

list
cea tech

Aix*Marseille
université
Initiative d'excellence

DE LA RECHERCHE À
L'INDUSTRIE

Exploration and conception of computing architecture of type *in-memory computing* based on emerging non volatile memories

23/10/18-22/10/21

PhD Director: Jean-Michel Portal

PhD Supervisor: Mathieu Moreau

CEA Supervisor: Jean-Philippe Noel



- ▶ **General context**
- ▶ **State of the art on different technologies and implementations**
- ▶ **CSRAM: our own design approach**
- ▶ **Software platforms and explorations**
- ▶ **Results and analysis on different memory computing architectures**
- ▶ **Conclusion and future works**

▶ General context

- Semiconductor industry and issues on current system architectures
- Introduction to memory systems and memory hierarchy
- Main thesis idea → memory computing

▶ State of the art on different technologies and implementations

▶ CSRAM: our own design approach

▶ Software platforms and explorations

▶ Results and analysis on different memory computing architectures

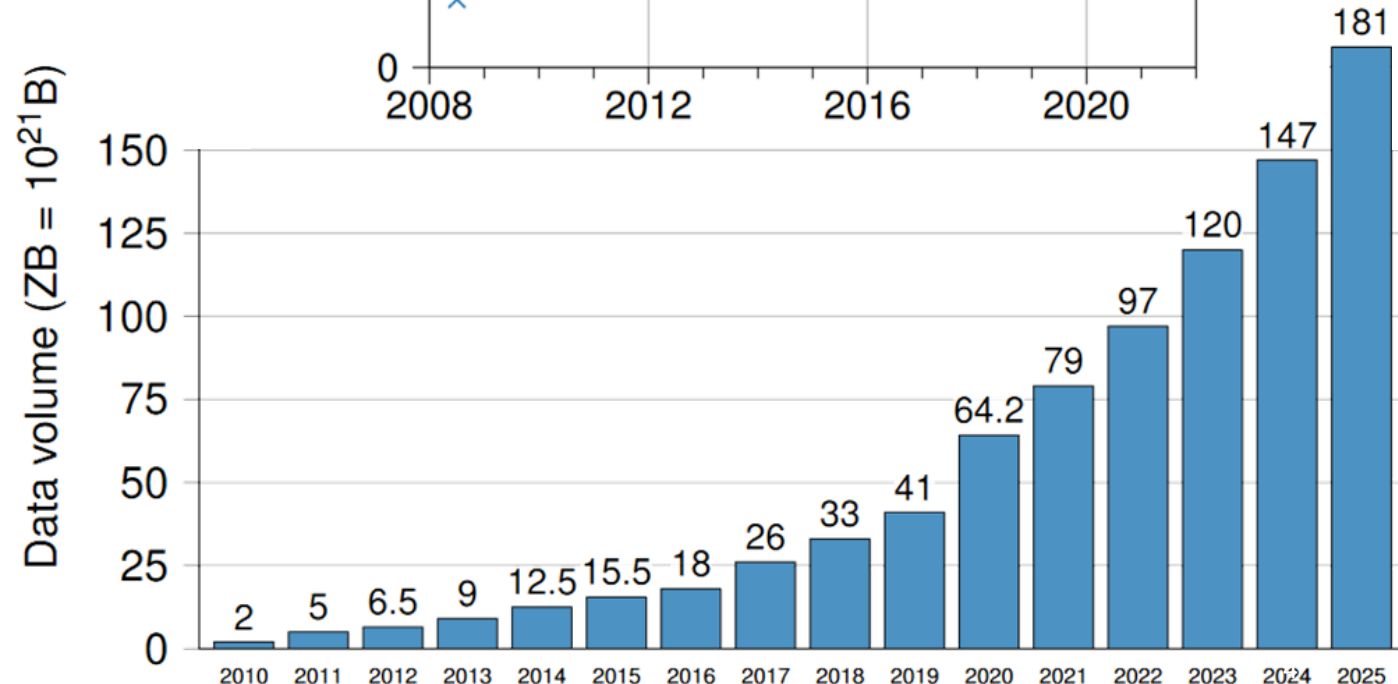
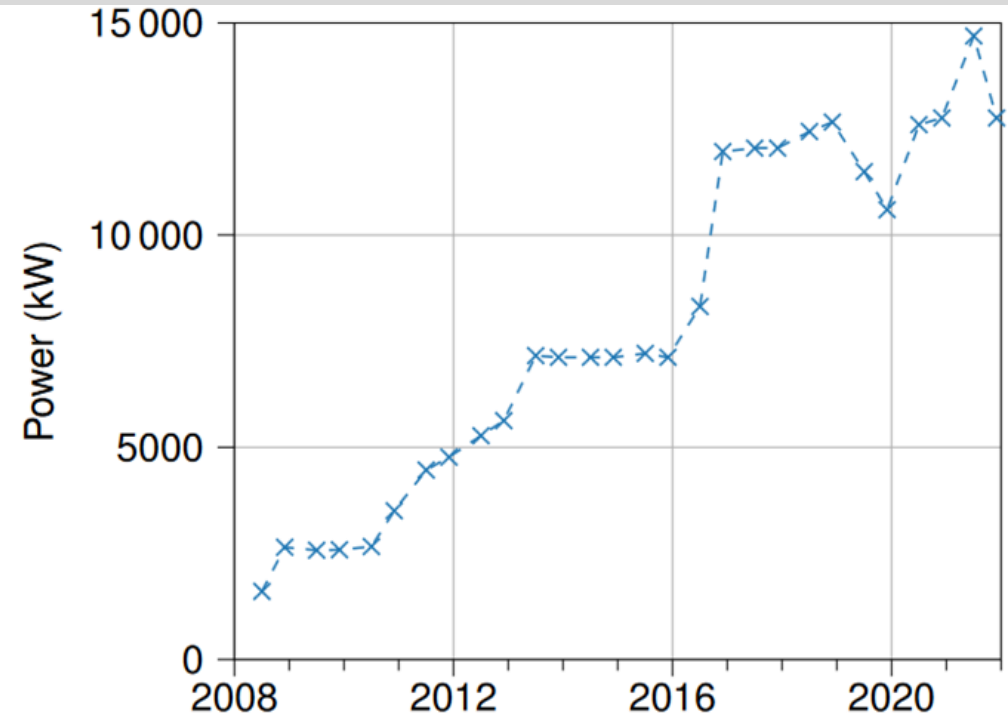
▶ Conclusion and future works

► Transition to digital services consumes more and more power

- Goes against energy transition
- Already 3% of world electric energy consumption (data centers + networks)

► And creates more and more data

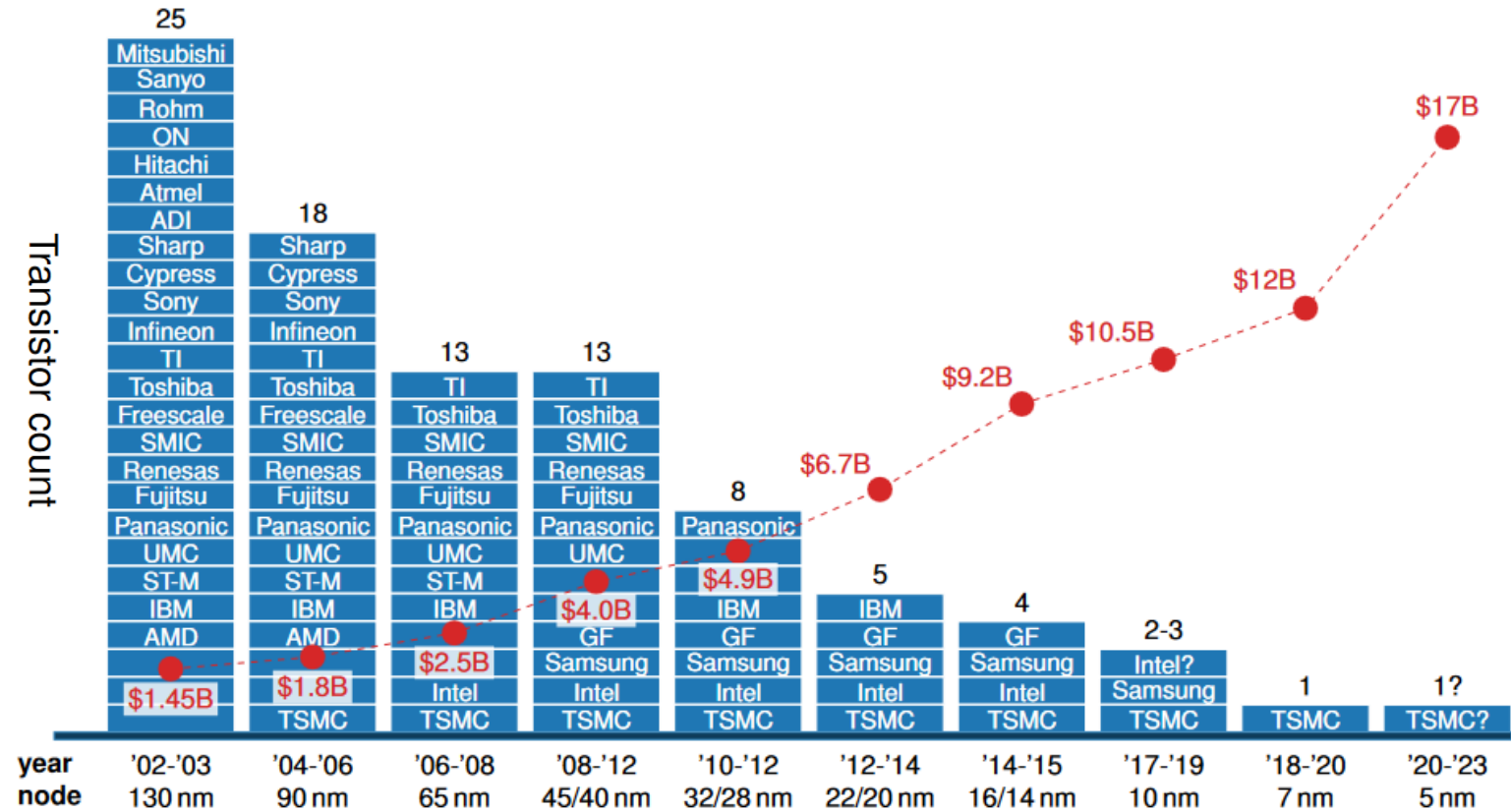
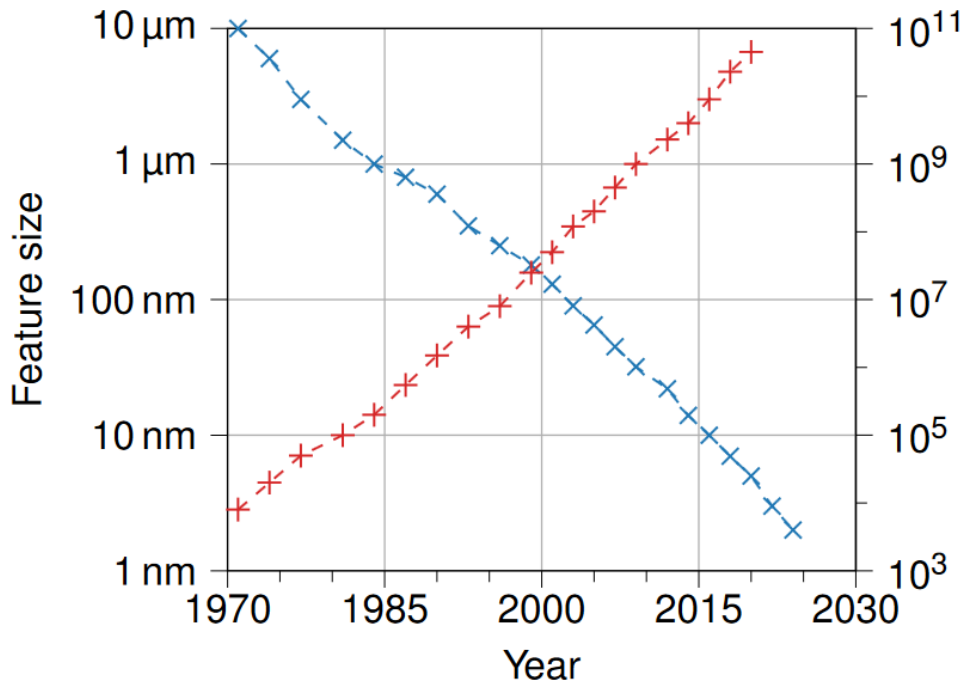
- Data that needs to be treated and stored
- 40% growth rate per year



General context : technology faces a dead end

► Feature size reduction reaches an end

► Cost of more advanced node is skyrocketing



ITRS roadmap, 2022, www.itrs2.net

V. Bertacco. "Re-Imagining Scalable System Design", 2018, VLSI-SoC

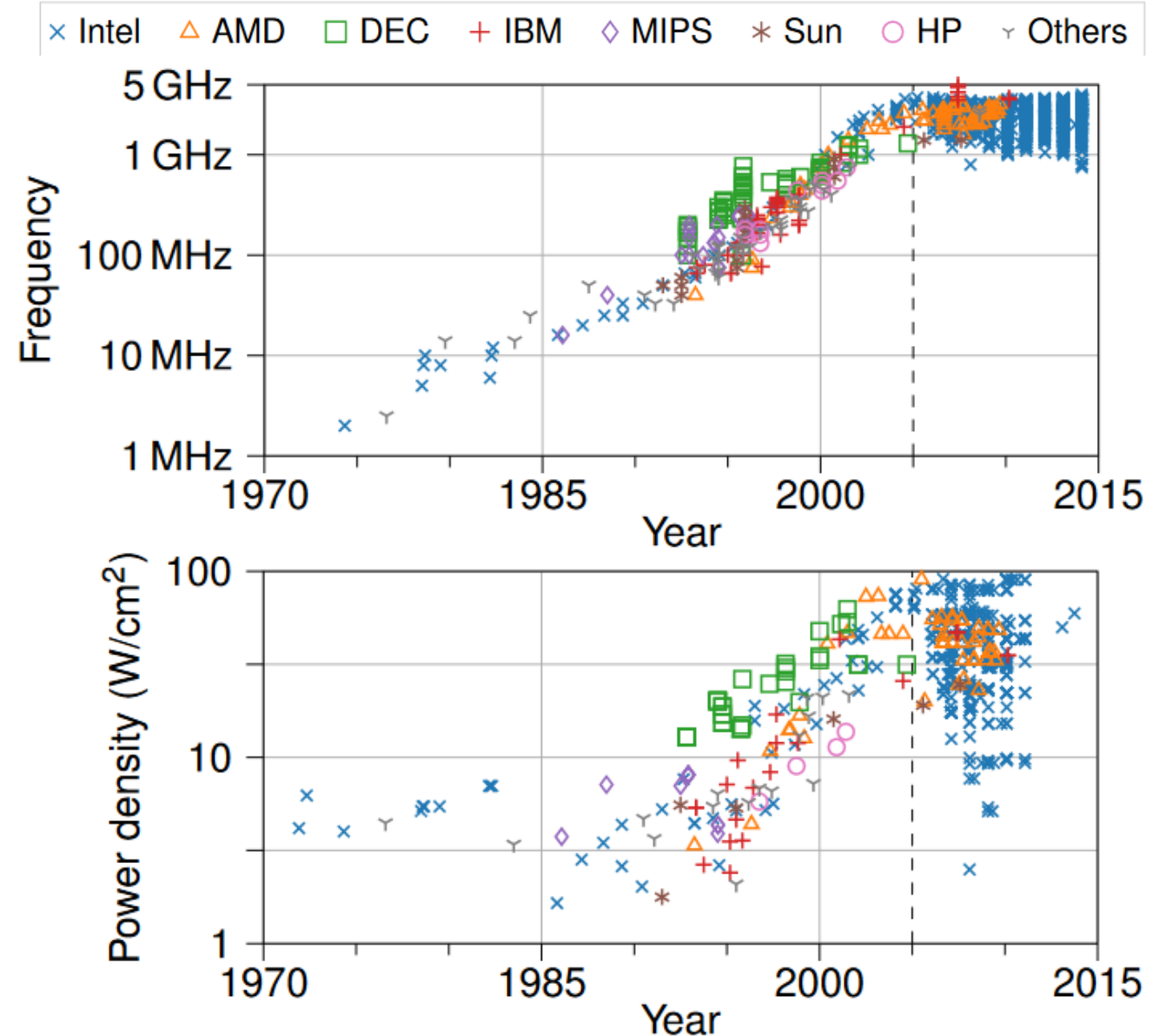


► Dennard's scaling broke in 2005

- Frequency stall in all CPUs around 3GHz

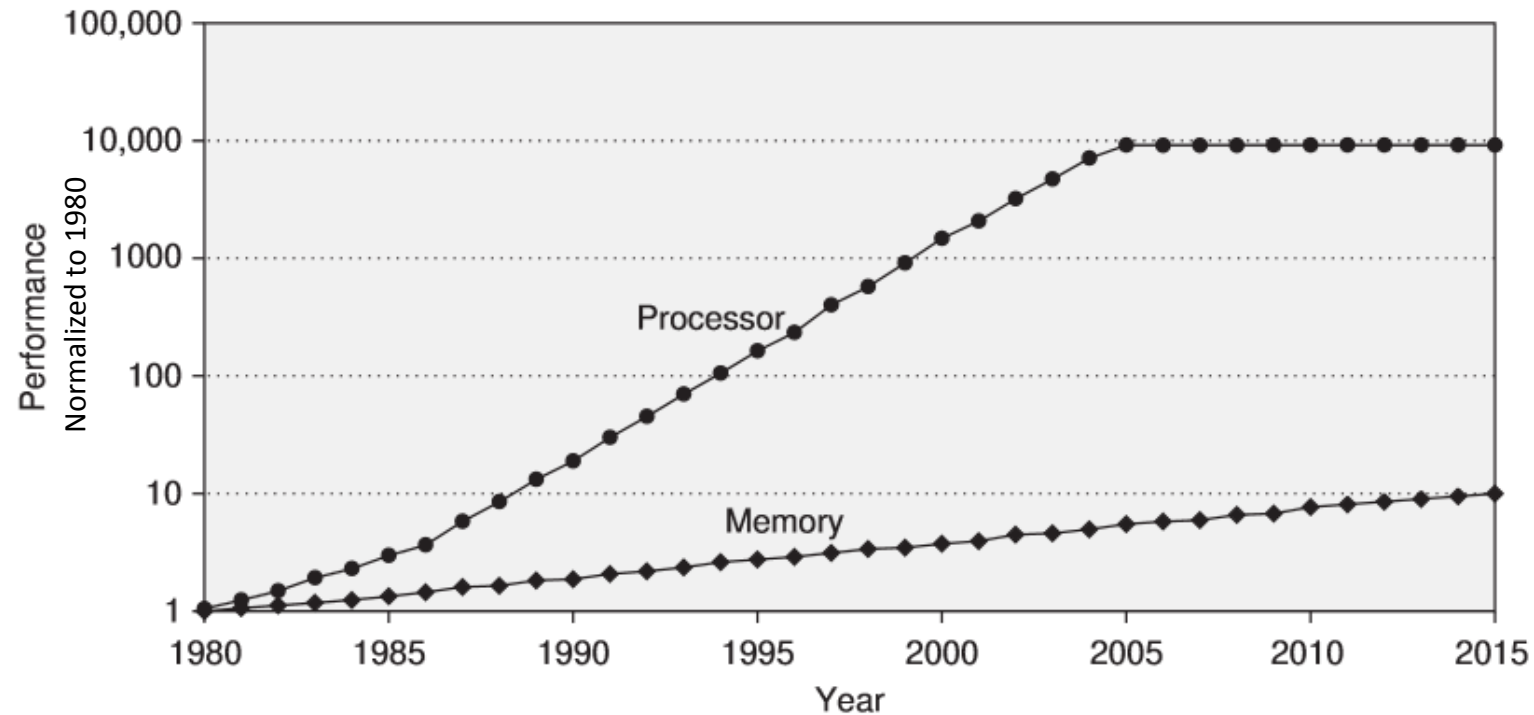
$$P_{dyn} = CV^2 f$$

- Max power density capped at 100W/cm² due to thermal constraint
- Dark Silicon



► Processor performances improvements exceeds memory's

- Different foundry processes
- Limited hardware innovation in memories

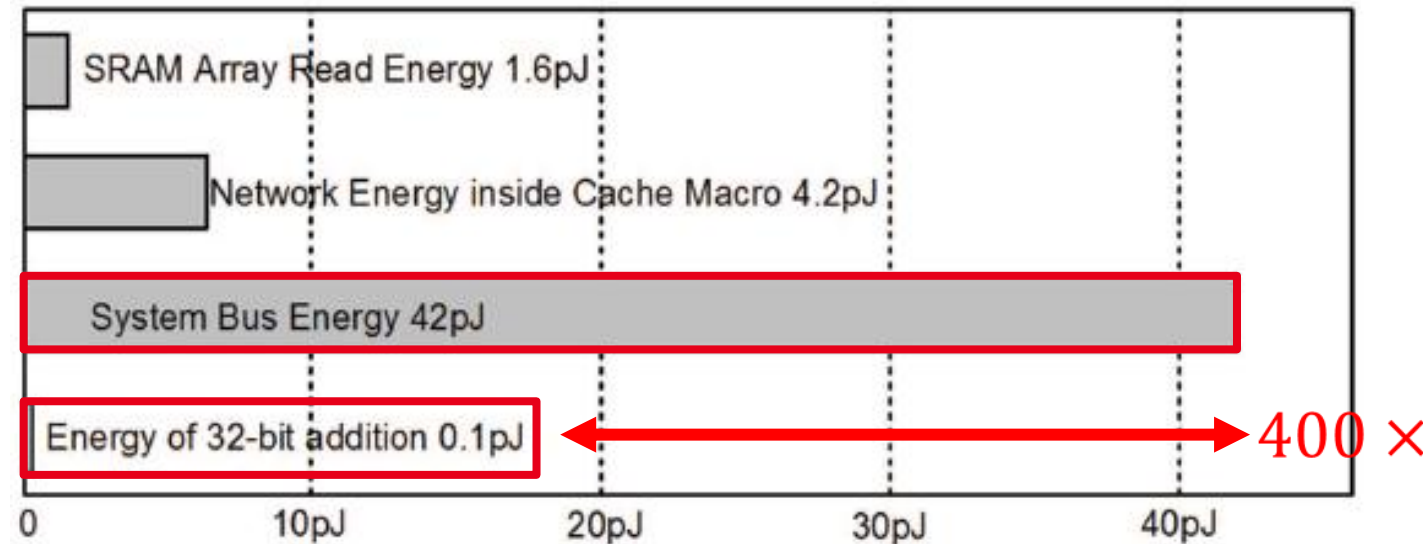
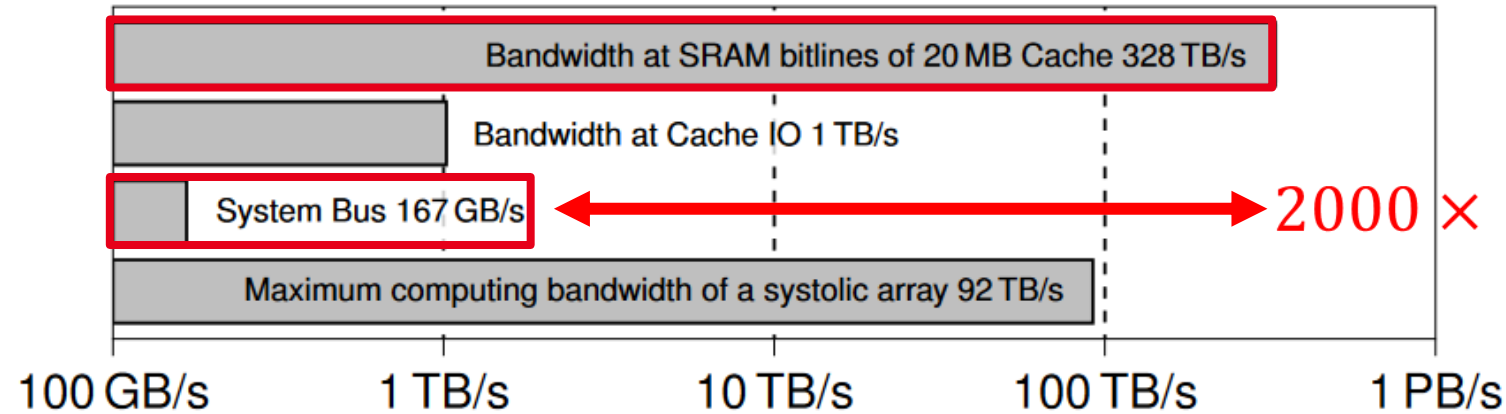


► Internal bandwidth is decimated by IO interfaces

- SA (analog to digital)
- Bus interfaces

► Memories are underutilized and waste 50% to 80% of the energy

- Full line read (64B) just to access a single byte

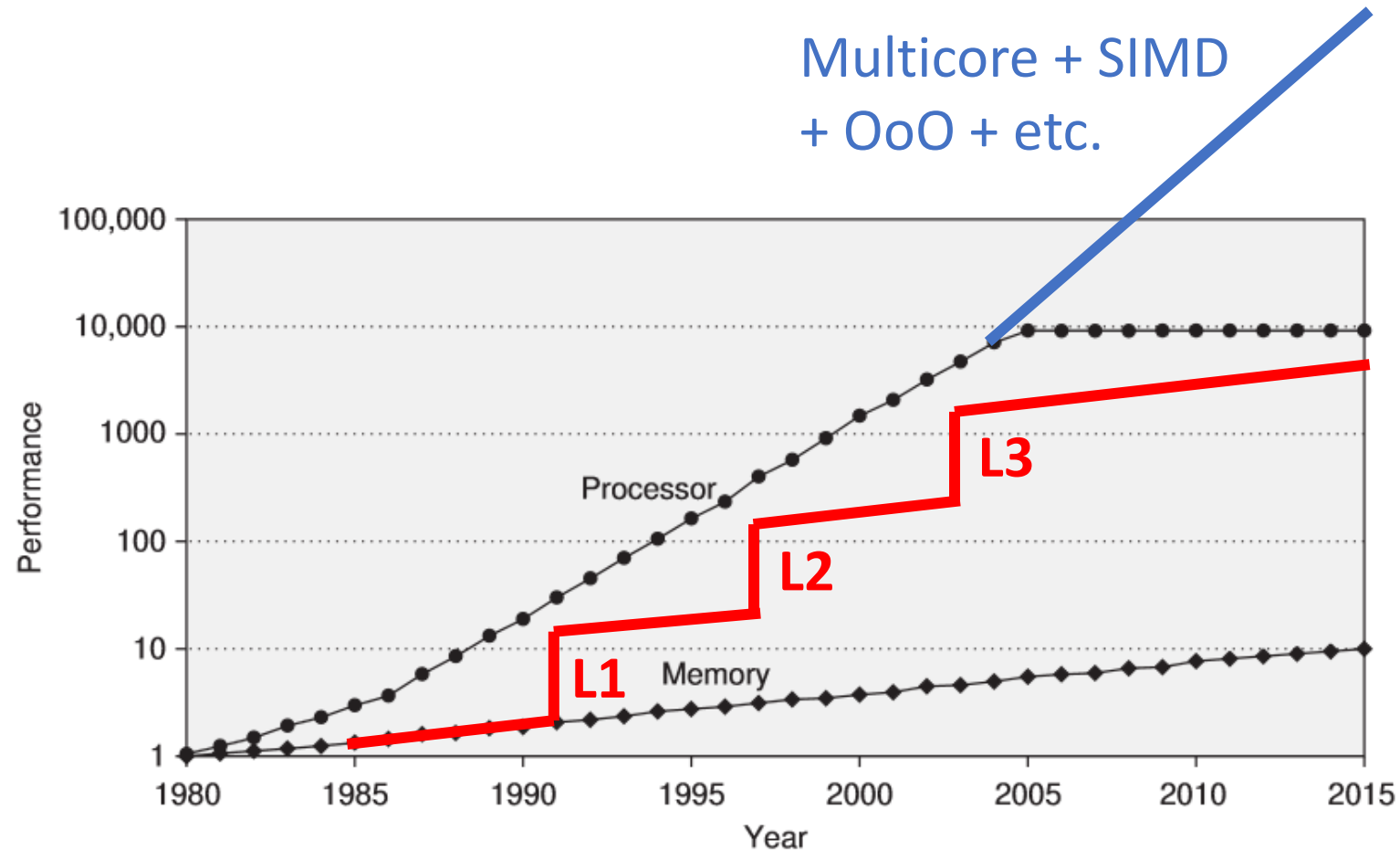


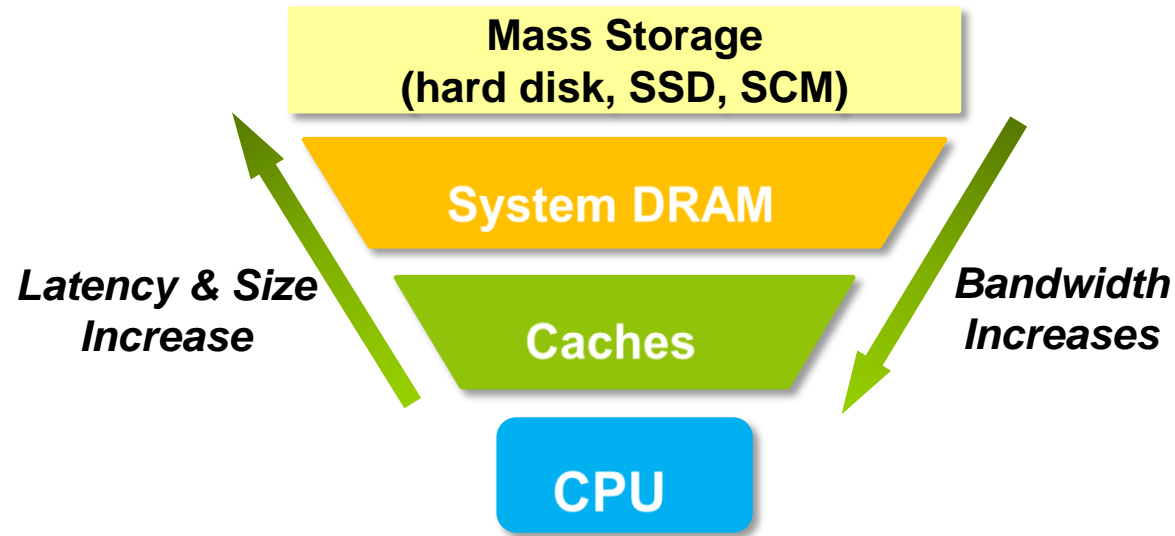
► To bring memory to the performance level of processors, a memory hierarchy is instated

- Prevents the processor from waiting doing nothing
- Every time processor performance scales up, a new memory is inserted

► Instruction centric point of view with lots of hardware innovation on CPU side

- Single Instruction Multiple Data (SIMD)
- Branch prediction
- Out of Order execution
- Speculative execution





Price (\$/GiB)

