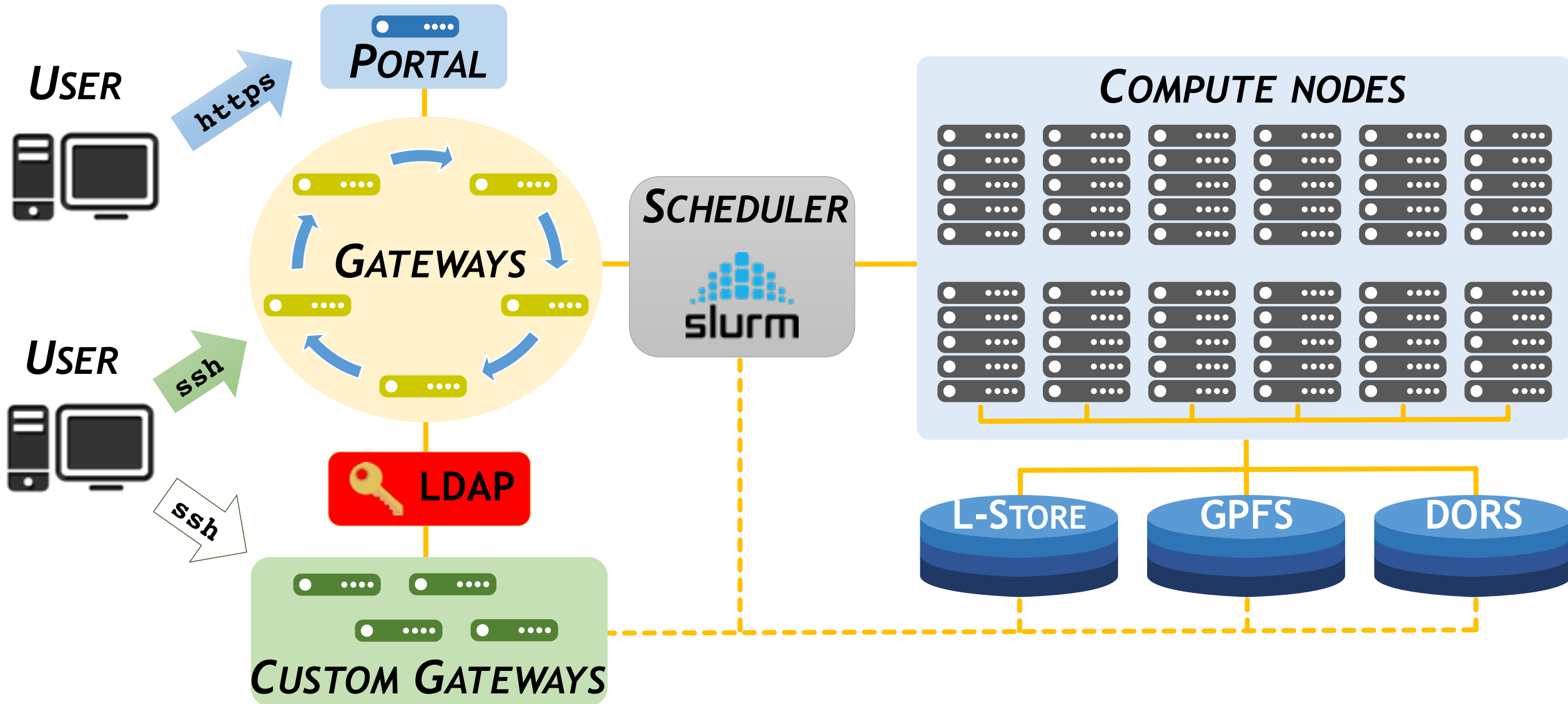# Using MPI at ACCRE

# ACCRE OVERVIEW

# ACCRE Overview

- ACCRE - Advanced Computing Center For Research and Education
- Centralized computing infrastructure for Vanderbilt researchers
- Operates as a co-op in which researchers share hardware
- ~10k CPU cores
- ~200 GPUs
- ~10PB disk storage + tape backups
- Optimized Scientific Software Stack
- Batch Job Scheduler
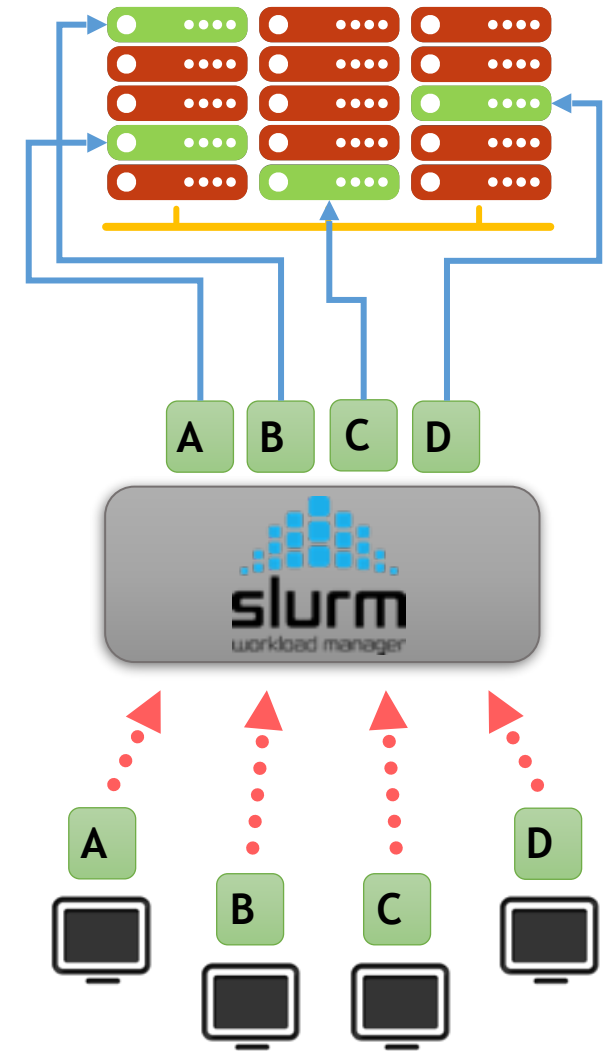- Interactive resources (Jupyter, etc.)
- Staff of ~10

- Using your own hardware:

  - can use all resources immediately

  - have to set up software, system, and networking yourself

  - full administrative access (root)

- Using ACCRE:

  - must schedule resource requirements

  - can "burst" to use more resources than you own

  - dedicated staff maintain system and software stack

  - no administrative access (regular user)

# ACCRE ARCHITECTURE

**User**

https

**Portal**

**Gateways**

ssh

ssh

LDAP

**Custom Gateways**

**Scheduler**

slurm

**Compute nodes**

L-Store    GPFS    DORS

**1** Execute user's workloads in the right priority order

**2** Provide requested resources on compute nodes

**3** Optimize cluster utilization

- Different Memory Configurations and CPU Core Counts
  - Nodes with 64GB, 128GB, 192GB, 256GB, and 384GB
  - Between 8 and 32 CPU-cores per node
