



<https://hao-ai-lab.github.io/cse234-w25/>

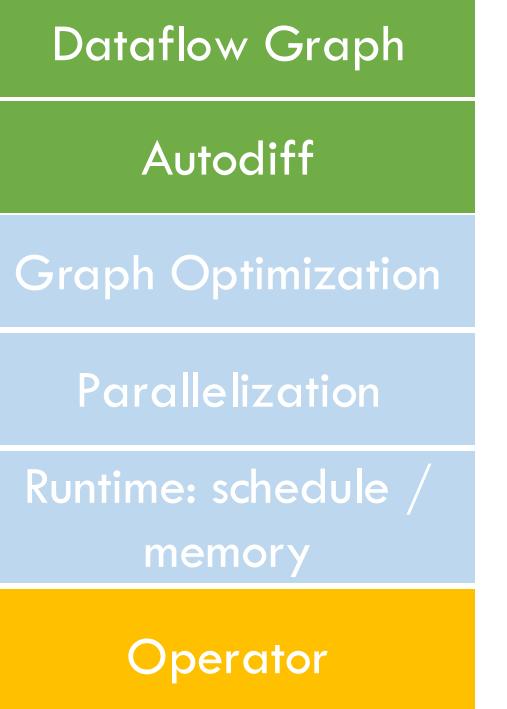
CSE 234: Data Systems for Machine Learning

Winter 2025

LLMSys

Optimizations and Parallelization

MLSys Basics



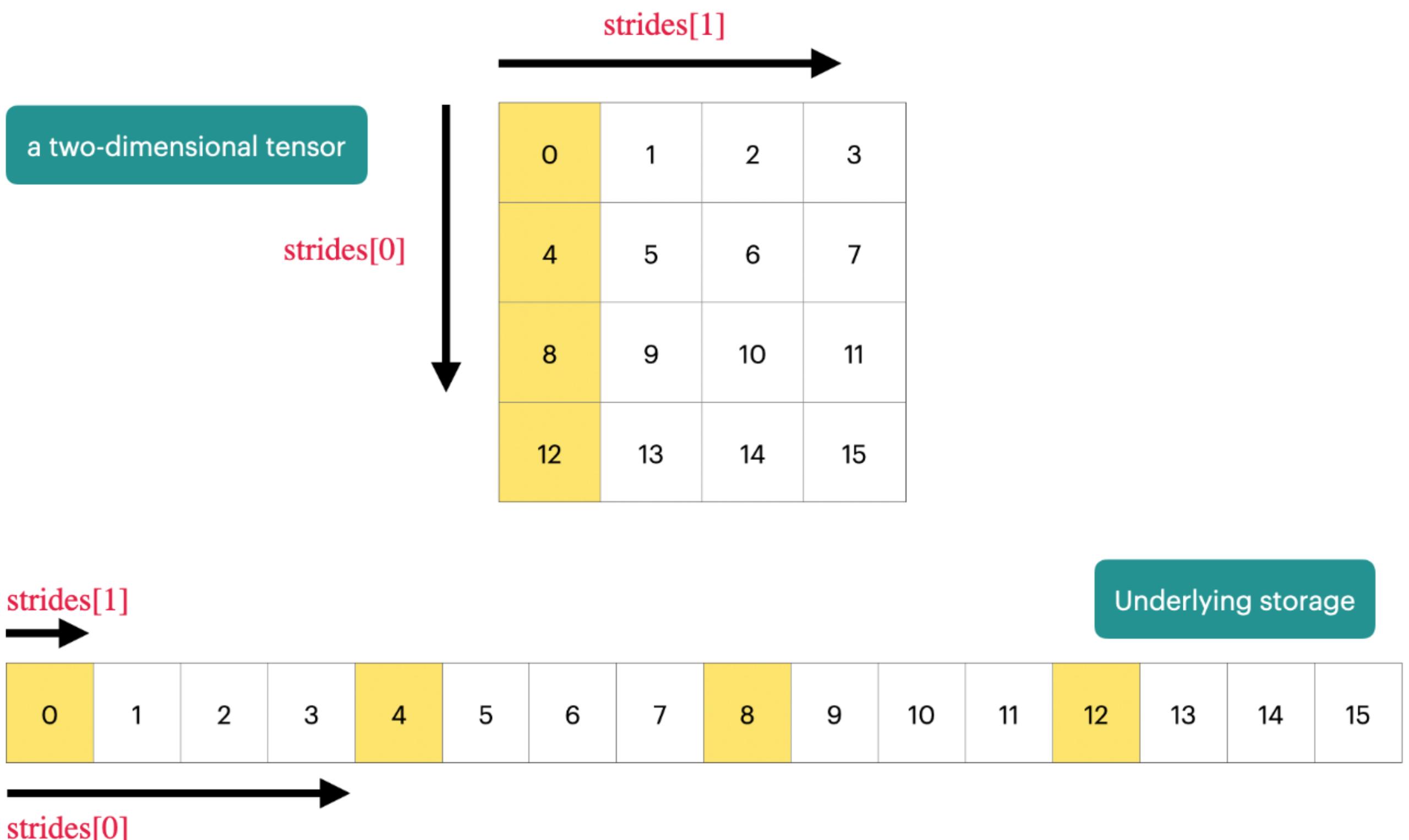
Today's Learning Goals

- How to make operators fast in general?
 - Vectorize
 - Data layout
 - Parallelization (at the operator level)
- Matmul-specific optimization
- GPUs and accelerators
 - High-level Idea
 - The accelerator market

Recap: Strides format

Python

```
1 A[i0][i1][i2]... = A_internal[  
2     stride_offset  
3     + i0 * A.strides[0]  
4     + i1 * A.strides[1]  
5     + i2 * A.strides[2]  
6     + ...  
7     + in-1 * A.strides[n-1]  
8 ]
```



Why we bother saving “strides” when saving tensors

- Strides can separate **the underlying storage** and **the view of the tensor**

torch.Tensor.view

Tensor.view(*shape) → Tensor

```
>>> x = torch.randn(4, 4)
>>> x.size()
torch.Size([4, 4])
>>> y = x.view(16)
>>> y.size()
torch.Size([16])
>>> z = x.view(-1, 8)  #
>>> z.size()
torch.Size([2, 8])
```

- Enable **zero-copy** of some very frequently used operators

Operator: Slice

- Change the `offset` by +1
- Reduce the `shape` to [3, 2]
- Note: zero copy

0	1	2	3	4
5	6	7	8	9
10	11	12	13	14
15	16	17	18	19

Operator: Transpose

- How to do Transpose?

▼ Python

```
1 print(t.stride())
2 # (24, 12, 4, 1)
3
4 print(t.permute((1, 2, 3, 0)).is_contiguous())
5 # True
6
7 print(t.permute((1, 2, 3, 0)).stride())
8 # (12, 4, 1, 24)
9
10 print_internal(t.permute((1, 2, 3, 0)))
11 # tensor([0, 1, 2, 3,
12 #         4, 5, 6, 7,
13 #         8, 9, 10, 11,
14 #         12, 13, 14, 15,
15 #         16, 17, 18, 19,
16 #         20, 21, 22, 23])
```

Original Matrix A

j	0	1	2
i	3	4	5
0	0	1	2
1	3	4	5

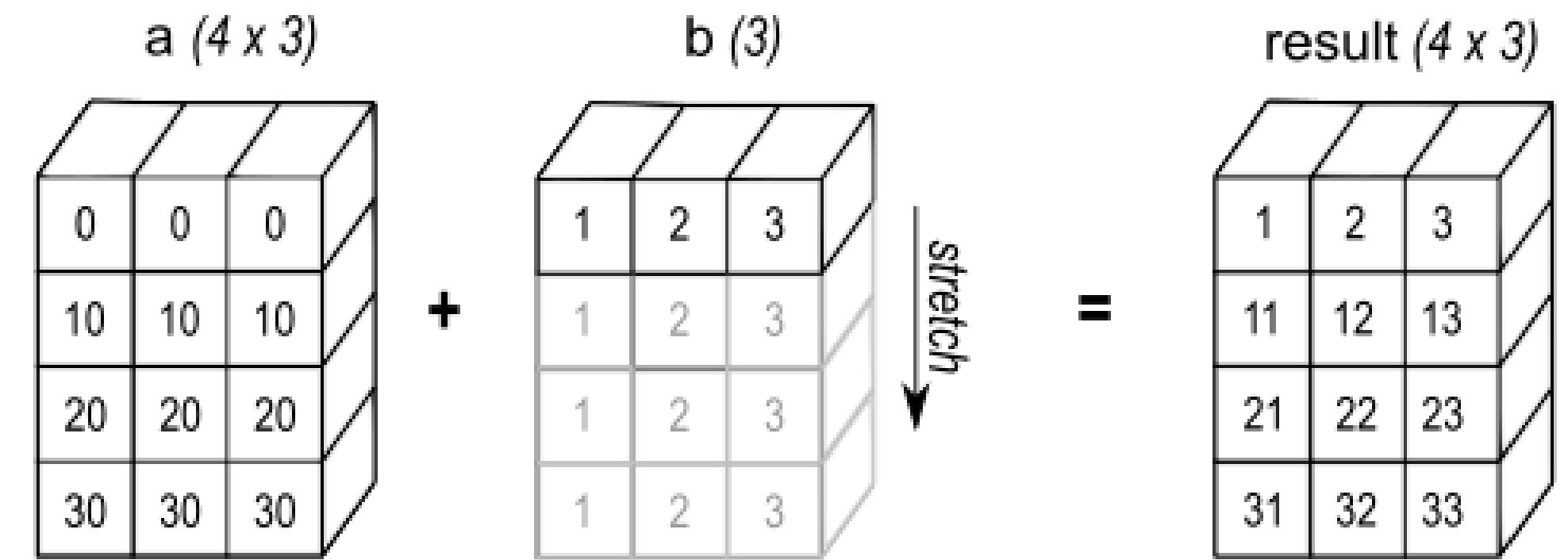
Transposed Matrix A^T

i	0	3
j	1	4
0	1	4
1	2	5

- Note 1: zero copy
- Note 2: underlying storage is unchanged

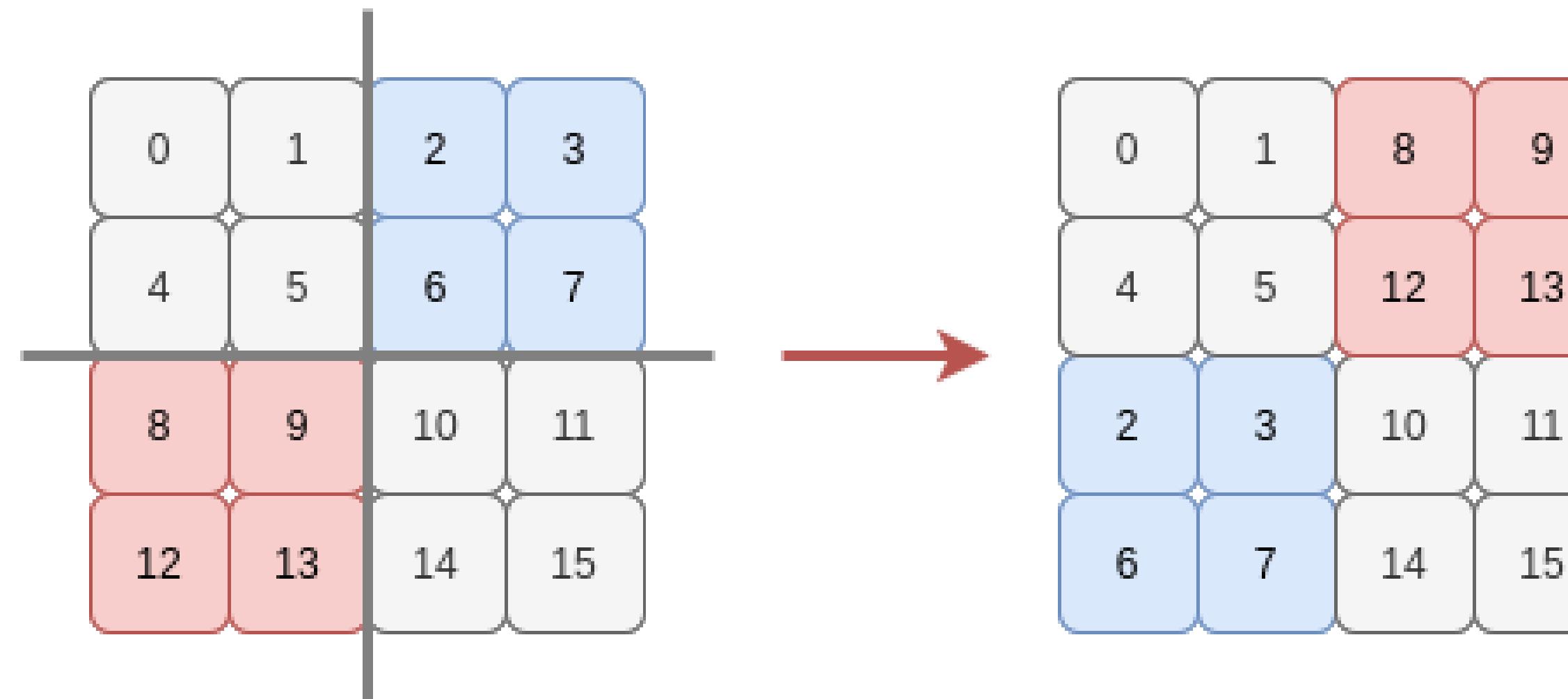
Broadcast

- Question: how to do broadcast?



- `b.strides=[1], b.shape=[1,3], b.data=[1, 2, 3]`
- After broadcast: `b.shape=[4,3], b.data=[1, 2, 3], b.strides=?`
- Recall the def of strides: $A[i, j] = A.data[offset + i * A.strides[0] + j * A.strides[1]]$

Home Exercise: Swapping tiles



```
>>> a = np.arange(16).reshape(4, 4)
>>> a
array([[ 0,  1,  2,  3],
       [ 4,  5,  6,  7],
       [ 8,  9, 10, 11],
       [12, 13, 14, 15]])
>>> np.vstack((np.hstack((a[0:2, 0:2], a[2:4, 0:2])), np.hstack((a[0:2, 2:4], a[2:4, 2:4]))))
array([[ 0,  1,  8,  9],
       [ 4,  5, 12, 13],
       [ 2,  3, 10, 11],
       [ 6,  7, 14, 15]])
```

Problems of Strides

- Memory Access may become not continuous
 - Many vectorized ops requires continuous storage
 - What's the underlying storage after slice?

0	1	2	3	4
5	6	7	8	9
10	11	12	13	14
15	16	17	18	19

TORCH.TENSOR.CONIGUOUS

`Tensor.contiguous(memory_format=torch.contiguous_format) → Tensor`

Returns a contiguous in memory tensor containing the same data as `self` tensor. If `self` tensor is already in the specified memory format, this function returns the `self` tensor.

Parameters

memory_format (`torch.memory_format`, optional) – the desired memory format of returned Tensor. Default: `torch.contiguous_format`.

