

CMPT 450/750: Computer Architecture

Fall 2023

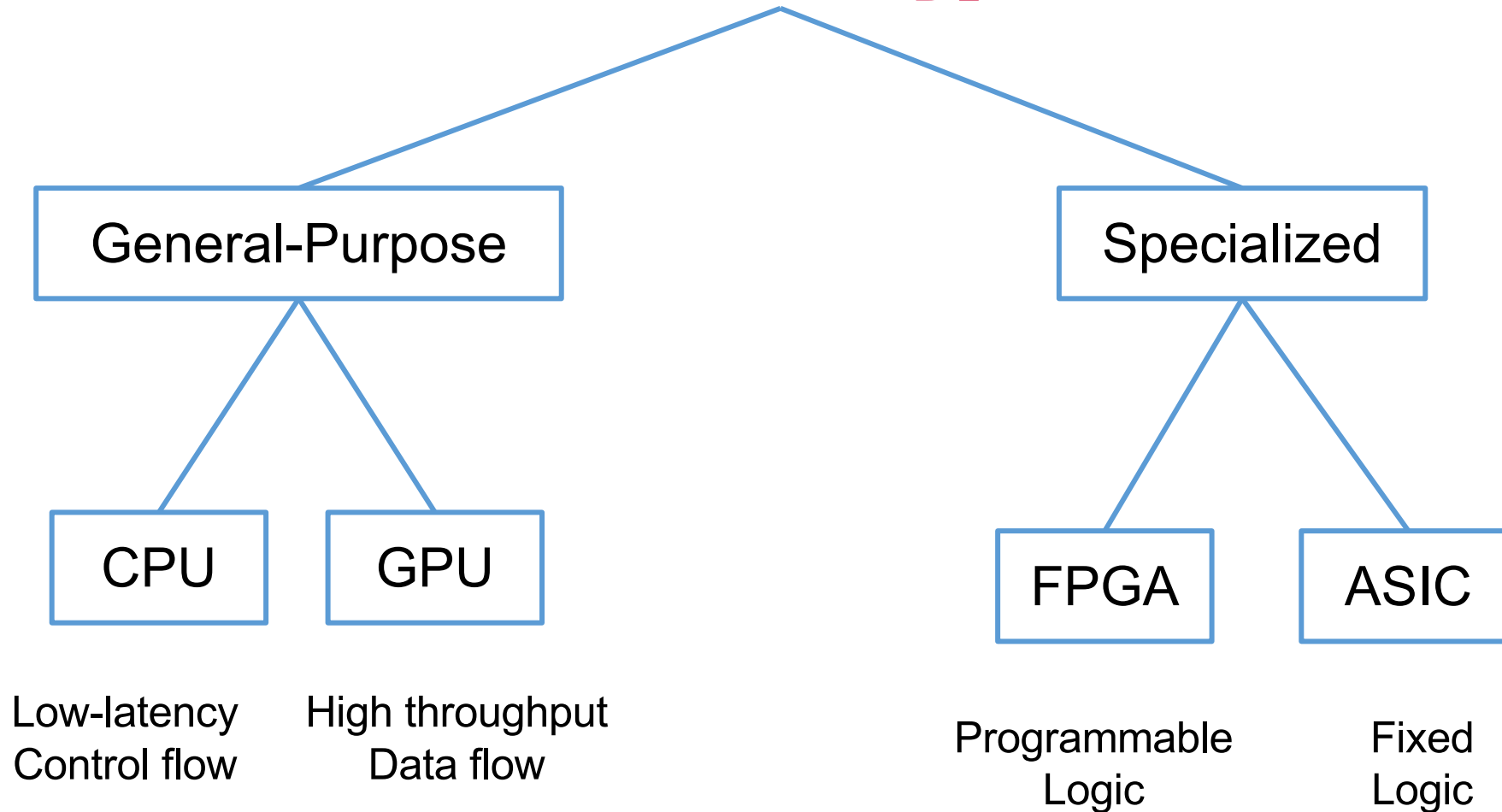
Domain-Specific Architecture I

How did we get here ?

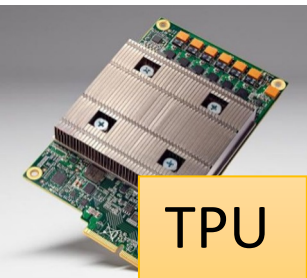
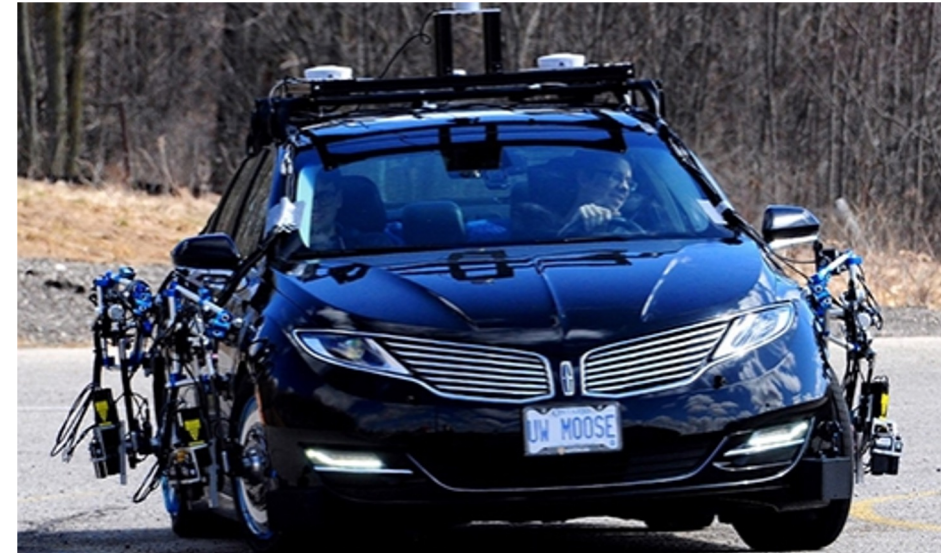
What are they ?

Alaa Alameldeen & Arrvindh Shriraman

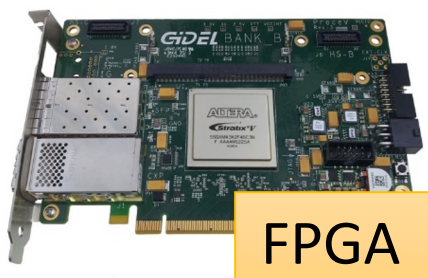
Hardware Types



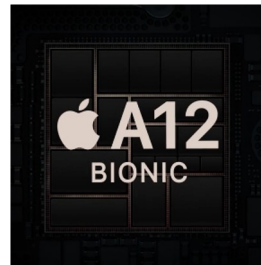
Specialized Hardware



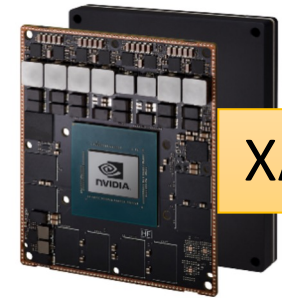
TPU



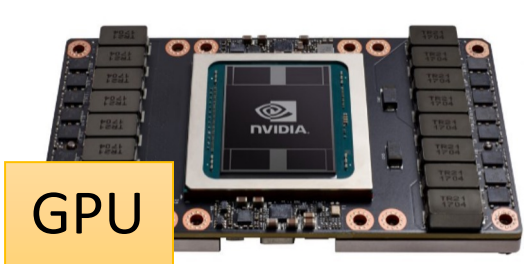
FPGA



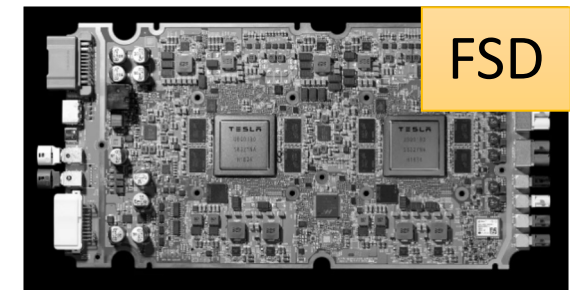
NPU



XAVIER



GPU



FSD

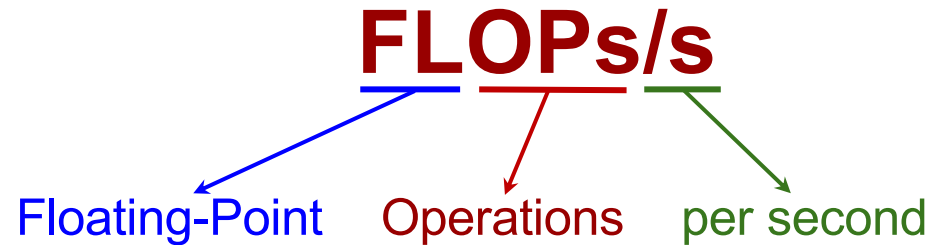
Learning Objectives

By the end of this lecture, you should be able to:

1. Calculate important performance metrics for hardware
2. Optimize the compute and memory efficiency of hardware
3. Analyze emerging hardware architectures

Hardware Metrics & Roofline

Compute Performance Metrics



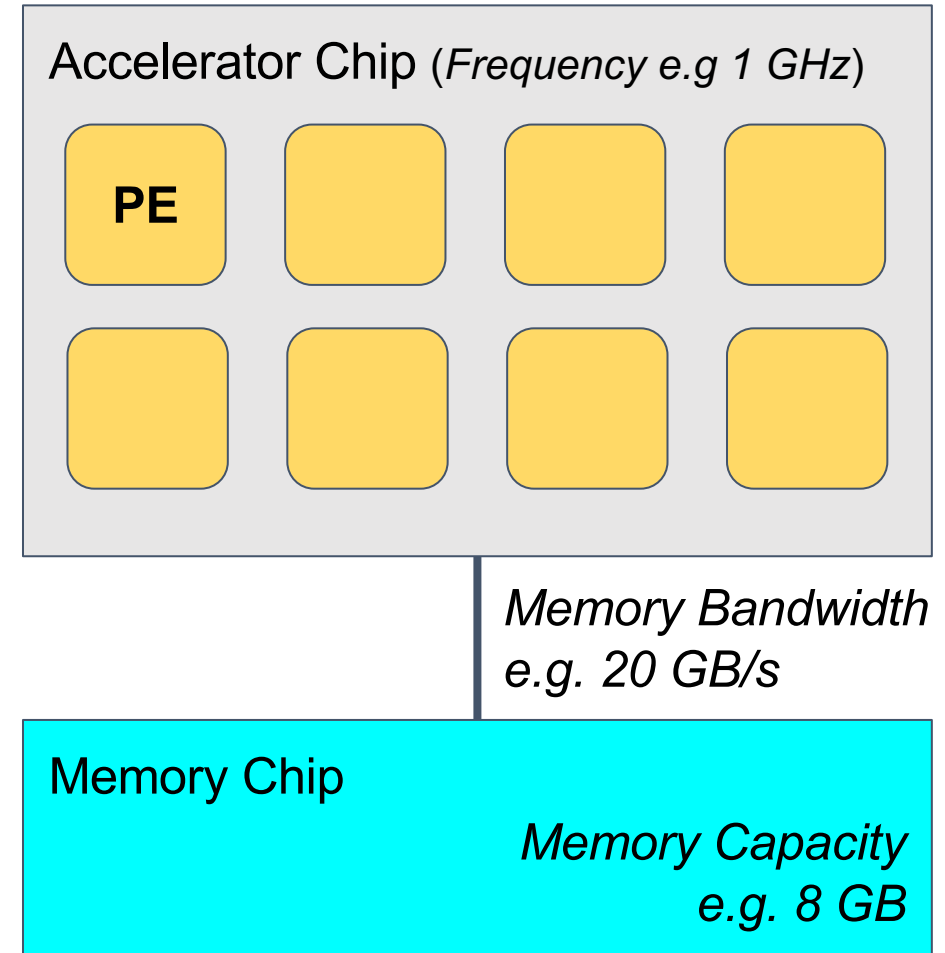
- MACs/s: Multiply-accumulate Ops/s
 - Half FLOPs/s
- OPs/s: for non floating-point operations
- Chips are often labeled with “peak FLOPs/s”
 - Not achievable under normal workloads
 - Very rough indication of performance



$$\frac{\text{operations}}{\text{second}} = \underbrace{\left(\frac{1}{\frac{\text{cycles}}{\text{operation}}} \times \frac{\text{cycles}}{\text{second}} \right)}_{\text{for a single PE}} \times \text{number of PEs}$$

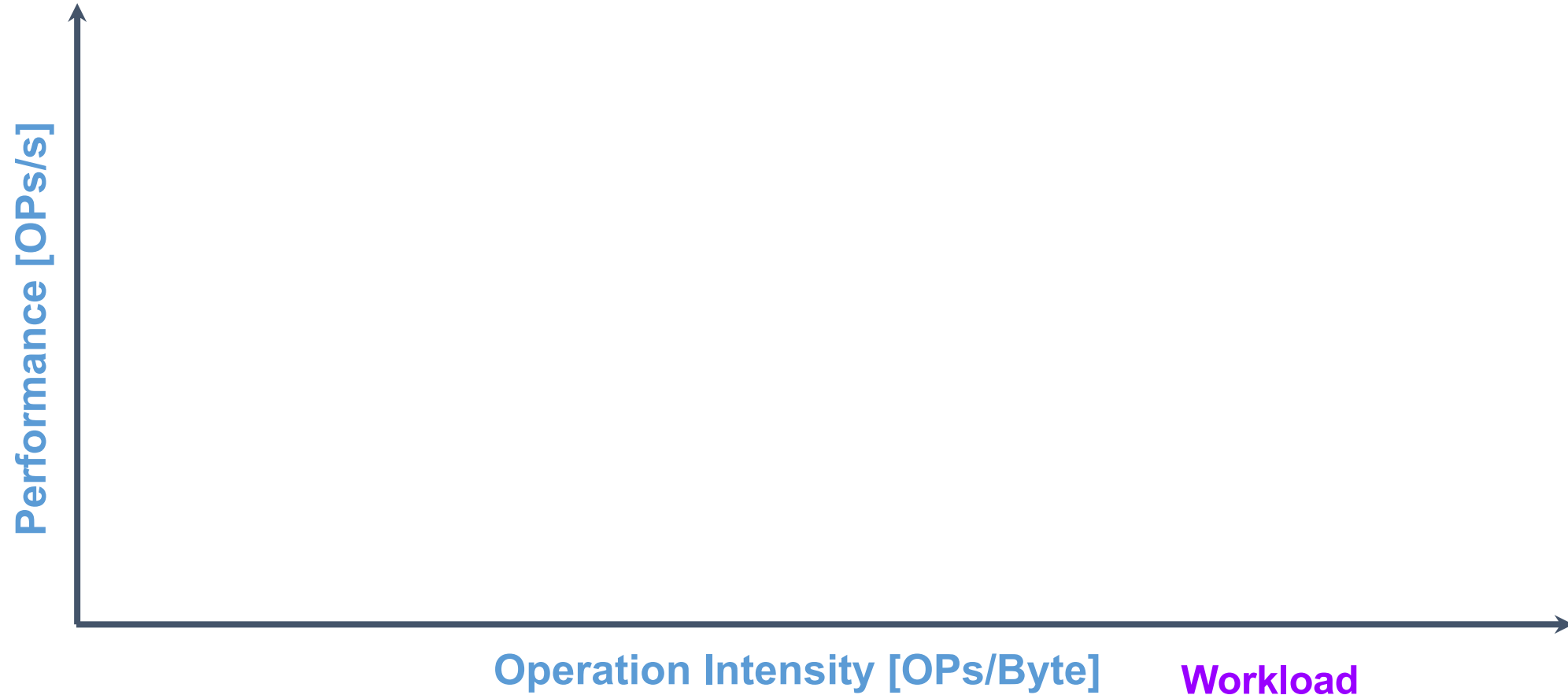
Memory Performance Metrics

- Memory capacity [GB]
- Memory bandwidth [GB/s]
 - Transfer speed from memory chip to compute chip
- More complicated because there is a *memory hierarchy*
 - Showing “external”/”main” memory
 - Can have caches, local memory, registers with much higher bandwidth



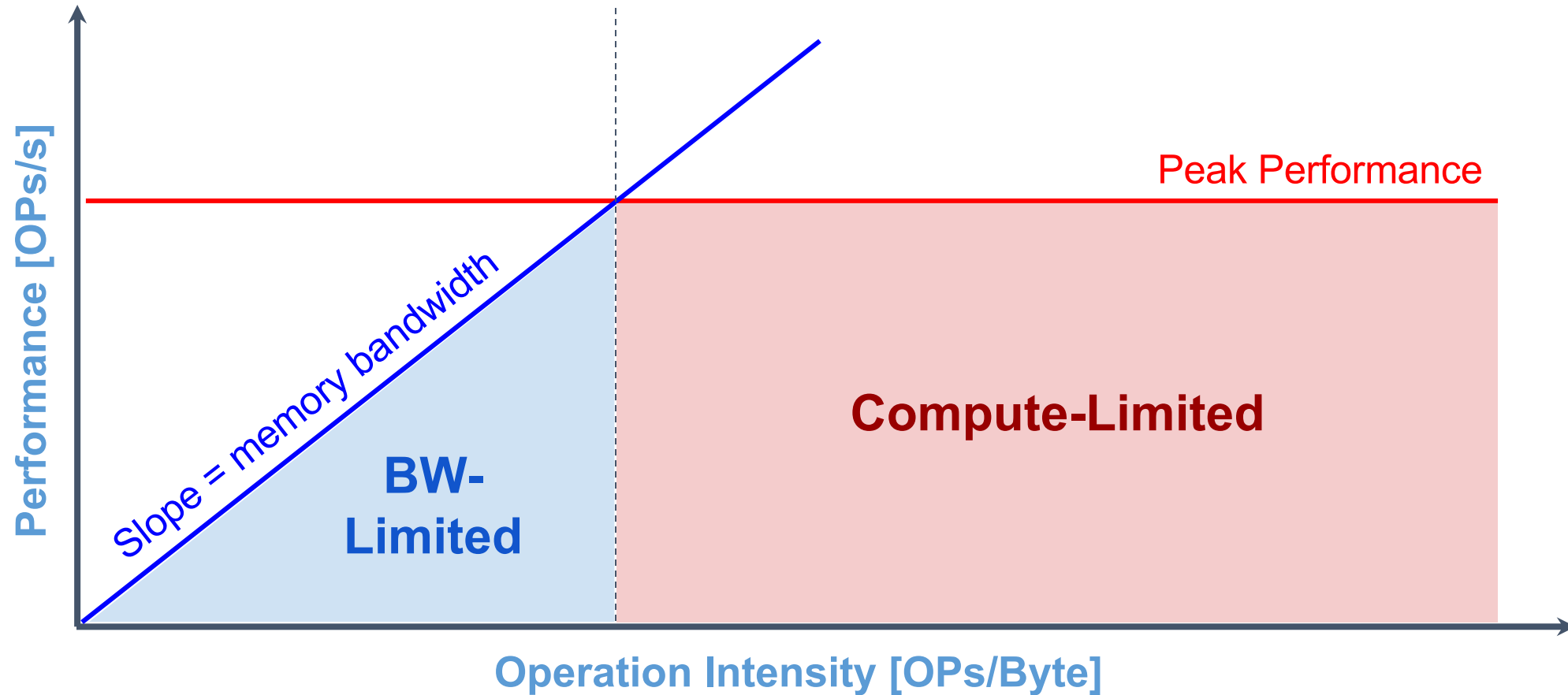
Roofline Plot

Characterize the performance of a given hardware device across different workloads



