Intermediate Level Introduction to Computing at CARC

Matthew Fricke

Version 0.2

Goals

- 1) SLURM scheduler literacy
- 2) Ability to employ coupled parallel solutions

• We won't cover file transfer, storage systems, module system, conda, PBS. (These are all covered in depth in the video tutorials)

Agenda

- HPC and Parallelism
- HPC Schedulers
- SLURM
- Coupled Parallelism
 - Message Passing
 Interface
 - StarCCM with MPI









We will have one 15 minute break. Opportunity to see the Sample Footer Text machine room.

Logging into Hopper



First login to the Linux **workstation** in front of you.

Use your CARC username and password.

Jacob, Keven, Tannor, and Jose can help you login if you have trouble.

This is an "important step" so don't let me move on until you have logged in

Logging into Hopper





Should prompt you for a password...

Don't let me move on until you are able to login.

Welcome to Hopper

Be sure to review the "Acceptable Use" guidelines posted on the CARC website.

For assistance using this system email help@carc.unm.edu.

Tutorial videos can be accessed through the CARC website: Go to http://carc.unm.edu, select the "New Users" menu and then click "Introduction to Computing at CARC".

Warning: By default home directories are world readable. Use the chmod command to restrict access.

Don't forget to acknowledge CARC in publications, dissertations, theses and presentations that use CARC computational resources:

"We would like to thank the UNM Center for Advanced Research Computing, supported in part by the National Science Foundation, for providing the research computing resources used in this work."

Please send citations to publications@carc.unm.edu.

There are three types of slurm partitions on Hopper: 1) General - this partition is accessible by all CARC users.

2) Condo - preemtable scavenger queue available to all condo users. Your job must use checkpointing to use this queue or you will lose any work you have done if it is preempted by the partition's owner.

Logging into Hopper



[vanilla@hopper~]\$ qgrok queues free busy offline jobs nodes CPUs GPUs

_____ ____ ____

general 1907103200debug 20002640totals: 1907103200

mvanilla@hopper:~ \$ qgrok
queues free busy offline jobs nodes CPUs GPUs

_____ ____

```
general 1
        9
          0
             7 10 320 0
debug 2 0 0
             0
               2
                   64 0
condo 18 19 1 7 38 1216 8
bugs 0 2 0
           0 2 64 0
pcnc 1 1 0
           0 2 64
                    0
pathogen 1 0 0
              0 1 32 0
          3 10 320 0
   5 5 0
tc
gold
   2 0 0
           0 2
                  64 0
fishgen 0 1 0 0 1 32 0
neuro-hsc 8 6 0
            0 14 448 0
cup-ecs 0 2 0 2 2 64 4
  0 1 0
          0 1 32 2
tid
biocomp 0 1 0
                1 32 1
              0
         0
chakra 1 0
            0
                 32
               1
                     1
             0 1 32 0
    0
      0 1
pna
      28 1
totals: 19
             14
                 48 1536 8
```

____ ____

vanilla@hopper:~ \$ qgrok
queues free busy offline jobs nodes CPUs GPUs

Open partitions for use by everyone with a CARC account.

Purchased by the Office for the Vice President for Research.

vanilla@hopper:~ \$ qgrok
queues free busy offline jobs nodes CPUs GPUs

_____ ____

```
general 1
       9
           0
              7
                 10
                    320 0
debug 2 0 0
              0
                 2
                   64 0
condo 18 19 1 7 38 1216 8
     0 2 0
            0 2 64 0
bugs
   1 1 0 0 2 <u>64</u> 0
pcnc
pathogen 1 0 0
               0 1 32 0
          3 10 320 0
    5 5 0
tc
gold
   2 0 0 0 2
                   64
                     0
fishgen 0 1 0 0 1 32 0
neuro-hsc 8 6 0
             0 14 448 0
cup-ecs 0 2 0 2 2 64 4
    0 1 0
tid
           0 1 32 2
biocomp 0 1 0
               0
                1 32 1
chakra 1
         0
                1 32 1
       0
              0
       0 1
             0 1 32 0
    0
pna
totals: 19
       28
              14
                 48 1536 8
          1
```

___ ____

vanilla@hopper:~ \$ qgrok queues free busy offline jobs nodes CPUs GPUs



Private partitions

Reserved for use by the purchaser.

