



Principles of Machine Learning Systems

4: Hardware Acceleration

Roadmap for Today

- HW enabling Deep Learning
- Performance Metrics
- Where does Energy Go?
- Hardware Efficiency Options
- Hardware Case Studies



Mohamed
Abdelfattah

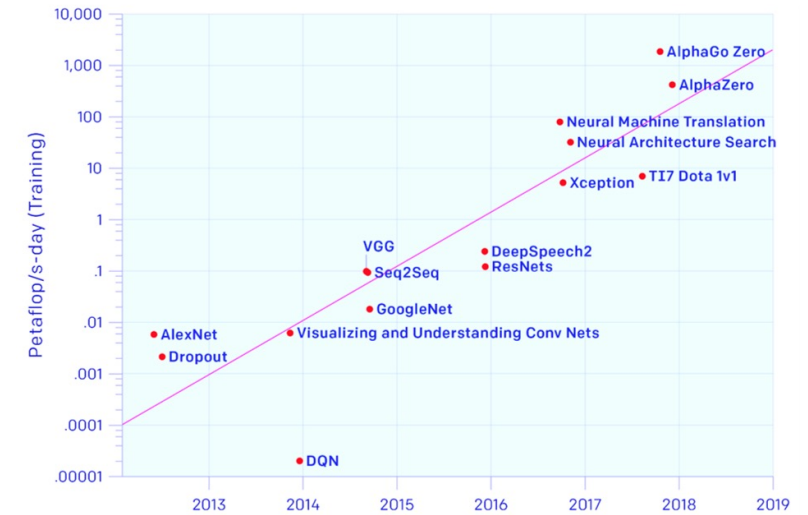
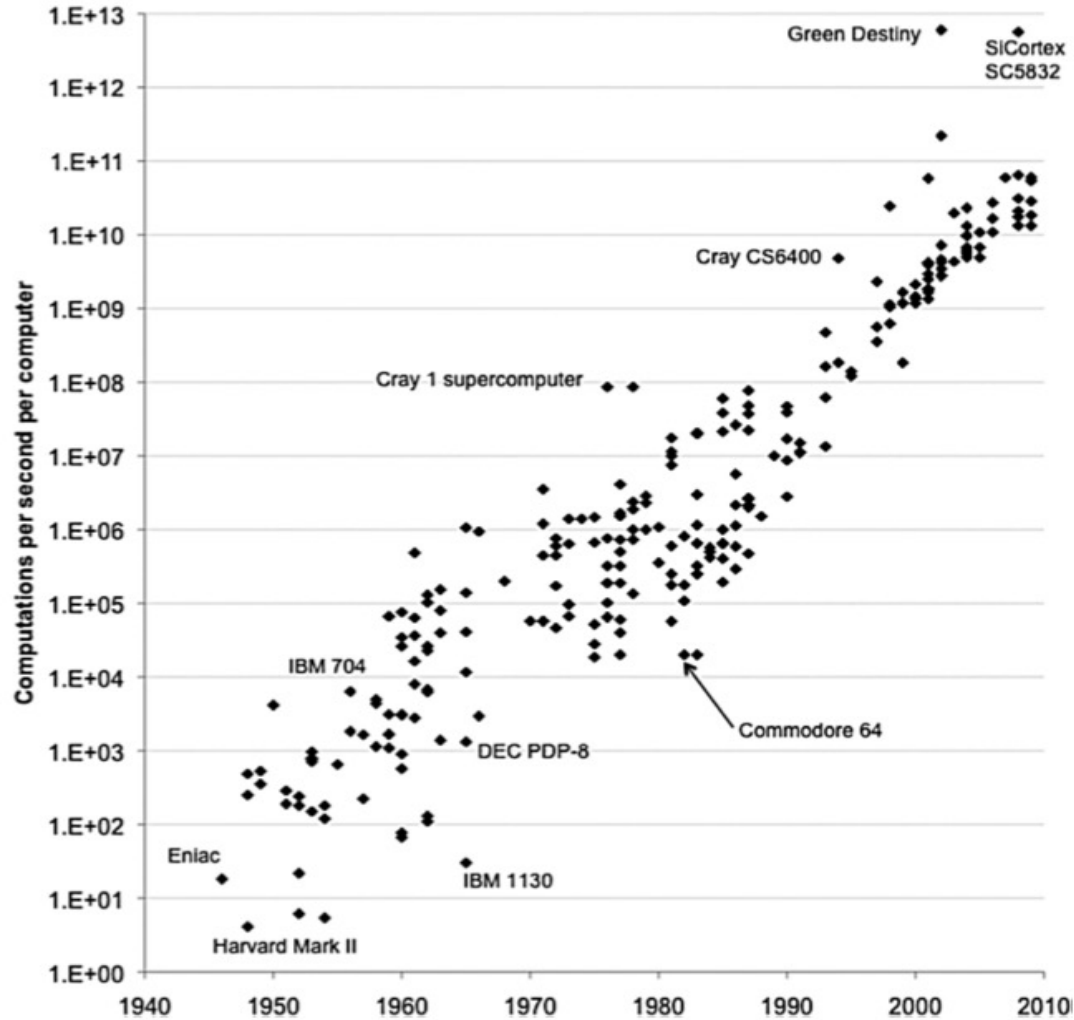


Roadmap for Today

- **HW enabling Deep Learning**
- Performance Metrics
- Where does Energy Go?
- Hardware Efficiency Options
- Hardware Case Studies



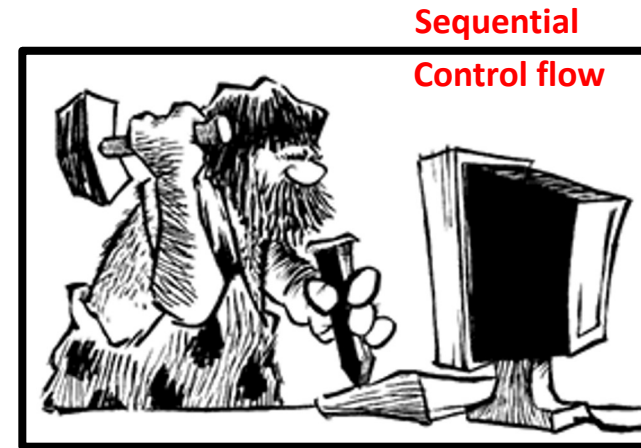
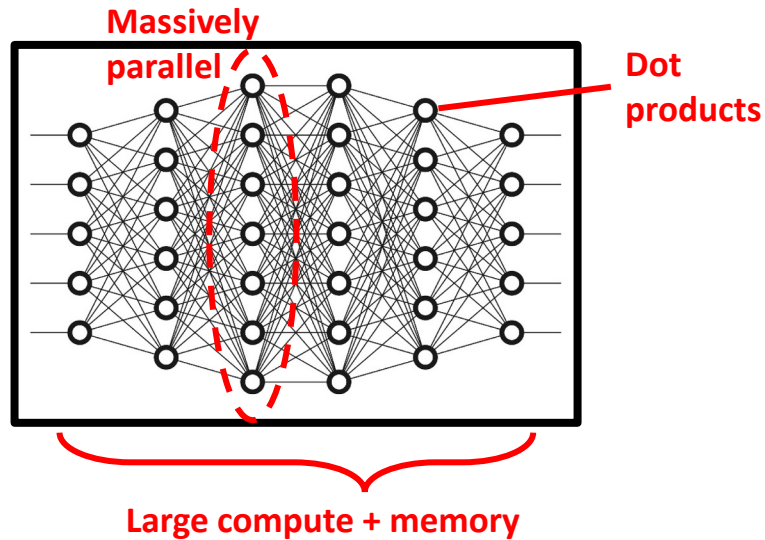
HW enables Deep Learning



Source: OpenAI



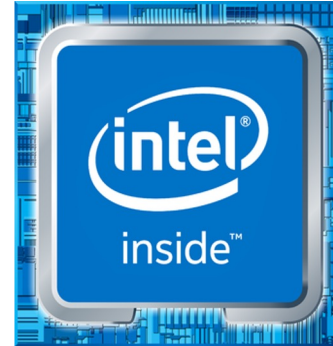
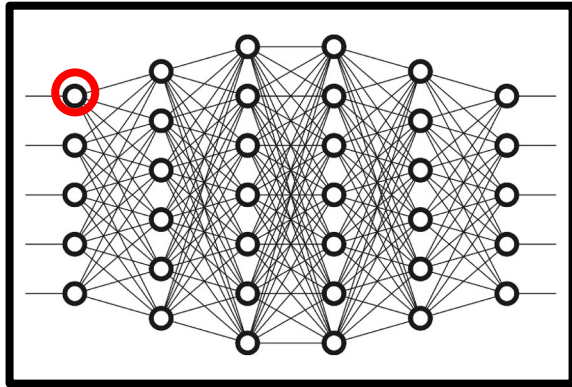
HW & Deep Learning Basics



- 1986: Backpropagation published



HW & Deep Learning Basics

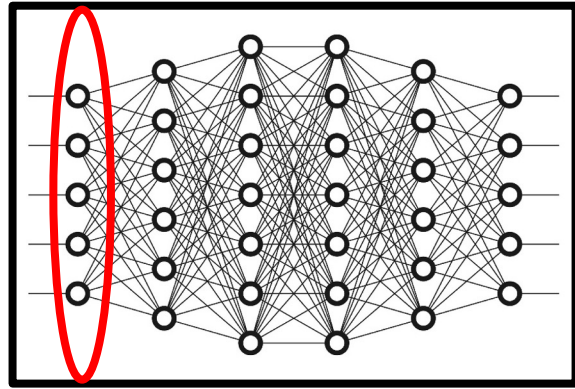


Sequential
Complex instruction set
Great for control flow

- 1986: Backpropagation published
- ~30 years of trying this on CPUs



HW & Deep Learning Basics

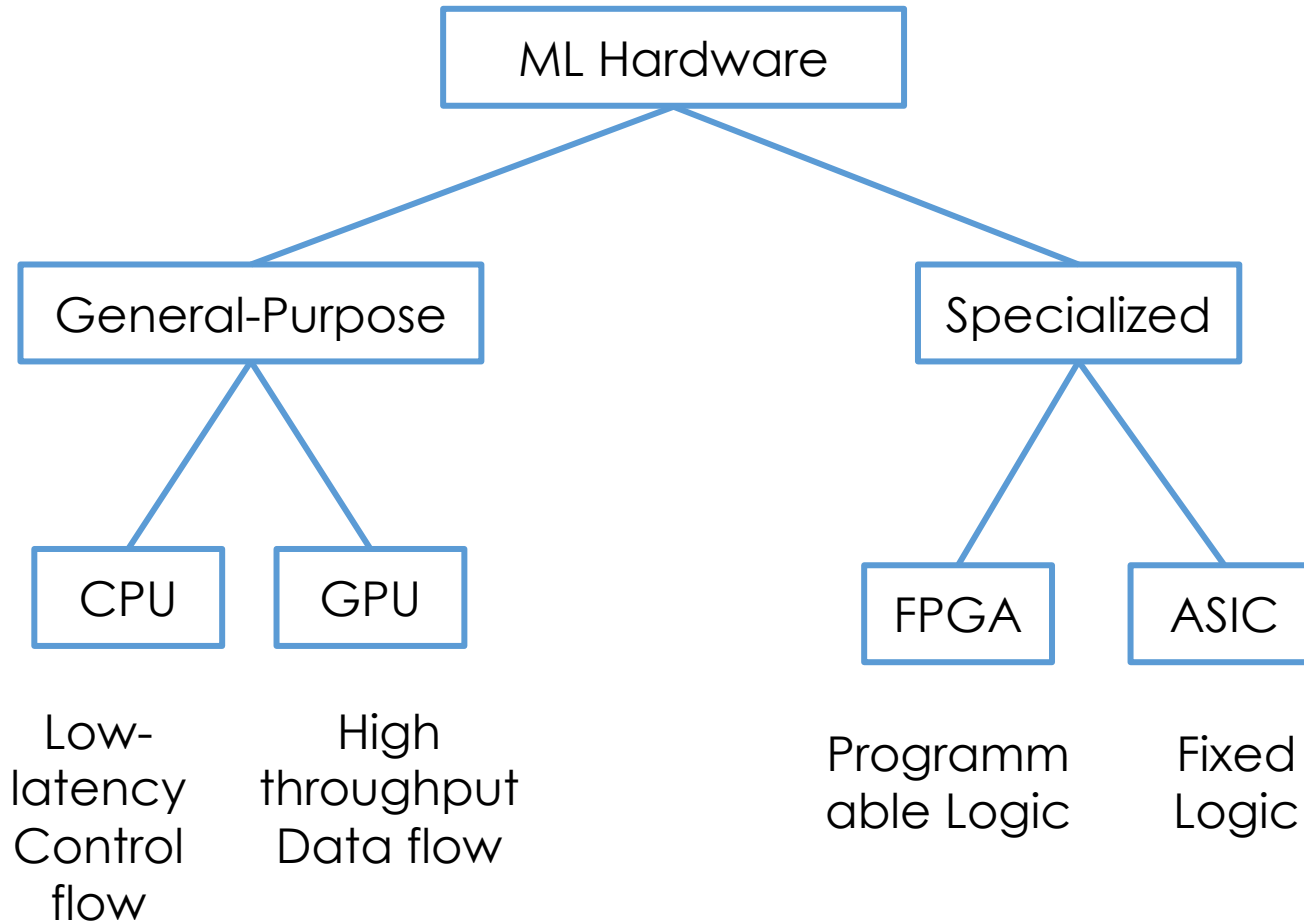


Parallel
Great for matrix math
Bad for control flow

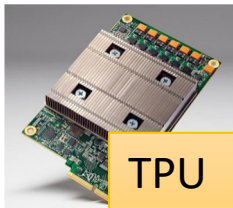
- 1986: Backpropagation published
- ~30 years of trying this on CPUs
- 2012: AlexNet paper \Rightarrow +10% accuracy \Rightarrow Deep learning explosion



