

Towards Memory Disaggregation via NVLink C2C: Benchmarking CPU-Requested GPU Memory Access

Felix Werner, **Marcel Weisgut**, Tilmann Rabl

Hasso Plattner Institute, University of Potsdam, Germany

Motivation

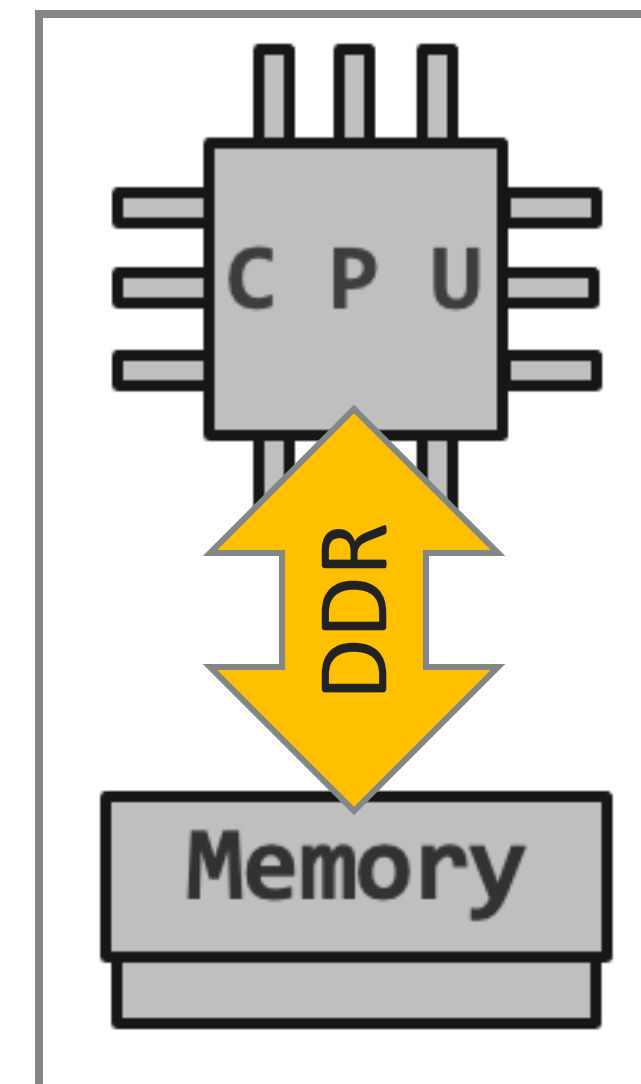
Separation of Resources

- ▶ Database systems have widely adapted separation of compute & storage
 - ▶ Elastic scaling of storage on demand
 - ▶ Avoid resource over-provisioning → cost reduction

Motivation

Separation of Resources

- ▶ Database systems have widely adapted separation of compute & storage
 - ▶ Elastic scaling of storage on demand
 - ▶ Avoid resource over-provisioning → cost reduction
- ▶ Compute & memory tightly coupled
 - ▶ Result: stranded memory



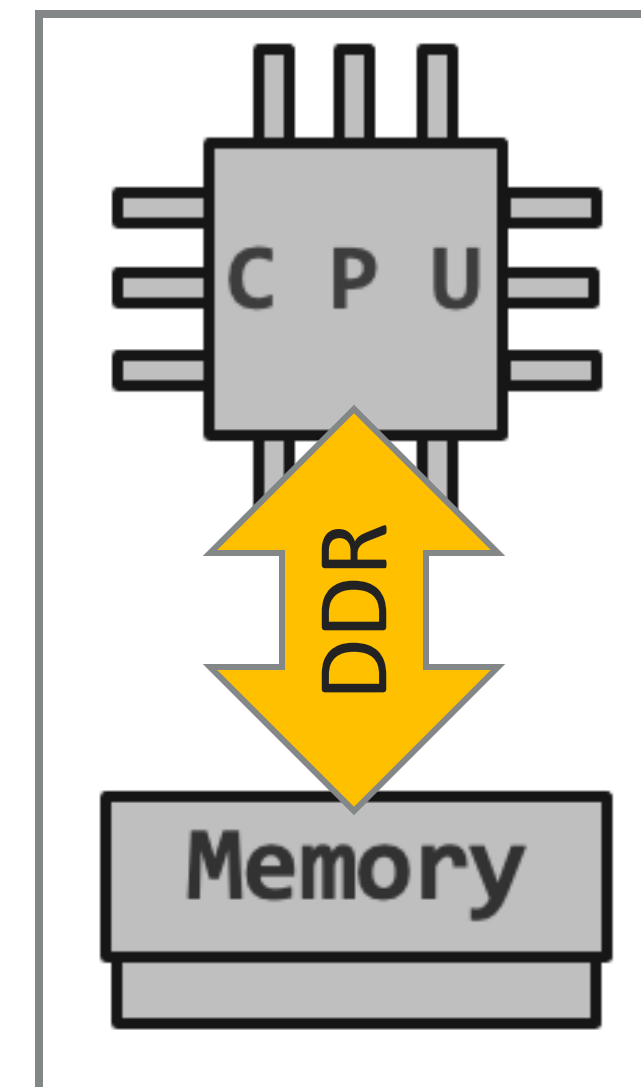
Motivation

Separation of Resources

- ▶ Database systems have widely adapted separation of compute & storage
 - ▶ Elastic scaling of storage on demand
 - ▶ Avoid resource over-provisioning → cost reduction
- ▶ Compute & memory tightly coupled
 - ▶ Result: stranded memory

Memory Disaggregation

- ▶ Initiatives towards memory disaggregation to separate memory from compute resources
 - ▶ Compute Express Link (CXL) as a recent option



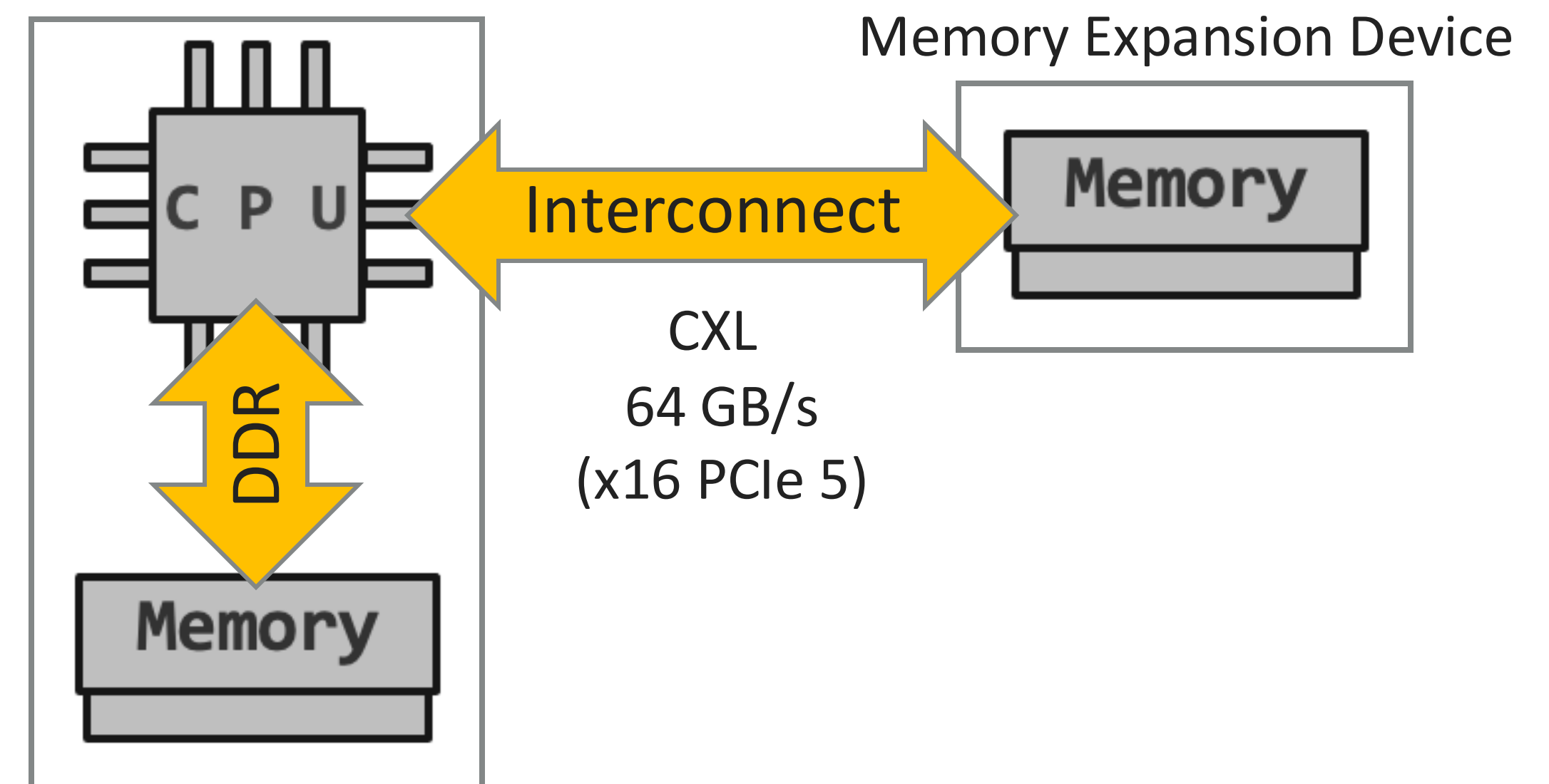
Motivation

Separation of Resources

- ▶ Database systems have widely adapted separation of compute & storage
 - ▶ Elastic scaling of storage on demand
 - ▶ Avoid resource over-provisioning → cost reduction
- ▶ Compute & memory tightly coupled
 - ▶ Result: stranded memory

Memory Disaggregation

- ▶ Initiatives towards memory disaggregation to separate memory from compute resources
 - ▶ Compute Express Link (CXL) as a recent option



Motivation

Separation of Resources

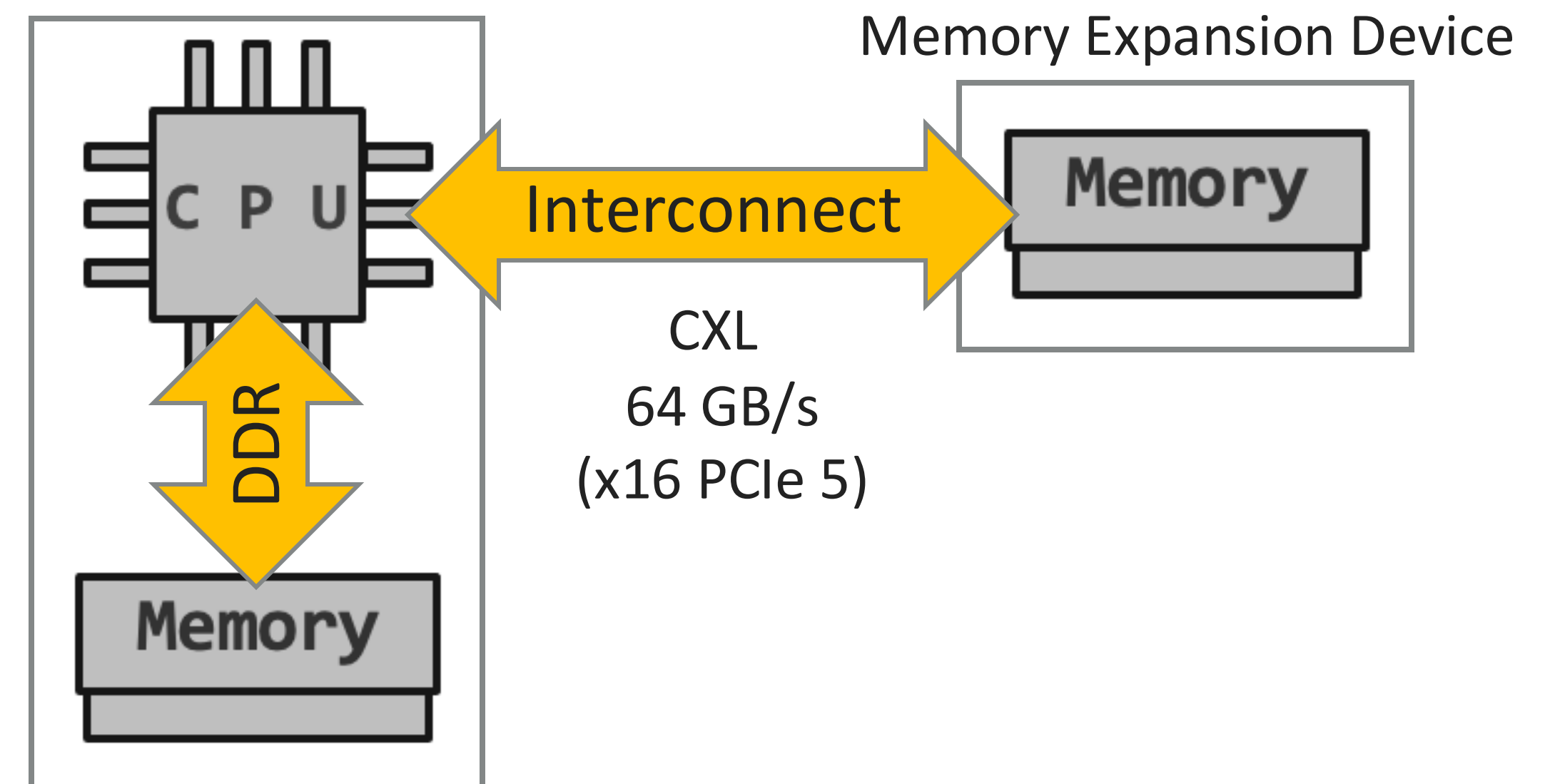
- ▶ Database systems have widely adapted separation of compute & storage
 - ▶ Elastic scaling of storage on demand
 - ▶ Avoid resource over-provisioning → cost reduction
- ▶ Compute & memory tightly coupled
 - ▶ Result: stranded memory