 Request access by emailing <u>support@carc.unm.edu</u> and CC the partition owner. mfricke@hopper:~\$qgrok queues free busy offline jobs nodes CPUs GPUs

Condo "scavenger" partition

Allows you to use compute nodes purchased by another group that are currently idle.

May be interrupted at any time if the owners start to use it.

[vanilla@hopper ~]\$ sinfo --partition debug PARTITION AVAIL TIMELIMIT NODES STATE NODELIST debug up 4:00:00 2 idle hopper[011-012]

sinfo reports information about partitions

[vanilla@hopper ~]\$ sinfo --partition debug PARTITION AVAIL TIMELIMIT NODES STATE NODELIST debug up 4:00:00 2 idle hopper[011-012]

The debug queues are intended for testing your programs.

And for interactive jobs.

[vanilla@hopper~]\$ sinfo --partition debug PARTITION AVAIL TIMELIMIT NODES STATE NODELIST debug up 4:00:00 2 idle hopper[011-012]



[vanilla@hopper ~]\$ sinfo --partition debug PARTITION AVAIL TIMELIMIT NODES STATE NODELIST debug up 4:00:00 2 idle hopper[011-012]

You can run a "job" for up to 4 hrs.

[vanilla@hopper ~]\$ sinfo --partition debug PARTITION AVAIL TIMELIMIT NODES STATE NODELIST debug up 4:00:00 2 idle hopper[011-012]

There are two nodes in this partition.

[vanilla@hopper~]\$ sinfo --partition debug PARTITION AVAIL TIMELIMIT NODES STATE NODELIST debug up 4:00:00 2 idle hopper[011-012]

The names of the nodes in the partition

[vanilla@hopper~]\$ sinfo --partition debug PARTITION AVAIL TIMELIMIT NODES STATE NODELIST debug up 4:00:00 2 idle hopper[011-012]

The names of the nodes in the partition

[vanilla@hopper ~]\$ sinfo --partition general PARTITION AVAIL TIMELIMIT NODES STATE NODELIST general* up 2-00:00:00 9 alloc hopper[001-009] general* up 2-00:00:00 1 idle hopper010

Hopper001 through 009 are running jobs.

Hopper010 is waiting to be used.

[vanilla@hopper ~]\$ hostname hopper [vanilla@hopper ~]\$

Running on the Head Node. The head node's name is "hopper". [vanilla@hopper ~]\$ hostname hopper [vanilla@hopper ~]\$ man hostname [vanilla@hopper ~]\$ hostname hopper [vanilla@hopper ~]\$ man hostname ('q' to quit)

[vanilla@hopper ~]\$ man man ('q' to quit)

[vanilla@hopper ~]\$ man sinfo

sinfo(1)

Slurm Commands

sinfo(1)

NAME

sinfo - View information about Slurm nodes and partitions.

SYNOPSIS sinfo [<u>OPTIONS</u>...]

DESCRIPTION

sinfo is used to view partition and node information for a system running Slurm

OPTIONS

-a, --all

Display information about all partitions. This causes information to be displayed about partitions that are configured as hidden and partitions that are unavailable to the user's group.

[vanilla@hopper ~]\$ sinfo --all

PARTITION AVAIL TIMELIMIT NODES STATE NODELIST

general* up 2-00:00:00 9 alloc hopper[001-009]
general* up 2-00:00:00 1 idle hopper010
debug up 4:00:00 2 idle hopper[011-012]
condo up 2-00:00:00 1 down* hopper045
condo up 2-00:00:00 3 mix hopper[018-020]
condo up 2-00:00:00 16 alloc hopper[013-015,028-036,049-052]
condo up 2-00:00:00 18 idle hopper[016-017,021-027,037-044,053]
bugs up 7-00:00:00 2 alloc hopper[013-014]
pcnc up 7-00:00:00 1 alloc hopper015
pcnc up 7-00:00:00 1 idle hopper016
pathogen up 7-00:00:00 1 idle hopper017
tc up 7-00:00:00 3 mix hopper[018-020]
tc up 7-00:00:00 2 alloc hopper[029-030]
tc up 7-00:00:00 5 idle hopper[021-025]
gold up 7-00:00:00 2 idle hopper[026-027]
fishgen up 7-00:00:00 1 alloc hopper028
neuro-hsc up 7-00:00:00 6 alloc hopper[031-036]
neuro-hsc up 7-00:00:00 8 idle hopper[037-044]
cup-ecs up 7-00:00:00 2 alloc hopper[049-050]
tid up 7-00:00:00 1 alloc hopper051
biocomp up 7-00:00:00 1 alloc hopper052
chakra up 7-00:00:00 1 idle hopper053
pna up 7-00:00:00 1 down* hopper045



Tell slurm to run a program on a compute node...

