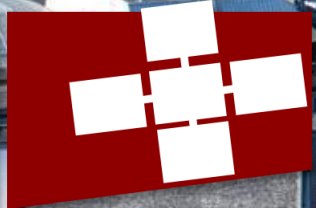
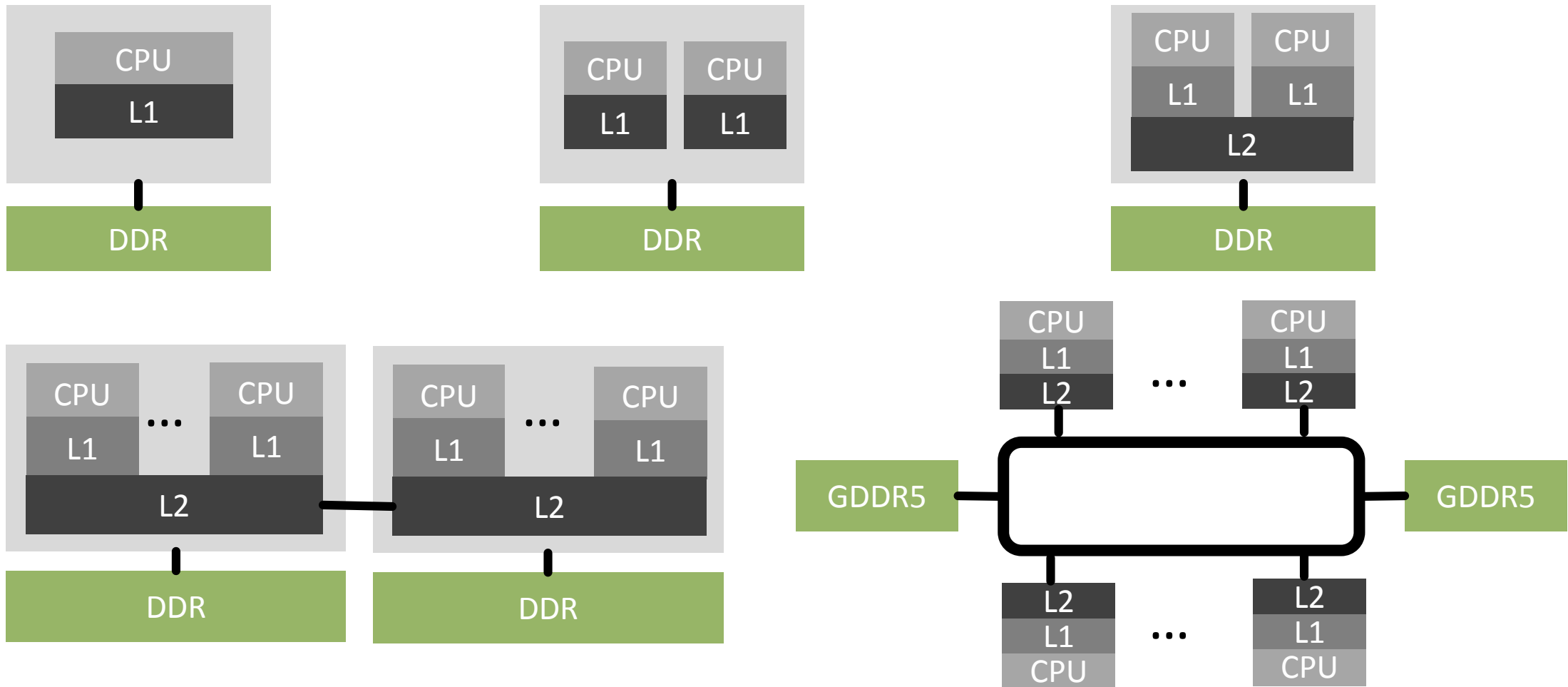


SABELA RAMOS, TORSTEN HOEFLER

Capability Models for Manycore Memory Systems: A Case-Study with Xeon Phi KNL



Microarchitectures are becoming more and more complex



How to optimize codes for these complex architectures?

- **Performance engineering:** “encompasses the set of roles, skills, activities, practices, tools, and deliverables applied at every phase of the systems development life cycle which ensures that a solution will be designed, implemented, and operationally supported to meet the non-functional requirements for performance (such as throughput, latency, or memory usage).”
- **Manually** profile codes and **tune** them to the given architecture
 - Requires highly-skilled performance engineers
 - Need familiarity with
 - NUMA (topology, bandwidths etc.)*
 - Caches (associativity, sizes etc.)*
 - Microarchitecture (number of outstanding loads etc.)*

Trust me, I'm an engineer!

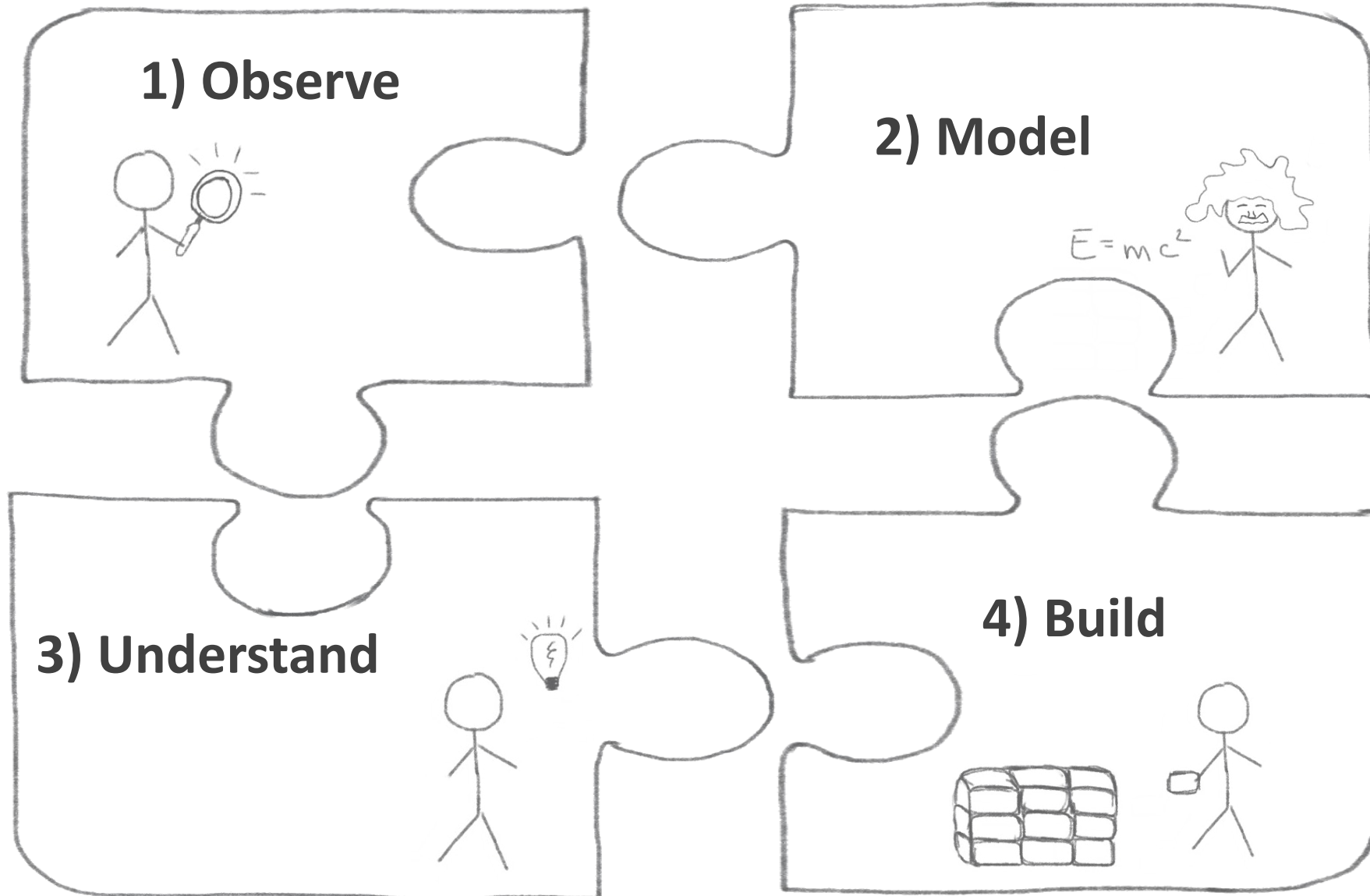


<title>code ninja</title>

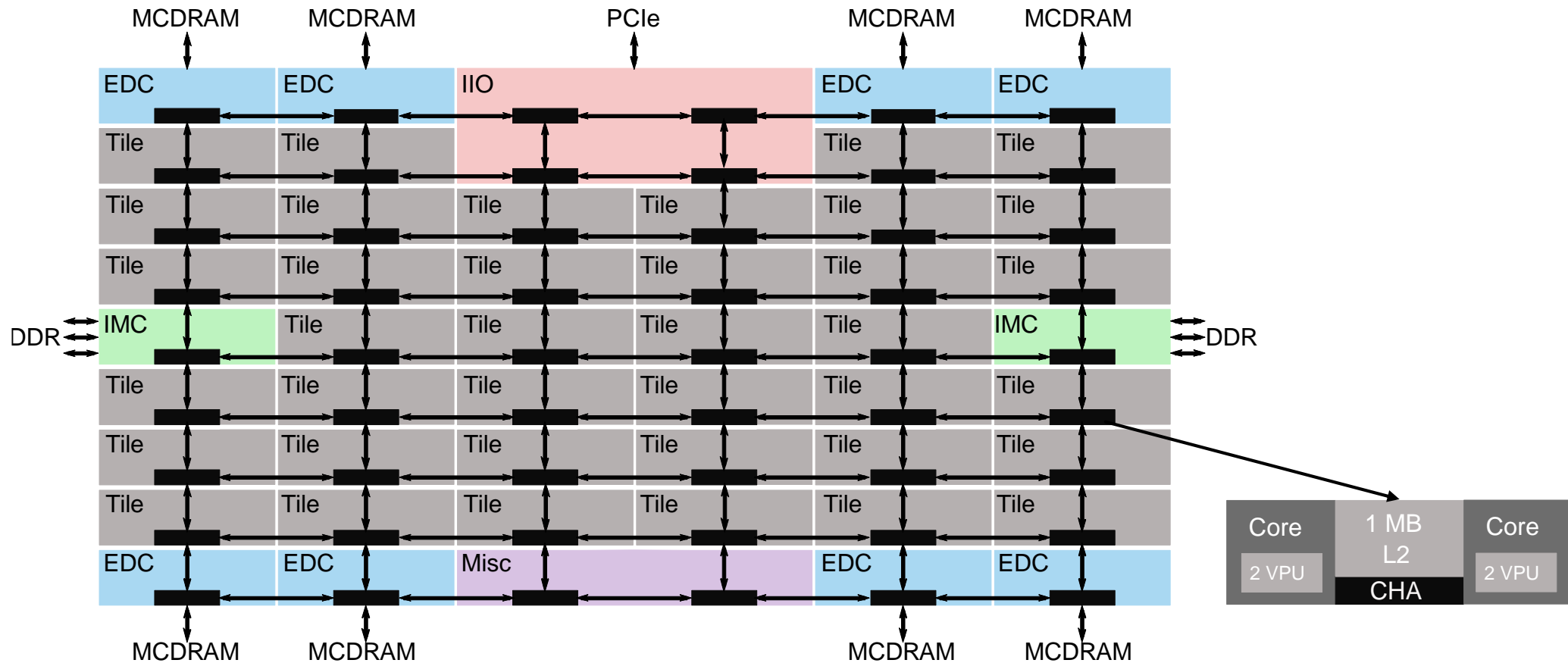
An engineering example – Tacoma Narrows Bridge



