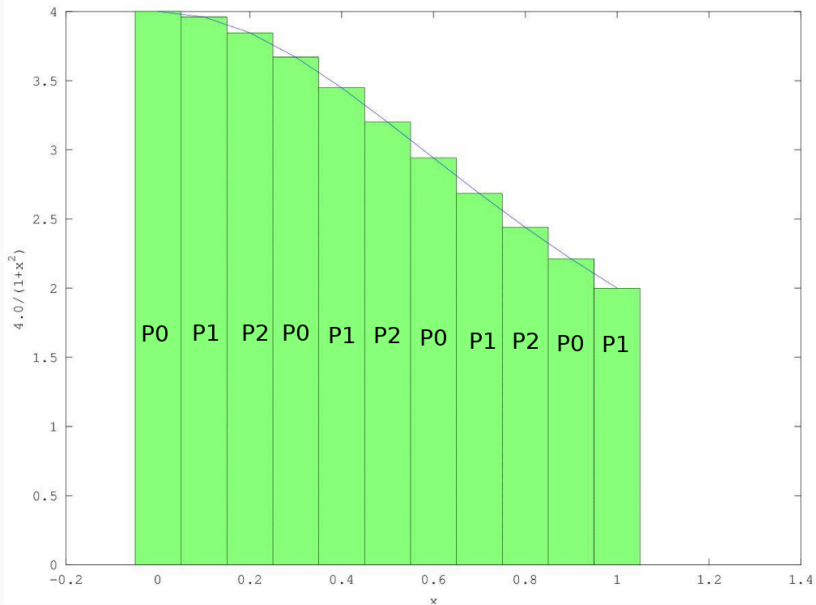
```
import time
def Pi(num_steps): # Function to calculate pi
    step = 1.0/num_steps
    sum = 0
    for i in xrange(num_steps):
        x = (i+0.5)*step
        sum = sum + 4.0/(1.0+x*x)
    pi = step * sum
    return pi

if __name__ == '__main__': # Main function
    start = time.time() # Start timing
    num_steps = 100000000
    pi=Pi(num_steps)
    end = time.time() #Stop timing
    # If this is the root process print the result
    print "Pi=%f(calculated_in_%fsecs)" %(pi, end-start)
```

Serial Program to Calculate π



Parallel Program to Calculate π



LMod Environment Modules

As before we will need to load an environment module to get access to the MPI libraries we need.

To load the MPI environment enter:

```
module load openmpi-3.1.1-intel-18.0.2-vde2j7x
```

Other useful commands:

Find modules: **module spider <search string>**

List all modules on the system: **module avail**

List loaded modules: **module list**

Unload a module: **module unload <module name>**

Anaconda Environment

We will create a conda environment that gives us access to mpi4py.

To create an environment called “wheeler_mpi_py2” that includes the python packages for numerical computing, scientific computing, and MPI enter:

```
module load anaconda to load the anaconda module  
conda create --name wheeler_mpi_py2 python=2 mpi4py numpy scipy
```

Once the python packages have finished installing enter:

```
source activate wheeler_mpi_py2
```

Parallel Program to Calculate π

calcPiMPI.py

```
from mpi4py import MPI
import time

# Get MPI variables
comm = MPI.COMM_WORLD      # Communication framework
root = 0                   # Root process
rank = comm.Get_rank()     # Rank of this process
num_procs = comm.Get_size() # Total number of processes

# Distributed function to calculate pi
def Pi(num_steps):
    step = 1.0/num_steps
    sum = 0
    for i in xrange(rank, num_steps, num_procs):
        x = (i+0.5)*step
        sum = sum + 4.0/(1.0+x*x)
    mypi = step * sum
    pi = comm.reduce(mypi, MPI.SUM, root)
    return pi

# Main function
if __name__ == '__main__':
    start = time.time() # Start timing
    num_steps = 10000000
    # Broadcast number of steps to use to the other processes
    comm.bcast(num_steps, root);
    pi=Pi(num_steps)
    end = time.time() #Stop timing
    # If this is the root process print the result
    if (rank==root): print "Pi=%f_(calculated_in_%f_secs)" %(pi, end-start)
```

Torque Job Scheduler

Write and compile your code on the wheeler head node. Run your programs on the compute nodes.

We use the **torque** system to schedule jobs on the compute nodes.

Some useful scheduler commands:

Show current jobs: **qstat -a**

Queue information: **qstat -q** or **qgrok**

Show jobs of a particular user: **qstat -u <username>**

Submit a job: **qsub <pbs² script>**

Delete a job: **qdel <job ID>**

²portable batch system

PBS Script

calc_pi.pbs

```
#!/bin/bash
```

```
#PBS -l nodes=2:ppn=8
```

```
#PBS -l walltime=00:05:00
```

```
#PBS -N calc_pi
```

```
#PBS -S /bin/bash
```

```
#PBS -j oe
```

```
#PBS -M youremailaddress@unm.edu
```

```
#PBS -V
```

```
module load openmpi-3.1.1-intel-18.0.2-vde2j7x
```

```
module load anaconda
```

```
source activate wheeler_mpi_py2
```

Need to add to mpirun arguments:—machinefile \$PBS_NODEFILE

```
mpirun -n $PBS_NP python calcPiMPI.py
```

File systems

Home directory. Limited space. Limited speed. Store your code here.
Backed up. Path: `~`

Scratch. Fast. Lot's of space. Not-backed up. Store data that you can regenerate here.
Path: `~/wheeler_scratch`

Temp (on Wheeler these are RAM drives). Very fast but decreases available RAM. Path: `/tmp`

Final Thoughts

During this tutorial you have learned about:

- Environment modules
- Writing PBS scripts and submission of jobs to the TORQUE/MAUI system
- Embarrassingly parallel tasks and GNU Parallel
- More tightly coupled parallelism with the Message Passing Interface (MPI)
- Scratch vs user space

For assistance send an email to help@carc.unm.edu