Memory Disaggregation

- ▶ Initiatives towards memory disaggregation to separate memory from compute resources
 - ▶ Compute Express Link (CXL) as a recent option

NVLink versions

Version	Data rate per link [GB/s]	# Lanes per link	# Links	Theoretical bandwidth	Architecture
1	20	8	4	80	Pascal
2	25	8	6	150	Volta
3	25	4	12	300	Ampere
4 / C2C	25	4	18	450	Hopper



Motivation

Separation of Resources

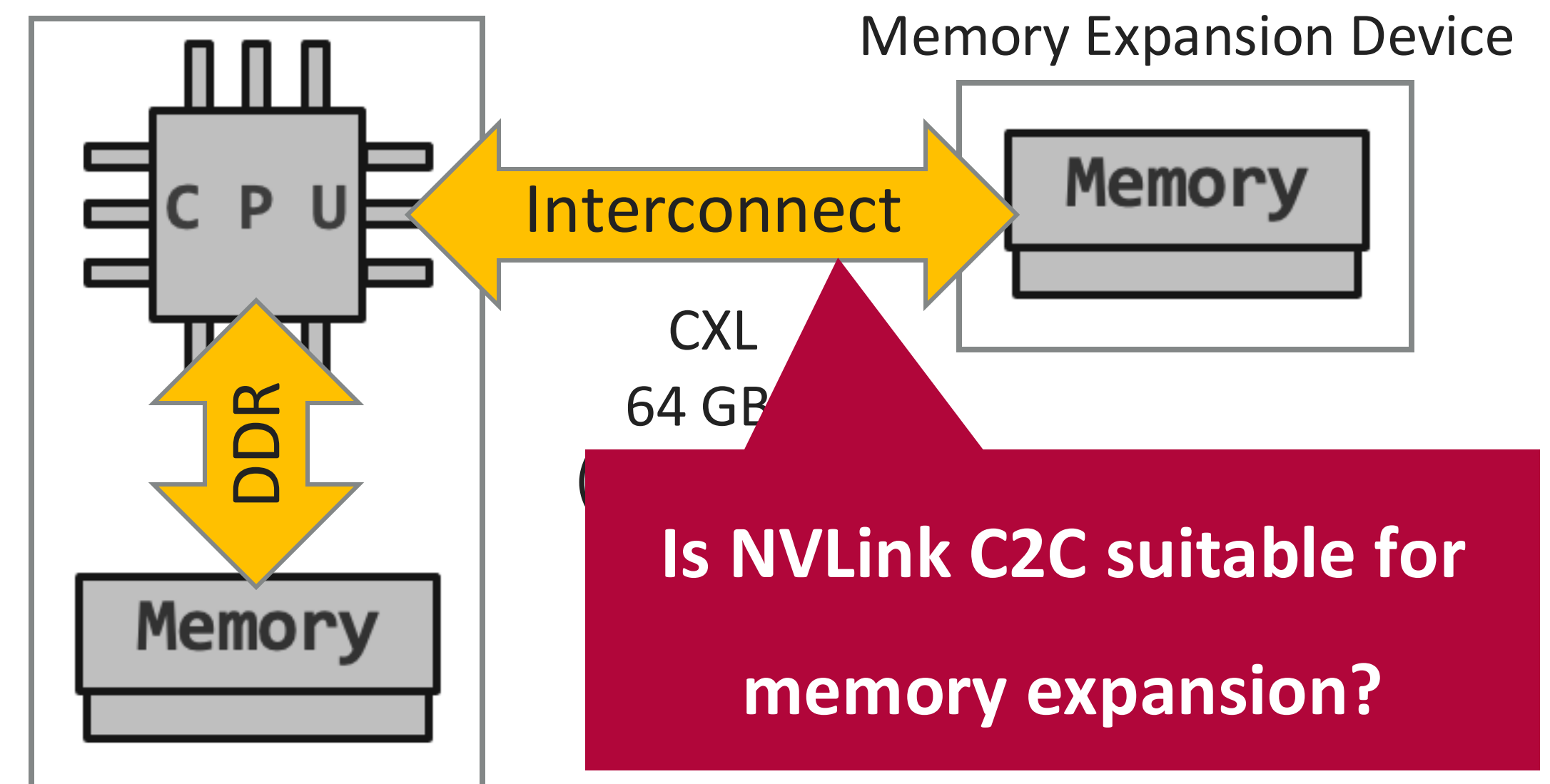
- ▶ Database systems have widely adapted separation of compute & storage
 - ▶ Elastic scaling of storage on demand
 - ▶ Avoid resource over-provisioning → cost reduction
- ▶ Compute & memory tightly coupled
 - ▶ Result: stranded memory

Memory Disaggregation

- ▶ Initiatives towards memory disaggregation to separate memory from compute resources
 - ▶ Compute Express Link (CXL) as a recent option

NVLink versions

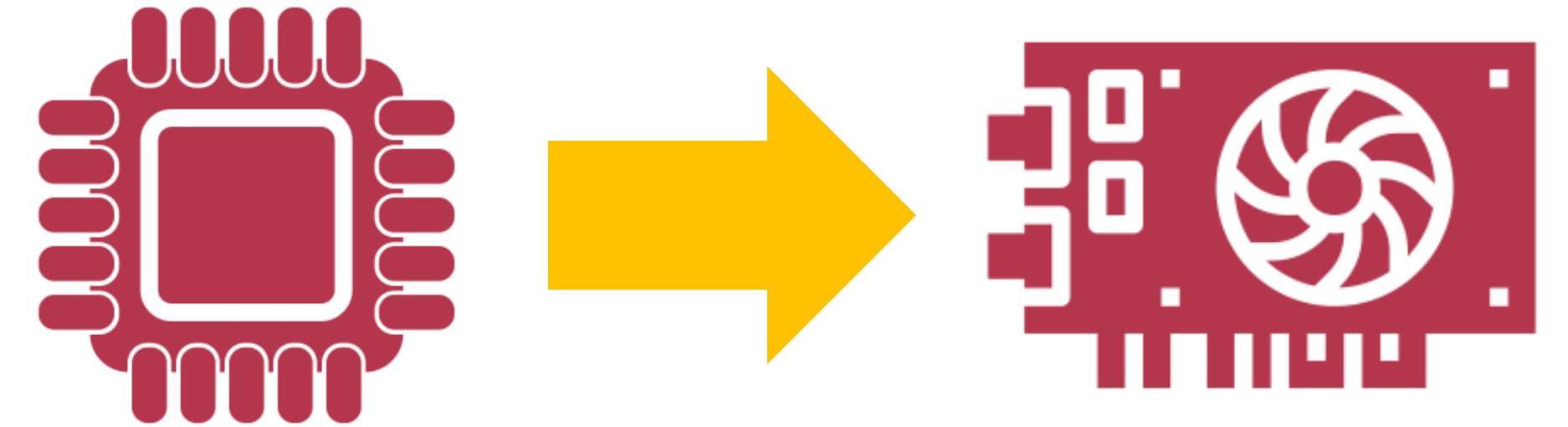
Version	Data rate per link [GB/s]	# Lanes per link	# Links	Theoretical bandwidth	Architecture
1	20	8	4	80	Pascal
2	25	8	6	150	Volta
3	25	4	12	300	Ampere
4 / C2C	25	4	18	450	Hopper



Contributions

- ▶ Performance evaluation of memory access performance over NVLink C2C

- ▶ CPU → GPU memory: throughput, latency
- ▶ CPU → CPU & GPU memory: throughput expansion

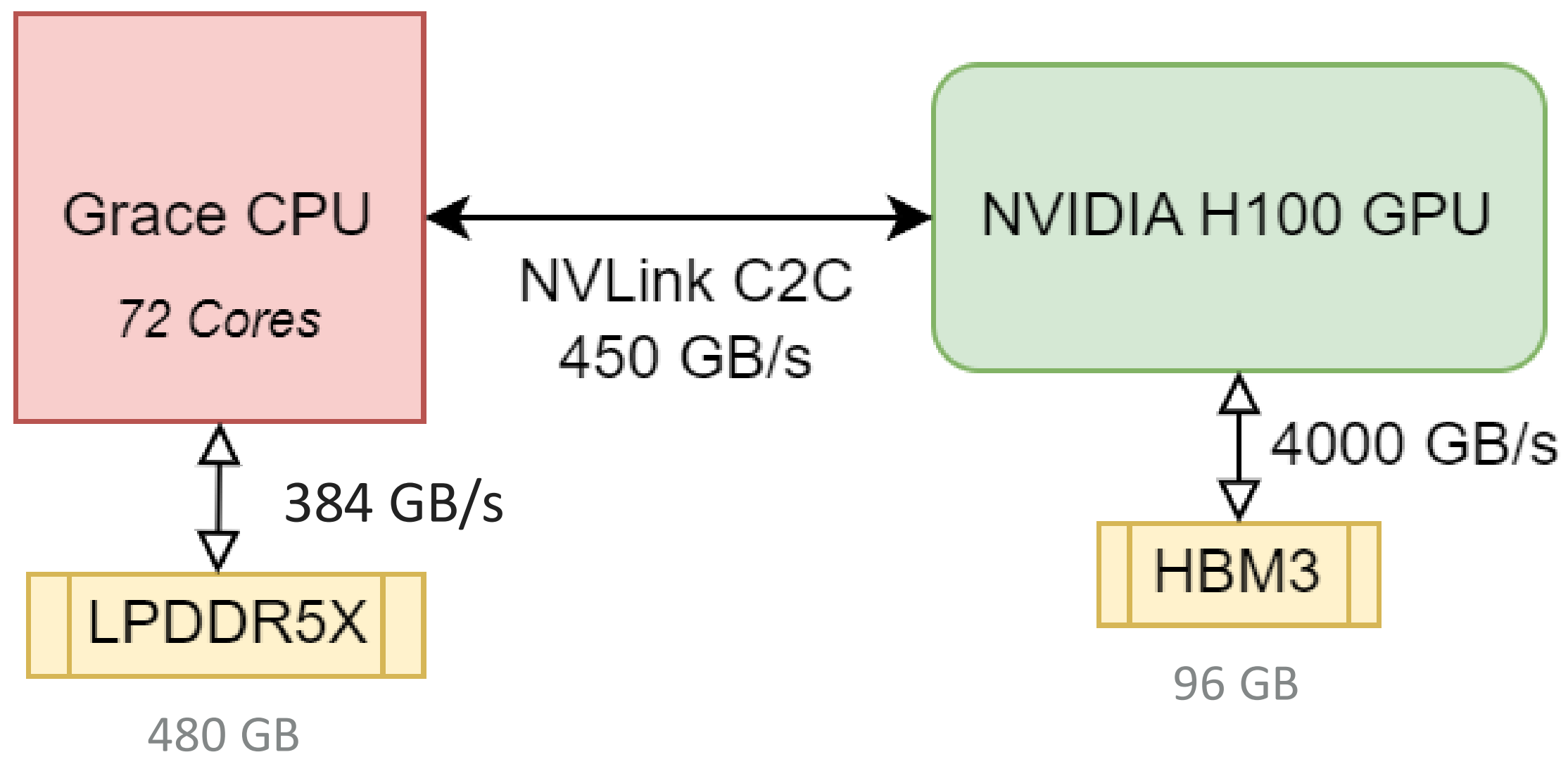


- ▶ Impact of storing data in GPU memory on database operations

- ▶ Discuss suitability of NVLink-attached memory for memory expansion

Hardware Setup

NVIDIA Grace-Hopper Superchip



Bandwidth & Latency

Sustained Bandwidth [GB/s]