- Different Intel CPU Architecture Families
  - Variable clock speed, L1/2/3 Cache Memory
  - Additional Instruction Sets on Newer CPUs
- Specialized Accelerated Nodes
  - Nvidia 4x GPU Nodes (Maxwell, Pascal, Turing)

- Heterogeneity is fine for "embarrassingly" parallel tasks
  - Don't care if jobs start at the same time
  - Don't care if some jobs are slower than others
- Not so great for distributed simulations
  - Need to synchronize between iterations
  - Limited to speed of slowest nodes
- ACCRE cluster users can choose to run on all architectures or specific architectures

- Westmere (~2010)
  - 64-bit extensions, MMX, SSE, SSE2, SSE3, SSSE3, SSE4.1, SSE4.2, POPCNT, AES and PCLMUL

- Sandy Bridge (~2011)
  - 64-bit extensions, MMX, SSE, SSE2, SSE3, SSSE3, SSE4.1, SSE4.2, POPCNT, **AVX**, AES and PCLMUL

- Haswell (~2014)
  - 64-bit extensions, **MOVBE**, MMX, SSE, SSE2, SSE3, SSSE3, SSE4.1, SSE4.2, POPCNT, AVX, **AVX2**, AES, PCLMUL, **FSGSBASE**, **RDRND**, **FMA**, **BMI**, **BMI2** and **F16C**

- Skylake (~2015)
  - 64-bit extensions, MOVBE, MMX, SSE, SSE2, SSE3, SSSE3, SSE4.1, SSE4.2, POPCNT, **PKU**, AVX, AVX2, AES, PCLMUL, FSGSBASE, RDRND, FMA, BMI, BMI2, F16C, **RDSEED**, **ADCX**, **PREFETCHW**, **CLFLUSHOPT**, **XSAVEC**, **XSAVES**, **AVX512F**, **CLWB**, **AVX512VL**, **AVX512BW**, **AVX512DQ** and **AVX512CD**

