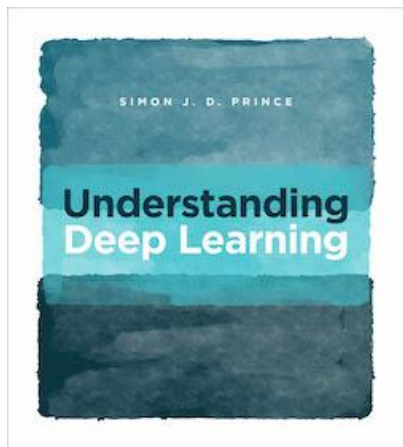


# Transformers for Vision

Computer Vision (SJK02)

Universitat Jaume I



## [Understanding Deep Learning \(MIT, 2023\)](#)

### Chapter 12: Transformers



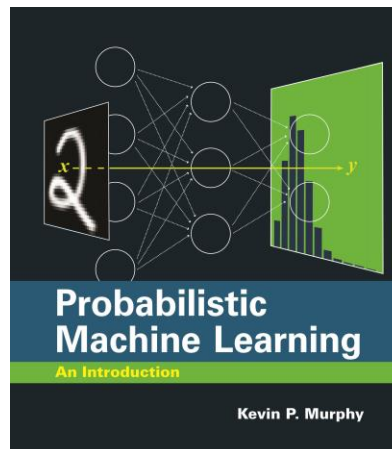
**Luis Serrano**

[Attention mechanisms](#)

[Attention with maths](#)

[Transformers](#)

[Transformer models](#)



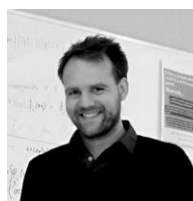
## [Probabilistic Machine Learning \(MIT 2022\)](#)

### Chapter 15: NN for sequences



**Jay Alammar**

[The Illustrated Transformer](#)



**Peter Bloem**

[Transformers from scratch](#)


# Why are they called “transformers”?

Hypothesis 1: Just marketing  
Hypothesis 2: Paradigm shift  
Hypothesis 3: Inputs are transformed  
Other?

# **A bottom-up explanation**

# **Attention mechanism**

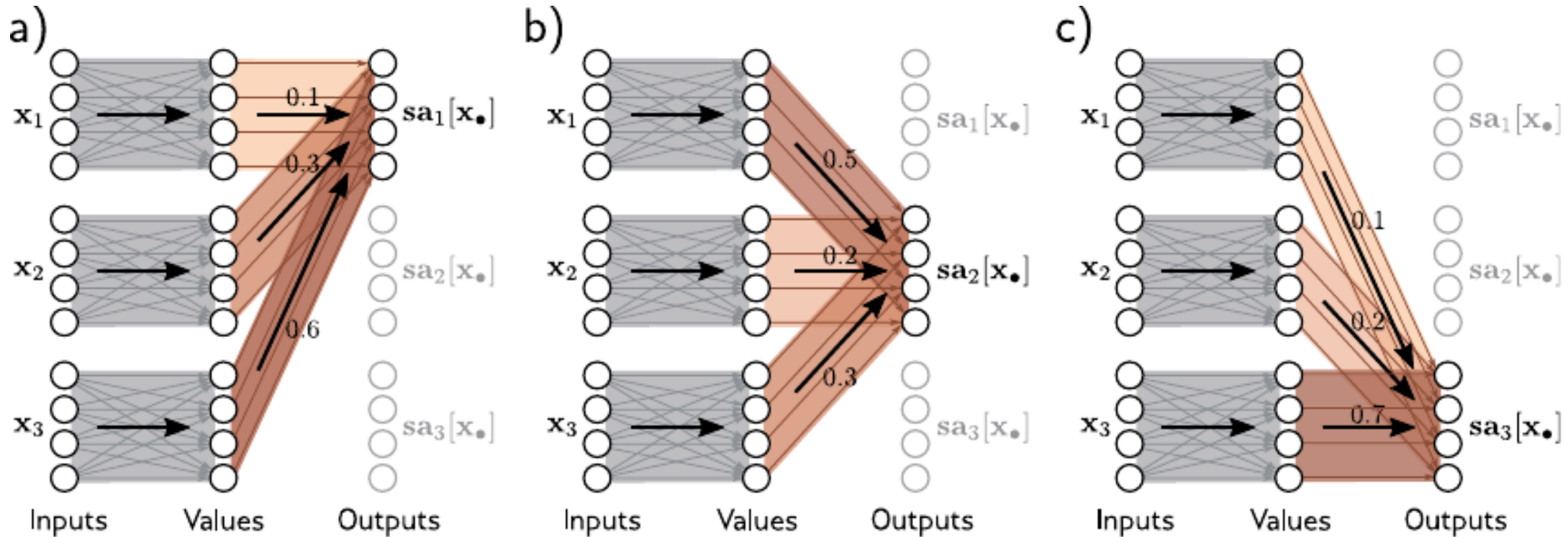
# Standard layer: fixed operation

$$z = \sigma(Wv)$$


fixed

What about something more **flexible**, input-dependent?

# Self-attention as “routing” of values



$$V_i = A_v \cdot x_i + b_v$$

$$sa_i = \sum_{j=1}^N \alpha_{ij} \cdot v_j$$

How many attention weights needed for input sequence of  $N$  elements?

# Attention = soft dictionary lookup

$q$  = query

$V$  = values (set of feature vectors)

