

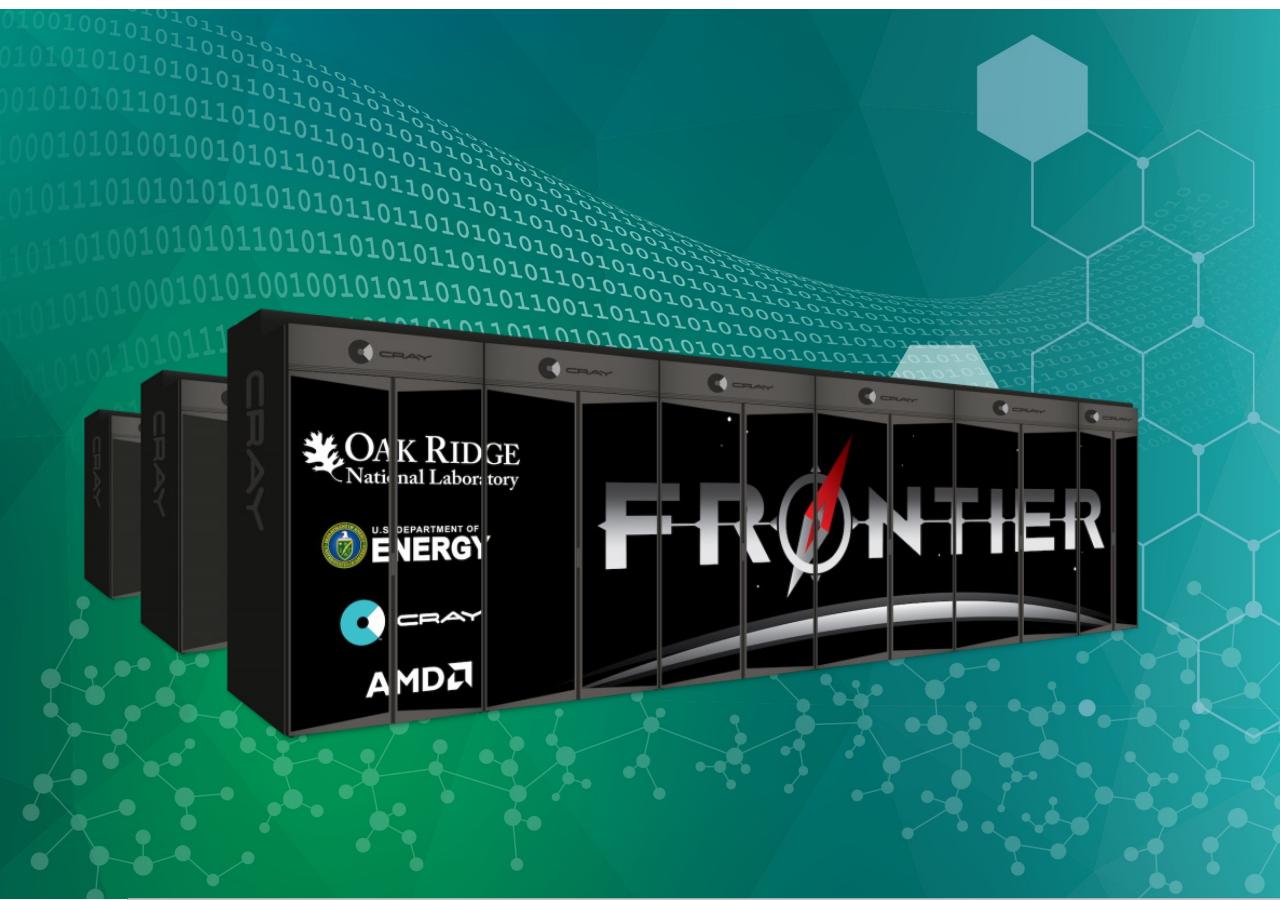
# Slurm on Frontier

Tom Papatheodore

HPC Engineer

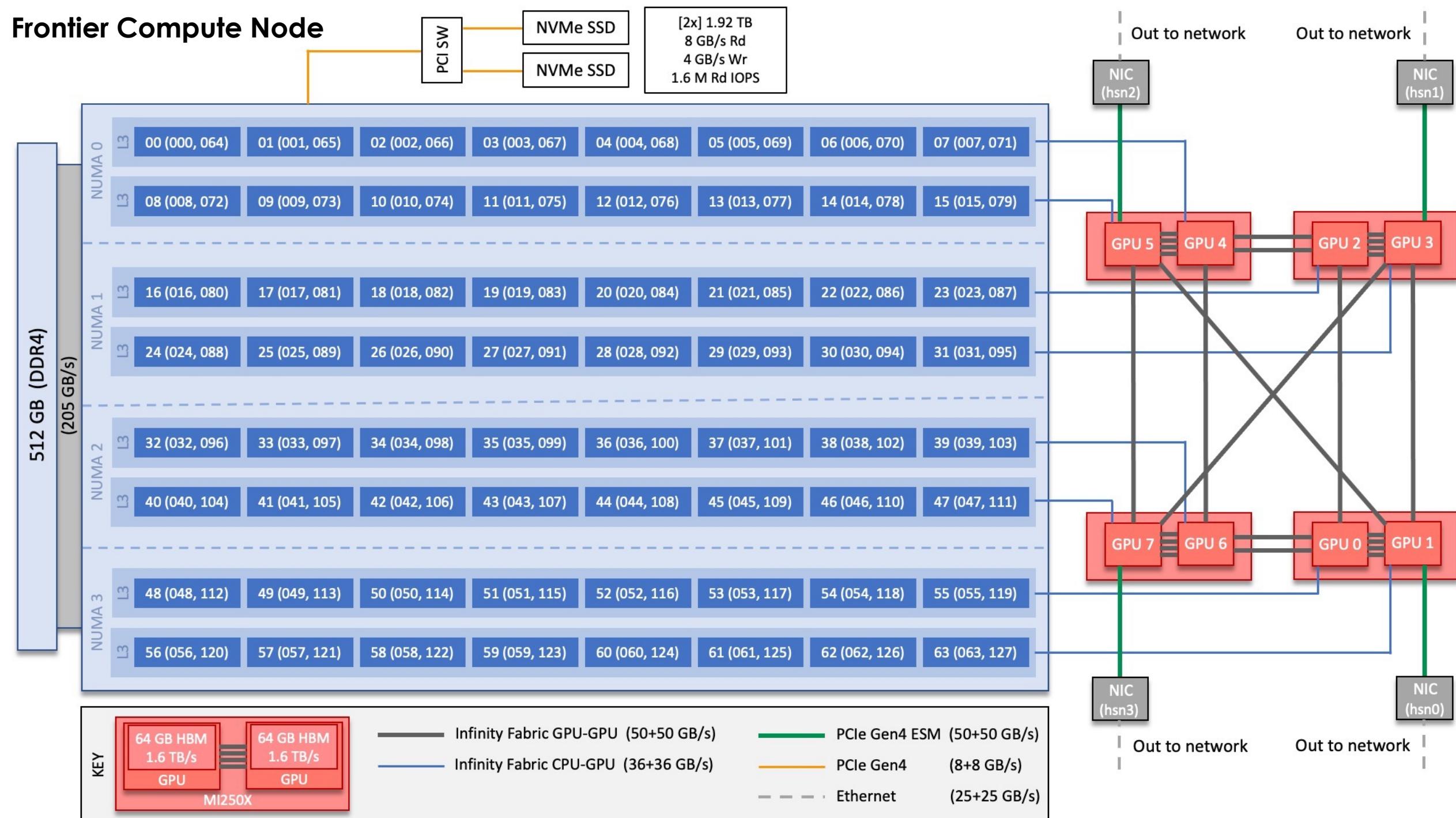
System Acceptance & User Environment Group  
Oak Ridge Leadership Computing Facility (OLCF)  
Oak Ridge National Laboratory (ORNL)