CPU → CPU Memory		CPU → GPU Memory	
Seq Read	Seq Write	Seq Read	Seq Write
362	370	130	163
		36%	
		CPU → CPU	
Rnd Read	Rnd Write	Rnd Read	Rnd Write
345	370	128	168
			45%
			CPU → CPU

Bandwidth & Latency

Sustained Bandwidth [GB/s]



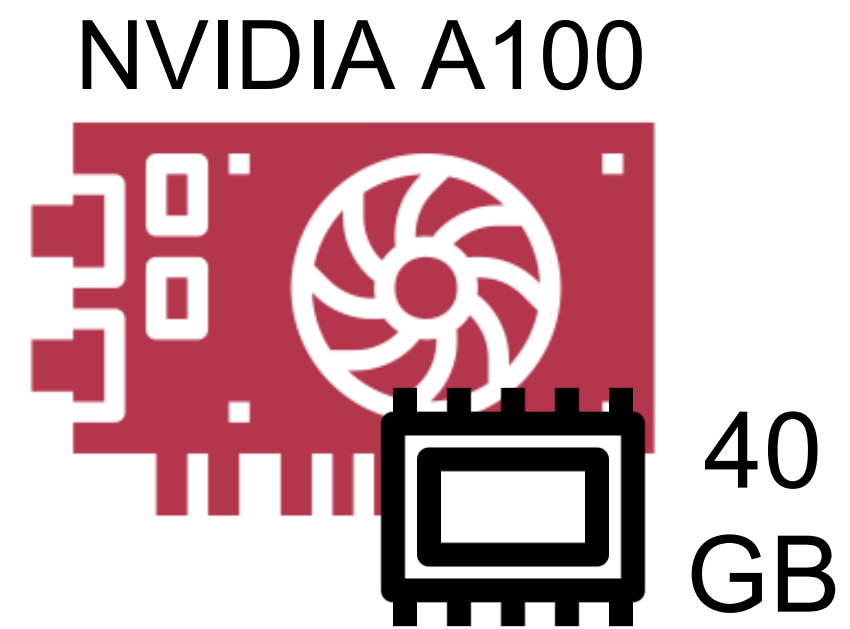
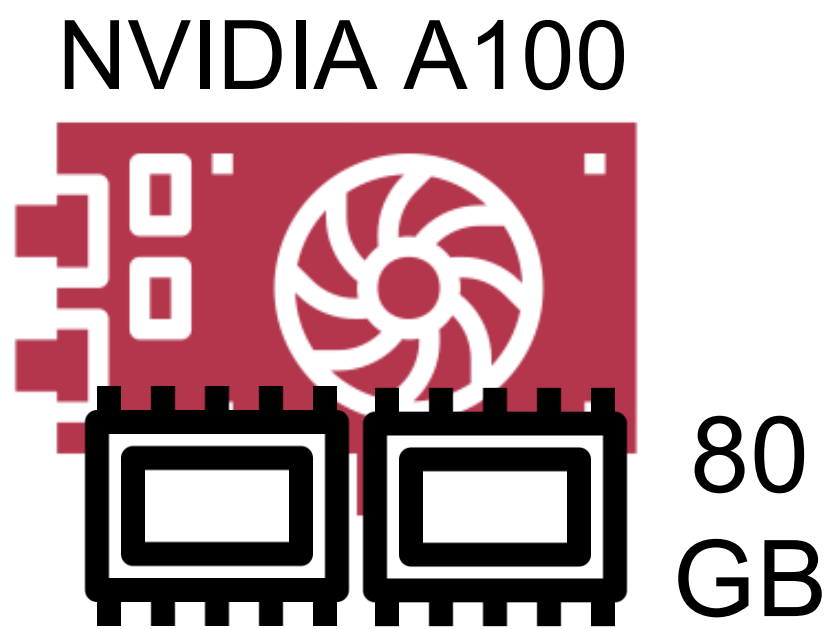
Latency [ns] (random reads)



	CPU→CPU Memory		CPU→GPU Memory	
	Seq Read	Seq Write	Seq Read	Seq Write
	362	370	130	163
			36%	
			CPU→CPU	
	Rnd Read	Rnd Write	Rnd Read	Rnd Write
	345	370	128	168
				45%
				CPU→CPU
	Idle	Idle	Idle	Idle
	220		810 (3.7x)	
	Loaded	Loaded	Loaded	Loaded
	440		1020 (2.3x)	

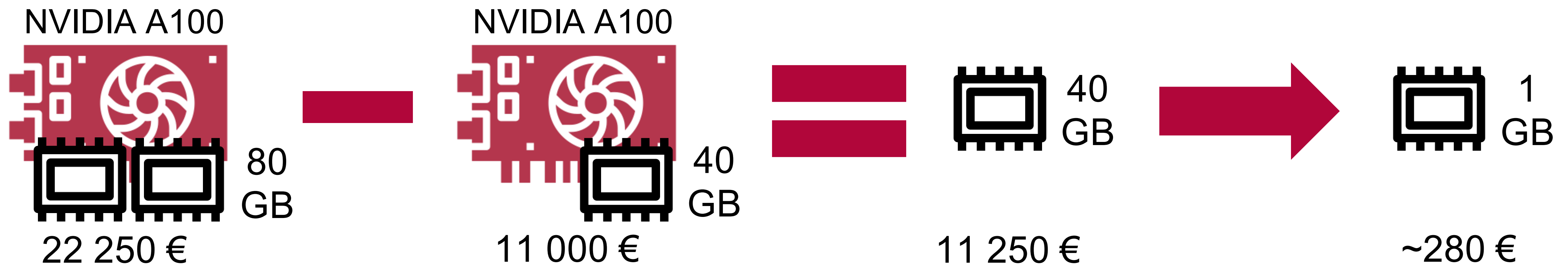
Economic Viability

- ▶ HBM is on-package memory → not modular like DRAM DIMMs
- ▶ Price estimation attempt: compare GPU models with different memory capacities



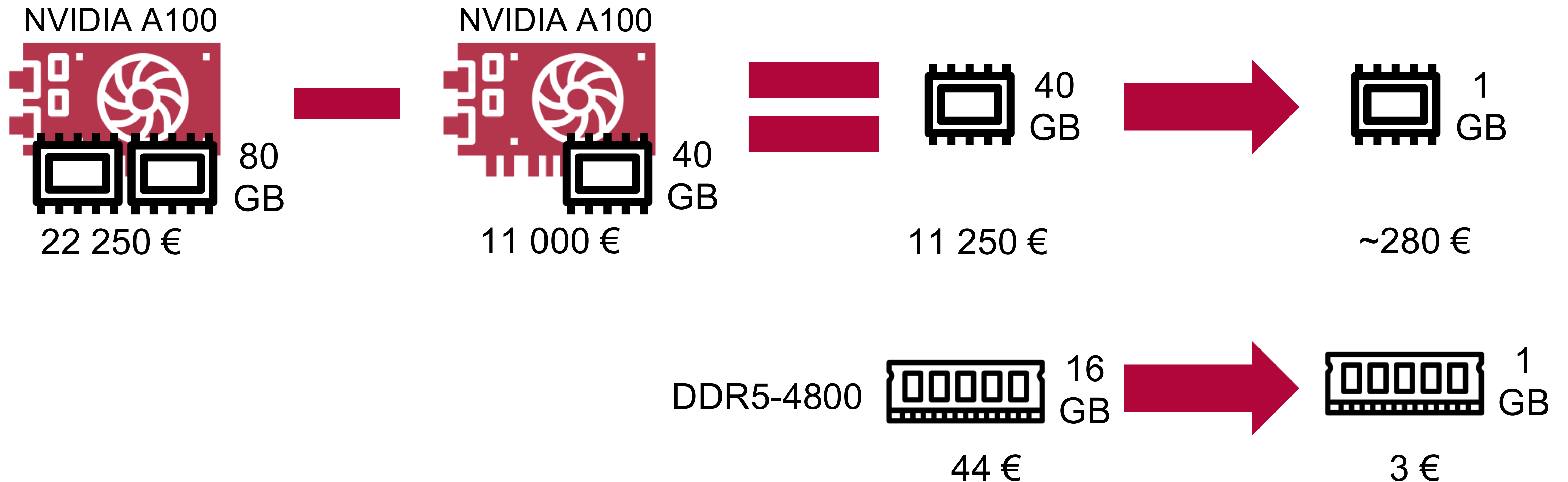
Economic Viability

- ▶ HBM is on-package memory → not modular like DRAM DIMMs
- ▶ Price estimation attempt: compare GPU models with different memory capacities

