



# Overcoming the Gap Between Compute and Memory Bandwidth in Modern GPUs

Lingqi Zhang<sup>1</sup>, Mohamed Wahib<sup>2</sup>, Toshio Endo<sup>1</sup>, Satoshi Matsuoka<sup>2,1</sup>  
1. Tokyo Institute of Technology, Japan 2. RIKEN Center for Computational Science, Japan

## OBSERVATIONS

### EASIER DEVICE SATURATION

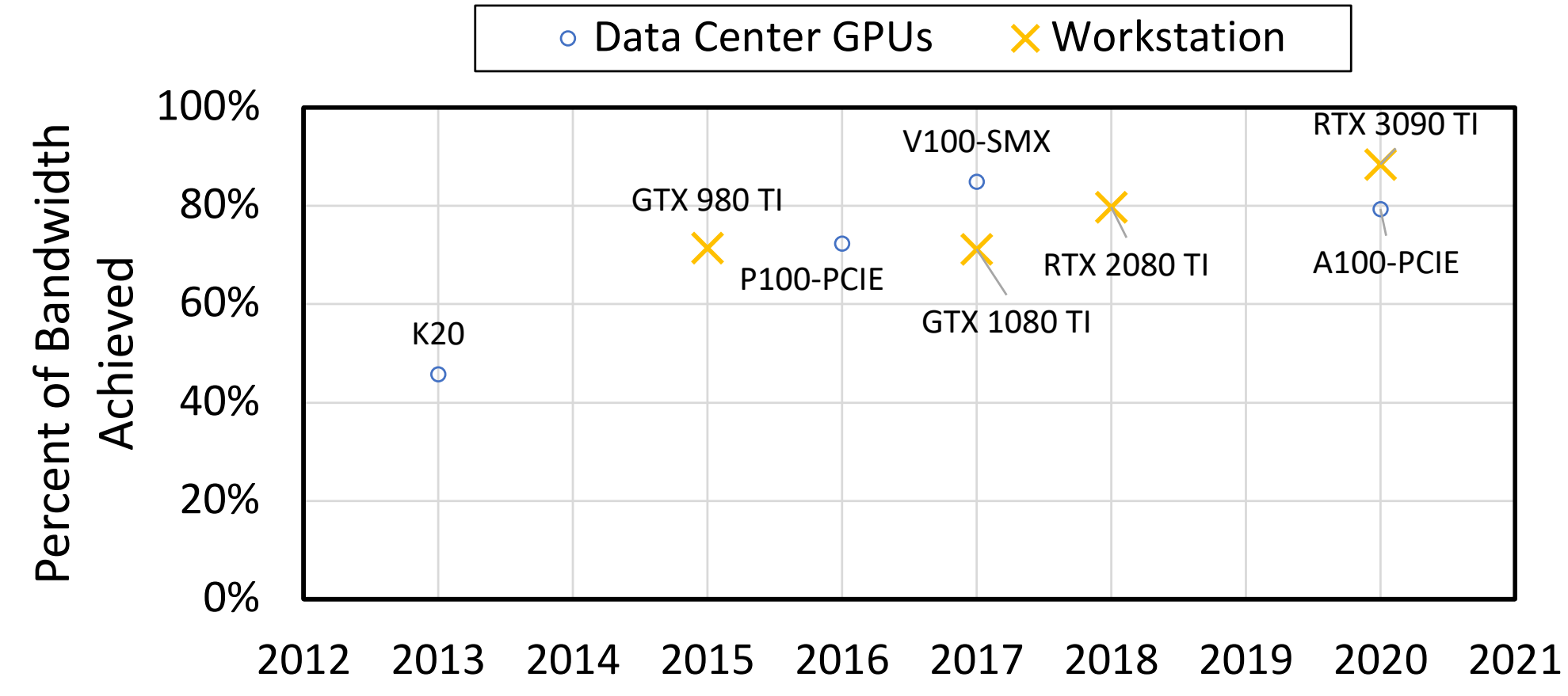
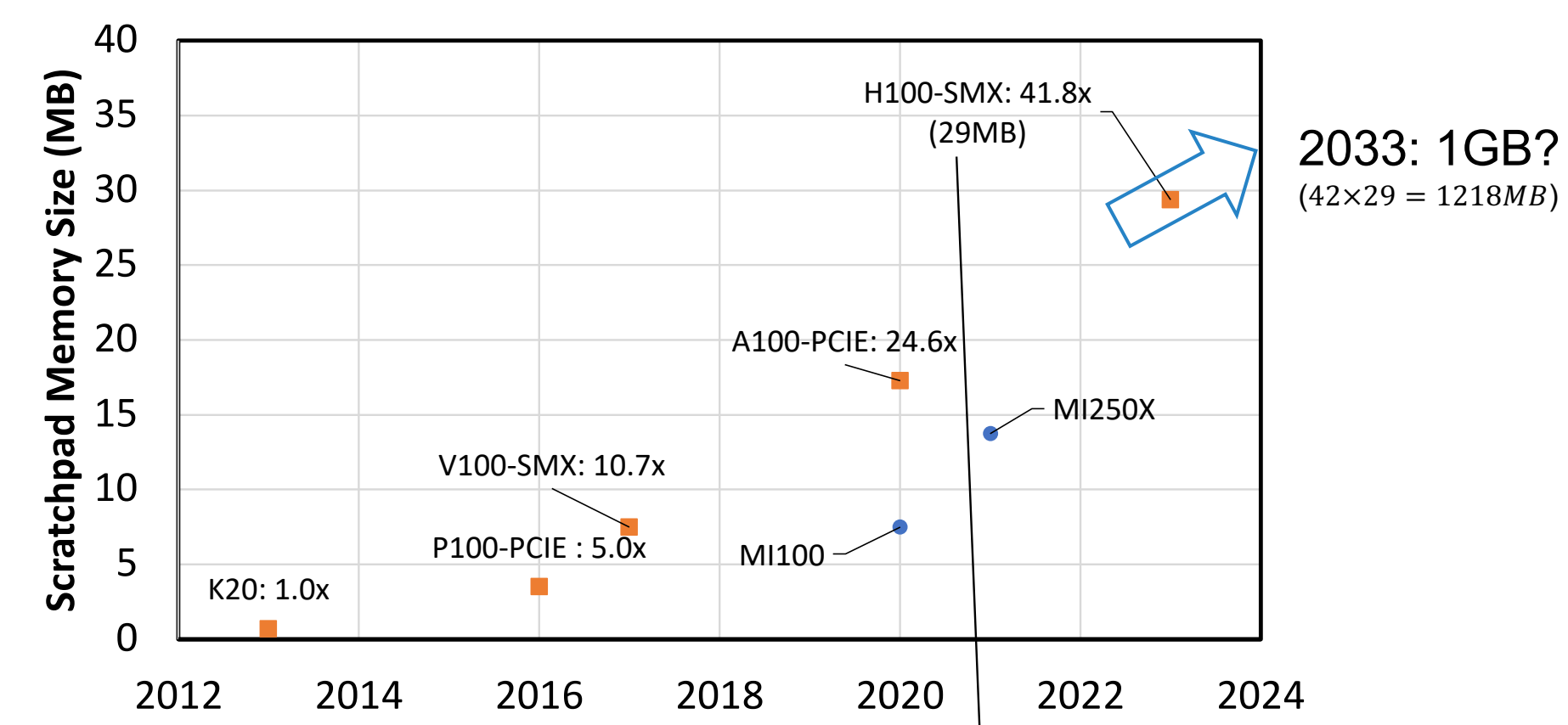


