

CS846

Machine Learning for Software Engineering

Pengyu Nie

Sequence-to-Sequence Models and Transformers

Tokenizer

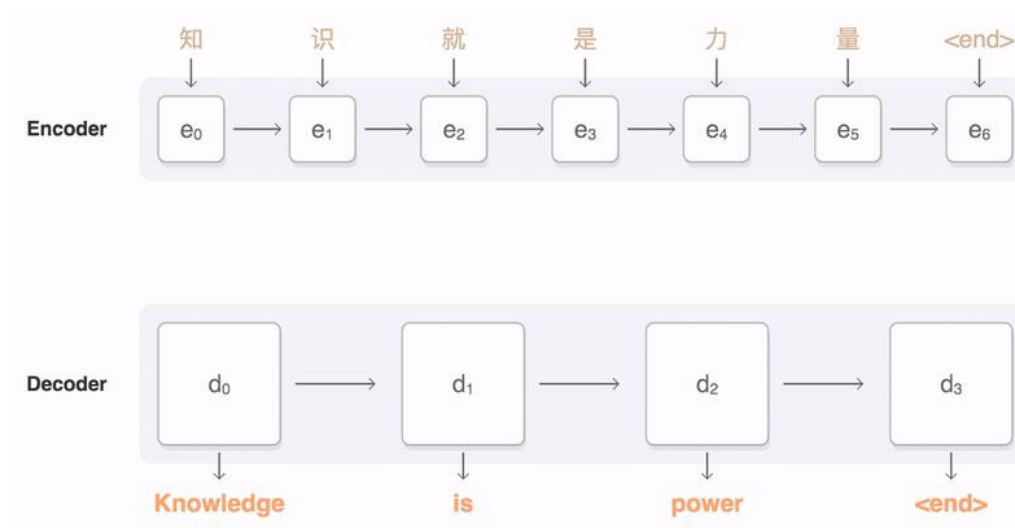
Embedding

Attention

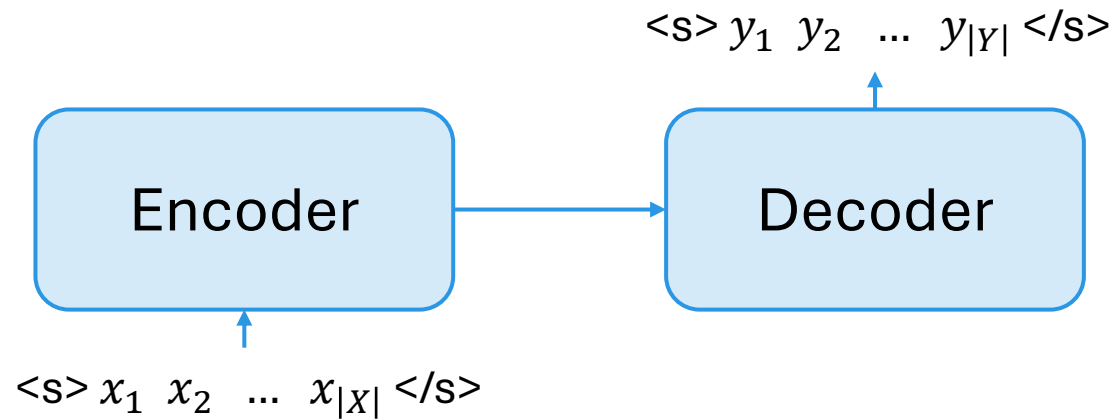
Motivation

- Language model: $P(W)$ or $P(w_i|w_1, \dots, w_{i-1})$
- Tasks with inputs and outputs
 - Neural machine translation (NMT)
 - Summarization
 - Examples in Software Engineering?
- $P(Y|X)$ or $P(y_i|y_1, \dots, y_{i-1}, x_1, \dots, x_{|X|})$

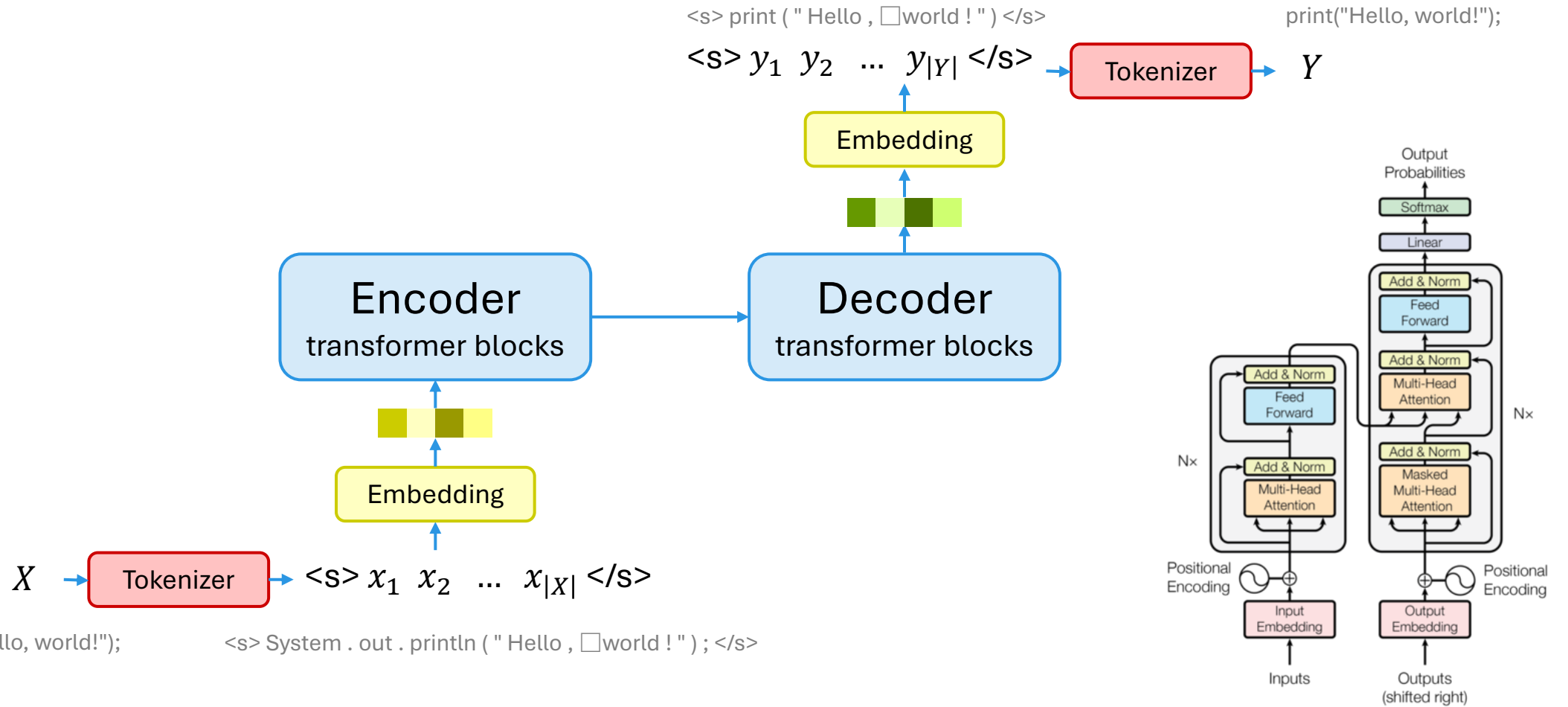
X	Y
Chinese sentence	English sentence
news article text	title



Architecture of Seq2Seq Model



Components



Tokenizer

- Used to be...