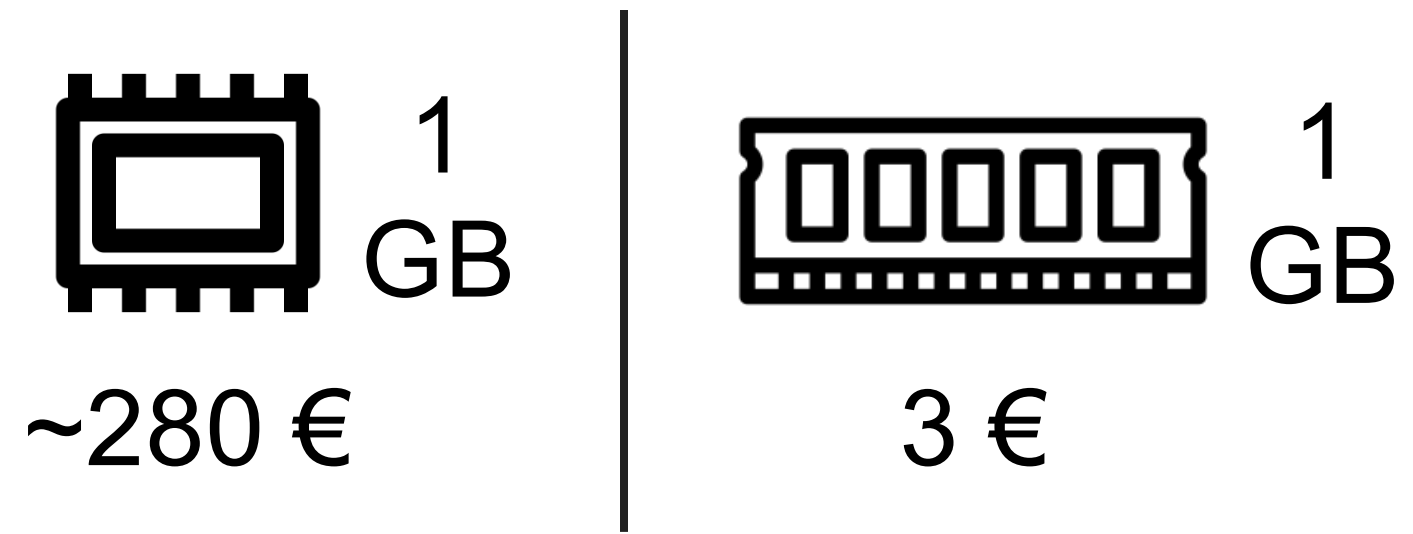


Economic Viability

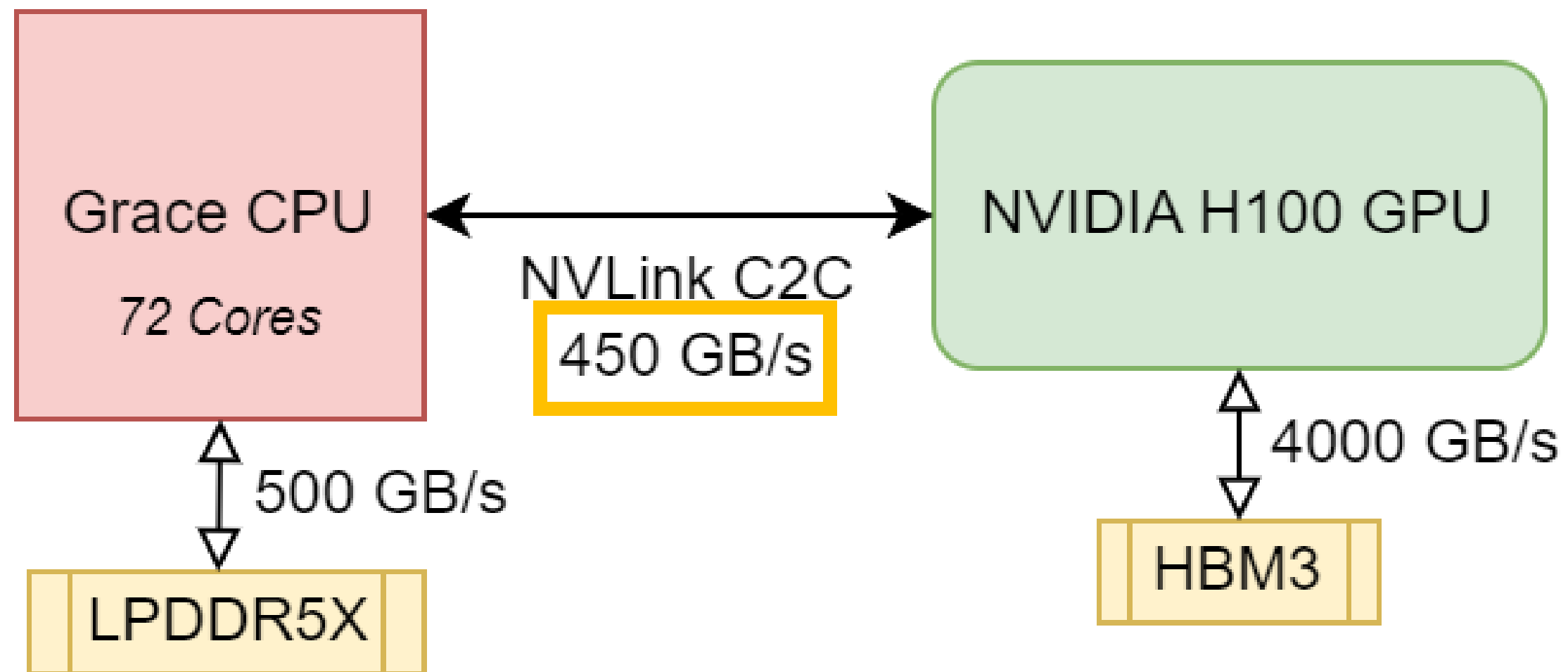
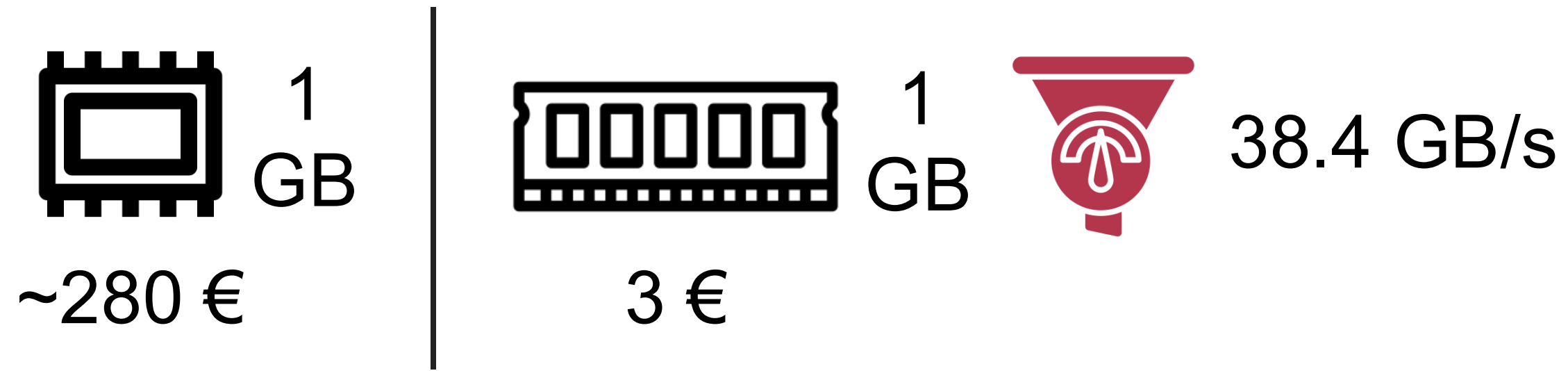
- ▶ HBM is on-package memory → not modular like DRAM DIMMs
- ▶ Price estimation attempt: compare GPU models with different memory capacities



