

Raise the roofline — hardware accelerators

Hung-Wei Tseng

Recap: from the software perspective

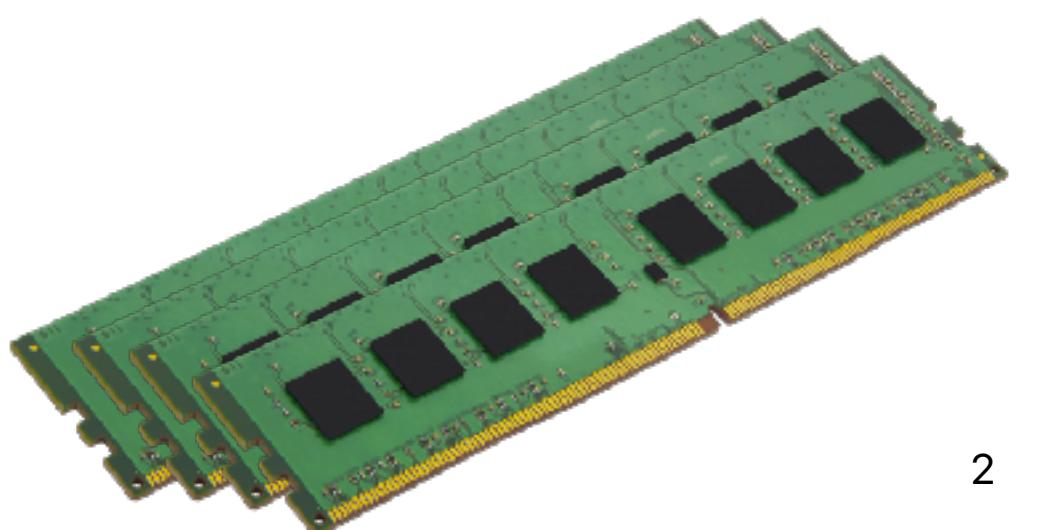


Should I use NDP?

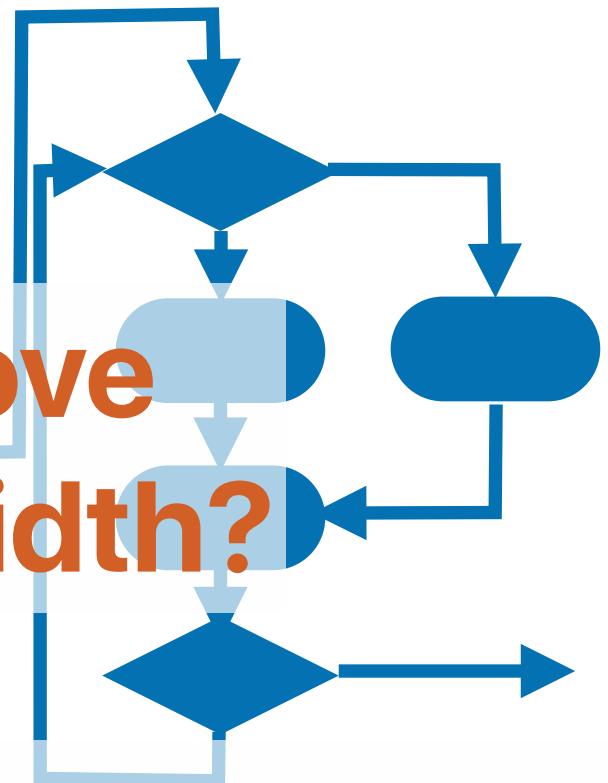
Source "data"



"Data" structures



Should I improve
memory bandwidth?



Should I use
accelerators?
Algorithms



Result

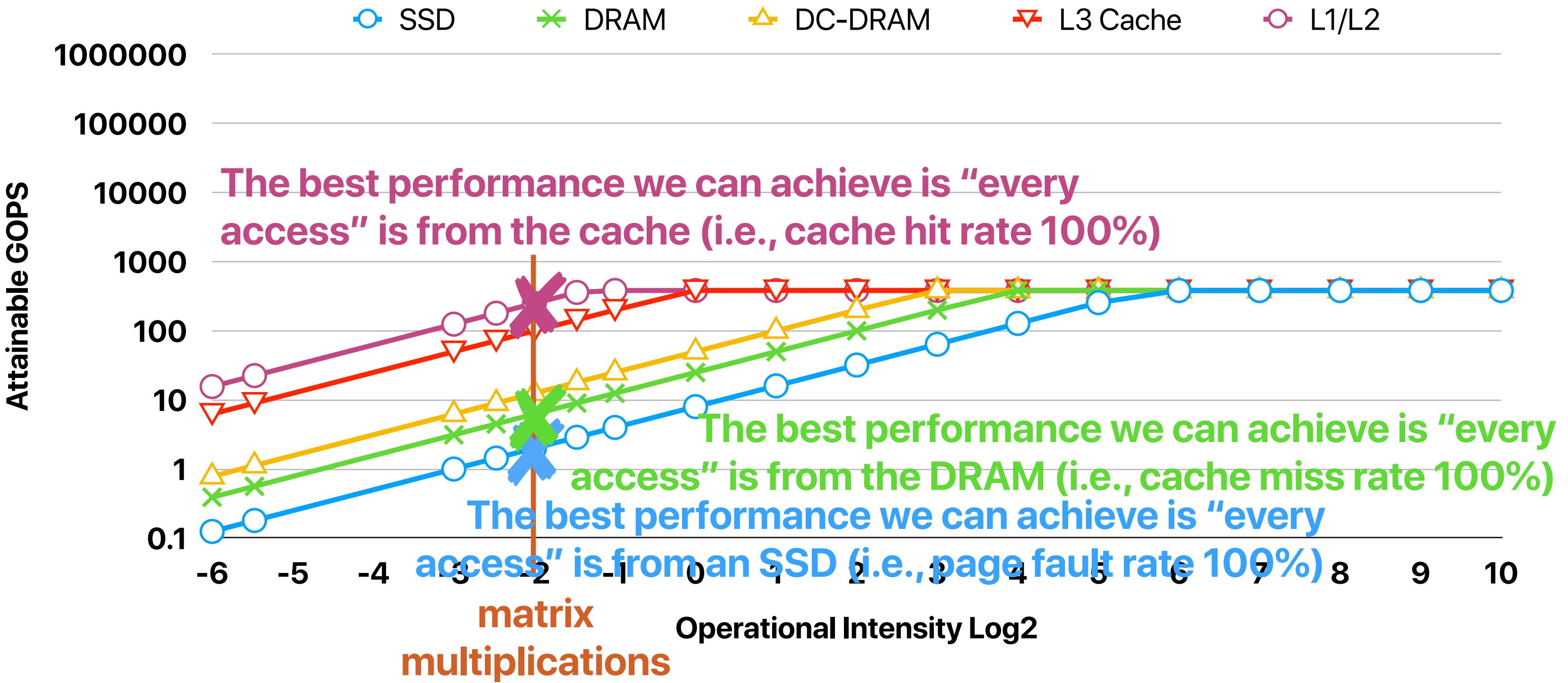


Recap: the roofline model

- The speed of computation — determined by the “OPS” (operations per second) of the processing unit
- The speed of data supply — determined by the “effective bandwidth” of the data path
- At any point, the slower one determines the performance

- $$\text{Attainable GFlops/sec} = \min \left\{ \frac{\text{Peak Floating-Point Performance}}{\text{Peak Memory Bandwidth} \times \text{Operational Intensity}} \right\}$$

Recap: the roofline of a CPU-based system



Outline

- Raise the roof
 - GPU
 - Tensor processing units
 - Tensor cores and ray tracing units

What's the operational intensity of modern workloads?

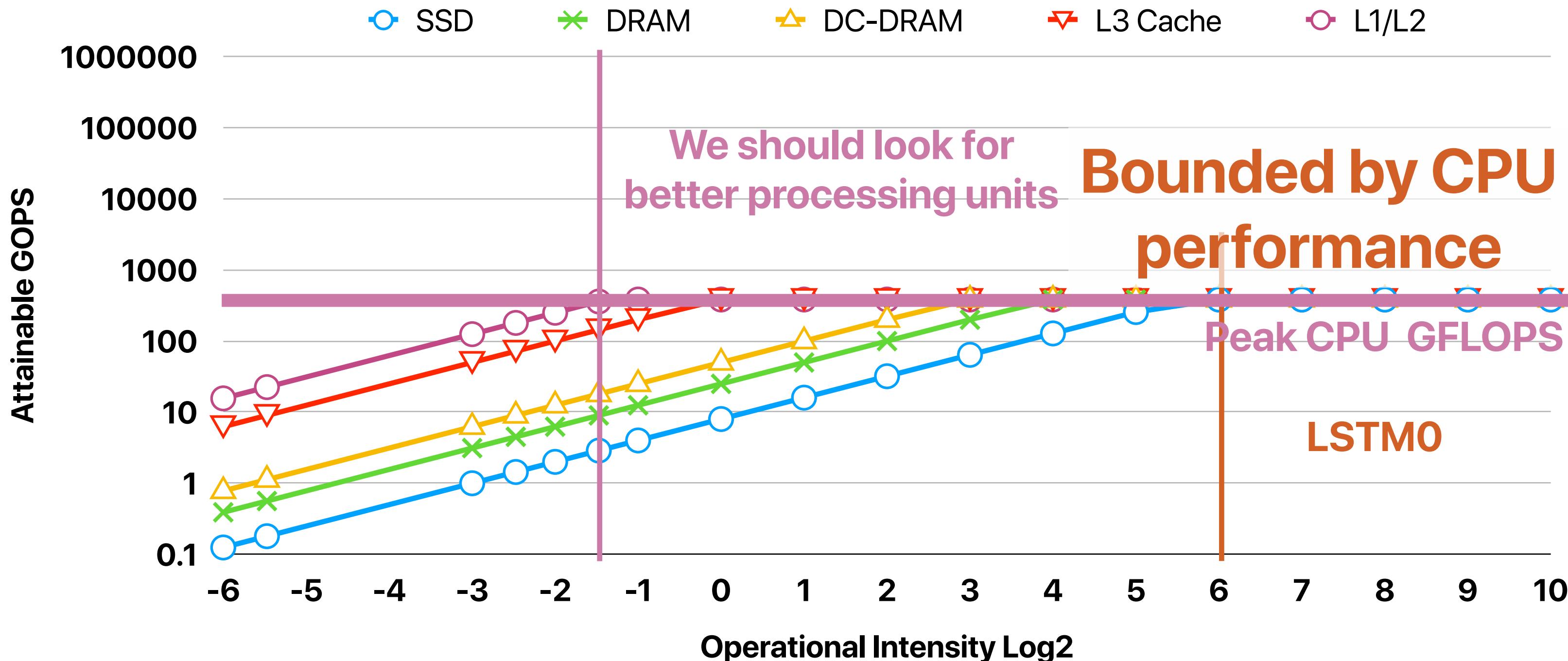
Table 1. Six DNN applications (two per DNN type). MLP stands for multi-layer perceptron, and LSTM stands for long short-term memory.

Name	Lines of Code	Weights	TPU Ops / Weight Byte	Operational Intensity Log2
MLP0	100	20 M	200	7.64385618977472
MLP1	1,000	5 M	168	7.39231742277876
LSTM0	1,000	52 M	64	6
LSTM1	1,500	34 M	96	6.58496250072116
CNN0	1,000	8 M	2,888	11.4958550268872
CNN1	1,000	100 M	1,750	10.7731392067197

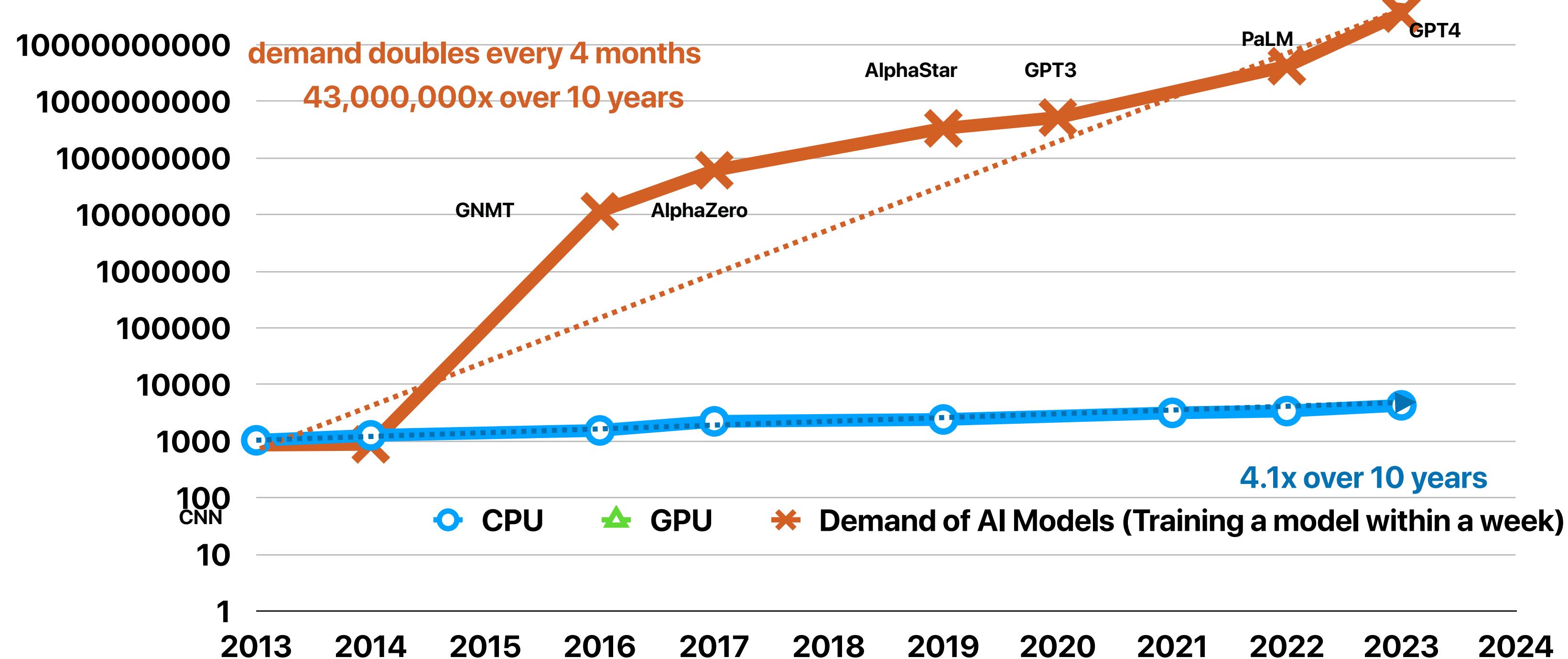
N. Jouppi, C. Young, N. Patil and D. Patterson, "Motivation for and Evaluation of the First Tensor Processing Unit," in *IEEE Micro*, vol. 38, no. 3, pp. 10-19, May./Jun. 2018, doi: 10.1109/MM.2018.032271057.

Recap: on conventional computers with CPUs and DDR DRAM as the main processing and runtime data storage, are AI/ML workloads compute-bound or memory-bound?

Where is LSTM0 in the roofline?



Mis-matching AI/ML demand and general-purpose processing



<https://ourworldindata.org/grapher/artificial-intelligence-training-computation>

Dennardian Broken

- Given a scaling factor S

Parameter	Relation	Classical Scaling	Leakage Limited
Power Budget		1	1
Chip Size		1	1
Vdd (Supply Voltage)		1/S	1
Vt (Threshold Voltage)	1/S	1/S	1
tex (oxide thickness)		1/S	1/S
W, L (transistor dimensions)		1/S	1/S
Cgate (gate capacitance)	WL/tox	1/S	1/S
I_{sat} (saturation current)	WVdd/tox	1/S	1
F (device frequency)	$I_{sat}/(C_{gate}V_{dd})$	S	S
D (Device/Area)	$1/(WL)$	S^2	S^2
p (device power)	$I_{sat}V_{dd}$	$1/S^2$	1
P (chip power)	D_p	1	S^2
U (utilization)	$1/P$	1	$1/S^2$

Recap: Power consumption to light on all transistors

Chip							
1	1	1	1	1	1	1	1
1	1	1	1	1	1	1	1
1	1	1	1	1	1	1	1
1	1	1	1	1	1	1	1
1	1	1	1	1	1	1	1
1	1	1	1	1	1	1	1
1	1	1	1	1	1	1	1
1	1	1	1	1	1	1	1

=49W

Dennardian Scaling

Chip							
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5

=50W

Dennardian Broken

Chip							
1	1	1	1	1	1	1	1
1	1	1	1	1	1	1	1
1	1	1	1	1	1	1	1
1	1	1	1	1	1	1	1
1	1	1	1	1	1	1	1
1	1	1	1	1	1	1	1
1	1	1	1	1	1	1	1
1	1	1	1	1	1	1	1

On ~ 50W
Off ~ 0W
Dark!