- Tokenize by whitespace / regex

```
String inputPath = args[1];
```

- Use PL's tokenizer

```
String inputPath = args [ 1 ] ;
```

- Sub-tokenize (important for PLs)

```
String input Path = args [ 1 ] ;
```

- Tokens not seen in the training set = <UNK>

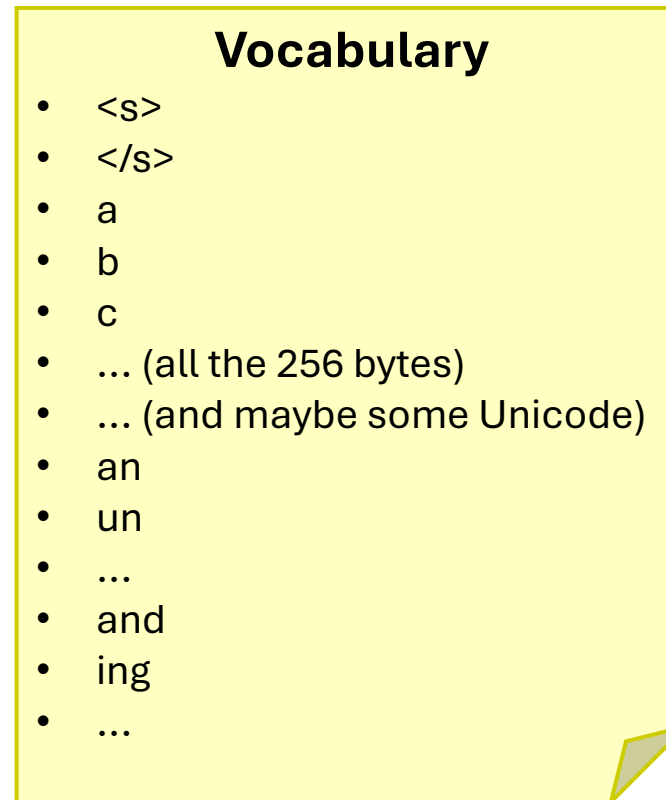
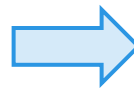
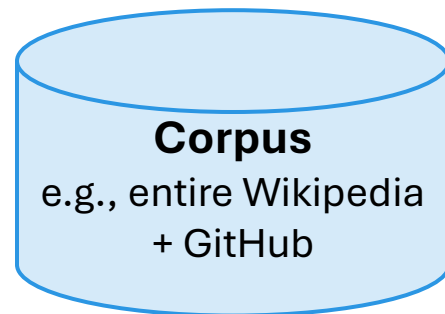
- Data-driven approach: byte-pair encoding (BPE)