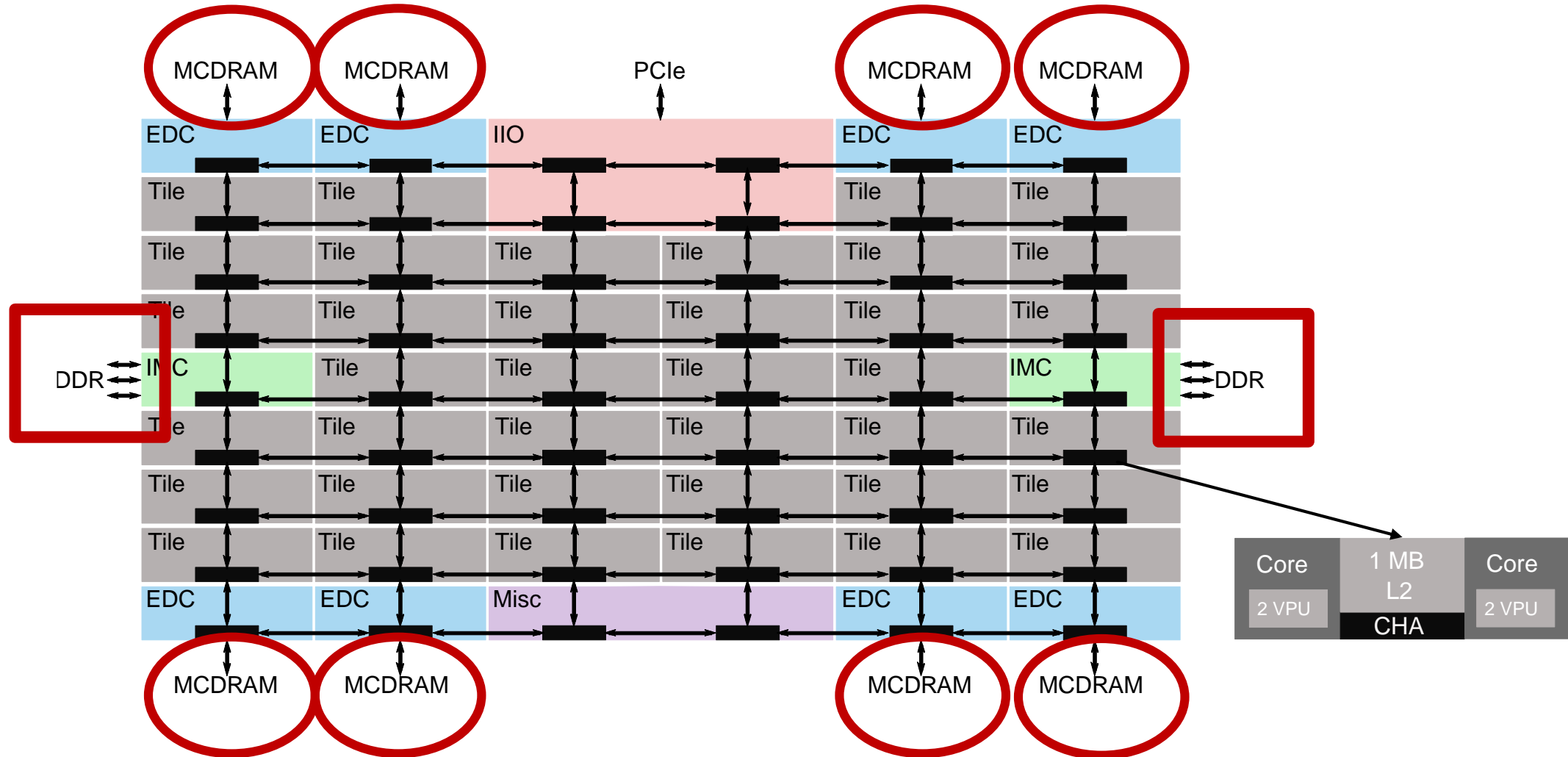
Scientific **Performance** Engineering



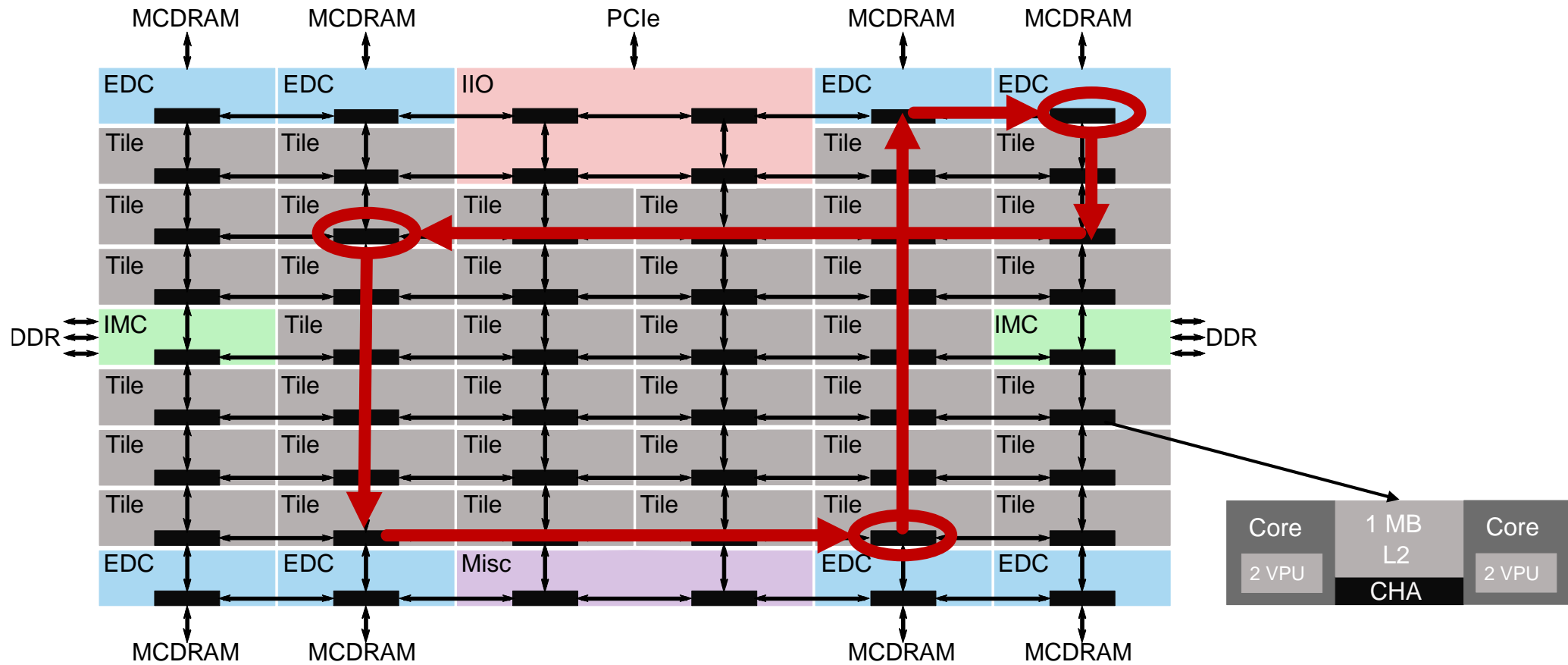
Modeling by example: KNL Architecture (mesh)