Run the program on a compute node in the debug partition.

The program to run.

srun: Account not specified in script or ~/.default_slurm_account, using latest project

You have not been allocated GPUs. To request GPUs, use the -G option in your submission script.

hopper011

[vanilla@hopper ~]\$ squeue



PD means programs that are waiting their turn.

2 (QUSMaxCpuPeruserLimit)

Shows you what the slurm scheduler is doing right now.

Here we can see that user 'erowland' has a lot of programs waiting to run.

- 2 (QOSMaxCpuPerUserLimit) 2 (QOSMaxCpuPerUserLimit)
- 2 (QOSMaxCpuPerUserLimit)

```
0.00
```

[vanilla@hopper ~]\$ squeue

4001

gene

JOBID	PARTITION	NAME	USER	ST	TIME	E NODES	• NODELIST	(REAS	SON)
	4314	general		PRE	erowland	PD	0:00	2	(QOSMaxCpuPerUserLimit)
	4315	general		PRE	erowland	PD	0:00	2	(QOSMaxCpuPerUserLimit)
	4317	general		PRE	erowland	PD	0:00	2	(QOSMaxCpuPerUserLimit)
	4318	general		PRE	erowland	PD	0:00	2	(QOSMaxCpuPerUserLimit)
								2	(QOSMaxCpuPerUserLimit)

The reason these jobs are not running is that 'erowland' is already using the maximum number of CPUs they are allowed

	2	(QOSMaxCpuPerUserLimit)
	2	(QOSMaxCpuPerUserLimit)
•	2	(QOSMaxCpuPerUserLimit)
	2	(QOSMaxCpuPerUserLimit)
00	2	(QOSMaxCpuPerUserLimit)

[vanilla@hopper ~]\$ squeue -t R --all

JOBID	PARTITIC	N NAN	IE US	SER ST	TIME NODES NODELIST(REASO	N)
4405	condo	2ndMA	mfrick	e R 1-07	7:48:30 6 hopper[031-036]	
5208	condo	NN	kgu R	5:48:49	9 1 hopper015	
5210	condo	NN	kgu R	6:30:13	3 1 hopper014	
5209	condo	NN	kgu R	6:31:13	3 1 hopper013	
5206	condo	NN	kgu R	6:32:13	3 1 hopper051	
5207	condo	NN	kgu R	6:32:13	3 1 hopper052	
5205	condo	NN	kgu R	6:32:43	3 1 hopper028	
4595	cup-ecs	golConfi	aalasan	d R 2-06	6:51:59 1 hopper050	
4594	cup-ecs	golConfi	aalasan	d R 2-06	6:52:03 1 hopper049	
5120	general	jupyterh	jacobr	n R 11:4	:45:47 1 hopper007	
4313	general	PRE er	owlanc	R 1:17	7:29 2 hopper[003-004]	
5111	general	1stMA	mfricke	e R 11:1	15:28 2 hopper[005-006]	
5025	general	c2n j	xzuo R	1:50	1 hopper001	
5024	general	c2n j	xzuo R	31:28	3 1 hopper002	
5203	general	NN	kgu R	6:37:50	0 1 hopper009	
5201	general	NN	kgu R	6:38:14	4 1 hopper008	
4390	tc UC	sTpCyd le	epluart	R 2-15:1	l8:18 3 hopper[018-020]	
5198	tc	NN kg	u R 6	:40:19	1 hopper030	
5196	tc	NN kg	ju R 6	:40:31	1 hopper029	

srun: Account not specified in script or ~/.default_slurm_account, using latest project

You have not been allocated GPUs. To request GPUs, use the -G option in your submission script. hopper011

hopper011

srun: Account not specified in script or ~/.default_slurm_account, using latest project

You have not been allocated GPUs. To request GPUs, use the -G option in your submission script. hopper011

hopper011

You ran two copies of your program.

ntasks is the number of copies to run.
[vanilla@hopper ~]\$ srun --partition debug --ntasks 8 hostname srun: Account not specified in script or ~/.default_slurm_account, using latest project hopper011 hopper011 You have not been allocated GPUs. To request GPUs, use the -G option in your submission script. hopper011 hopper011

hopper011

hopper011

hopper011

You ran eight copies of your program.

ntasks is the number of copies to run.