Parallelization (on CPUs)

We'll com
back to this
later

- How to parallelize the loop?

```
for (int i = 0; i < 64; ++i) {  
    float4 a = load_float4(A + i*4);  
    float4 b = load_float4(B + i*4);  
    float4 c = add_float4(a, b);  
    store_float4(C + i* 4, c);
```

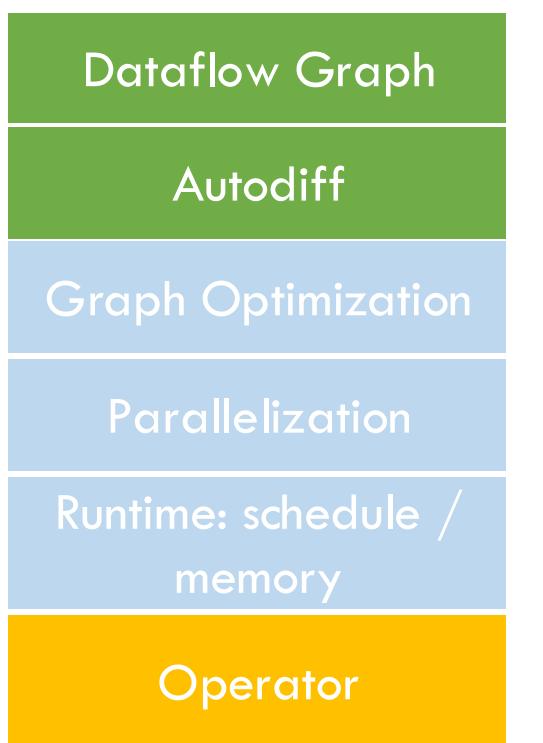
```
#pragma omp parallel for  
for (int i = 0; i < 64; ++i) {  
    float4 a = load_float4(A + i*4);  
    float4 b = load_float4(B + i*4);  
    float4 c = add_float4(a, b);  
    store_float4(C * 4, c);
```

vectorized

Vectorized &
parallelized

Summary

- Vectorization
 - Leverage platform-specific vectorized functions
 - reduce seek time
- Data layout
 - Stride format
 - Zero copy
 - Enable fast array-manipulation: slice, transpose, broadcast, etc.
- Parallelization on CPUs



Next

- How to make operators fast in general?
 - Vectorize
 - Data layout
 - Parallelization (at the operator level)
- Matmul-specific optimization
- GPUs and accelerators
 - High-level Idea
 - The accelerator market

What is Matmul in Code?

Compute $C = \text{dot}(A, B.T)$

```
float A[n][n], B[n][n], C[n][n];

for (int i = 0; i < n; ++i)
    for (int j = 0; j < n; ++j) {
        C[i][j] = 0;
        for (int k = 0; k < n; ++k) {
            C[i][j] += A[i][k] * B[j][k];
        }
    }
```

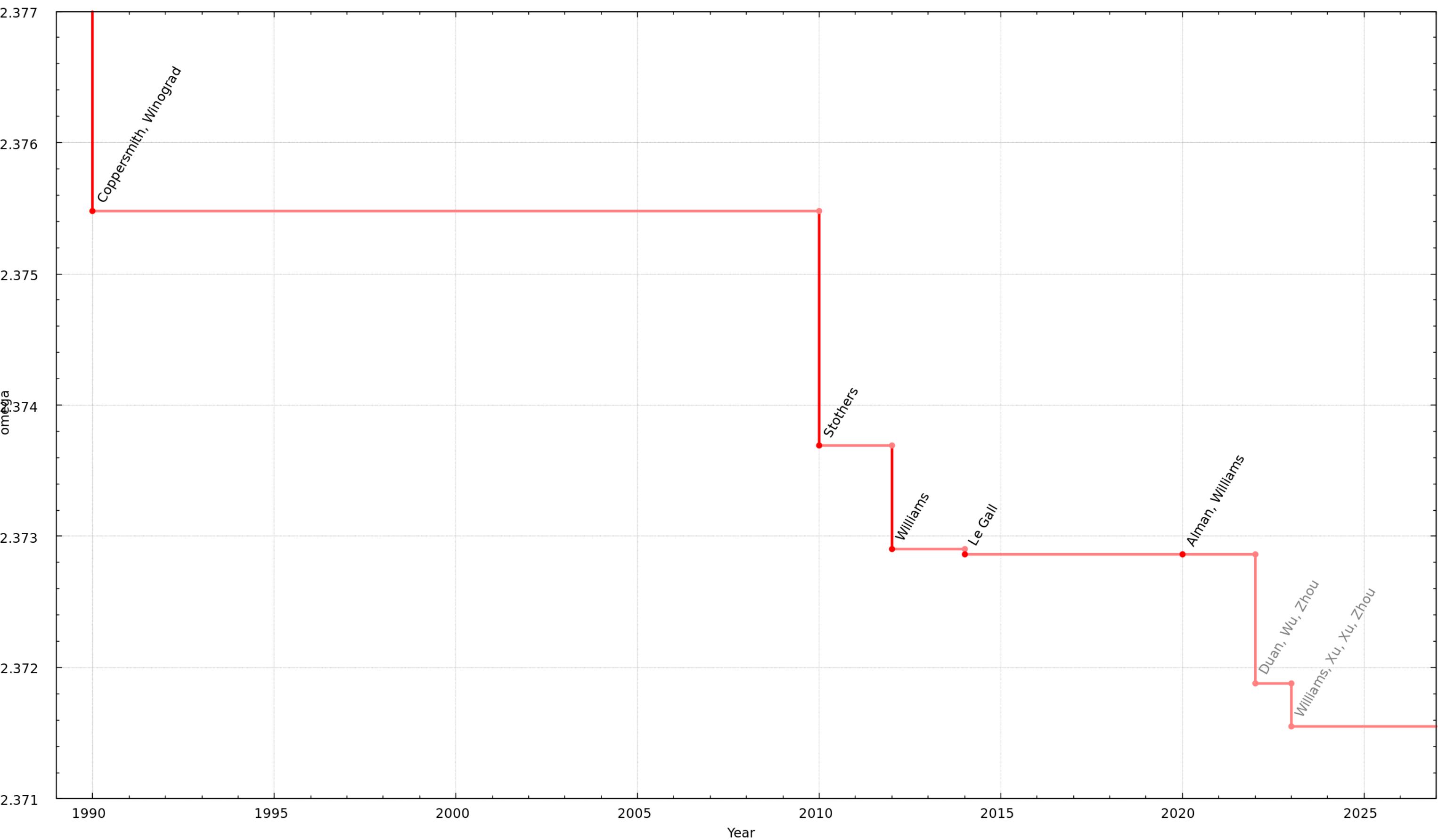
- What is the time complexity of 2D matmul?
- $O(n^3)$
- What is the best complexity we can achieve?
- $O(n^{2.371552})$

Matmul Complexity

Not a good area
to do research 😊

Timeline of matrix multiplication exponent

Year	Bound on omega	Authors
1969	2.8074	Strassen ^[1]
1978	2.796	Pan ^[10]
1979	2.780	Bini, Capovani ^[it] , Romani ^[11]
1981	2.522	Schönhage ^[12]
1981	2.517	Romani ^[13]
1981	2.496	Coppersmith, Winograd ^[14]
1986	2.479	Strassen ^[15]
1990	2.3755	Coppersmith, Winograd ^[16]
2010	2.3737	Stothers ^[17]
2012	2.3729	Williams ^{[18][19]}
2014	2.3728639	Le Gall ^[20]
2020	2.3728596	Alman, Williams ^{[21][22]}
2022	2.371866	Duan, Wu, Zhou ^[23]
2024	2.371552	Williams, Xu, Xu, and Zhou ^[2]



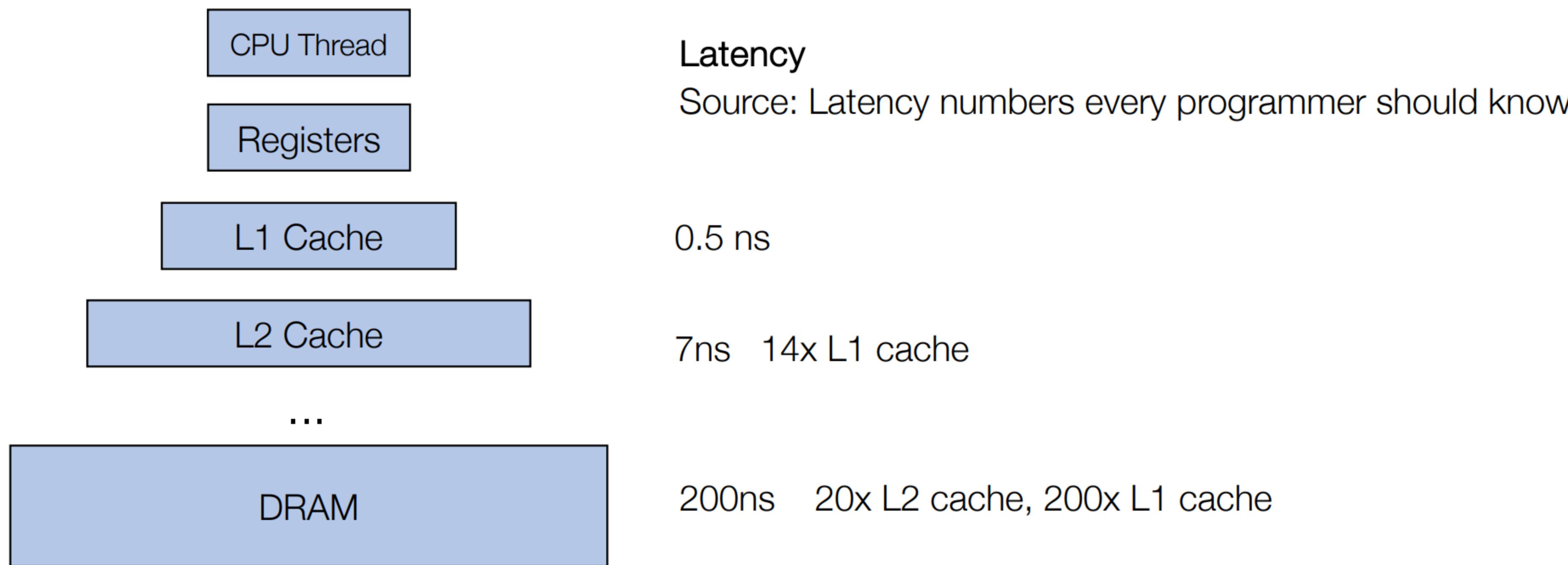
How to Make Matmul Fast

$$\max \text{AI} = \#ops / \#bytes$$

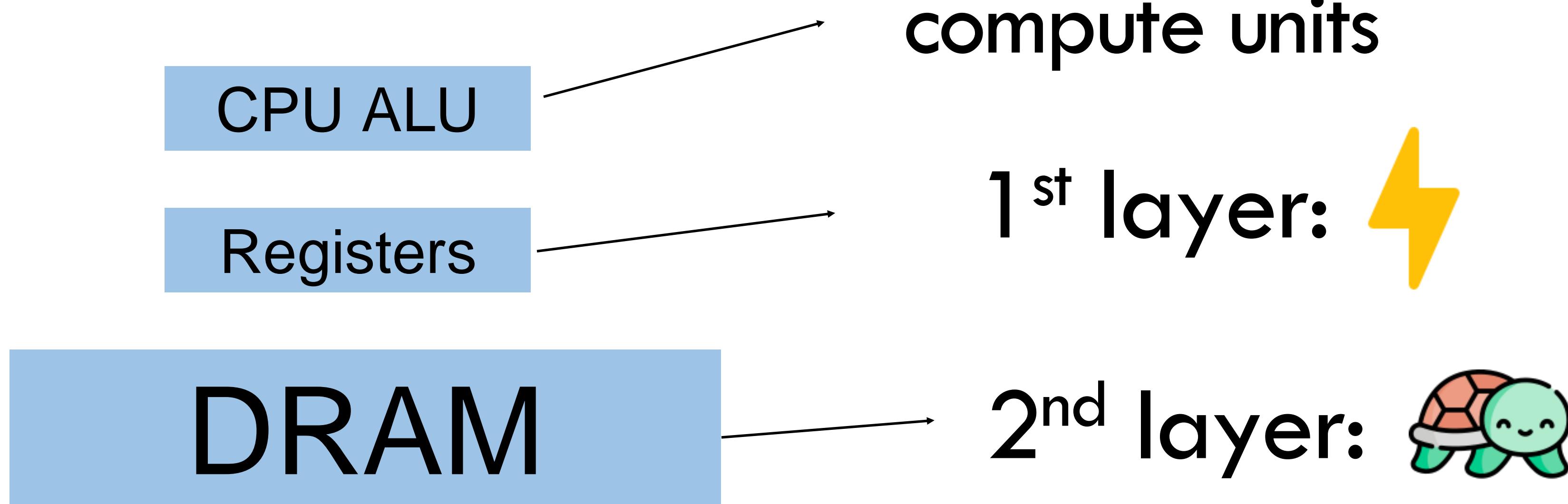


Recall: Memory Hierarchy

- Ideally: we want everything to be local to processors (In registers)
- But registers are expensive and small, hence memory hierarchy



Simplify It a bit



Recall How to Estimate Al: count loads

```
float* A, *B, *C, *D, *E, *tmp1, *tmp2;  
// assume arrays are allocated here  
// compute E = D + ((A + B) * C)  
add(n, A, B, tmp1);  
mul(n, tmp1, C, tmp2);  
add(n, tmp2, D, E);
```

Overall arithmetic intensity = 1/3

```
void fused(int n, float* A, float* B, float* C, float* D,  
          float* E) {  
    for (int i=0; i<n; i++)  
        E[i] = D[i] + (A[i] + B[i]) * C[i];  
}  
// compute E = D + (A + B) * C  
fused(n, A, B, C, D, E);
```

Four loads, one store per 3 math ops
arithmetic intensity = 3/5

Review Matmul loop

```
dram float A[n][n], B[n][n], C[n][n];
for (int i = 0; i < n; ++i) {
    for (int j = 0; j < n; ++j) {
        register float c = 0;
        for (int k = 0; k < n; ++k) {
            register float a = A[i][k];
            register float b = B[j][k];
            c += a * b;
        }
        C[i][j] = c;
    }
}
```

