



<https://haoailab.com/cse291-s26/>

# CSE/DSC 291: Deep Learning Systems Spring 2026

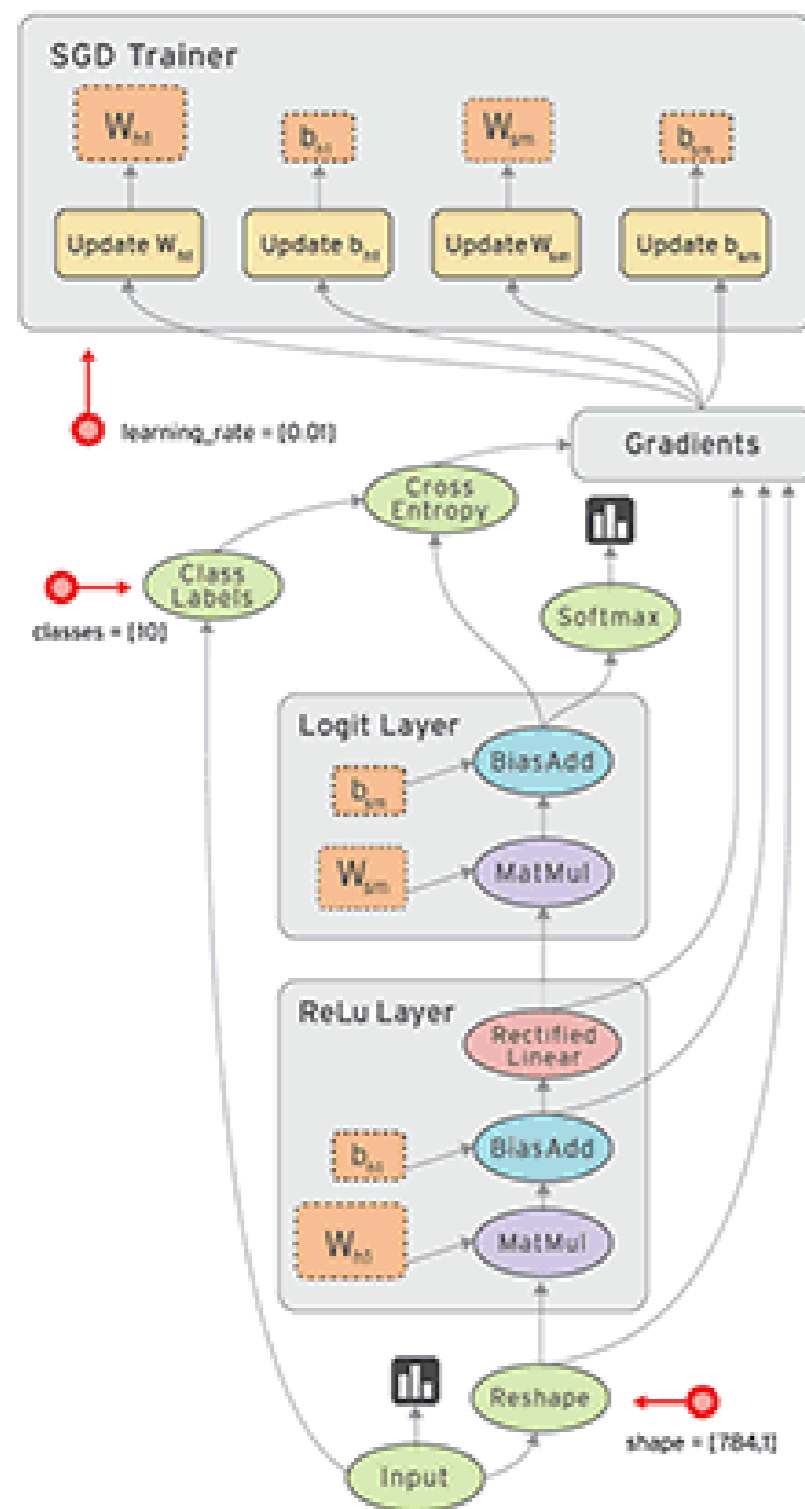
---

LLM, diffusion, and case studies

Optimizations and Parallelization

Basics

# ML System Overview



Dataflow Graph

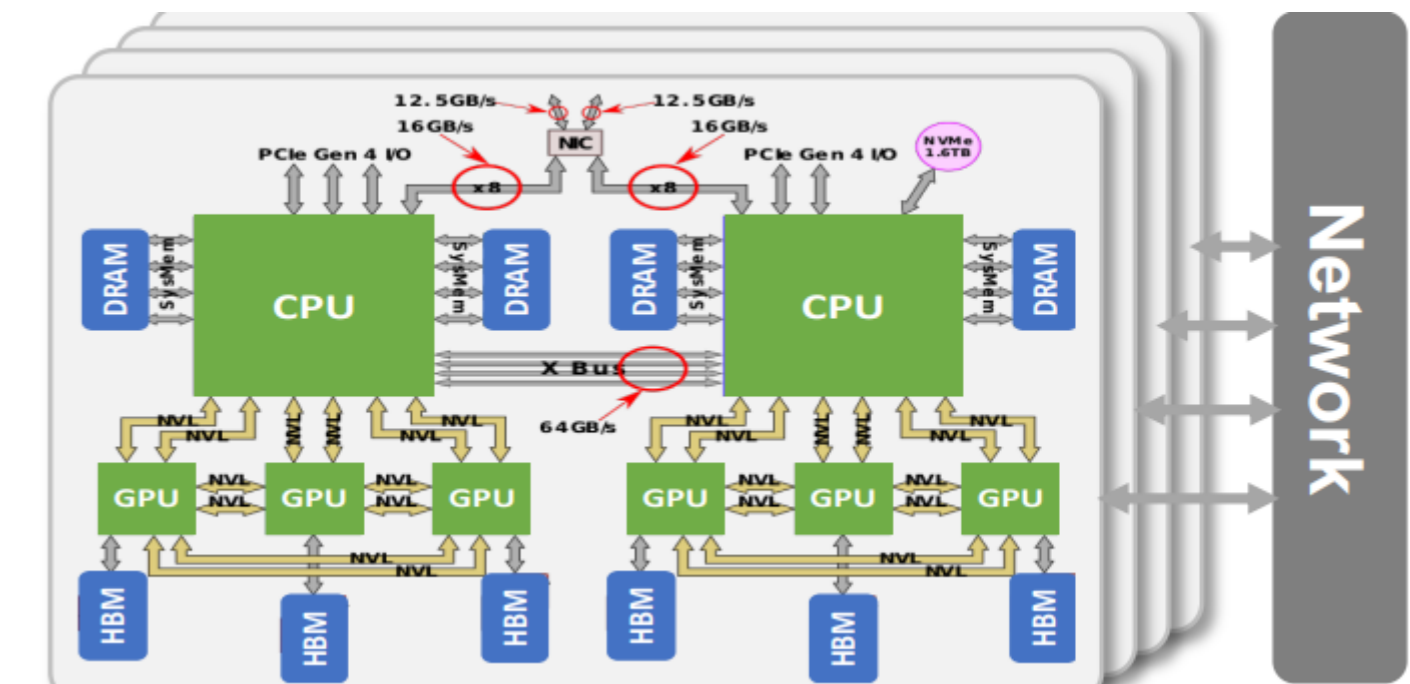
Autodiff

Graph Optimization

Parallelization

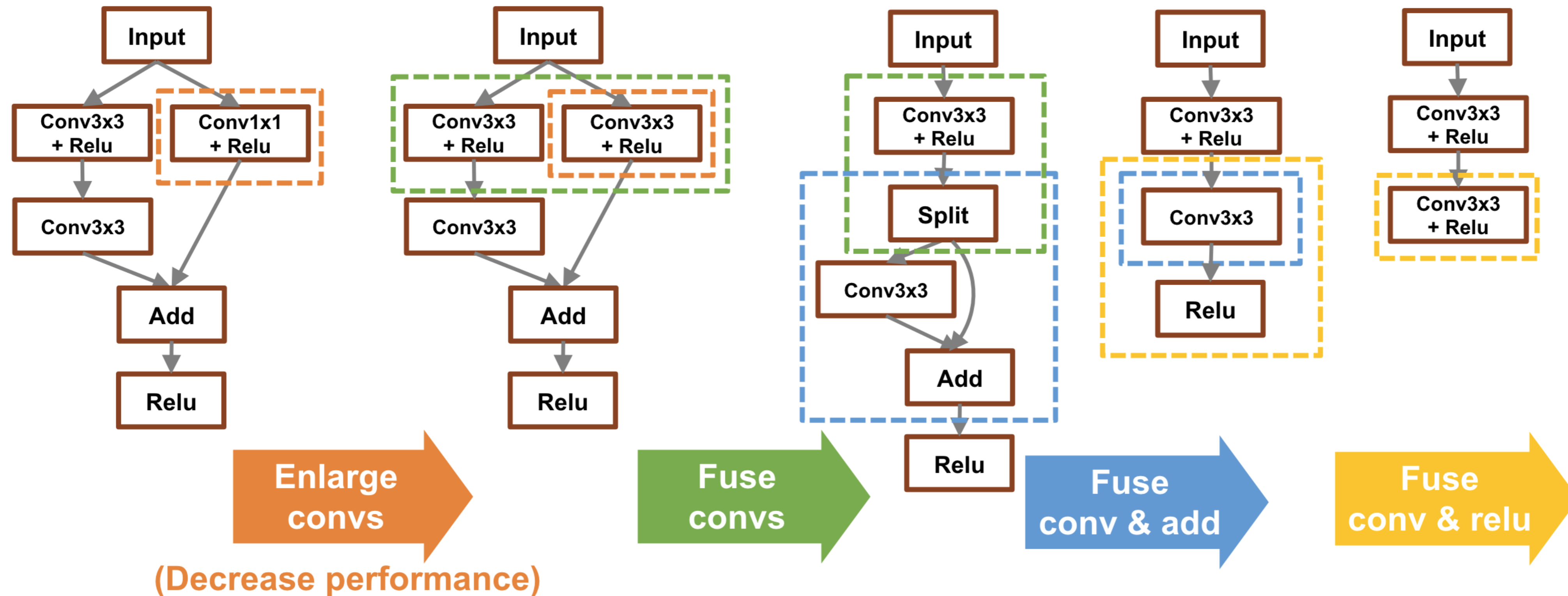
Runtime: schedule / memory

Operator optimization/compilation



Dataflow Graph
Autodiff
Graph Optimization
Parallelization
Runtime: schedule / memory
Operator