[vanilla@hopper ~]\$ srun --partition debug --ntasks 8 hostname srun: Account not specified in script or ~/.default_slurm_account, using latest project hopper011 hopper011 hopper011 You have not been allocated GPUs. To request GPUs, use the -G option in your submission script.

hopper011

hopper011

hopper011

hopper011

hopper011

By default, each task (copy of your program) is allowed to use one CPU.

Many programs are able to use more than one CPU at a time.

[vanilla@hopper ~]\$ srun --partition debug --ntasks 2 --cpus-per-task 2 hostname srun: Account not specified in script or ~/.default_slurm_account, using latest project You have not been allocated GPUs. To request GPUs, use the -G option in your submission script. hopper011 hopper011

> Here we are telling SLURM to run 2 copies of our program and let each copy of our program use 2 CPUs.

[vanilla@hopper ~]\$ srun --partition debug --nodes 2 --ntasks-per-node 4 hostname srun: Account not specified in script or ~/.default_slurm_account, using latest project hopper012

You have not been allocated GPUs. To request GPUs, use the -G option in your submission script.

hopper012

hopper011

hopper011

hopper012

hopper012

hopper011

hopper011

Here we are telling SLURM to run 4 copies of our program on 2 different compute nodes.

This is useful when our programs need a bigger share of the compute node.

- [vanilla@hopper ~]\$ srun --partition debug --nodes 2
- --ntasks-per-node 2 --cpus-per-task 2 hostname
- srun: Account not specified in script or ~/.default_slurm_account, using latest project
- hopper011
- You have not been allocated GPUs. To request GPUs, use the -G option in your submission script.
- hopper011
- hopper012
- hopper012

And we can combine all three.

[vanilla@hopper ~]\$ srun --partition debug --mem 4G --nodes 2 --ntasks-per-node 2 --cpus-per-task 2 hostname srun: Account not specified in script or ~/.default_slurm_account, using latest project

hopper012 hopper012 You have not been allo submission script. hopper011 Hopper011

And we can specify how much memory we want.

--mem 4G means give me 4 gigabytes of memory per node. [vanilla@hopper ~]\$ srun --partition debug --mem 4G --nodes 2 --ntasks-per-node 2 --cpus-per-task 2 hostname srun: Account not specified in script or ~/.default_slurm_account, using latest project

hopper012 hopper012 You have not been allo submission script. hopper011 Hopper011

Why does all this matter?

The purpose of SLURM is to provide you the hardware your programs need.

So you have to understand what those requirements are really well.

[vanilla@hopper ~]\$ srun --partition debug --mem 4G --nodes 2 --ntasks-per-node 2 --cpus-per-task 2 hostname srun: Account not specified in script or ~/.default slurm account, using latest project

hopper012 hopper012

You have not been all submission script. hopper011 Hopper011 3)

Can my program use multiple 1) **CPUs**? 2) How much memory does my program need? Can my program use multiple compute nodes (MPI*, GNU Parallel*)? Can my program use GPUs? 4)

[vanilla@hopper ~]\$ srun --partition debug --mem 4G --nodes 2 --ntasks-per-node 2 --cpus-per-task 2 hostname srun: Account not specified in script or ~/.default slurm account, using latest project hopper012 hopper012 This command is getting pretty long. You have not been alloca submission script. We can use salloc to avoid asking for hopper011 Hopper011 the same resources every time we

use srun.

[vanilla@hopper ~]\$ salloc --partition debug --nodes 2 --ntasks-per-node 2 salloc: Account not specified in script or ~/.default_slurm_account, using latest project

salloc: Granted job allocation 5251

salloc: Waiting for resource configuration

salloc: Nodes hopper[011-012] are ready for job

[vanilla@hopper ~]\$

This command is getting pretty long. We can use salloc to avoid asking for the same resources every time we use srun. [vanilla@hopper ~]\$ srun hostname
hopper012
hopper012
hopper011
You have not been allocated GPUs. To request GPUs, use the -G option in your
submission script.

hopper011

[vanilla@hopper ~]\$ srun hostname

hopper012

hopper011

hopper012

hopper011

[vanilla@hopper~]\$

Now we can use srun over and over without having to ask for a new hardware allocation each time.