=100W!

Alternative computing resources

GPUs as alternative computing resource

DDR4

25 GB/sec

DRAM

CPU-
Memory
BUS

Memory
Controller



PCI
EXPRESS®

PCIe Root
Complex

GPU



PCIe 4.0 x 4 ~
8GB/sec

HDD

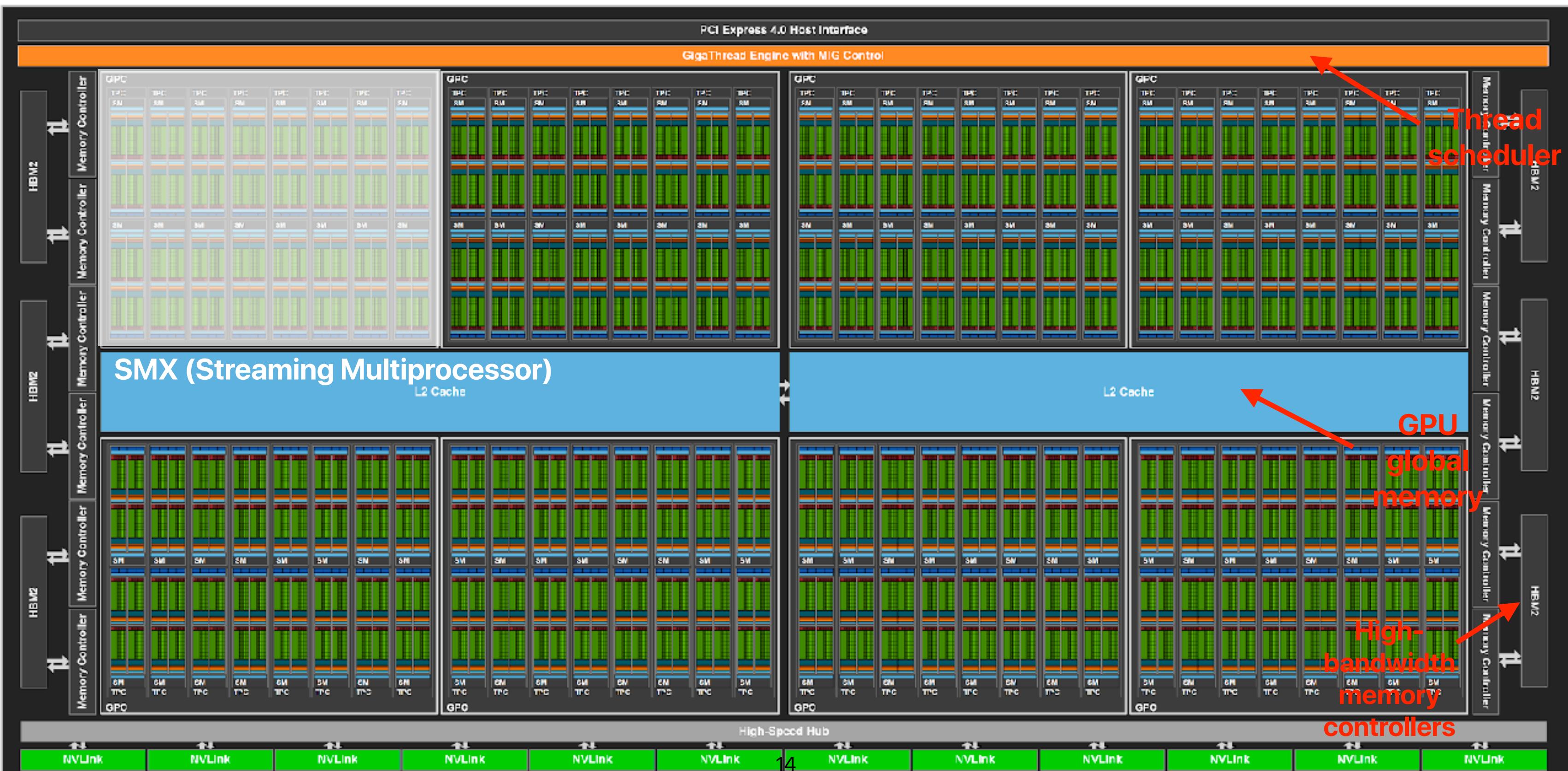


SSD



NIC

GPU Architecture



Inside an SM



A total of $16 \times 4 = 64$ FP32 cores

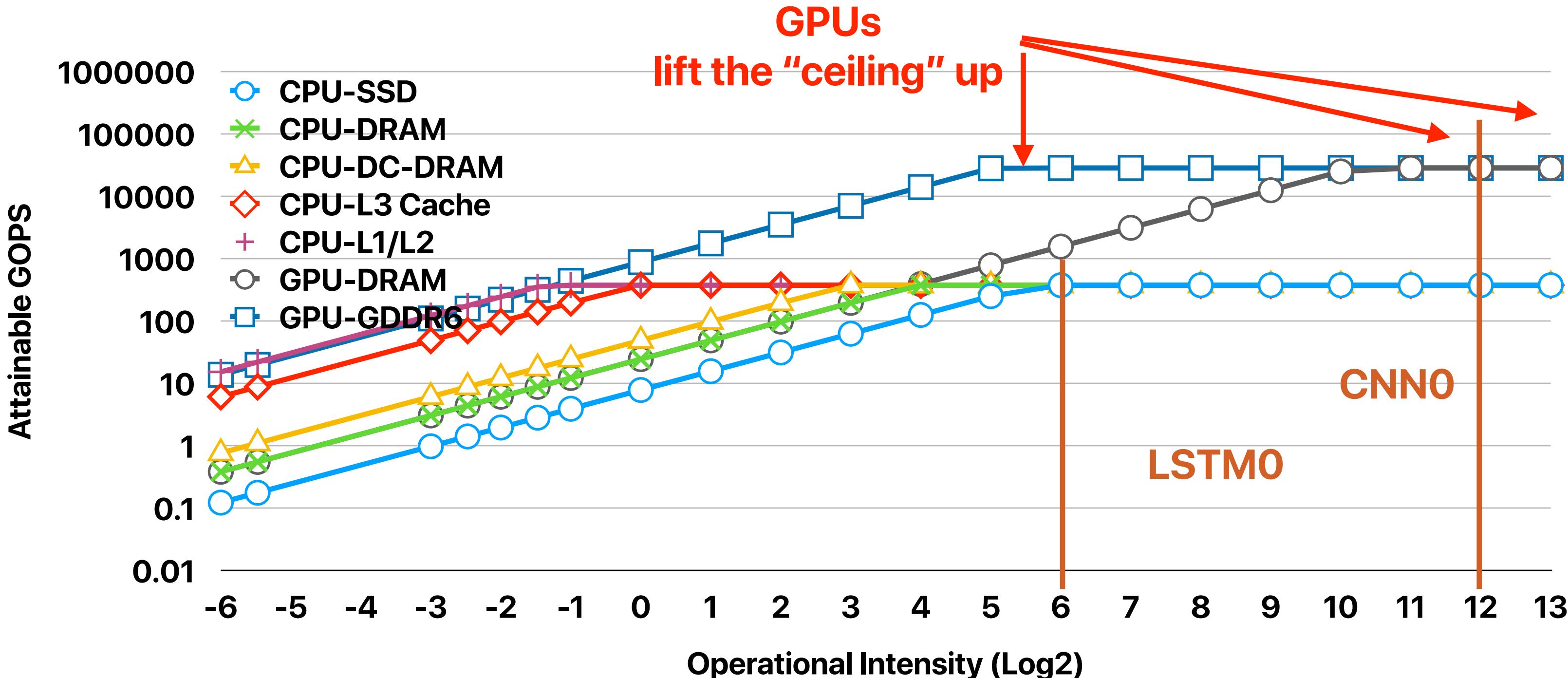
A total of $16 \times 4 = 64$ INT32 cores

A total of $16 \times 4 = 16$ FP64 cores

- All of these can only perform the same operation at the same time, but each of these is named as a "thread" in CUDA
- You can only use either FP32, FP64, INT32 and "Tensor Cores" at the same time

How will GPUs change the roofline model?

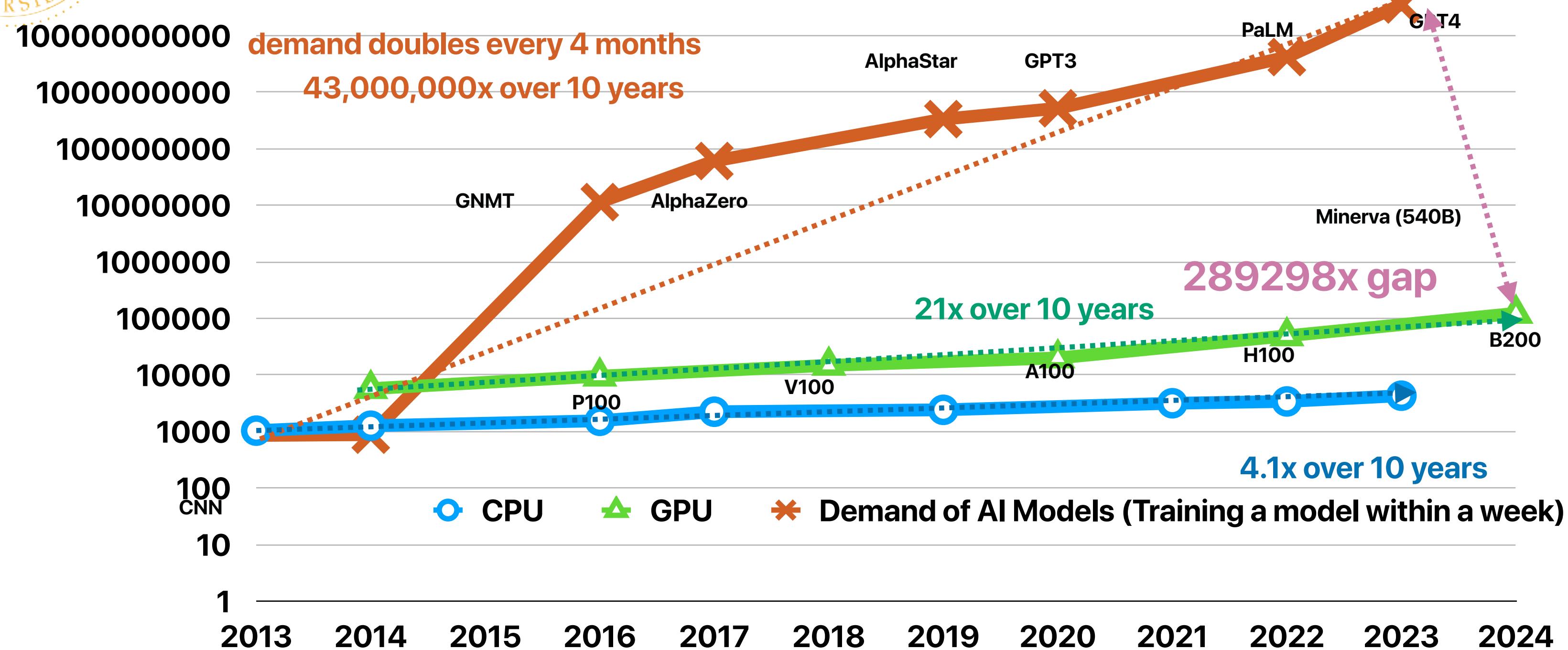
Placing GPUs in the roofline



For future AI/ML applications, should we look for more powerful computing resources? Why or why not?

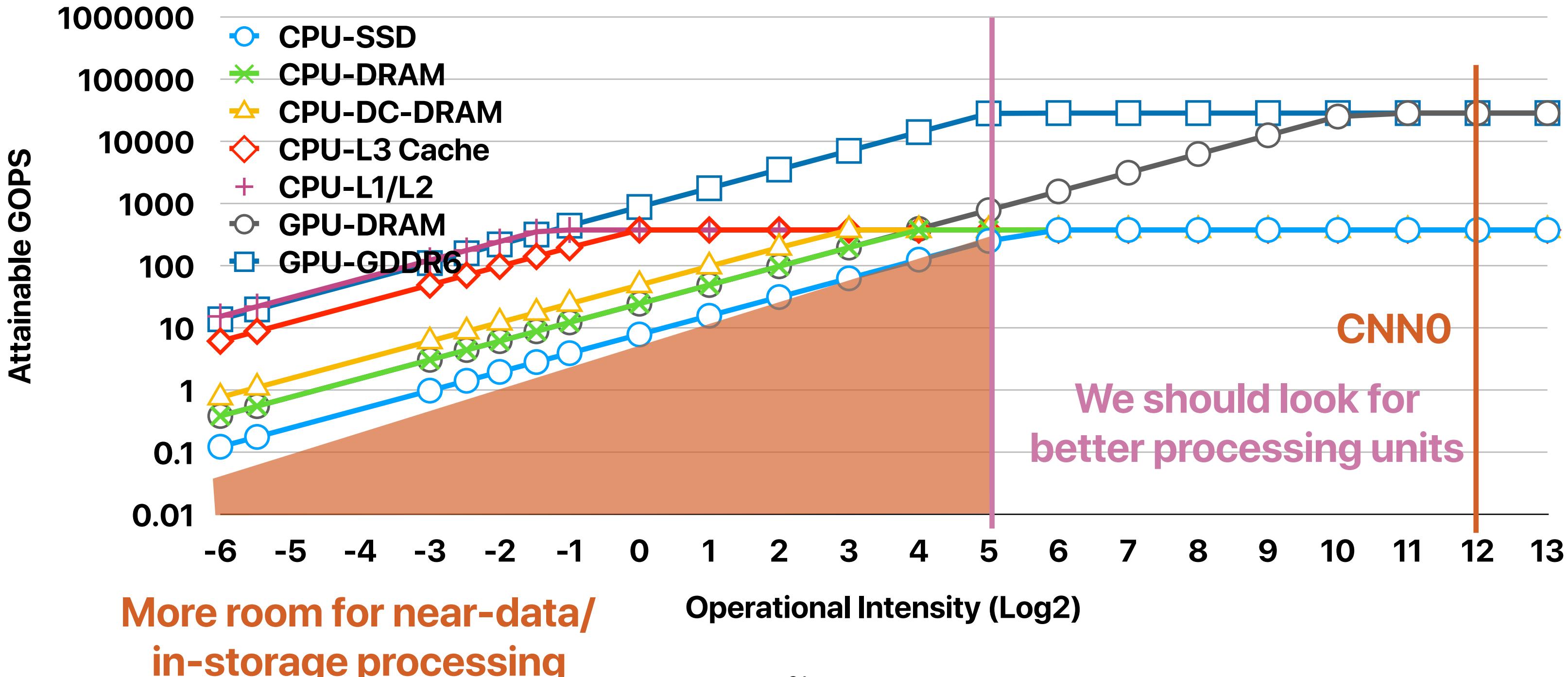


Mis-matching AI/ML demand and general-purpose processing



<https://ourworldindata.org/grapher/artificial-intelligence-training-computation>