$K$  = keys

If  $q$  is most similar to key  $i$ , then use more value  $i$



$$\text{Attn}(q, \{(k_i, v_i)\}) = \sum \alpha_i \cdot v_i$$

$$0 \leq \alpha_i \leq 1$$

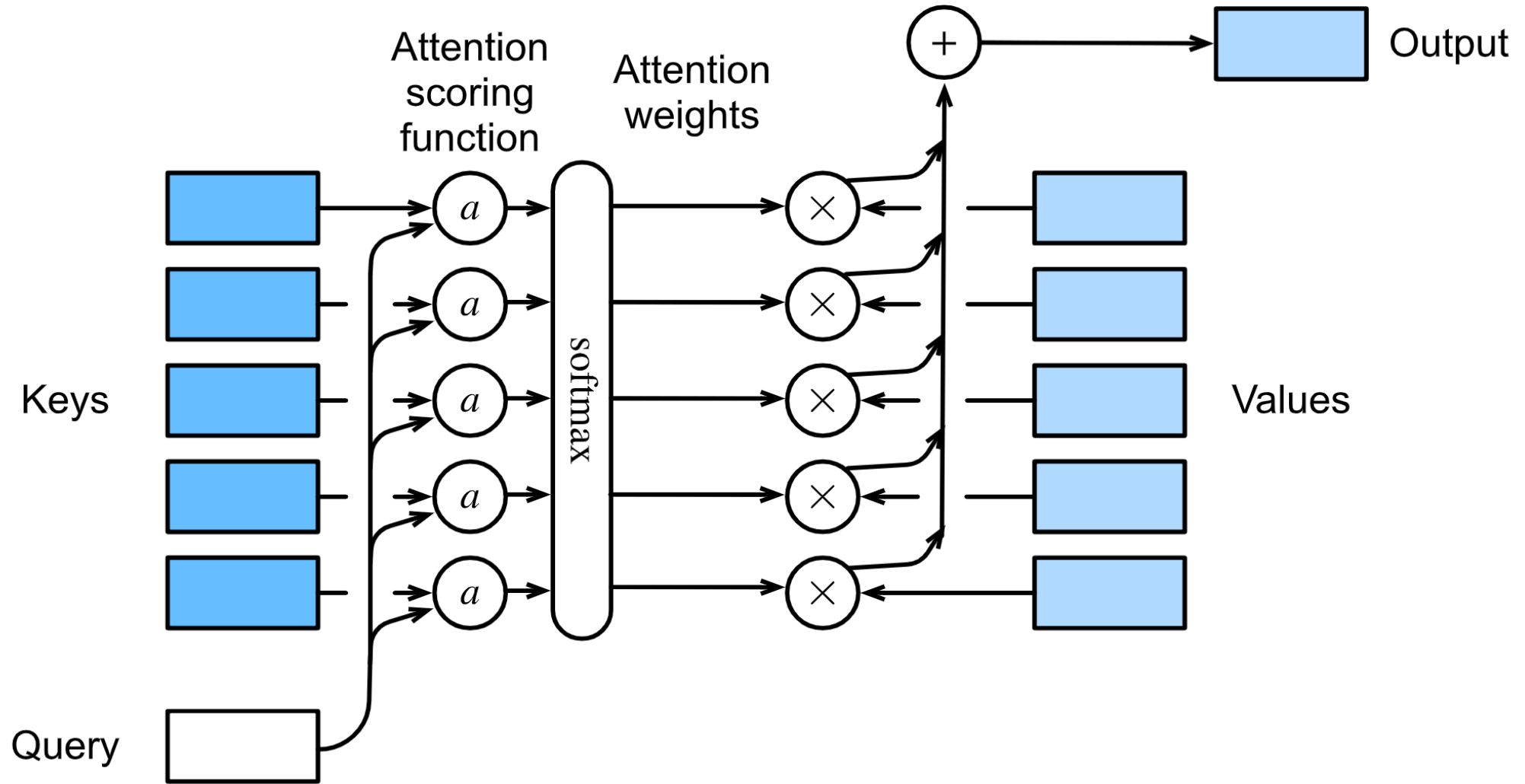
$$\sum \alpha_i = 1$$

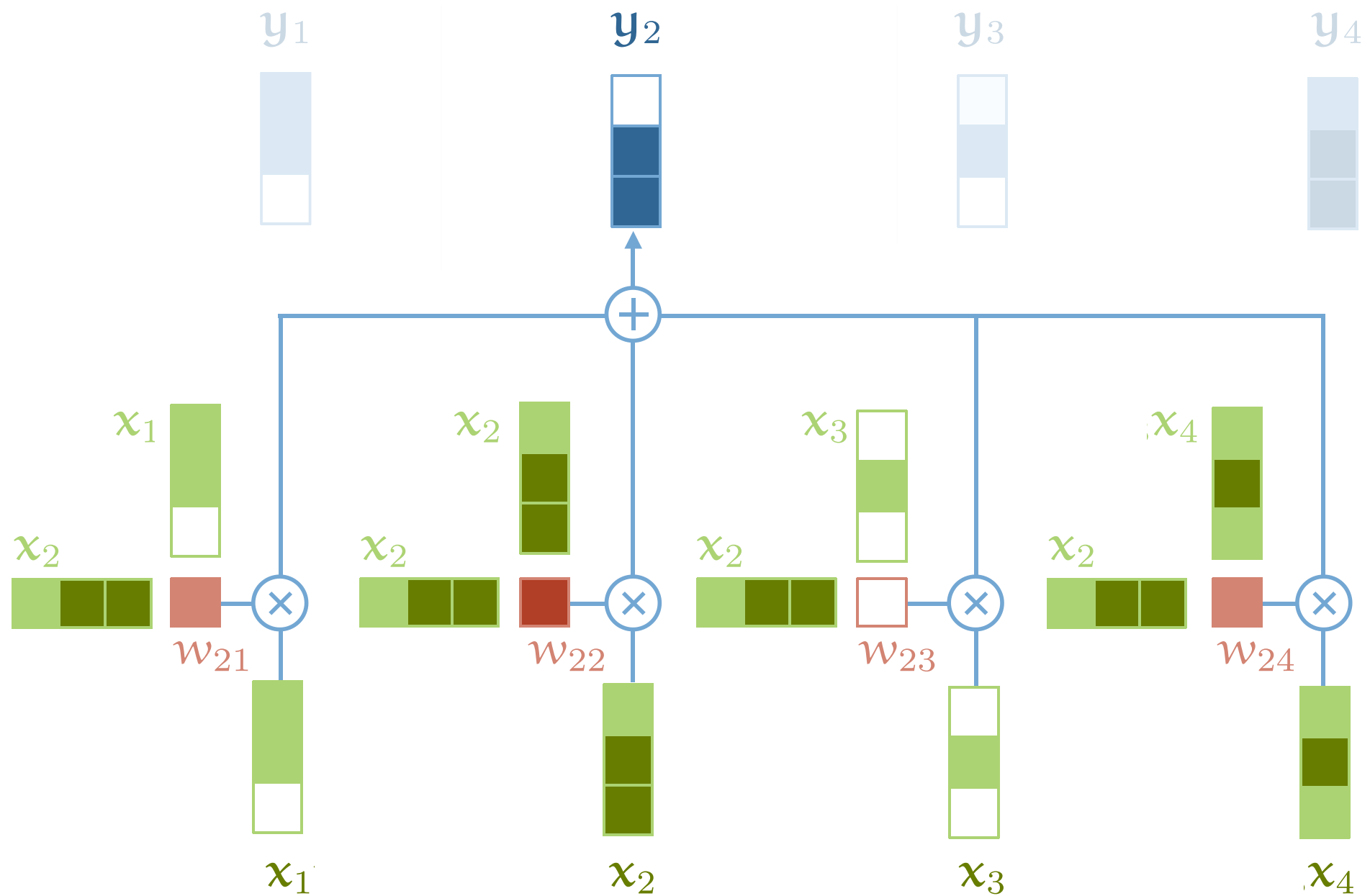
attention weight

attention  
score

$$\alpha_i = \frac{e^{a(q, k_i)}}{\sum_j e^{a(q, k_j)}}$$

Differentiable lookup





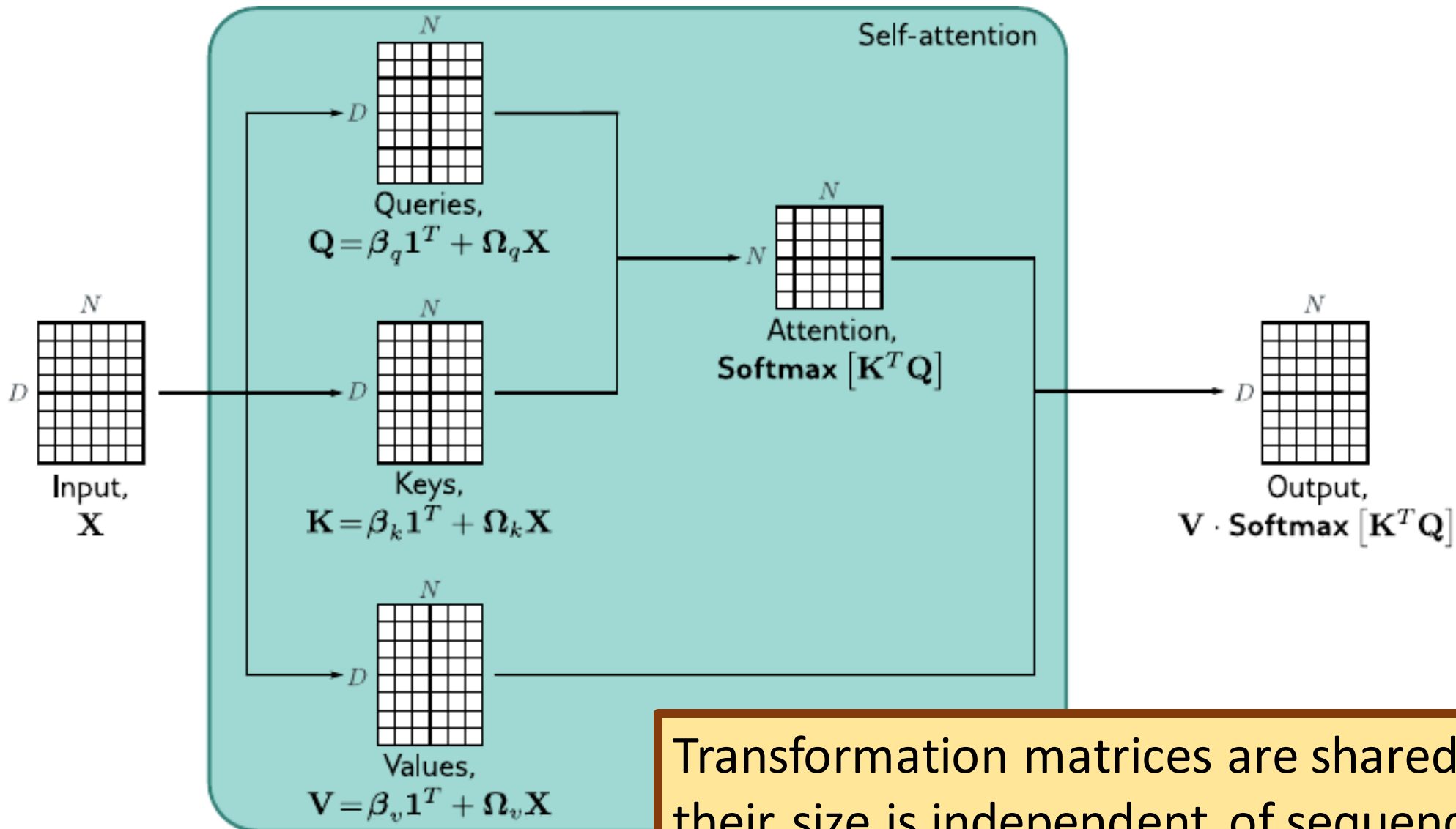
# Scaled dot-product attention

$$q, k \in \mathbb{R}^d \quad a(q, k) = \frac{q^T k}{\sqrt{d}}$$

$$\text{Attn}(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d}}\right)V$$

$$Q \in \mathbb{R}^{n \times d} \quad K \in \mathbb{R}^{m \times d} \quad V \in \mathbb{R}^{m \times d}$$

# Block diagram



Transformation matrices are shared among inputs; their size is independent of sequence length!

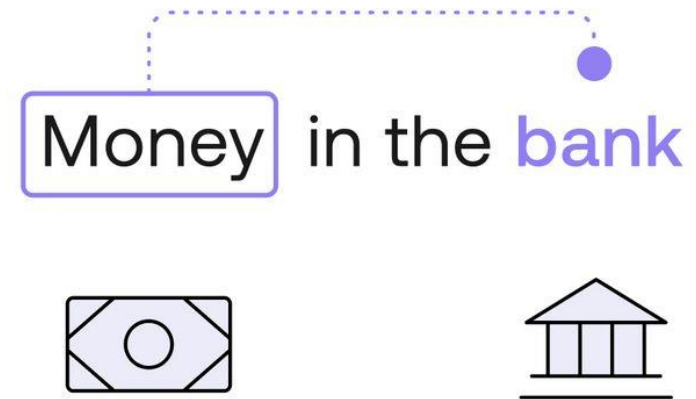
# Visual explanation

Sentence 1: The bank of the river.

Sentence 2: Money in the bank.

## Attention:

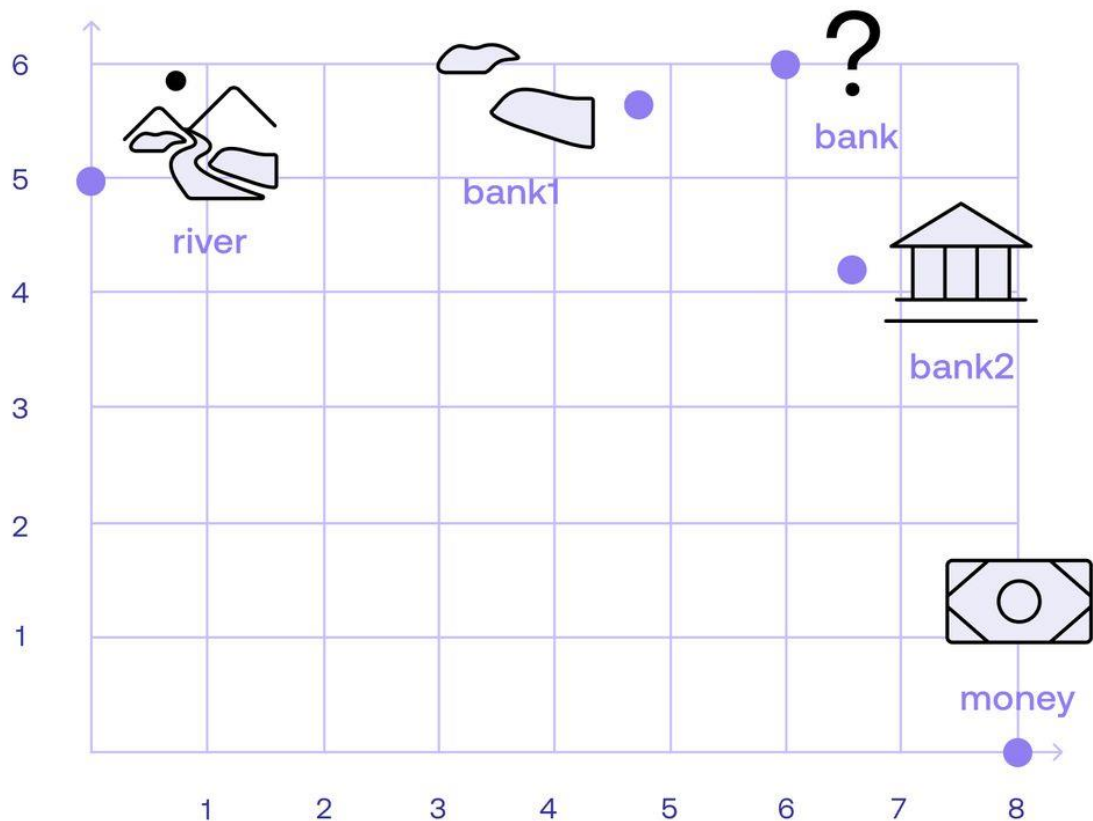
Telling context in words



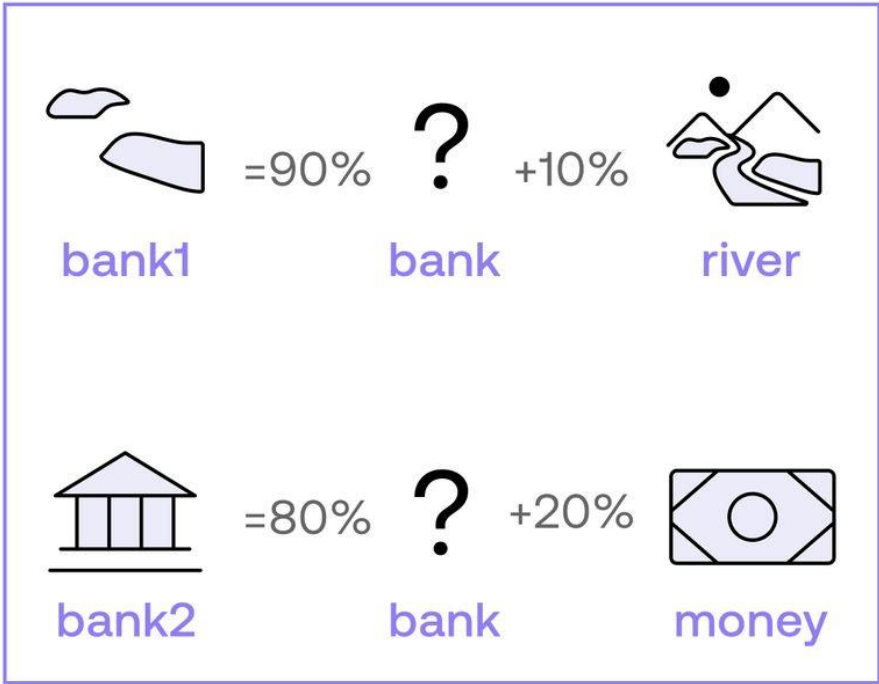
Modified sentence 1: The **bank1** of the river.

Modified sentence 2: Money in the **bank2**.

Embedding

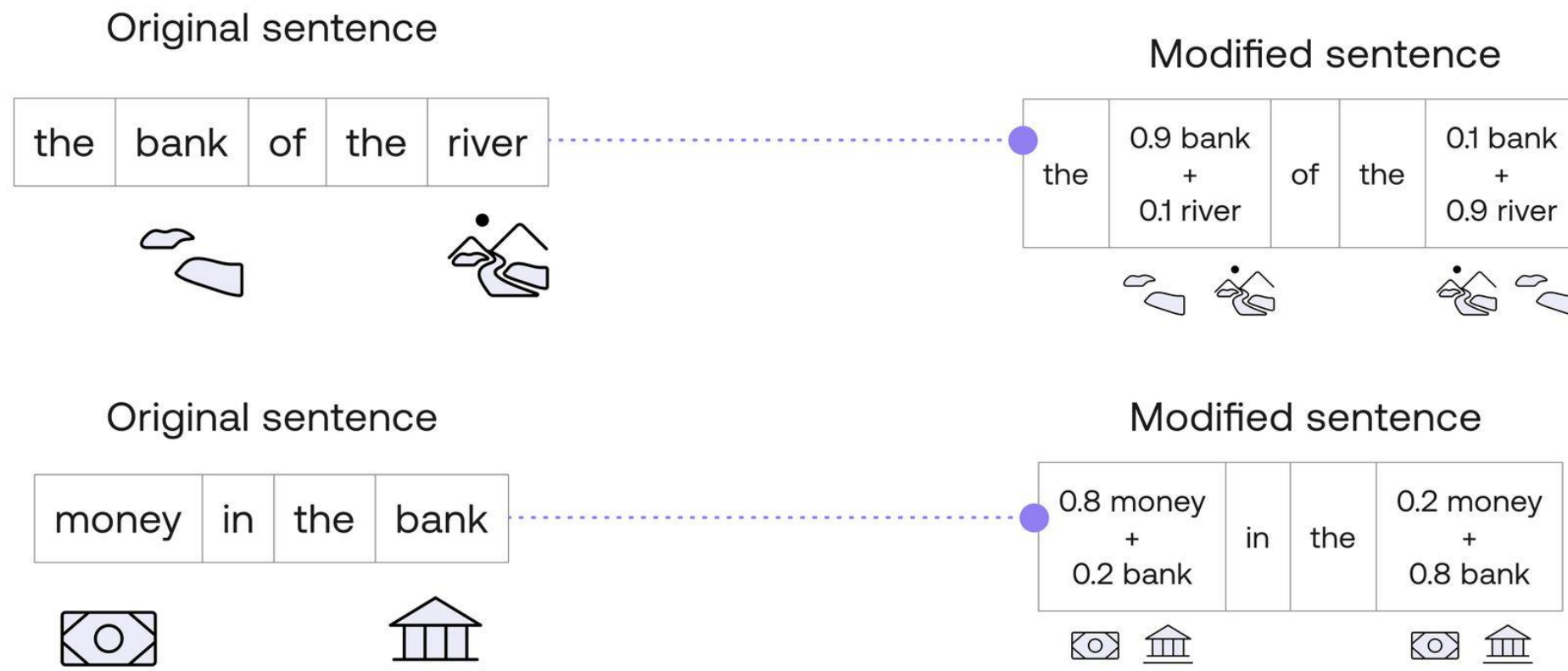


Equations



Bank1 = 0.9\*Bank + 0.1\*River  
Bank2 = 0.8\*Bank + 0.2\*Money

We find these "weights"  
through mechanisms such  
as similarity and attention





# Similarity matrices

The **bank** of the river

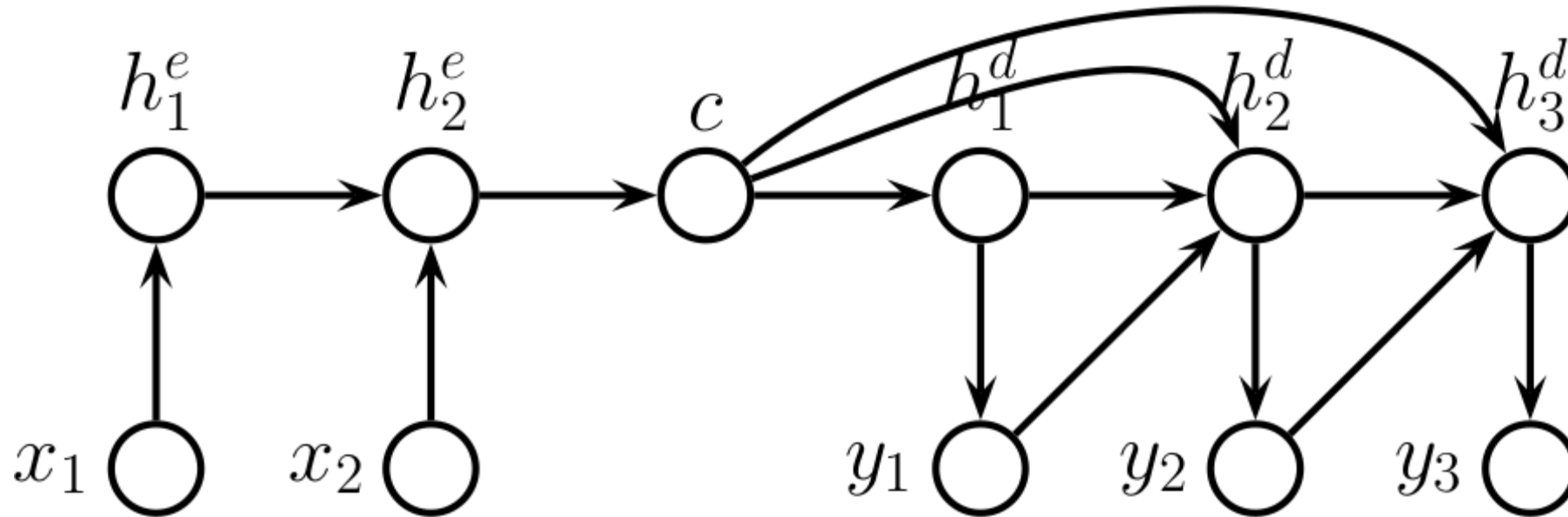
	the	bank	of	river
the	1	0	0	0
bank	0	1	0	0.11
of	0	0	1	0
river	0	0.11	0	1