Density

Latency

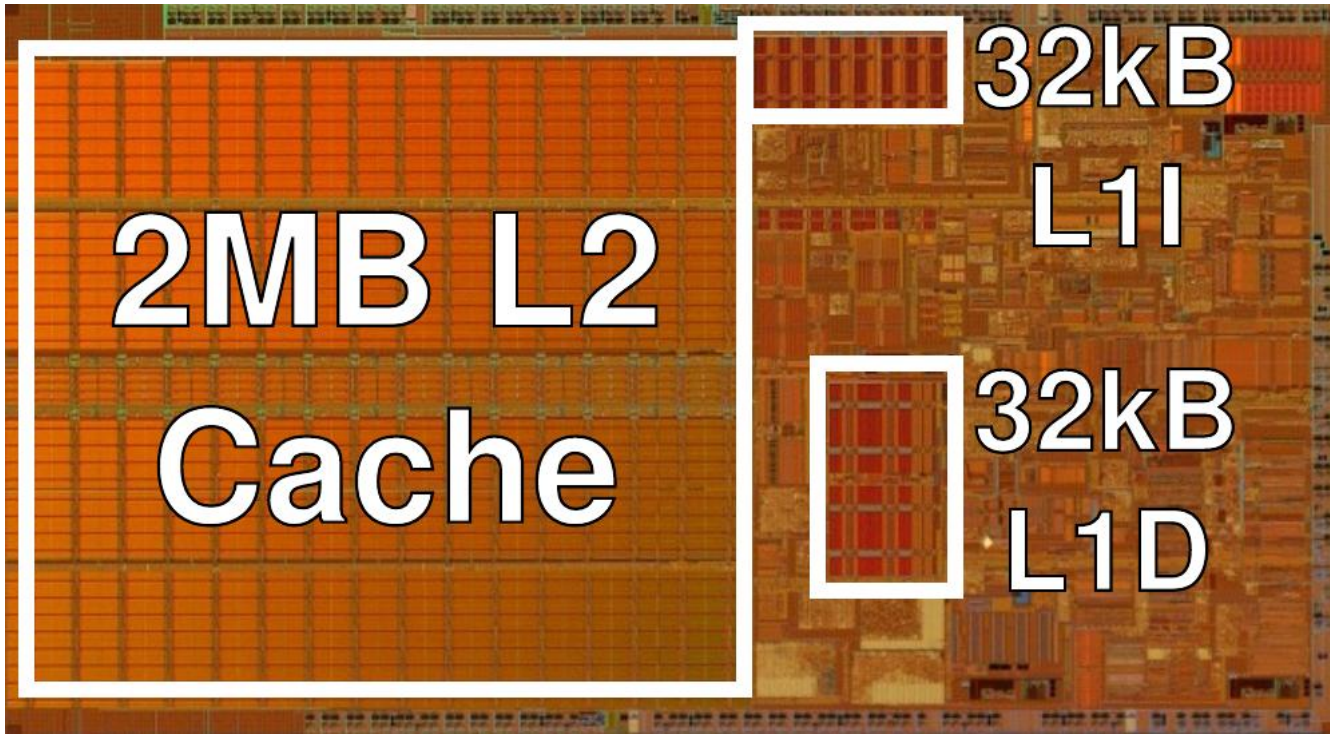
Bandwidth

Persistence

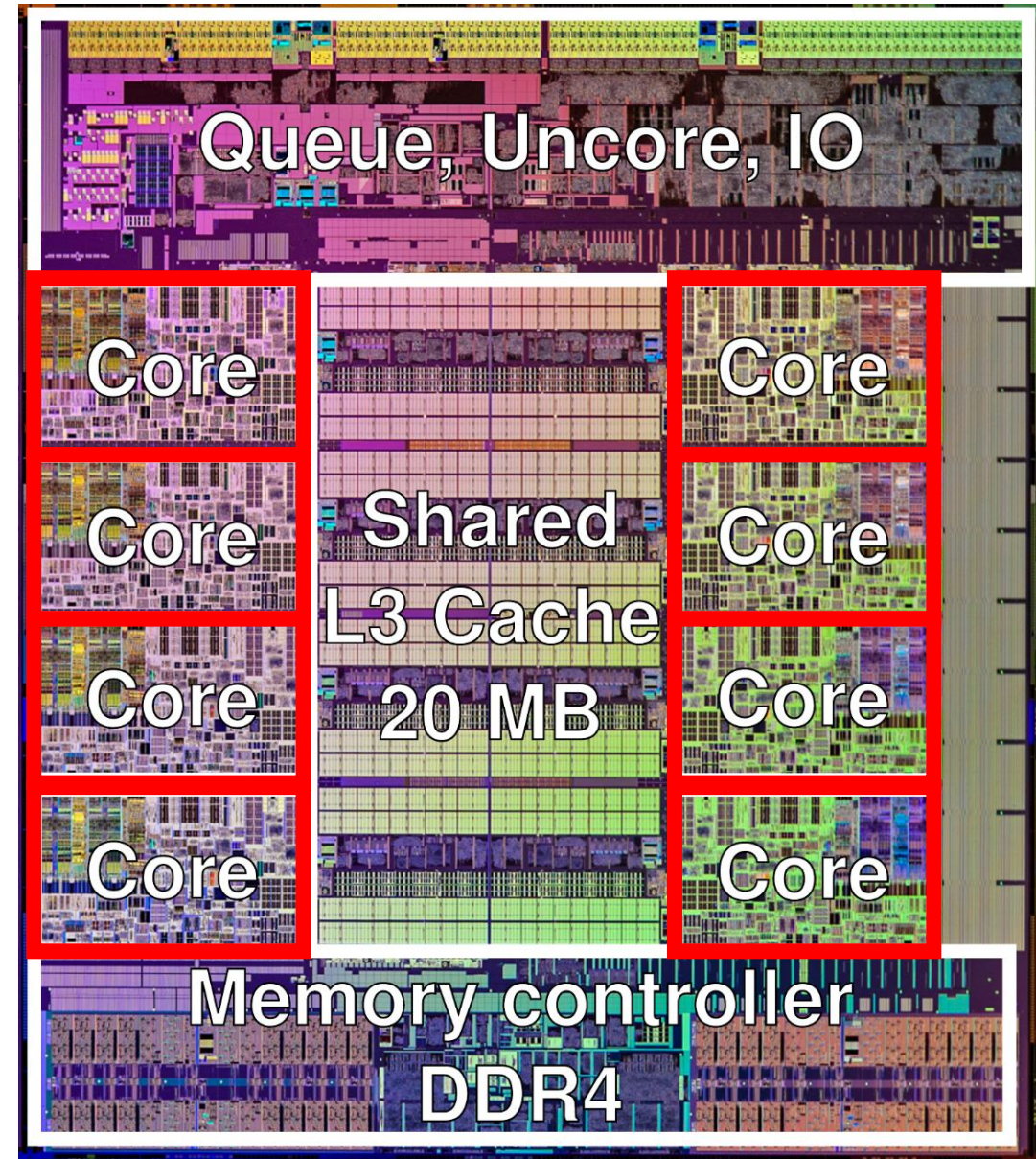
Largest size

	Price (\$/GiB)	Density	Latency	Bandwidth	Persistence	Largest size
SRAM	5000	120 F ² or 2 Gbit/cm ²	1–50 ns	1 TiB/s	10 μs (Power off)	10–100 MiB
DRAM	20	8 F ² or 25 Gbit/cm ²	20–400 ns	10–100 GiB/s	64 ms	100–1000 GiB
Flash	4	<1 F ² or >100 Gbit/cm ²	1–10 μs	1 GiB/s	10–20 yr	10–100 TiB
HDD	0.1	100 Gbit/cm ²	5–20 ms	100 MiB/s	10–100 yr	10–100 TiB

Illustration of low SRAM density



130nm Intel Pentium M (2004)



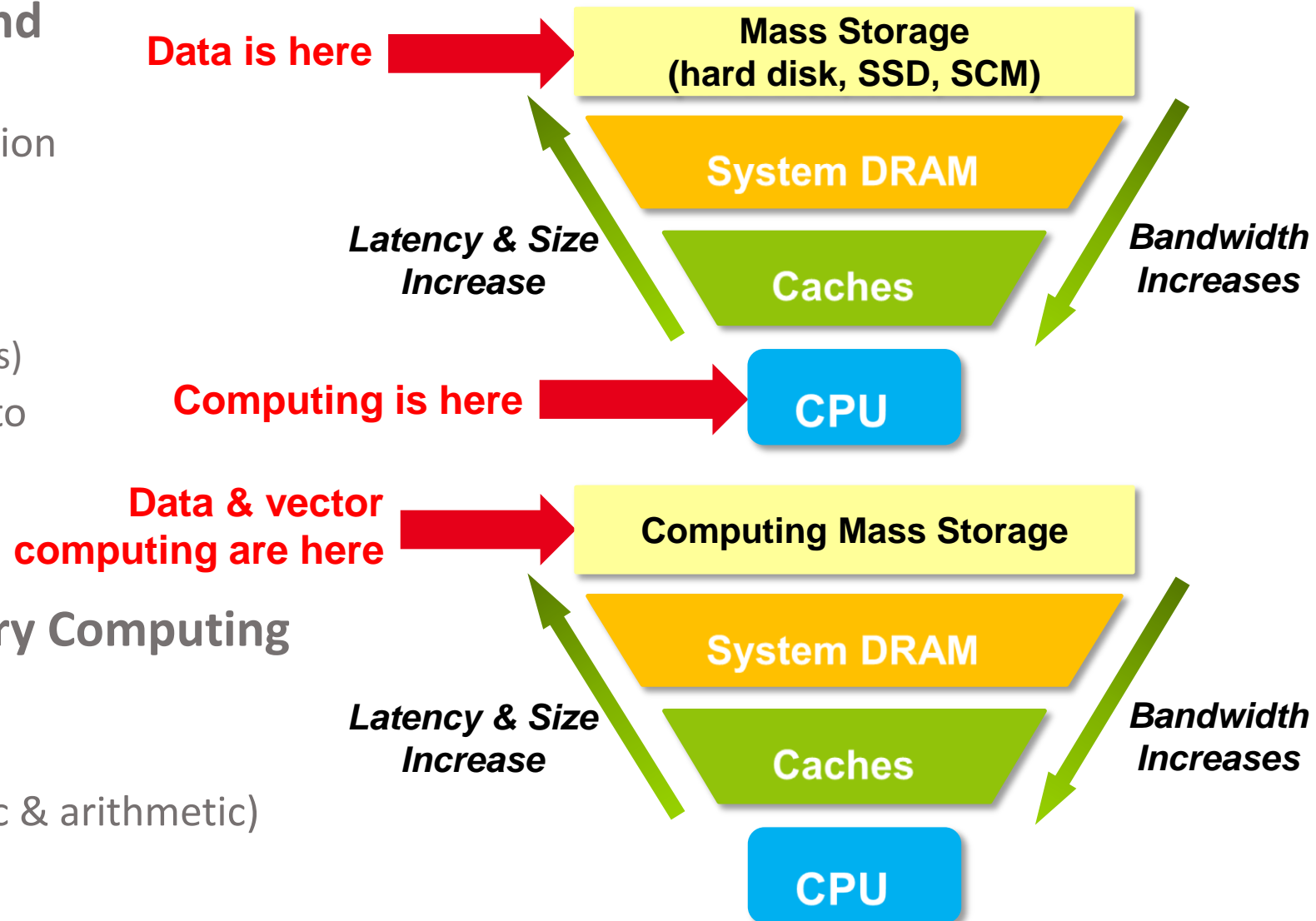
22nm Intel Haswell 8 cores (2014)

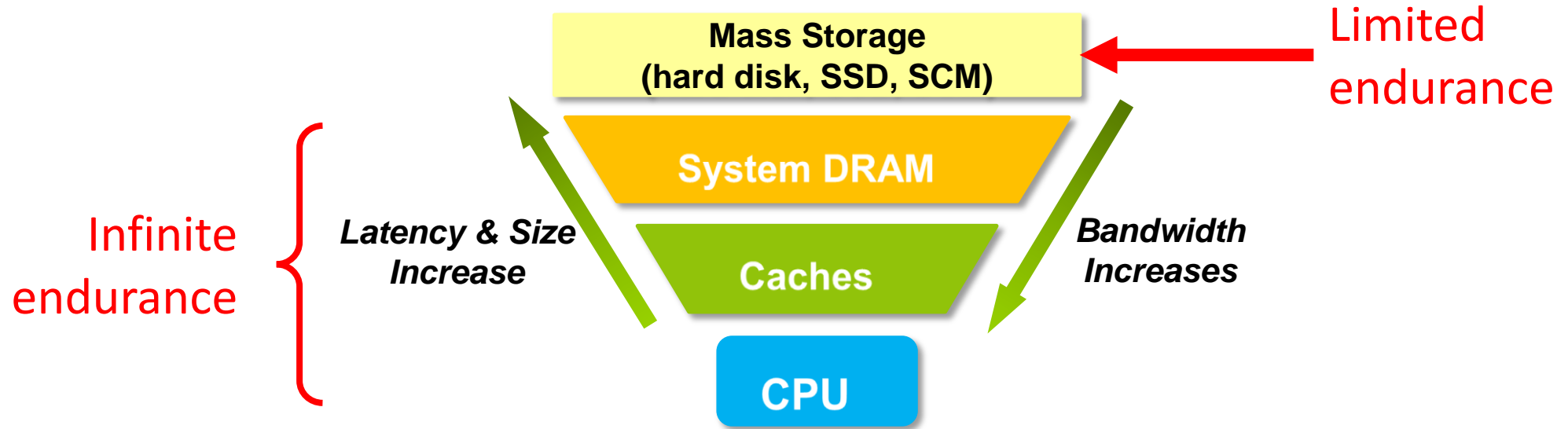
► von-Neumann Architecture and Bottleneck

- Linked to data intensive application (memory bound)
 - Big data era
 - Neural networks
 - Irregular access patterns (graphs)
- Data keeps moving around due to coarse grained blocks

► Proposed solution: In-Memory Computing

- Data is in memory
- Bigger parallelism
- But only simple operations (logic & arithmetic)





	Price (\$/GiB)	Density	Latency	Bandwidth	Persistence	Largest size
SRAM	5000	120 F ² or 2 Gbit/cm ²	1–50 ns	1 TiB/s	10 μs (Power off)	10–100 MiB
DRAM	20	8 F ² or 25 Gbit/cm ²	20–400 ns	10–100 GiB/s	64 ms	100–1000 GiB
Flash	4	<1 F ² or >100 Gbit/cm ²	1–10 μs	1 GiB/s	10–20 yr	10–100 TiB
HDD	0.1	100 Gbit/cm ²	5–20 ms	100 MiB/s	10–100 yr	10–100 TiB

- ▶ General context
- ▶ **State of the art on different technologies and implementations**
 - Taxonomy
- ▶ CSRAM: our own design approach
- ▶ Software platforms and explorations
- ▶ Results and analysis on different memory computing architectures
- ▶ Conclusion and future works

Memory Bandwidth & Size



Interface levels

bit-cells array
(custom bitcell)

closer to **Memory**



closer to **Processor**

Computation complexity

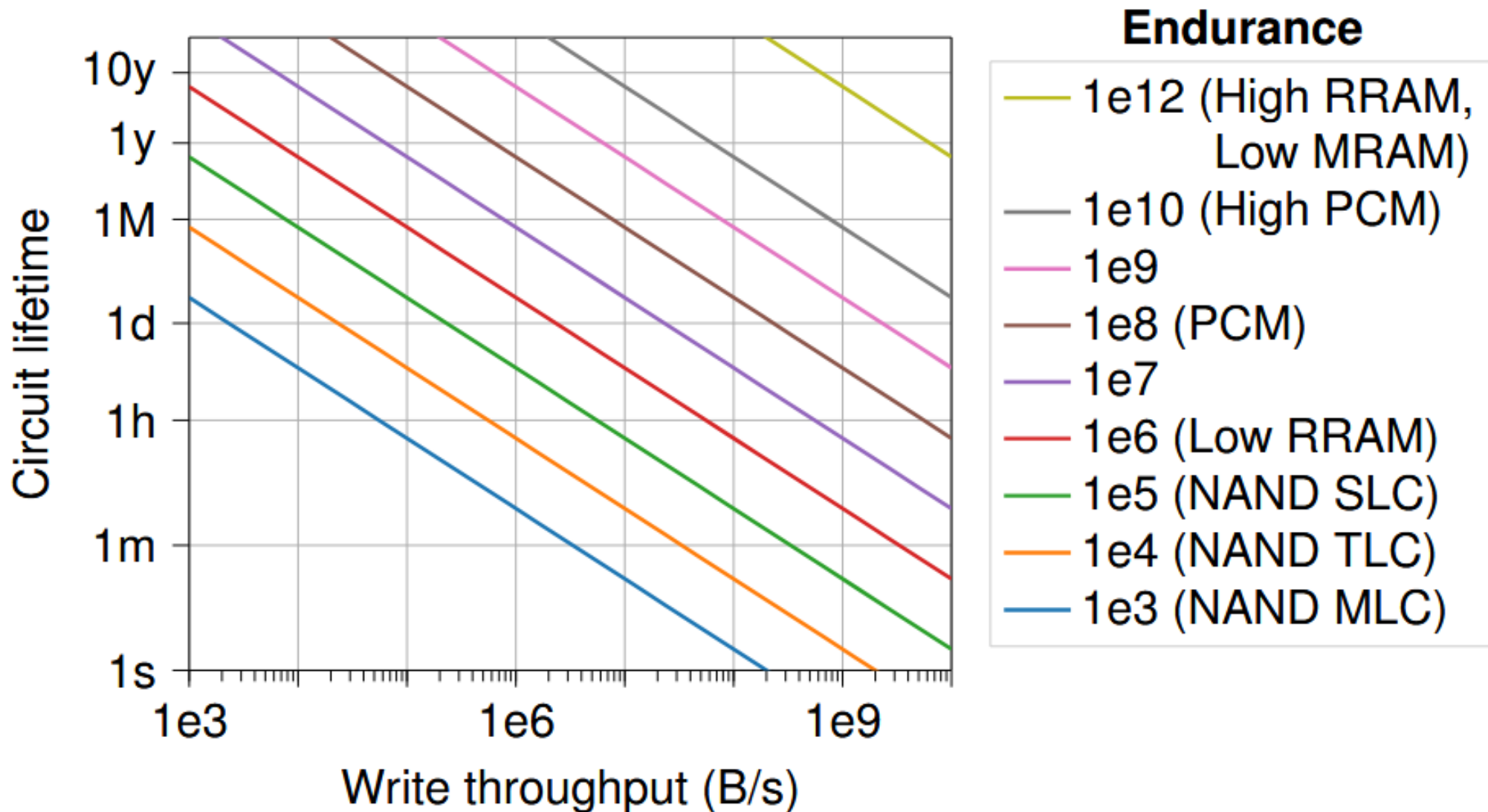
nor,
nand

Limited
enduranceHigh power
consumptionData is not
naturally here

	Memory Technology	Architecture	Key Results
[1] ReRAM CIM	<input checked="" type="checkbox"/> Circuit is demonstrated physically with ReRAM	<input checked="" type="checkbox"/> Standalone circuit	<input checked="" type="checkbox"/> No mention of endurance <input checked="" type="checkbox"/> Logic operations only
[2] NAND Flash CIM	<input checked="" type="checkbox"/> Analogic computation (cu	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/> No write to Flash 40 TOPS/W on MAC
[3] Ambit	<input checked="" type="checkbox"/> DR	<input checked="" type="checkbox"/>	Logic operations only 32x speedup 35x energy saving on database
[4] UPMEM	<input checked="" type="checkbox"/> Har logi with	<input checked="" type="checkbox"/>	No vector processing 20x speedup 10x better power efficiency
[5] Compute Caches	<input checked="" type="checkbox"/> SRAM with standard bitcell	<input checked="" type="checkbox"/> HPC Architecture <input checked="" type="checkbox"/> Stay in cache <input checked="" type="checkbox"/> Not concerned about DRAM consumption	<input checked="" type="checkbox"/> Logic operations <input checked="" type="checkbox"/> 1.9x speedup <input checked="" type="checkbox"/> 2.4x energy reduction
[6] Blade	<input checked="" type="checkbox"/> SRAM with non standard organization	<input checked="" type="checkbox"/> MCU Architecture	<input checked="" type="checkbox"/> Logic operation and addition <input checked="" type="checkbox"/> 6x speedup vs NEON

- No paper with full architecture & data coming from SCM
- No paper with a mix of SRAM & SCM to preserve the latter's endurance
- Analog computation in SCM gives very few freedom in operators and has lots of constraints

Endurance is the limiting factor in NVMs

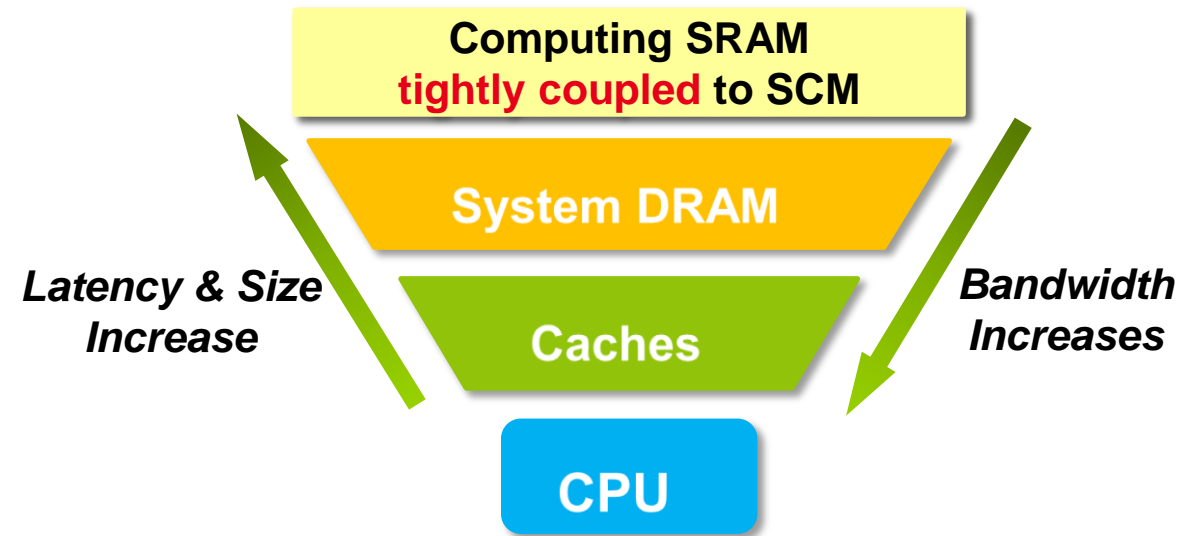


► **Merge computing and data in the same circuit**

- ~~Low energy efficiency and slow~~ high speed
- Endurance is a ~~problem~~ (10^9)
- Use existing row buffer and repurpose it for computing

► DRAM is still used as a write buffer

► Caches can be reduced (Instructions & stack)

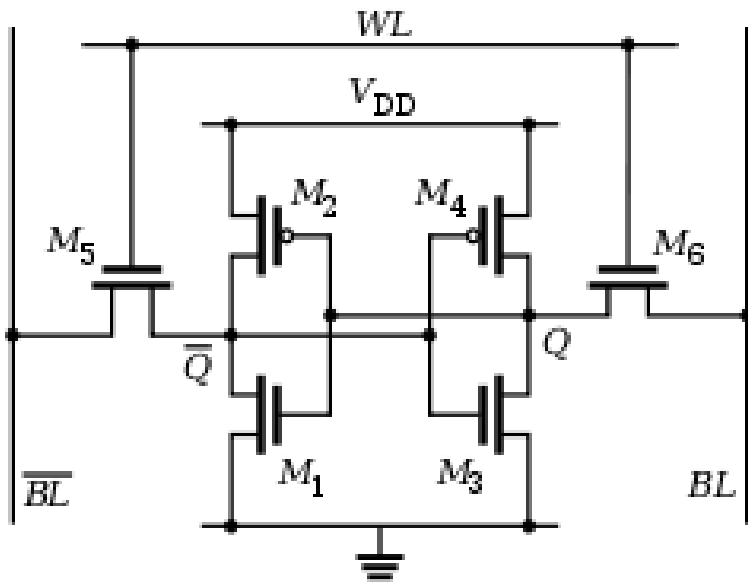


Let's move computing to the top with tightly coupled SRAM to get the best of both worlds!

- ▶ General context
- ▶ State of the art on different technologies and implementations
- ▶ **CSRAM: our own design approach**
 - SRAM based design
 - Automated design approach
 - Comparison with state of the art
- ▶ Software platforms and explorations
- ▶ Results and analysis on different memory computing architectures
- ▶ Conclusion and future works

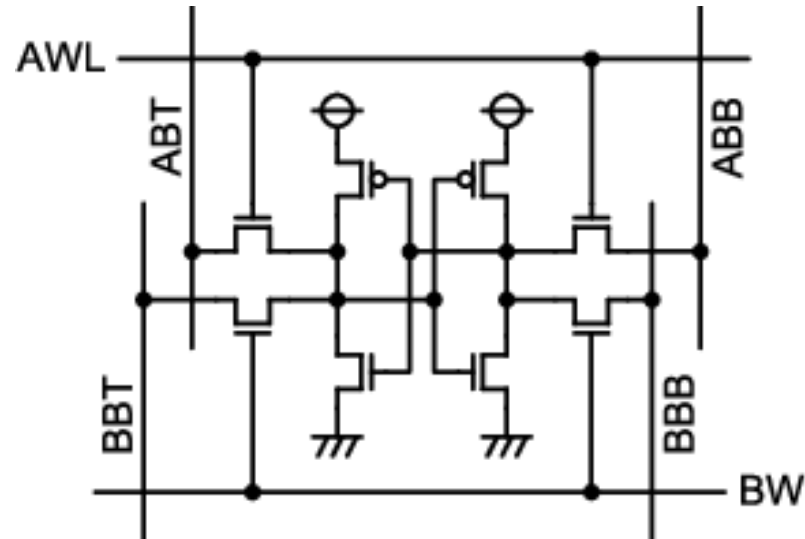
► Single Port (1RW)

- 6 transistors



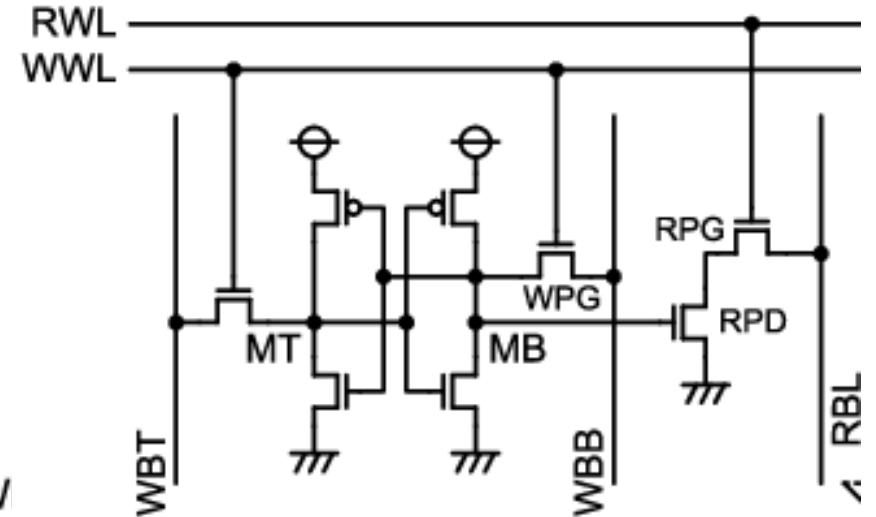
► Dual Port (2RW)

- 8 transistors



► Two Port (1R1RW)

- 8 transistors

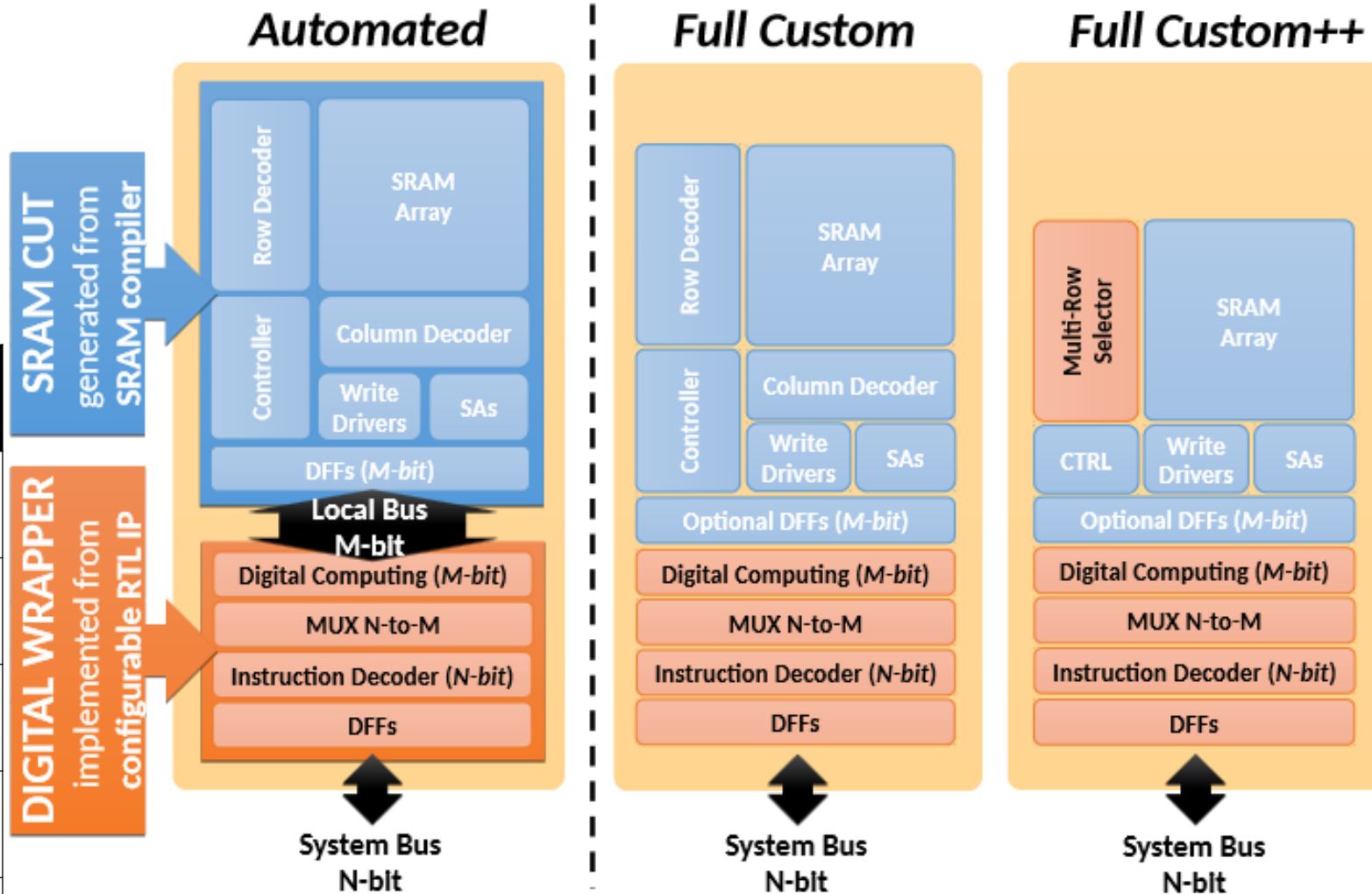


Legend:

MEMORY PART

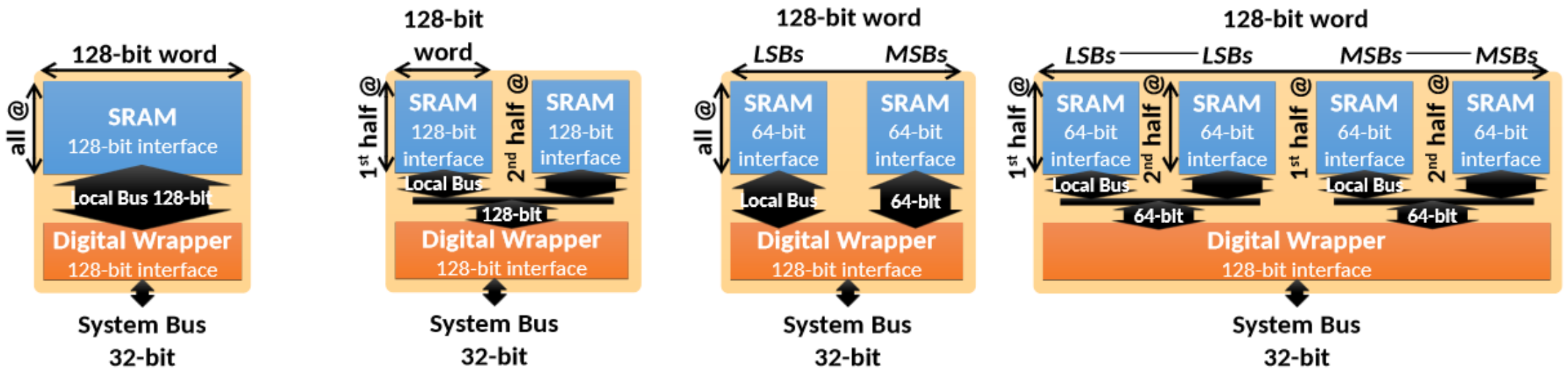
COMPUTING PART

Design Methodology	Full Custom ++
Computing classification	NMC/IMC
Computing area overhead	low
Form factor Flexibility	low
Development effort*	very high
Silicon qualification effort	very high