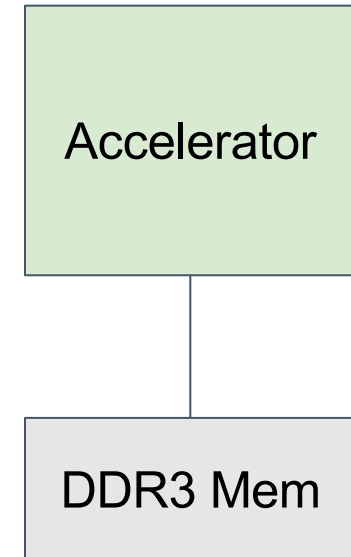
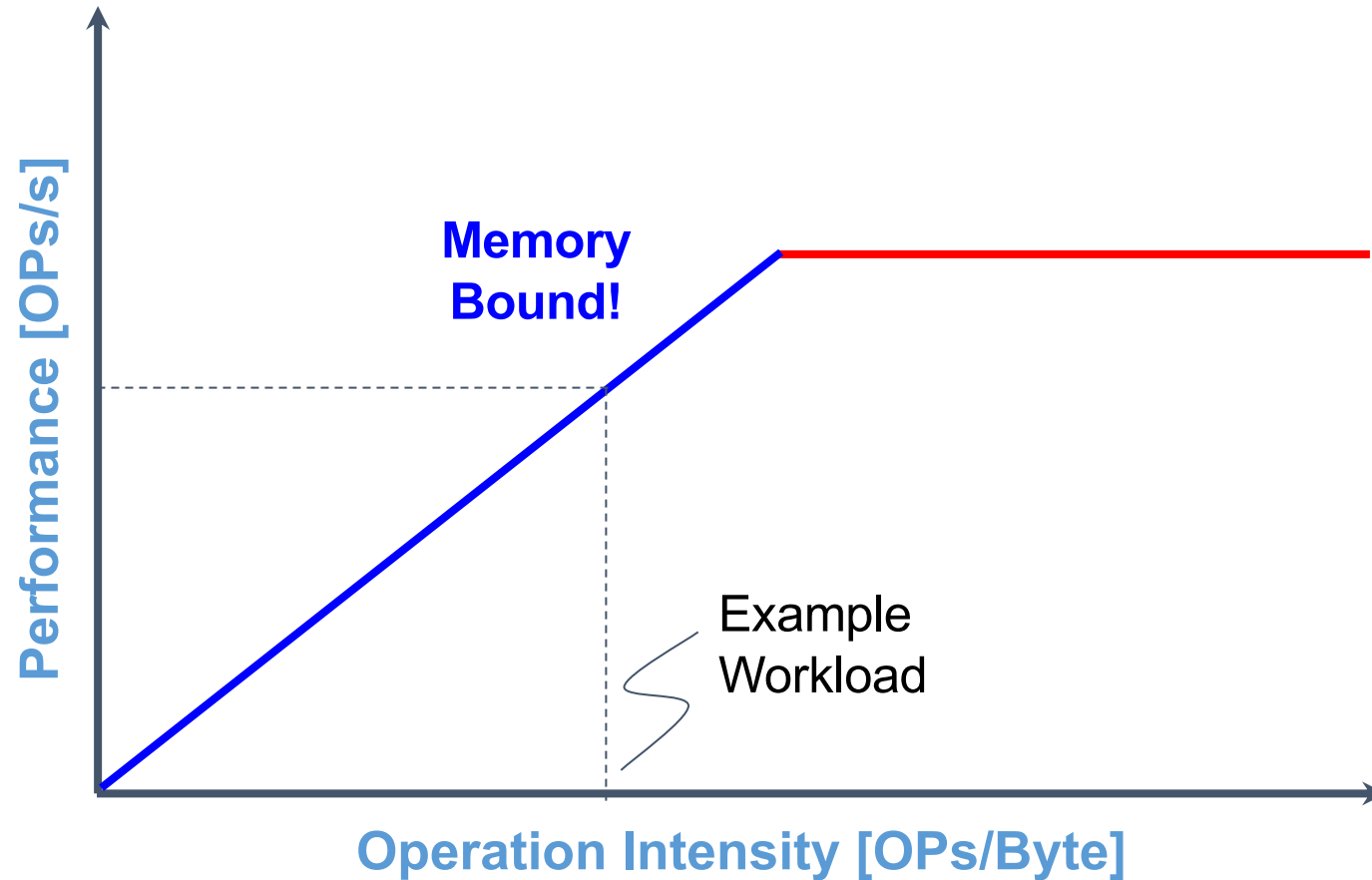
Roofline Plot

Characterize the performance of a given hardware device across different workloads



Roofline Plot

Characterize the performance of a given hardware device across different workloads



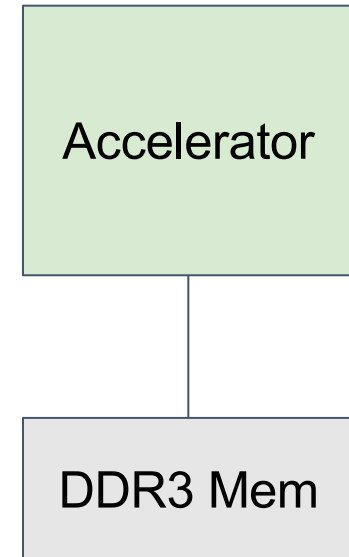
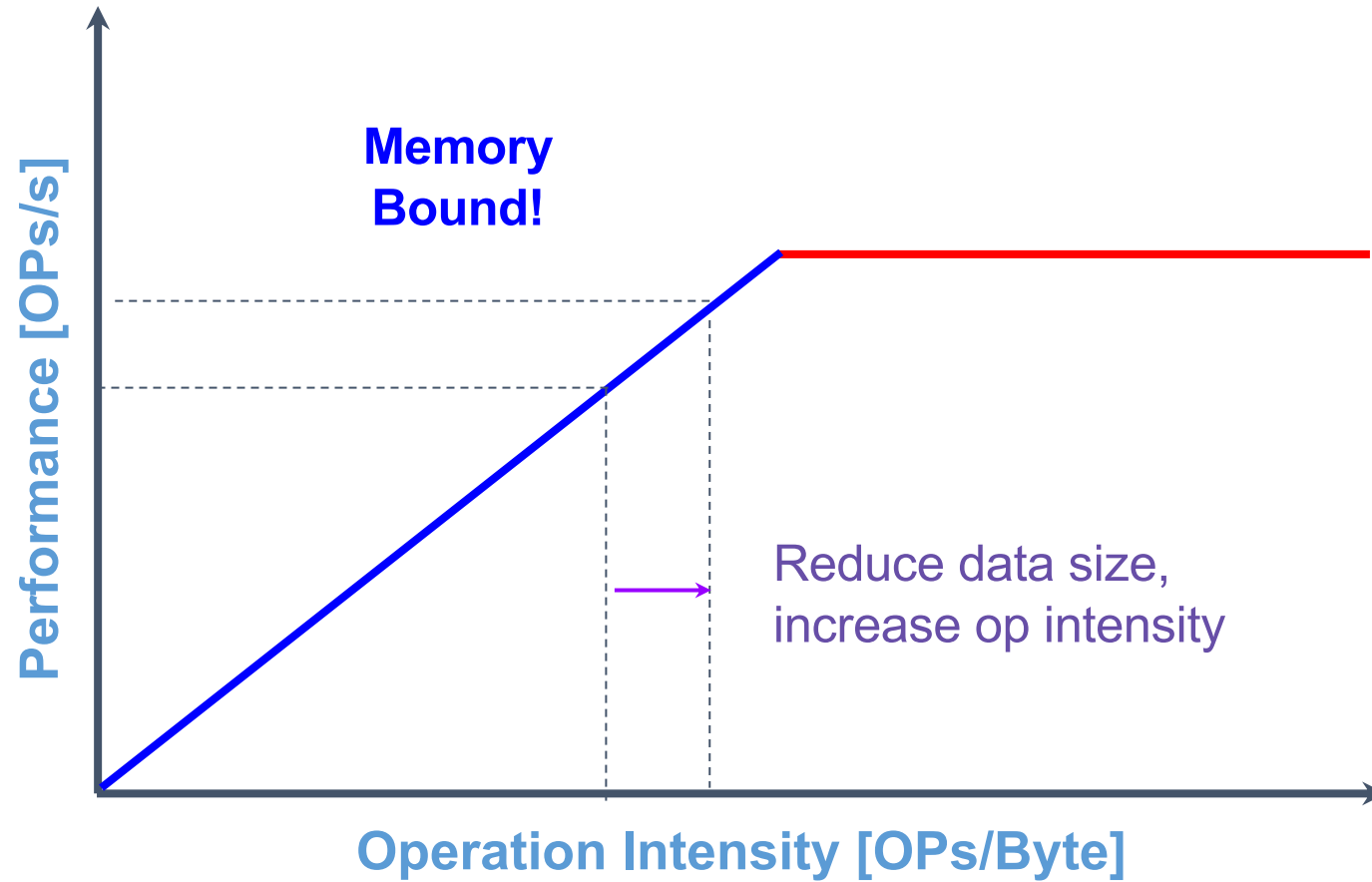
What is OPs/Byte of a DNN?

- Operational intensity [OPs/Byte] quantifies the ratio of computations to memory footprint of a DNN
- Total number of operations = multiplications + additions
- Total memory footprint = size of parameters + size of activations

$$\text{Operational Intensity} = \frac{\text{Total number of operations}}{\text{Total memory footprint}}$$

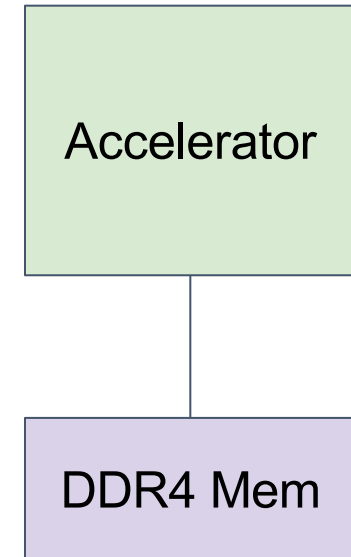
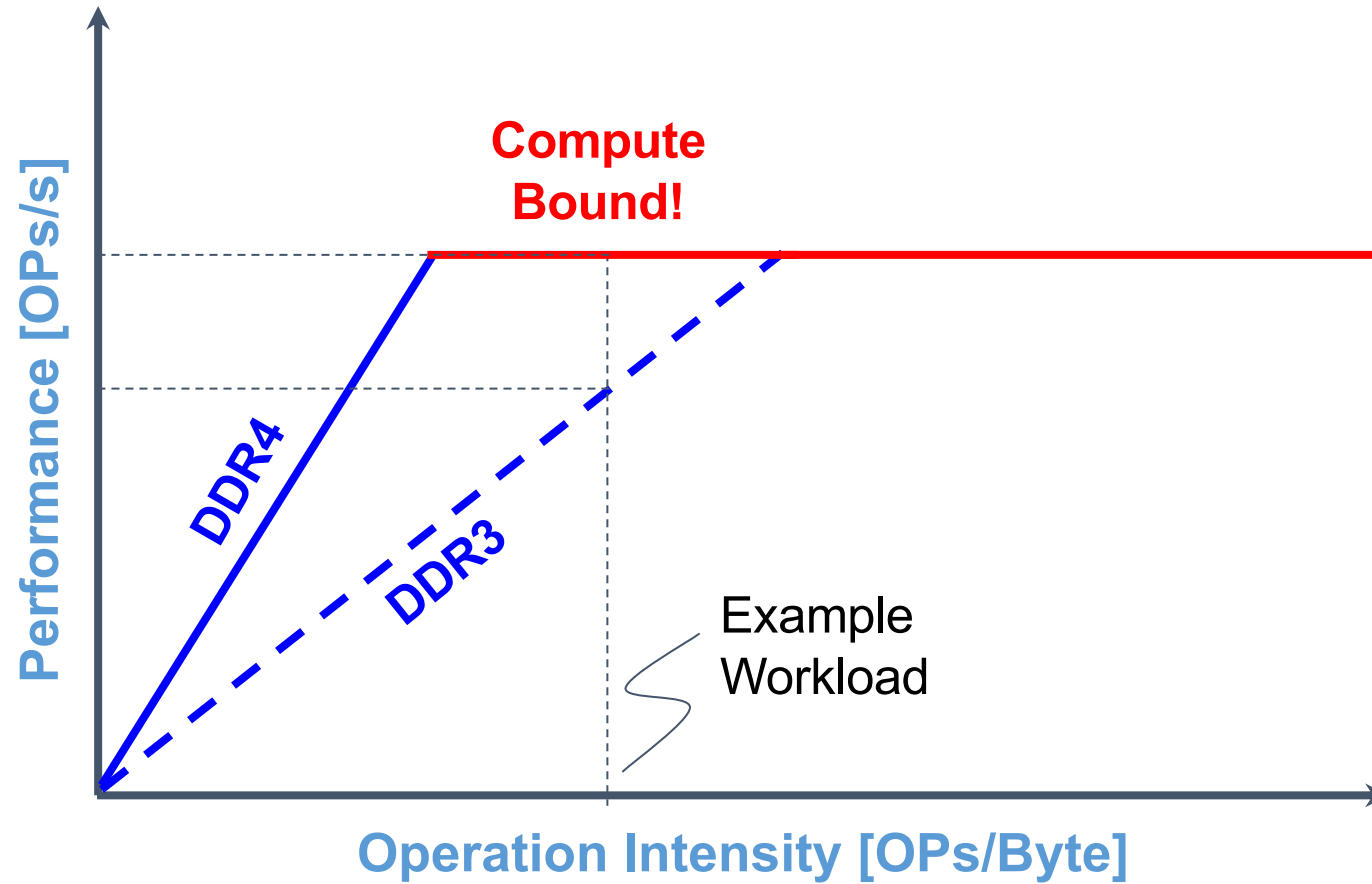
Roofline Plot

Compressed data format e.g. reduced precision



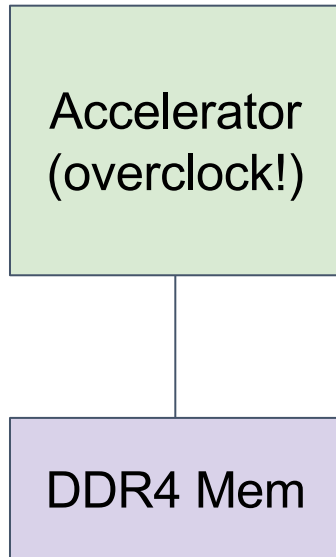
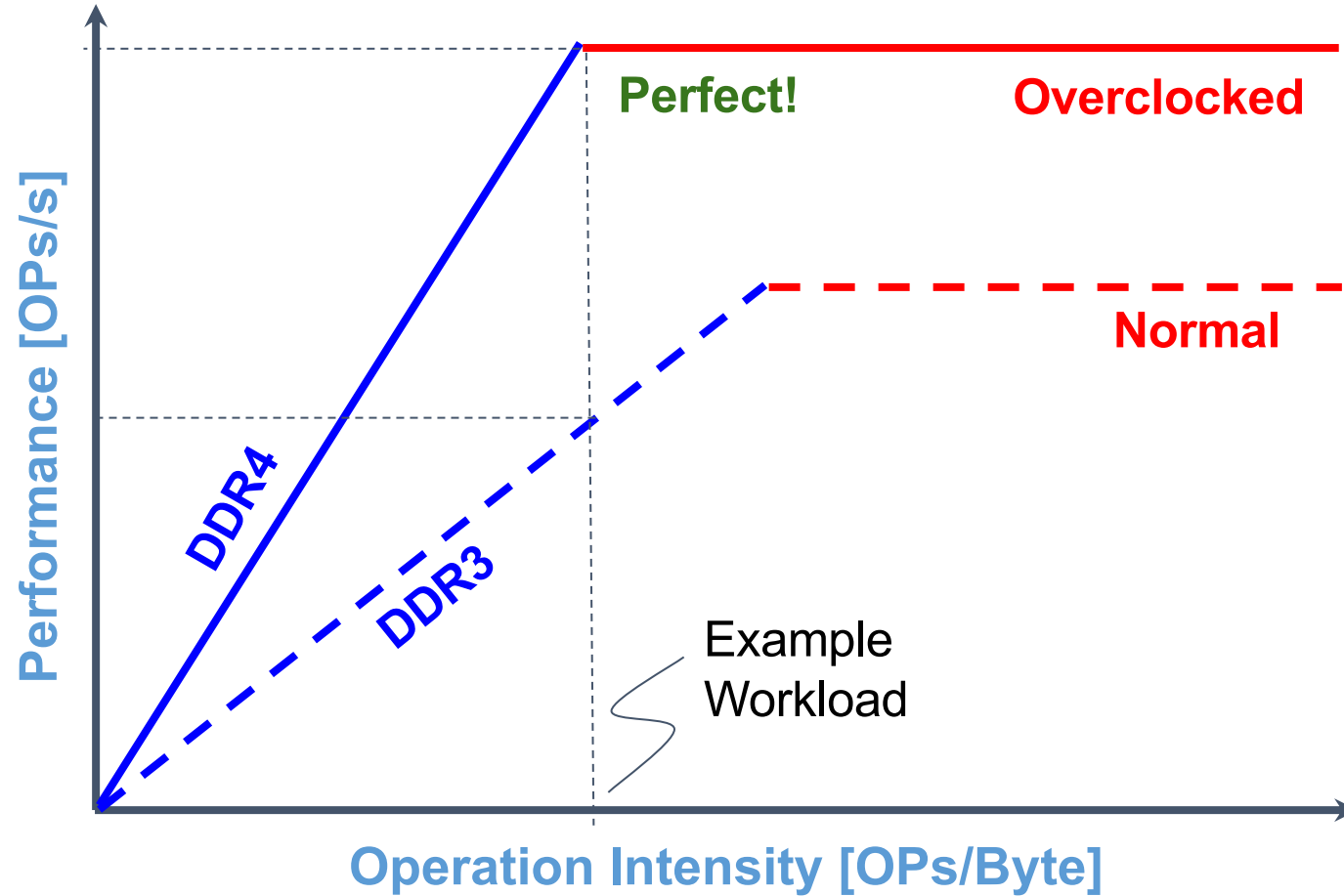
Roofline Plot

Faster memory chip increases slope of roofline



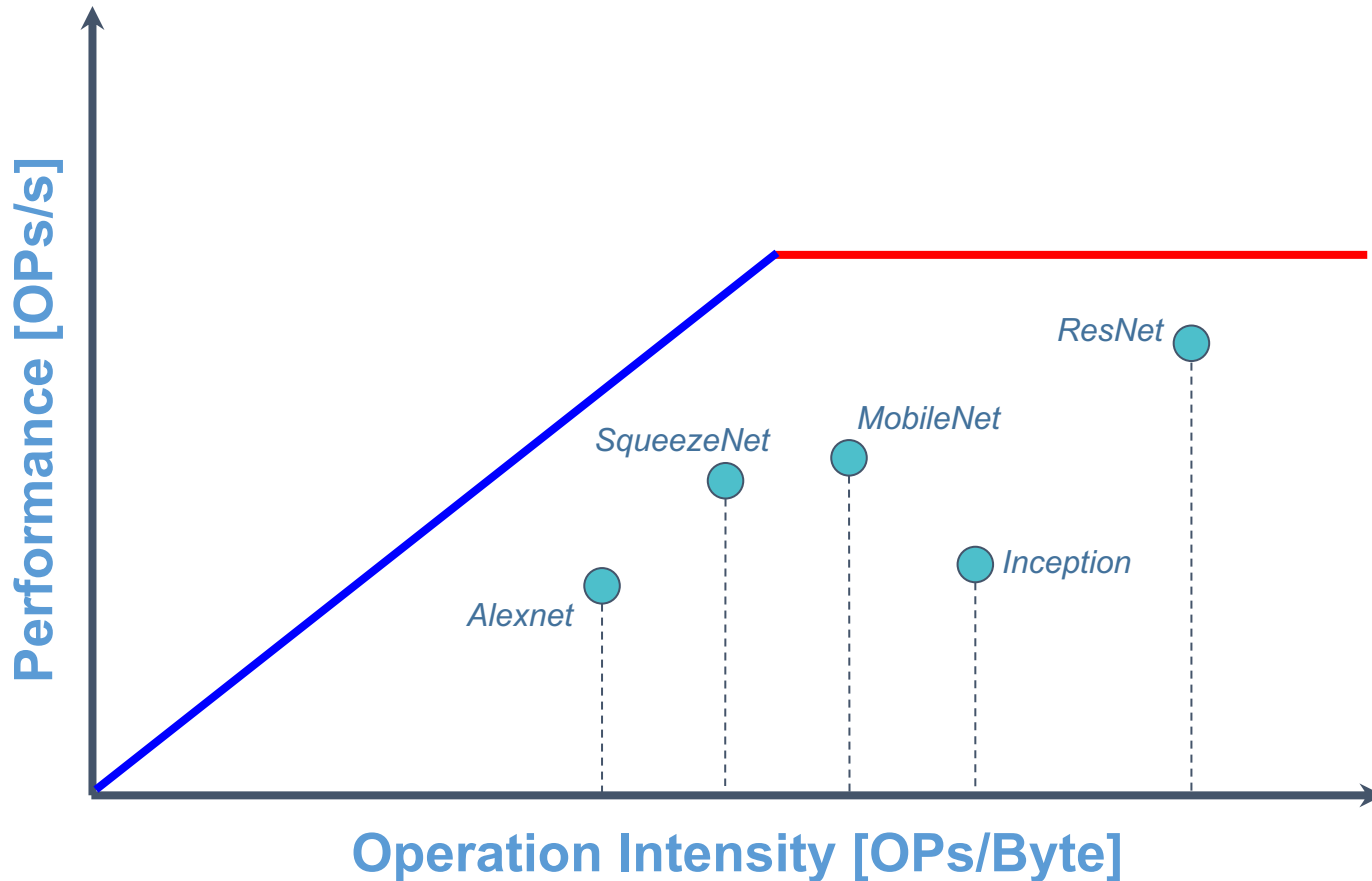
Roofline Plot

Raise the roofline by increasing the speed/throughput of compute



Roofline Example

Measured performance is (by definition) below the roofline.