- Cascade Lake (~2019)
  - 64-bit extensions, MOVBE, MMX, SSE, SSE2, SSE3, SSSE3, SSE4.1, SSE4.2, POPCNT, PKU, AVX, AVX2, AES, PCLMUL, FSGSBASE, RDRND, FMA, BMI, BMI2, F16C, RDSEED, ADCX, PREFETCHW, CLFLUSHOPT, XSAVEC, XSAVES, AVX512F, CLWB, AVX512VL, AVX512BW, AVX512DQ, AVX512CD and **AVX512VNNI**
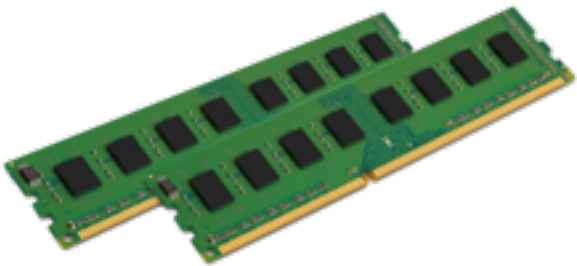
# ACCRE CLUSTER COMPUTE NODES

## Regular nodes

*Dual multicore CPUs*

*Random Access Memory*

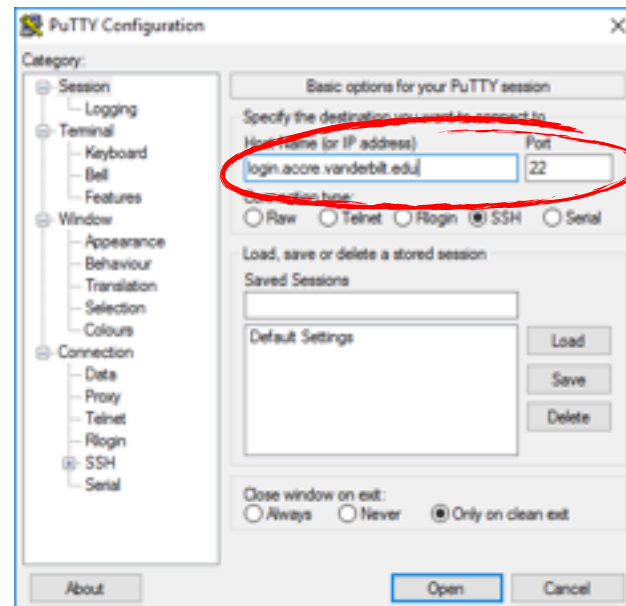| Family | No. of cores | RAM / GB | No. of nodes |
|--------|:------------:|:--------:|:------------:|
| Skylake | 16 | 256 | 41 |
| | 24 | 128 | 52 |
| Haswell | 12 | 128 | 41 |
| | 16 | 128 | 120 |
| | | 256 | 50 |
| Sandy Bridge | 12 | 64 | 31 |
| | | 96 | 2 |
| | | 128 | 193 |
| | | 256 | 4 |
| | 16 | 128 | 3 |
| Westmere | 8 | 128 | 22 |
| | 12 | 48 | 16 |
| **Total** | **8,292** | **82,432** | **575** |

Newer ↑ Older

# USING ACCRE

Users can log into the cluster with a Secure Shell (**ssh**) client.