Money in the **bank**

	money	in	the	bank
money	1	0	0	0.25
in	0	1	0	0
the	0	0	1	0
bank	0.25	0	0	1

# Seq2Seq model

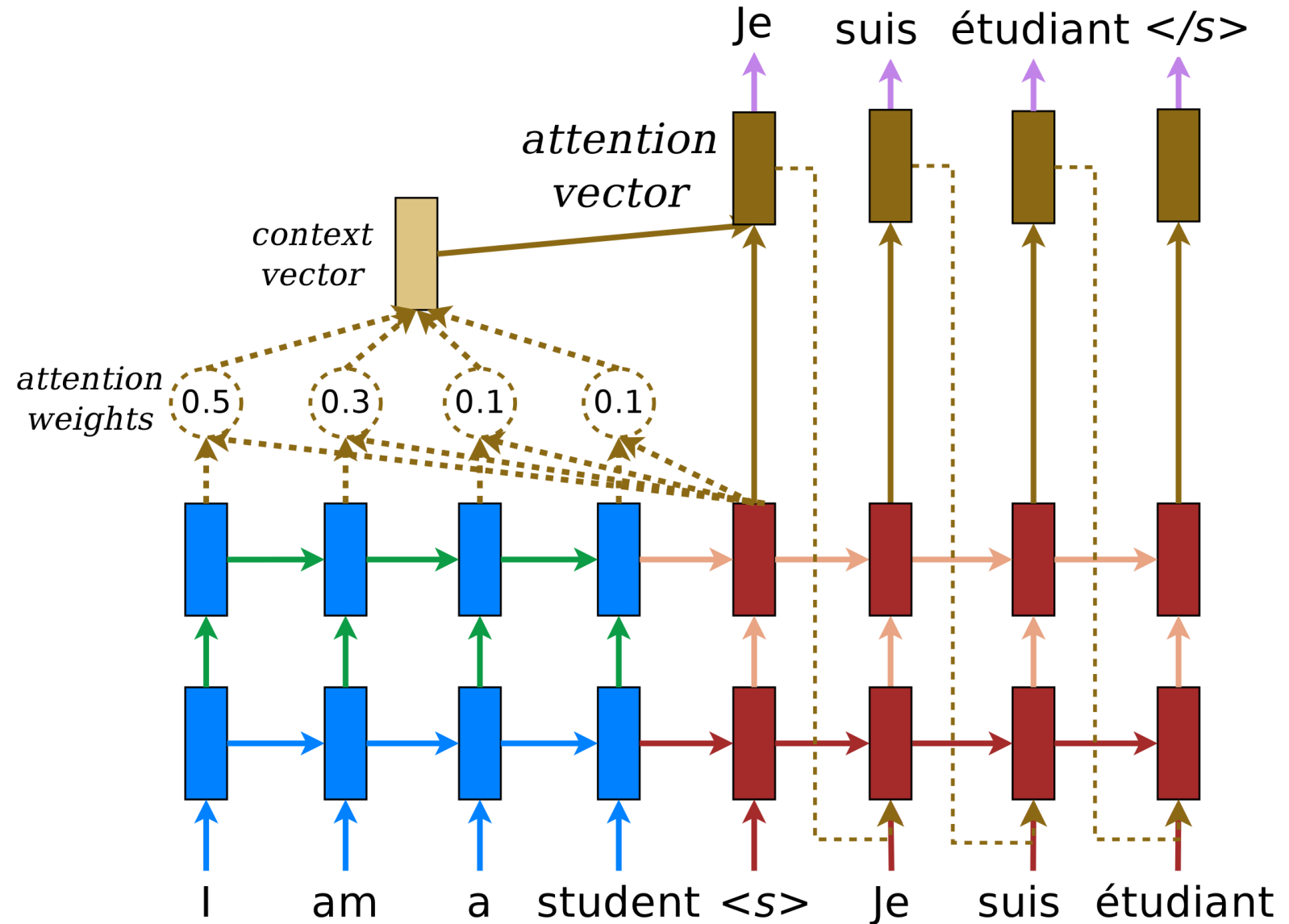
$$E \rightarrow C \rightarrow D$$



Context vector  $c$  encodes the *whole* input sequence  $x$   
Output has no access to input tokens  
If we had access, which token should we "pay attention to"?

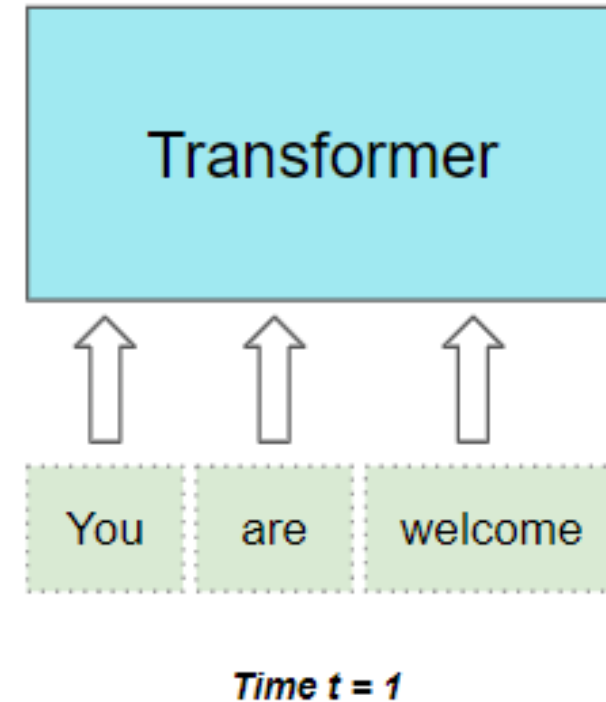
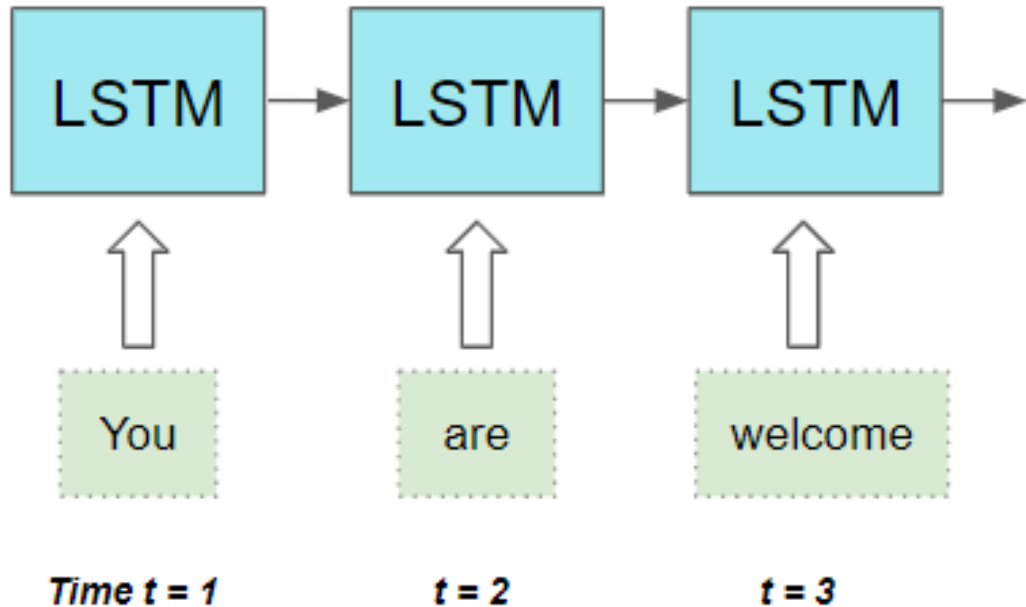
# Seq2Seq with attention

$$q = h_{t-1}^d$$
$$k_i = h_i^e$$
$$v_i = h_i^e$$



# Are transformers better than RNNs?

- Transformers deal better with long-range dependencies
- Transformers process the input at same time (not step-by-step)





# Transformer

# Transformers

Transformer = Seq2Seq model with attention at both E and D

Do we really need both *encoder* and decoder?

"Any architecture designed to process a connected set of units—such as the tokens in a sequence or the pixels in an image—where the only interaction between units is through self-attention."

Don't we need also cross-attention?

# Self attention (at encoder)

Input attends to itself!

$$y_i = \text{Attn}(x_i, \{x_j, x_j\})$$

$$q_i = x_i$$

$$k_j = x_j$$

$$v_j = x_j$$

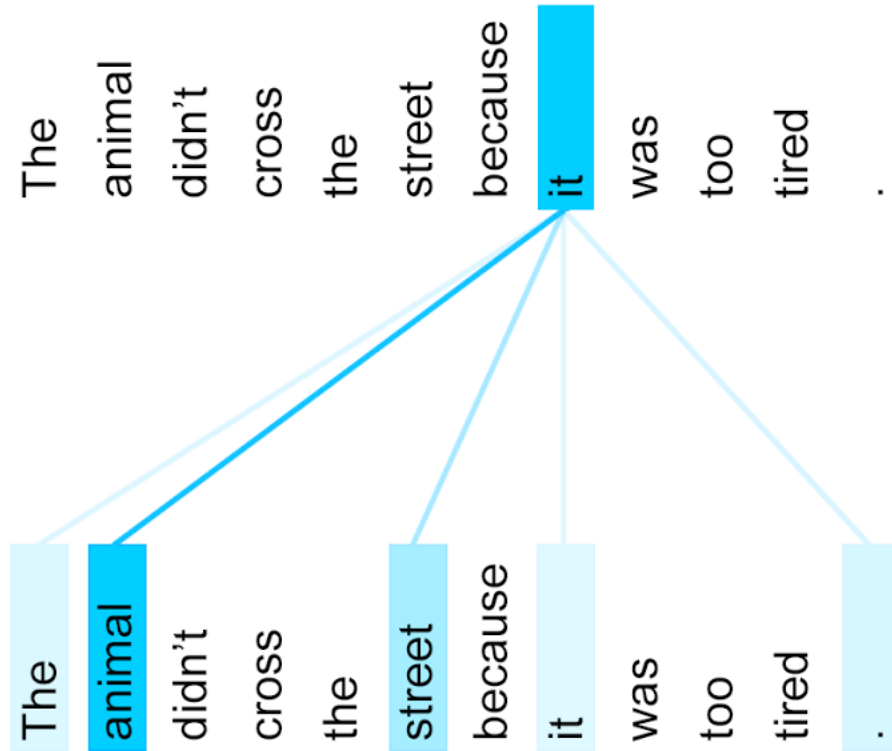


# At decoder: masked self-attention

$$y_i = \text{Attn} \left( \underset{q_i}{y_{i-1}}, \left\{ \underset{k_i}{(y_1, y_1)} \cdots \underset{v_i}{(y_{i-1}, y_{i-1})} \right\} \right)$$

# Example: Translation English-French

**x** = The animal didn't cross the street because it was too tired



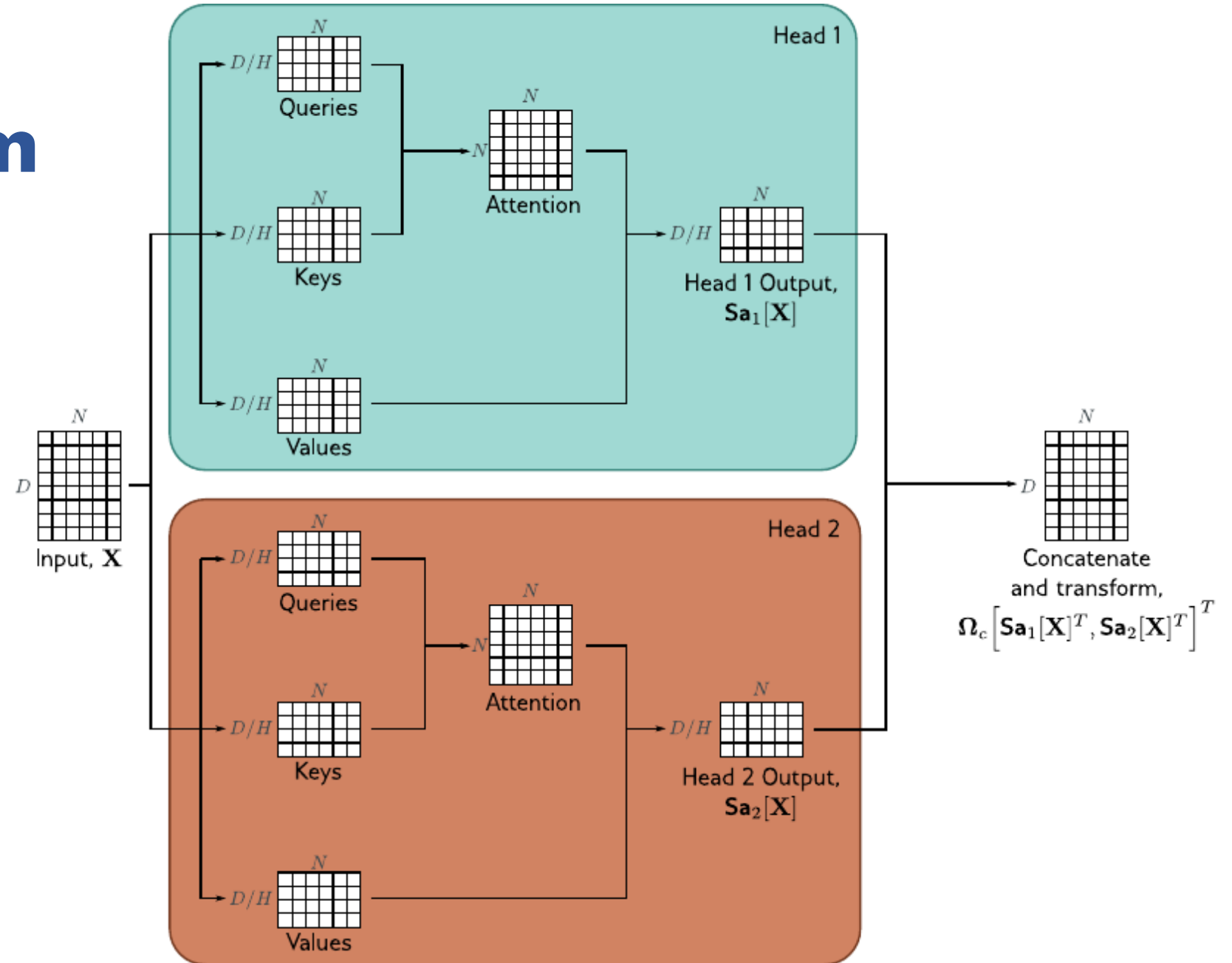
# Multi-head attention

Enrich the representation: having several notions of similarity

$$h_i = \text{Attn} \left( W_i^{(q)} q, \{W_i^{(k)} k_j, W_i^{(v)} v_j\} \right)$$

$$h = \text{MHA} (q, \{(k_j, v_j)\}) = W_o \begin{bmatrix} h_1 \\ h_2 \\ \vdots \\ h_n \end{bmatrix}$$

# Multihead – block diagram



# Positional encoding

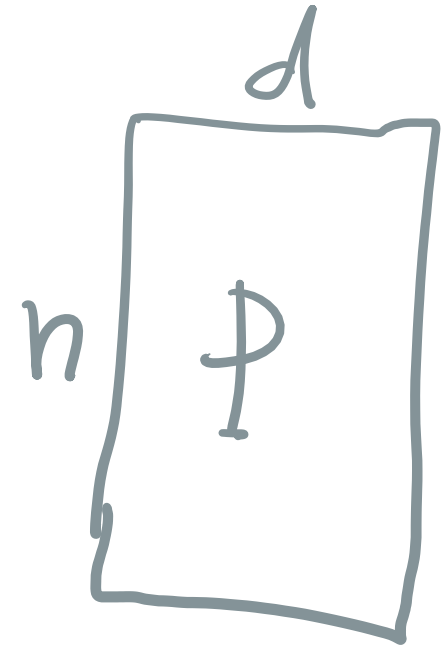
Sequences have order, but self-attention is permutation invariant!

