

ARGONNE  
**ATPESC2025**  
EXTREME - SCALE COMPUTING

# Quick Start on ATPESC Computing Resources

**JaeHyuk Kwack**

Argonne National Laboratory



Argonne Leadership  
Computing Facility

[extremecomputingtraining.anl.gov](http://extremecomputingtraining.anl.gov)



# Outline



**The DOE Leadership Computing Facility**



**ALCF Aurora/Polaris System**



**OLCF Odo System**



**NERSC Perlmutter System**



**Hands-on**

# The DOE Leadership Computing Facility

- Collaborative, multi-lab, DOE/SC initiative ranked top national priority in *Facilities for the Future of Science: A Twenty-Year Outlook*.
- Mission: Provide the computational and data science resources required to solve the most important scientific & engineering problems in the world.
- Highly competitive user allocation program (INCITE, ALCC).
- Projects receive 100x more hours than at other generally available centers.
- LCF centers partner with users to enable science & engineering breakthroughs (Liaisons, Catalysts).



# Leadership Computing Facility System

	Argonne LCF		Oak Ridge LCF
System	HPE	HPE	HPE
Name	Polaris	Aurora	Frontier
Compute nodes	560	10,624	9,408
Node architecture	1 x AMD Milan CPU + 4x NVIDIA A100 GPU	2 x Intel Xeon SPR + 6 x Intel PVC GPU	1 x AMD EYPC CPU + 4 x AMD MI250x GPU
Processing Units	560 CPUs + 2,240 GPUs	21,248 CPUs + 63,744 GPUs	9,408 CPUs + 37,362 GPUs
Memory per node, (gigabytes)	512 GB DDR4 + 160 GB HBM2 + 1600 GB SSD	128 GB HBM2e on CPU + 1024 GB DDR5 on CPU + 768 GB HBM2e on GPU	512 GB DDR4 + 512 GB HBM2e
Peak performance, (petaflops)	44	> 2 Exaflops DP	1.6 Exaflop DP

# AVAILABLE COMPUTING RESOURCES FOR ATPESC

## ALCF Systems

Intel CPUs + Intel PVC GPUs (Aurora)  
AMD CPUs + NVIDIA A100 GPUs (Polaris)

## OLCF

AMD CPUs + AMD MI-250x GPUs (Odo)

## NERSC

AMD CPUs + NVIDIA A100 GPUs (Perlmutter)

# ALCF Aurora/Polaris System

# Aurora Exascale Compute Blade

## NODE CHARACTERISTICS

**6** GPUs - Intel Data Center GPU Max Series

**2** CPUs - Intel Xeon CPU Max Series

**768 GB** GPU HBM Memory

**19.66 TB/s** Peak GPU HBM BW

**128 GB** CPU HBM Memory

**2.87 TB/s** Peak CPU HBM BW

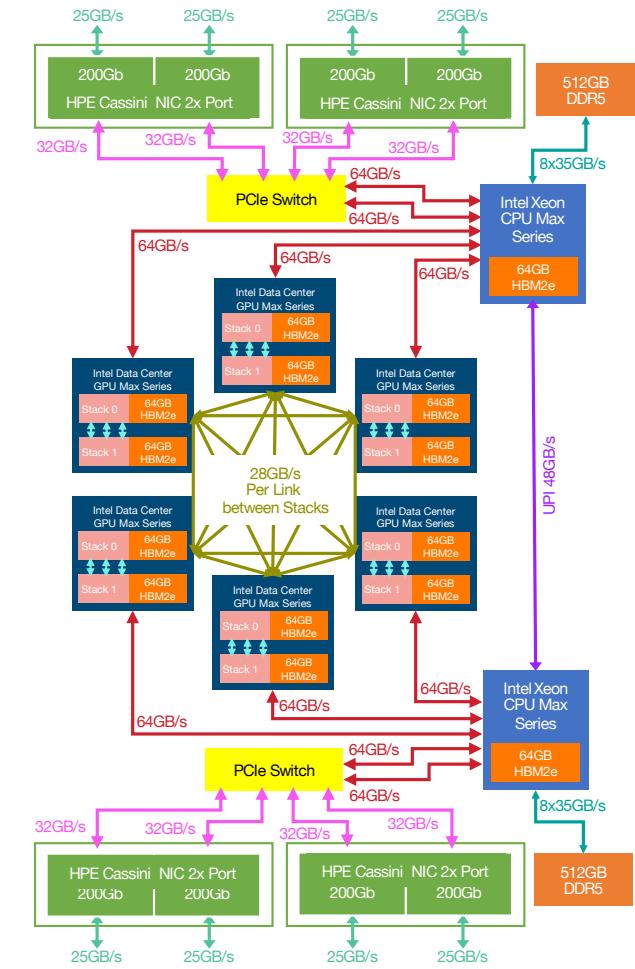
**1024 GB** CPU DDR5 Memory

**0.56 TB/s** Peak CPU DDR5 BW

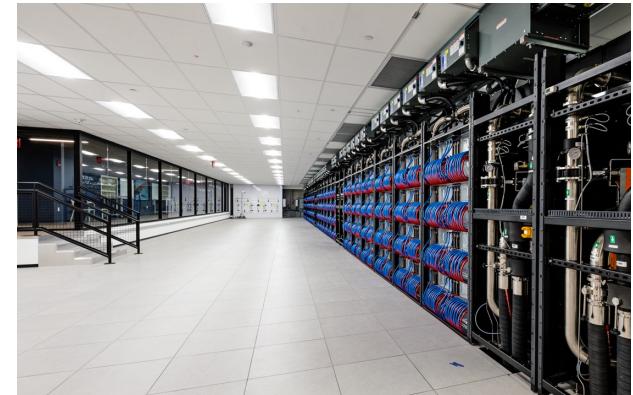
**≥ 130 TF** Peak Node DP FLOPS

**200 GB/s** Max Fabric Injection

**8** NICs



# Aurora System Configuration



**Intel® Data Center GPU  
Max 1550 (PVC)**

**4<sup>th</sup> Gen Intel XEON Max  
Series CPU with *High  
Bandwidth Memory***

Platform  
**HPE Cray-Ex**

**Racks - 166**  
**Nodes - 10,624**  
CPUs - 21,248  
GPUs – 63,744

**Interconnect**  
HPE Slingshot 11  
Dragonfly topology with adaptive routing  
Cassini NIC, 200 Gb/s (25 GB/s), 8 per node  
Network Switch:  
25.6 Tb/s per switch (64 200 Gb/s ports)  
Links with 25 GB/s per direction

**Peak FP64 Performance  
≥ 2 exaFLOPS**

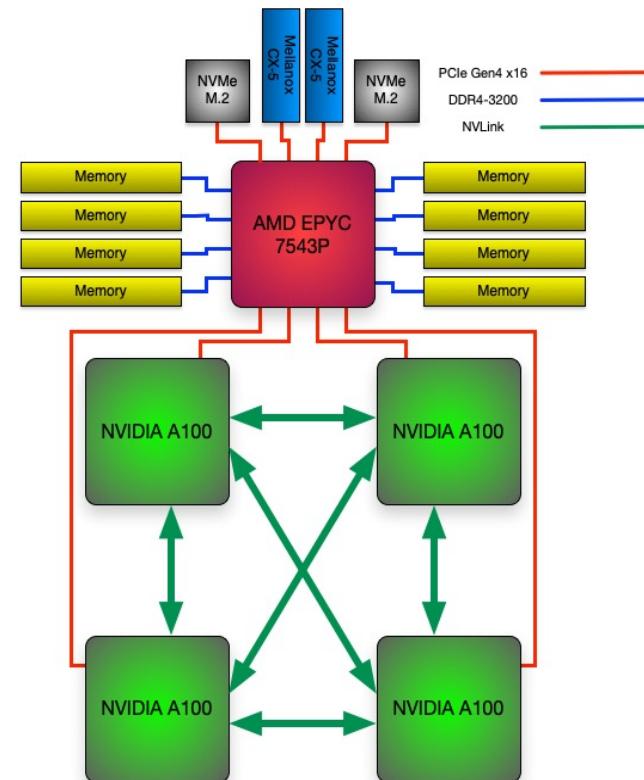
**Memory**  
**10.9PiB of DDR @ 5.95 PB/s**  
**1.36PiB of CPU HBM @ 30.5 PB/s**  
**8.16PiB of GPU HBM @ 208.9 PB/s**