Byte-Pair Encoding

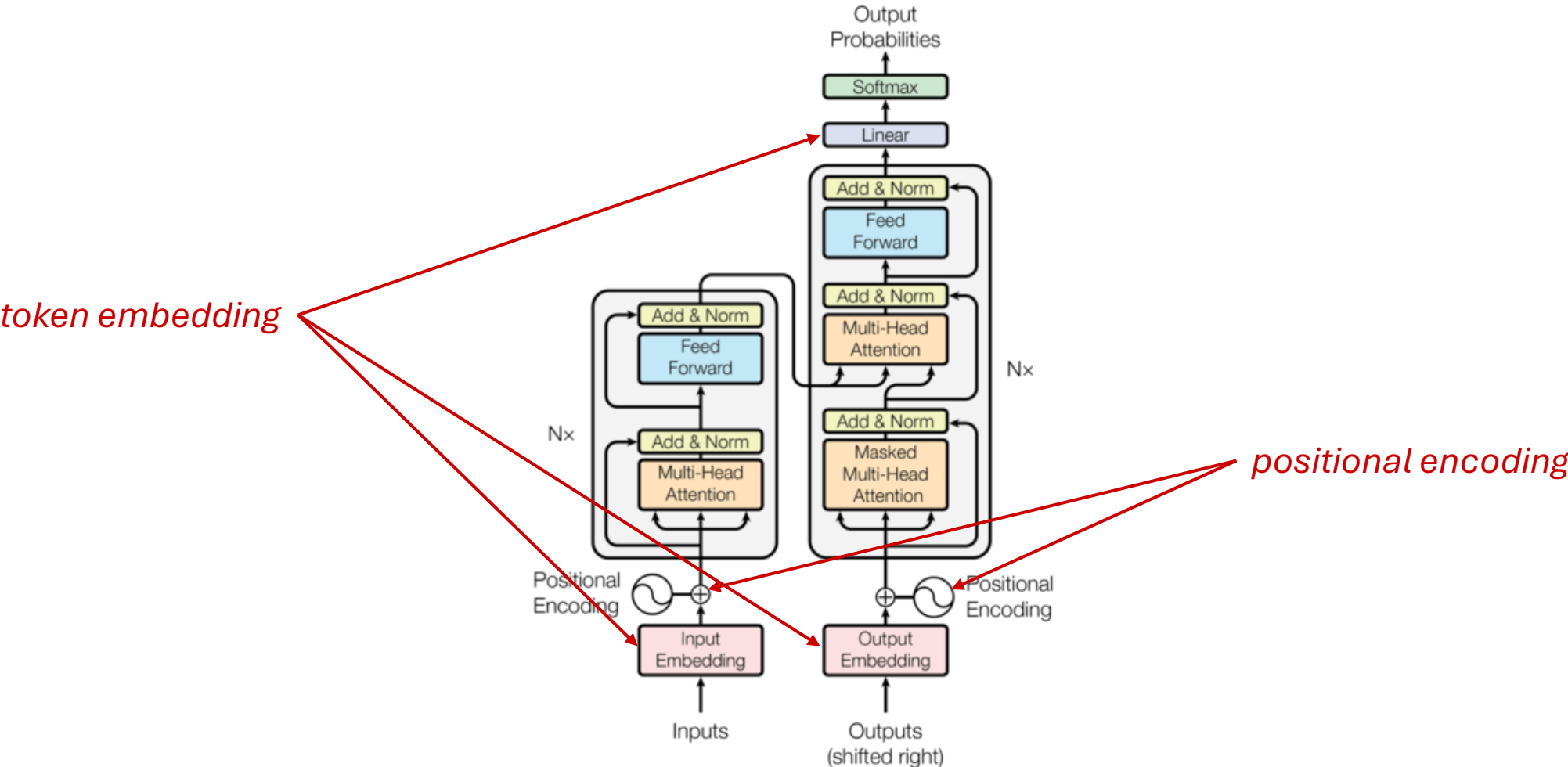
inputs: corpus, vocabulary size v

outputs: the vocabulary

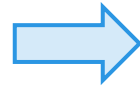
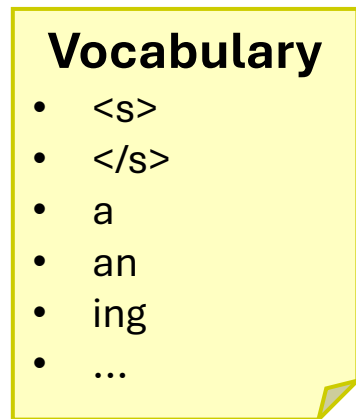
- initialize the vocabulary with base tokens
- while $|\text{vocabulary}| < v$
 - find the most common bigram in corpus
 - add that bigram as a new token



Embedding



Token/Word Embedding



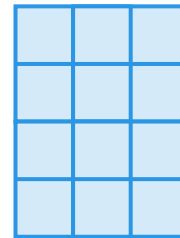
one-hot vector

(1, 0, 0, 0, 0, ...)
(0, 1, 0, 0, 0, ...)
(0, 0, 1, 0, 0, ...)
(0, 0, 0, 1, 0, ...)
(0, 0, 0, 0, 1, ...)

shape (n, v)
n: sequence length
v: vocabulary size

×

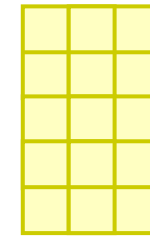
embedding matrix



shape (v, d)
d: hidden state dimension



token embedding



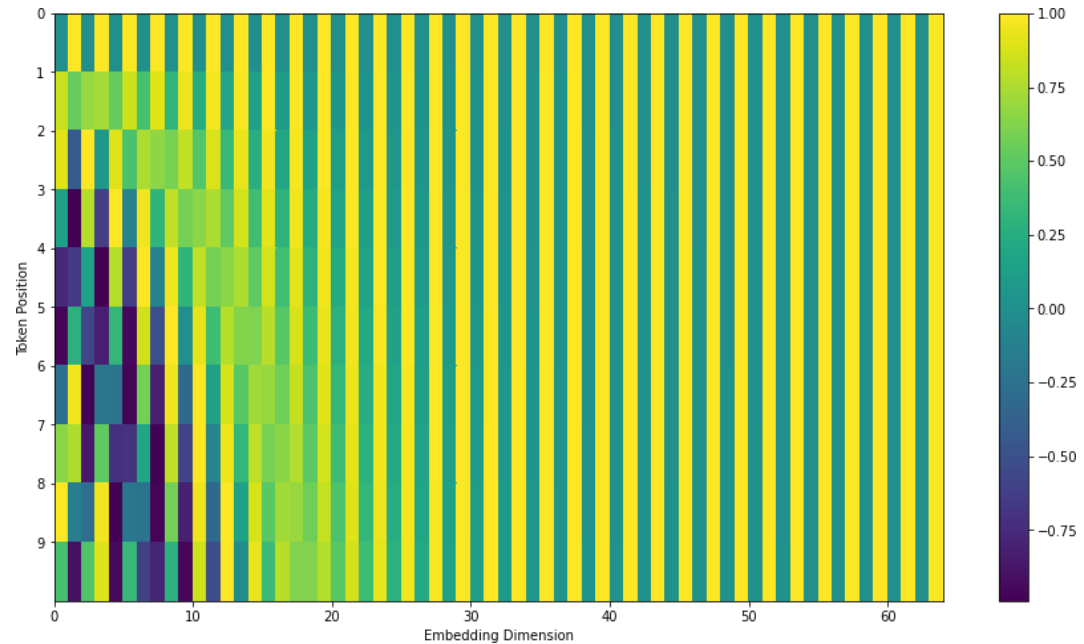
shape (n, d)

Positional Encoding

- Transformers \neq recurrent neural network;
we need something to represent the order of tokens

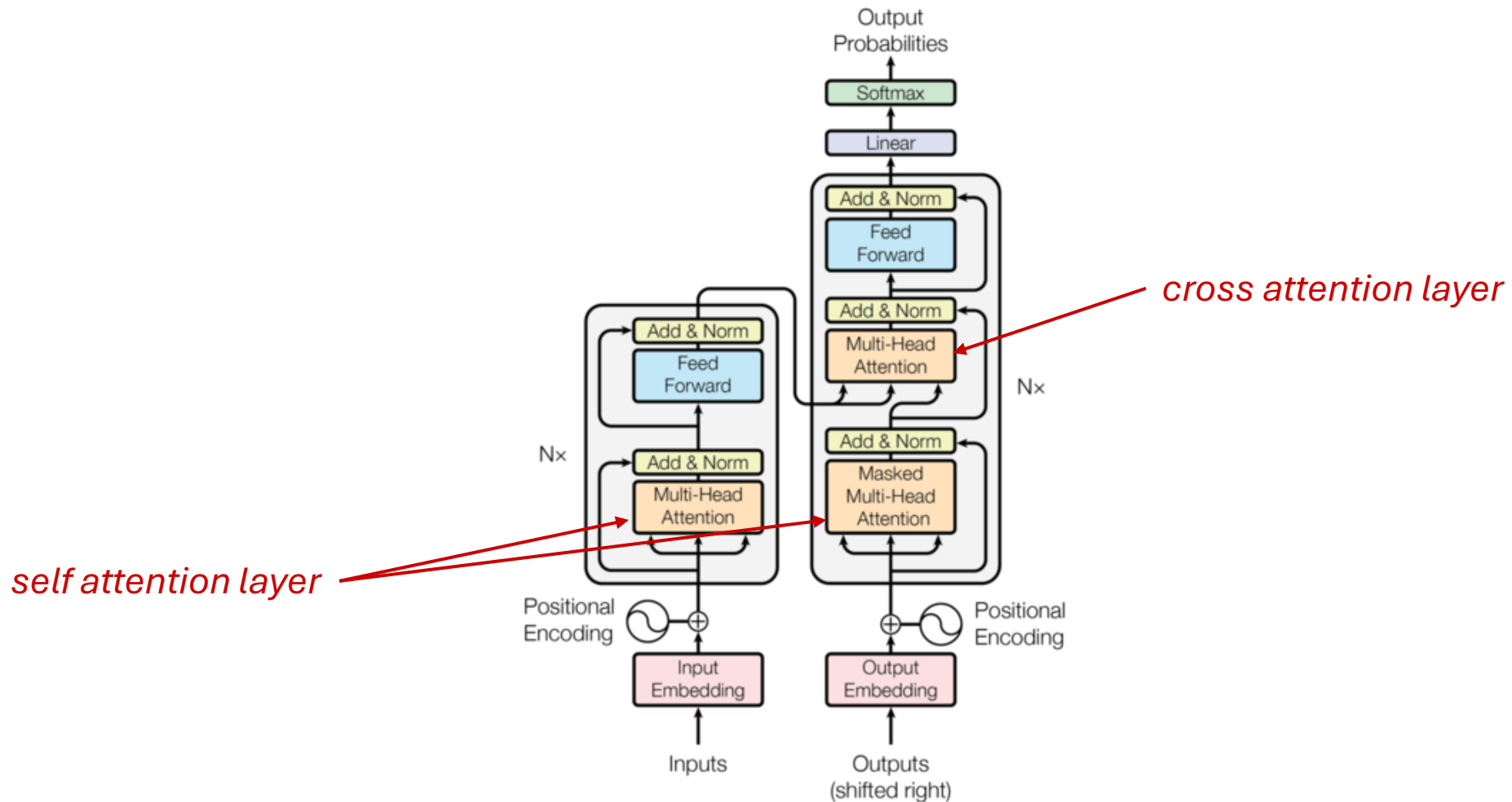
$$PE(pos, 2i) = \sin\left(\frac{pos}{10000^{\frac{2i}{d}}}\right)$$