$$x_1, x_2, \dots, x_n \quad x_i \in \mathbb{R}^d$$

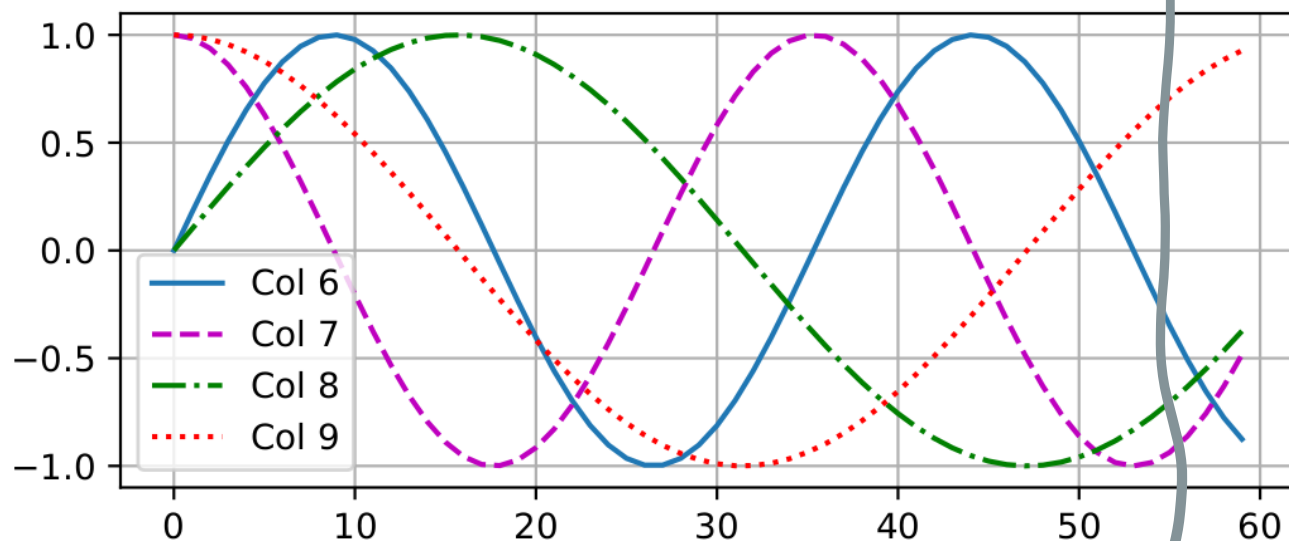
$$P_{i,j} = \sin\left(\frac{i}{C^{2j/d}}\right) \quad \text{if } j \text{ is even}$$

$$P_{i,j} = \cos\left(\frac{i}{C^{2j/d}}\right) \quad \text{if } j \text{ is odd}$$

$C = \text{max seq length}$



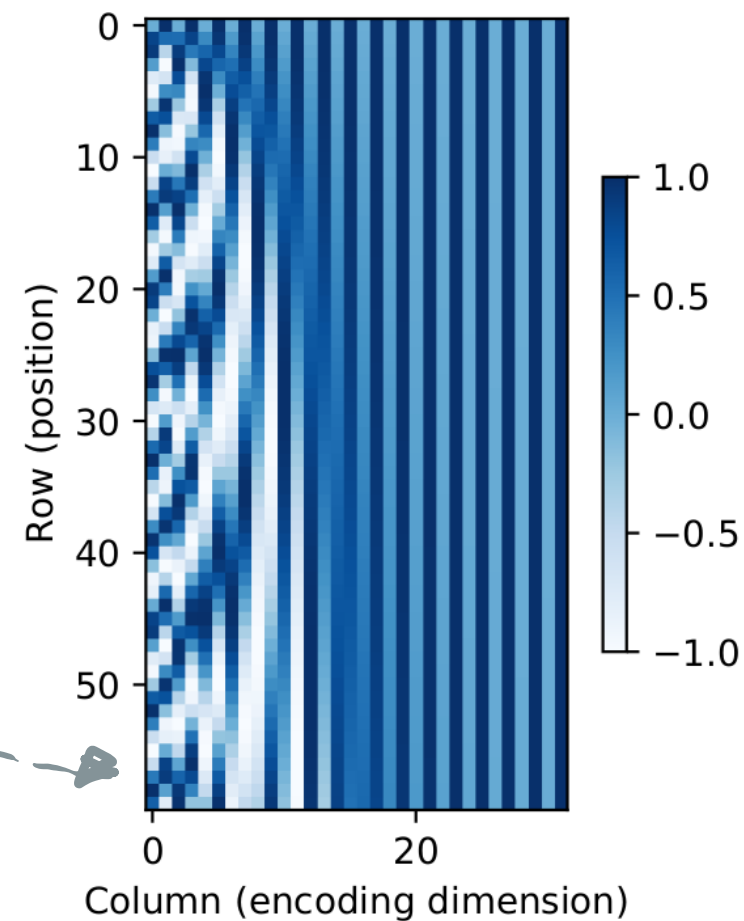
$n=60, d=32$



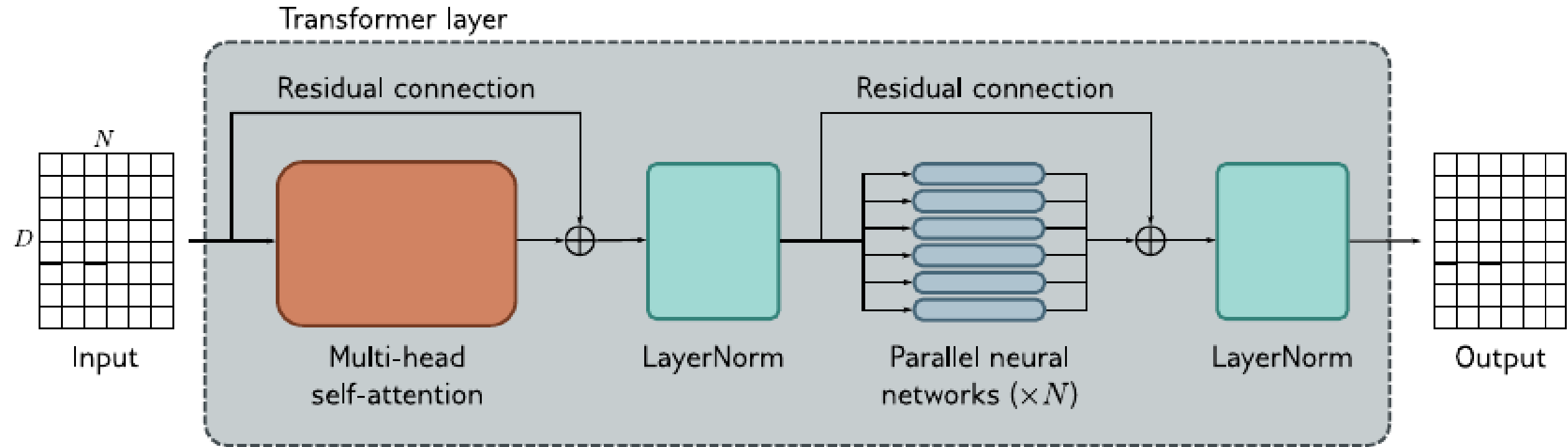
$X$  = original embeddings

$P$  = positional embeddings

new embeddings:  $X + P$



# Transformer: MH SA + something else



# Encoder [+ decoder] transformers

## Self-attention:

- input tokens attend to one another

## Masked self-attention

- output tokens attend to *previous* output tokens

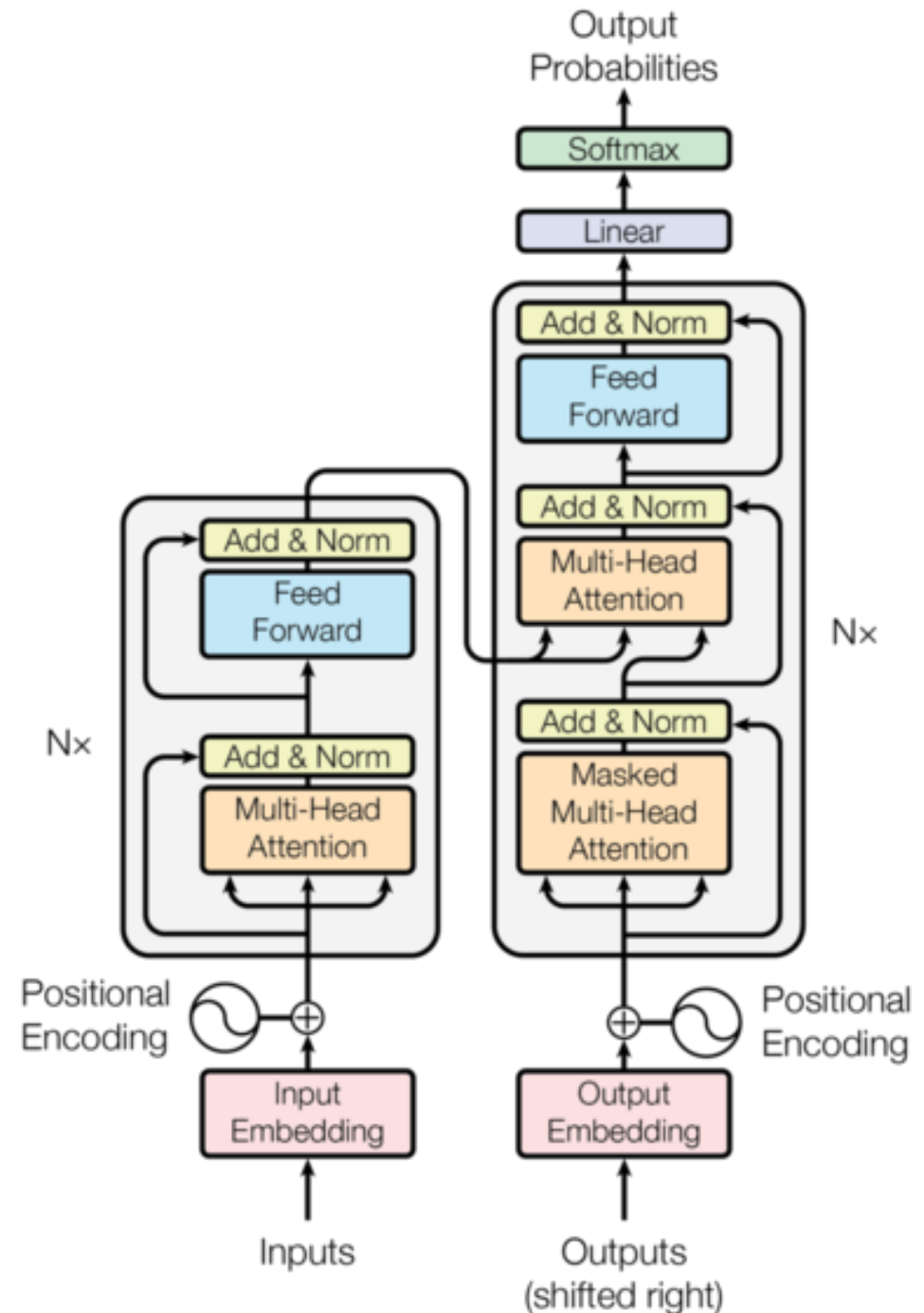
## Cross-attention

- output tokens attend to input tokens



Let's identify the (q,k,v) at the three different attention blocks

$k_i$   $v_i$   $q_i$



# Pseudocode

Encoder = series of  $M$  encoder blocks

```
def EncoderBlock(X) :  
    Z = LayerNorm(MHA(Q=X, K=X, V=X) + X) # note the residual connection  
    E = LayerNorm(FeedForward(Z) + Z)  
    return E  
  
def Encoder(X, M) : # M = number of layers  
    E = POS(Embed(X))  
    for _ in range(M) :  
        E = EncoderBlock(E)  
    return E
```

Decoder has access to:

- encoder (via another MHA)
- Previously generated output

```
def DecoderBlock(X, E):  
    Z = LayerNorm(MHA(Q=X, K=X, V=X) + X)  
    Z' = LayerNorm(MHA(Q=Z, K=E, V=E) + Z)  
    D = LayerNorm(FeedForward(Z') + Z')  
    return D
```

```
def Decoder(X, E, N):  
    D = POS(Embed(X))  
    for _ in range(N):  
        E = DecoderBlock(D, E)  
    return D
```

# Attention is not explanation

Or is it?

# Many variations proposed

Conformer (conv layers inside Transformer)

Reformer

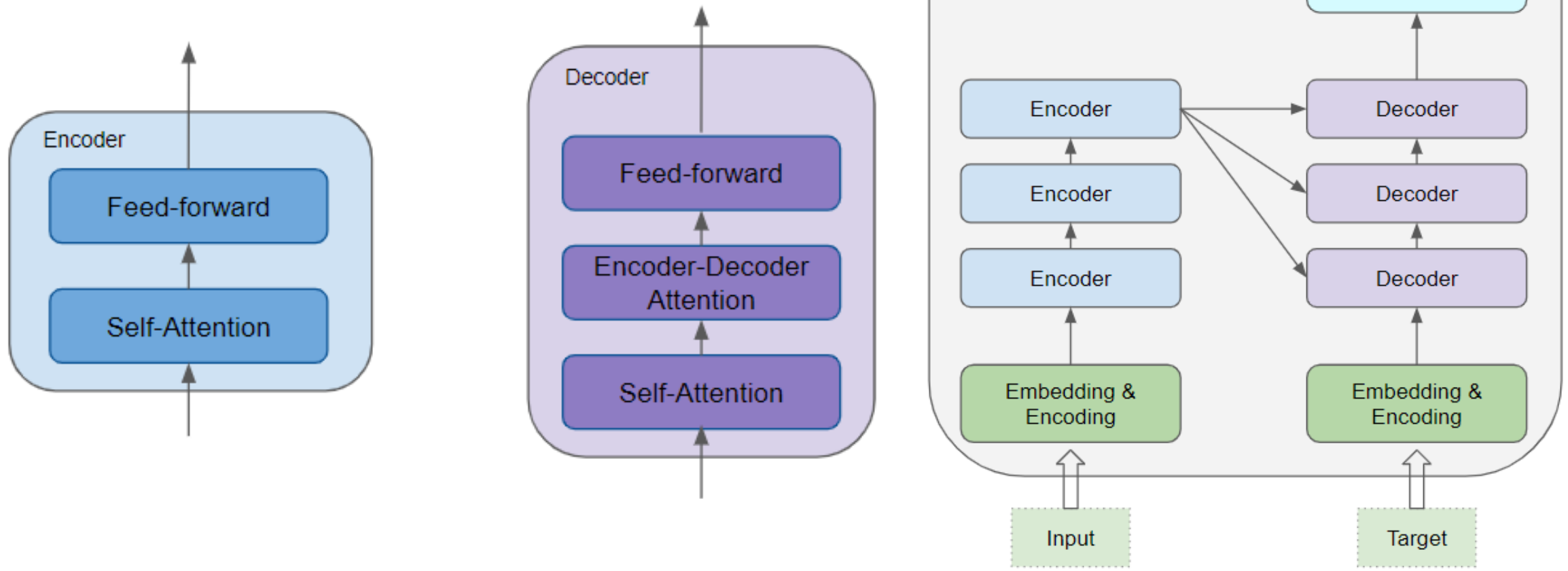
Linformer

Performer

...

