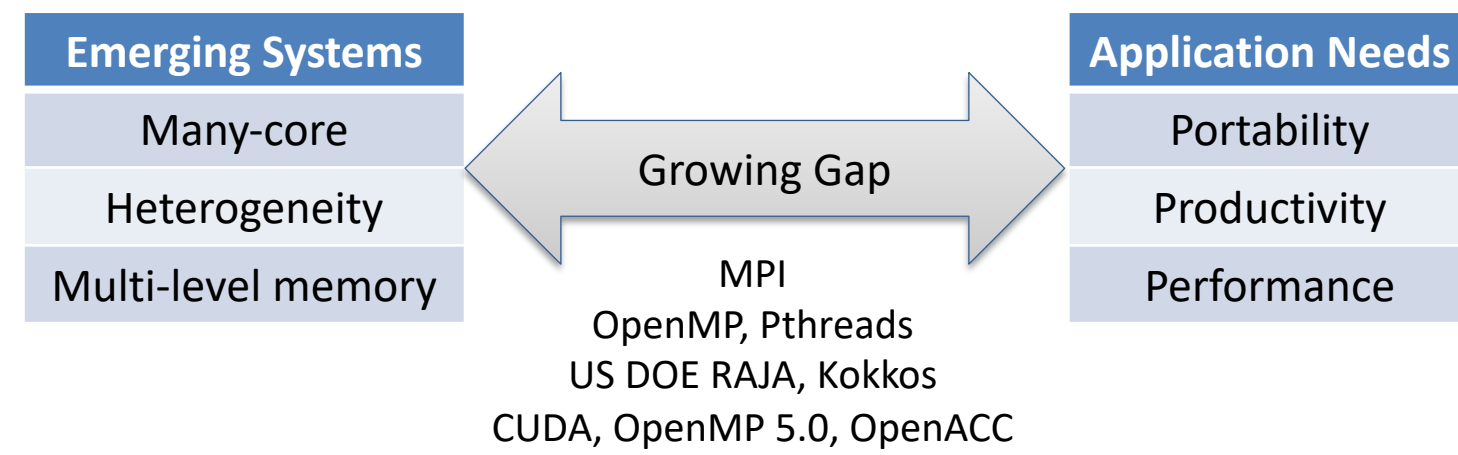


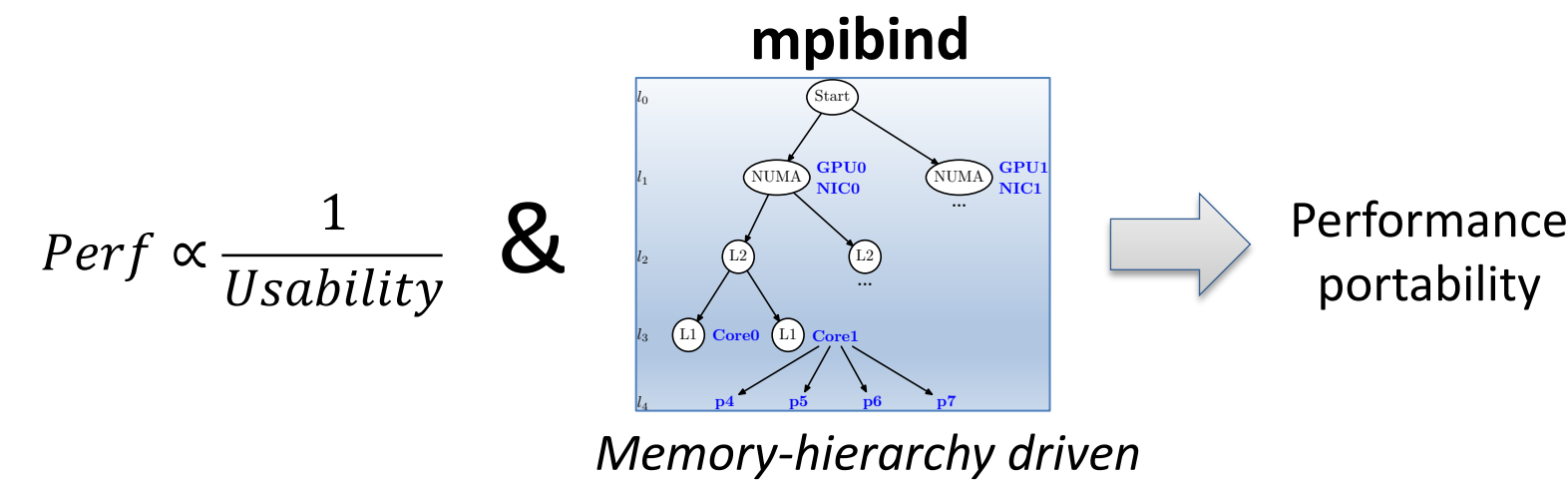
Cross-Architecture Affinity of Supercomputers