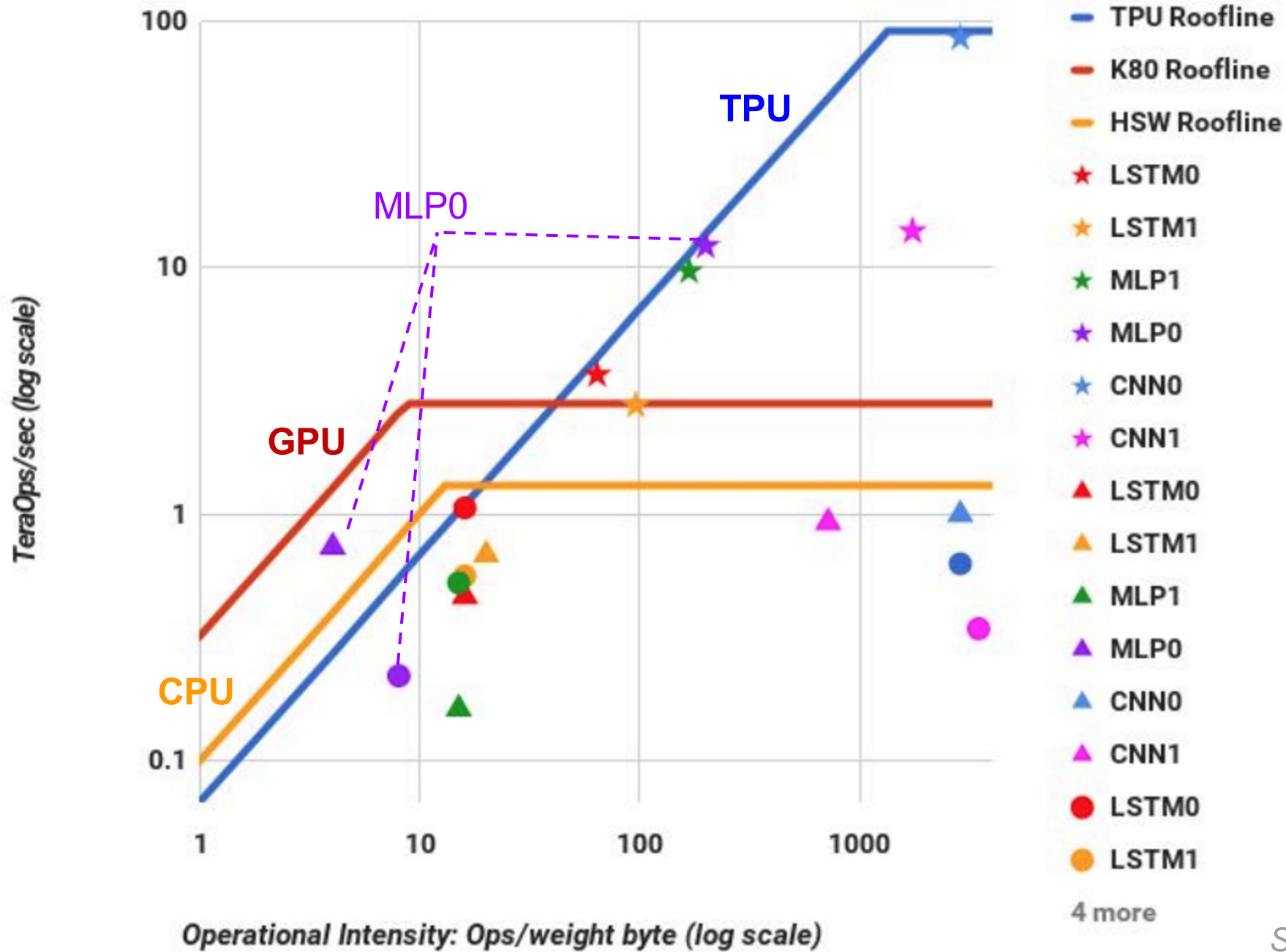


Achieved Performance can be limited by:

- Memory access efficiency
 - E.g.: uncoalesced reads - most DRAM chips require successive reads, each of a specific width to use maximum bandwidth.
- Compute utilization
 - E.g.: In DNN, MAC array hardcoded to 16 channels per tile but first layer has 3 channels
 - Overhead of control logic
- Complexity
 - Control flow and data hazards may stall execution even if the hardware is available

note: points are not plotted in their correct place and are just for illustrative purposes

Google TPU Roofline



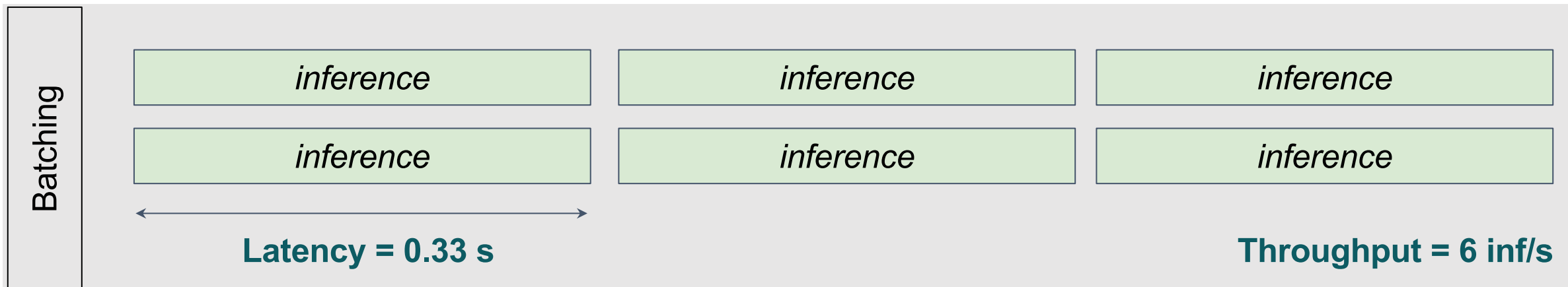
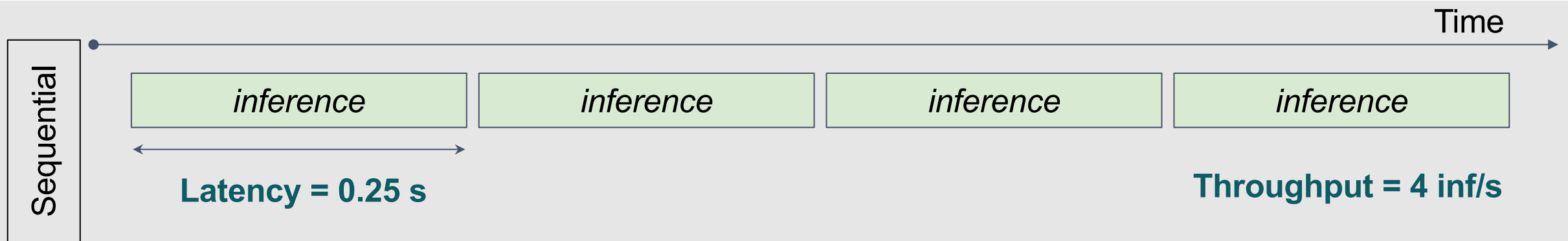
Source: Google

Metrics Summary (so far)

Metric	Hardware	DNN
Peak Performance [OPs/s]	●	
Memory Bandwidth [GB/s]	●	
Operational Intensity [OP/B]		●
HW Utilization	●	●
Throughput [OPs/s]	●	●
Latency [seconds]	●	●

Throughput and Latency

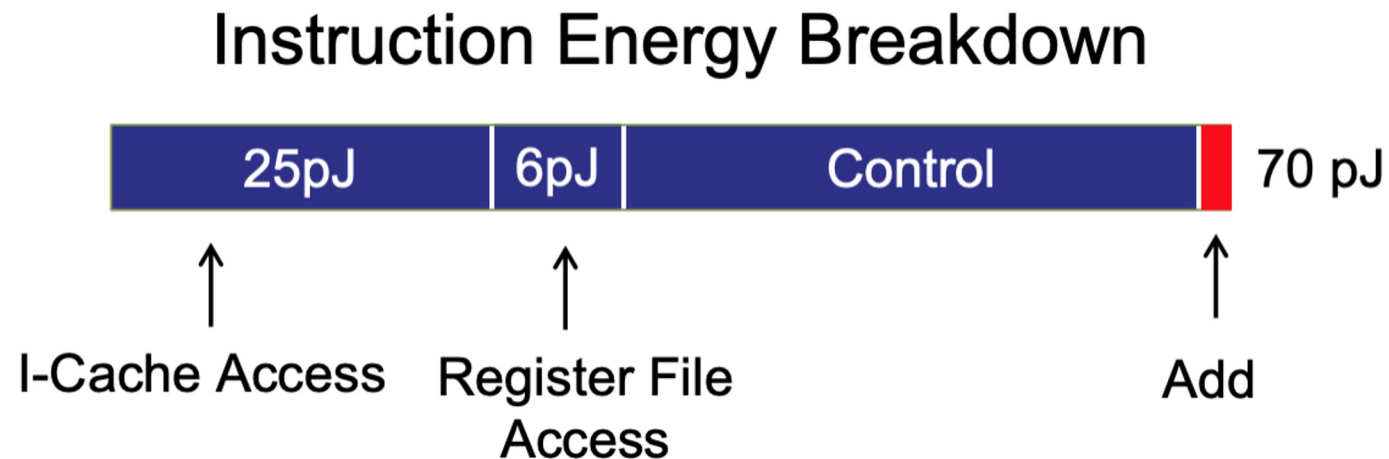
- **Latency:** Number of seconds per inference (unit = seconds)
- **Throughput:** Number of inferences per second (unit = inference/second)



2. Hardware Efficiency

Where does the Energy go?

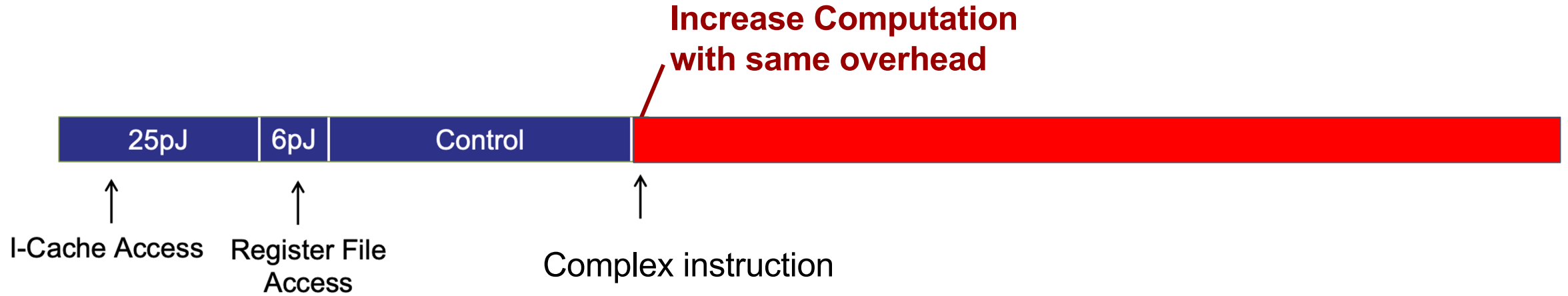
- Energy breakdown of an add instruction in a 45nm CPU
- How can we optimize this?



Hardware Efficiency Approaches

1. **Arithmetic**
 - Specialized Instructions: To amortize overhead.
 - Lower precision (Quantization)
2. **Memory**
 - Locality: Move data to inexpensive on-chip memory.
 - Reuse: To avoid expensive memory fetches.
3. **Ineffectual Operations**
 - Sparsity: Skip useless operations

Amortize Overhead



Half-precision
Fused Multiply-Add

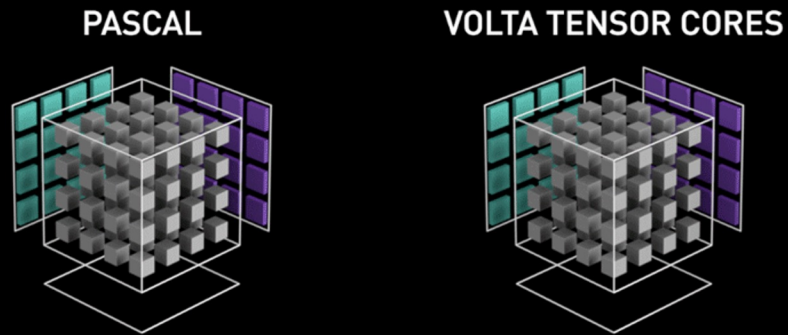
4-way dot-product

16x16 matrix
multiplication

Operation	Energy**	Overhead*
HFMA	1.5pJ	2000%
HDP4A	6.0pJ	500%
HMMA	110pJ	27%

“Special” Instruction Examples

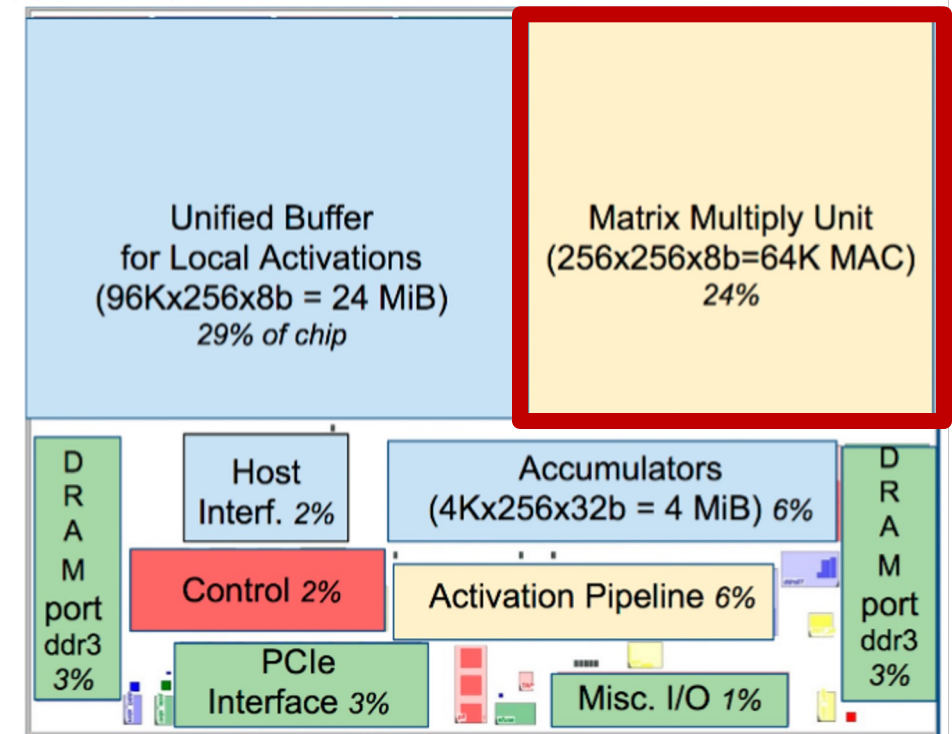
GPU



$$16 \times 16 = 256^* \text{ MAC/cycle}$$

*~ 500 tensor cores per GPU

ASIC (TPUv1)



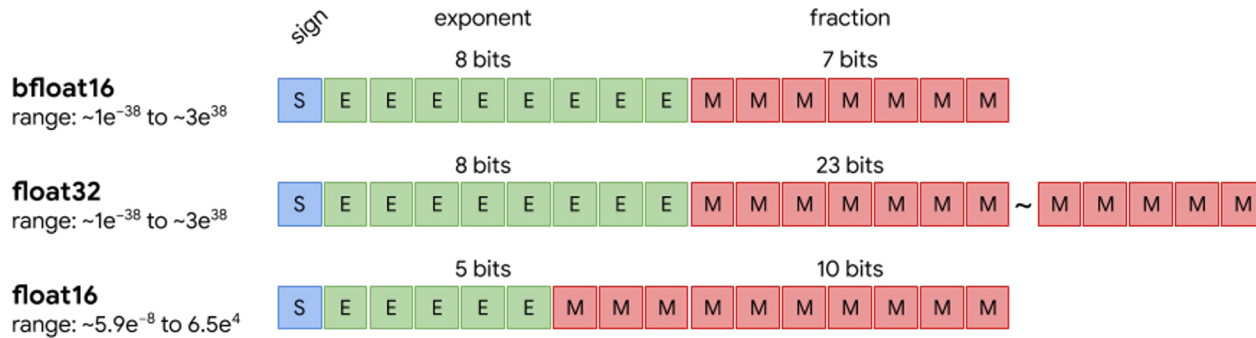
$$256 \times 256 = 64 \text{ kMAC/cycle}$$