Hardware types



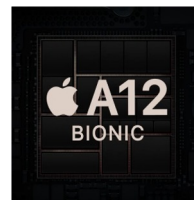
Specialized Hardware



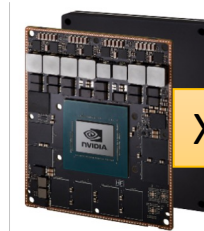
TPU



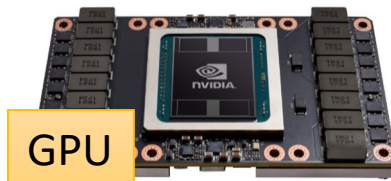
FPGA



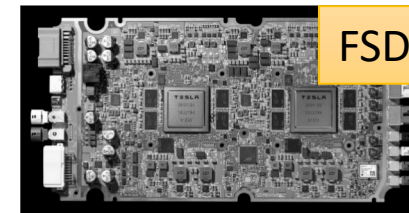
NPU



XAVIER



GPU



FSD

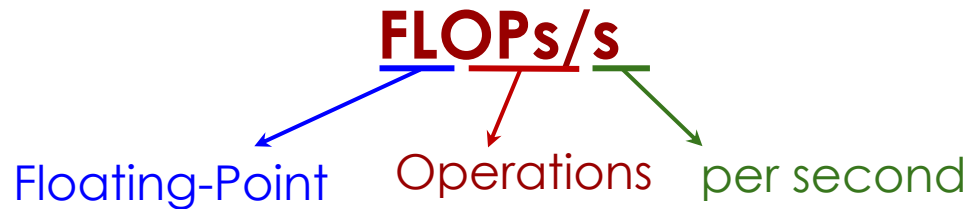


Roadmap for Today

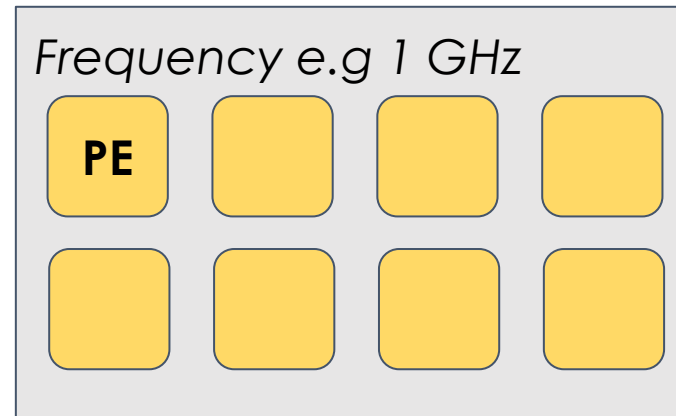
- HW enabling Deep Learning
- **Performance Metrics**
- Where does Energy Go?
- Hardware Efficiency Options
- Hardware Case Studies



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

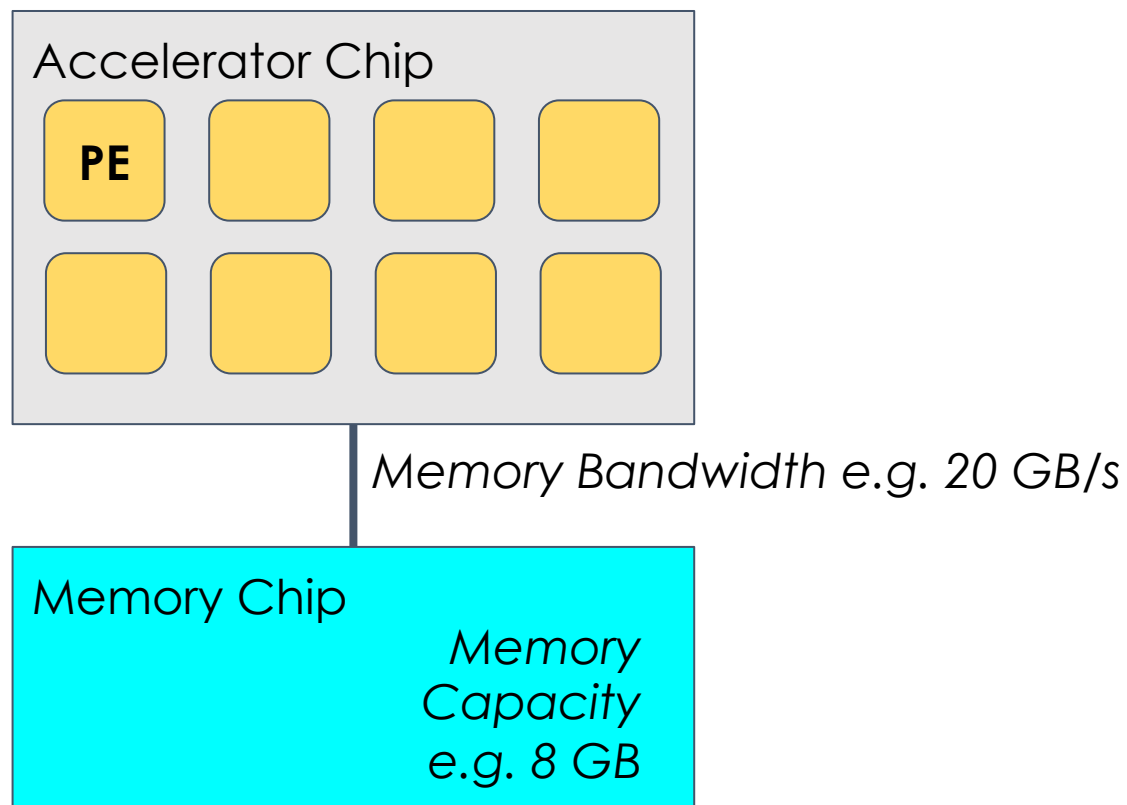


$$\frac{\text{operations}}{\text{second}} = \underbrace{\left(\frac{1}{\frac{\text{cycles}}{\text{operation}}} \times \frac{\text{cycles}}{\text{second}} \right)}_{\text{for a single PE}} \times \text{number of PEs}$$



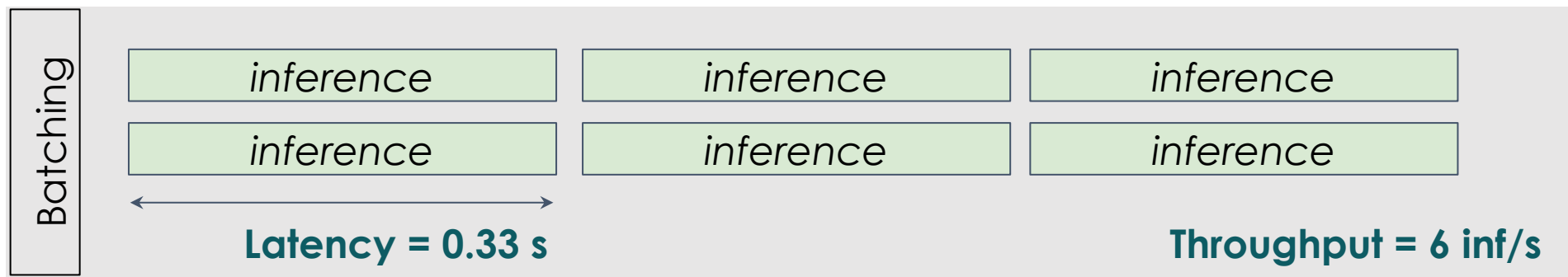
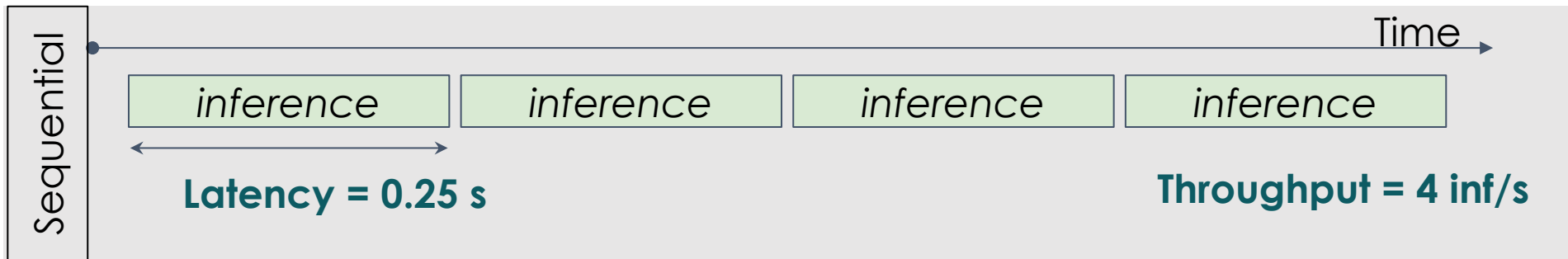
Memory Performance Metrics

- 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



DNN Performance

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



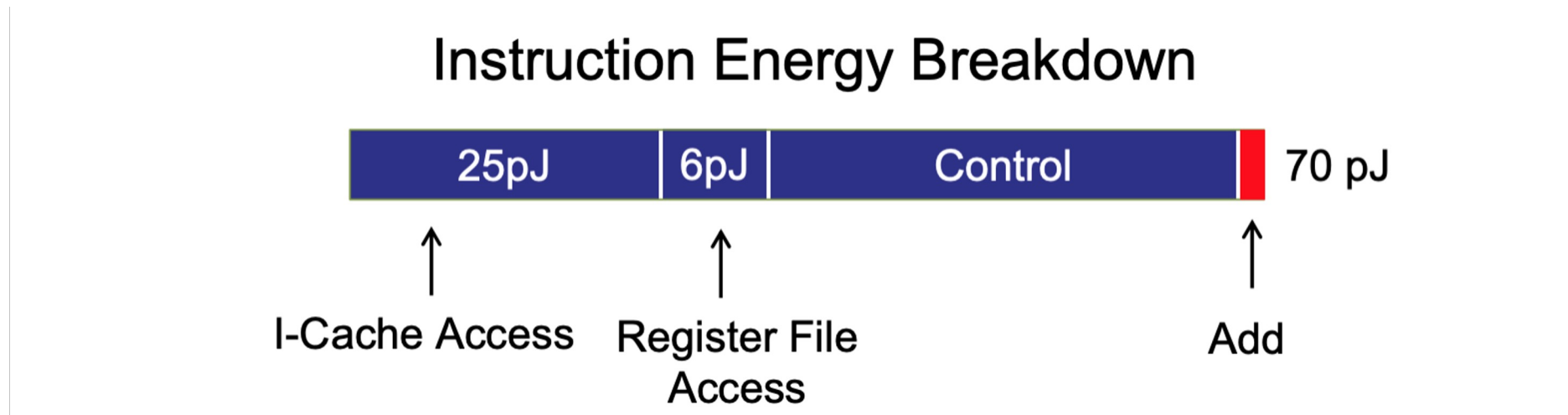
Roadmap for Today

- HW enabling Deep Learning
- Performance Metrics
- **Where does Energy Go?**
- Hardware Efficiency Options
- Hardware Case Studies



Where does the energy go?