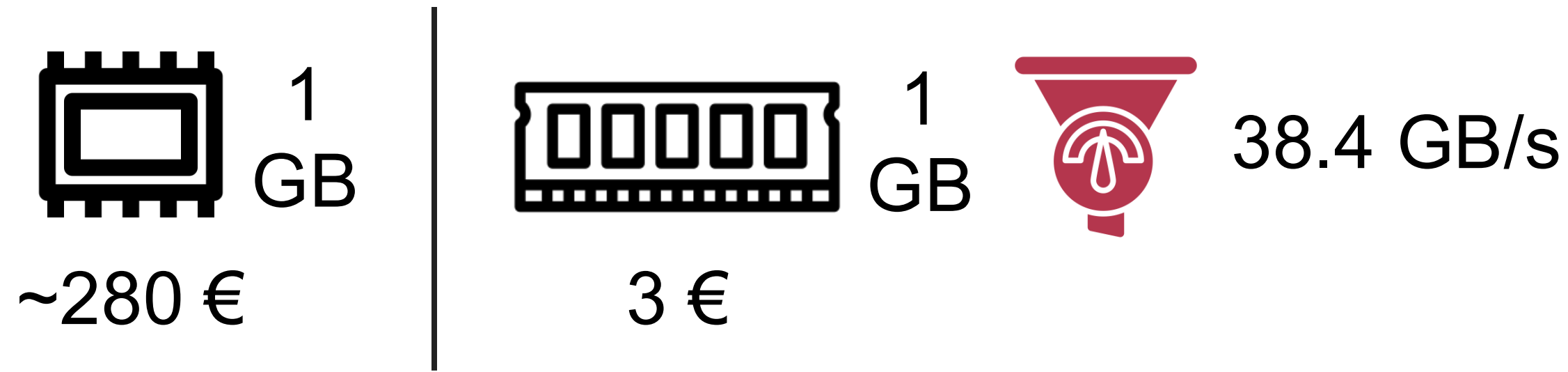
Economic Viability



Economic Viability

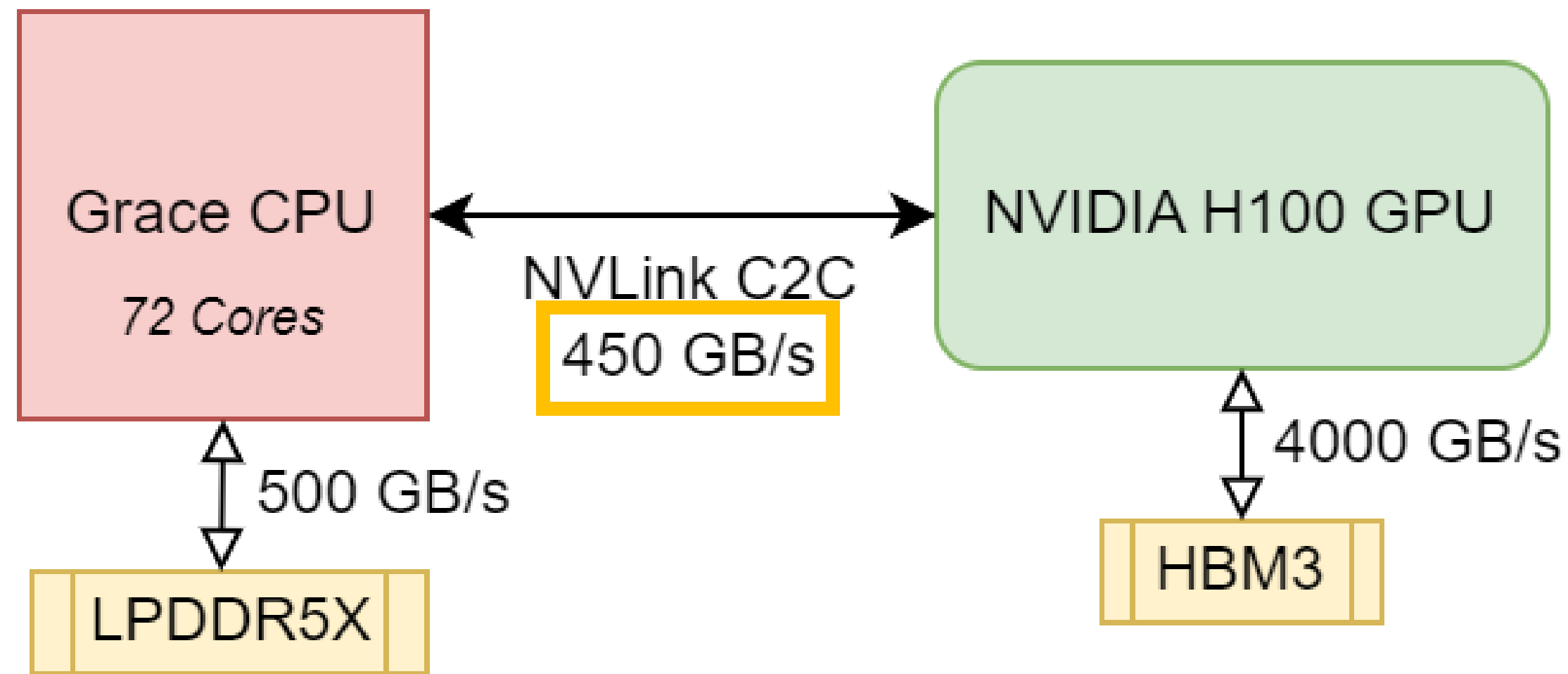


Economic Viability

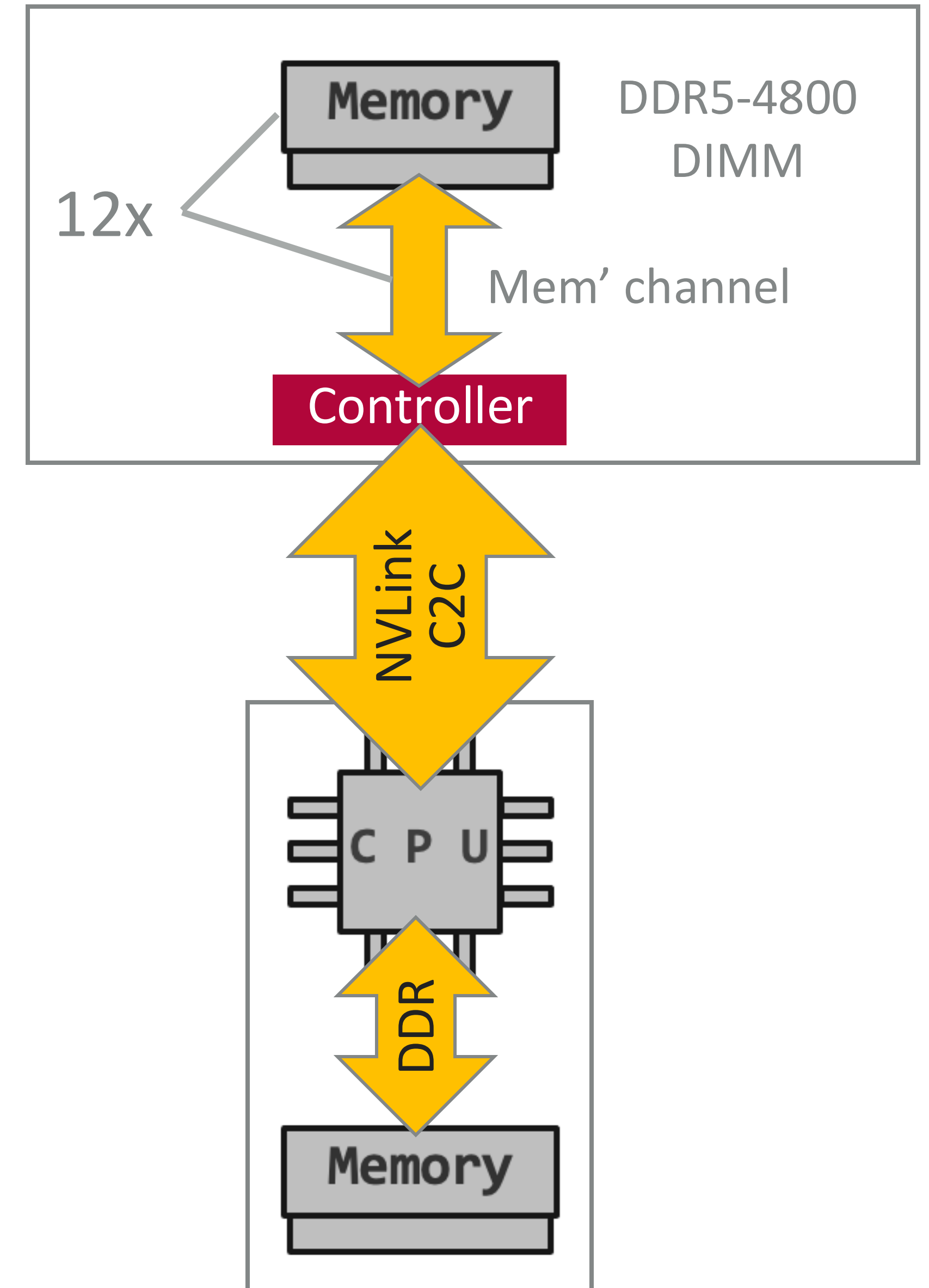


DIMM setup

- 460.8 GB/s
- 192 GB
- 528 €



Memory Expansion Device



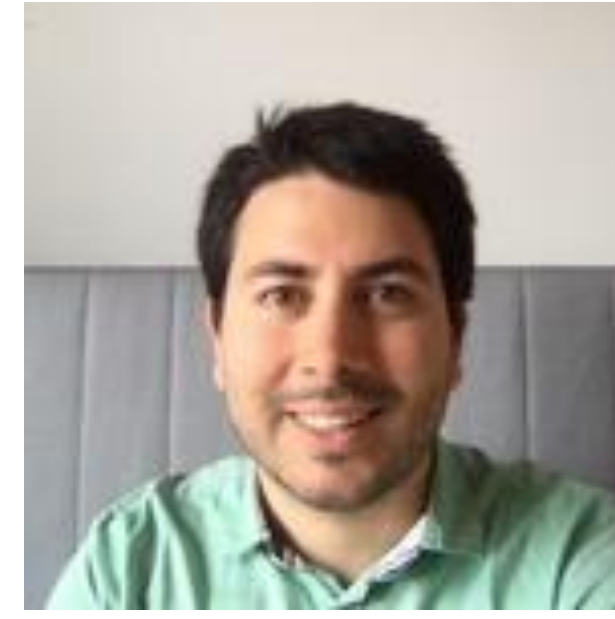
Summary

- ▶ CPU→GPU achieves 36% and 45% of the CPU→CPU read/write throughput
- ▶ CPU→GPU has access latencies of between ~800 ns to ~1000 ns (3.7x / 2.3x of CPU→CPU)
- ▶ Bandwidth expansion study shows improvements between of 1.3x for reads and 1.7x for writes
- ▶ Multiple DDR DIMMs more cost-efficient for pure memory expansion than GPU HBM