Source: Google

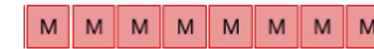
Source: Nvidia

Numerical Format and Precision

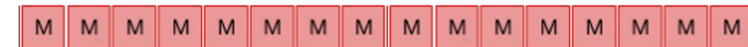
Floating Point



Integer



8-bit

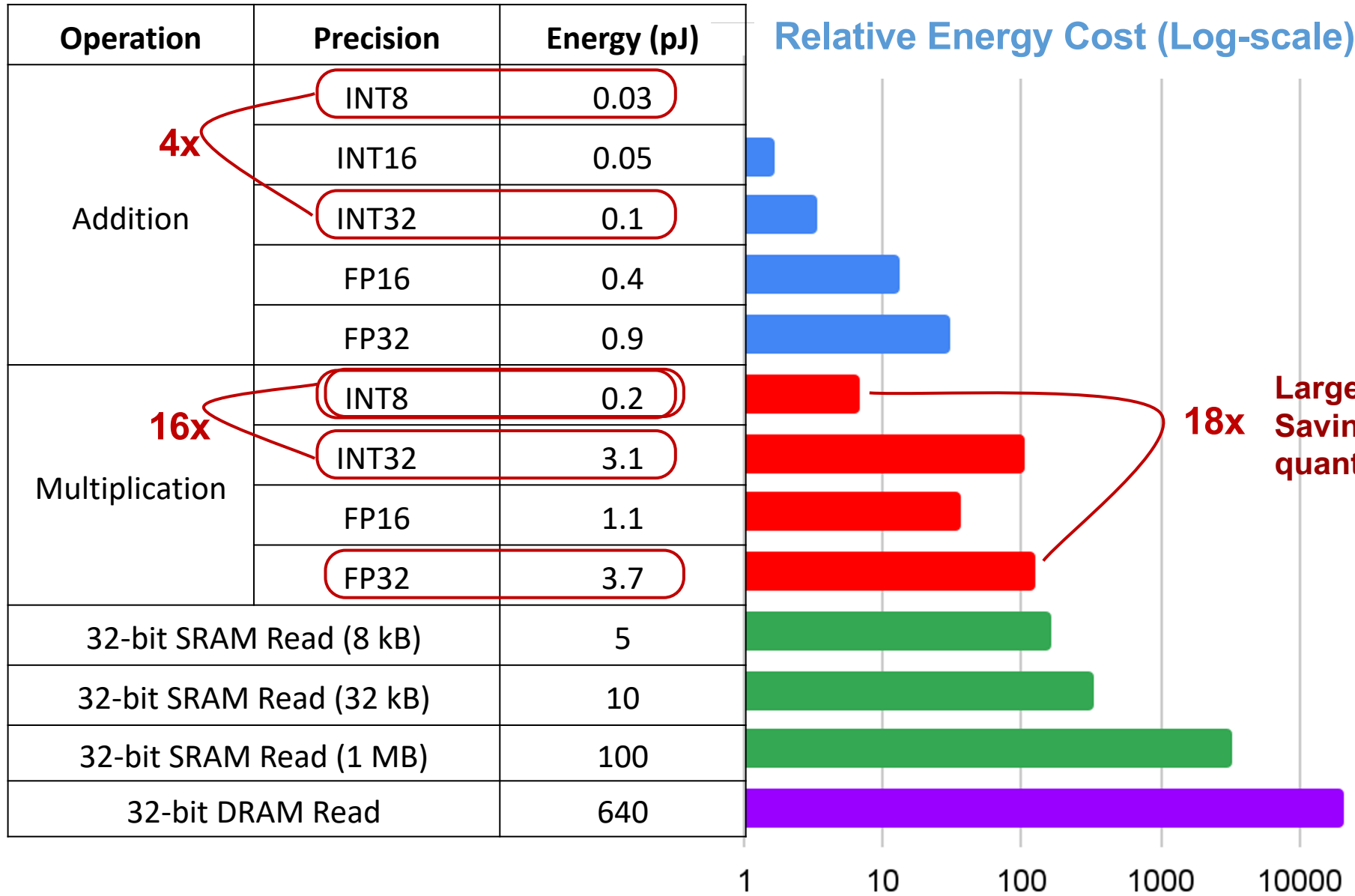


16-bit

- IEEE standard includes FP32 and FP16
- Many exotic FP numbers in DNN
 - E.g. bfloat, minifloat

- Whole numbers only
- (typically) much cheaper circuit area and power

Cost of Arithmetic Operations



Adapted from Mark Horowitz "Computing's Energy Problem (and What we can do about it)" ISSCC 2014


QUESTION

Why is floating-point add so expensive compared to integer add?

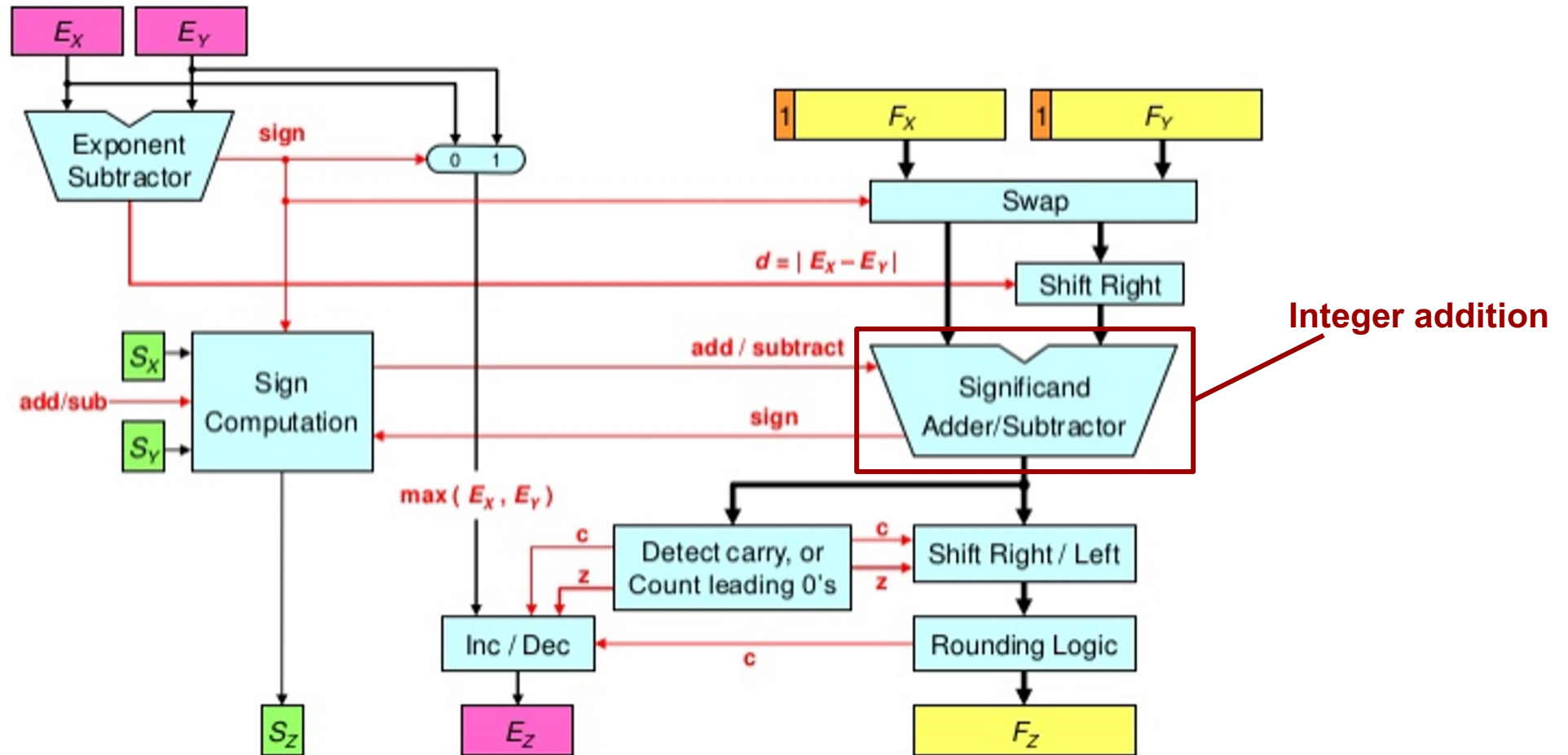
Operation	Precision	Energy (pJ)
Addition	INT8	0.03
	INT16	0.05
	INT32	0.1
	FP16	0.4
	FP32	0.9
Multiplication	INT8	0.2
	INT32	3.1
	FP16	1.1
	FP32	3.7

9x

1.2x



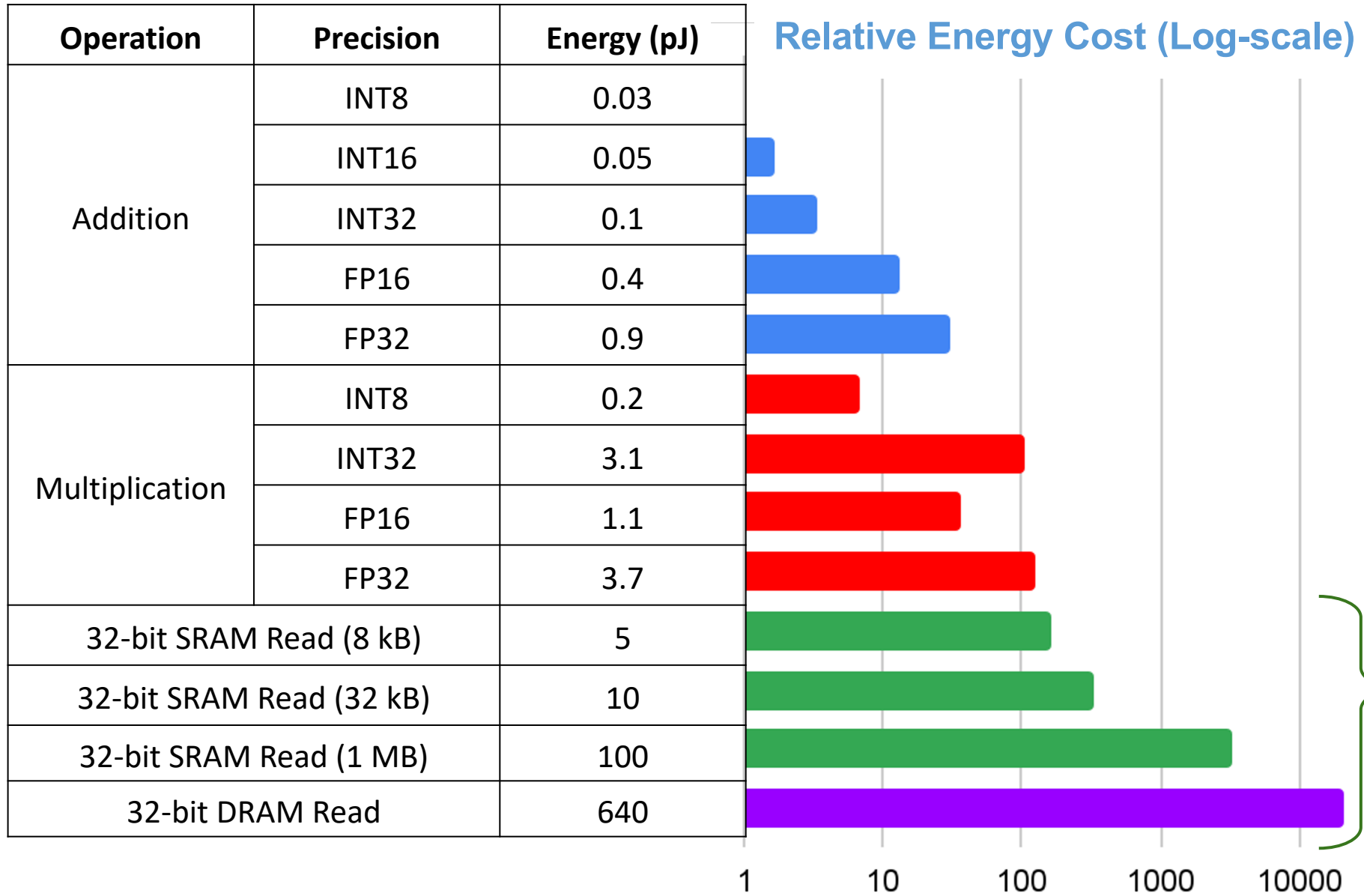
Floating-Point Addition



2. Hardware Efficiency

1. Arithmetic
 - Specialized Instructions: To amortize overhead. ✓
 - Lower precision (Quantization) ✓
2. Memory
 - Locality: Move data to inexpensive on-chip memory.
 - Reuse: To avoid expensive memory fetches.
3. Ineffectual Operations
 - Sparsity: Skip useless operations
 - Compressed Sparse Column (CSC) Format

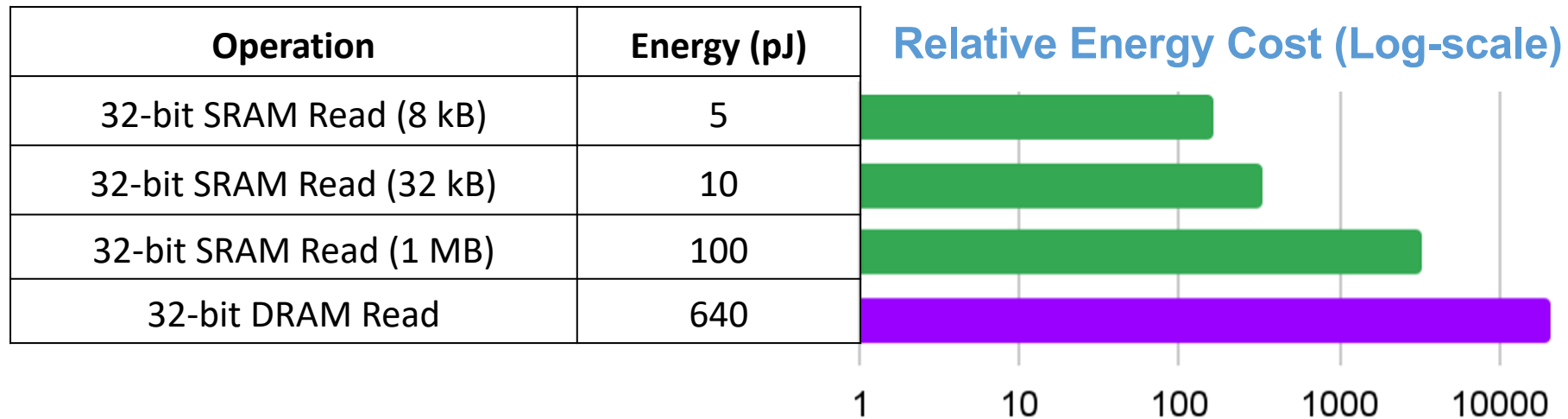
Cost of Arithmetic Operations



Memory is the issue!

Memory Hierarchy Optimizations

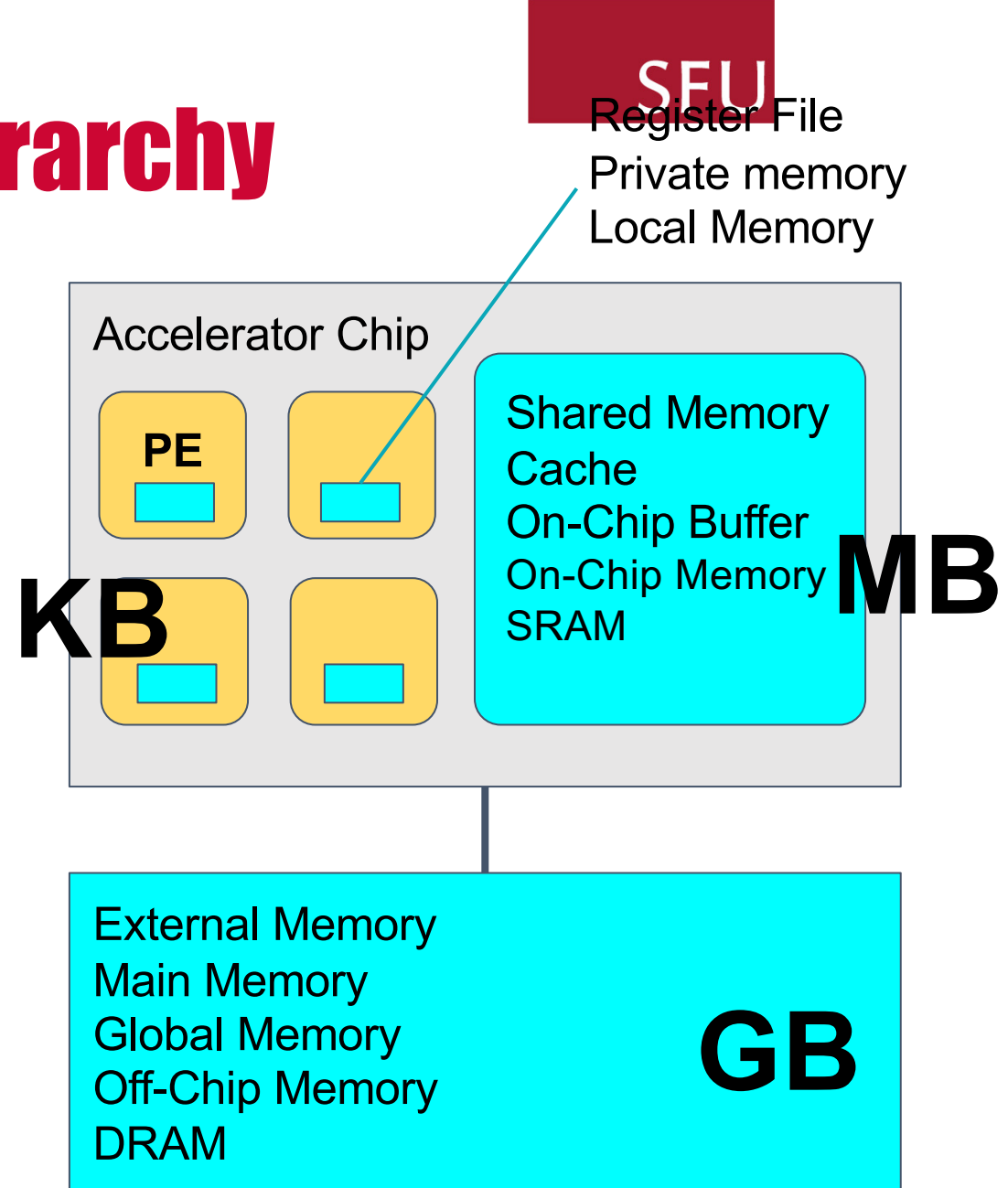
1. Get data close to the computation. (**LOCALITY**)
2. Once data is close - perform all computations with this data. (**REUSE**)



Memory Hierarchy

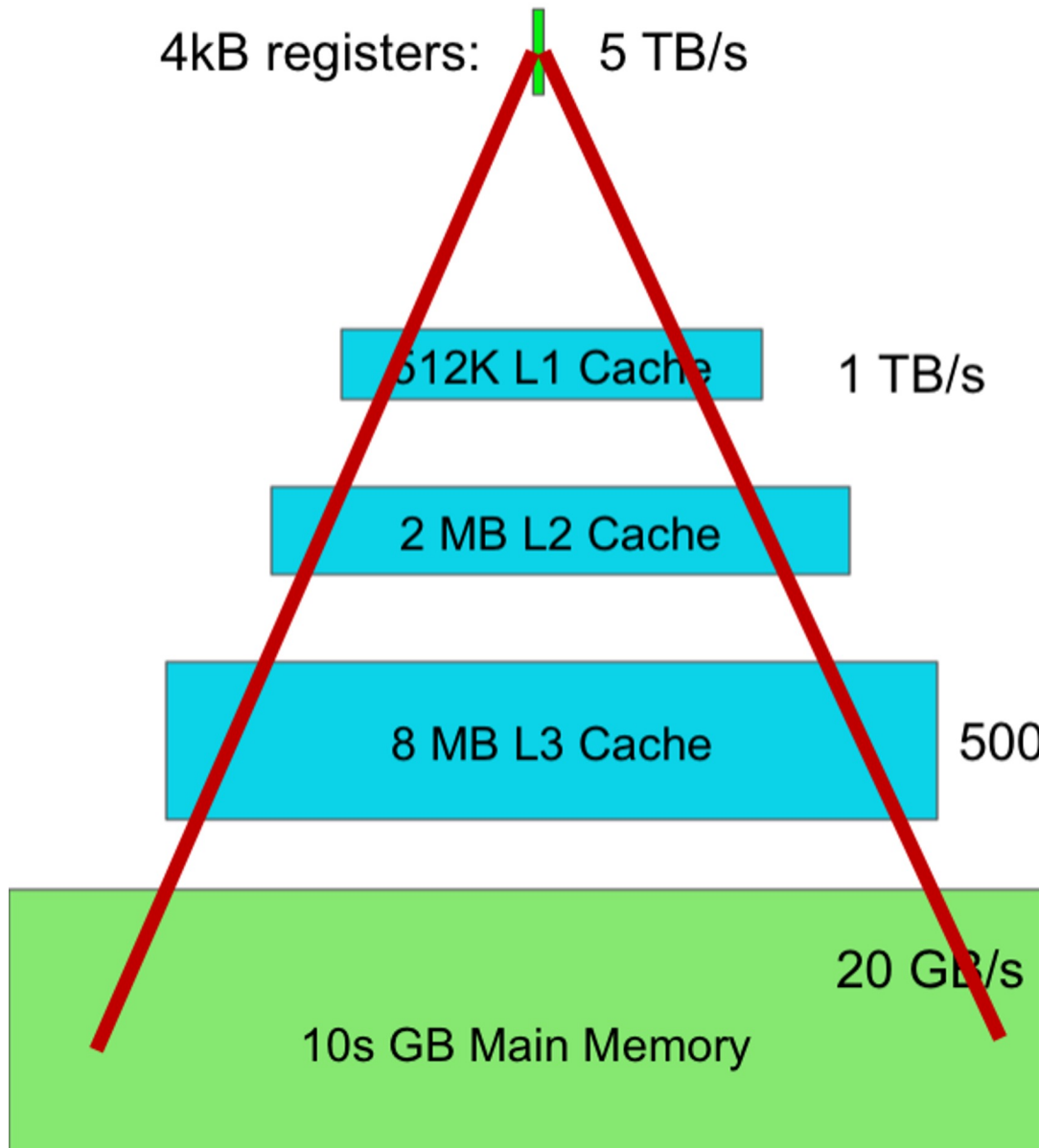
Why do we have a memory hierarchy?