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

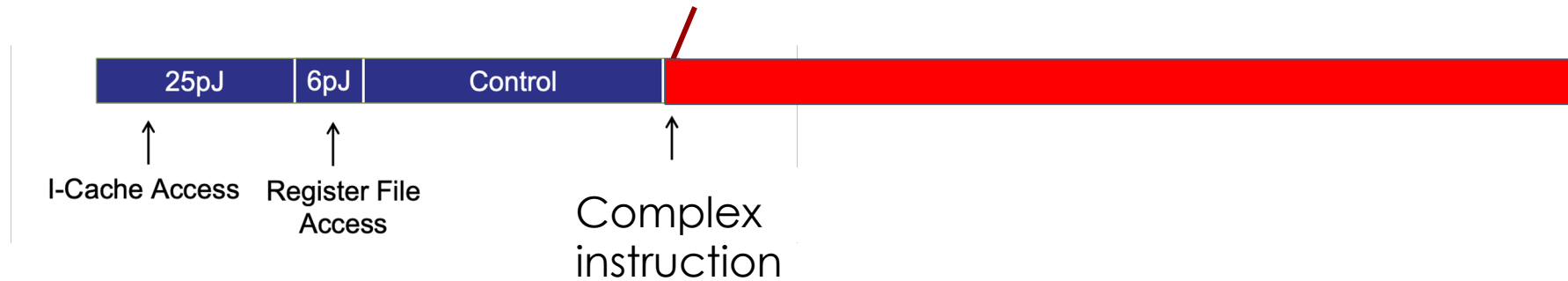


Source: Horowitz



Amortize Overhead

Increase Computation with same overhead



Half-precision Fused Multiply-Add

4-way dot-product

16x16 matrix multiplication

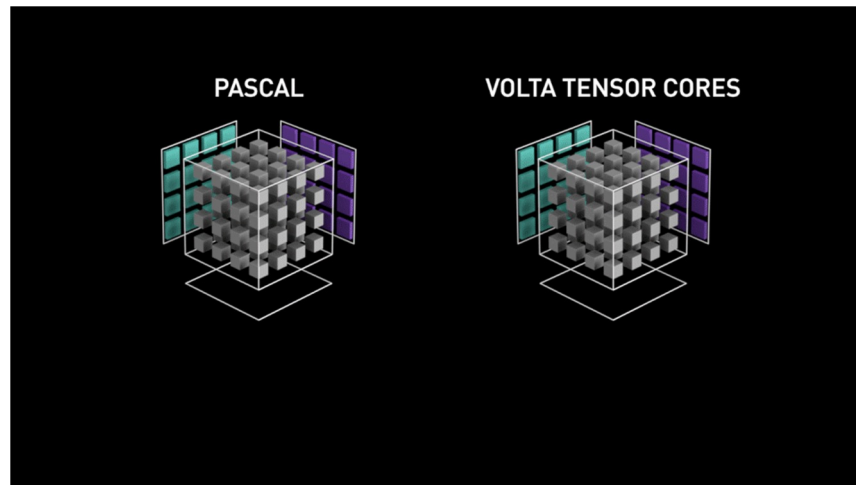
Operation	Energy**	Overhead*
HFMA	1.5pJ	2000%
HDP4A	6.0pJ	500%
HMMA	110pJ	27%

Source: Dally



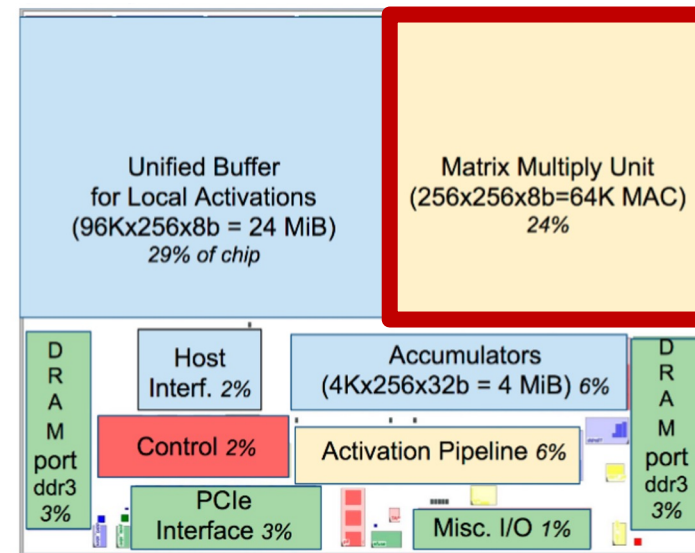
"Special" Instruction Examples

GPU



$16 \times 16 = 256^*$
MAC/cycle **~ 500 tensor cores per GPU*

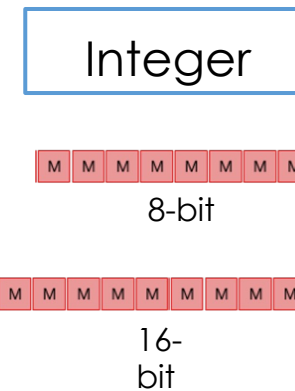
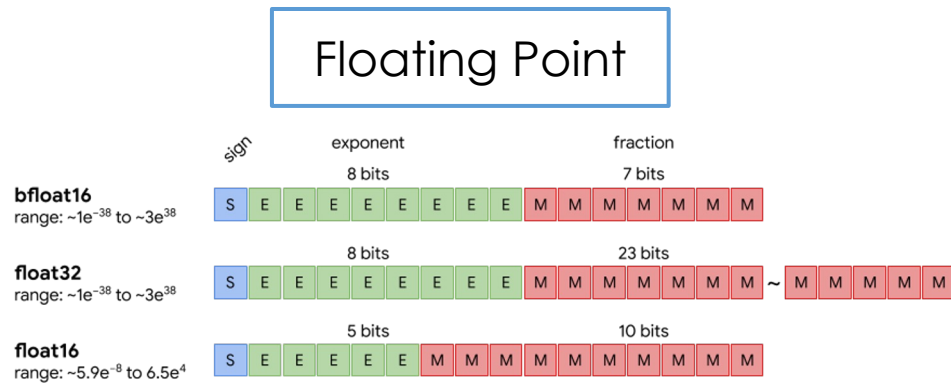
ASIC (TPUv1)



$256 \times 256 = 64$
kMAC/cycle



Numerical Format & Precision

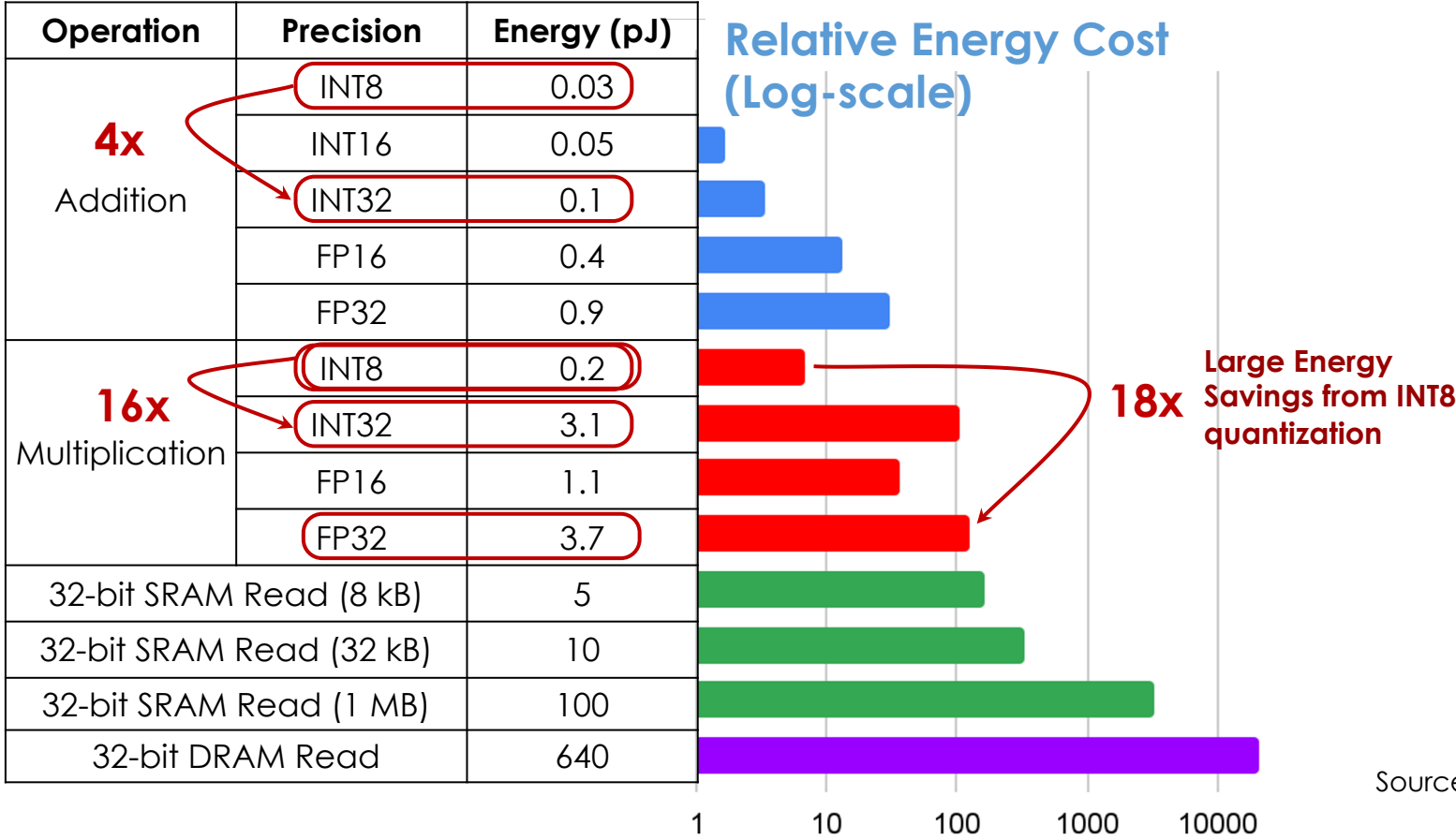


- IEEE standard includes FP32 and FP16
- Many exotic FP numbers in DNN
 - E.g. bfloat, minifloat

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



Cost of Arithmetic Operations



Source: Horowitz



Roadmap for Today

- HW enabling Deep Learning
- Performance Metrics
- Where does Energy Go?
- **Hardware Efficiency Options**
- Hardware Case Studies



Roadmap for Today

- HW enabling Deep Learning
- Performance Metrics
- Where does Energy Go?
- **Hardware Efficiency Options**
 - **Arithmetic**
 - **Memory**
 - Ineffectual Operation
- Hardware Case Studies

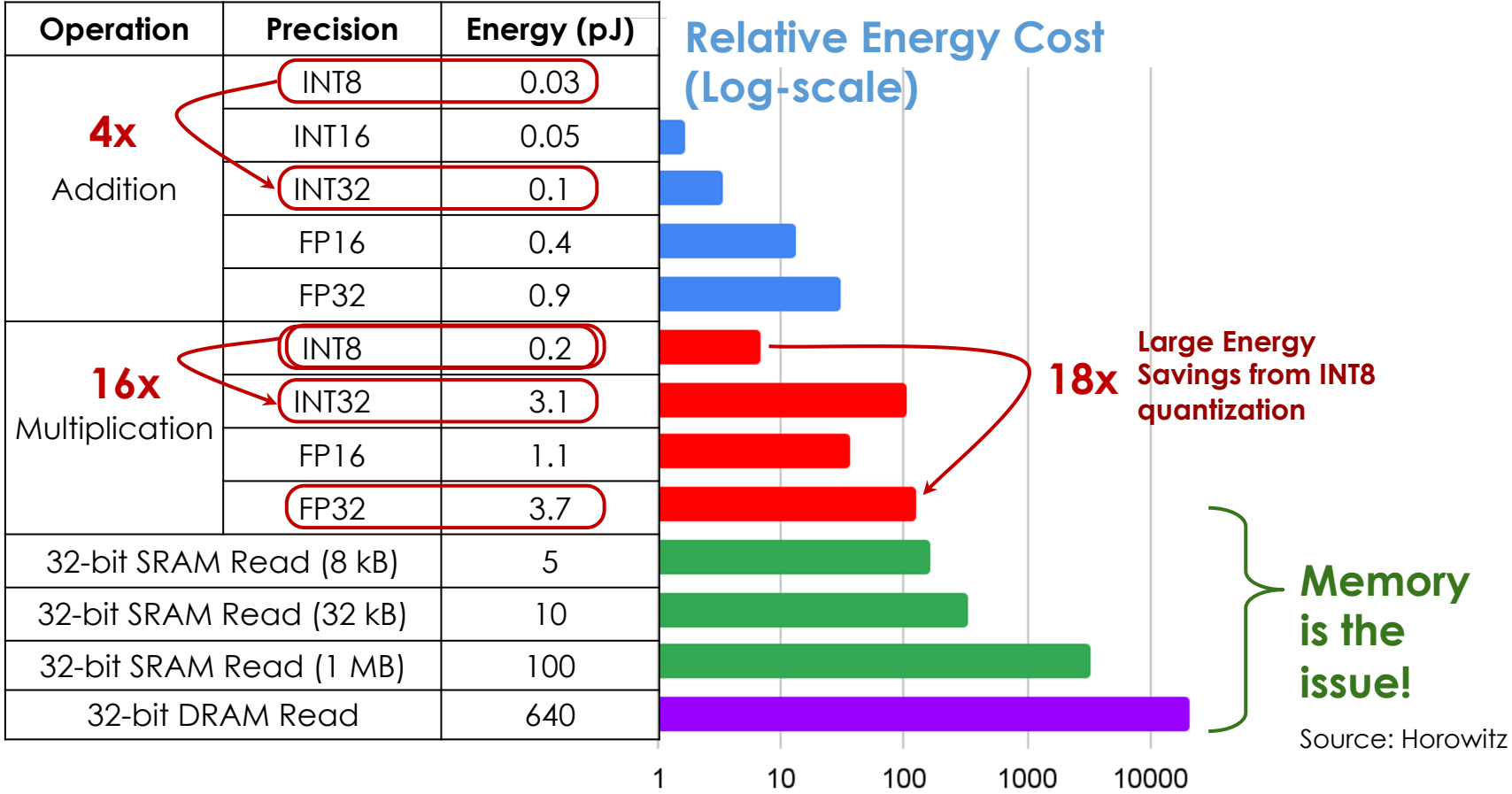


Roadmap for Today

- HW enabling Deep Learning
- Performance Metrics
- Where does Energy Go?
- **Hardware Efficiency Options**
 - Arithmetic
 - **Memory**
 - Ineffectual Operation
- Hardware Case Studies