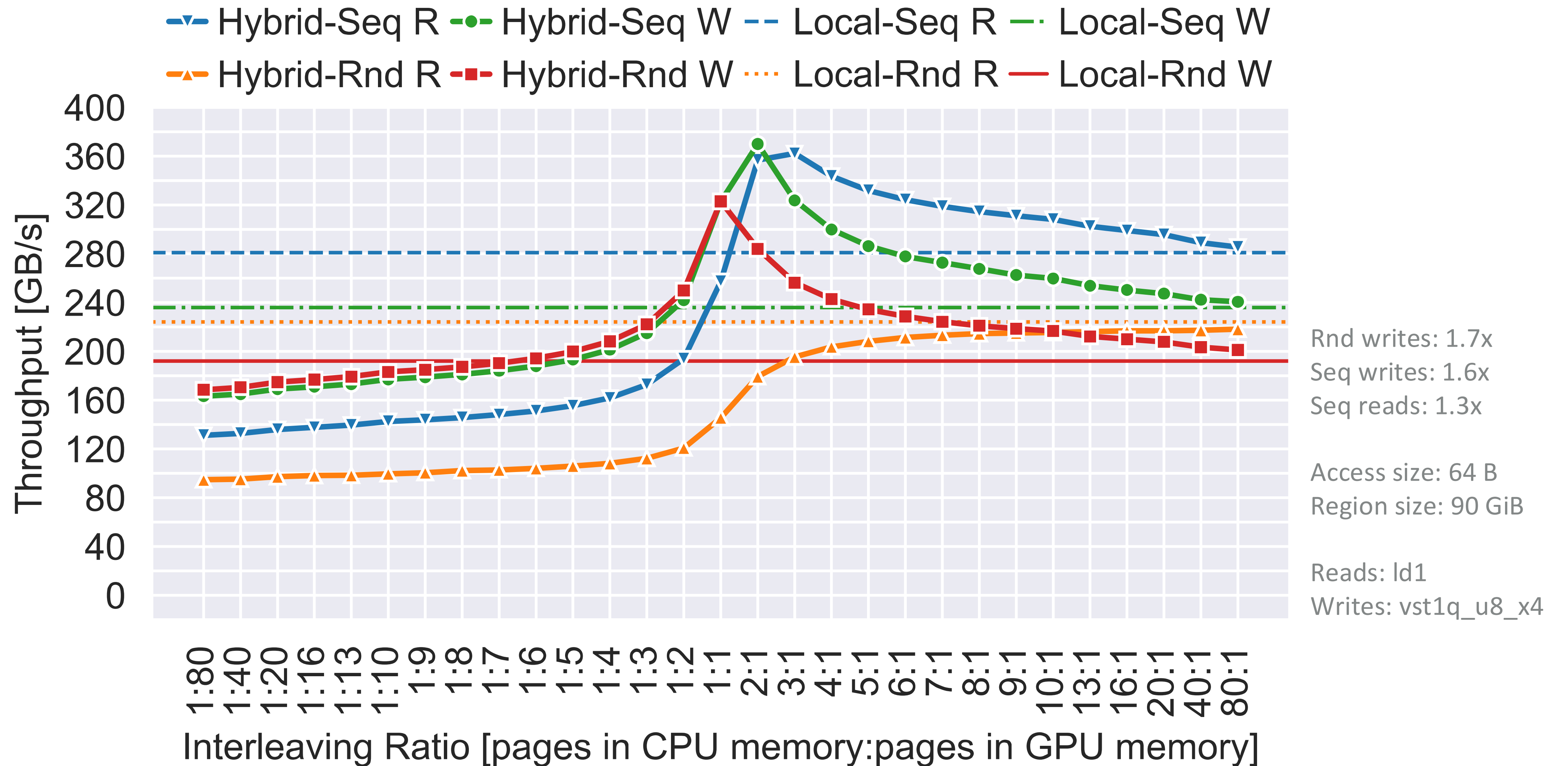
Data Engineering Systems Group @ HPI

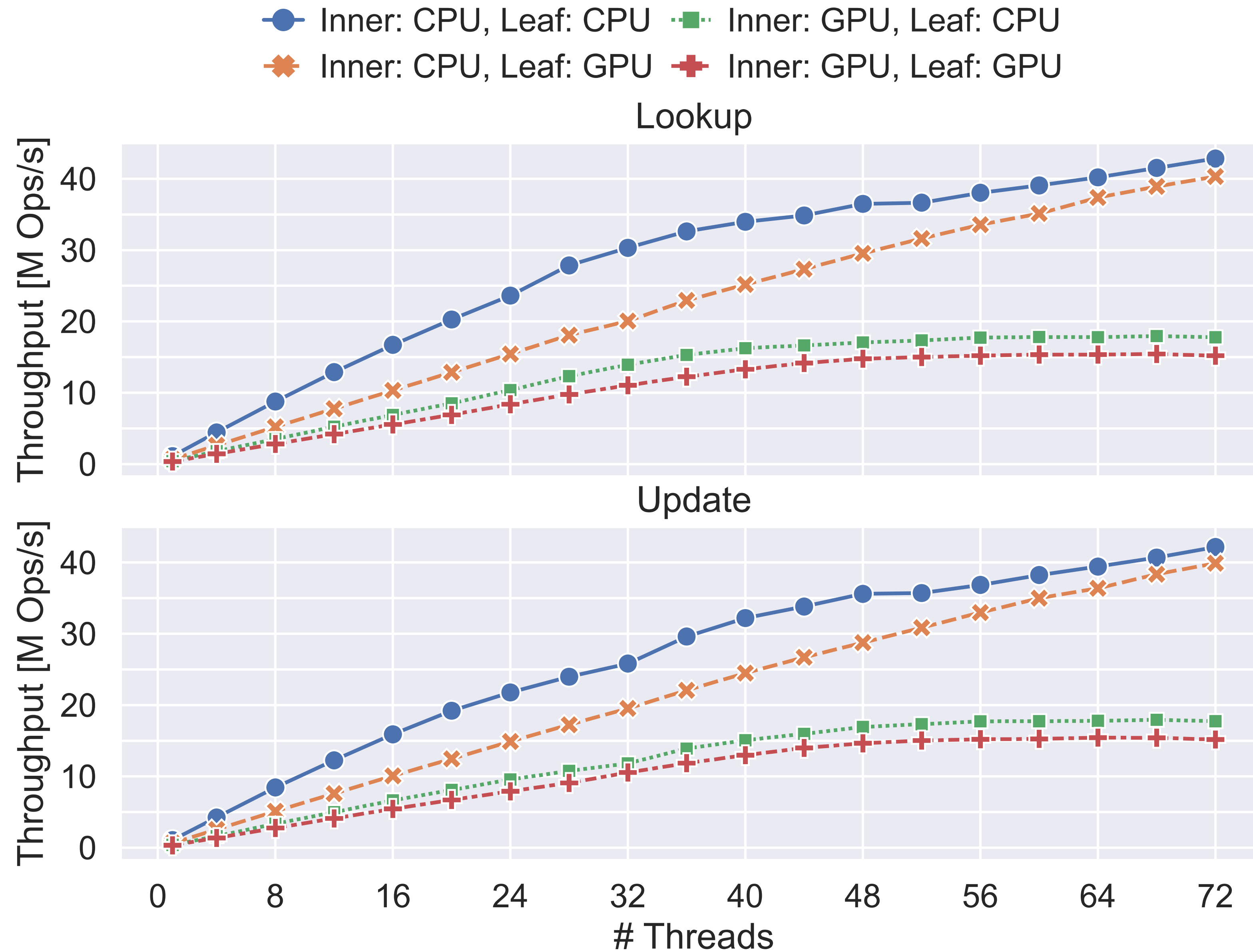
🔍 hpi.de/rabl

✉️ marcel.weisgut@hpi.de



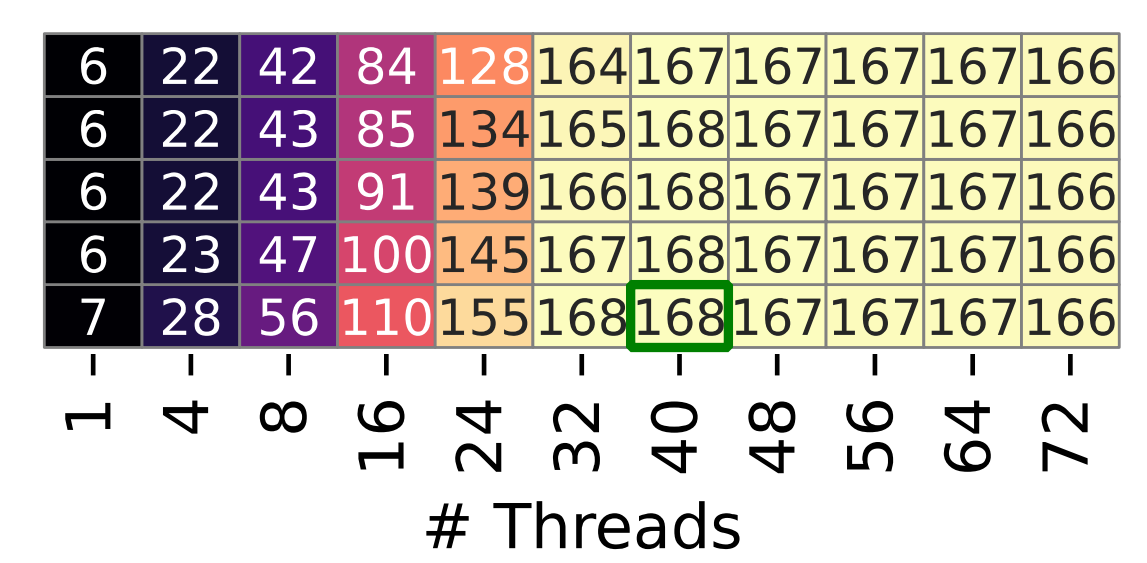
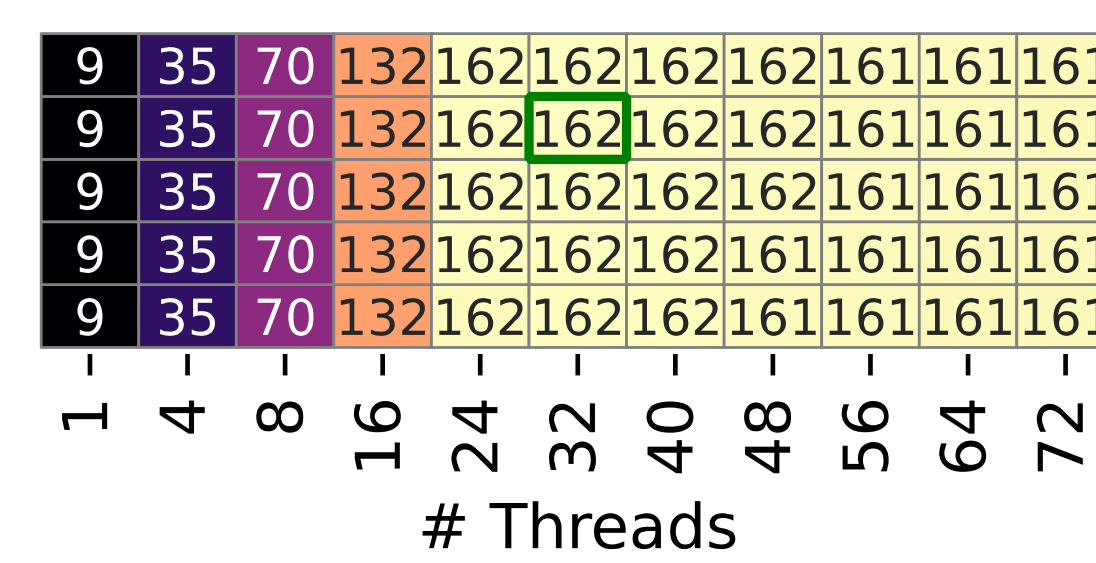
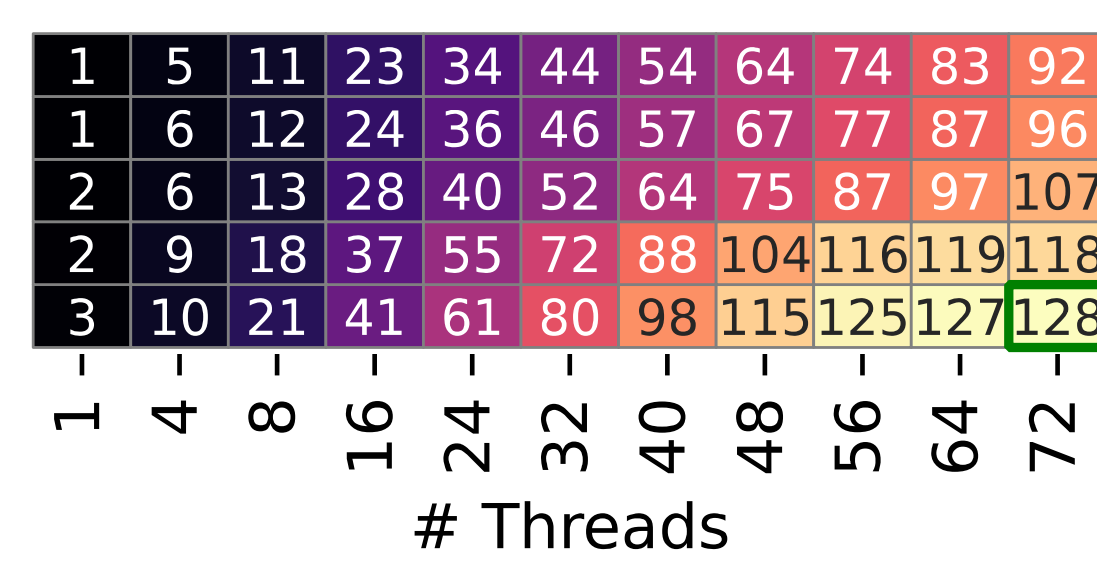
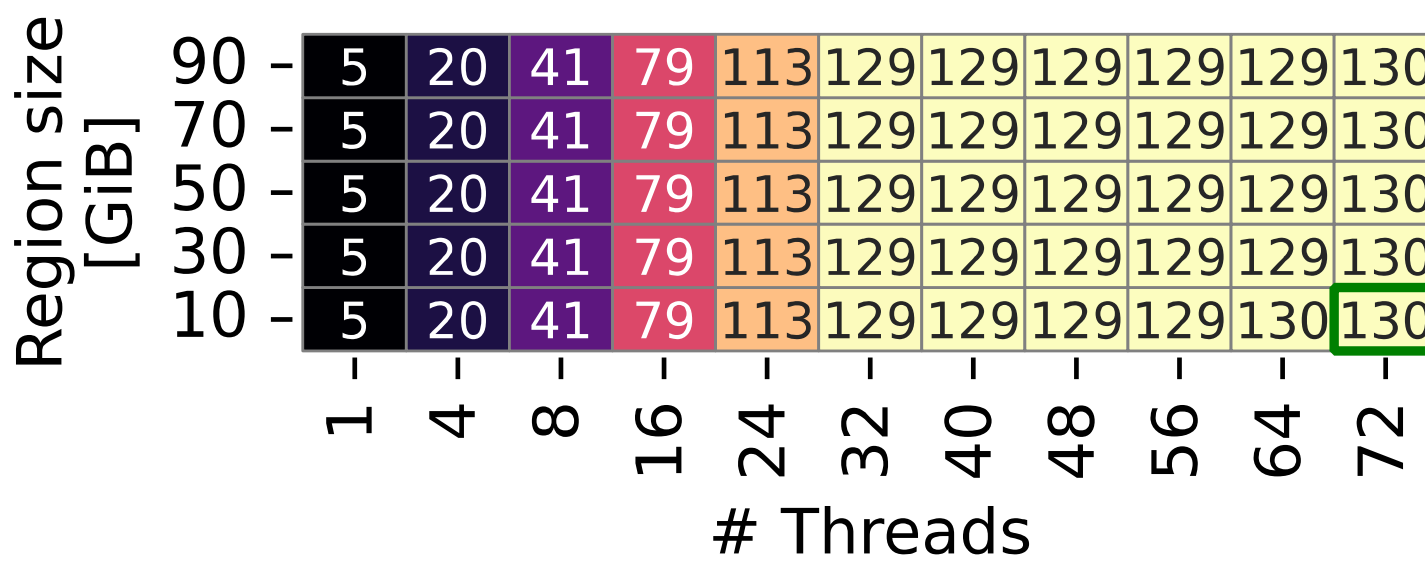
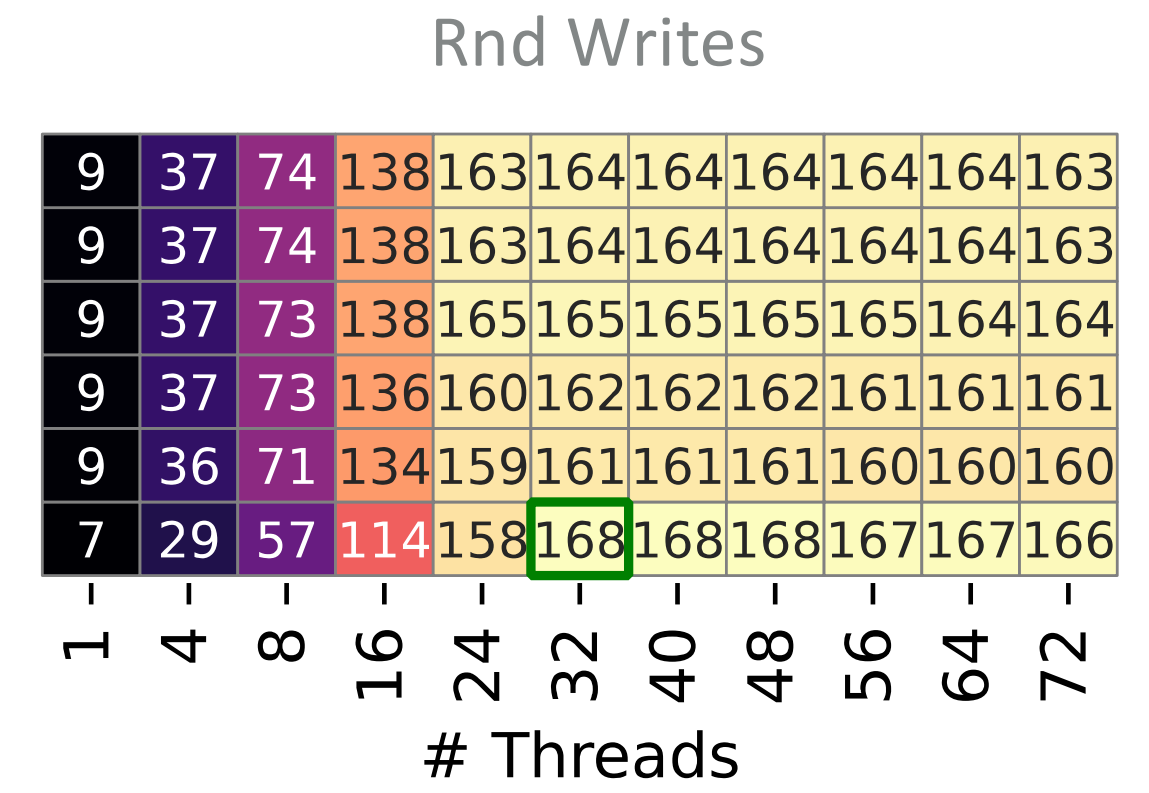
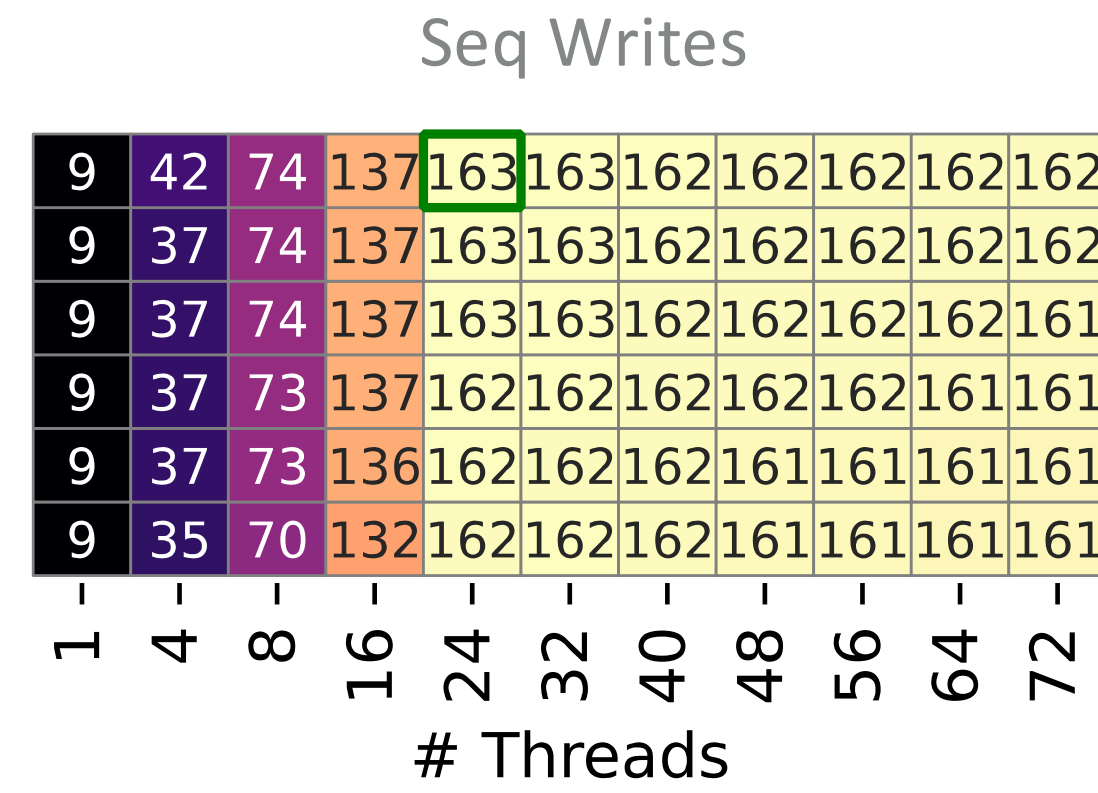
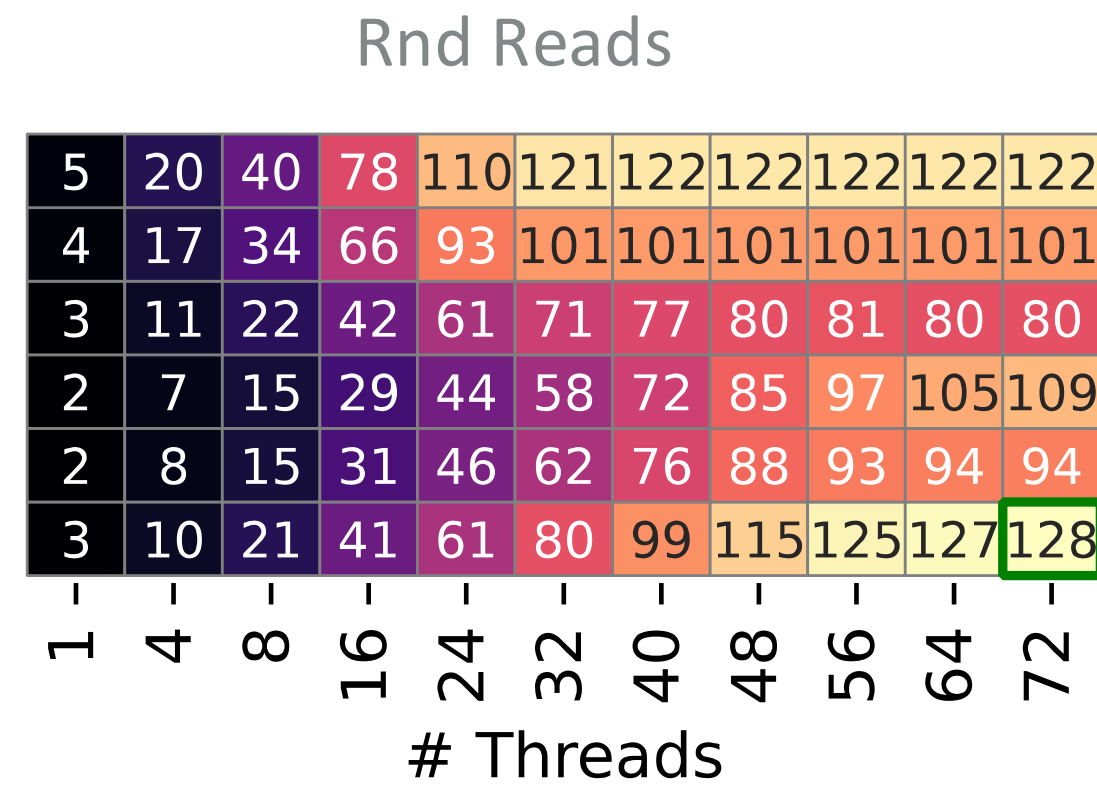
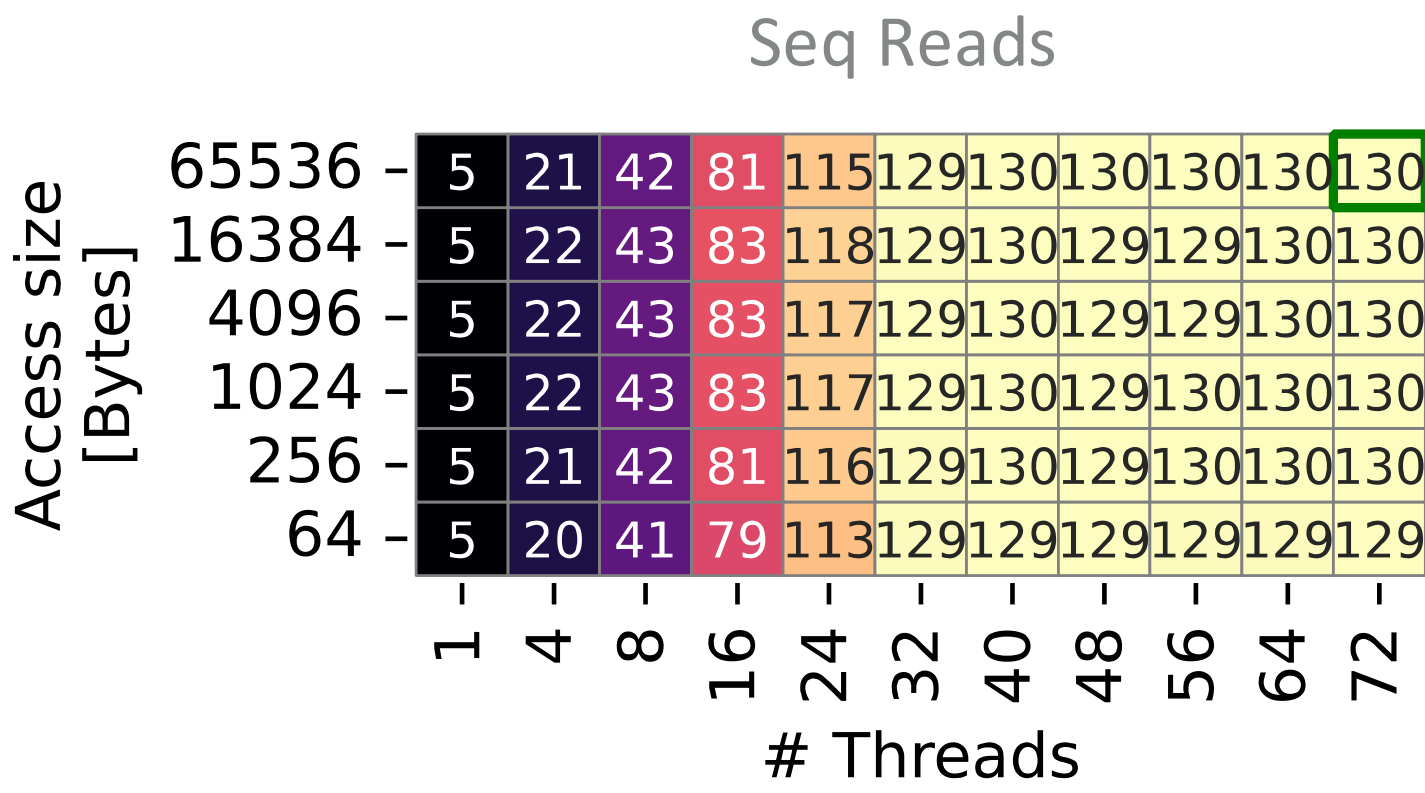
Bandwidth Expansion