Figure 1: Percent of peak performance of STREAM kernel in with low occupancy (occupancy=256 threads, ILP=4, 8 Byte per memory access).

### LARGER ON-CHIP RESOURCES

Scratchpad memory as an example:



The size of shared memory has increased by 41.8x over a decade (peak performance has increased by ~22x)

Figure 2: The capacity trend of scratchpad memory.

### MATURE DEVICE SYNCHRONIZATION [3]

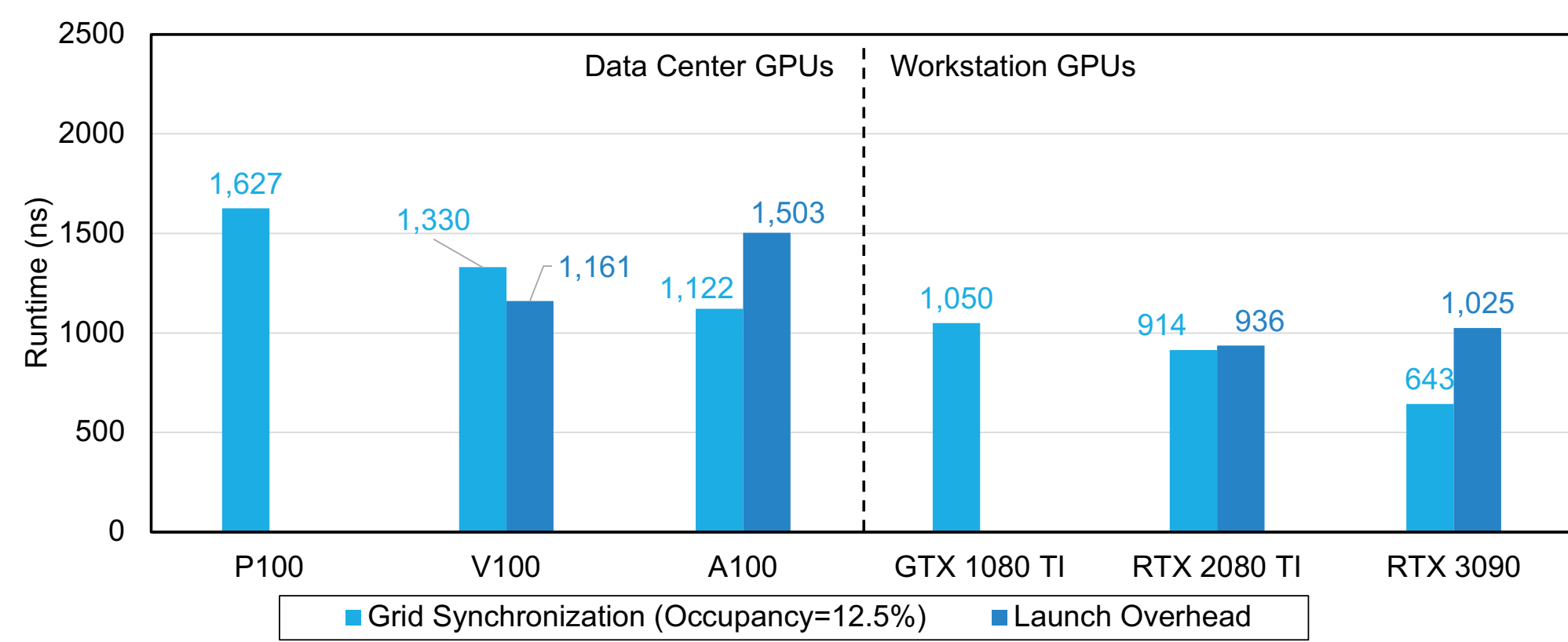


Figure 3: The overhead of device-wide synchronization (grid synchronization as explicit device-wide synchronization; kernel launch as implicit device-wide synchronization in the latest GPUs).

## REFERENCES

- L. Zhang, M. Wahib, P. Chen, J. Meng, X. Wang, T. Endo, and S. Matsuoka. Perks: A locality-optimized execution model for iterative memory-bound gpu applications. In *Proceedings of the 37th International Conference on Supercomputing, ICS '23*, page 167–179, New York, NY, USA, 2023. Association for Computing Machinery.
- L. Zhang, M. Wahib, P. Chen, J. Meng, X. Wang, T. Endo, and S. Matsuoka. Revisiting temporal blocking stencil optimizations. In *Proceedings of the 37th International Conference on Supercomputing, ICS '23*, page 251–263, New York, NY, USA, 2023. Association for Computing Machinery.
- L. Zhang, M. Wahib, H. Zhang, and S. Matsuoka. A study of single and multi-device synchronization methods in nvidia gpus. In *2020 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, pages 483–493, 2020.