# **A top-down explanation**

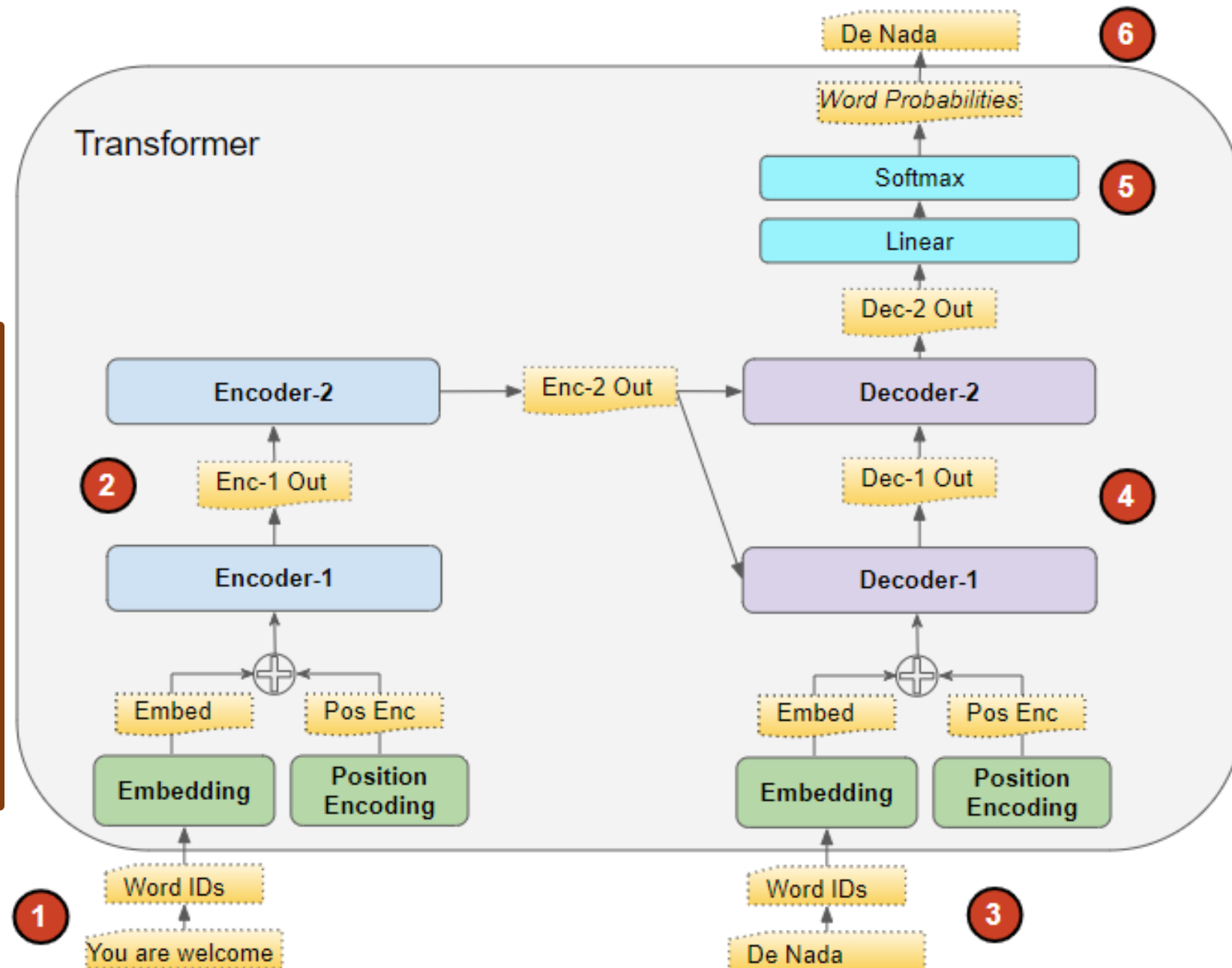
# Encoder, Decoder and how they relate



# At training...

Loss:  
compares output  
sequence with  
target sequence  
(known because it's  
in the training data)

Teacher Forcing

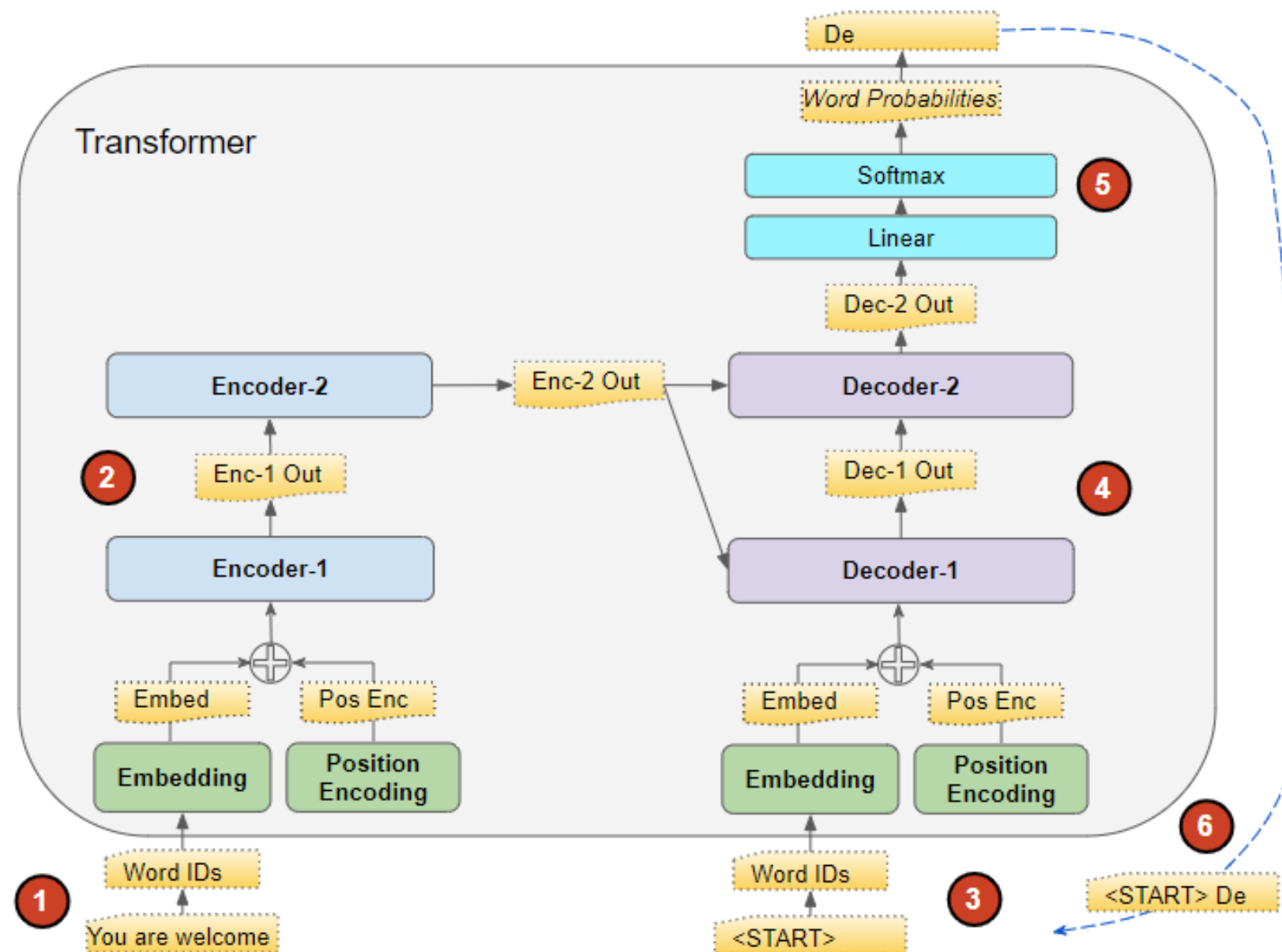




# At inference...

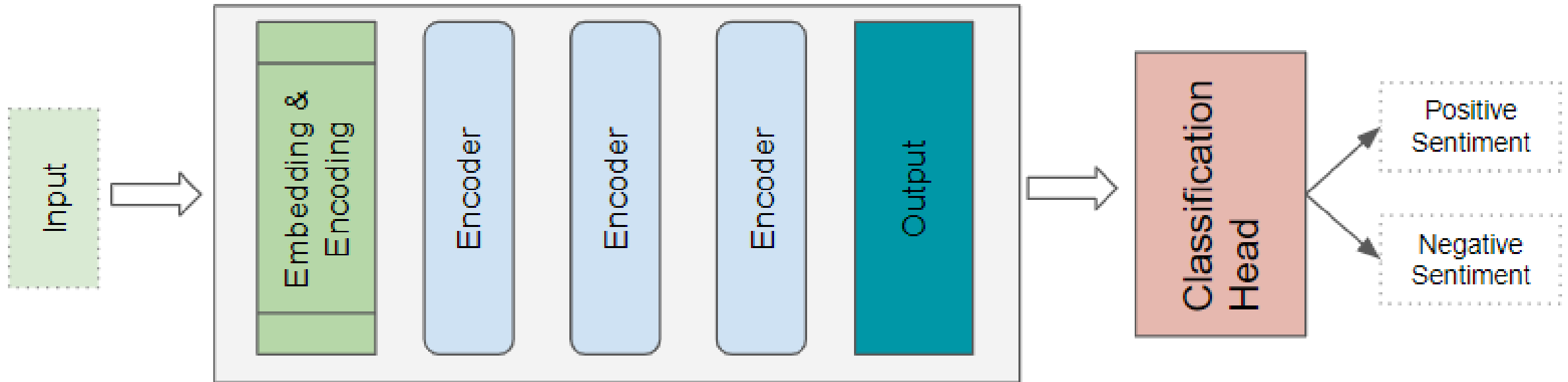
Last word of output sequence added at the end of the decoder input

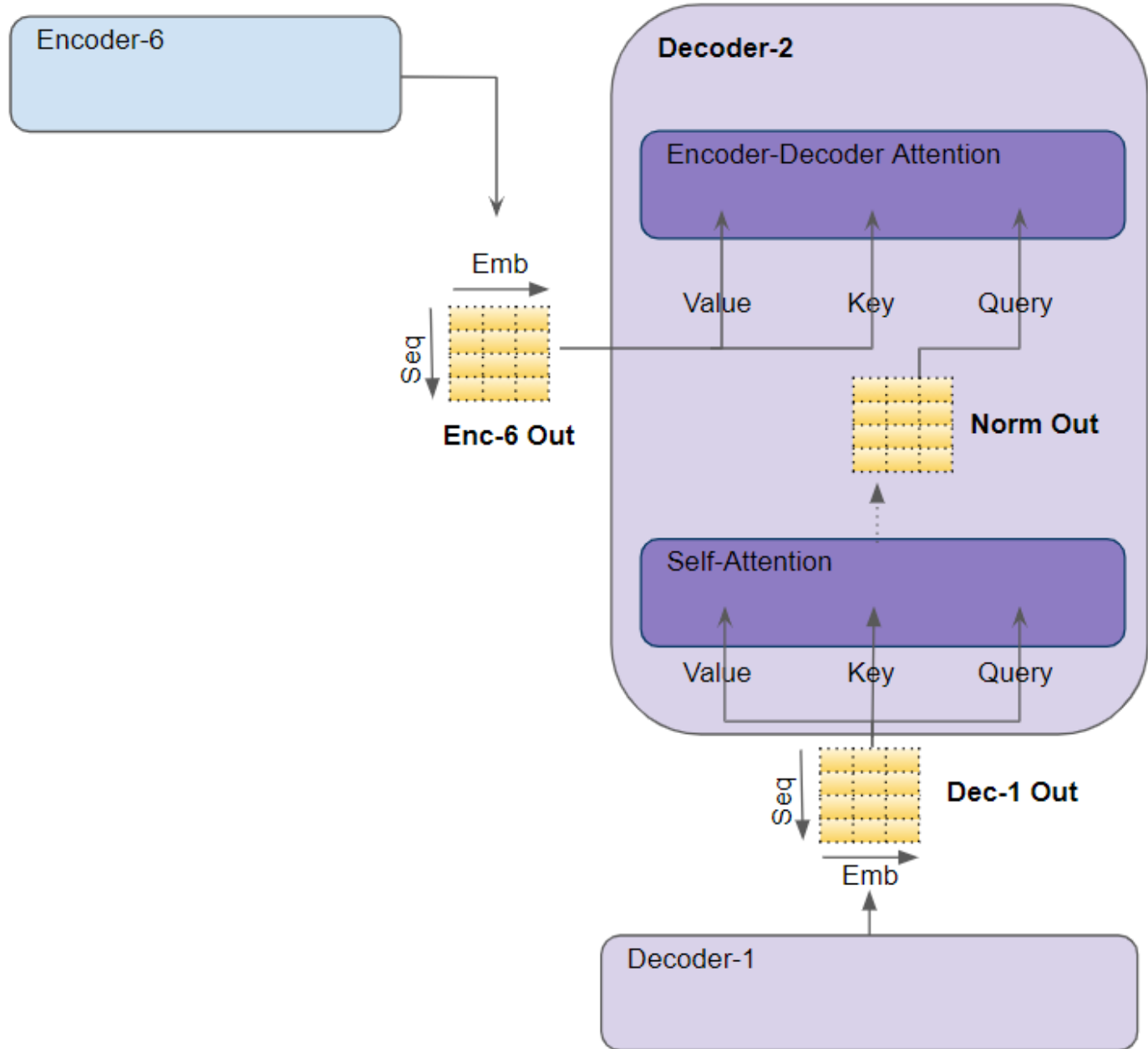
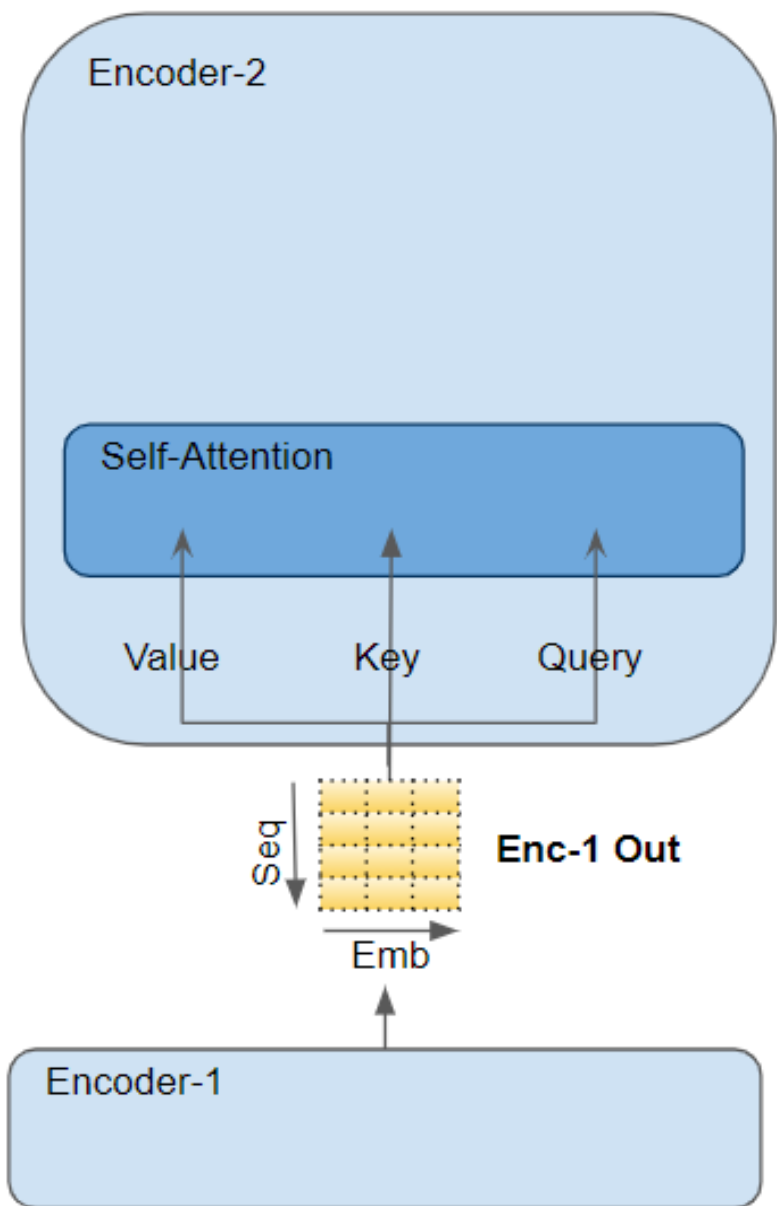
Repeat 3, 4, 5, 6 until "end of sequence"