[vanilla@hopper ~]\$ exit exit salloc: Relinquishing job allocation 5251

Always type exit when you are done with the hardware. Running salloc inside an allocation gets very confusing.

Interactive vs Batch Mode

Interactive Mode

• Everything so far has been interactive. You request hardware, run your program, and get the output on your screen right away.

Batch Mode

- Most programs at an HPC center are run in "batch" mode.
- Batch mode means we write a shell script that the SLURM scheduler runs for us. The script requests hardware just like we did with salloc and then runs the commands in the script.
- Whatever would have been written to the screen is saved to a file instead.

[vanilla@hopper ~]\$ git clone https://lobogit.unm.edu/CARC/workshops.git Cloning into 'workshops'... remote: Enumerating objects: 132, done. remote: Counting objects: 100% (75/75), done. remote: Compressing objects: 100% (43/43), done. remote: Total 132 (delta 33), reused 74 (delta 32), pack-reused 57 Receiving objects: 100% (132/132), 57.58 KiB | 3.60 MiB/s, done. Resolving deltas: 100% (51/51), done.

> Rather than make you write shell scripts lets just download some we wrote for this workshop...

[vanilla@hopper ~]\$ tree workshops

workshops/

— intro_workshop

- code
- ---- calcPiMPI.py
- ---- calcPiSerial.py

—– vecadd

– Makefile

---- vecadd_gpu.cu

-— vecadd_mpi_cpu

---- vecadd_mpi_cpu.c

---- vecaddmpi_cpu.sh

L-vecadd_mpi_gpu.c

— data

— H2O.gjf

---- step_sizes.txt

- slurm

- calc_pi_array.sh
- calc_pi_mpi.sh
- calc_pi_parallel.sh
- calc_pi_serial.sh

— gaussian.sh

- hostname_mpi.sh
- ---- vecadd_hopper.sh

---- vecadd_xena.sh

- workshop_example2.sh
- workshop_example3.sh

— workshop_example.sh

README.md

Run tree to see how the workshops directories are organized...

[vanilla@hopper ~]\$ tree workshops

intro_workshop code code calcPiMPI.py calcPiSerial.py vecadd Makefile vecadd_gpu.cu vecadd_mpi_cpu vecadd_mpi_cpu.c vecadd_mpi_cpu.sh vecadd_mpi_gpu.c data H2O.gjf step_sizes.txt

— slurm

workshops/

— calc_pi_array.sh

- calc_pi_mpi.sh

---- calc_pi_parallel.sh

— calc_pi_serial.sh

— gaussian.sh

- hostname_mpi.sh
- --- vecadd_hopper.sh

--- vecadd_xena.sh

---- workshop_example2.sh

- workshop_example3.sh

└── workshop_example.sh

- README.md

Run tree to see how the workshops directories are organized...

The workshop files are divided into "code", "slurm", and "data" directories.

[vanilla@hopper \$] cd ~/workshops/intro workshop [vanilla@hopper intro workshop]\$ pwd /users/vanilla/workshops/intro workshop [vanilla@hopper intro workshop]\$ cat slurm/workshop example.sh #!/bin/bash **#SBATCH** --partition debug **#SBATCH** --ntasks 4 #SBATCH --time 00:05:00 #SBATCH --job-name ws_example **#SBATCH** --mail-user your username@unm.edu **#SBATCH** --mail-type ALL

srun hostname

Let's take a look at the workshop_example.sh script in the slurm directory... [vanilla@hopper intro_workshop]\$ sbatch slurm/workshop_example.sh sbatch: Account not specified in script or ~/.default_slurm_account, using latest project Submitted batch job 5252 [vanilla@hopper intro_workshop]\$

> We submit our slurm shell script with the sbatch command.

[vanilla@hopper intro_workshop]\$ sbatch slurm/workshop_example1.sh sbatch: Account not specified in script or ~/.default_slurm_account, using latest project Submitted batch job 5252 [vanilla@hopper intro_workshop]\$

Notice that the only output we get is a job id.

This indicates that the script was successfully sent to the scheduler.

The commands in the script will run as soon as the hardware requested is available. We submit our slurm shell script with the sbatch command.