## OVERVIEW

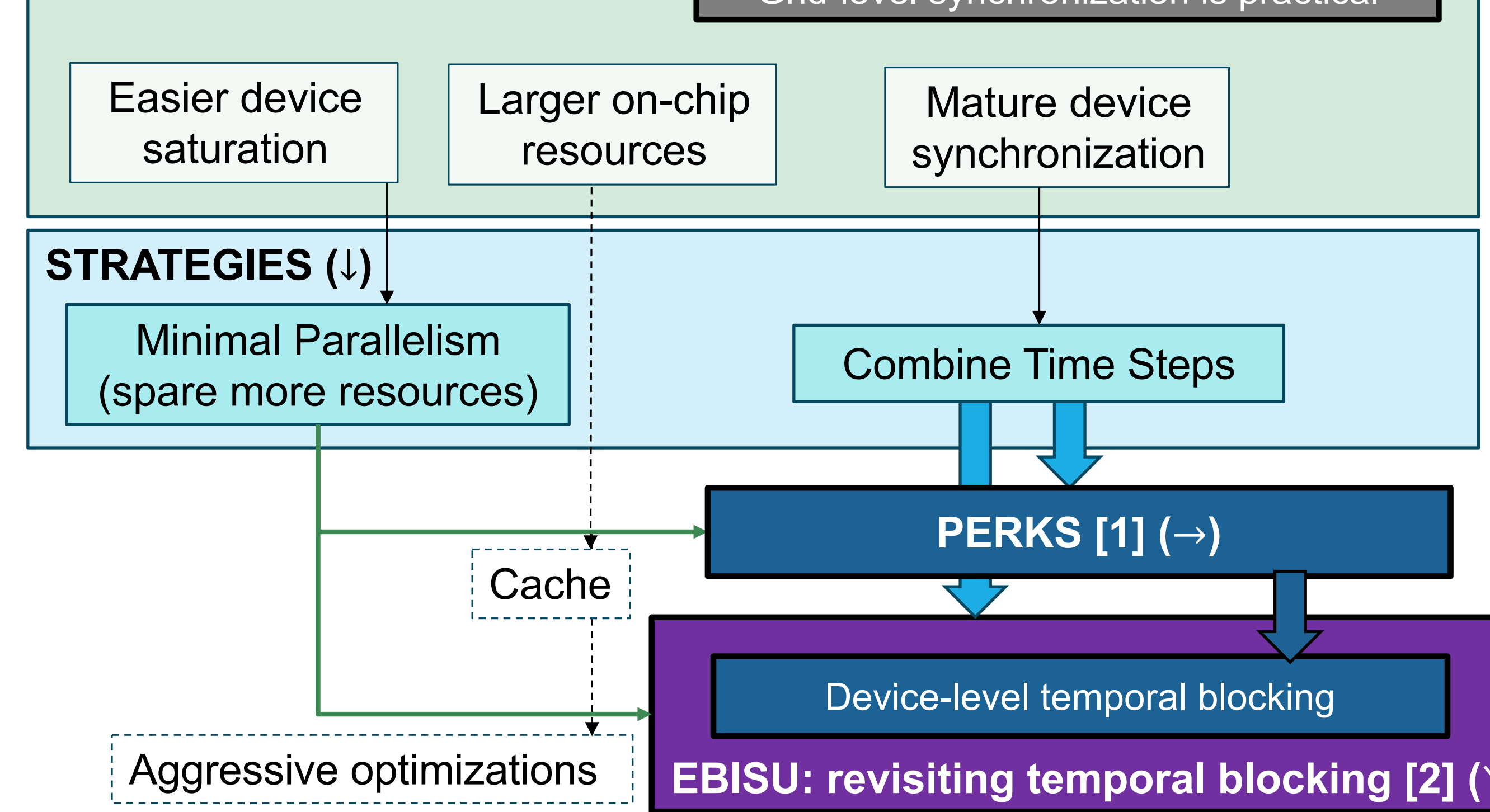
### MOTIVATION

- Compute is faster than memory bandwidth.
- The gap is still widening.
- As GPU architecture evolves, **can we leverage the new features to better overcome the gap?**
- We showcase our methodology with memory-bound kernels (e.g., stencils).