Placing GPUs in the roofline

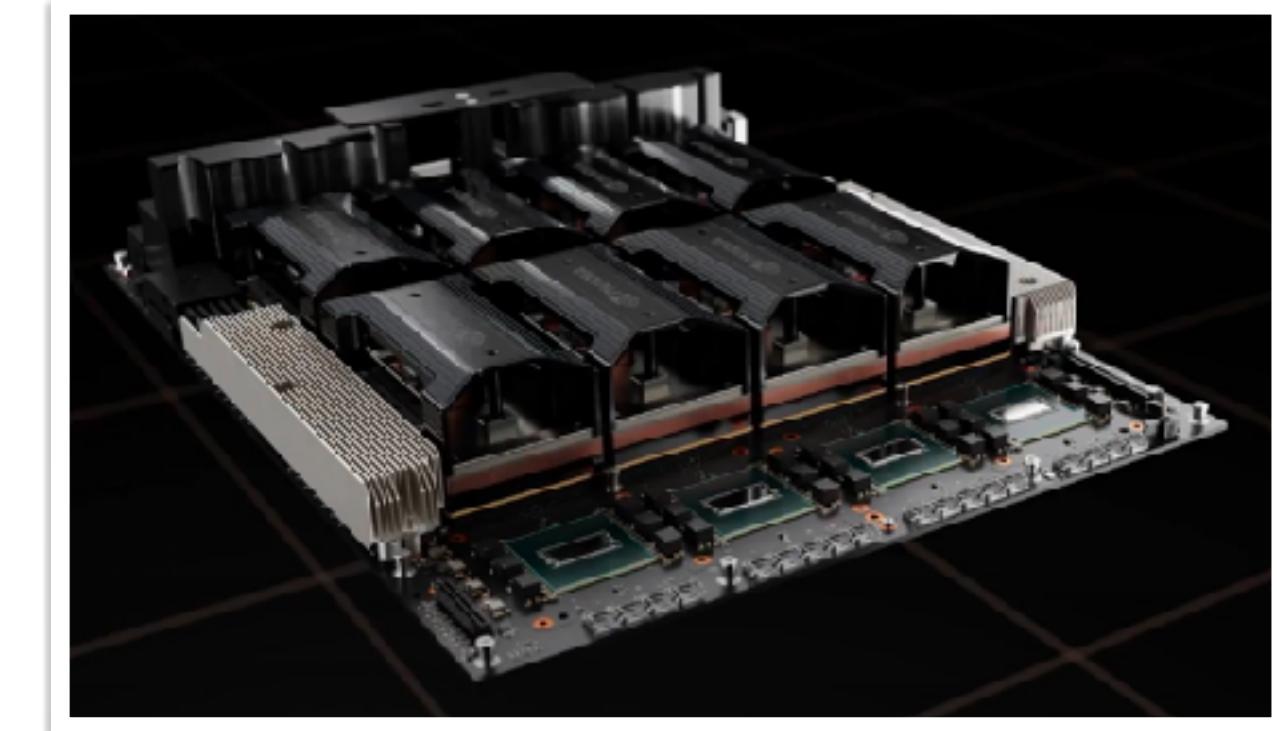


Power/energy is the main issue

NVIDIA Accelerator Specification Comparison			
	H100	A100 (80GB)	V100
FP32 CUDA Cores	16896	6912	5120
Tensor Cores	528	432	640
Boost Clock	~1.78GHz (Not Finalized)	1.41GHz	1.53GHz
Memory Clock	4.8Gbps HBM3	3.2Gbps HBM2e	1.75Gbps HBM2
Memory Bus Width	5120-bit	5120-bit	4096-bit
Memory Bandwidth	3TB/sec	2TB/sec	900GB/sec
VRAM	80GB	80GB	16GB/32GB
FP32 Vector	60 TFLOPS	19.5 TFLOPS	15.7 TFLOPS
FP64 Vector	30 TFLOPS	9.7 TFLOPS (1/2 FP32 rate)	7.8 TFLOPS (1/2 FP32 rate)
INT8 Tensor	2000 TOPS	624 TOPS	N/A
FP16 Tensor	1000 TFLOPS	312 TFLOPS	125 TFLOPS
TF32 Tensor	500 TFLOPS	156 TFLOPS	N/A
FP64 Tensor	60 TFLOPS	19.5 TFLOPS	N/A
Interconnect	NVLink 4 18 Links (900GB/sec)	NVLink 3 12 Links (600GB/sec)	NVLink 2 6 Links (300GB/sec)
GPU	GH100 (814mm ²)	GA100 (826mm ²)	GV100 (815mm ²)
Transistor Count	80B	54.2B	21.1B
TDP	700W	400W	300W/350W
Manufacturing Process	TSMC 4N	TSMC 7N	TSMC 12nm FFN
Interface	SXM5	SXM4	SXM2/SXM3
Architecture	Hopper	Ampere	Volta



<https://www.workstationspecialist.com/product/nvidia-tesla-a100/>



<https://www.servethehome.com/wp-content/uploads/2022/03/NVIDIA-GTC-2022-H100-in-HGX-H100.jpg>

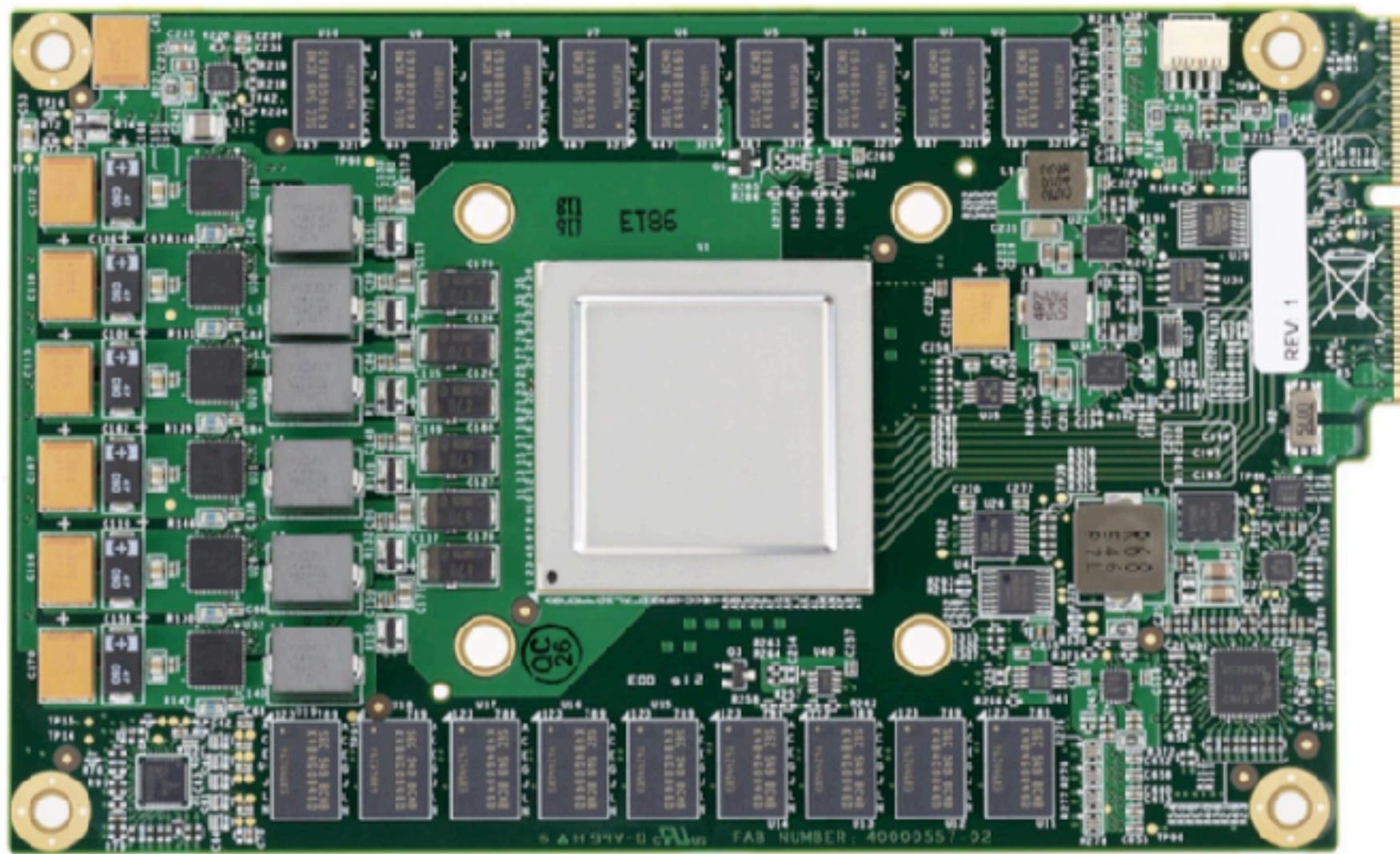
**Can we be fast but also power/
energy efficient?**

Example of a hardware accelerator

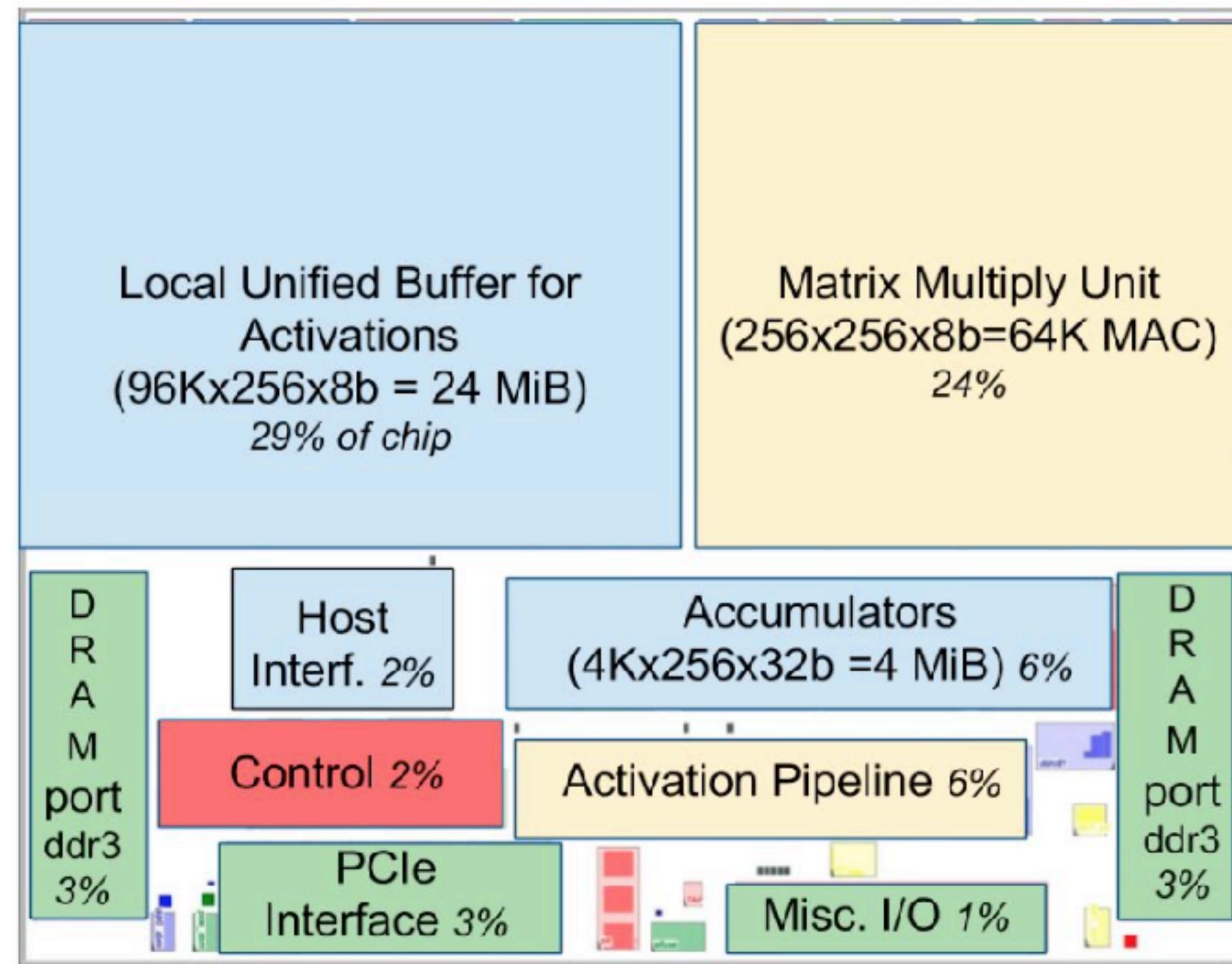
— Tensor Processing Unit

- N. Jouppi, C. Young, N. Patil and D. Patterson. Motivation for and Evaluation of the First Tensor Processing Unit. In IEEE Micro, vol. 38, no. 3, pp. 10-19. 2018
- Norman P. Jouppi, Cliff Young, Nishant Patil, David Patterson, Gaurav Agrawal, Raminder Bajwa, Sarah Bates, Suresh Bhatia, Nan Boden, Al Borchers, Rick Boyle, Pierre-luc Cantin, Clifford Chao, Chris Clark, Jeremy Coriell, Mike Daley, Matt Dau, Jeffrey Dean, Ben Gelb, Tara Vazir Ghaemmaghami, Rajendra Gottipati, William Gulland, Robert Hagmann, C. Richard Ho, Doug Hogberg, John Hu, Robert Hundt, Dan Hurt, Julian Ibarz, Aaron Jaffey, Alek Jaworski, Alexander Kaplan, Harshit Khaitan, Daniel Killebrew, Andy Koch, Naveen Kumar, Steve Lacy, James Laudon, James Law, Diemthu Le, Chris Leary, Zhuyuan Liu, Kyle Lucke, Alan Lundin, Gordon MacKean, Adriana Maggiore, Maire Mahony, Kieran Miller, Rahul Nagarajan, Ravi Narayanaswami, Ray Ni, Kathy Nix, Thomas Norrie, Mark Omernick, Narayana Penukonda, Andy Phelps, Jonathan Ross, Matt Ross, Amir Salek, Emad Samadiani, Chris Severn, Gregory Sizikov, Matthew Snelham, Jed Souter, Dan Steinberg, Andy Swing, Mercedes Tan, Gregory Thorson, Bo Tian, Horia Toma, Erick Tuttle, Vijay Vasudevan, Richard Walter, Walter Wang, Eric Wilcox, and Doe Hyun Yoon.
[In-Datacenter Performance Analysis of a Tensor Processing Unit](#). In ISCA '17. 2017.

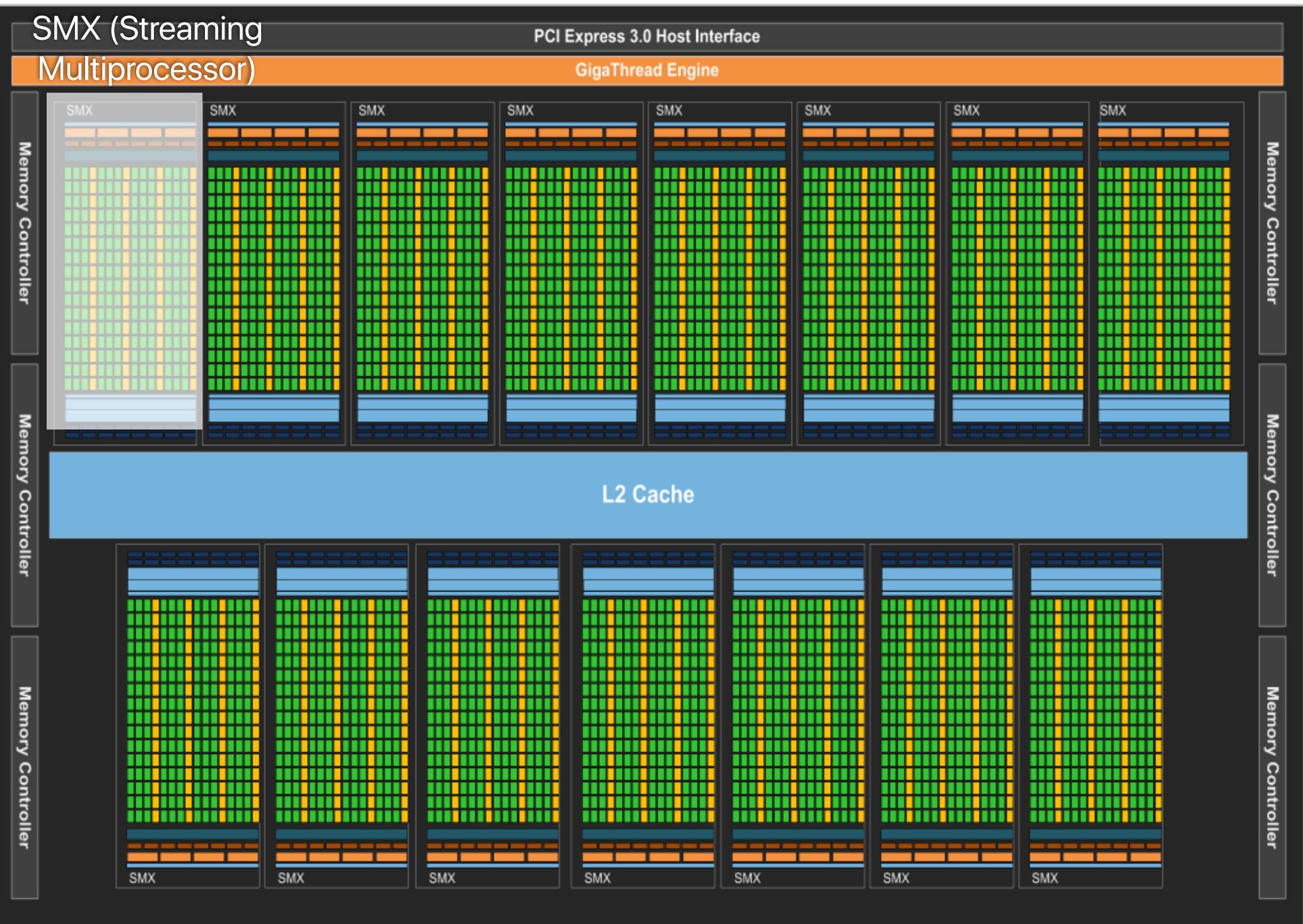
What the “first” generation of TPU looks like



