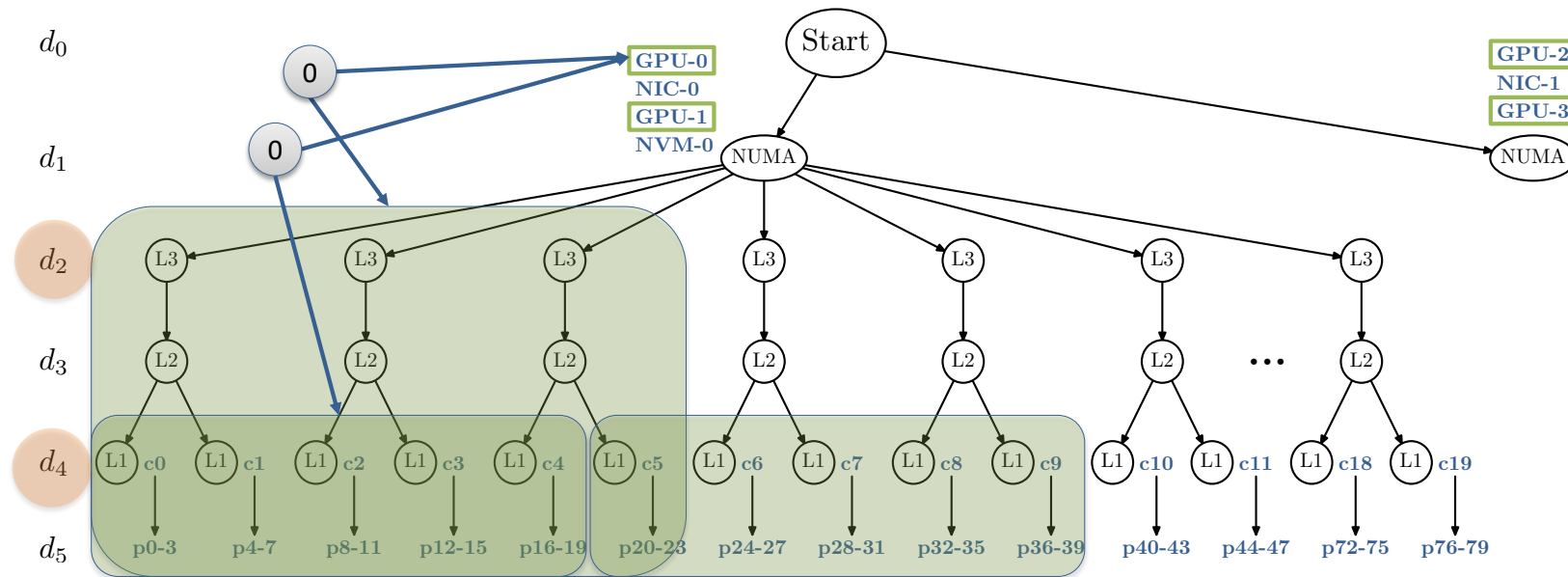


# Bringing in Robust, Memory-Driven Affinity to Slurm

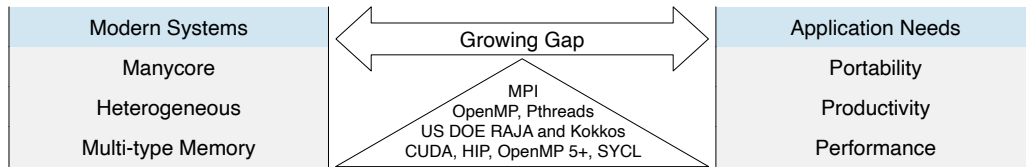
*mpibind*

Edgar A. León  
Livermore Computing

Slurm User Group  
12 Sep 2024, University of Oslo



# HPC users face complex computing architectures



AMD Instinct MI300A

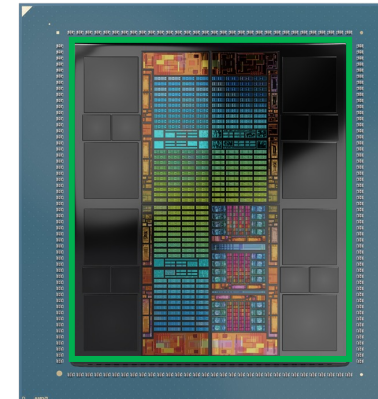