Frontier Training Workshop – February 15, 2023

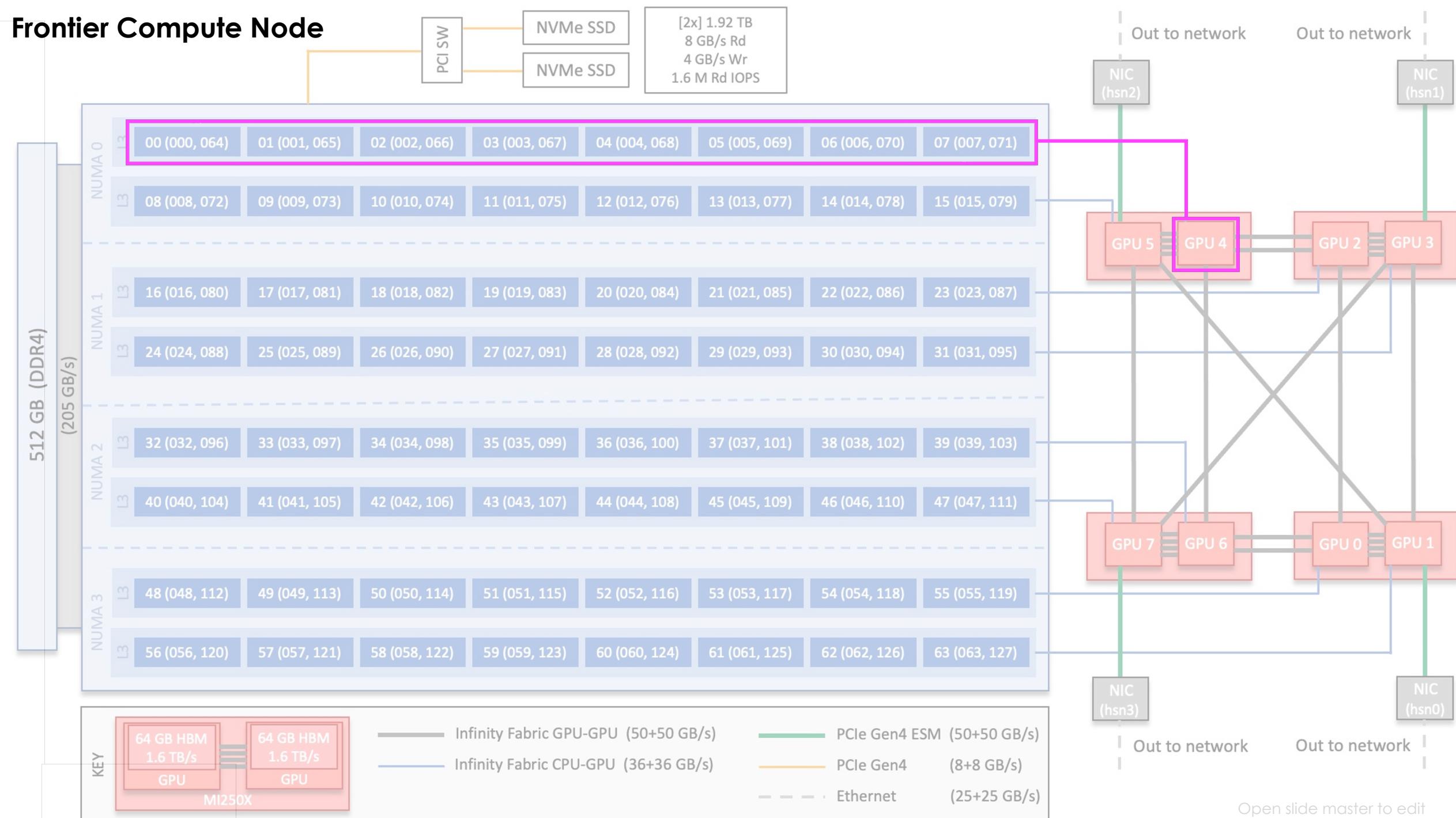


ORNL is managed by UT-Battelle LLC for the US Department of Energy

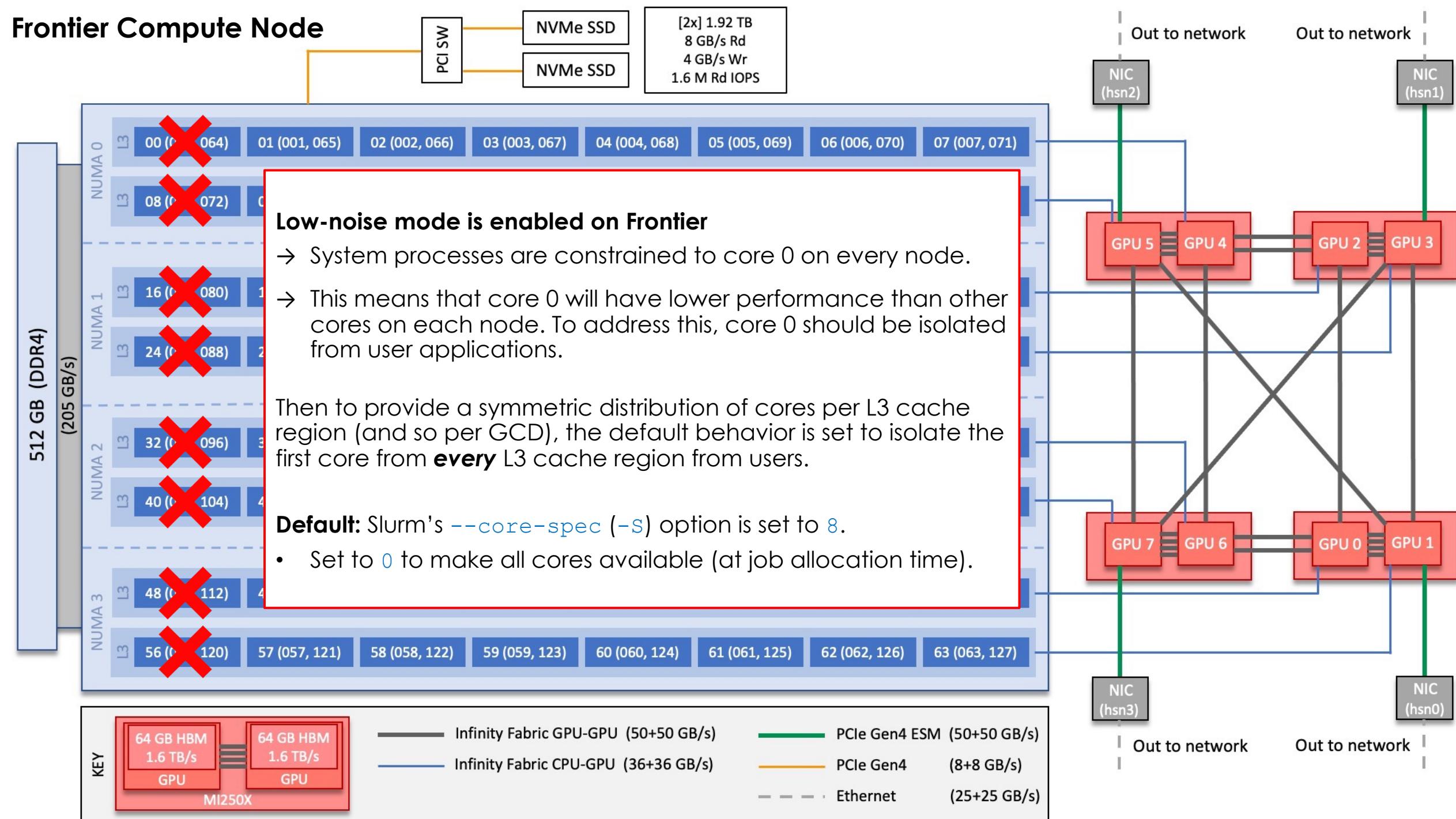
# Frontier Compute Node



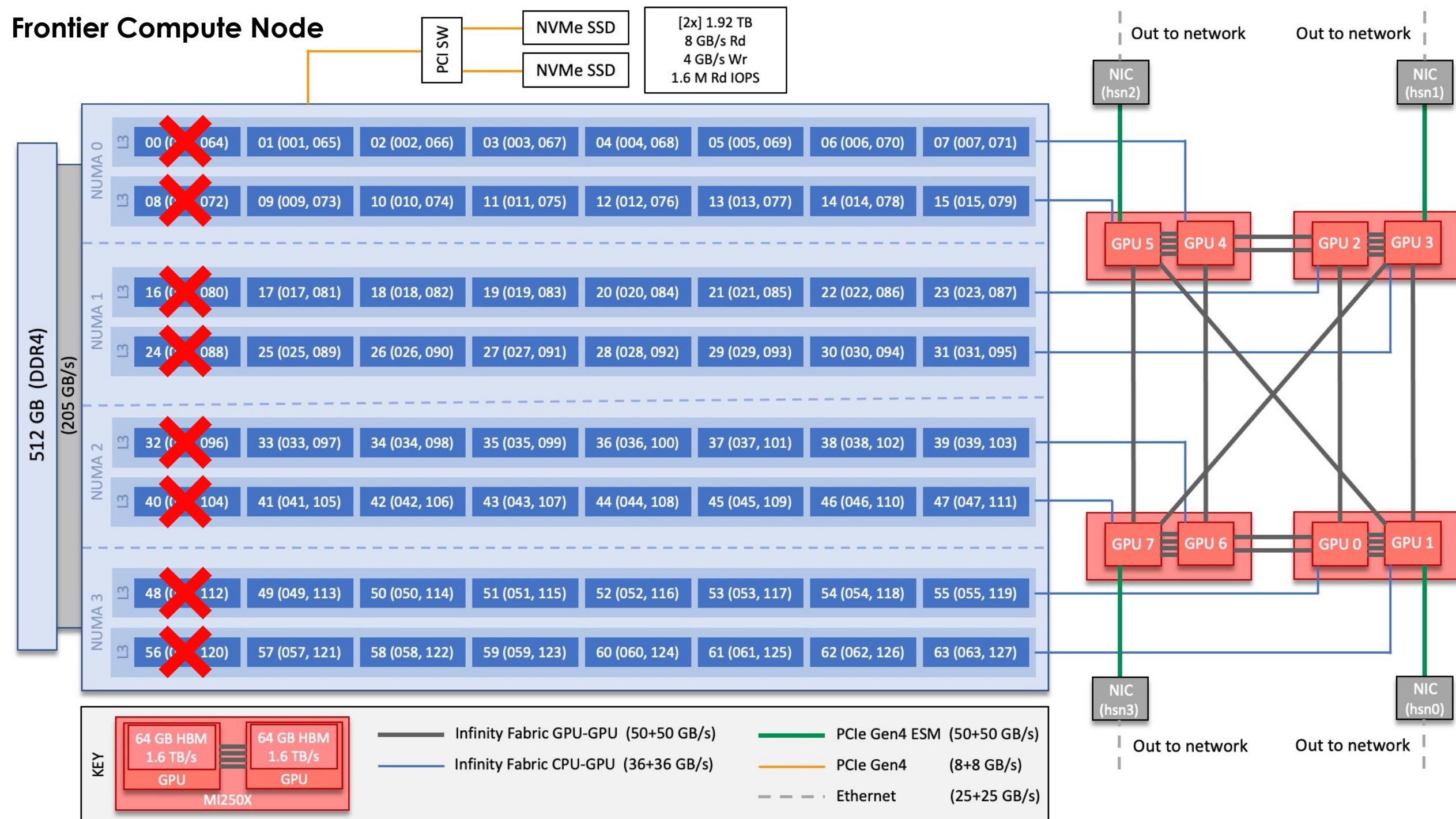
# Frontier Compute Node



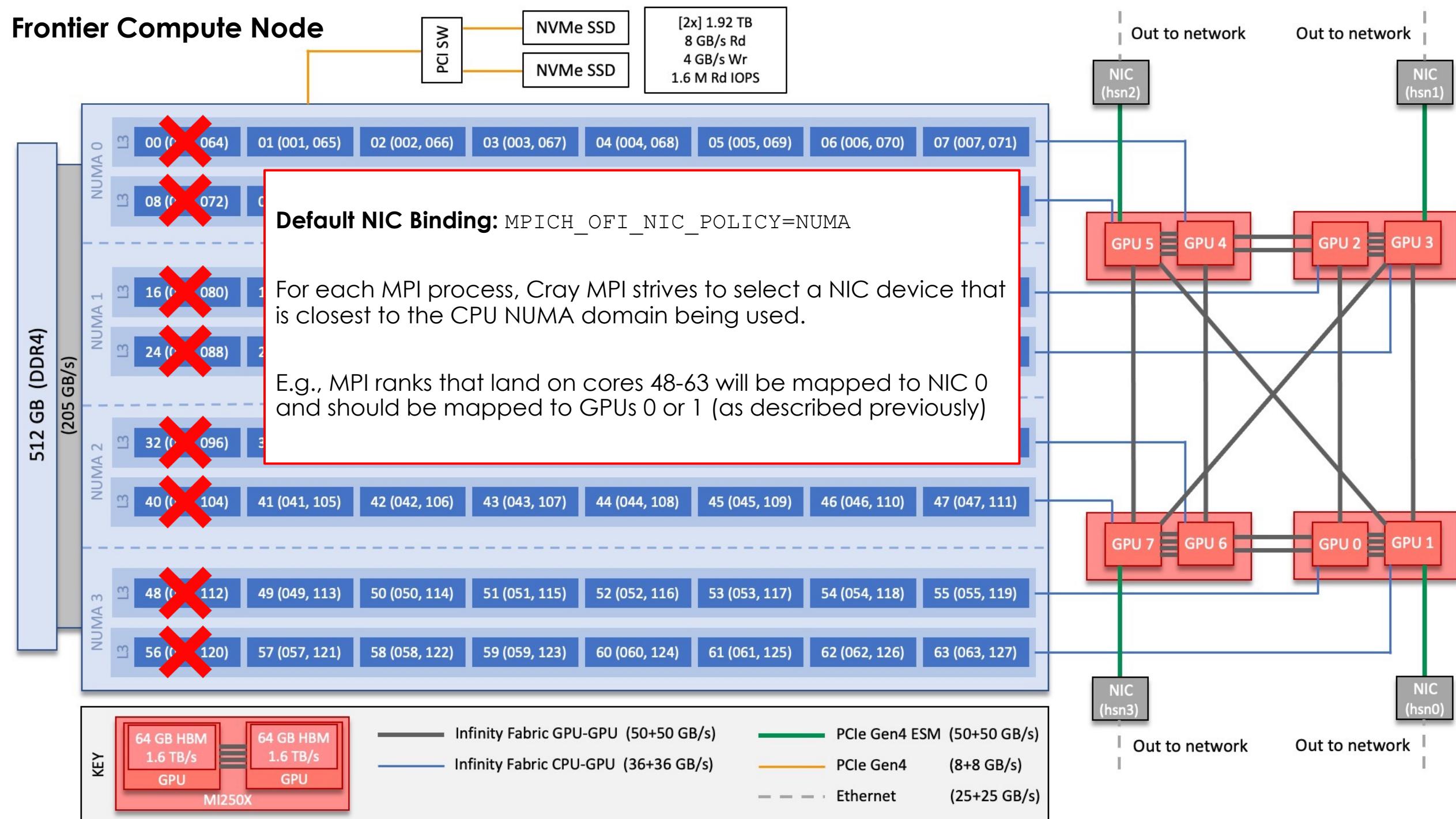
# Frontier Compute Node



# Frontier Compute Node



# Frontier Compute Node



# Slurm on Frontier and Crusher

Slurm provides 3 ways of submitting and launching jobs on Frontier/Crusher compute nodes: **batch scripts, interactive, and single-command**.

**NOTE:** Jobs launch on the allocated compute nodes, with the first compute node in the allocation serving as head-node.

Slurm commands associated with the 3 methods:

<code>sbatch</code>	Used to submit a batch script to allocate a Slurm job allocation. The script contains options preceded with <code>#SBATCH</code> .
<code>salloc</code>	Used to allocate an interactive Slurm job allocation, where one or more job steps (i.e., <code>srun</code> commands) can then be launched on the allocated resources (i.e., nodes).
<code>srun</code>	Used to run a parallel job (job step) on the resources allocated with <code>sbatch</code> or <code>salloc</code> , or used to create a resource allocation and run a job step in a single command.

# Batch Scripts

A **batch script** can be used to submit a job to run on the compute nodes at a later time. In this case, stdout and stderr will be written to a file(s) that can be opened after the job completes.