► Form factor flexibility

C-SRAM Macro = *Memory Cut Partitioning w/ Digital Wrapper*



a) WORD & BIT # < physical limits

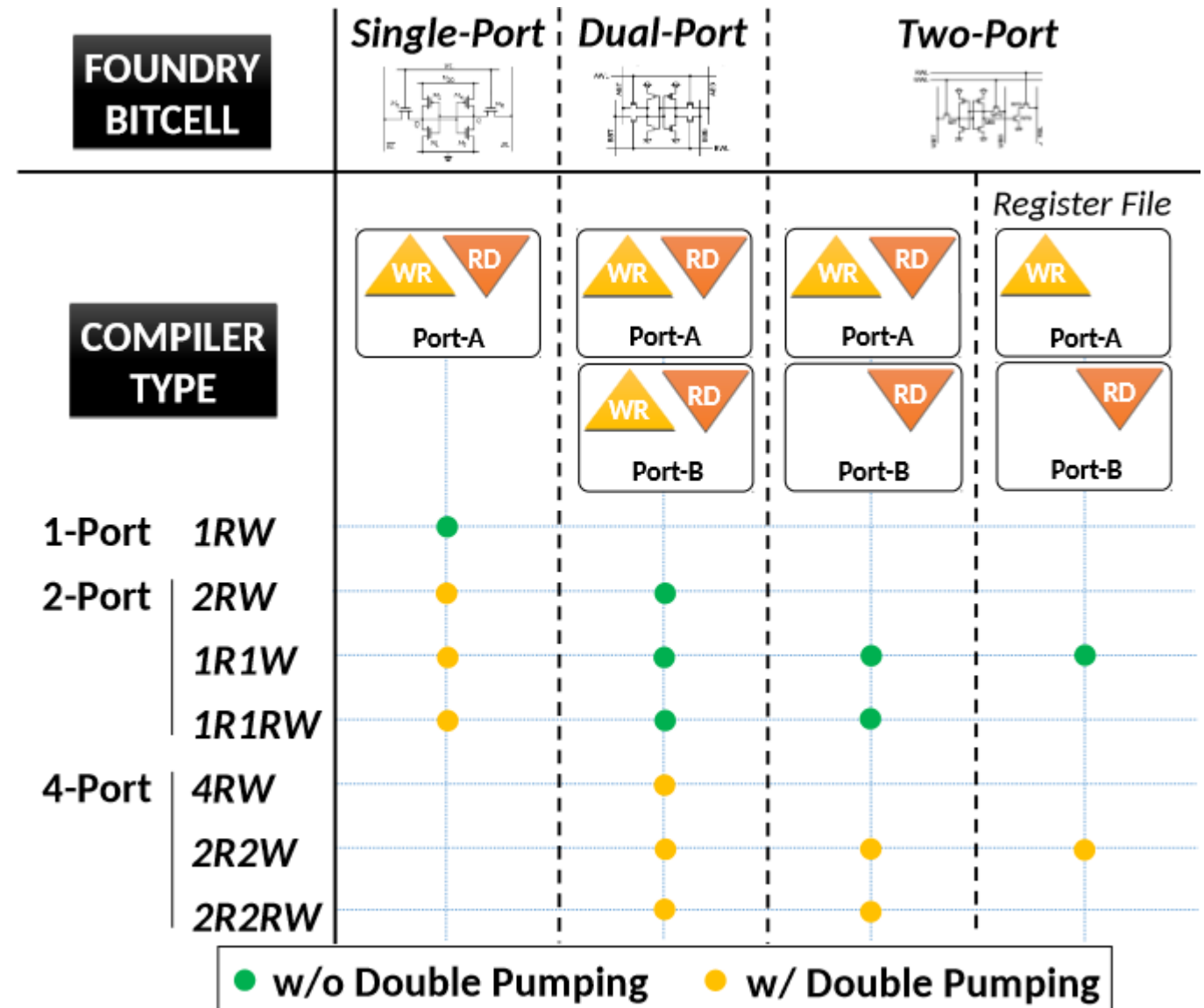
b) WORD # > physical limits

c) BIT # > physical limits

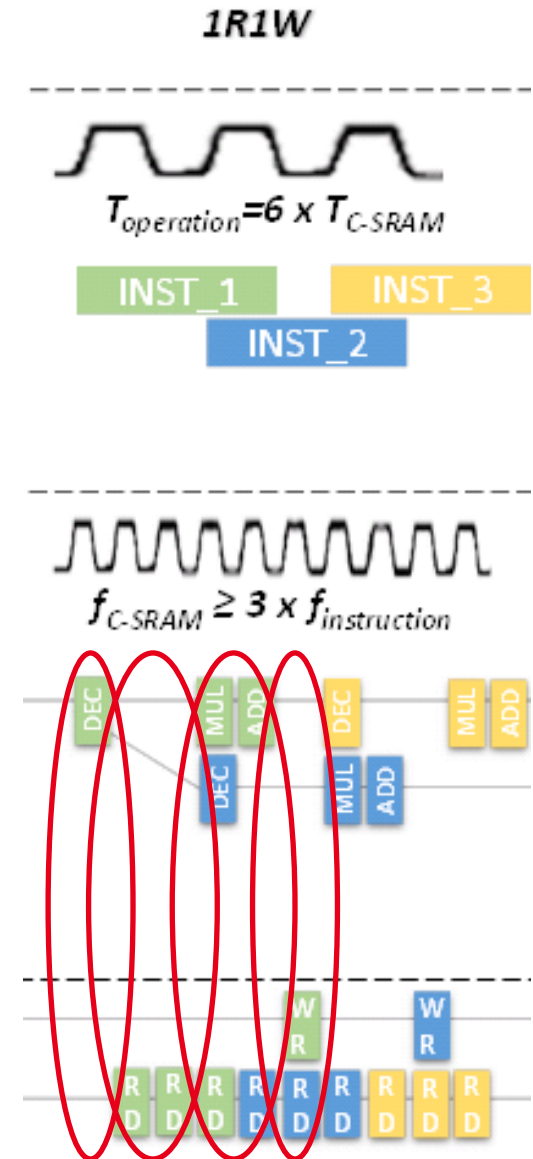
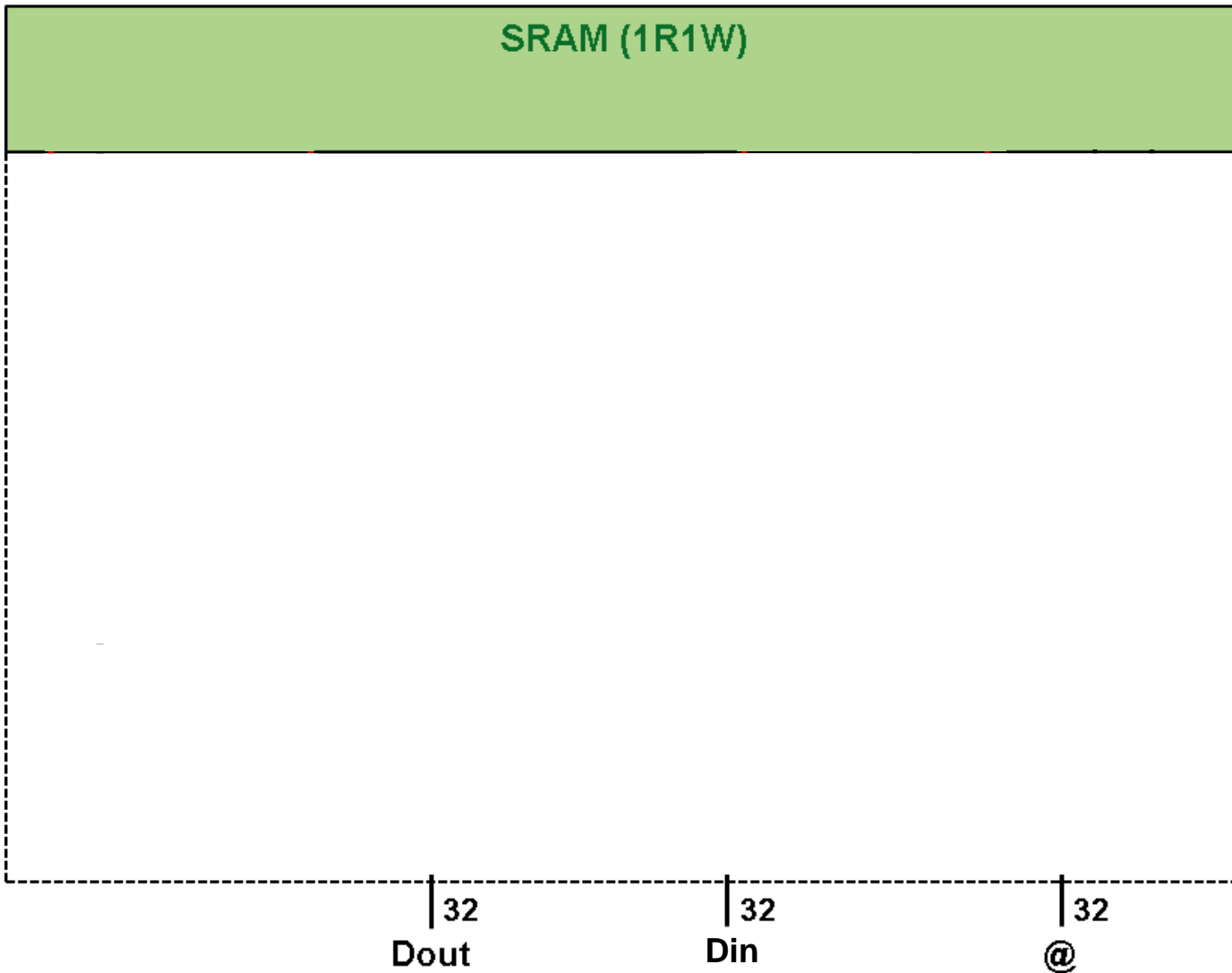
d) WORD & BIT # > physical limits

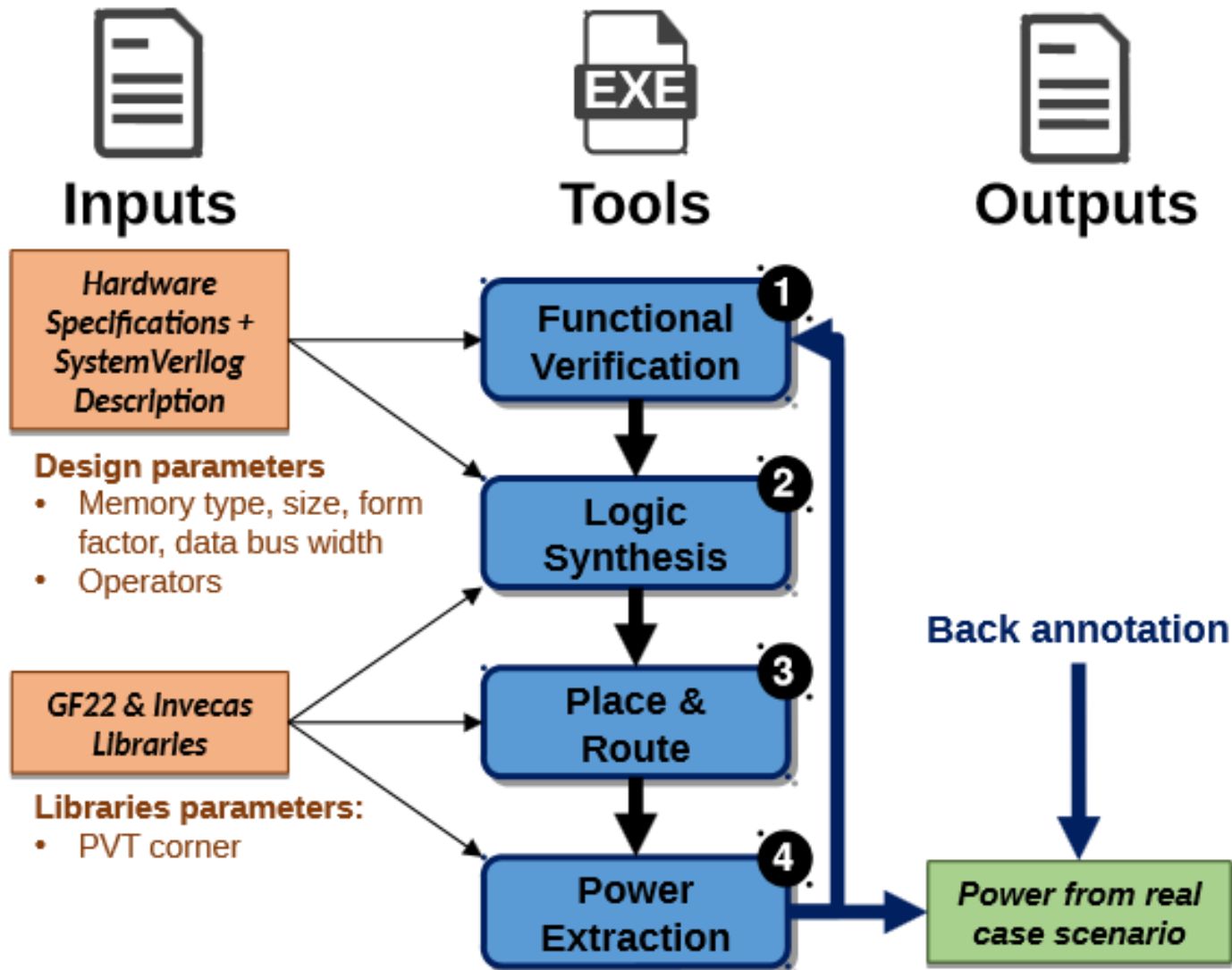
► Double pumping

- Makes memory active on both rising and falling edges
- Slightly decreases frequency



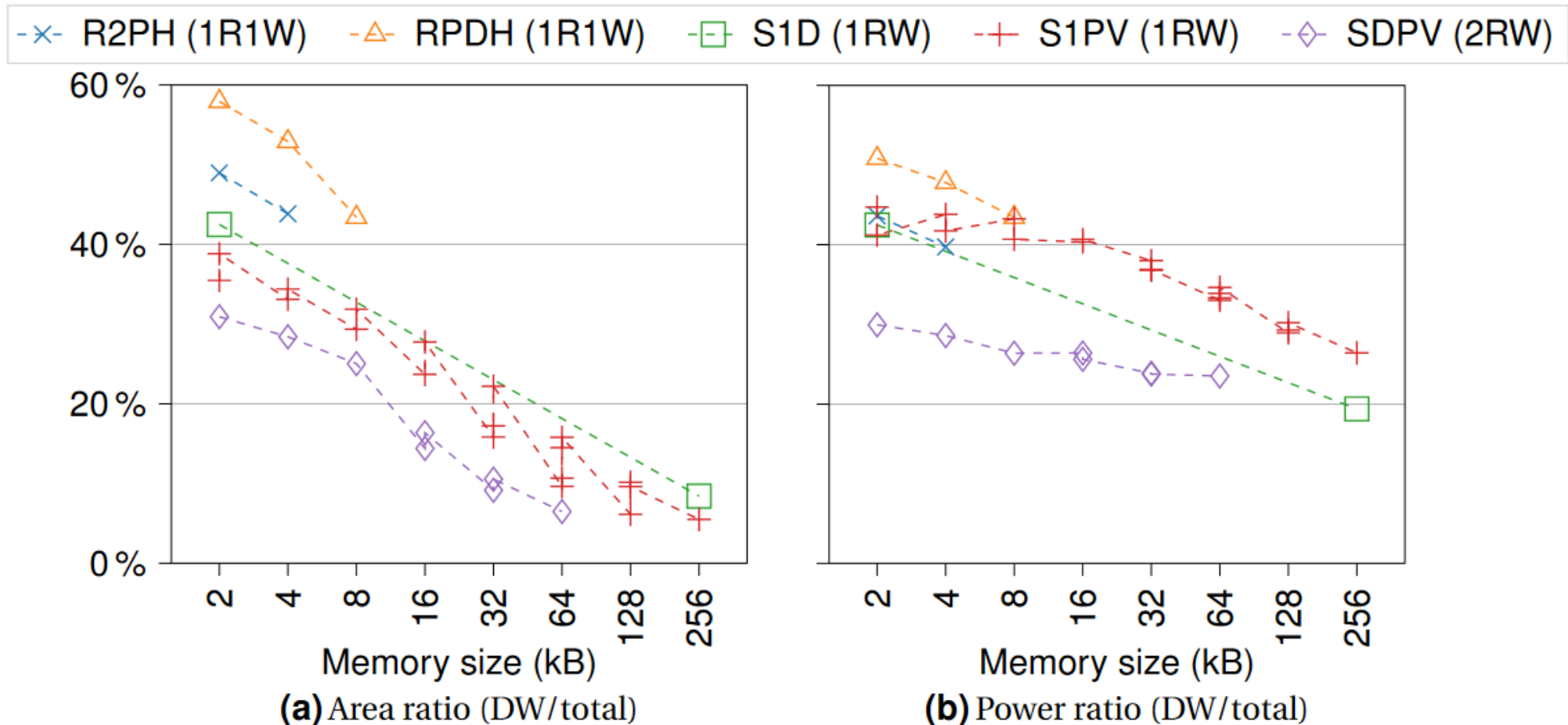
- ▶ **To be generic, we need 8, 16 and 32 bits operators**
- ▶ **Logical operators**
 - AND, OR, XOR, NOT and their complement → 8 instructions
- ▶ **Advanced logical operators**
 - Shift Right, Shift Left and Broadcast → 3 sizes × 3 operations = 9 instructions
- ▶ **Arithmetic operators**
 - Addition, Subtraction, Comparison and Shift Right Arithmetic
→ 3 sizes × (3 operations + 6 comparisons) = 27 instructions
- ▶ **Advanced arithmetic operators (8 bits only)**
 - Multiplication and Multiply-And-Accumulate (MAC)
→ 3 multiplications for different signedness and 1 MAC = 4 instructions
- ▶ **Other operators**
 - Horizontal swap = 2 instructions

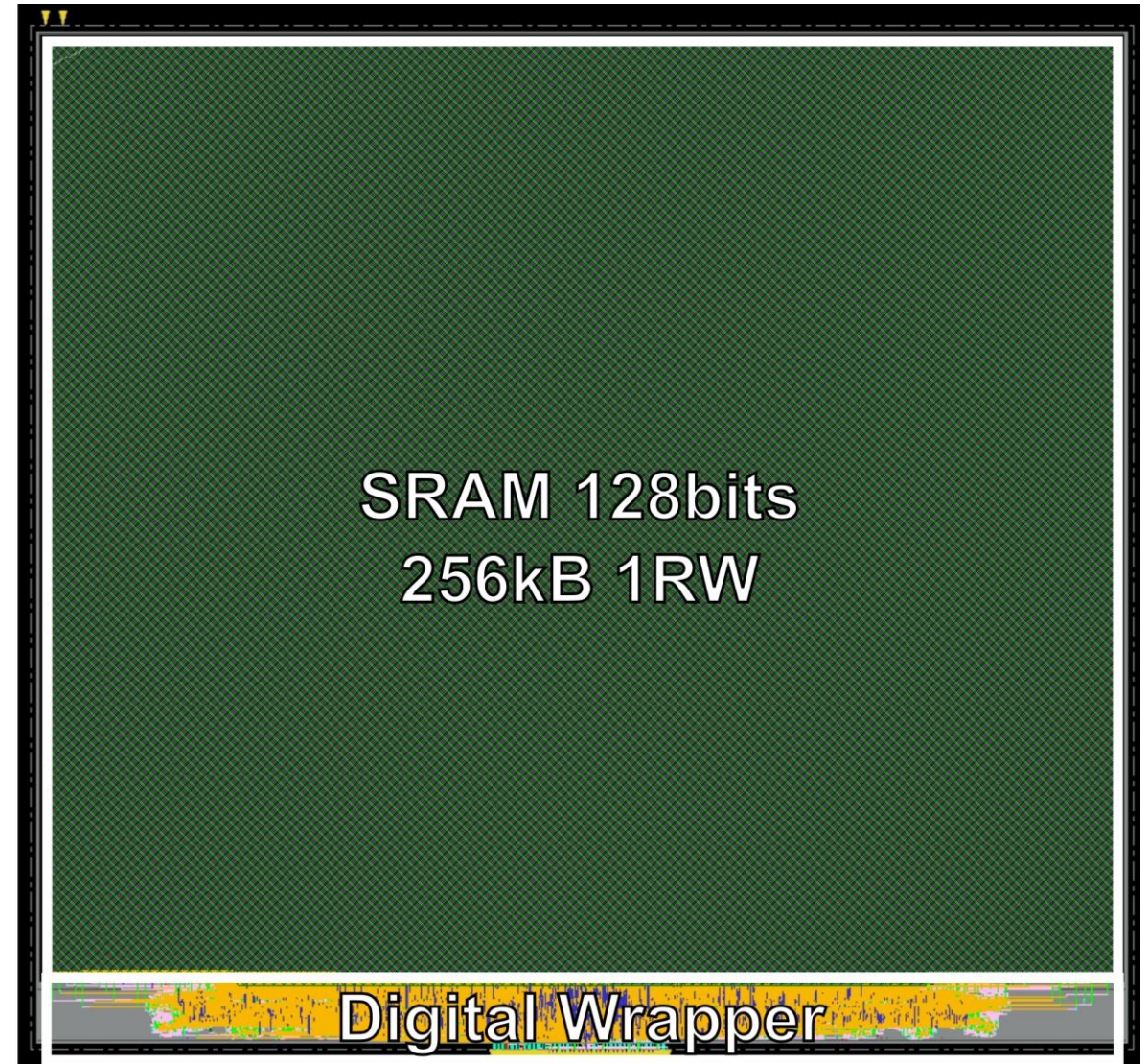
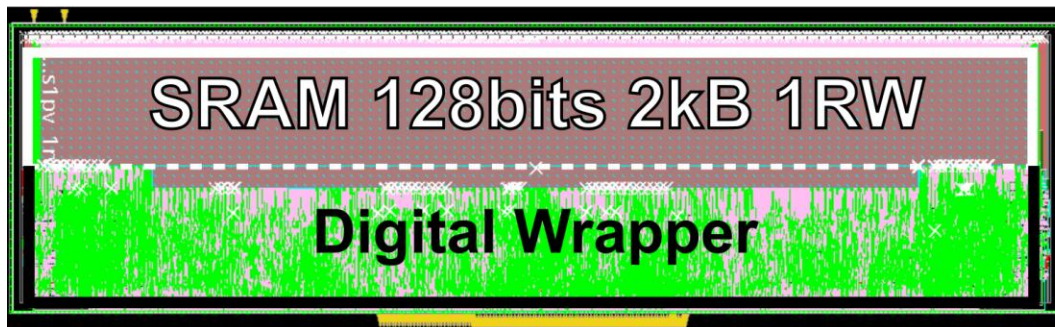
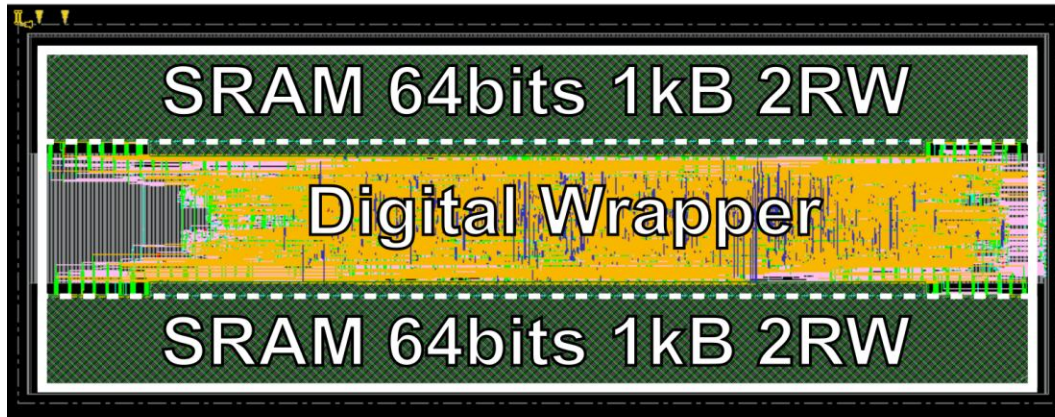
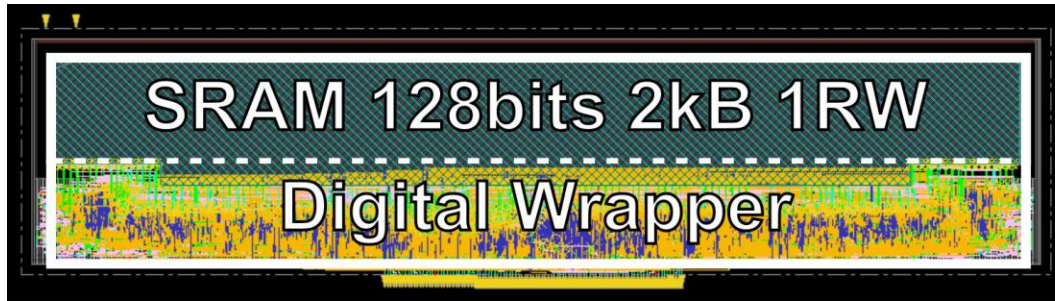




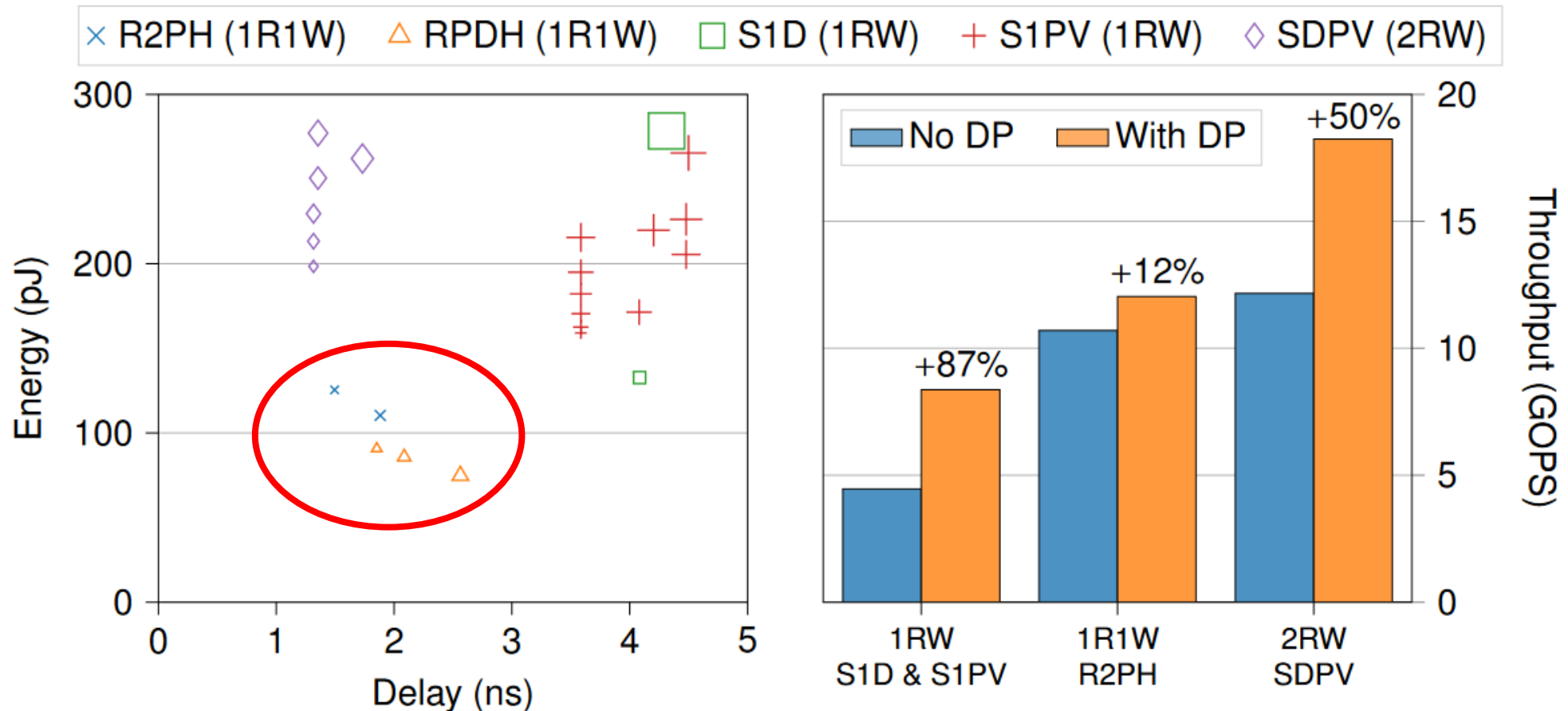
	Tool	Vendor	Version
①	ModelSim	MentorGraphics	10.5c
②	DesignCompiler	Synopsys	n-2017.09-sp5
③	Innovus	Cadence	19.13
④	Star-RC	Synopsys	o-2018.06-sp2
	PrimeTime Power	Synopsys	o-2018.06-sp5

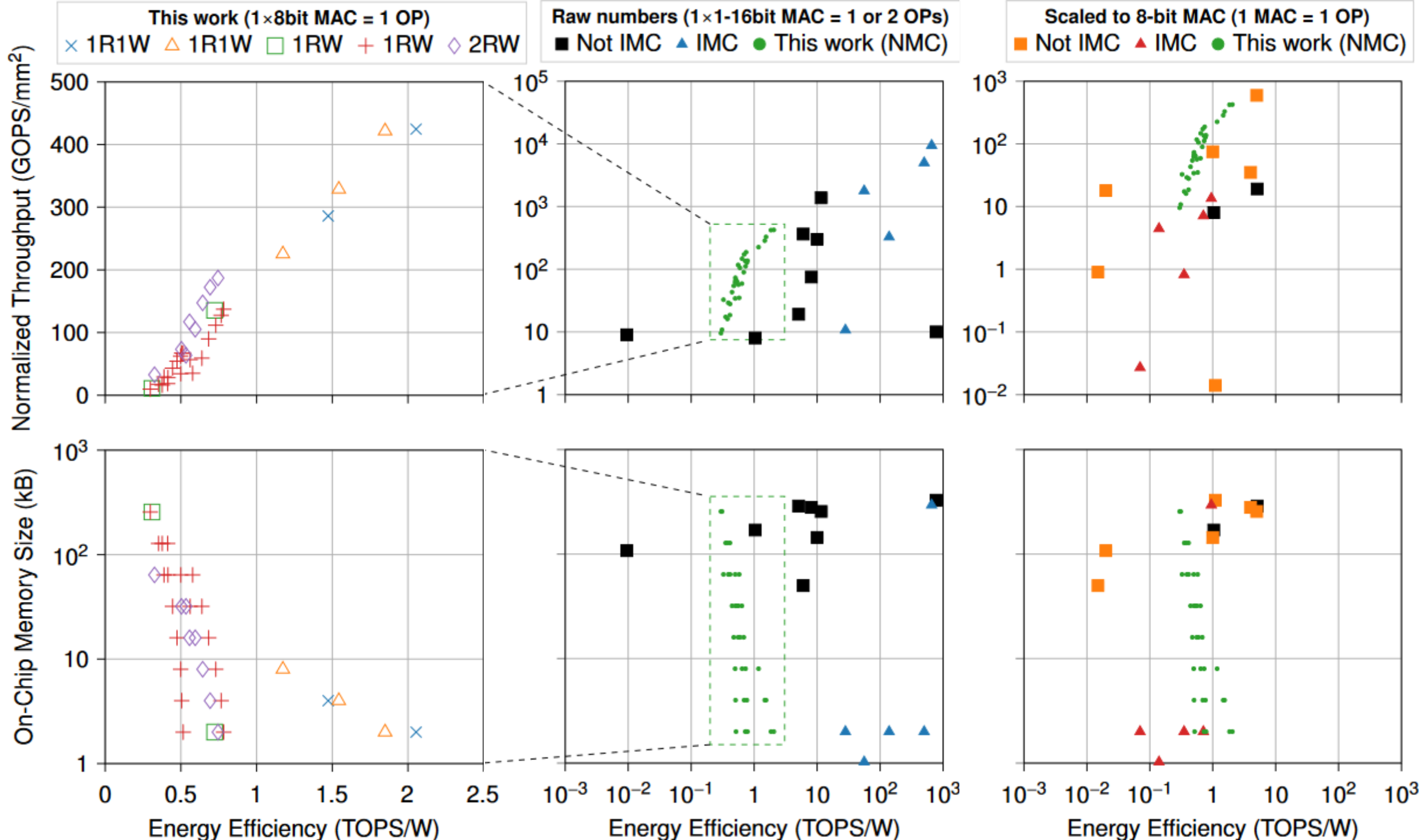
- ▶ Area overhead ranging from almost 60% for 2kB memories to 5% for >128kB memories
- ▶ Power overhead decreasing more gradually from 50% down to 20%





- ▶ Pareto front for optimal in either energy or delay
 - 1R1W memories are the best choice but limited to small size (8kB max)
- ▶ Doublepump grants +87% throughput for 1RW memories





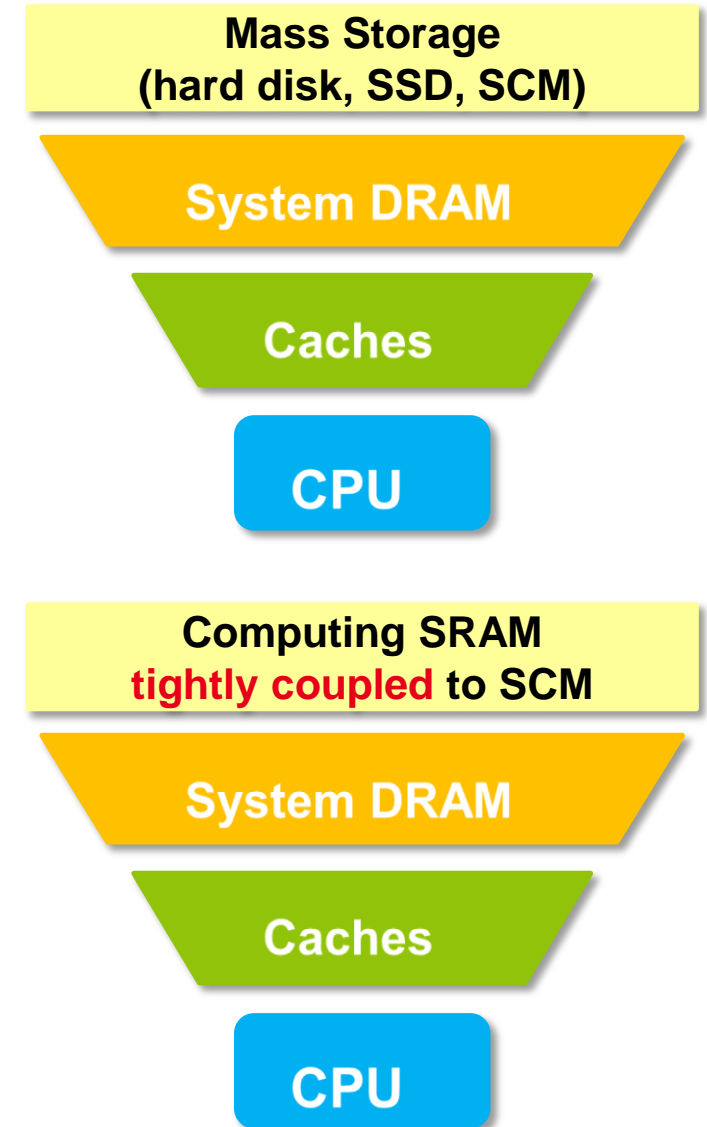
- ▶ **SRAM based with automated approach**
 - Fast and versatile hardware development
 - Standard pipeline design
 - Acceptable area and power overhead
- ▶ **Better than analog IMC implementation**
 - More throughput and efficiency on 8-bit MAC
 - More memory size
- ▶ **This work has been published in DATE 2020**

- ▶ General context
- ▶ State of the art on different technologies and implementations
- ▶ CSRAM: our own design approach
- ▶ **Software platforms and explorations**
 - Architecture evaluation
 - Memory access cost evaluation
 - Workflow
- ▶ Results and analysis on different memory computing architectures
- ▶ Conclusion and future works