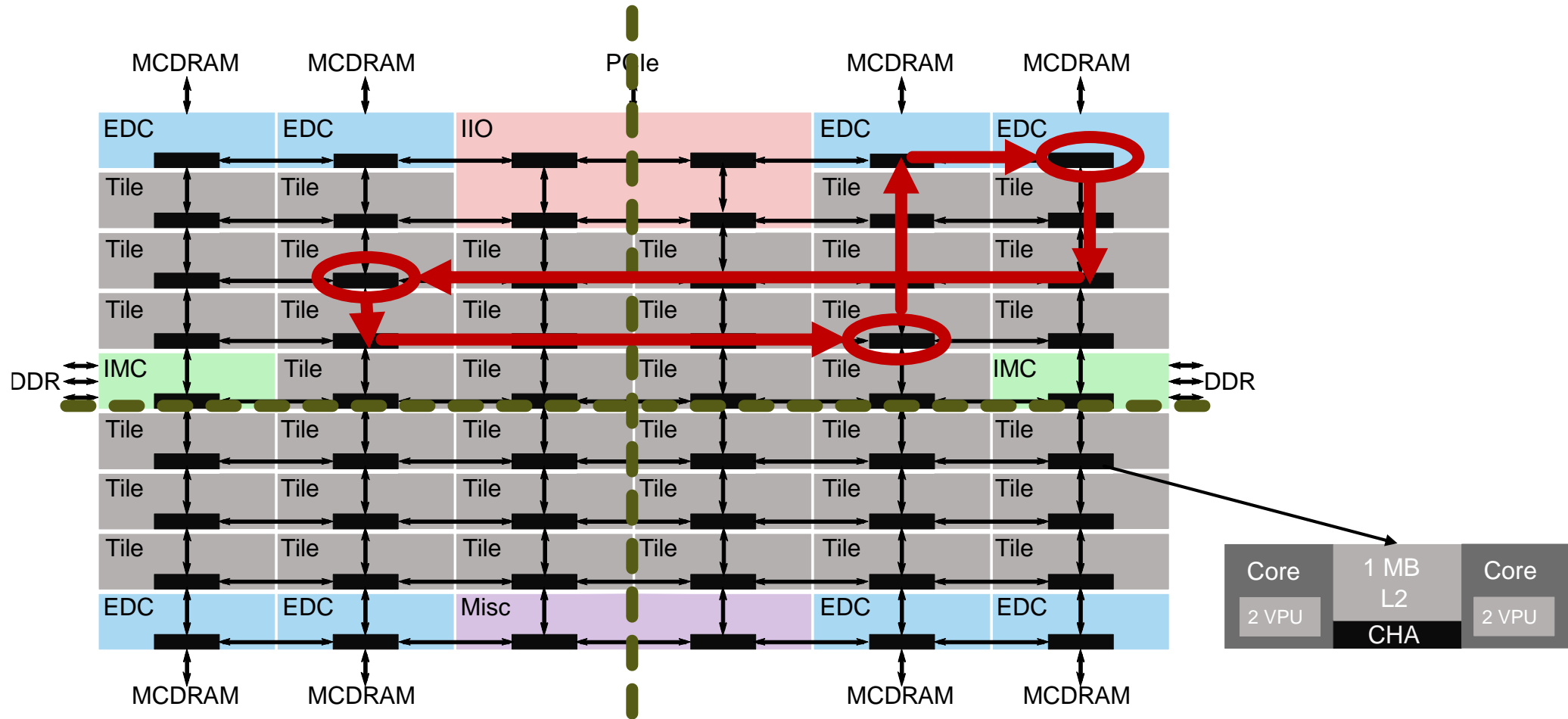
KNL Architecture (memory: Flat & Cache)



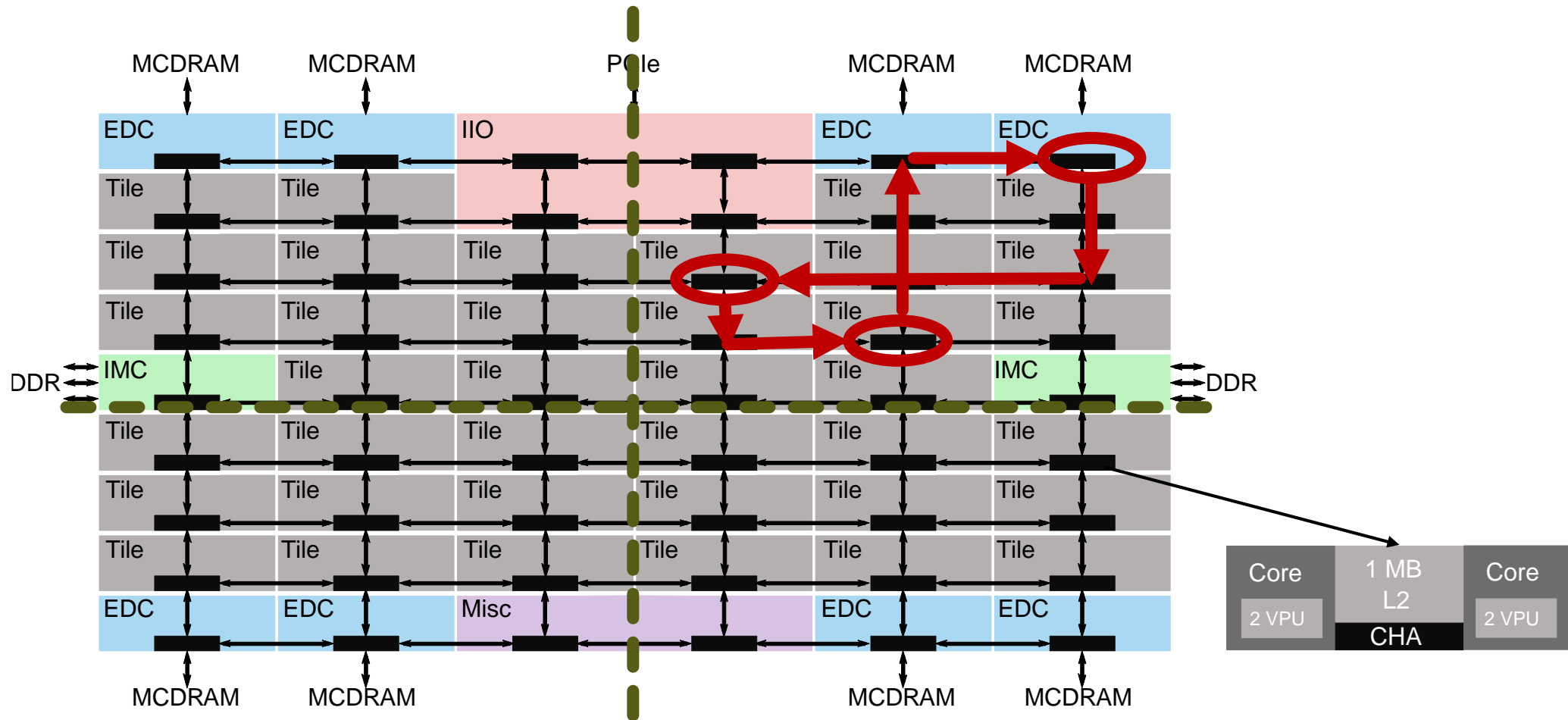
KNL Architecture (all to all mode)



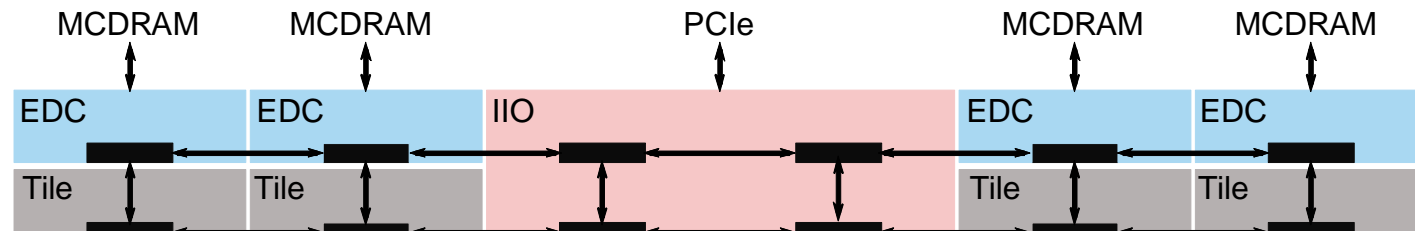
KNL Architecture (Quadrant or Hemisphere)



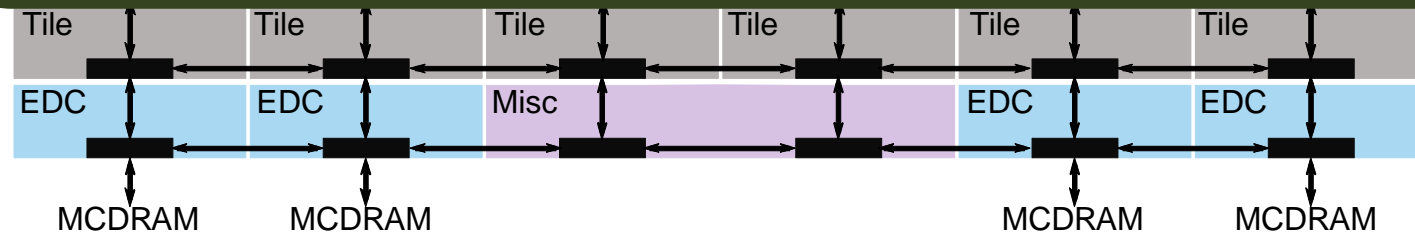
KNL Architecture (SNC-4 or SNC-2)



KNL Architecture



How much does this all matter?
 What is the real cost of accessing cache?
 What is the cost of accessing memory?



Step 1: Understand core-to-core transfers – MESIF cache coherence

Write back overhead

Location: only 5-15% difference

Contention effects?

That is curious!

Step 2: Understand core-to-memory transfers – DRAM and MCDRAM

MCDRAM 20% slower!

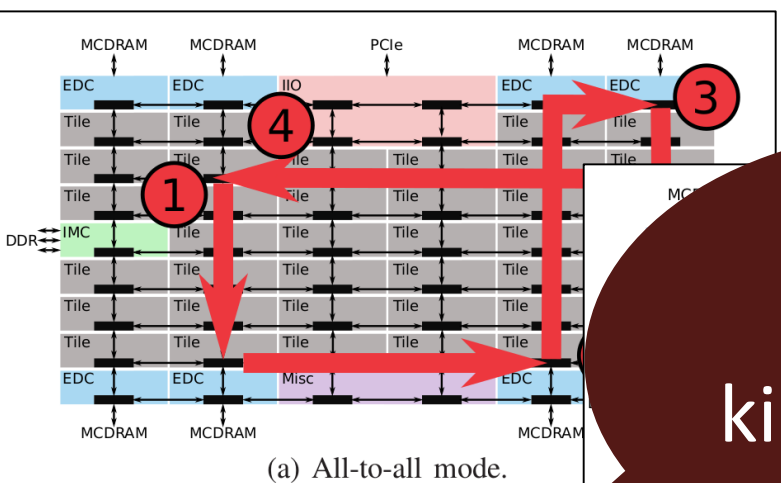
MCDRAM 4-6x faster!