- The closer you get to compute, the more \$\$ and scarce the memory resource becomes
- In *most* cases, the DNN parameters live off chip and are fetched layer-by-layer or tile-by-tile
- Data locality: how to get data close to the PEs (to keep them fully utilized)

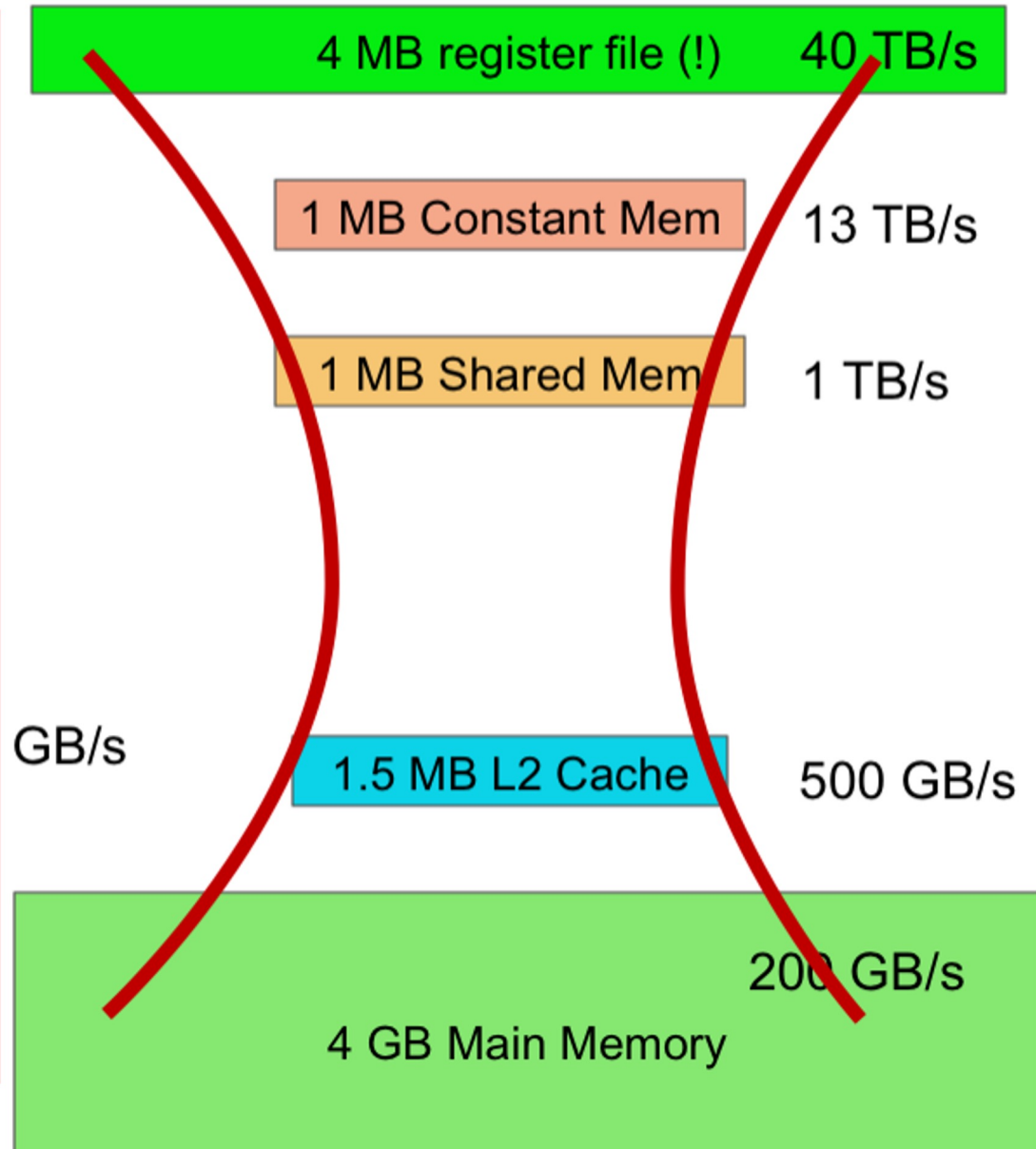


Memory Hierarchy Examples

Intel® 8 core Sandy Bridge CPU



NVIDIA® GK110 GPU



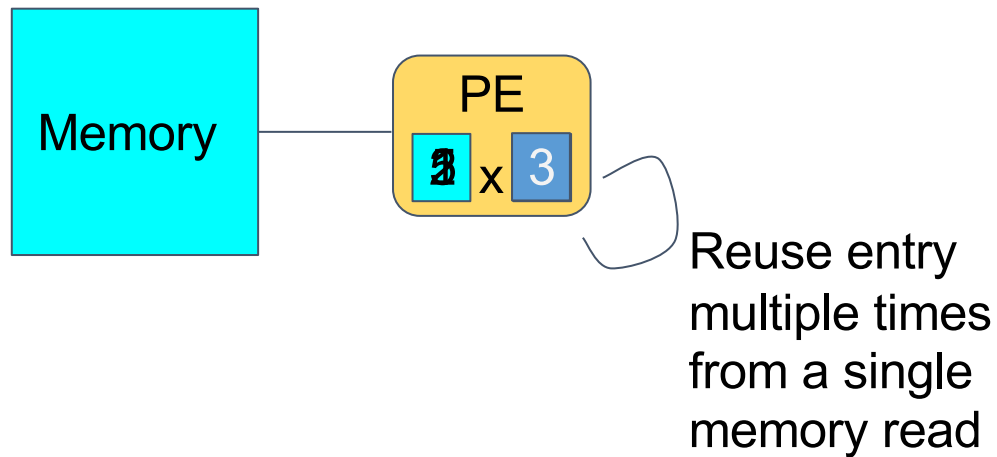
Source: Nvidia

Data Reuse

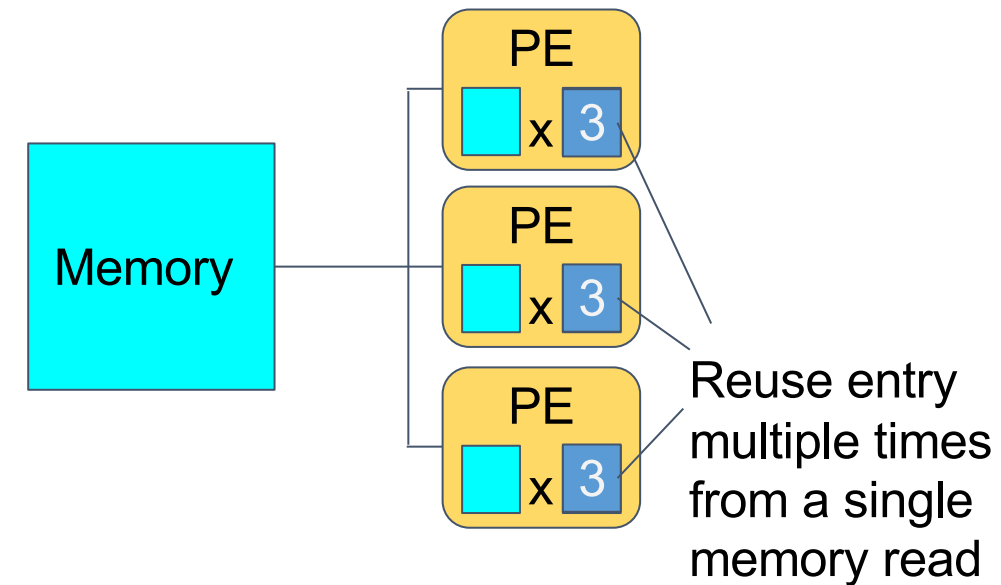
Temporal Reuse

Spatial Reuse

Read once from memory, use same data multiple times by same PE



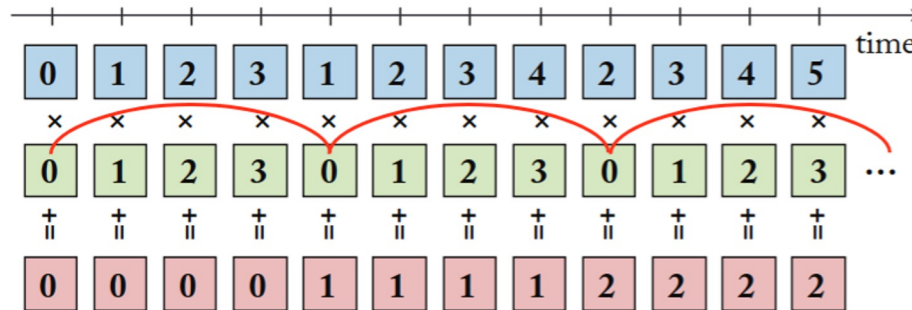
Read once from memory, use same data multiple times by multiple PEs



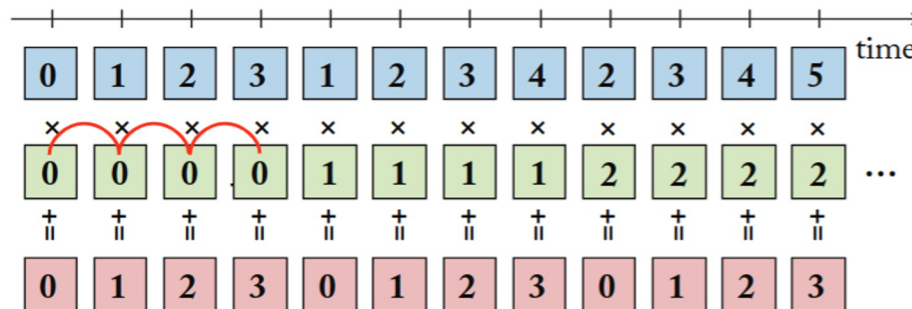
Temporal Reuse Example

$$\begin{bmatrix} 0 & 1 & 2 & 3 & 4 & 5 & 6 \end{bmatrix} * \begin{bmatrix} 0 & 1 & 2 & 3 \end{bmatrix} = \begin{bmatrix} 0 & 1 & 2 & 3 \end{bmatrix}$$

(a) Example 1D convolution



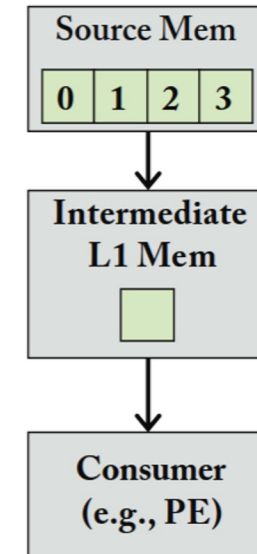
(b) Operation ordering 1: weight reuse distance = 4



(c) Operation ordering 2: weight reuse distance = 1

Reuse Distance

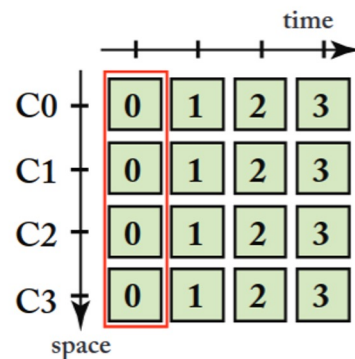
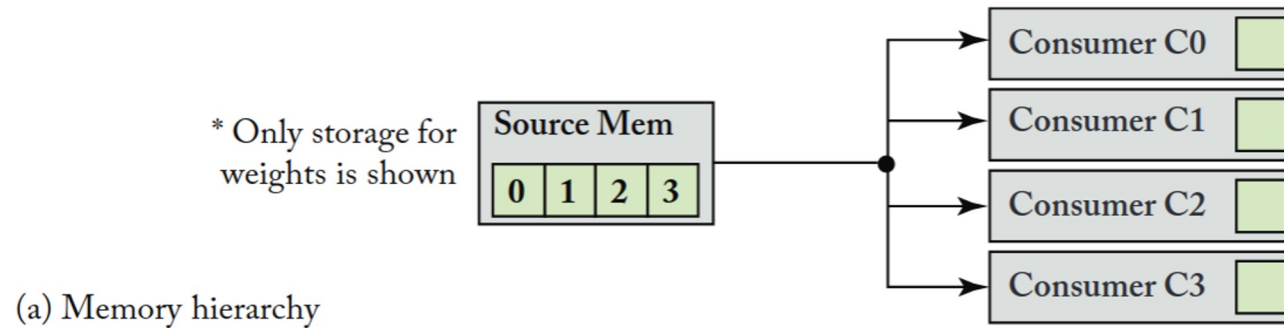
- Number of memory accesses
- Size of your on-chip memory
- Smaller is better



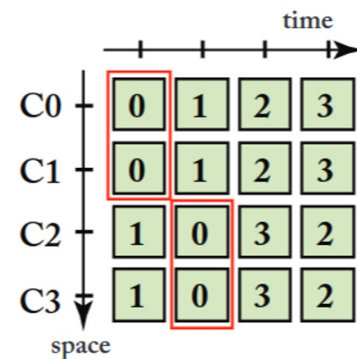
* Only storage for weights is shown

(d) Memory Hierarchy

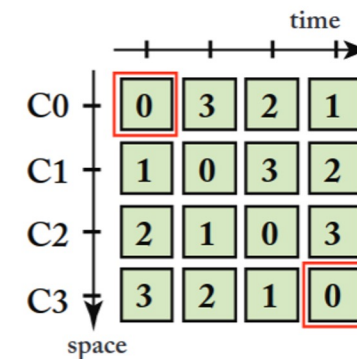
Spatial Reuse Example



(b) Operation ordering 1



(c) Operation ordering 2

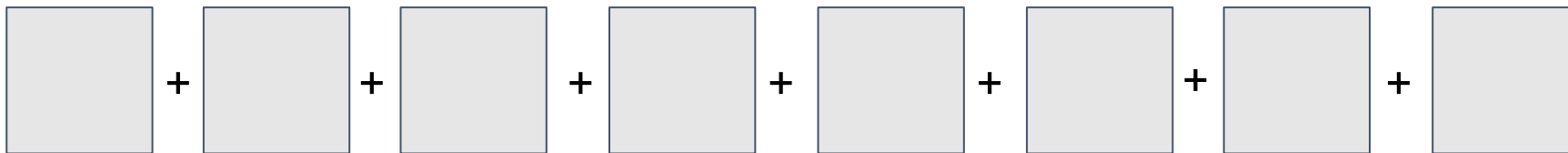
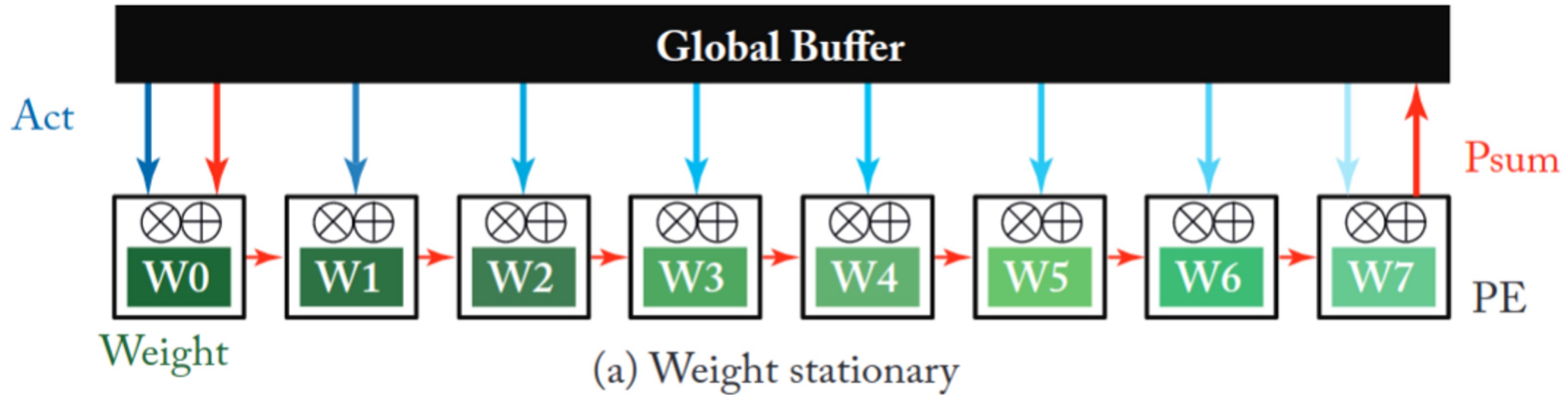


(d) Operation ordering 3

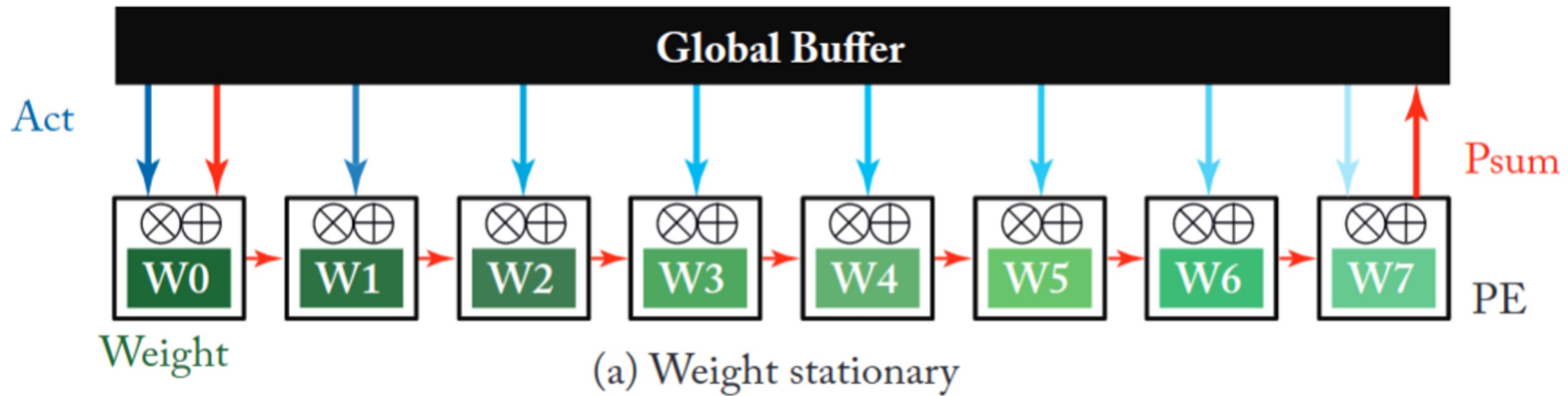
Why a smaller reuse distance with spatial reuse?