► Need to monitor all memories accesses

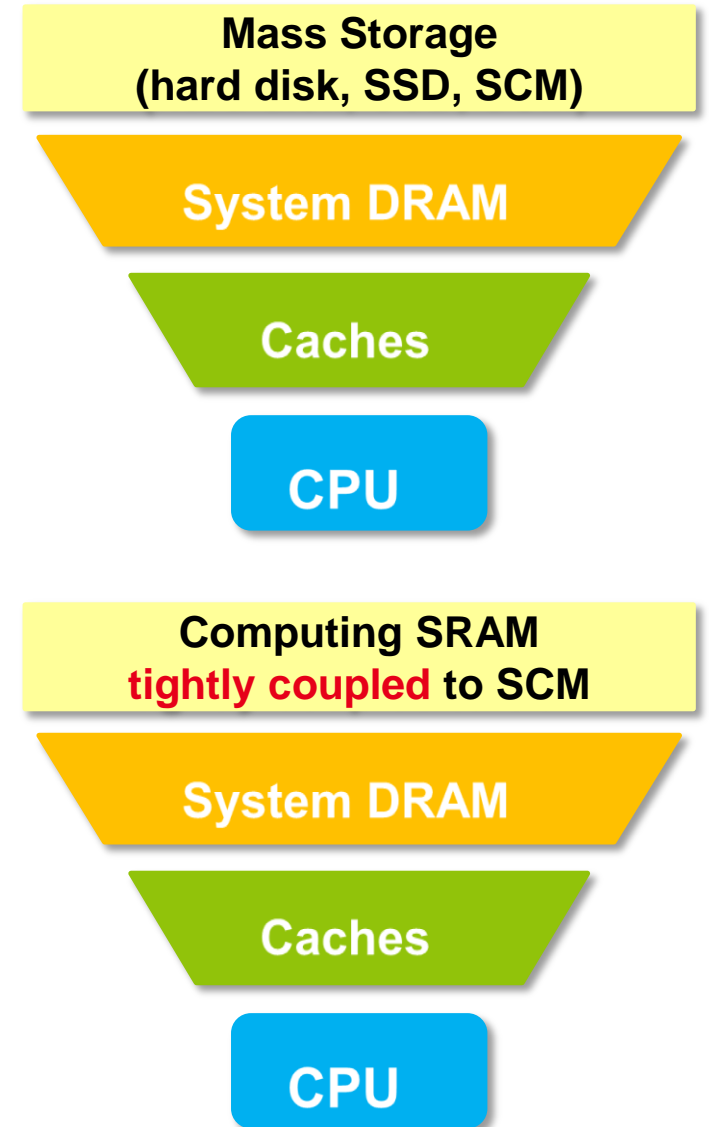
- Analytic models cannot describe exchange CPU ↔ CSRAM
- Hardware counter not reliable for L3, DRAM and above
- Use existing simulators such as gem5 or QEMU

Platform	Memory interfaces	Accuracy	Speed	Development effort
QEMU	Callback functions	Instruction	High	High
gem5	Ruby	Cycle	Slow	High
ArchSim	TLM sockets	Transaction	Medium	Moderate
ZSim	None	Instruction	Very high	Moderate
Sniper	N/A	Cycle	Slow	Moderate
LLVM	None	Instruction	Medium	Low
Pin	Callback functions	Instruction	High	Low

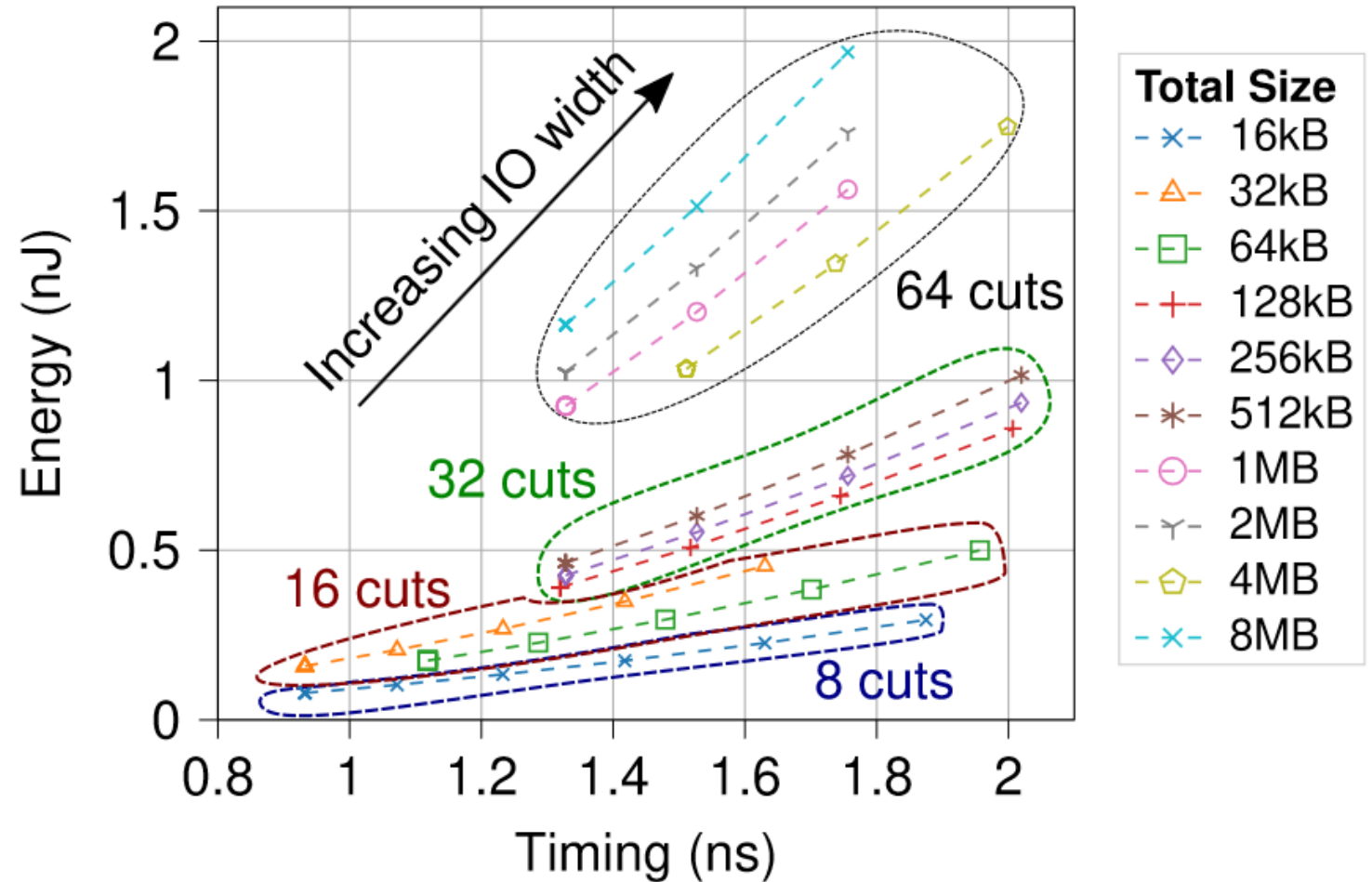


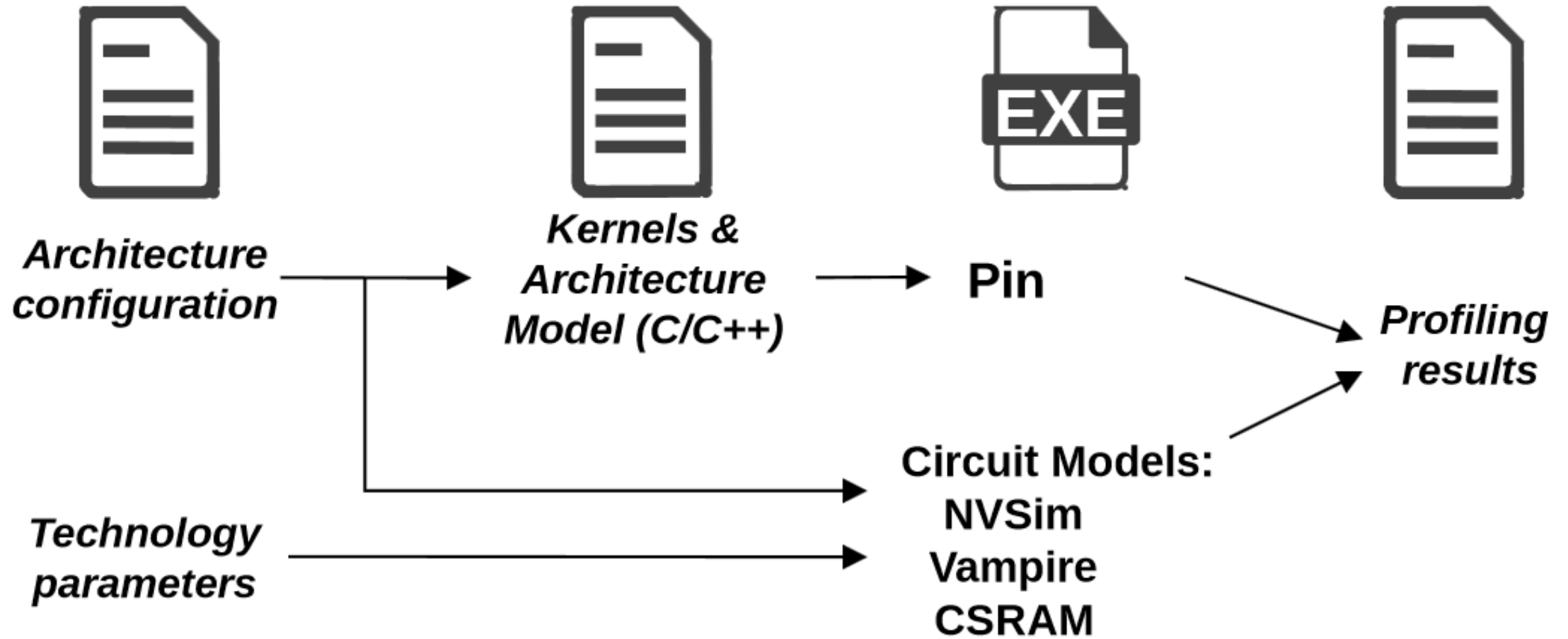
► Use specialized hardware model tools

- NVSim (based on cacti) for caches and NVM
 - PCM based SCM with technology parameters from state of the art
- DRAM evaluation is tricky
 - Calibration done with VAMPIRE
- CSRAM evaluation made in-house



- ▶ L3 cache like size to contain working dataset
- ▶ Tiling and widening 1RW memories
 - From 16B wide IO up to 4kB
 - From 2kB to 256kB total size up to 8MB

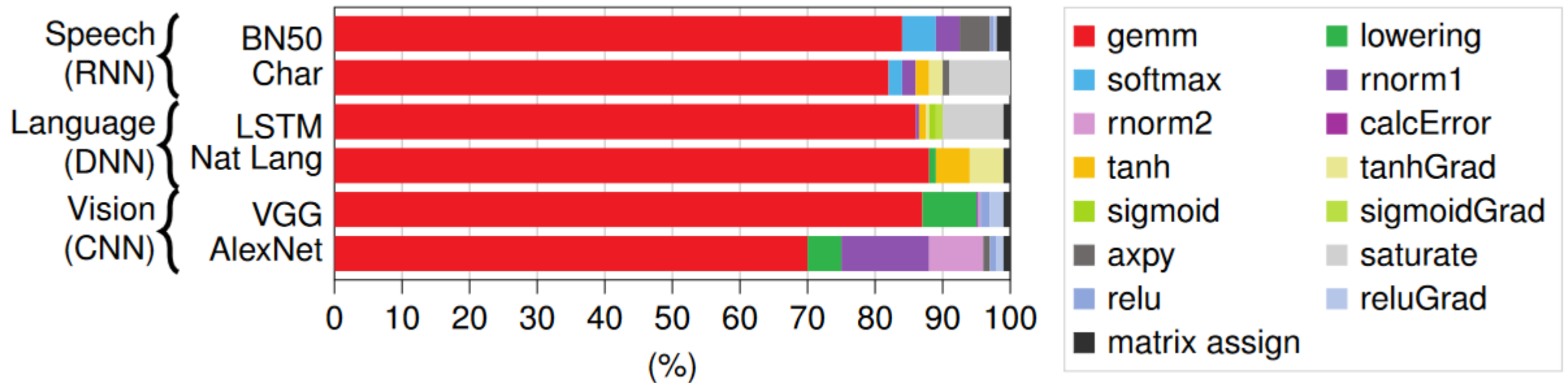




- ▶ General context
- ▶ State of the art on different technologies and implementations
- ▶ CSRAM: our own design approach
- ▶ Software platforms and explorations
- ▶ **Results and analysis on different memory computing architectures**
 - Benchmarks used
 - Computing at the top
- ▶ Conclusion and future works

► Big data era with neural networks and sensors everywhere

- Neural networks mostly use matrix multiplication
- Convolution can be turned into matrix multiplication through *matrix assign* operation



▶ Linear kernels

- Hamming weights \rightarrow information theory and also BNN with popcnt operation
- shift-or \rightarrow pattern matching in DNA sequence

Compute bound

- AXPY ($y = ax + y$) \rightarrow Basic linear algebra, very high disk access

▶ Quadratic kernels

- atax ($y = A^T(Ax)$) \rightarrow linear solver and signal processing
- gesummv ($y = \alpha Ax + \beta Bx$) \rightarrow signal processing

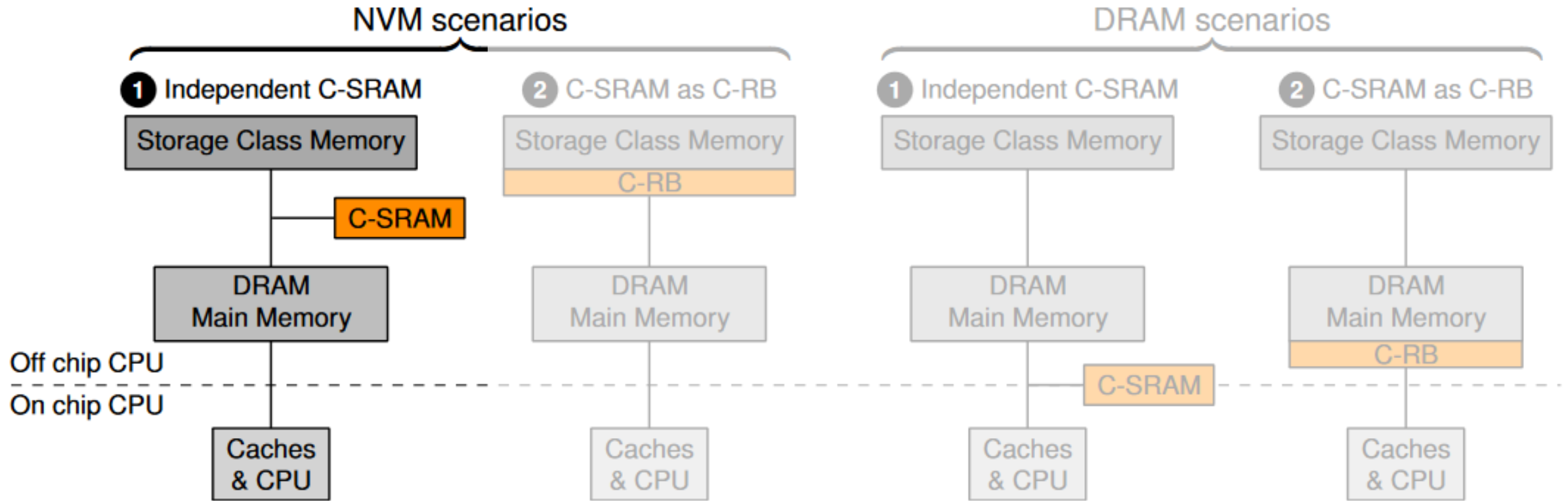
▶ Cubic kernels

- gemm ($C = \alpha AB + \beta C$) \rightarrow matrix multiplication, image processing, neural networks

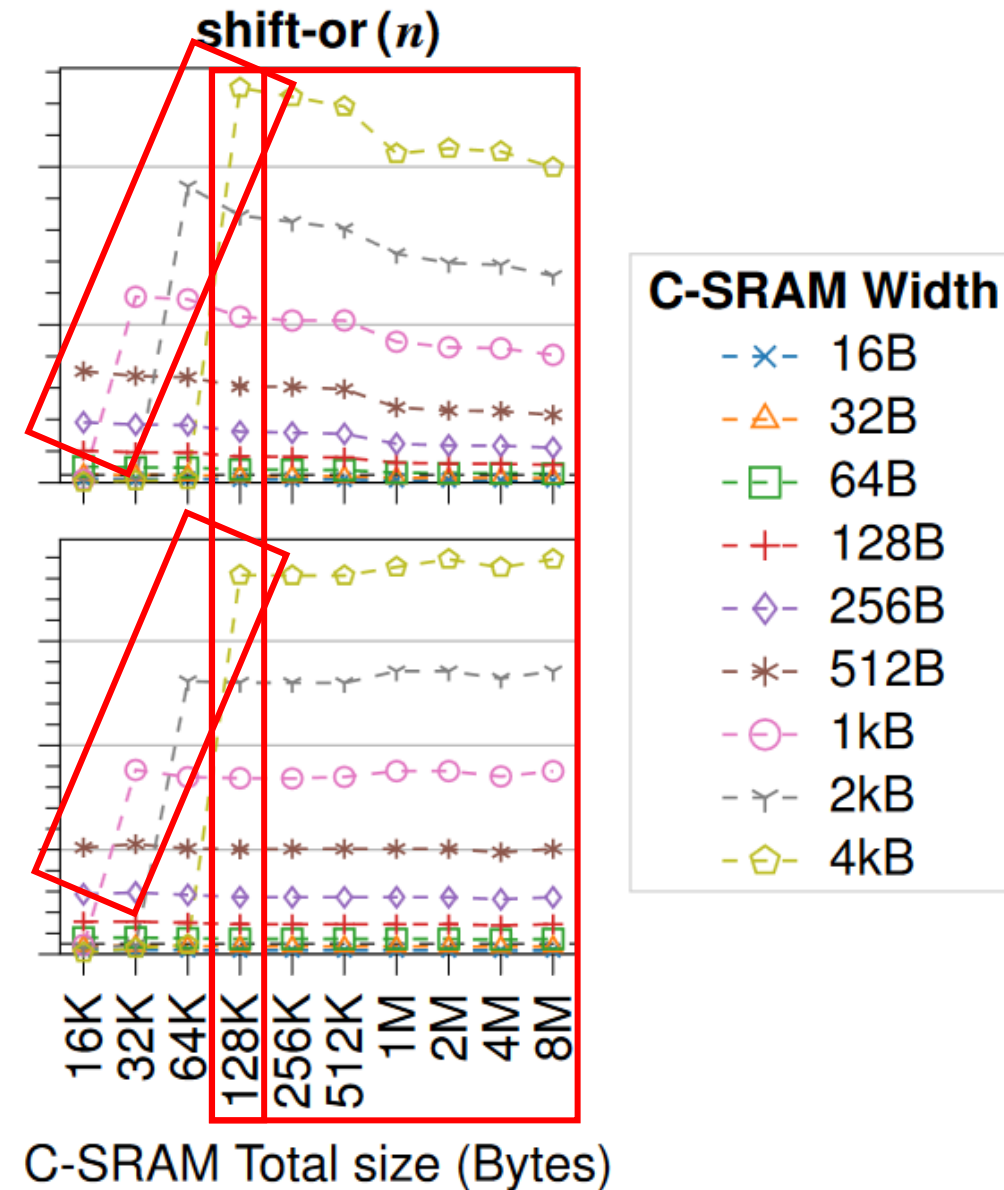
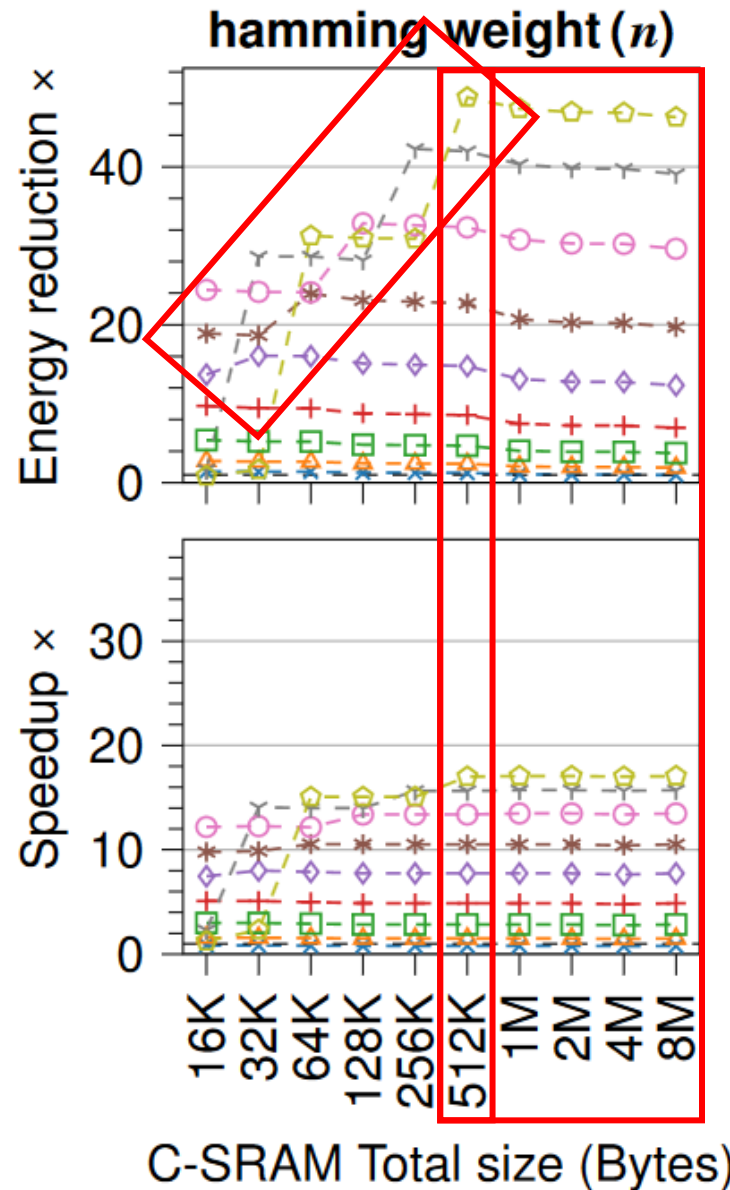
▶ Real case application

- darknet \rightarrow implementation of CNN for image classification

Memory bound



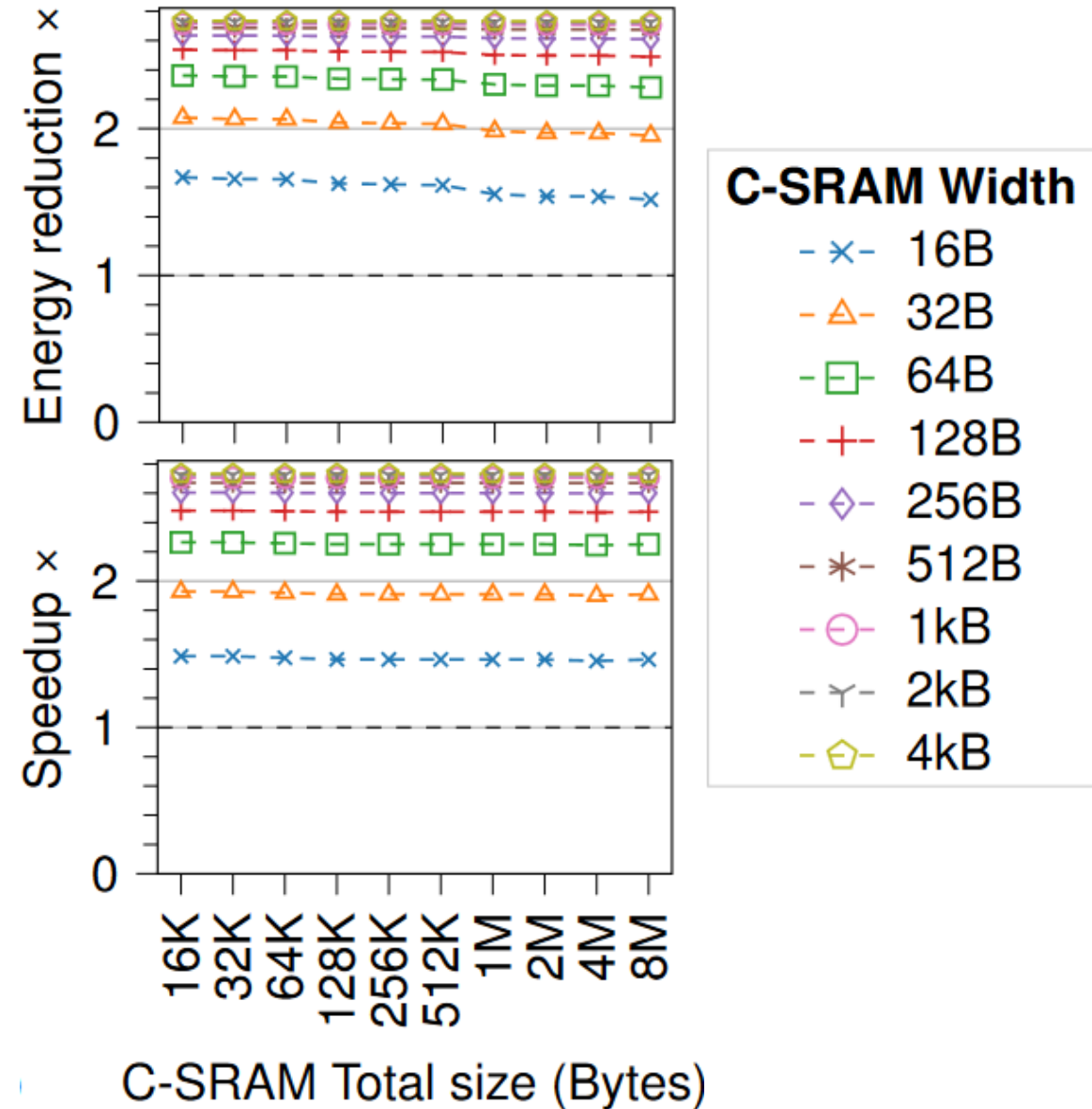
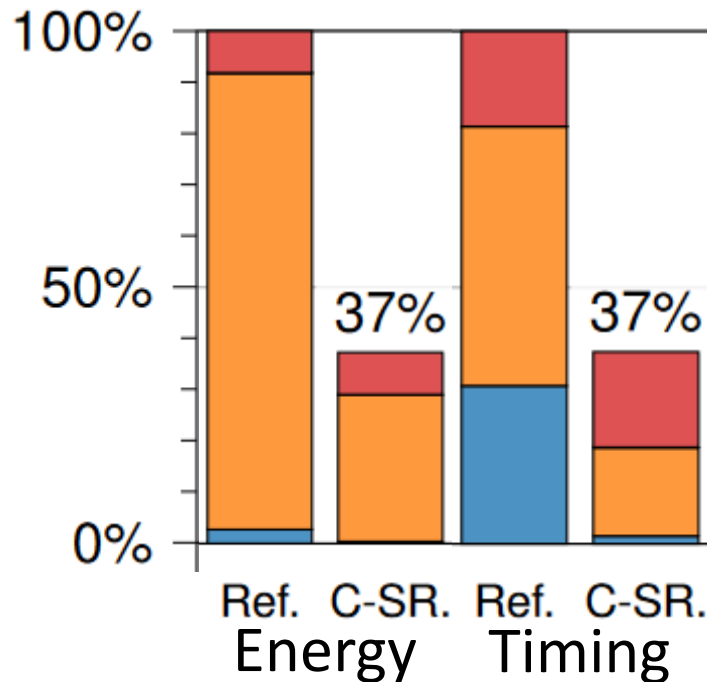
- ▶ Optimal size between 128kB and 512kB
- ▶ Larger vector width leads to more improvements
- ▶ Larger than optimal total size decreases energy reduction due to more costly accesses
- ▶ Stair like shape due to minimal working dataset size



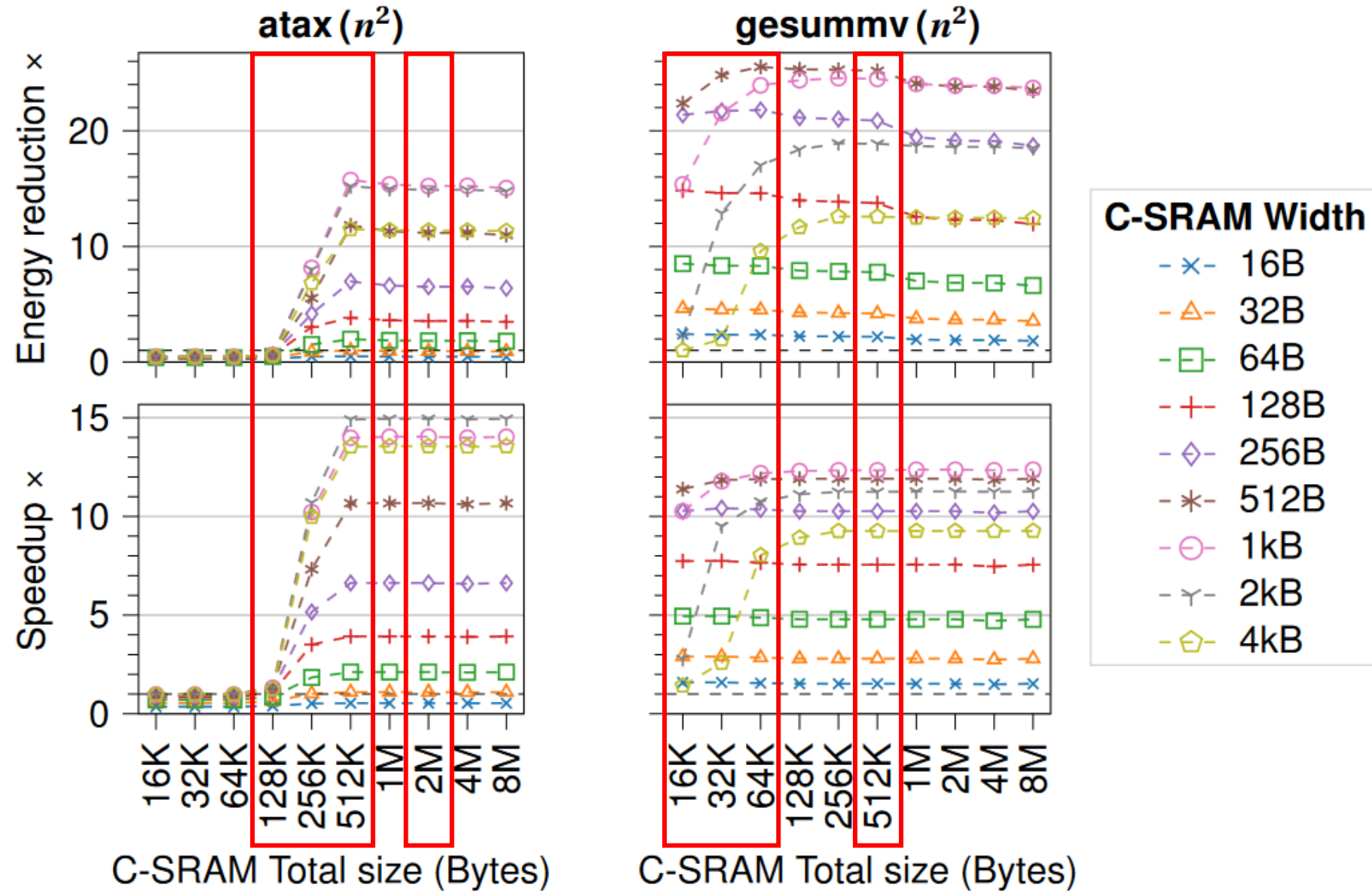
► Total size has almost no impact

- Working dataset is just 4 data

► Vector width shows very fast saturation

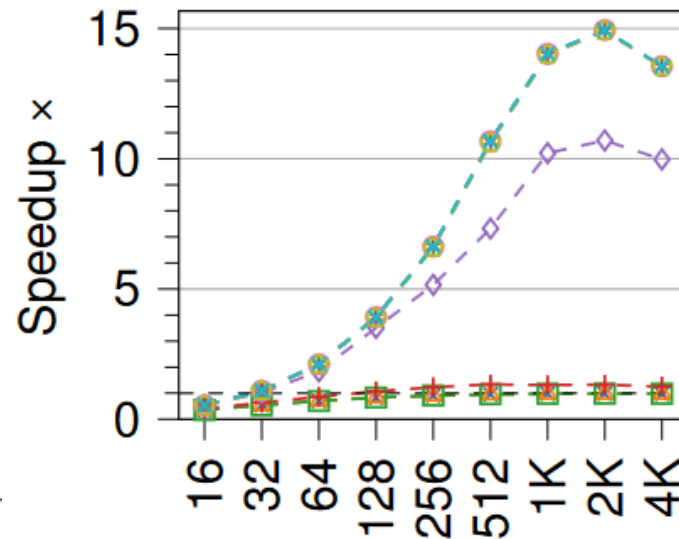
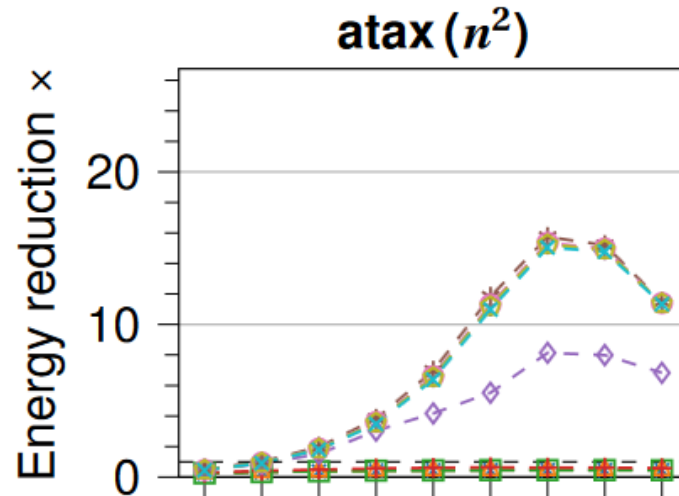
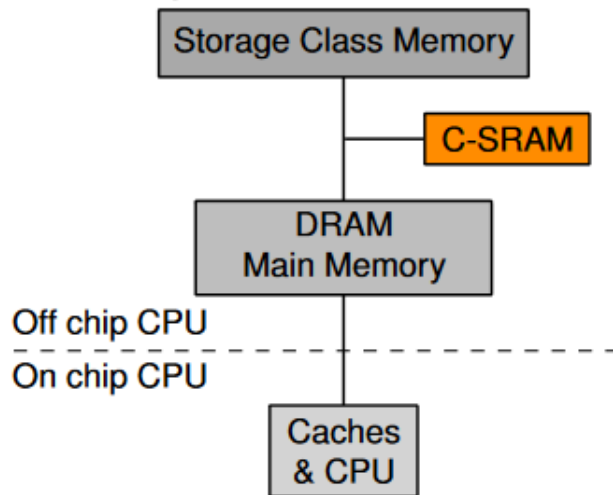