Need to read and write for full bandwidth

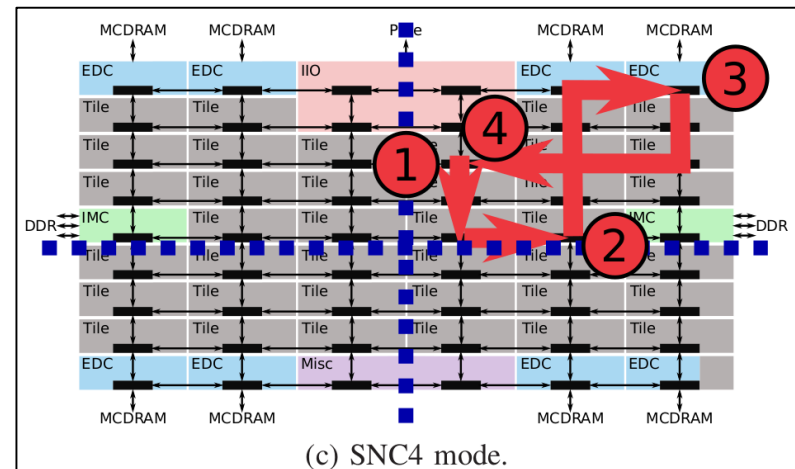
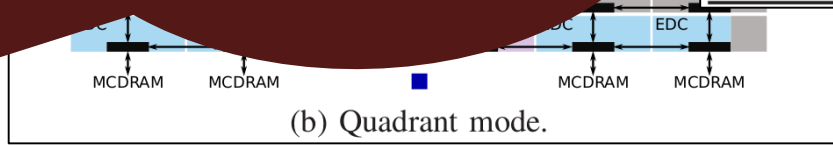
Cache mode >20% slower

Bandwidth suffers a bit

Performance engineers optimize your code!



Are you kidding me?



		Software NUMA		Software UMA		
		SNC4	SNC2	QUAD	HEM	A2A
Latency [ns] (Copy/BenchIT)	Local (L1)	3.8	3.8	3.8	3.8	3.8
	Tile (L2)	34 (M) 17 (E) 14 (S,F)	34 (M) 18 (E) 14 (S,F)	34 (M) 18 (E) 14 (S,F)	34 (M) 18 (E) 14 (S,F)	34 (M) 18 (E) 14 (S,F)
	Remote	107-122 (M) 98-114 (E) 96-118 (S,F)	111-125 (M) 104-117 (E) 104-118 (S,F)	119 (M) 116 (E) 107-117 (S,F)	120 (M) 116 (E) 107-117 (S,F)	122 (M) 116 (E) 109-117 (S,F)
	Bandwidth [GB/s] (Read)	2.5	2.5	2.5	2.5	2.5
Bandwidth [GB/s] (Copy)	Tile	6.7 (M) 7.6 (E)	6.7 (M) 6.7 (E)	7.5 (M) 9.2 (E)	7.4 (M) 9.2 (E)	7.5 (M) 9.2 (E)
	Remote	7.7	6.7	7.5	7.5	7.5
Congestion (P2P pairs)				None		
Contention [ns] (1:N copy)		α	200	200	200	200
Linear, $\mathcal{T}_C(N) = \alpha + \beta \cdot N$		β	34	34	34	34

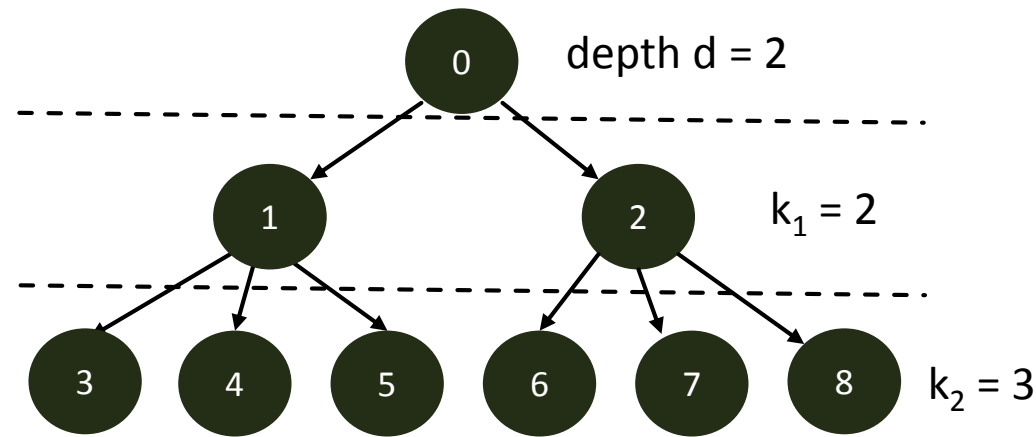
		Software NUMA		Software UMA			
		SNC4	SNC2	QUAD	HEM	A2A	
Flat Mode	Latency [ns] (BenchIT)	DRAM	130-140	134-146	140	140	139
		MCDRAM	160-175	160-170	167	167	168
	Bandwidth [GB/s] (Copy NT / STREAM Copy)	DRAM	69 / 77	69 / 77	70 / 77	71 / 77	71 / 77
		MCDRAM	342 / 418	333 / 388	333 / 415	315 / 372	306 / 359
	Bandwidth [GB/s] (Read)	DRAM	71	71	77	77	77
		MCDRAM	243	288	314	314	314
	Bandwidth [GB/s] (Write)	DRAM	33	34	36	36	36
		MCDRAM	147	163	171	165	161
	Bandwidth [GB/s] (Triad NT / STREAM Triad)	DRAM	71 / 82	71 / 82	74 / 82	73 / 82	73 / 82
		MCDRAM	371 / 448	347 / 441	340 / 441	332 / 434	325 / 427
Cache Mode	Latency [ns] (BenchIT)		158-178	161-171	166	168	172
	Bandwidth [GB/s] (Copy NT / STREAM Copy)		150 / 252	130 / 252	175 / 255	134 / 237	132 / 233
	Bandwidth [GB/s] (Read)		87	95	124	128	118
			56	56	72	72	68
	Bandwidth [GB/s] (Triad NT / STREAM Triad)		296 / 292	246 / 294	296 / 309	273 / 274	264 / 269



<title>code ninja</title>

A principled approach to designing cache-to-cache broadcast algorithms

Multi-ary tree example



Tree depth

Level size

$$\mathcal{T}_{tree} = \sum_{i=1}^d \mathcal{T}_C(k_i) = \sum_{i=1}^d (c \cdot k_i + b)$$

Tree cost

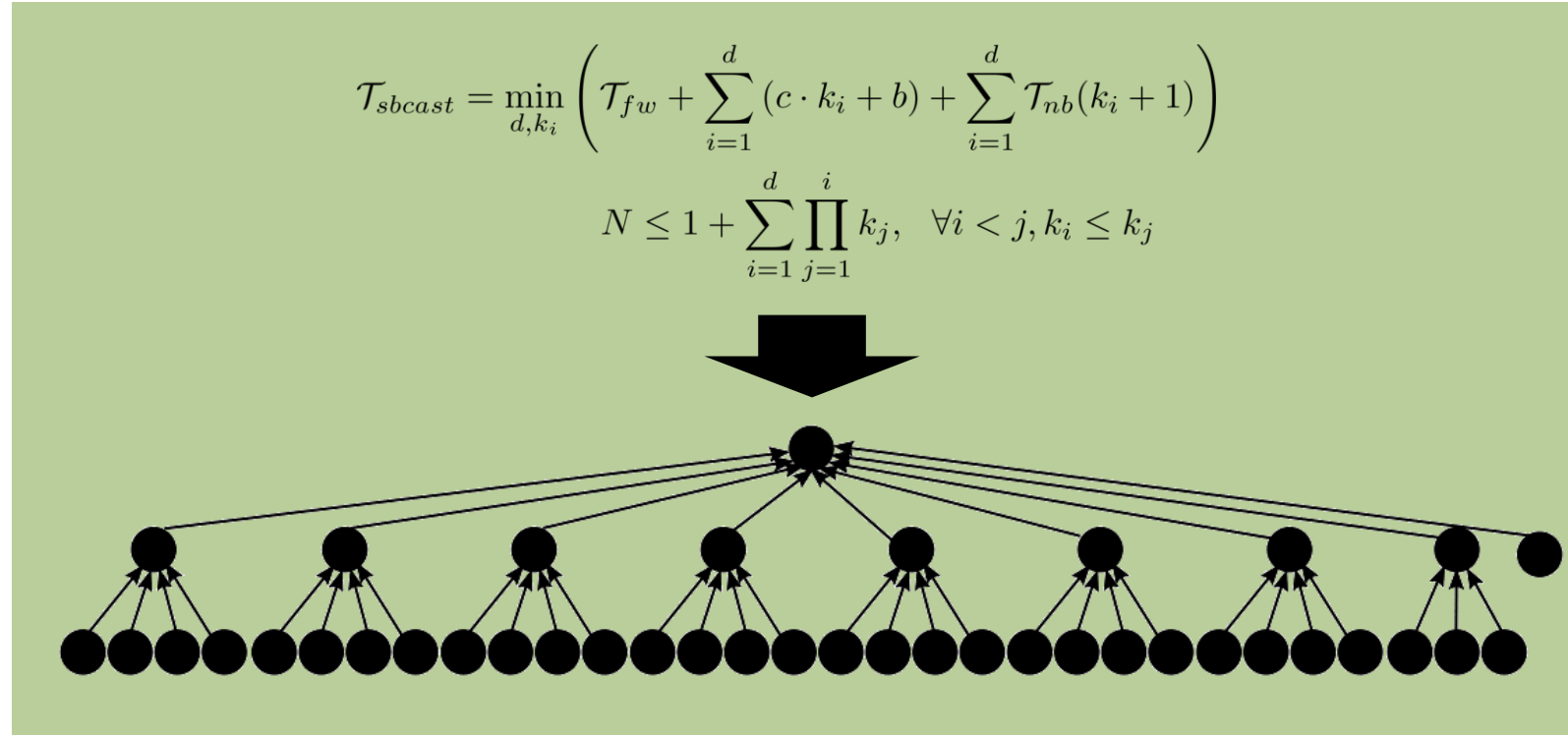
$$= \sum_{i=1}^d (R_R + R_L + c \cdot (k_i - 1))$$

$$\mathcal{T}_{sbcast} = \min_{d, k_i} \left(\mathcal{T}_{fw} + \sum_{i=1}^d (c \cdot k_i + b) + \sum_{i=1}^d \mathcal{T}_{nb}(k_i + 1) \right)$$

