



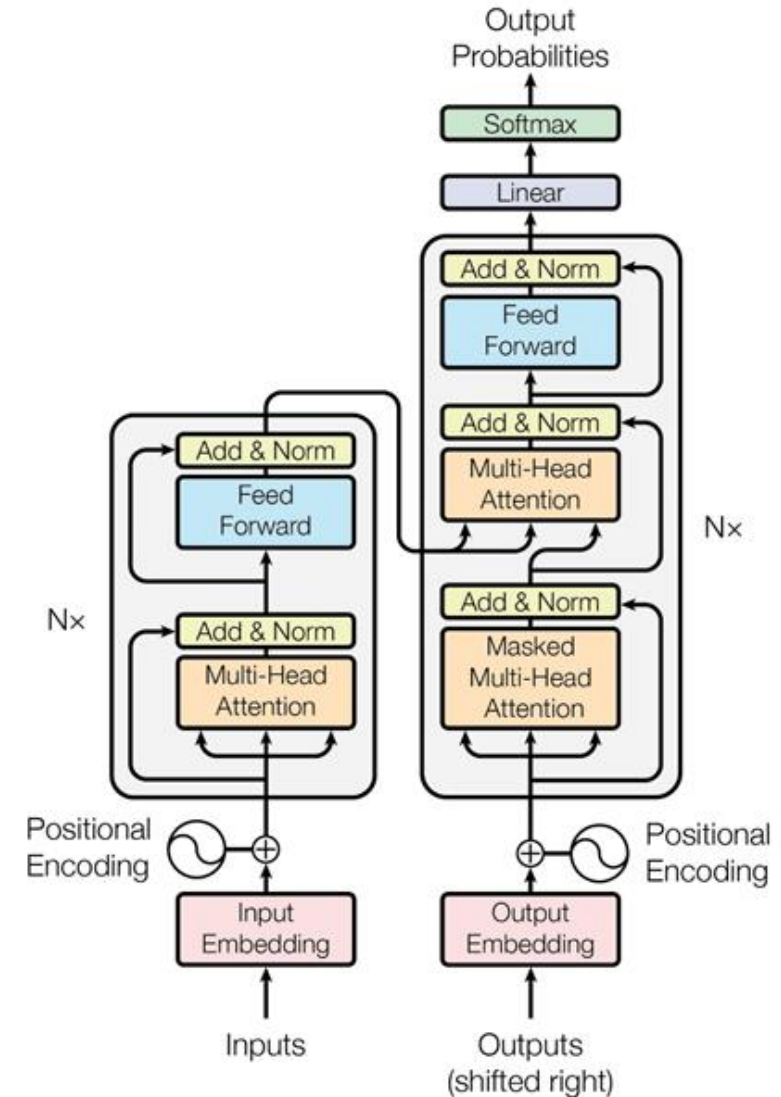
# CS 498: Machine Learning System Spring 2025

Minjia Zhang

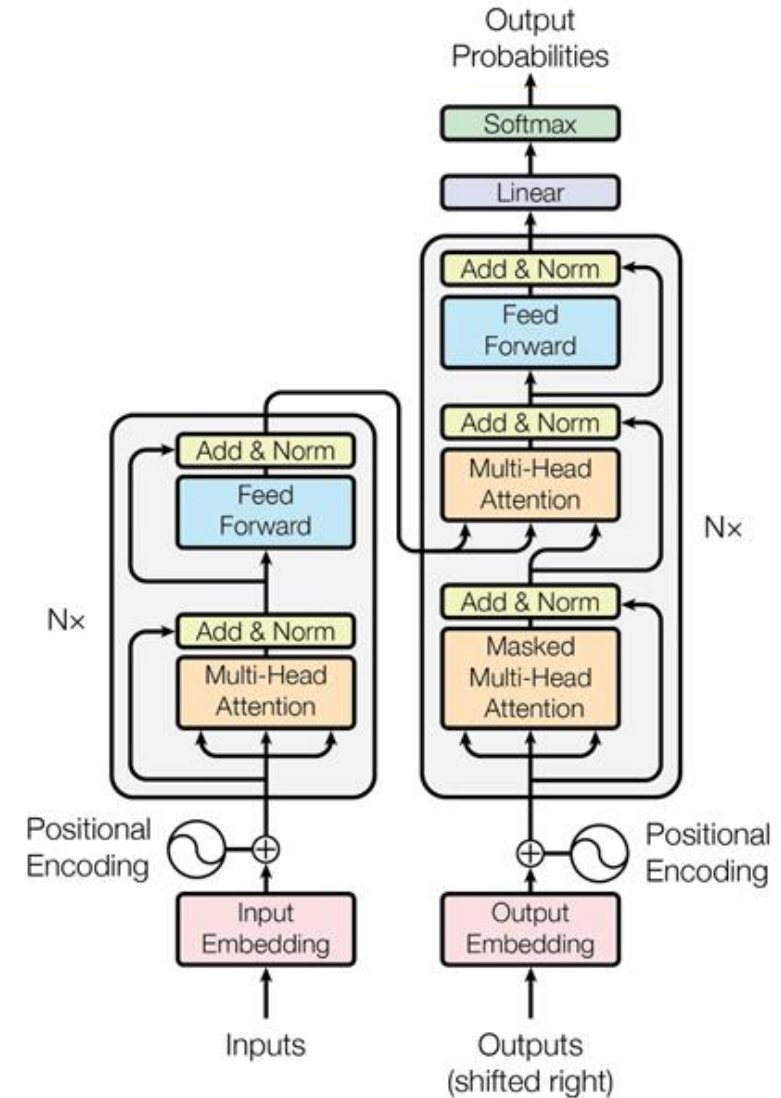
The Grainger College of Engineering

- **Transformers Deep Dive**
- **Arithmetic Intensity**

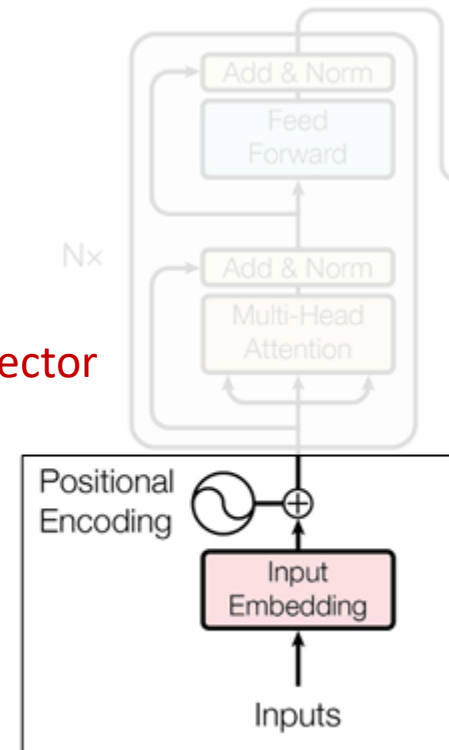
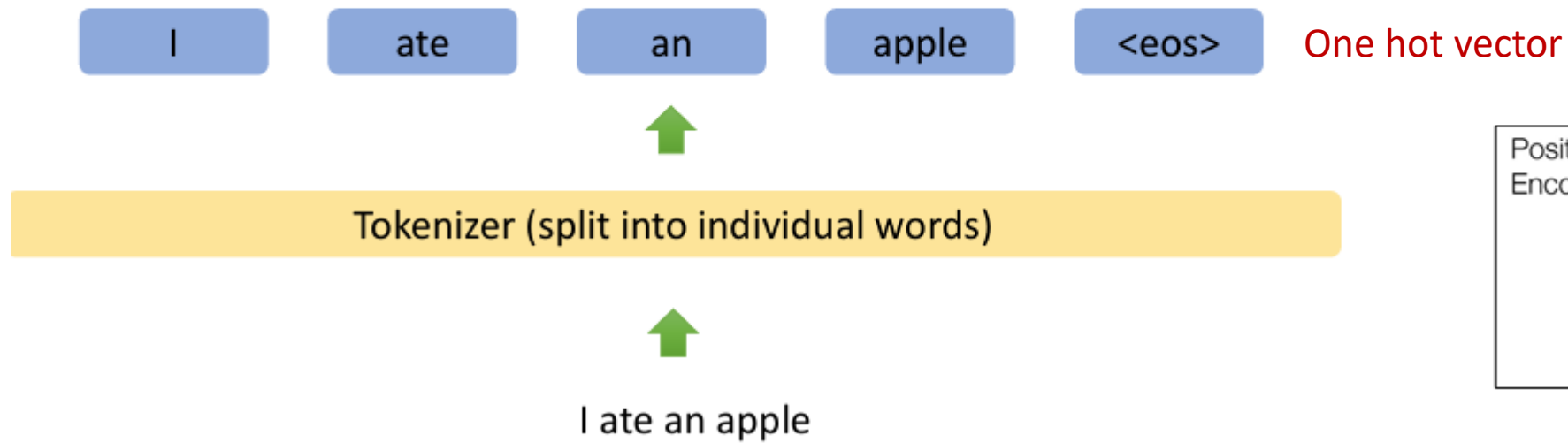
- Tokenization
- Input embeddings
- Position encodings
- Query, Key, Value
- Attention
- Self-attention
- Multi-head attention
- Feed forward
- Add & norm
- Residual
- Masked attention
- Causal attention
- Linear
- Softmax
- Encoders
- Decoder
- Encoder-decoder



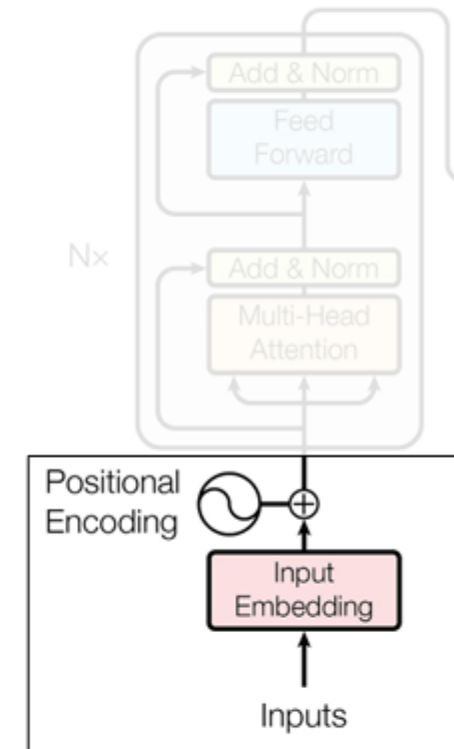
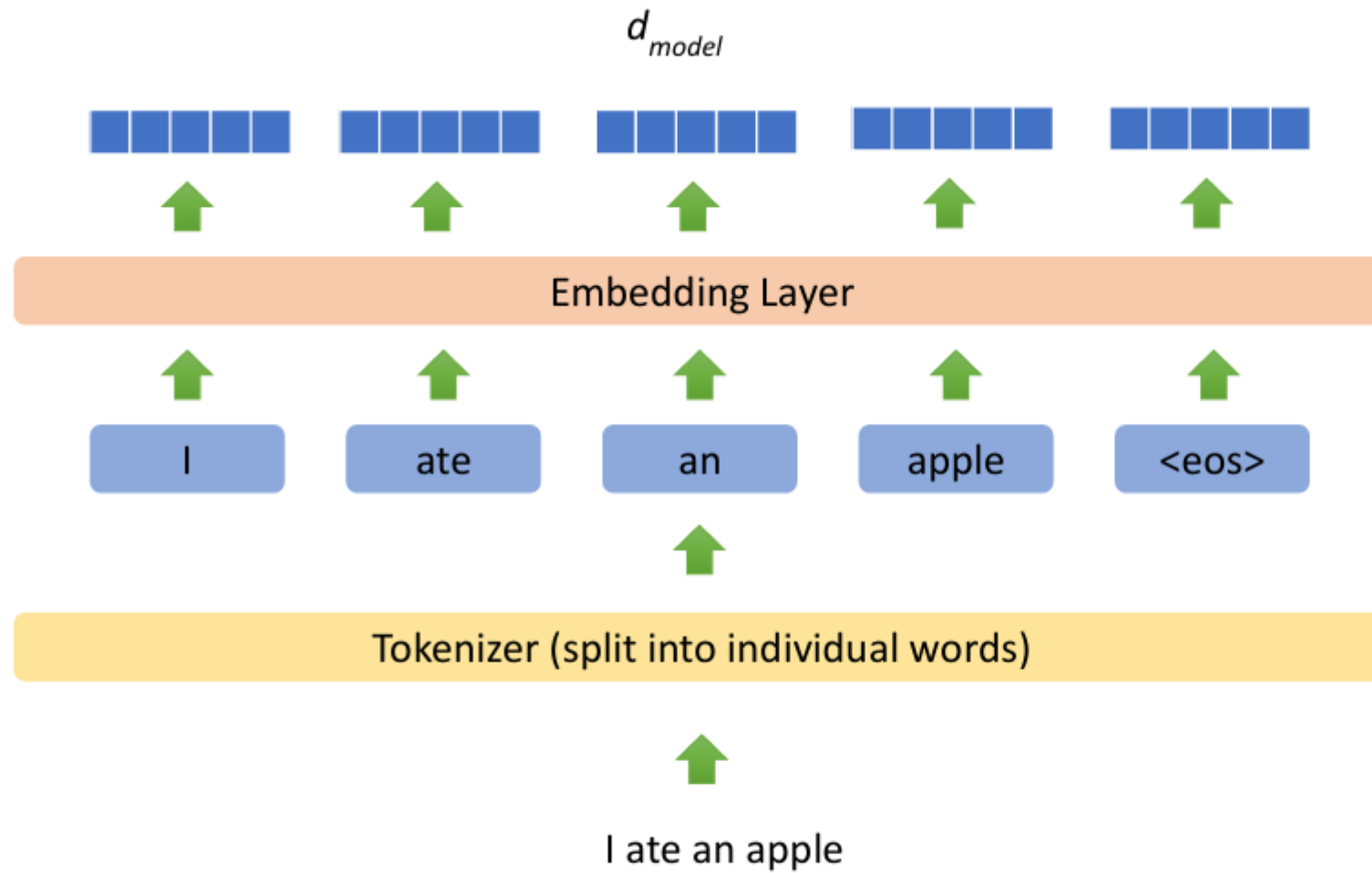
# Machine Translation