- Atax needs minimal total size of 512kB
- Gesummv benefits from small memories
- Both kernels, larger vector width has negative impact
 - Reduction loop performed by the CPU

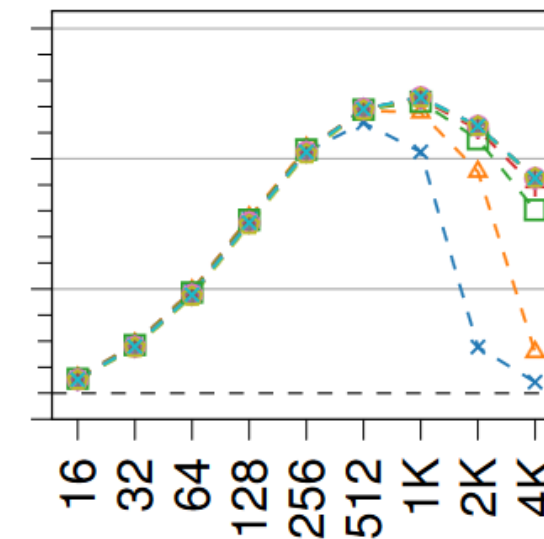
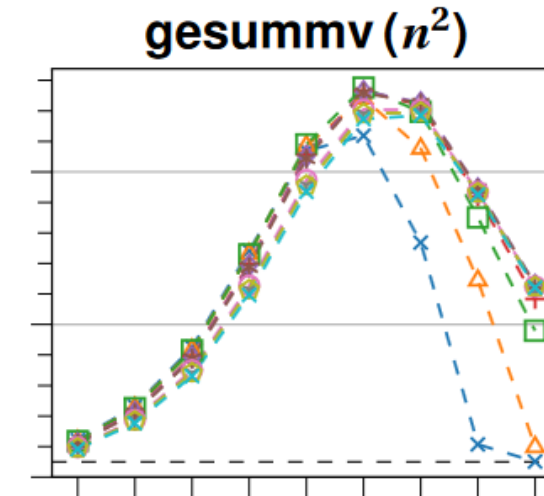


- ▶ Atax needs minimal total size of 512kB
- ▶ Gesummv benefits from small memories
- ▶ Both kernels, larger vector width has negative impact

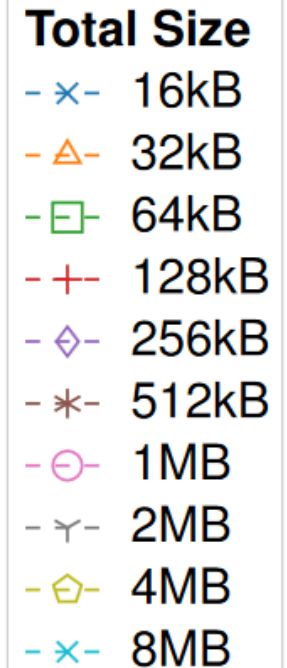
- Reduction loop performed by the CPU



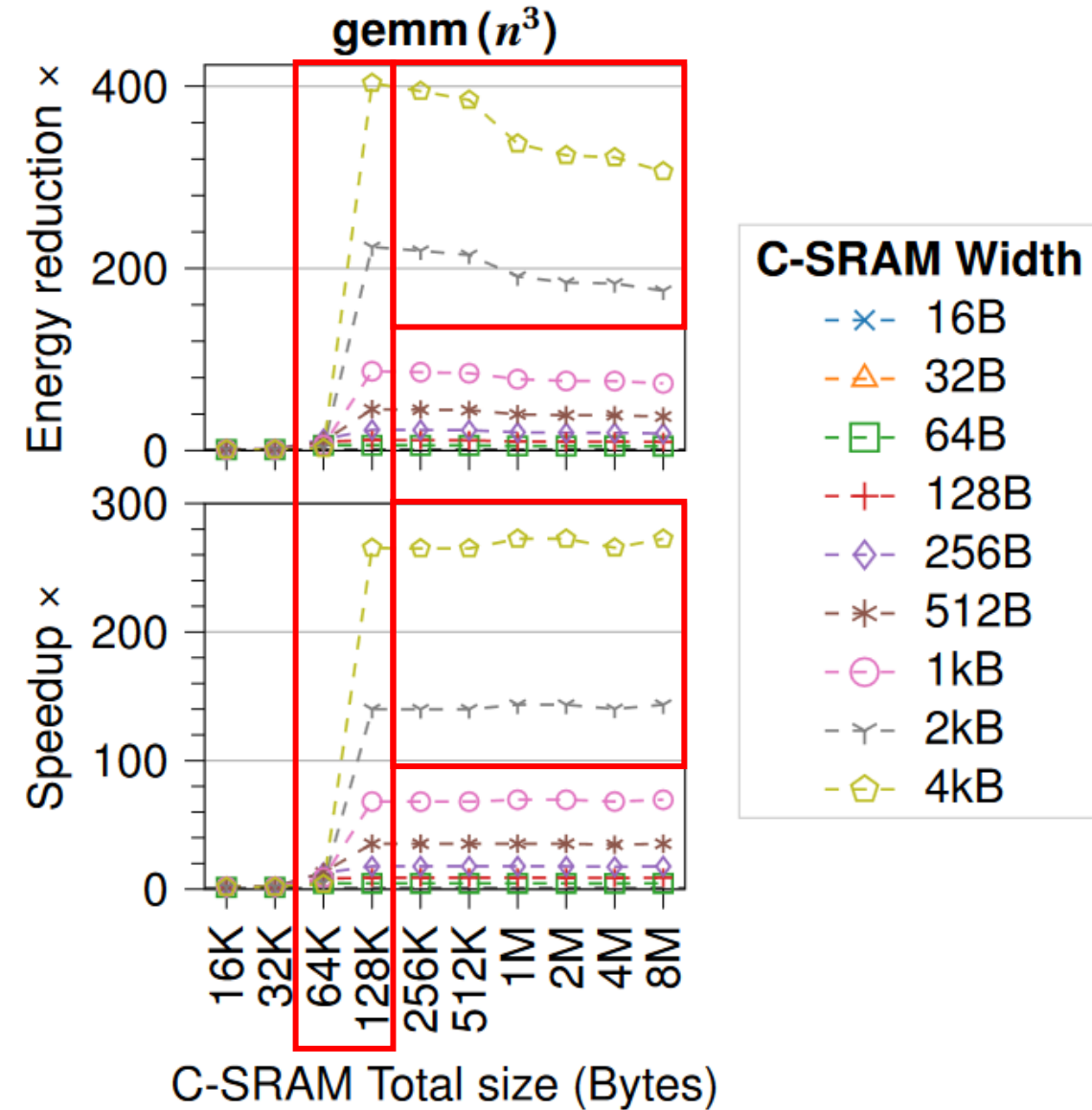
C-SRAM Vector Width (Bytes)

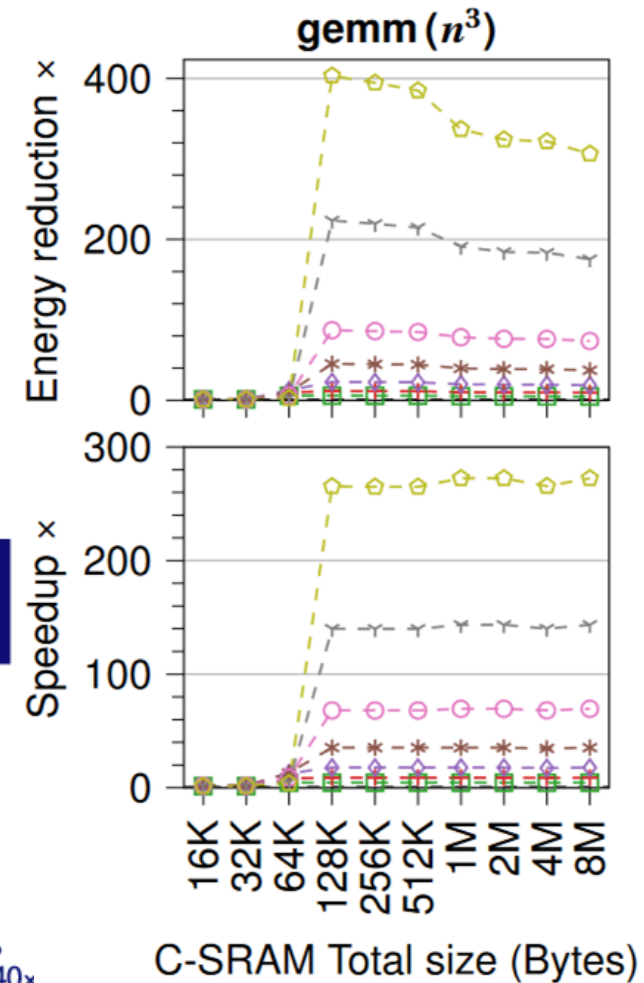
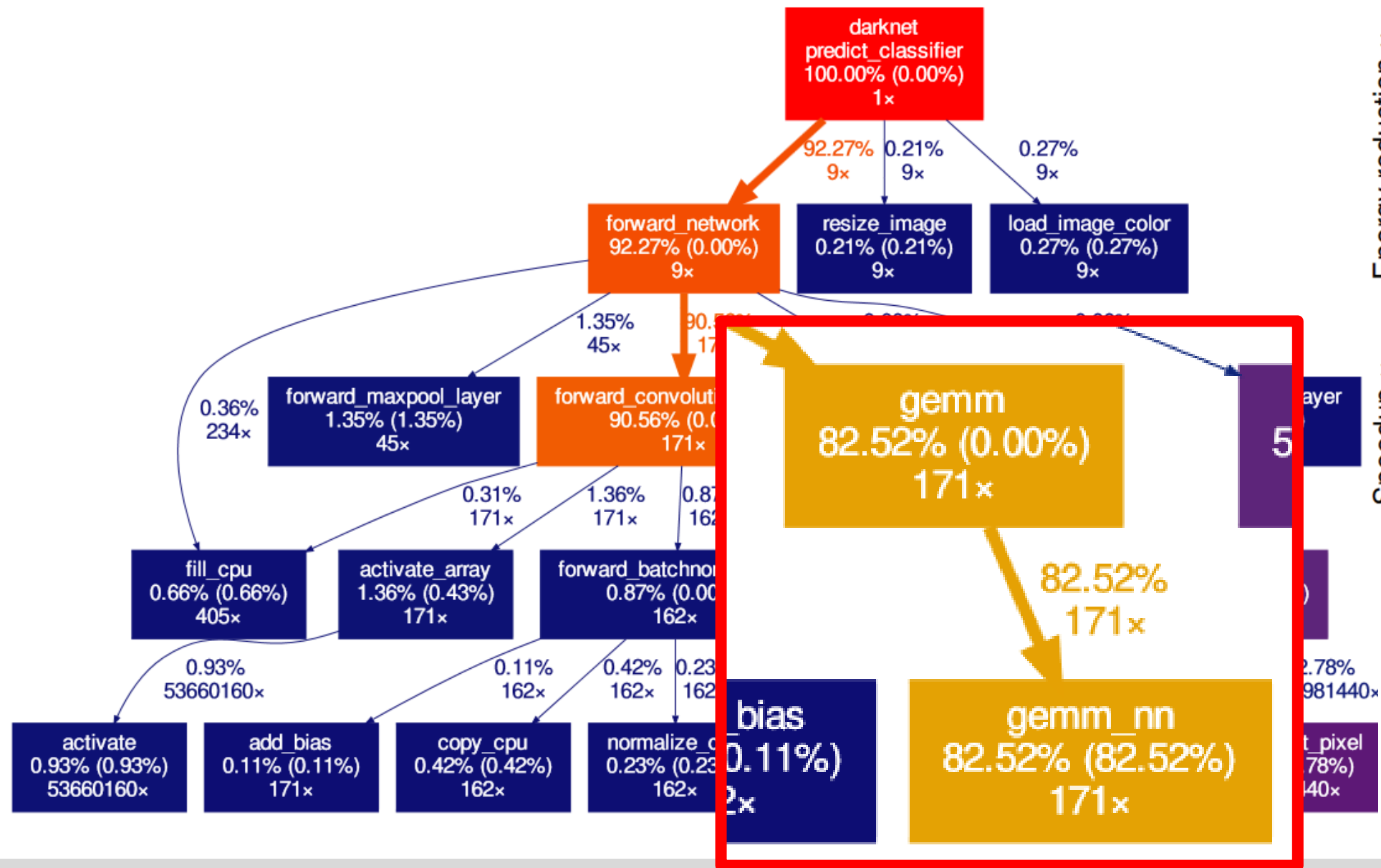


C-SRAM Vector Width (Bytes)



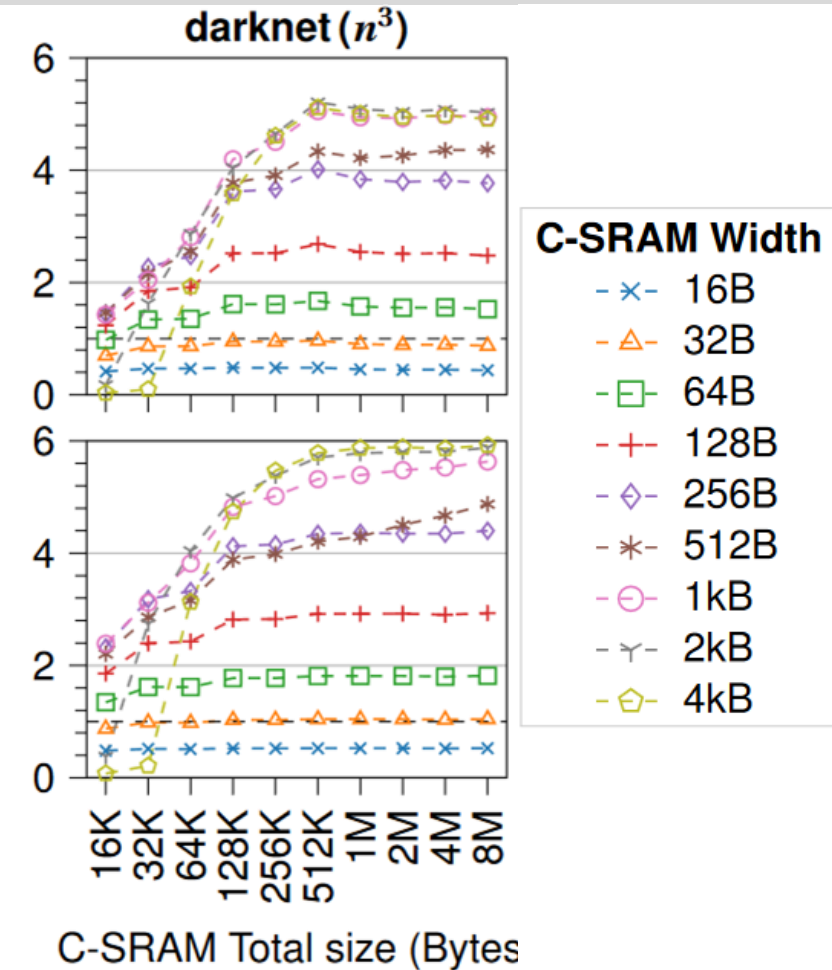
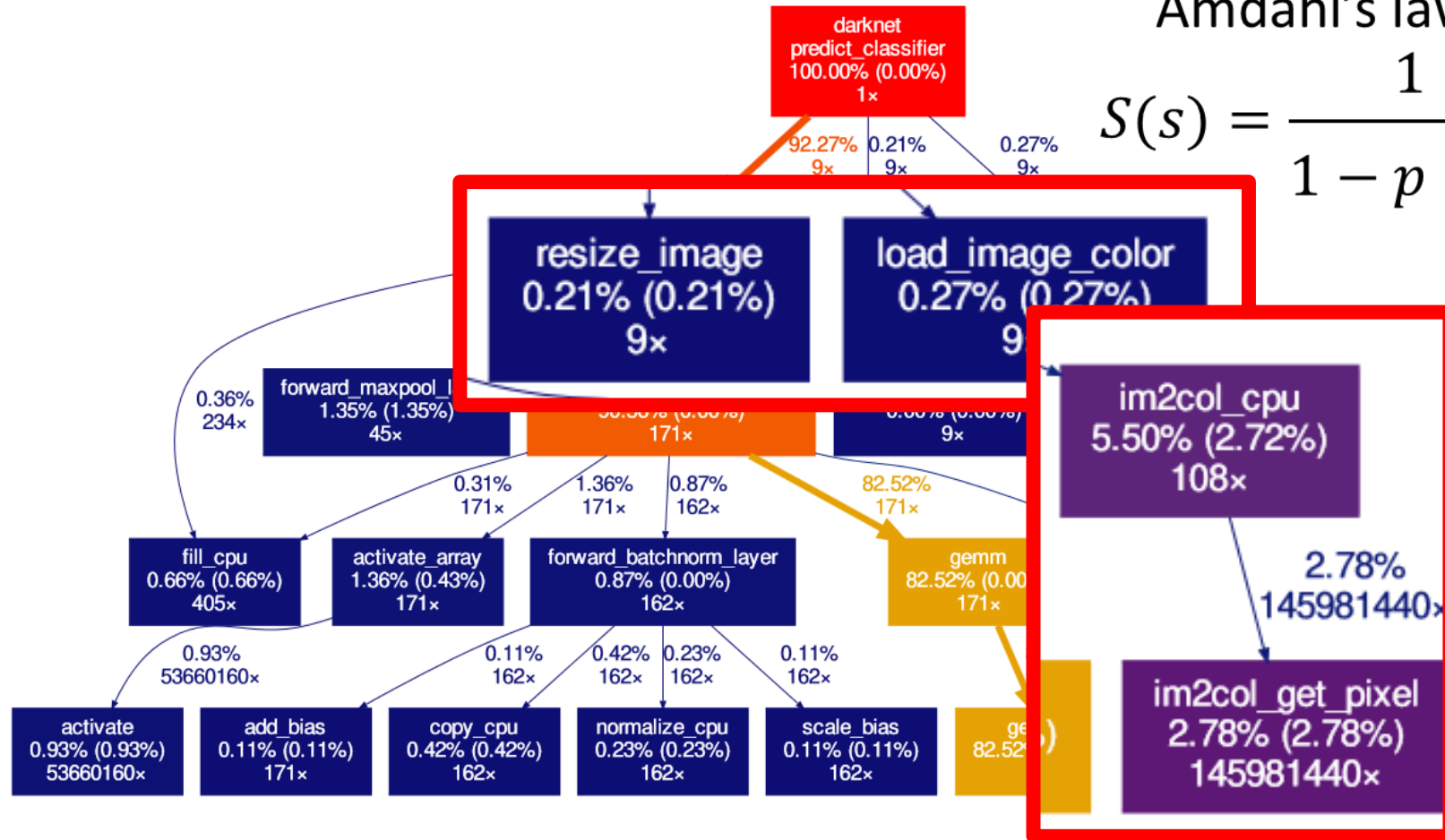
- ▶ Similar to atax, gemm requires minimal total size of 128kB
- ▶ More than optimal total size costs more similar as linear kernels
 - Timing unaffected due to less eviction
- ▶ Best case of 403x energy reduction and 271x speedup





Amdahl's law

$$S(s) = \frac{1}{1 - p + \frac{p}{s}}$$



- Data are stored on disks but must be decoded (jpg, png, mp4, etc.) before applying algorithms
- Need to layout data for neural network to turn convolution into gemm

- ▶ **No more room for performance improvements**
 - **Technology** reaches an **end**
 - Energy wall & dark silicon
 - CPU **architecture** has been **pushed** to its **maximum**
 - Memory wall & von Neumann bottleneck
 - **Energy efficiency** of current solutions is quite **low**
- ▶ **Need for a more energy efficient and data oriented architecture**
 - Big data era : neural networks and social networks
- ▶ → **In-Memory computing is promising solution**
 - **Tackle von Neumann bottleneck** at its source
 - **Computes** where data is and severely **reduces data movement**
- ▶ **SOTA is limited**
 - Analog computing
 - NVM endurance is often not discussed

► CSRAM design

- SRAM based and **automated** design approach
- **Competes** with state-of-the-art solutions

► Results

- Moving **computing closer** to where **data** live seems to be the **best option**
 - **Computing** in the **SCM** will **kill** it in less than a **day**
 - **Data** are stored on **disks** but must be **decoded** (jpg, png, mp4, etc.) **before** applying **algorithms**
 - Need to **layout data** for **neural network** to turn **convolution** into **gemm**
 - **DRAM** is not yet evinced as it serves as a **write buffer**
 - **Cache** hierarchy can be reduced but is still **used** for **instructions, stack**
- ## ► Average 17.4x energy reduction and 12.9x speedup
- **Optimal** memory size around **512kB**

► Future works

- CPU should be used for **parallel computing**
- **Compare** against other **architectures** such as **GPUs**
- **Silicon implementation**

► Perspectives

- **Take all application into account, not just kernel**
- **Watch out for rebound effect**

► Publications

- “Shuffle operator for matrix multiplication in in-memory computing architecture”, **Author**, Poster in COMPAS, 06/19, and also in ACACES HiPEAC, 07/19
- “Computational SRAM Design Automation using Pushed-Rule Bitcells for Energy-Efficient Vector Processing”, **2nd author**, DATE 2020 (special session)
- “Reconfigurable Tiles of Computing-In-Memory SRAM Architecture for Scalable Vectorization”, **2nd author**, ISLPED 2020
- “Storage Class Memory with Computing Row Buffer: A Design Space Exploration”, **author**, DATE 2021

► Patents

- “Outil de selection automatique de mémoires pour circuits de type *in-* ou *near-memory computing*”, brevet **co-inventeur**, 2020
- “Tampon de lignes mémoire volatile calculant pour mémoires non volatiles ”, **inventeur**, 2021

The logo for the Commissariat à l'énergie atomique et aux énergies alternatives (CEA), featuring the lowercase letters 'cea' in white on a red square background, with a thin green horizontal line underneath.

Thank you for your attention

The logo for Aix-Marseille université, featuring a blue curved line on the left, the text 'Aix*Marseille' in black, 'université' in blue, and 'Initiative d'excellence' in a smaller blue font below.

Backup

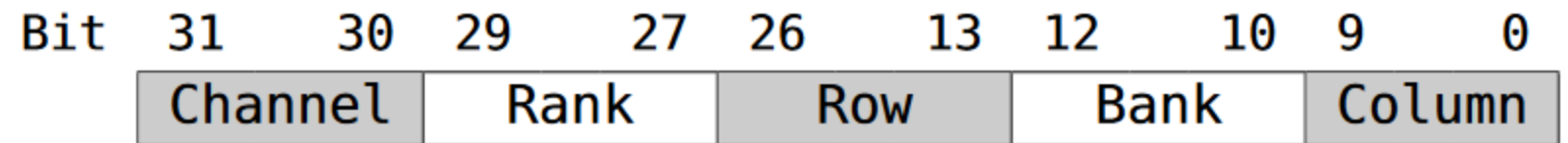
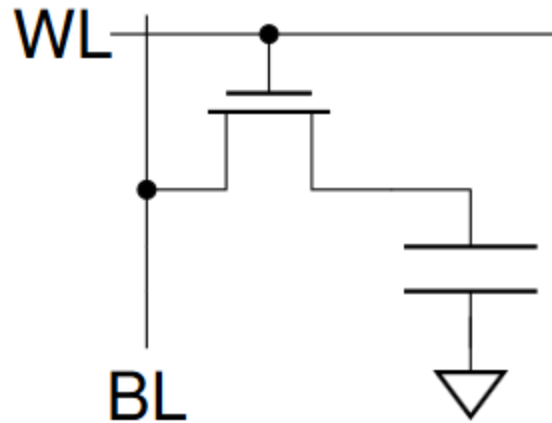
- Frequency scaling of an Intel Xeon Silver 4116, a 12 cores chip, function of active cores and active SIMD extension.

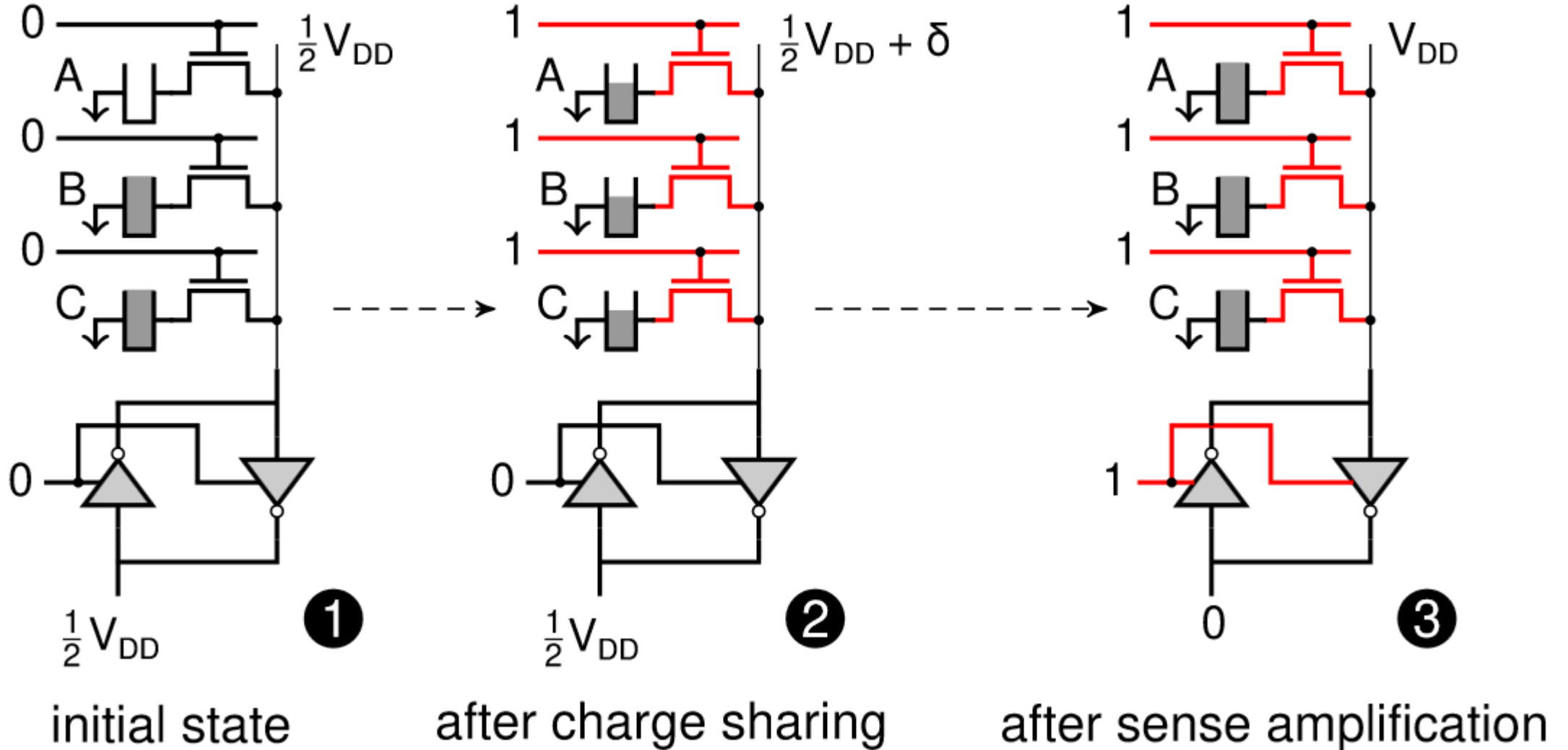
Mode	Base	Turbo Frequency/Active Cores											
		1	2	3	4	5	6	7	8	9	10	11	12
Normal	2.1 GHz	3.0 GHz	3.0 GHz	2.8 GHz	2.8 GHz	2.7 GHz	2.7 GHz	2.7 GHz	2.7 GHz	2.4 GHz	2.4 GHz	2.4 GHz	2.4 GHz
AVX2	1.7 GHz	2.9 GHz	2.9 GHz	2.7 GHz	2.7 GHz	2.4 GHz	2.4 GHz	2.4 GHz	2.4 GHz	2.1 GHz	2.1 GHz	2.1 GHz	2.1 GHz
AVX512	1.1 GHz	1.8 GHz	1.8 GHz	1.6 GHz	1.6 GHz	1.5 GHz	1.5 GHz	1.5 GHz	1.5 GHz	1.4 GHz	1.4 GHz	1.4 GHz	1.4 GHz

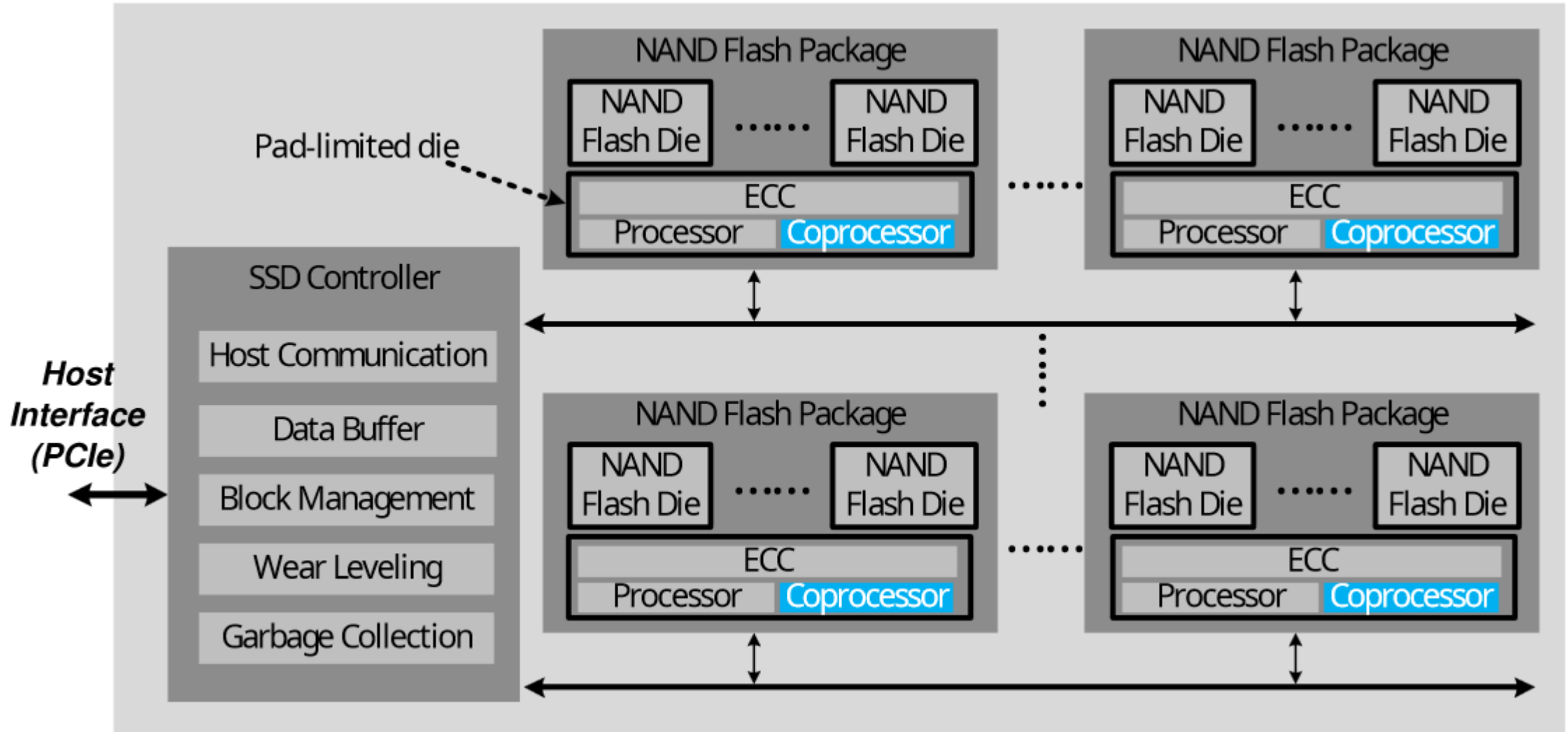
Integer		Floating Point		Memory	
Addition				Cache	64 bit access
8 bit	0.03pJ	16 bit	0.4pJ	8kB	10pJ
32 bit	0.1pJ	32 bit	0.9pJ	32kB	20pJ
Multiplication				1MB	100pJ
8 bit	0.2pJ	16 bit	1.1pJ	DRAM	1.3-2.6nJ
32 bit	3.1pJ	32 bit	3.7pJ		