From a terminal:    `ssh` *vunetid*`@login.accre.vanderbilt.edu`



**PuTTY**

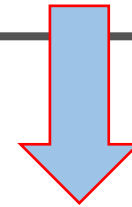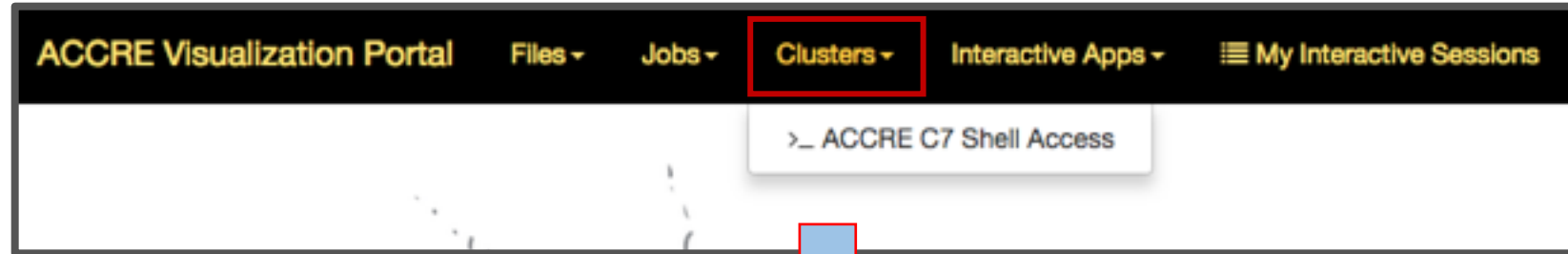**Bash on Windows**

Full Ubuntu-based Bash shell

To install:

*https://goo.gl/tAsj8U*

⚠ Windows 10 only

# SHELL ACCESS TO THE CLUSTER – PORTAL

Opens in a new browser tab

**GPFS** and **DORS** are distributed parallel filesystems that allow users to get access to the same set of directories on <u>all nodes</u> and <u>all gateways</u> on the cluster.

| | Nightly backup | Included with account | For purchase |
|---|:---:|:---:|:---:|
| **GPFS** | | | |
| /home | ✔ | ✔ | ✖ |
| /scratch | ✖ | ✔ | ✔ |
| /data | ✔ | ✖ | ✔ |

**DORS**

/dors     Managed by *Center for Structural Biology*, supported by ACCRE.

Provides easy access to data from both desktops and cluster.

**QUOTA:** When exceeded the user receives a warning message.
Usage has to return below the quota within the **GRACE PERIOD**.

**LIMIT:** Cannot be exceeded.
Automatically set to the actual quota usage when grace period expires.

**GPFS**

/home

/scratch

/data

| | Data size | | Number of files | | GRACE PERIOD |
|---|---|---|---|---|---|
| | QUOTA | LIMIT | QUOTA | LIMIT | |
| /home | 15 GB | 20 GB | 200,000 | 300,000 | 7 days |
| /scratch | 50 GB | 200 GB | 200,000 | 1,000,000 | 14 days |
| /data | | Can be purchased at 1 TB increments. | | | |

Note that groups may also purchase additional scratch quota in 1 TB increments.

**How can I check my current quota usage?**

`accre_storage`

- Shows the current usage for all quotas associated with the user.

```
                           Space                       Files
                 Usage     Quota     Limit    Usage    Quota    Limit

Home (user):     12.41G    15G       20G      120304   200000   300000
Scratch (user):  36.23G    50G       200G     180276   200000   1000000
```