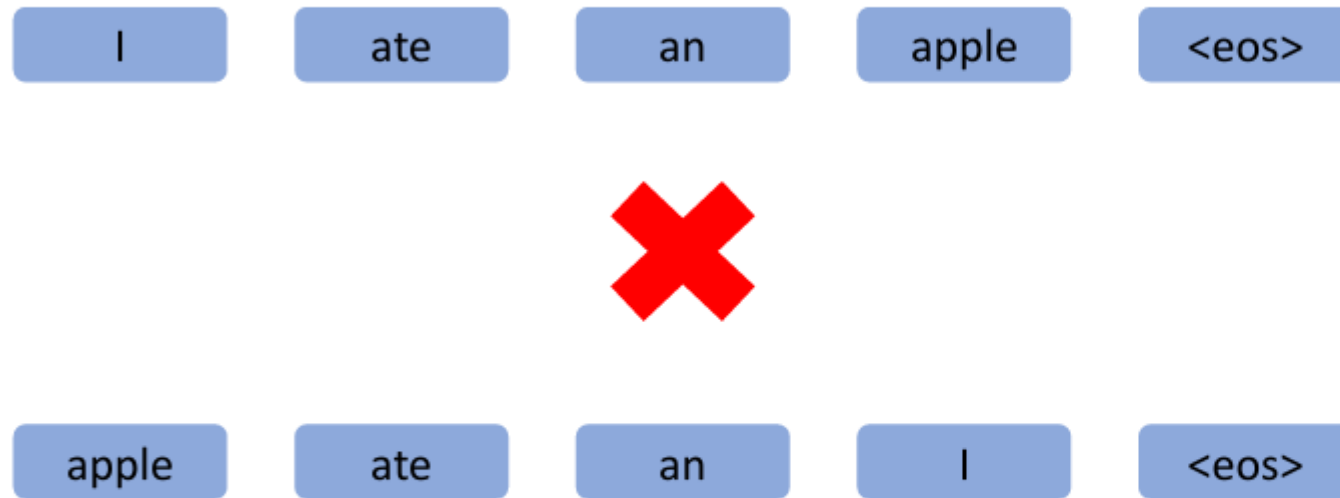
# Tokenization



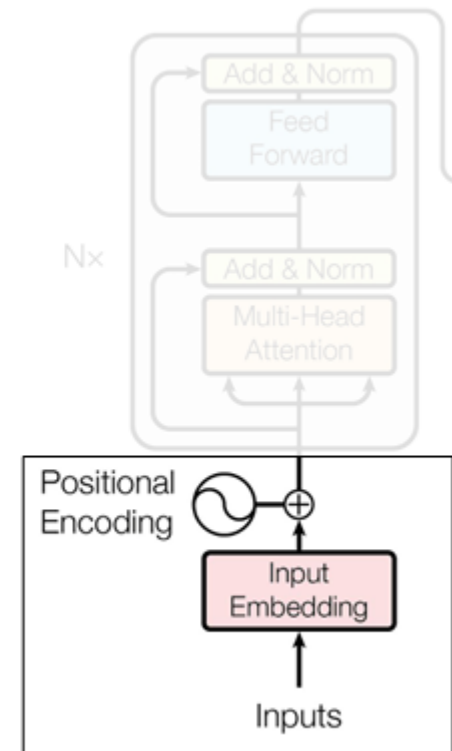
# Input Embeddings



Generate Input Embeddings



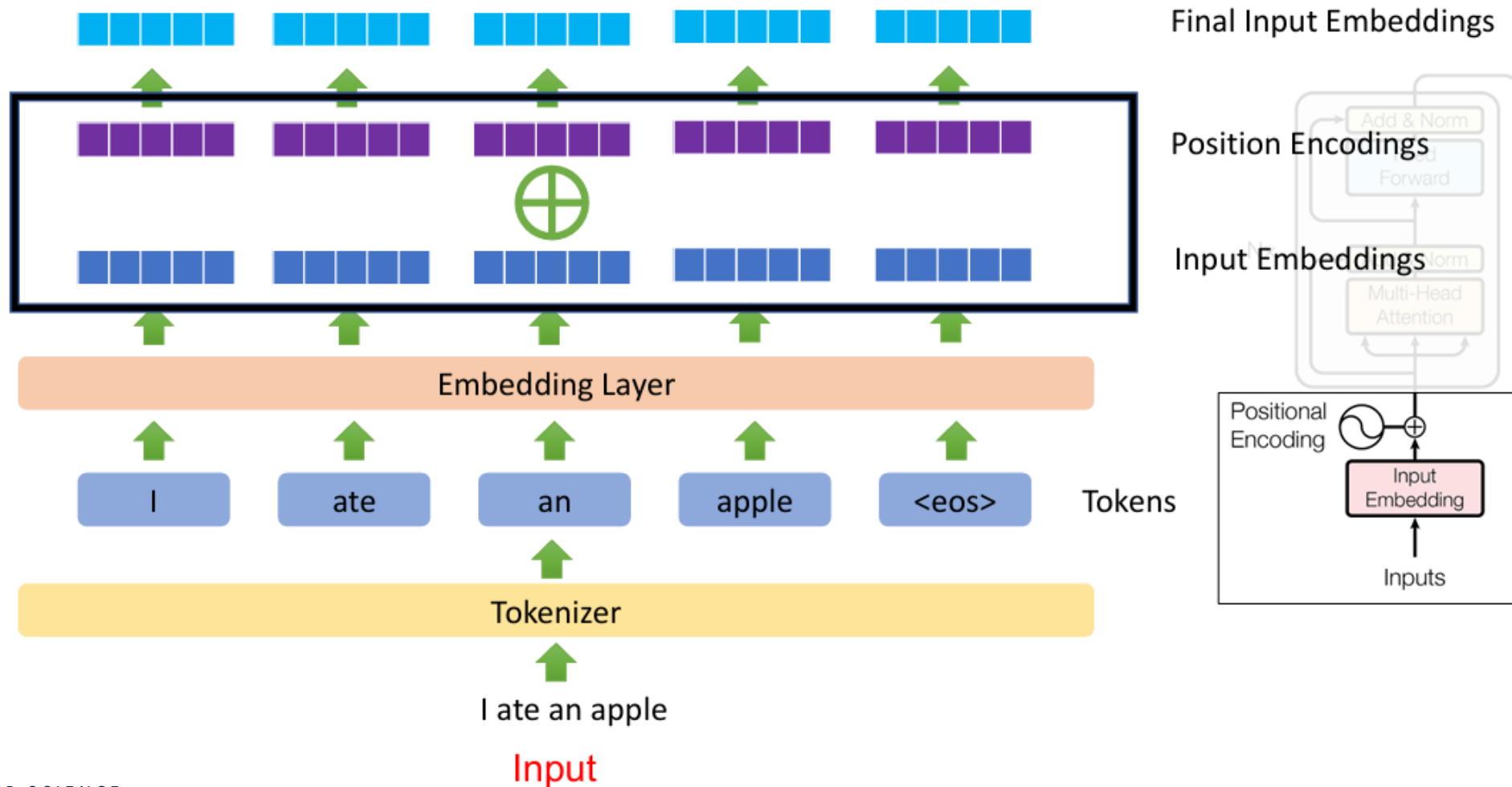
- The position of words in a sentence carries information!
- **Idea:** Add some information to the representation at the beginning that indicates where it is in the sentence



# Position Encodings




Absolute positional encoding (original Transformer): The token embeddings and the absolute position embedding are **added** together element-wise.





- The original position embedding in Transformer is not a great idea, because absolute position is less important than relative position

I walk my dog every day  

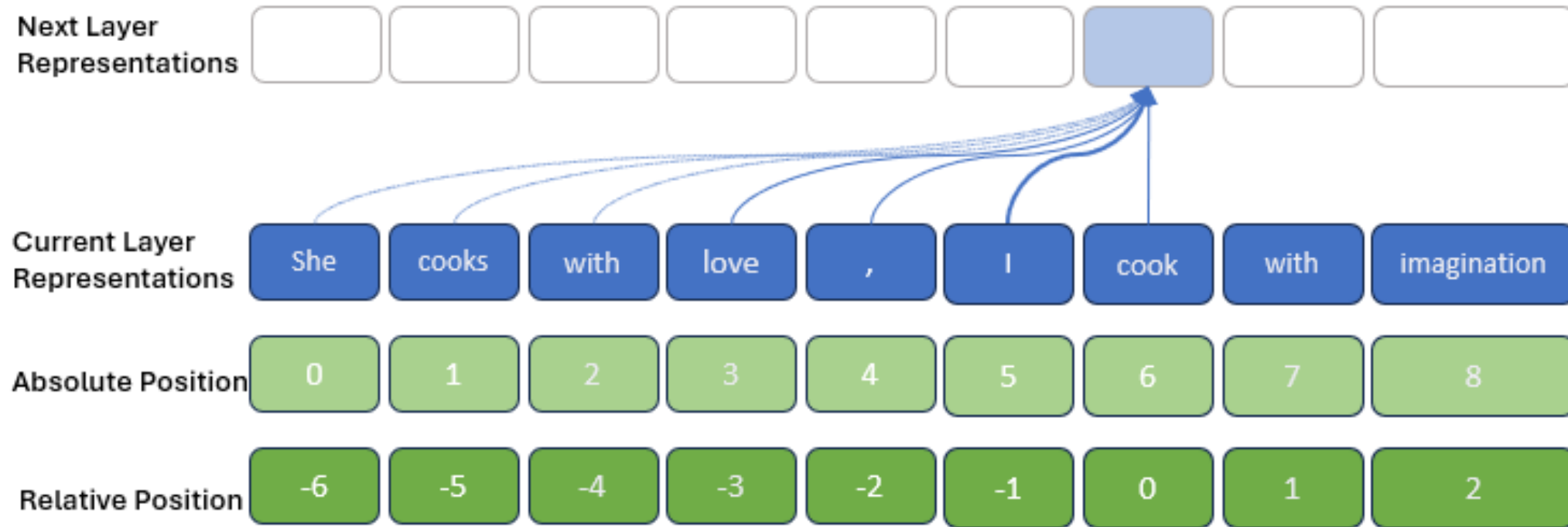

every single day I walk my dog  


The fact that “my dog” is right after “I walk” is the important part, not its absolute position

# Relative Position Encodings



- Translation invariance to position information
- Lead to performance improvements



- Most recent LLMs adopt RoPE

---

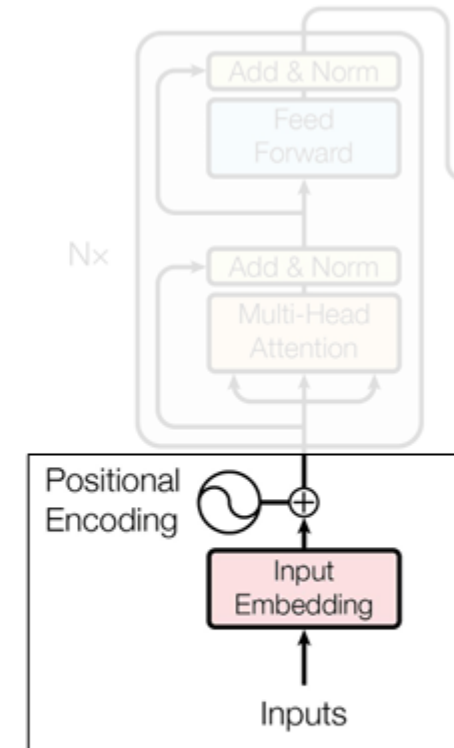
## ROFORMER: ENHANCED TRANSFORMER WITH ROTARY POSITION EMBEDDING

---

$$f_{\{q,k\}}(\mathbf{x}_m, m) = \begin{pmatrix} \cos m\theta & -\sin m\theta \\ \sin m\theta & \cos m\theta \end{pmatrix} \begin{pmatrix} W_{\{q,k\}}^{(11)} & W_{\{q,k\}}^{(12)} \\ W_{\{q,k\}}^{(21)} & W_{\{q,k\}}^{(22)} \end{pmatrix} \begin{pmatrix} x_m^{(1)} \\ x_m^{(2)} \end{pmatrix}$$

Table 2: Comparing RoFormer and BERT by fine tuning on downstream GLEU tasks.