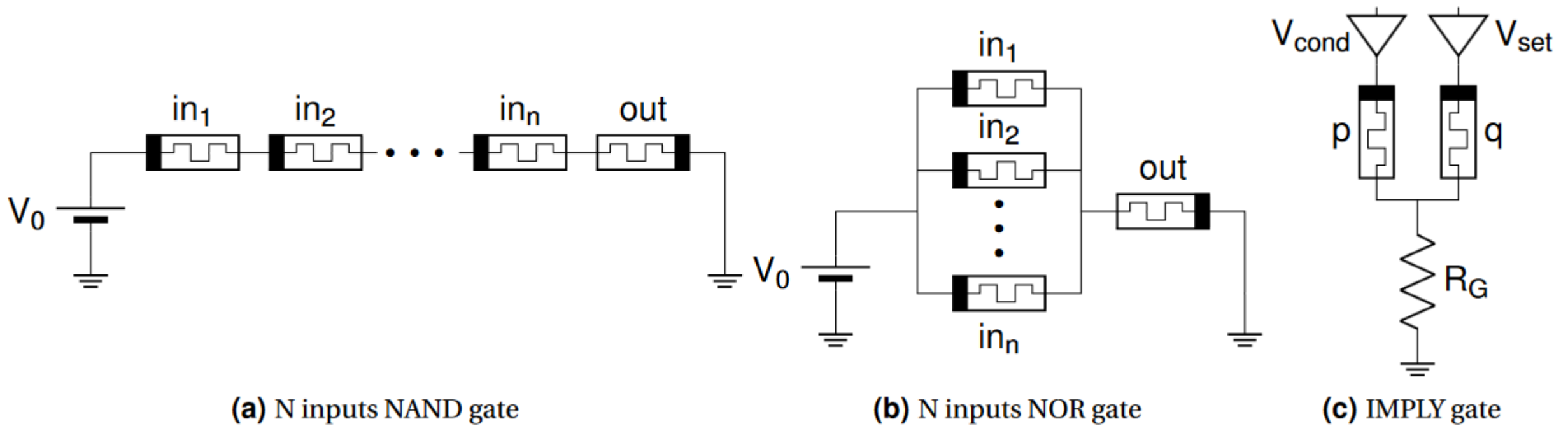


	DRAM	SRAM	NAND Flash	STT-RAM	PCM	RRAM
Cell Size	6 F ²	120–200 F ²	<4 F ² (3D)	6–50 F ²	4–30 F ²	4–12 F ²
Multibit	1	1	3	2	4	2
Read latency	10–50 ns	0.2–2 ns	15–35 μs	2–35 ns	10–60 ns	5–20 ns
Write latency	10–50 ns	0.2–2 ns	100–1000 μs	3–50 ns	20–150 ns	10–50 ns
Write energy (/bit)	10 fJ	1 fJ	10 nJ	10 fJ–1 pJ	10–500 pJ	0.1–10 pJ
Endurance	∞	∞	1e4–1e5	1e12–1e15	1e7–1e10	1e6–1e12
Cons	Read destructive + Refresh	High leakage	Erase block (4kB) before write	Process incompatibility	Heat transfer + high current	Sneak current + IR drop
Maturity	Mature	Mature	Mature	Testchips	Com.	Testchips





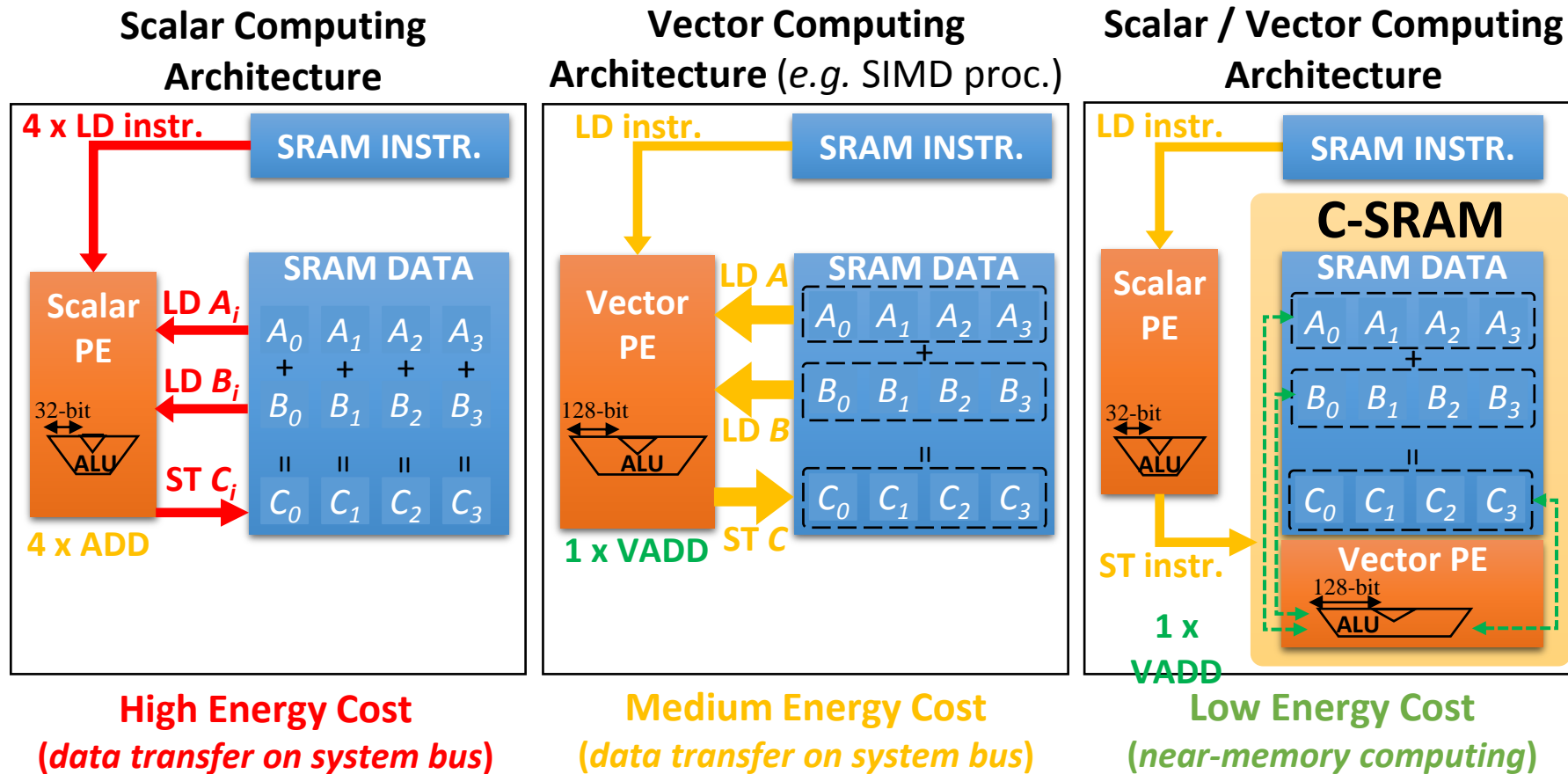




(a) N inputs NAND gate

(b) N inputs NOR gate

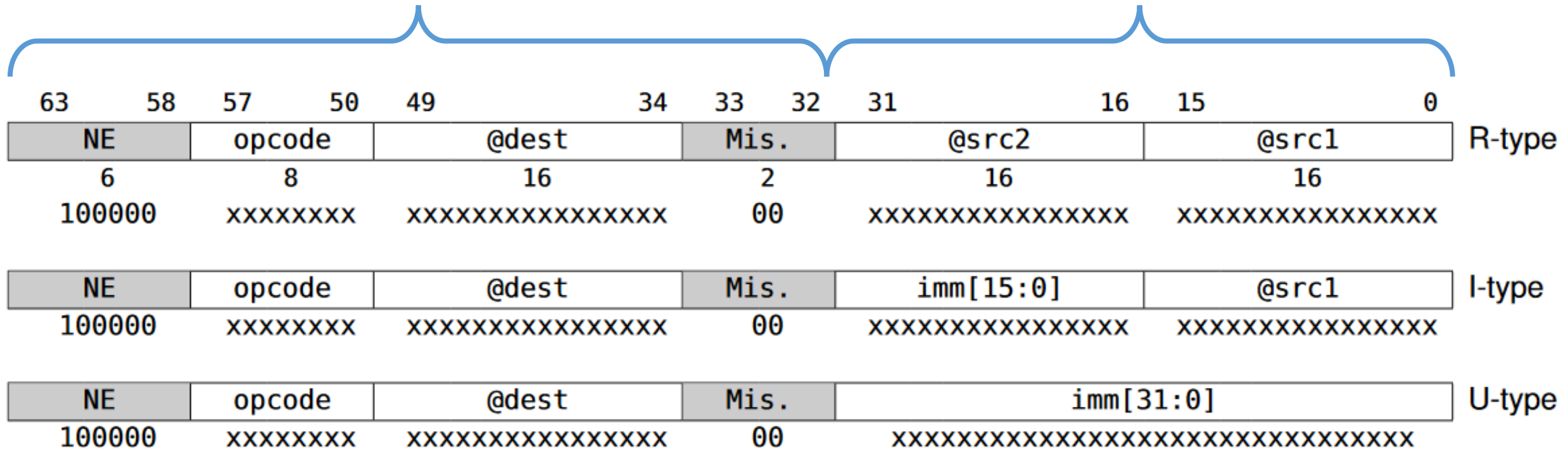
(c) IMPLY gate



Defined 32 bits ISA (32 address bits + 32 data bits)

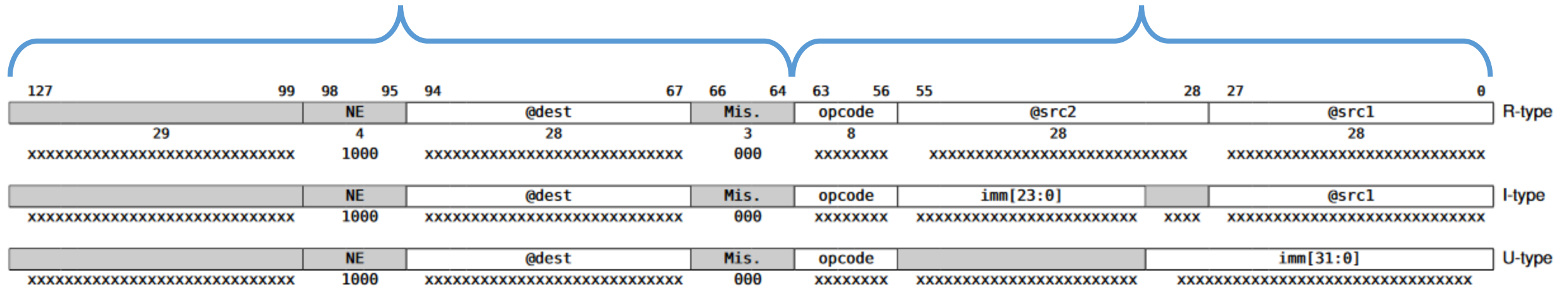
32 address bits

32 data bits

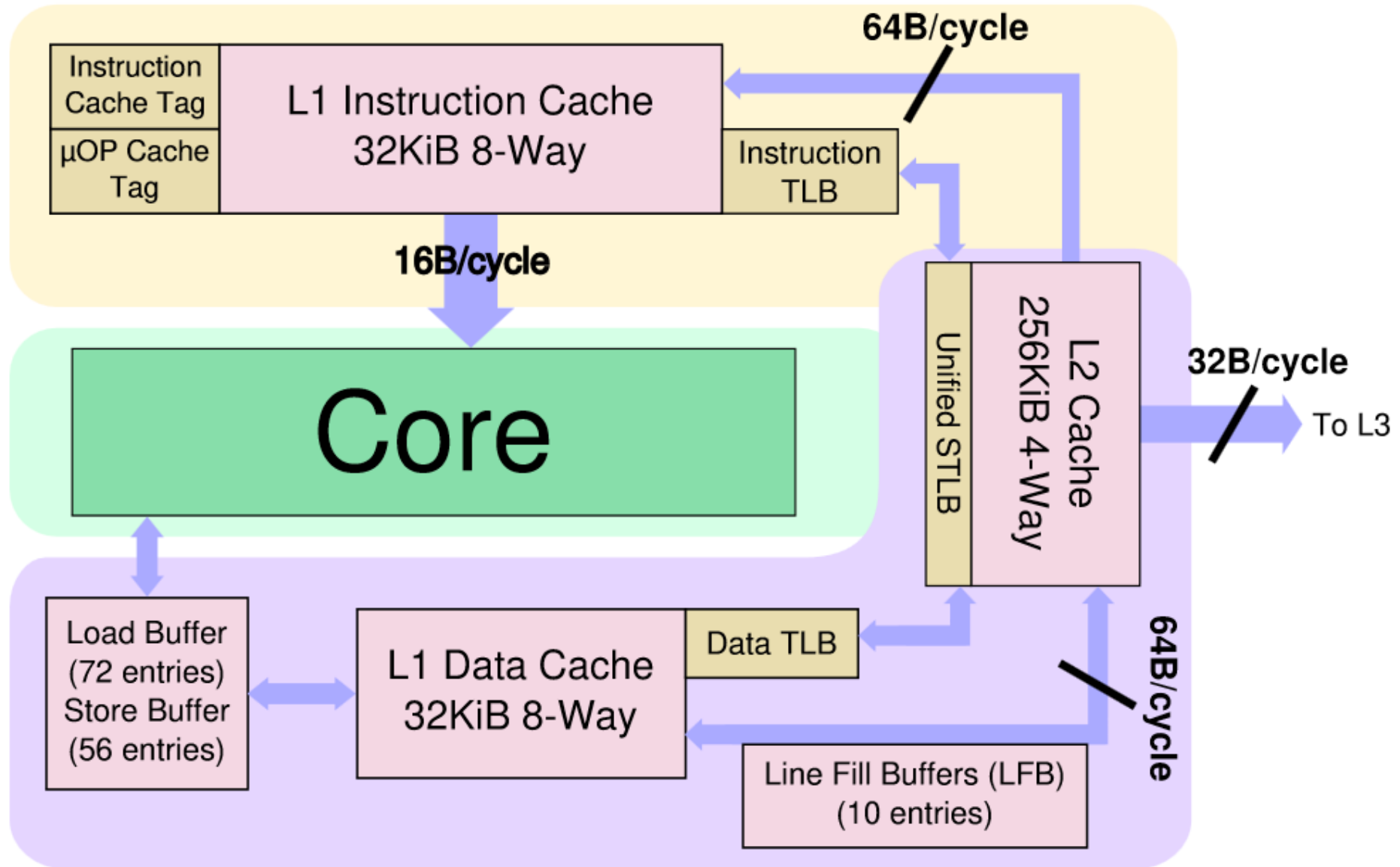


64 address bits

64 data bits



Simplified Intel Skylake memory system



► Scaling only applies to routing

(a) Original 128 bits IO 128 MB PCRAM

	Energy	Latency
Read	0.875 nJ	142 ns
Write	4.88 nJ	512 ns

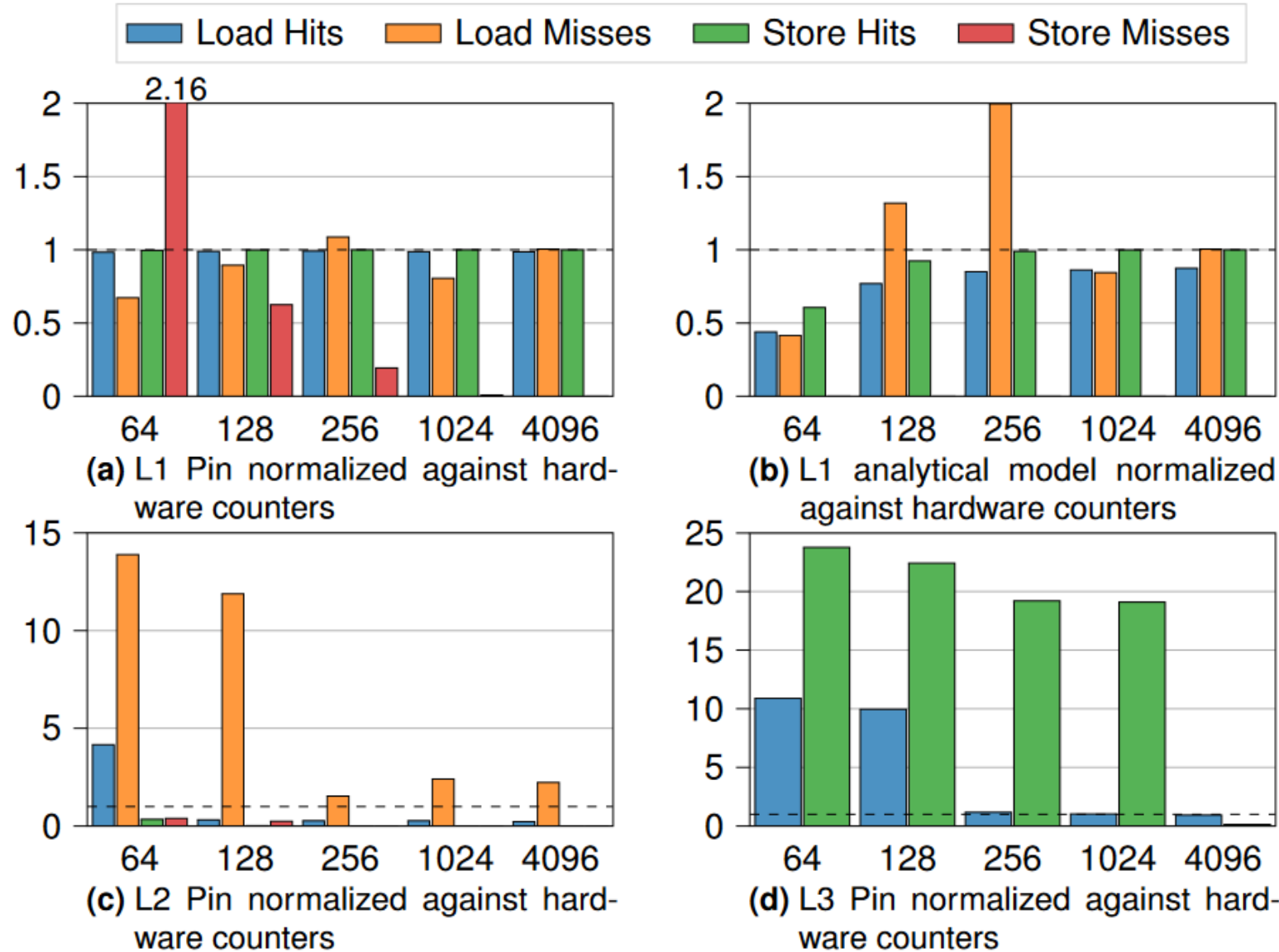
(b) Scaled 4 GB PCRAM

IO Width	Read		Write		Widen Factor
	Energy	Latency	Energy	Latency	
512 B	38.5 nJ	263 ns	215 nJ	949 ns	0
1 kB	77.0 nJ	447 ns	429 nJ	1.61 μ s	1
2 kB	154 nJ	759 ns	858 nJ	2.74 μ s	2
4 kB	308 nJ	1.29 μ s	1.72 μ J	4.66 μ s	3

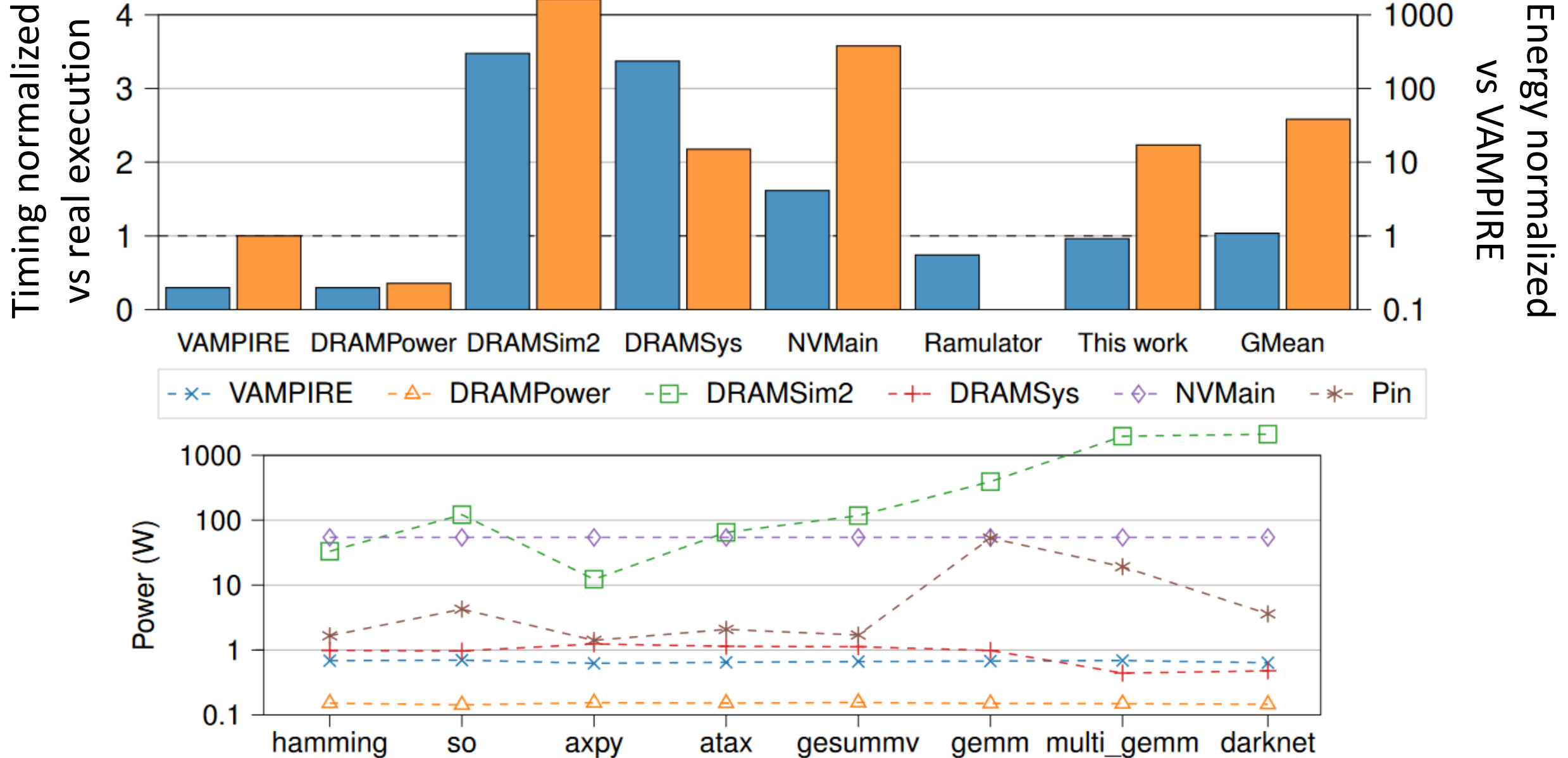
Size	16 kB	32 kB	64 kB	128 kB	256 kB	512 kB	1–8 MB
N_{cuts}	8	16	16	32	32	32	64
E_{tf}	10 %	10 %	19 %	19 %	37 %	37 %	37 %
T_{tf}	30 %	30 %	56 %	56 %	85 %	85 %	85 %

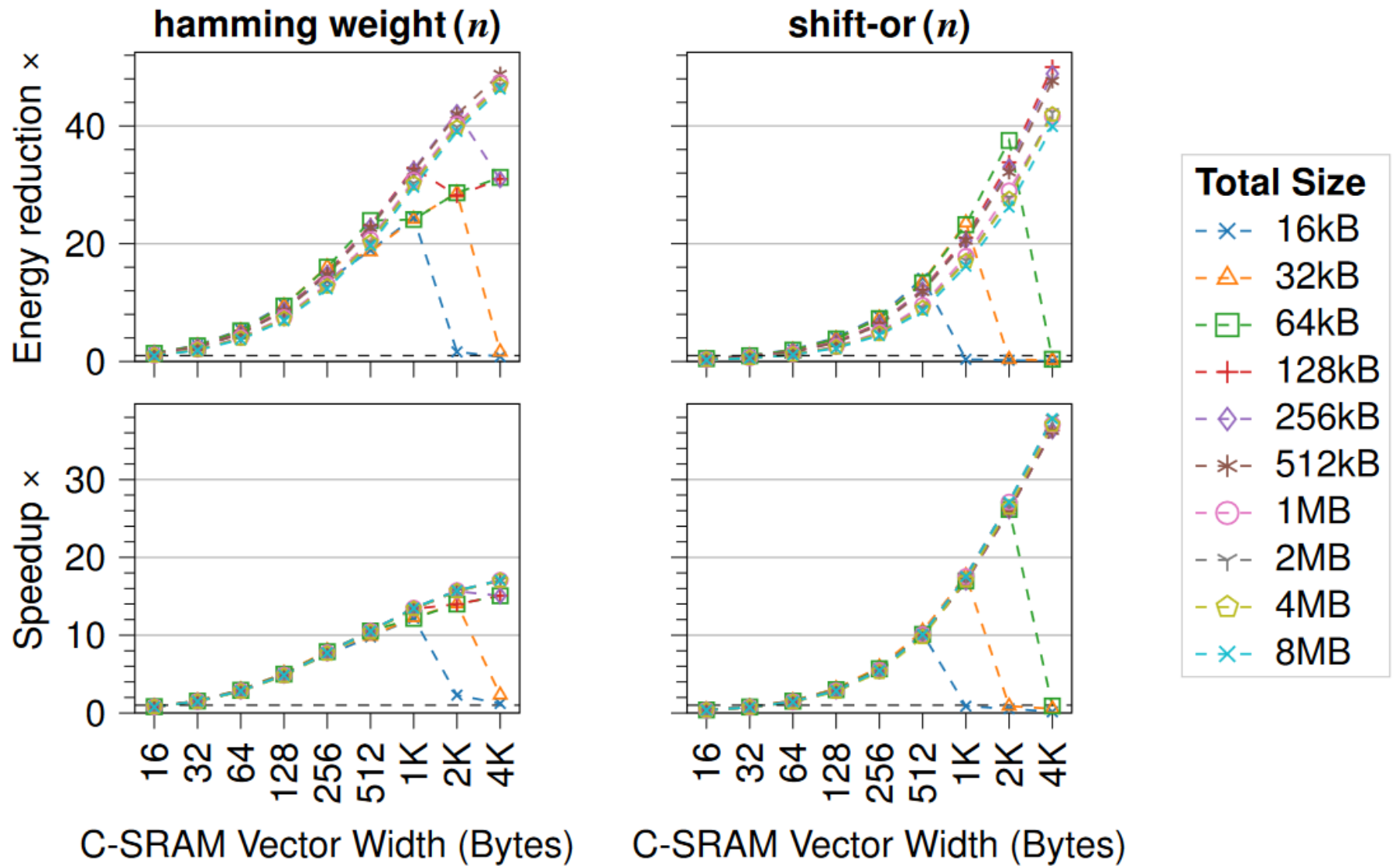
$$E = (E_{base} \times E_{W_f}^{W_f}) \times E_{tf} \times N_{cuts}$$

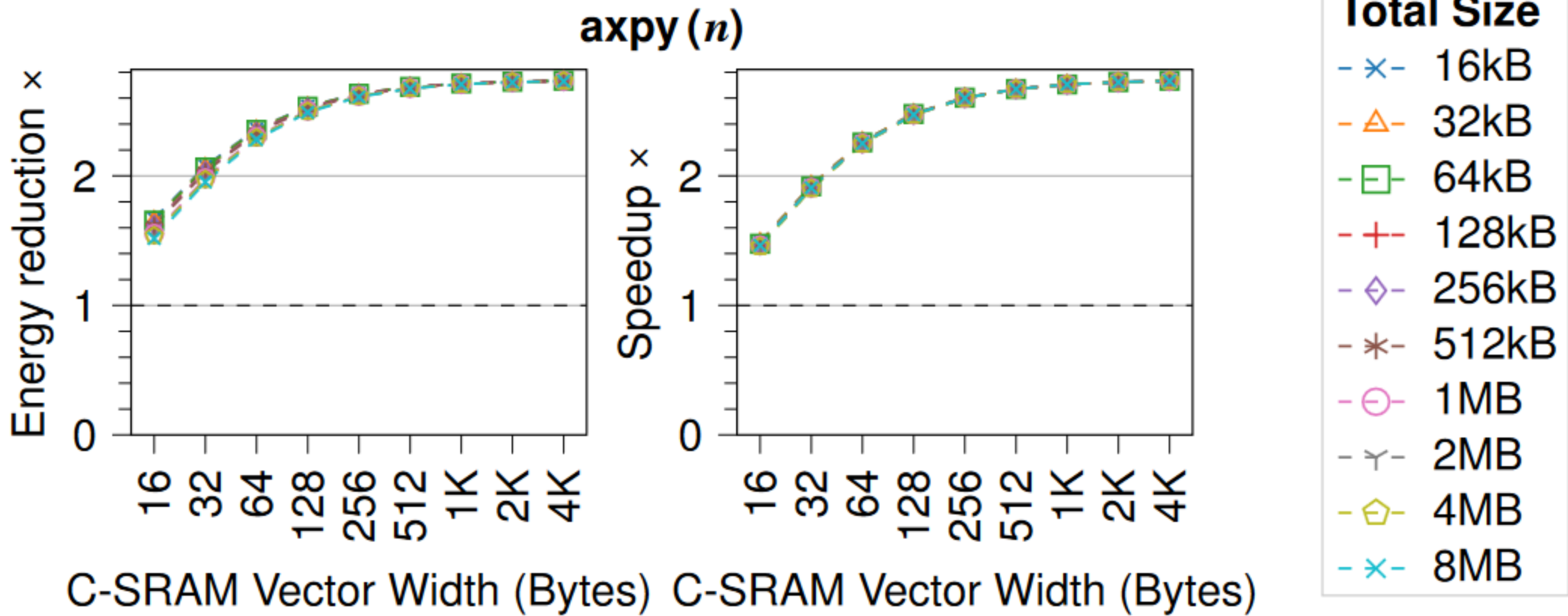
$$T = (T_{base} \times T_{W_f}^{W_f}) \times T_{tf}$$