Reached threads

$$\dots N \leq 1 + \sum_{i=1}^d \prod_{j=1}^i k_j, \quad \forall i < j, k_i \leq k_j$$

Model-driven performance engineering for broadcast

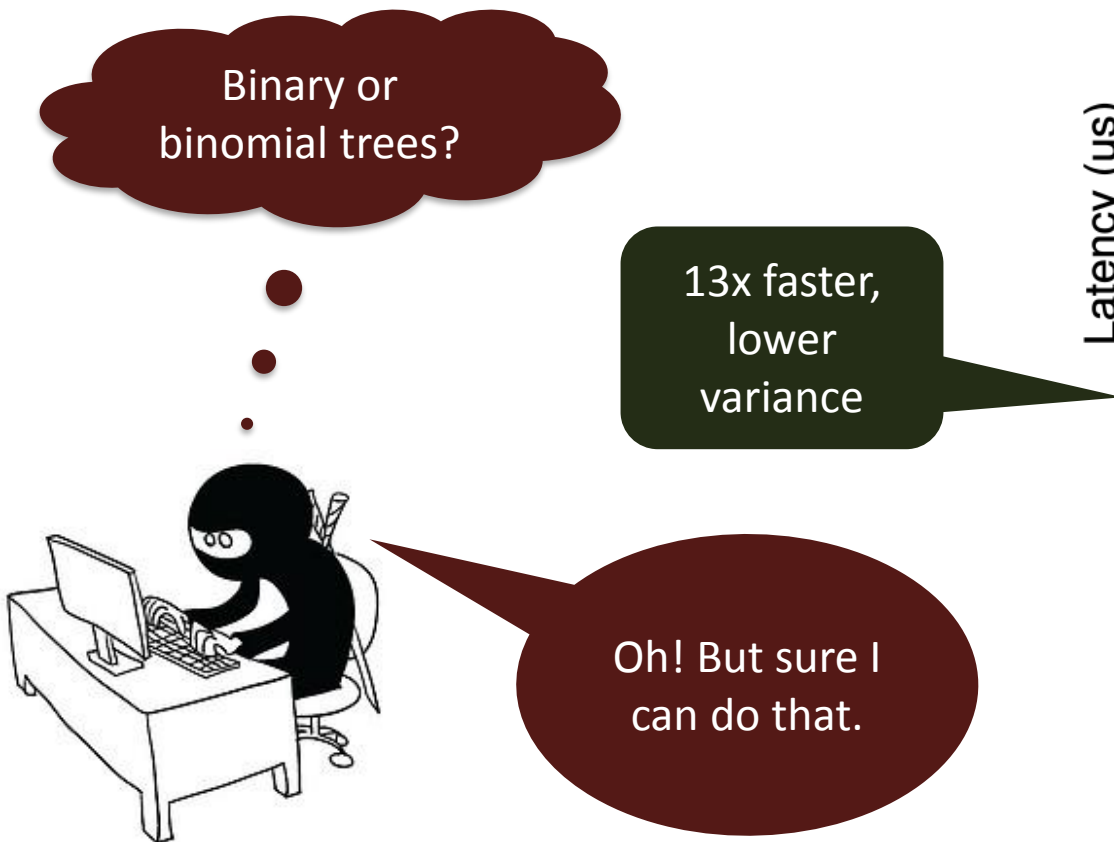


Binary or binomial trees?

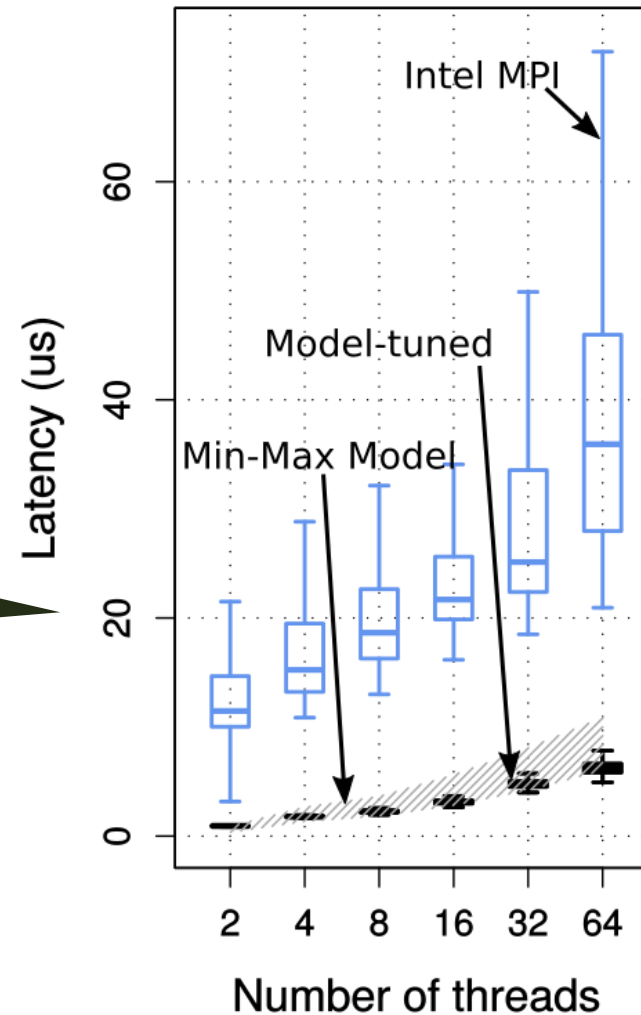


Oh! But sure I can do that.

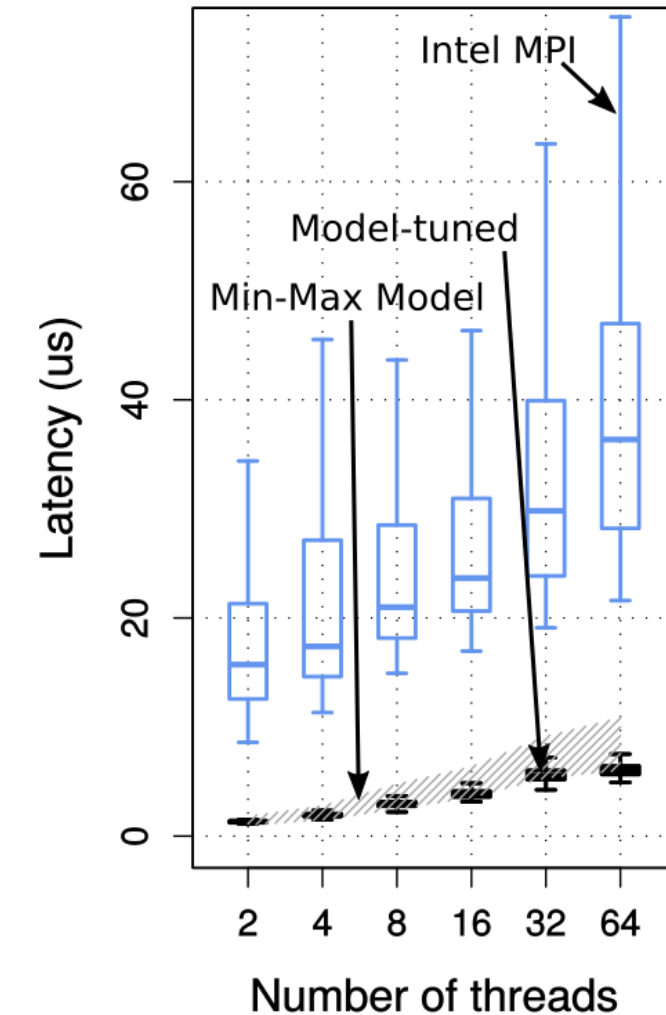
Model-driven performance engineering for broadcast



<title>code ninja</title>

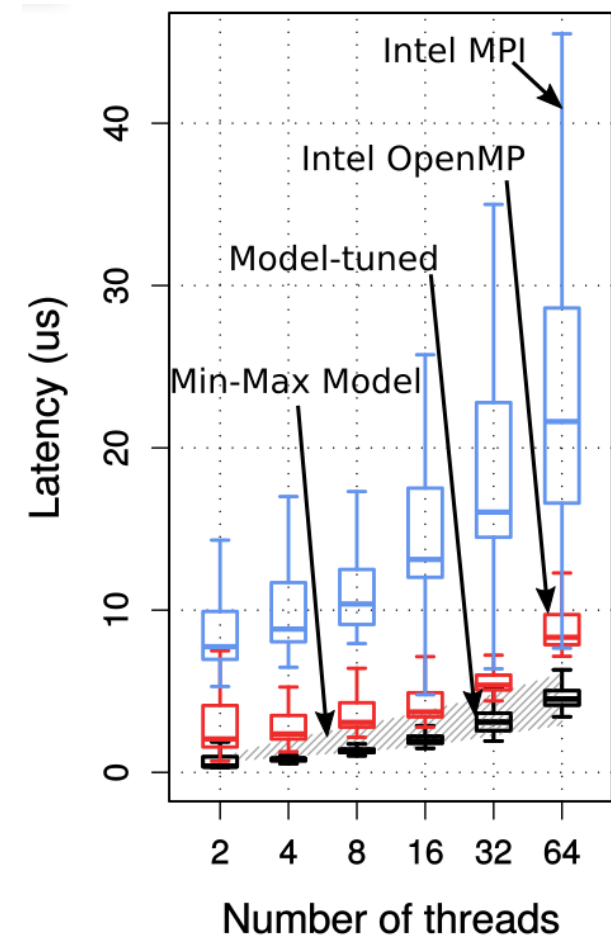


(a) Filling Tiles.