Read a	n^3
Read b	n^3
Write c	n^2

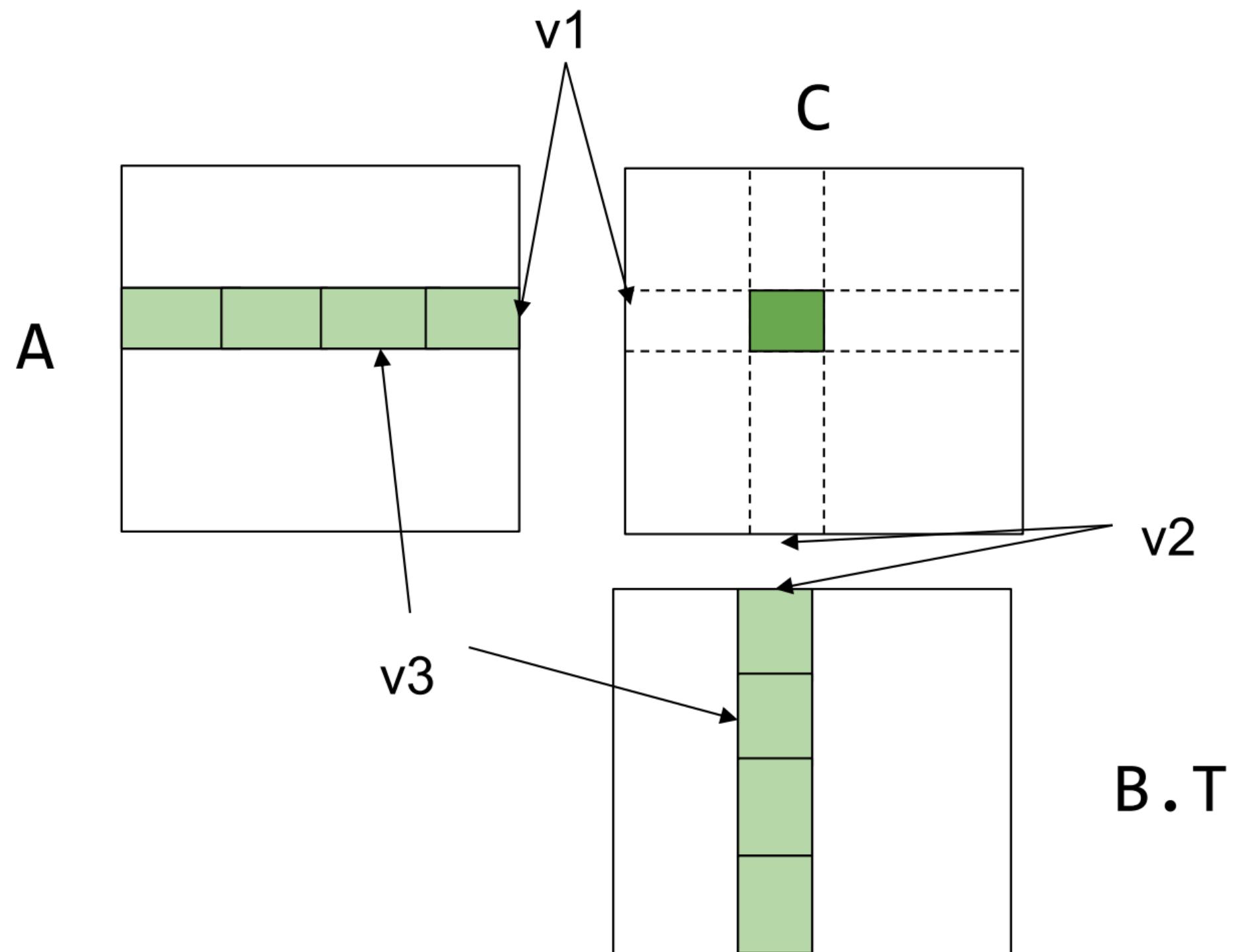
#Registers needed:
 $1 + 1 + 1 = 3$

Read cost:
 $2 * n^3 * \text{speed(dram -> register)}$

Register Tiled Matrix Multiplication

```
dram float A[n/v1][n/v3][v1][v3];
dram float B[n/v2][n/v3][v2][v3];
dram float C[n/v1][n/v2][v1][v2];

for (int i = 0; i < n/v1; ++i) {
    for (int j = 0; j < n/v2; ++j) {
        register float c[v1][v2] = 0;
        for (int k = 0; k < n/v3; ++k) {
            register float a[v1][v3] = A[i][k];
            register float b[v2][v3] = B[j][k];
            c += dot(a, b.T);
        }
        C[i][j] = c;
    }
}
```



Register Tiled Matrix Multiplication

```
dram float A[n/v1][n/v3][v1][v3];
dram float B[n/v2][n/v3][v2][v3];
dram float C[n/v1][n/v2][v1][v2];

for (int i = 0; i < n/v1; ++i) {
    for (int j = 0; j < n/v2; ++j) {
        register float c[v1][v2] = 0;
        for (int k = 0; k < n/v3; ++k) {
            register float a[v1][v3] = A[i][k];
            register float b[v2][v3] = B[j][k];
            c += dot(a, b.T);
        }
        C[i][j] = c;
    }
}
```

Read a

N^3 / v_2

Read b

N^3 / v_1

Write c

N^2

#Registers needed:

$v1*v3 + v2*v3 + v1*v2$

Read cost:

$(N^3/v2 + N^3/v1) * \text{speed(dram -> register)}$

Register Tiled Matrix Multiplication

- Q: is the load cost related to v3?

Read a

N^3 / v_2

- Q: How to set v1 / v2?

Read b

N^3 / v_1

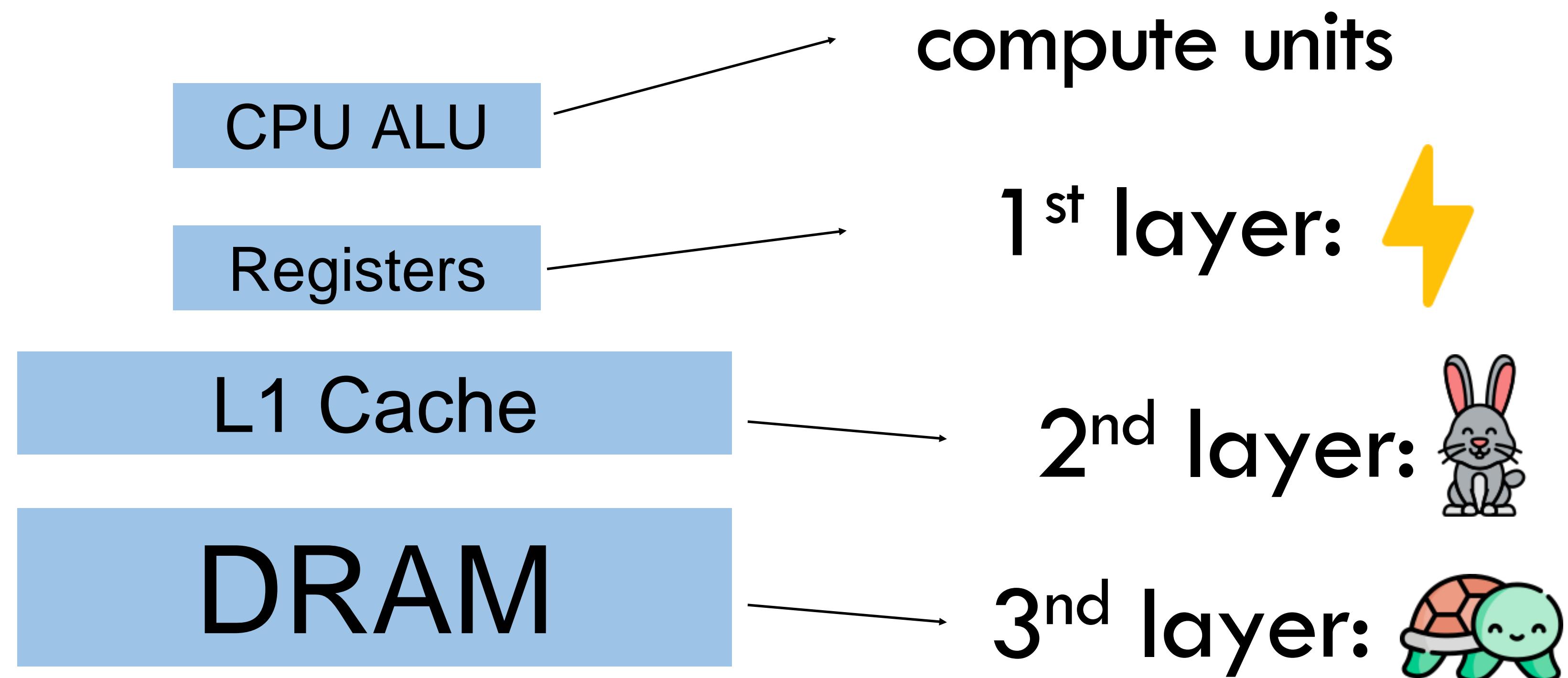
- What are the constraints?

#registers needed:

$v1*v3 + v2*v3 + v1*v2$

- Q: Why essentially can tiling reduce read cost?

Make it more complicated: Consider L1 cache

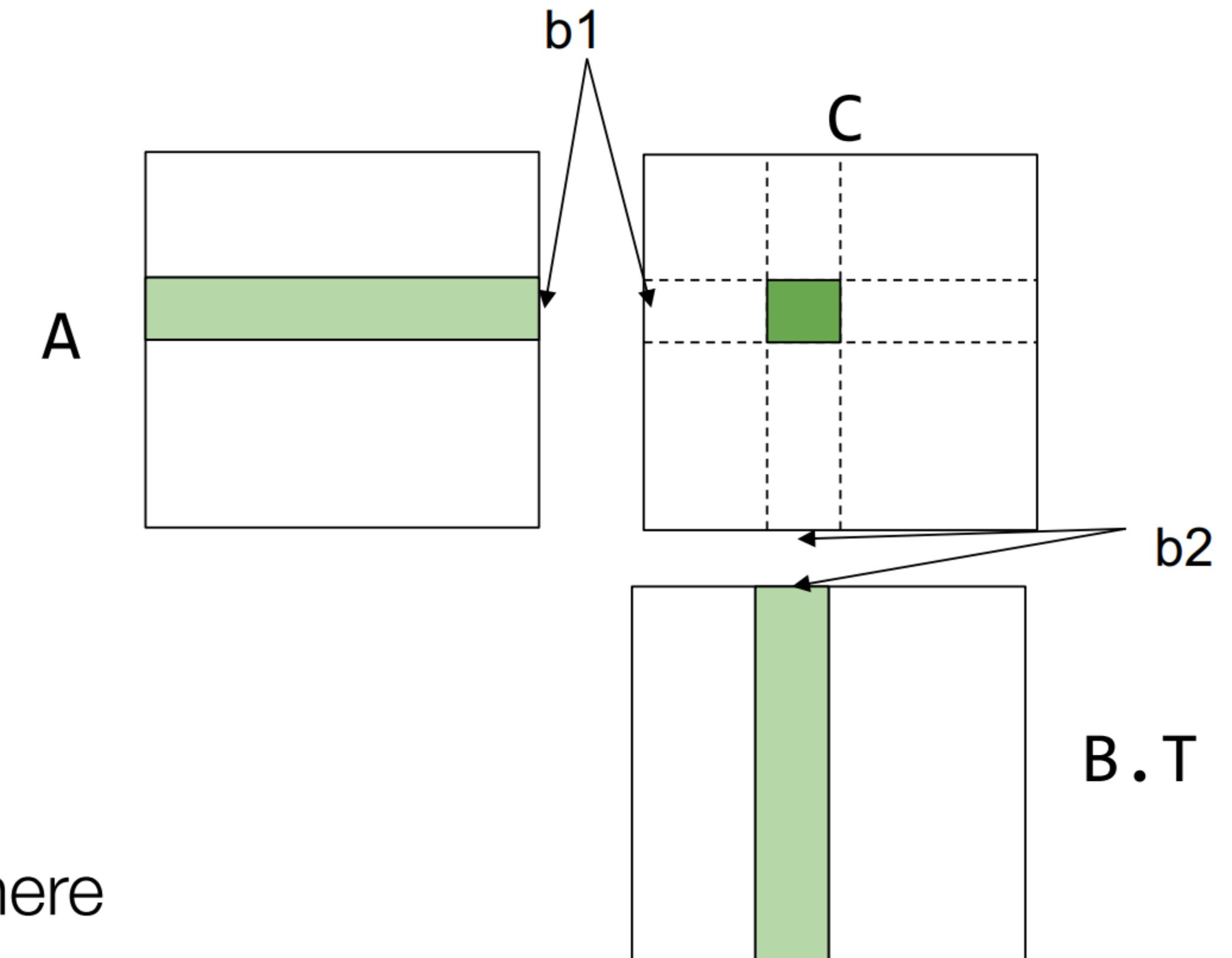


Cache-aware tiling

- We can further tile $[b_1][n]$ / $[b_2][n]$ using at the L1->register level
- What's the required condition?

```
dram float A[n/b1][b1][n];
dram float B[n/b2][b2][n];
dram float C[n/b1][n/b2][b1][b2];
for (int i = 0; i < n/b1; ++i) {
    l1cache float a[b1][n] = A[i];
    for (int j = 0; j < n/b2; ++j) {
        l1cache float b[b2][n] = B[j];
        C[i][j] = dot(a, b.T);
    }
}
```

Sub-procedure, can apply register tiling here



Cache-aware tiling

```
dram float A[n/b1][b1][n];
dram float B[n/b2][b2][n];
dram float C[n/b1][n/b2][b1][b2];
for (int i = 0; i < n/b1; ++i) {
    l1cache float a[b1][n] = A[i];
    for (int j = 0; j < n/b2; ++j) {
        l1cache b[b2][n] = B[j];
        C[i][j] = dot(a, b.T);
    }
}
```

Later we apply
register tiling
here

Data movement path:

- 1.Dram
- 2.Dram -> L1cache (cache tiling)
- 3.L1cache -> register (reg tiling)

Cache-aware tiling

```
dram float A[n/b1][b1][n];
dram float B[n/b2][b2][n];
dram float C[n/b1][n/b2][b1][b2];
for (int i = 0; i < n/b1; ++i) {
    l1cache float a[b1][n] = A[i];
    for (int j = 0; j < n/b2; ++j) {
        l1cache b[b2][n] = B[j];
        C[i][j] = dot(a, b.T);
    }
}
```

A's dram -> L1 cost:

$$n / b1 * n * b1 = n^2$$

B's dram -> L1 time cost:

$$n / b1 * n / b2 * b2 * n = n^3 / b1$$

Vs. previous untiled version?

s.t.

- $b1 * n + b2 * n < \text{L1 cache size}$

Putting Things Together

```
dram float A[n/b1][b1/v1][n][v1];  
dram float B[n/b2][b2/v2][n][v2];  
  
for (int i = 0; i < n/b1; ++i) {  
    l1cache float a[b1/v1][n][v1] = A[i];  
    for (int j = 0; j < n/b2; ++j) {  
        l1cache float b[b2/v2][n][v2] = B[j];  
        for (int x = 0; x < b1/v1; ++x)  
            for (int y = 0; y < b2/v2; ++y) {  
                register float c[v1][v2] = 0;  
                for (int k = 0; k < n; ++k) {  
                    register float ar[v1] = a[x][k][:];  
                    register float br[v2] = b[y][k][:];  
                    c += dot(ar, br.T);  
                }  
            }  
    }  
}
```

We set v3 = 1 (we know it does not matter)

Outside: cache tiling
Inside: register tiling

- Cache tiling using b1 and b2
- DRAM → L1 cache reads here
- Register tiling using v1 and v2
- L1 → register cache reads here

Putting Things Together

```
dram float A[n/b1][b1/v1][n][v1];
dram float B[n/b2][b2/v2][n][v2];

for (int i = 0; i < n/b1; ++i) {
    l1cache float a[b1/v1][n][v1] = A[i];
    for (int j = 0; j < n/b2; ++j) {
        l1cache float b[b2/v2][n][v2] = B[j];
        for (int x = 0; x < b1/v1; ++x)
            for (int y = 0; y < b2/v2; ++y) {
                register float c[v1][v2] = 0;
                for (int k = 0; k < n; ++k) {
                    register float ar[v1] = a[x][k][:];
                    register float br[v2] = b[y][k][:];
                    C += dot(ar, br.T)
                }
            }
    }
}
```

Outside: cache tiling
Inside: register tiling

Cost: dram -> l1

- $n/b1 * b1/v1 * n * v1 = n^2$
- $n/b1 * n/b2 * b2/v2 * n * v2 = n^2/v1$

Cost: l1 -> register:

- $n / b1 * n / b2 * b1 / v1 * b2 / v2 * n * v1 = n^3 / v2$
- $n^3 / v1$

In practice

- On CPUs: We have disk -> dram -> L2 -> L1 -> Register
- How to choose v1, v2, b1, b2, c1, c2?
- While we are reading from dram -> L2, can we concurrently read:
 - L2 -> L1
 - L1 -> register
- S.t. sizes of L2, L1, registers

Why tiling works: **reuse** loading

Access of A is independent of the dimension of j

```
float A[n][n];
float B[n][n];
float C[n][n];

C[i][j] = sum(A[i][k] * B[j][k], axis=k)
```

Tile the j dimension by v1 enables reuse of A for v1 times

Homework Candidate?