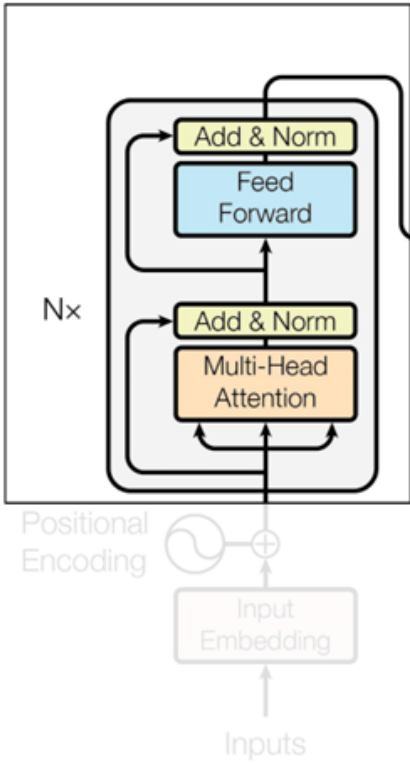
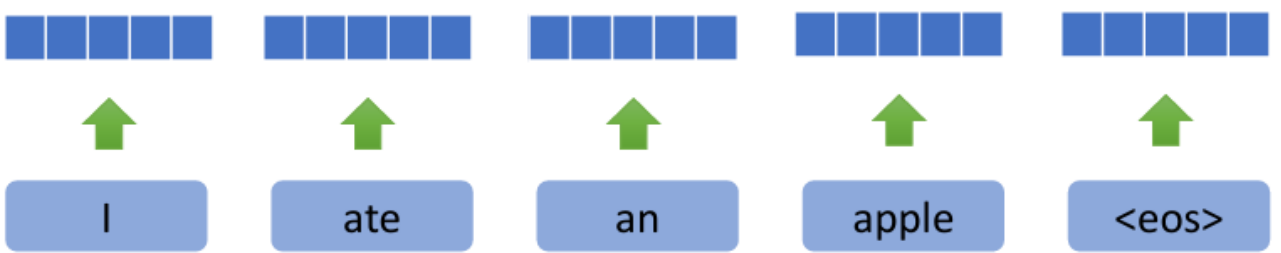
| Model                     | MRPC        | SST-2 | QNLI | STS-B       | QQP         | MNLI(m/mm) |
|---------------------------|-------------|-------|------|-------------|-------------|------------|
| BERT Devlin et al. [2019] | 88.9        | 93.5  | 90.5 | 85.8        | 71.2        | 84.6/83.4  |
| RoFormer                  | <b>89.5</b> | 90.7  | 88.0 | <b>87.0</b> | <b>86.4</b> | 80.2/79.8  |

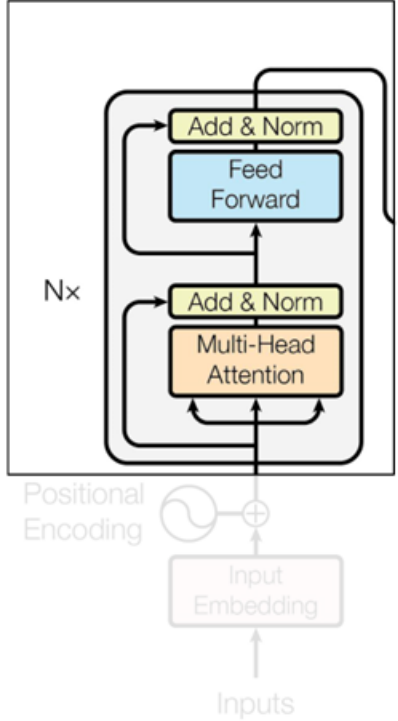
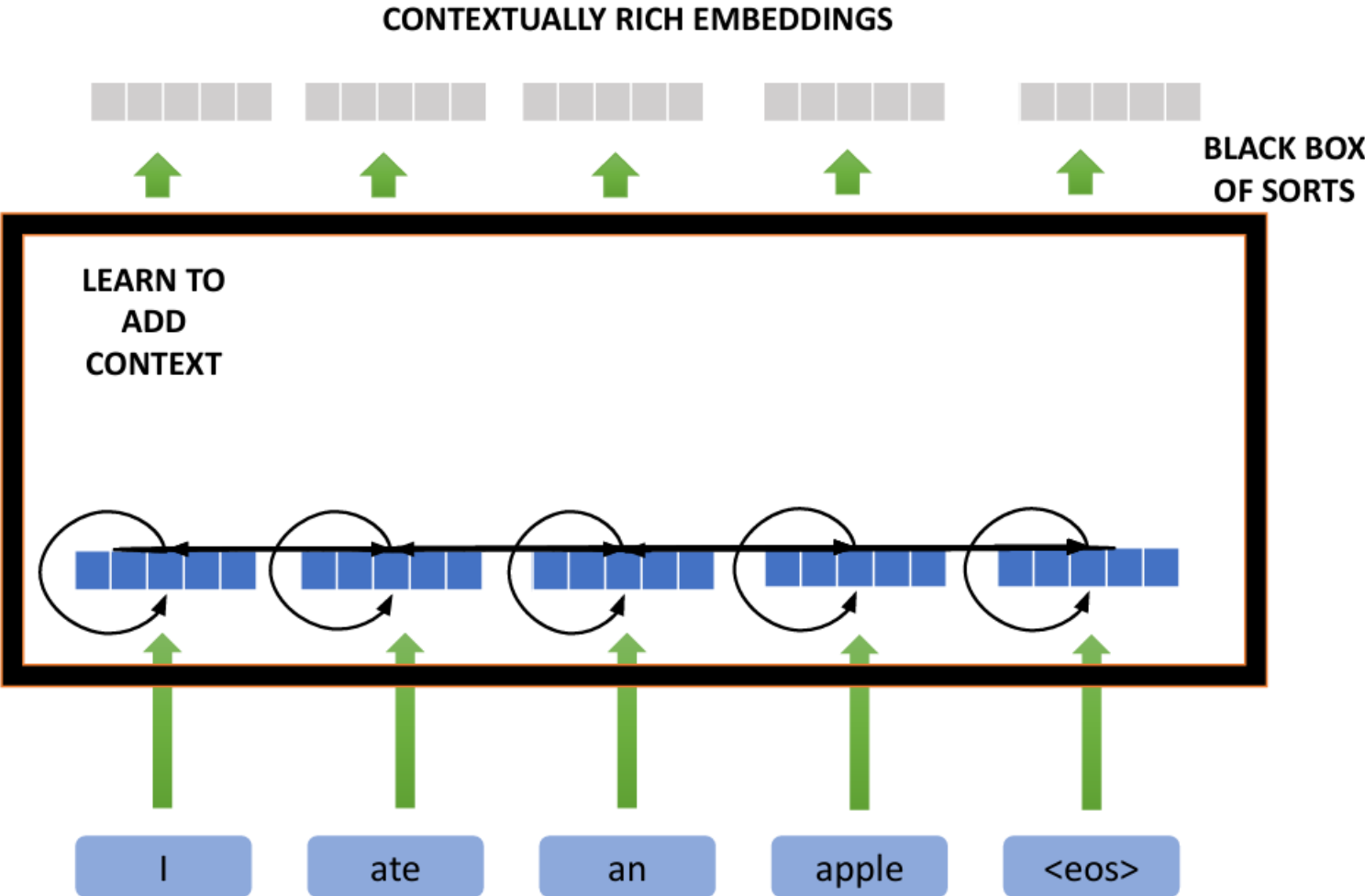


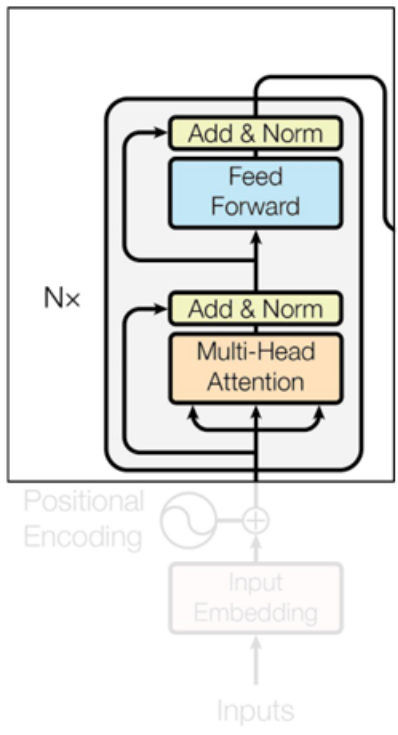
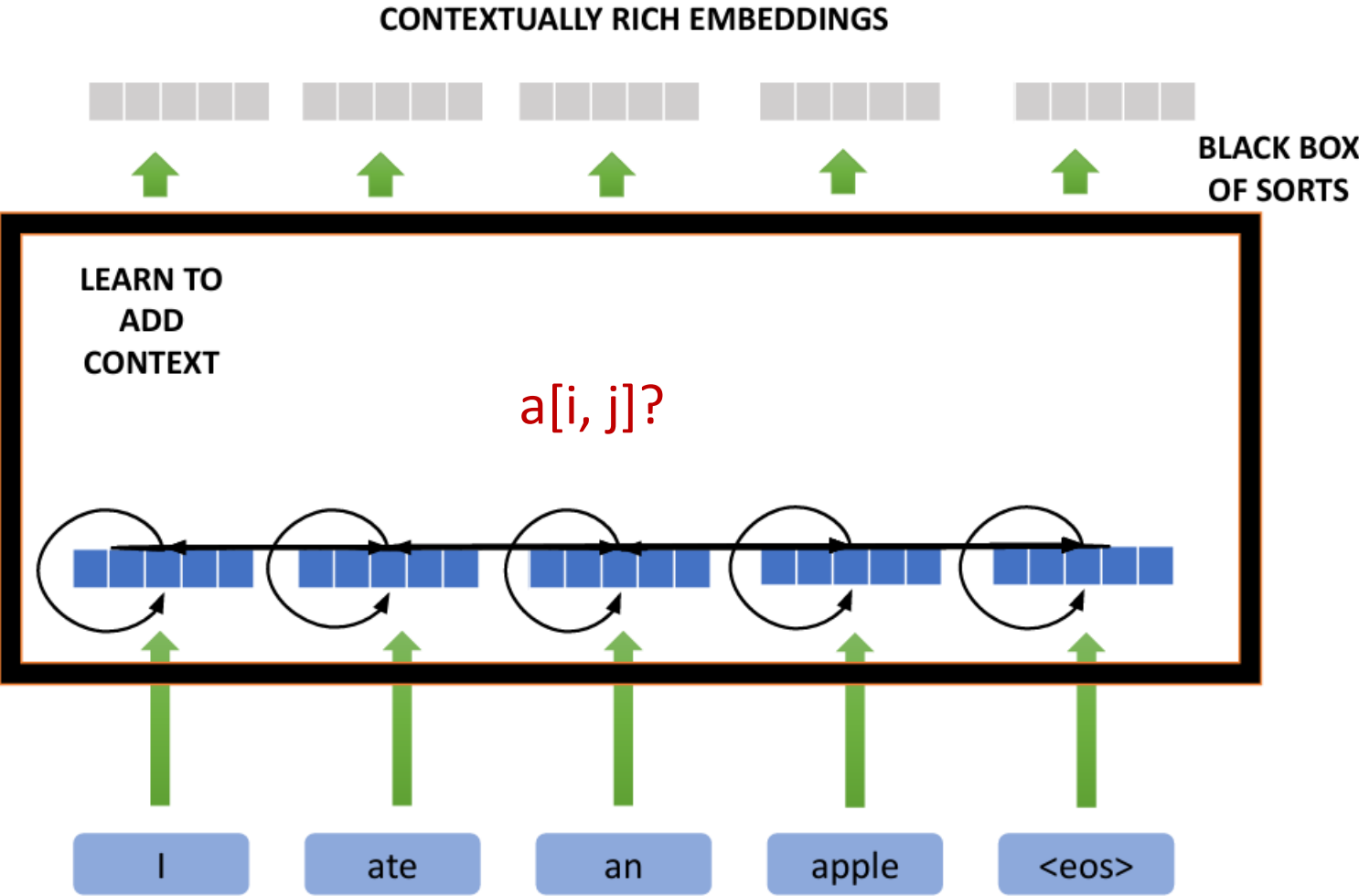
[REF: Rotary Position Embeddings](#)