Image credit: AMD

AMD Instinct MI300X

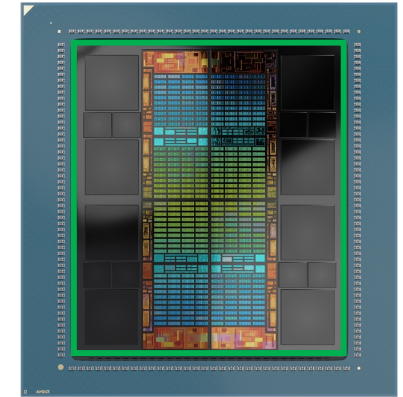
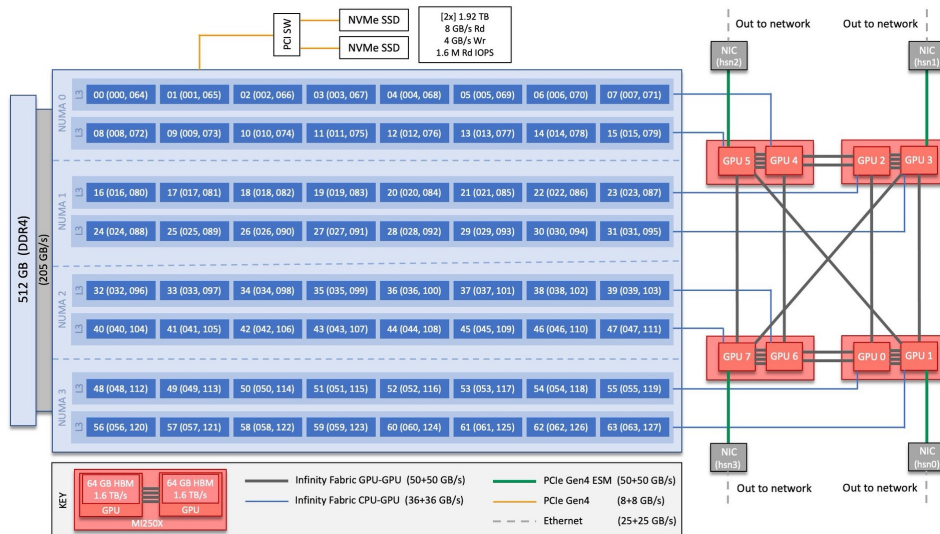


Image credit: AMD

## AMD 3rd Gen EPYC CPU + AMD Instinct MI250X GPUs



[https://docs.olcf.ornl.gov/systems/crusher\\_quick\\_start\\_guide.html](https://docs.olcf.ornl.gov/systems/crusher_quick_start_guide.html)