- **Demo Cluster Login**

- **Demo Portal Usage**

- **Demo ACCRE Storage**

# A Simple Example Program with Lmod

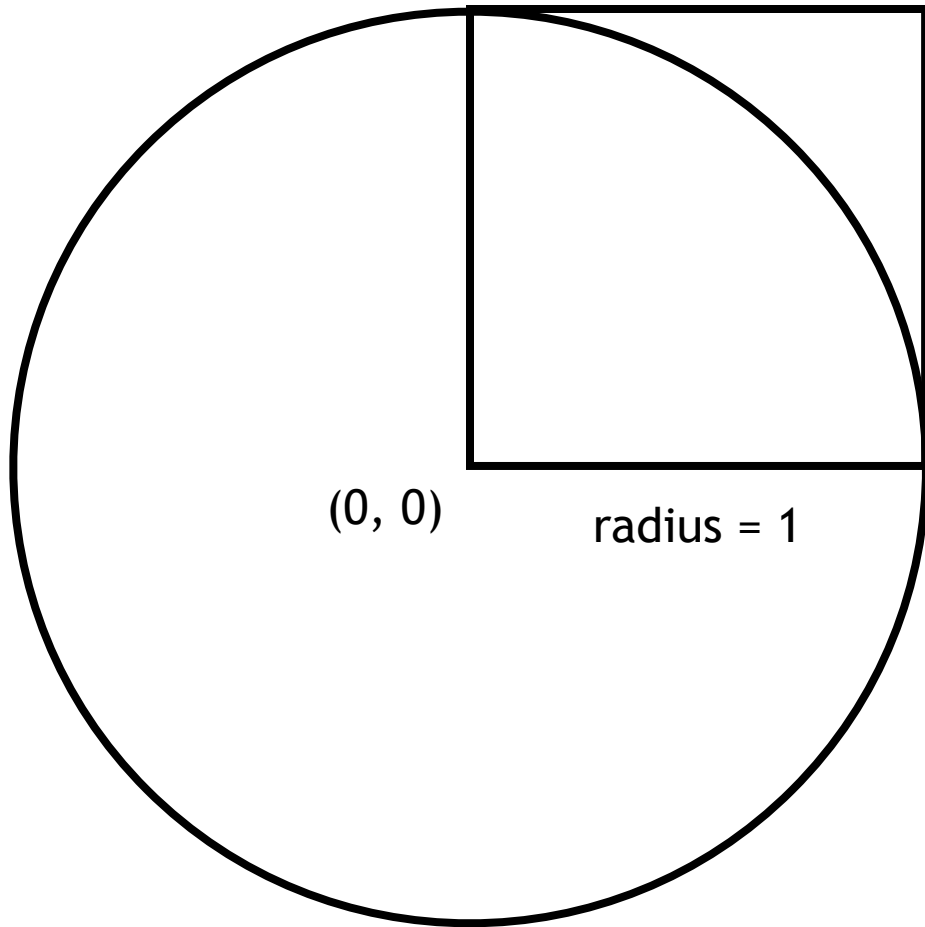# An Inefficient Way to Calculate Pi

(1, 1)

(0, 0)

radius = 1

$$\frac{\text{Area Quarter Circle}}{\text{Area Square}} = \frac{\text{pi} / 4}{1} = \text{pi} / 4$$

Pick lots of random coordinates from (0, 0) to (1, 1)

$$\frac{\text{coordinates in the quarter circle}}{\text{coordinates in the square}} \approx \frac{\text{Area Quarter Circle}}{\text{Area Square}} = \text{pi} / 4$$

# An Inefficient Way to Calculate Pi

(1, 1)

$$\frac{\text{coordinates in the quarter circle}}{\text{coordinates in the square}} \approx \text{pi} / 4$$

(0, 0)

radius = 1

```
attempts = 0
hits = 0

for i = 0; i < N; i++
    x = random(0, 1)
    y = random(0, 1)
    attempts++
    if x*x + y*y < 1
        hits++

return 4 * hits / attempts
```

Researchers should focus on the **science** behind the software they use.

Easily search available software

Effortless access to selected software
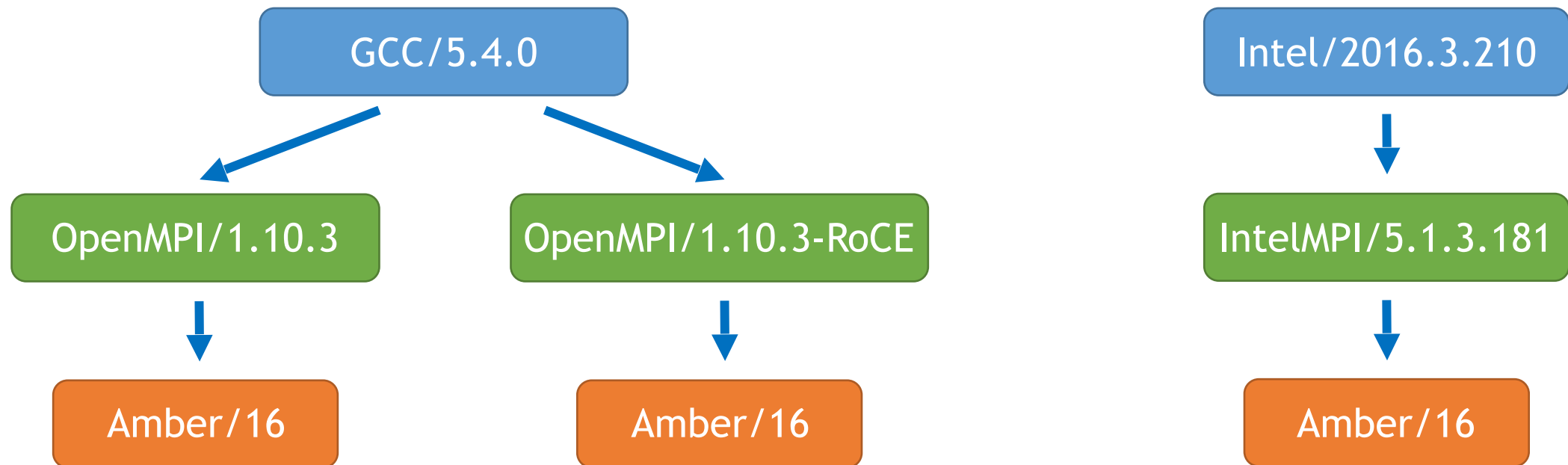
Seamless shell initialization

Incompatible software conflict resolution

# LMOD

Lua-based framework that provides a convenient way to customize user's environment through software modules.