- Fewer memory read ports
 - Less area
- Fewer reads
 - Less power

What to Keep Stationary? Weights

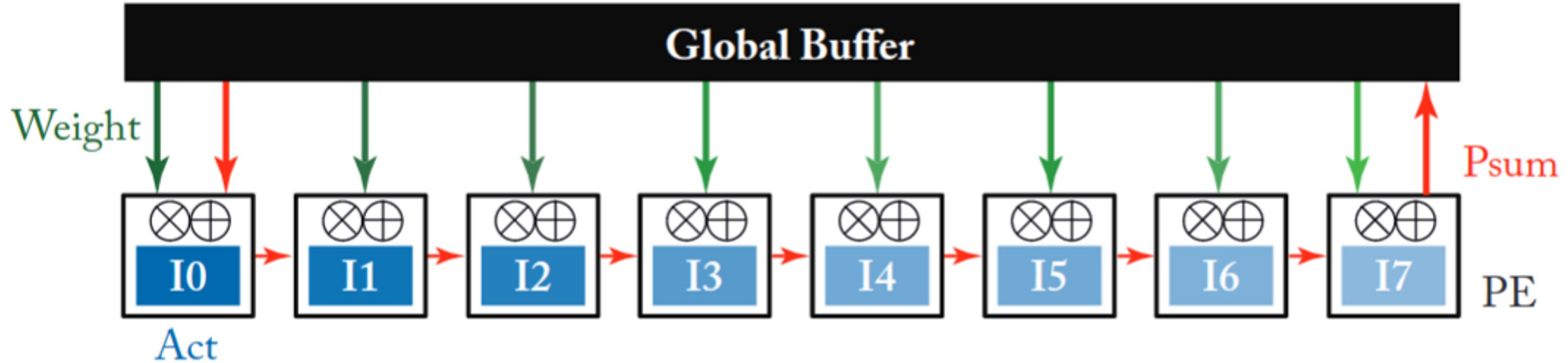


What to Keep Stationary? Weights

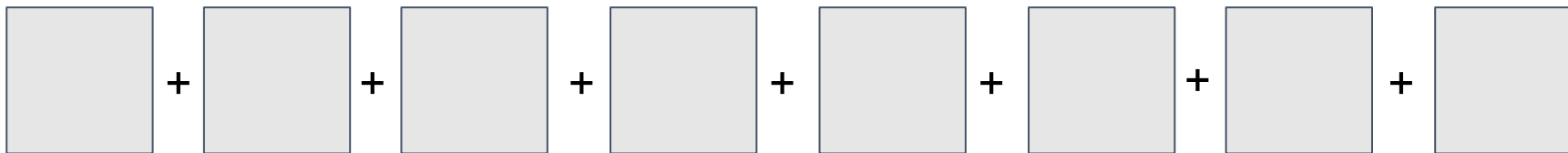


$$a_0w_0 + a_1w_1 + a_2w_2 + a_3w_3 + a_4w_4 + a_5w_5 + a_6w_6 + a_7w_7$$

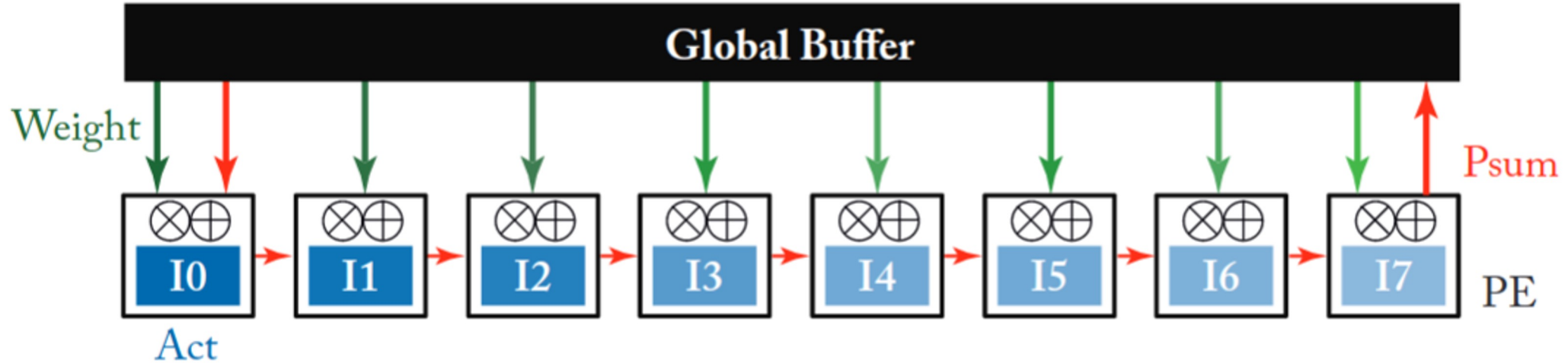
What to Keep Stationary? Inputs



(c) Input stationary



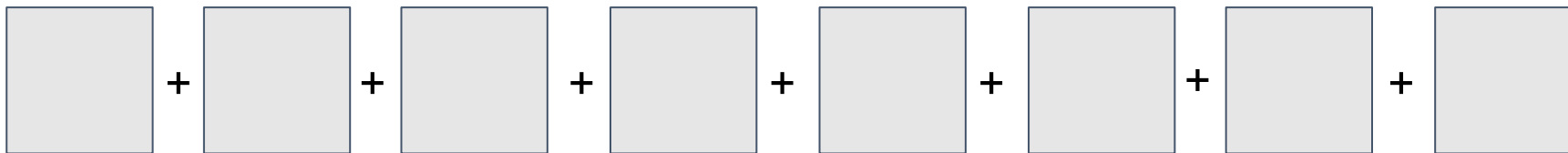
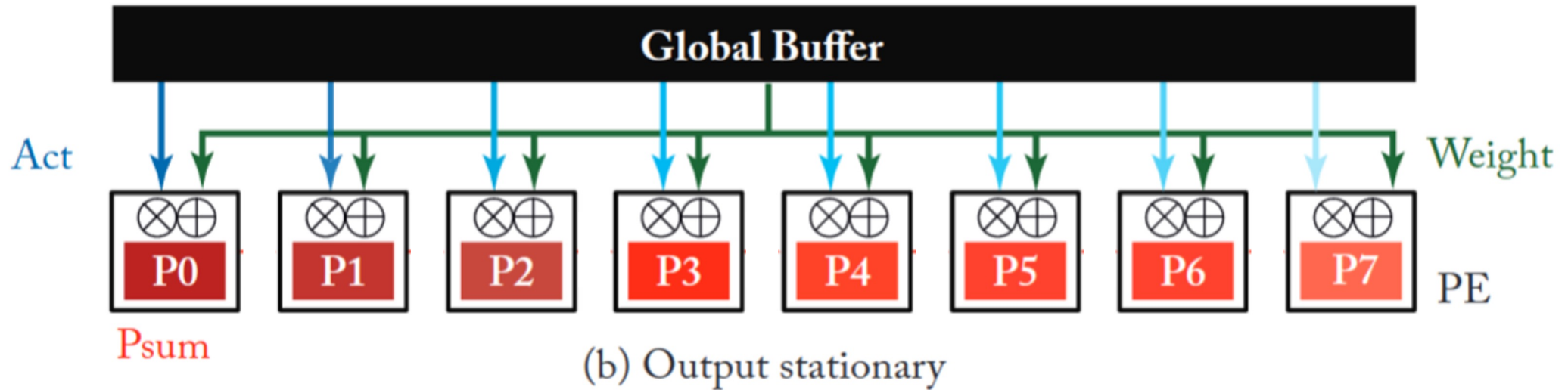
What to Keep Stationary? Inputs



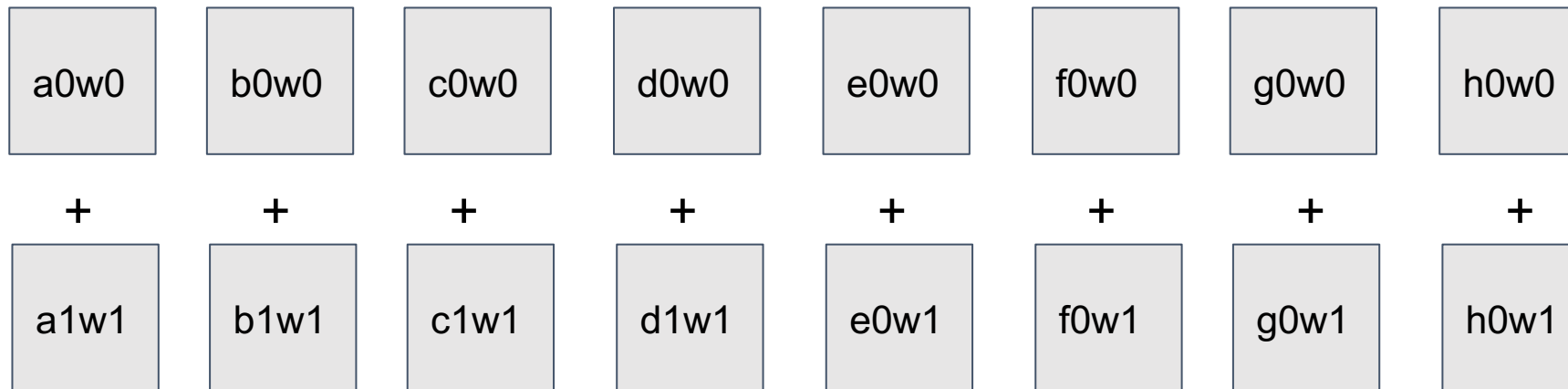
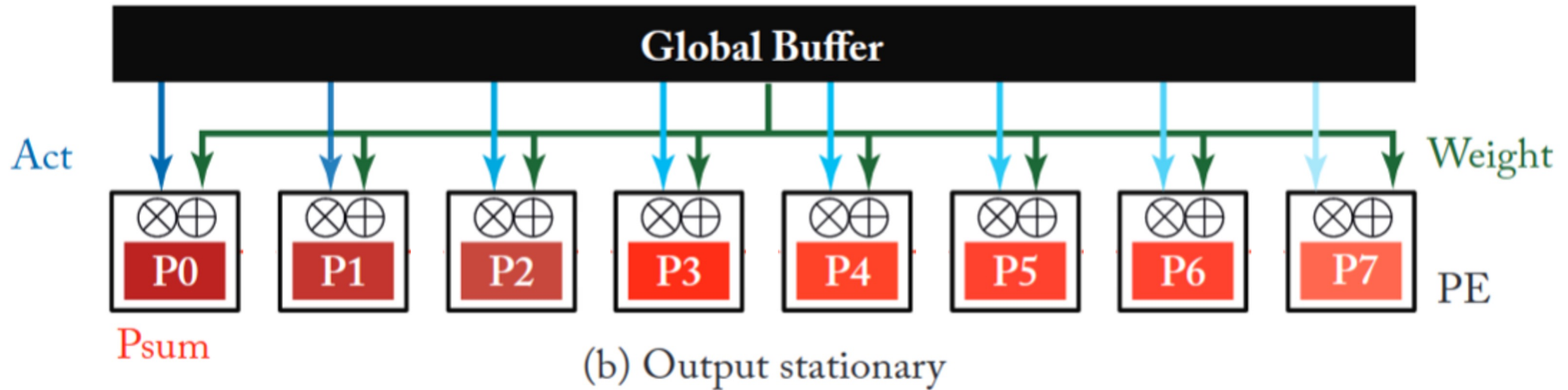
(c) Input stationary

$$a_0w_0 + a_1w_1 + a_2w_2 + a_3w_3 + a_4w_4 + a_5w_5 + a_6w_6 + a_7w_7$$

What to Keep Stationary? Outputs

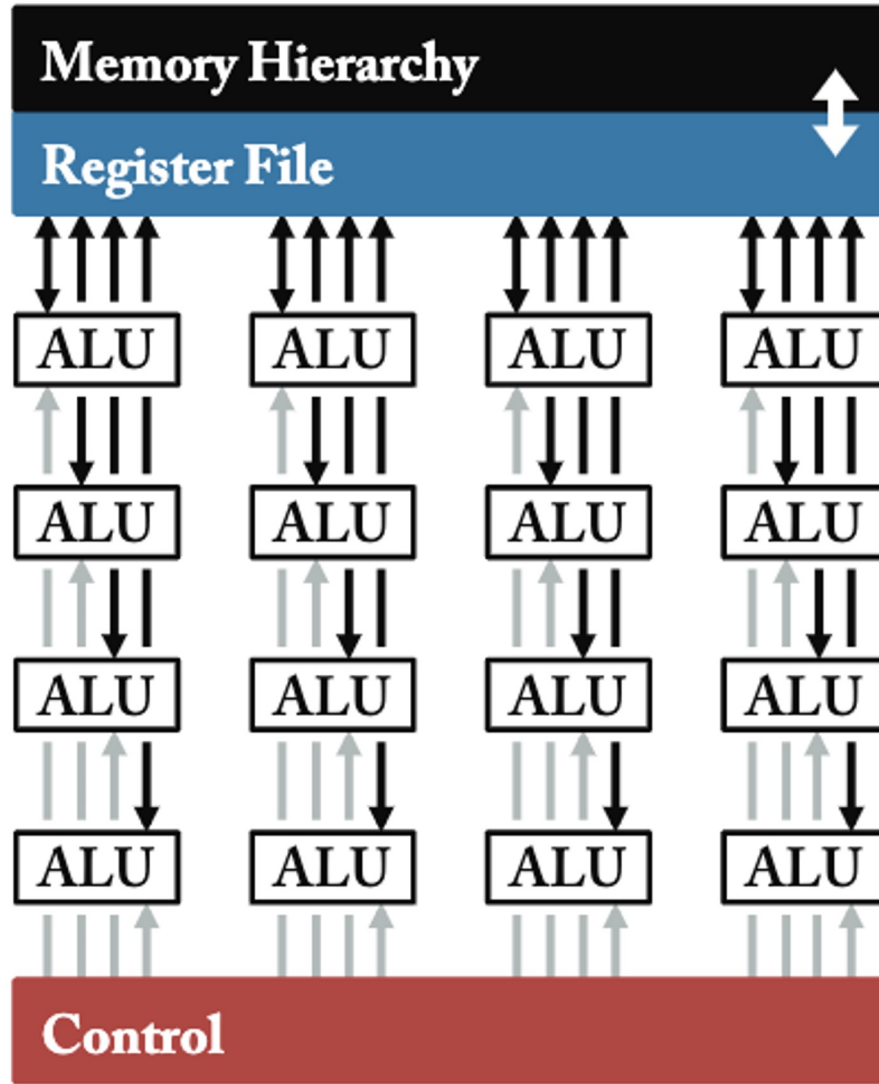


What to Keep Stationary? Outputs

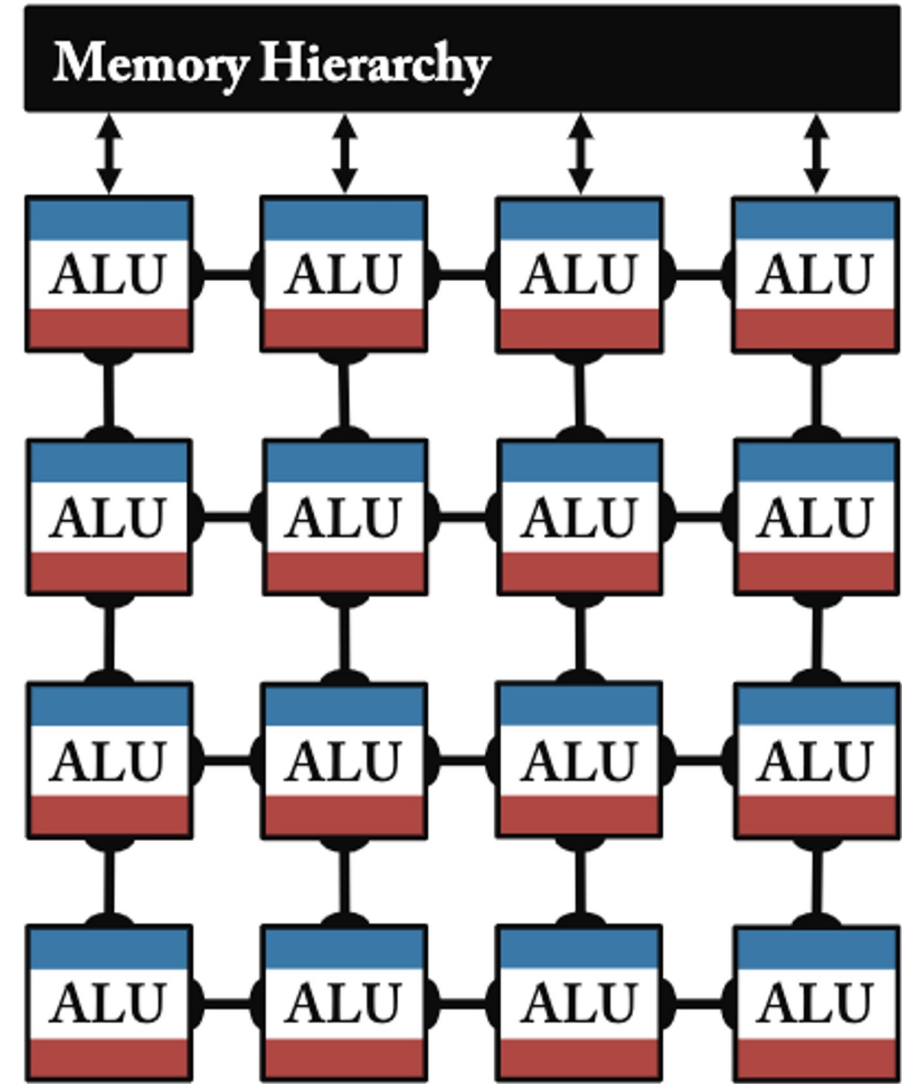


Highly Parallel Architectures

Temporal Architecture (SIMD/SIMT)

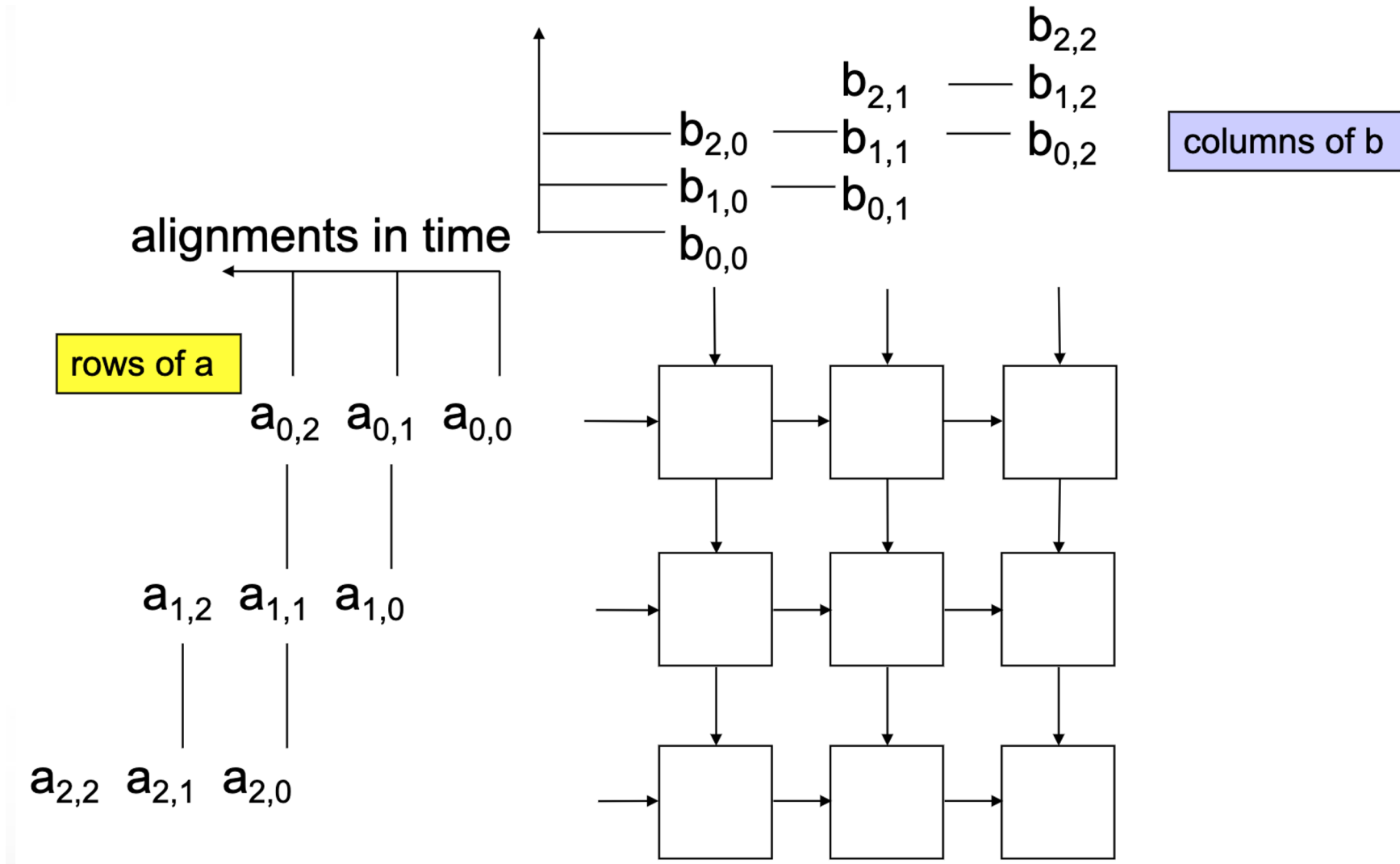


Spatial Architecture (Dataflow Processing)