Here is an example of a simple batch script:

Line	Actual batch script	Description
1	<code>#!/bin/bash</code>	[optional] shell interpreter line
2	<code>#SBATCH -A &lt;project_id&gt;</code>	OLCF project to charge
3	<code>#SBATCH -J &lt;job_name&gt;</code>	Job name
4	<code>#SBATCH -o %x-%j.out</code>	stdout file name ( <code>%x</code> is job name, <code>%j</code> is job id)
5	<code>#SBATCH -t 00:05:00</code>	Walltime requested (HH:MM:SS)
6	<code>#SBATCH -p &lt;partition&gt;</code>	Batch queue
7	<code>#SBATCH -N 2</code>	Number of computed nodes requested
8		Blank line
9	<code>srun -N2 -n4 --ntasks-per-node=2 ./a.out</code>	srun command to launch parallel job

# Interactive Jobs

An **interactive job** allows multiple job steps (i.e., multiple `srun` commands) to be launched (interactively) on compute nodes allocated with the `salloc` command.

Here is an example of such a job:

```
$ salloc -A <project_id> -J <job_name> -t 00:05:00 -p <partition> -N 2
salloc: Granted job allocation 4258
salloc: Waiting for resource configuration
salloc: Nodes frontier[10469-10470] are ready for job

$ srun -n 4 --ntasks-per-node=2 ./a.out
<output printed to terminal>

$ srun -n 2 --ntasks-per-node=1 ./a.out
<output printed to terminal>
```

Here, `salloc` is used to request 2 compute nodes for 5 minutes, and once the compute nodes become available, job steps can be launched on them using `srun`.

# Single Command (non-interactive)

A **single command** job allows job allocation and a job step to be run in the same `srun` command.

Here is an example of such a job:

```
$ srun -A <project_id> -t 00:05:00 -p <partition> -N 2 -n 4 --ntasks-per-node=2 ./a.out  
<output printed to terminal>
```

Here, `srun` is used to request 2 compute nodes for 5 minutes and launch a job step with 4 tasks (2 on each node) – all in the same single command.

# Common Slurm Options and Commands

## Common Slurm Submission Options:

<code>-A &lt;project_id&gt;</code>	Project ID to charge	see <a href="#">man sbatch</a> for more info
<code>-J &lt;job_name&gt;</code>	Name of job	
<code>-p &lt;partition&gt;</code>	Partition / batch queue	
<code>-t &lt;time&gt;</code>	Wall clock time <HH:MM:SS> (or you can just give minutes)	
<code>-N &lt;number_of_nodes&gt;</code>	Number of compute nodes	
<code>-o &lt;file_name&gt;</code>	Standard output file name	
<code>-e &lt;file_name&gt;</code>	Standard error file name	
<code>--threads-per-core=&lt;threads&gt;</code>	Number of active hardware threads per core [1 (default) or 2]	

## Common Slurm Submission Options:

<code>sinfo</code>	Used to view partition and node information.	see individual <a href="#">man</a> pages for more info
<code>squeue</code>	Used to view job and job step information for jobs in the scheduling queue.	
<code>sacct</code>	Used to view accounting data for jobs and job steps in the accounting log (in queue or recently completed).	
<code>scancel</code>	Used to signal or cancel jobs or job steps.	
<code>scontrol</code>	Used to view or modify job configuration.	

# Examples

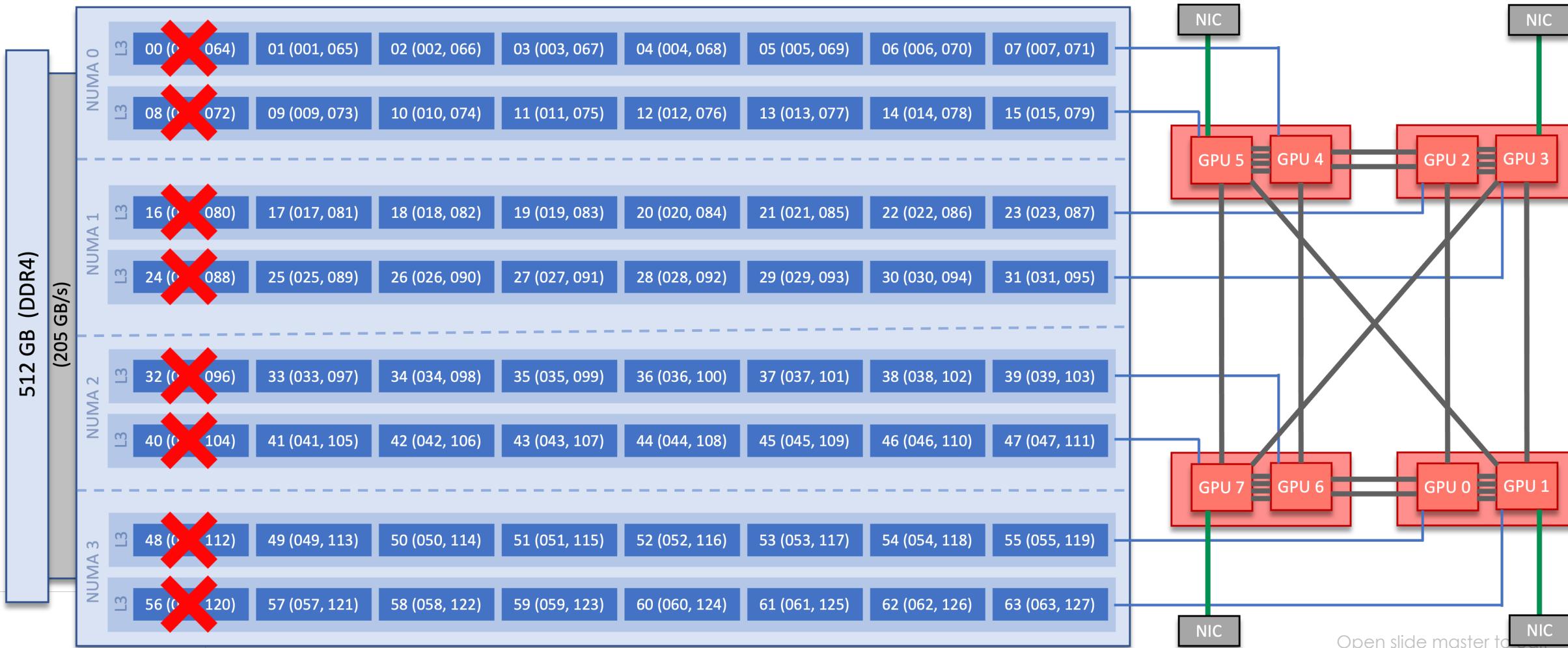


# Slurm GPU Bindings

(Recall: MPI ranks should target the GPU associated with their L3 cache region)

Always check process/thread/GPU bindings: [https://code.ornl.gov/olcf/hello\\_jobstep](https://code.ornl.gov/olcf/hello_jobstep)

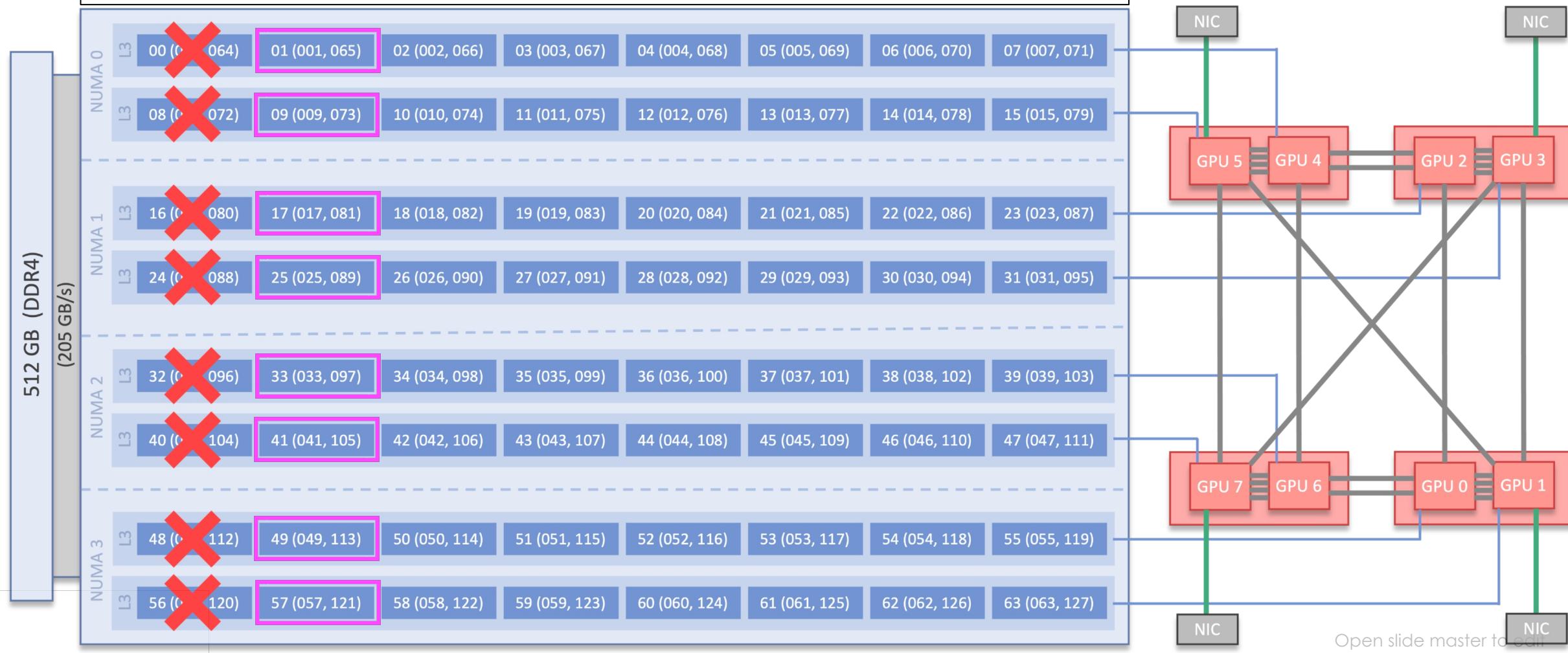
- **GPU\_ID** – node-level (or global) GPU ID  
read from `ROCR_VISIBLE_DEVICES`.
- **RT\_GPU\_ID** – HIP runtime GPU ID  
(as reported from `hipGetDevice`).
- **Bus\_ID** – physical bus ID associated with the GPUs.  
Comparing bus IDs shows that different GPUs are being used.



```
$ OMP_NUM_THREADS=1 srun -N1 -n8 -c7 --gpus-per-node=8 ./hello_jobstep | sort
```

MPI	OMP	HWT	Node	frontier08303	RT_GPU_ID	GPU_ID	Bus_ID
000	000	001		0,1,2,3,4,5,6,7	0,1,2,3,4,5,6,7	c1,c6,c9,ce,d1,d6,d9,de	
001	000	009		0,1,2,3,4,5,6,7	0,1,2,3,4,5,6,7	c1,c6,c9,ce,d1,d6,d9,de	
002	000	017		0,1,2,3,4,5,6,7	0,1,2,3,4,5,6,7	c1,c6,c9,ce,d1,d6,d9,de	
003	000	025		0,1,2,3,4,5,6,7	0,1,2,3,4,5,6,7	c1,c6,c9,ce,d1,d6,d9,de	
004	000	033		0,1,2,3,4,5,6,7	0,1,2,3,4,5,6,7	c1,c6,c9,ce,d1,d6,d9,de	
005	000	041		0,1,2,3,4,5,6,7	0,1,2,3,4,5,6,7	c1,c6,c9,ce,d1,d6,d9,de	
006	000	049		0,1,2,3,4,5,6,7	0,1,2,3,4,5,6,7	c1,c6,c9,ce,d1,d6,d9,de	
007	000	057		0,1,2,3,4,5,6,7	0,1,2,3,4,5,6,7	c1,c6,c9,ce,d1,d6,d9,de	