$$PE(pos, 2i + 1) = \cos\left(\frac{pos}{10000^{\frac{2i}{d}}}\right)$$



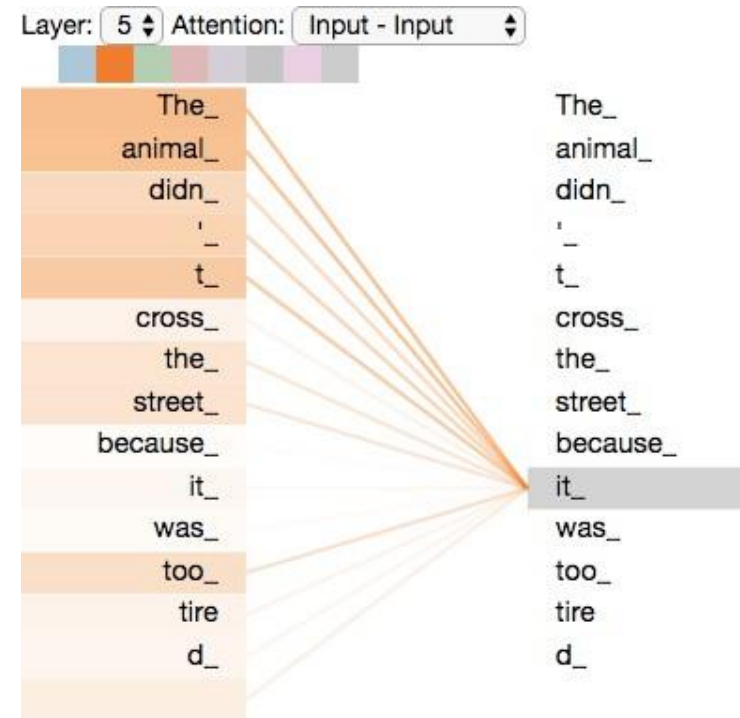
Alternative approaches: learned position embedding, attention bias (<https://arxiv.org/pdf/2108.12409>)

Attention

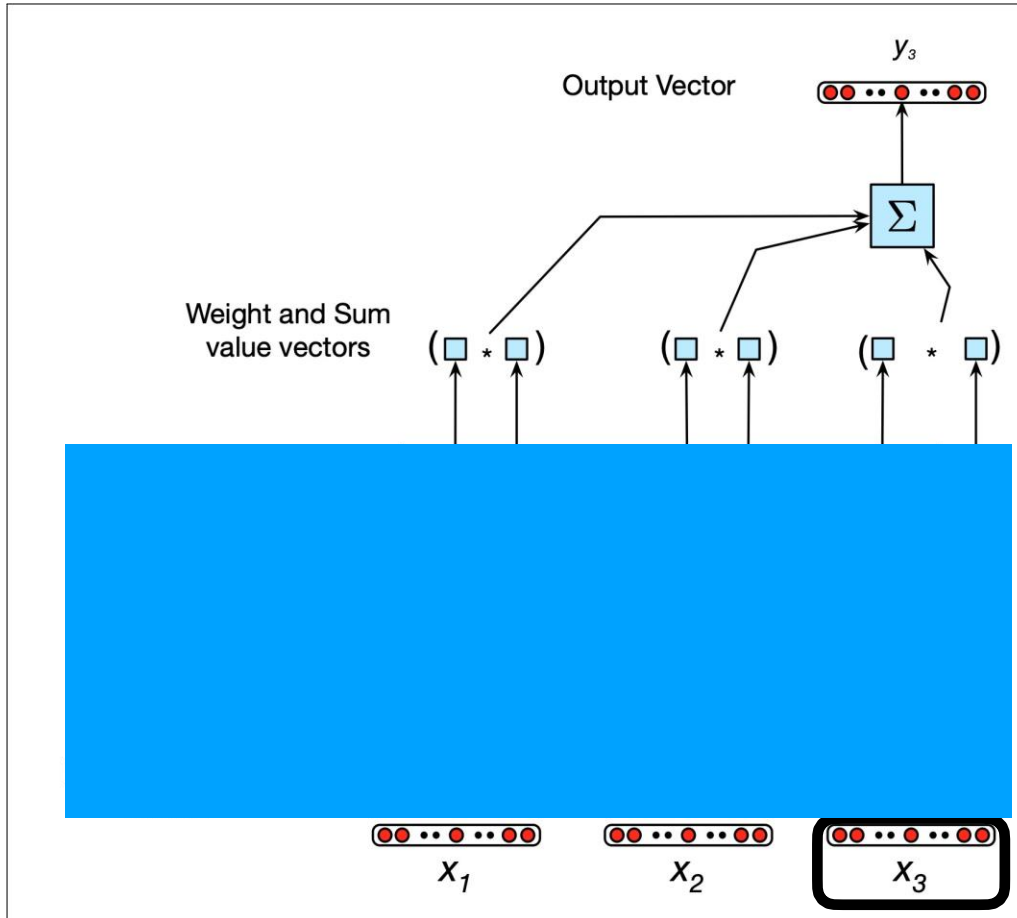


Self-attention

- Consider: “The animal didn't cross the street because it was too tired”
- What meaning should we associate with the word “it”?