**Network**  
**2.12 PB/s Peak Injection BW**  
**0.69 PB/s Peak Bisection BW**

**Storage**  
**230PB DAOS Capacity**  
**31 TB/s DAOS Bandwidth**

# Polaris Single Node Configuration

# of AMD EPYC 7543P CPUs	1
# of NVIDIA A100 GPUs	4
Total HBM2 Memory	160 GB
HBM2 Memory BW per GPU	1.6 TB/s
Total DDR4 Memory	512 GB
DDR4 Memory BW	204.8 GB/s
# OF NVMe SSDs	2
Total NVMe SSD Capacity	3.2 TB
# of Mellanox NICs	2
Total Injection BW (w/ Cassini)	25 (50) GB/s
PCIe Gen4 BW	64 GB/s
NVLink BW	600 GB/s
Total GPU DP Tensor Core Flops	78 TF



# Polaris System Configuration

# of River Compute racks	40
# of Apollo Gen10+ Chassis	280
# of Nodes	560
# of AMD EPYC 7543P CPUs	560
# of NVIDIA A100 GPUs	2240
Total GPU HBM2 Memory	87.5TB
Total CPU DDR4 Memory	280 TB
Total NVMe SSD Capacity	1.75 PB
Interconnect	HPE Slingshot
# of Cassini NICs	1120
# of Rosetta Switches	80
Total Injection BW (w/ Cassini)	28 TB/s 13 TB/s
Total GPU DP Tensor Core Flops	44 PF
Total Power	1.8 MW



**Apollo 6500 Gen10+**

[extremecomputingtraining.anl.gov](http://extremecomputingtraining.anl.gov)

Argonne  
NATIONAL LABORATORY

# Aurora/Polaris Filesystems

- Lustre
  - Home directories (/home)
    - Default quota 50GiB
    - Your home directory is backed up
  - Project directory locations
    - Aurora: **/flare/ATPESC2025**
      - **CREATE A SUBDIRECTORY /flare/ATPESC2025/usr/your\_username**
    - Polaris: **/eagle/ATPESC2025**
      - **CREATE A SUBDIRECTORY /eagle/ATPESC2025/usr/your\_username**
      - Access controlled by unix group of your project
      - Default quota 1TiB
      - Project directories are NOT backed up
  - With large I/O on Lustre, be sure to consider **stripe width**

# Aurora/Polaris Modules

- A tool for managing a user's environment
  - Sets your PATH to access desired front-end tools
  - *Your compiler version can be changed here*
- *module commands*
  - *help*
  - *list* ← what is currently loaded
  - *avail*
  - *load*
  - *unload*
  - *switch|swap*
  - *use* ← add a directory to MODULEPATH
  - *display|show*

# Aurora Compiling

- Intel oneAPI Environment
  - The oneAPI programming environment is currently the single environment for building and running software to maximally use the available hardware resources. The oneAPI environment is loaded by default for users and is principally defined by the following set of modules and related variants.
    - oneapi: Intel oneAPI HPC toolkit.
    - mpich: MPI libraries.
  - Compiler wrappers
    - mpicc -> icx
    - mpicxx -> icpx
    - mpifort -> ifx
    - Support SYCL and OpenMP target offload
- Libraries found in
  - /opt/aurora/24.347.0/oneapi/
  - /opt/cray

# Polaris Compiling

- Cray Programming Environment (PE)
  - HPE provides compiler wrappers by default which includes various libraries (including MPI libraries)
    - Integrates with modules environment
    - HPE provided modules will add headers/libraries/compiler+linker options to compiler
    - -craype-verbose to show actual compile/link command
  - PrgEnv-nvhpc (default)
    - cc -> nvc
    - CC -> nvc++
    - ftn -> nvfortran
    - Support CUDA and OpenMP target offload
    - nvcc still available but not used by wrappers
  - PrgEnv-gnu
    - cc -> gcc
    - CC -> g++
    - ftn -> gfortran
- Libraries found in
  - /opt/nvidia
  - /opt/cray

# Aurora/Polaris Running MPI Applications

- Jobs run directly on the compute nodes. The `mpiexec` command runs applications using the Parallel Application Launch Service (PALS)
- `mpiexec`
  - Execute MPI applications on compute nodes using `mpiexec`
    - `-n` Total number of MPI ranks
    - `-ppn` Total number of MPI ranks per node
    - `--cpu-bind` CPU binding for application
    - `--depth` Number of CPUs per rank
    - `--env` Set environment variables (e.g., `OMP_NUM_THREADS=nthreads`)
    - `--hostfile` Indicate file with hostname
  - Full list of options available from the man page
    - <https://docs.alcf.anl.gov/aurora/running-jobs-aurora/>
    - <https://docs.alcf.anl.gov/polaris/running-jobs/>

# Aurora Jobs

- Two parts for running jobs
  - Interacting with scheduler
  - Launching job using mpiexec
- Shell script
  - describes parameters for scheduler
  - Commands to run included mpiexec to launch
  - Runs on ‘head’ node of your job
    - Permissible to run computation in your shell script
  - Need to load any of your non-default modules which provide library paths
- qsub -q prod ./run.sh
  - Will return the jobid
  - Output and error logs are in submission directory

```
#!/bin/bash
#PBS -A $PROJECT
#PBS -l walltime=01:00:00
#PBS -l select=4
#PBS -l filesystems=home:flare

rpn=12 # assume 1 process per GPU
procs=$((PBS_NODES*rpn))

# job to "run" from your submission directory
cd $PBS_O_WORKDIR

module load <something>

set +x # report all commands to stderr
env
mpiexec -n $procs -ppn $rpn --cpu-bind core -
genvall ./bin <opts>
```

# Aurora Interactive job

- Useful for short tests or debugging
- Submit the job with `-I` (letter I for Interactive)
  - Debug queue:
    - `qsub -I -l select=1 -l walltime=1:00:00 -q debug -A ATPESC2025 -l filesystems=home:flare`
  - Using reservation:
    - `qsub -I -l select=1 -l walltime=1:00:00 -q ATPESC -A ATPESC2025 -l filesystems=home:flare`
    - Check the reservation queues with `pbs_rstat`
- Wait for job's shell prompt
  - Exit this shell to end your job
- From job's shell prompt, run just like in a script job,
  - `mpiexec -n 4 -ppn 4 ./a.out`
- After job expires, `mpiexec` will fail. *Check `qstat $PBS_JOBID`*

# Polaris Jobs

- Two parts for running jobs
  - Interacting with scheduler
  - Launching job using mpiexec
- Shell script
  - describes parameters for scheduler
  - Commands to run included mpiexec to launch
  - Runs on ‘head’ node of your job
    - Permissible to run computation in your shell script
  - Need to load any of your non-default modules which provide library paths
- qsub -q prod ./run.sh
  - Will return the jobid
  - Output and error logs are in submission directory

```
#!/bin/bash
#PBS -A $PROJECT
#PBS -l walltime=01:00:00
#PBS -l select=4
#PBS -l filesystems=home:eagle

rpn=4 # assume 1 process per GPU
procs=$((PBS_NODES*rpn))

# job to "run" from your submission directory
cd $PBS_O_WORKDIR

module load <something>

set +x # report all commands to stderr
env
mpiexec -n $procs -ppn $rpn --cpu-bind core -
genvall ./bin <opts>
```