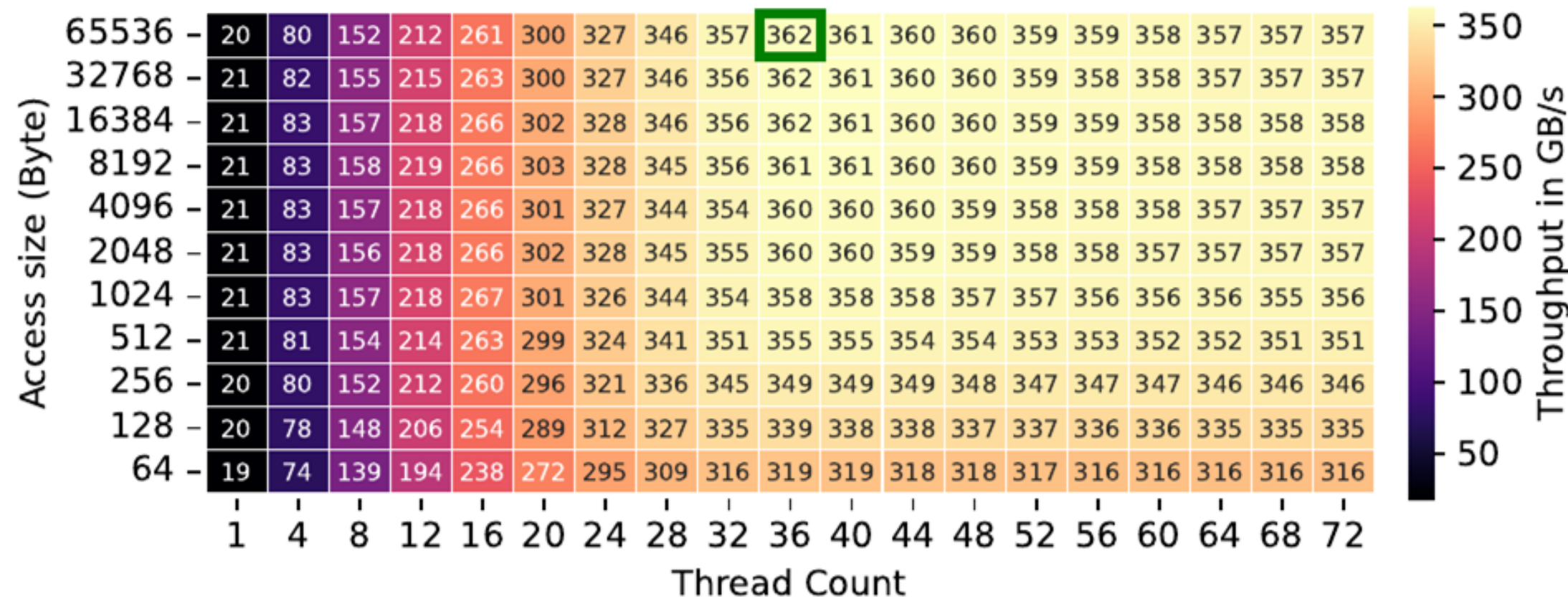
CXL-Bench

- ▶ Support for fundamental and chains of memory access operations
- ▶ Memory region preparation
 - ▶ Pin memory region to NUMA node(s) (`mbind`) – such as CXL and (modern) GPU memory
 - ▶ Multiple regions, partitions, and NUMA nodes supported
 - ▶ Round-robin & weighted page interleaving supported
- ▶ Task execution
 - ▶ Thread pinning to set of cores (`pthread_setaffinity_np`)
 - ▶ Memory access via scalar & SIMD instructions (using compiler vector intrinsics)
- ▶ Configurable with YAML files