Self-attention



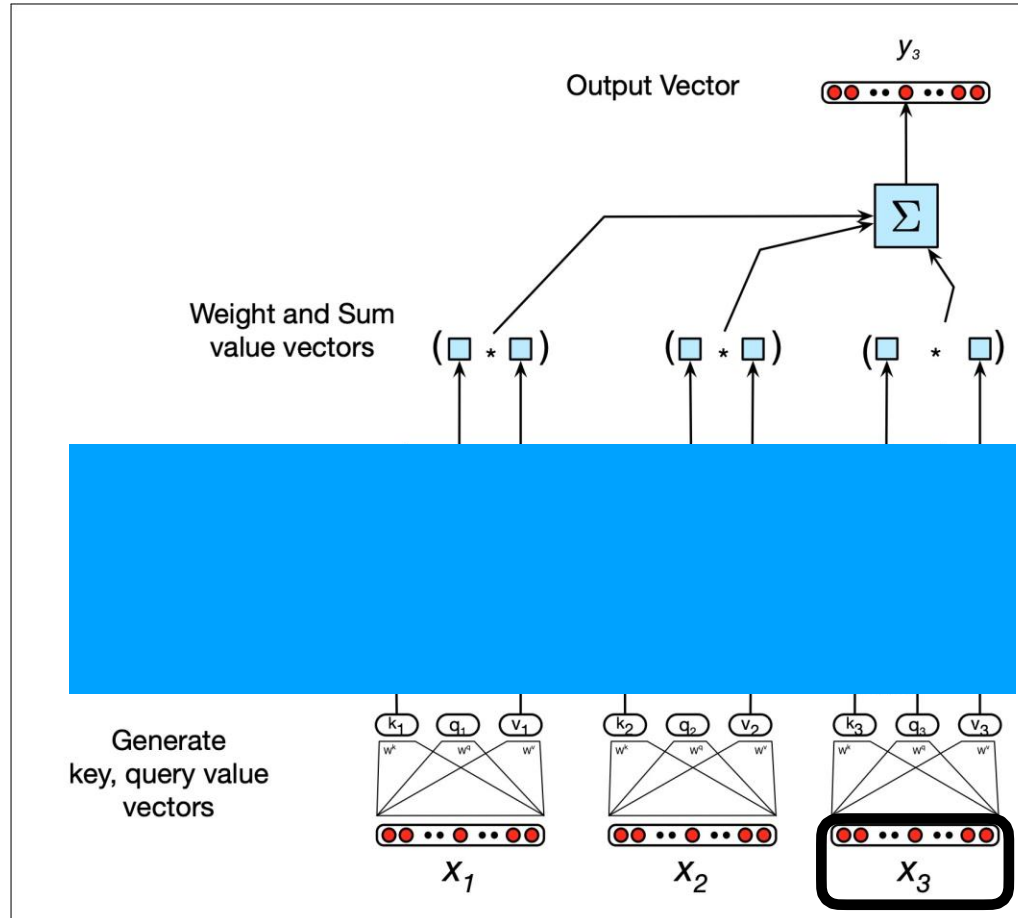
$$y_i = \sum_{j \leq i} \alpha_{ij} x_j$$

$$\begin{aligned} \alpha_{ij} &= \text{softmax}(\text{score}(x_i, x_j)) \\ &= \frac{\exp(\text{score}(x_i, x_j))}{\sum_{k=1}^i \exp(\text{score}(x_i, x_k))} \end{aligned}$$

$$\text{score}(x_i, x_j) = x_i \cdot x_j$$

But there is nothing to learn here?

Self-attention

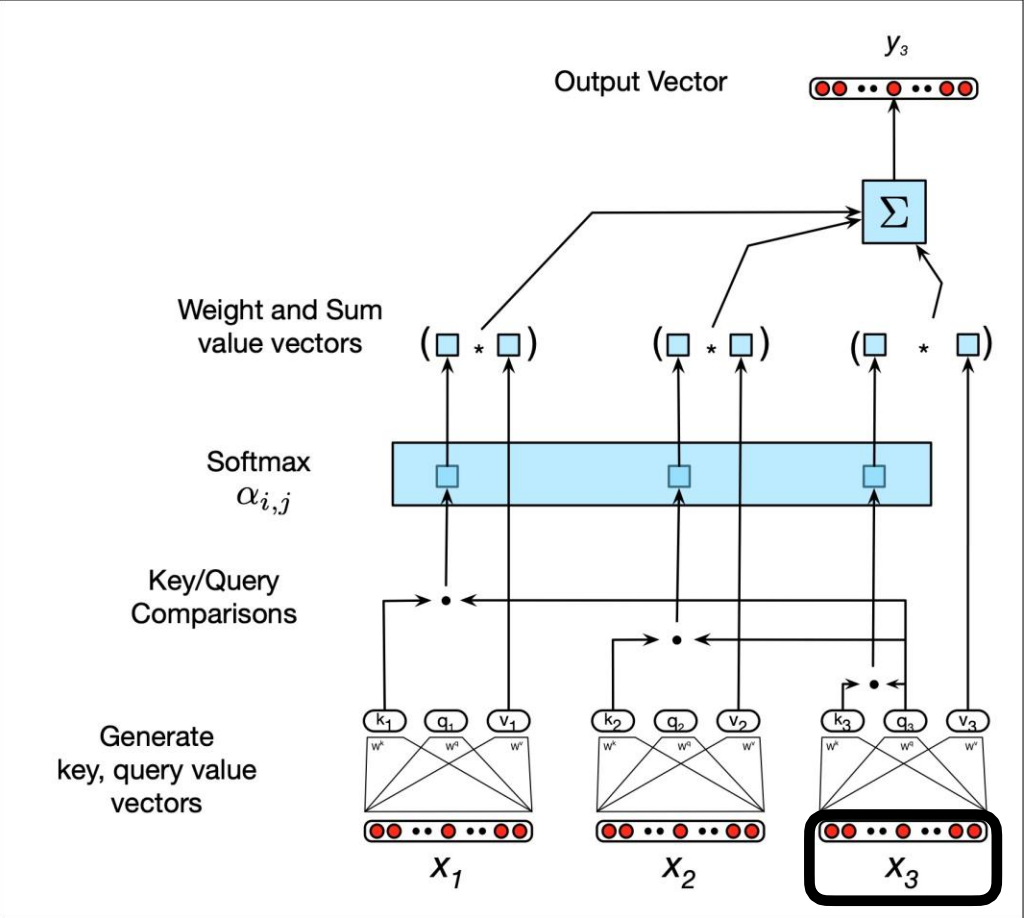


Introducing 3 types of weights, corresponding to 3 roles of each word w :

- **Query**: w is the current word under question
- **Key**: w is the word in context being compared with
- **Value**: learnable weights for the output

$$q_i = W^Q x_i; \quad k_i = W^K x_i; \quad v_i = W^V x_i$$

Self-attention



$$y_i = \sum_{j \leq i} \alpha_{ij} v_j$$

also multiply with v

$$\alpha_{ij} = \text{softmax}(\text{score}(x_i, x_j))$$

normalize

$$\text{score}(x_i, x_j) = q_i \cdot k_j$$

$$q_i = W^Q x_i; k_i = W^K x_i; v_i = W^V x_i$$

Self-attention in one graph

$$\begin{matrix} \mathbf{X} \\ \text{2x4 grid} \end{matrix} \times \begin{matrix} \mathbf{W}^Q \\ \text{3x4 grid} \end{matrix} = \begin{matrix} \mathbf{Q} \\ \text{2x3 grid} \end{matrix}$$

$$\begin{matrix} \mathbf{X} \\ \text{2x4 grid} \end{matrix} \times \begin{matrix} \mathbf{W}^K \\ \text{3x4 grid} \end{matrix} = \begin{matrix} \mathbf{K} \\ \text{2x3 grid} \end{matrix}$$