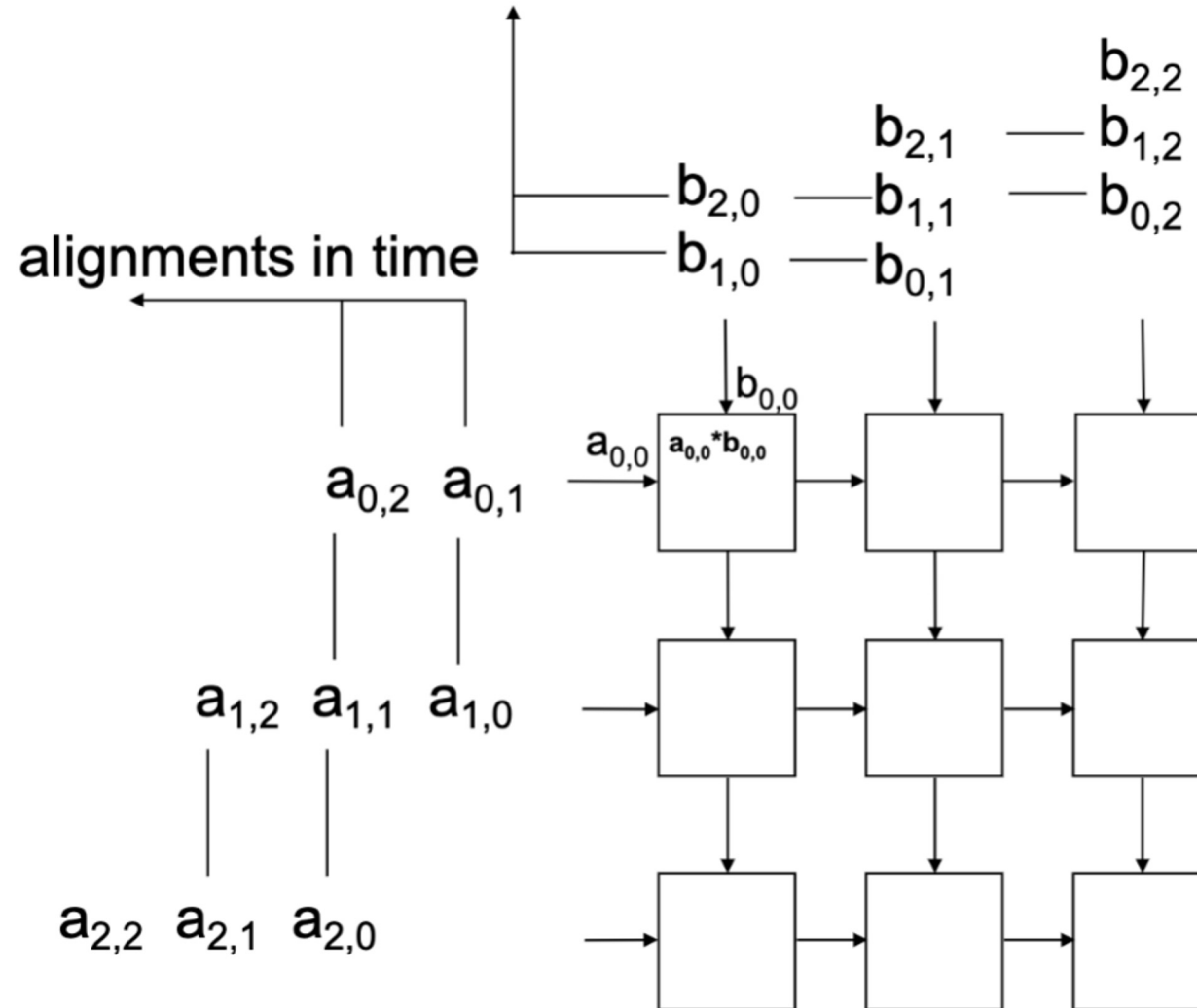


Source: V. Sze et al. "Efficient Processing of Deep Neural Networks" 2020

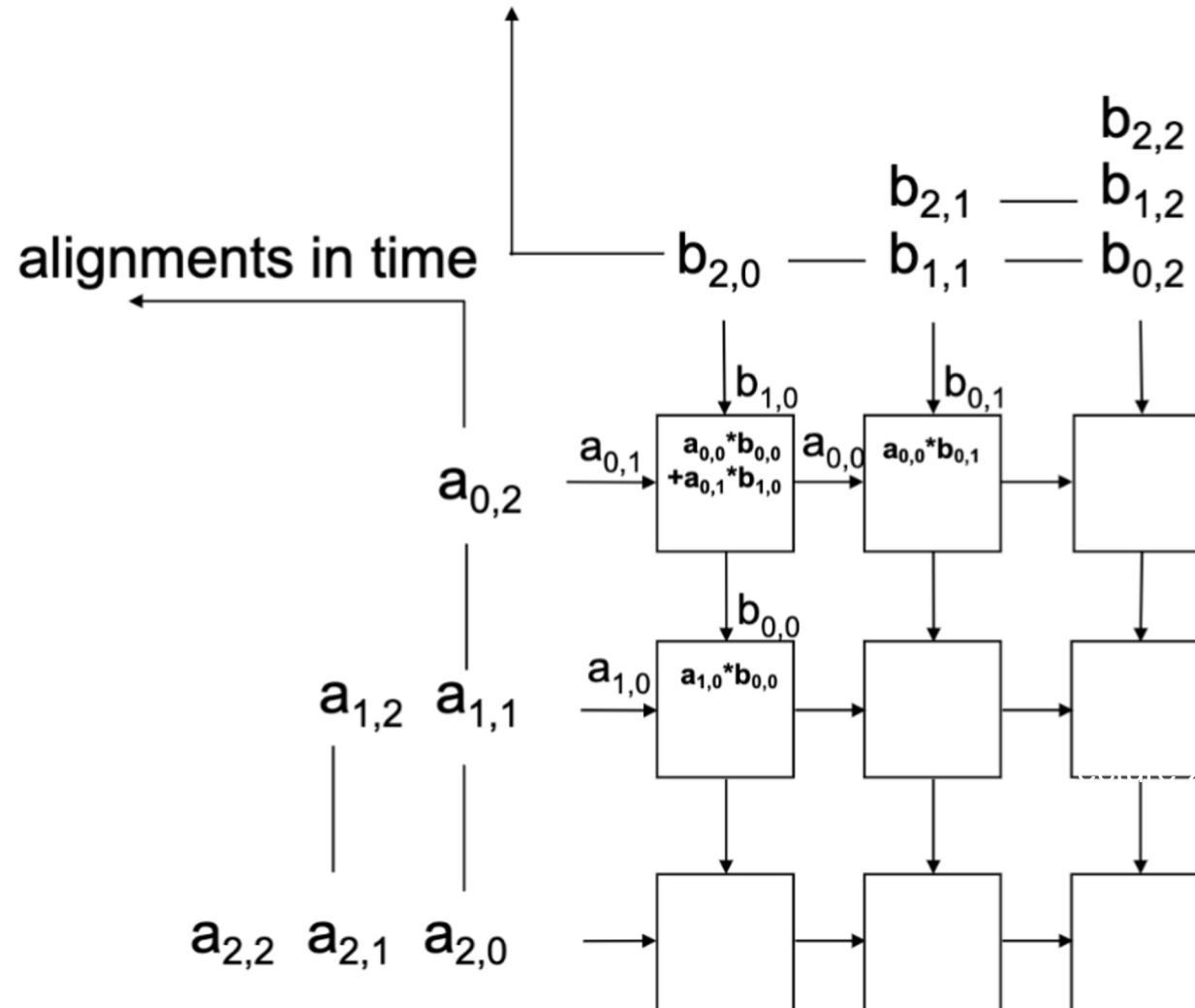
Systolic Array: Matrix Multiply Example



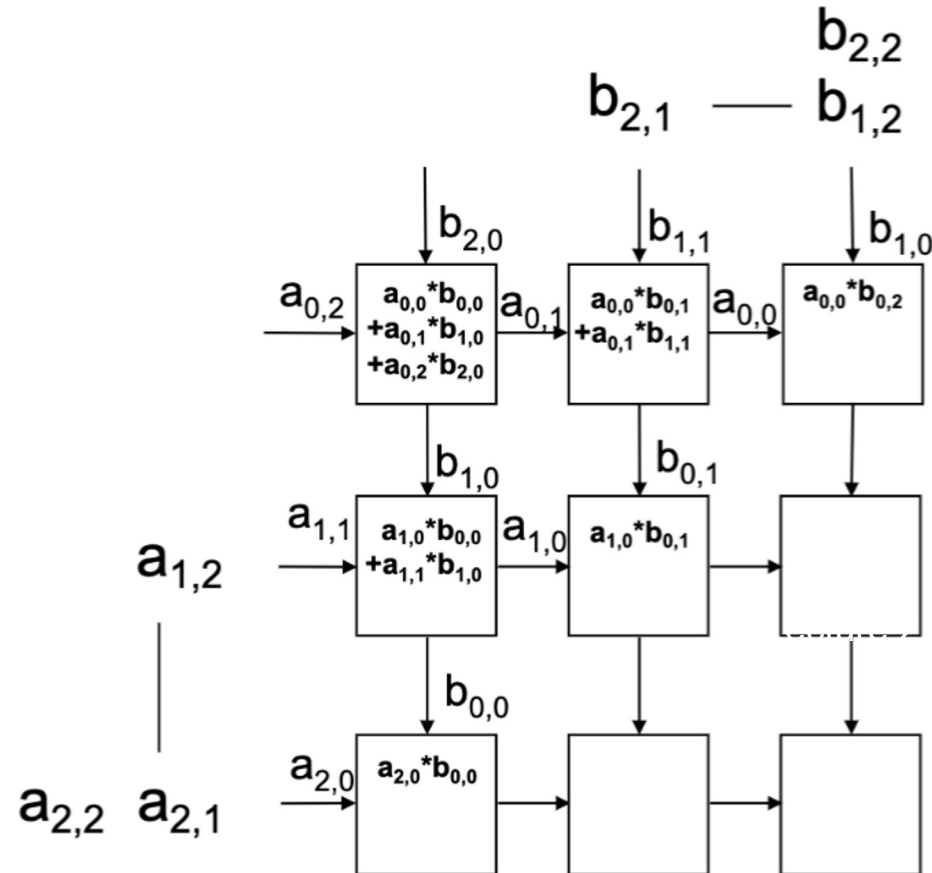
Systolic Array: Matrix Multiply Example



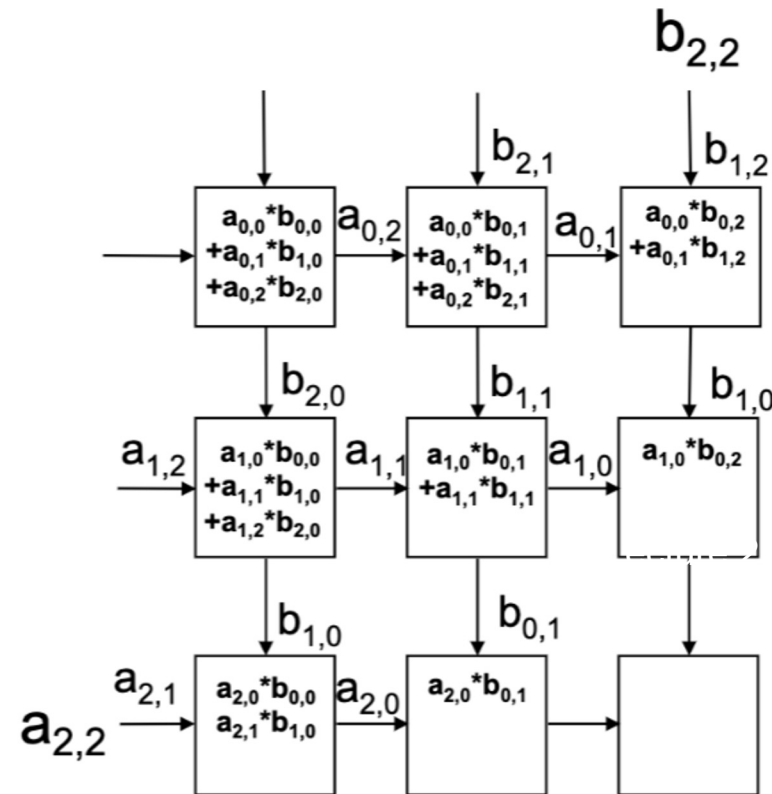
Systolic Array: Matrix Multiply Example



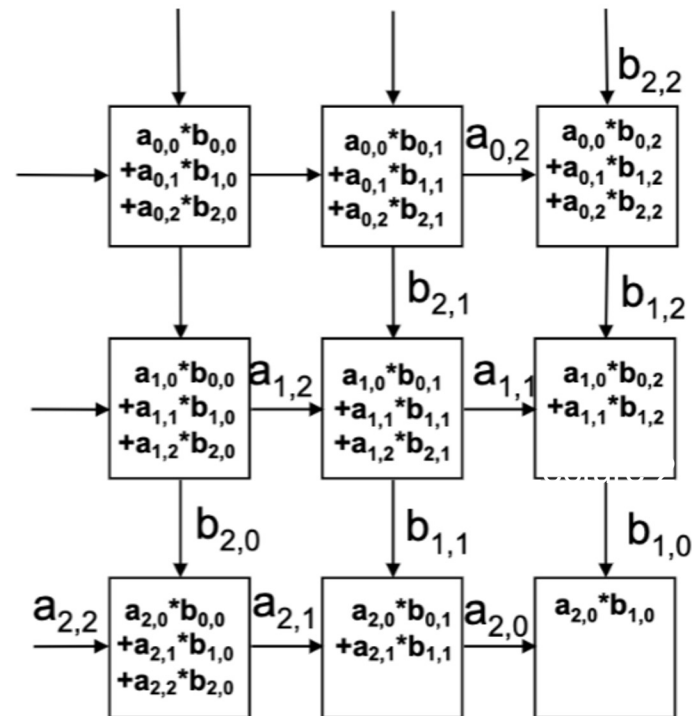
Systolic Array: Matrix Multiply Example



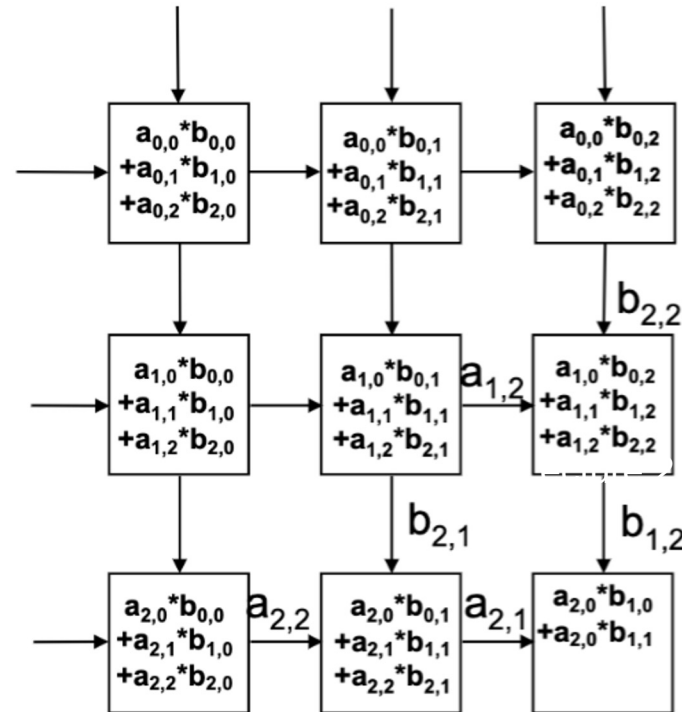
Systolic Array: Matrix Multiply Example



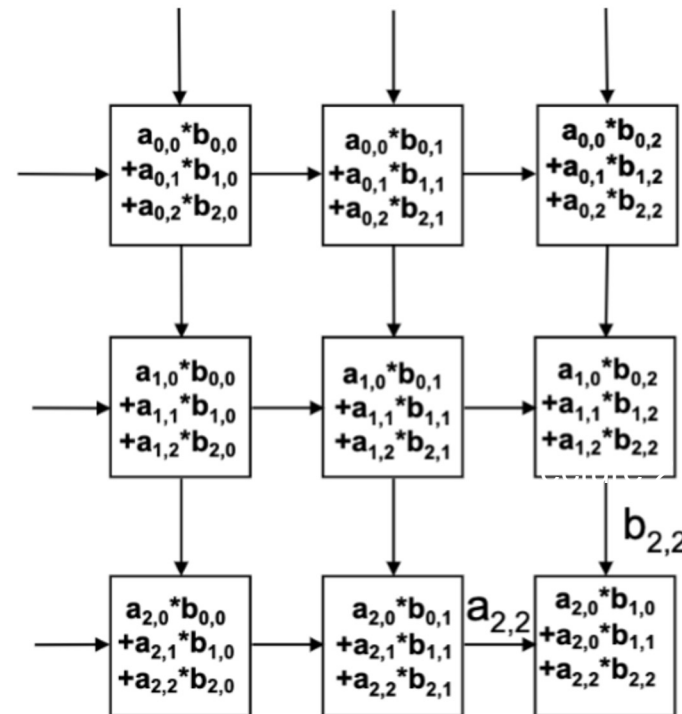
Systolic Array: Matrix Multiply Example



Systolic Array: Matrix Multiply Example



Systolic Array: Matrix Multiply Example



_____ stationary?

2. Hardware Efficiency

1. Arithmetic

- Specialized Instructions: To amortize overhead. ✓
- Lower precision (Quantization) ✓

2. Memory

- Locality: Move data to inexpensive on-chip memory. ✓
- Reuse: To avoid expensive memory fetches. ✓ ✓

3. Ineffectual Operations

- Sparsity: Skip useless operations
- Compressed Sparse Column (CSC) Format

Kinds of Sparsity

Activation

5	0	1	2
3	1	0	1
0	8	4	4
9	0	0	1

Activation Sparsity



Sparse activation functions (e.g. ReLU)

 \times

Weight

2	0	1	2
-4	-1	3	0
0	0	3	2
0	0	-5	7

Weight Sparsity

Block Sparsity

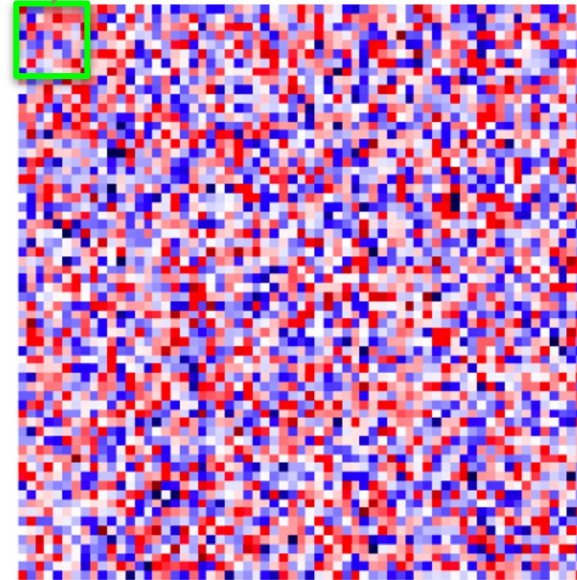


Pruning (covered in later lectures)

