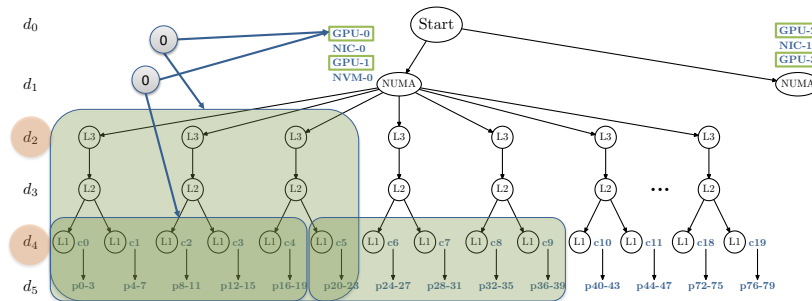


# Mapping Applications to the Hardware Portably and Transparently

Edgar A. León  
Livermore Computing

*mpibind*

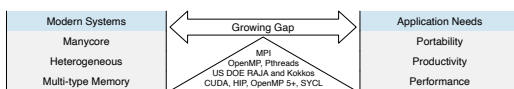
FOSDEM, Brussels  
2 February 2025



Prepared by LLNL under contract DE-AC52-07NA27344. LLNL-PRES-872166.



## HPC users face complex computing architectures



AMD Instinct MI300A

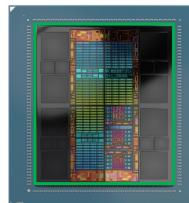


Image credit: AMD

AMD Instinct MI300X

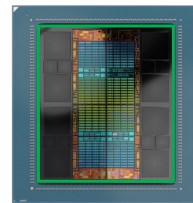
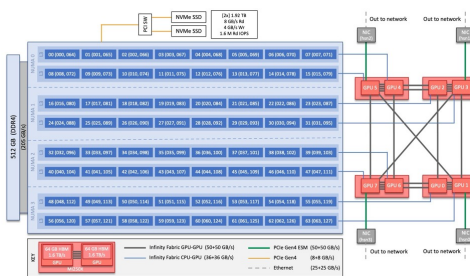