# Polaris Interactive job

- Useful for short tests or debugging
- Submit the job with `-I` (letter I for Interactive)
  - Debug queue:
    - `qsub -I -l select=1 -l walltime=1:00:00 -q debug -A ATPESC2025 -l filesystems=home:eagle`
  - Using reservation:
    - `qsub -I -l select=1 -l walltime=1:00:00 -q ATPESC -A ATPESC2025 -l filesystems=home:eagle`
    - Check the reservation queues with `pbs_rstat`
- Wait for job's shell prompt
  - Exit this shell to end your job
- From job's shell prompt, run just like in a script job,
  - `mpiexec -n 4 -ppn 4 ./a.out`
- After job expires, `mpiexec` will fail. *Check `qstat $PBS_JOBID`*

# Aurora/Polaris Scheduler – PBS Professional

- Primary commands
  - qsub
    - Request resources and start your script on the head node
    - -A Allocation
    - -l Options
    - -I Interactive mode
    - -q Which queue to submit otherwise default queue
  - qstat
    - Check on the status of requests
    - -Q List queues
    - -f <jobid> Detailed information about a job
    - -x <jobid> Information about a completed job
  - qalter
    - Update your requests
  - qdel
    - Cancel/delete jobs
  - pbs\_rstat
    - Check reservations

# Aurora Queues

- Aurora had 5 main queues
  - <https://docs.alcf.anl.gov/aurora/running-jobs-aurora/>
  - debug
    - 2 nodes max
    - 1 hour max
    - 5 minutes min
  - debug-scaling
    - 31 nodes max
    - 1 hour max
    - 5 minutes min
  - prod
    - 1 nodes min
    - 2048 nodes max
    - 5 minutes min
    - 18 hours max
  - prod-large (*by request only*)
    - 2048 nodes min
    - 10624 nodes max
    - 5 minutes min
    - 24 hours max
  - Visualization (*by request only*)
    - 1 nodes min
    - 32 nodes max
    - 5 minutes min
    - 8 hours max

# Polaris Queues

- Polaris had 3 main queues
  - <https://docs.alcf.anl.gov/polaris/running-jobs/>
  - debug
    - 2 nodes max
    - 1 hour max
    - 5 minutes min
  - debug-scaling
    - 10 nodes max
    - 1 hour max
    - 5 minutes min
  - prod
    - 10 nodes min
    - 496 nodes max
    - 5 minutes min
    - 24 hours max

# Cryptocard tips for ALCF systems

- Tokens
  - Physical Token: The displayed value is a hex string. Type your PIN followed by all letters as CAPITALS.
  - Mobile Token: Type the displayed numbers from your mobile app
- If you fail to authenticate the first time, you may have typed it incorrectly
  - Try again with the **same crypto string** (do NOT refresh button again)
- If you fail again, try a different ALCF host with a fresh crypto #
  - A successful login resets your count of failed logins
- Too many failed logins → your account locked
  - Symptom: You get password prompt but login denied even if it is correct
- Too many failed logins from a given IP → the IP will be blocked
  - Symptom: connection attempt by ssh or web browser will just time out

# ALCF References

- Sample files
  - `/flare/ATPESC2025/EXAMPLES/track-0-getting-started/` #on Aurora
  - `/eagle/ATPESC2025/EXAMPLES/track-0-getting-started/` #on Polaris
- Online docs
  - <https://www.alcf.anl.gov/support-center>
  - <https://docs.alcf.anl.gov/aurora/getting-started-on-aurora/>
  - <https://docs.alcf.anl.gov/polaris/getting-started/>
  - ALCF Beginners Guide (Instructions, examples, and videos)
    - <https://github.com/argonne-lcf/ALCFBeginnersGuide/tree/master/aurora>
    - <https://github.com/argonne-lcf/ALCFBeginnersGuide/tree/master/polaris>
    - <https://www.alcf.anl.gov/events/getting-started-aurora-and-polaris-bootcamp>

# OLCF Odo System

# Frontier System overview

- HPE Cray EX Supercomputer architecture
- 74 cabinets, 128 nodes per cabinet (9408 nodes)
- 3<sup>rd</sup> Gen AMD EPYC 64-core CPU
- 4 AMD Instinct MI250X GPUs
- HPE Slingshot interconnect
- Peak 1.206 Exaflops on HPL benchmark
- Cray, AMD, and GNU software stacks
- Orion - 679 PB multi-tier Lustre filesystem
- NFS storage (/ccs/home, /ccs/proj)



# Odo System overview

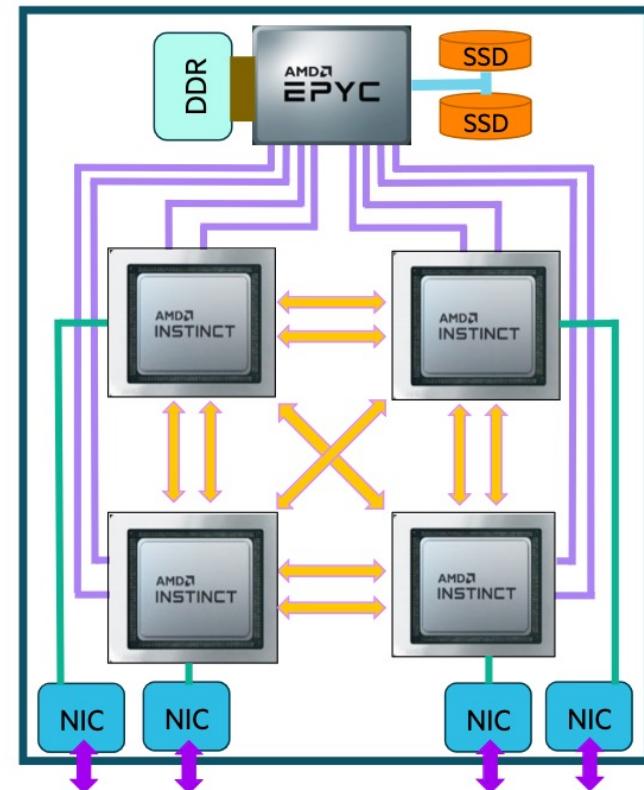
- HPE Cray EX Supercomputer architecture
- 74 cabinets, 128 nodes per cabinet (9408 nodes)

Odo is 30 Frontier nodes, and uses the GPFS filesystem (/gpfs/alpine2), and NFS on the Open Enclave for home areas (/ccsopen/home, /ccsopen/proj)

- Orion - 679 PB multi-tier Lustre filesystem
- NFS storage (/ccs/home, /ccs/proj)

# Odo Compute Node Configuration

- 1x AMD Optimized 3rd Gen EPYC 64 core processor
  - 2 hardware threads per physical core,
  - 2.0GHz base clock, 3.7GHz boost clock
- 512 GB DDR4 memory with 205 GB/s peak bandwidth
- 2x NVMe 2TB SSDs, peak 8 GB/s R, 4 GB/s W, >1.5M IOPs
- 4x AMD MI250X Instinct GPUs
  - 128 GB High-Bandwidth Memory (HBM2E)
  - 3.2 TB/s peak bandwidth
  - 53 TFLOPS double-precision peak for modeling & simulation
  - 2 Graphic Compute Dies (GCDs)
- AMD Infinity Fabric between CPU and GPUs
  - Peak host-to-device (H2D) and device-to-host (D2H) data transfers of 36+36 GB/s per link
- AMD Infinity Fabric between MI250Xs
  - Peak device-to-device bandwidth of 50+50 GB/s per link, low latency
- 4x HPE Slingshot Interconnect 200 GbE NICs
  - Provides 100 GB/s to other nodes, 25 GB/s per port



# Odo Module Commands

Command	Description
<code>module -t list</code>	Shows a terse list of the currently loaded modules
<code>module avail</code>	Shows a table of the currently available modules
<code>module help &lt;modulename&gt;</code>	Shows help information about <code>&lt;modulename&gt;</code>
<code>module show &lt;modulename&gt;</code>	Shows the environment changes made by the <code>&lt;modulename&gt;</code> modulefile
<code>module spider &lt;string&gt;</code>	Searches all possible modules according to <code>&lt;string&gt;</code>
<code>module load &lt;modulename&gt; [...]</code>	Loads the given <code>&lt;modulename&gt;</code> (s) into the current environment
<code>module use &lt;path&gt;</code>	Adds <code>&lt;path&gt;</code> to the modulefile search cache and <code>MODULESPATH</code>
<code>module unuse &lt;path&gt;</code>	Removes <code>&lt;path&gt;</code> from the modulefile search cache and <code>MODULESPATH</code>
<code>module purge</code>	Unloads all modules
<code>module reset</code>	Resets loaded modules to system defaults
<code>module update</code>	Reloads all currently loaded modules

# Odo: Compiling

- Cray, AMD, and GCC compilers are provided through modules. The Cray and AMD compilers are both based on LLVM/Clang. There is also a system/OS versions of GCC available in /usr/bin. The table below lists details about each of the module-provided compilers.