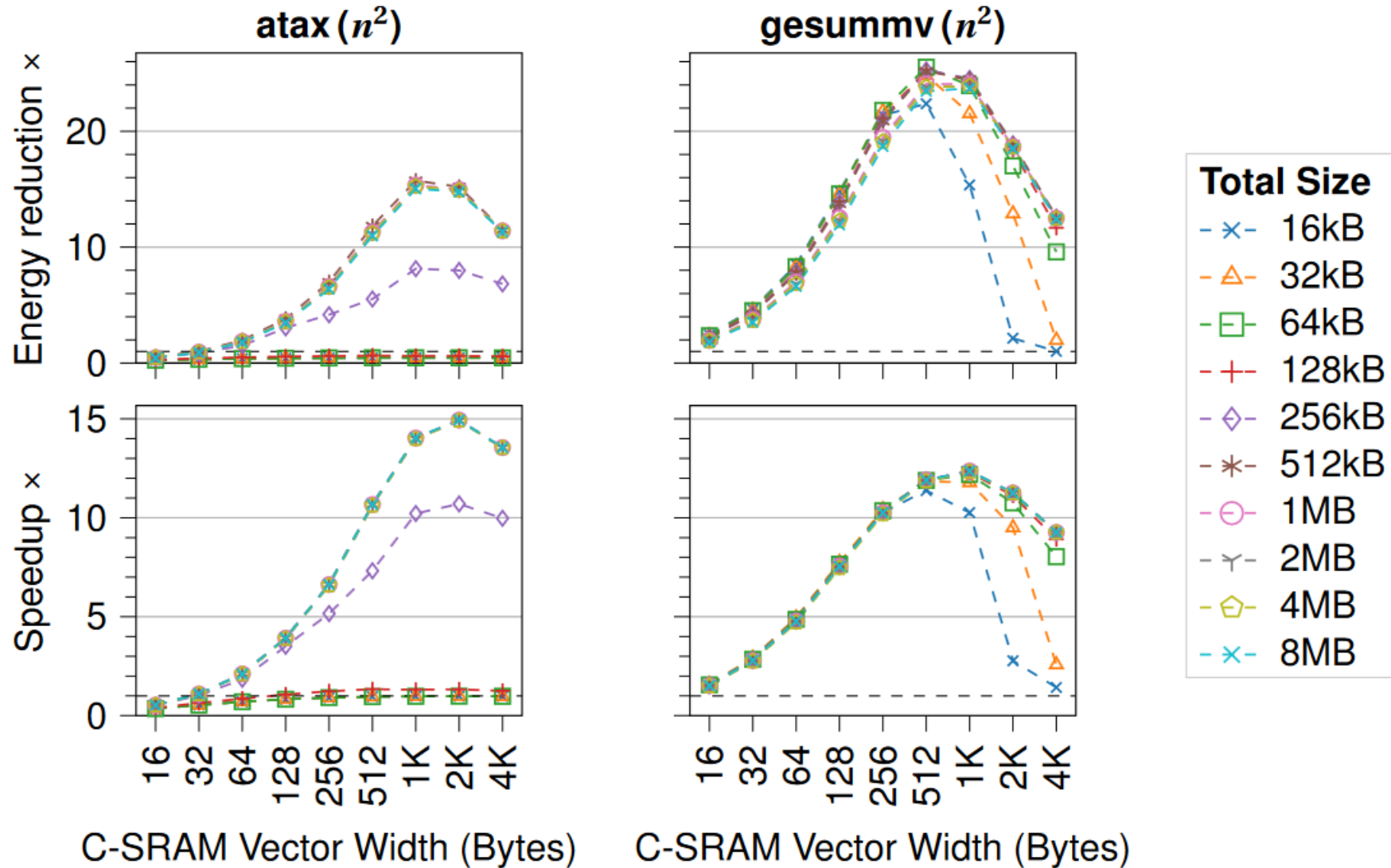
Platform validation: DRAM accesses cost

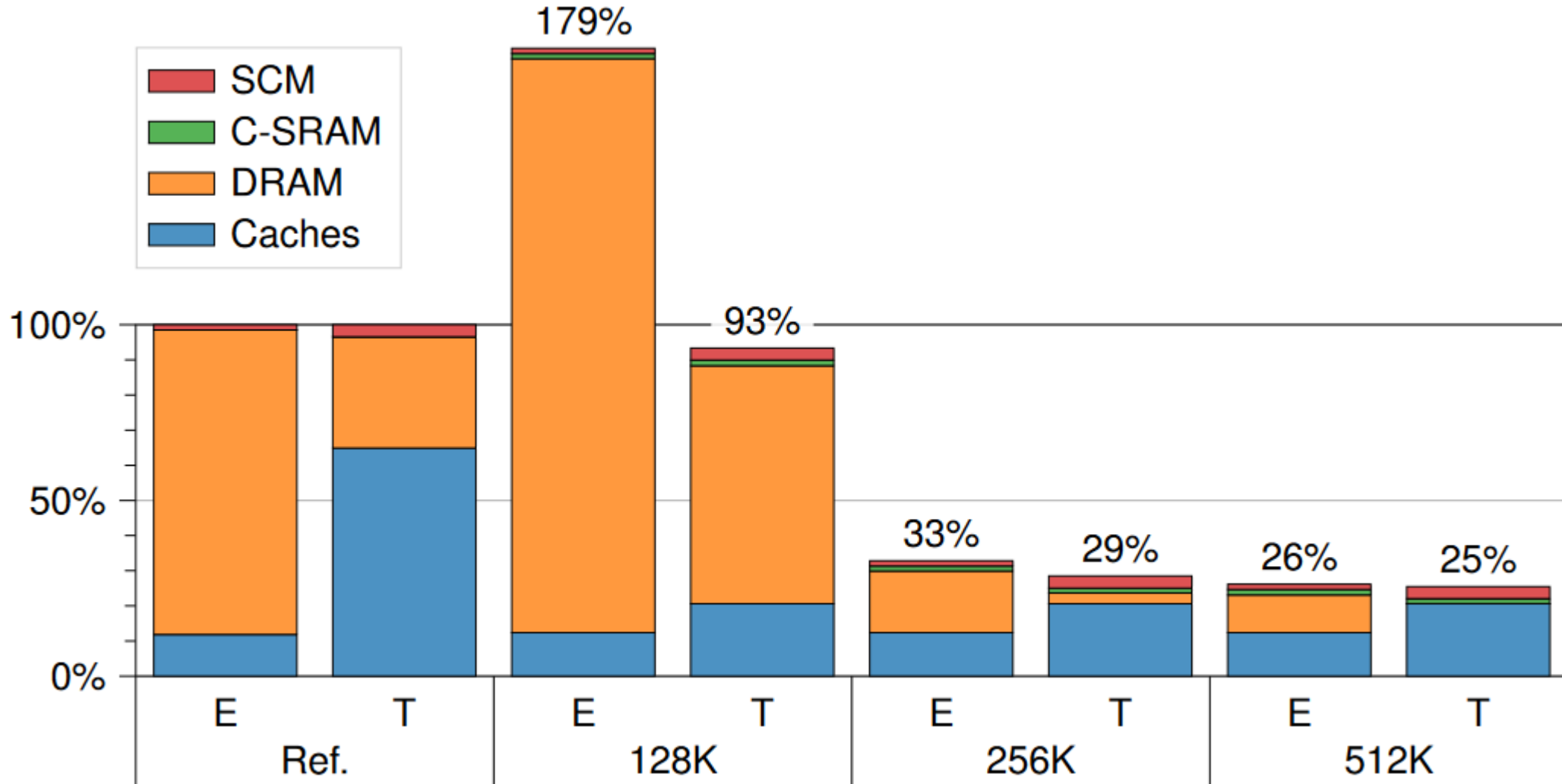




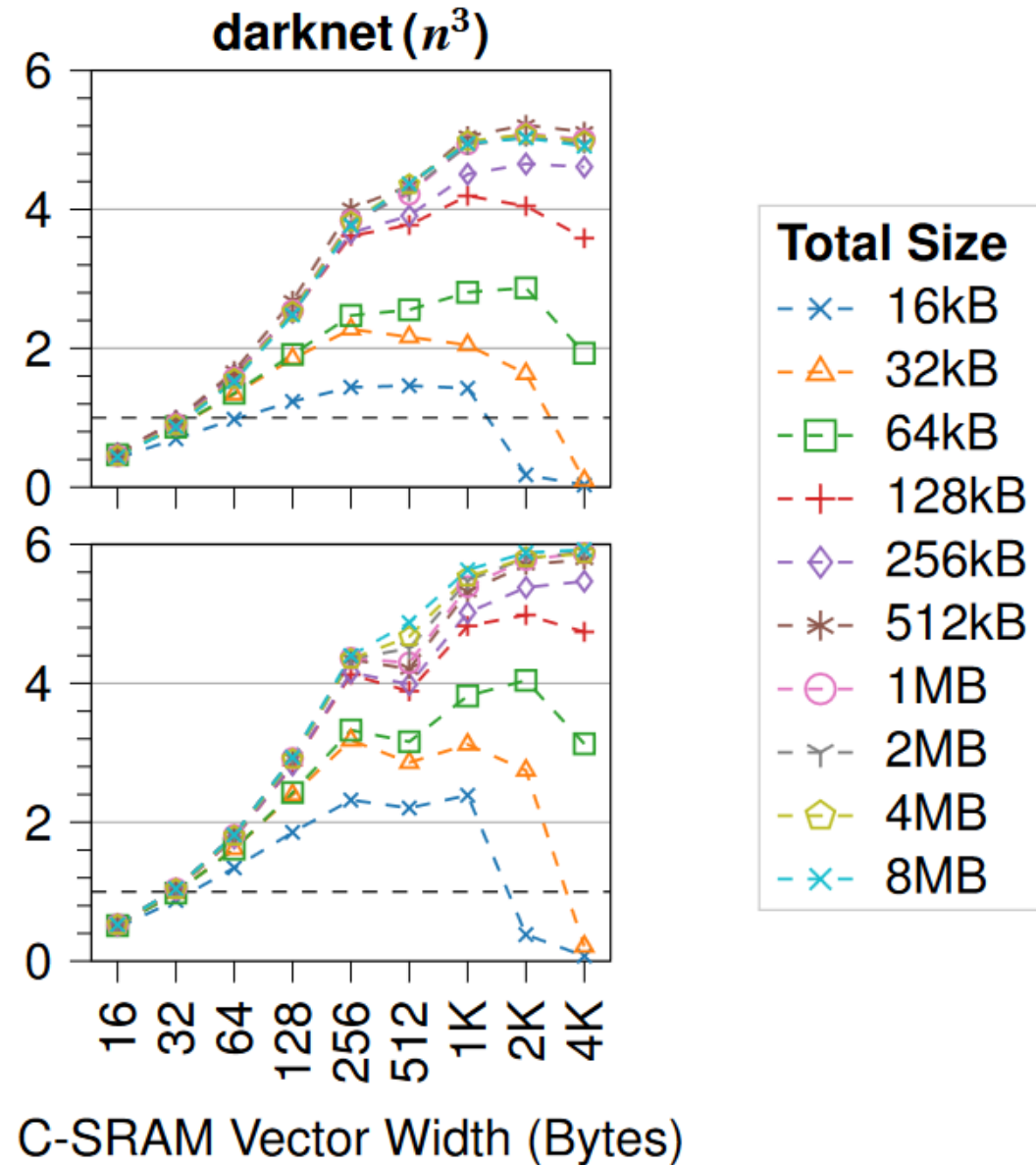
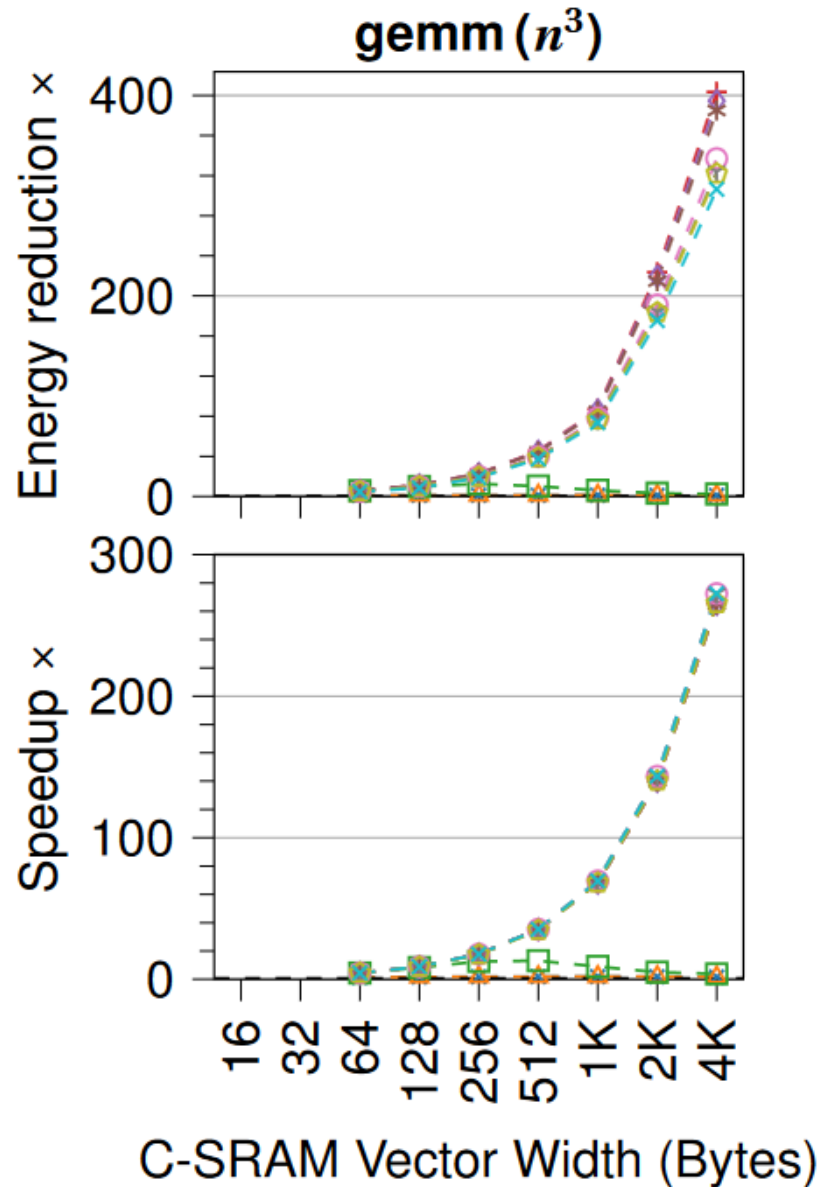


Quadratic kernels: atax & gesummv

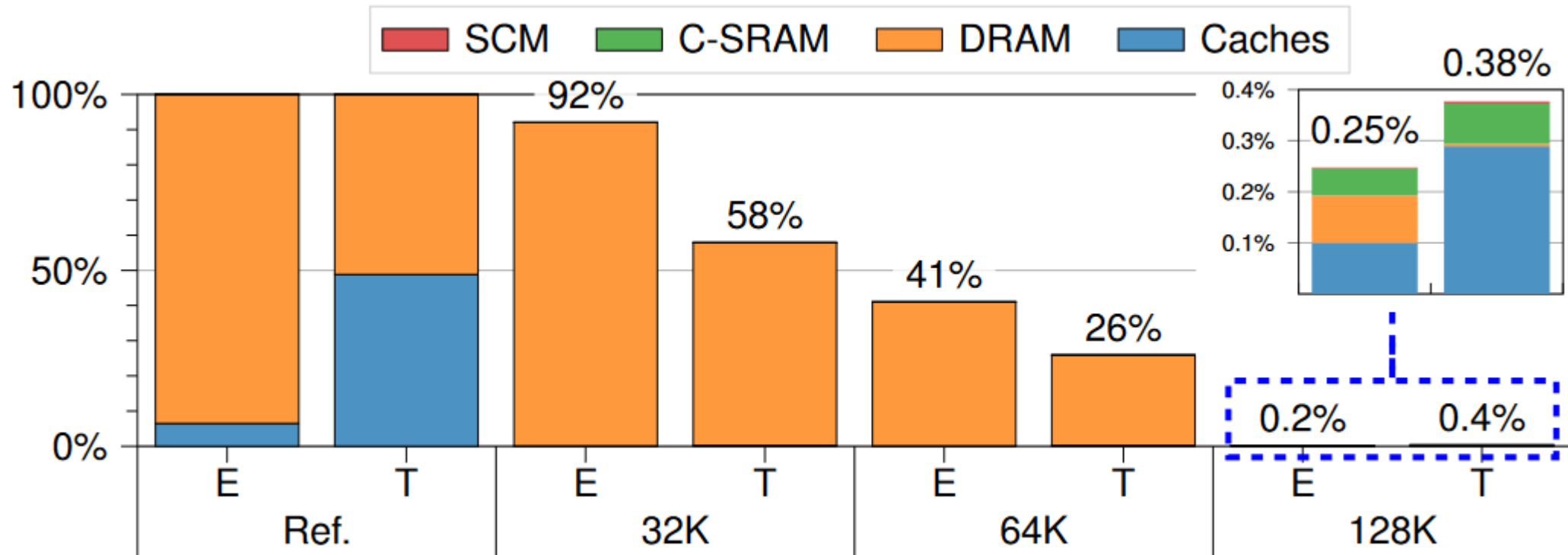




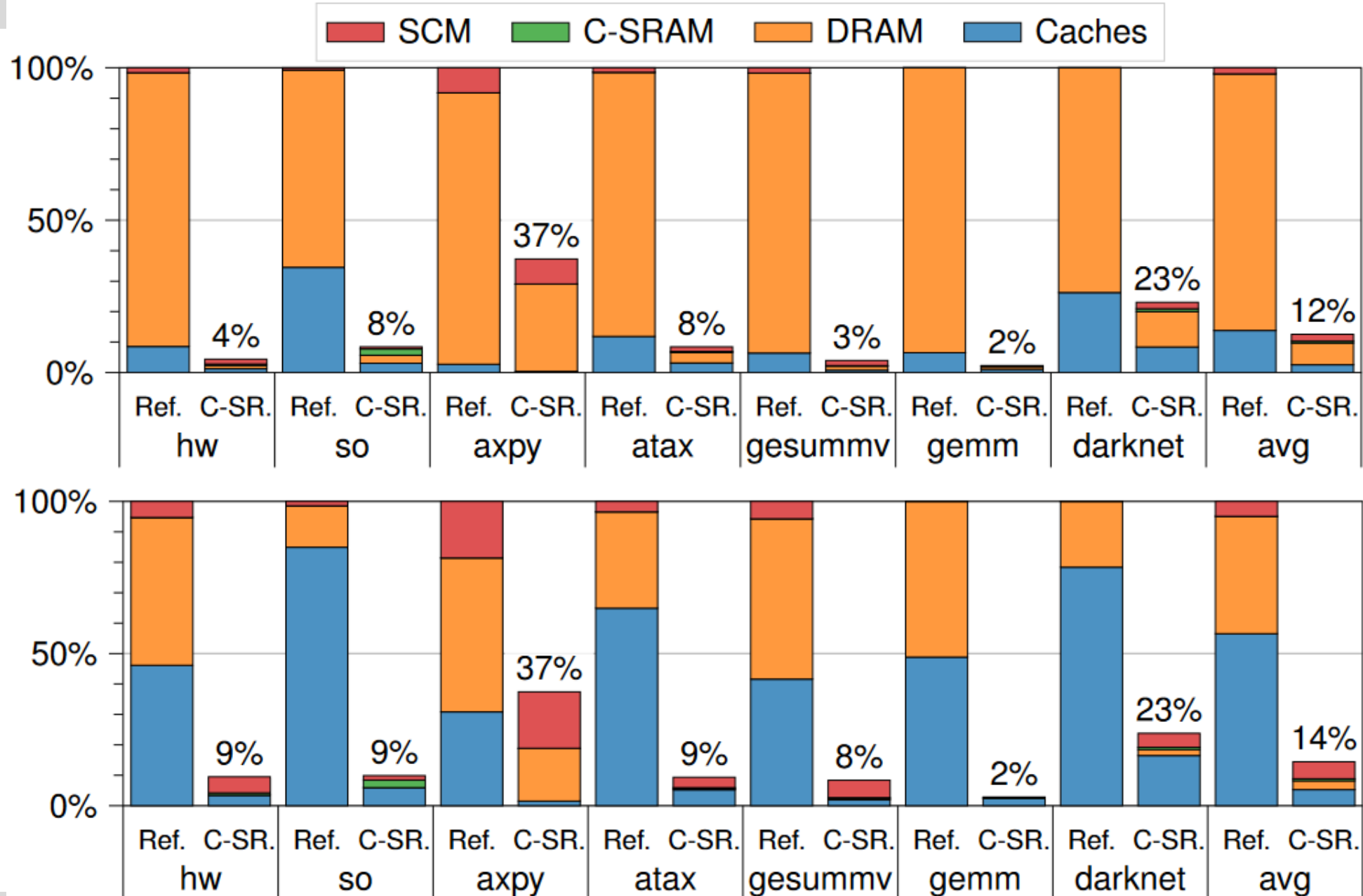
Cubic kernels: gemm & darknet



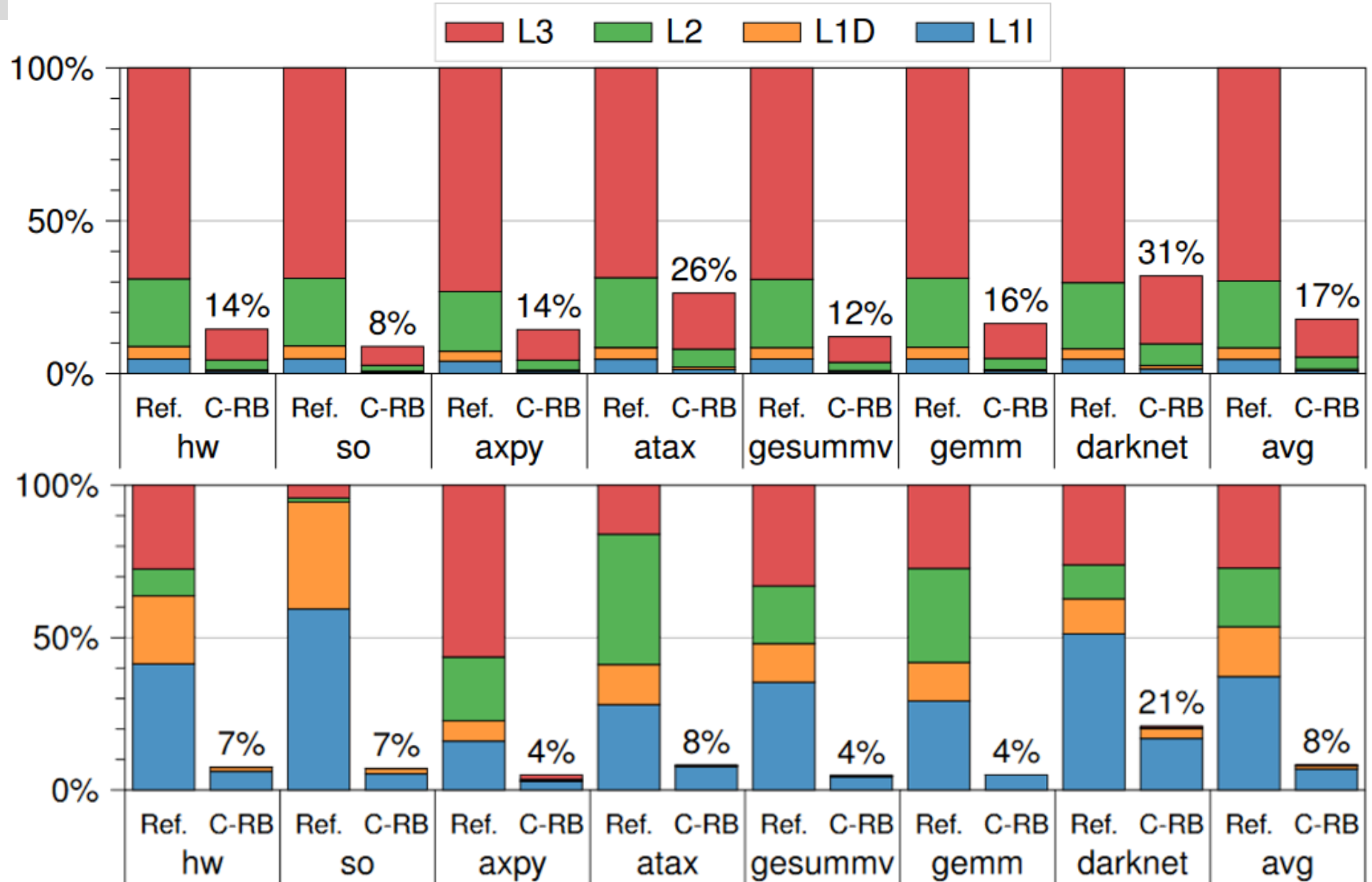
Cubic kernels: gemm sudden jump

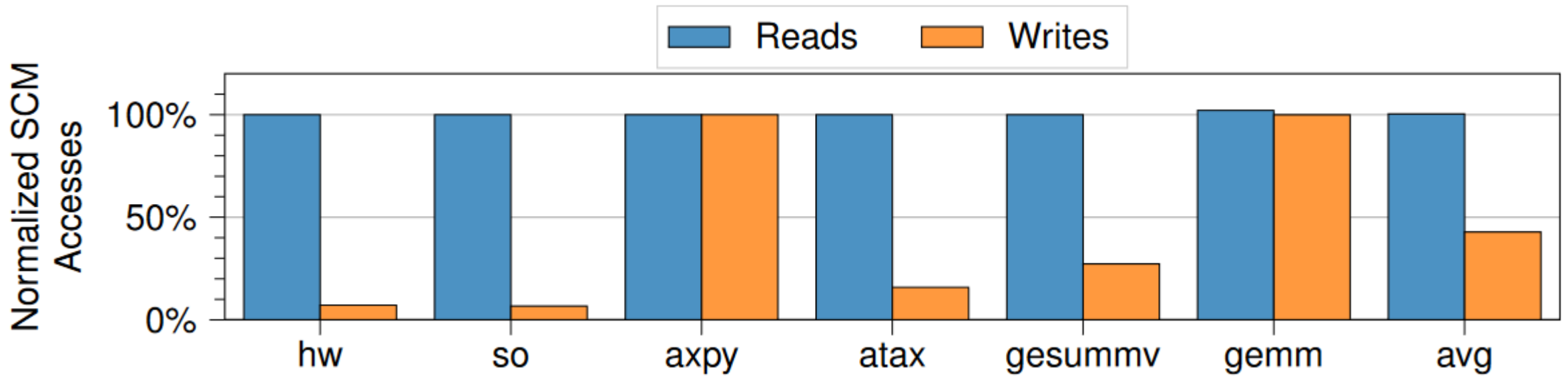


All benchmarks distribution



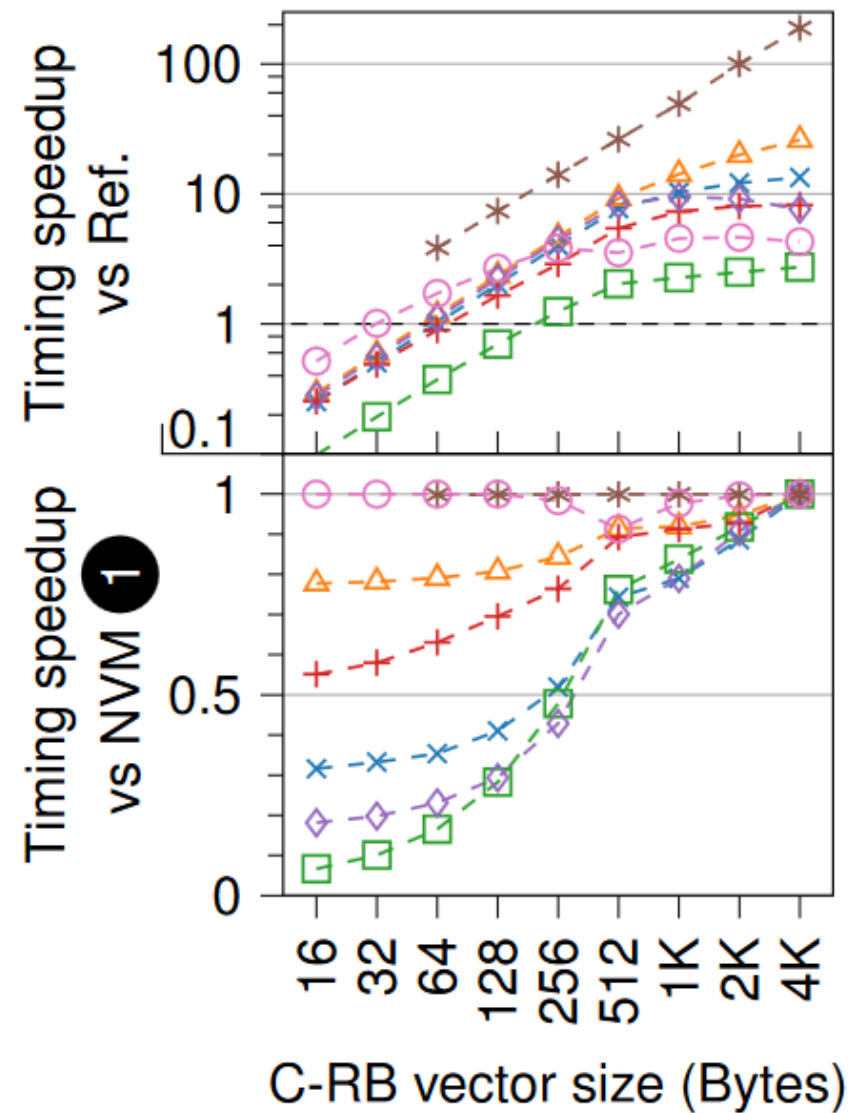
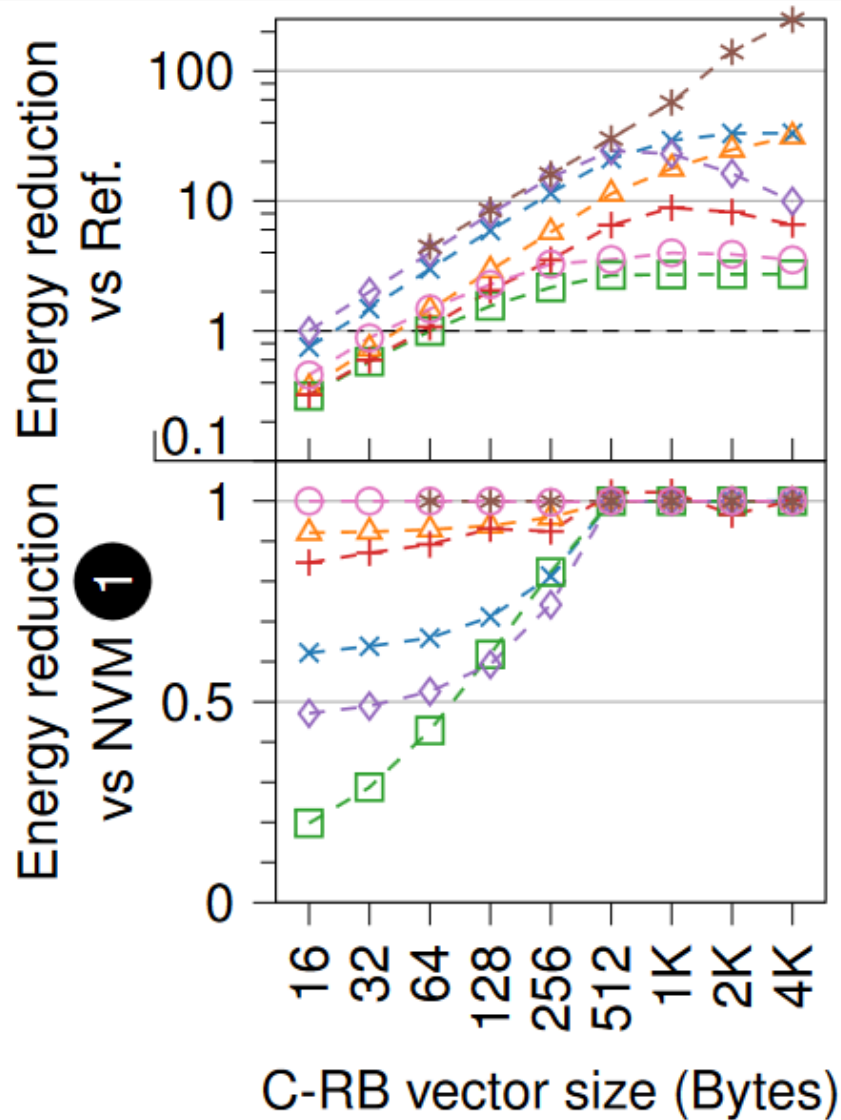
All benchmarks caches distribution





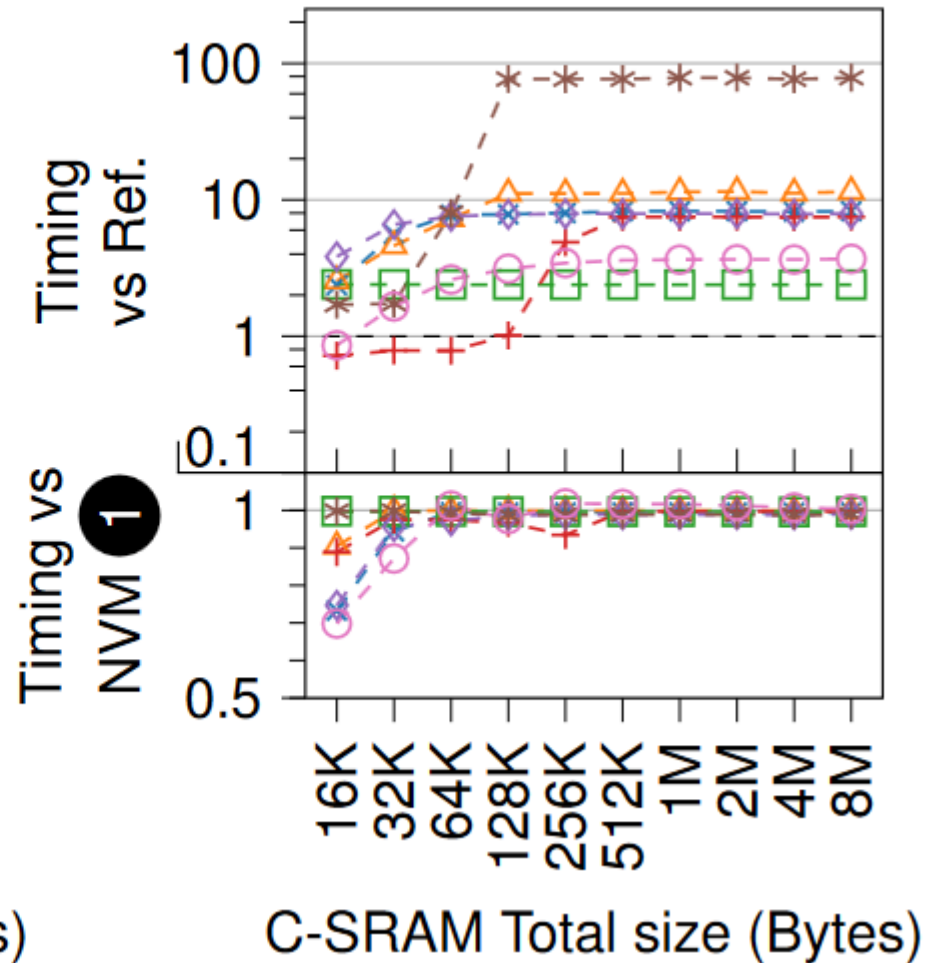
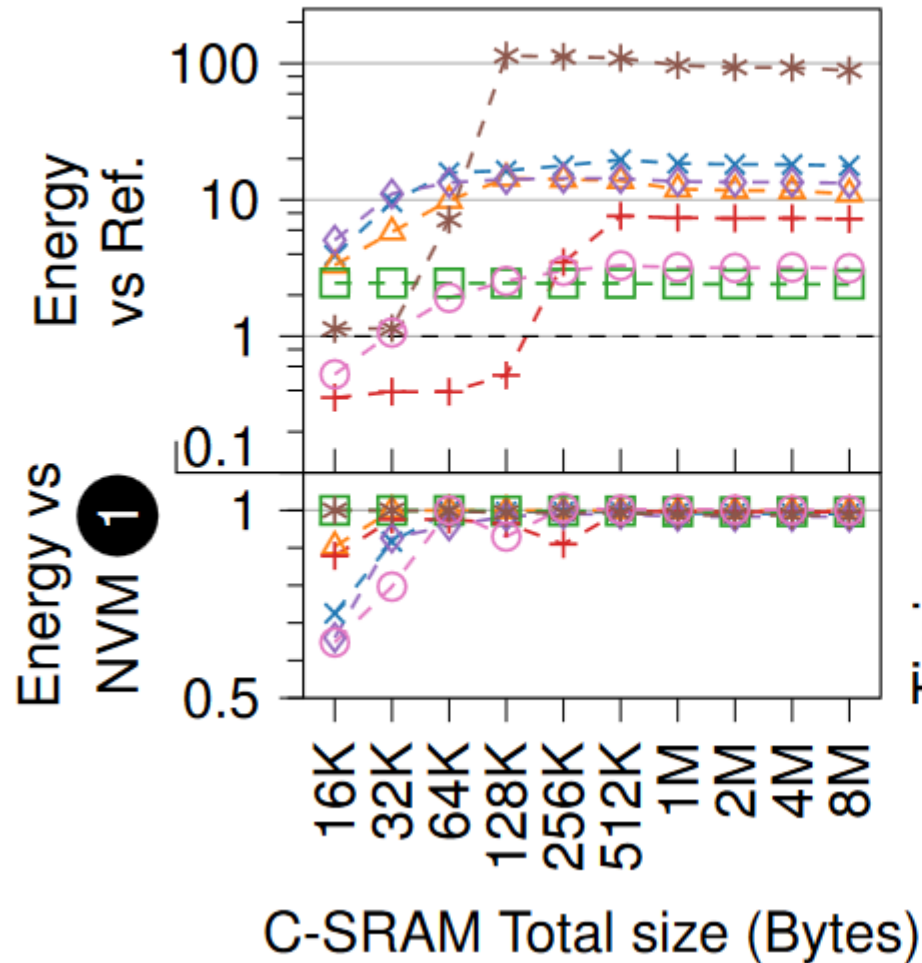
Row buffer scenario