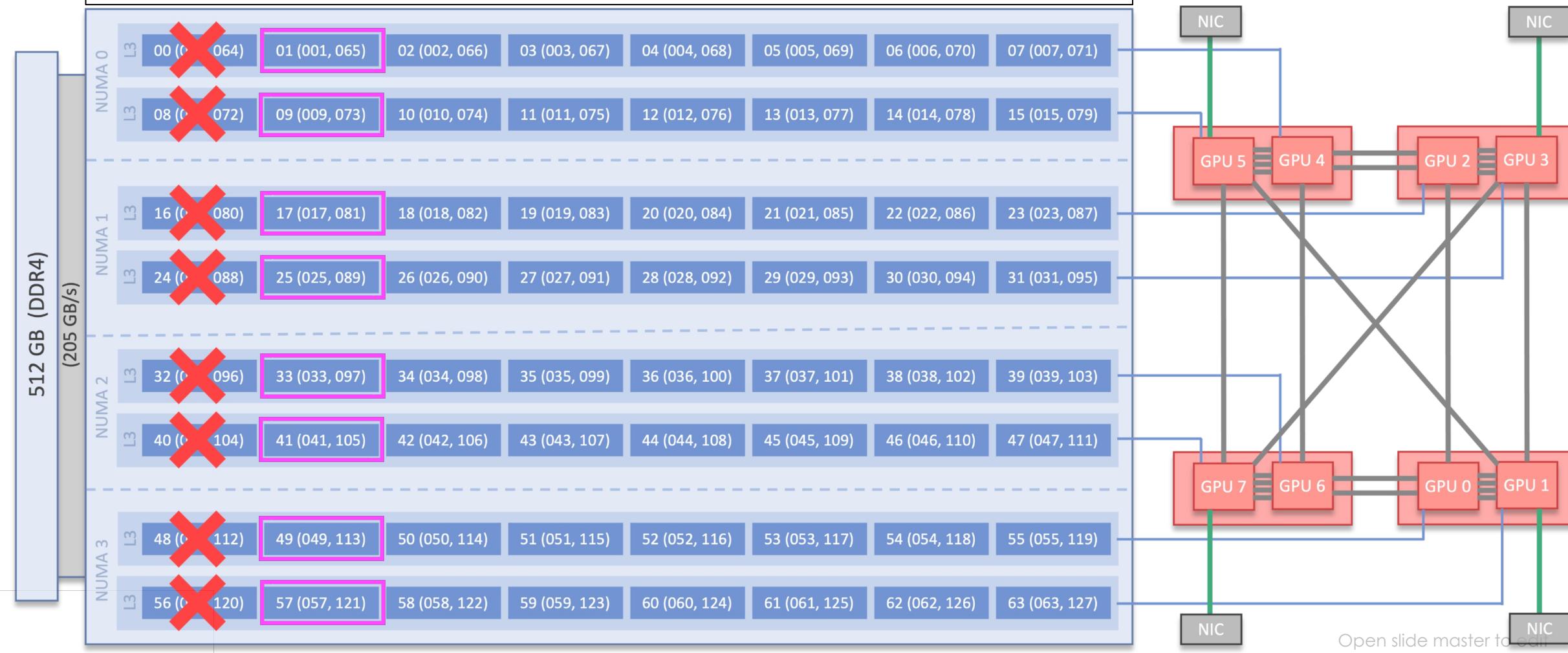
8 tasks per node, each with 7 CPU cores and 1 GPU



```
$ OMP_NUM_THREADS=1 srun -N1 -n8 -c7 --gpus-per-node=8 --gpu-bind=closest ./hello_jobstep | sort
```

MPI	OMP	HWT	Node	frontier08303	RT_GPU_ID	GPU_ID	Bus_ID
000	000	001		0	0	4	d1
001	000	009		0	0	5	d6
002	000	017		0	0	2	c9
003	000	025		0	0	3	ce
004	000	033		0	0	6	d9
005	000	041		0	0	7	de
006	000	049		0	0	0	c1
007	000	057		0	0	1	c6

8 tasks per node, each with 7 CPU cores and 1 GPU

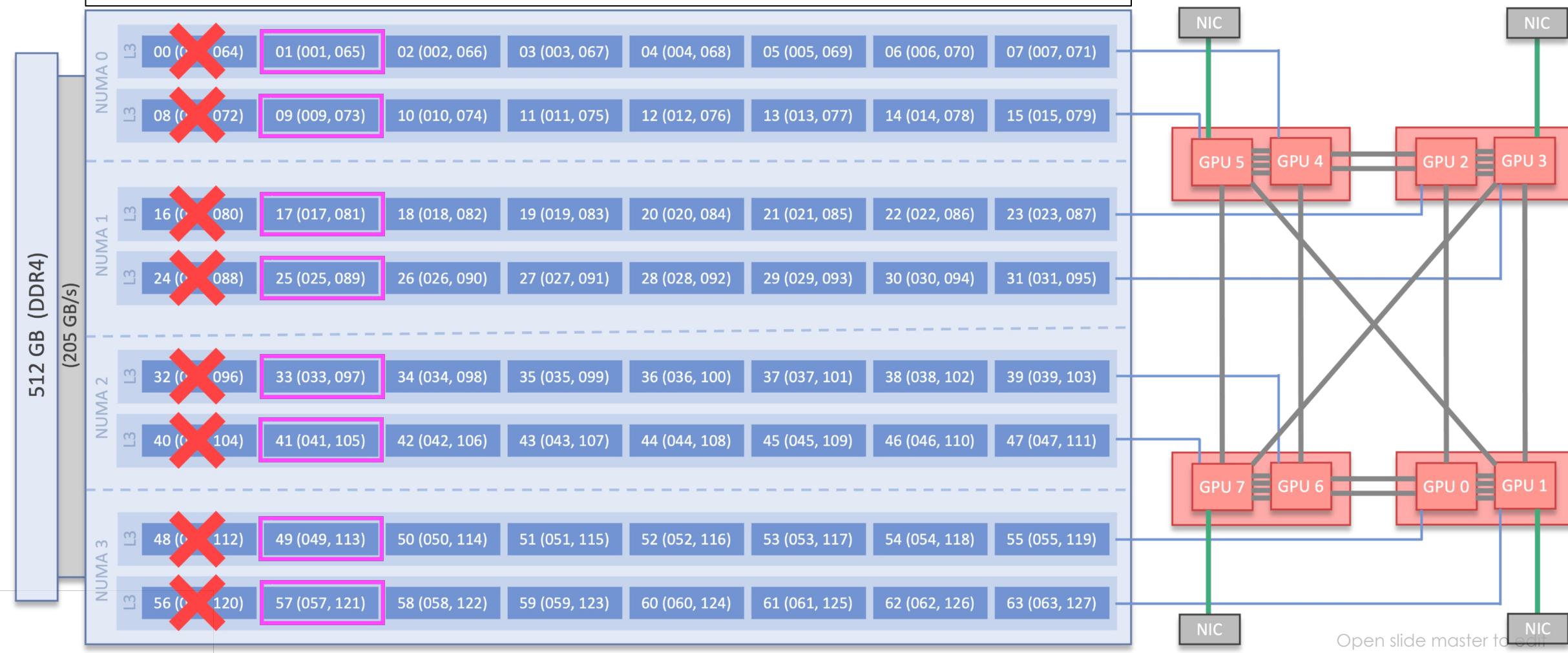


\$ OMP\_NUM\_THREADS=1 srun -N1 -n8 -c7 --gpus-per-task=1 ./hello jobstep | sort

```
MPI 000 - OMP 000 - HWT 001 - Node frontier08303 - RT_GPU_ID 0 - GPU_ID 0 - Bus_ID c1
MPI 001 - OMP 000 - HWT 009 - Node frontier08303 - RT_GPU_ID 0 - GPU_ID 1 - Bus_ID c6
MPI 002 - OMP 000 - HWT 017 - Node frontier08303 - RT_GPU_ID 0 - GPU_ID 2 - Bus_ID c9
MPI 003 - OMP 000 - HWT 025 - Node frontier08303 - RT_GPU_ID 0 - GPU_ID 3 - Bus_ID ce
MPI 004 - OMP 000 - HWT 033 - Node frontier08303 - RT_GPU_ID 0 - GPU_ID 4 - Bus_ID d1
MPI 005 - OMP 000 - HWT 041 - Node frontier08303 - RT_GPU_ID 0 - GPU_ID 5 - Bus_ID d6
MPI 006 - OMP 000 - HWT 049 - Node frontier08303 - RT_GPU_ID 0 - GPU_ID 6 - Bus_ID d9
MPI 007 - OMP 000 - HWT 057 - Node frontier08303 - RT_GPU_ID 0 - GPU_ID 7 - Bus_ID de
```

Incorrect mapping!

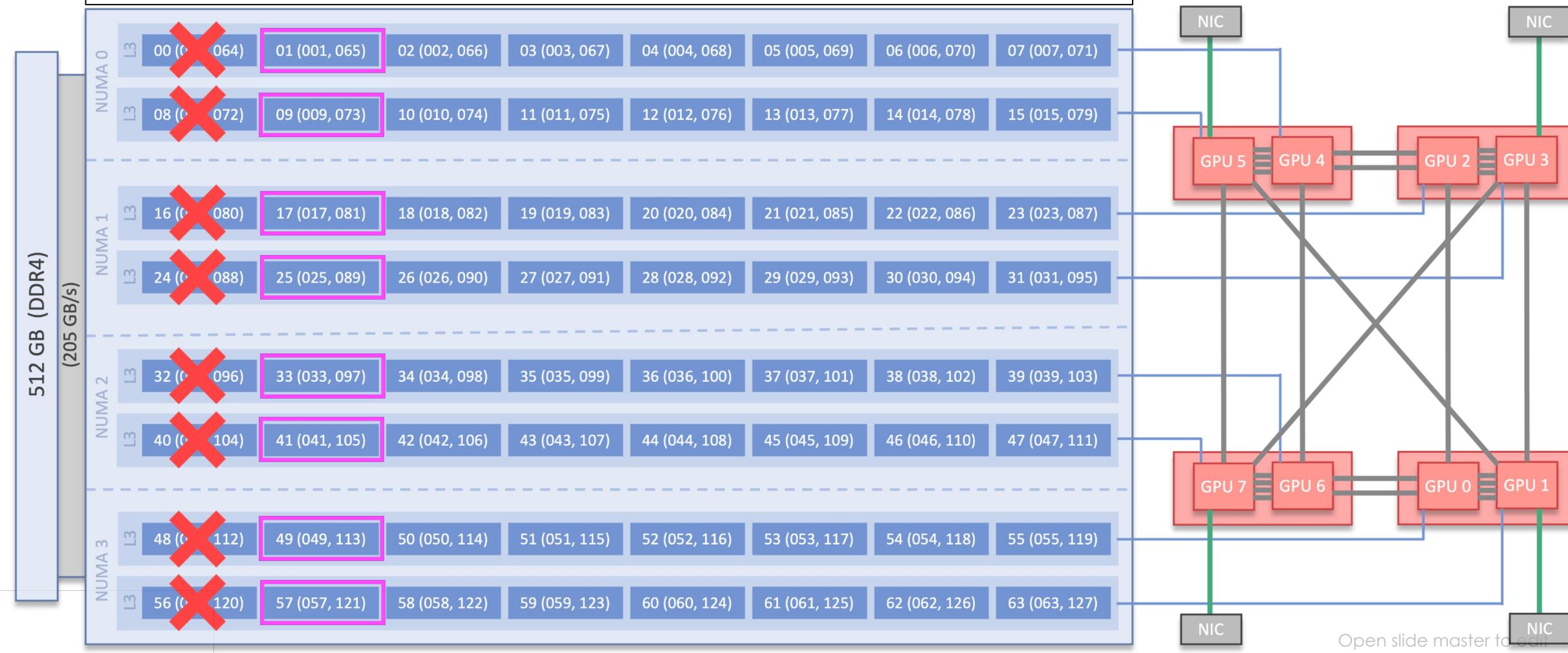
8 tasks per node, each with 7 CPU cores and 1 GPU



```
$ OMP_NUM_THREADS=1 srun -N1 -n8 -c7 --gpus-per-task=1 --gpu-bind=closest ./hello_jobstep | sort
```

MPI	OMP	HWT	Node	frontier08303	RT_GPU_ID	GPU_ID	Bus_ID
000	000	001		0	0	4 ✓	d1
001	000	009		0	0	5	d6
002	000	017		0	0	2	c9
003	000	025		0	0	3	ce
004	000	033		0	0	6	d9
005	000	041		0	0	7	de
006	000	049		0	0	0	c1
007	000	057		0	1	1	c6