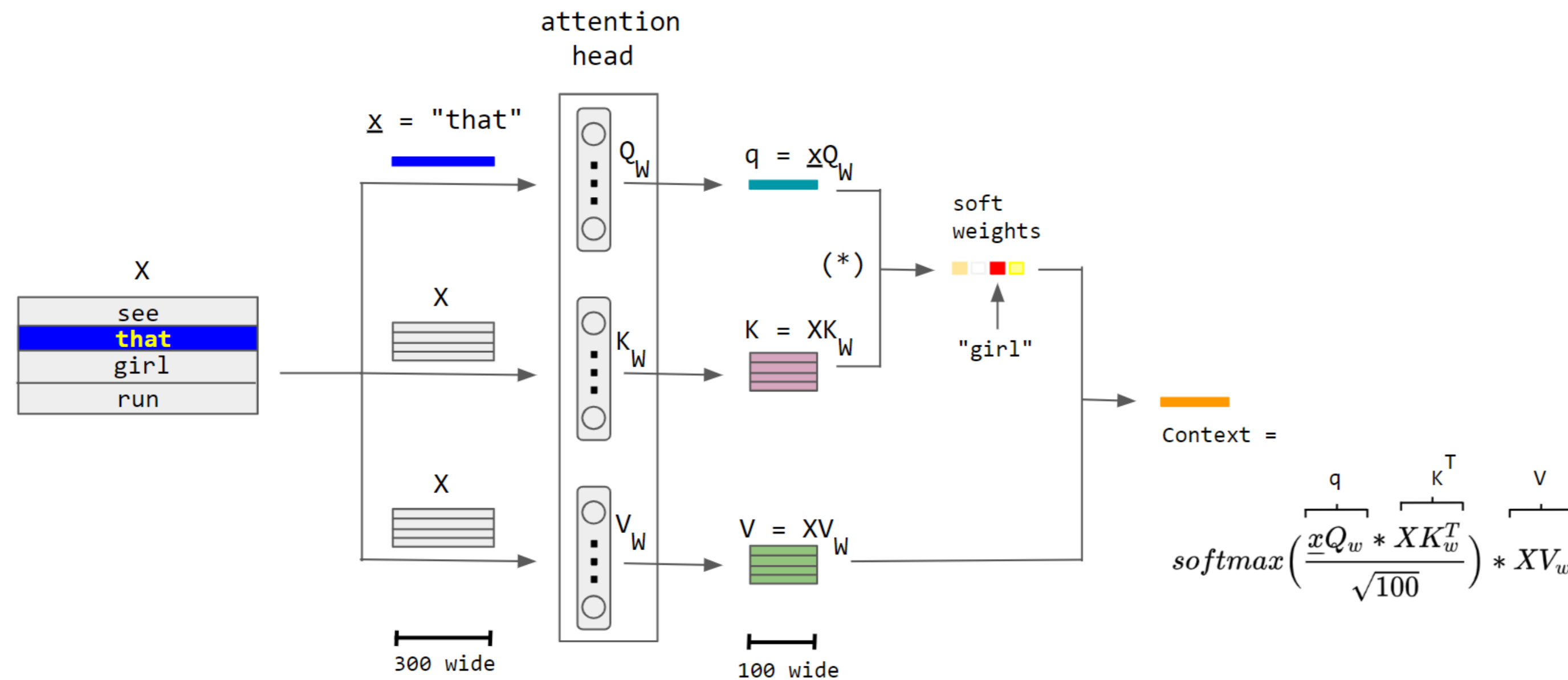
# Motivating Example: we can go further



- Does each step become faster than previous step?
- How does it perf on different hardware?

Dataflow Graph
Autodiff
Graph Optimization
Parallelization
Runtime: schedule / memory
Operator

# Motivating Example 2: Attention



# Original

$$Q = \text{matmul}(W_q, h)$$

$$K = \text{matmul}(W_k, h)$$

$$V = \text{matmul}(W_v, h)$$

# Merged QKV

$$QKV = \text{matmul}(\text{concat}(W_q, W_k, W_v), h)$$

$$\text{Context} = \text{softmax}\left(\frac{\overbrace{xQ_w}^q * \overbrace{XK_w^T}^{K^T}}{\sqrt{100}}\right) * \overbrace{XV_w}^v$$

- Why merged QKV is faster?

# Arithmetic Intensity

$$AI = \#ops / \#bytes$$

# Arithmetic Intensity

```
void add(int n, float* A, float* B, float* C){  
    for (int i=0; i<n; i++)  
        C[i] = A[i] + B[i];  
}
```

Two loads, one store per math op  
(arithmetic intensity = 1/3)

1. Read A[i]
2. Read B[i]
3. Add A[i]+B[i]
4. Store C[i]

# Which program performs better? Program 1

```
void add(int n, float* A, float* B, float* C) {  
    for (int i=0; i<n; i++)  
        C[i] = A[i] + B[i];  
}
```

```
void mul(int n, float* A, float* B, float* C) {  
    for (int i=0; i<n; i++)  
        C[i] = A[i] * B[i];  
}
```

```
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);
```

Two loads, one store per math op  
(arithmetic intensity = 1/3)

Two loads, one store per math op  
(arithmetic intensity = 1/3)

Overall arithmetic intensity = 1/3

# Which program performs better? Program 2

```
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);

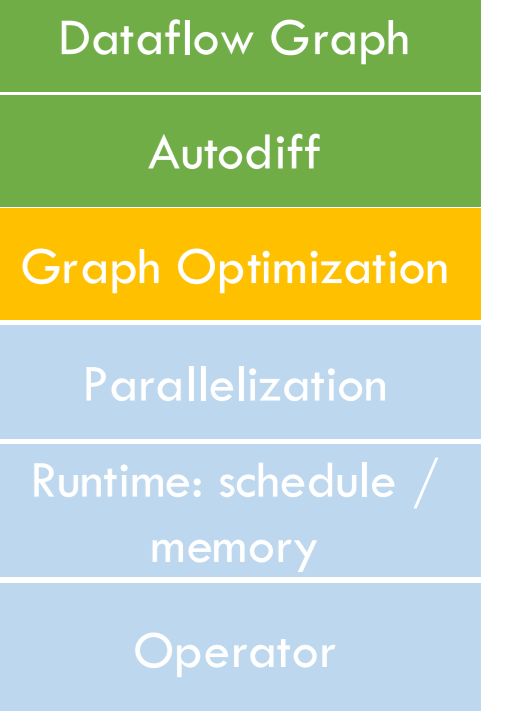
void fused(int n, float* A, float* B, float* C, float* D,
float* E) {
    for (int i=0; i<n; i++)
        E[i] = D[i] + (A[i] + B[i]) * C[i];
}
// compute E = D + (A + B) * C
fused(n, A, B, C, D, E);
```