Vendor	Programming Environment	Compiler Module	Language	Compiler Wrapper	Compiler
Cray	PrgEnv-cray	cce	C	cc	craycc
			C++	CC	craycxx or crayCC
			Fortran	ftn	crayftn
AMD	PrgEnv-amd	amd	C	cc	am clang
			C++	CC	am clang++
			Fortran	ftn	am flang
GCC	PrgEnv-gnu	gcc-native or gcc (<12.3)	C	cc	gcc
			C++	CC	g++
			Fortran	ftn	gfortran

- [https://docs.olcf.ornl.gov/systems/frontier\\_user\\_guide.html#compiling](https://docs.olcf.ornl.gov/systems/frontier_user_guide.html#compiling)

# Odo: Slurm Workload Manager

- Slurm Commands (vs. LSF commands)

Command	Action/Task	LSF Equivalent
<code>squeue</code>	Show the current queue	<code>bjobs</code>
<code>sbatch</code>	Submit a batch script	<code>bsub</code>
<code>salloc</code>	Submit an interactive job	<code>bsub -Is \$SHELL</code>
<code>srun</code>	Launch a parallel job	<code>jrun</code>
<code>sinfo</code>	Show node/partition info	<code>bqueues</code> or <code>bhosts</code>
<code>sacct</code>	View accounting information for jobs/job steps	<code>bacct</code>
<code>scancel</code>	Cancel a job or job step	<code>bkill</code>
<code>scontrol</code>	View or modify job configuration.	<code>bstop</code> , <code>bresume</code> , <code>bmod</code>

- [https://docs.olcf.ornl.gov/systems/frontier\\_user\\_guide.html#slurm](https://docs.olcf.ornl.gov/systems/frontier_user_guide.html#slurm)

# Odo: Batch Scripts

- To submit a batch script,
  - Use the command, `sbatch myjob.sl`
- Description of the example

Line	Description
2	OLCF project to charge
3	Job name
4	Job standard output file (%x: the job name, %j: the Job ID)
5	Walltime requested (in HH:MM:SS format).
6	Partition (queue) to use
7	Number of compute nodes requested
9	Change into the run directory
10	Copy the input file into place
11	Run the job ( add layout details )
12	Copy the output file to an appropriate location.

```
1 #!/bin/bash
2 #SBATCH -A ABC123
3 #SBATCH -J RunSim123
4 #SBATCH -o %x-%j.out
5 #SBATCH -t 1:00:00
6 #SBATCH -p batch
7 #SBATCH -N 1024
8
9 cd $MEMBERWORK/abc123/Run.456
10 cp $PROJWORK/abc123/RunData/Input.456 ./Input.456
11 srun ...
12 cp my_output_file $PROJWORK/abc123/RunData/Output.456
```

- [https://docs.olcf.ornl.gov/systems/frontier\\_user\\_guide.html#batch-scripts](https://docs.olcf.ornl.gov/systems/frontier_user_guide.html#batch-scripts)

# Odo: Interactive job

- Useful for short tests or debugging
- Submit the job with `salloc`
  - `salloc -A trn038 -J RunSim123 -t 1:00:00 -p batch -N 1`
- Wait for job's shell prompt
  - Exit this shell to end your job
- From job's shell prompt, run just like in a script job,
  - `srun -N 1 -n 8 --ntasks-per-node=8 ./a.out`
- [https://docs.olcf.ornl.gov/systems/frontier\\_user\\_guide.html#interactive-jobs](https://docs.olcf.ornl.gov/systems/frontier_user_guide.html#interactive-jobs)

# Odo: Monitoring and Modifying Jobs

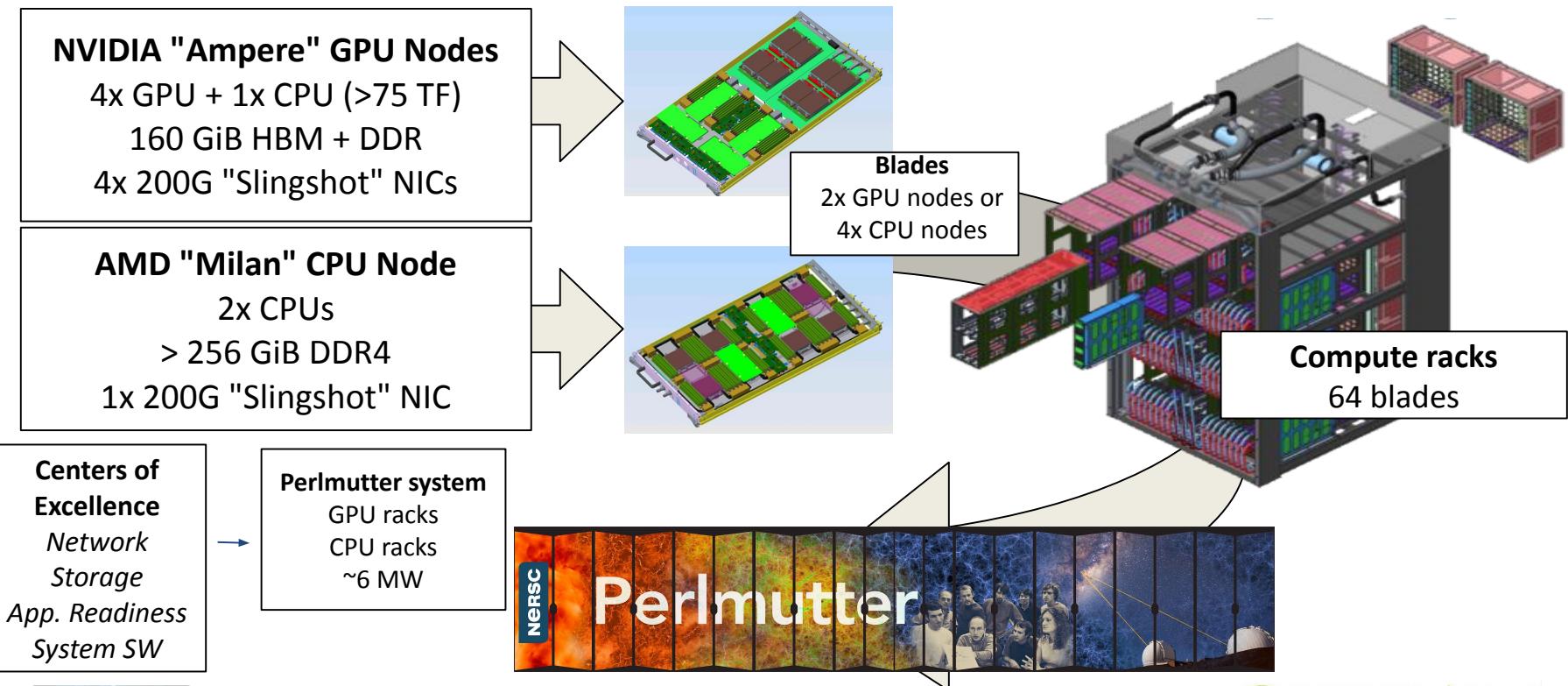
- Primary commands
  - `scancel`: Cancel or Signal a Job
  - `squeue`: View the queue
    - `squeue -l` : Show all jobs currently in the queue
    - `squeue -l -u $USER` : Show all of your jobs currently in the queue
  - `scontrol show job`: get Detailed Job information
- [https://docs.olcf.ornl.gov/systems/frontier\\_user\\_guide.html#monitoring-and-modifying-batch-jobs](https://docs.olcf.ornl.gov/systems/frontier_user_guide.html#monitoring-and-modifying-batch-jobs)