TPU Floorplan

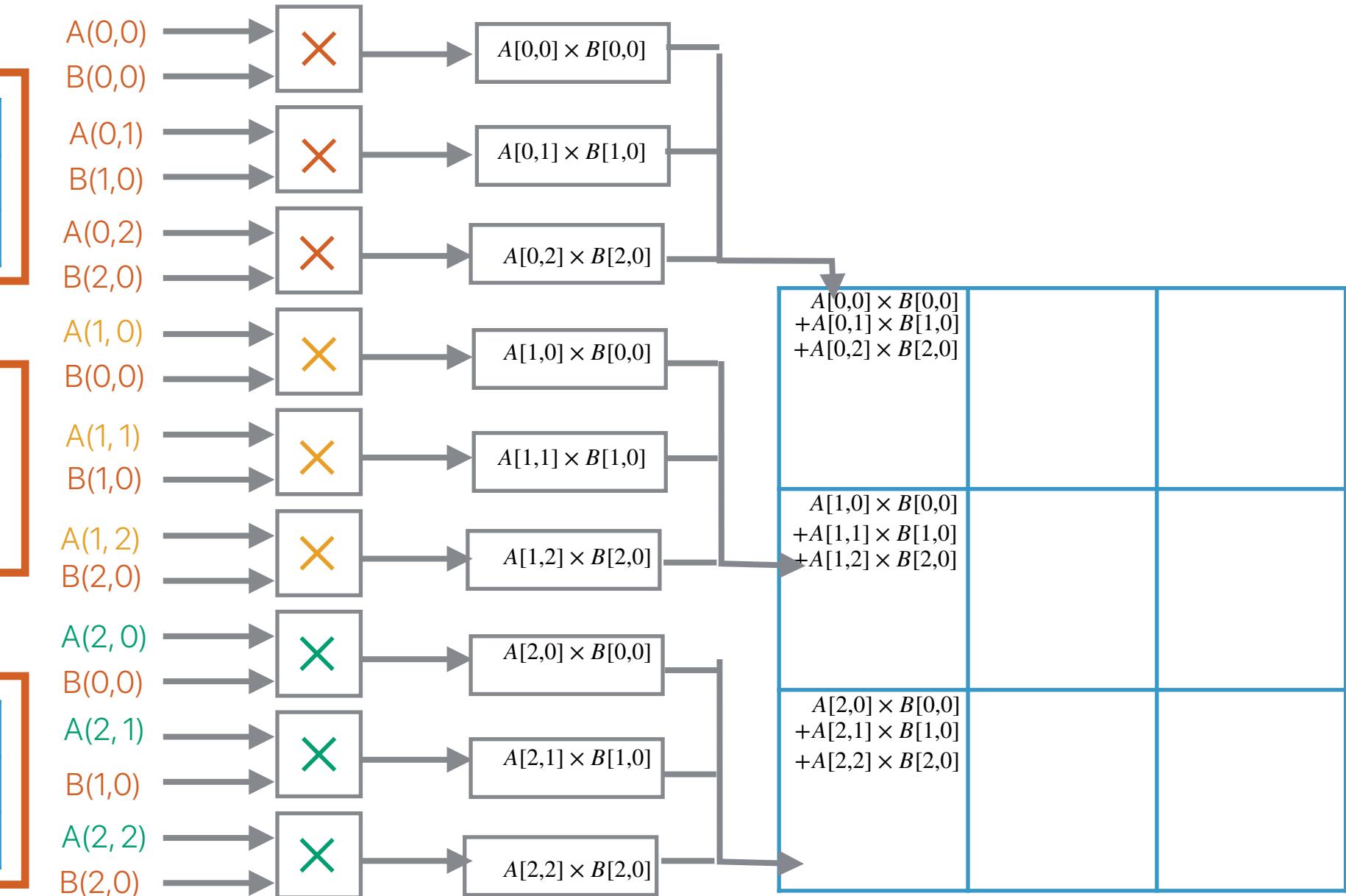
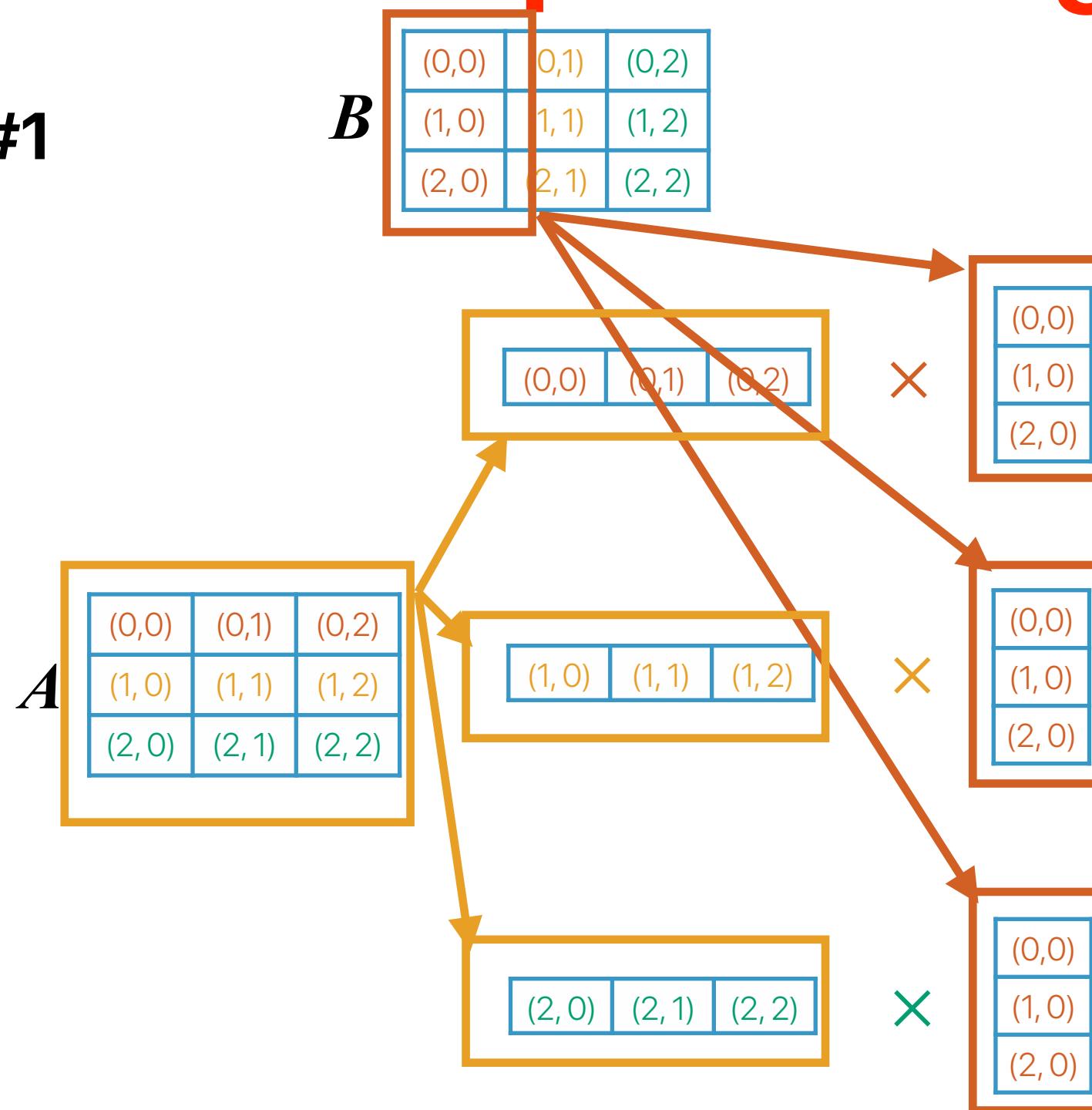


Vector processing model by GPUs



Vector processing for MM using N^2 PEs

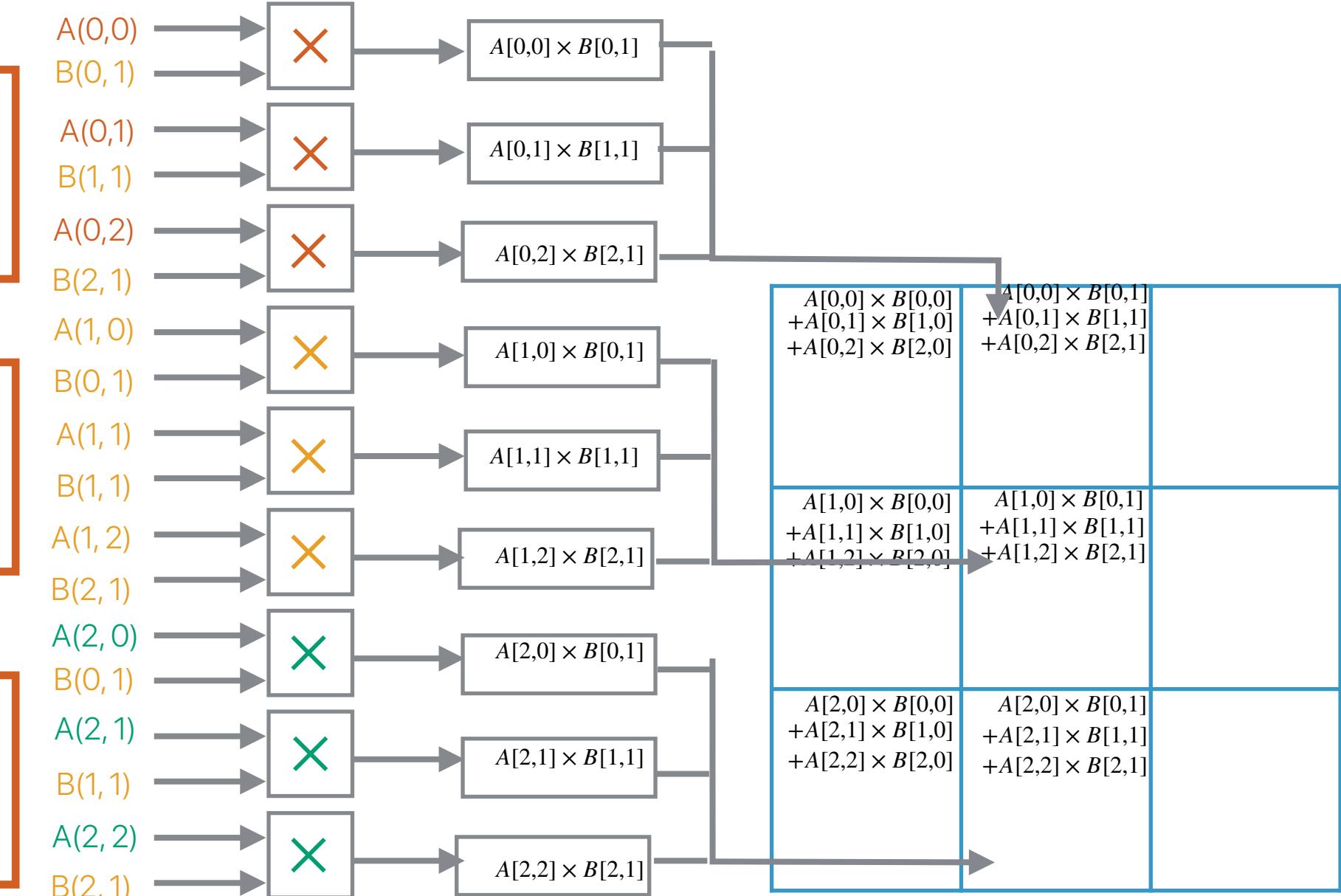
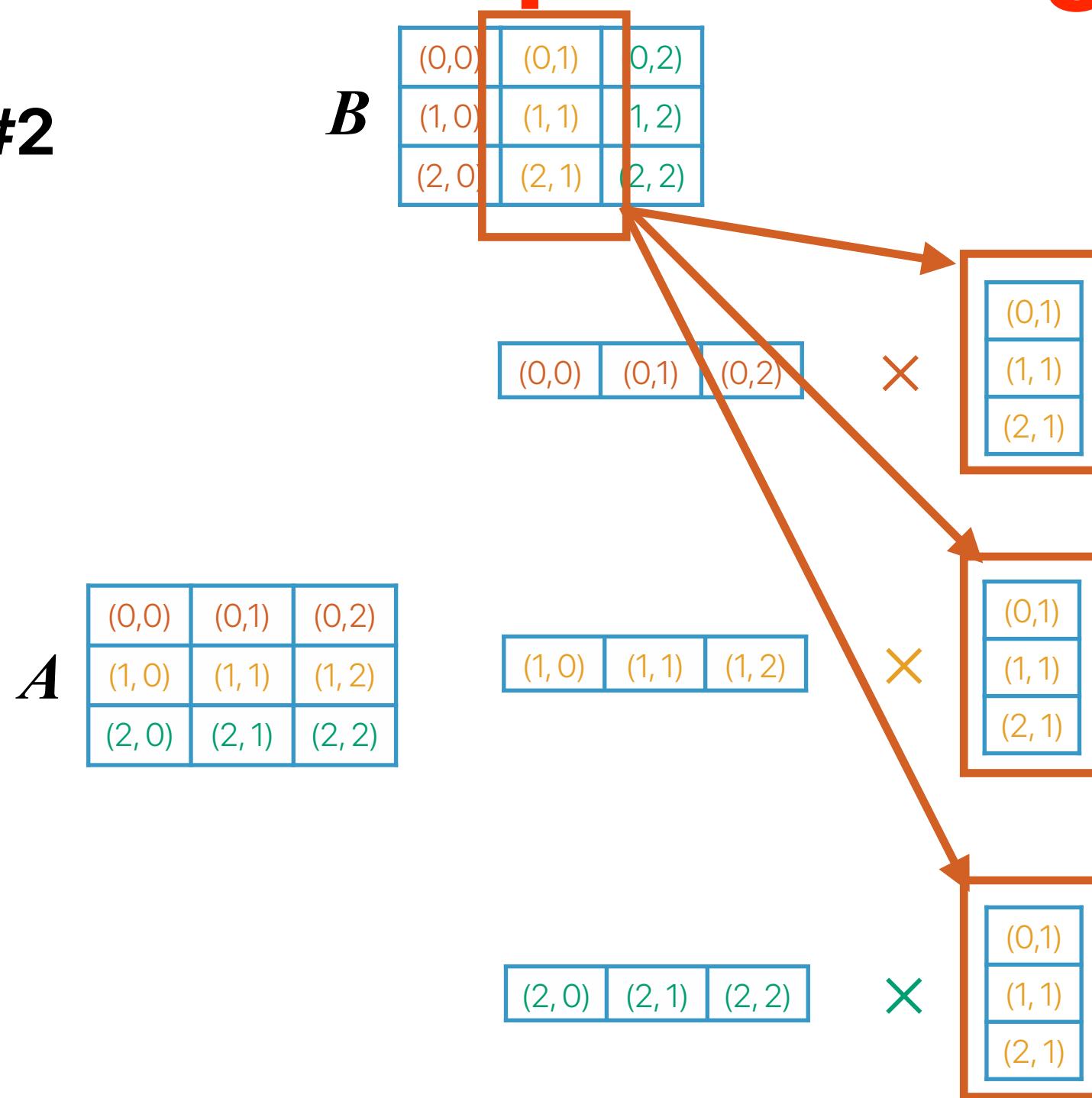
#1