- Q: How to tile?

```
for n in range(0, N):  
    for co in range(0, CO):  
        for h in range(0, H):  
            for w in range(0, W):  
                for ci in range(0, CI):  
                    for kh in range(0, KH):  
                        for kw in range(0, KW):  
                            C[n,co,h,w] += A[n,co,h+kh,w+kw] x B[kh,kw,co,ci]
```

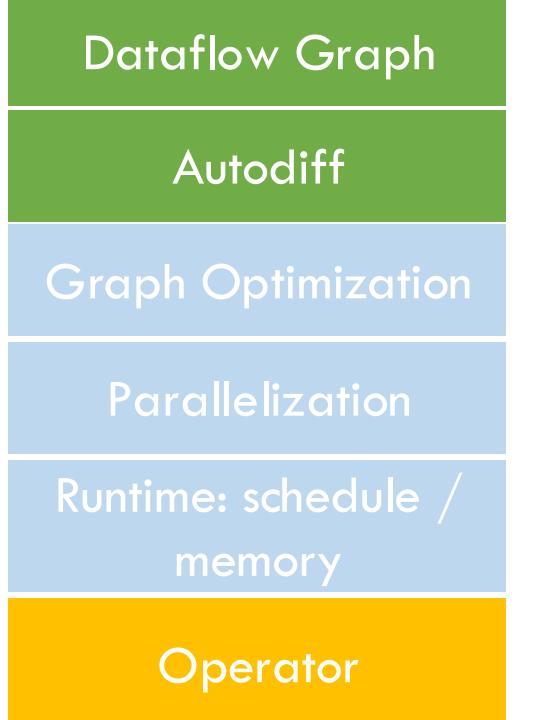


Simple spatial loops.

Stencil computation loops.

Reduction loop.

Reduction loops. But usually too small (≤ 5) for parallelization.



Where we are

- How to make operators fast in general?
 - Vectorize
 - Data layout
 - Parallelization (at the operator level)
 - Matmul-specific optimization
 - ? GPUs and accelerators
 - High-level Idea
 - The accelerator market

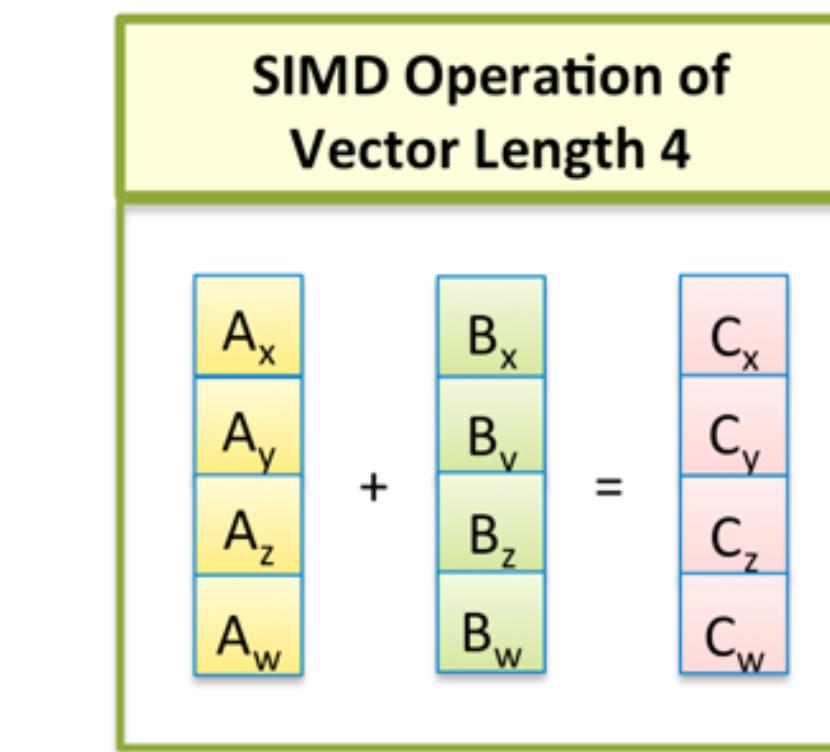
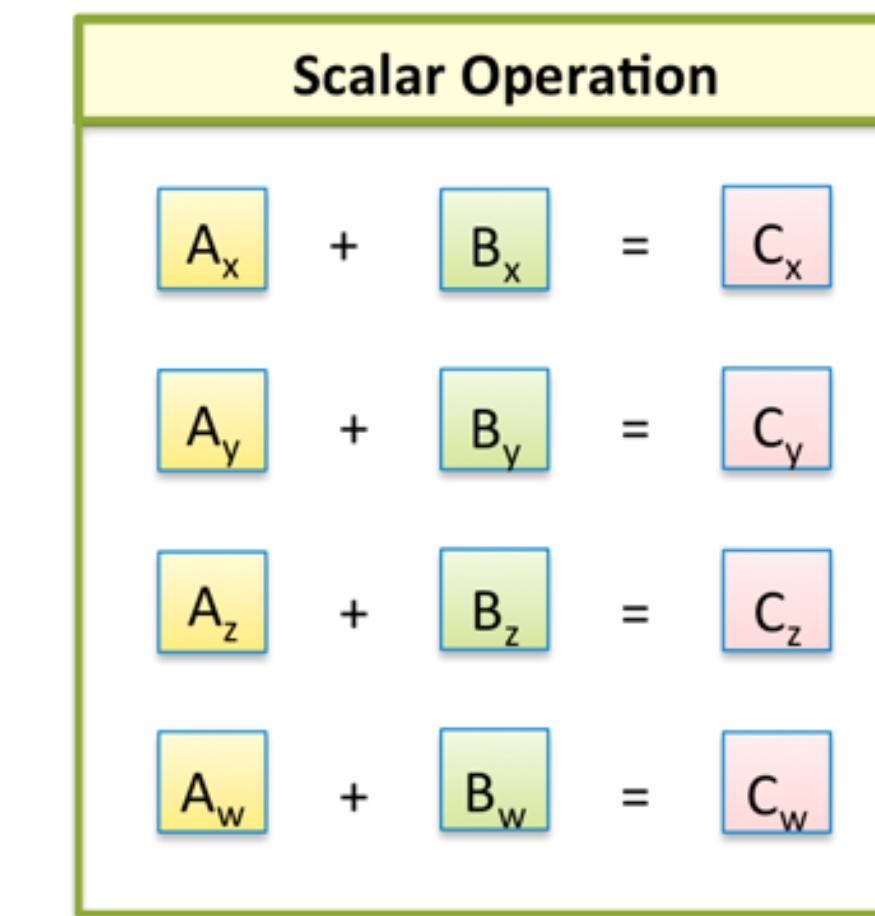
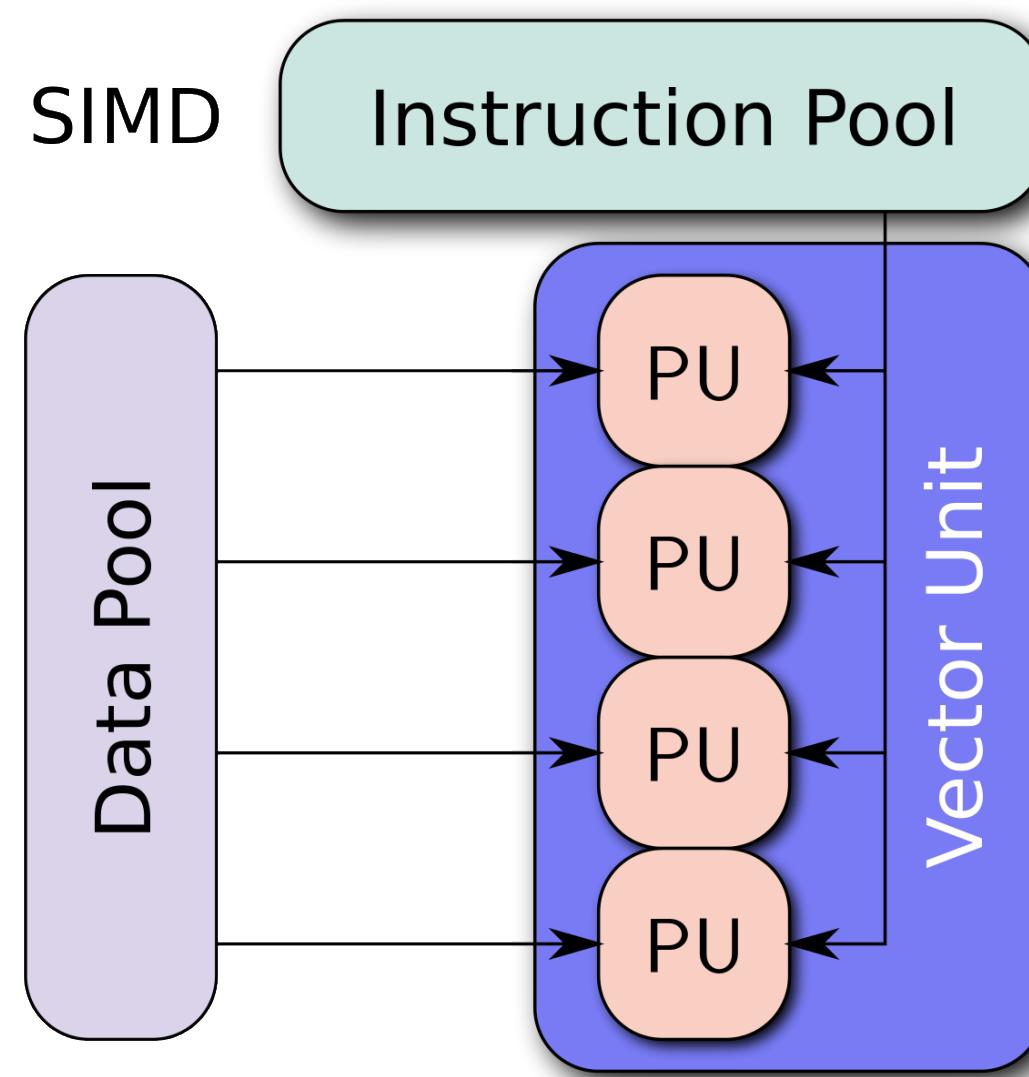
Recap CPU parallelization

- We can parallelize this loop using CPU threads
- == using many concurrent cores

```
#pragma omp parallel for
for (int i = 0; i < 64; ++i) {
    float4 a = load_float4(A + i*4);
    float4 b = load_float4(B + i*4);
    float4 c = add_float4(a, b);
    store_float4(C * 4, c);
}
```

Vectorized &
parallelized

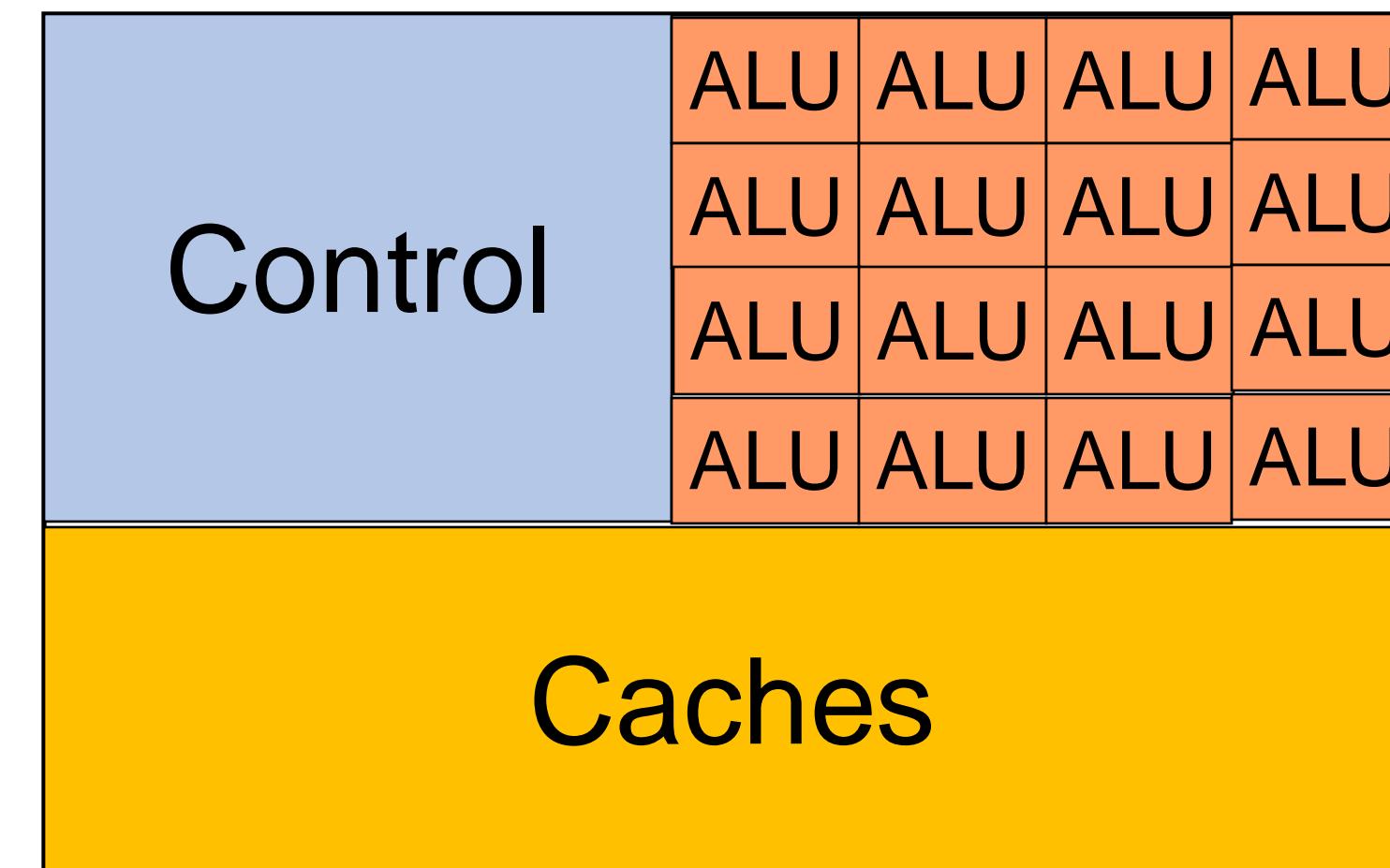
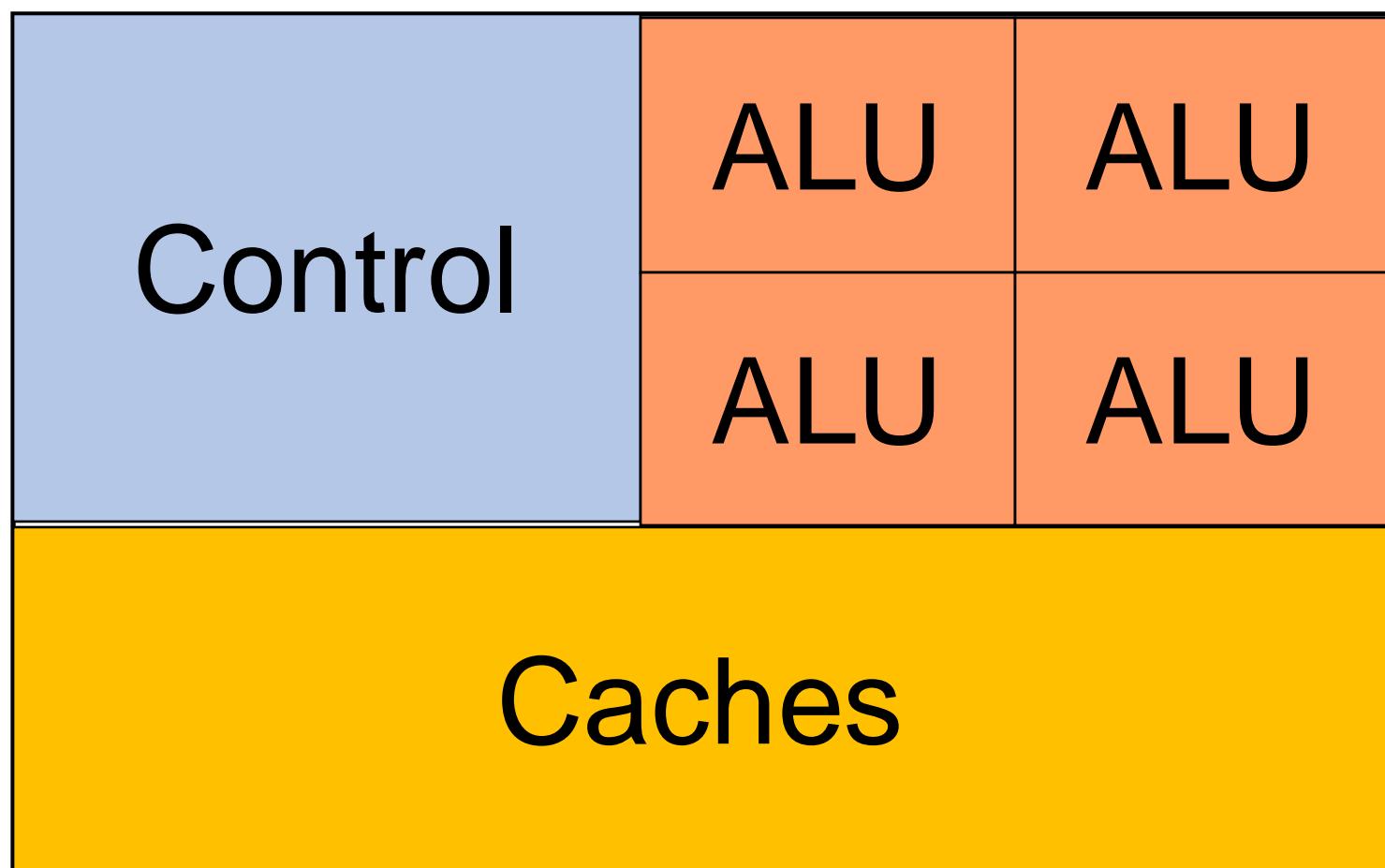
Single-Instruction Multiple-Data