Overall arithmetic intensity = 1/3

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

# How to perform graph optimization?

- Writing rules / template
- Auto discovery

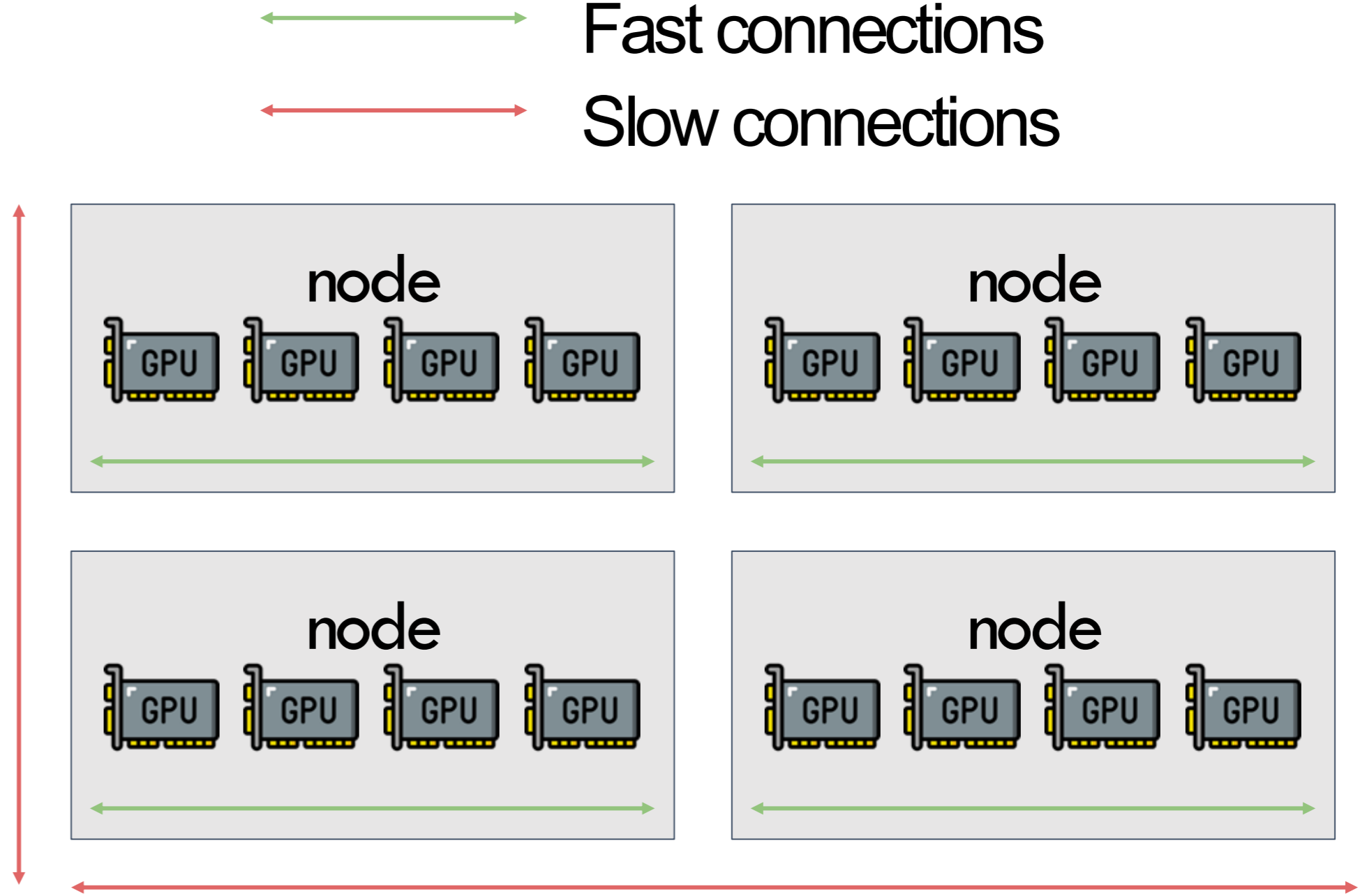
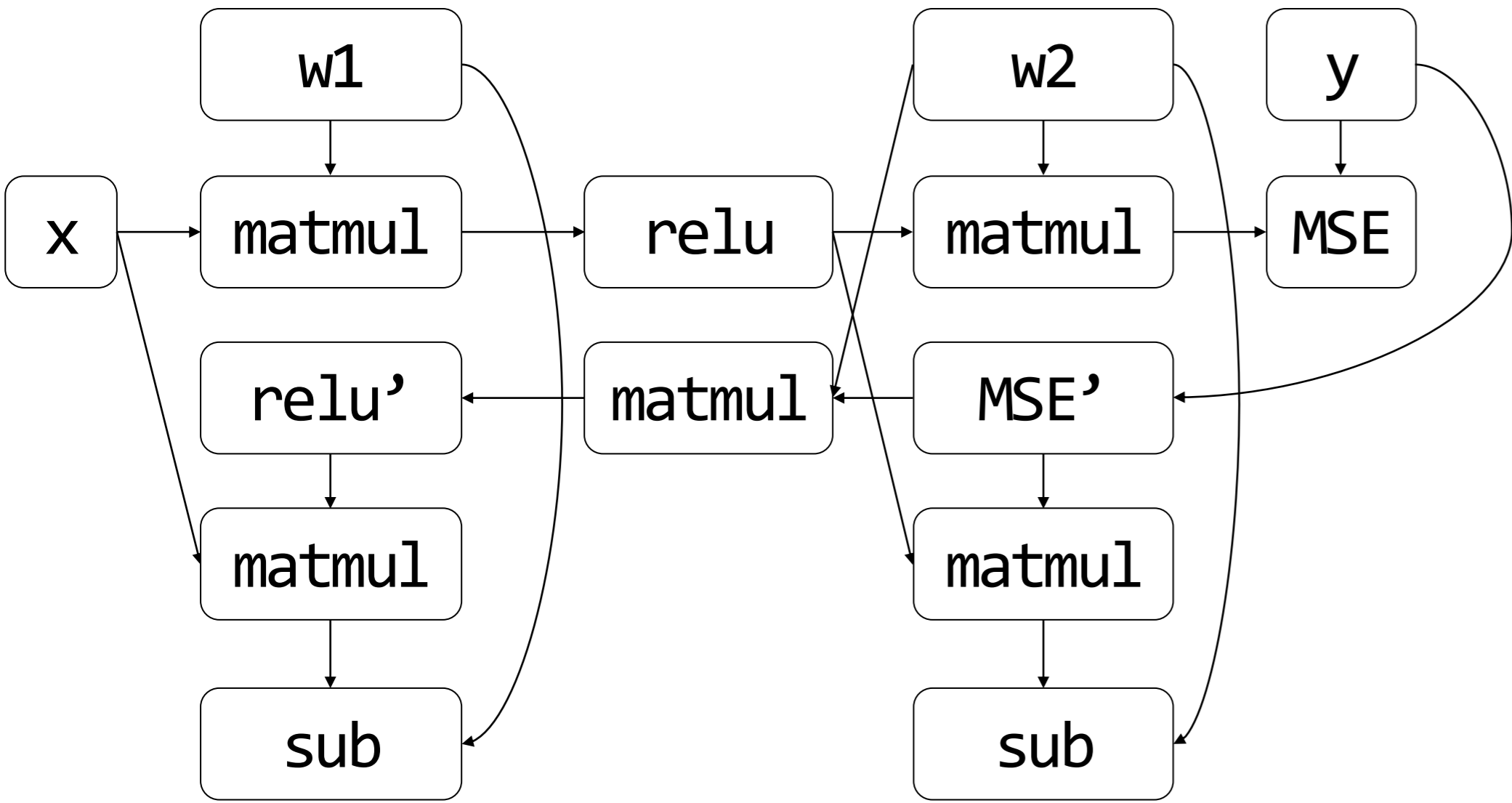


Dataflow Graph
Autodiff
Graph Optimization
Parallelization
Runtime: schedule / memory
Operator

# Parallelization

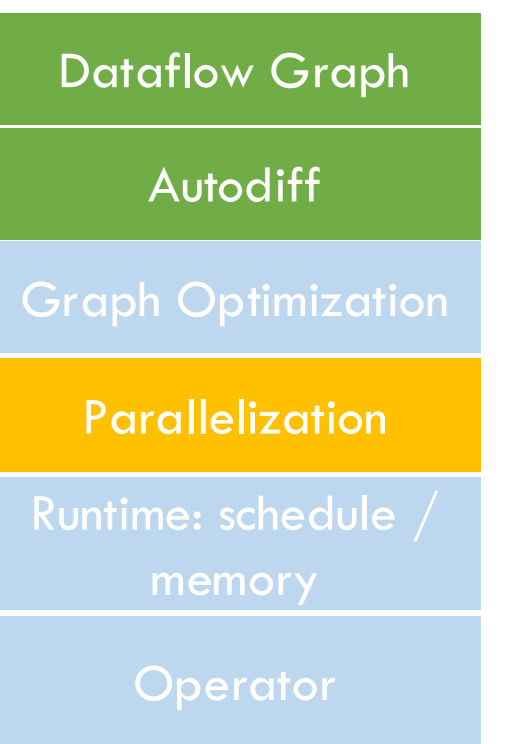
- Goal: parallelize the graph compute over multiple devices

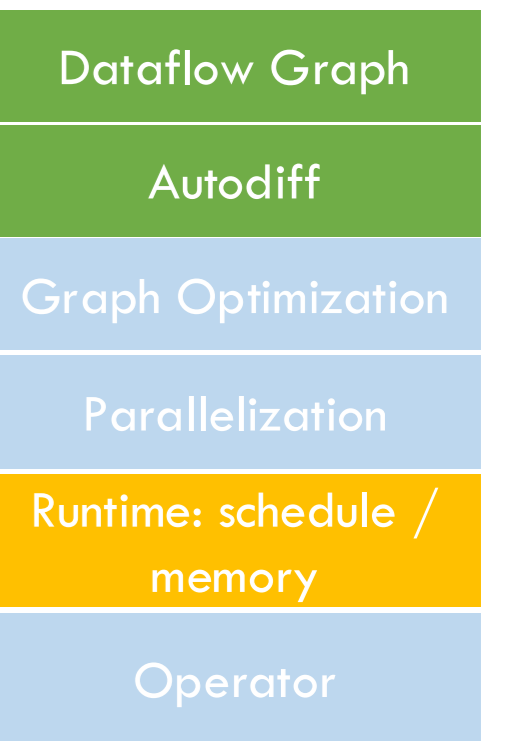
How to partition the computational graph on the device cluster?



# Parallelization Problems

- How to partition
- How to communicate
- How to schedule
- Consistency
- How to auto-parallelize?