8 tasks per node, each with 7 CPU cores and 1 GPU

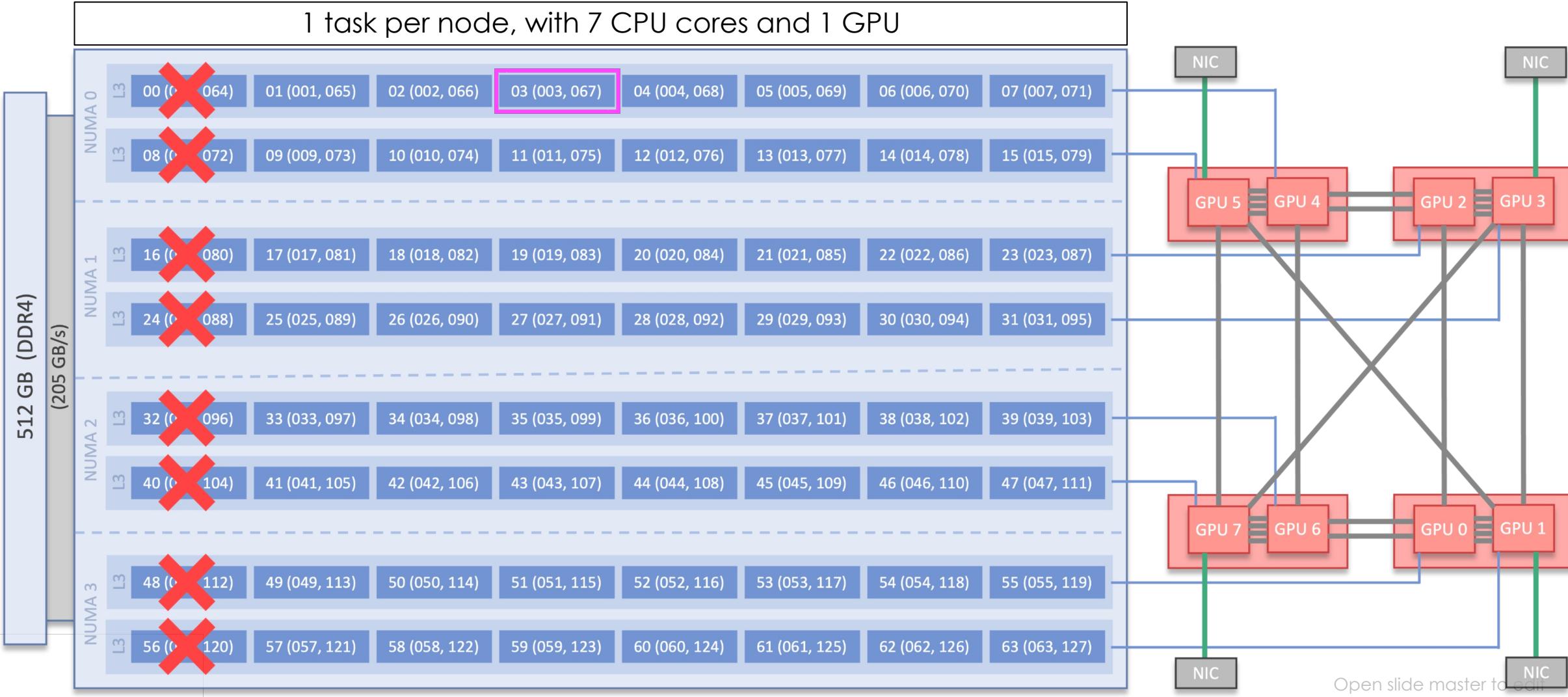


```
$ OMP_NUM_THREADS=1 srun -N1 -n1 -c7 --gpus-per-node=8 --gpu-bind=closest ./hello jobstep | sort
MPI 000 - OMP 000 - [HWT 003] - Node frontier08303 - RT_GPU_ID 0,1,2,3,4,5,6,7 - [GPU_ID 0,1,2,3,4,5,6,7] - Bus_ID c1,c6,c9,ce,d1,d6,d9,de
```

--gpu-bind=[verbose,]<type>

Bind tasks to specific GPUs. By default every spawned task can access every GPU allocated to the step. If "verbose," is specified before <type>, then print out GPU binding debug information to the stderr of the tasks. GPU binding is ignored if there is only one task.

**Incorrect mapping!**



```
$ OMP_NUM_THREADS=1 srun -N1 -n1 -c7 --gpus-per-task=1 --gpu-bind=closest ./hello_jobstep | sort
```

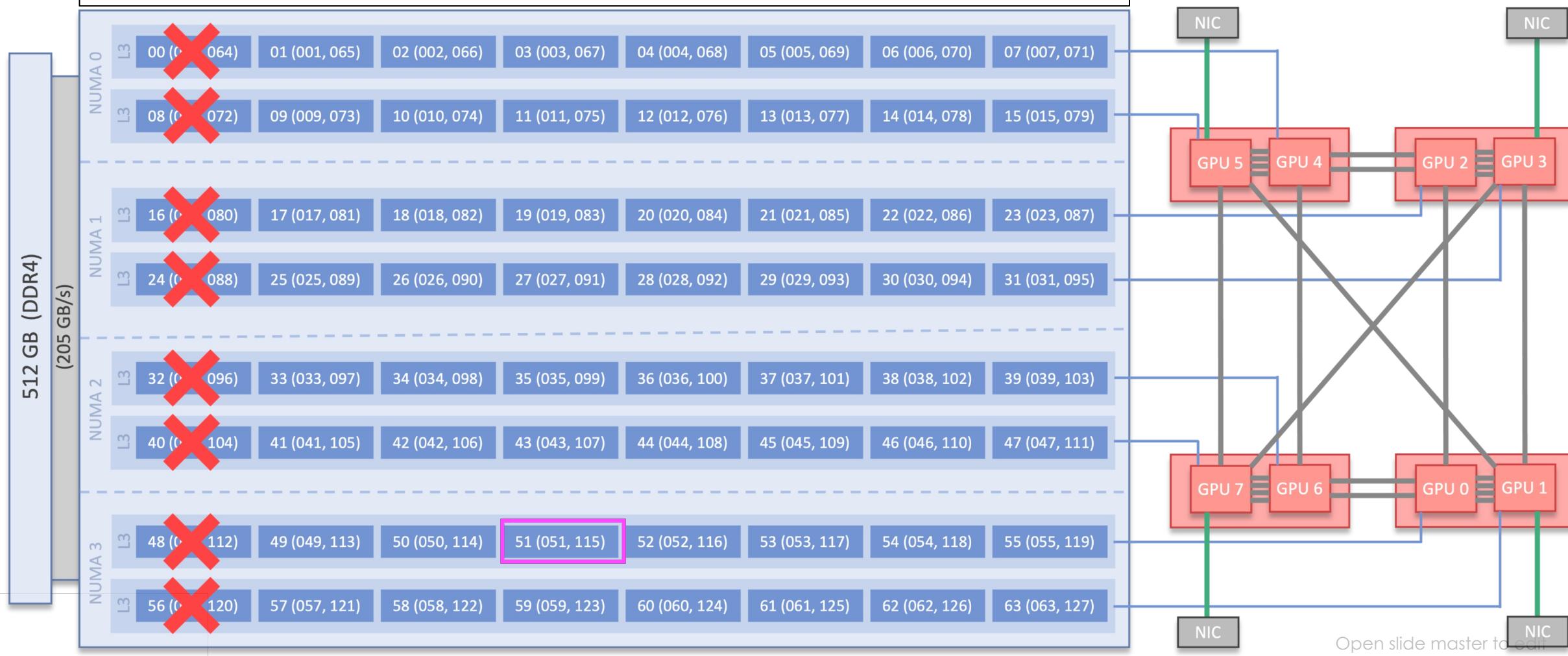
MPI 000 - OMP 000 - HWT 051 - Node frontier08303 - RT\_GPU\_ID 0 - GPU\_ID 0 - Bus\_ID c1



--gpu-bind=[verbose, ]<type>

Bind tasks to specific GPUs. By default every spawned task can access every GPU allocated to the step. If "verbose," is specified before <type>, then print out GPU binding debug information to the stderr of the tasks. GPU binding is ignored if there is only one task.

### 1 task per node, with 7 CPU cores and 1 GPU



```
$ OMP_NUM_THREADS=1 srun -N2 -n2 -c7 --gpus-per-node=8 --gpu-bind=closest ./hello jobstep | sort
MPI 000 - OMP 000 - HWT 002 - Node frontier10227 - RT_GPU_ID 0,1,2,3,4,5,6,7 - GPU_ID 0,1,2,3,4,5,6,7 - Bus_ID c1,c6,c9,ce,d1,d6,d9,de
MPI 001 - OMP 000 - HWT 002 - Node frontier10228 - RT_GPU_ID 0,1,2,3,4,5,6,7 - GPU_ID 0,1,2,3,4,5,6,7 - Bus_ID c1,c6,c9,ce,d1,d6,d9,de
```

--gpu-bind=[verbose,]<type>

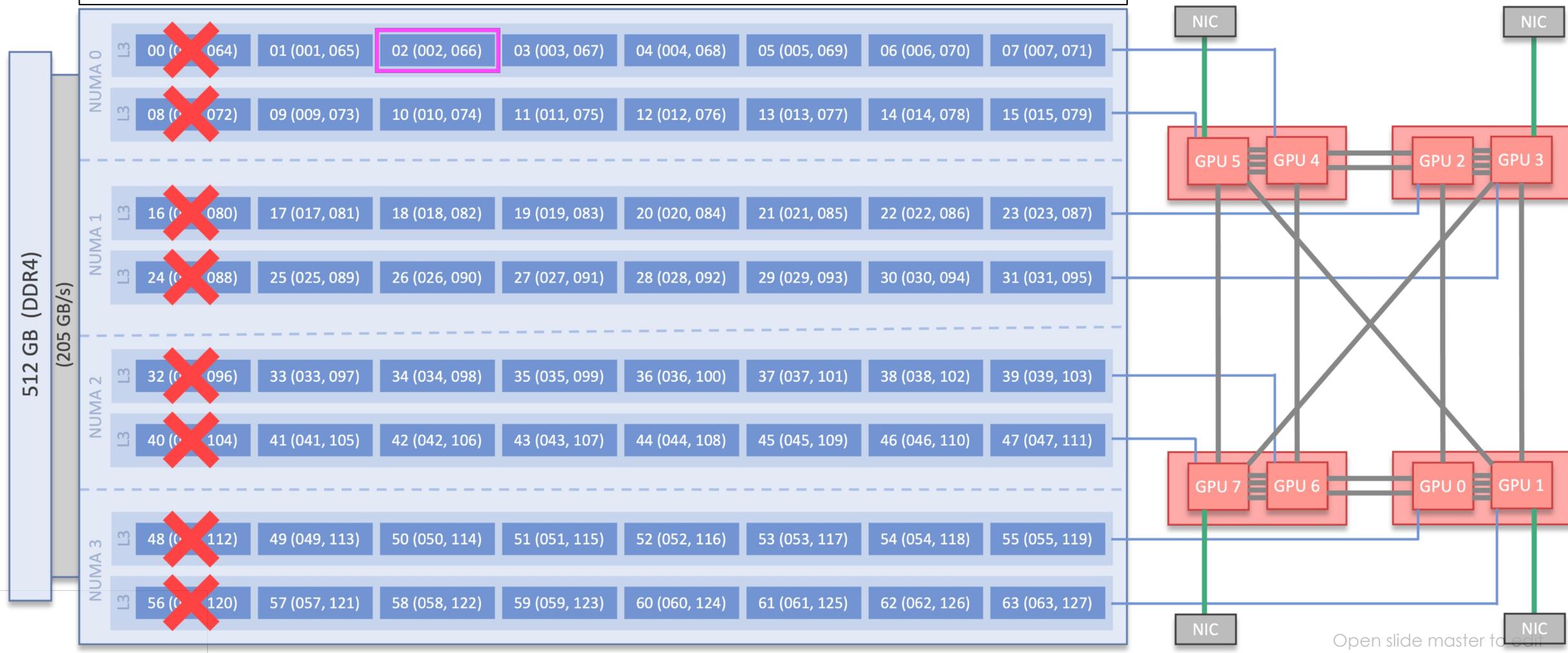


Incorrect mapping!

Bind tasks to specific GPUs. By default every spawned task can access every GPU allocated to the step. If "verbose," is specified before <type>, then print out GPU binding debug information to the stderr of the tasks. GPU binding is ignored if there is only one task.