Must feed $2N^2$ elements ($8N^2$ bytes) before
each vector MAC operation

Vector processing for MM using N^2 PEs

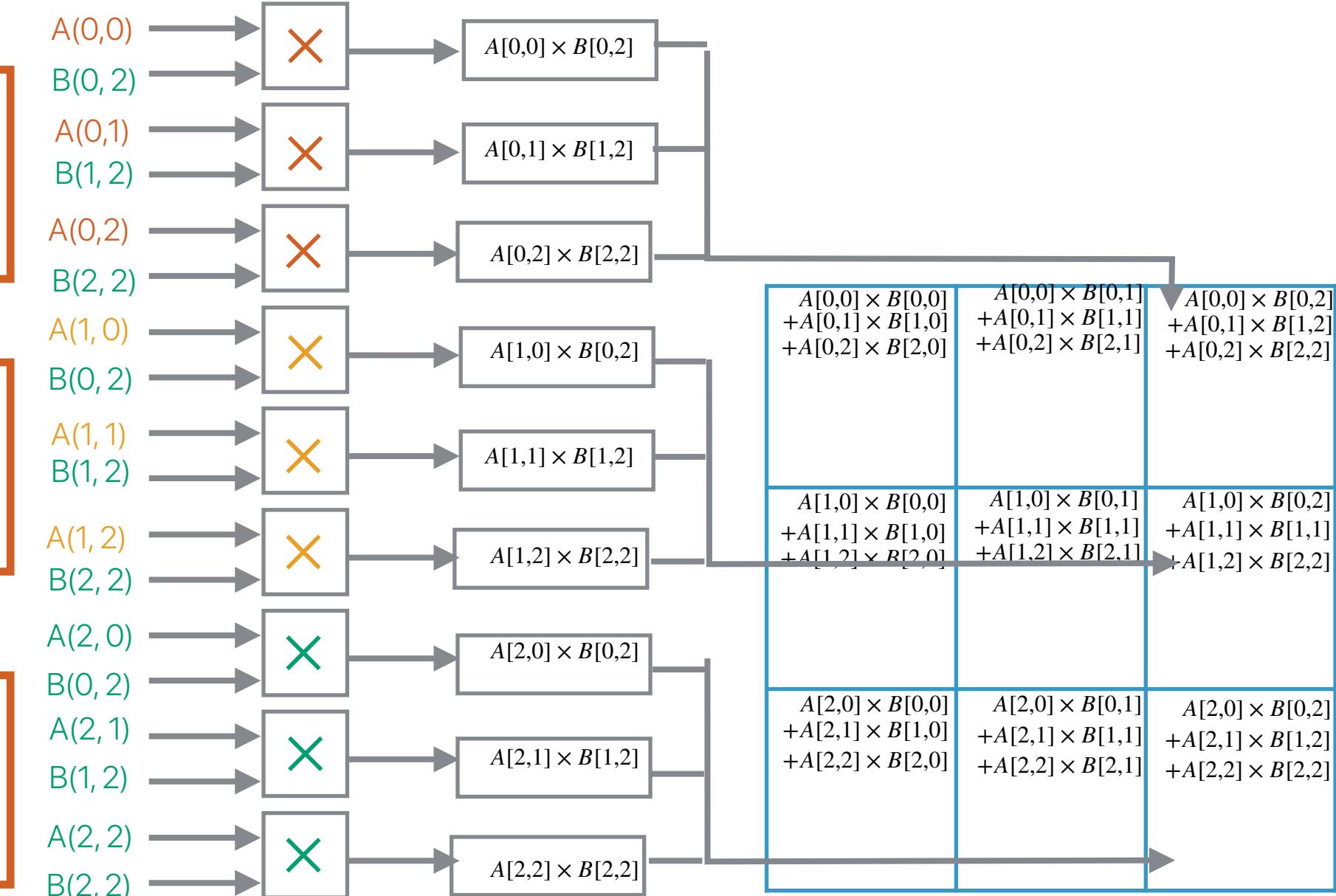
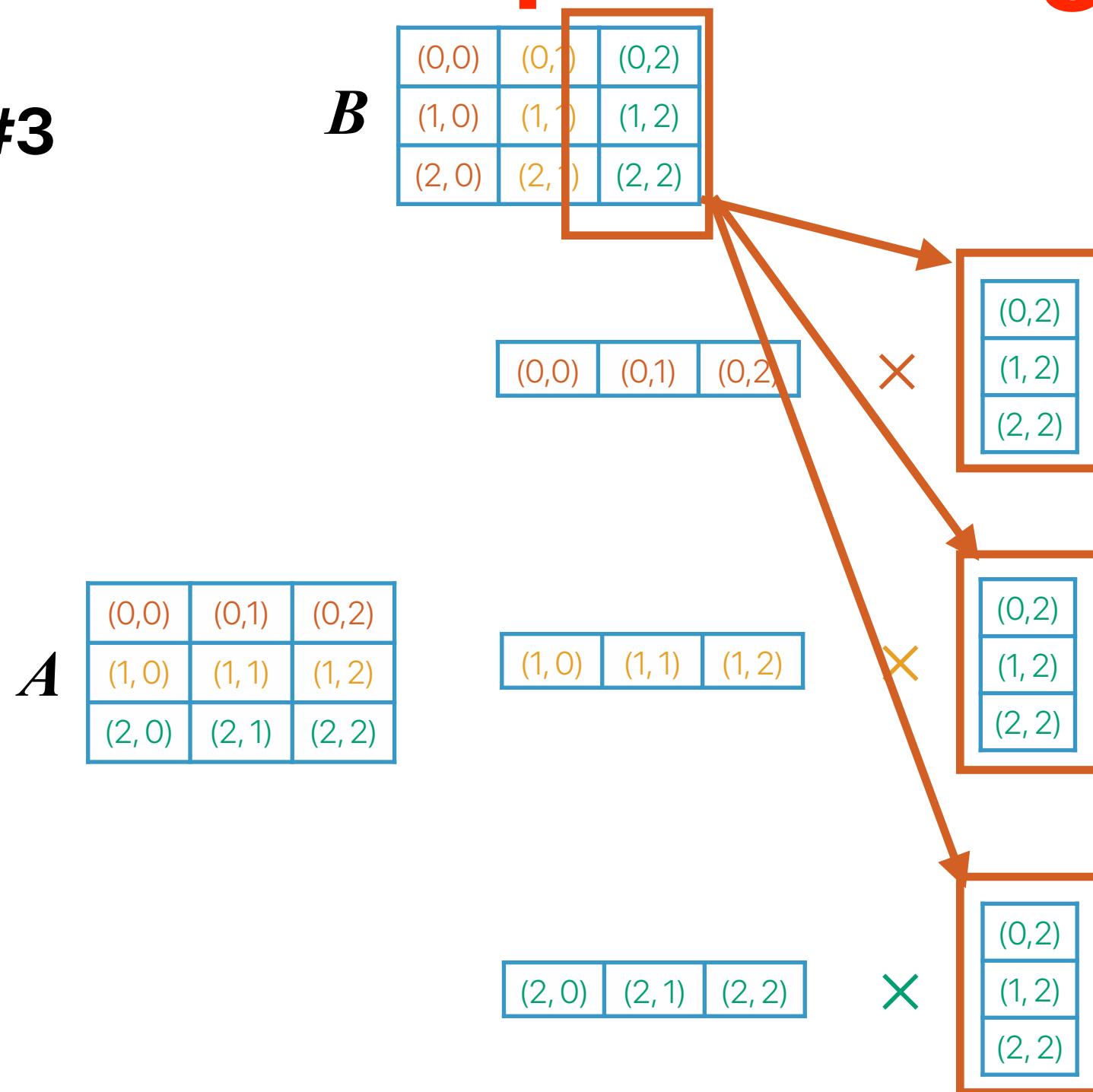
#2



Must feed N^2 elements ($4N^2$ bytes) before each vector MAC operation

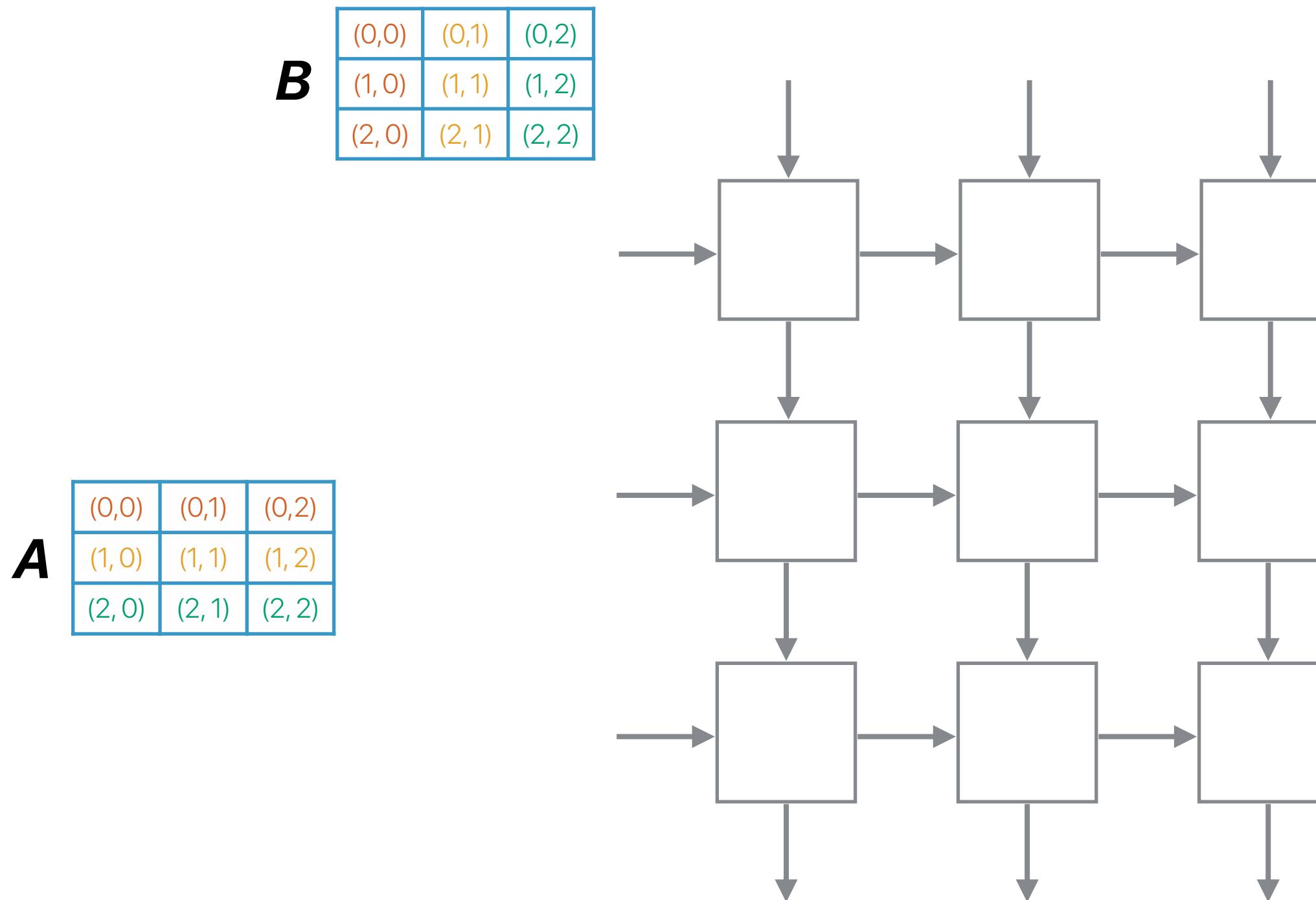
Vector processing for MM using N^2 PEs

#3



Must feed N^2 elements (4 N^2 bytes) before each vector MAC operation

Systolic Arrays used by AI/ML Accelerators



H.T. Kung, "Why systolic architectures?", in IEEE Computer, vol. 15, no. 1, pp. 37-46, Jan. 1982, doi: 10.1109/MC.1982.1653825.

Systolic Arrays used by AI/ML Accelerators

#1

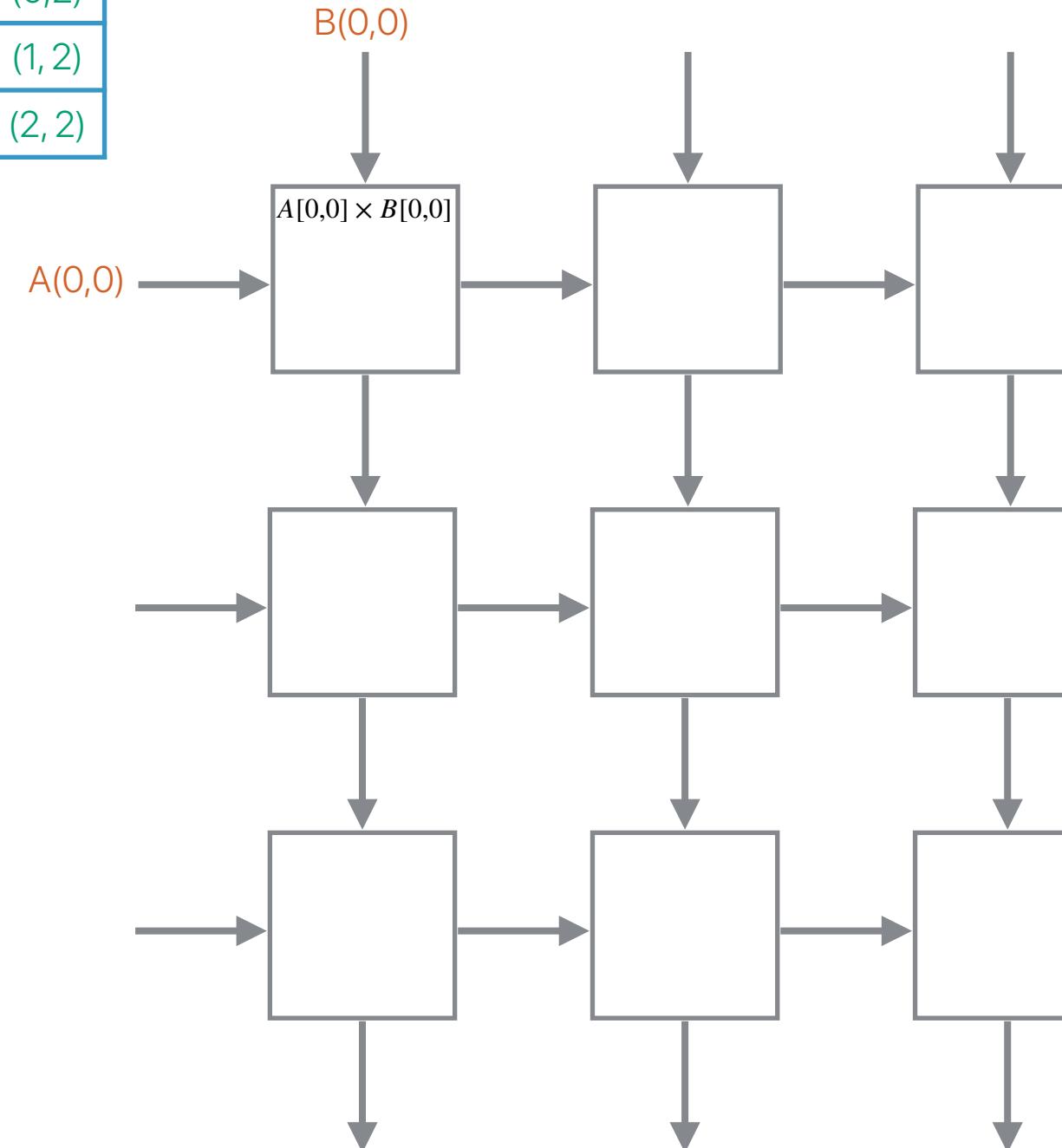
B

(0,0)	(0,1)	(0,2)
(1,0)	(1,1)	(1,2)
(2,0)	(2,1)	(2,2)

Only need to fetch 2 elements (8 bytes) before we start

A

(0,0)	(0,1)	(0,2)
(1,0)	(1,1)	(1,2)
(2,0)	(2,1)	(2,2)



H.T. Kung, "Why systolic architectures?", in IEEE Computer, vol. 15, no. 1, pp. 37-46, Jan. 1982, doi: 10.1109/MC.1982.1653825.