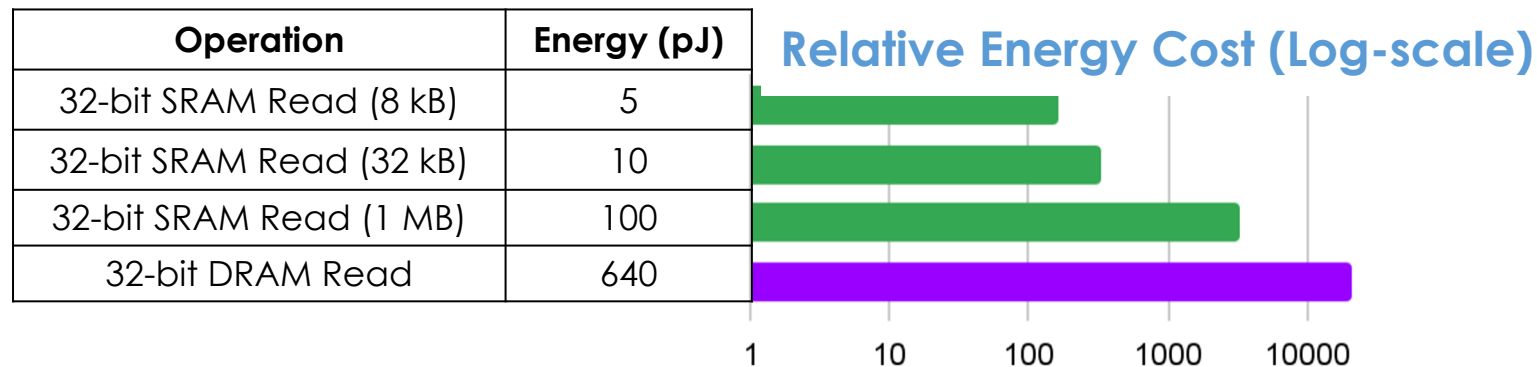


Memory is the bottleneck



Mem. Hierarchy Optimizations

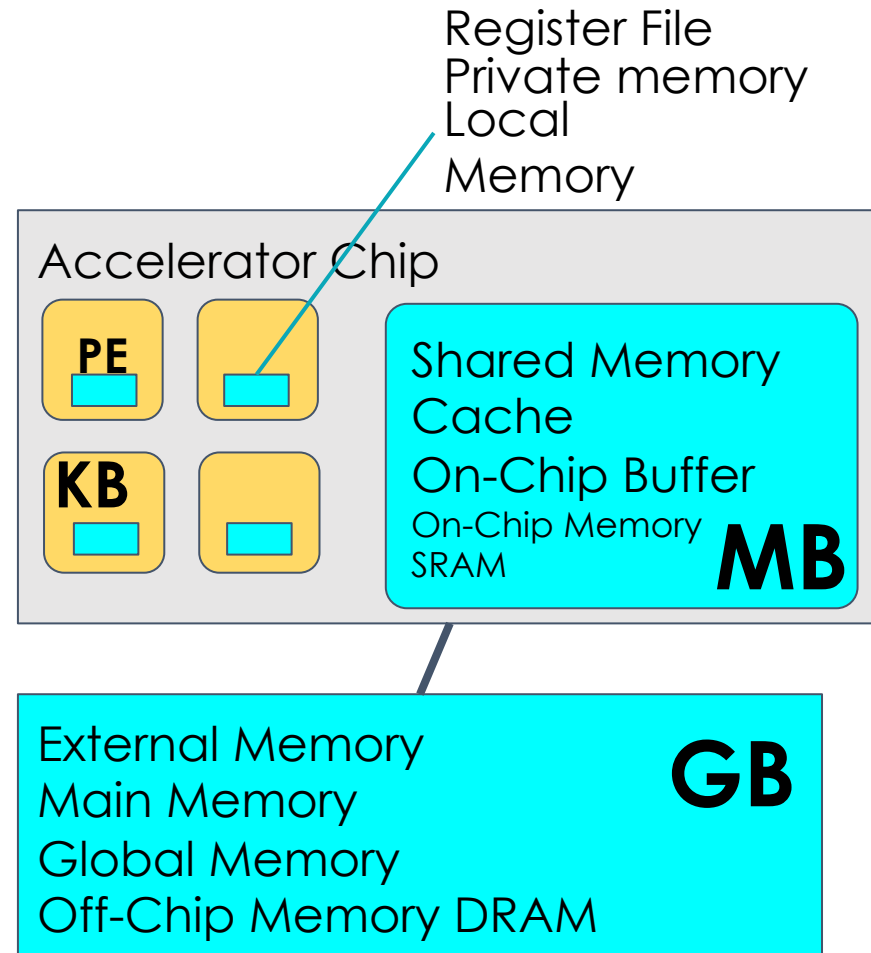
1. Get data close to the computation. (LOCALITY)
2. Once data is close - perform all computations with this data. (REUSE)



Memory Hierarchy

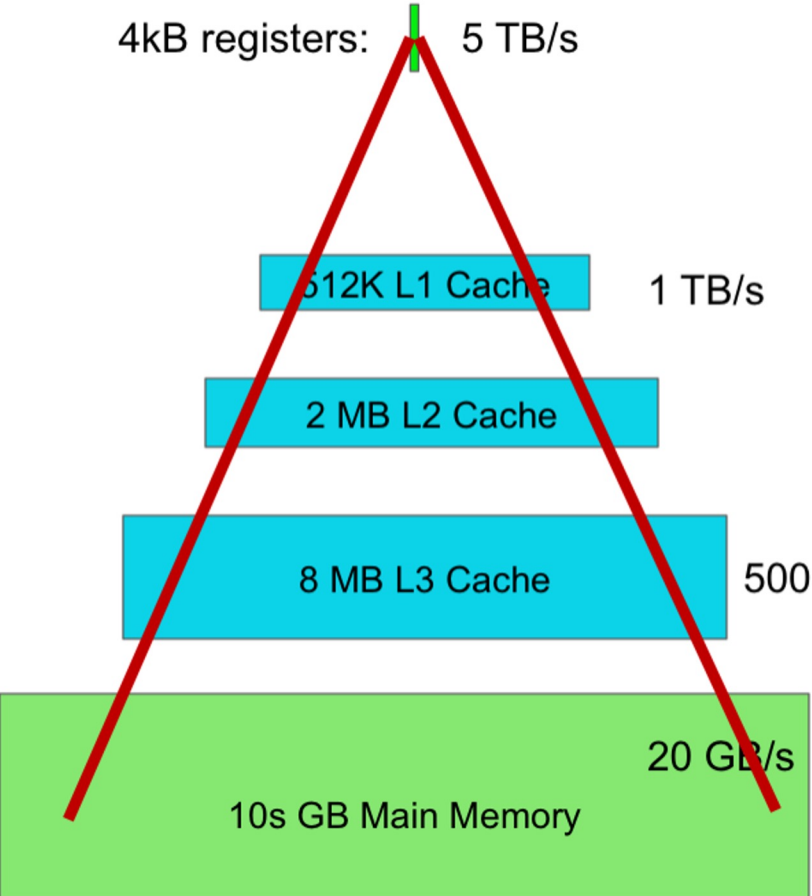
Why do we have a memory hierarchy?

- 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)