1 task per node, with 7 CPU cores and 1 GPU (on 2 nodes)



```
$ OMP_NUM_THREADS=1 srun -N2 -n2 -c7 --gpus-per-task=1 --gpu-bind=closest ./hello_jobstep | sort
```

MPI 000 - OMP 000 - HWT 051 - Node frontier10227 - RT\_GPU\_ID 0 - GPU\_ID 0 - Bus\_ID c1

MPI 001 - OMP 000 - HWT 049 - Node frontier10228 - RT\_GPU\_ID 0 - GPU\_ID 0 - Bus\_ID c1

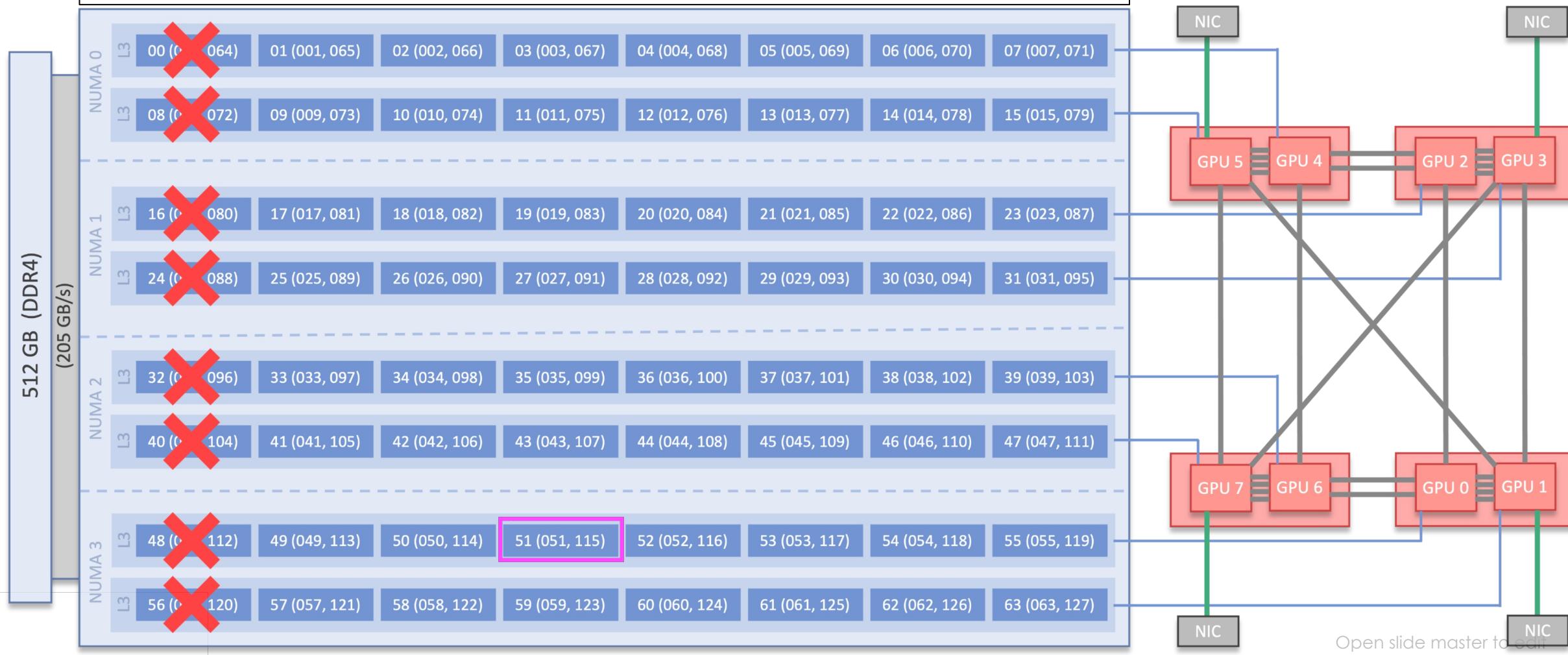
--gpu-bind=[verbose,]<type>



Bind tasks to specific GPUs. By default every spawned task can access every GPU allocated to the step. If "verbose," is specified before <type>, then print out GPU binding debug information to the stderr of the tasks. GPU binding is ignored if there is only one task.



### 1 task per node, with 7 CPU cores and 1 GPU (on 2 nodes)



# Distribution of tasks to node, sockets, and CPUs

--distribution=<value>[:<value>][:<value>][,{Pack | NoPack}]

Specifies the distribution of MPI ranks across compute nodes, sockets (L3 regions on Crusher), and cores, respectively. The default values are block:cyclic:cyclic

\* Can be used to refer to the default value

```
-m, --distribution=(*|block|cyclic|arbitrary|plane=<size>)[:{*|block|cyclic|fcyclic}[:{*|block|cyclic|fcyclic}]][,{Pack|NoPack}]
```

This option controls the distribution of tasks to the nodes on which resources have been allocated, and the distribution of those resources to tasks for binding (task affinity).

The **first distribution method** (before the first ":") controls the distribution of tasks to nodes.

[default method for distributing tasks to nodes (block)]

The **second distribution method** (after the first ":") controls the distribution of allocated CPUs across sockets for binding to tasks.

[default method for distributing CPUs across sockets (cyclic)]

**first distribution method** (before first ":") controls **the distribution of tasks to nodes**.

[default method for distributing tasks to nodes (block)]

```
$ OMP_NUM_THREADS=1 srun -N2 -n16 -c7 --gpus-per-task=1 --gpu-bind=closest ./hello_jobstep | sort
MPI 000 - OMP 000 - HWT 001 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 4 - Bus_ID d1
MPI 001 - OMP 000 - HWT 009 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 5 - Bus_ID d6
MPI 002 - OMP 000 - HWT 017 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 2 - Bus_ID c9
MPI 003 - OMP 000 - HWT 025 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 3 - Bus_ID ce
MPI 004 - OMP 000 - HWT 033 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 6 - Bus_ID d9
MPI 005 - OMP 000 - HWT 041 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 7 - Bus_ID de
MPI 006 - OMP 000 - HWT 049 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 0 - Bus_ID c1
MPI 007 - OMP 000 - HWT 057 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 1 - Bus_ID c6
MPI 008 - OMP 000 - HWT 001 - Node frontier10228 - RT_GPU_ID 0 - GPU_ID 4 - Bus_ID d1
MPI 009 - OMP 000 - HWT 009 - Node frontier10228 - RT_GPU_ID 0 - GPU_ID 5 - Bus_ID d6
MPI 010 - OMP 000 - HWT 017 - Node frontier10228 - RT_GPU_ID 0 - GPU_ID 2 - Bus_ID c9
MPI 011 - OMP 000 - HWT 025 - Node frontier10228 - RT_GPU_ID 0 - GPU_ID 3 - Bus_ID ce
MPI 012 - OMP 000 - HWT 033 - Node frontier10228 - RT_GPU_ID 0 - GPU_ID 6 - Bus_ID d9
MPI 013 - OMP 000 - HWT 041 - Node frontier10228 - RT_GPU_ID 0 - GPU_ID 7 - Bus_ID de
MPI 014 - OMP 000 - HWT 049 - Node frontier10228 - RT_GPU_ID 0 - GPU_ID 0 - Bus_ID c1
MPI 015 - OMP 000 - HWT 057 - Node frontier10228 - RT_GPU_ID 0 - GPU_ID 1 - Bus_ID c6
```

block distribution  
(default)

```
$ OMP_NUM_THREADS=1 srun -N2 -n16 -c7 --gpus-per-task=1 --gpu-bind=closest -m cyclic ./hello_jobstep | sort
MPI 000 - OMP 000 - HWT 001 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 4 - Bus_ID d1
MPI 001 - OMP 000 - HWT 001 - Node frontier10228 - RT_GPU_ID 0 - GPU_ID 4 - Bus_ID d1
MPI 002 - OMP 000 - HWT 009 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 5 - Bus_ID d6
MPI 003 - OMP 000 - HWT 009 - Node frontier10228 - RT_GPU_ID 0 - GPU_ID 5 - Bus_ID d6
MPI 004 - OMP 000 - HWT 017 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 2 - Bus_ID c9
MPI 005 - OMP 000 - HWT 017 - Node frontier10228 - RT_GPU_ID 0 - GPU_ID 2 - Bus_ID c9
MPI 006 - OMP 000 - HWT 025 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 3 - Bus_ID ce
MPI 007 - OMP 000 - HWT 025 - Node frontier10228 - RT_GPU_ID 0 - GPU_ID 3 - Bus_ID ce
MPI 008 - OMP 000 - HWT 033 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 6 - Bus_ID d9
MPI 009 - OMP 000 - HWT 033 - Node frontier10228 - RT_GPU_ID 0 - GPU_ID 6 - Bus_ID d9
MPI 010 - OMP 000 - HWT 041 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 7 - Bus_ID de
MPI 011 - OMP 000 - HWT 041 - Node frontier10228 - RT_GPU_ID 0 - GPU_ID 7 - Bus_ID de
MPI 012 - OMP 000 - HWT 049 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 0 - Bus_ID c1
MPI 013 - OMP 000 - HWT 049 - Node frontier10228 - RT_GPU_ID 0 - GPU_ID 0 - Bus_ID c1
MPI 014 - OMP 000 - HWT 057 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 1 - Bus_ID c6
MPI 015 - OMP 000 - HWT 057 - Node frontier10228 - RT_GPU_ID 0 - GPU_ID 1 - Bus_ID c6
```

cyclic distribution

**second distribution method** (after first ":") controls **the distribution of allocated CPUs across sockets for binding to tasks.**

[default method for distributing CPUs across sockets (cyclic)]

```
$ OMP_NUM_THREADS=1 srun -N1 -n16 -c3 --gpu-bind=closest --ntasks-per-gpu=2 ./hello_jobstep | sort
```

```
MPI 000 - OMP 000 - HWT 001 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 4 - Bus_ID d1  
MPI 001 - OMP 000 - HWT 009 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 5 - Bus_ID d6  
MPI 002 - OMP 000 - HWT 017 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 2 - Bus_ID c9  
MPI 003 - OMP 000 - HWT 025 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 3 - Bus_ID ce  
MPI 004 - OMP 000 - HWT 033 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 6 - Bus_ID d9  
MPI 005 - OMP 000 - HWT 041 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 7 - Bus_ID de  
MPI 006 - OMP 000 - HWT 049 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 0 - Bus_ID c1  
MPI 007 - OMP 000 - HWT 057 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 1 - Bus_ID c6  
MPI 008 - OMP 000 - HWT 004 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 4 - Bus_ID d1  
MPI 009 - OMP 000 - HWT 012 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 5 - Bus_ID d6  
MPI 010 - OMP 000 - HWT 020 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 2 - Bus_ID c9  
MPI 011 - OMP 000 - HWT 028 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 3 - Bus_ID ce  
MPI 012 - OMP 000 - HWT 036 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 6 - Bus_ID d9  
MPI 013 - OMP 000 - HWT 044 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 7 - Bus_ID de  
MPI 014 - OMP 000 - HWT 052 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 0 - Bus_ID c1  
MPI 015 - OMP 000 - HWT 060 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 1 - Bus_ID c6
```

cyclic distribution  
(default)

```
$ OMP_NUM_THREADS=1 srun -N1 -n16 -c3 --gpu-bind=closest --ntasks-per-gpu=2 -m block:block ./hello_jobstep | sort
```

```
MPI 000 - OMP 000 - HWT 001 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 4 - Bus_ID d1  
MPI 001 - OMP 000 - HWT 004 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 4 - Bus_ID d1  
MPI 002 - OMP 000 - HWT 007 - Node frontier10227 - RT_GPU_ID 0,1 - GPU_ID 4,5 - Bus_ID d1,d6  
MPI 003 - OMP 000 - HWT 011 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 5 - Bus_ID d6  
MPI 004 - OMP 000 - HWT 017 - Node frontier10227 - RT_GPU_ID 0,1 - GPU_ID 2,5 - Bus_ID c9,d6  
MPI 005 - OMP 000 - HWT 018 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 2 - Bus_ID c9  
MPI 006 - OMP 000 - HWT 021 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 2 - Bus_ID c9  
MPI 007 - OMP 000 - HWT 025 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 3 - Bus_ID ce  
MPI 008 - OMP 000 - HWT 028 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 3 - Bus_ID ce  
MPI 009 - OMP 000 - HWT 033 - Node frontier10227 - RT_GPU_ID 0,1 - GPU_ID 3,6 - Bus_ID ce,d9  
MPI 010 - OMP 000 - HWT 035 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 6 - Bus_ID d9  
MPI 011 - OMP 000 - HWT 041 - Node frontier10227 - RT_GPU_ID 0,1 - GPU_ID 6,7 - Bus_ID d9,de  
MPI 012 - OMP 000 - HWT 042 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 7 - Bus_ID de  
MPI 013 - OMP 000 - HWT 045 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 7 - Bus_ID de  
MPI 014 - OMP 000 - HWT 049 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 0 - Bus_ID c1  
MPI 015 - OMP 000 - HWT 052 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 0 - Bus_ID c1
```

block distribution

??