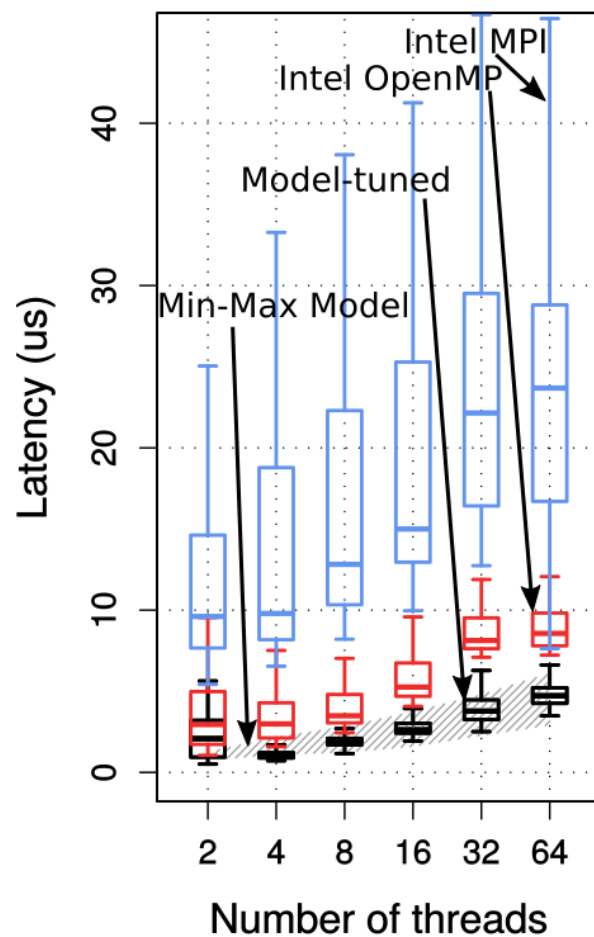


(b) Scatter.

Easy to generalize to similar algorithms

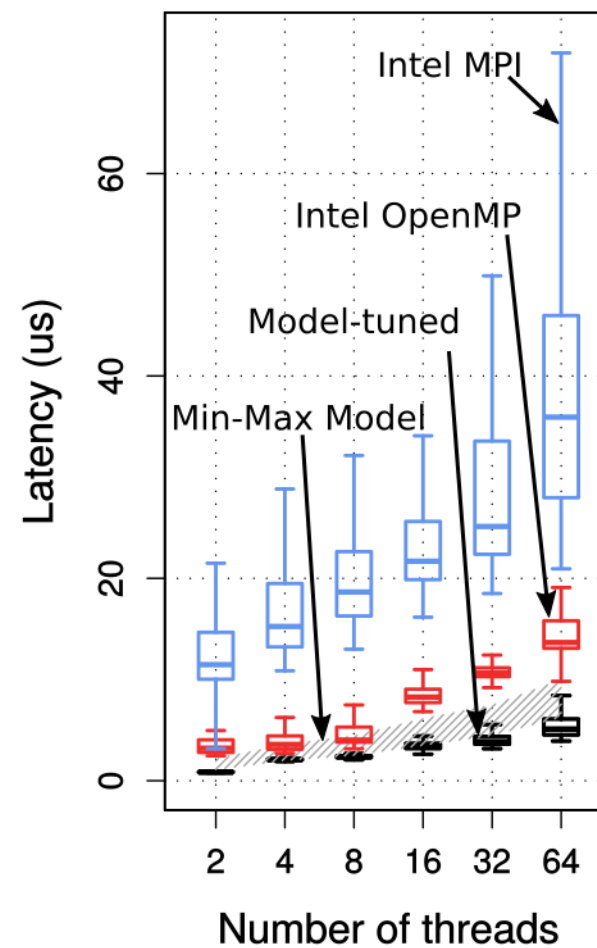


(a) Filling Tiles.

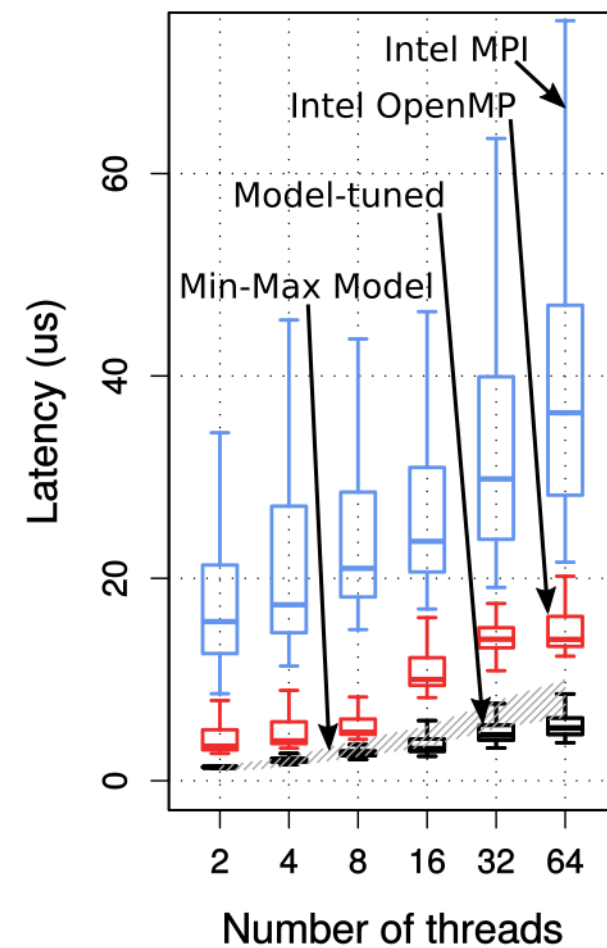


(b) Scatter.

Barrier (7x faster than OpenMP)



(a) Filling Tiles.



(b) Scatter.

Reduce (5x faster than OpenMP)

Sorting in complex memories

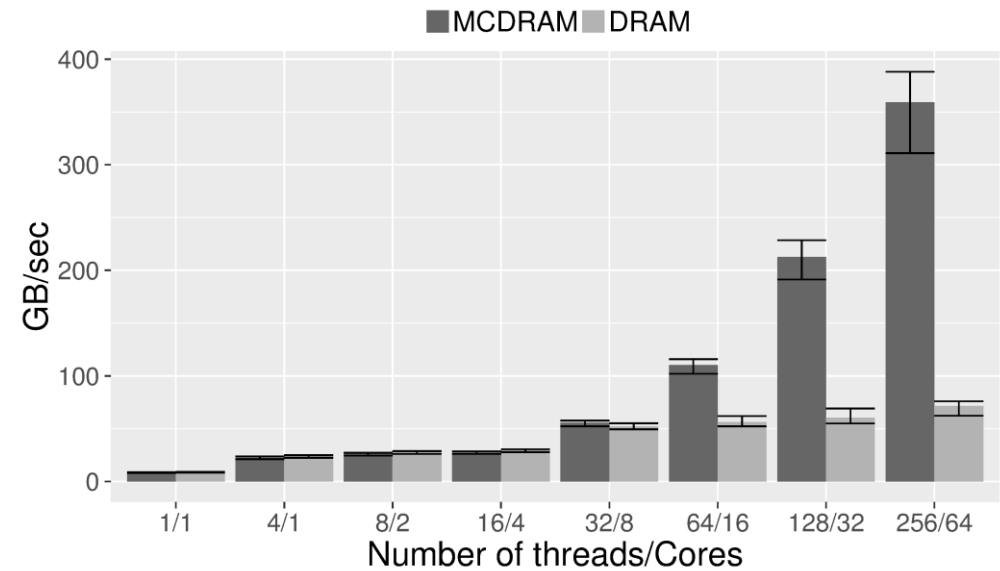
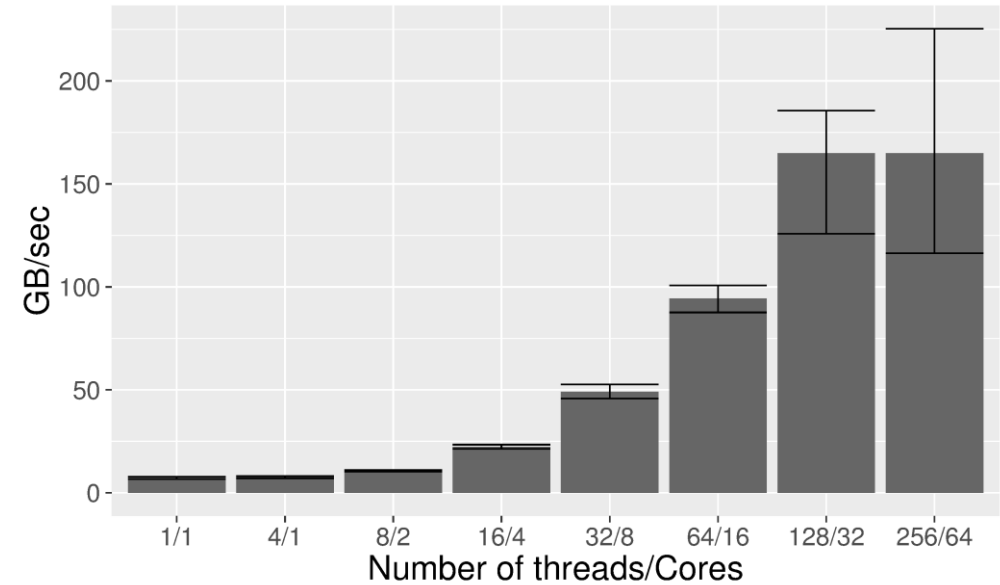
Triad
Cache Mode
SNC4

Triad
Flat Mode
SNC4

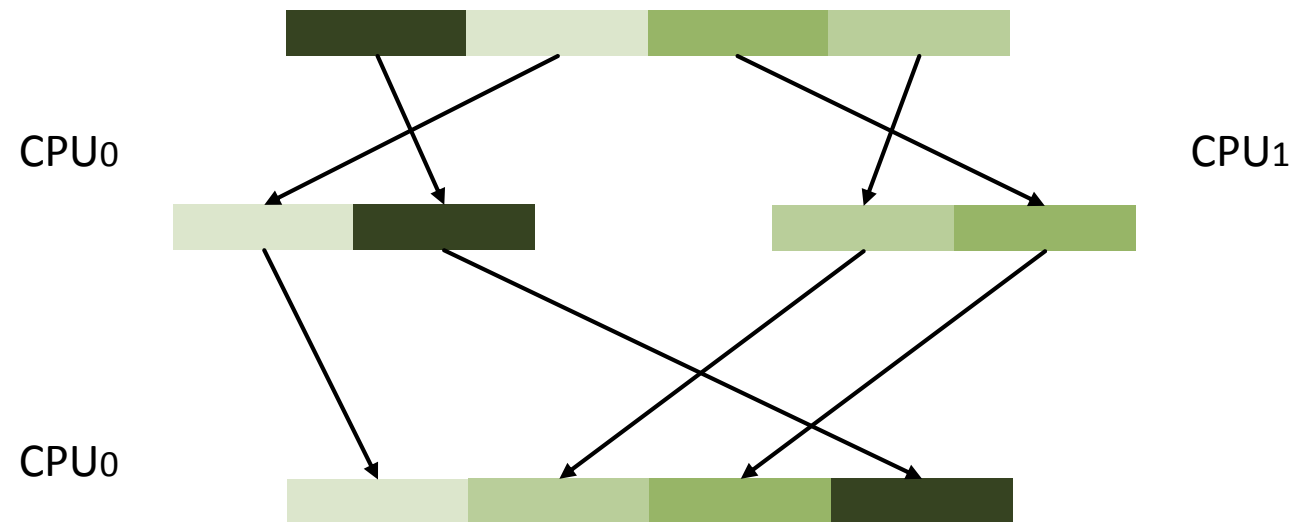
Check! Let's do something "Big Data"!