Platform	Launched Memory (Year)	Memory (GB/s)	Compute (TFLOPS/s)	Balance Flops/Bytes
Fujitsu (A64FX)	2019	1024	3.4	3.32
Intel (Platinum 8368)	2021	204.8	1.094	5.34
AMD (7773X)	2022	204.8	2.253	11
Nvidia (A100)	2020	1555	9.7	6.24
Nvidia (H100-SMX)	2022	3000	30	10
AMD (MI250X)	2021	3200	47.9	14.97

### OBSERVATIONS (←)

Observing the GPU trends



## STRATEGIES

### Minimal Parallelism

Orchestrating parallelism.

Little's Law (Hardware)	
C	Concurrency;
L	Latency;
THR	Throughput;
Parallelisms (Software)	
PAR	Parallelism of a program;
ILP	Instruction Level Parallelism;
TLP	Thread Level Parallelism (thread per Stream Multiprocessor);

$$\begin{aligned} &\text{minimize } PAR(TLP, ILP) \\ &\text{subject to } PAR \geq C \\ &PAR = TLP \times ILP \\ &C = L \times THR \end{aligned} \quad (1)$$

### Combine Time Steps

Reducing/Eliminating memory traffic.

Roofline	
I	Operation Intensity;
W	Number of works;
Q	Number of bytes of memory traffic;
Additional Parameters	
t	Combined time steps;
R	Memory traffic reduced in between time steps;

$$I = \frac{W \times t}{Q + (t-1) \times (1-R) \times Q} \quad (2)$$

Increasing combined time steps  $t$  and percent of memory traffic cached  $R$  to increase Operation Intensity  $I$ . As such, the memory-bound kernel becomes closer to compute-bound. The kernel becomes temporal blocking if  $R = 1$ .

## PERKS [1]

### MOTIVATION

Stencil:

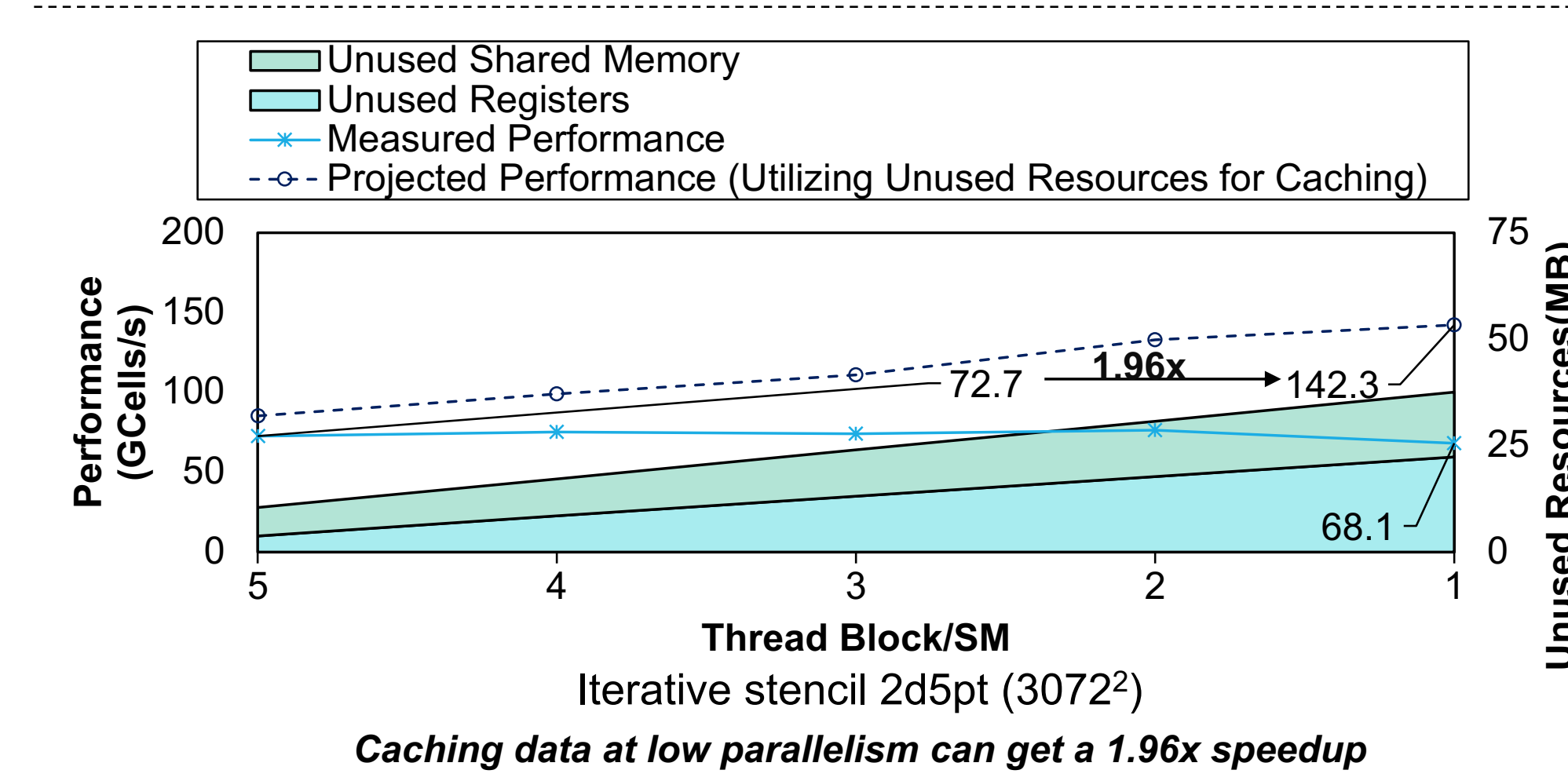
Machine Learning

PDE applications:

- Weather modeling
- Fluid dynamics simulation

$$A^t(s) = \sum_{a \in \mathcal{N}} w_a \times A^{t-1}(s+a) + c$$