# Odo: Running MPI Applications

- Jobs run directly on the compute nodes. The `srun` command is used to execute an MPI binary on one or more compute nodes in parallel.
- `srun`
  - `-N` Number of nodes
  - `-n` Total number of MPI tasks
  - `-c` Logical cores per MPI task
  - `--ntasks-per-node=<ntasks>` A maximum count of tasks per node
  - `--gpus` Specify the number of GPUs
  - `--gpu-per-node` Specify the number of GPUs per node
- [https://docs.olcf.ornl.gov/systems/frontier\\_user\\_guide.html#srun](https://docs.olcf.ornl.gov/systems/frontier_user_guide.html#srun)

# NERSC Perlmutter System

# Perlmutter System Configuration

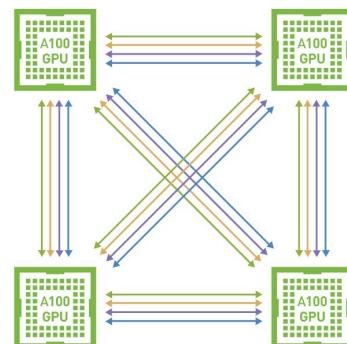
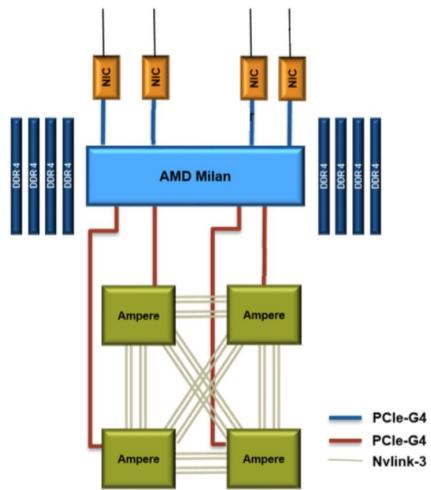


# Perlmutter Nodes

## GPU Nodes:

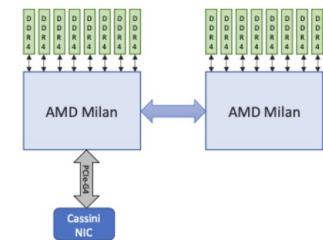
- Single [AMD EPYC 7763 \(Milan\)](#) CPU
- 64 cores per CPU
- Four [NVIDIA A100](#) (Ampere) GPUs
- PCIe 4.0 GPU-CPU connection
- PCIe 4.0 NIC-CPU connection
- 4 [HPE Slingshot 11](#) NICs
- 256 GB of DDR4 DRAM
- 40 GB of HBM per GPU with
- 1555.2 GB/s GPU memory bandwidth
- 204.8 GB/s CPU memory bandwidth
- 12 third generation NVLink links between each pair of gpus
- 25 GB/s/direction for each link

Data type	GPU TFLOPS
FP32	19.5
FP64	9.7
TF32 (tensor)	155.9
FP16 (tensor)	311.9
FP64 (tensor)	19.5



## CPU Nodes:

- 2x [AMD EPYC 7763 \(Milan\)](#) CPUs
- 64 cores per CPU
- AVX2 instruction set
- 512 GB of DDR4 memory total
- 204.8 GB/s memory bandwidth per CPU
- 1x [HPE Slingshot 11](#) NIC
- PCIe 4.0 NIC-CPU connection
- 39.2 GFlops per core
- 2.51 TFlops per socket
- 4 NUMA domains per socket (NPS=4)



# Perlmutter Modules Environment

- LMod is used to manage the user environment
  - <https://docs.nersc.gov/environment/lmod/>

module	
list	To list the modules in your environment
spider <name>	To list available modules with <name> as substring, and how to load
load/unload ..	To load or unload module
swap ... ..	To swap modules
show/display ..	To see what a module loads, what env a module sets
whatis ..	Display the module file information
help ..	General help: \$module help Information about a module: \$ module help PrgEnv-cray

# Perlmutter Software Environment

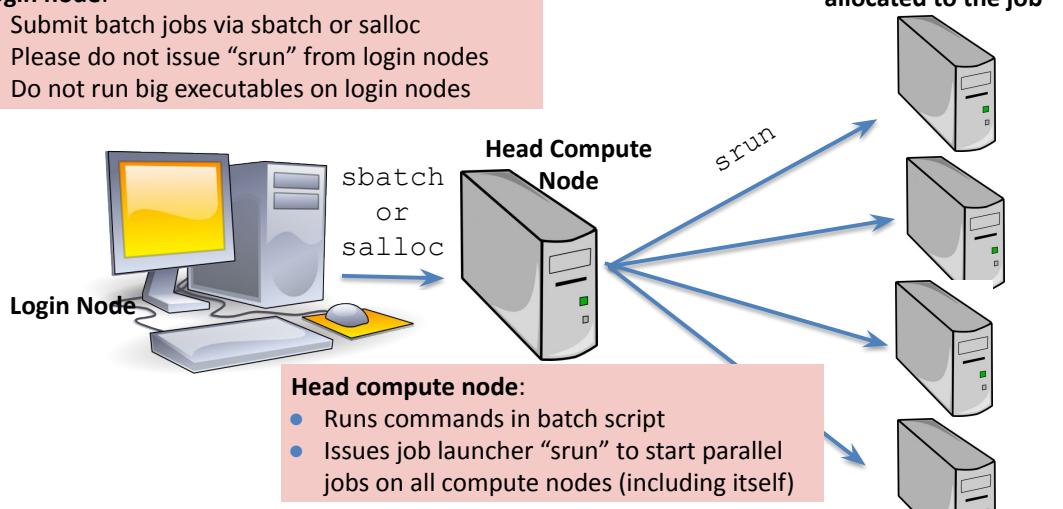
- Available compilers: GNU, Nvidia, CCE, (and Intel, in progress)
- It calls native compilers for each compiler (such as gfortran, gcc, g++, etc.) underneath.
  - Do not use native compilers directly
  - ftn for Fortran codes: **ftn my\_code.f90**
  - cc for C codes: **cc my\_code.c**
  - CC for C++ codes: **CC my\_code.cc**
- Compiler wrappers add header files and link in MPI and other loaded Cray libraries by default
  - Builds applications dynamically by default.