## NVIDIA GH200 Grace Hopper Superchip

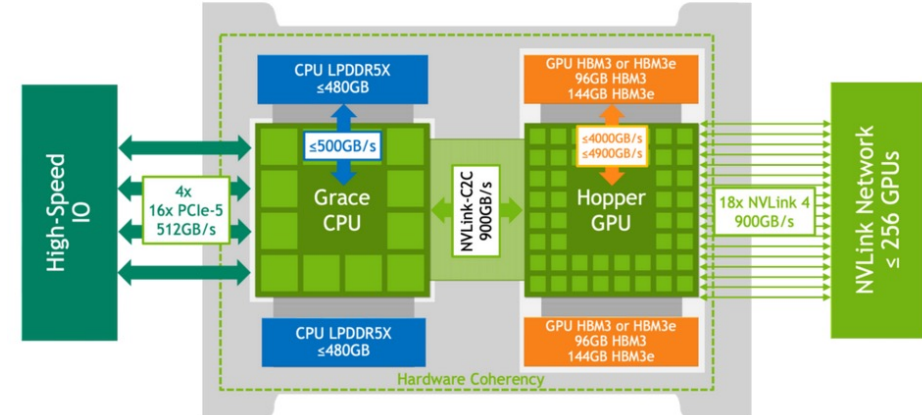
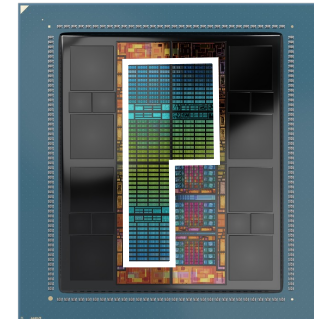
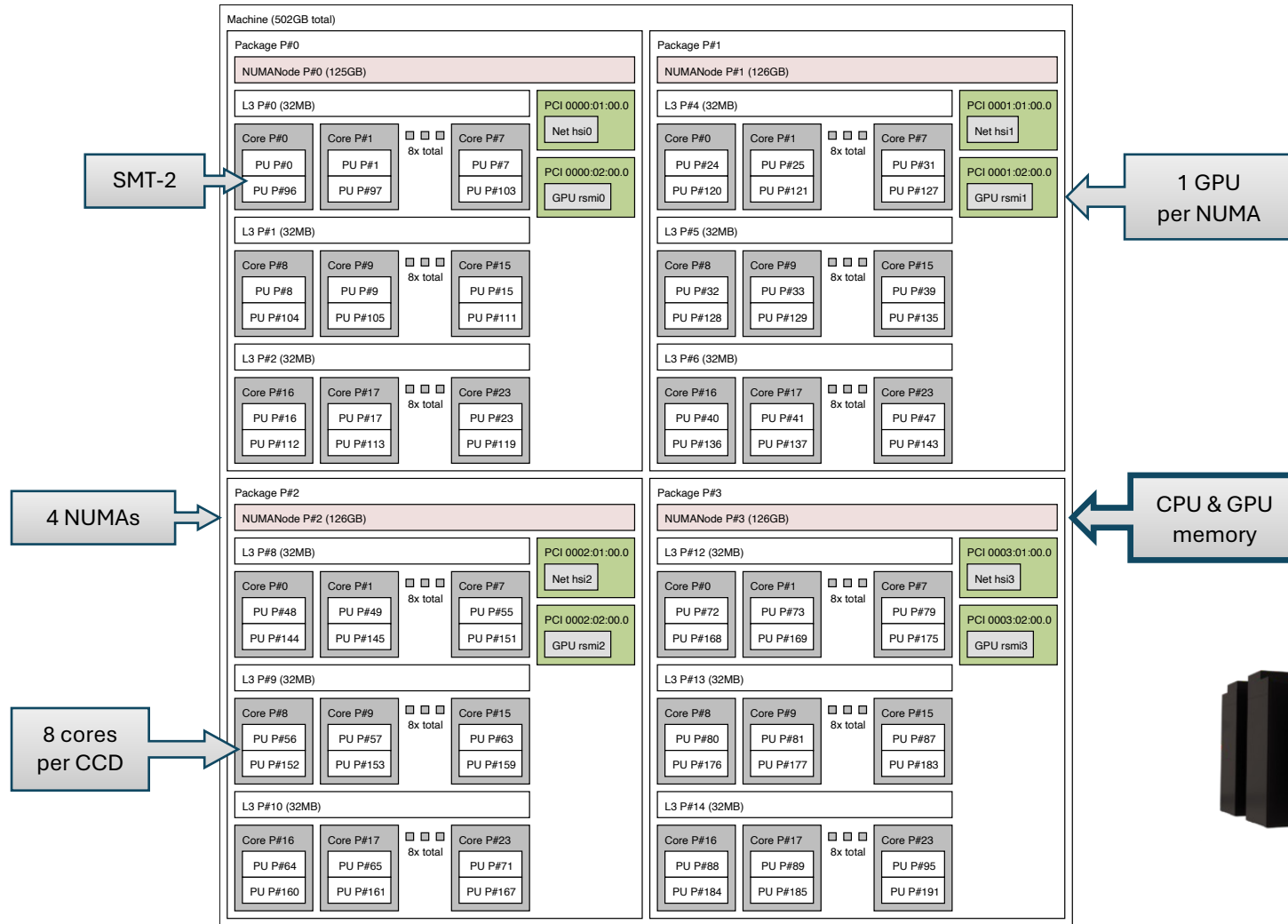