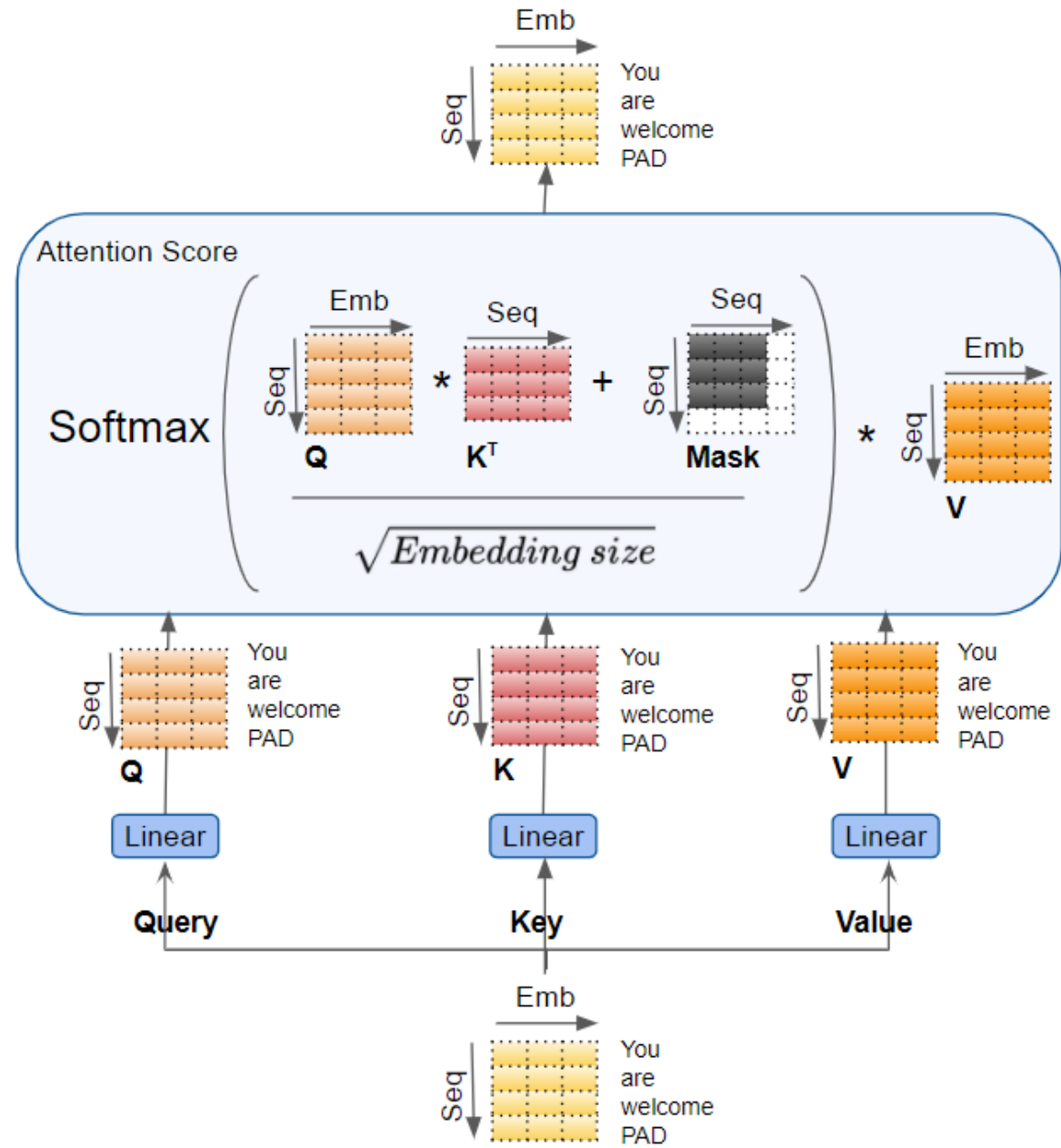


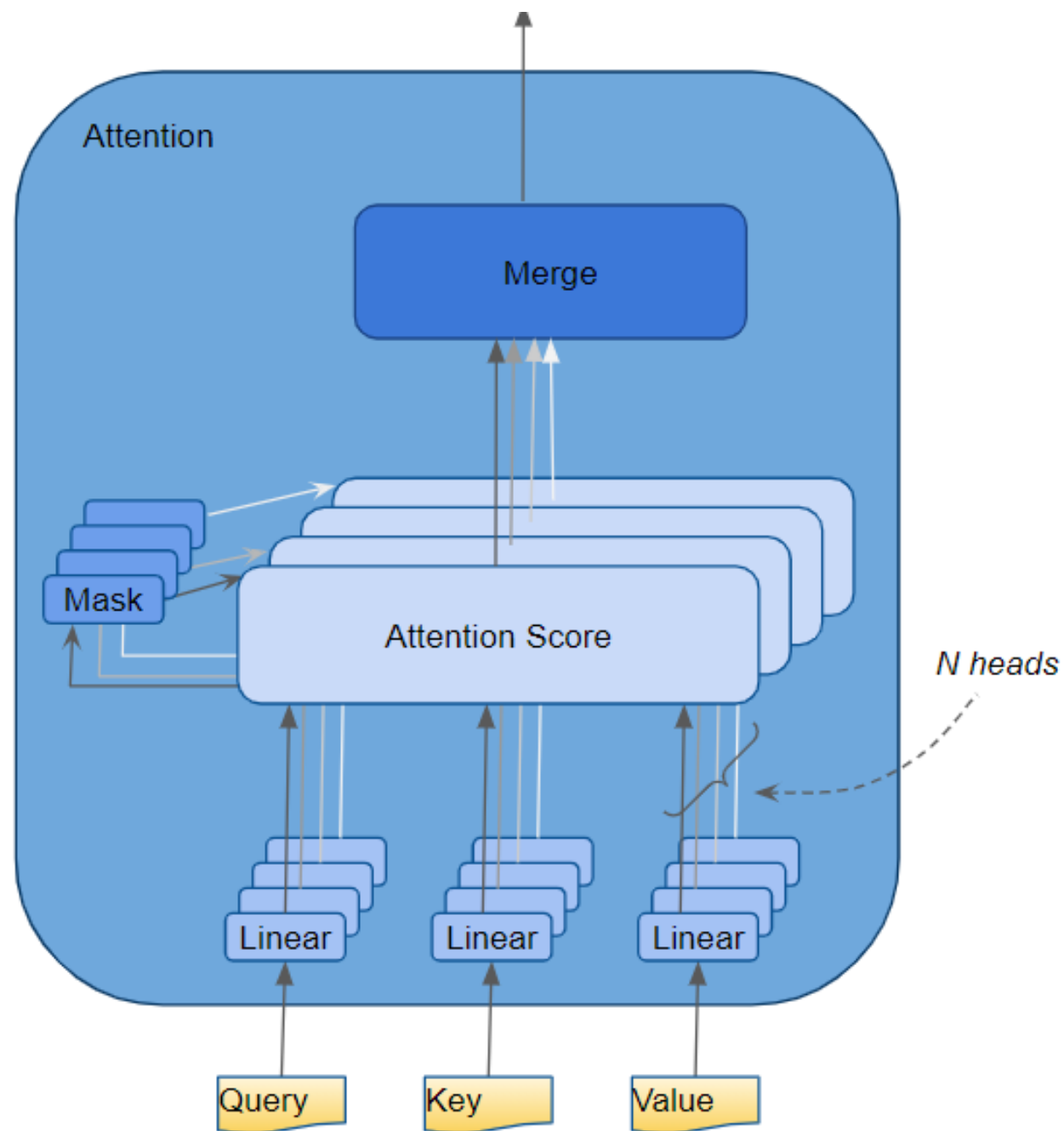
# We don't always have a decoder...

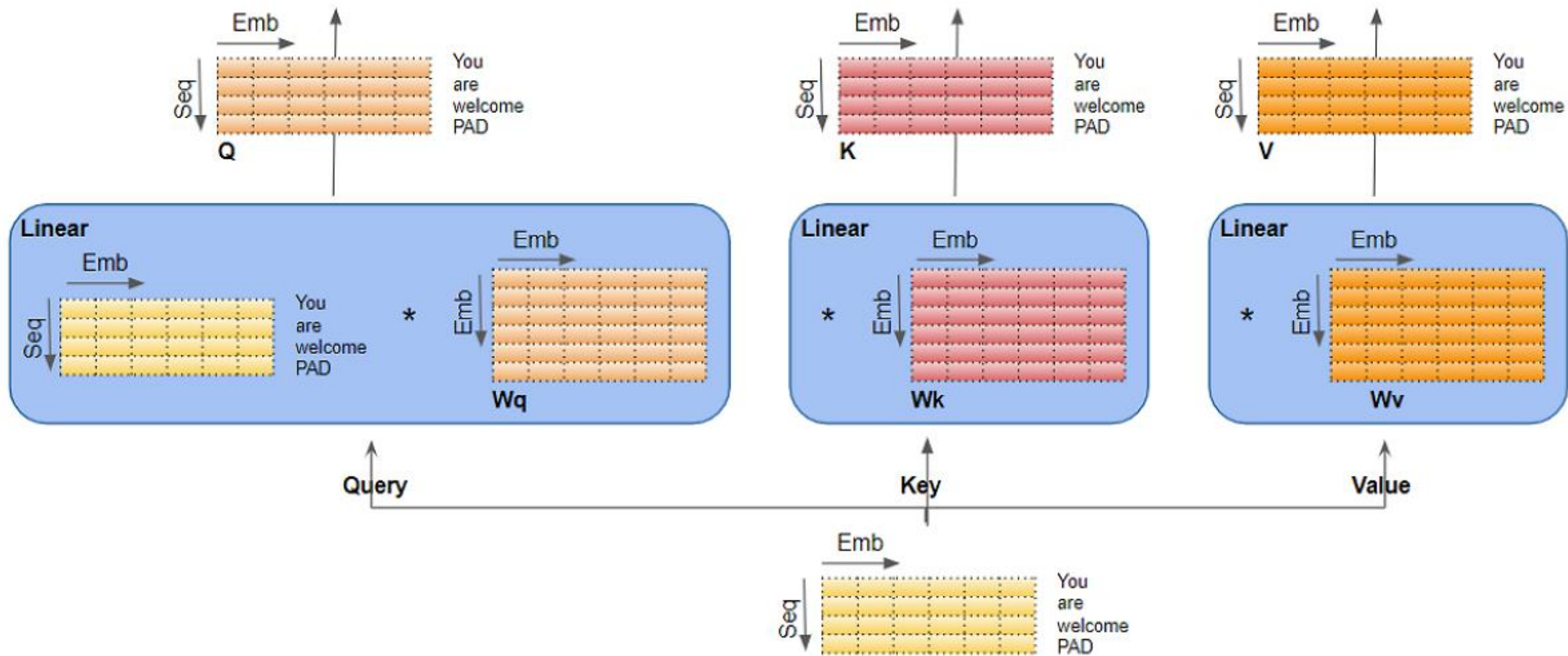
Encoder as a building block for different applications