Compute Node 01

Compute Node 02

Compute Node 03

Compute Node 04

Compute Node 05









[vanilla@hopper intro_workshop]\$ ls code data pbs slurm slurm-5252.out

The hostname command is very fast so everyone's job should finish in a few seconds.

When it is finished you will have a new file named slurm-{your job id}.out. [vanilla@hopper intro_workshop]\$ ls code data pbs slurm slurm-5252.out

> When it is finished you will have a new file named slurm-{your job id}.out.

[vanilla@hopper intro_workshop]\$ cat slurm-5252.out hopper011 [vanilla@hopper intro_workshop]\$ ls code data pbs slurm slurm-5252.out

> When it is finished you will have a new file named slurm-{your job id}.out.

[vanilla@hopper intro_workshop]\$ cat slurm-5252.out What do you see?

Serial Program to Calculate π



[vanilla@hopper intro_workshop]\$ module load miniconda3
[vanilla@hopper intro_workshop]\$ conda create -n numpy numpy

Wait a while – introduce yourselves to your neighbor...

Conda allows you to install software into your home directory. In this case we need the numerical python libraries for calcPiSerial.py

Let's experiment with a program that does slightly more than print the hostname. [vanilla@hopper intro_workshop]\$ source activate numpy [vanilla@hopper intro_workshop]\$ srun --partition debug python code/calcPiSerial.py 10 srun: Using account 2016199 from ~/.default_slurm_account You have not been allocated GPUs. To request GPUs, use the -G option in your submission script.

Pi = 3.14242598500109870940, (Diff=0.00083333141130559341) (calculated in 0.000005 secs with 10 steps)

Activate the numpy environment and Run calcPiSerial.py on a compute node.

For our example program the more steps it takes the more accurate it is, but the longer it takes.

[vanilla@hopper intro_workshop]\$ sbatch slurm/calc_pi_serial.sh sbatch: Using account 2016199 from ~/.default_slurm_account Submitted batch job 5263 vanilla@hopper:~/workshops/intro_workshop\$ squeue --me JOBID PARTITION NAME USER ST TIME NODES NODELIST(REASON) 5263 debug calc_pi_ vanilla R 0:44 1 hopper011

> Edit slurm/calc_pi_serial.sh. Change the email address to your address and submit the script.

Then enter squeue -- me to see the job status.

Take a look at the job output.

Parallelism – Coupled Parallelism

- Coupled problems are those where the CPUs need to work together to solve a problem by communicating with each other.
- Many commercial and research programs designed to run on HPC systems like CARC use a library called the message passing interface (MPI) to do this.
- We have written an MPI version of our python pi calculator to demonstrate.



Serial Program to Calculate π



Parallel Program to Calculate π



MPI: Message Passing Interface

When programs need to run on many processors but also communicate with one another.

Here the parallel version of calcPi needs to communicate the partial sums computed by each process so they can all be added up.

To communicate we will use the MPI library:

module load minconda3 conda create –n mpi_numpy mpi mpi4py numpy

import time import sys import numpy as np # Value of PI to compare to

```
#Distributed function to calculate pi
def Pi(num_steps):
    step = 1.0 / num_steps
    sum = 0
    for i in range(rank, num_steps, num_procs): # Divide sum among
processes
    x = (i + 0.5) * step
    sum = sum + 4.0 / (1.0 + x * x)
```

```
sum = sum + 4.0
mypi = step * sum
```

Get that partial sums from all the processes, add them up, and give to the root process

```
pi = comm.reduce(mypi, MPI.SUM, root)
return pi
```

#Main function
Check that the caller gave us the number of steps to use
if len(sys.argv) != 2:
 print("Usage: ", sys.argv[0], " <number of steps>")
 sys.exit(1)

num_steps = int(sys.argv[1],10);

#Broadcast number of steps to use to the other processes
comm.bcast(num_steps, root)

Call function to calculate pi
start = time.time() #Start timing
pi = Pi(num_steps) # Call the function that calculates pi
end = time.time() # End timing

If we are the root process then print our estimation of pi, # the difference from numpy's value, and how long it took print("Pi = %.20f, (Diff=%.20f) (calculated in %f secs with %d steps)" %(pi, pi-np.pi, end - start, num_steps))
```
#!/bin/bash
#SBATCH --partition debug
#SBATCH --nodes 2
#SBATCH --ntasks-per-node 4
#SBATCH --time 00:05:00
#SBATCH --job-name calc_pi_mpi
#SBATCH --mail-user your_username@unm.edu
#SBATCH --mail-type ALL
```

module load miniconda3
source activate mpi_numpy

cd \$SLURM_SUBMIT_DIR srun --mpi=pmi2 python code/calcPiMPI.py 100000000

sbatch slurm/calc_pi_mpi.sh

srun --mpi=pmi2 python code/calcPiMPI.py 100000000

srun understands MPI programs!

