

RETHINKING DATABASE ALGORITHMS FOR PCM

Shimin Chen (HP Labs), Phillip Gibbons (Intel Labs), Suman Nath (Microsoft Research)

PHASE CHANGE MEMORY (PCM)

Emerging byte-addressable non-volatile memory technology

	DRAM	PCM*	NAND Flash
Page size	64B	64B	4KB
Page read latency	20-50ns	~ 50ns	~ 25 μ s
Page write latency	20-50ns	~ 1 μ s	~ 500 μ s
Write bandwidth	~GB/s	50-100 MB/s	5-40 MB/s
	per die	per die	per die
Erase latency	N/A	N/A	~ 2 ms
Endurance	∞	$10^6 - 10^8$	$10^4 - 10^5$
Read energy	0.8 J/GB	1 J/GB	1.5 J/GB [28]
Write energy	1.2 J/GB	6 J/GB	17.5 J/GB [28]
Idle power	~100 mW/GB	~1 mW/GB	1-10 mW/GB
Density	1x	2 - 4x	4x

Sources (all public): [Doller'09] [Lee et al. '09] [Qureshi et al.'09] [Tseng et al. '06]. *PCM data is forecasted data

- Compared to DRAM, PCM has better density and scalability; PCM has similar read latency but longer write latency
- Compared to NAND Flash, PCM is byte-addressable, has orders of magnitude lower latency and higher endurance

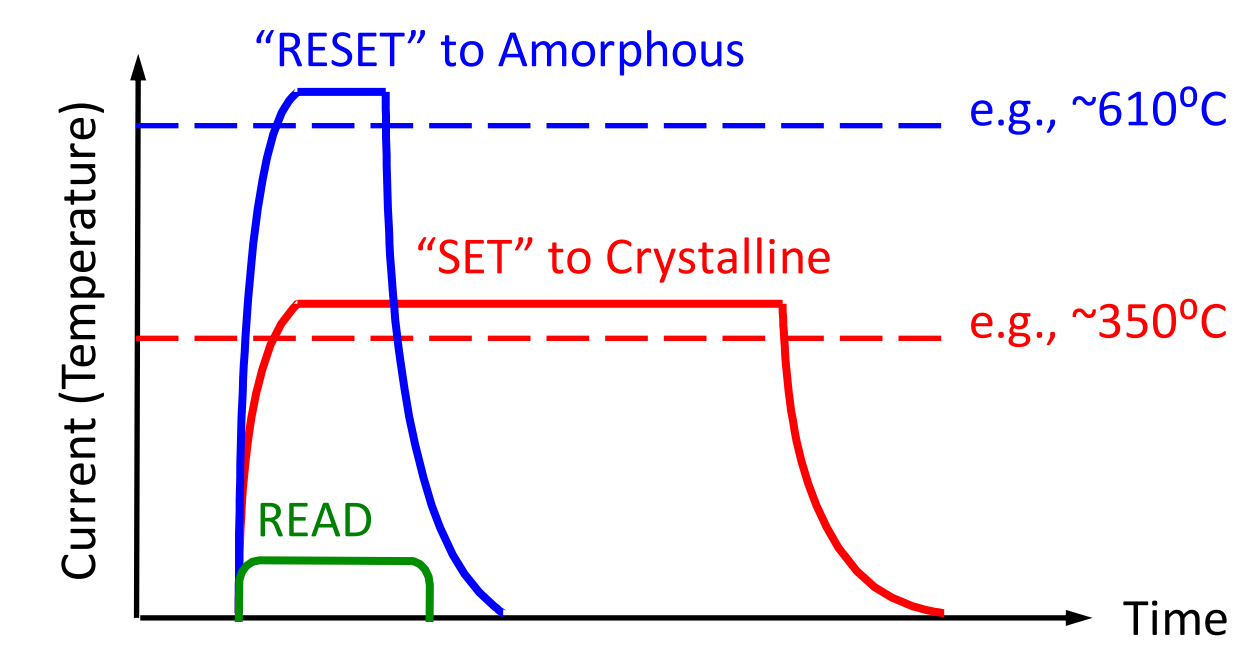
GOAL: DATABASE MINIMIZES PCM WRITES

PCM-DB: Database systems exploiting PCM as primary main memory

- Initial focus on key database algorithms:
 - B⁺-Tree Index & Hash Joins**
 - Structures reside in PCM main memory
 - Optional DRAM is another (transparent or explicit) cache
- Key analytical metrics:
 - Total Wear, Energy, Total PCM Access Latency
- Experimental Setup: PTLSSim extended with PCM support
 - B⁺-Tree: Node size 8 cache lines, 50M entries, 75% full; Inserting / Deleting / Searching 500K random keys
 - Hash Join: 500MN joins 100MB; varying record size from 20B-100B