Image credit: AMD

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

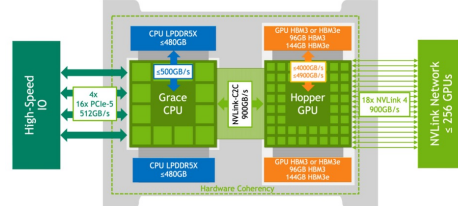
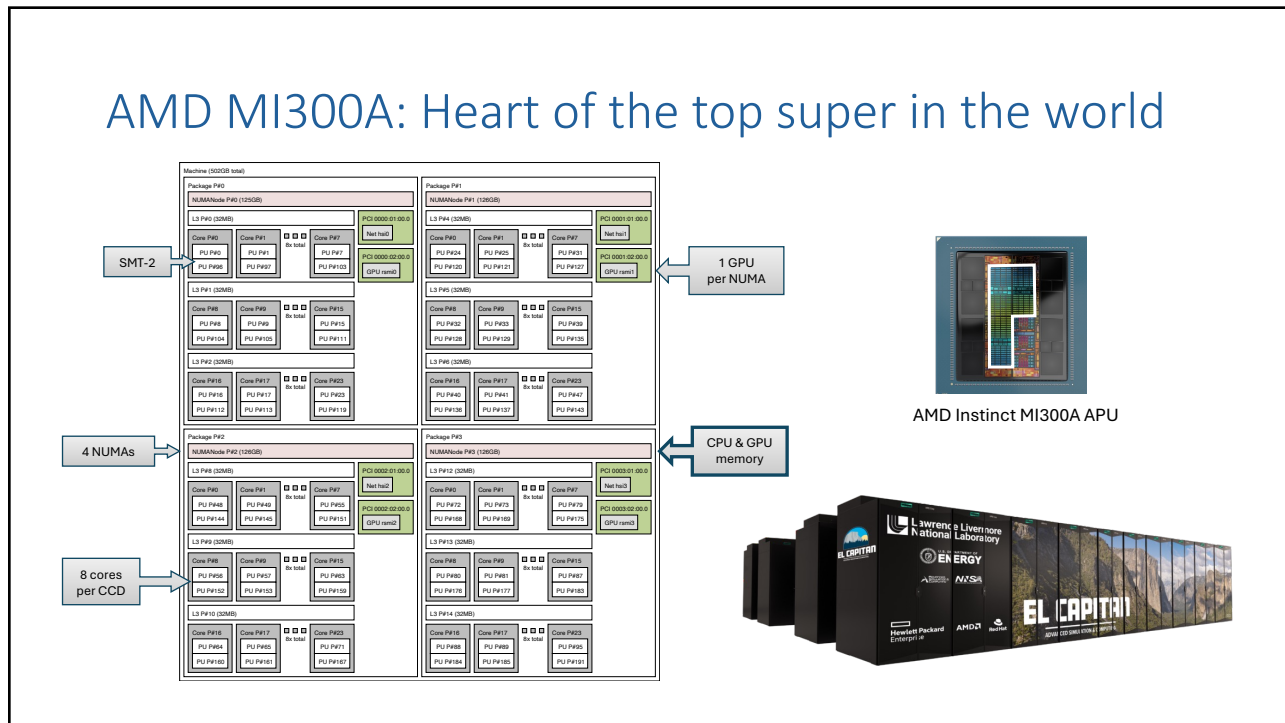
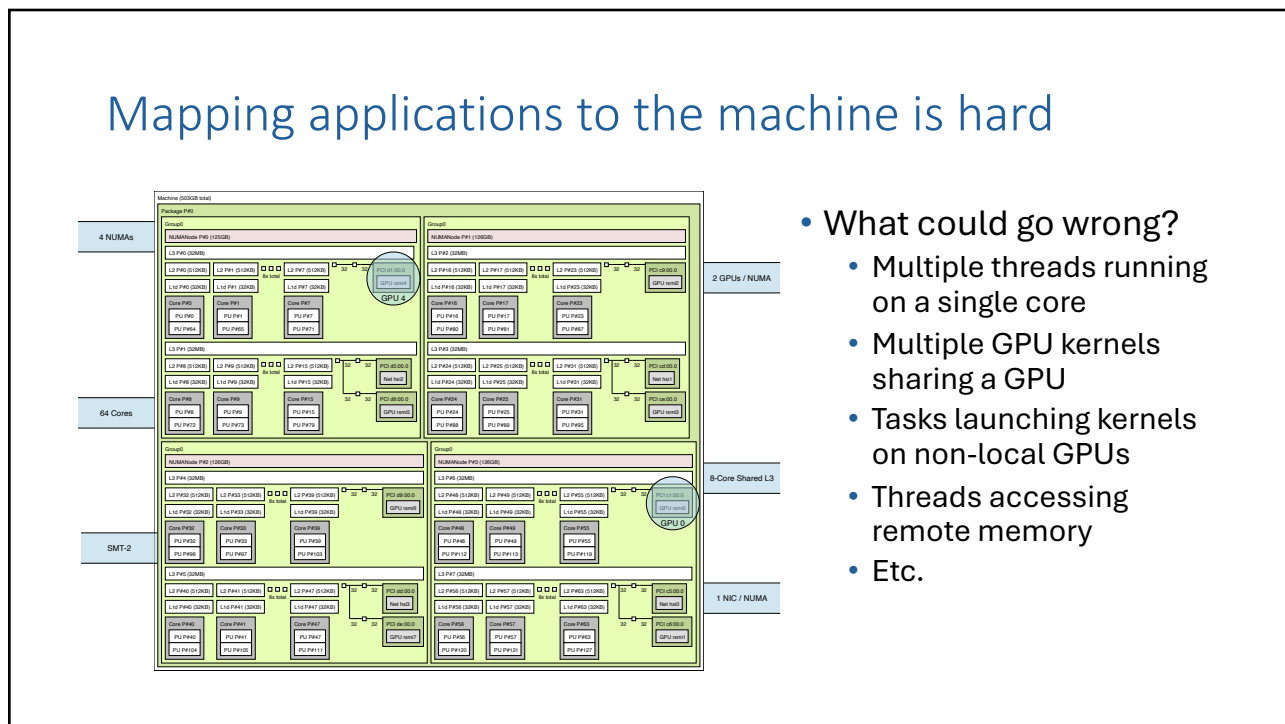


Image credit: NVIDIA

# AMD MI300A: Heart of the top super in the world



# 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

- Memory-system driven
- Locality based
- Portable across architectures



- Properties

1. Provides a simple interface
2. Requires minimal user input
3. Provides cross-system portability
4. Minimizes rem. memory accesses
5. Maximizes cache per worker
6. Leverages compute/mem locality
7. Enables system-noise mitigation
8. Enables job co-scheduling