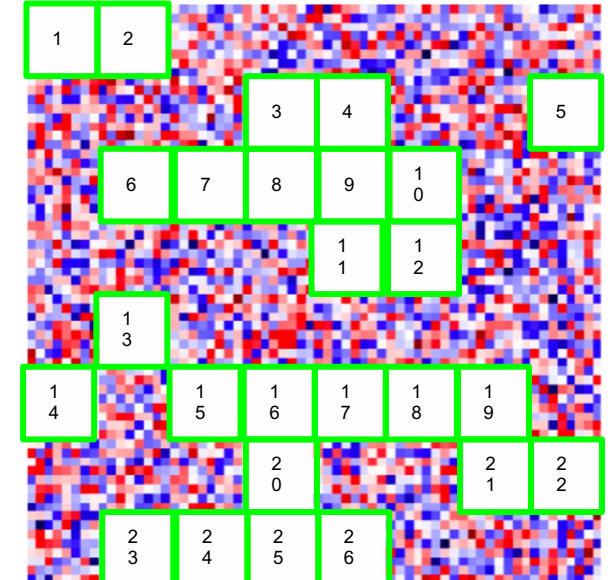
Coarse-grained “Block” Sparsity

- All DNN accelerators are parallel
 - Multiple MACs/cycle
- The smallest unit of computation that can be skipped is a large block (*recall [amortized overhead](#)*)
- Example:
 - Systolic array with 64 MACs/cycle
 - 8x8 pattern
 - 64x64 matrix = 4096 MACs
 - Total # cycles = 64 cycles
 - Block sparsity pattern needs to skip blocks of 8x8
 - Speedup = $64/(64-26) = 1.7X$ faster

64 MACs/cycle



Dense weights

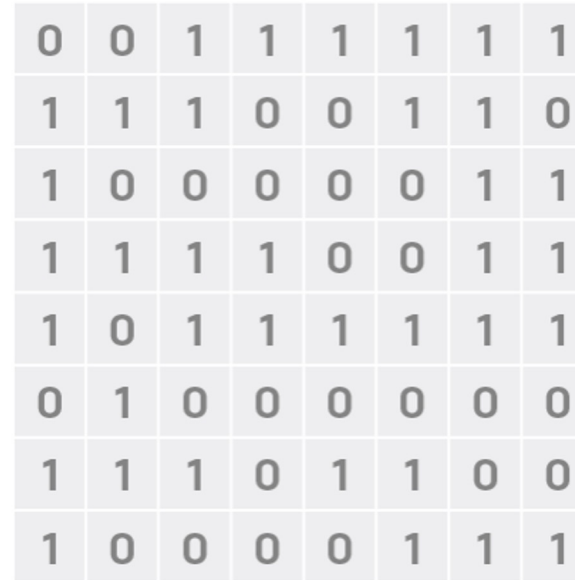


Block-sparse weights

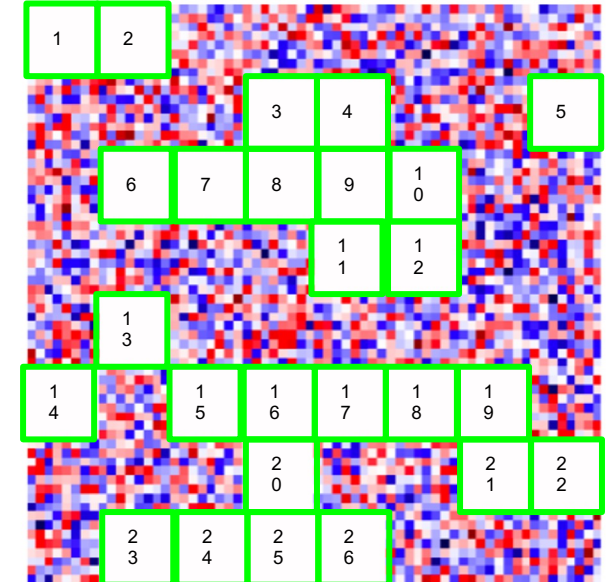
Coarse-grained “Block” Sparsity

Source: Open AI

- All DNN accelerators are parallel
 - Multiple MACs/cycle
- The smallest unit of computation that can be skipped is a large block (*recall [amortized overhead](#)*)
- Example:
 - Systolic array with 64 MACs/cycle
 - 8x8 pattern
 - 64x64 matrix = 4096 MACs
 - Total # cycles = 64 cycles
 - Block sparsity pattern needs to skip blocks of 8x8
 - Speedup = $64/(64-26) = 1.7X$ faster



Corresponding sparsity pattern



Block-sparse weights

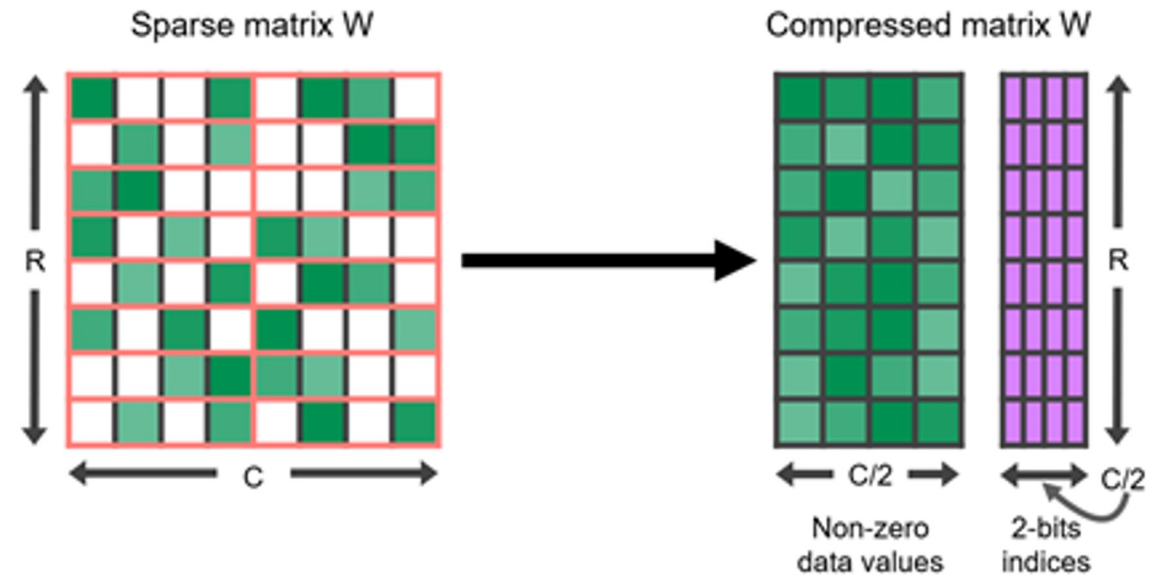
Simplest way to leverage sparsity with low overhead

⇒ Single bit per 8x8 block ($1/64 = 1.6\%$ overhead)

⇒ Simple control logic because entire block is skipped

Fine-grained Sparsity in Ampere GPUs

- Very recently, fine-grained sparsity was added to Tensor Cores on Nvidia GPUs
- 2 elements for every block of 4 elements can be zero
- Requires retraining to regain accuracy
- Overhead?
 - 2 bits per 8-bit element
 - 12.5% memory overhead
 - Control logic? Performance improvement? Power savings?



QUESTION

What is the performance improvement of 50% fine-grained sparsity on Nvidia GPUs?

1. 2.0 X
2. 1.5 X
3. 1.2 X
4. 0.5 X

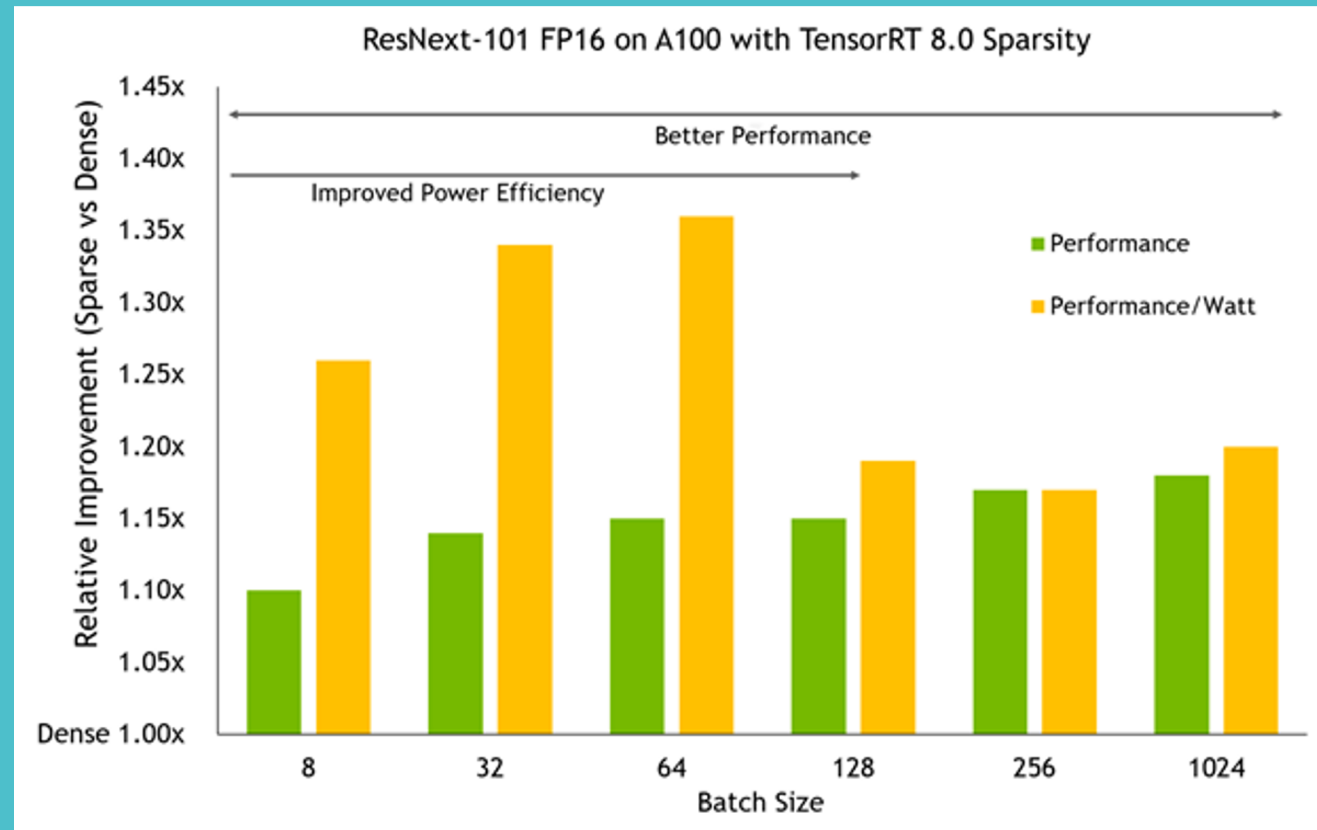


QUESTION

What is the performance improvement of 50% fine-grained sparsity on Nvidia GPUs?

1. 2.0 X
2. 1.5 X
3. 1.2 X
4. 0.5 X

Even though we skip half the computations, there is overhead to support sparsity, like figuring out where all the zeroes are to be able to skip



2. Hardware Efficiency

1. Arithmetic

- Specialized Instructions: To amortize overhead. ✓
- Lower precision (Quantization) ✓

2. Memory

- Locality: Move data to inexpensive on-chip memory. ✓
- Reuse: To avoid expensive memory fetches. ✓

3. Ineffectual Operations

- Sparsity: Skip useless operations ✓
- Compressed Sparse Column (CSC) Format

Compressed Sparse Column (CSC) Format

Virtual Weight	a	b	c	d	e	f	g	h	i	j	k	l	m
Relative Row Index	0	1	0	1	0	2	0	0	0	2	0	2	0
Column Pointer	0	3	4	6	6	8	10	11	13				

Compressed Sparse Column (CSC) Format

Virtual Weight	a	b	c	d	e	f	g	h	i	j	k	l	m
Relative Row Index	0	1	0	1	0	2	0	0	0	2	0	2	0
Column Pointer	0	3	4	6	6	8	10	11	13				

a	0	e	0	g	i	k	0
0	d	0	0	h	0	0	0
b	0	0	0	0	0	0	l
c	0	f	0	0	j	0	m

Interleaving

$$\vec{a} \begin{pmatrix} 0 & 0 & a_2 & 0 & a_4 & a_5 & 0 & a_7 \end{pmatrix} \times \begin{matrix} PE0 \\ PE1 \\ PE2 \\ PE3 \end{matrix} \begin{pmatrix} w_{0,0} & 0 & w_{0,2} & 0 & w_{0,4} & w_{0,5} & w_{0,6} & 0 \\ 0 & w_{1,1} & 0 & w_{1,3} & 0 & 0 & w_{1,6} & 0 \\ 0 & 0 & w_{2,2} & 0 & w_{2,4} & 0 & 0 & w_{2,7} \\ 0 & w_{3,1} & 0 & 0 & 0 & w_{0,5} & 0 & 0 \\ 0 & w_{4,1} & 0 & 0 & w_{4,4} & 0 & 0 & 0 \\ 0 & 0 & 0 & w_{5,4} & 0 & 0 & 0 & w_{5,7} \\ 0 & 0 & 0 & 0 & w_{6,4} & 0 & w_{6,6} & 0 \\ w_{7,0} & 0 & 0 & w_{7,4} & 0 & 0 & w_{7,7} & 0 \\ w_{8,0} & 0 & 0 & 0 & 0 & 0 & 0 & w_{8,7} \\ w_{9,0} & 0 & 0 & 0 & 0 & 0 & w_{9,6} & w_{9,7} \\ 0 & 0 & 0 & 0 & w_{10,4} & 0 & 0 & 0 \\ 0 & 0 & w_{11,2} & 0 & 0 & 0 & 0 & w_{11,7} \\ w_{12,0} & 0 & w_{12,2} & 0 & 0 & w_{12,5} & 0 & w_{12,7} \\ w_{13,0} & w_{13,2} & 0 & 0 & 0 & 0 & w_{13,6} & 0 \\ 0 & 0 & w_{14,2} & w_{14,3} & w_{14,4} & w_{14,5} & 0 & 0 \\ 0 & 0 & w_{15,2} & w_{15,3} & 0 & w_{15,5} & 0 & 0 \end{pmatrix} = \begin{pmatrix} b_0 \\ b_1 \\ -b_2 \\ b_3 \\ -b_4 \\ b_5 \\ b_6 \\ -b_7 \\ -b_8 \\ -b_9 \\ b_{10} \\ -b_{11} \\ -b_{12} \\ b_{13} \\ b_{14} \\ -b_{15} \end{pmatrix} \xRightarrow{ReLU} \begin{pmatrix} b_0 \\ b_1 \\ 0 \\ b_3 \\ 0 \\ b_5 \\ b_6 \\ 0 \\ 0 \\ 0 \\ b_{10} \\ 0 \\ 0 \\ b_{13} \\ b_{14} \\ 0 \end{pmatrix}$$

2. Hardware Efficiency

1. Arithmetic

- Specialized Instructions: To amortize overhead. ✓
- Lower precision (Quantization) ✓

2. Memory

- Locality: Move data to inexpensive on-chip memory. ✓
- Reuse: To avoid expensive memory fetches. ✓

3. Ineffectual Operations

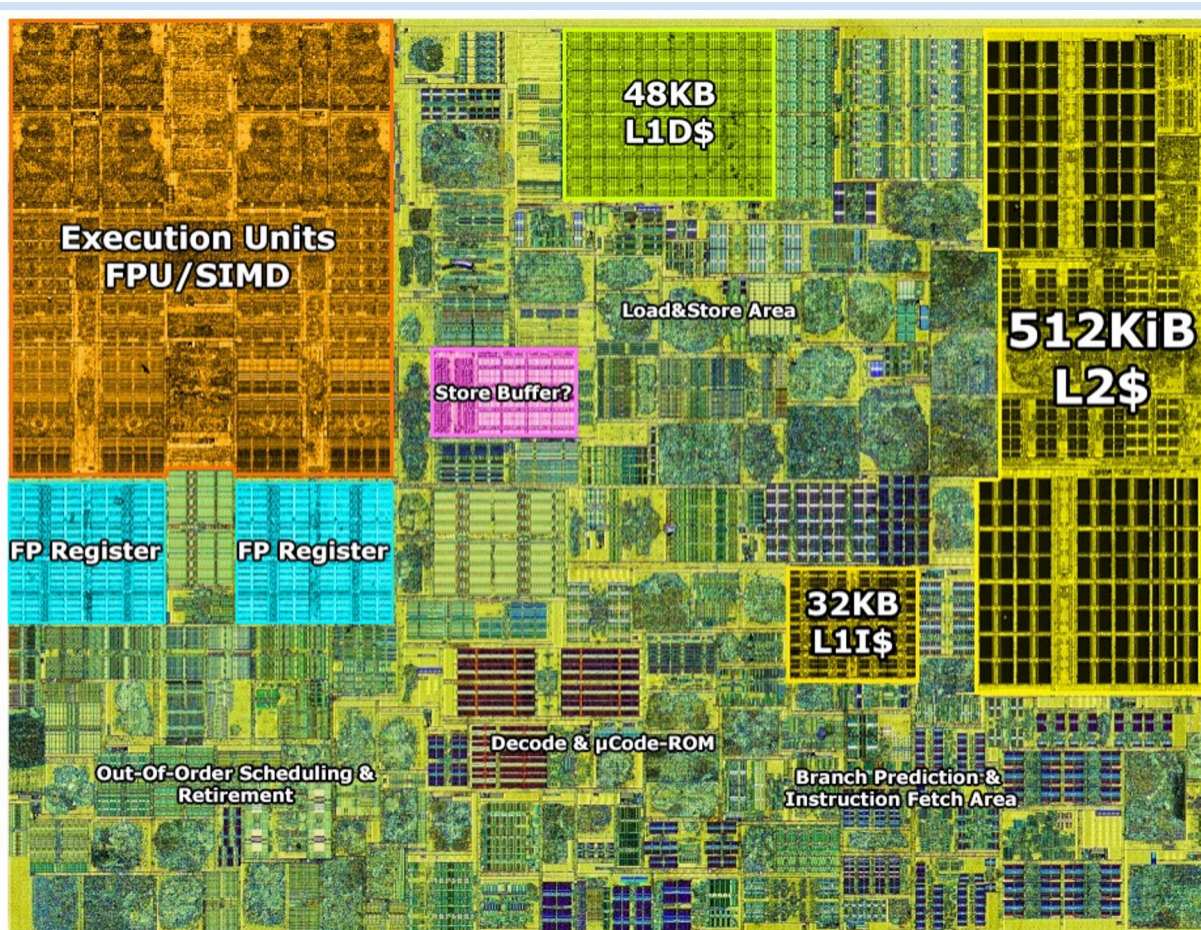
- Sparsity: Skip useless operations ✓
- Compressed Sparse Column (CSC) Format ✓

Guidelines Domain-Specific Architectures (DSAs)

1. Use dedicated memory to minimize the distance over which data is moved.
2. Invest the resources saved from dropping advanced microarchitectural optimizations into more arithmetic units or bigger memories.
3. Use the easiest form of parallelism that matches the domain.
4. Reduce data size and type to the simplest needed for the domain.
5. Use a domain-specific programming language to port code to the DSA.



Putting it all together



On-Chip Memory

MAC Array

Other Stuff
- Control Logic
- Memory Controllers
- etc.

Layer-Specific
Arithmetic Units