<https://docs.nersc.gov/development/compilers/wrappers/>

# Perlmutter: Launching Parallel Jobs with Slurm

## Login node:

- Submit batch jobs via sbatch or salloc
- Please do not issue “srun” from login nodes
- Do not run big executables on login nodes



Head Compute Node

Other Compute Nodes  
allocated to the job

Login Node

## Head compute node:

- Runs commands in batch script
- Issues job launcher “srun” to start parallel jobs on all compute nodes (including itself)

## my\_batch\_script:

```
#!/bin/bash
#SBATCH -q debug
#SBATCH -N 2
#SBATCH -t 10:00
#SBATCH -C cpu
##SBATCH -L SCRATCH
##SBATCH -J myjob
srun -n 64 ./helloWorld
```

## To run via batch queue

% sbatch my\_batch\_script

## To run via interactive batch

% salloc -N 2 -q interactive -C cpu -t 10:00

<wait\_for\_session\_prompt. Land on a compute node>

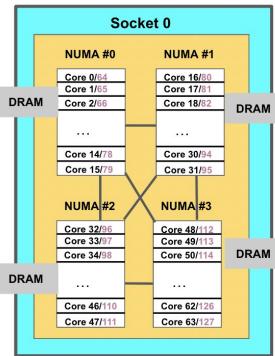
% srun -n 64 ./helloWorld

# Perlmutter: CPU and GPU Compute Nodes Affinity

	Perlmutter CPU	CPU on Perlmutter GPU
Physical cores	128	64
Logical CPUs per physical core	2	2
Logical CPUs per node	256	128
NUMA domains	8	4
-c value for srun	$2 * \text{floor}(128/\text{tpn})$	$2 * \text{floor}(64/\text{tpn})$

tpn = Number of MPI tasks per node

CPU on Perlmutter GPU



- Correct process, thread and memory affinity is critical for getting optimal performance on Perlmutter CPU and GPU
  - Process Affinity: bind MPI tasks to CPUs
  - Thread Affinity: bind threads to CPUs allocated to its MPI process
  - Memory Affinity: allocate memory from specific NUMA domains
- Both **-c xx** and **--cpu-bind=cores** are essential, otherwise multiple processes may land on the same core, while other cores are idle, hurting performance badly
- <https://docs.nersc.gov/jobs/affinity/>

# Perlmutter: Shared QOS for the reserved nodes

The “shared” QOS allows multiple executables from different users to share a node

To use nodes to be shared by multiple users, ATPESC attendees can request with salloc or sbatch with flags such as:

```
-C gpu -q shared -A trn015 -N 1 -c 32 -G 1 -t 60:00
```

Please notice the `-q shared` and `-c 32 -G 1` options. It will get each user 1/4 of node CPU and 1 GPU. And users can run CPU or GPU jobs in this allocation.

<https://docs.nersc.gov/jobs/examples/#shared>

To check reservations during ATPESC, use scontrol:

```
scontrol show res
```

To use a reservation during ATPESC, use the following:

```
--reservation=atpesc_xxxxx
```

Perlmutter Job script generator:

[https://my.nersc.gov/script\\_generator.php](https://my.nersc.gov/script_generator.php)

The screenshot shows the 'Jobscript Generator' page on the 'my.nersc.gov/script\_generator.php' website. The left sidebar contains links like Dashboard, Jobs, Center Status, File Browser, Service Tickets, Data Dashboard, PI Toolbox, Jupyter Hub, NERSC Homepage, Documentation Portal, and Accounts Portal. The main right panel has sections for 'Job Information', 'Machine' (set to 'Perlmutter - CPU'), 'Application Name' (set to 'myapp.x'), 'Job Name' (empty), 'Email Address' (empty), 'Quality of Service' (set to 'regular'), and 'Wallclock Time'. A large red circle highlights the 'Command Line' section which displays a bash script template:

```
#!/bin/bash
#SBATCH -N 4
#SBATCH -C cpu
#SBATCH -q regular
#SBATCH -t 01:30:00

#OpenMP settings:
export OMP_NUM_THREADS=4
export OMP_PLACES=threads
export OMP_PROC_BIND=spread

#run the application:
srun -n 32 -c 32 --cpu_bind=cores myapp.x
```

# Perlmutter GPU Queue Policy

QOS	Max nodes	Max time (hrs)	Min time (hrs)	Submit limit	Run limit	Priority
regular	-	48	-	5000	-	medium
interactive <sup>1</sup>	4	4	-	2	2	high
shared_interactive	0.5	4	-	2	2	high
jupyter	4	6	-	1	1	high
debug	8	0.5	-	5	2	medium
shared <sup>2</sup>	0.5	48	-	5000	-	medium
preempt <sup>3</sup>	128	48 (preemptible after two hours)	2	5000	-	medium
debug_preempt	2	0.5 (preemptible after five minutes)	-	5	2	medium

# Perlmutter: Monitoring your Jobs

- Jobs are waiting in the queue until resources are available
- Overall job priorities are a combination of QOS, queue wait time, job size, wall time request, etc.
- You can monitor with
  - **squeue**: Slurm native command
  - **sqsh**: NERSC custom wrapper script
  - **sacct**: Query Completed and Pending Jobs
  - <https://docs.nersc.gov/jobs/monitoring/>

# Cheat Sheets

# ATPESC Resources

## ALCF – Aurora/ Polaris

-Project name: **ATPESC2025**

-**Note:** use your ALCF Username.

-**Support:** ALCF staff available to help you [via slack!!](#) and [support@alcf.anl.gov](mailto:support@alcf.anl.gov)

-**Reservations:** Please check the details of the reservations directly on Aurora/Polaris (**command:** *pbs\_rstat*)

-**Queue**

    -**check *pbs\_rstat*, or **default** for running without reservation**

-**User guide:**

    -<https://www.alcf.anl.gov/support-center>

    -<https://docs.alcf.anl.gov/aurora/getting-started-on-aurora/>

    -<https://docs.alcf.anl.gov/polaris/getting-started/>

# ATPESC Resources

## OLCF – Odo

- Project name: [trn038](#)
- Support: [help@olcf.ornl.gov](mailto:help@olcf.ornl.gov)
- Queue: **running without reservation**
- User guide:
  - Frontier User Guide: [https://docs.olcf.ornl.gov/systems/frontier\\_user\\_guide.html#frontier-user-guide](https://docs.olcf.ornl.gov/systems/frontier_user_guide.html#frontier-user-guide)
  - Odo User Guide: [https://docs.olcf.ornl.gov/systems/odo\\_user\\_guide.html](https://docs.olcf.ornl.gov/systems/odo_user_guide.html)

# ATPESC Resources

## NERSC – Perlmutter

- Project name: [trn015](#)
- Support: [help desk](#)
- Reservations: Please check the details of the reservations directly on Perlmutter (**command:** `scontrol show res`)
- Queue:
  - check `scontrol show res`, or `default` for running without reservation
- User guide: <https://docs.nersc.gov/>

# Questions?

- *Use this presentation as a reference during ATPESC!*
- Supplemental info will be posted as well

# Hands-on Exercises

## Hands-on exercise

- Aurora hands-on
- Polaris hands-on
- Odo hands-on
- Perlmutter hands-on

# Hands-on exercise: Aurora

```
$ ssh -Y {your_username}@aurora.alcf.anl.gov # Login to Polaris
$ module avail # See available modules
$ module list # See loaded modules
[jkwack@aurora-uan-0012:~] module li

Currently Loaded Modules:
 1) gcc-runtime/13.3.0-ghotoln (H)  5) gcc/13.3.0          9) hwloc/2.11.3-mpich-level-zero   13) cray-pals/1.4.0
 2) gmp/6.3.0-mtokfaw      (H)  6) oneapi/release/2025.0.5  10) yaksa/0.3-7ks5f26           (H)  14) cray-libpals/1.4.0
 3) mpfr/4.2.1-gkcdl5w      (H)  7) libiconv/1.17-jjpbb4sl  (H)  11) mpich/opt/develop-git.6037a7a
 4) mpc/1.3.1-rdrllvsl      (H)  8) libxml2/2.13.5        12) libfabric/1.22.0

$ qstat -u ${USER} # To see your jobs
$ pbs_rstat # Check reservation
```

```
jkwack@aurora-uan-0012:~] pbs_rstat
Resv ID Queue User State Start / Duration / End
-----
M6457232.aurora M6457232 bsallen@ DG Tue 09:00 / 32400 / Tue 18:00
S6703362.aurora S6703362 richp@au CO Wed 10:00 / 14400 / Wed 14:00
R6710521.aurora R6710521 richp@au CO Today 14:30 / 16200 / Today 19:00
R6710528.aurora R6710528 richp@au CO Mon 15:30 / 12600 / Mon 19:00
R6710533.aurora R6710533 richp@au CO Tue 08:30 / 37800 / Tue 19:00
R6710543.aurora R6710543 richp@au CO Wed 08:30 / 46800 / Wed 21:30
R6710664.aurora R6710664 richp@au CO Thu 09:00 / 43200 / Thu 21:00
R6710673.aurora R6710673 richp@au CO Fri 08:30 / 36000 / Fri 18:30
R6710677.aurora R6710677 richp@au CO Wed Aug 06 08:00 / 39600 / Wed Aug 06 19:00
R6710682.aurora R6710682 richp@au CO Thu Aug 07 09:00 / 43200 / Thu Aug 07 21:00
R6711188.aurora R6711188 richp@au CO Wed 21:00 / 21600 / Thu 03:00
R6711189.aurora R6711189 richp@au CO Thu 21:00 / 21600 / Fri 03:00
R6711190.aurora R6711190 richp@au CO Fri 21:00 / 32400 / Sat 06:00
R6711192.aurora R6711192 richp@au CO Tue Aug 05 21:00 / 21600 / Wed Aug 06 03:00
R6711193.aurora R6711193 richp@au CO Wed Aug 06 21:00 / 32400 / Thu Aug 07 06:00
R6711197.aurora R6711197 richp@au CO Thu Aug 07 21:00 / 32400 / Fri Aug 08 06:00
```