# Runtime and Scheduling

- Goal: schedule the compute/communication/memory in a way that
  - As fast as possible
  - Overlap communication with compute
  - Subject to memory constraints

Dataflow Graph
Autodiff
Graph Optimization
Parallelization
Runtime: schedule / memory
Operator

# Operator Implementation

- Goal: get the fastest possible implementation of
  - Matmul
  - Conv2d?
- For different hardware: V100, A100, H100, phone, TPU
- For different precision: fp32, fp16, fp8, fp4
- For different shape: conv2d\_3x3, conv2d\_5x5, matmul2D, 3D, attention

# High-level Picture

Data

?  $\{x_i\}_{i=1}^n$

Model



Math primitives  
(mostly matmul)

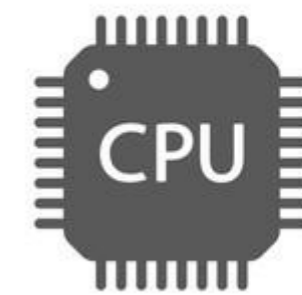
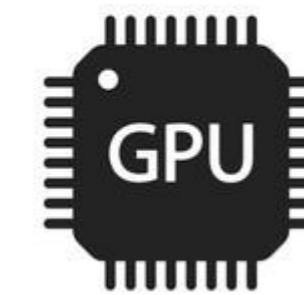
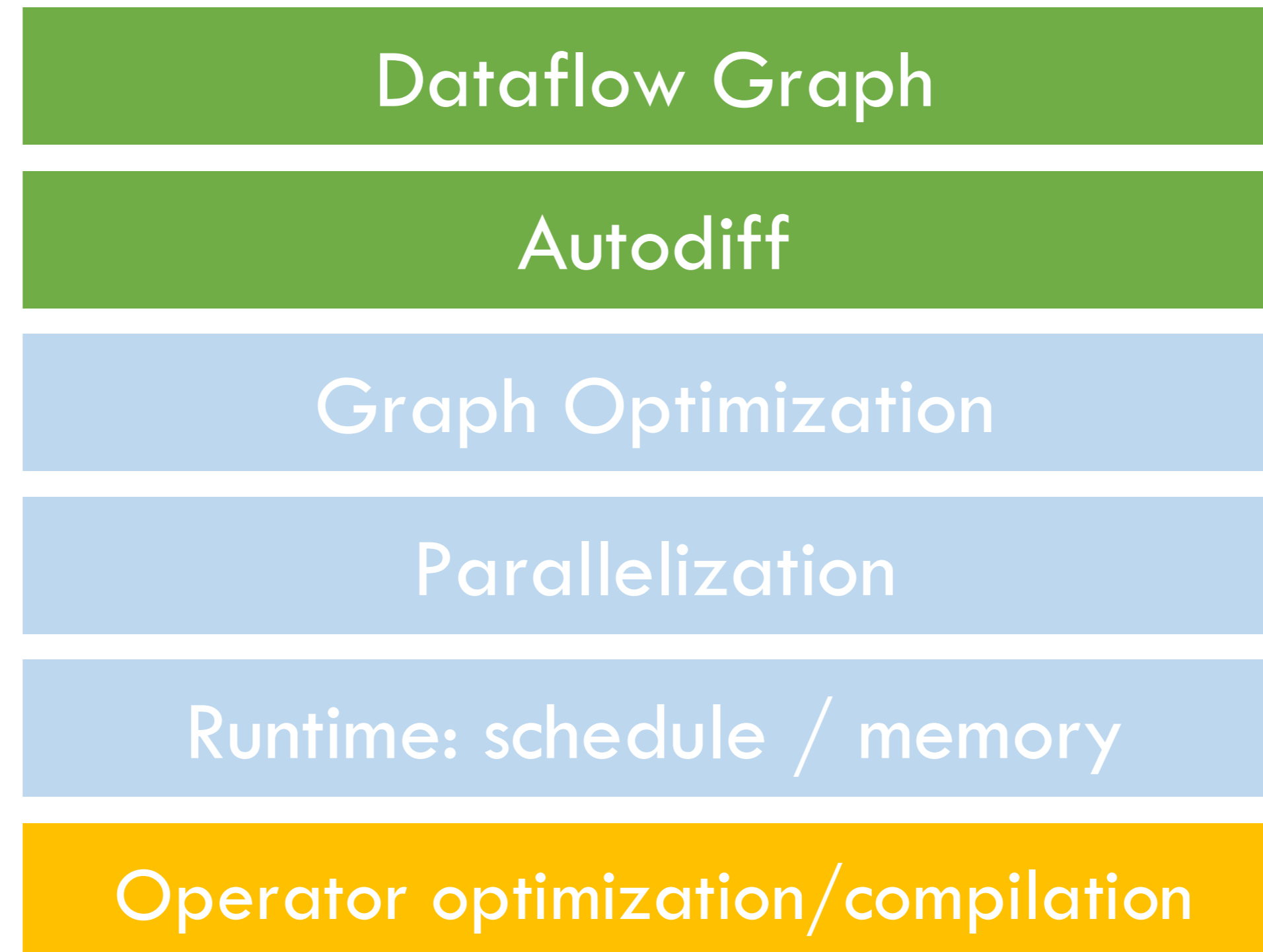
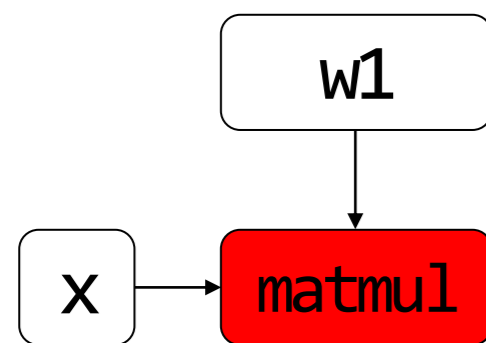


A repr that expresses the  
computation using primitives

Compute

? Make them run on (clusters  
of ) different kinds of  
hardware

# Next: How to make operators run (fast) on devices?



Our Goal in This Layer: Maximize Arithmetic Intensity

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

# Next

- How we can make operator fast in general
- Case study: Matmul
- GPU architecture and programming

# How we can make operators fast in general

- Vectorization
- Data layout
- Parallelization

Using vectorized operations: array add

## Why vectorized is faster than unvectorized?

```
Float A[256], B[256], C[256]
For (int i = 0; i < 256; ++i) {
    C[i] = A[i] + B[i]
}
```

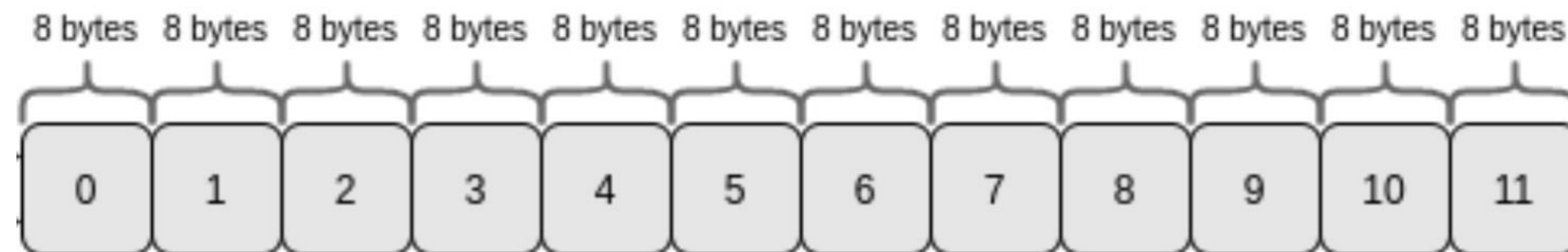
unvectorized

```
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);
}
```

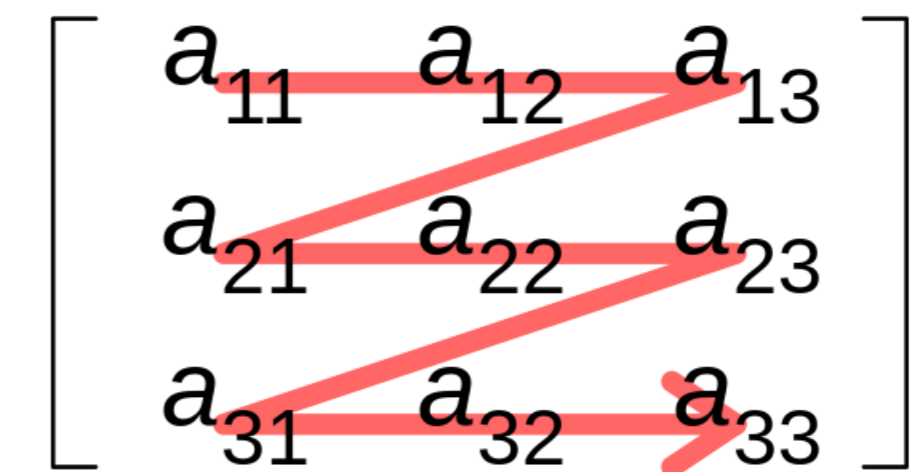
vectorized

# Data Layout: make read/write faster

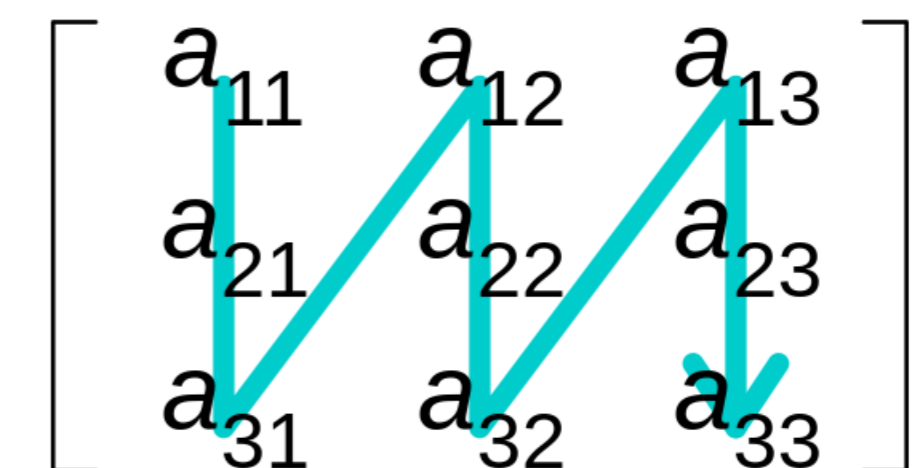
- How to store a matrix in memory
  - Data in memory are stored sequentially (no tensor awareness)
- Row Major:  $A[i, j] = A.data[i * A.shape[1] + j]$
- Column major:  $A[i, j] = A.data[j * A.shape[0] + i]$