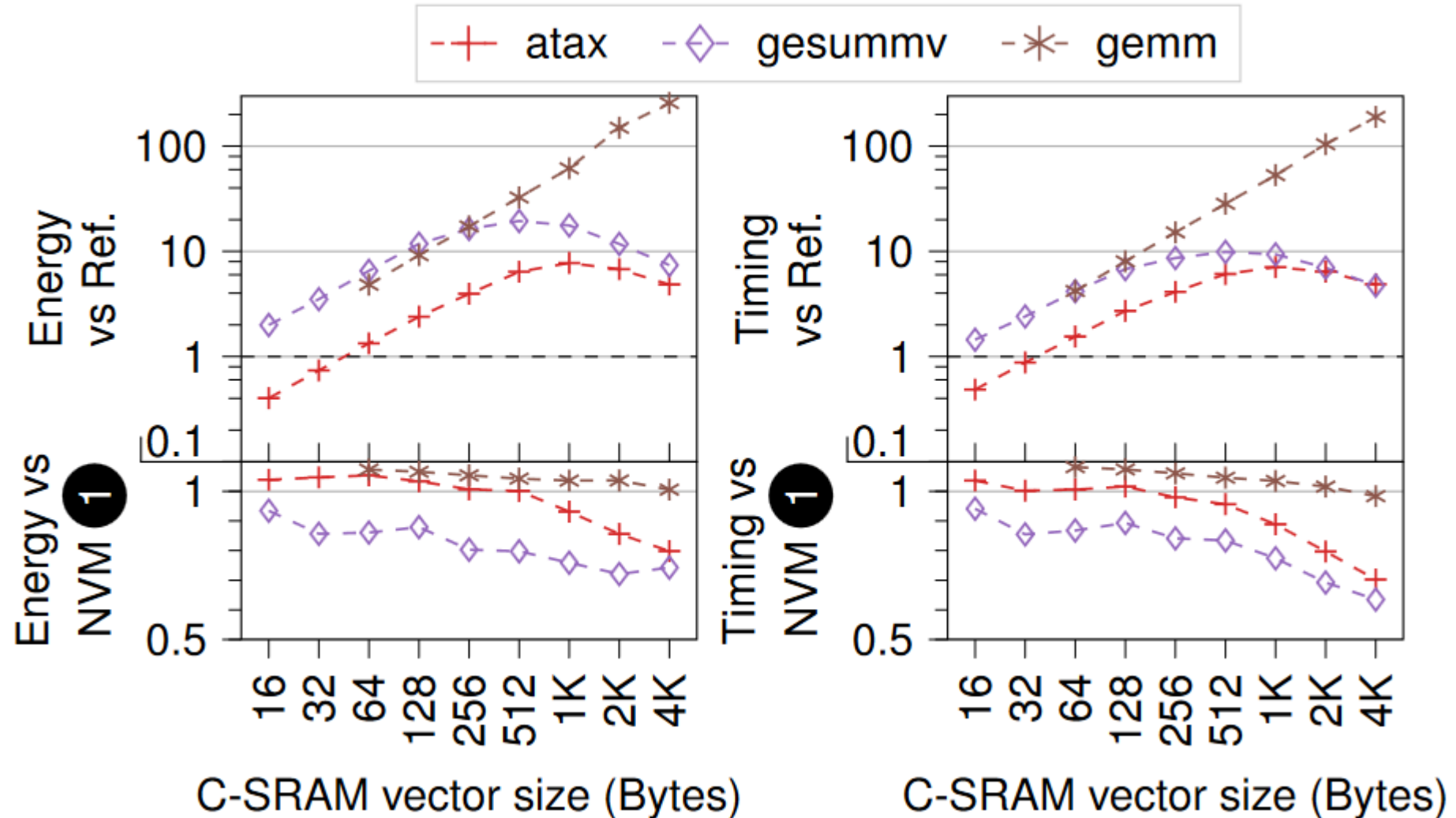
-x- hw -△- so -□- AXPY -+ atax -◇ gesummv -* gemm -○ darknet

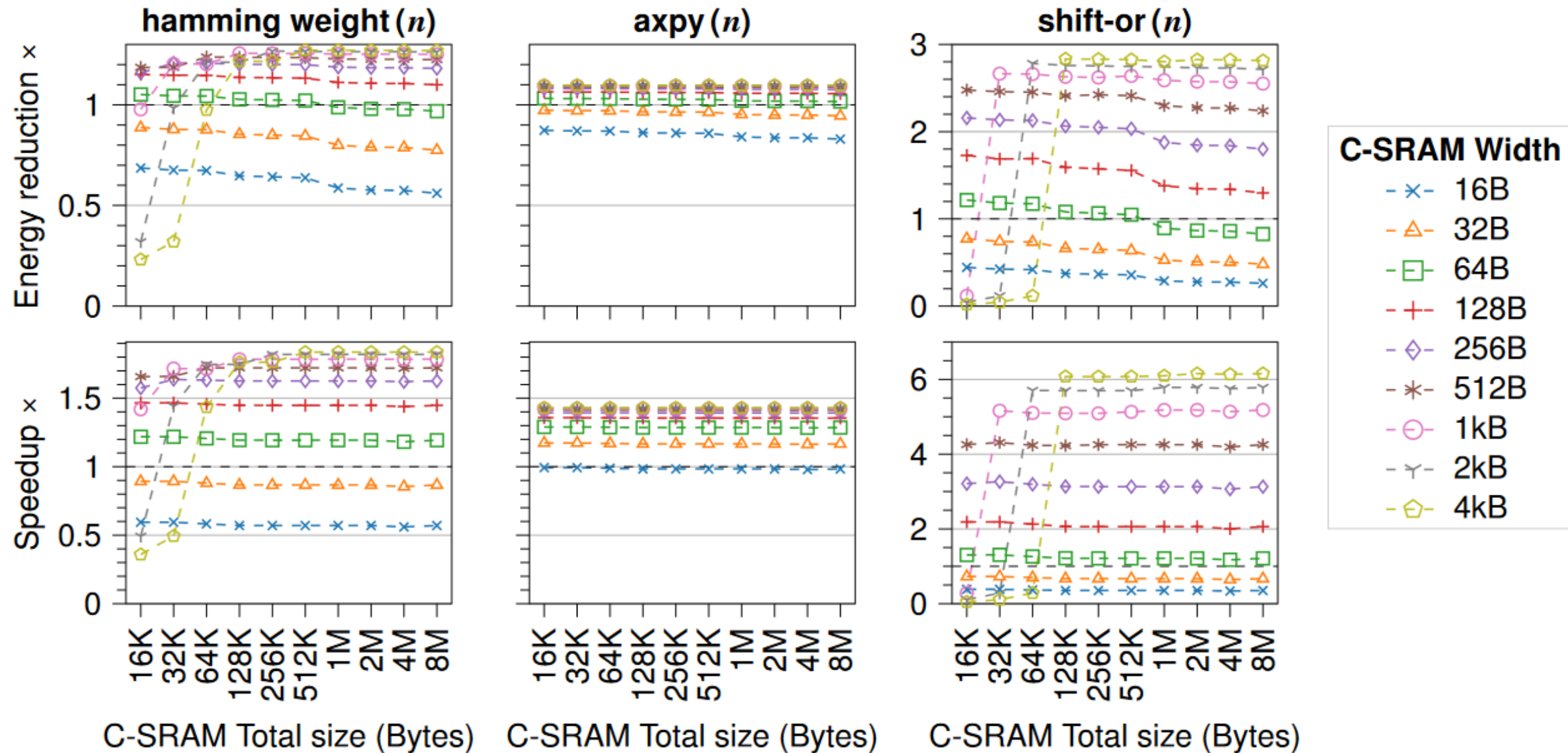


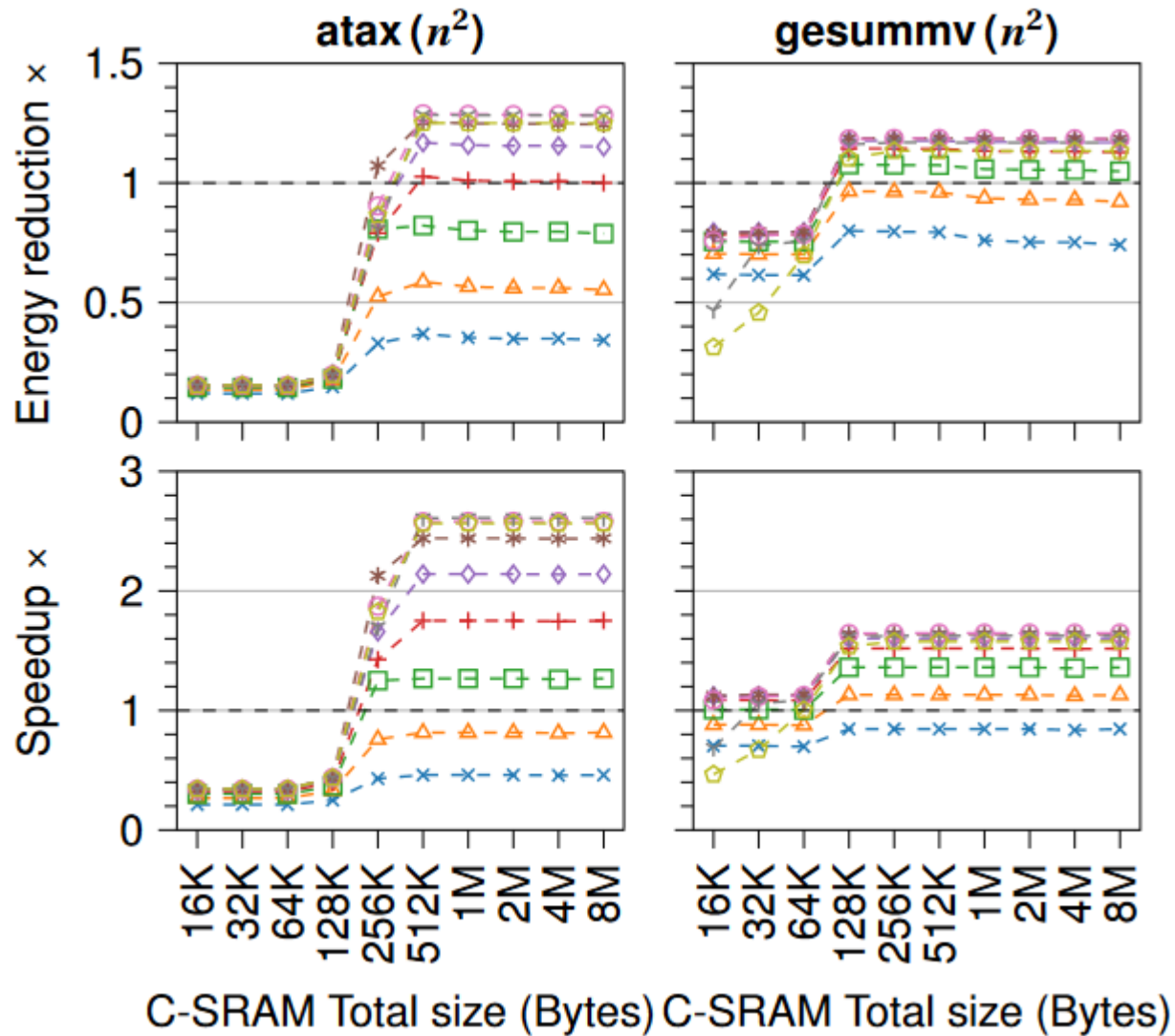
Base scenario but with paged size transfer

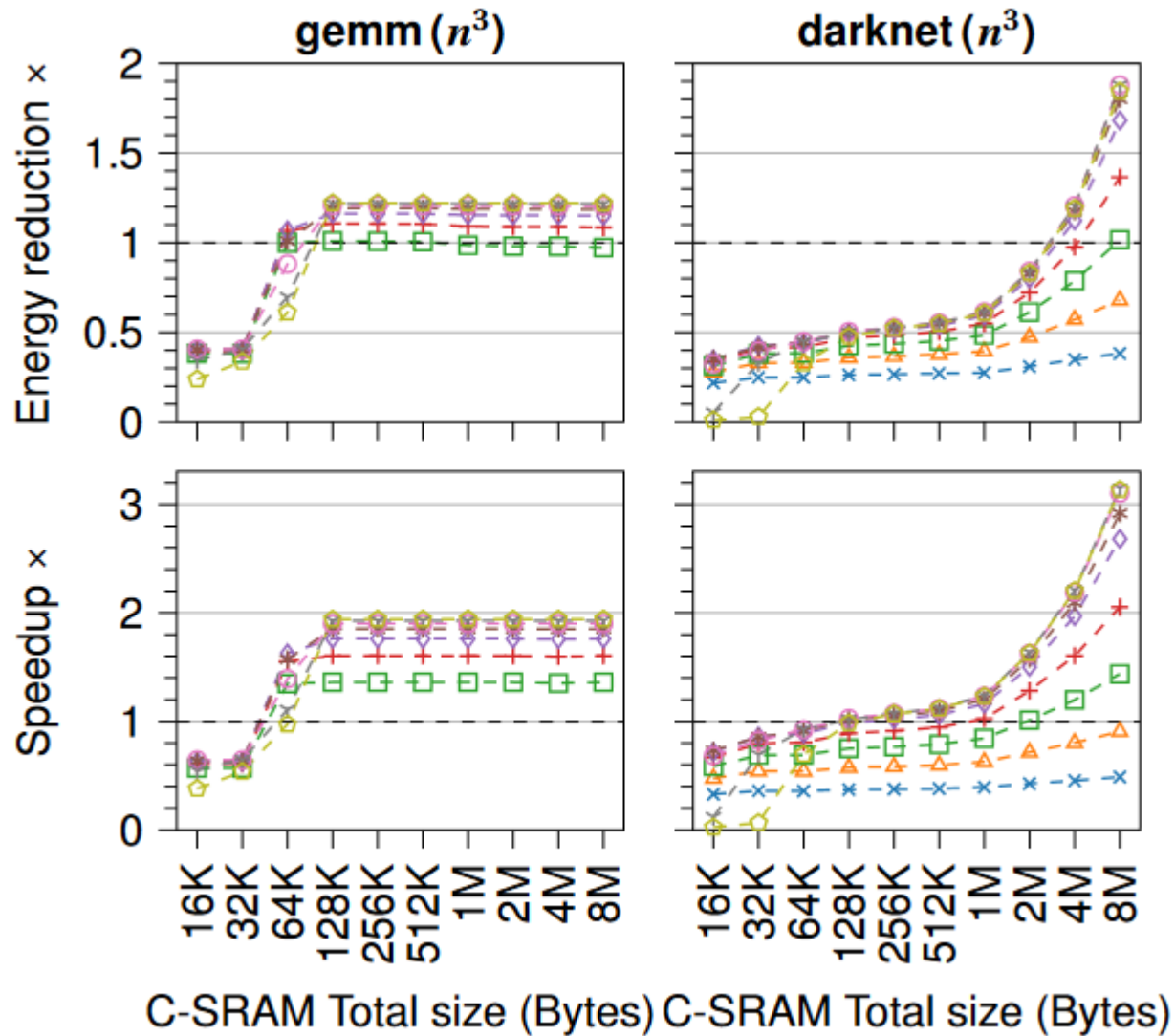
-x- hw -△- so -□- AXPY -+ atax -◇- gesummv -* gemm -○ darknet



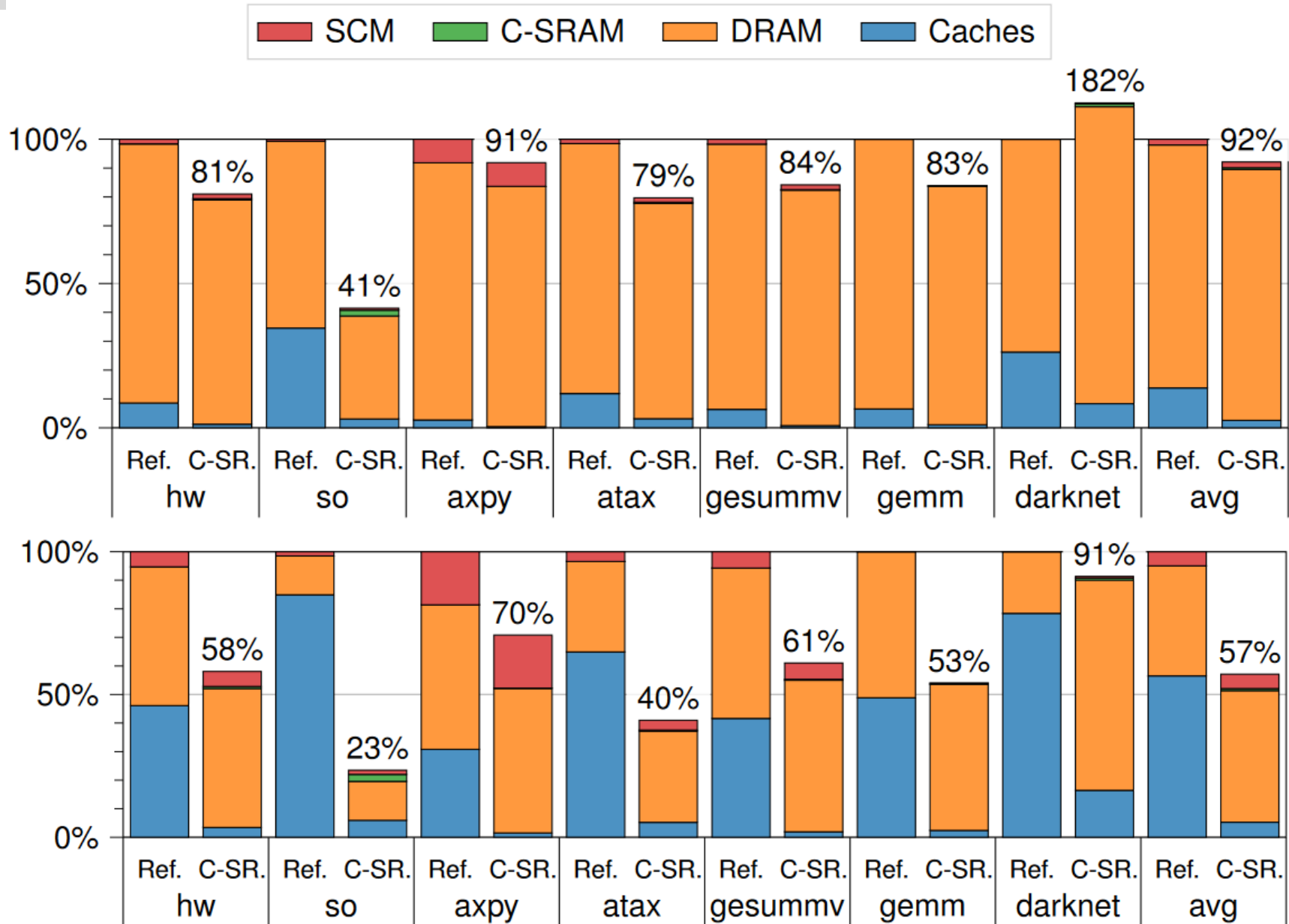


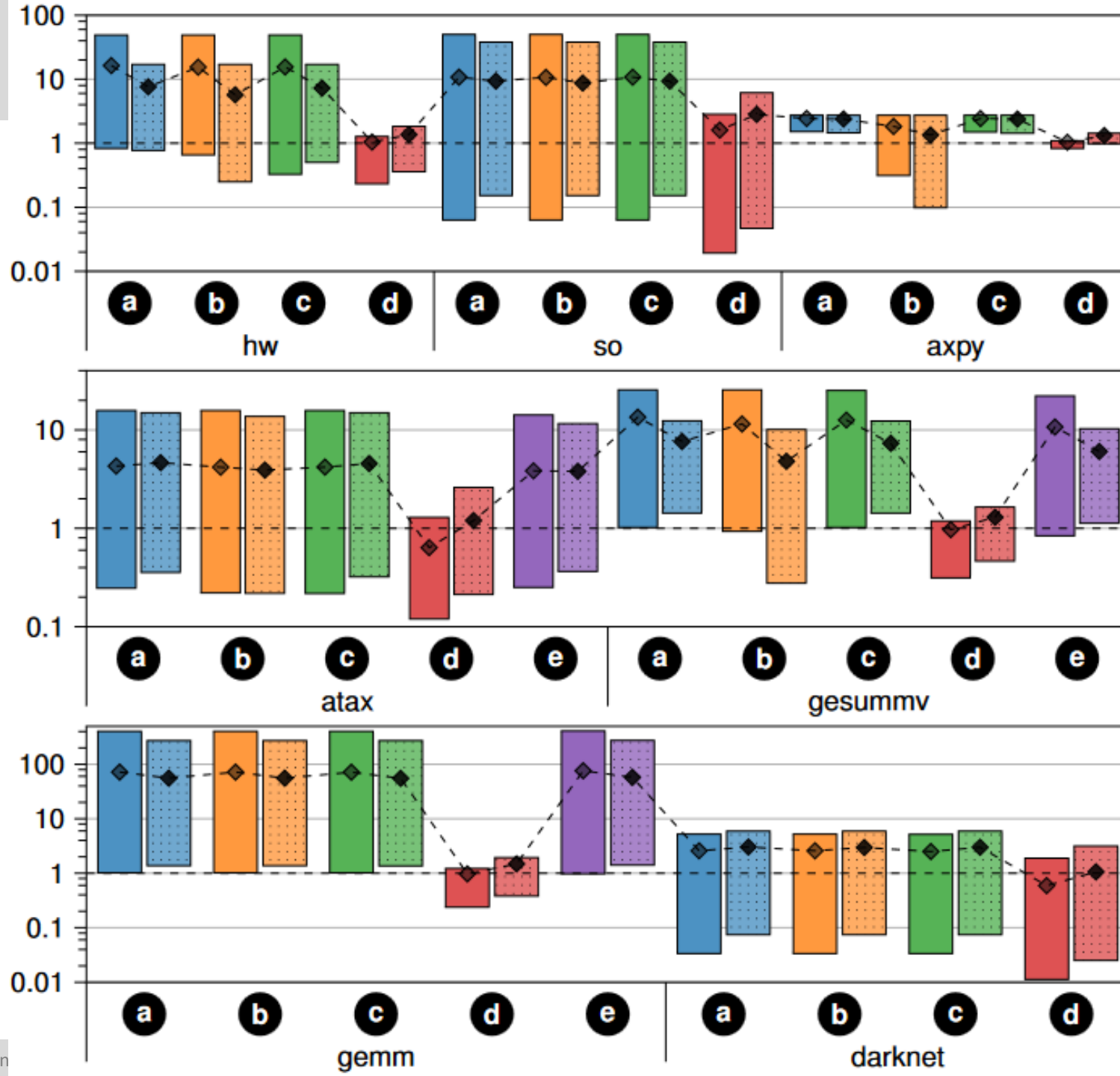




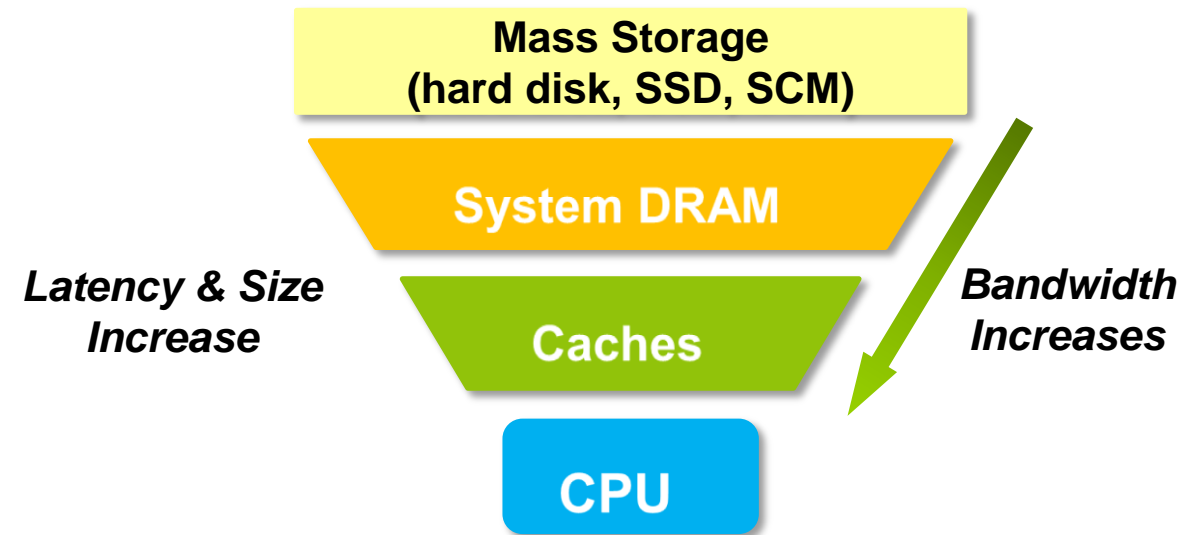


DRAM scenario distribution

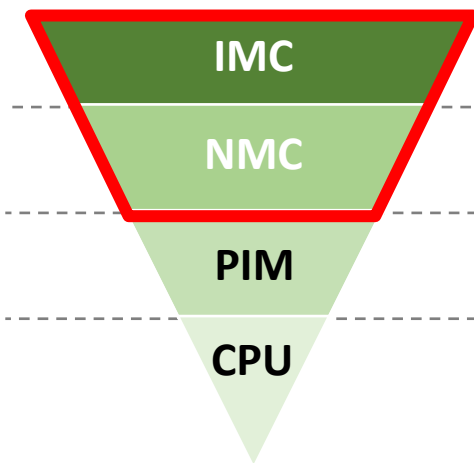




Copy paste



Memory Bandwidth & Size

closer to **Memory**

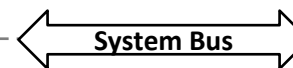
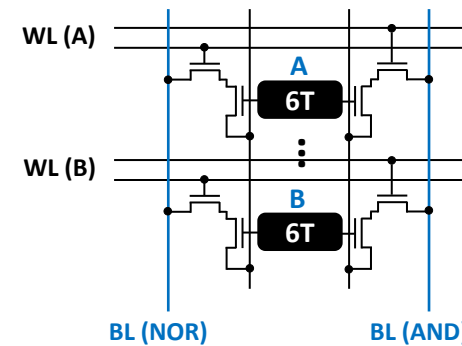
Interface levels

bit-cells array
(custom bitcell)

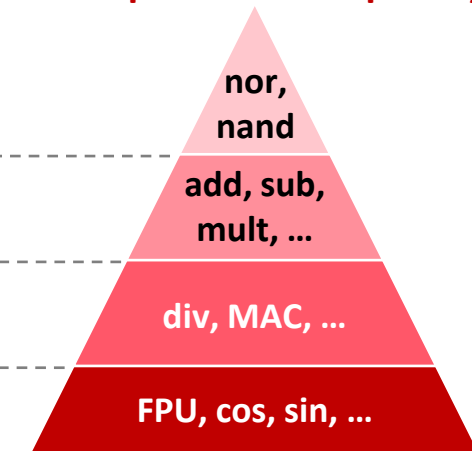
word-line IOs
(digital wrapper)

logic tight integration
(RISC CPU)

system integration

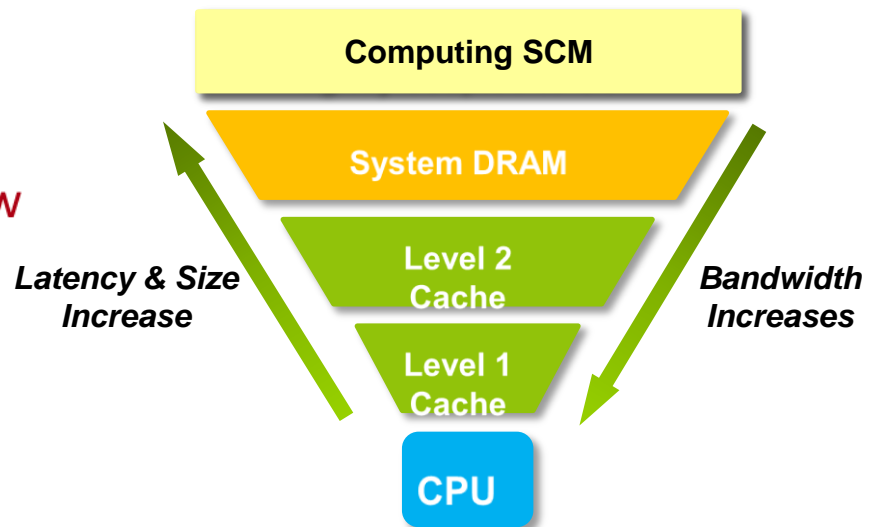
closer to **Processor**

Computation complexity



- Merge computing and data in the same circuit

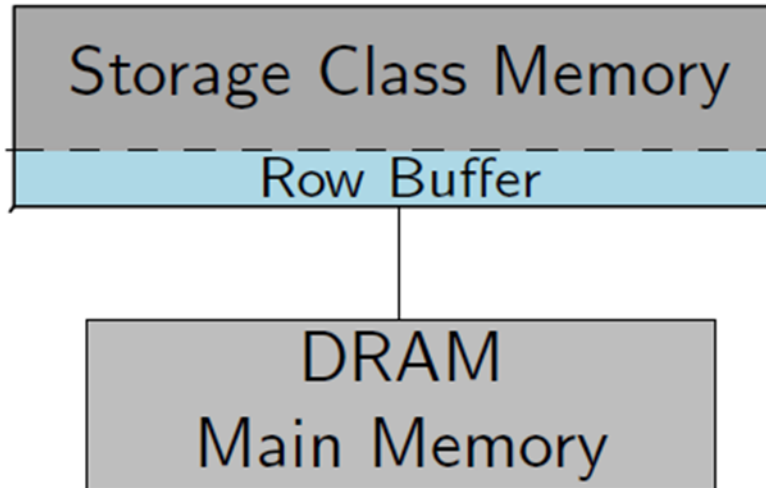
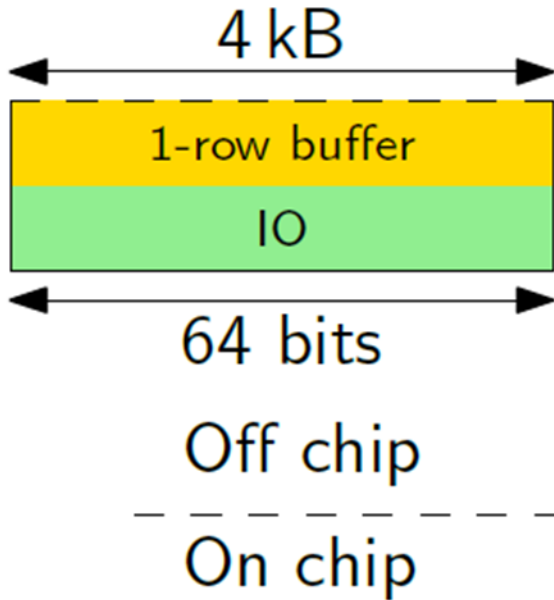
- ~~Low energy efficiency and low~~
- ~~High speed is a problem (10^9)~~
- Endurance is unlimited
- Use existing row buffer and repurpose it for computing



- DRAM is still used as a write buffer
- Caches can be reduced (Instructions & stack)

Let's move computing at the top with tightly coupled SRAM to get the best of both worlds!

1 Reference Architecture



2 C-RB Architecture