CHALLENGE: PCM WRITES

- Limited endurance**
 - Wear out quickly for hot spots
- High energy consumption**
 - 6-10X more energy than a read
- High latency & low bandwidth**
 - SET/RESET time > READ time
 - Limited instantaneous electric current level, requires multiple rounds of writes: Writes 20X slower than reads

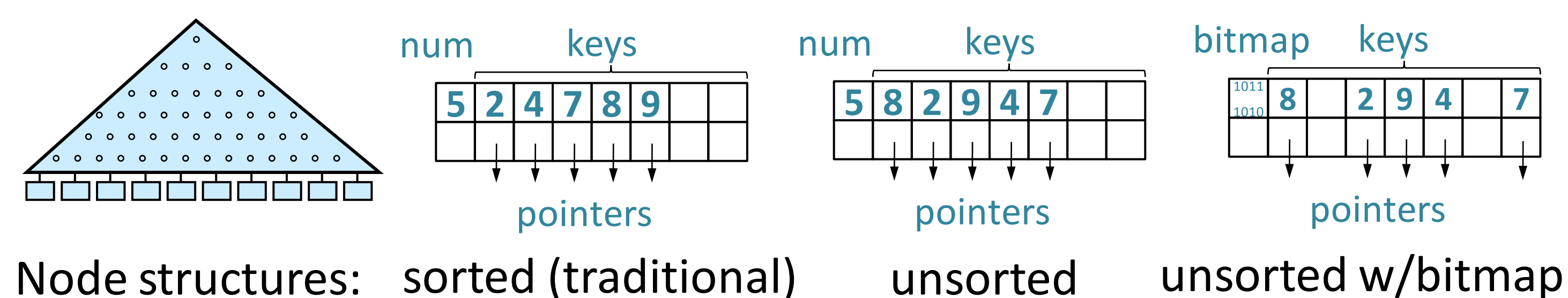


[Cho&Lee'09] [Lee et al '09] [Yang et al'07] [Zhou et al'09]

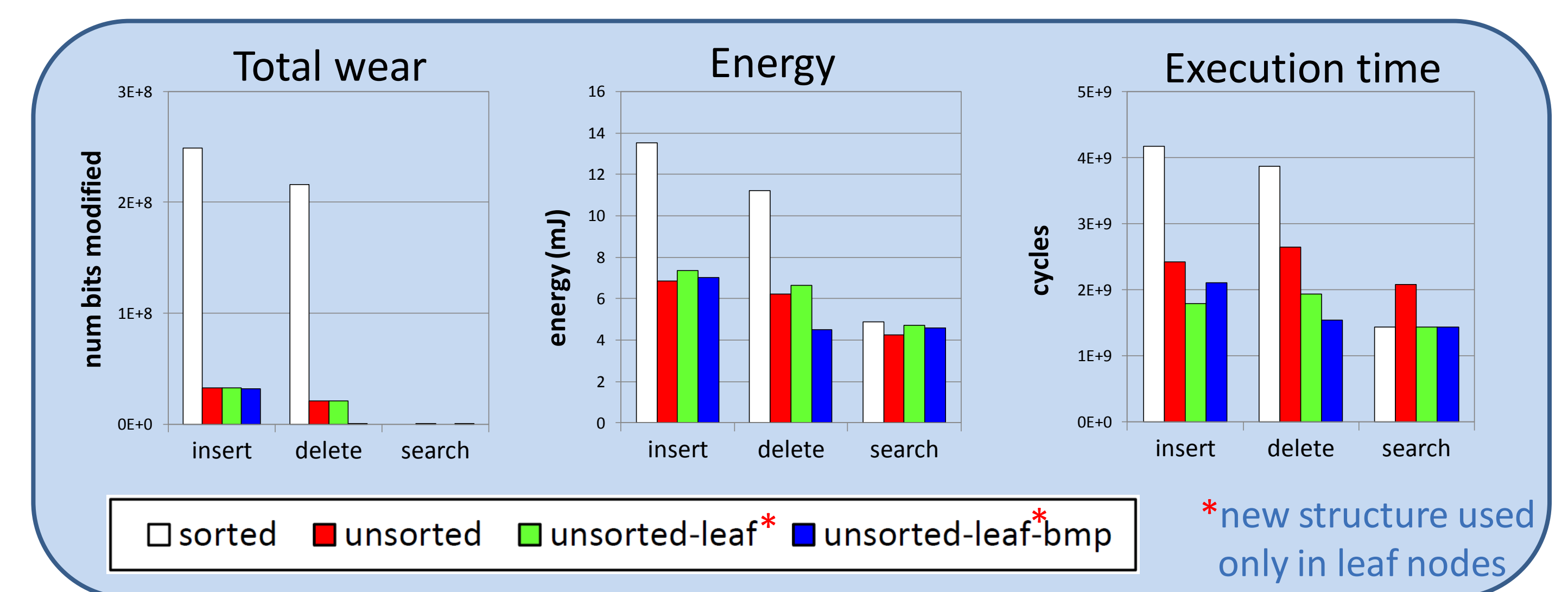
PCM Write Optimizations in literature:

- Baseline: several rounds of writes for a cache line
 - Which bits in which rounds are hard wired
- Optimization: data comparison write:
 - Goal: write only modified bits rather than entire cache line
 - Approach: read-compare-write
- Skipping rounds with no modified bits

B⁺-TREE ALGORITHMS & RESULTS

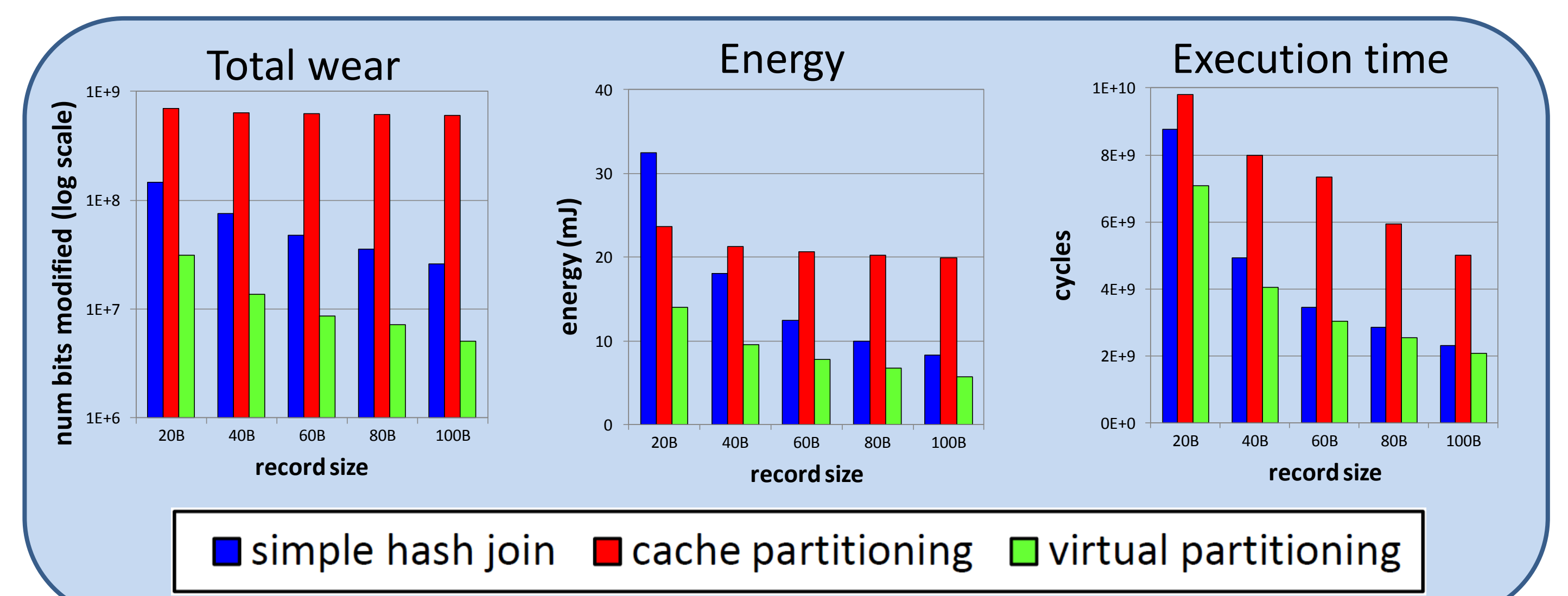
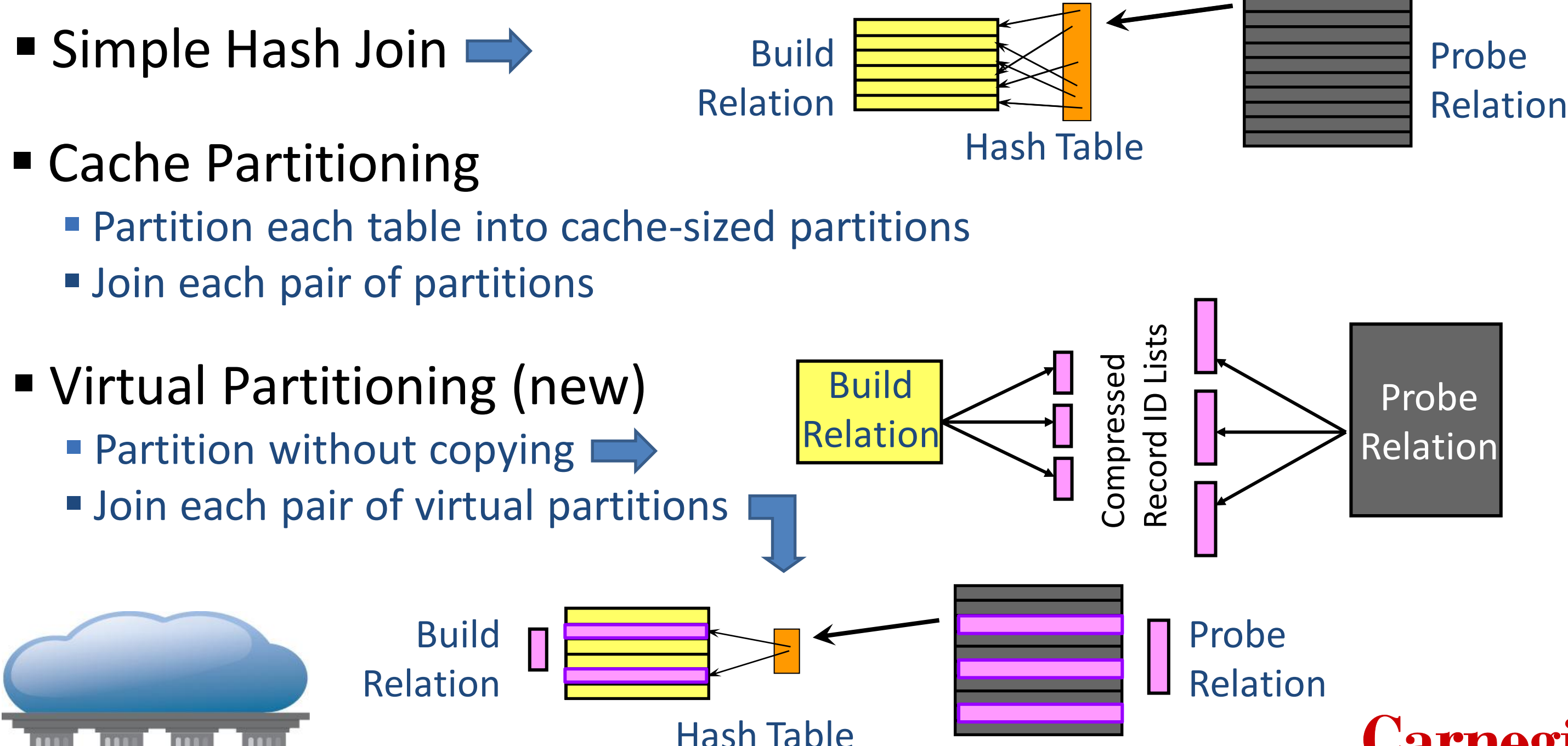


Node structures: sorted (traditional) unsorted unsorted w/bitmap



Unsorted leaf schemes achieve the best performance

IN-MEMORY HASH JOIN ALGORITHMS & RESULTS



Virtual Partitioning achieves the best performance



[Paper in CIDR'11]