Systolic Arrays used by AI/ML Accelerators

#2

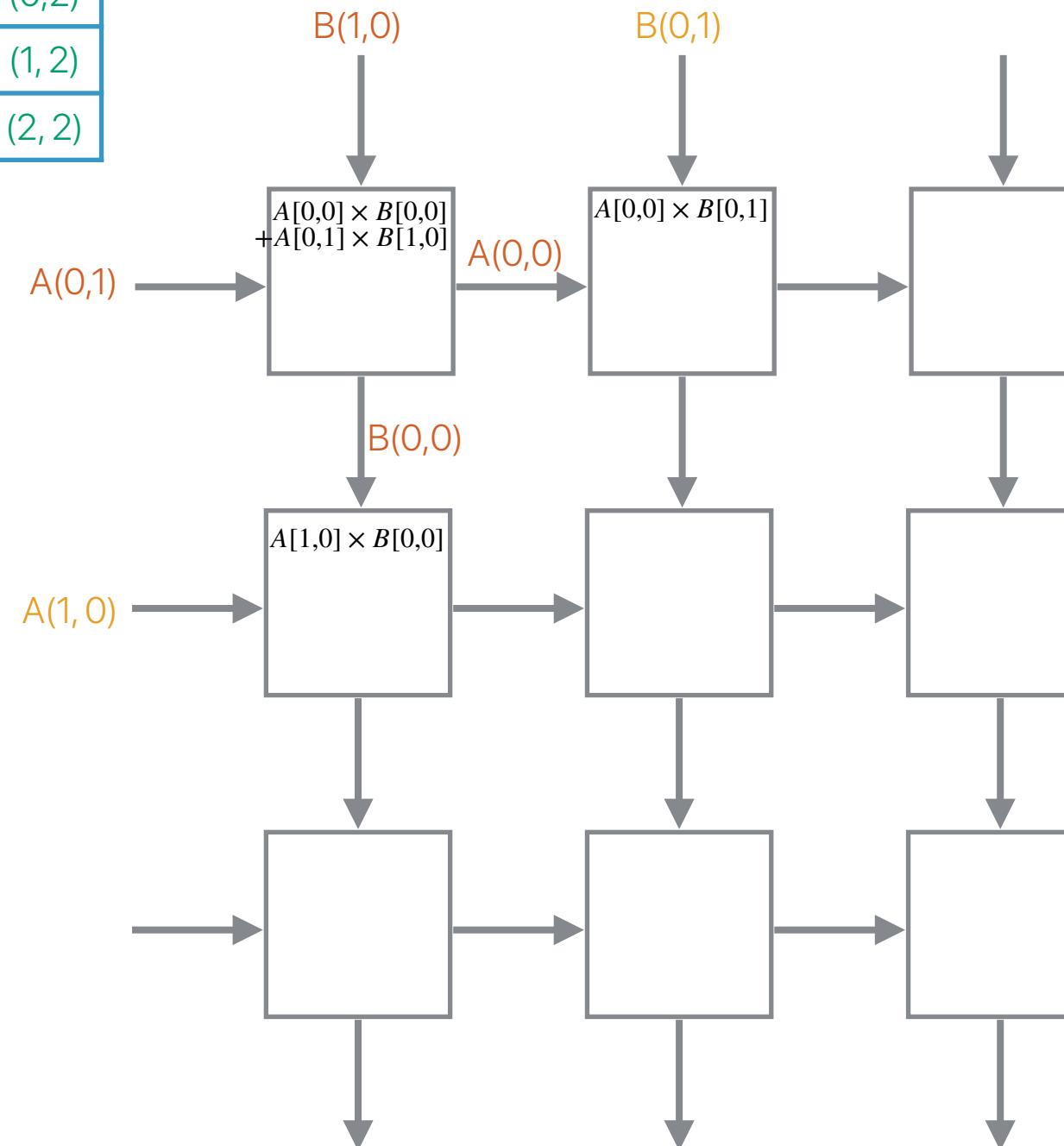
B

(0,0)	(0,1)	(0,2)
(1,0)	(1,1)	(1,2)
(2,0)	(2,1)	(2,2)

Only need to fetch 4 elements (16 bytes) before we start

A

(0,0)	(0,1)	(0,2)
(1,0)	(1,1)	(1,2)
(2,0)	(2,1)	(2,2)



H.T. Kung, "Why systolic architectures?", in IEEE Computer, vol. 15, no. 1, pp. 37-46, Jan. 1982, doi: 10.1109/MC.1982.1653825.

Systolic Arrays used by AI/ML Accelerators

#3

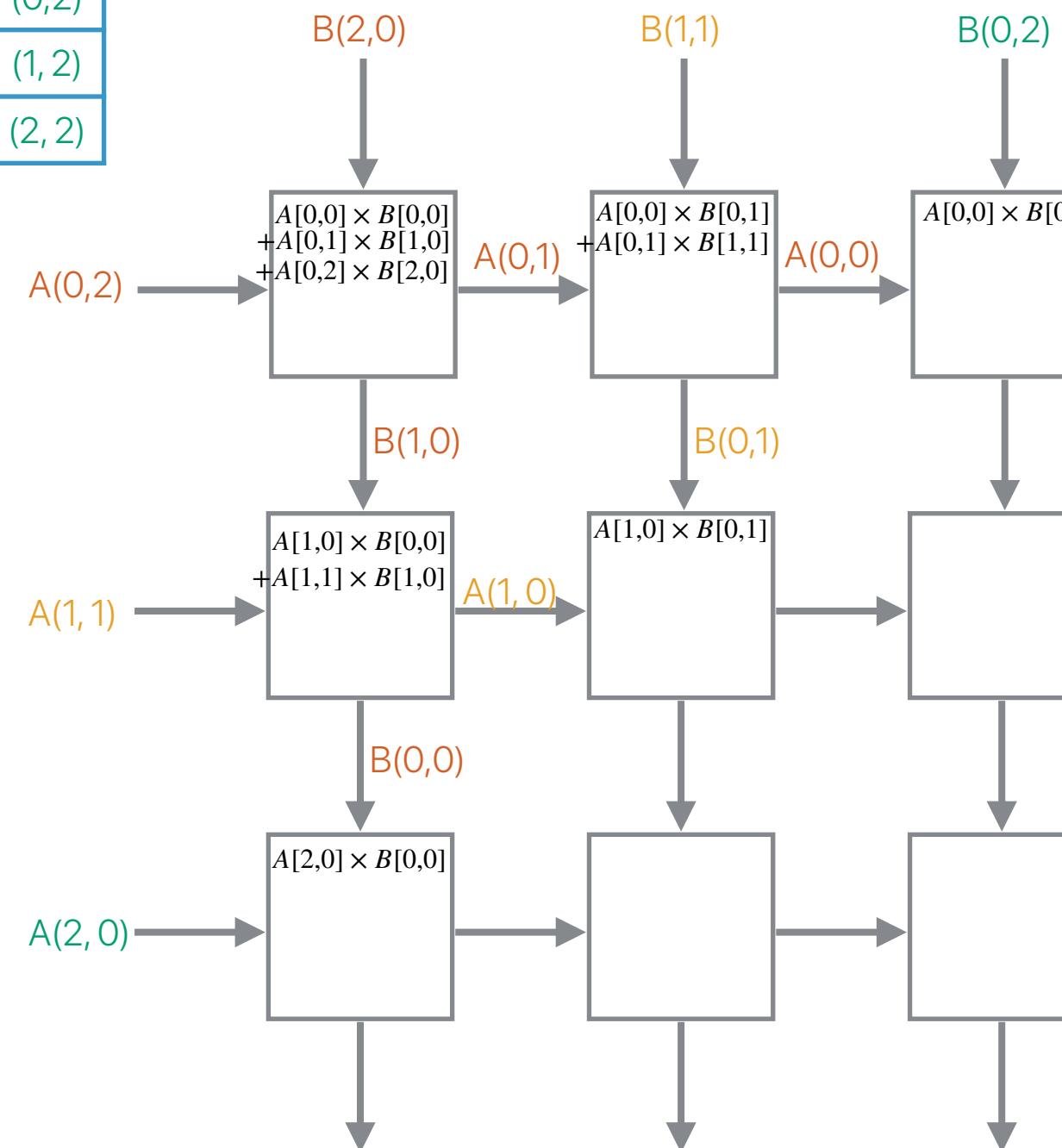
B

(0,0)	(0,1)	(0,2)
(1,0)	(1,1)	(1,2)
(2,0)	(2,1)	(2,2)

Only need to fetch $2N$ elements ($8N$ bytes) before we start

A

(0,0)	(0,1)	(0,2)
(1,0)	(1,1)	(1,2)
(2,0)	(2,1)	(2,2)



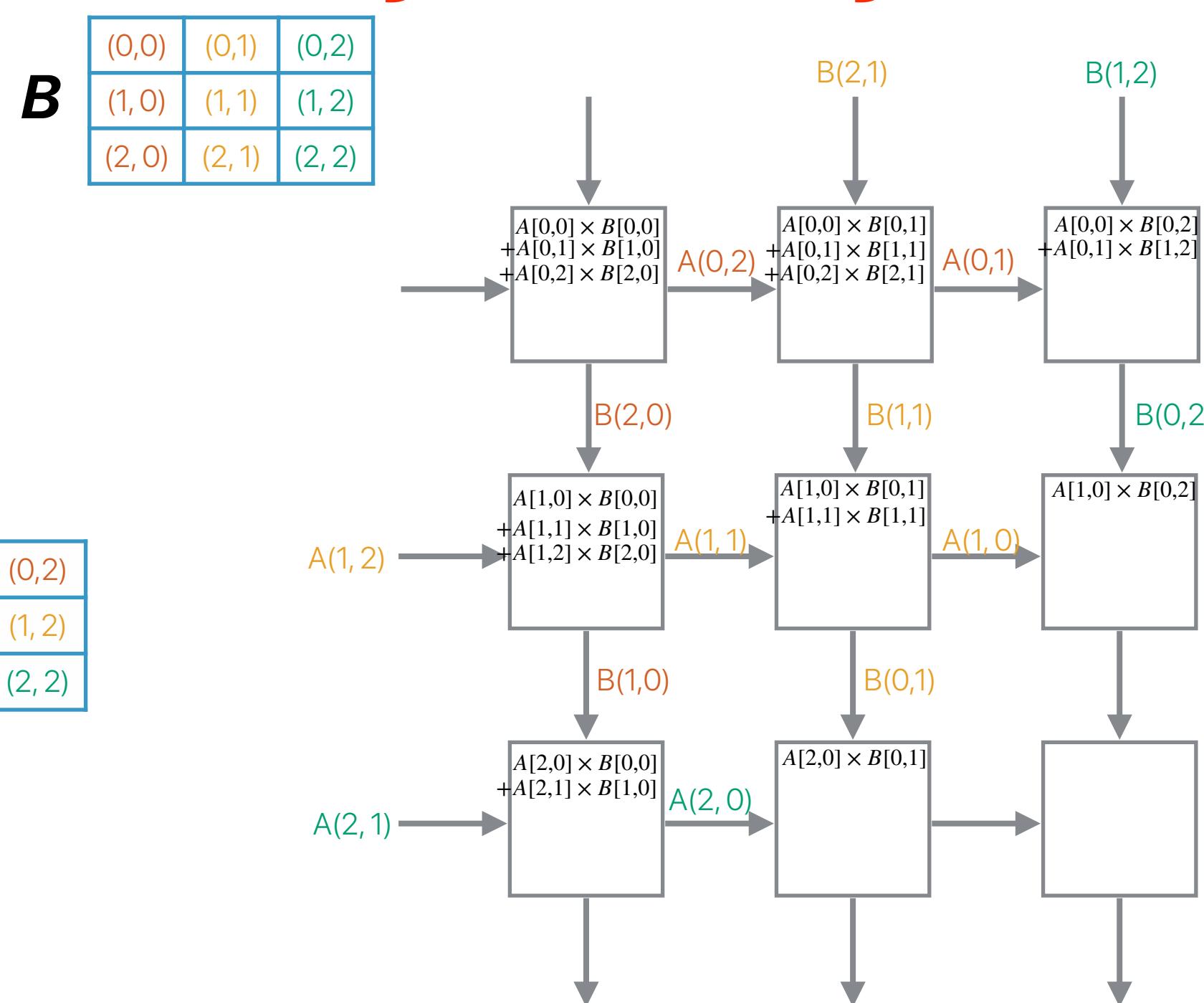
H.T. Kung, "Why systolic architectures?", in IEEE Computer, vol. 15, no. 1, pp. 37-46, Jan. 1982, doi: 10.1109/MC.1982.1653825.

Systolic Arrays used by AI/ML Accelerators

#4

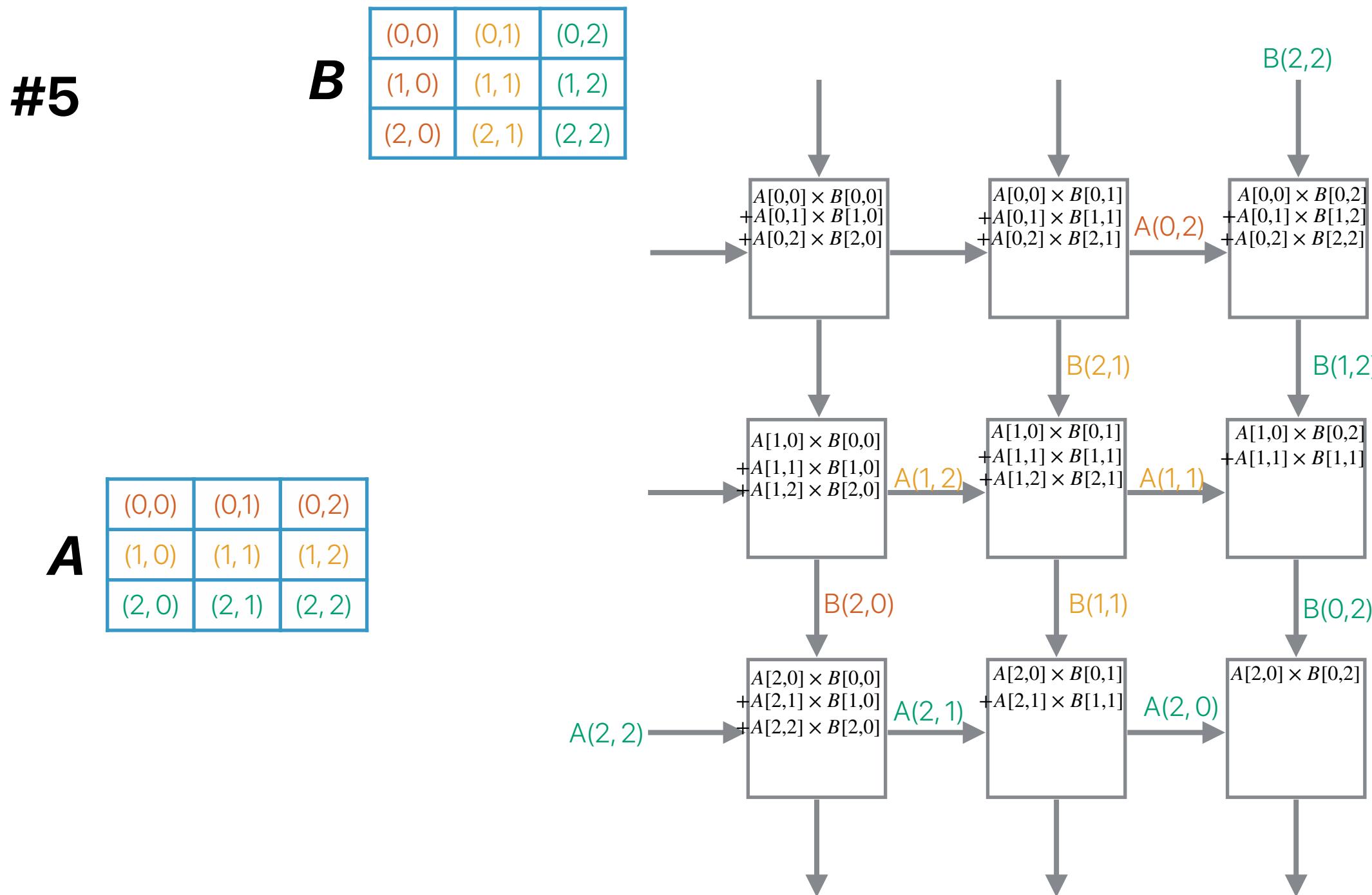
A

(0,0)	(0,1)	(0,2)
(1,0)	(1,1)	(1,2)
(2,0)	(2,1)	(2,2)



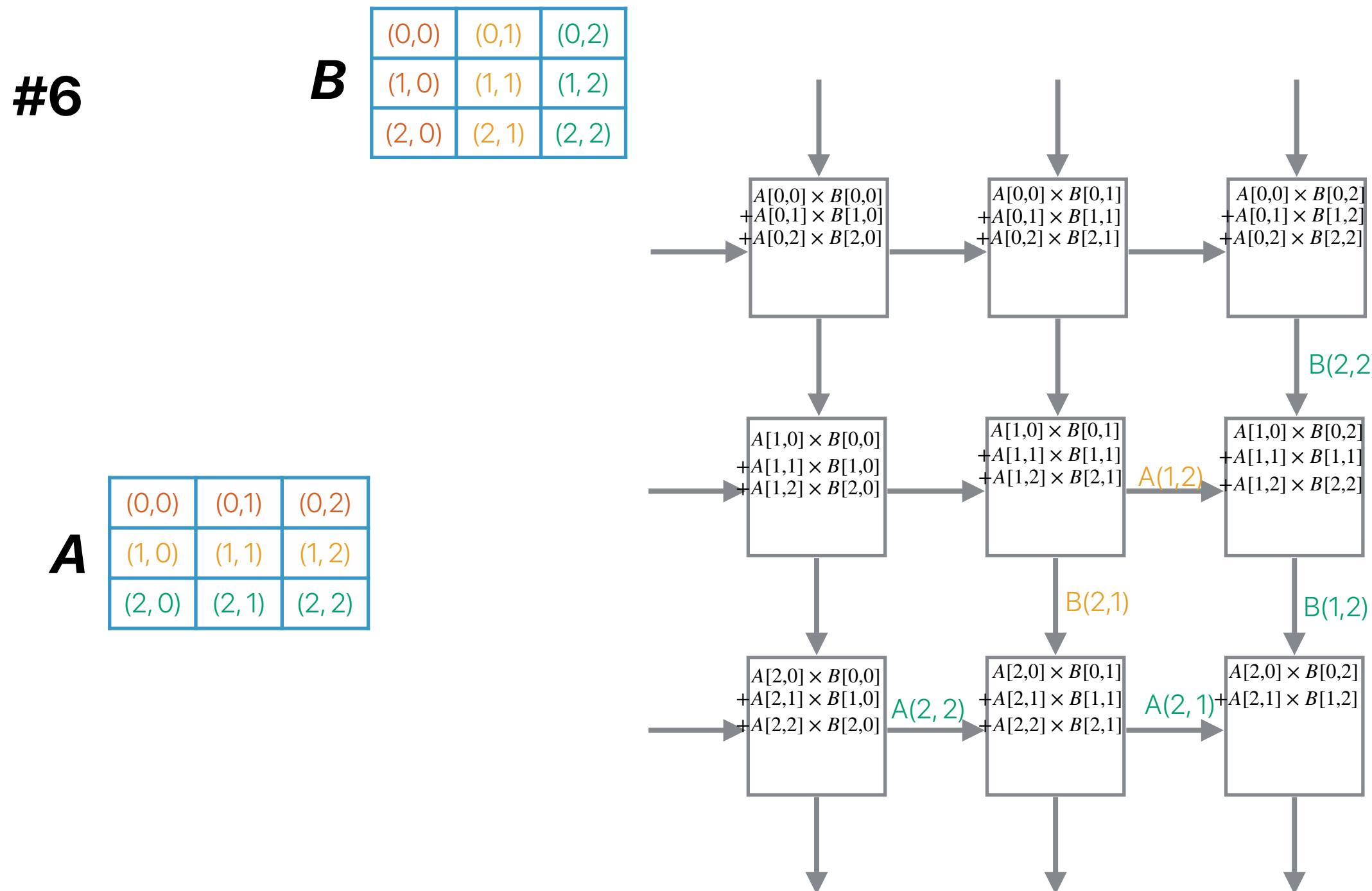
H.T. Kung, "Why systolic architectures?", in IEEE Computer, vol. 15, no. 1, pp. 37-46, Jan. 1982, doi: 10.1109/MC.1982.1653825.