Image credit: NVIDIA

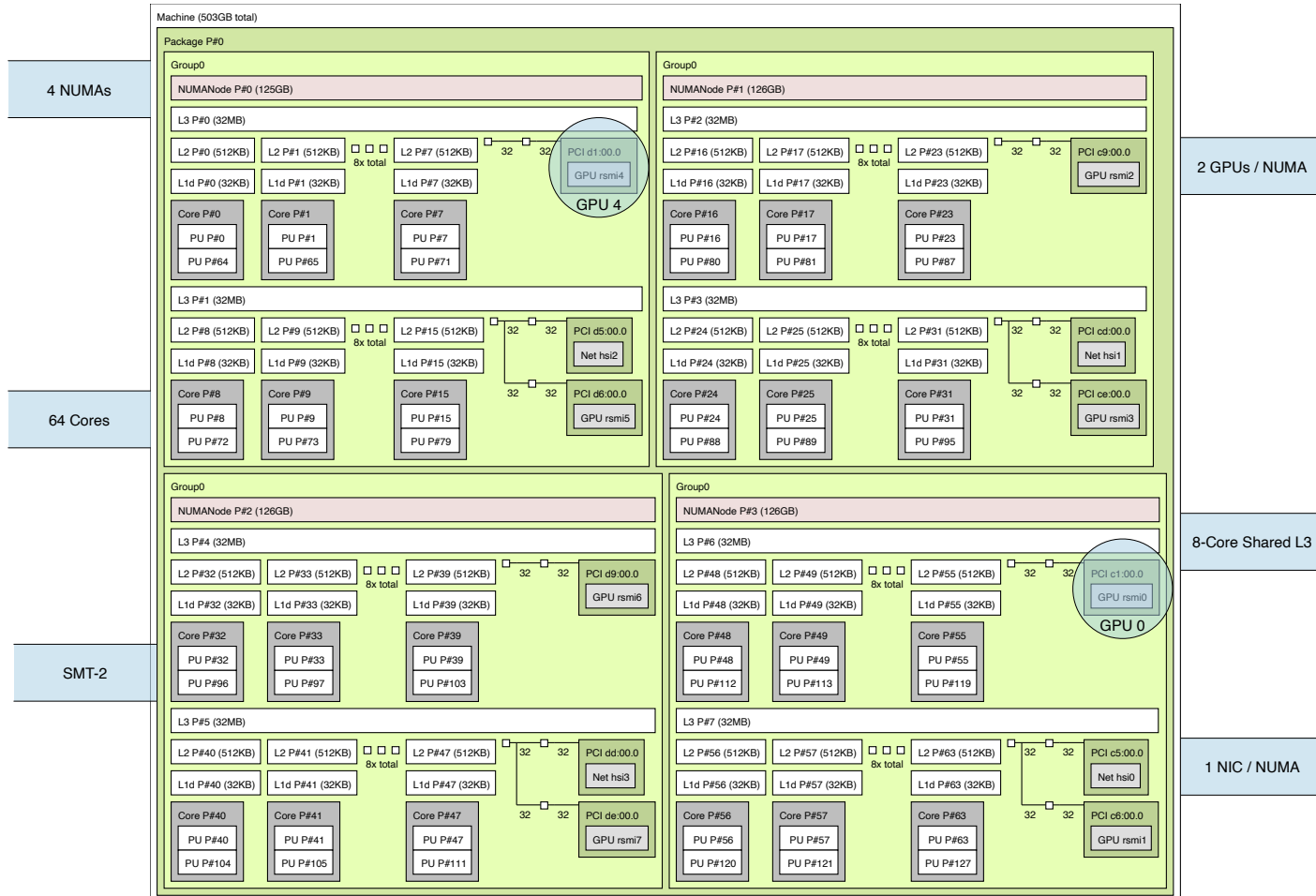
# El Capitan's building block is AMD's MI300A APU



AMD Instinct MI300A APU



# Mapping applications to the machine is hard



- What could go wrong?
  - Multiple threads running on a single core
  - Multiple GPU kernels sharing a GPU
  - Tasks launching kernels on non-local GPUs
  - Threads accessing remote memory
  - Etc.