We'll need a parallel sort on all cores, right?



Memory model: Bitonic Mergesort



- Slices of 16 elements go through a bitonic network.
- Communication: CPU₀ accesses data from local and remote caches.
- Synchronization: CPU₀ waits for CPU₁.
- Memory accesses: latency vs. bandwidth.

Modeling Bitonic Sorting

L1 access cost

memory access cost

$$C_{L1}(n) = [\log_2(n) - 1]2n \cdot \text{cost}_{L1} + 2n \cdot \text{cost}_{mem}$$

elements in L1

$$C_{L2}(n) = \frac{n}{n_{L1}} C_{L1}(n_{L1}) +$$

$$+ [\log_2(n) - \log_2(n_{L1})]2n \cdot \text{cost}_{L2}$$

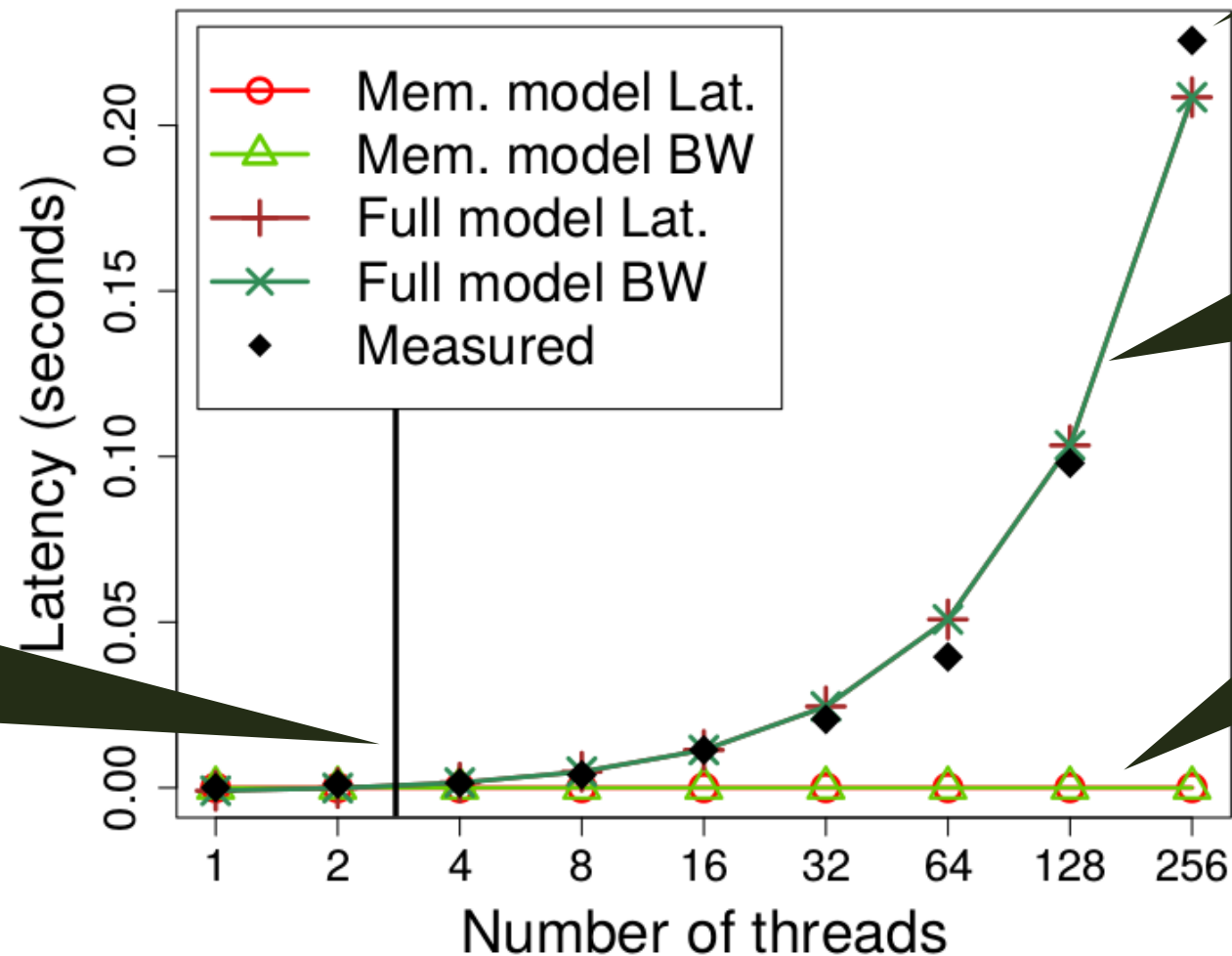
L2 access cost

elements in L2

$$C_{mem}(n) = \frac{n}{n_{L2}} C_{L2}(n_{L2}) +$$

$$+ [\log_2(n) - \log_2(n_{L2})]2n \cdot \text{cost}_{mem}$$

Bitonic Sort of 1 kiB



measurement

model including synchronization cost

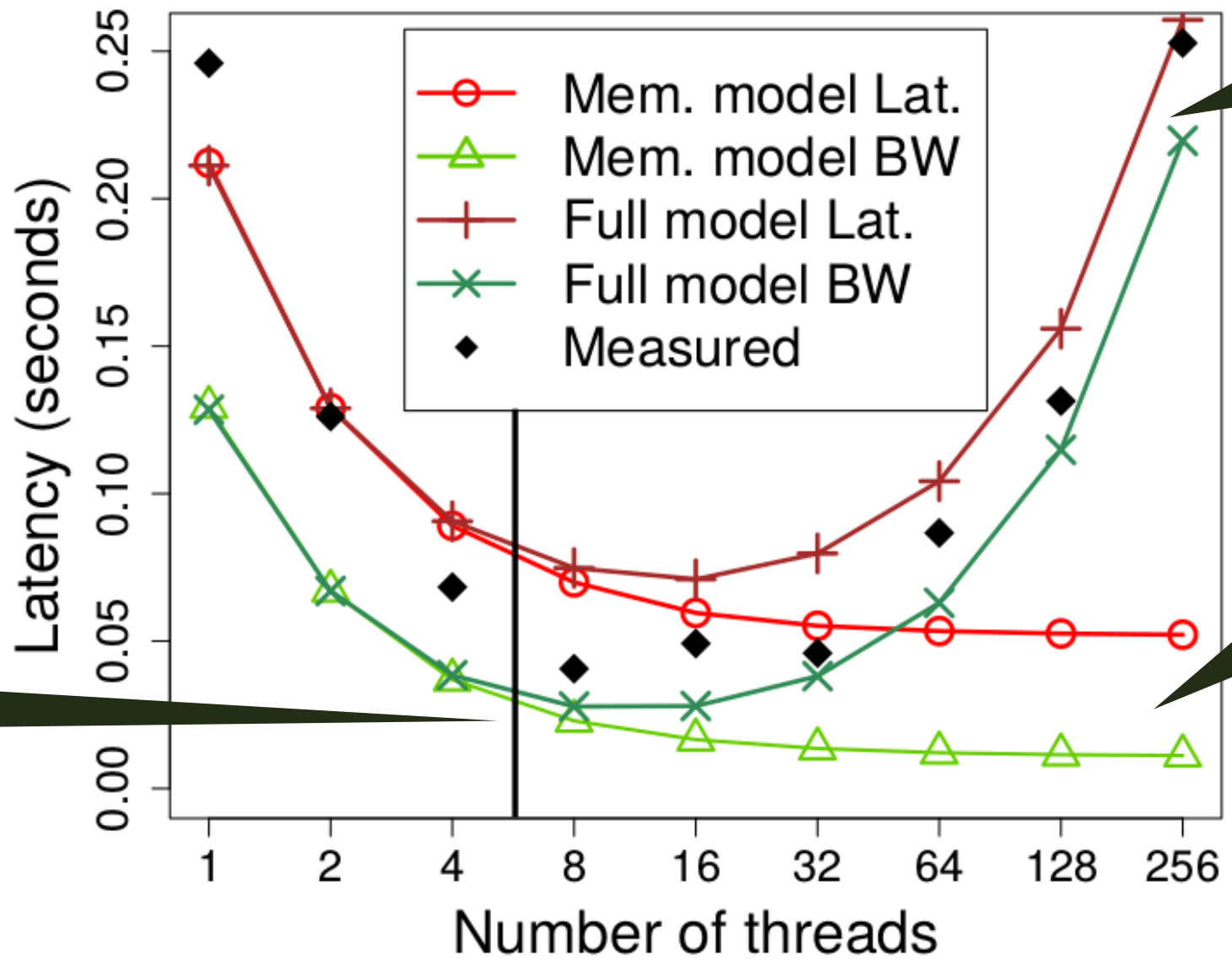
model (just) memory costs

Synchronization outweighs memory costs for small data!

So don't parallelize too much!

(a) Sorting 1 KB of integers.