# Hands-on exercise: Aurora

```
$ qsub -I -l select=1 -l walltime=00:30:00 -l filesystems=home:flare -A ATPESC2025 -q ATPESC
jkwack@aurora-uau-0012:~> qsub -I -l select=1 -l walltime=00:30:00 -l filesystems=home:flare -A ATPESC2025 -q debug
qsub: waiting for job 6785771.aurora-pbs-0001.hostmgmt.cm.aurora.alcf.anl.gov to start
qsub: job 6785771.aurora-pbs-0001.hostmgmt.cm.aurora.alcf.anl.gov ready
jkwack@x4204c1s3b0n0:~>
```

```
$ cd /flare/ATPESC2025/usr/
$ mkdir $USER
$ cd $USER
$ cp -rf /flare/ATPESC2025/EXAMPLES/track-0-getting-started .
$ cd track-0-getting-started/aurora/
$ more Makefile
...
CC=mpicc

hellompi: hellompi.c
    which $(CC)
    $(CC) -g -O0 -o hellompi hellompi.c
...
```

# Hands-on exercise: Aurora

```
$ cat submit.sh
#!/bin/bash
#PBS -l select=1
#PBS -l walltime=00:30:00
##PBS -q debug
#PBS -l filesystems=home:flare
#PBS -A ATPESC2025
#PBS -q ATPESC

cd $PBS_O_WORKDIR

mpiexec -n 4 --ppn 4 ./hellompi
status=$?

echo "mpiexec status is $status"
exit $status
```

# Hands-on exercise: Aurora

```
$ mpicc -o hellompi hellompi.c          # Build the example  
$ make clean; make                         # Another way to build the example
```

```
$ mpiexec -n 4 --ppn 4 ./hellompi  
jkwack@x4204c1s3b0n0:~/ATPESC2025/track-0-getting-started/aurora> mpirun -n 4 --ppn 4 ./hellompi  
0: Hello!  
3: Hello!  
2: Hello!  
1: Hello!
```

```
$ module load xpu-smi  
$ xpu-smi discovery
```

## More references for Polaris

<https://github.com/argonne-lcf/ALCFBeginnersGuide/tree/master/aurora>  
<https://www.alcf.anl.gov/events/getting-started-aurora-and-polaris-bootcamp>  
<https://docs.alcf.anl.gov/aurora/getting-started-on-aurora/>

# Hands-on exercise: Polaris

```
$ ssh -Y {your_username}@polaris.alcf.anl.gov # Login to Polaris  
$ module avail # See available modules  
$ module list # See loaded modules
```

```
| jkwack@polaris-login-04:~> module list  
  
Currently Loaded Modules:  
 1) nvhpc/23.9      4) cray-mpich/8.1.28  7) cray-libpals/1.3.4  10) libfabric/1.15.2.0      13) darshan/3.4.4  
 2) craype/2.7.30    5) cray-pmi/6.1.13    8) craype-x86-milan   11) craype-network-ofi   14) xalt/3.0.2-202408282050  
 3) cray-dsmml/0.2.2  6) cray-pals/1.3.4    9) PrgEnv-nvhpc/8.5.0  12) perftools-base/23.12.0
```

```
$ qstat -u ${USER} # To see your jobs  
$ pbs_rstat # Check reservation
```

```
| jkwack@polaris-login-04:~> pbs_rstat  
Resv ID     Queue     User     State          Start / Duration / End  
-----  
R5594657.polari R5594657  mluczkow CO  Sun Jul 27 14:30 / 16200 / Sun Jul 27 19:00  
R5594658.polari R5594658  mluczkow CO  Mon Jul 28 15:30 / 12600 / Mon Jul 28 19:00  
R5594659.polari R5594659  mluczkow CO  Tue Jul 29 08:30 / 37800 / Tue Jul 29 19:00  
R5594660.polari R5594660  mluczkow CO  Wed Jul 30 08:30 / 46800 / Wed Jul 30 21:30  
R5594661.polari R5594661  mluczkow CO  Thu Jul 31 09:00 / 43200 / Thu Jul 31 21:00  
R5594662.polari R5594662  mluczkow CO  Fri Aug 01 08:30 / 36000 / Fri Aug 01 18:30  
R5594663.polari R5594663  mluczkow CO  Mon Aug 04 17:30 / 12600 / Mon Aug 04 21:00  
R5594664.polari R5594664  mluczkow CO  Tue Aug 05 08:30 / 37800 / Tue Aug 05 19:00  
R5594665.polari R5594665  mluczkow CO  Wed Aug 06 08:00 / 39600 / Wed Aug 06 19:00  
R5594666.polari R5594666  mluczkow CO  Thu Aug 07 09:00 / 43200 / Thu Aug 07 21:00  
R5594667.polari R5594667  mluczkow CO  Wed Jul 30 21:00 / 21600 / Thu Jul 31 03:00  
R5594668.polari R5594668  mluczkow CO  Thu Jul 31 21:00 / 21600 / Fri Aug 01 03:00  
R5594669.polari R5594669  mluczkow CO  Fri Aug 01 21:00 / 32400 / Sat Aug 02 06:00  
R5594670.polari R5594670  mluczkow CO  Mon Aug 04 21:00 / 21600 / Tue Aug 05 03:00  
R5594671.polari R5594671  mluczkow CO  Tue Aug 05 21:00 / 21600 / Wed Aug 06 03:00  
R5594672.polari R5594672  mluczkow CO  Wed Aug 06 21:30 / 30600 / Thu Aug 07 06:00
```

# Hands-on exercise: Polaris

```
$ qsub -I -l select=1 -l walltime=00:30:00 -l filesystems=home:grand -A ATPESC2025 -q ATPESC
[jkwack@polaris-login-04:~> qsub -I -l select=1 -l walltime=00:30:00 -l filesystems=home:grand -A ATPESC2025 -q debug
qsub: waiting for job 5684646.polaris-pbs-01.hsn.cm.polaris.alcf.anl.gov to start
qsub: job 5684646.polaris-pbs-01.hsn.cm.polaris.alcf.anl.gov ready
```

```
Currently Loaded Modules:
 1) nvhpc/23.9      4) cray-mpich/8.1.28   7) cray-libpals/1.3.4 10) libfabric/1.15.2.0      13) darshan/3.4.4
 2) craype/2.7.30    5) cray-pmi/6.1.13     8) craype-x86-milan   11) craype-network-ofi    14) xalt/3.0.2-202408282050
 3) cray-dsml/0.2.2   6) cray-pals/1.3.4    9) PrgEnv-nvhpc/8.5.0 12) perftools-base/23.12.0
```

```
$ cd /eagle/ATPESC2025/usr/
$ mkdir $USER
$ cd $USER
$ cp -rf /eagle/ATPESC2025/EXAMPLES/track-0-getting-started .
$ cd track-0-getting-started/polaris/
$ more Makefile
...
CC=cc

hellompi: hellompi.c
which $(CC)
$(CC) -g -O0 -o hellompi hellompi.c
...
```