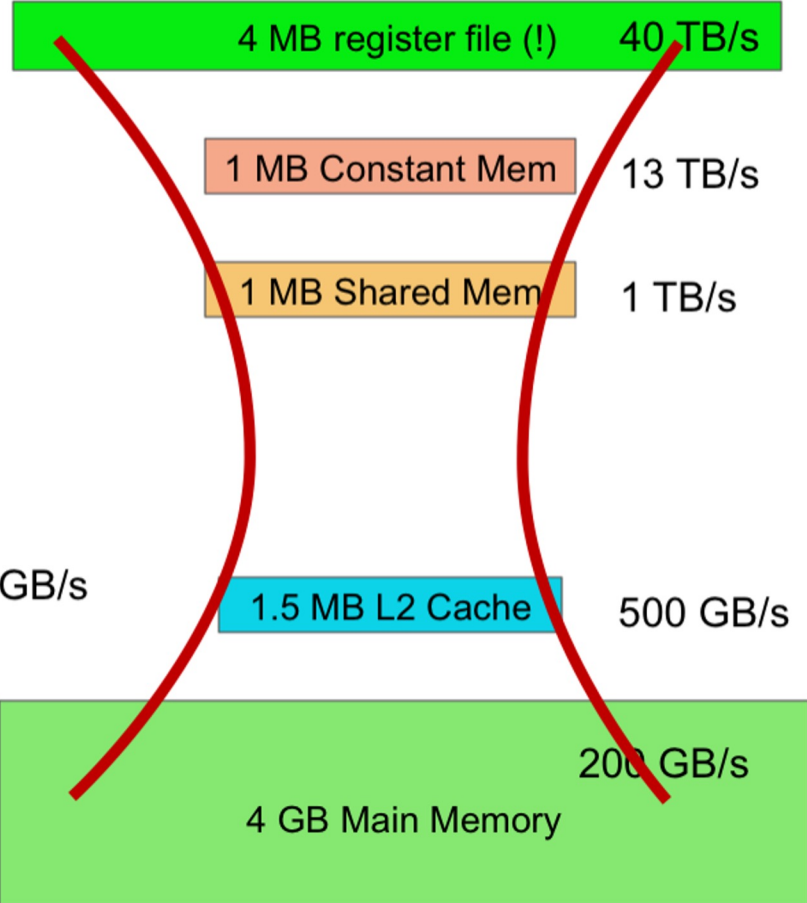


Memory Hierarchy Examples

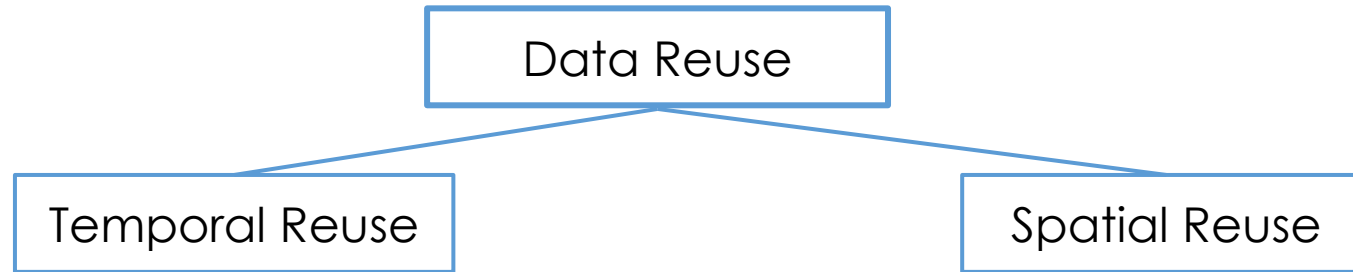
Intel® 8 core Sandy Bridge CPU



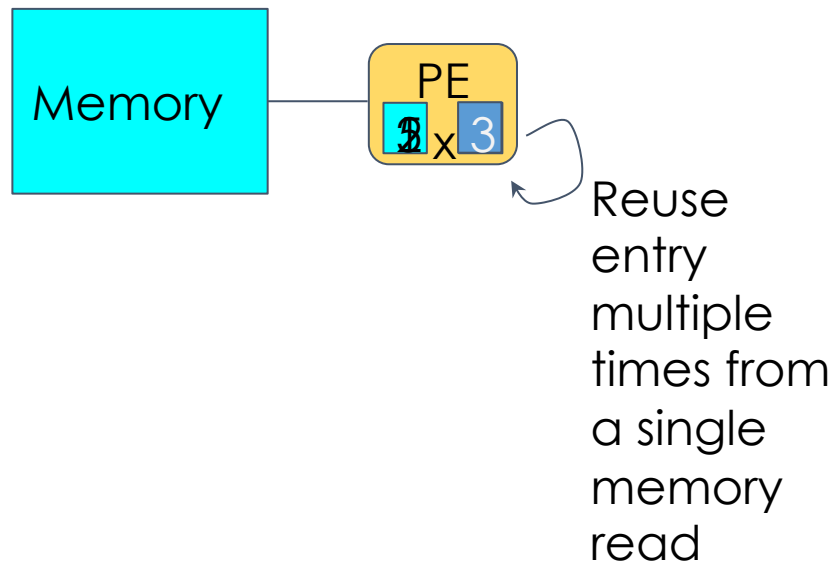
NVIDIA® GK110 GPU



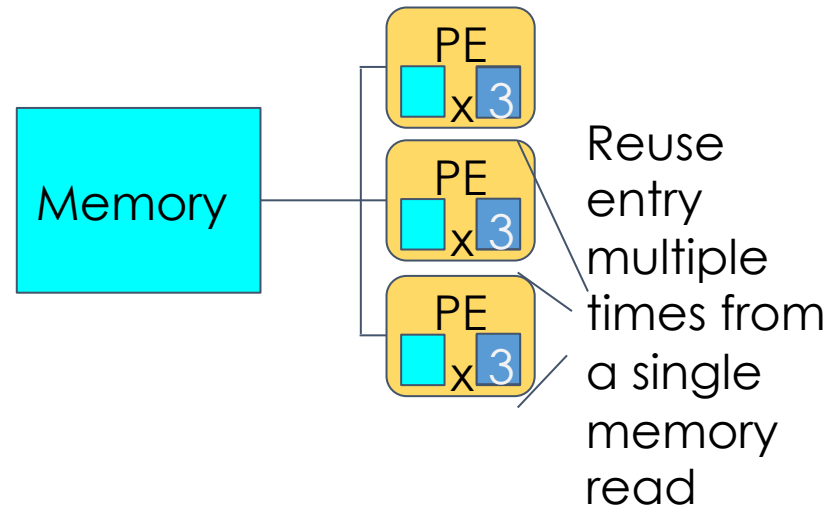
Data Reuse



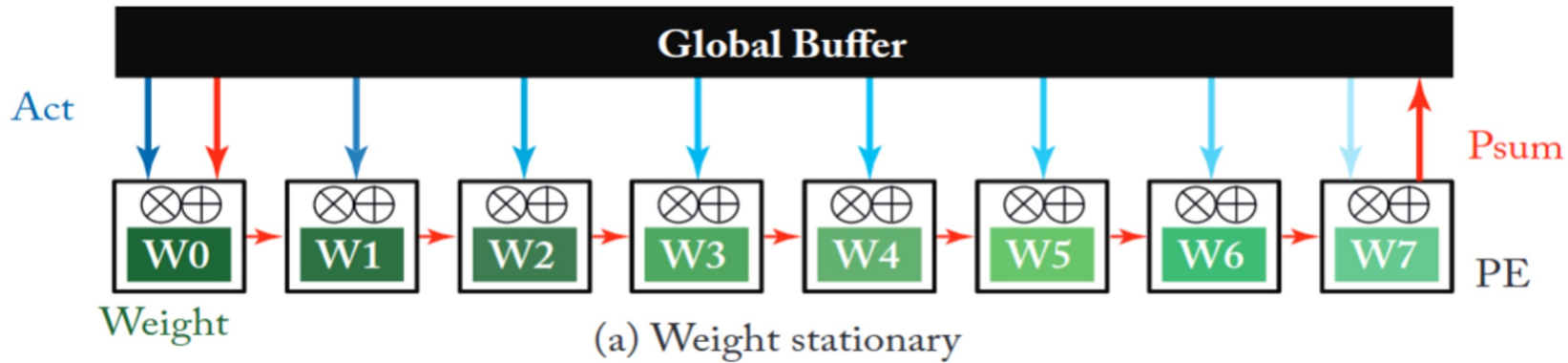
Read once from memory, use same data multiple times by same PE



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



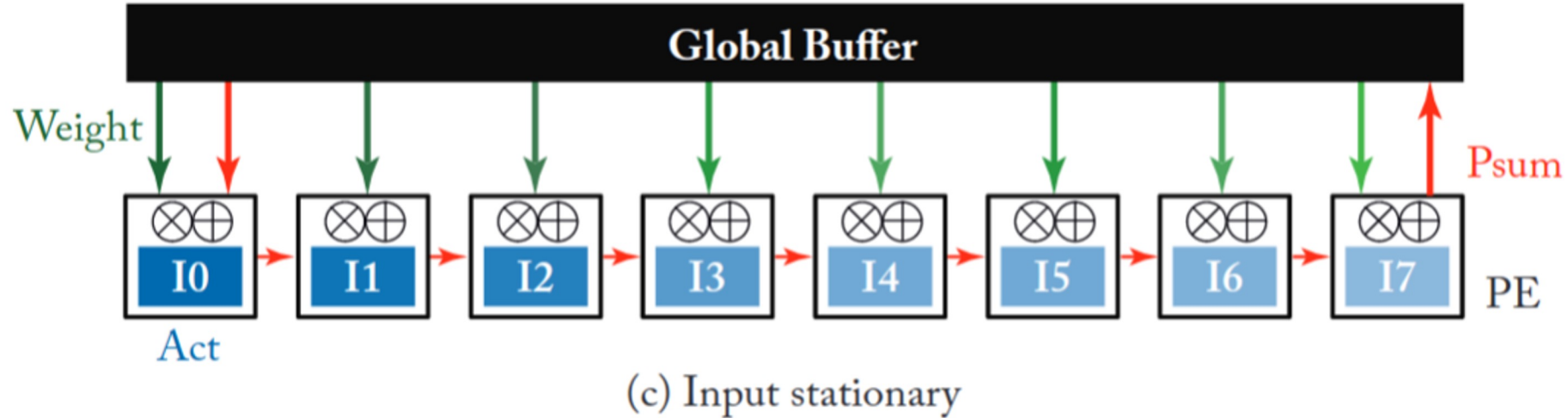
Stationary Weights?



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



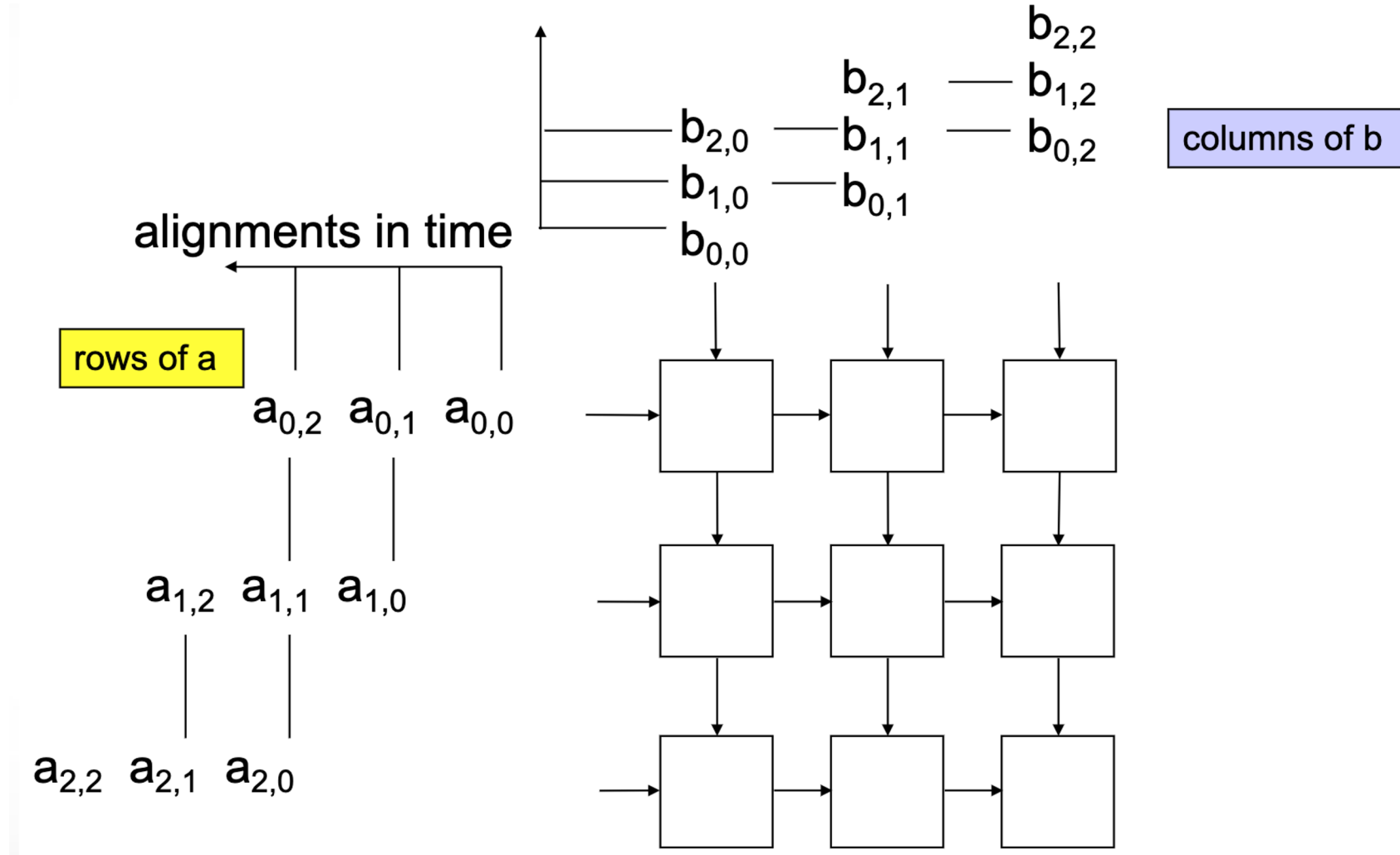
Stationary Inputs?



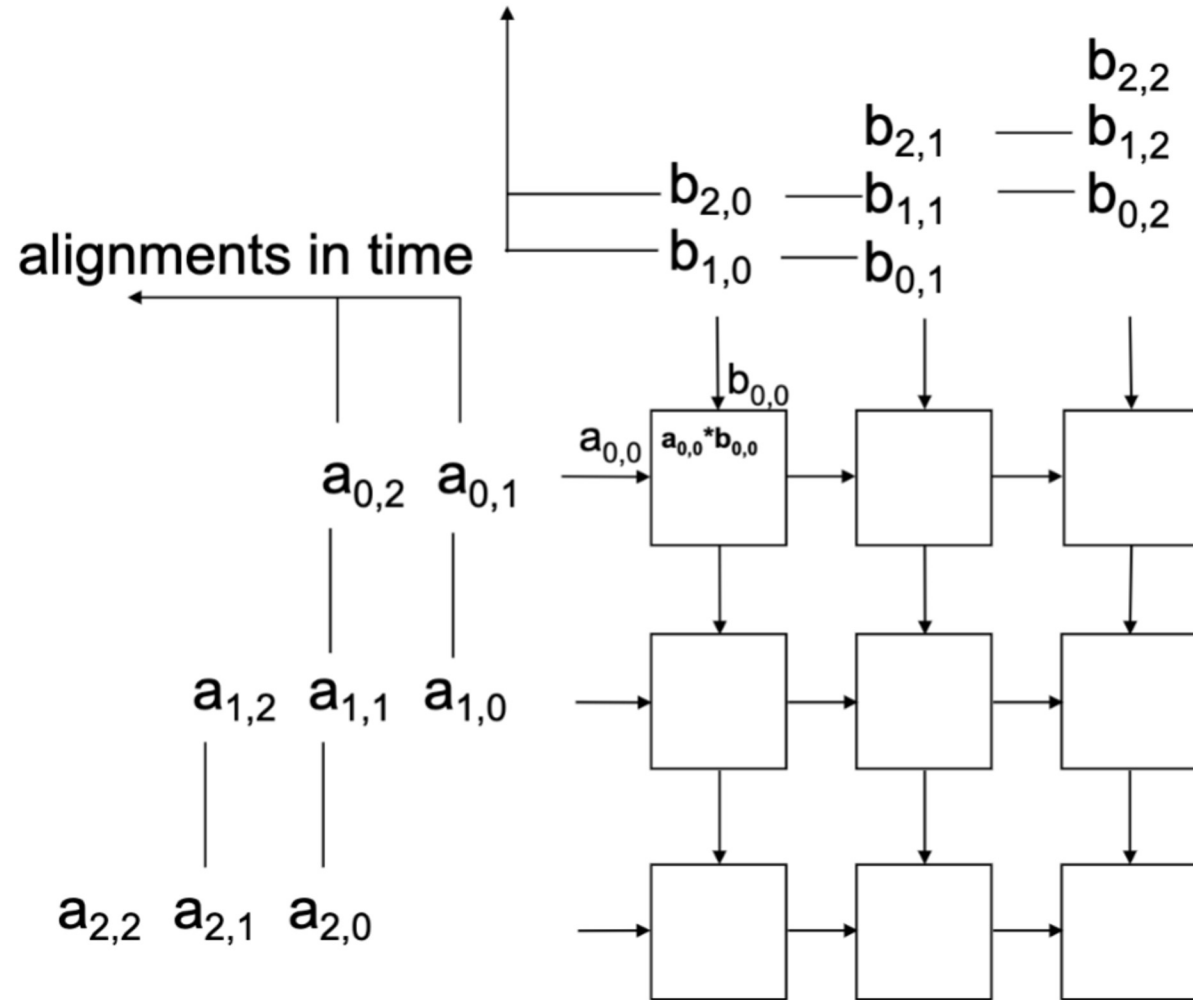
$$a_0w_0 + a_1w_1 + a_2w_2 + a_3w_3 + a_4w_4 + a_5w_5 + a_6w_6 + a_7w_7$$