Row-major order



Column-major order



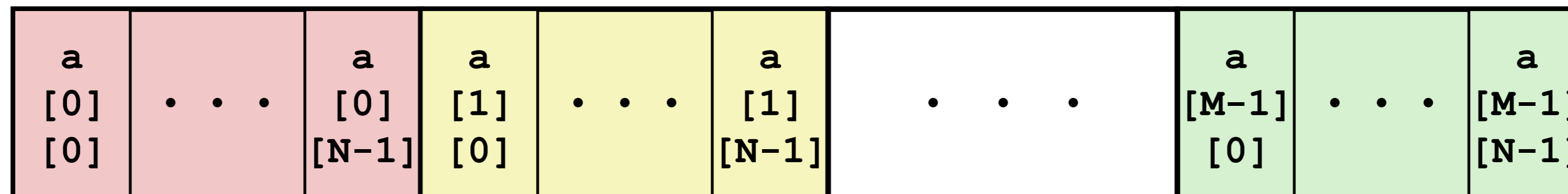
# Be aware of your data layout

**Assuming row-major  
array**

```
int sum_array_rows(int a[M][N])
{
    int i, j, sum = 0;

    for (j = 0; j < N; j++)
        for (i = 0; i < M; i++)
            sum += a[i][j];

    return sum;
}
```



How to improve the above program?

# MCQ Time

## Data

?  $\{x_i\}_{i=1}^n$

ML Systems Store Data in:

A. Row major

B. Col major

C. Strides format:  $A[i, j] = A.data[offset + i * A.strides[0] + j * A.strides[1]]$

# Strides in High-dimension

**Offset:** the offset of the tensor relative to the underlying storage

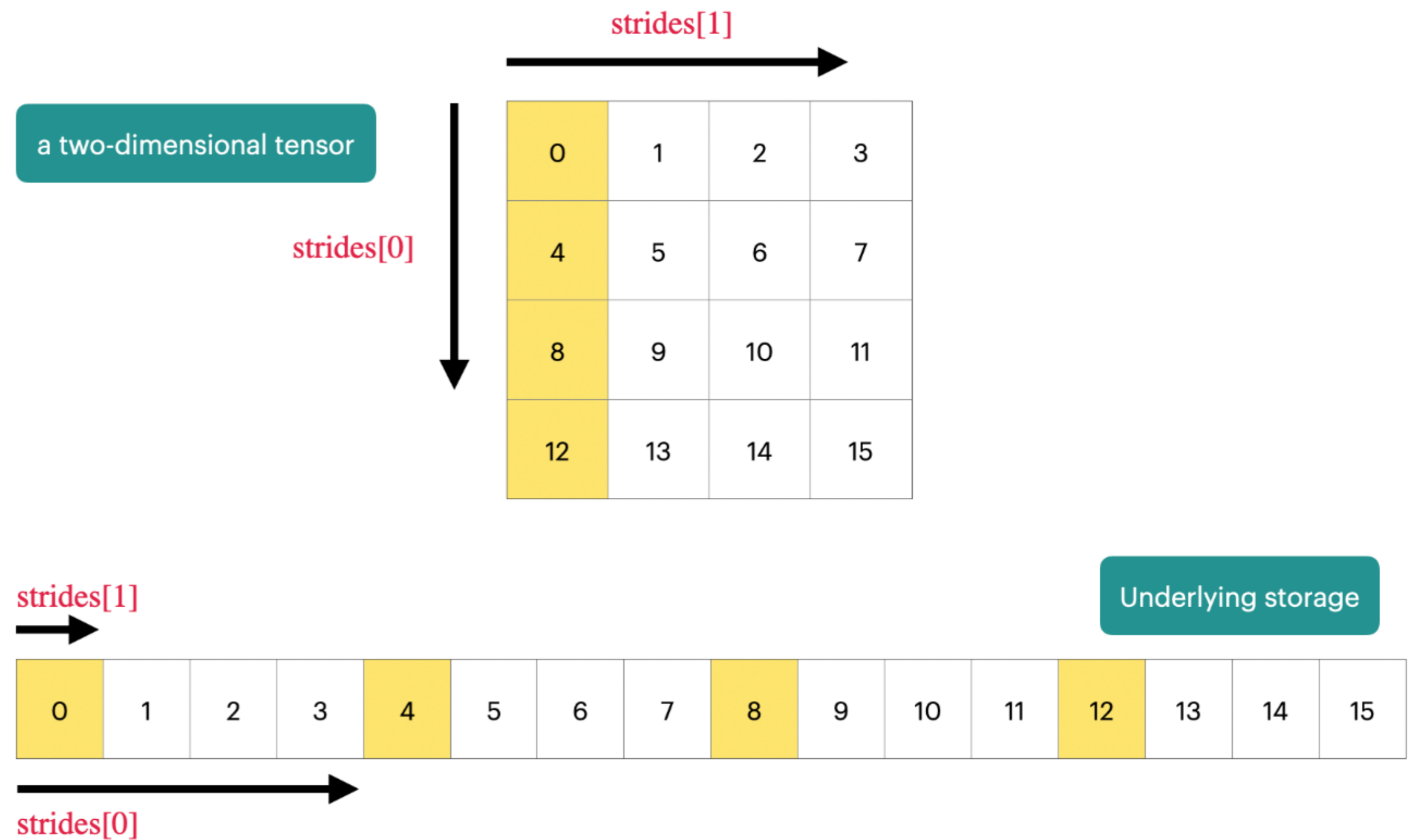
**Strides:** `strides[i]` indicates how many “elements” need to be skipped in memory to move “one unit” in the *i*-th dimension of the tensor