# mpibind is a memory-first, user-friendly algorithm

- Design principles

- Driven by the memory system
- Based on locality
- Portable across architectures

*mpibind*

- Guiding principles

- Provide a simple interface
- Require minimal user input
- Minimize remote memory accesses
- Maximize cache per worker
- Leverage compute/memory locality
- Leverage architecture's features



# Provides a simple interface & Requires minimal user input

JSM

```
# 8 MPI tasks across all cores and all GPUs on a CORAL node  
jsrun -a 2 -c 10 -g 1 -r 4 -d packed -b packed:5 <prog>
```

mpibind

```
# 8 MPI tasks across all cores and all GPUs on a CORAL node  
srun -n8 <prog>  
srun -n8 -mpibind=on <prog>
```

<b>Required</b>	Number of tasks
<b>Optional</b>	Number of threads
<b>User options</b>	Greedy (true or false)
	GPU optimized (true or false)
	SMT (1 to num. HW threads per core)

*mpibind*

# Provides portability across systems

- Counterexample

<b>Intel MPI</b>	I_MPI_PIN_DOMAIN	core, sock, numa, node, cache
	I_MPI_PIN_ORDER	range, scatter, compact, spread, bunch
<b>MVAPICH2</b>	MV2_CPU_BINDING_LEVEL	core, socket, numanode
	MV2_CPU_BINDING_POLICY	bunch, scatter, hybrid
	MV2_HYBRID_BINDING_POLICY	bunch, scatter, linear, compact, spread
<b>OpenMPI</b>	--bind-to, --rank-by	slot, hwthread, core, cache, socket, numa, board
	--map-by	+ node, sequential, distance, ppr:n:unit:pe=n
<b>IBM Spectrum MPI</b>	-aff	bandwidth, latency, cycle:unit



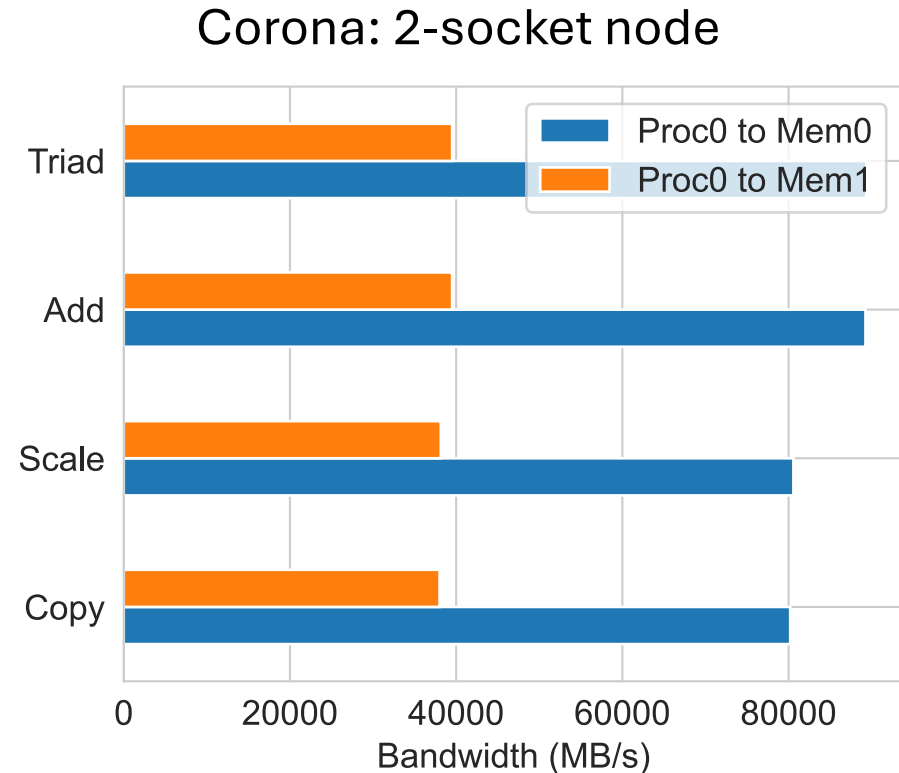
# Provides portability across systems

- Relies on abstract memory-compute tree
  - Portable Hardware Locality (hwloc)
- Affinity algorithm is separate from applying affinity
  - C interface
    - `{Task} → {CPUs, GPUs, Thread mapping}`
  - **Slurm plugin**
    - Use job info as input to mpibind
    - Get mpibind mapping
    - Bind tasks
    - Set environment variables