Edgar A. León, Lawrence Livermore National Laboratory

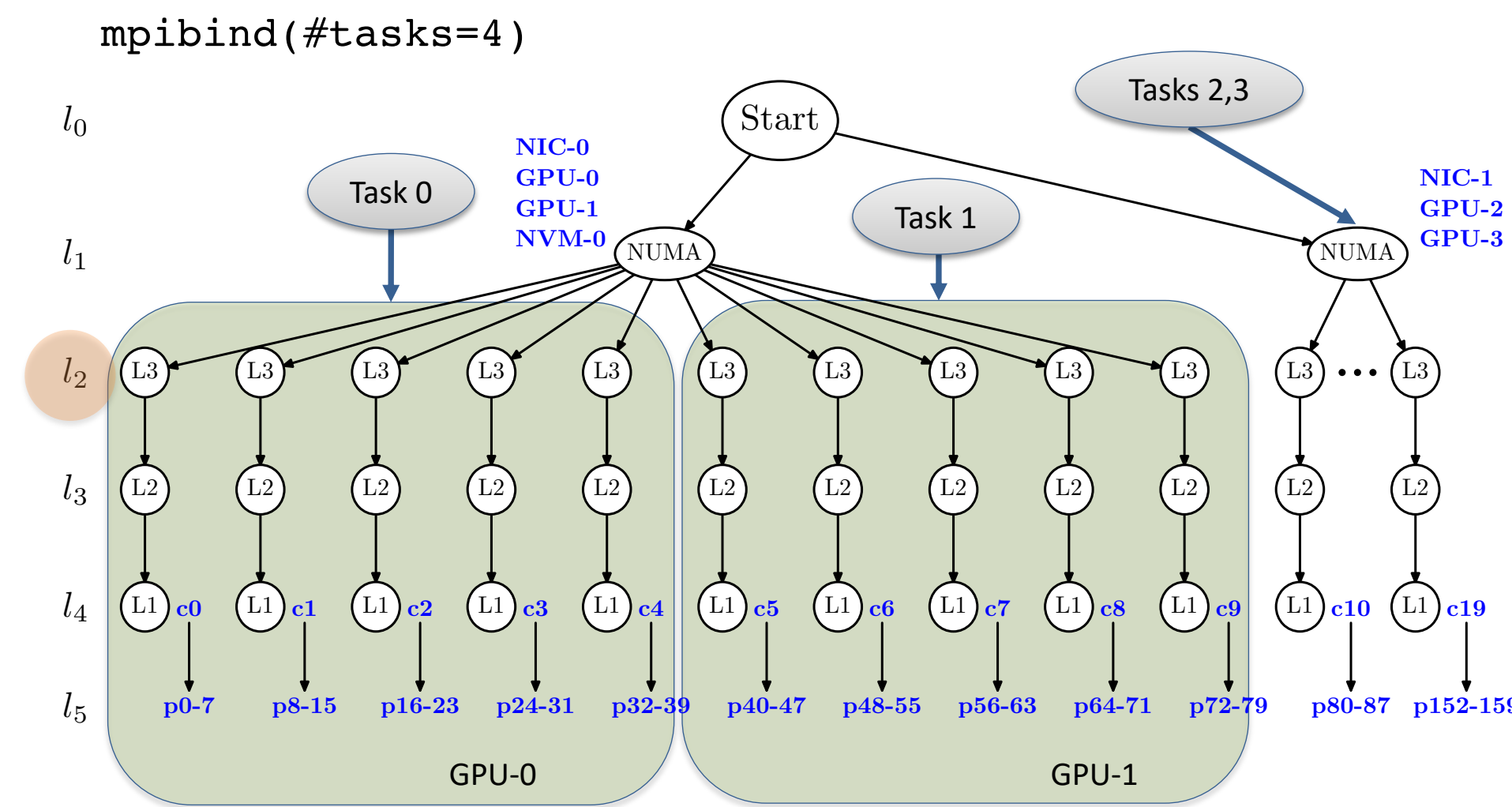
How to map hybrid applications portably to complex machines?