## ▼ 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  ]
```

# Strides format

- What we have when:
  - $A.\text{strides}[0] = 1,$
  - $A.\text{strides}[1] = A.\text{shape}[0]?$
- What we have when:
  - $A.\text{strides}[0] = A.\text{shape}[1]$
  - $A.\text{strides}[1] = 1,$
- Strides offers more flexibility



# Questions

- If a tensor of shape [1, 2, 3, 4] is stored contiguous in memory following row Major, write down its strides?

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

# 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

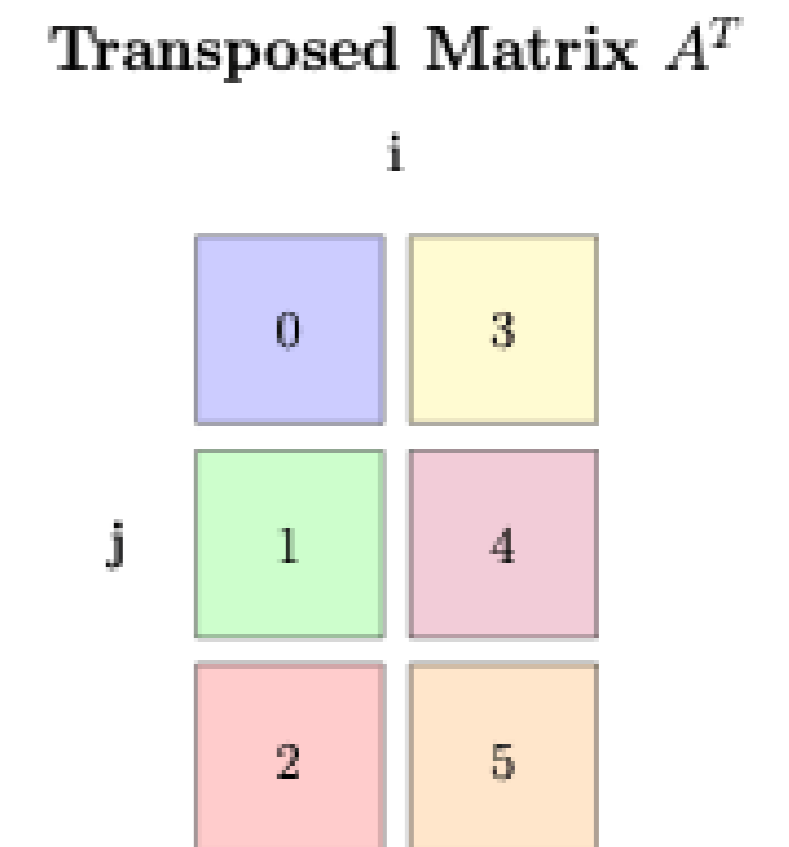
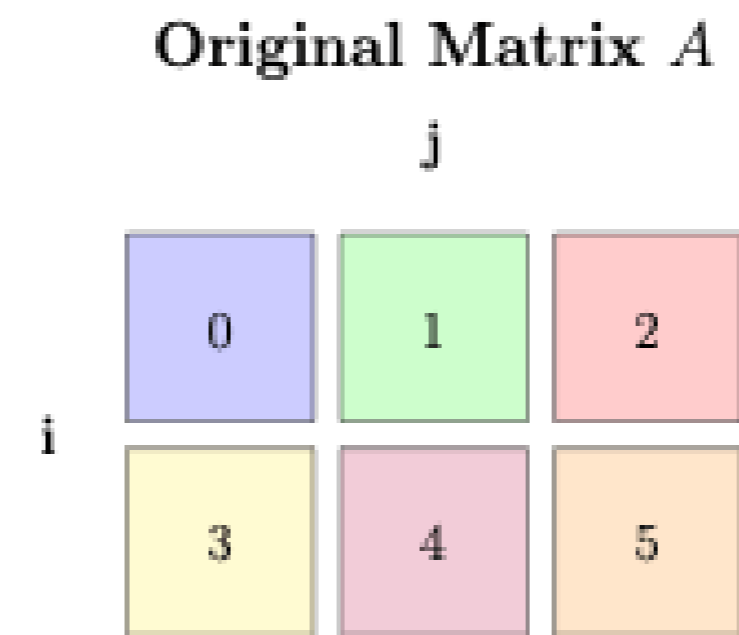
- Change the **offset** by +1
- Reduce the **shape** to [3, 2]
  
- Note: zero copy

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?

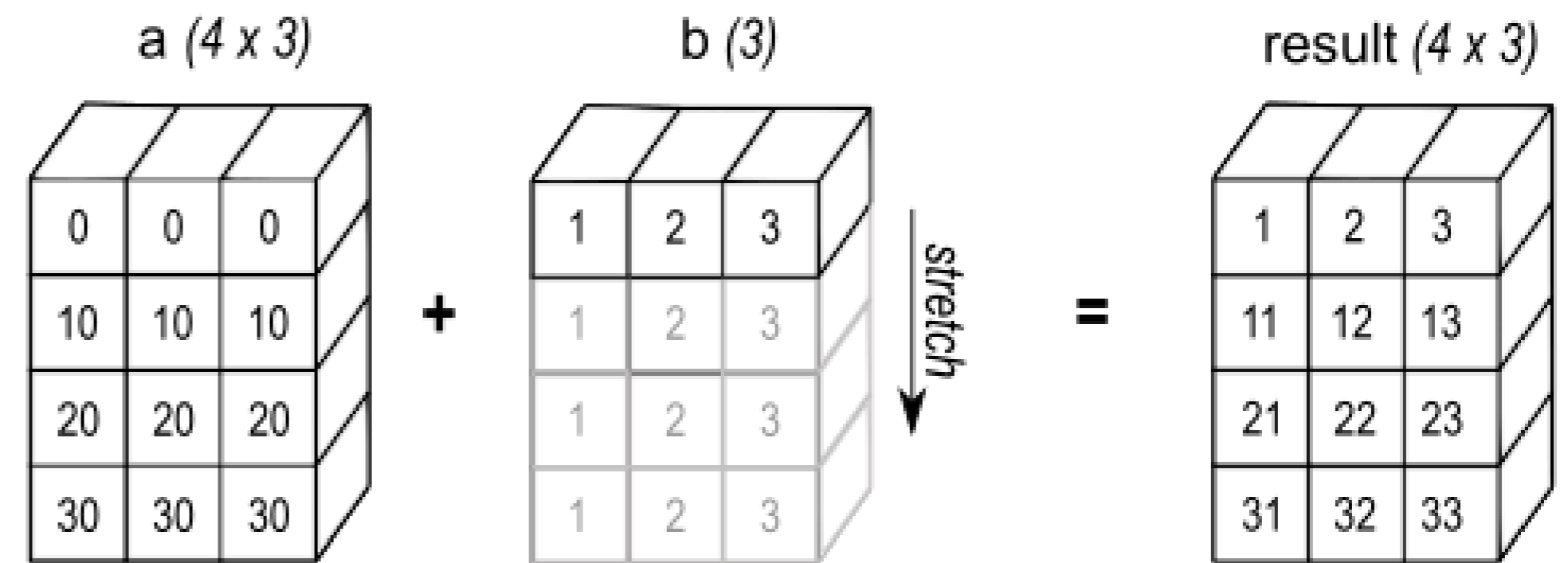
```
Python
1 print(t.stride())
2 # (24, 12, 4, 1)
3
4 print(t.permute((1, 2, 3, 0)).is_contiguous())
5 # True
6
7 print(t.permute((1, 2, 3, 0)).stride())
8 # (12, 4, 1, 24)
9
10 print_internal(t.permute((1, 2, 3, 0)))
11 # tensor([0, 1, 2, 3,
12 #         4, 5, 6, 7,
13 #         8, 9, 10, 11,
14 # 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



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

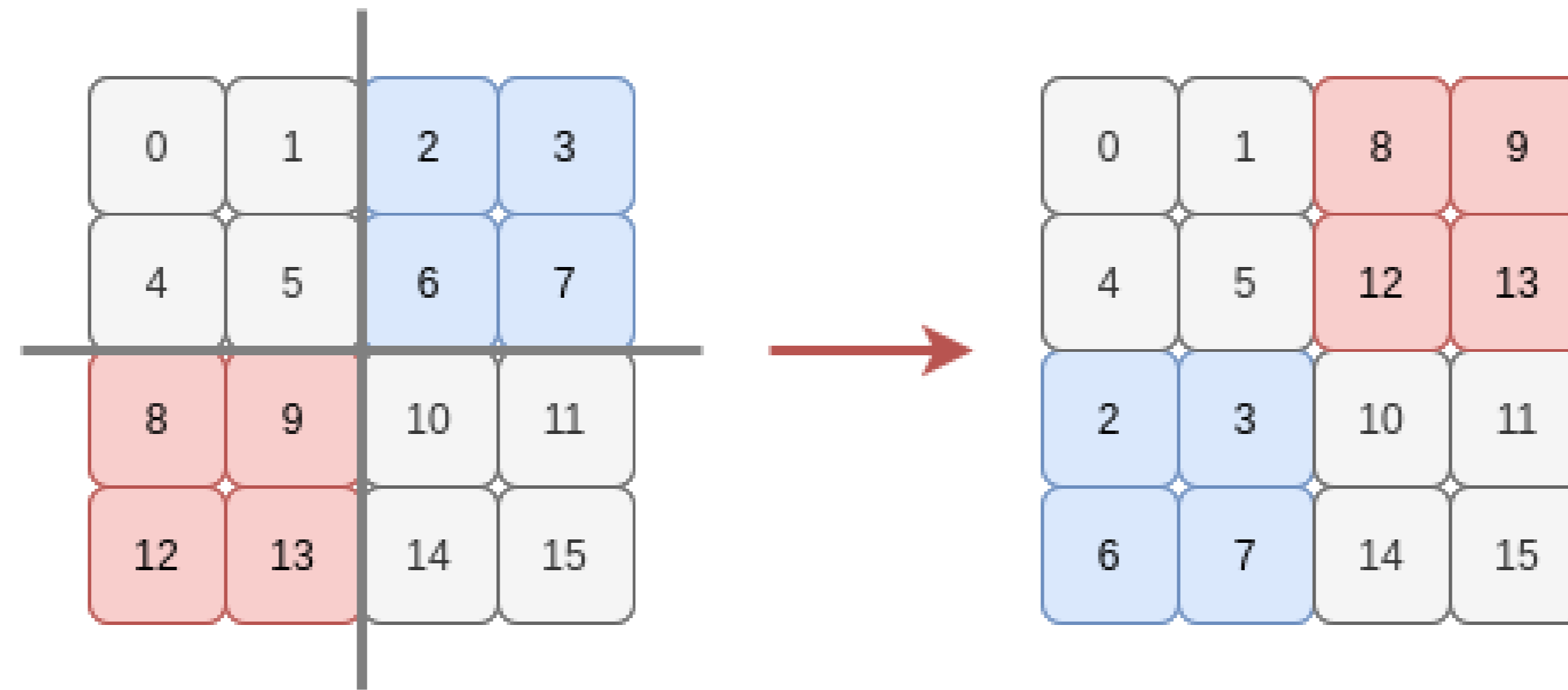
# Broadcast

- 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:  $A[i, j] = A.data[\text{offset} + i * A.strides[0] + j * A.strides[1]]$

# 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.CONTIGUOUS

`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 (Elementise)

- 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);  
}
```

vectorized

We'll com  
back to this  
later