Intel® Architecture currently has SIMD operations of vector length 4, 8, 16

Chip Design Trajectory: SIMD

If we're able to reduce size of ALU
(transistors) while keeping its power

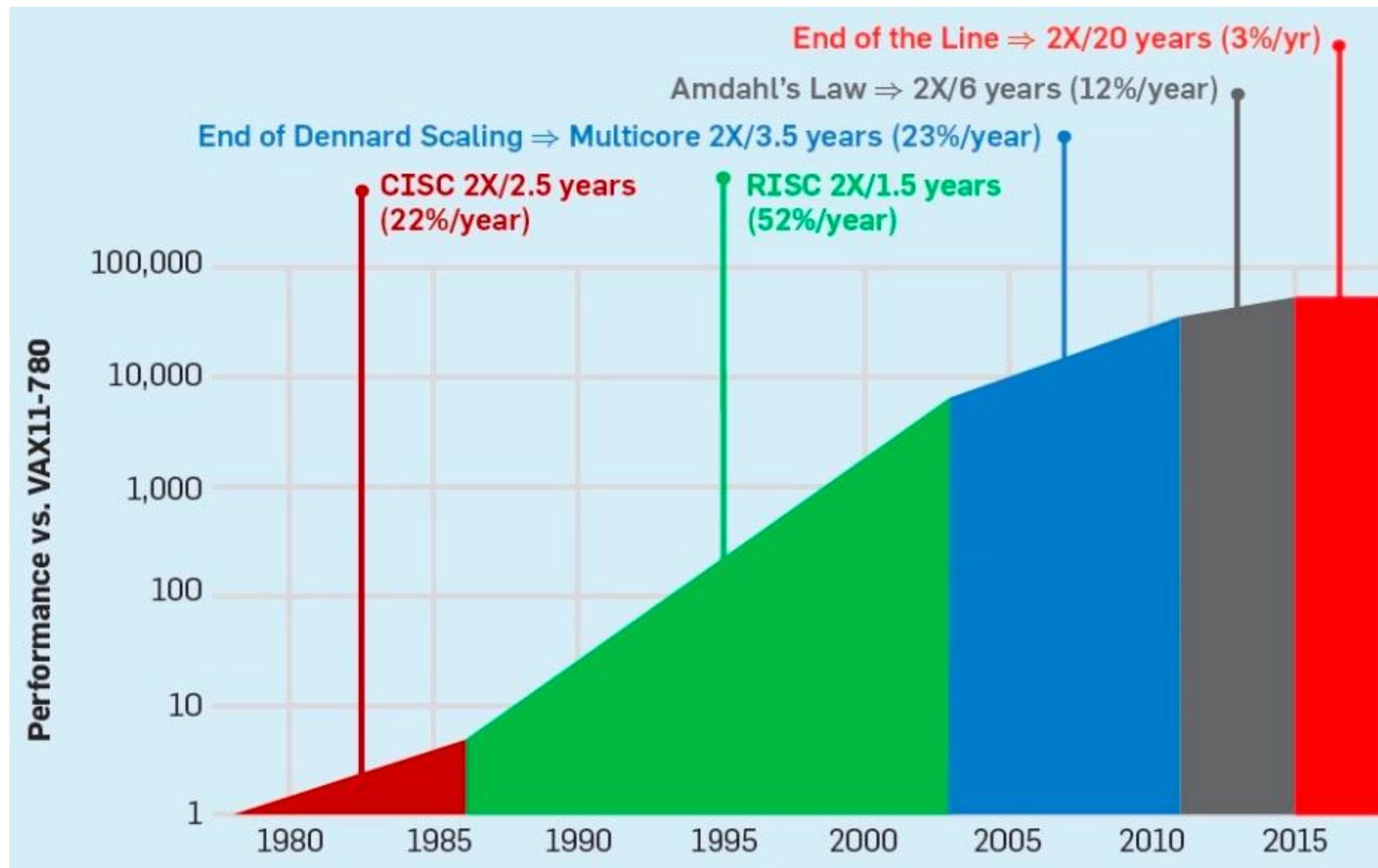


Chip Industry: 70nm -> 60nm -> 50nm -> ... ->?

- Problem: this is not substantiable; there are also power/heat issues when you put more ALUs in a limited area (s.t. physics limitations)



Chip Industry: Moore's Law Comes to an End

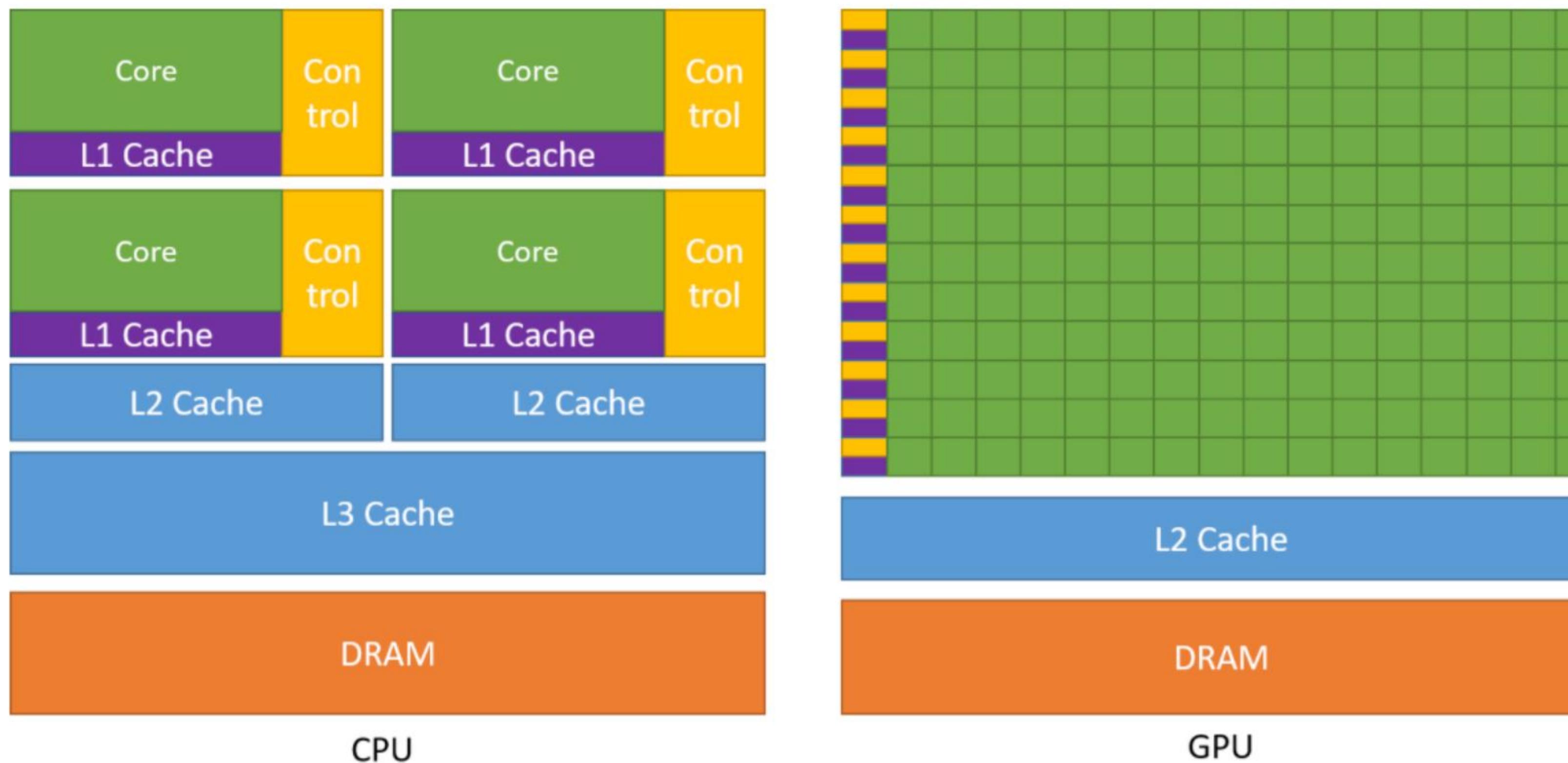


Option 1: Go to the quantum world



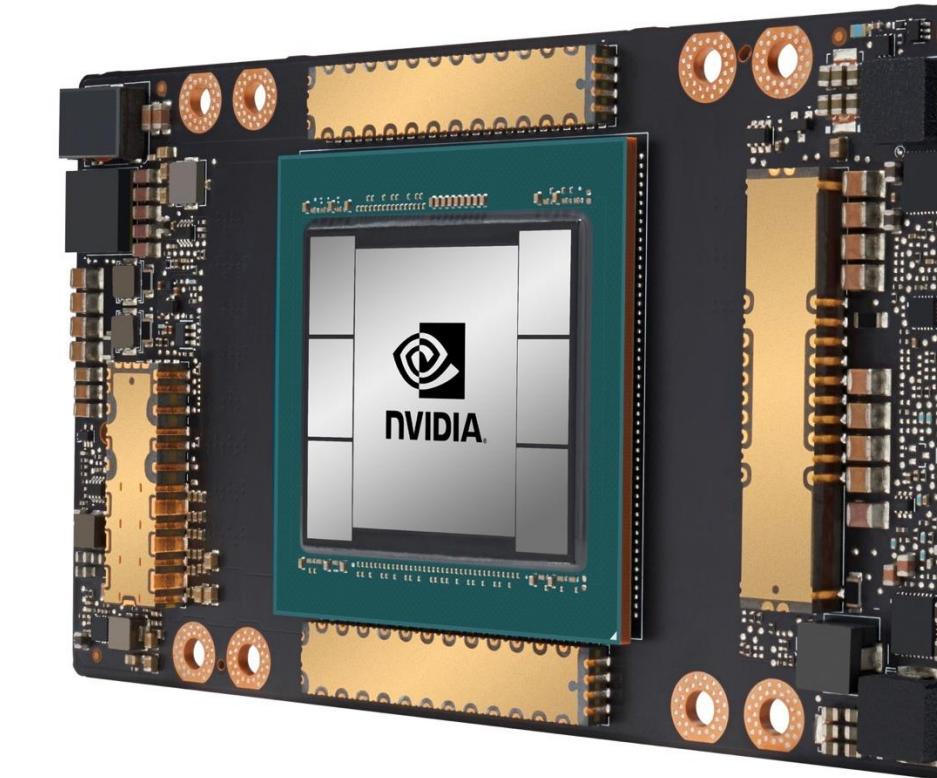
Option 2: Specialized hardware

Idea: How about we use a lot of weak/specialized cores



Hardware Accelerators: GPUs

- **Graphics Processing Unit (GPU)**: Tailored for matrix/tensor ops
- Basic idea: Use tons of ALUs (but weak and more specialized); massive data parallelism (SIMD on steroids);
- Popularized by NVIDIA in early 2000s for video games, graphics, and multimedia; now ubiquitous in DL
- CUDA released in 2007; later wrapper APIs on top: CuDNN, CuSparse, CuDF (RapidsAI), NCCL, etc.



Other Hardware Accelerators

- E.g.
 - Tensor Processing Unit (TPU)
 - An “application-specific integrated circuit” (ASIC) created by Google in mid 2010s; used for AlphaGo
 - E.g.
 - B200 (projected release 2025): fp4 / fp8 Tensorcore
 - E.g.
 - M3 max: mixing tensorcore and normal core



What Does It Mean by “Specialized” In accelerator world

In General:

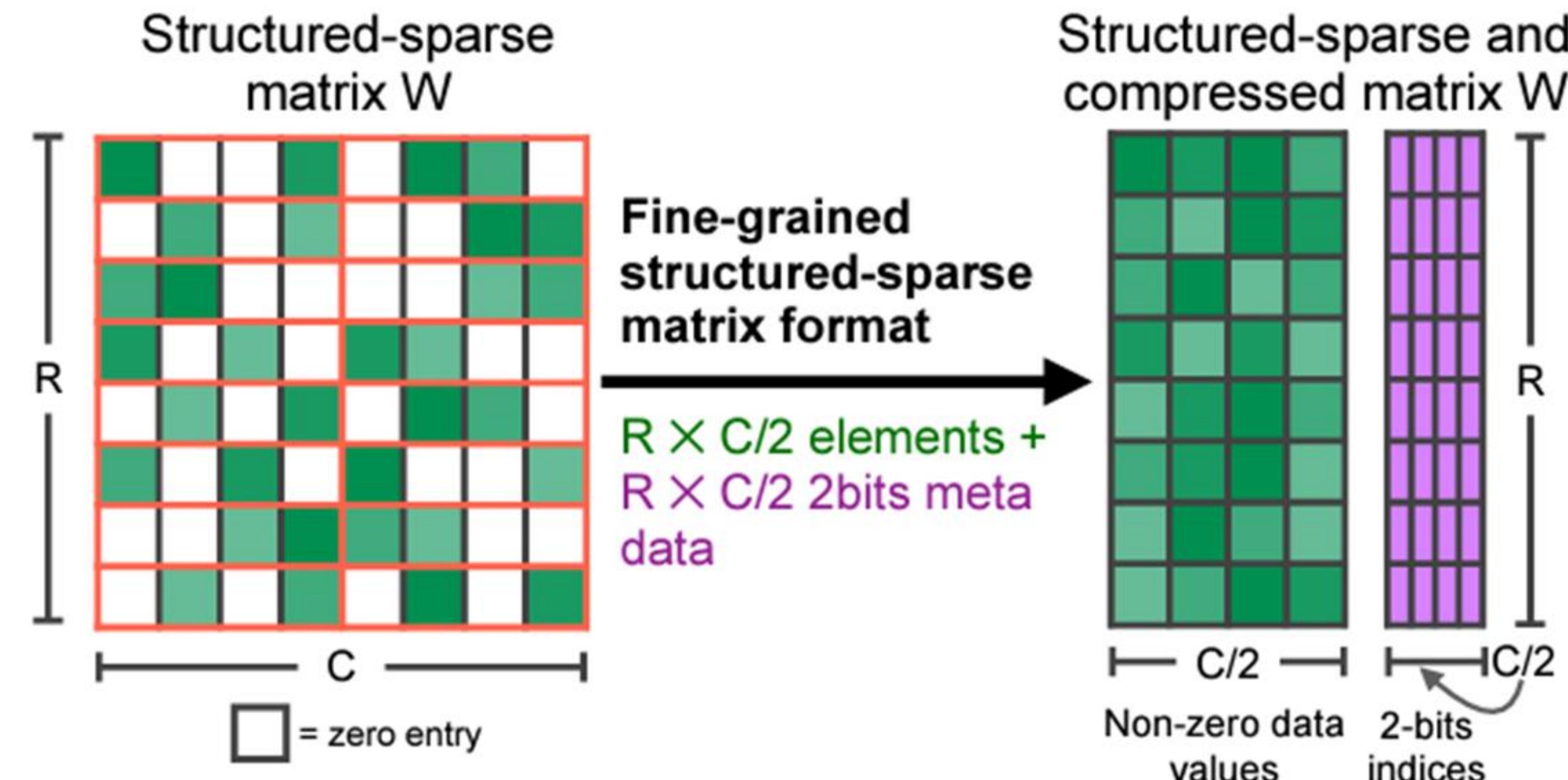
- Functionality-specialized:
 - Can only compute certain computations: matmul, sparsity
 - Mixing specialized cores with versatile cores
- Reduce precision
 - Floating point operations: fp32, fp16, fp8, int8, int4, ...
- Tune the distribution of different components for specific workloads
 - SRAM, cache, registers, etc.