# Hands-on exercise: Polaris

```
$ cat submit.sh
#!/bin/bash
#PBS -l select=1
#PBS -l walltime=00:30:00
##PBS -q debug
#PBS -l filesystems=home:eagle
#PBS -A ATPESC2025
#PBS -q ATPESC

cd $PBS_O_WORKDIR

mpiexec -n 4 --ppn 4 ./hellompi
status=$?

echo "mpiexec status is $status"
exit $status
```

# Hands-on exercise: Polaris

```
$ cc -o hellompi hellompi.c          # Build the example  
$ make clean; make                   # Another way to build the example  
  
$ mpiexec -n 4 --ppn 4 ./hellompi  
[jkwack@x3004c0s25b1n0:/eagle/ATPESC2025/usr/jkwack/track-0-getting-started/polaris> mpiexec -n 4 --ppn 4 ./hellompi  
3: Hello!  
0: Hello!  
1: Hello!  
2: Hello!
```

```
$ nvidia-smi
```

## More references for Polaris

- <https://github.com/argonne-lcf/ALCFBeginnersGuide/tree/master/polaris>
- <https://www.alcf.anl.gov/events/getting-started-aurora-and-polaris-bootcamp>
- <https://docs.alcf.anl.gov/polaris/getting-started/>

# Hands-on exercise: Odo

```
$ ssh -Y {your_username}@odo.olcf.ornl.gov          # Login to Odo
$ module avail                                     # See available modules
$ module list                                      # See loaded modules
[jkwack@login1.odo ~]$ module list
Currently Loaded Modules:
 1) craype-x86-trento      5) xpmem/2.11.3-1.3_gdbda01a1eb3d  9) cray-dsml/0.3.0    13) Core/25.01
 2) libfabric/1.22.0        6) cray-pmi/6.1.15            10) cray-mpich/8.1.31   14) tmux/3.4
 3) craype-network-ofi     7) cce/18.0.1              11) cray-libsci/24.11.0  15) darshan-runtime/3.4.6-mpi (E4S)
 4) perftools-base/24.11.0  8) craype/2.7.33           12) PrgEnv-cray/8.6.0   16) DefApps
Where:
E4S: E4S: Extreme-scale Scientific Software Stack (E4S) https://e4s.io/index.html
$ squeue -l -u $USER                                # To see your jobs
```

## Hands-on exercise: Odo

```
$ salloc -A trn038 -t 1:00:00 -p batch -N 1
[jkwack@login1.odo odo]$ salloc -A trn038 -t 1:00:00 -p batch -N 1
salloc: Granted job allocation 12243
salloc: Waiting for resource configuration
salloc: Nodes odo01 are ready for job
```

```
$ cp -r /ccsopen/proj/trn038/track-0-getting-started .
$ cd track-0-getting-started/odo/
$ more Makefile
...
CC=cc

hellompi: hellompi.c
    which $(CC)
    $(CC) -g -O0 -o hellompi hellompi.c
...
```

## Hands-on exercise: Odo

```
$ more submit.sh
#!/bin/bash
#SBATCH -p batch
#SBATCH -A trn038
#SBATCH -N 1
#SBATCH -t 60:00

srun -n 8 ./hellompi
status=$?

echo "mpiexec status is $status"
exit $status

$ cc -o hellompi hellompi.c          # Build the example
$ make clean; make                   # Another way to build the example
$ srun -n 4 ./hellompi
[jkwack@odo01:~/track-0-getting-started/odo> srun -n 4 ./hellompi
0: Hello!
1: Hello!
3: Hello!
2: Hello!
```

- More references for Odo and Frontier
  - [https://docs.olcf.ornl.gov/systems/odo\\_user\\_guide.html](https://docs.olcf.ornl.gov/systems/odo_user_guide.html)
  - [https://docs.olcf.ornl.gov/systems/frontier\\_user\\_guide.html](https://docs.olcf.ornl.gov/systems/frontier_user_guide.html)

# Hands-on exercise: Perlmutter

```
$ ssh -Y {your_username}@perlmutter.nersc.gov # Login to Perlmutter
$ module avail # See available modules
$ module list # See loaded modules
[jkwack2@perlmutter:login34:~> module li
Currently Loaded Modules:
 1) craype-x86-milan           6) cray-dsmmml/0.3.0    11) perftools-base/24.07.0  16) sqs/2.0
 2) libfabric/1.20.1            7) cray-libsci/24.07.0   12) cpe/24.07          17) darshan/default
 3) craype-network-ofi          8) cray-mpich/8.1.30   13) cudatoolkit/12.4
 4) xpmem/2.9.7-1.1_20250411150514__g191b5f8bea4c 9) craype/2.7.32      14) craype-accel-nvidia80
 5) PrgEnv-gnu/8.5.0           10) gcc-native/13.2    15) gpu/1.0

$ sqs # To see your jobs
$ salloc -q shared -C gpu -A trn015 -c 32 -G 1 -N 1 -t 30:00 --reservation=atpesc_Jul27
train580@perlmutter:login37:~> salloc -q shared -C gpu -A ntrain5 -c 32 -G 1 -N 1 -t 30:00
salloc: Pending job allocation 28681463
salloc: job 28681463 queued and waiting for resources
salloc: job 28681463 has been allocated resources
salloc: Granted job allocation 28681463
salloc: Waiting for resource configuration
salloc: Nodes nid002808 are ready for job
train580@nid002808:~>
```

## Hands-on exercise: Perlmutter

```
$ cp -r /global/homes/j/jkwack2/track-0-getting-started/perlmutter .
$ cd perlmutter/

$ more Makefile
...
CC=cc

hellompi: hellompi.c
    which $(CC)
    $(CC) -g -O0 -o hellompi hellompi.c
...
```

# Hands-on exercise: Perlmutter

```
$ more submit.sh
#!/bin/bash
#SBATCH -N 1
#SBATCH -C gpu
#SBATCH -A trn015
#SBATCH -q shared
#SBATCH --reservation=atpesc_Jul27
#SBATCH -c 32
#SBATCH -G 1
#SBATCH -t 60:00

srun -n 4 -c 4 ./hellompi
status=$?

echo "mpiexec status is $status"
exit $status

$ sbatch submit.sh
```

## Hands-on exercise: Perlmutter

```
$ cc -o hellompi hellompi.c          # Build the example  
$ make clean; make                   # Another way to build the example  
  
$ srun -n 4 -c 4 ./hellompi  
[train580@nid002808:~/perlmutter> srun -n 4 -c 4 ./hellompi  
0: Hello!  
1: Hello!  
3: Hello!  
2: Hello!  
  
$ nvidia-smi
```

More references for Perlmutter

- <https://docs.nersc.gov/jobs/>

**Thank you!**

## Supplemental Info

•

**ATPESC2025**

[extremecomputingtraining.anl.gov](http://extremecomputingtraining.anl.gov)



ARGONNE  
**ATPESC2025**  
EXTREME - SCALE COMPUTING

## ARGONNE TRAINING PROGRAM ON EXTREME-SCALE COMPUTING

Produced by Argonne National Laboratory, a U.S. Department of Energy Laboratory managed by UChicagoArgonne, LLC under contract DE-AC02-06CH11357.

Special thanks to the National Energy Research Scientific Computing Center (NERSC) and Oak Ridge Leadership Computing Facility (OLCF) for the use of their resources during the training event.

The U.S. Government retains for itself and others acting on its behalf a nonexclusive, royalty-free license in this video, with the rights to reproduce, to prepare derivative works, and to display publicly.



Argonne Leadership  
Computing Facility

[extremecomputingtraining.anl.gov](http://extremecomputingtraining.anl.gov)