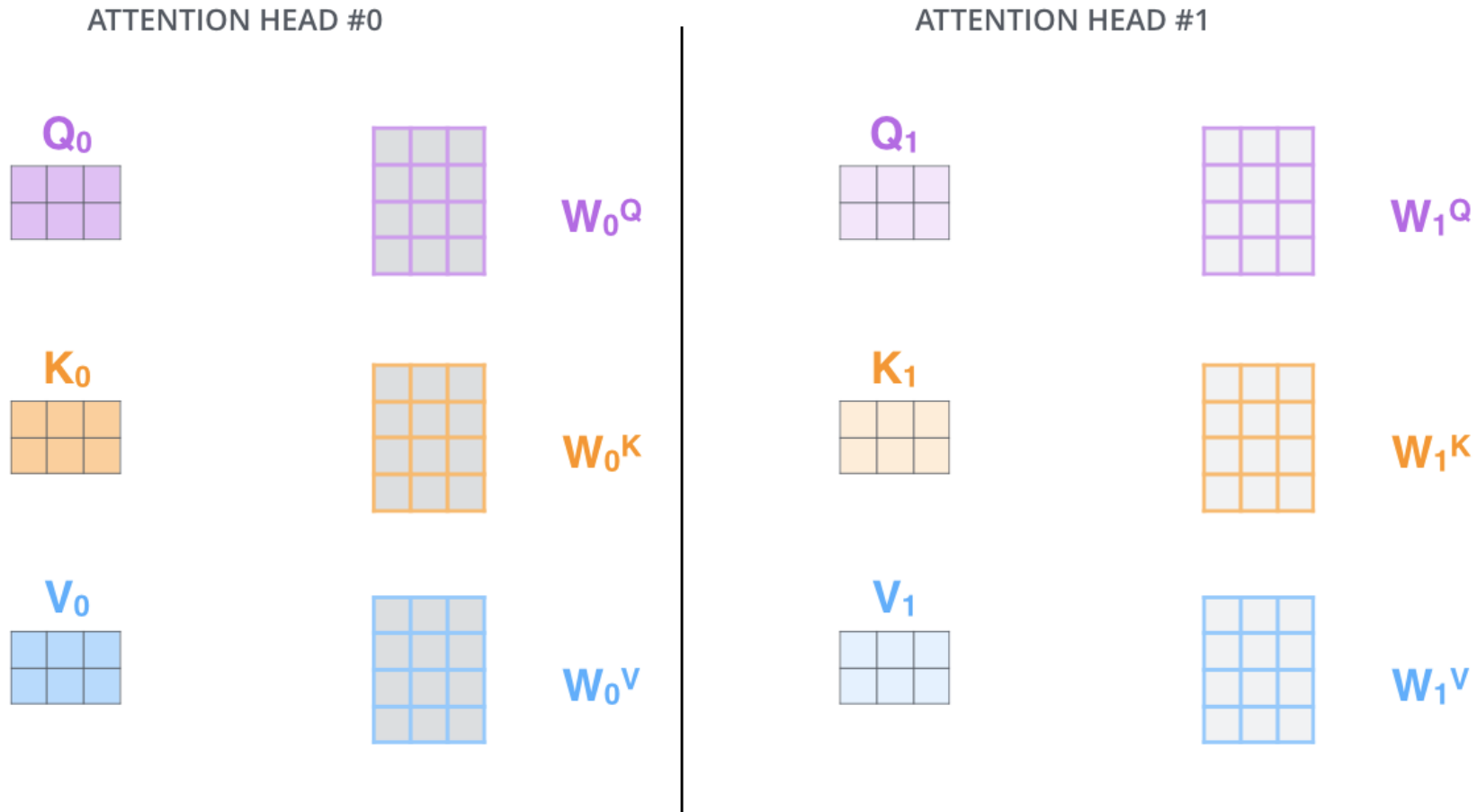
$$\begin{matrix} \mathbf{X} \\ \text{2x4 grid} \end{matrix} \times \begin{matrix} \mathbf{W}^V \\ \text{3x4 grid} \end{matrix} = \begin{matrix} \mathbf{V} \\ \text{2x3 grid} \end{matrix}$$

$$\text{softmax} \left(\frac{\begin{matrix} \mathbf{Q} \\ \text{2x3 grid} \end{matrix} \times \begin{matrix} \mathbf{K}^T \\ \text{3x2 grid} \end{matrix}}{\sqrt{d_k}} \right) \begin{matrix} \mathbf{V} \\ \text{2x3 grid} \end{matrix}$$
$$= \begin{matrix} \mathbf{Z} \\ \text{2x3 grid} \end{matrix}$$