**WHERE IS THE  
CONTEXT ?**



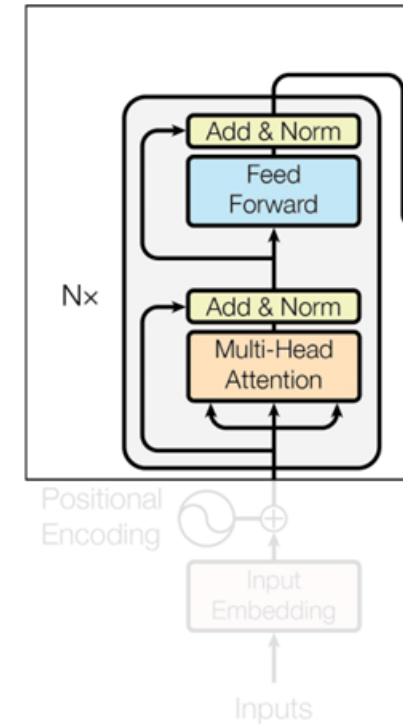




- Before “multi-head” attention, what is “single head” attention?

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

- Query
- Key
- Value



- Database {key, value store}

{Query: "Order details of **order\_104**"}

OR

{Query: "Order details of **order\_106**"}

```
{ "order_100": { "items": "a1", "delivery_date": "a2", ... } },  
{ "order_101": { "items": "b1", "delivery_date": "b2", ... } },  
{ "order_102": { "items": "c1", "delivery_date": "c2", ... } },  
{ "order_103": { "items": "d1", "delivery_date": "d2", ... } },  
{ "order_104": { "items": "e1", "delivery_date": "e2", ... } },  
{ "order_105": { "items": "f1", "delivery_date": "f2", ... } },  
{ "order_106": { "items": "g1", "delivery_date": "g2", ... } },  
{ "order_107": { "items": "h1", "delivery_date": "h2", ... } },  
{ "order_108": { "items": "i1", "delivery_date": "i2", ... } },  
{ "order_109": { "items": "j1", "delivery_date": "j2", ... } },  
{ "order_110": { "items": "k1", "delivery_date": "k2", ... } }
```

## Query

1. Search for info

## Key

1. Interacts directly with Queries
2. Distinguishes one object from another
3. Identify which object is the most relevant and by how much

## Value

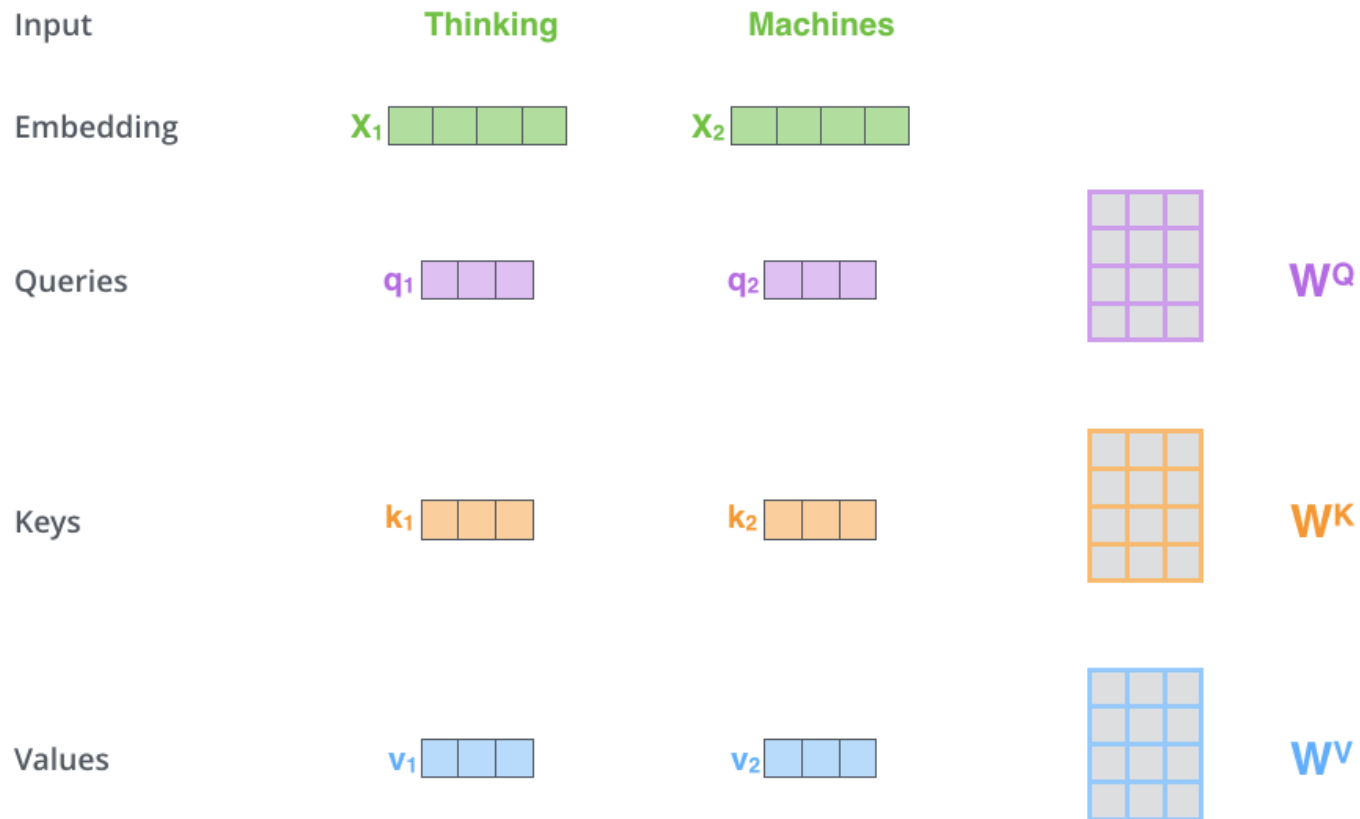
1. Actual details of the object
2. More fine grained



# Self-Attention



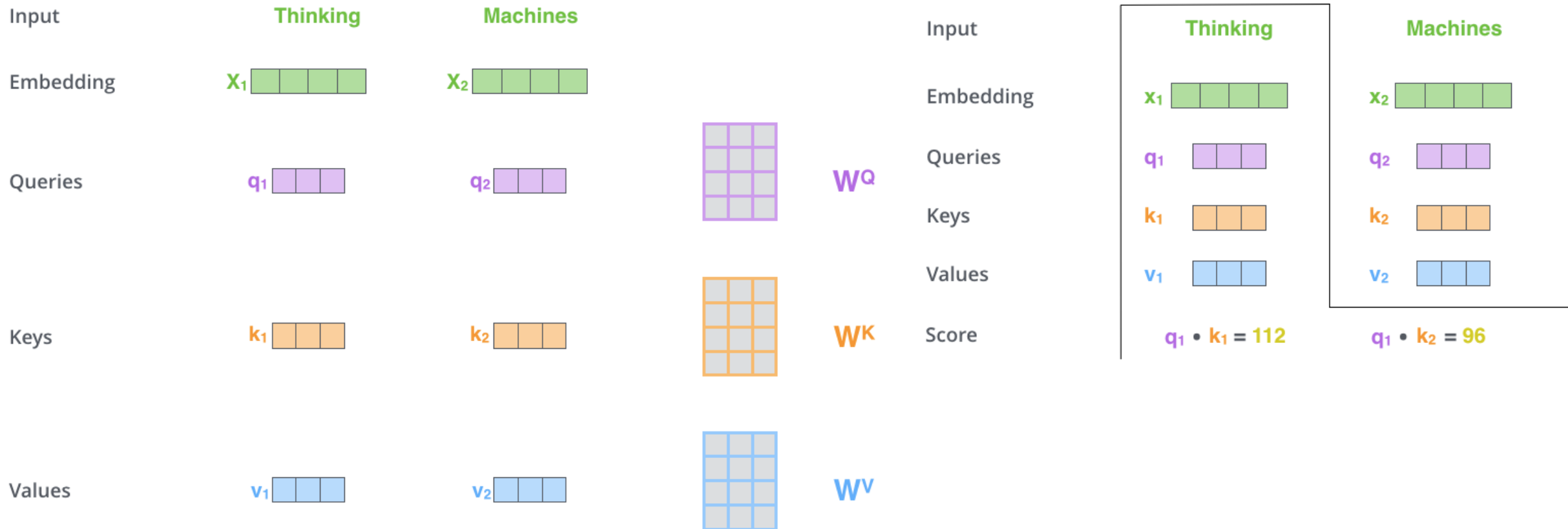
- **Step 1:** compute “key, value, query” embedding for each input token



# Self-Attention



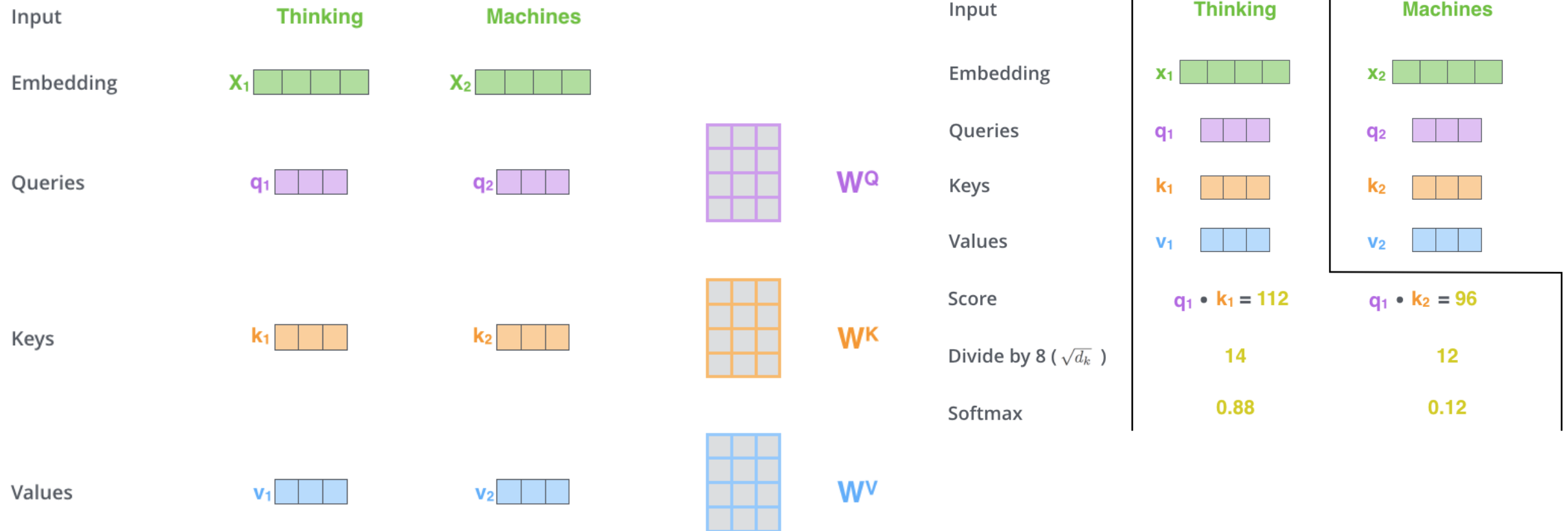
- **Step 1:** compute “key, value, query” embedding for each input token
- **Step 2:** compute scores between pairs of tokens