Case Study 1: Nvidia GPU Specification

FP64 Tensor Core	67 teraFLOPS
FP32	67 teraFLOPS
TF32 Tensor Core*	989 teraFLOPS
BFLOAT16 Tensor Core*	1,979 teraFLOPS
FP16 Tensor Core*	1,979 teraFLOPS
FP8 Tensor Core*	3,958 teraFLOPS
INT8 Tensor Core*	3,958 TOPS
GPU Memory	80GB
GPU Memory Bandwidth	3.35TB/s
Decoders	7 NVDEC 7 JPEG
Max Thermal Design Power (TDP)	Up to 700W (configurable)
Multi-Instance GPUs	Up to 7 MIGS @ 10GB each
Form Factor	SXM
Interconnect	NVIDIA NVLink™: 900GB/s PCIe Gen5: 128GB/s
Server Options	NVIDIA HGX H100 Partner and NVIDIA-Certified Systems™ with 4 or 8 GPUs NVIDIA DGX H100 with 8 GPUs
NVIDIA AI Enterprise	Add-on

* With sparsity

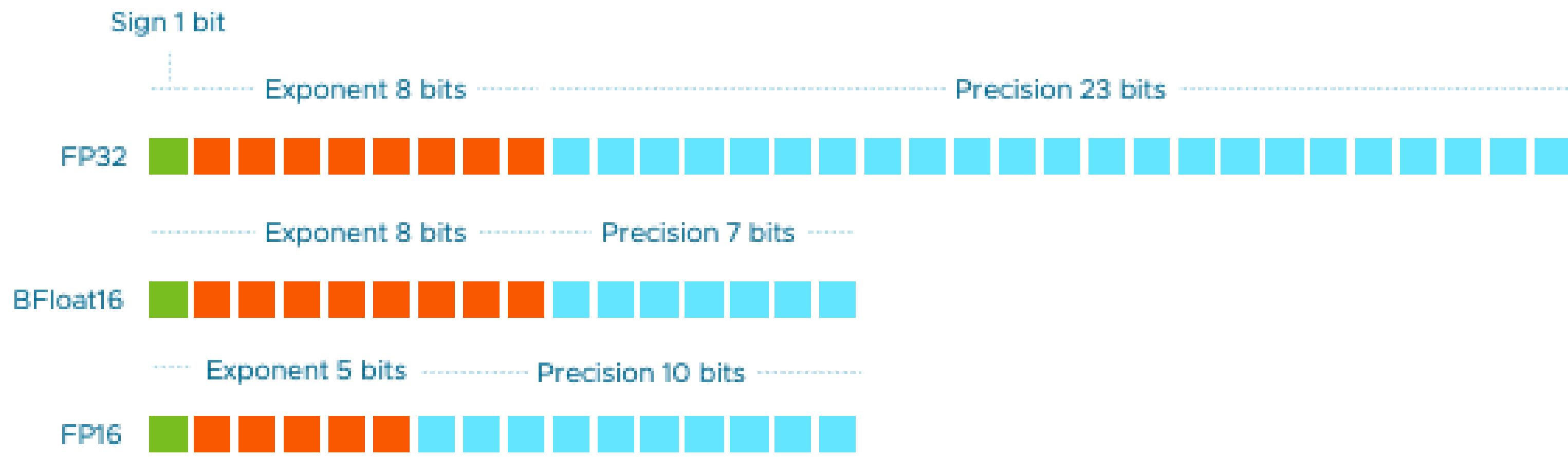
Case Study 1: Nvidia GPU Specification



Case Study 1: Nvidia GPU Specification

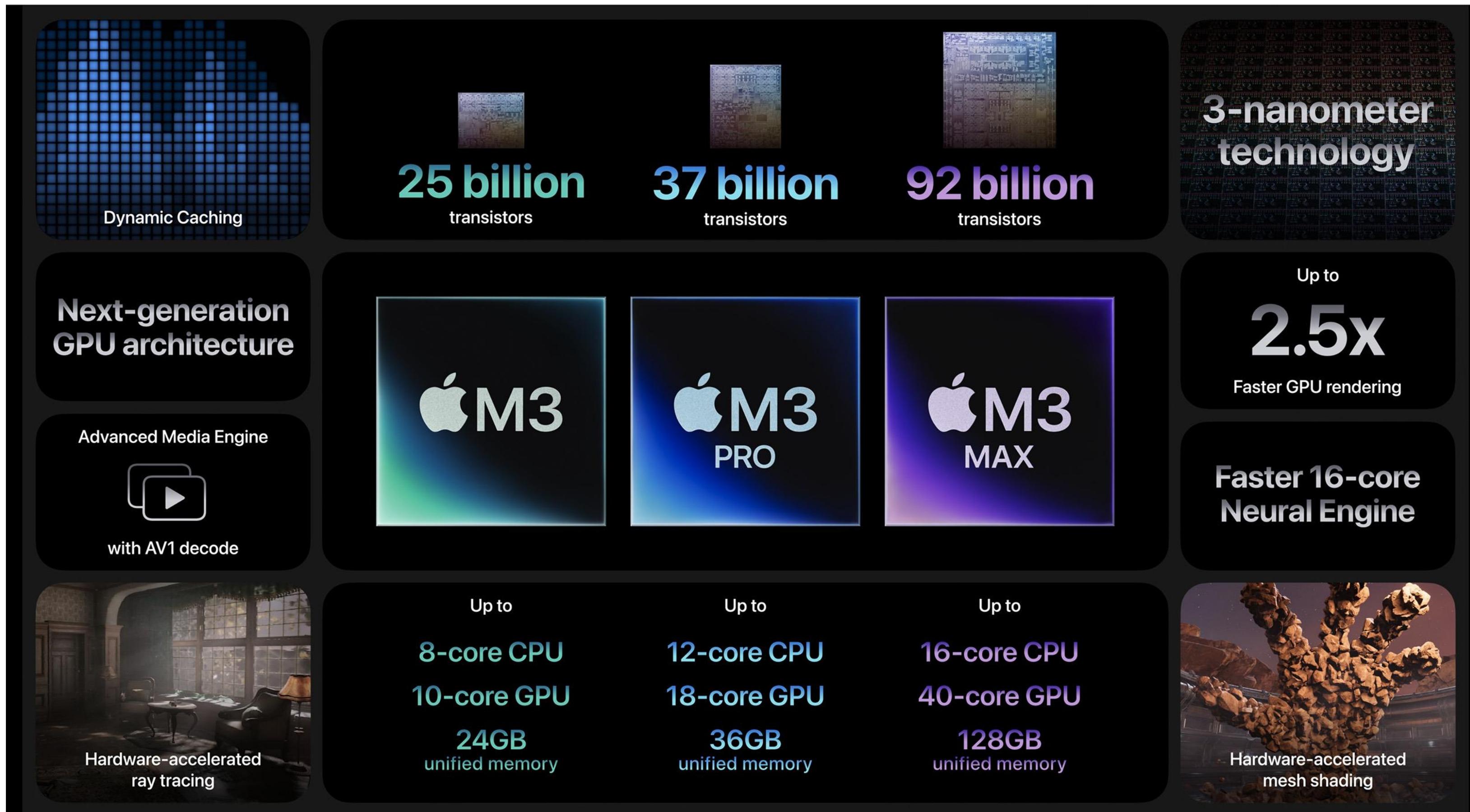
FP64 Tensor Core	67 teraFLOPS
FP32	67 teraFLOPS
TF32 Tensor Core*	989 teraFLOPS
BFLOAT16 Tensor Core*	1,979 teraFLOPS
FP16 Tensor Core*	1,979 teraFLOPS
FP8 Tensor Core*	3,958 teraFLOPS
INT8 Tensor Core*	3,958 TOPS
GPU Memory	80GB
GPU Memory Bandwidth	3.35TB/s
Decoders	7 NVDEC 7 JPEG
Max Thermal Design Power (TDP)	Up to 700W (configurable)
Multi-Instance GPUs	Up to 7 MIGS @ 10GB each
Form Factor	SXM
Interconnect	NVIDIA NVLink™: 900GB/s PCIe Gen5: 128GB/s
Server Options	NVIDIA HGX H100 Partner and NVIDIA-Certified Systems™ with 4 or 8 GPUs NVIDIA DGX H100 with 8 GPUs
NVIDIA AI Enterprise	Add-on

* With sparsity

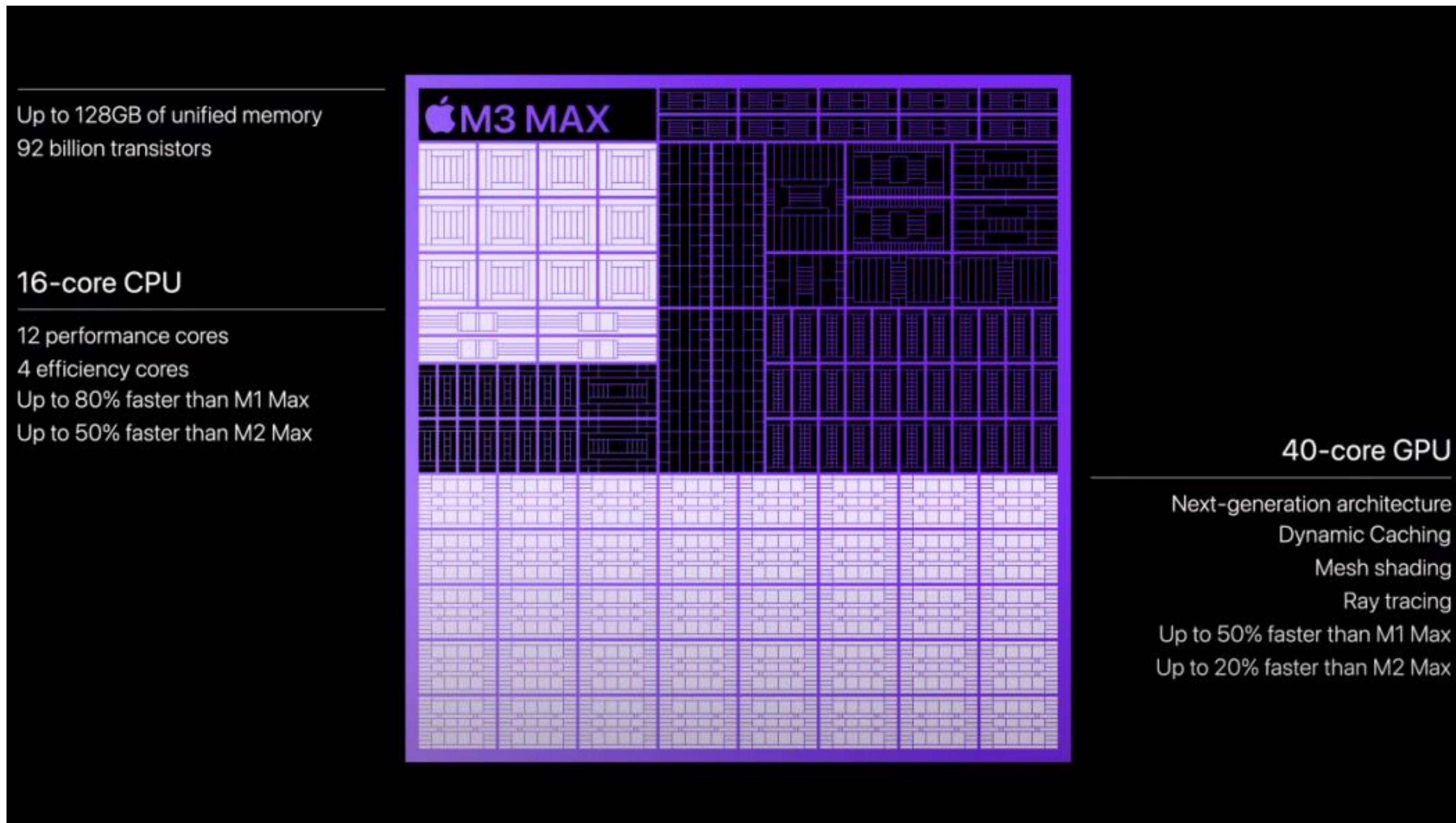


Question: why this could work in ML programs?

Case Study 2: Apple Silicon



Case Study 2: Apple Silicon Revealed

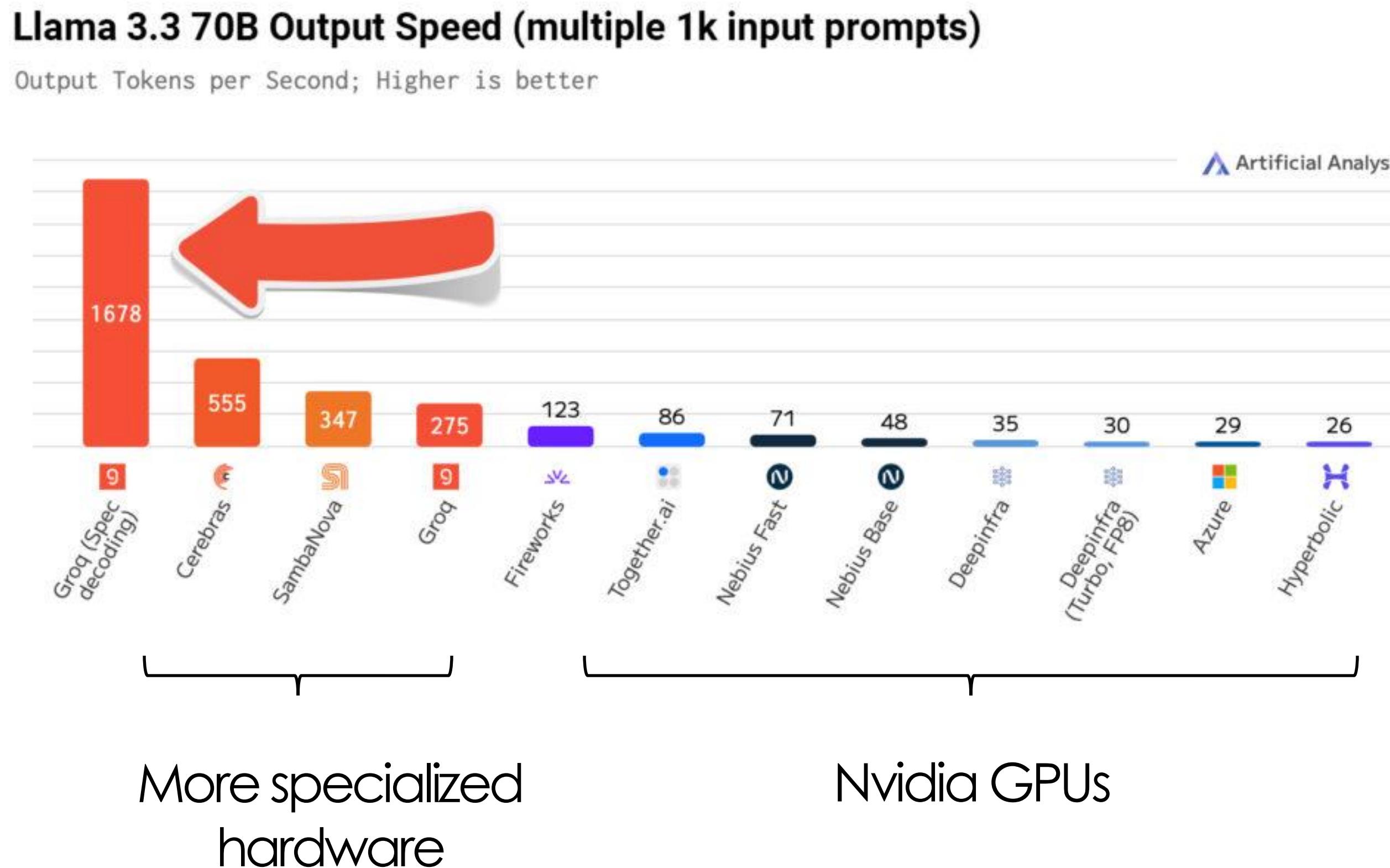


Case Study 3: Leading Chip Startups



Case Study 3: Groq

Question: How did Groq achieve that?