```
#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);  
}
```

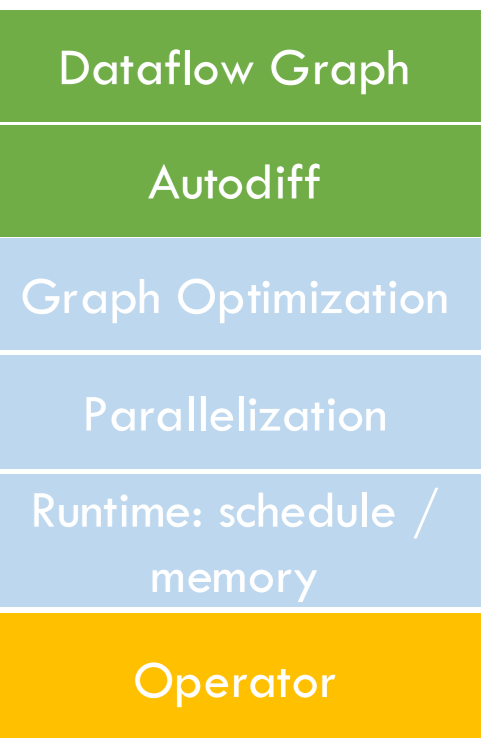
Vectorized &  
parallelized

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

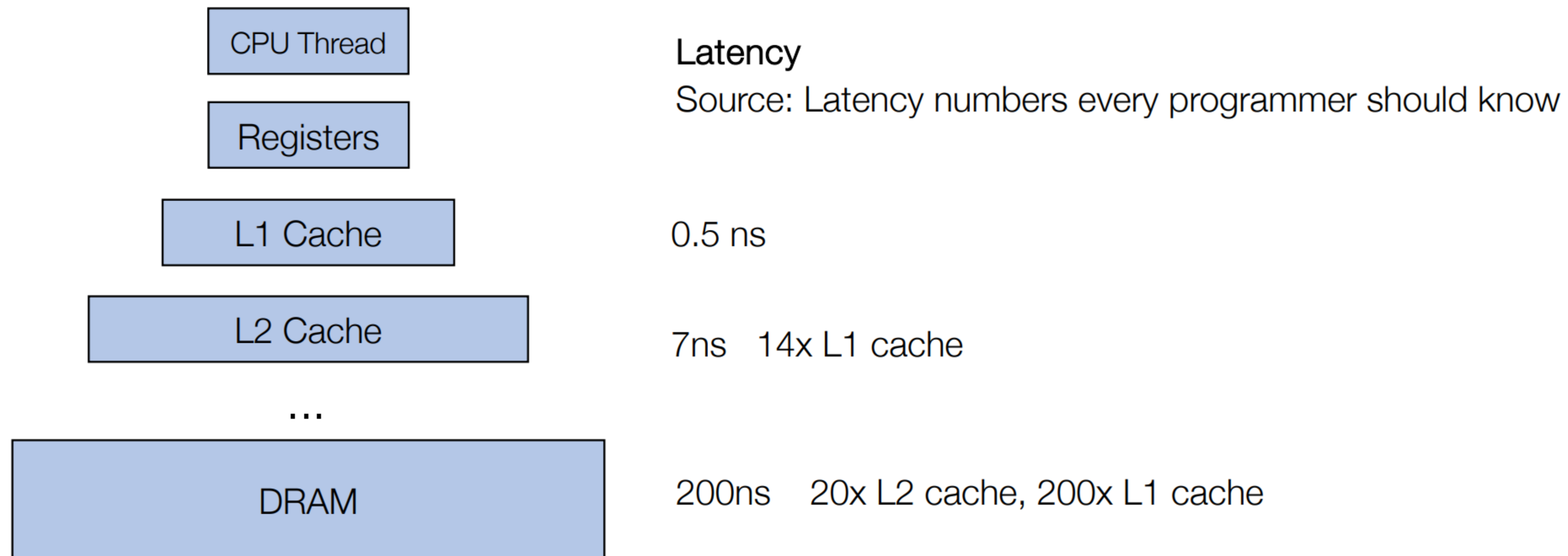


# How to Make Matmul Fast

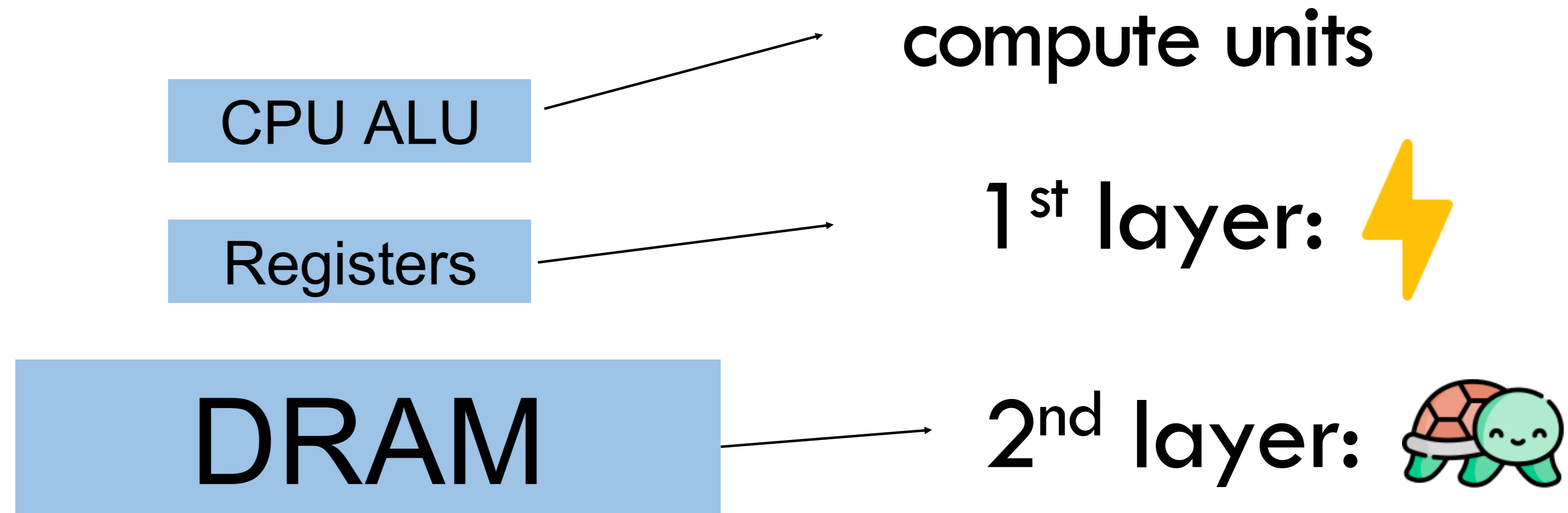
$$\text{max AI} = \#ops / \#bytes$$
A diagram illustrating the components of the equation. A green arrow points upwards from the text "#ops" in the equation above. A red arrow points downwards from the text "#bytes" in the equation above.

# 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 AI: 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);

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);
```

Overall arithmetic intensity = 1/3

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:

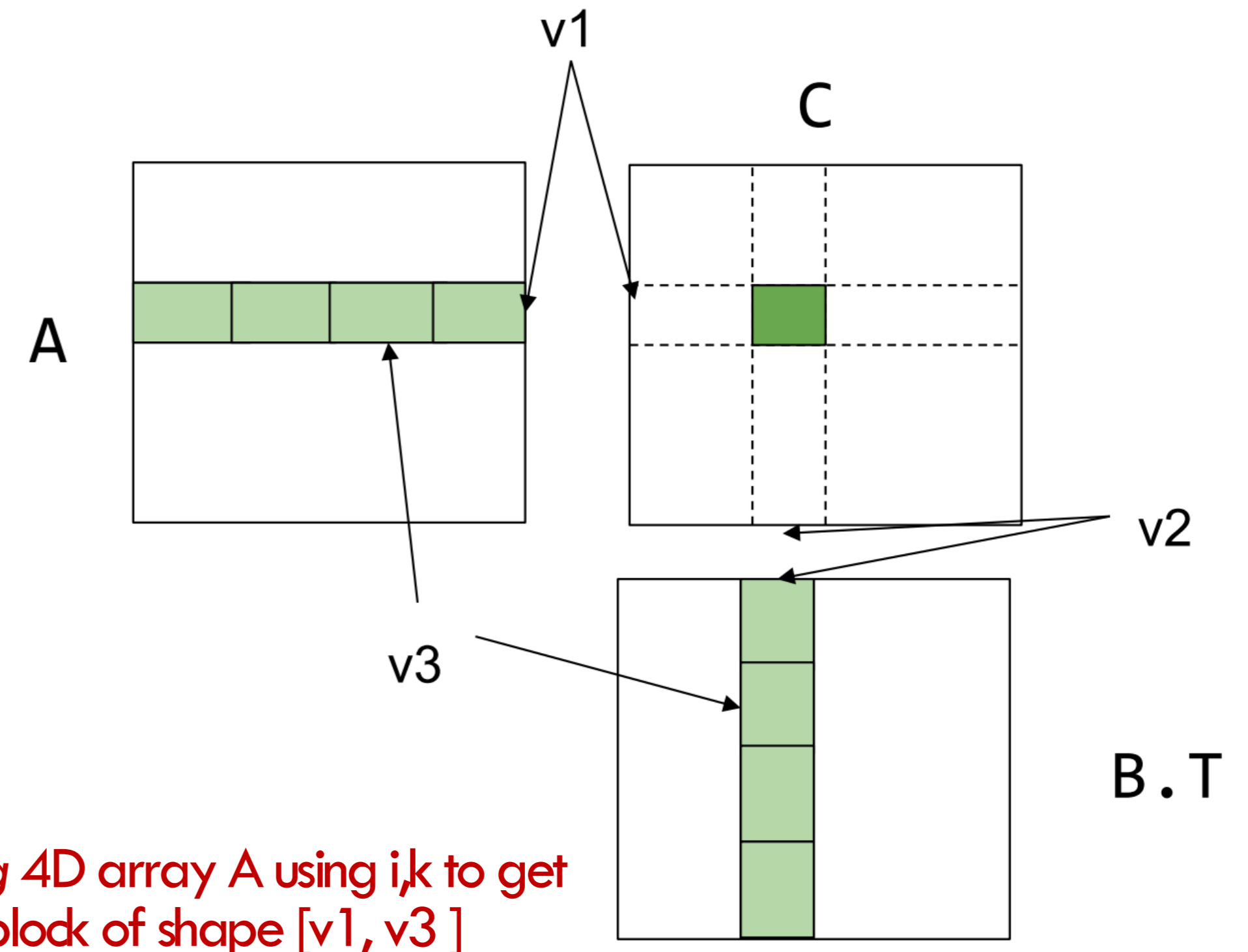
$$2 * n^3 * \text{speed}(\text{dram} \rightarrow \text{register})$$