# Minimizes remote memory accesses

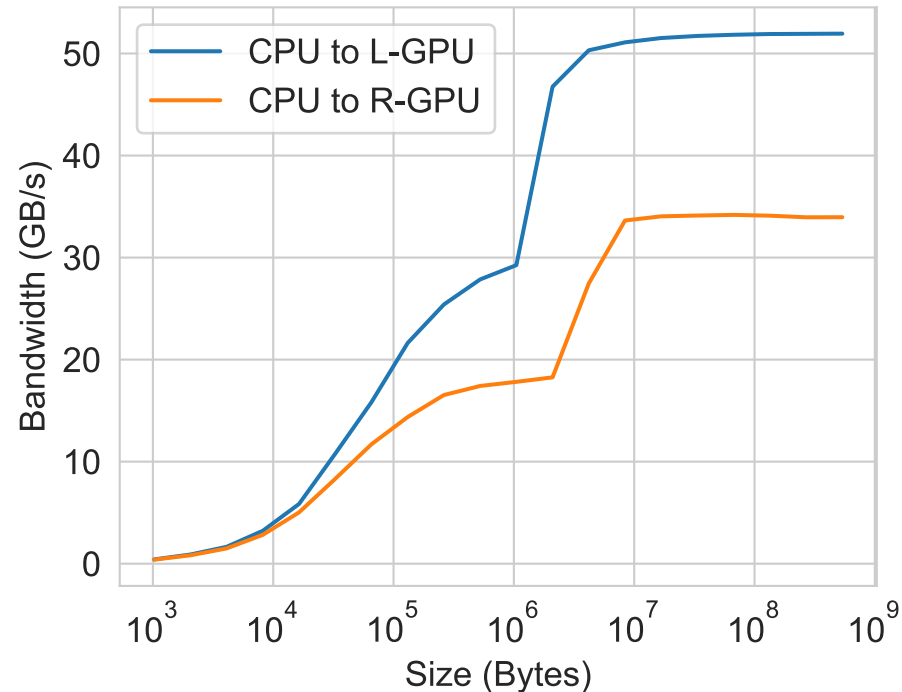
- Each task restricted to CPUs within a NUMA domain
- Leverage local memory
  - Spillover if necessary



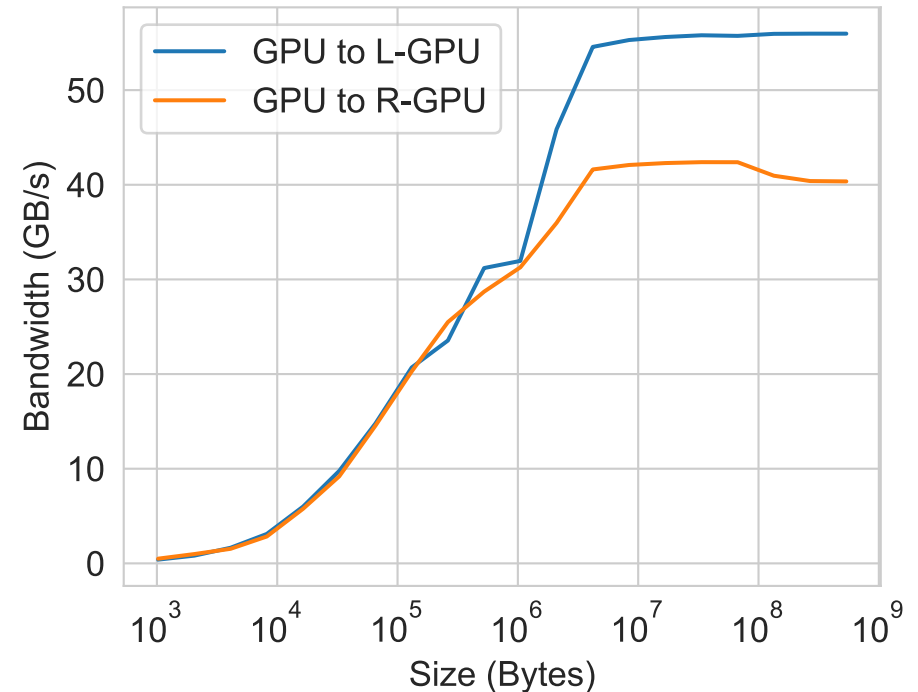
# Selects local CPU-GPU sets

- Does it matter?