**second distribution method** (after first ":") controls **the distribution of allocated CPUs across sockets for binding to tasks.**

[default method for distributing CPUs across sockets (cyclic)]

```
$ OMP_NUM_THREADS=1 srun -N1 -n16 -c3 --gpu-bind=closest --ntasks-per-gpu=2 ./hello_jobstep | sort
MPI 000 - OMP 000 - HWT 002 - Node crusher131 - RT_GPU_ID 0 - GPU_ID 4 - Bus_ID d1
MPI 001 - OMP 000 - HWT 008 - Node crusher131 - RT_GPU_ID 0 - GPU_ID 5 - Bus_ID d6
MPI 002 - OMP 000 - HWT 017 - Node crusher131 - RT_GPU_ID 0 - GPU_ID 2 - Bus_ID c9
MPI 003 - OMP 000 - HWT 025 - Node crusher131 - RT_GPU_ID 0 - GPU_ID 3 - Bus_ID ce
MPI 004 - OMP 000 - HWT 032 - Node crusher131 - RT_GPU_ID 0 - GPU_ID 6 - Bus_ID d9
MPI 005 - OMP 000 - HWT 040 - Node crusher131 - RT_GPU_ID 0 - GPU_ID 7 - Bus_ID de
MPI 006 - OMP 000 - HWT 049 - Node crusher131 - RT_GPU_ID 0 - GPU_ID 0 - Bus_ID c1
MPI 007 - OMP 000 - HWT 057 - Node crusher131 - RT_GPU_ID 0 - GPU_ID 1 - Bus_ID c6
MPI 008 - OMP 000 - HWT 003 - Node crusher131 - RT_GPU_ID 0 - GPU_ID 4 - Bus_ID d1
MPI 009 - OMP 000 - HWT 013 - Node crusher131 - RT_GPU_ID 0 - GPU_ID 5 - Bus_ID d6
MPI 010 - OMP 000 - HWT 019 - Node crusher131 - RT_GPU_ID 0 - GPU_ID 2 - Bus_ID c9
MPI 011 - OMP 000 - HWT 027 - Node crusher131 - RT_GPU_ID 0 - GPU_ID 3 - Bus_ID ce
MPI 012 - OMP 000 - HWT 035 - Node crusher131 - RT_GPU_ID 0 - GPU_ID 6 - Bus_ID d9
MPI 013 - OMP 000 - HWT 043 - Node crusher131 - RT_GPU_ID 0 - GPU_ID 7 - Bus_ID de
MPI 014 - OMP 000 - HWT 053 - Node crusher131 - RT_GPU_ID 0 - GPU_ID 0 - Bus_ID c1
MPI 015 - OMP 000 - HWT 061 - Node crusher131 - RT_GPU_ID 0 - GPU_ID 1 - Bus_ID c6
```

cyclic distribution  
(default)

```
$ OMP_NUM_THREADS=1 srun -N1 -n16 -c3 --gpu-bind=closest --ntasks-per-gpu=2 -m block:block ./hello_jobstep | sort
PE 4:
MPICH ERROR: Unable to use a NIC_POLICY of 'NUMA'. Rank 4 is not confined to a single NUMA node. There are 4 numa_nodes detected
(rc=0).
MPICH ERROR [Rank 0] [job id 264018.27] [Sun Feb 12 18:20:30 2023] [crusher109] - Abort(2664847) (rank 0 in comm 0): Fatal error in
PMPI_Init: Other MPI error, error stack:
MPIR_Init_thread(171).....
MPI_DInit(495).....
MPIIDI_OFI_mpi_init_hook(623).....
open_fabric(1455).....
MPIIDI_CRAY_ofi_nic_assign_policy(3206):
MPIIDI_CRAY_ofi_get_nic_index(1701)....: OFI invalid value for environment variable
```

??

block distribution

On Crusher

# A brief aside...

Another way to check for resources available to each process

Prepend stdout lines with process ID

```
$ srun -l <srun-options> /bin/bash -c 'echo $(hostname) $(grep Cpus_allowed_list /proc/self/status) GPUS: $ROCR_VISIBLE_DEVICES' | sort
```

```
$ srun -l -N1 -n8 -c7 --gpus-per-task=1 --gpu-bind=closest /bin/bash -c 'echo $(hostname) $(grep Cpus_allowed_list /proc/self/status) GPUS: $ROCR_VISIBLE_DEVICES' | sort
```

```
0: crusher109 Cpus_allowed_list: 1-7 GPUS: 4
1: crusher109 Cpus_allowed_list: 9-15 GPUS: 5
2: crusher109 Cpus_allowed_list: 17-23 GPUS: 2
3: crusher109 Cpus_allowed_list: 25-31 GPUS: 3
4: crusher109 Cpus_allowed_list: 33-39 GPUS: 6
5: crusher109 Cpus_allowed_list: 41-47 GPUS: 7
6: crusher109 Cpus_allowed_list: 49-55 GPUS: 0
7: crusher109 Cpus_allowed_list: 57-63 GPUS: 1
```

Tells you which CPU cores and GPUs are available to each process, but not where they actually ran.

**second distribution method** (after first ":") controls **the distribution of allocated CPUs across sockets for binding to tasks.**

[default method for distributing CPUs across sockets (cyclic)]

```
$ OMP_NUM_THREADS=1 srun -N1 -n16 -c3 --gpu-bind=closest --ntasks-per-gpu=2 -m block:block ./hello_jobstep | sort
MPI 000 - OMP 000 - HWT 001 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 4 - Bus_ID d1
MPI 001 - OMP 000 - HWT 004 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 4 - Bus_ID d1
MPI 002 - OMP 000 - HWT 007 - Node frontier10227 - RT_GPU_ID 0,1 - GPU_ID 4,5 - Bus_ID d1,d6
MPI 003 - OMP 000 - HWT 011 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 5 - Bus_ID d6
MPI 004 - OMP 000 - HWT 017 - Node frontier10227 - RT_GPU_ID 0,1 - GPU_ID 2,5 - Bus_ID c9,d6
MPI 005 - OMP 000 - HWT 018 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 2 - Bus_ID c9
MPI 006 - OMP 000 - HWT 021 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 2 - Bus_ID c9
MPI 007 - OMP 000 - HWT 025 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 3 - Bus_ID ce
MPI 008 - OMP 000 - HWT 028 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 3 - Bus_ID ce
MPI 009 - OMP 000 - HWT 033 - Node frontier10227 - RT_GPU_ID 0,1 - GPU_ID 3,6 - Bus_ID ce,d9
MPI 010 - OMP 000 - HWT 035 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 6 - Bus_ID d9
MPI 011 - OMP 000 - HWT 041 - Node frontier10227 - RT_GPU_ID 0,1 - GPU_ID 6,7 - Bus_ID d9,de
MPI 012 - OMP 000 - HWT 042 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 7 - Bus_ID de
MPI 013 - OMP 000 - HWT 045 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 7 - Bus_ID de
MPI 014 - OMP 000 - HWT 049 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 0 - Bus_ID c1
MPI 015 - OMP 000 - HWT 052 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 0 - Bus_ID c1

$ OMP_NUM_THREADS=1 srun -l -N1 -n16 -c3 --gpu-bind=closest --ntasks-per-gpu=2 -m block:block /bin/bash -c 'echo $(hostname) $(grep
Cpus_allowed_list /proc/self/status) GPUS: $(ROCR_VISIBLE_DEVICES)' | sort -g
→ 0: frontier10227 Cpus_allowed_list: 1-3 GPUS: 4
→ 1: frontier10227 Cpus_allowed_list: 4-6 GPUS: 4
→ 2: frontier10227 Cpus_allowed_list: 7,9-10 GPUS: 4,5
3: frontier10227 Cpus_allowed_list: 11-13 GPUS: 5
4: frontier10227 Cpus_allowed_list: 14-15,17 GPUS: 2,5
5: frontier10227 Cpus_allowed_list: 18-20 GPUS: 2
6: frontier10227 Cpus_allowed_list: 21-23 GPUS: 2
7: frontier10227 Cpus_allowed_list: 25-27 GPUS: 3
8: frontier10227 Cpus_allowed_list: 28-30 GPUS: 3
9: frontier10227 Cpus_allowed_list: 31,33-34 GPUS: 3,6
10: frontier10227 Cpus_allowed_list: 35-37 GPUS: 6
11: frontier10227 Cpus_allowed_list: 38-39,41 GPUS: 6,7
12: frontier10227 Cpus_allowed_list: 42-44 GPUS: 7
13: frontier10227 Cpus_allowed_list: 45-47 GPUS: 7
14: frontier10227 Cpus_allowed_list: 49-51 GPUS: 0
15: frontier10227 Cpus_allowed_list: 52-54 GPUS: 0
```

Due to core isolation:

2 tasks (each with 3 cores) can fit on a node, but there is still 1 more core left over that will be given to the next task.

This can be resolved by removing the core isolation and using 4 threads per task.

[default method for distributing CPUs across sockets (cyclic)]

```
$ OMP_NUM_THREADS=1 srun -N1 -n56 -c1 --gpu-bind=closest --ntasks-per-gpu=7 -m block:block ./hello_jobstep | sort
MPI 000 - OMP 000 - HWT 001 - Node frontier02816 - RT_GPU_ID 0 - GPU_ID 4 - Bus_ID d1
MPI 001 - OMP 000 - HWT 002 - Node frontier02816 - RT_GPU_ID 0 - GPU_ID 4 - Bus_ID d1
MPI 002 - OMP 000 - HWT 003 - Node frontier02816 - RT_GPU_ID 0 - GPU_ID 4 - Bus_ID d1
MPI 003 - OMP 000 - HWT 004 - Node frontier02816 - RT_GPU_ID 0 - GPU_ID 4 - Bus_ID d1
MPI 004 - OMP 000 - HWT 005 - Node frontier02816 - RT_GPU_ID 0 - GPU_ID 4 - Bus_ID d1
MPI 005 - OMP 000 - HWT 006 - Node frontier02816 - RT_GPU_ID 0 - GPU_ID 4 - Bus_ID d1
MPI 006 - OMP 000 - HWT 007 - Node frontier02816 - RT_GPU_ID 0 - GPU_ID 4 - Bus_ID d1
MPI 007 - OMP 000 - HWT 009 - Node frontier02816 - RT_GPU_ID 0 - GPU_ID 5 - Bus_ID d6
MPI 008 - OMP 000 - HWT 010 - Node frontier02816 - RT_GPU_ID 0 - GPU_ID 5 - Bus_ID d6
MPI 009 - OMP 000 - HWT 011 - Node frontier02816 - RT_GPU_ID 0 - GPU_ID 5 - Bus_ID d6
MPI 010 - OMP 000 - HWT 012 - Node frontier02816 - RT_GPU_ID 0 - GPU_ID 5 - Bus_ID d6
MPI 011 - OMP 000 - HWT 013 - Node frontier02816 - RT_GPU_ID 0 - GPU_ID 5 - Bus_ID d6
MPI 012 - OMP 000 - HWT 014 - Node frontier02816 - RT_GPU_ID 0 - GPU_ID 5 - Bus_ID d6
MPI 013 - OMP 000 - HWT 015 - Node frontier02816 - RT_GPU_ID 0 - GPU_ID 5 - Bus_ID d6
MPI 014 - OMP 000 - HWT 017 - Node frontier02816 - RT_GPU_ID 0 - GPU_ID 2 - Bus_ID c9
MPI 015 - OMP 000 - HWT 018 - Node frontier02816 - RT_GPU_ID 0 - GPU_ID 2 - Bus_ID c9
MPI 016 - OMP 000 - HWT 019 - Node frontier02816 - RT_GPU_ID 0 - GPU_ID 2 - Bus_ID c9
MPI 017 - OMP 000 - HWT 020 - Node frontier02816 - RT_GPU_ID 0 - GPU_ID 2 - Bus_ID c9
MPI 018 - OMP 000 - HWT 021 - Node frontier02816 - RT_GPU_ID 0 - GPU_ID 2 - Bus_ID c9
MPI 019 - OMP 000 - HWT 022 - Node frontier02816 - RT_GPU_ID 0 - GPU_ID 2 - Bus_ID c9
MPI 020 - OMP 000 - HWT 023 - Node frontier02816 - RT_GPU_ID 0 - GPU_ID 2 - Bus_ID c9
MPI 021 - OMP 000 - HWT 025 - Node frontier02816 - RT_GPU_ID 0 - GPU_ID 3 - Bus_ID ce
MPI 022 - OMP 000 - HWT 026 - Node frontier02816 - RT_GPU_ID 0 - GPU_ID 3 - Bus_ID ce
MPI 023 - OMP 000 - HWT 027 - Node frontier02816 - RT_GPU_ID 0 - GPU_ID 3 - Bus_ID ce
MPI 024 - OMP 000 - HWT 028 - Node frontier02816 - RT_GPU_ID 0 - GPU_ID 3 - Bus_ID ce
MPI 025 - OMP 000 - HWT 029 - Node frontier02816 - RT_GPU_ID 0 - GPU_ID 3 - Bus_ID ce
MPI 026 - OMP 000 - HWT 030 - Node frontier02816 - RT_GPU_ID 0 - GPU_ID 3 - Bus_ID ce
MPI 027 - OMP 000 - HWT 031 - Node frontier02816 - RT_GPU_ID 0 - GPU_ID 3 - Bus_ID ce
MPI 028 - OMP 000 - HWT 033 - Node frontier02816 - RT_GPU_ID 0 - GPU_ID 6 - Bus_ID d9
MPI 029 - OMP 000 - HWT 034 - Node frontier02816 - RT_GPU_ID 0 - GPU_ID 6 - Bus_ID d9
MPI 030 - OMP 000 - HWT 035 - Node frontier02816 - RT_GPU_ID 0 - GPU_ID 6 - Bus_ID d9
MPI 031 - OMP 000 - HWT 036 - Node frontier02816 - RT_GPU_ID 0 - GPU_ID 6 - Bus_ID d9
MPI 032 - OMP 000 - HWT 037 - Node frontier02816 - RT_GPU_ID 0 - GPU_ID 6 - Bus_ID d9
...
...
```