Case Study 3: Groq

GroqCard™



Card Specifications

Form Factor

Dual width, full height, 3/4 length PCI Express Gen4 x16 adapter

Performance

Up to 750 TOPs, 188 TFLOPs (INT8, FP16 @900 MHz)

Memory

230 MB SRAM per chip

Up to 80 TB/s on-die memory bandwidth

Chip Scaling

Up to 9 RealScale™ chip-to-chip connectors

Numerics

INT8, INT16, INT32 & TruePoint™ technology

MXM: FP32

VXM: FP16, FP32

Power

Max: 375W; TDP: 275 ; Typical: 240W

Data Center GPU	NVIDIA Tesla V100	NVIDIA A100	NVIDIA H100
GPU Architecture	NVIDIA Volta	NVIDIA Ampere	NVIDIA Hopper
Compute Capability	7.0	8.0	9.0
Threads / Warp	32	32	32
Max Warps / SM	64	64	64
Max Threads / SM	2048	2048	2048
Max Thread Blocks (CTAs) / SM	32	32	32
Max Thread Blocks / Thread Block Clusters	NA	NA	16
Max 32-bit Registers / SM	65536	65536	65536
Max Registers / Thread Block (CTA)	65536	65536	65536
Max Registers / Thread	255	255	255
Max Thread Block Size (# of threads)	1024	1024	1024
FP32 Cores / SM	64	64	128
Ratio of SM Registers to FP32 Cores	1024	1024	512
Shared Memory Size / SM	Configurable up to 96 KB	Configurable up to 164 KB	Configurable up to 228 KB

Case Study 3: Groq

- Recall

```
dram float A[n][n], B[n][n], C[n][n];
for (int i = 0; i < n; ++i) {
    for (int j = 0; j < n; ++j) {
        register float c = 0;
        for (int k = 0; k < n; ++k) {
            register float a = A[i][k];
            register float b = B[j][k];
            c += a * b;
        }
        C[i][j] = c;
    }
}
```

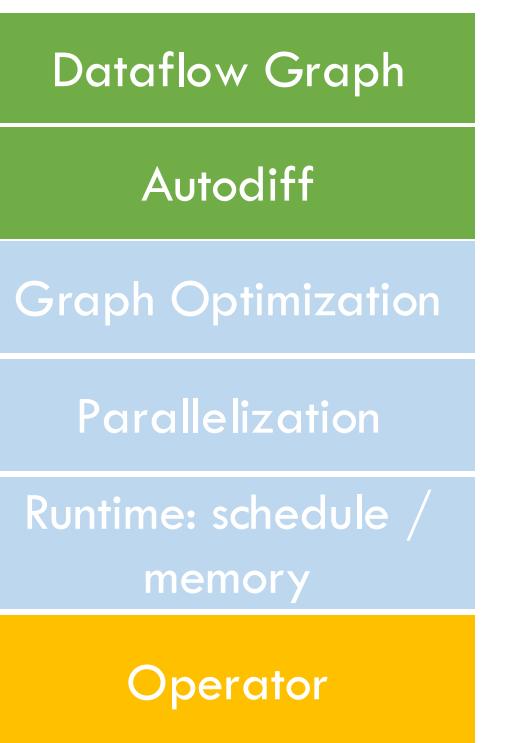
Take-home Exercise

- Study B100 specification and compare it to H100
 - How nvidia claims another 2x from H100 -> B100?
 - How about B200?

Economic Question

Question: What is Nvidia's Moat?

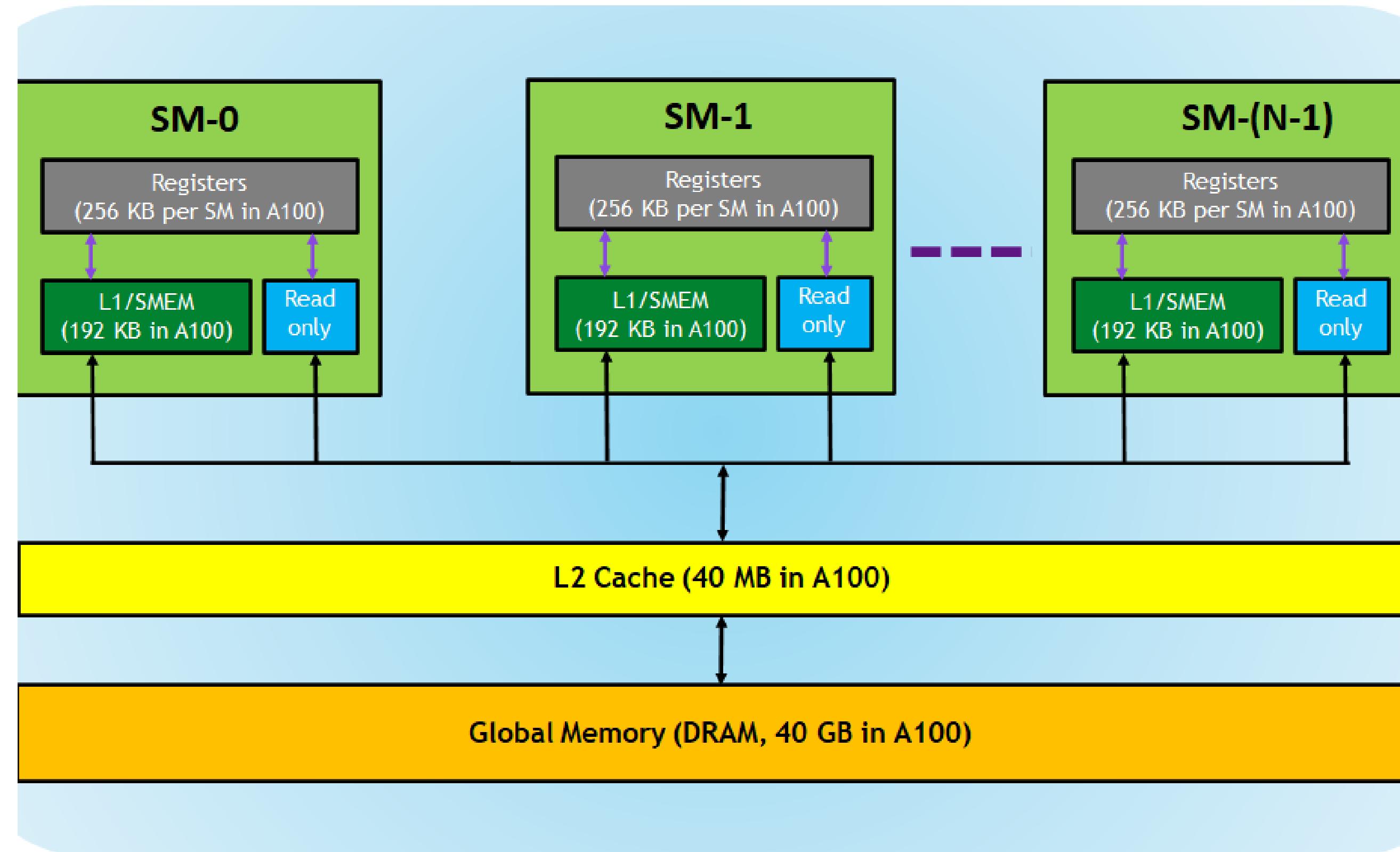




Next: GPU and CUDA

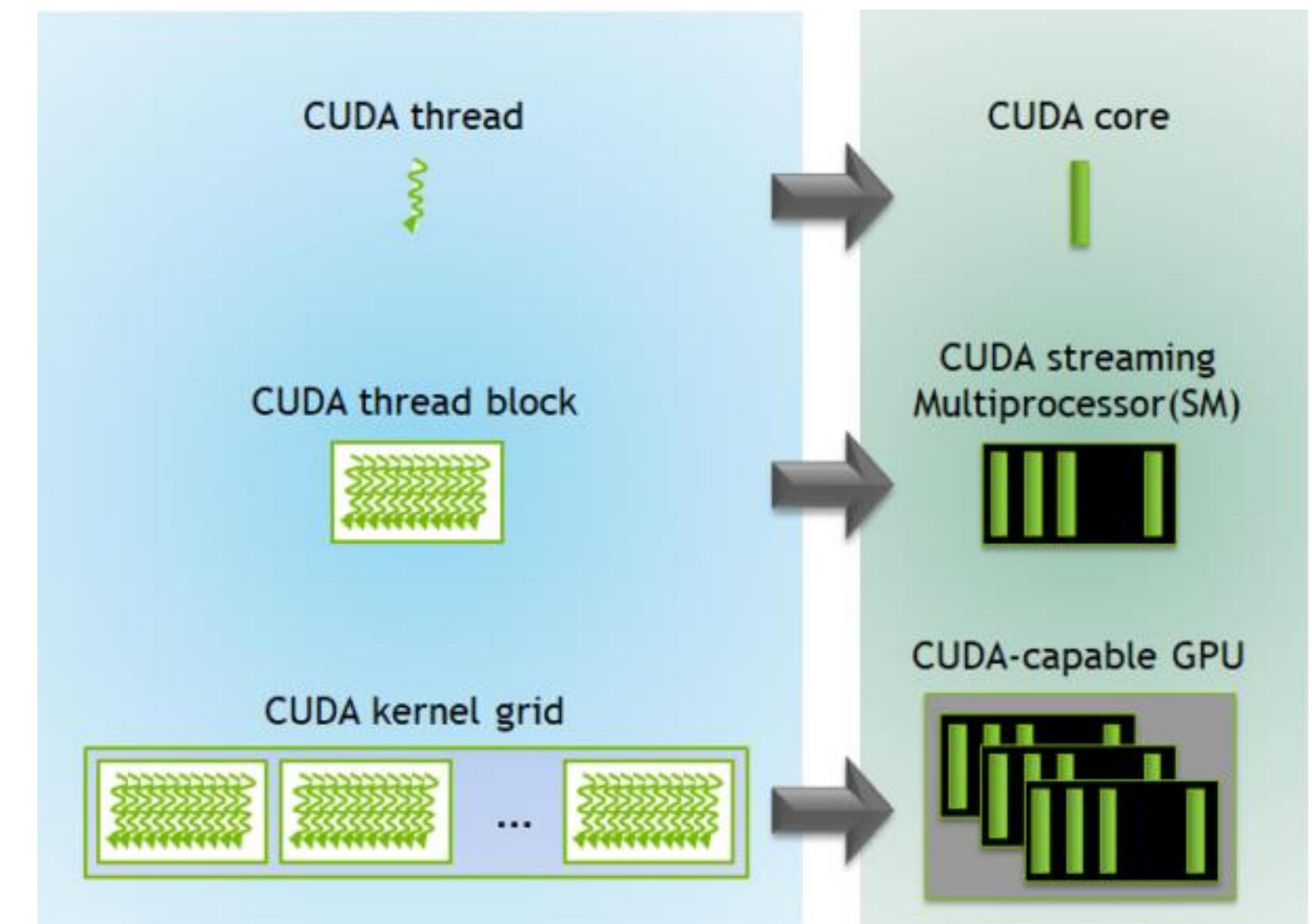
- Basic concepts and Architecture
- Programming abstraction
- Case study: Matmul

GPU Overview



Threads, Blocks, Grids

- Threads: smallest units to process a chunk of data
- Blocks: A group of threads that share memory
- Grid: A collection of blocks that execute the same kernel



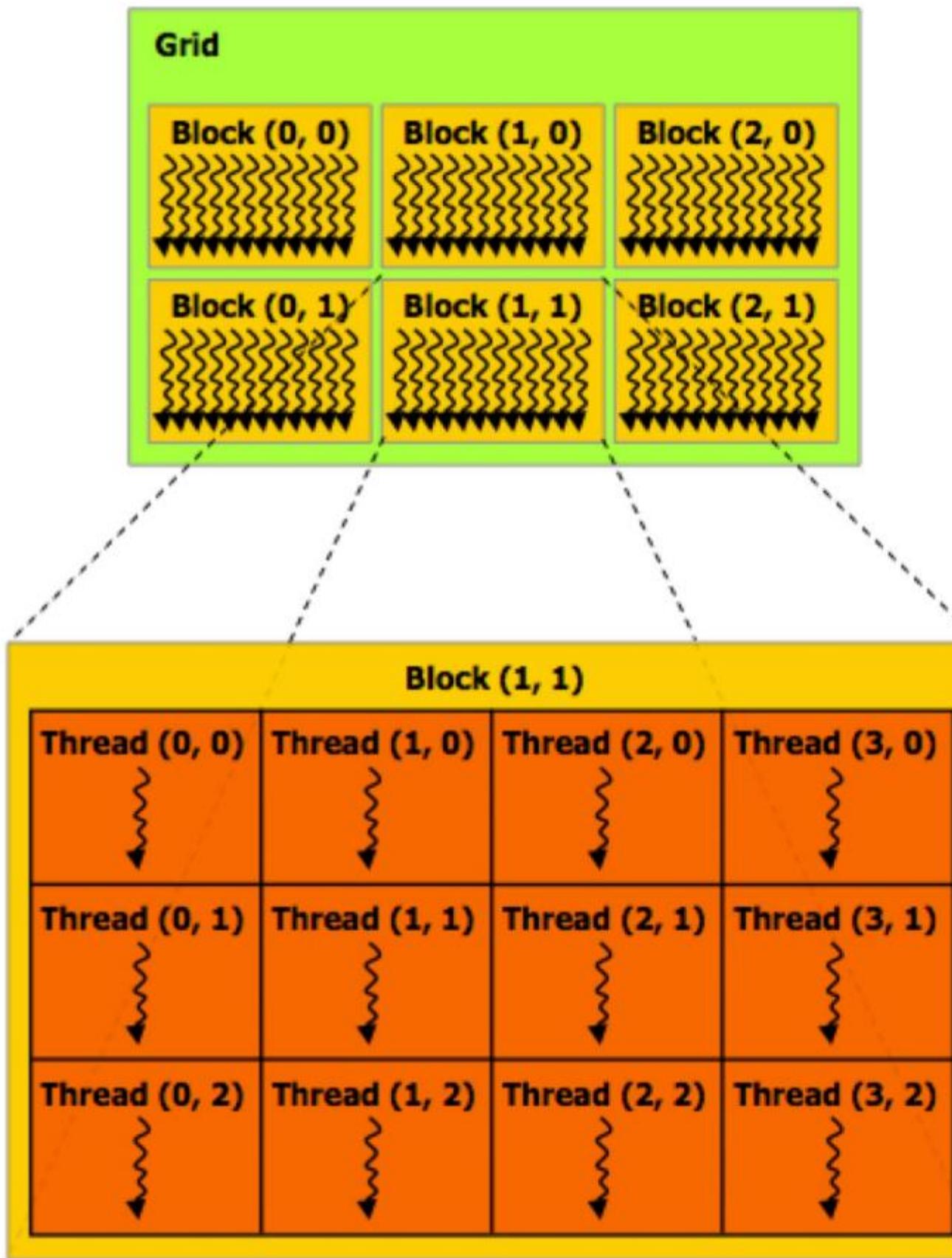
How many SMs/Threads we have?