$s$ : address,  $\mathcal{N}$ : concerned neighbours  
 $w$ : weight,  $c$ : constant



### EVALUATION

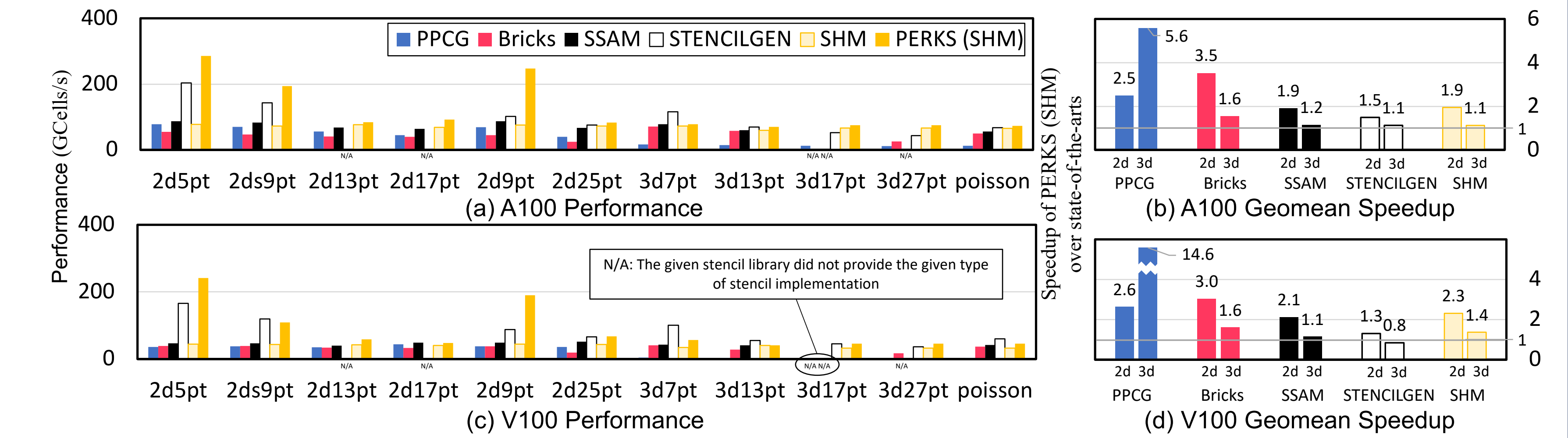
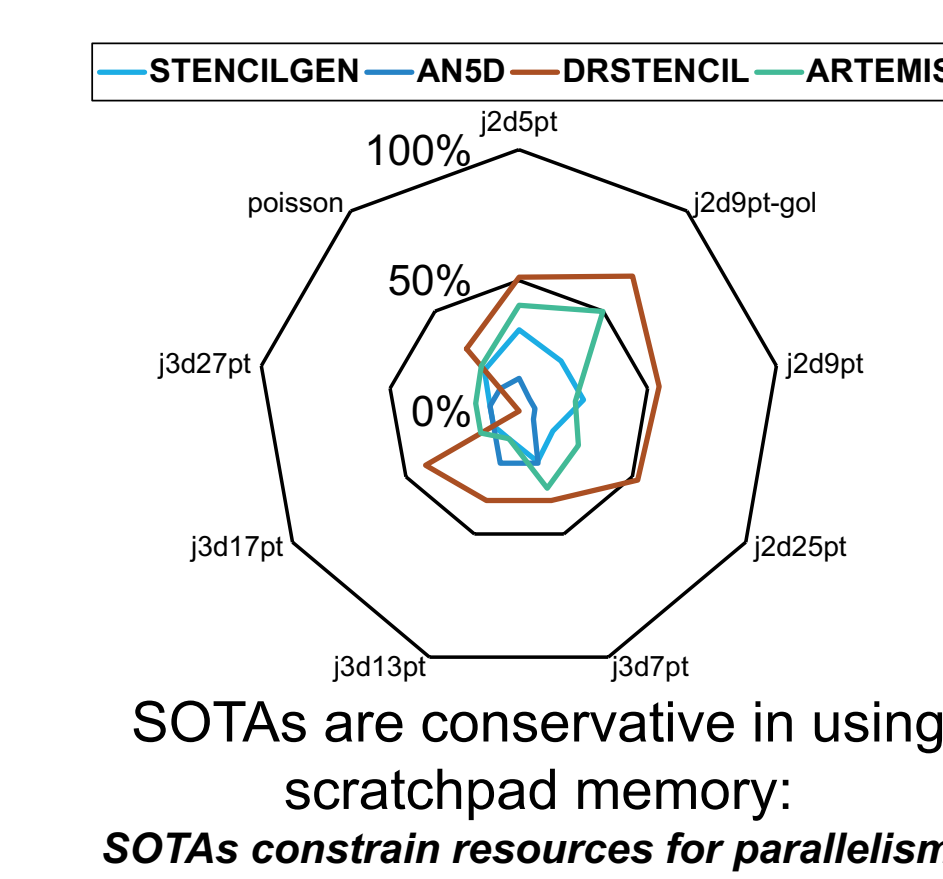