Bitonic Sort of 4 MiB

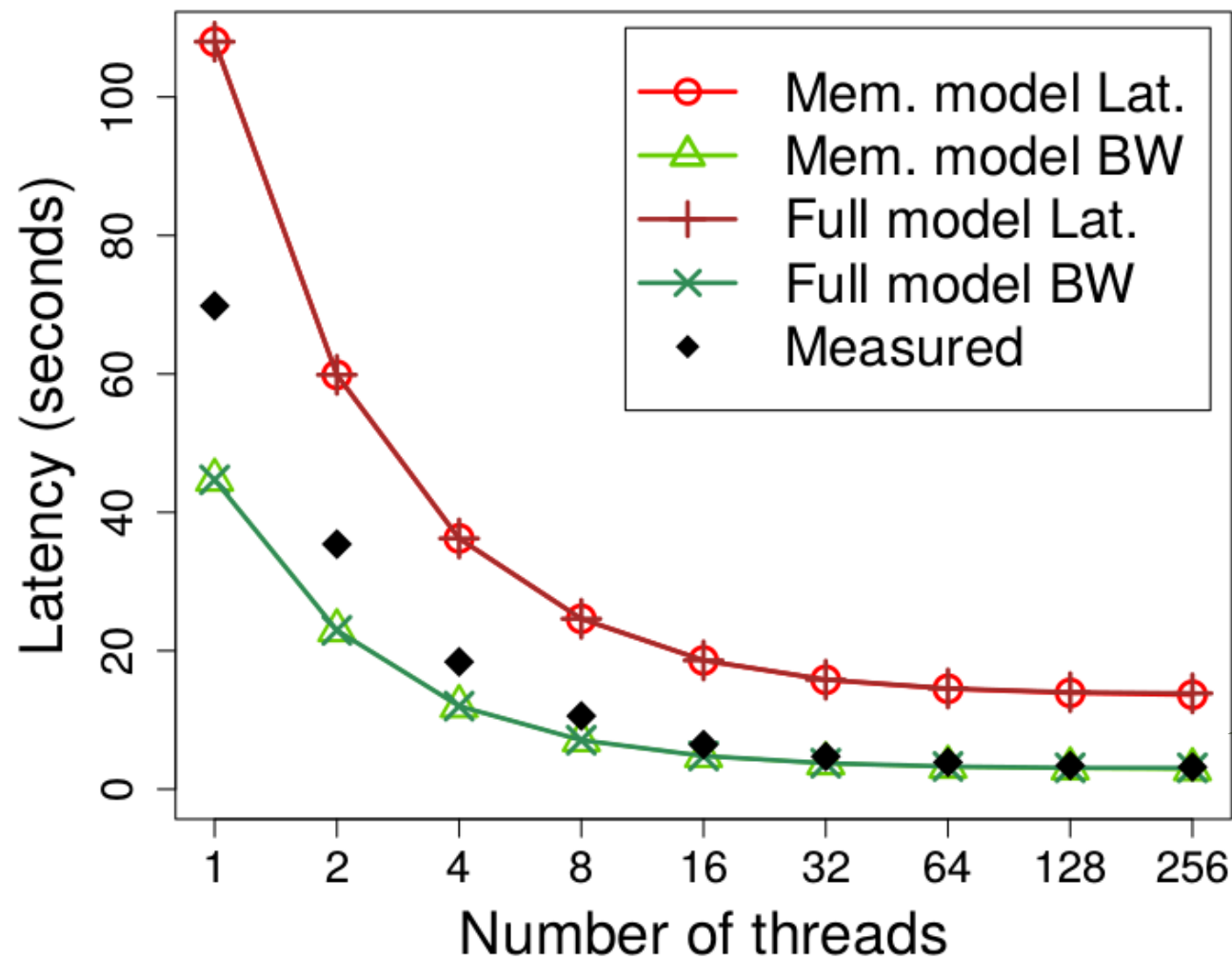


model including synchronization cost

"parallelization boundary"

model (just) memory costs

Bitonic Sort of 1 GiB



Always use all cores!

synchronization negligible

The most surprising result last ...

Hey, but which memory? DRAM or MCDRAM?



<title>code ninja</title>

The model (and practical measurements) indicate that it does not matter.

Thesis: the higher bandwidth of MCDRAM did not help due to the higher latency ($\log^2 n$ depth).

Disclaimer: this is NOT the best sorting algorithm for Xeon Phi KNL. It is the best we found with limited effort. We suspect that a combination of algorithms will perform best.

Scientific performance engineering for complex memory systems

Step 1: Understand core-to-core transfers – MESIF cache coherence

		Software NUMA		Software UMA		
		SNC4	SNC2	QUAD	HEM	A2A
Latency [ns] (Copy/BenchIT)	Local (L1)	3.8	3.8	3.8	3.8	3.8
	Tile (L2)	34 (M)	34 (M)	34 (M)	34 (M)	34 (M)
		17 (E)	18 (E)	18 (E)	18 (E)	18 (E)
		14 (S,F)	14 (S,F)	14 (S,F)	14 (S,F)	14 (S,F)
	Remote	107-122 (M)	111-125 (M)	119 (M)	120 (M)	122 (M)
	98-114 (E)	104-117 (E)	116 (E)	116 (E)	116 (E)	
	96-118 (S,F)	104-118 (S,F)	107-117 (S,F)	107-117 (S,F)	109-117 (S,F)	
Bandwidth [GB/s] (Read)		2.5	2.5	2.5	2.5	2.5
Bandwidth [GB/s] (Copy)	Tile	6.7 (M)	7.6 (E)			
	Remote	7.7				
Congestion (P2P pairs)						
Contention [ns] (1:N copy)	Linear, $\mathcal{T}_C(N) = \alpha + \beta \cdot N$	α	200	β	34	

Write back overhead

Only 5-15% difference

Contention effects?

All values are medians within 10% of the 95% nonparametric CI. Cf. TH, RB: "Scie"

Step 2: Understand core-to-memory transfers – DRAM and MCDRAM

Flat Mode			Software NUMA		Software UMA		
			SNC4	SNC2	QUAD	HEM	A2A
Latency [ns] (BenchIT)	DRAM	130-140	134-146	140	140	139	
	MCDRAM	160-175	160-170	167	167	168	
Bandwidth [GB/s] (Copy NT / STREAM Copy)	DRAM	69 / 77	69 / 77	70 / 77	71 / 77	71 / 77	
	MCDRAM	342 / 418	333 / 388	333 / 415	315 / 372	306 / 359	
Bandwidth [GB/s] (Read)	DRAM	71	71	77	77	77	
	MCDRAM	243	288	314	314	314	
Bandwidth [GB/s] (Write)	DRAM	33	34	36	36	36	
	MCDRAM	147	163	171	165	161	
Cache Mode	Latency [ns] (BenchIT)	71 / 82	74 / 82	73 / 82	73 / 82	73 / 82	
	Bandwidth [GB/s] (Copy NT / STREAM Copy)	347 / 441	340 / 441	332 / 434	325 / 427		
Bandwidth [GB/s] (Read)	DRAM	161-171	166	168	172		
	MCDRAM	130 / 252	175 / 255	134 / 237	132 / 233		
Bandwidth [GB/s] (Write)	DRAM	95	124	128	118		
	MCDRAM	56	72	72	68		
Bandwidth [GB/s] (Read)	DRAM	246 / 294	296 / 309	273 / 274	264 / 269		
	MCDRAM						

MCDRAM 20% slower!

MCDRAM 4-6x faster!

Need to read and write for full bandwidth

Modeling Communication in Cache-Coherent SMP Systems - A Case-Study with Xeon Phi, ACM HPDC'13

S. Ramos, TH: "Cache line aware optimizations for cclUMA systems (IEEE TPDS'17)"

Questions/Discussions?

Performance engineers optimize your code!

Are you kidding me?

		Software NUMA		Software UMA		
		SNC4	SNC2	QUAD	HEM	A2A
Latency [ns] (Copy/BenchIT)	Local (L1)	3.8	3.8	3.8	3.8	3.8
	Tile (L2)	34 (M)	34 (M)	34 (M)	34 (M)	34 (M)
		17 (E)	18 (E)	18 (E)	18 (E)	18 (E)
		14 (S,F)	14 (S,F)	14 (S,F)	14 (S,F)	14 (S,F)
	Remote	107-122 (M)	111-125 (M)	119 (M)	120 (M)	122 (M)
	98-114 (E)	104-117 (E)	116 (E)	116 (E)	116 (E)	
	96-118 (S,F)	104-118 (S,F)	107-117 (S,F)	107-117 (S,F)	109-117 (S,F)	
Bandwidth [GB/s] (Read)		2.5	2.5	2.5	2.5	2.5
Bandwidth [GB/s] (Copy)	Tile	6.7 (M)	6.7 (M)	7.5 (M)	7.4 (M)	7.5 (M)
	Remote	7.6 (E)	6.7 (E)	9.2 (E)	9.2 (E)	9.2 (E)
		7.7	6.7	7.5	7.5	7.5
Congestion (P2P pairs)		None				
Contention [ns] (1:N copy)	Linear, $\mathcal{T}_C(N) = \alpha + \beta \cdot N$	α	200	β	34	

(a) All-to-all mode.

(b) Quadrant mode.

(c) SNC4 mode.

<title>code ninja</title>

A principled approach to designing cache-to-cache broadcast algorithms

Multi-ary tree example

Tree depth

Level size

Tree cost

Reached threads

$$T_{tree} = \sum_{i=1}^d \mathcal{T}_C(k_i) = \sum_{i=1}^d (c \cdot k_i + b)$$

$$= \sum_{i=1}^d (R_R + R_L + c \cdot (k_i - 1))$$

$$T_{sbcast} = \min_{d, k_i} \left(T_{fw} + \sum_{i=1}^d (c \cdot k_i + b) + \sum_{i=1}^d \mathcal{T}_{nb}(k_i + 1) \right)$$

$$\dots N \leq 1 + \sum_{i=1}^d \prod_{j=1}^i k_j, \forall i < j, k_i \leq k_j$$

Modeling Communication in Cache-Coherent SMP Systems - A Case-Study with Xeon Phi, ACM HPDC'13

S. Ramos, TH: "Cache line aware optimizations for cclUMA systems (IEEE TPDS'17)"