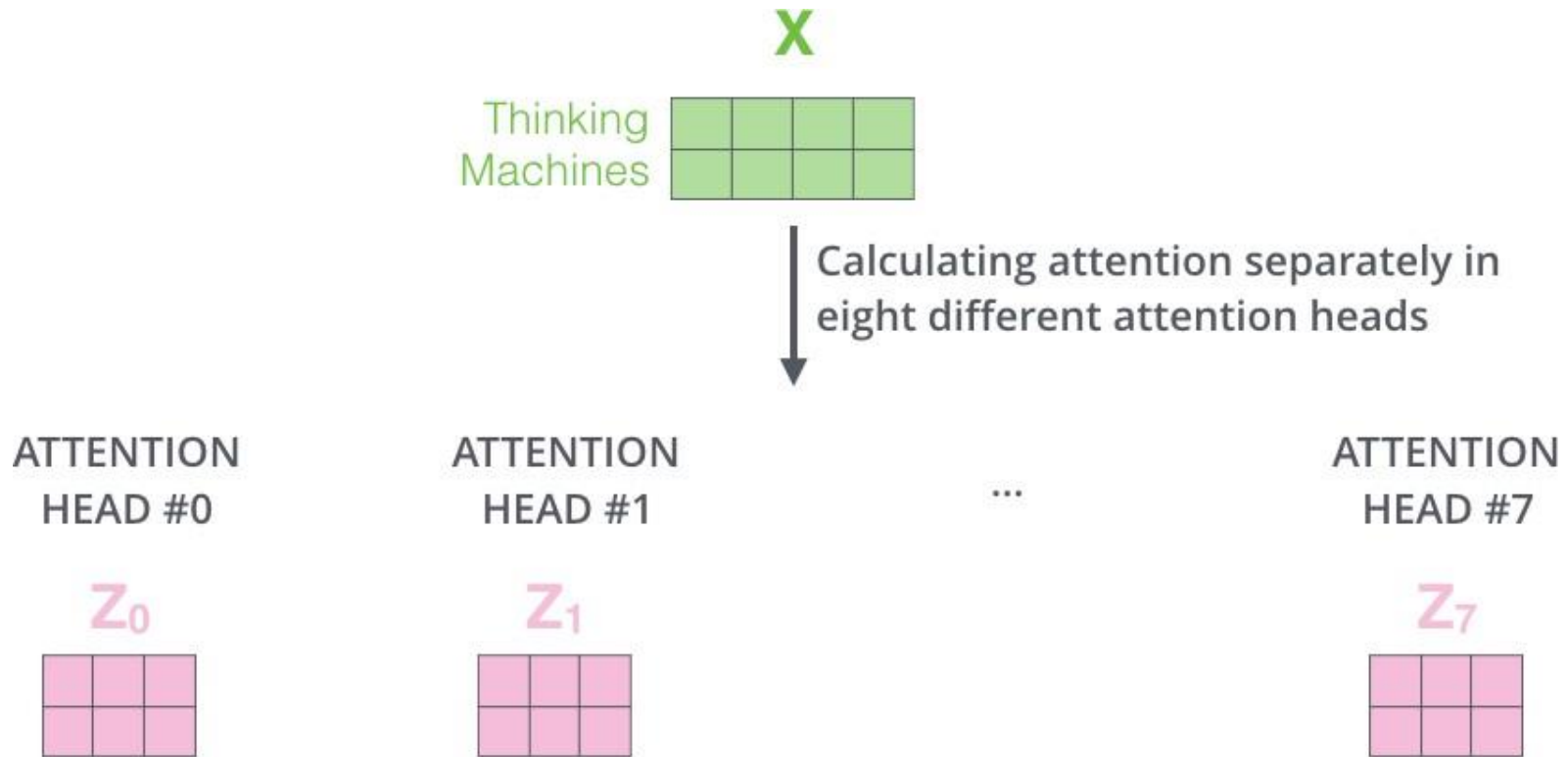
This is called one attention "head"...

Multi-headed attention

- Or, the beast with many heads

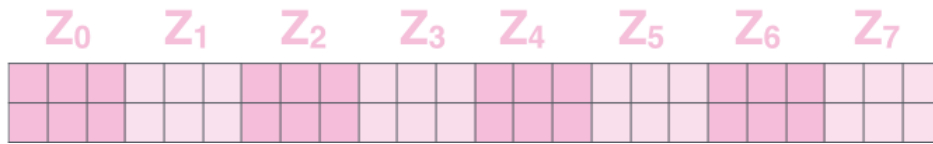


Multi-headed attention



Multi-headed attention

1) Concatenate all the attention heads



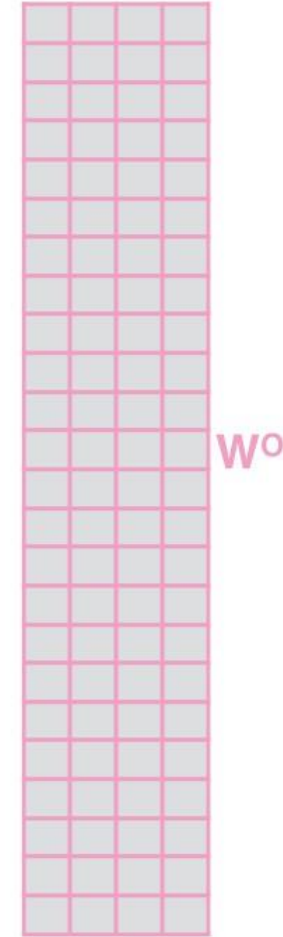
Multi-headed attention

1) Concatenate all the attention heads



2) Multiply with a weight matrix W^O that was trained jointly with the model

x



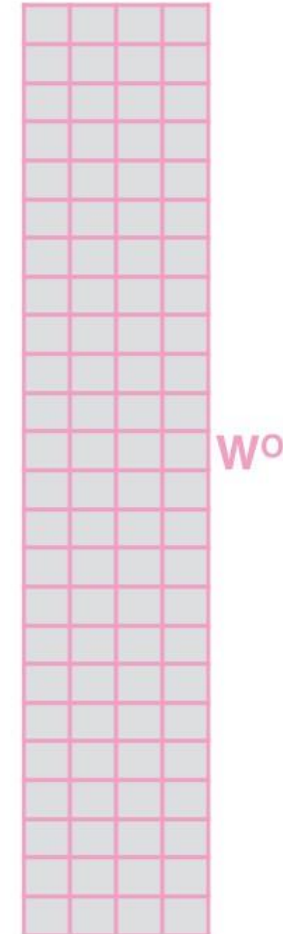
Multi-headed attention

1) Concatenate all the attention heads

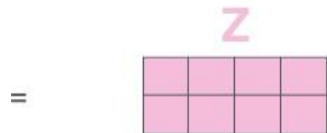


2) Multiply with a weight matrix W^O that was trained jointly with the model

x



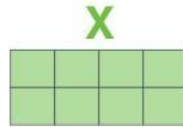
3) The result would be the Z matrix that captures information from all the attention heads. We can send this forward to the FFNN



Multi-headed attention

- 1) This is our input sentence*
- 2) We embed each word*

Thinking
Machines



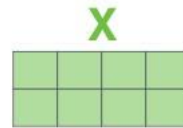
* In all encoders other than #0,
we don't need embedding.
We start directly with the output
of the encoder right below this one

Multi-headed attention

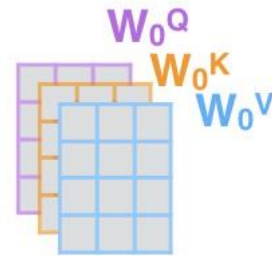
1) This is our input sentence*

Thinking
Machines

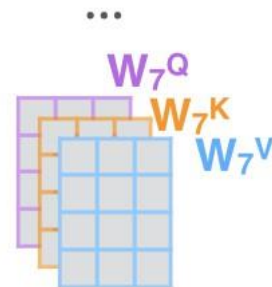
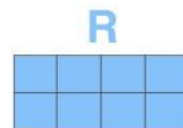
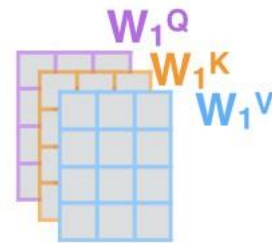
2) We embed each word*



3) Split into 8 heads.
We multiply X or R with weight matrices



* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one

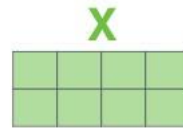


Multi-headed attention

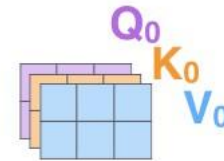
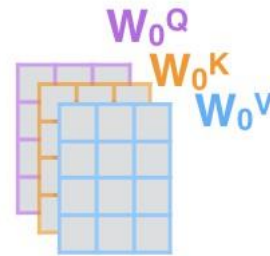
1) This is our input sentence*