Supports all shells available on the cluster.

Software is organized in a tree structure and displayed accordingly to the loaded dependencies.

`module avail <mod>`

- If no module is passed, print a list of all modules that are available to be loaded.

- If a module is specified, show all available modules with that name.

`module load mod1 mod2 …`

- Load the specified modules.

`module unload mod1 mod2 …`

- Unload the specified modules.

`module list`

- Show all modules loaded in the current environment.

`module purge`

- Remove all loaded modules from the environment.

- Show Program Code in editor

- Use Lmod to load GCC and compile

- Compile and run single-cpu program

# Submitting Batch Jobs

# Creating a Batch Job Script

A **batch job** consists of a sequence of commands listed in a file with the purpose of being interpreted as shell commands.

## *SHEBANG*

- Specify the script interpreter (Bash)
- Must be the first line!

## *SLURM DIRECTIVES*

- Start with "#SBATCH":
  Parsed by Slurm but ignored by Bash.
- Can be separated by spaces.
- Comments between and after directives are allowed.
- Must be before actual commands!

## *SCRIPT COMMANDS*

- Commands you want to execute on the compute nodes.

**myjob.slurm**

```
#!/bin/bash

#SBATCH --nodes=1
#SBATCH --ntasks=1
#SBATCH --mem=1G
#SBATCH --time=1-06:30:00
#SBATCH --job-name=myjob
#SBATCH --output=myjob.out


# Just a comment
module load GCC Python
python myscript.py
```

# CREATE A BATCH JOB SCRIPT

`--nodes=N`

- Request *N* nodes to be allocated. (*Default*: *N*=1)

`--ntasks=N`

- Request *N* tasks to be allocated. (*Default*: *N*=1)
- Unless otherwise specified, one task maps to one CPU core.

`--mem=NG`

- Request *N* gigabytes of memory <u>per node</u>. (*Default*: *N*=1)

`--time=d-hh:mm:ss`

- Request *d* days, *hh* hours, *mm* minutes and *ss* seconds. (*Default*: 00:15:00)

`--output=<file_name>`

- Write the batch script's standard output in the specified file.
- If not specified the output will be saved in the file: `slurm-<jobid>.out`

- Show batch script for single-cpu job

- Schedule job

- Monitor job and look at results

# Choosing Your Architecture

```
--constraint=<arch>
```

- Currently westmere, sandybridge, haswell, and skylake are available

# COMPILING ON A COMPUTE NODE

**salloc** *options*

- Obtain job allocation with shell access.
- Accepts all the same *options* previously seen for `sbatch`.

Gateway

```
[vanzod@vmps10 ~]$ salloc --nodes=1 --ntasks=4 --mem=16G --time=1:00:00
```

👍

Recommended for compiling, debugging and benchmarking sessions.

- **Use Salloc to run on a skylake node**

- Compile program for skylake architecture

- Try enabling avx512 instructions

- Compare results to un-optimized

# Using MPI

- Show mpi version of program in editor

- Load OpenMPI with LMod

- Run MPI program with 2, 4, 8 tasks interactively

- Compare results to single-cpu job

# CREATE A BATCH JOB SCRIPT FOR MPI

**`--nodes=`**$N$

- Request $N$ nodes to be allocated. (*Default*: $N$=1)

**`--ntasks=`**$N$

- Request $N$ tasks to be allocated. (*Default*: $N$=1)
- Unless otherwise specified, one task maps to one CPU core.

**`--tasks-per-node=`**$N$

- Request specifically $N$ tasks to run on each <u>node</u>.

**`--cpus-per-task=`**$N$

- Request $N$ cpu cores to be allocated to each task. (*Default*: $N$=1)

- Show mpi version slurm batch script

- Launch a scheduled MPI job with 48+ tasks

- Launch a scheduled MPI job with 8 tasks on each node

- Compare results