# Self-Attention



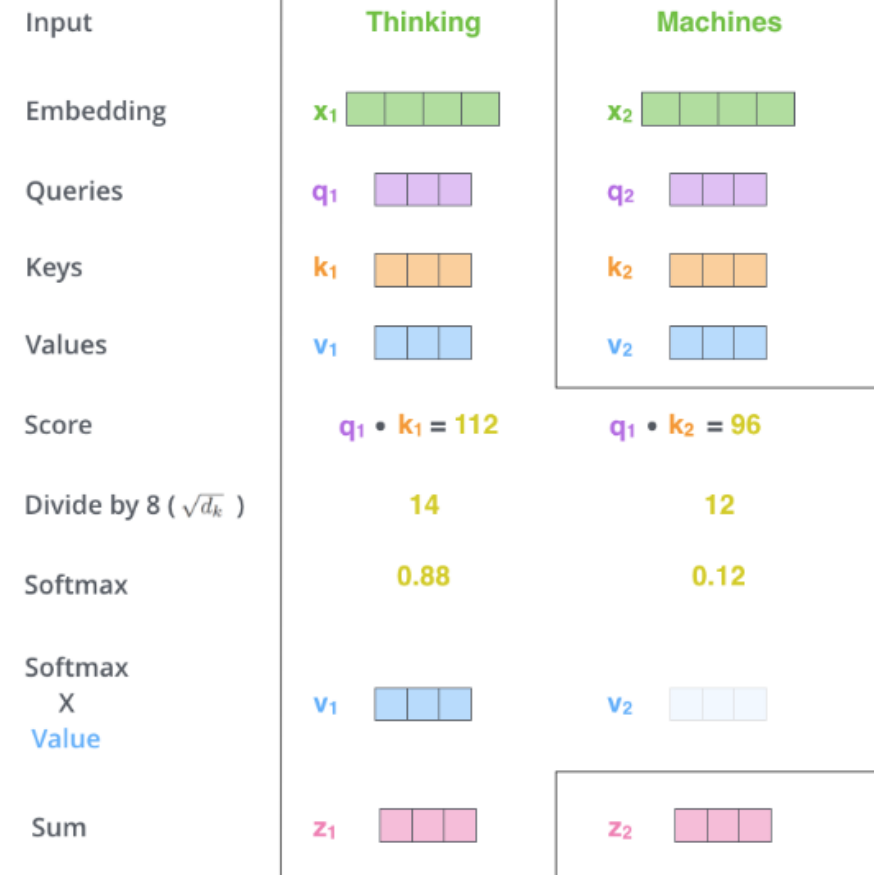
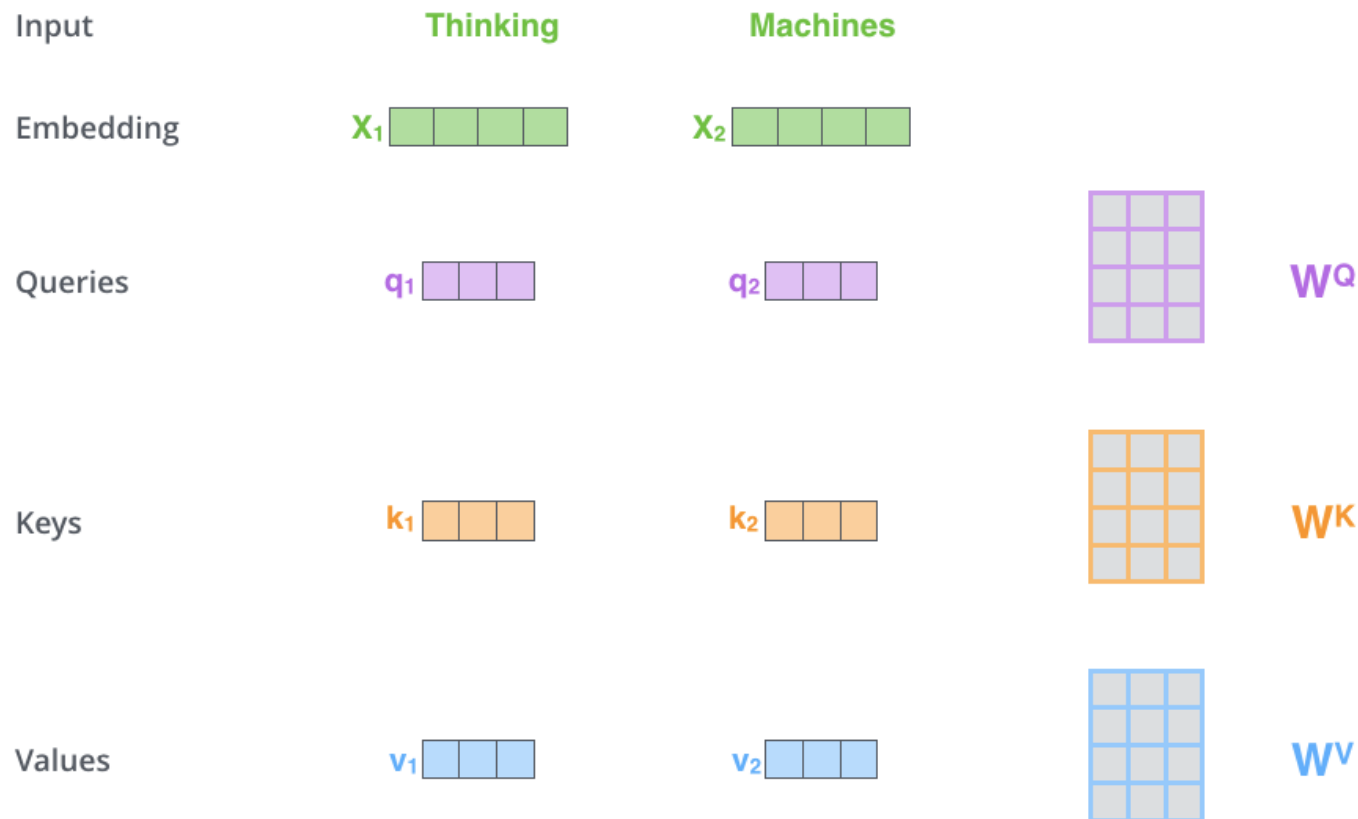
- **Step 1:** compute “key, value, query” embedding for each input token
- **Step 2:** compute scores between pairs of tokens
- **Step 3:** compute normalized attention scores



# Self-Attention



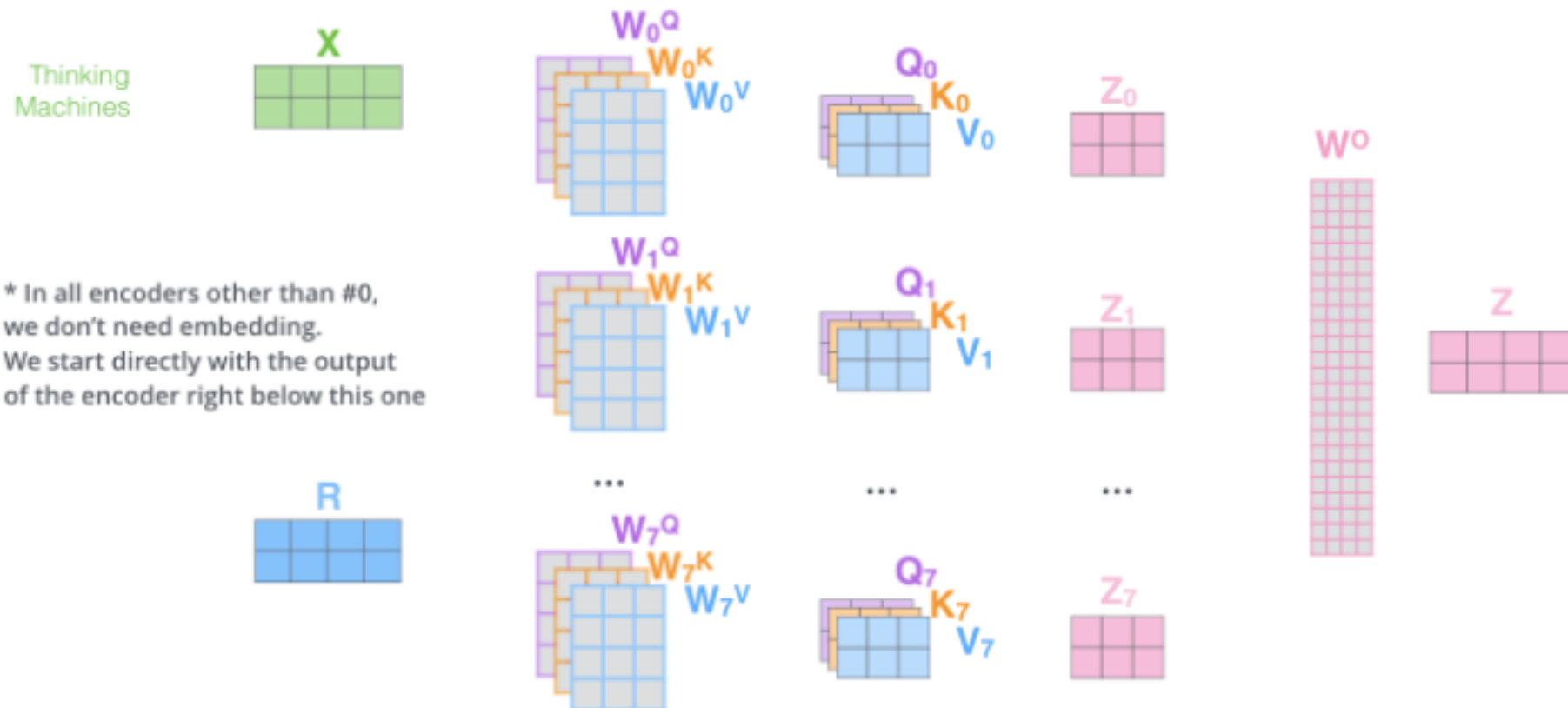
- **Step 1:** compute “key, value, query” embedding for each input token
- **Step 2:** compute scores between pairs of tokens
- **Step 3:** compute normalized attention scores
- **Step 4:** get new representation by weighted sum of values