If you ever used mpirun or mpiexec you had to provide a lot of parameters to describe how many compute nodes you had and what their names are, etc.

But srun is part of SLURM so it already knows all that.

The only thing you have to specify is the communication library to use. In our case "pmi2".

Run calc_pi_mpi.sh.

Vary the number of tasks it uses.

Use squeue to monitor the state of your job,

Look at your output files.

What is the relationship between the number of tasks and how cast it calculates pi?







STAR-CCM+°

[vanilla@hopper intro_workshop]\$ cd ~/workshops/starccm/ [vanilla@hopper starccm]\$ nano starccm_example.slurm

#!/bin/bash -l
#SBATCH --job-name starccm-example
#SBATCH --partition general
#SBATCH --time=6:00:00
#SBATCH --nodes 2
#SBATCH --nodes 2
#SBATCH --ntasks-per-node 32
#SBATCH --mem=90GB
#SBATCH --mail-user your@email.address
#SBATCH --mail-type all

SIM_PATH=/carc/scratch/users/\$USER/starccm_example/L22_AE_FullCar_Cornering_v2.sim

module load starccm starccm+ -batchsystem **slurm** -batch **\$SIM_PATH** -time -batch-report -licpath 1717@10.101.164.203

Copy example sim file to your personal scratch directory

[vanilla@hopper intro_workshop]\$ rsync --info=progress2 /projects/shared/workshops/starccm/starccm_example/L22_AE_FullCar_Cornering_v2.sim /carc/scratch/users/\$USER/starccm_example/

Queue up the job

[vanilla@hopper starccm]\$ sbatch starccm_example.slurm

Check the job status

[vanilla@hopper starccm]\$ squeue --me JOBID PARTITION NAME USER ST TIME NODES NODELIST(REASON) 2746424 general starccm- vanilla R 0:05 2 hopper[001-002]

Queue up the job

[vanilla@hopper starccm]\$ tail -f slurm-YOURJOBID.out

Starting local server: /opt/local/starccm/18.06.006/STAR-CCM+18.06.006/star/bin/starccm+ -batchsystem slurm -licpath 1717@10.101.164.203 -server -xencoded-session L2NhcmMvc2NyYXRjaC91c2Vycy9tZnJpY2tlL3N0YXJjY21fZXhhbXBsZS9MMjJfQUVfRnVsbENhcl9Db3JuZXJpbmdfdjluc2lt

Starting parallel server MPI Distribution : Open MPI-4.0.3 Host 0 -- hopper001 -- Ranks 0-31 Host 1 -- hopper002 -- Ranks 32-63 Process rank 0 hopper001 609981 Total number of processes : **64**

 Iteration
 Continuity
 X-momentum
 Y-momentum
 Z-momentum
 Tke
 Sdr

 1
 1.000000e+00
 1.00

After the job finishes check to see how efficiently the resources were used

[vanilla@hopper starccm]\$ Seff 2746310 Job ID: 2746310 **Cluster: hopper** User/Group: mfricke/users State: COMPLETED (exit code 0) Nodes: 2 Cores per node: 32 CPU Utilized: 5-18:44:49 CPU Efficiency: 99.02% of 5-20:07:28 core-walltime Job Wall-clock time: 02:11:22 Memory Utilized: 29.03 GB Memory Efficiency: 16.13% of 180.00 GB

Timing on Hopper

12 GB Racecar Model Input Hopper 1 node 32 tasks – about 4 hours. Hopper 2 nodes 64 tasks – 1 hour 50 minutes. Hopper 4 nodes 128 tasks – Under 1 hour.

Useful Slurm Commands

squeue --me --long squeue --me --start scancel jobid scancel --u \$USER sacct seff jobid shows information about jobs you submitted shows when slurm expects your job to start cancels a job cancels all your jobs shows your job history shows how efficiently the hardware was used