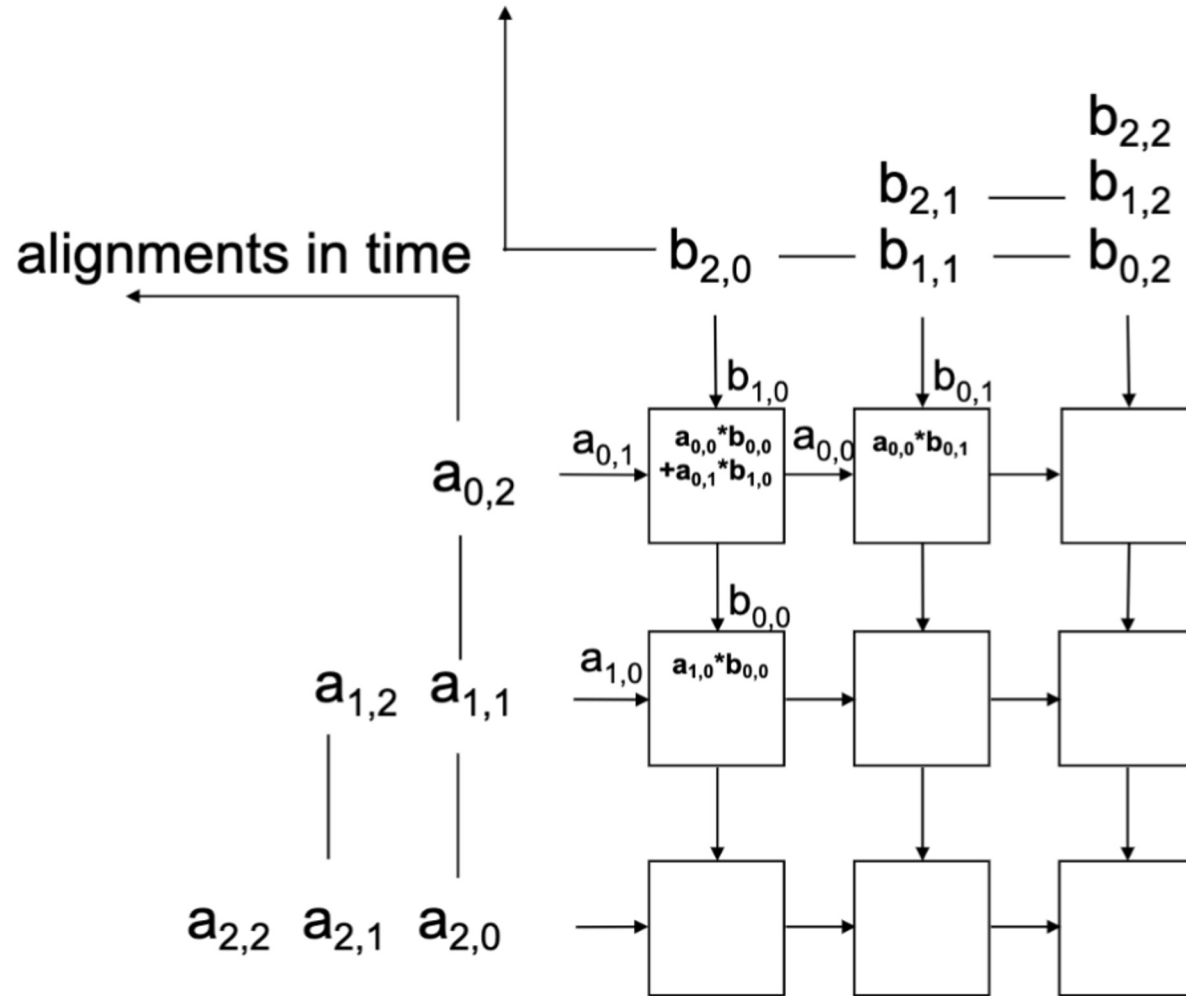
Systolic Array: Matrix Mults.



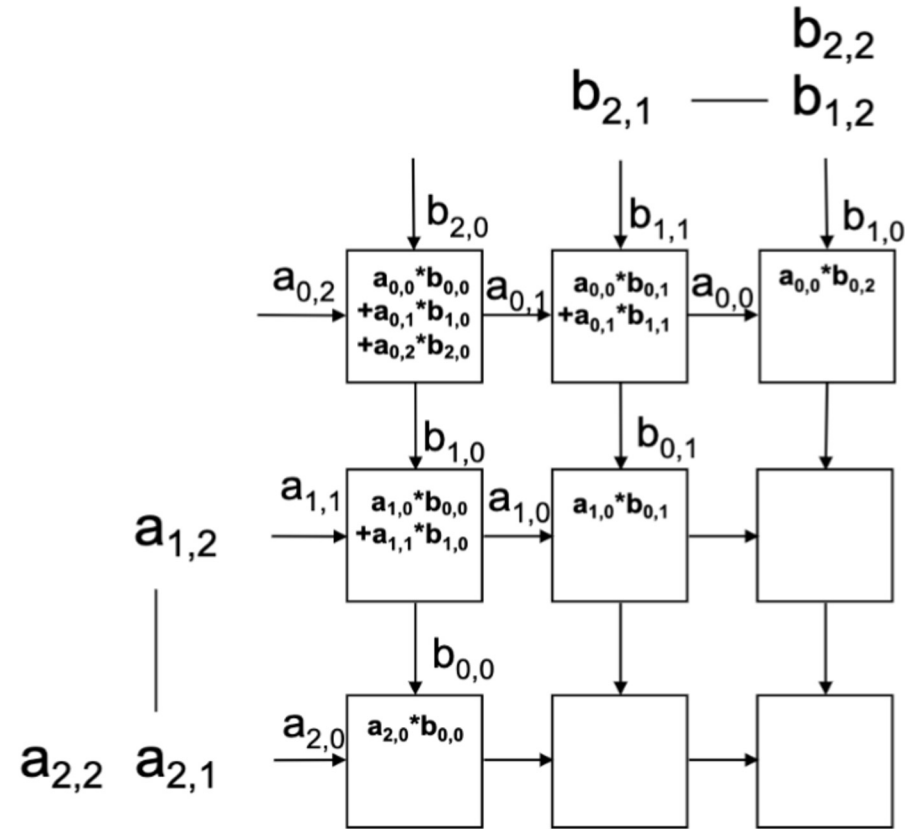
Systolic Array: Matrix Mults.



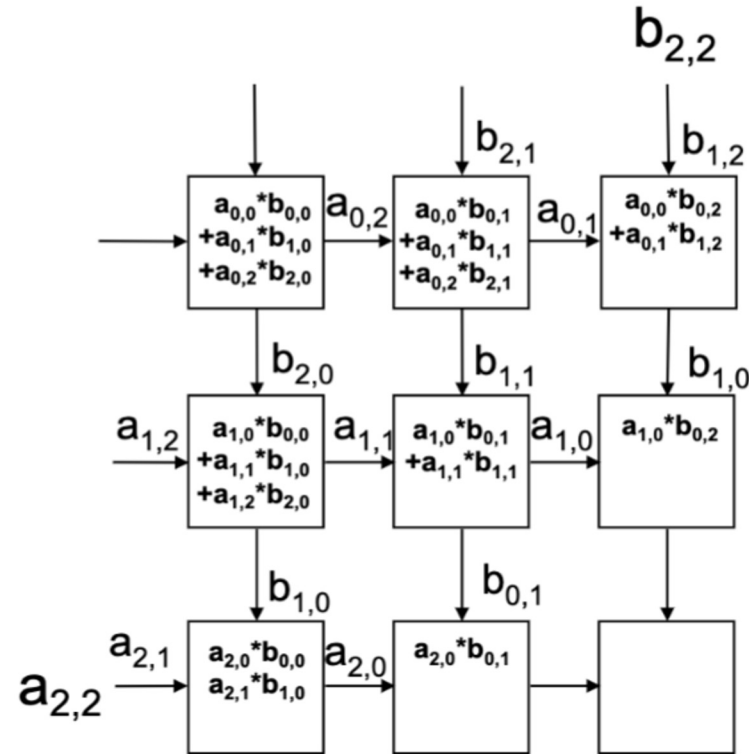
Systolic Array: Matrix Mults.



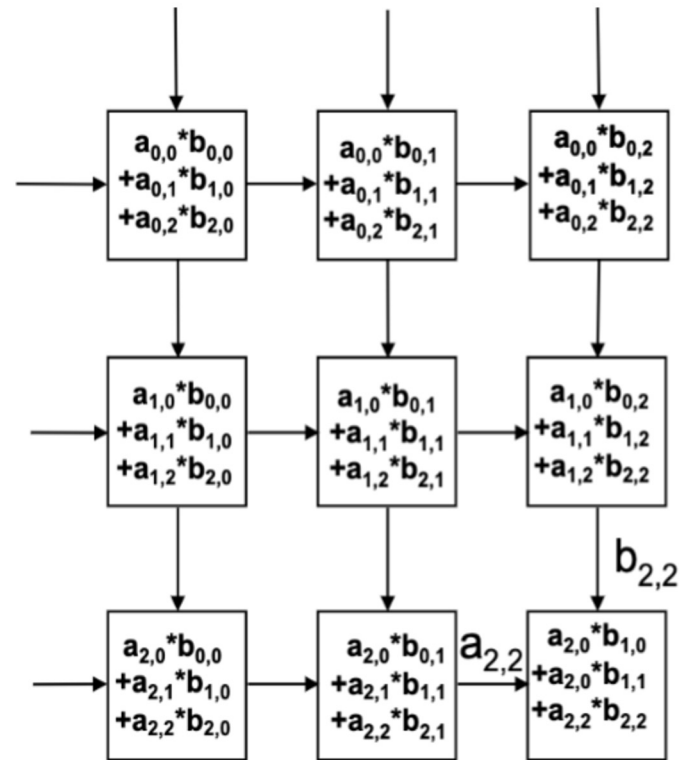
Systemic Array: Matrix Mults.



Systemic Array: Matrix Mults.



Systemic Array: Matrix Mults.

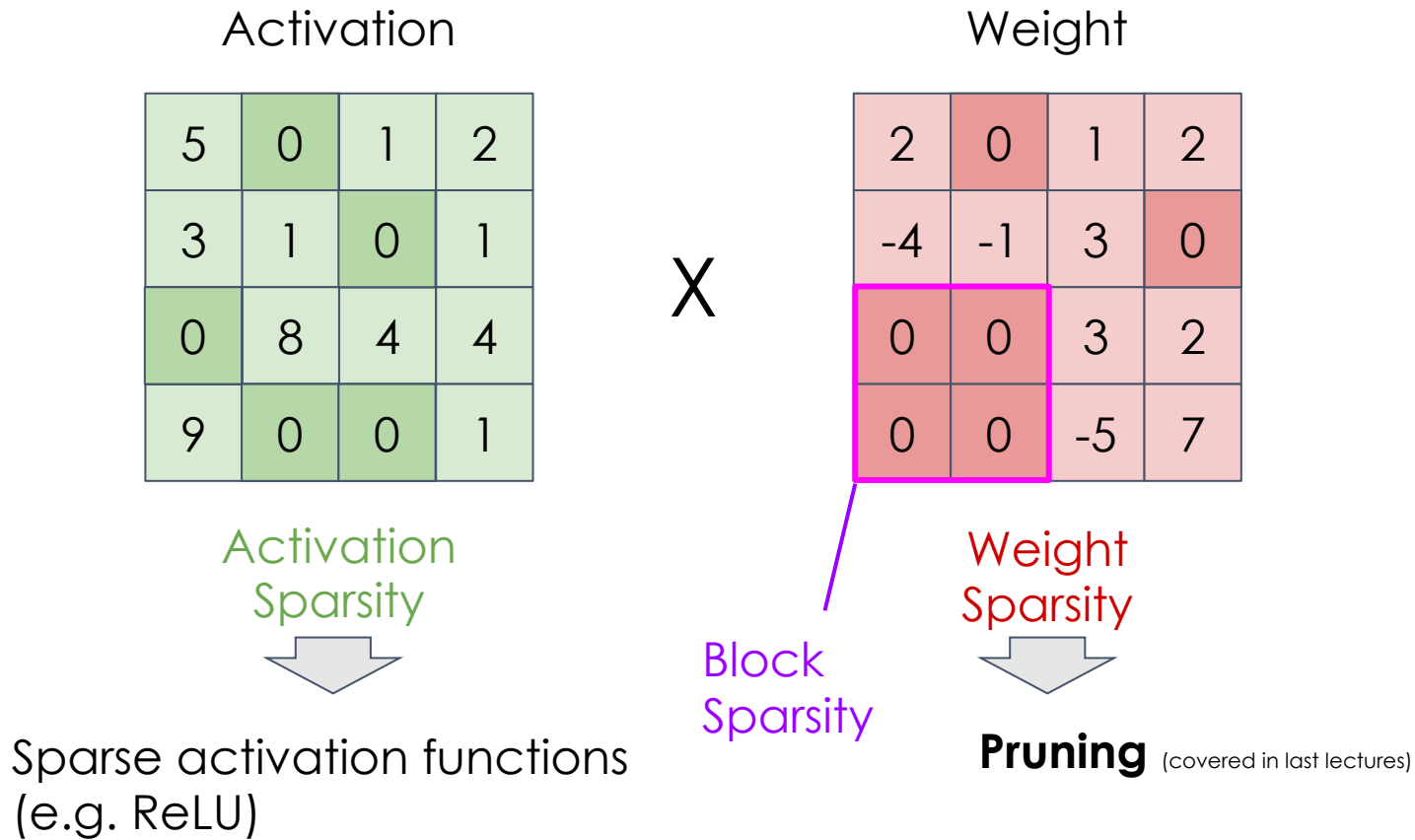


Roadmap for Today

- HW enabling Deep Learning
- Performance Metrics
- Where does Energy Go?
- **Hardware Efficiency Options**
 - Arithmetic
 - Memory
 - **Ineffectual Operation**
- Hardware Case Studies