Systolic Arrays used by AI/ML Accelerators



H.T. Kung, "Why systolic architectures?", in IEEE Computer, vol. 15, no. 1, pp. 37-46, Jan. 1982, doi: 10.1109/MC.1982.1653825.

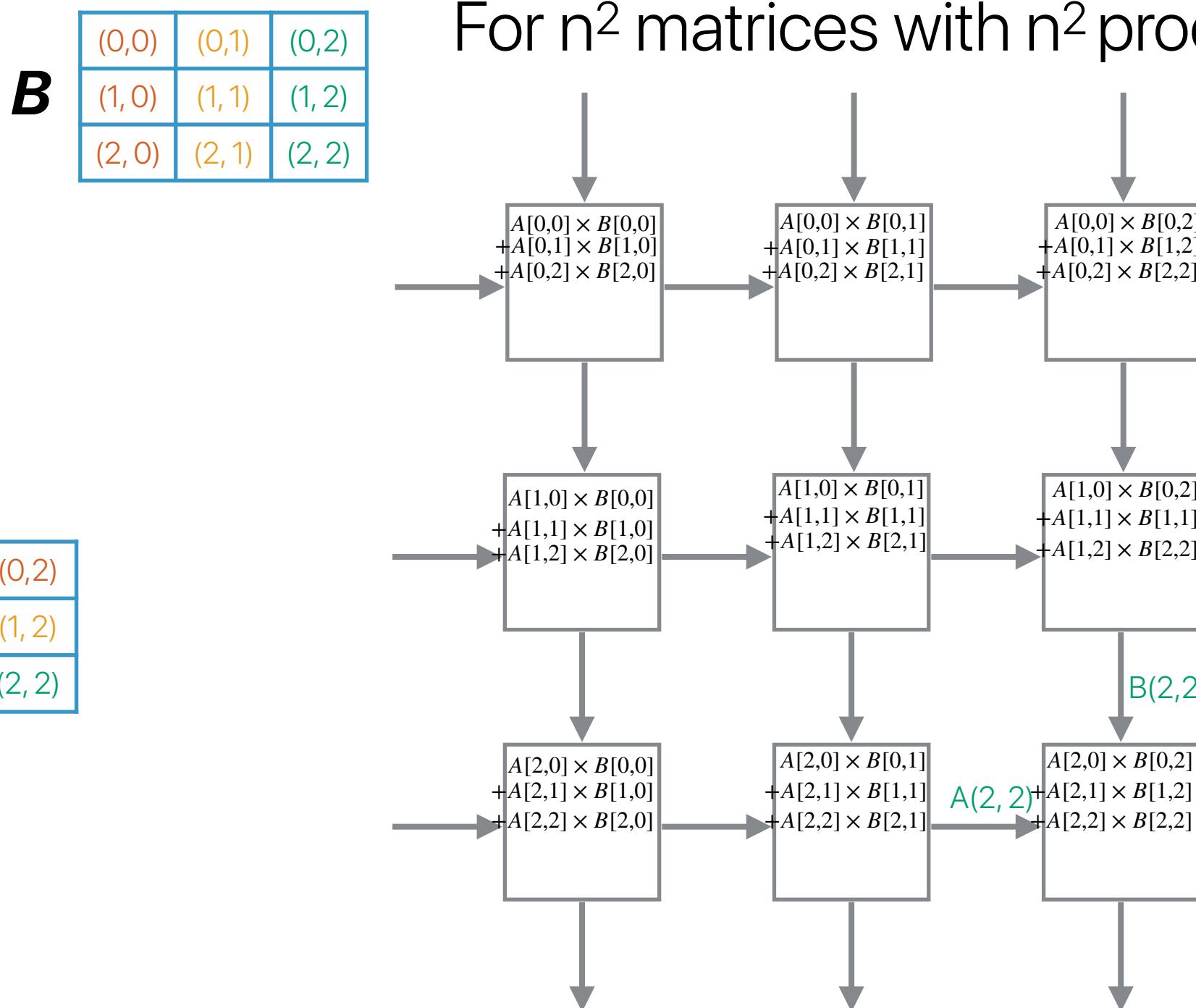
Systolic Arrays used by AI/ML Accelerators



H.T. Kung, "Why systolic architectures?", in IEEE Computer, vol. 15, no. 1, pp. 37-46, Jan. 1982, doi: 10.1109/MC.1982.1653825.

Systolic Arrays used by AI/ML Accelerators

#7



H.T. Kung, "Why systolic architectures?", in IEEE Computer, vol. 15, no. 1, pp. 37-46, Jan. 1982, doi: 10.1109/MC.1982.1653825.

Maxwell (2015)'s memory bandwidth

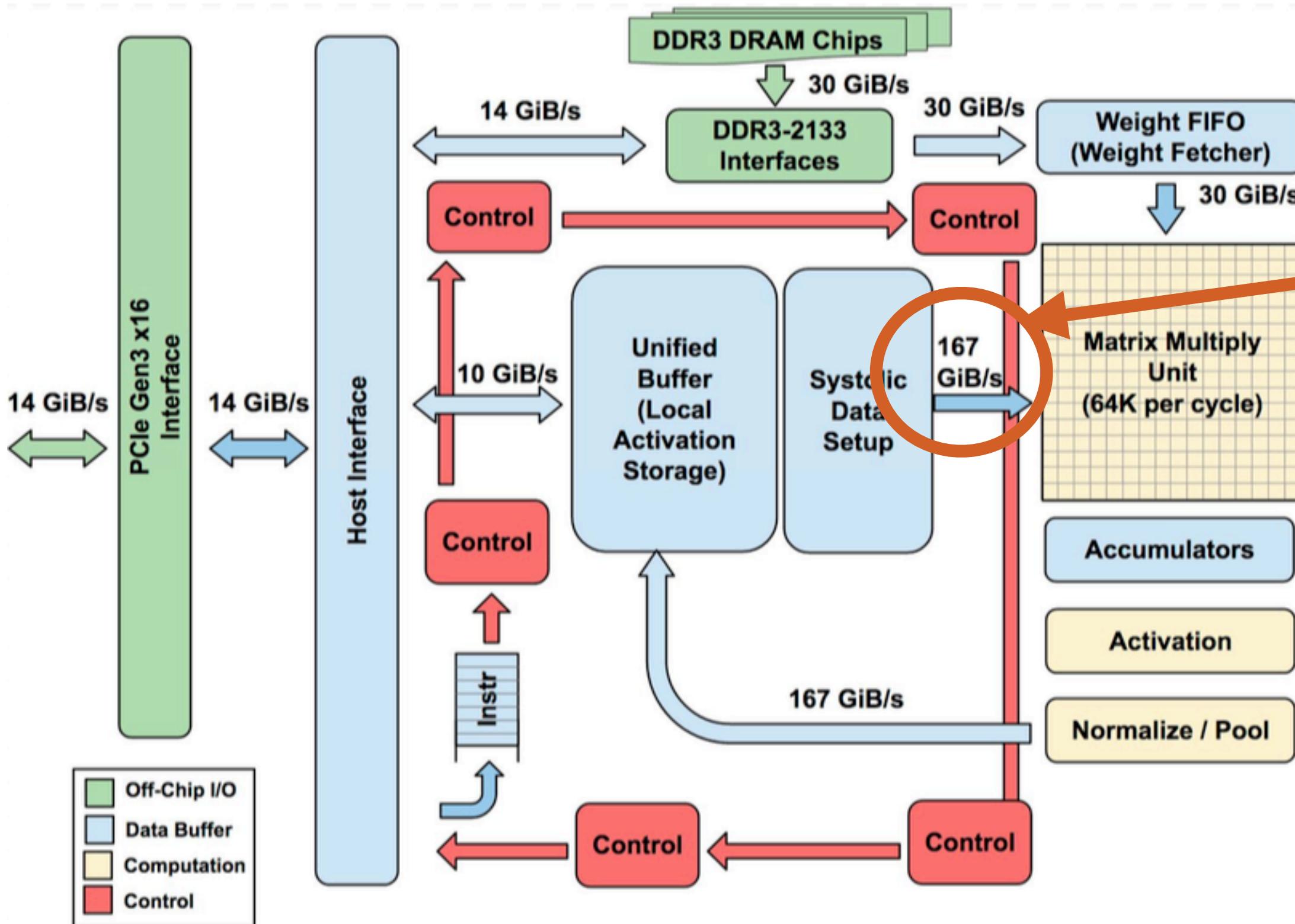
(Kepler GK104) and the New GeForce GTX 980 (Maxwell GM204).

GPU	GeForce GTX 680 (Kepler GK104)	GeForce GTX 980 (Maxwell GM204)
CUDA Cores	1536	2048
Base Clock	1006 MHz	1126 MHz
GPU Boost Clock	1058 MHz	1216 MHz
GFLOPs	3090	4612
Compute Capability	3.0	5.2
SMs	8	16
Shared Memory / SM	48KB	96KB
Register File Size / SM	256KB	256KB
Active Blocks / SM	16	32
Texture Units	128	128
Texel fill-rate	128.8 Gigatexels/s	144.1 Gigatexels/s
Memory	2048MB	4096MB
Memory Clock	6008 MHz	7010 MHz
Memory Bandwidth	192.3 GB/sec	224.3 GB/sec

$$2048 \times 1126 \times 10^6 = 2306.05 \times 10^9$$

$$\frac{2048 \times 1126 \times 10^6 \times 4 \times 2}{2^{30}} = 17181 \text{GB/sec}$$

TPU Block diagram

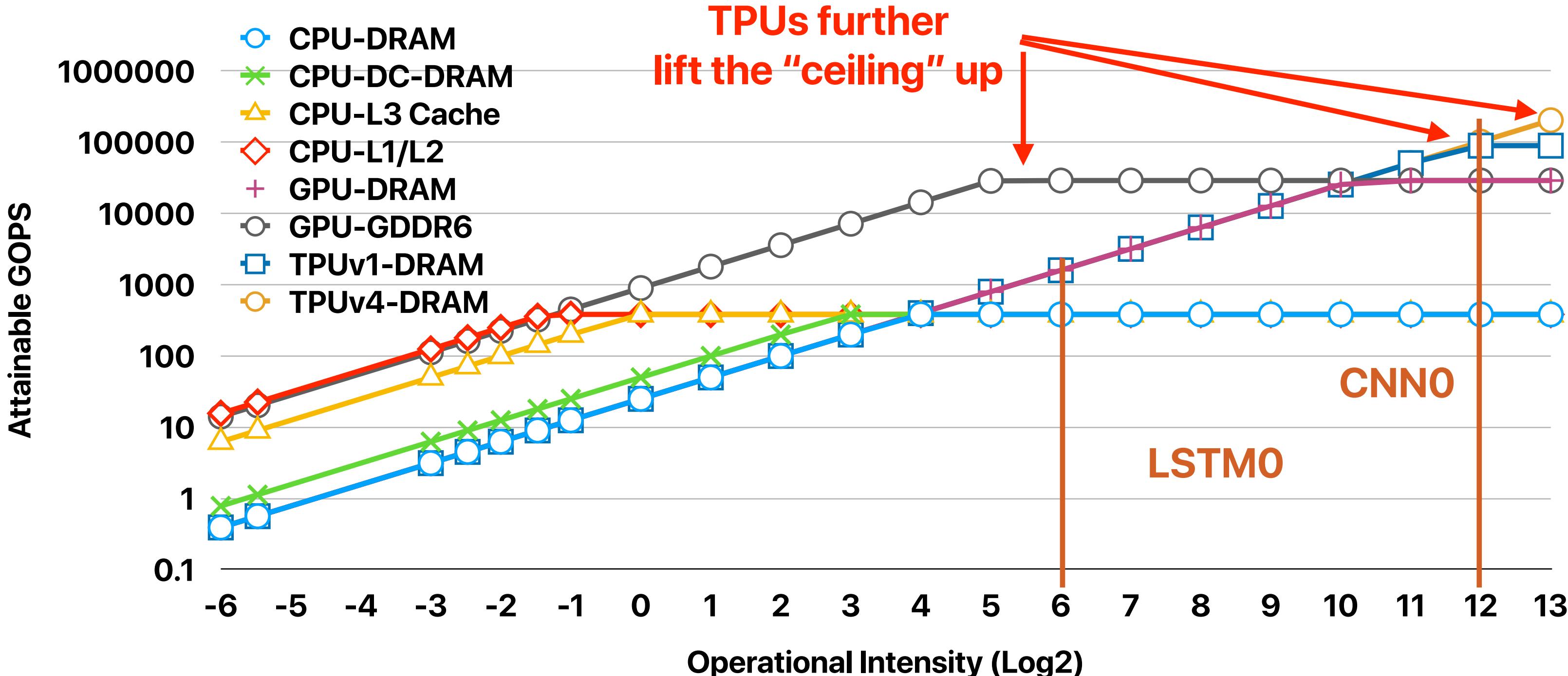


Smaller than the demand of GPUs

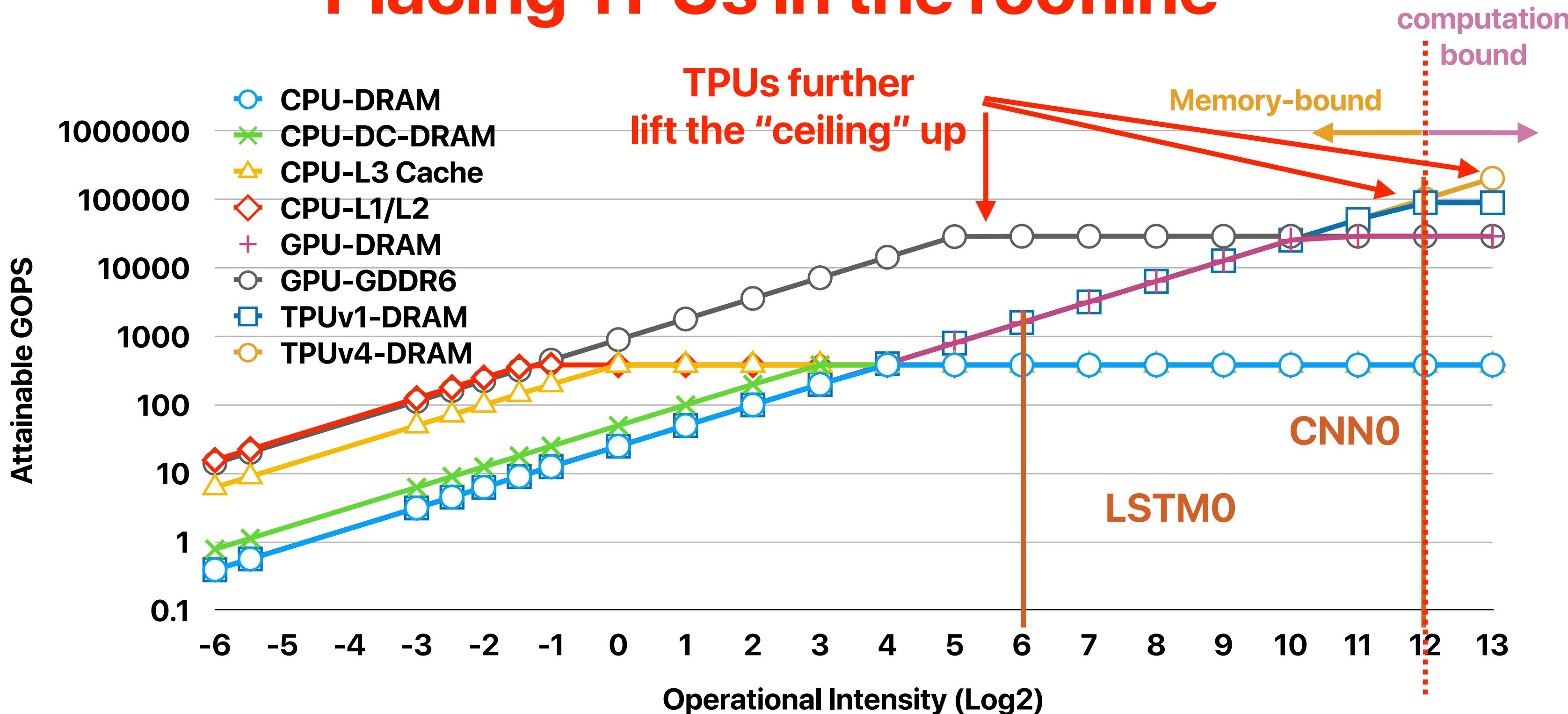
Why can accelerator help?