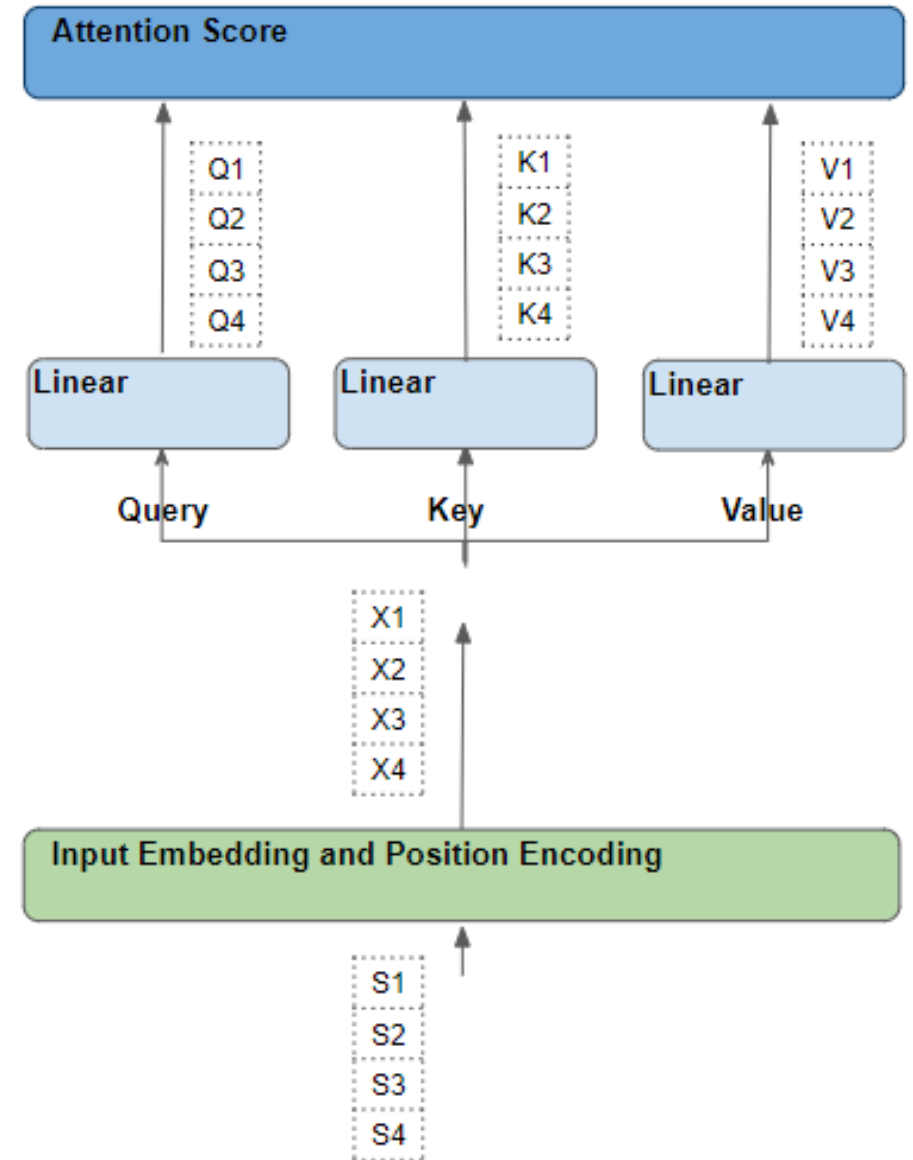




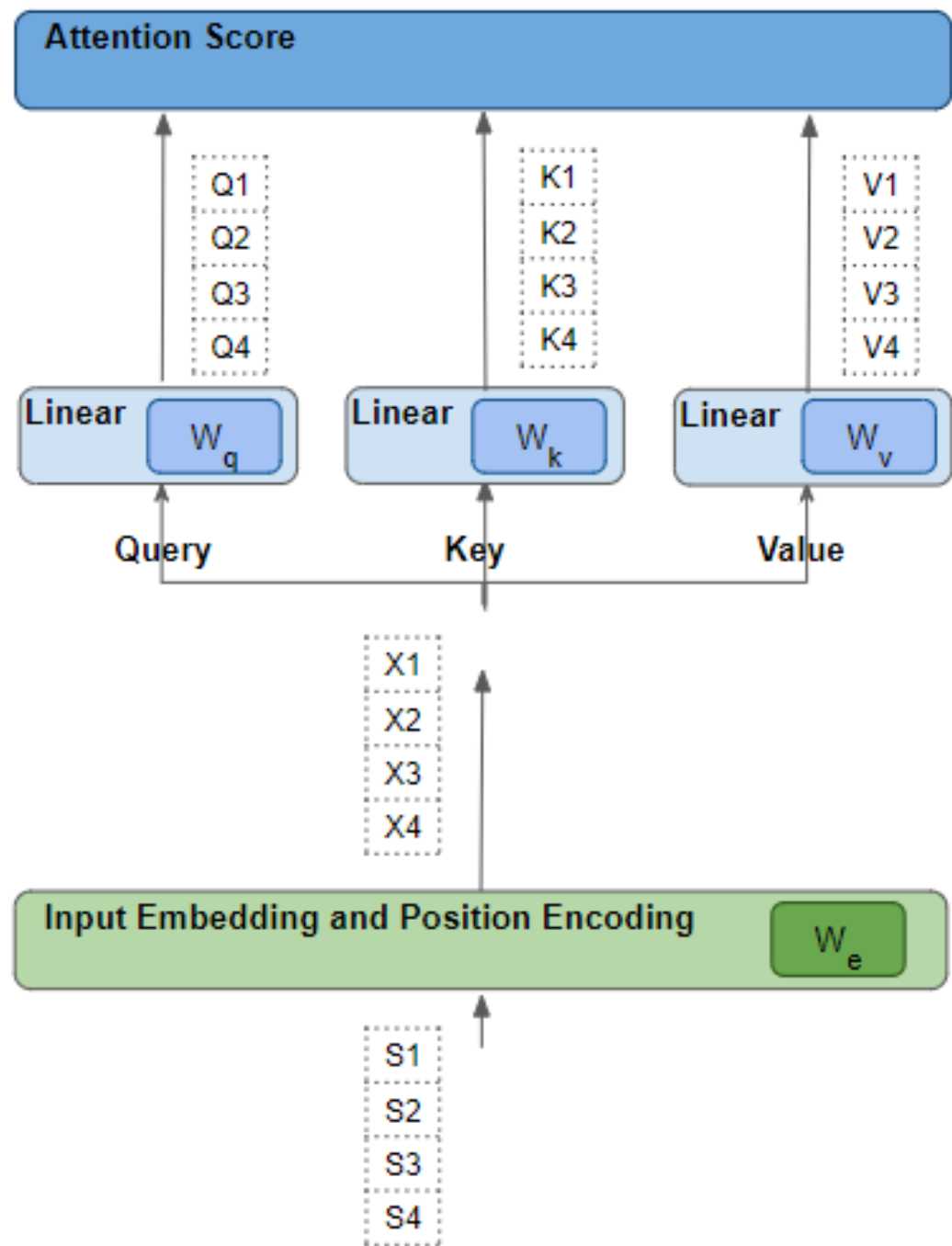




**Words (sequence items)  
go through a series  
of transformations**



# Transformations are learnable!





# Role of Q, K, V

The **Query** word: word **for which** we are calculating Attention

The **Key** and **Value** word: word ***to which*** we are paying attention  
(how relevant is that word to the Query word)

# And how is "relevance" learned?

The more "aligned" (similar embeddings) two words are, the higher their attention score

So, the transformer learns to produce low or high scores

And how are such scores found?  
From the training instances

# **Vision Transformers**

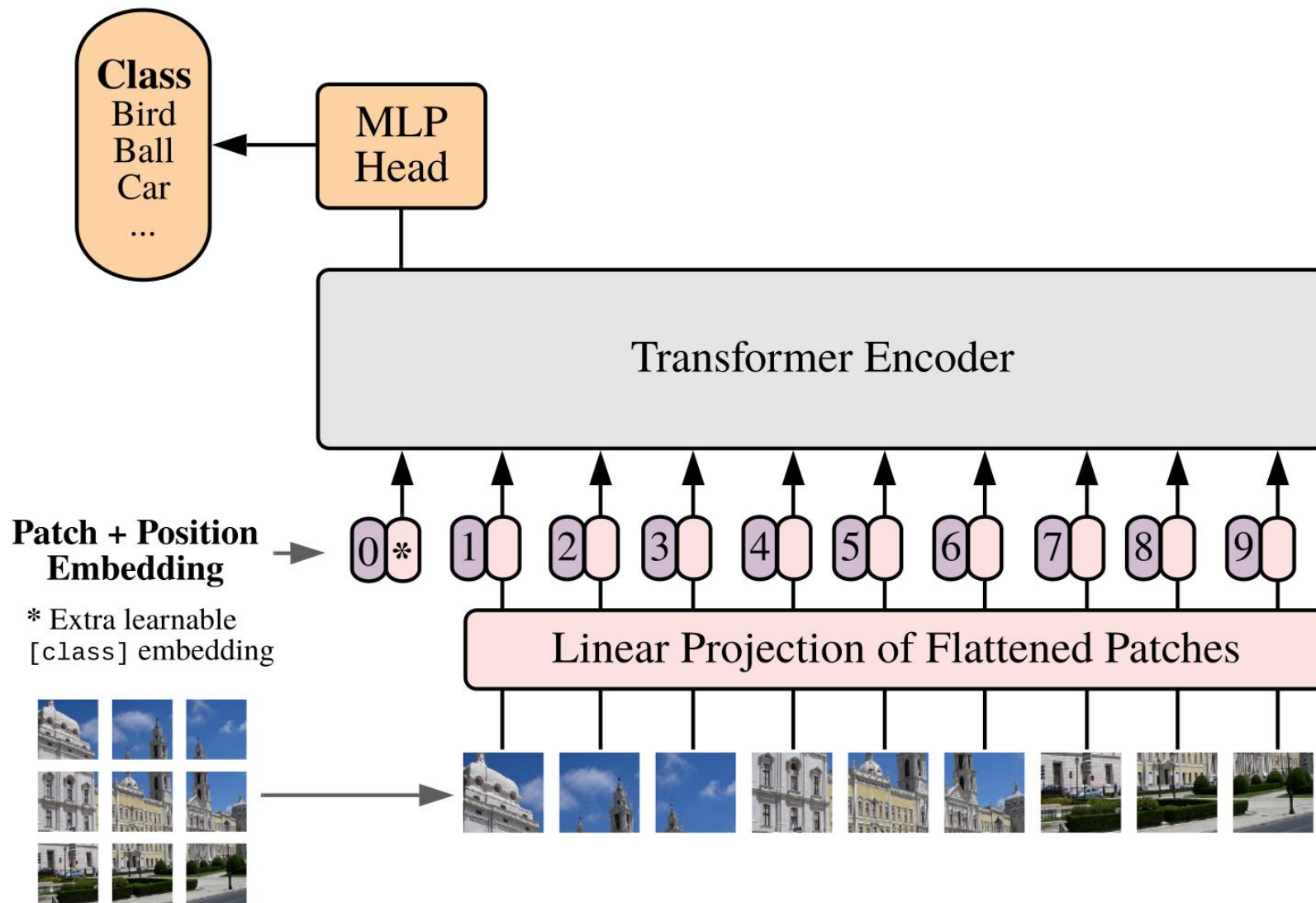
# CNNs vs Transformers for images

CNNs have built-in **inductive bias**:

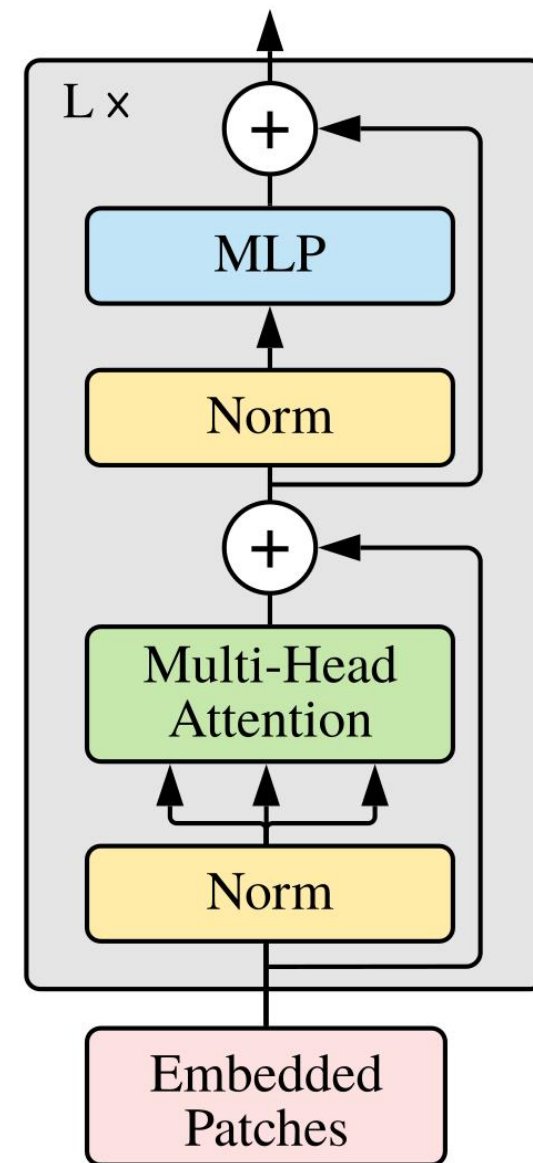
- **locality** (small kernel)
- **equivariance** (weight sharing)
- **location invariance** (pooling)

Transformers don't seem a good fit for 2D data...

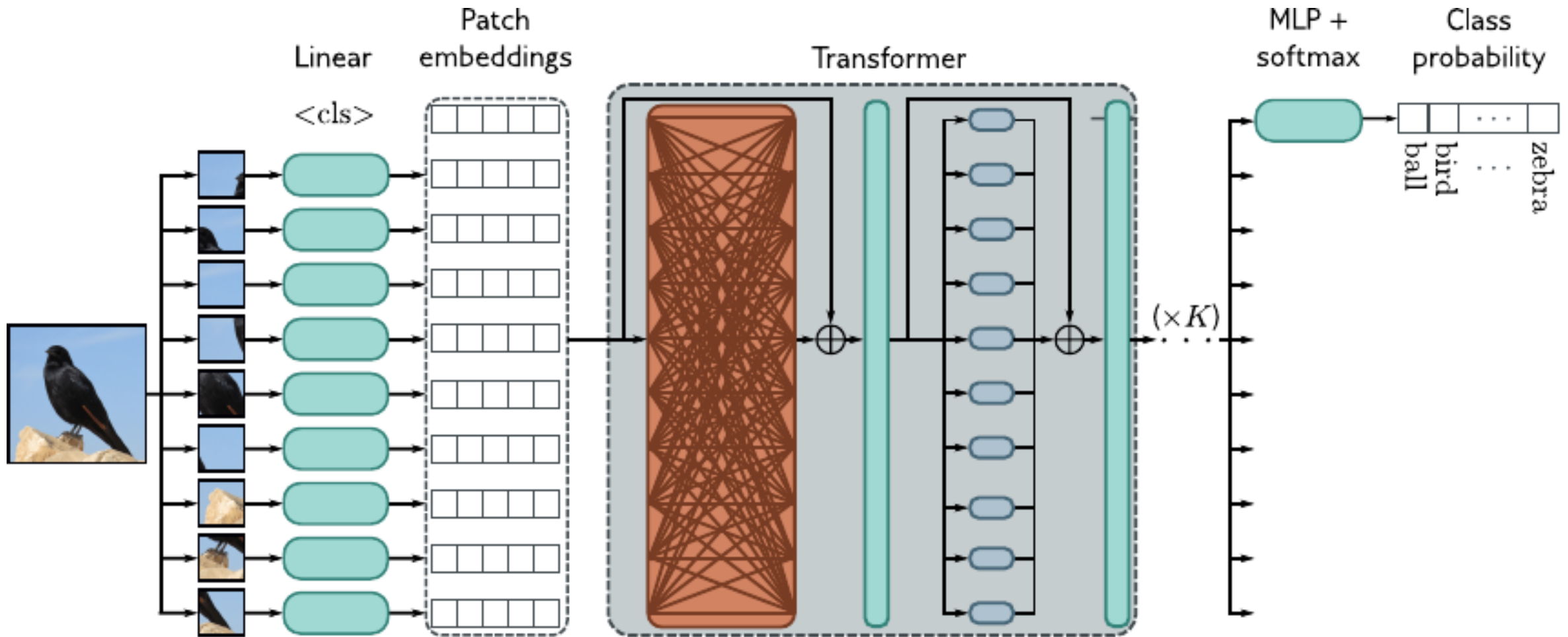
## Vision Transformer (ViT)