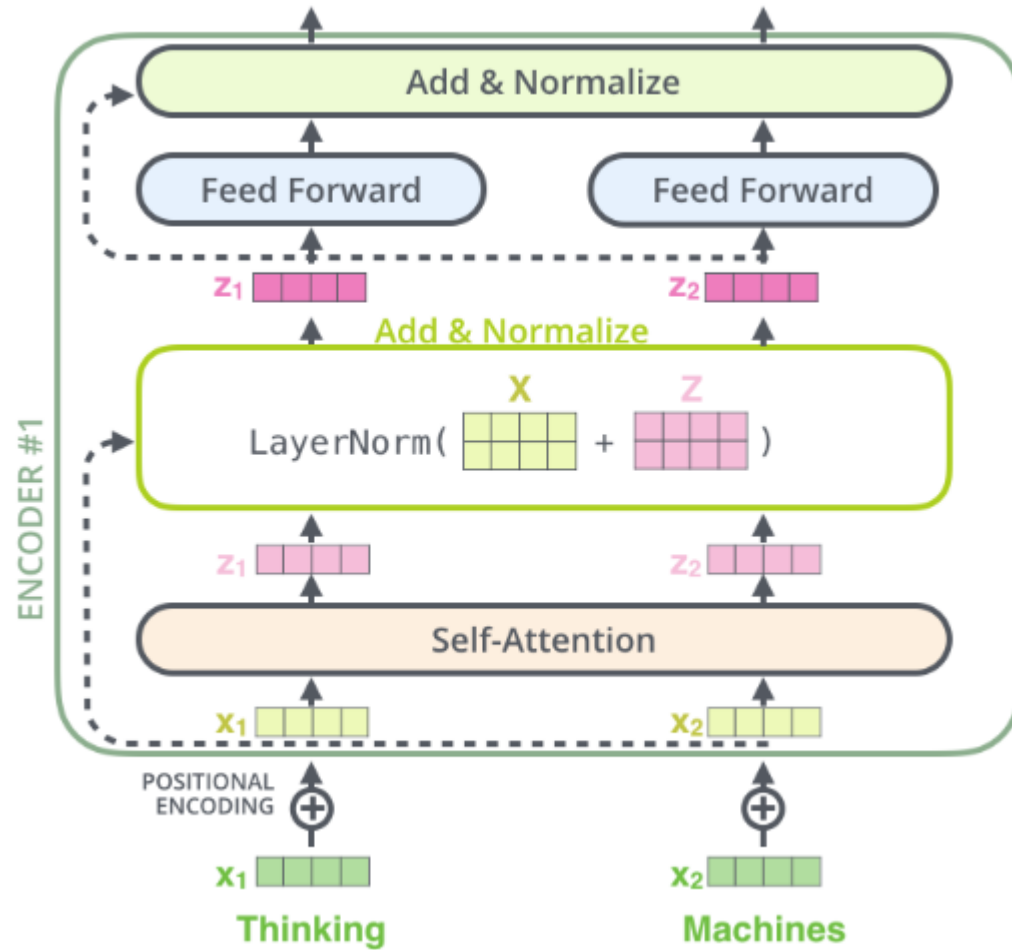
# Multi-head Attention



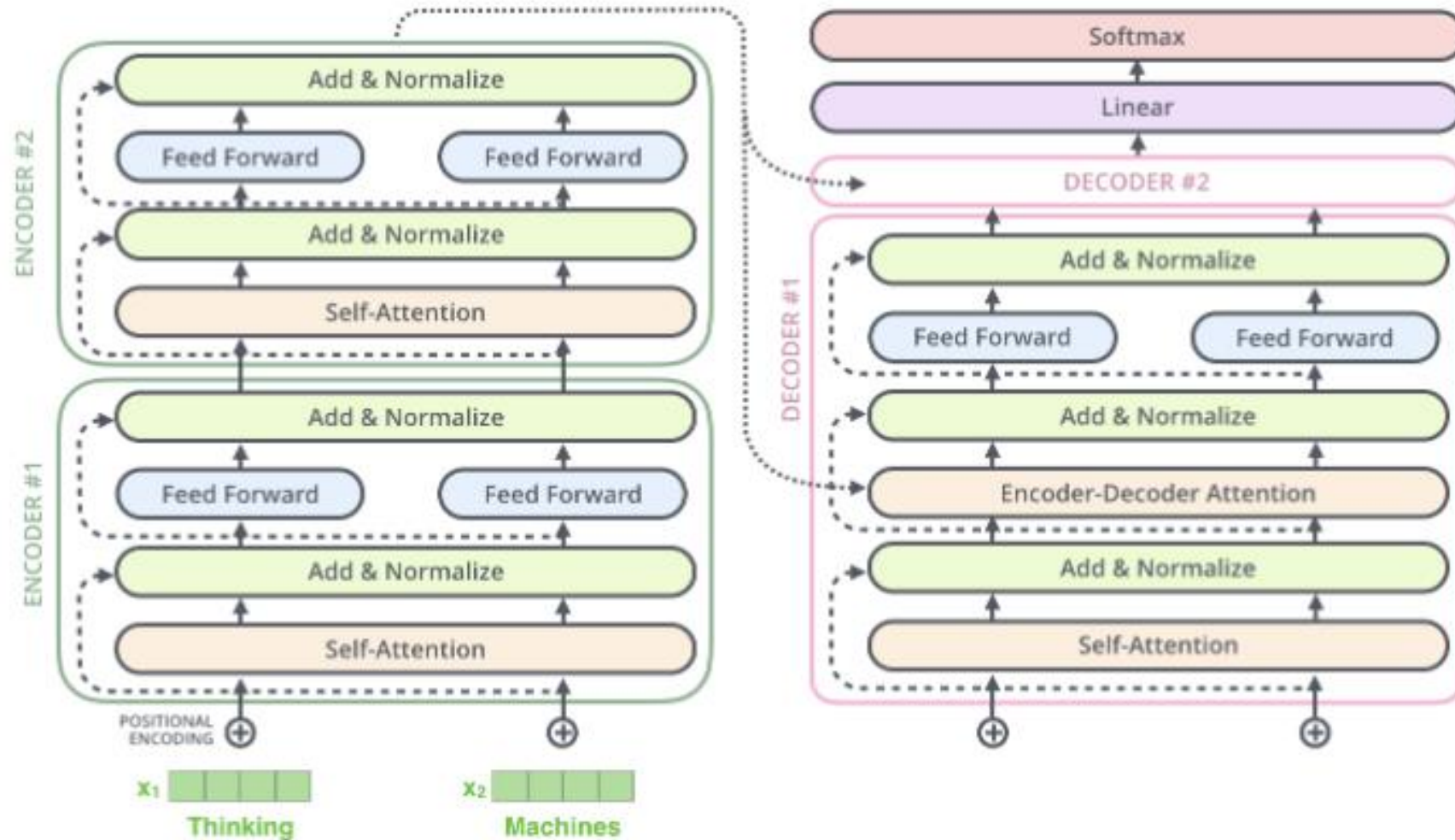
- Do many attention head calculation in parallel, and combine
- Each head has its own set of parameters
- Different heads can learn different “interactions” between inputs



# The Residuals and LayerNorm



# The Decoder Side

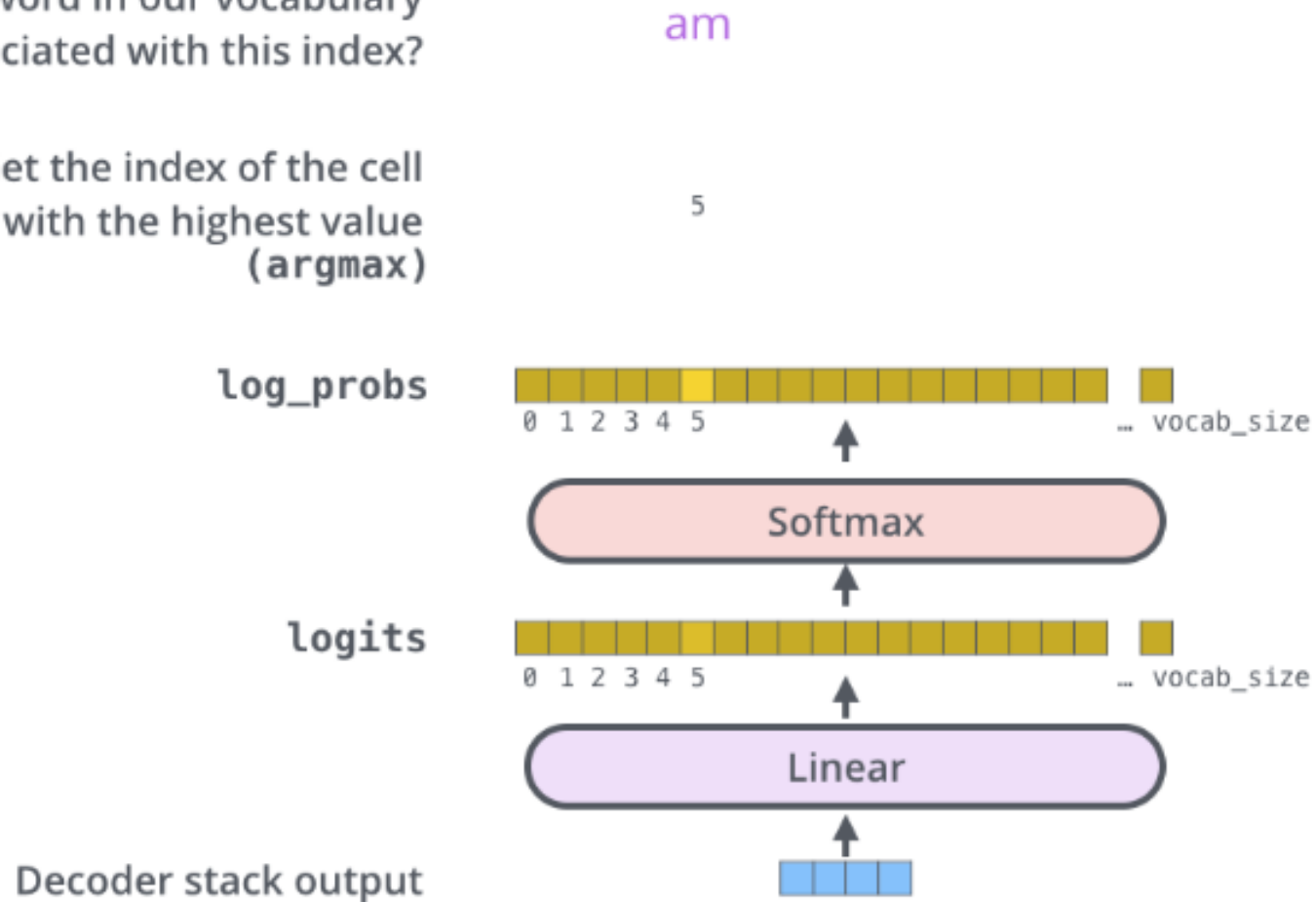


# The Final Linear and Softmax Layer



Which word in our vocabulary  
is associated with this index?

Get the index of the cell  
with the highest value  
(**argmax**)





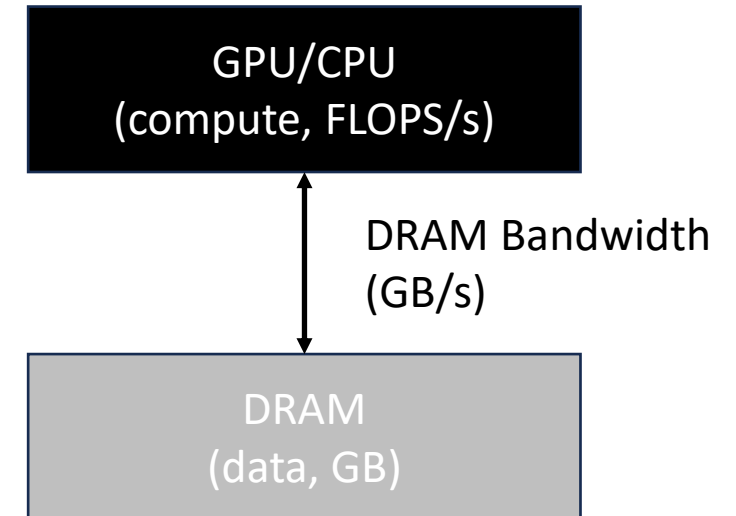
$$AI = \#ops / \#bytes$$

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

Used to evaluate the efficiency of computational algorithms

System performance is bound by

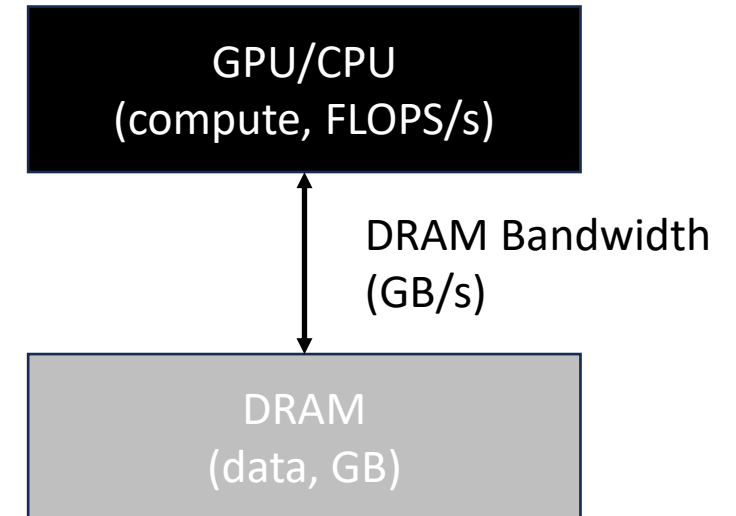
- 1) The peak compute TFLOPS
- 2) The memory bandwidth