Thinking
Machines

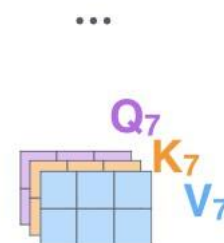
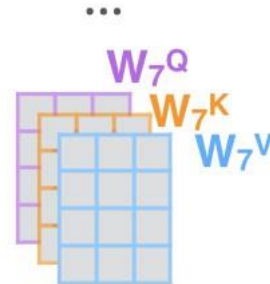
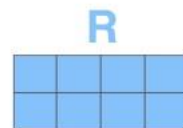
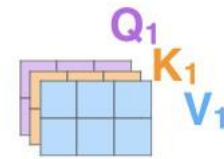
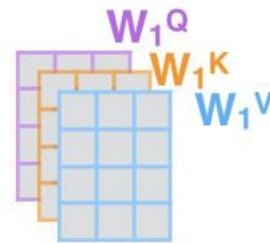
2) We embed each word*



3) Split into 8 heads.
We multiply X or R with weight matrices



* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one



Multi-headed attention

1) This is our input sentence*

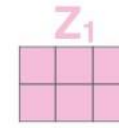
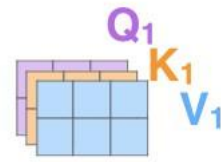
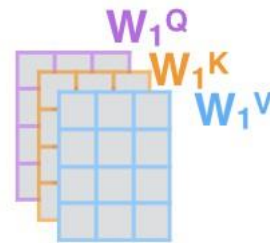
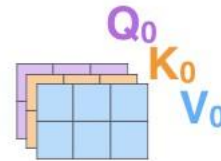
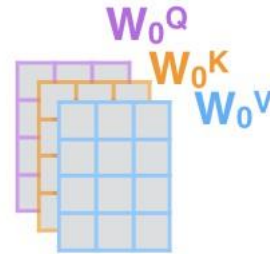
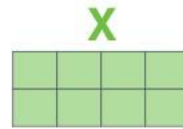
2) We embed each word*

3) Split into 8 heads. We multiply X or R with weight matrices

4) Calculate attention using the resulting $Q/K/V$ matrices

5) Concatenate the resulting Z matrices, then multiply with weight matrix W^O to produce the output of the layer

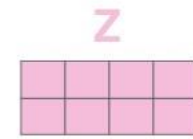
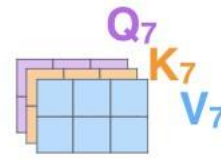
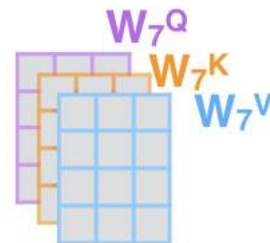
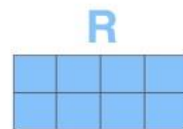
Thinking
Machines



...

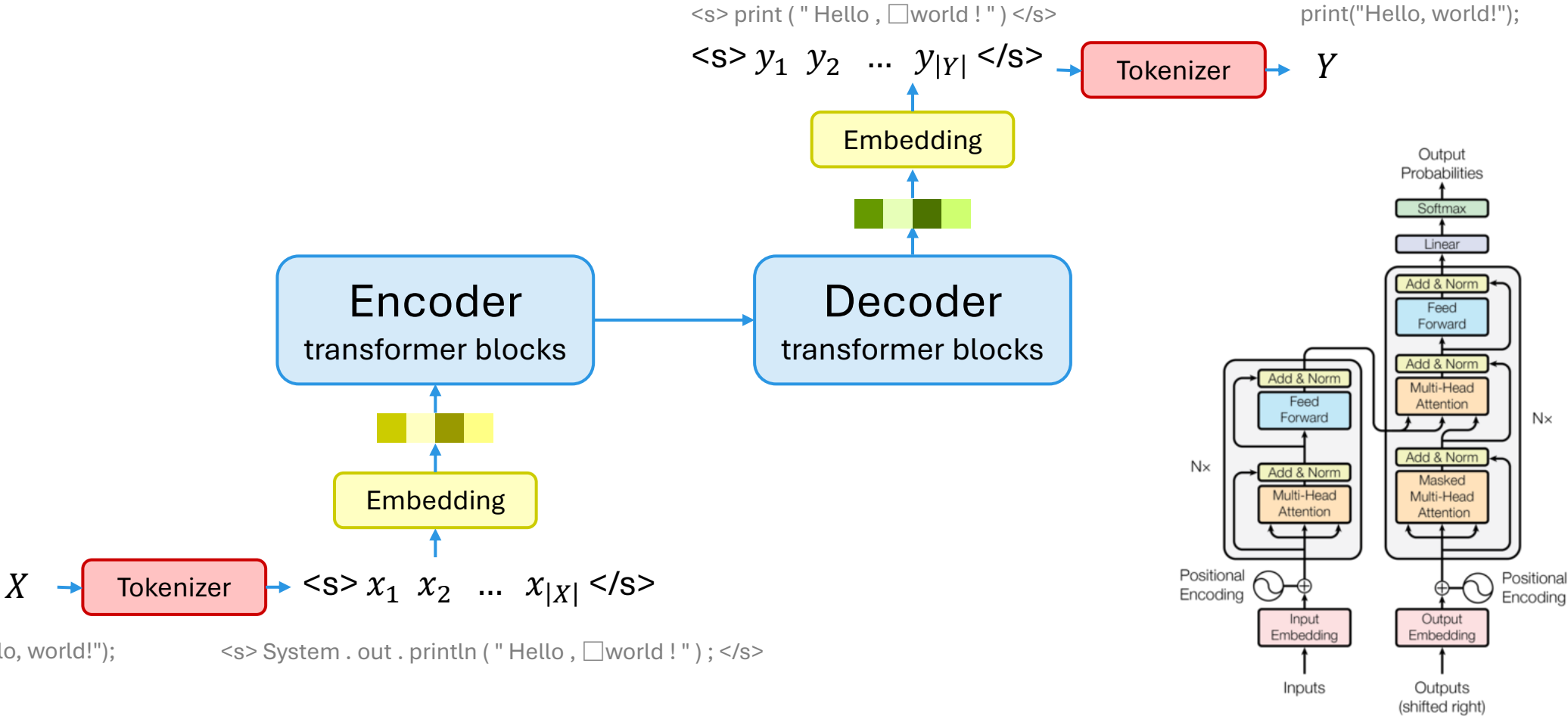
...

...



* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one

Recap



Variants: Encoder-only and Decoder-only

- Encoder-decoder: T5, BART
 - good for transduction tasks
- Encoder-only: BERT
 - good for classification tasks
- Decoder-only: GPT, Llama
 - good for generation tasks