- Performance-wise $ET = IC \times CPI \times CT$
 - The accelerator itself reduces the amount of “instructions” used to implement a program
 - The accelerator’s logic is simpler than that of a general-purpose pipeline (potentially higher clock rate)
- Power-wise $P = \alpha CV^2f$
 - Simpler logic reduces the waste of gates
- Energy-wise $Energy = P \times ET$
 - Energy is power times execution time
 - If you improve both, it’s better anyway

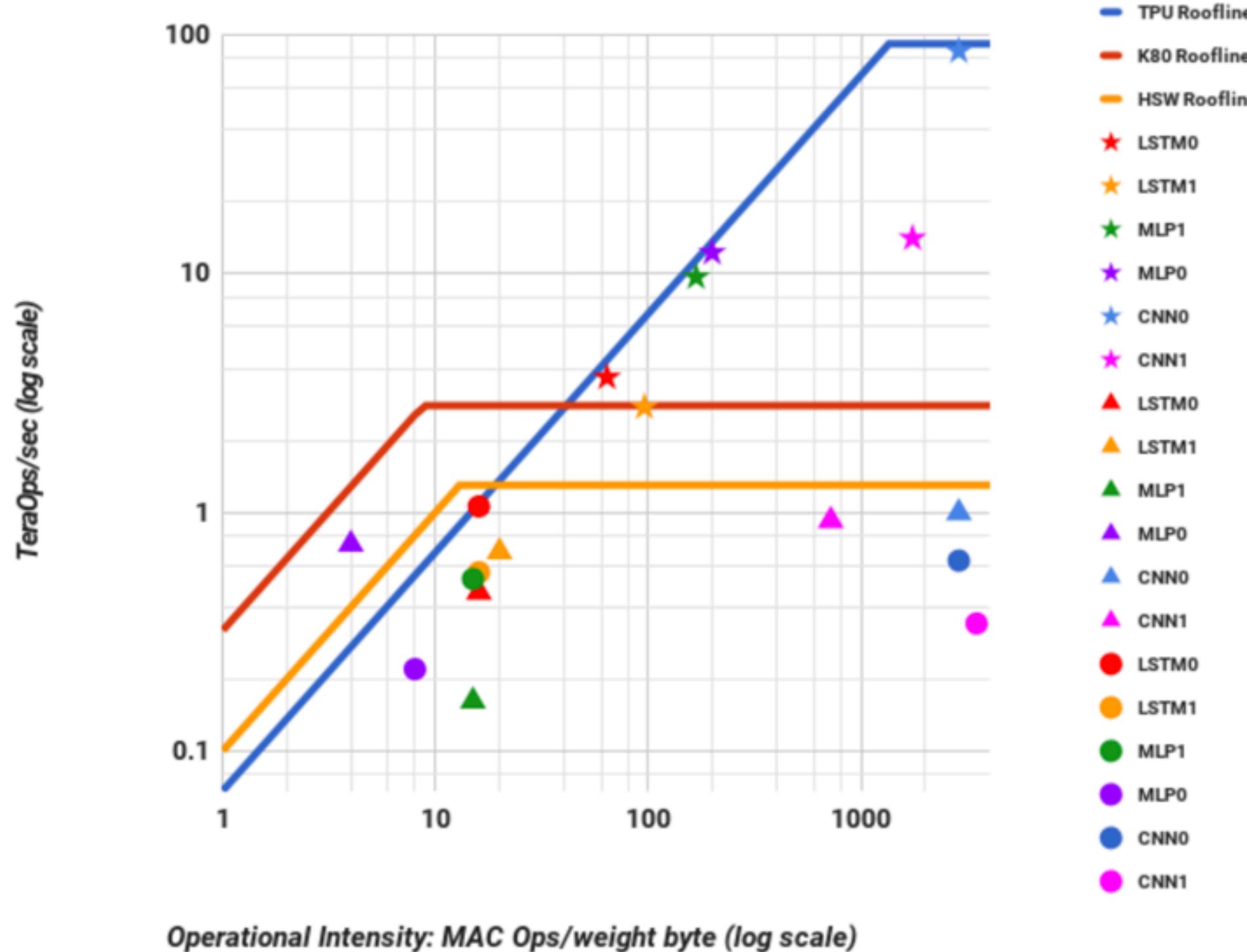
Placing TPUs in the roofline



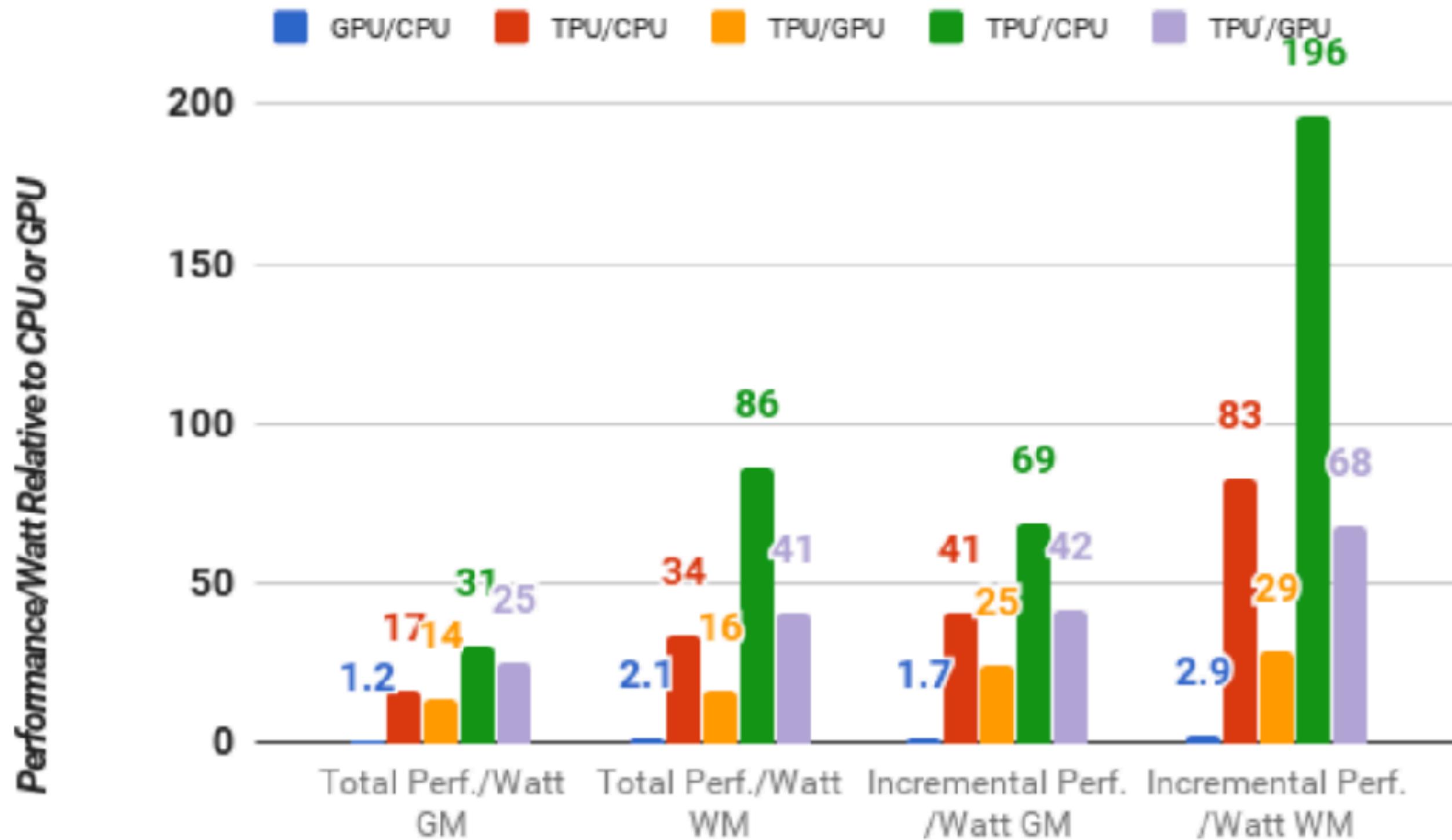
Placing TPUs in the roofline



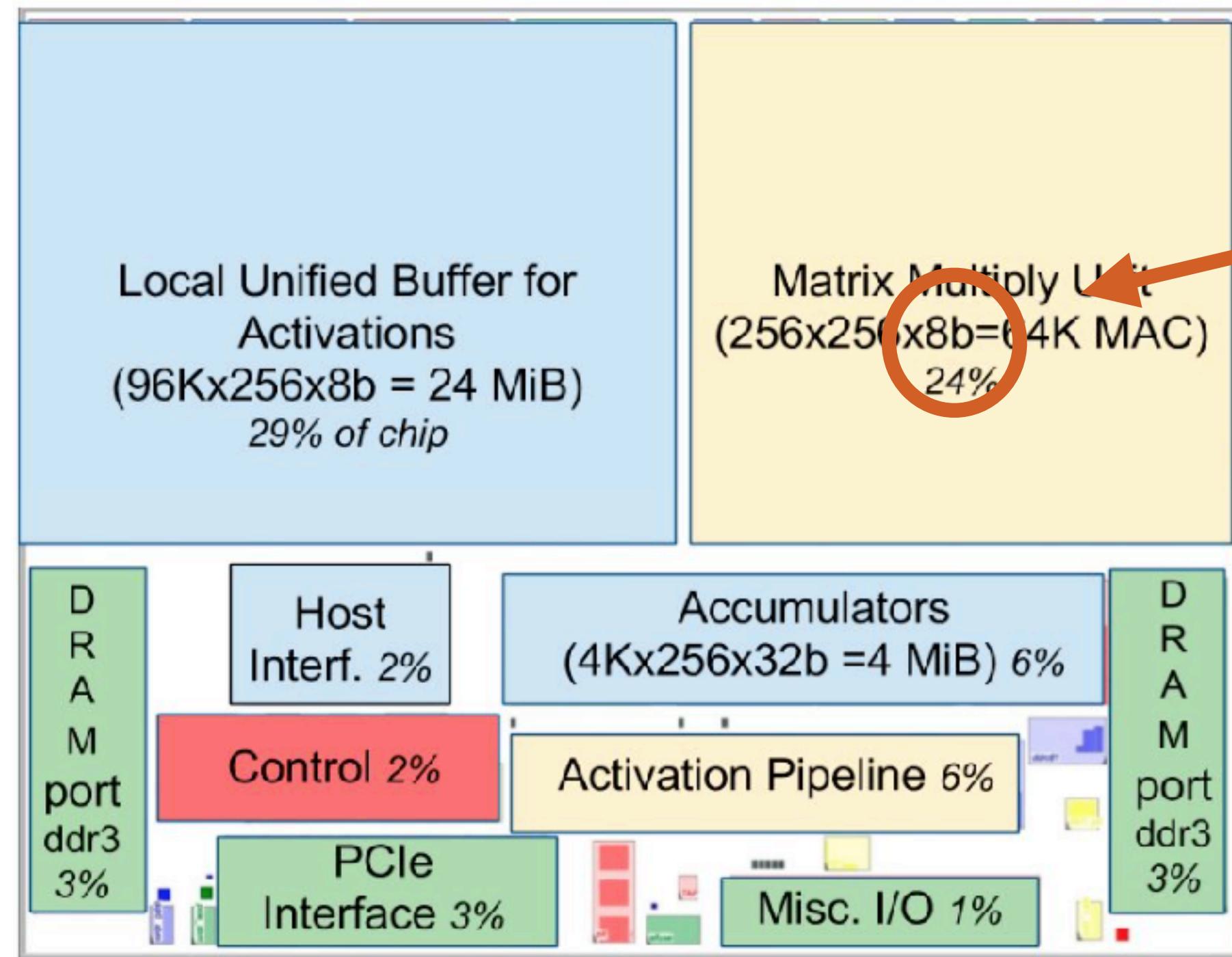
Most workloads surpass the roofline of GPUs



Energy efficiency is better! — reduced cost!



TPU Floorplan



Why 8-bit?

The benefit of application-specific design

- We only need to provide “just enough” hardware support rather than complete support
 - We can deliver higher throughput within the same area
 - Support 2x bits in precision 4x the area
 - We can have more “values” sent using the same memory bandwidth
 - We can reduce the total power/energy consumption

Table 2. Benchmarked servers use Haswell CPUs, K80 GPUs, and TPUs. The GPU and TPU use the Haswell server as host. The TPU die is less than half the Haswell die size.

Model	mm ²	nm	MHz	Thermal Design Power (TDP) per Chip	Cores	TeraOps/s		Memory Gbyte/s	Chips/Server	TDP per Server
						8 bit	FP32			
Haswell	662	22	2,300	145 W	18	2.6	1.3	51	2	504 W
Nvidia K80	561	28	560	150 W	13	--	2.8	160	8	1,838 W
TPU	< 331	28	700	75 W	1	92	--	34	4	861 W

From the TPU paper

Pitfall: Being Ignorant of Architecture History when Designing a DSA.

Ideas that didn't fly for general-purpose computing might be ideal for DSAs. For the TPU, three important architectural ideas date from the early 1980s: systolic arrays,⁵ decoupled-access/execute,⁴ and CISC instructions. The first reduced the area and power of the large matrix multiply unit, the second fetches weights concurrently during operation of the matrix multiply unit, and the third better utilizes the limited bandwidth of the PCIe bus for delivering instructions. History-aware architects could have a competitive edge in this DSA era.

TPU Programming model — domain specific languages

Cloud TPU

Product overview

Introduction to Cloud TPU

Quickstarts

All quickstarts

Run TensorFlow on Cloud TPU VM

Run JAX on Cloud TPU VM

Run PyTorch on Cloud TPU VM

How-to guides

Concepts

Tutorials

Colab notebooks

★ Note: For single TPUs, Slices, and Pods, pass your TPU name to `TPUClusterResolver()`.

```
import tensorflow as tf
print("Tensorflow version " + tf.__version__)

tpu = tf.distribute.cluster_resolver.TPUClusterResolver(tpu='your(tpu-name') # TPU detection
print('Running on TPU ', tpu.cluster_spec().as_dict()['worker'])

tf.config.experimental_connect_to_cluster(tpu)
tf.tpu.experimental.initialize_tpu_system(tpu)
strategy = tf.distribute.experimental.TPUStrategy(tpu)

@tf.function
def add_fn(x,y):
    z = x + y
    return z

x = tf.constant(1.)
y = tf.constant(1.)
z = strategy.run(add_fn, args=(x,y))
print(z)
```

Roogle Project Meetings

- Make an appointment through the Google Calendar
- Available resources
 - SmartSSD, Edge TPUs, CUDA GPUs, Tensor Cores, or something else
- Project ideas
 - Accelerating applications through AI/ML accelerators
 - Accelerating applications through intelligent storage devices
 - Accelerating applications through innovative parallel programming models that hardware accelerators enable
 - Anything related to what we discussed in this class!



Electrical Computer Science Engineering

277

つづく