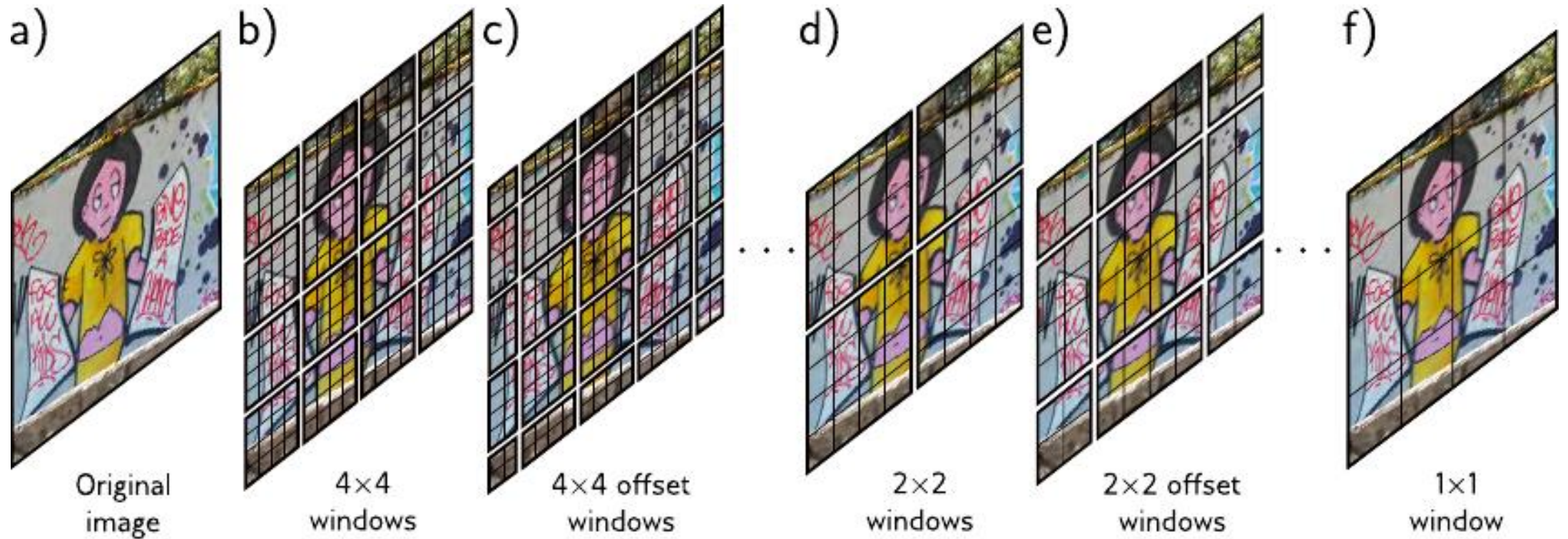
## Transformer Encoder



# ViT: transformer encoder



# SWin (shifted-window): multi-scale transformer



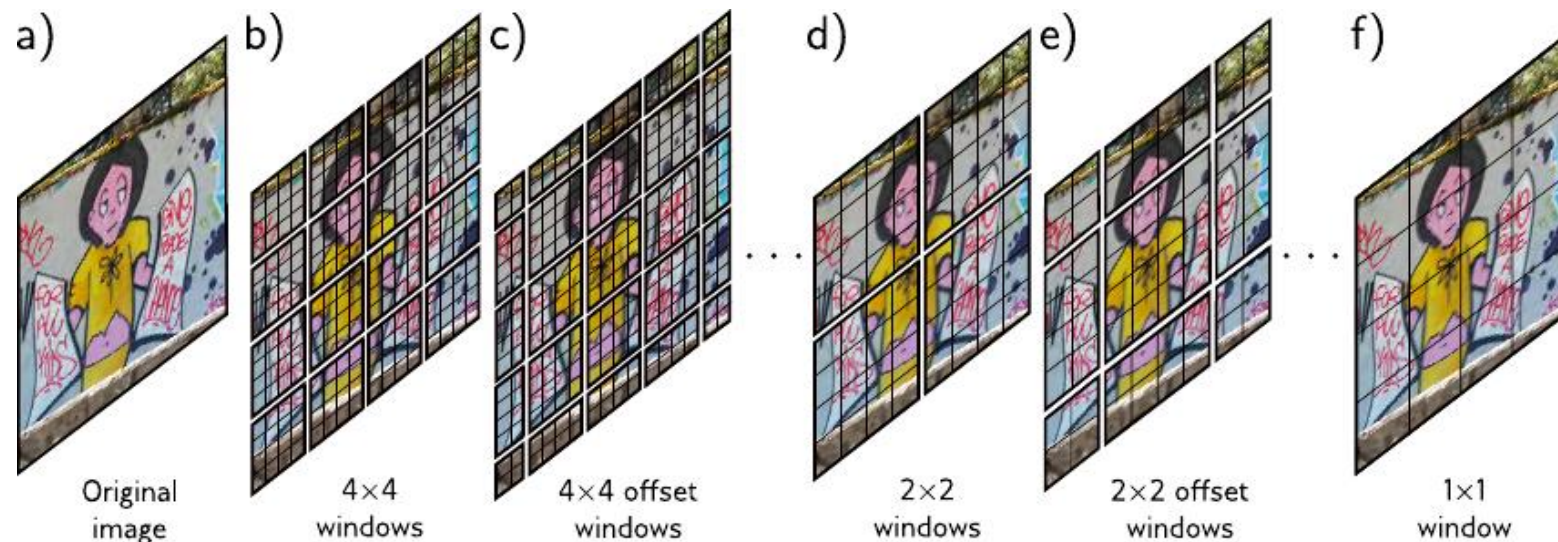


## Efficiency:

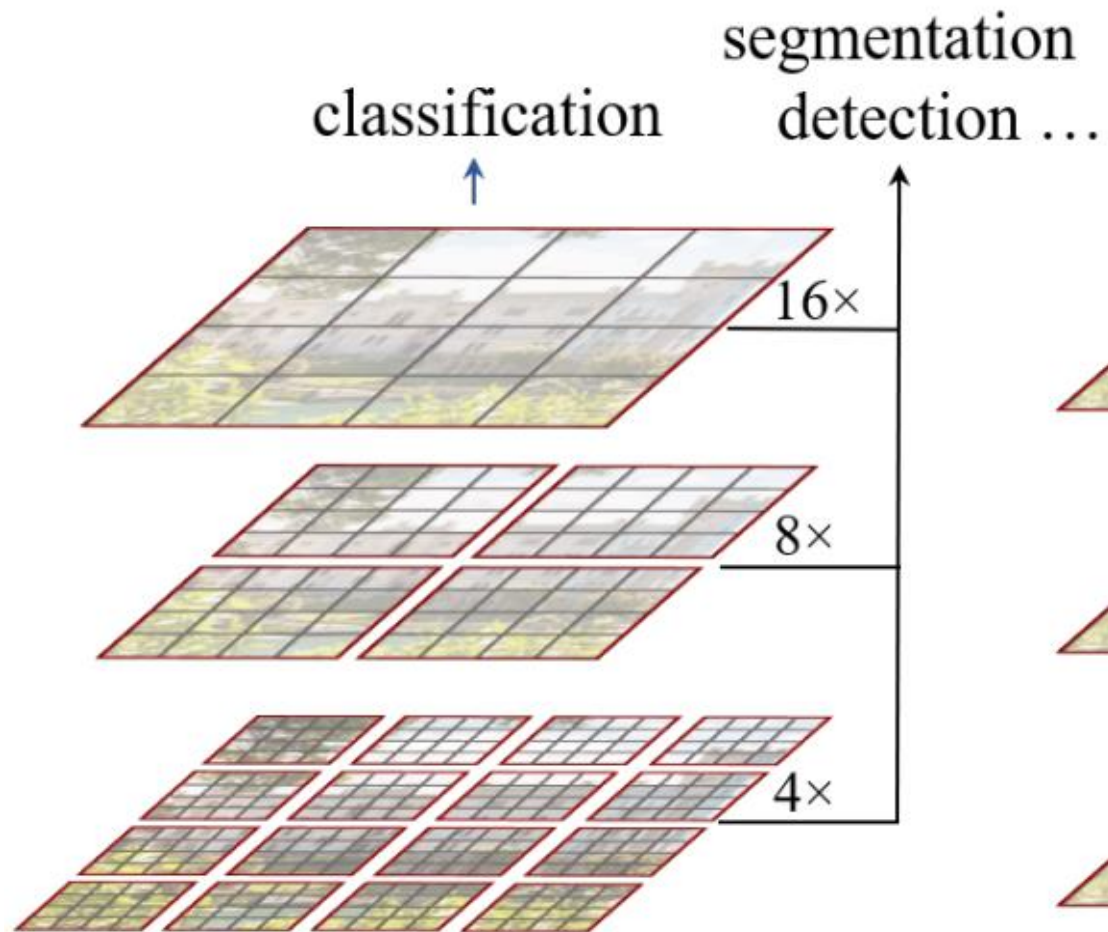
- Limiting self-attention to non-overlapping **local** windows

## Hierarchical architecture:

- Flexibility to model at various **scales**
- **Linear** computational complexity with respect to image size

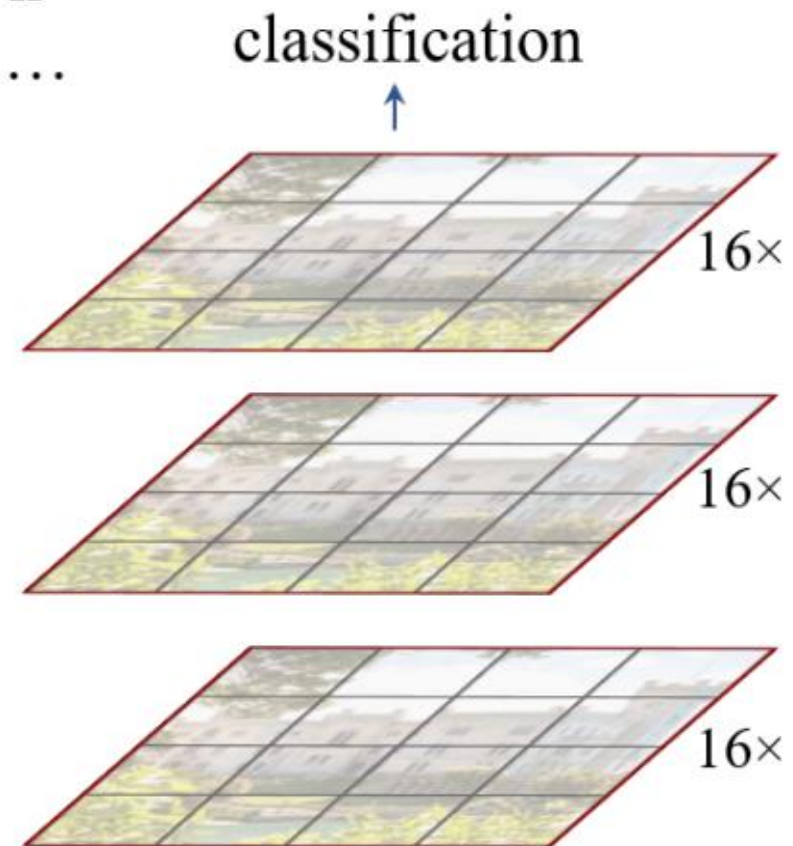






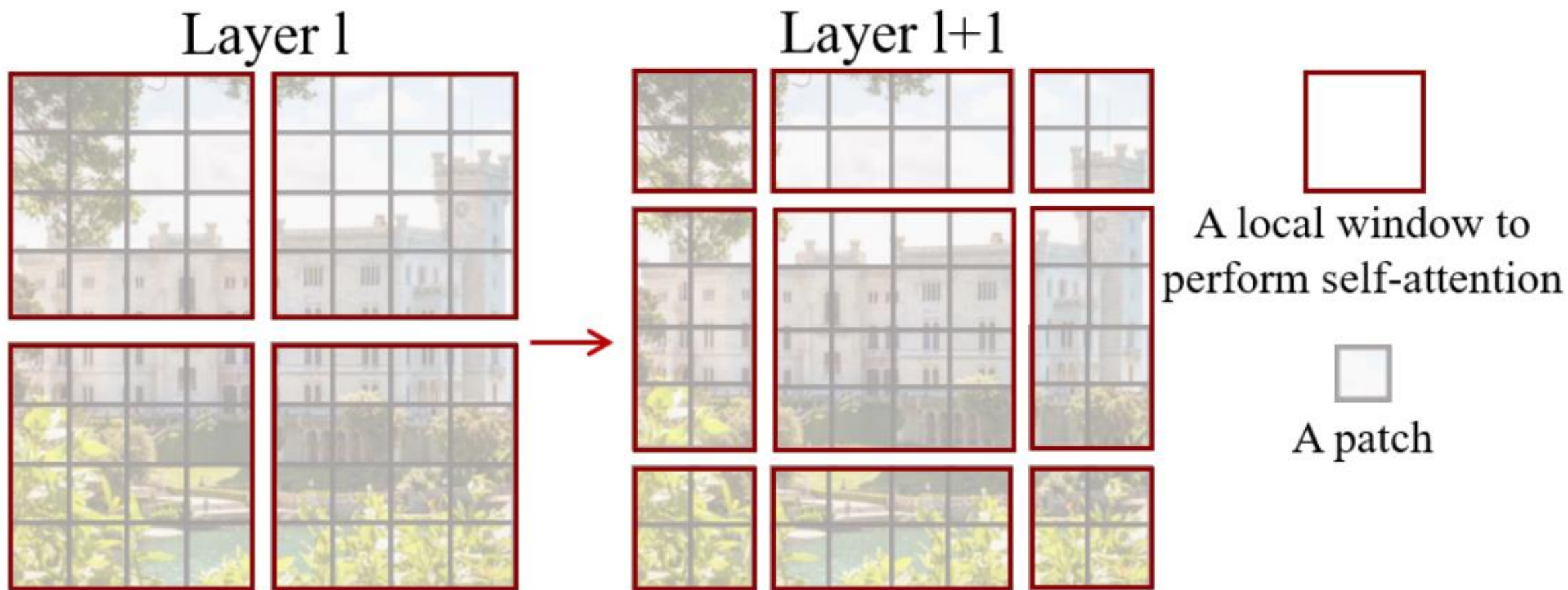
(a) Swin Transformer (ours)

Local attention  
Hierarchical  
 $O(N)$



(b) ViT

Global attention  
Single and coarse level  
 $O(N^2)$



Window partitioning is shifted in alternate layers  
--> this changes the interacting patches (and connect windows)

# Ablation study (SWin)

	ImageNet		COCO		ADE20k
	top-1	top-5	AP <sup>box</sup>	AP <sup>mask</sup>	mIoU
w/o shifting	80.2	95.1	47.7	41.5	43.3
shifted windows	<b>81.3</b>	<b>95.6</b>	<b>50.5</b>	<b>43.7</b>	<b>46.1</b>
no pos.	80.1	94.9	49.2	42.6	43.8
abs. pos.	80.5	95.2	49.0	42.4	43.2
abs.+rel. pos.	81.3	95.6	50.2	43.4	44.0
rel. pos. w/o app.	79.3	94.7	48.2	41.9	44.1
rel. pos.	<b>81.3</b>	<b>95.6</b>	<b>50.5</b>	<b>43.7</b>	<b>46.1</b>

# Findings

Transformers also work well (despite the lack of inductive bias!)

- If trained with **enough data** (to compensate this lack)

Expensive training: 30 days on ImageNet-21k on Google Cloud  
TPUx3 with 8 cores

# Video Transformers

Tokenization ► Embedding ► Positioning