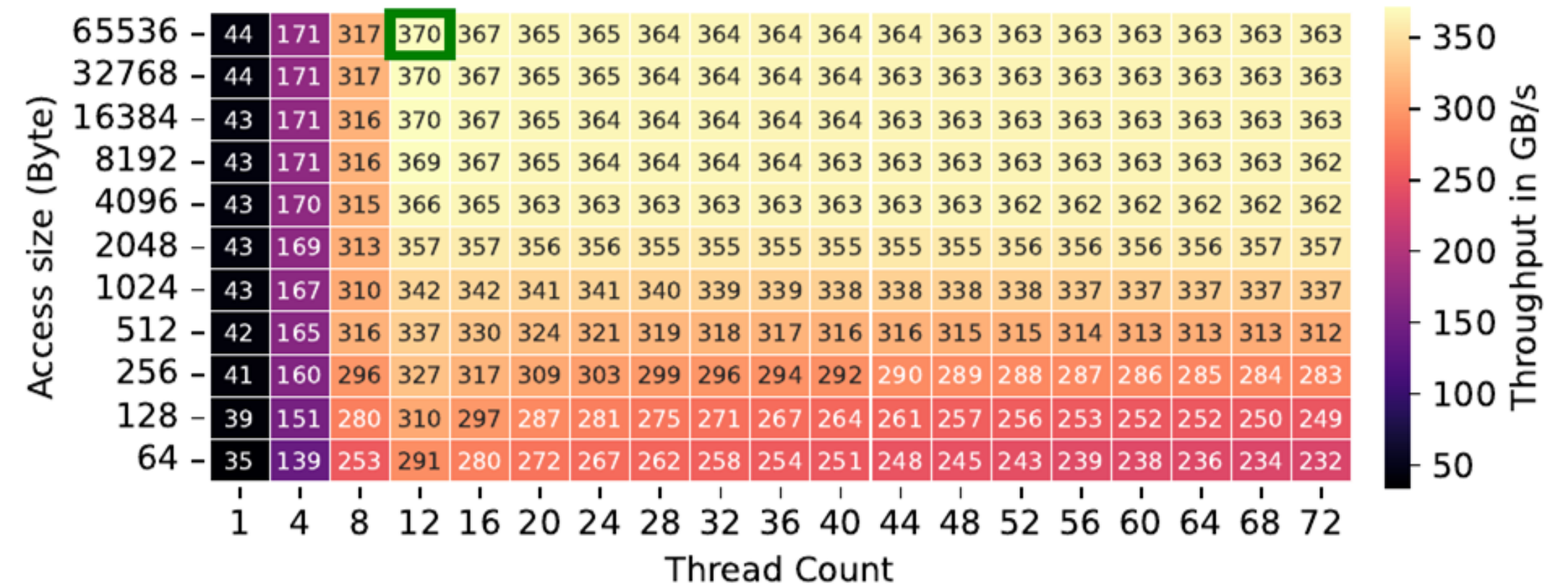


LOCAL MEMORY ACCESS THROUGHPUT

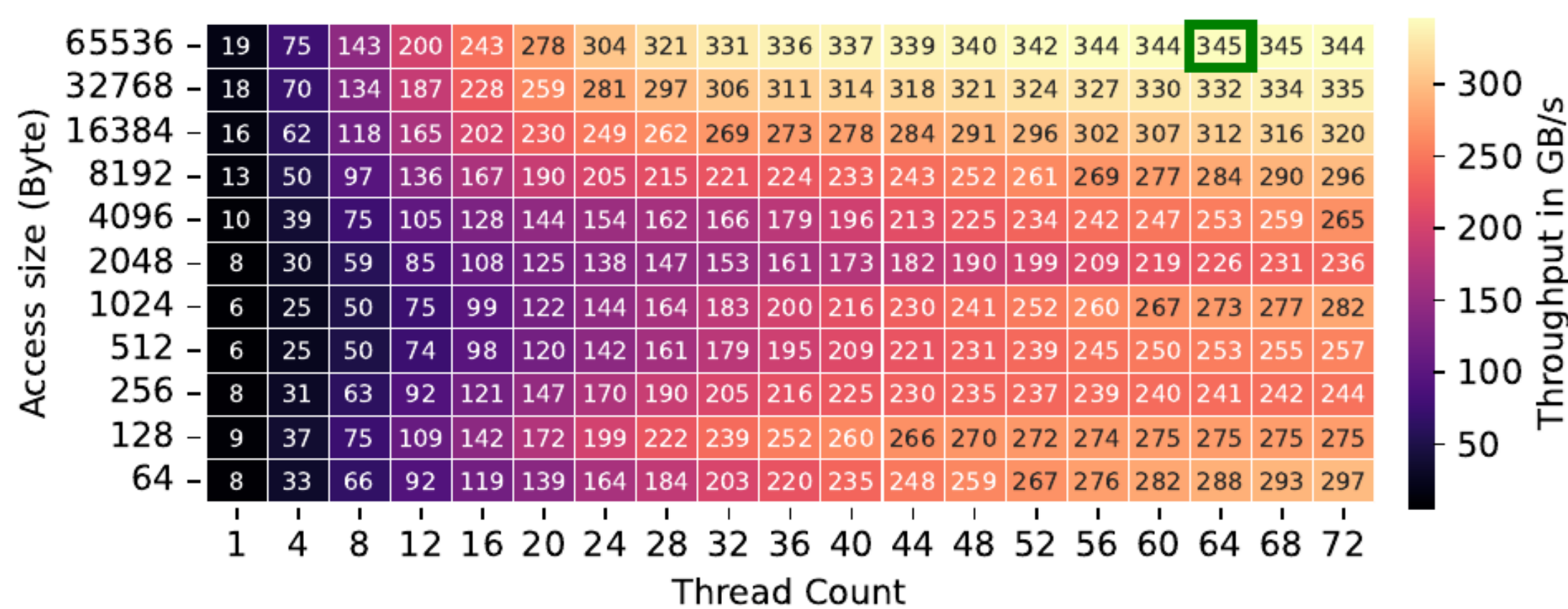
Sequential Reads



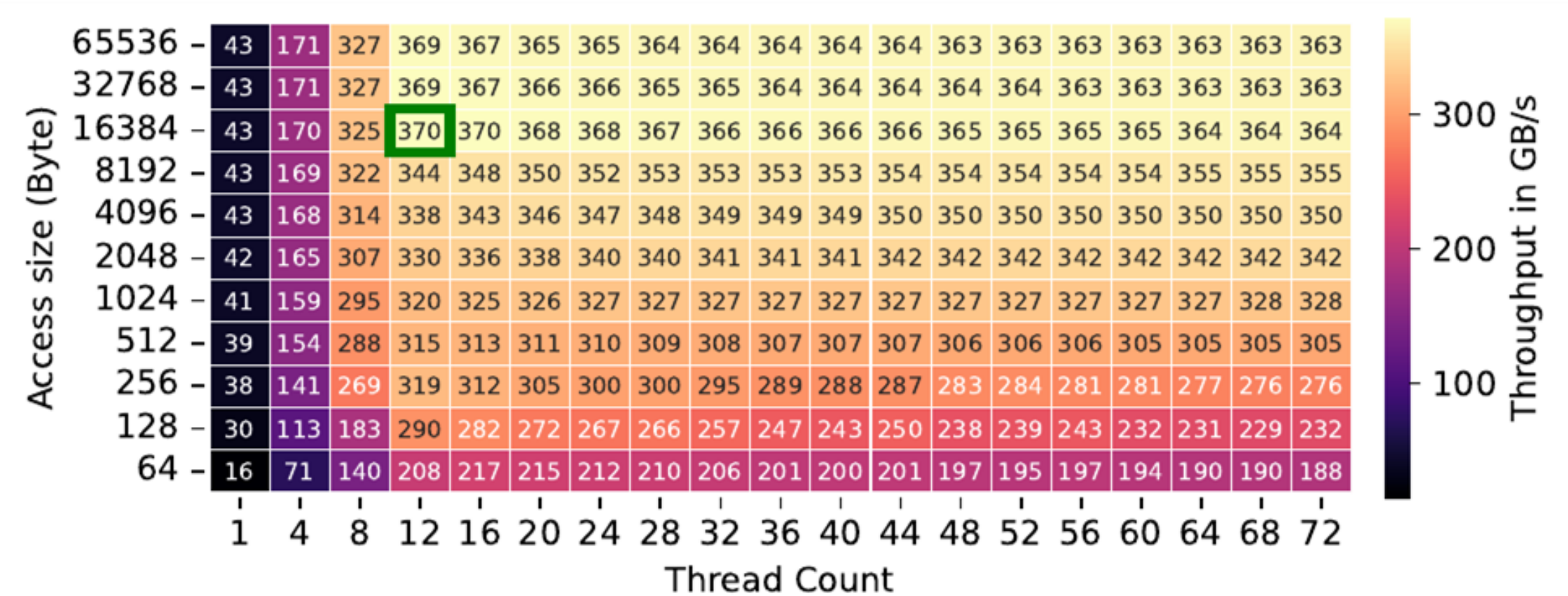
Sequential Writes

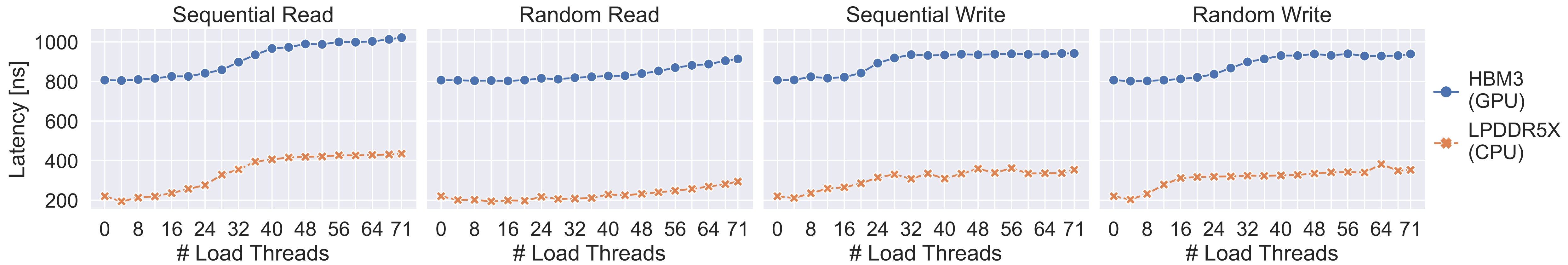


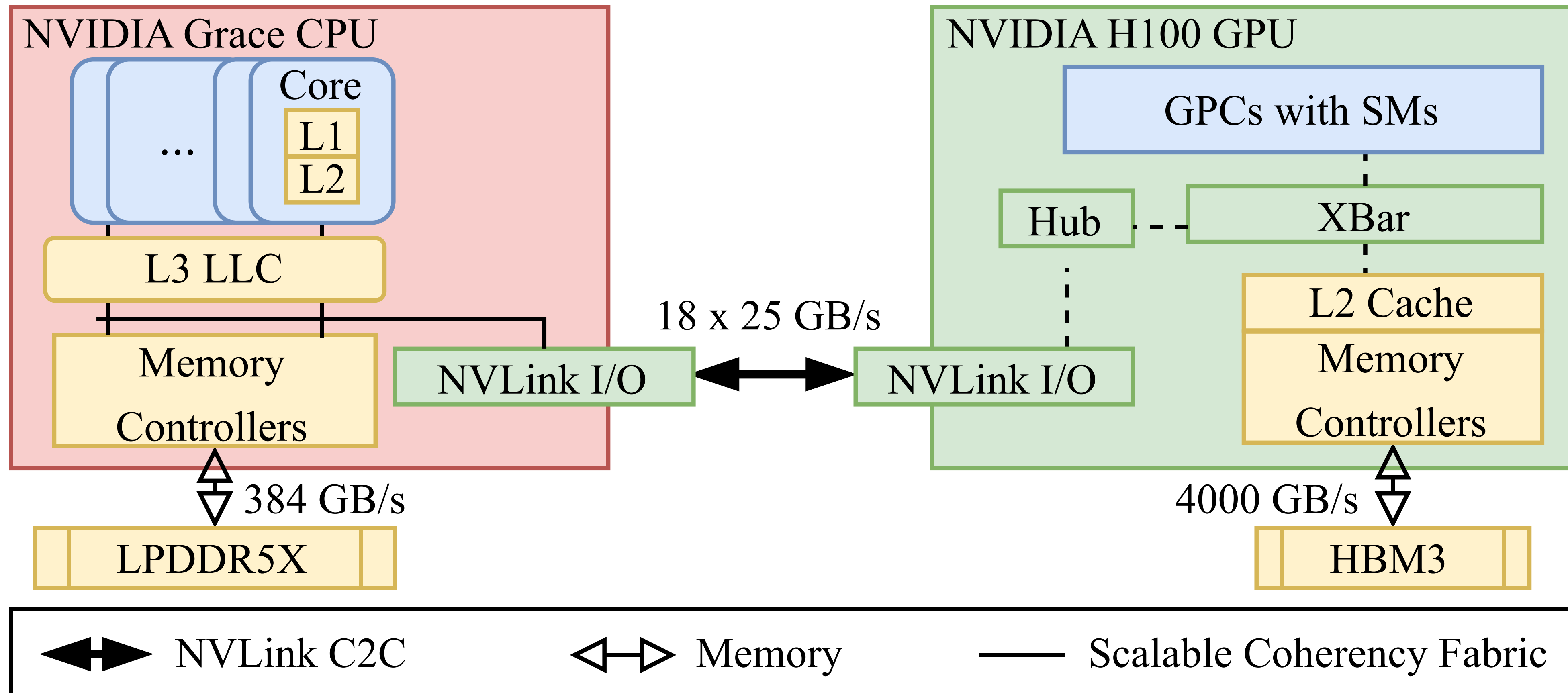
Random Reads



Random Writes







NVLink Evolution

Version	Data rate per link [GB/s]	# Lanes per link	# Links	Theoretical bandwidth	Architecture
1	20	8	4	80	Pascal
2	25	8	6	150	Volta
3	25	4	12	300	Ampere
4 / C2C	25	4	18	450	Hopper

NVLINK C2C VS CPU INTERCONNECTS