## Follow the memory system and leverage locality

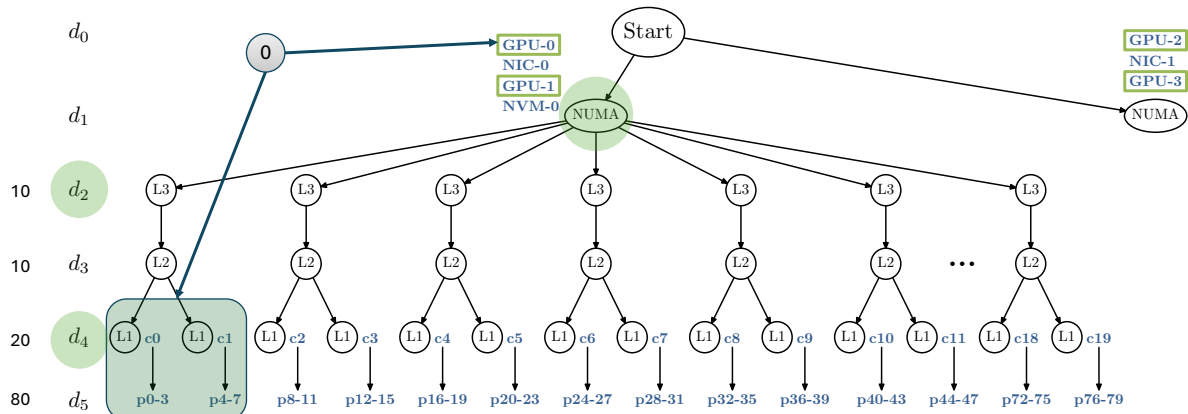
Tasks = 32

→ Tasks/NUMA = 16

→ Workers = 16  
→ Vertices( $d_2$ ) = 10



→ Workers = 16  
→ Vertices( $d_4$ ) = 20



## 1. Provides a simple interface & 2. Requires minimal user input

Coral Node  
2 NUMAs,  
40 cores, 4 GPUs

JSM	<p style="text-align: center; margin: 0;"># 8 MPI tasks across all cores and all GPUs on a CORAL node</p> <pre style="margin: 0; font-family: monospace;">jsrun -a 2 -c 10 -g 1 -r 4 -d packed -b packed:5 &lt;prog&gt;</pre>
mpibind	<p style="text-align: center; margin: 0;"># 8 MPI tasks across all cores and all GPUs on a CORAL node</p> <pre style="margin: 0; font-family: monospace;">srun -n8 &lt;prog&gt; srun -n8 -mpibind=on &lt;prog&gt;</pre>

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

## 3. 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

8

### 3. Provides portability across systems

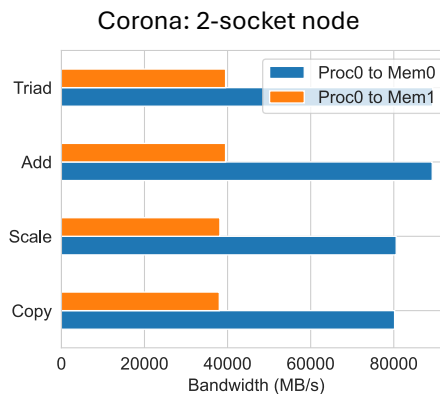
- Mapping algorithm relies on abstract memory-compute tree
  - Portable Hardware Locality (hwloc)
- Mapping algorithm is separate from applying affinity

C Interface	Slurm and Flux Plugins
{Task} → {CPUs, GPUs, Thread mapping}	Use job info as input to mpibind
	Get mpibind mapping
	Bind tasks
	Set environment variables



### 4. Minimizes remote memory accesses

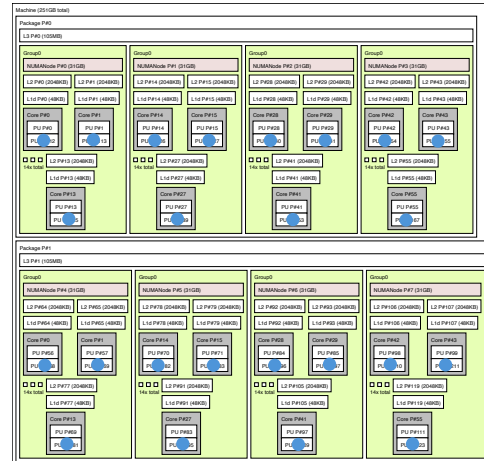
- Each task restricted to CPUs within a NUMA domain
- Leverage local memory
  - Spillover if necessary



10

## 7. Enables system noise mitigation

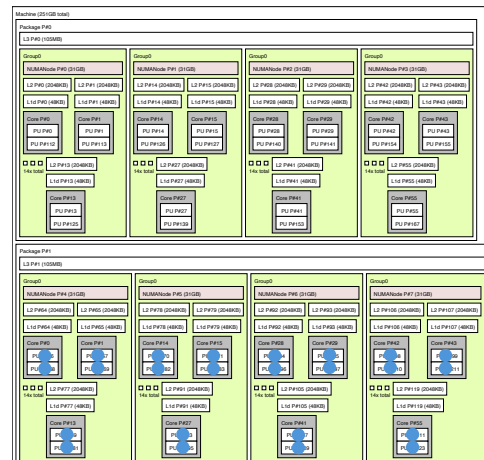
- CPU specialization
  - Application CPUs
  - System CPUs
- Noise mitigation
  - Place app threads on AppCPUs
  - Place system services on SysCPUs
- mpibind controls app placement
  - MPIBIND\_RESTRICT=<cpu|mem>
  - **Apps can still use all resources**