# 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;  
  }  
}
```

- Assigning the right block to a register block of shape  $[v1, v3]$
- It should be:  $a[0:v1, 0:v3] = A[i][k]$
- RHS is indexing, but LHS is not!



# 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:  
 $v_1 * v_3 + v_2 * v_3 + v_1 * v_2$

Read cost:  
 $(n^3/v_2 + n^3 / v_1) * \text{speed}(\text{dram} \rightarrow \text{register})$

# Register Tiled Matrix Multiplication

- Q: is the load cost related to  $v_3$ ?

$$\text{Read a} \quad N^3 / v_2$$

- Q: How to set  $v_1 / v_2$ ?

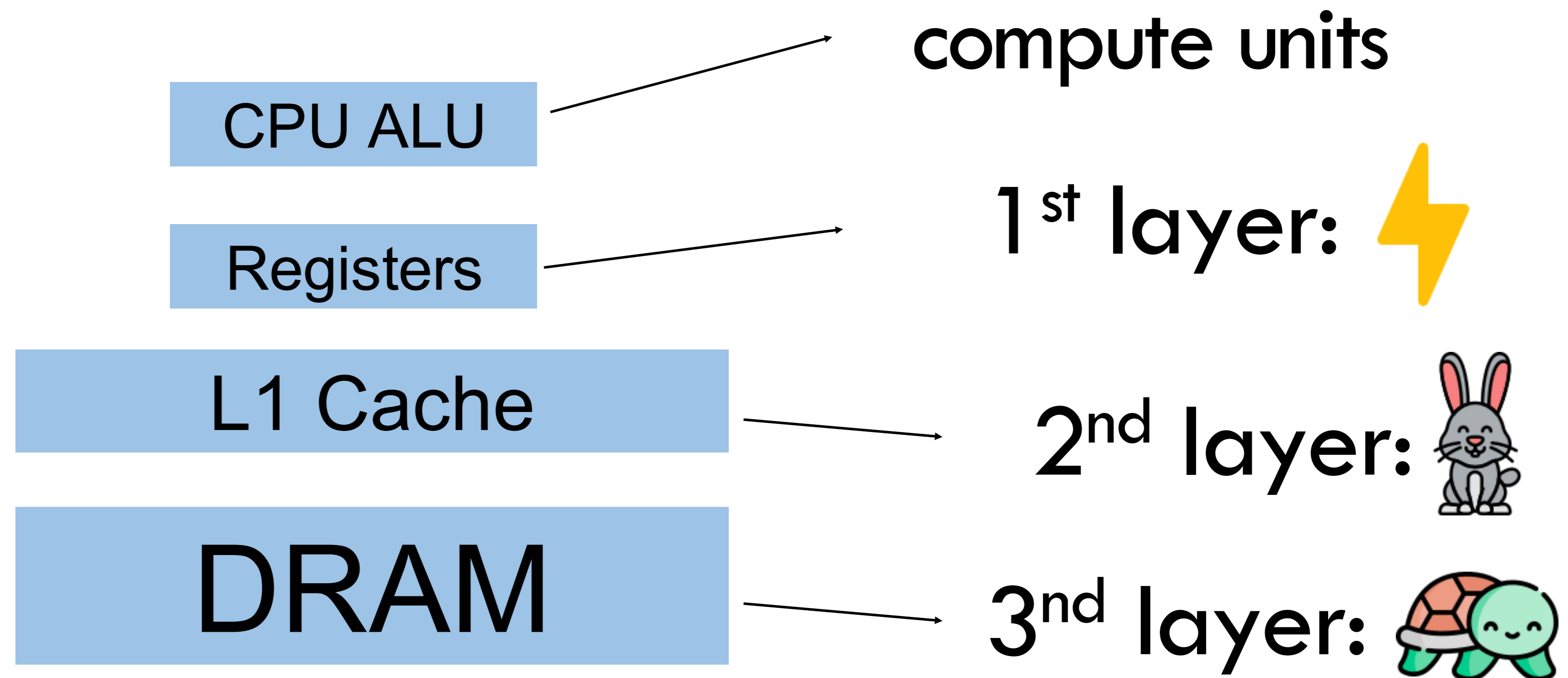
$$\text{Read b} \quad N^3 / v_1$$

- What are the constraints?

#registers needed:

- Q: Why essentially can tiling reduce read cost?  
 $v_1 * v_3 + v_2 * v_3 + v_1 * v_2$

Make it more complicated: Consider L1 cache

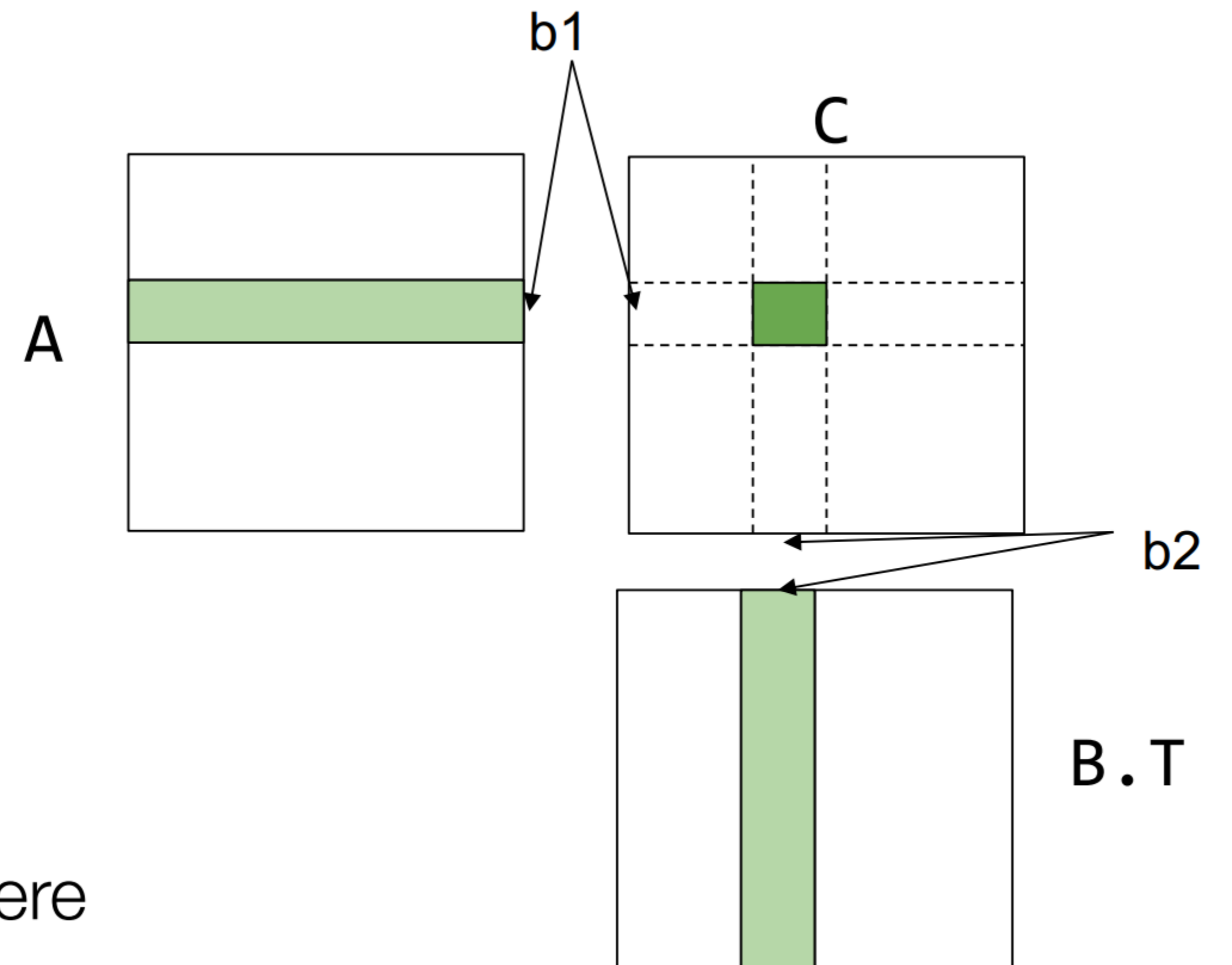


# Cache-aware tiling

- We can further tile the array  $[b1][n]$  or  $[b2][n]$  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 b[b2][n] = B[j];  
  
    C[i][j] = dot(a, b.T);  
  }  
}
```

Sub-procedure, can apply register tiling here

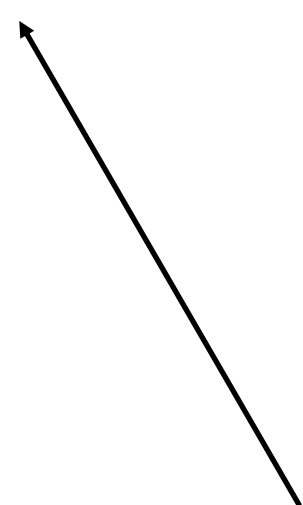


# Cache-aware tiling

```
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 -> l1 cache (cache tiling)
3. l1 cache -> 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 untilted version?

s.t.

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

# 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 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  $\rightarrow$  L1 cache reads here

- Register tiling using  $v1$  and  $v2$
- L1  $\rightarrow$  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 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**

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

**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  $\rightarrow$  dram  $\rightarrow$  L2  $\rightarrow$  L1  $\rightarrow$  Register
- How to choose  $v_1, v_2, b_1, b_2, c_1, c_2$ ?
- While we are reading from dram  $\rightarrow$  L2, can we concurrently read:
  - L2  $\rightarrow$  L1
  - L1  $\rightarrow$  register
- S.t. sizes of L2, L1, registers

## Why tiling works: **reuse** loading

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

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

Access of A is independent of the dimension of j

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 ( $\leq 5$ ) for parallelization.