Up to 35% penalty for remote transfers



Up to 28% penalty for remote transfers

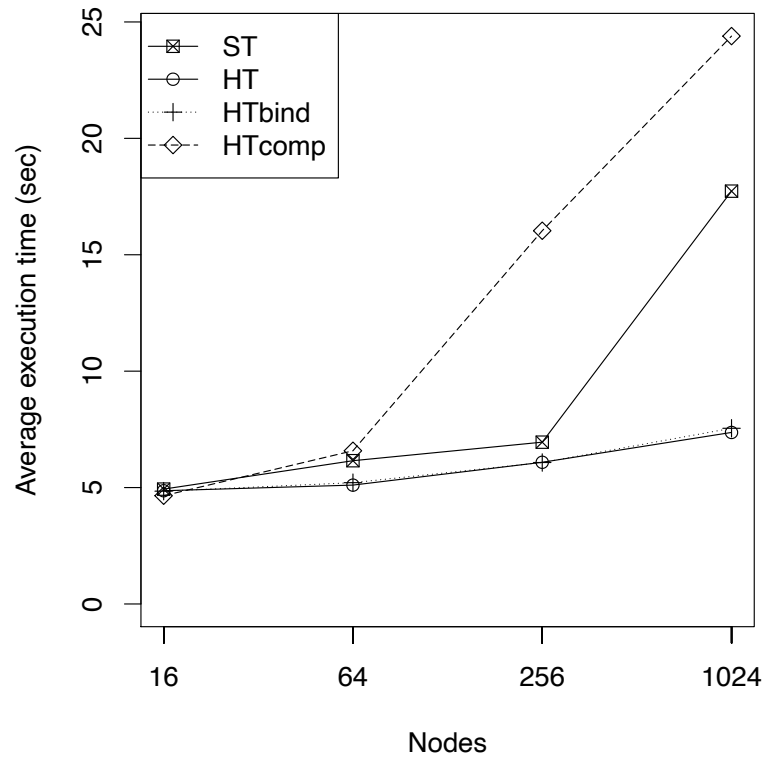
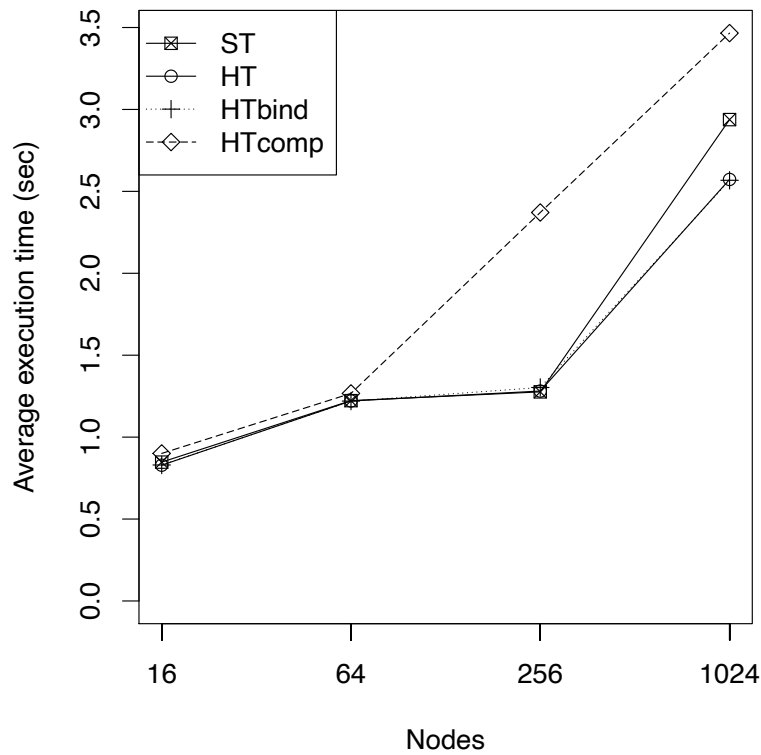