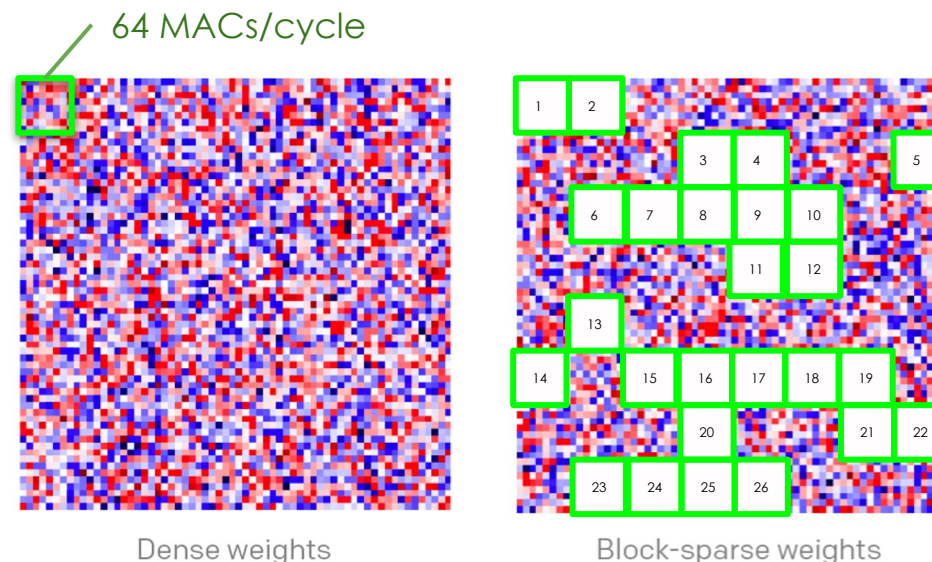


Kinds of Sparsity



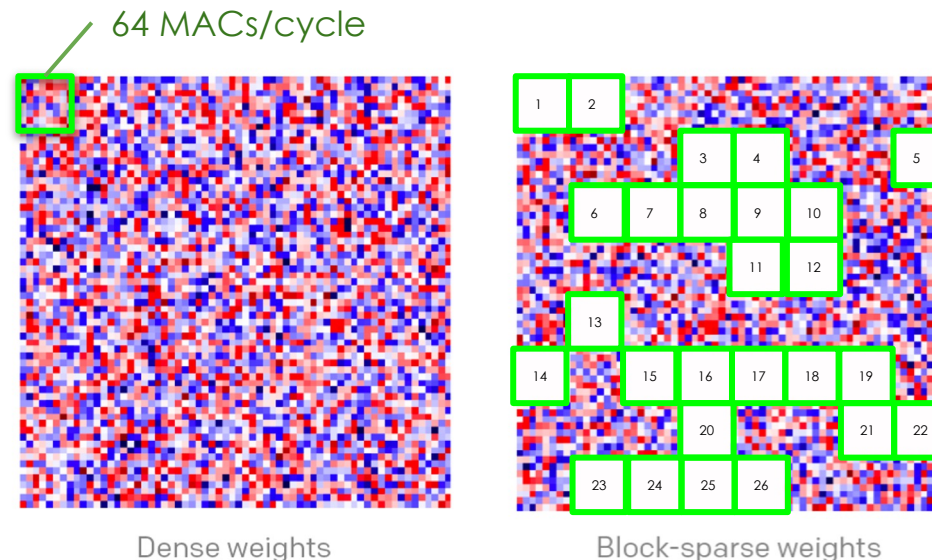
Coarse-grained “Block” Sparsity

- 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



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



Simplest way to leverage sparsity with low overhead

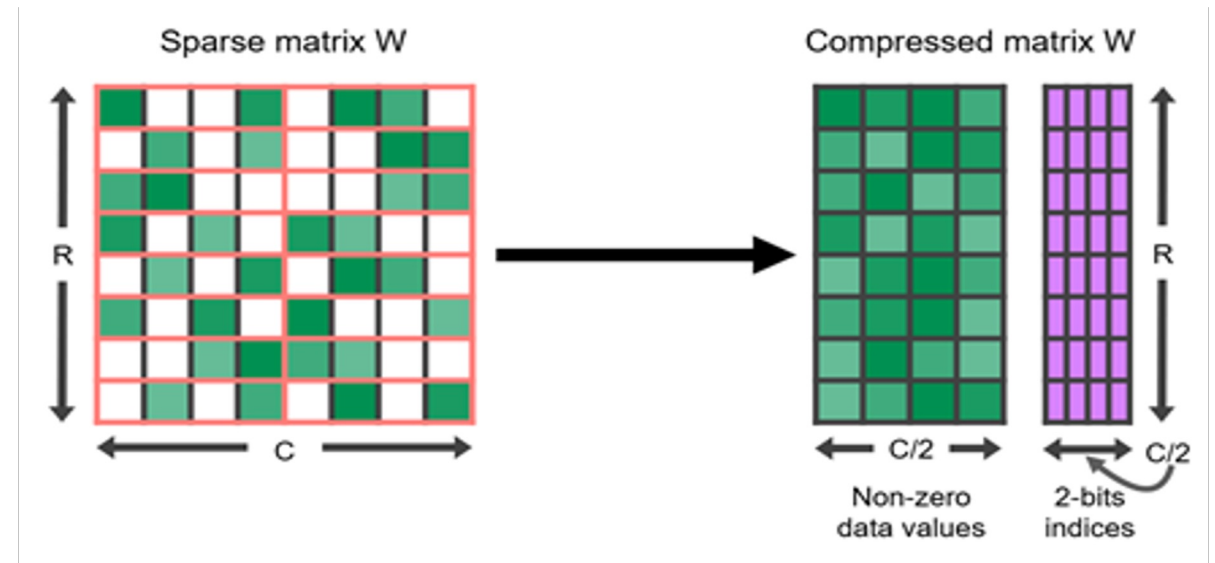
⇒ Single bit per 8x8 block ($1/64 = 1.6\%$ overhead)

⇒ Simple control logic because entire block is skipped



Fine-grained Sparsity (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?

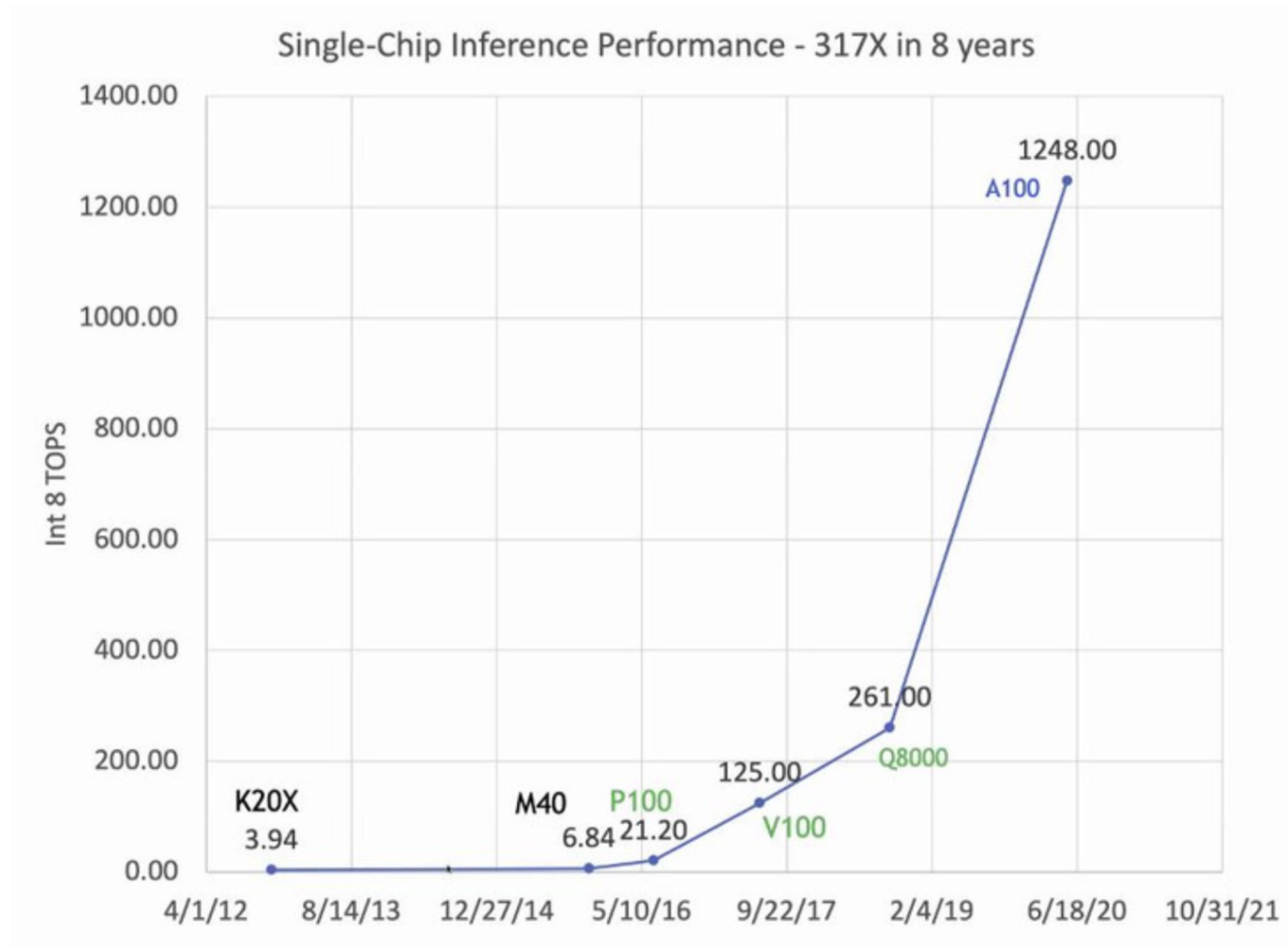


Roadmap for Today

- HW enabling Deep Learning
- Performance Metrics
- Where does Energy Go?
- Hardware Efficiency Options
- **Hardware Case Studies**