Bridging the gap with *mpibind*: A memory-driven approach to map scientific applications



Primary consideration: *The memory hierarchy*



Key takeaways mapping parallel applications

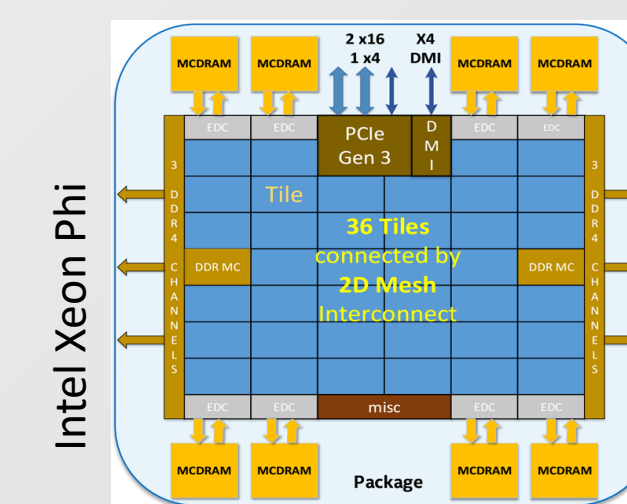
- A memory-driven approach is key for performance portability
 - Many applications are memory bound, not compute bound
 - Memory tree can be derived on most systems
- Maintain a simple interface to improve productivity
 - Computational scientists can focus on science & engineering
- Avoid using hardware topology at the application level
 - May improve performance but may break portability
- Provide a Rookie interface with Pro performance**
 - Provide a simple interface, but
 - Leverage advanced configurations to extract performance

`mpibind(#tasks, [#threads]):`
Mapping of MPI tasks, threads, GPU kernels to processing units and memory domains