# On LLNL systems with SMT we mitigate noise via thread specialization + mpibind

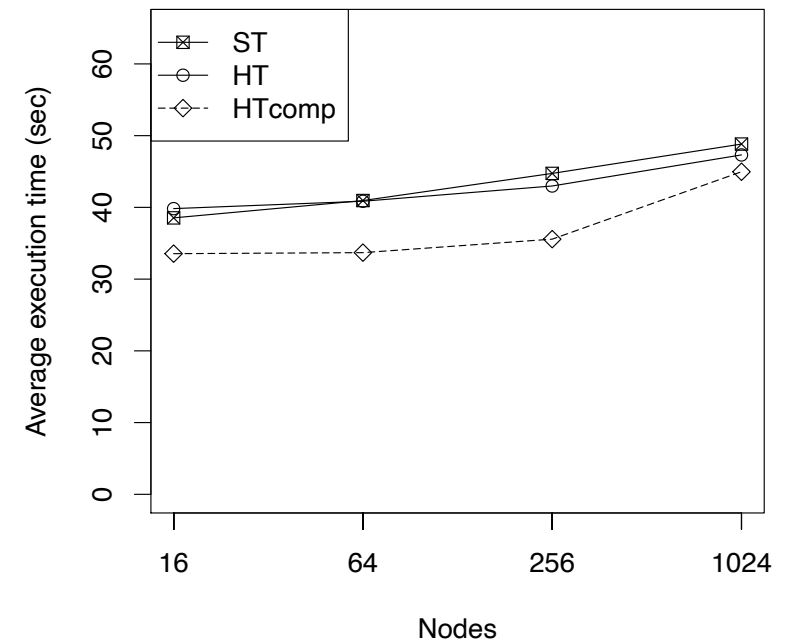
- Thread specialization
  - Application threads
  - System threads
- At boot time (TOSS)
  - Bind system processes to one HW thread of every core
  - *RHEL's tuna*
- At run time
  - Bind application to the other HW thread of every core
  - *LLNL's mpibind*
- But, applications have a choice!
  - `--mpibind=smt:<n>`
  - **Use all resources if needed**
- System noise studies
  - Must look at applications
  - Useful benchmarks
    - Parallel FWQ/FTQ
    - Synchronous collectives
  - Leon et al., IPDPS 2016, SC 2020

# It matters how SMT is used by applications!

AMG and BLAST do not take advantage of additional HW threads, but significantly benefit from using them for system processes



pF3D can leverage the additional HW threads for compute



# mpibind is available in Slurm as a SPANK plugin

- Open source

- MIT license
- Written in C
- Depends on hwloc

- Slurm SPANK plugin

```
$ grep mpibind /etc/slurm/plugstack.conf  
required /usr/lib64/mpibind/mpibind_slurm.so default_off
```



- Building mpibind

- GNU autotools

```
bootstrap  
configure  
make  
make install
```



- Spack

```
spack install mpibind+rocm  
spack install mpibind+cuda
```



<https://github.com/LLNL/mpibind>

# mpibind is an excellent initial policy, but...

- Not suited for apps with dynamically changing mappings
  - Static policy set at job start
- Not intended for benchmarking
  - Performance of remote memory, remote GPUs, etc.
- Does not replace custom mappings
  - Resource manager's affinity masks, etc.
- Advanced use cases will be covered by Quo Vadis  
<https://github.com/hpc/quo-vadis>

# Documentation is available

- Tutorials on mpibind and Slurm affinity  
<https://github.com/LLNL/mpibind/tree/master/tutorials>

- Articles

<b>SC 2020</b>	TOSS-2020: A commodity software stack for HPC
<b>MEMSYS 2018</b>	Achieving transparency mapping parallel applications: A memory hierarchy affair
<b>GTC 2018</b>	Mapping MPI+X applications to multi-GPU architectures: A performance-portable approach
<b>MEMSYS 2017</b>	mpibind: A memory-centric affinity algorithm for hybrid applications
<b>IPDPS 2016</b>	System Noise Revisited: Enabling Application Scalability and Reproducibility with SMT



# mpibind helps applications with performance, portability, and productivity

- Performance
  - Leverages local memory and local devices
  - Mitigates system noise via SMT
  - Maximizes cache per worker
- Productivity
  - Provides a simple interface
  - Requires minimal user input
- Portability
  - Same algorithm across system architectures, MPI libraries, and resource managers

*mpibind*