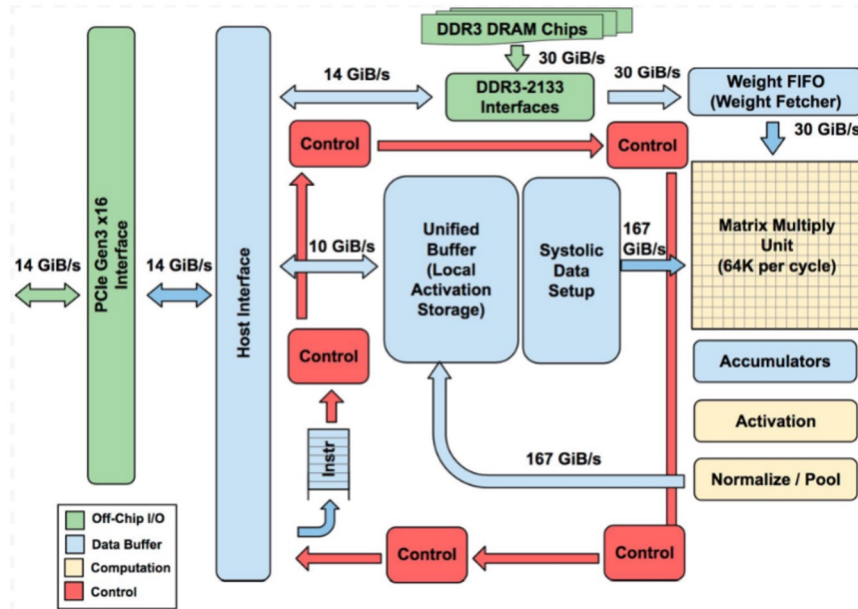
Nvidia GPU Progression



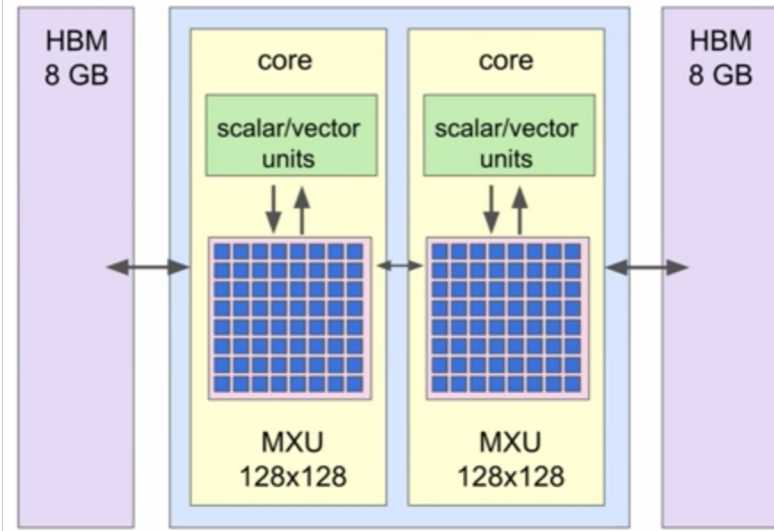
Source: Dally



Google TPU



v2



Source: Google

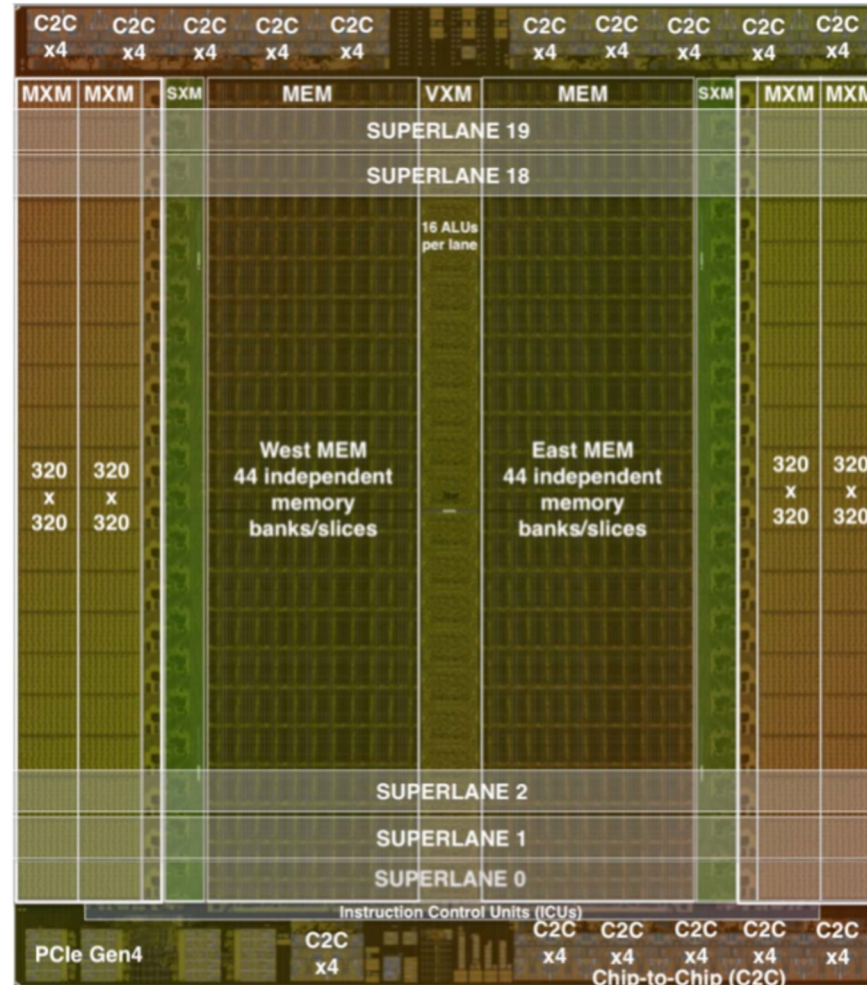
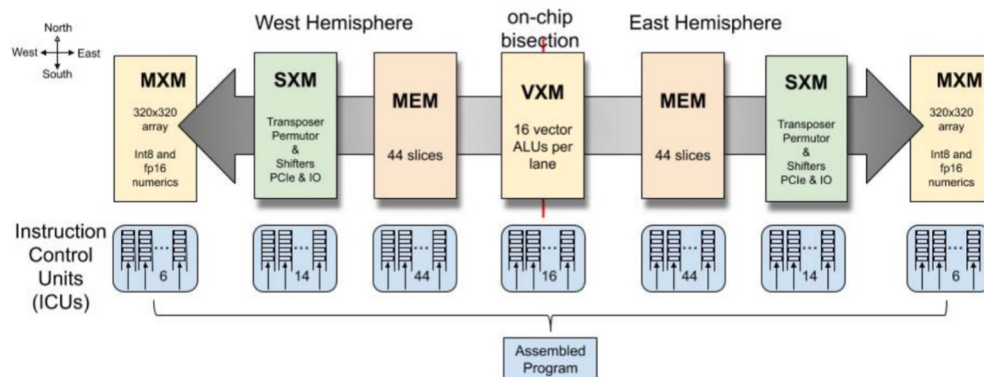
- 16 GB of HBM
- 600 GB/s mem BW
- Scalar/vector units: 32b float
- MXU: 32b float accumulation but reduced precision for multipliers
- 45 TFLOPS

Source: D. Harris



Groq

- Programmable dataflow architecture
- 1000 TOPs/s peak INT8 performance
- 200 MB on-chip SRAM (80 TB/s)
 - No external memory, scales by increasing number of chips
- FP16 and INT8 precision
- Philosophy: “unroll” a multicore architecture on-chip spatially to allow for custom instructions



Source: Groq



Groq

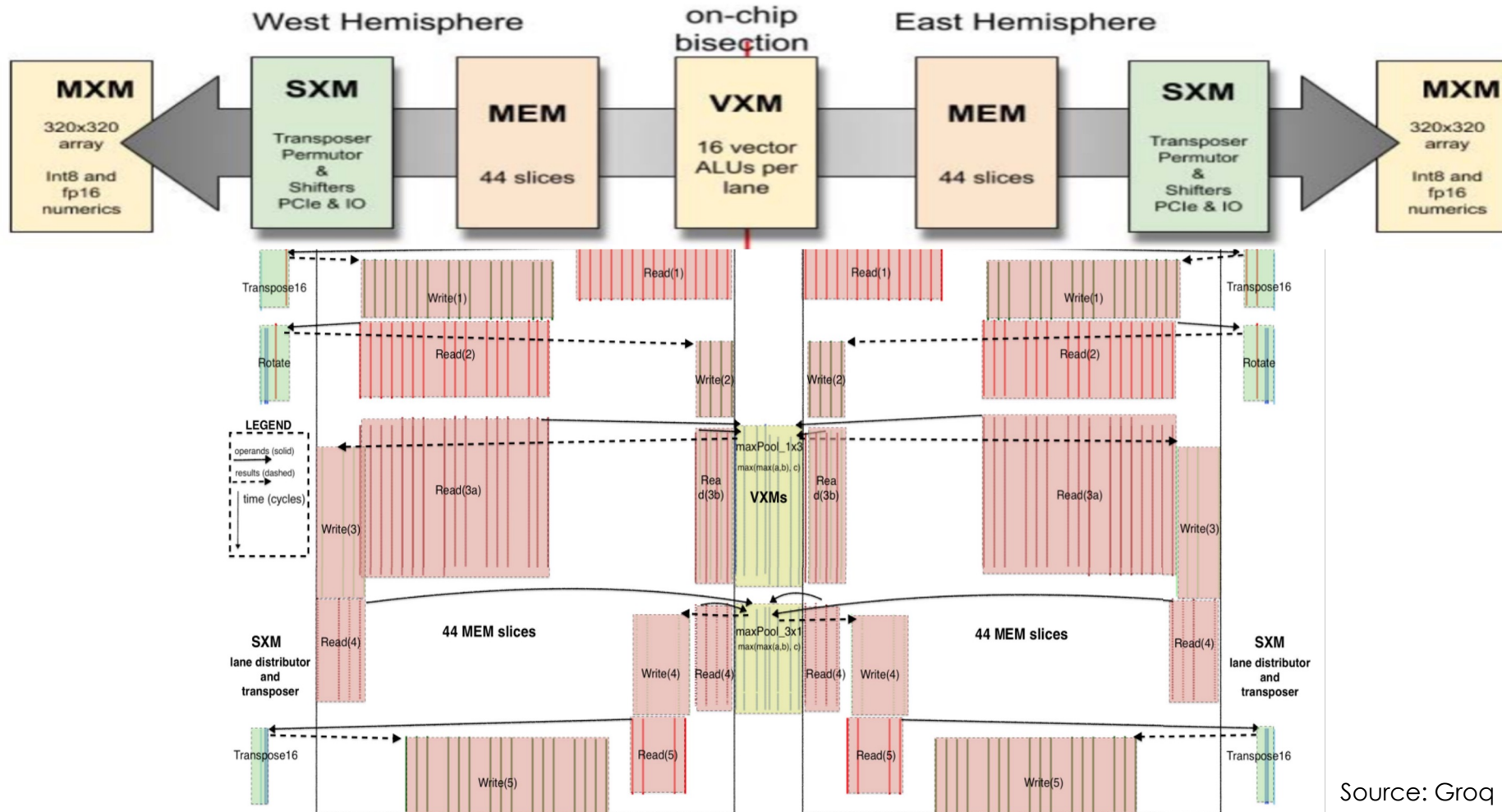
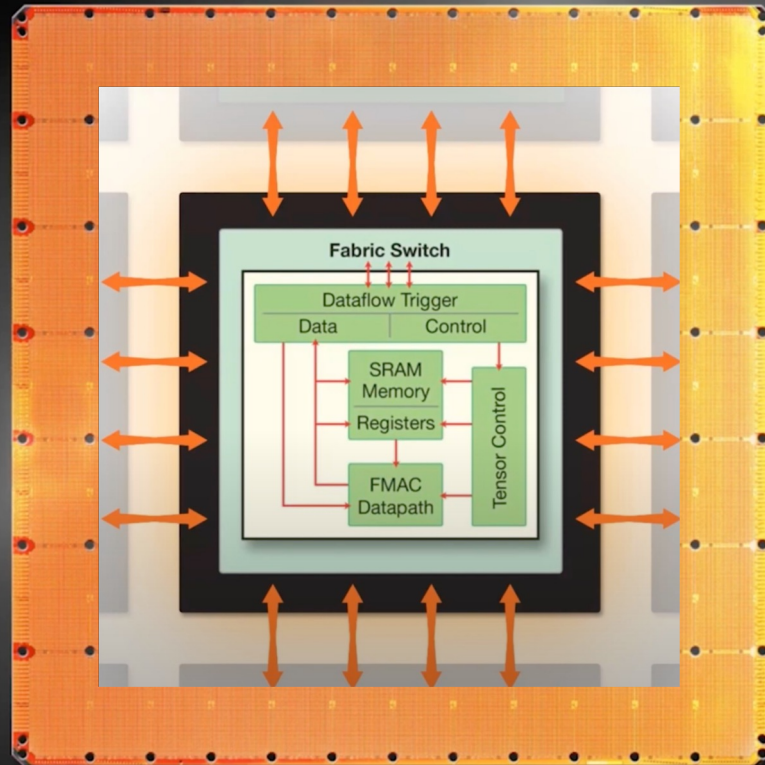


Fig. 11. Example instruction schedule for 3x3 max pool in ResNet50.



Cerebras



Cerebras WSE
1.2 Trillion Transistors
46.225 mm² Silicon

Largest Chip Ever Built

- 46,225 mm² silicon
- 1.2 trillion transistors
- 400,000 AI optimized cores
- 18 Gigabytes of On-chip Memory
- 9 PByte/s memory bandwidth
- 100 Pbit/s fabric bandwidth
- TSMC 16nm process



Largest GPU
21.1 Billion Transistors
815 mm² Silicon

Source:Cerebras



Summary of the Day

- HW enabling Deep Learning
- Performance Metrics
- Where does Energy Go?
- Hardware Efficiency Options
- Hardware Case Studies