- V100 (2018 - Now): 80 SMs, 2048 threads/SM
- A100 (2020 - Now): 108 SMs, 2048 threads/SM
- H100 (2022 - Now): 144 SMs, 2048 threads/SM
- B100 (2025 -): go surveying the numbers

CUDA

- Introduced in 2007 with NVIDIA Tesla architecture
- C-like languages for programming GPUs
- CUDA's design matches the grid/block/thread concepts in GPUs

CUDA Programs contain A Hierarchy of Threads

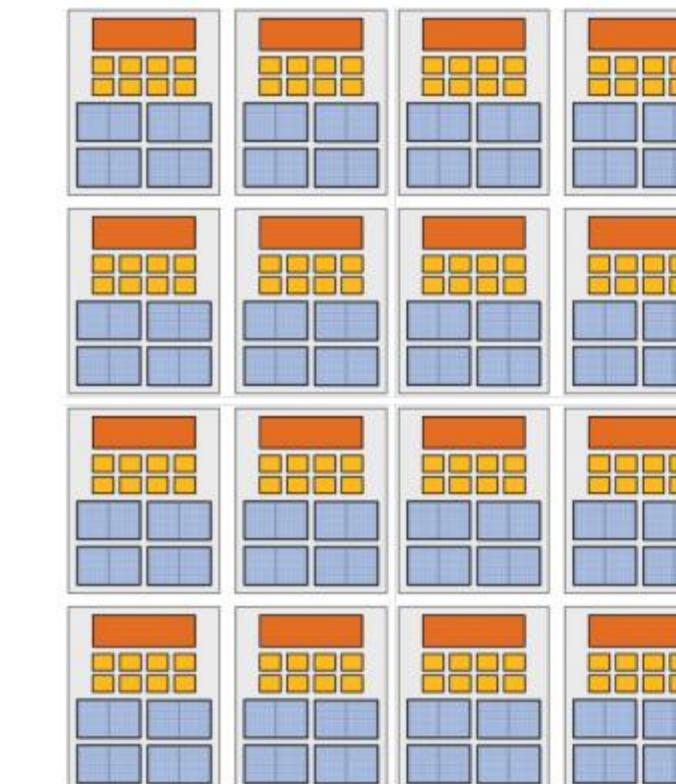


```
const int Nx = 12;
const int Ny = 6;

dim3 threadsPerBlock(4, 3, 1);
dim3 numBlocks(Nx/threadsPerBlock.x,
               Ny/threadsPerBlock.y, 1);

// assume A, B, C are allocated Nx x Ny float arrays
// this call will trigger execution of 72 CUDA threads:
// 6 thread blocks of 12 threads each
matrixAdd<<<numBlocks, threadsPerBlock>>>(A, B, C);
```

Run on
CPU



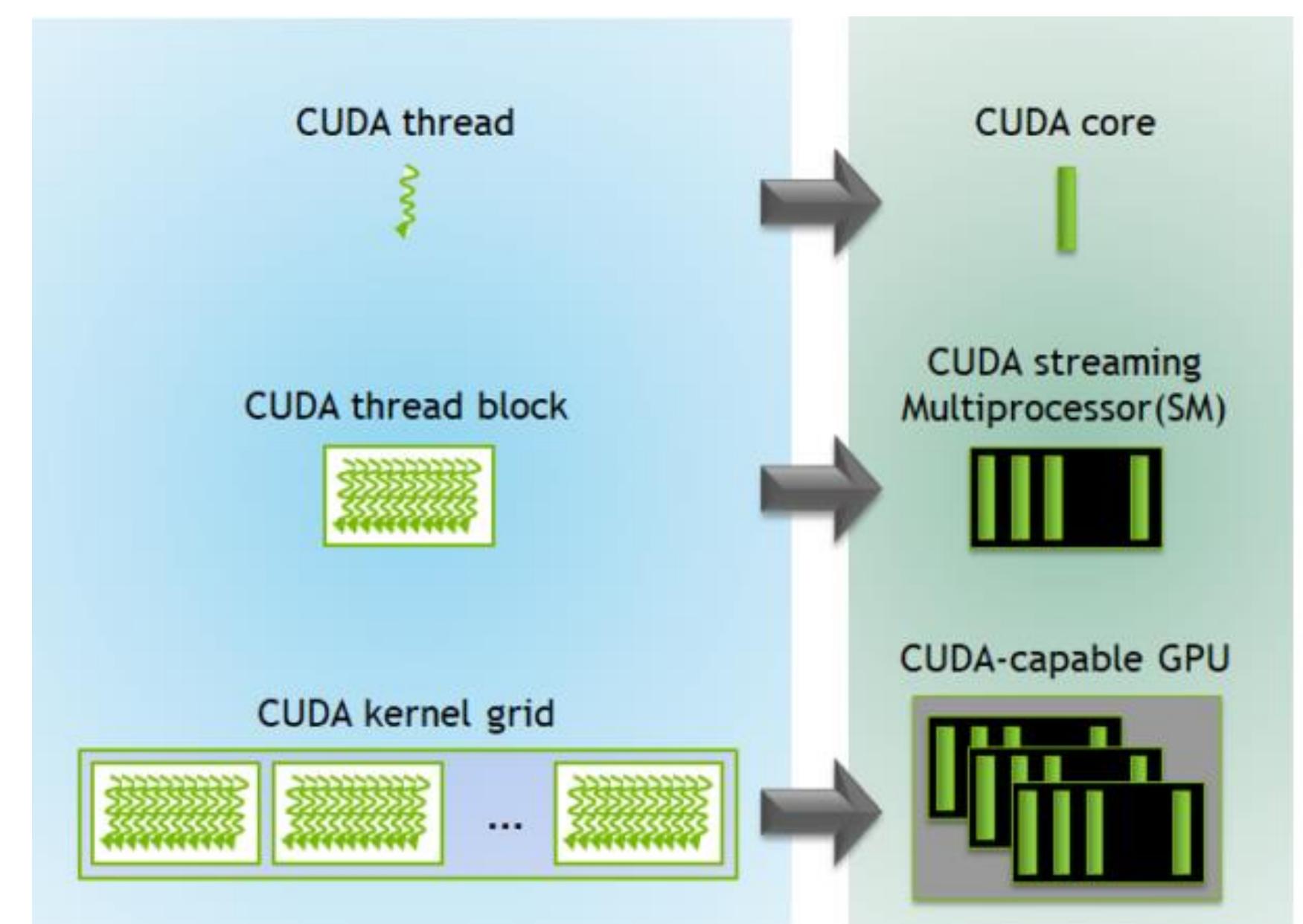
How Many threads/Blocks it runs on?

```
const int Nx = 12;
const int Ny = 6;

dim3 threadsPerBlock(4, 3, 1);
dim3 numBlocks(Nx/threadsPerBlock.x,
               Ny/threadsPerBlock.y, 1);

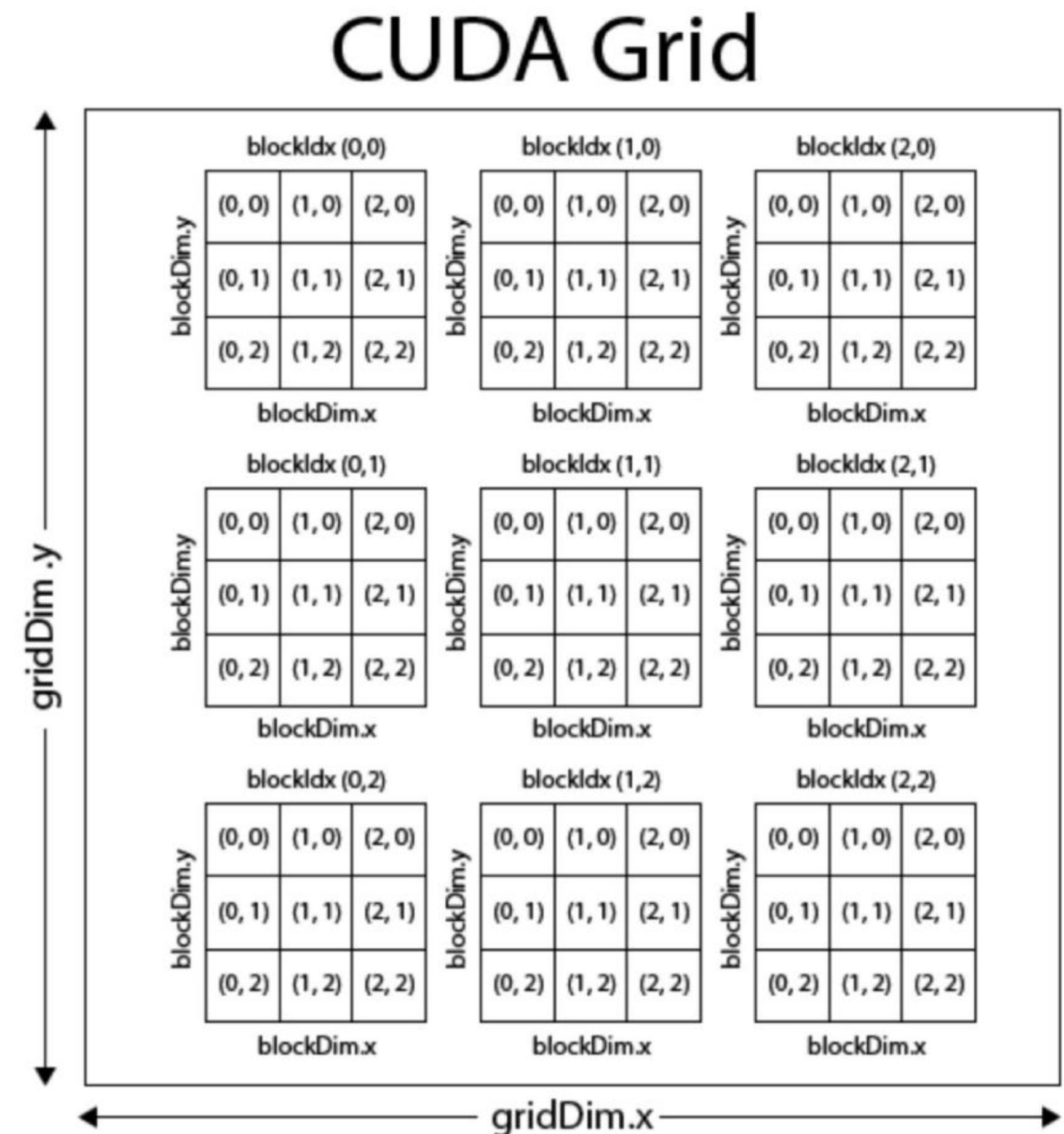
// assume A, B, C are allocated Nx x Ny float arrays

// this call will trigger execution of 72 CUDA threads:
// 6 thread blocks of 12 threads each
matrixAdd<<<numBlocks, threadsPerBlock>>>(A, B, C);
```



Grid, Block, and Thread

- GridDim: The dimensions of the grid
- blockIdx: The block index within the grid
- blockDim: The dimensions of a block
- threadIdx: The thread index within a block
- What About GridId?
- What about threadDim?



An Example CUDA Program

```
const int Nx = 12;
const int Ny = 6;

dim3 threadsPerBlock(4, 3, 1);
dim3 numBlocks(Nx/threadsPerBlock.x,
               Ny/threadsPerBlock.y, 1);

// assume A, B, C are allocated Nx x Ny float arrays

// this call will cause execution of 72 threads
// 6 blocks of 12 threads each
matrixAddDoubleB<<<numBlocks, threadsPerBlock>>>(A, B, C);
```

- “launch a grid of CUDA thread blocks” Call returns when all threads have terminated

```
__device__ float doubleValue(float x)
{
    return 2 * x;
}

// kernel definition
__global__ void matrixAddDoubleB(float A[Ny][Nx],
                                 float B[Ny][Nx],
                                 float C[Ny][Nx])
{
    int i = blockIdx.x * blockDim.x + threadIdx.x;
    int j = blockIdx.y * blockDim.y + threadIdx.y;

    C[j][i] = A[j][i] + doubleValue(B[j][i]);
}
```

- **__global__** denotes a CUDA kernel function runs on GPU
- Each thread indexes its data using blockIdx, blockDim, threadIdx and execute the compute

Separation CPU and GPU Execution

```
const int Nx = 12;
const int Ny = 6;

dim3 threadsPerBlock(4, 3, 1);
dim3 numBlocks(Nx/threadsPerBlock.x,
               Ny/threadsPerBlock.y, 1);

// assume A, B, C are allocated Nx x Ny float arrays

// this call will cause execution of 72 threads
// 6 blocks of 12 threads each
matrixAddDoubleB<<<numBlocks, threadsPerBlock>>>(A, B, C);
```

```
__device__ float doubleValue(float x)
{
    return 2 * x;
}

// kernel definition
__global__ void matrixAddDoubleB(float A[Ny][Nx],
                                 float B[Ny][Nx],
                                 float C[Ny][Nx])
{
    int i = blockIdx.x * blockDim.x + threadIdx.x;
    int j = blockIdx.y * blockDim.y + threadIdx.y;

    C[j][i] = A[j][i] + doubleValue(B[j][i]);
}
```

- Host code: serial execution on CPU
- Device code: SIMD parallel execution on GPUs

#Threads is Explicit and Static in Programs

```
const int Nx = 11; // not a multiple of threadsPerBlk.x
const int Ny = 5; // not a multiple of threadsPerBlk.y

dim3 threadsPerBlk(4, 3, 1);
dim3 numBlocks(3, 2, 1);

// assume A, B, C are allocated Nx x Ny float arrays
// this call will trigger execution of 72 CUDA threads:
// 6 thread blocks of 12 threads each
matrixAdd<<<numBlocks, threadsPerBlk>>>(A, B, C);
```

Developers to:

- To provide CPU/GPU code separation
- Statically declare blockDim, shapes.
- Map data to blocks/threads
- Check boundary conditions

```
// kernel definition
__global__ void matrixAdd(float A[Ny][Nx],
                          float B[Ny][Nx],
                          float C[Ny][Nx])
{
    int i = blockIdx.x * blockDim.x + threadIdx.x;
    int j = blockIdx.y * blockDim.y + threadIdx.y;
    // guard against out of bounds array access
    if (i < Nx && j < Ny)
        C[j][i] = A[j][i] + B[j][i];
}
```

SIMD Constraints: how to handle control flow?

SIMD requires all ALUs/Core Must proceed in the same pace

```
// kernel definition
__global__ void f(float A[N])
{
    int i = blockIdx.x * blockDim.x + threadIdx.x;
    float x = A[i];
    if (x > 0) {
        x = 2.0f * x;
    } else {
        x = exp(x, 5.0f);
    }
    A[i] = x;
}
```