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

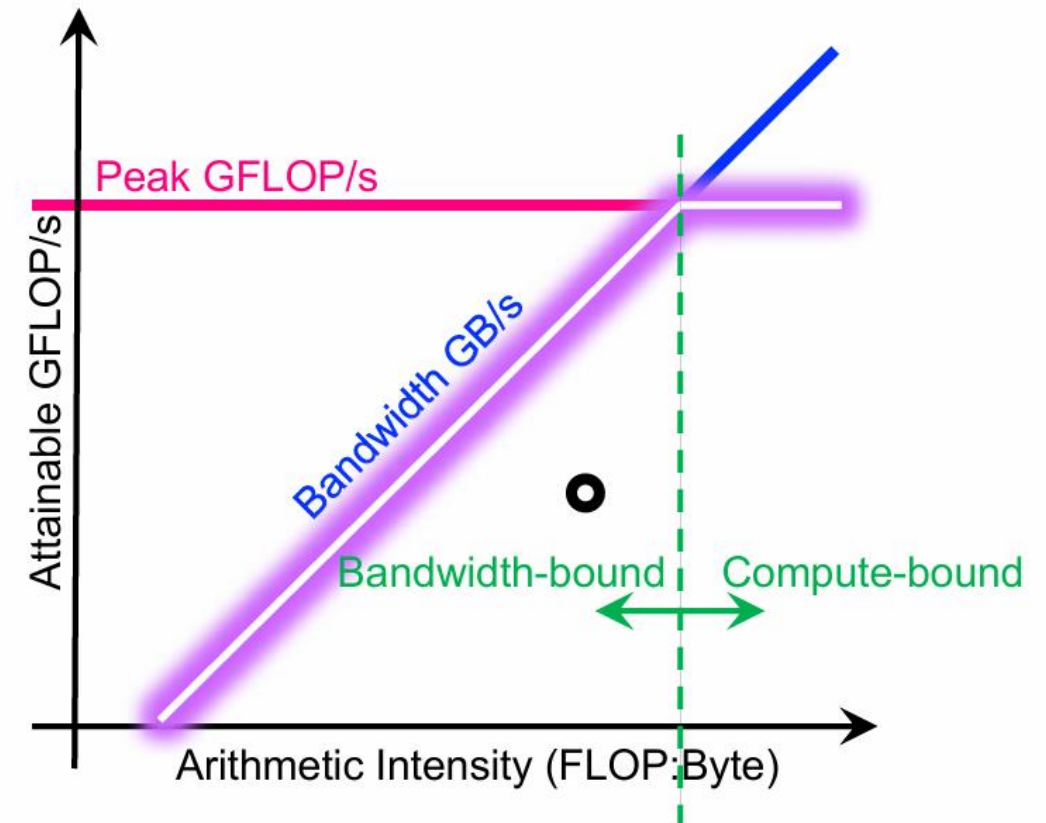
System performance is bound by

$$TFLOPS/s = \min(\text{Peak TFLOPS/s}, \text{Peak bw GB/s} * AI)$$



The Roofline model provides a relatively simple way for performance estimates based on the computation of workload and hardware characteristics

- High AI: Compute-bound
- Low AI: Memory bandwidth bound



# Why Roofline Model?

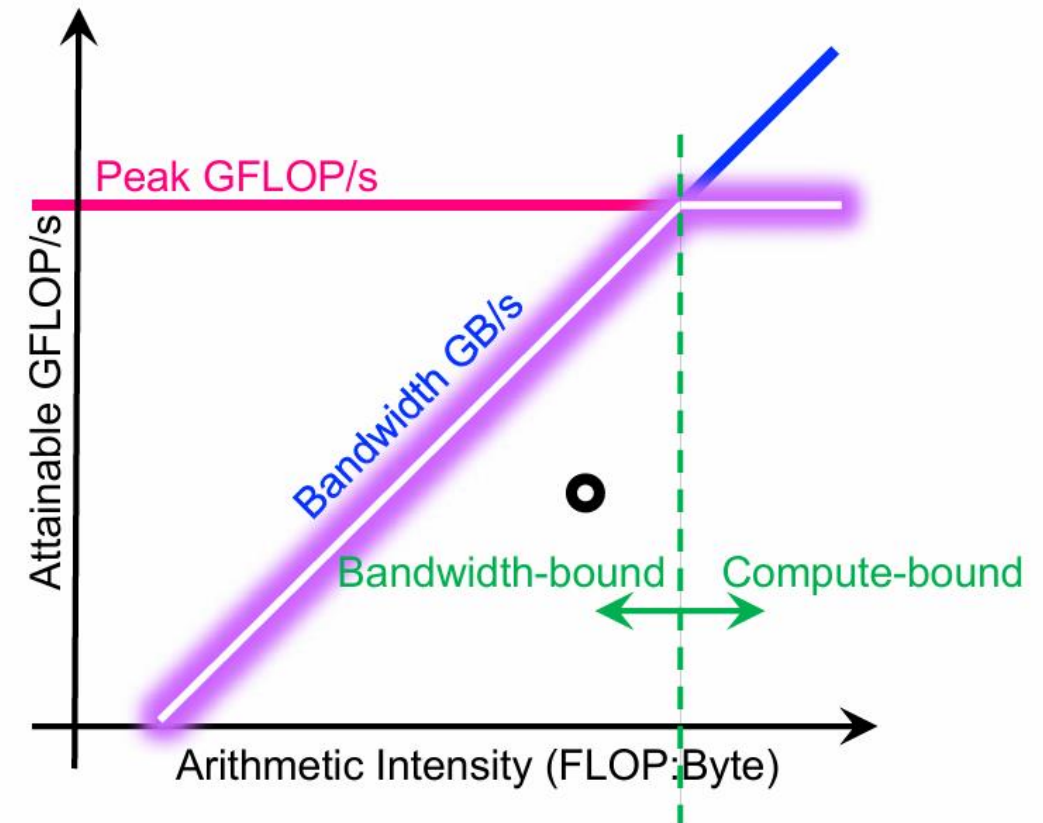


Helps identify the bottlenecks

Program performance depends on how well it fits the hardware architecture

Create optimizations to exhaust both compute and bandwidth **at the same time** (many times it is impossible)

The mode also tells you when to stop



$$\begin{aligned} \text{AI\_A100} &= \text{compute} / \text{memory\_bandwidth} \\ &= 312 \text{ TFLOPS} / 800 \text{ GB/s} \\ &= 390 \text{ ops/ byte} \end{aligned}$$

## NVIDIA A100 TENSOR CORE GPU SPECIFICATIONS (SXM4 AND PCIE FORM FACTORS)

|                                | A100<br>80GB PCIe  | A100<br>80GB SXM                     |
|--------------------------------|--|--------------------------------------|
| FP64                           | 9.7 TFLOPS   |                                      |
| FP64 Tensor Core               | 19.5 TFLOPS  |                                      |
| FP32                           | 19.5 TFLOPS  |                                      |
| Tensor Float 32 (TF32)         | 156 TFLOPS   312 TFLOPS*   |                                      |
| BFLOAT16 Tensor Core           | 312 TFLOPS   624 TFLOPS*   |                                      |
| FP16 Tensor Core               | 312 TFLOPS   624 TFLOPS*   |                                      |
| INT8 Tensor Core               | 624 TOPS   1248 TOPS*  |                                      |
| GPU Memory                     | 80GB HBM2e   | 80GB HBM2e                           |
| GPU Memory Bandwidth           | 1,935GB/s  | 2,039GB/s                            |
| Max Thermal Design Power (TDP) | 300W   | 400W***                              |
| Multi-Instance GPU             | Up to 7 MIGs @ 10GB  | Up to 7 MIGs @ 10GB                  |
| Form Factor                    | PCIe<br>dual-slot air cooled or<br>single-slot liquid cooled       | SXM                                  |
| Interconnect                   | NVIDIA® NVLink® Bridge for 2 GPUs: 600GB/s **<br>PCIe Gen4: 64GB/s | NVLink: 600GB/s<br>PCIe Gen4: 64GB/s |

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

```
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]



```
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)

```
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?



```
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);
```

# Which Program Performs Better?



```
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 =  $3/5$

# Understanding Transformer Calculations



**Input**  $b$ : batch size;  $s$ : sequence length

**Model**  $n$ : the number of attention heads;  $d$ : dimension of single attention head

$h$ : hidden dimension ( $h = n \times d$ )

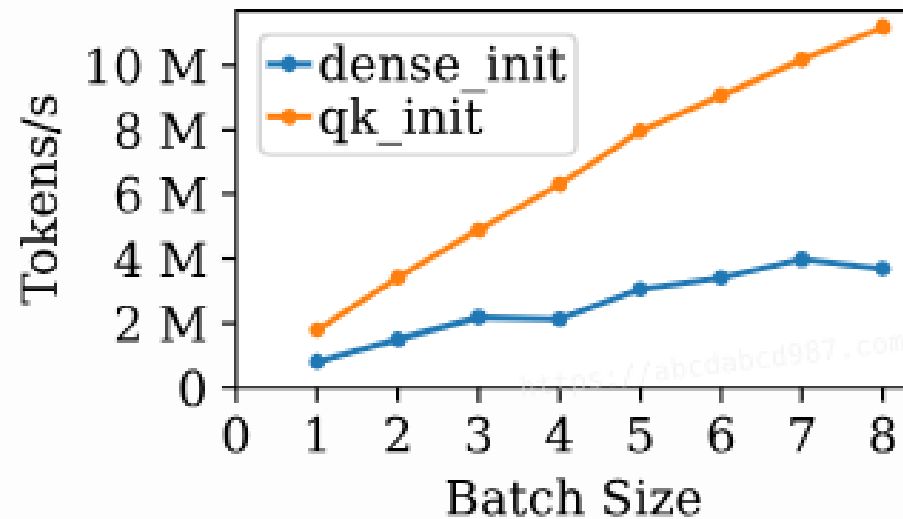
| Symbol                               | Definition                                   | Shape          | FLOP         | I/O                | FLOP:I/O            |
|--------------------------------------|--|----------------|--------------|--------------------|---------------------|
| Input                                |  |                |              |                    |                     |
| $X$                                  | Input for self attention                     | $(b, s, h)$    |              |                    |                     |
| Self-Attention                       |  |                |              |                    |                     |
| $Q^\ddagger, K^\ddagger, V^\ddagger$ | $XW_Q, XW_K, XW_V$                           | $(b, s, h)$    | $O(bsh^2)$   | $O(2bsh + h^2)$    | $O(1/(1/h + 1/bs))$ |
| $Q^\dagger, K^\dagger, V^\dagger$    | Reshape $Q^\ddagger, K^\ddagger, V^\ddagger$ | $(b, s, n, d)$ |              |                    |                     |
| $Q, K, V$                            | Transpose $Q^\dagger, K^\dagger, V^\dagger$  | $(b, n, s, d)$ |              |                    |                     |
| $K^T$                                | Transpose $K$                                | $(b, n, d, s)$ |              |                    |                     |
| $P$                                  | Softmax ( $QK^T/\sqrt{d}$ )                  | $(b, n, s, s)$ | $O(bs^2nd)$  | $O(2bsnd + bs^2n)$ | $O(1/(1/d + 1/s))$  |
| $A^\ddagger$                         | $PV$   | $(b, n, s, d)$ | $O(bs^2nd)$  | $O(2bsnd + bs^2n)$ | $O(1/(1/d + 1/s))$  |
| $A^\dagger$                          | Transpose $A^\ddagger$                       | $(b, s, n, d)$ |              |                    |                     |
| $A$                                  | Reshape $A^\dagger$                          | $(b, s, h)$    |              |                    |                     |
| $Y$                                  | $AW_O$                                       | $(b, s, h)$    | $O(bsh^2)$   | $O(2bsh + h^2)$    | $O(1/(1/h + 1/bs))$ |
| Feed-Forward Network                 |  |                |              |                    |                     |
| $Z$                                  | $\text{ReLU}(YW_1)W_2$                       | $(b, s, h)$    | $O(16bsh^2)$ | $O(2bsh + 8h^2)$   | $O(1/(1/h + 1/bs))$ |