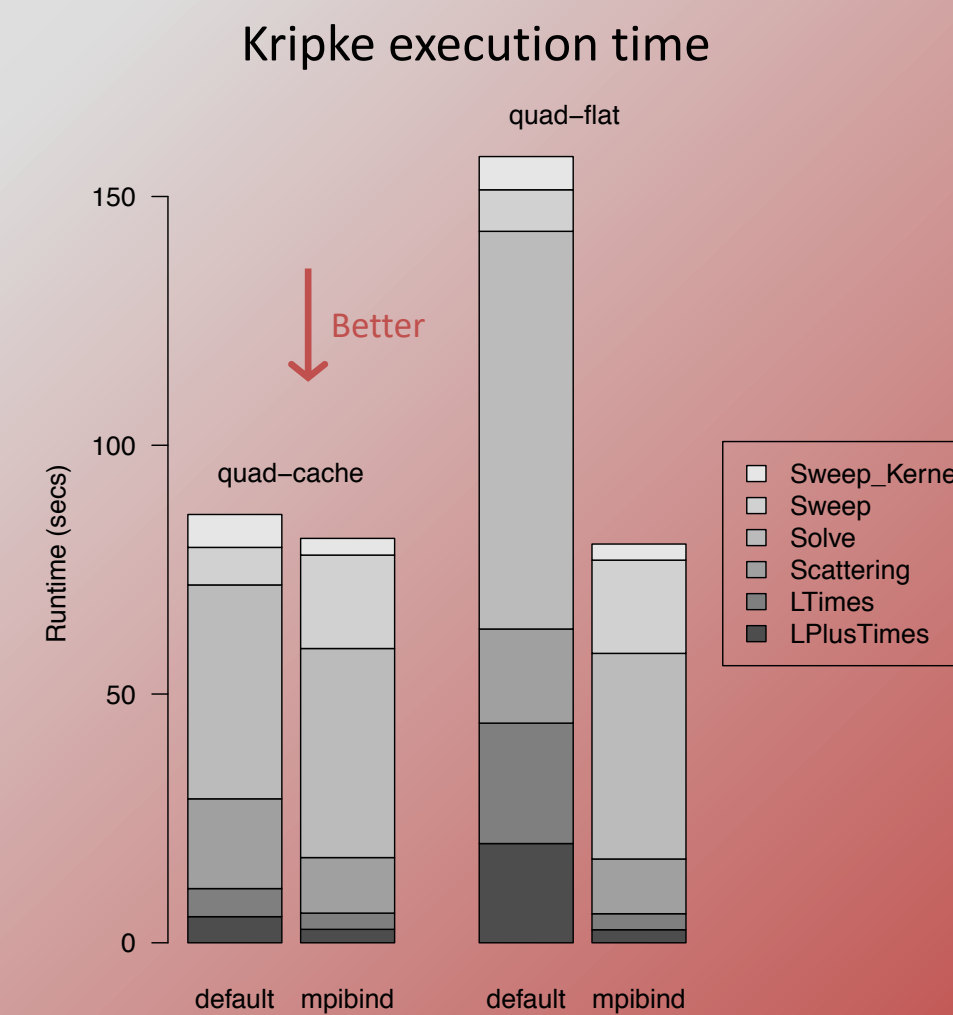
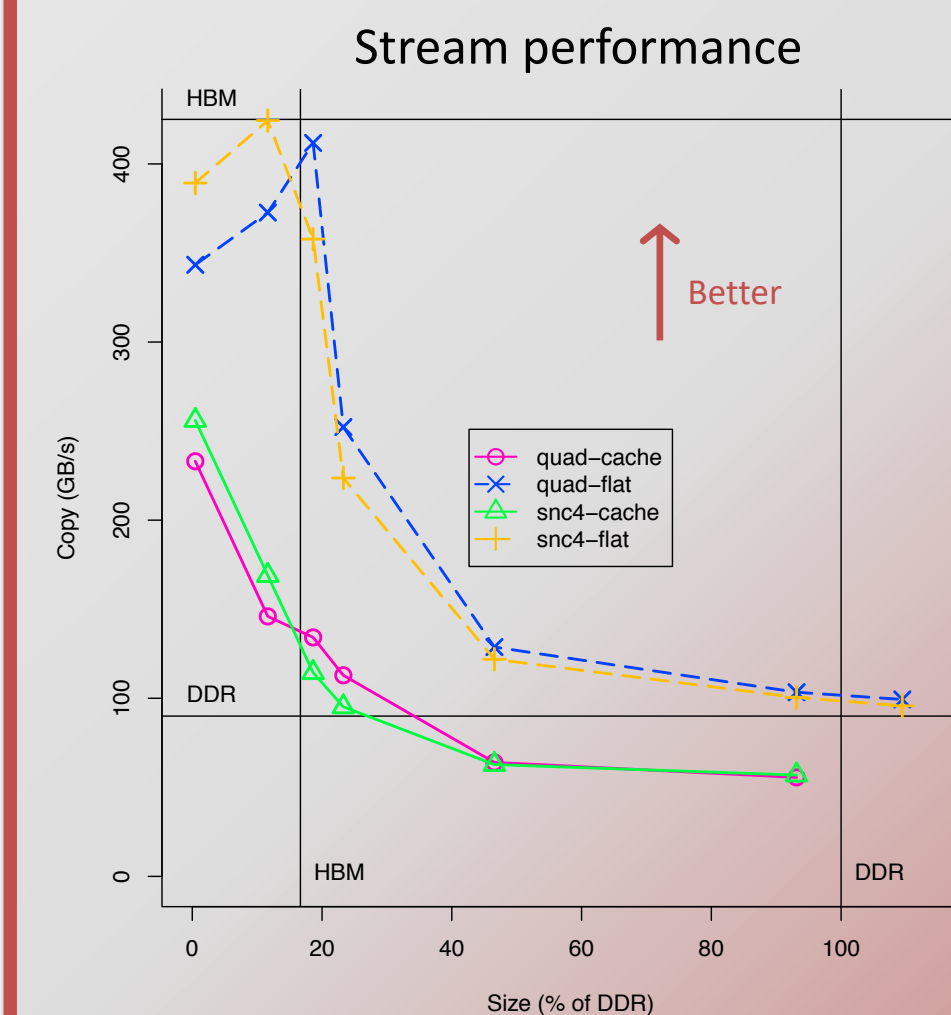
- Build memory tree
- Assign compute resources to memory vertices
- Distribute tasks over NUMA nodes
- Traverse tree to match user workers to memory subtrees and associated compute resources

Application productivity w/2-level memory—HBM+DDR:

- Bridge the gap between performance and ease-of-use
- Users can choose HW mode with exactly the same interface



	HW Mode	NUMA	Use	Perf.	Name
Quad	Cache	1	High	Low	Rookie
	Flat	1x2	Med	Med	Skilled
SNC4	Cache	4	Med	Med	Skilled
	Flat	4x2	Low	High	Pro



Evaluation—3 case studies

	Latency-Optimized Compute	Hybrid CPU + GPU	Many-Core Two-Level Memory
Architecture	IBM Power8+	NVIDIA Pascal GPUs w/ IBM Power8+	Intel Knights Landing
Memory	DDR	HBM + DDR	HBM + DDR
Affinity	Spectrum-MPI default	Spectrum-MPI default	OpenMPI default
	Spectrum-MPI latency	Spectrum-MPI latency	
Application Benchmarks	Spectrum-MPI bandwidth	Spectrum-MPI bandwidth	
	mpibind	mpibind	mpibind
Repetitions	Allreduce: MPI	SW4lite: MPI+OpenMP, MPI+RAJA	Stream: MPI
	1,000,000	10	10
Tasks per Node	4	4	64
Nodes	1, 2, 4, 8, 16	1	1
Case Study	System Noise Mitigation	Application Performance	Application Productivity