Figure 4: Comparison of PERKS(SHM) over a wide range of stencil libraries.

## EBISU [2]

### MOTIVATION



### EVALUATION

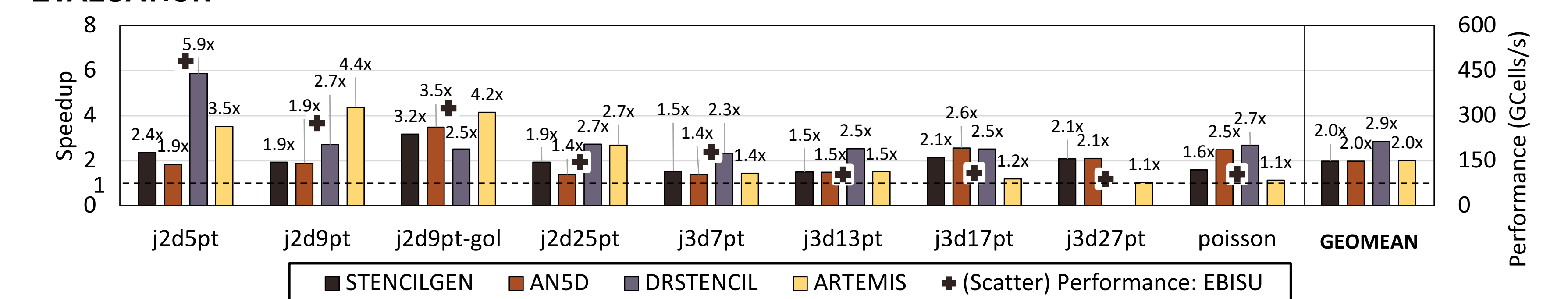


Figure 5: Speedup of EBISU over the state-of-the-art temporal blocking implementations. We also plot the performance of EBISU (right Y-axis plotted as '+' ticks).