# Transformer Performance: Varying Batch Sizes



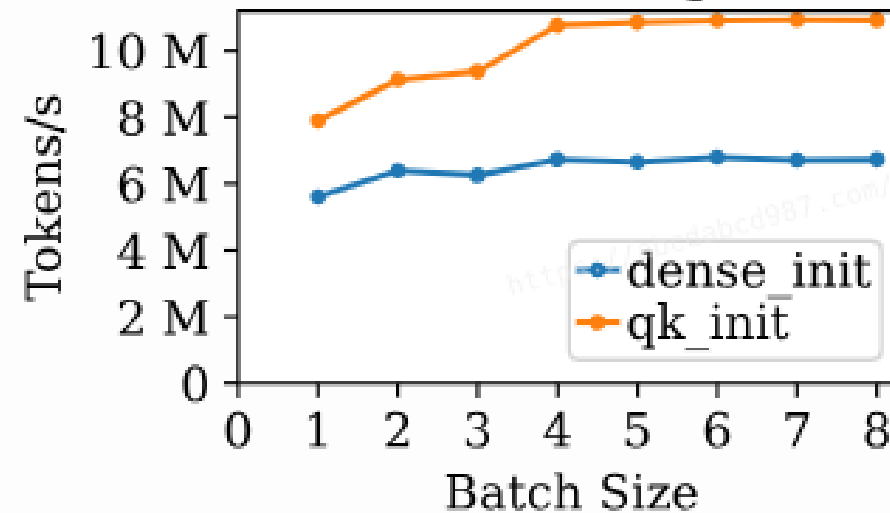
$h=4096$   $s=50$

Initial Stage

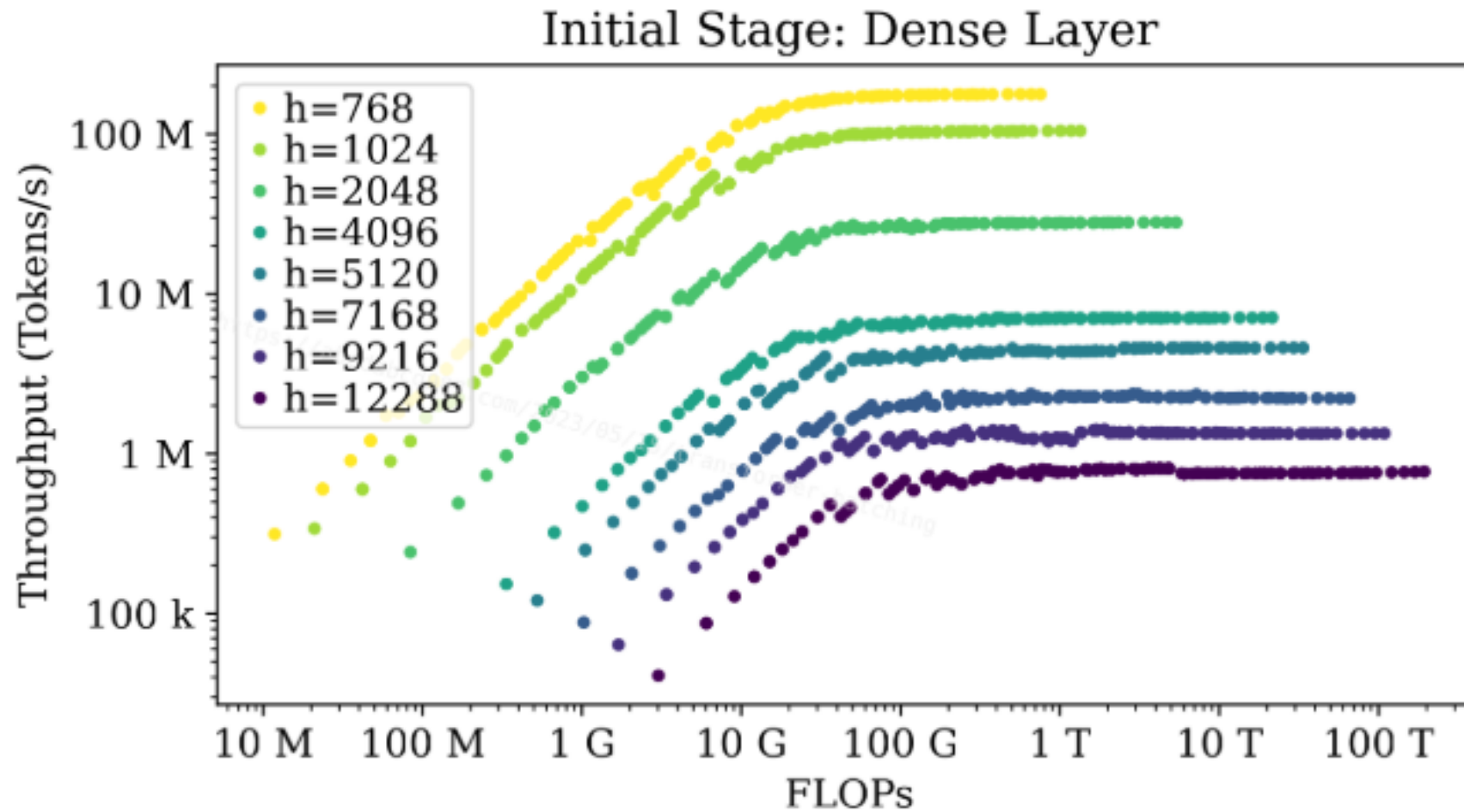


$h=4096$   $s=1000$

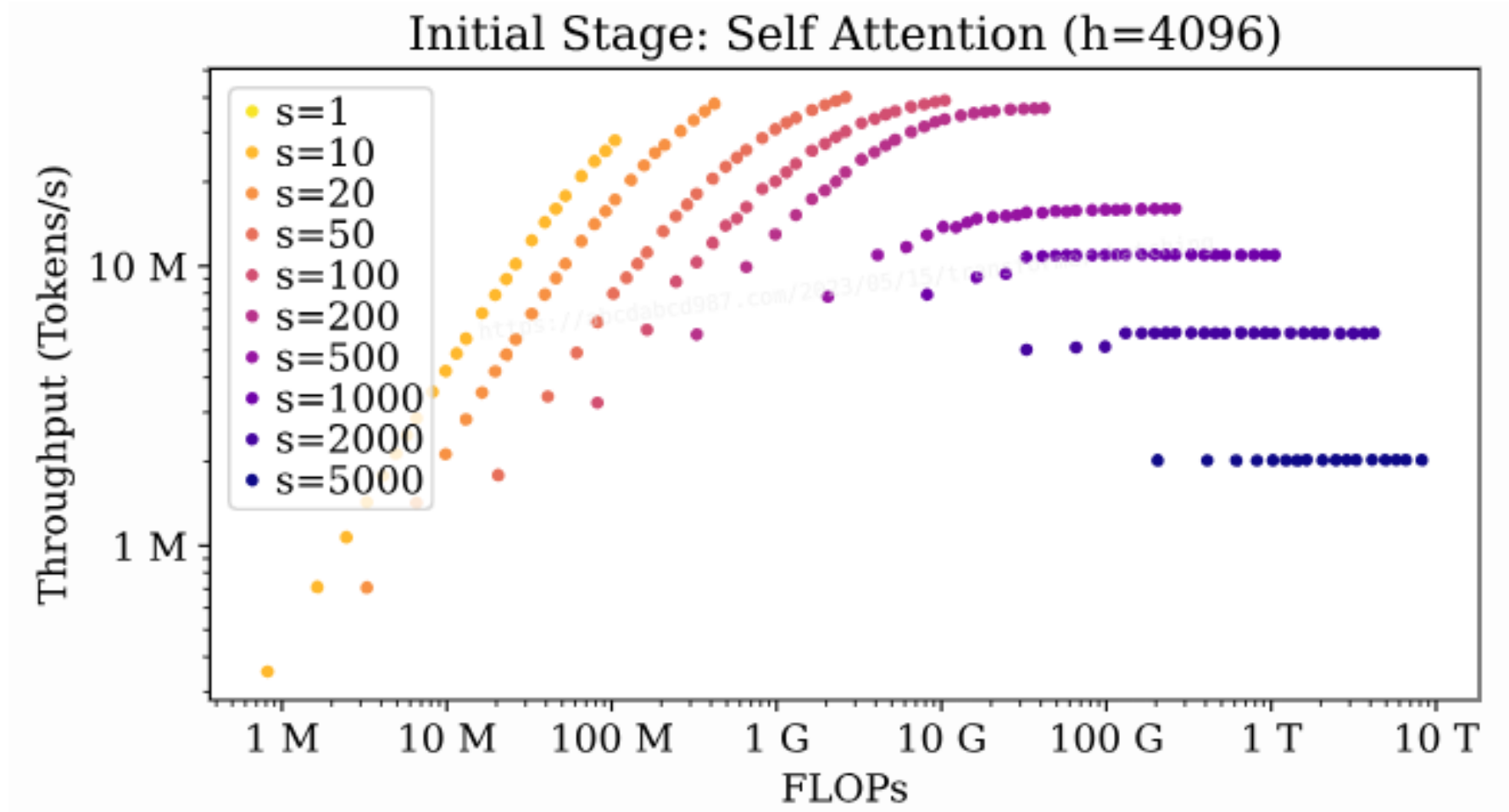
Initial Stage



# Transformer Performance: Batching of Dense Layer

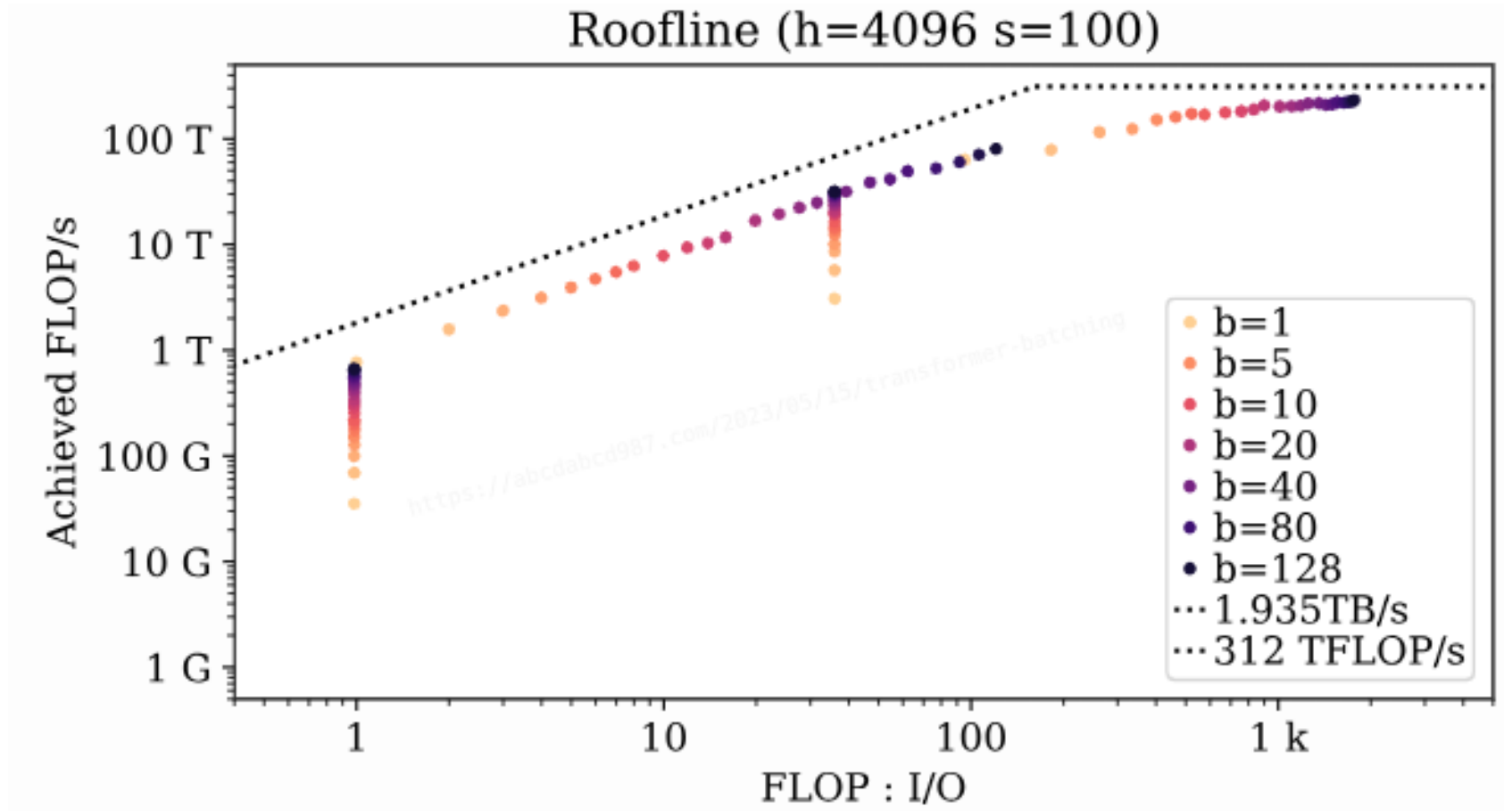


# Transformer Performance: Batching of Self-Attention





# Transformer Performance: Roofline



- After class:
  - Walk through the AI calculation of Transformers
  - How many floating-point operations in total for a 7B decoder-only model?

# Questions?