MPIBIND\_RESTRICT=112-223

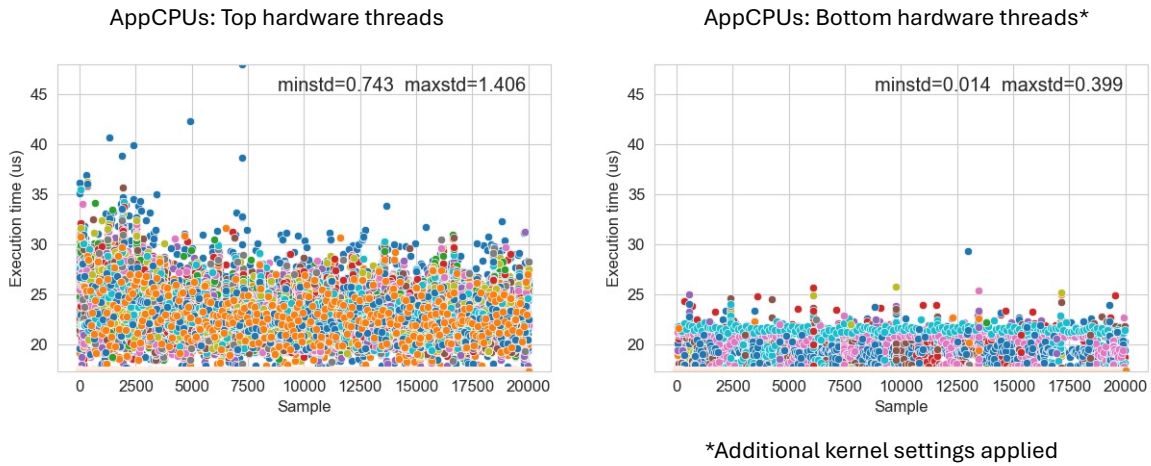
## 7. Enables system noise mitigation

- CPU specialization
  - Application CPUs
  - System CPUs
- Noise mitigation
  - Place app threads on AppCPUs
  - Place system services on SysCPUs
- mpibind controls app placement
  - MPIBIND\_RESTRICT=<cpu|mem>
  - **Apps can still use all resources**



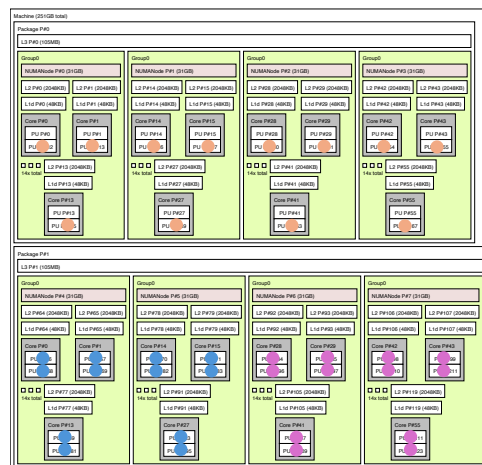
MPIBIND\_RESTRICT\_TYPE=mem  
MPIBIND\_RESTRICT=4-7

## 7. Enables system noise mitigation studies



## 8. Enables arbitrary placement of co-located, concurrent jobs

- **Job 1 on Mem 4-5**
  - MPIBIND\_RESTRICT\_TYPE=mem
  - MPIBIND\_RESTRICT=4-5
  - `srun -n4 -overlap -mpibind=on prog1 &`
- **Job 2 on Mem 6-7**
  - MPIBIND\_RESTRICT\_TYPE=mem
  - MPIBIND\_RESTRICT=6-7
  - `srun -n14 -overlap -mpibind=on prog2 &`
- **Job 3 on bottom CPUs of Mem 0-3**
  - MPIBIND\_RESTRICT\_TYPE=cpu
  - MPIBIND\_RESTRICT=112-167
  - `srun -n4 -overlap -mpibind=on prog3 &`



## LLNL using mpibind in production since 2015

- 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>

18

## 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

- Flux affinity, Slurm affinity, mpibind  
<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 users make better use of modern supercomputers

- Performance

- Leverages local memory and local devices
- Helps mitigate system noise
- Maximizes cache per worker

- Productivity

- Provides a simple interface
- Requires minimal user input

- Portability and vendor-neutrality

- Same algorithm across system architectures, MPI libraries, and resource managers