Note: There is no problem when running with 1 core per MPI rank (i.e., 7 ranks per GPU) because the task can't span multiple L3s

# Using all CPU cores on a Crusher/Frontier node

```
-S, --core-spec=<num>
```

Count of specialized cores per node reserved by the job for system operations and not used by the application.

```
$ salloc -Astdf016_frontier -N1 -t120 -S0 ← -S flag is used at job allocation time
salloc: Pending job allocation 1273083
salloc: job 1273083 queued and waiting for resources
salloc: job 1273083 has been allocated resources
salloc: Granted job allocation 1273083
salloc: Waiting for resource configuration
salloc: Nodes frontier09076 are ready for job
```

```
$ OMP_NUM_THREADS=8 srun -N1 -n8 -c8 --gpus-per-task=1 --gpu-bind=closest ./hello_jobstep | sort
MPI 000 - OMP 000 - HWT 000 - Node frontier09076 - RT_GPU_ID 0 - GPU_ID 4 - Bus_ID d1
MPI 000 - OMP 001 - HWT 001 - Node frontier09076 - RT_GPU_ID 0 - GPU_ID 4 - Bus_ID d1
MPI 000 - OMP 002 - HWT 002 - Node frontier09076 - RT_GPU_ID 0 - GPU_ID 4 - Bus_ID d1
MPI 000 - OMP 003 - HWT 003 - Node frontier09076 - RT_GPU_ID 0 - GPU_ID 4 - Bus_ID d1
MPI 000 - OMP 004 - HWT 004 - Node frontier09076 - RT_GPU_ID 0 - GPU_ID 4 - Bus_ID d1
MPI 000 - OMP 005 - HWT 005 - Node frontier09076 - RT_GPU_ID 0 - GPU_ID 4 - Bus_ID d1
MPI 000 - OMP 006 - HWT 006 - Node frontier09076 - RT_GPU_ID 0 - GPU_ID 4 - Bus_ID d1
MPI 000 - OMP 007 - HWT 007 - Node frontier09076 - RT_GPU_ID 0 - GPU_ID 4 - Bus_ID d1
MPI 001 - OMP 000 - HWT 008 - Node frontier09076 - RT_GPU_ID 0 - GPU_ID 5 - Bus_ID d6
MPI 001 - OMP 001 - HWT 009 - Node frontier09076 - RT_GPU_ID 0 - GPU_ID 5 - Bus_ID d6
MPI 001 - OMP 002 - HWT 010 - Node frontier09076 - RT_GPU_ID 0 - GPU_ID 5 - Bus_ID d6
MPI 001 - OMP 003 - HWT 011 - Node frontier09076 - RT_GPU_ID 0 - GPU_ID 5 - Bus_ID d6
MPI 001 - OMP 004 - HWT 012 - Node frontier09076 - RT_GPU_ID 0 - GPU_ID 5 - Bus_ID d6
MPI 001 - OMP 005 - HWT 013 - Node frontier09076 - RT_GPU_ID 0 - GPU_ID 5 - Bus_ID d6
MPI 001 - OMP 006 - HWT 014 - Node frontier09076 - RT_GPU_ID 0 - GPU_ID 5 - Bus_ID d6
MPI 001 - OMP 007 - HWT 015 - Node frontier09076 - RT_GPU_ID 0 - GPU_ID 5 - Bus_ID d6
...
...
```

But remember that core 0 will have lower performance than other cores on each node.

```

$ cat multiple_jobsteps.sh
#!/bin/bash

for idx in {1..8};
do
    OMP_NUM_THREADS=1 srun --exact -u -N1 -w <node_name> -n1 -c6 --gpus-per-task=1 --gpu-bind=closest ./hello_jobstep &
    date &
    sleep 1
done

wait

$ ./multiple_jobsteps.sh
Tue 14 Feb 2023 12:31:59 PM EST
MPI 000 - OMP 000 - HWT 051 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 0 - Bus_ID c1

Tue 14 Feb 2023 12:32:00 PM EST
MPI 000 - OMP 000 - HWT 059 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 1 - Bus_ID c6

Tue 14 Feb 2023 12:32:01 PM EST
MPI 000 - OMP 000 - HWT 019 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 2 - Bus_ID c9

Tue 14 Feb 2023 12:32:02 PM EST
MPI 000 - OMP 000 - HWT 027 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 3 - Bus_ID ce

Tue 14 Feb 2023 12:32:03 PM EST
MPI 000 - OMP 000 - HWT 003 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 4 - Bus_ID d1

Tue 14 Feb 2023 12:32:04 PM EST
MPI 000 - OMP 000 - HWT 011 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 5 - Bus_ID d6

Tue 14 Feb 2023 12:32:05 PM EST
MPI 000 - OMP 000 - HWT 035 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 6 - Bus_ID d9

Tue 14 Feb 2023 12:32:06 PM EST
MPI 000 - OMP 000 - HWT 042 - Node frontier10227 - RT_GPU_ID 0 - GPU_ID 7 - Bus_ID de

```

Find with \$SLURM\_NODELIST

# Multiple job steps on a single node

Bug? This will only work correctly for 1-6 cores when -S8 is set.

Not a problem for -S0

Added a `sleep(20)` to the bottom of `hello_jobstep` so it didn't return so quickly.

Questions?

Summit here

Frontier here



papatheodore@ornl.gov

OAK RIDGE  
National Laboratory

# Bonus Slides



# Slurm Tips – Flags and Completed Jobs

- u flag gives unbuffered output, which can be helpful when debugging.
- l flag prepends the task ID to lines of stdout.

To show completed jobs in a specific time period, specify a start (-S) and end (-E) time

- The default time window depends on other options (see [man sacct](#))

```
$ sacct --user=tppapathe -s 2022-05-05T00:01 -e 2022-05-05T23:59 | awk '$0 !~ /\./'
```

JobID	JobName	Partition	Account	AllocCPUS	State	ExitCode
107814	interacti+	batch	stf016	128	COMPLETED	0:0
107836	interacti+	batch	stf016	128	COMPLETED	0:0
107849	hello_job+	batch	stf016	128	COMPLETED	0:0
107858	hello_job+	batch	stf016	256	COMPLETED	0:0
107866	hello_job+	batch	stf016	256	COMPLETED	0:0
107867	hello_job+	batch	stf016	256	CANCELLED+	0:0
107868	hello_job+	batch	stf016	256	COMPLETED	0:0
107965	MatrixTra+	batch	stf016	128	COMPLETED	0:0
107966	rocprof	batch	stf016	128	FAILED	6:0
107969	rocprof	batch	stf016	128	FAILED	6:0

Then you can drill into more details about each job with the -j flag and customized output.

# Slurm Tips – Capture Job Information

```
$ cat submit.sh
#!/bin/bash

#SBATCH -A stf016
#SBATCH -J job-name
#SBATCH -o %x-%j.out
#SBATCH -t 5
#SBATCH -p batch
#SBATCH -N 1

scontrol show job ${SLURM_JOBID}

srun -n1 ./p2p/p2p --correct

sacct -j ${SLURM_JOBID} -o jobid%20,Start%20,elapsed%20
```

Also add...

```
module -t list
(lists currently-loaded environment modules)
```

```
ldd ./<your-executable>
#prints shared object dependencies)
```

```
$ cat job-name-108121.out
```

JobId=108121 JobName=job-name  
UserId=tpapathe(5987) GroupId=tpapathe(8654) MCS\_label=N/A  
Priority=200 Nice=0 Account=stf016 QOS=normal  
JobState=RUNNING Reason=None Dependency=(null)  
Requeue=0 Restarts=0 BatchFlag=1 Reboot=0 ExitCode=0:0  
RunTime=00:00:00 TimeLimit=00:05:00 TimeMin=N/A  
SubmitTime=2022-05-05T17:34:51 EligibleTime=2022-05-05T17:34:51  
AccrueTime=2022-05-05T17:34:51  
StartTime=2022-05-05T17:35:03 EndTime=2022-05-05T17:40:03 Deadline=N/A  
SuspendTime=None SecsPreSuspend=0 LastSchedEval=2022-05-05T17:35:03 Scheduler=Backfill  
Partition=batch AllocNode:Sid=login2:106343  
ReqNodeList=(null) ExcNodeList=(null)  
NodeList=crusher048  
BatchHost=crusher048  
NumNodes=1 NumCPUs=128 NumTasks=1 CPUs/Task=1 ReqB:S:C:T=0:0:0:1  
TRES=cpu=128,node=1,billing=128  
Socks/Node=\* NtasksPerN:B:S:C=0:0:0:1 CoreSpec=\*  
MinCPUsNode=1 MinMemoryNode=0 MinTmpDiskNode=0  
Features=(null) DelayBoot=00:00:00  
OverSubscribe=NO Contiguous=0 Licenses=(null) Network=(null)  
Command=/autofs/nccs-svml\_home1/tpapathe/capture-job-info/submit.sh  
WorkDir=/autofs/nccs-svml\_home1/tpapathe/capture-job-info  
StdErr=/autofs/nccs-svml\_home1/tpapathe/capture-job-info/job-name-108121.out  
StdIn=/dev/null  
StdOut=/autofs/nccs-svml\_home1/tpapathe/capture-job-info/job-name-108121.out

...<APPLICATION OUTPUT>...

JobID	Start	Elapsed
108121	2022-05-05T17:35:03	00:02:34
108121.batch	2022-05-05T17:35:03	00:02:34
108121.extern	2022-05-05T17:35:03	00:02:34
108121.0	2022-05-05T17:35:04	00:02:32

# Submitting Helpful User Support Tickets

Submit support tickets to  
[help@olcf.ornl.gov](mailto:help@olcf.ornl.gov)

Typical questions we ask users for:

- Job ID
- List of modules loaded when compiling and running
- Batch script used to launch the job
- Job stdout and stderr
- Output from `ldd` run on executable

Notice these are all things that are included from the previous slide



Other helpful things to include in your tickets

- Full errors that are generated
- Has the program run successfully before?
  - If so, did you change anything since then?
- Programming models used in your code