# Natural Language Processing

**LLMs, Transformer, PEFT, RAG, LLMs Project Life Cycle**

**KyawSwarTun**

**Most illustrations in these slides are from the internet, will be used for an educational purpose.**

# History of Large Language Models



**1967** ELiza

**1997** LSTM

**2010** Stanford CoreNLP

**2011** Google Brain

**2017** Google Transformer Architecture
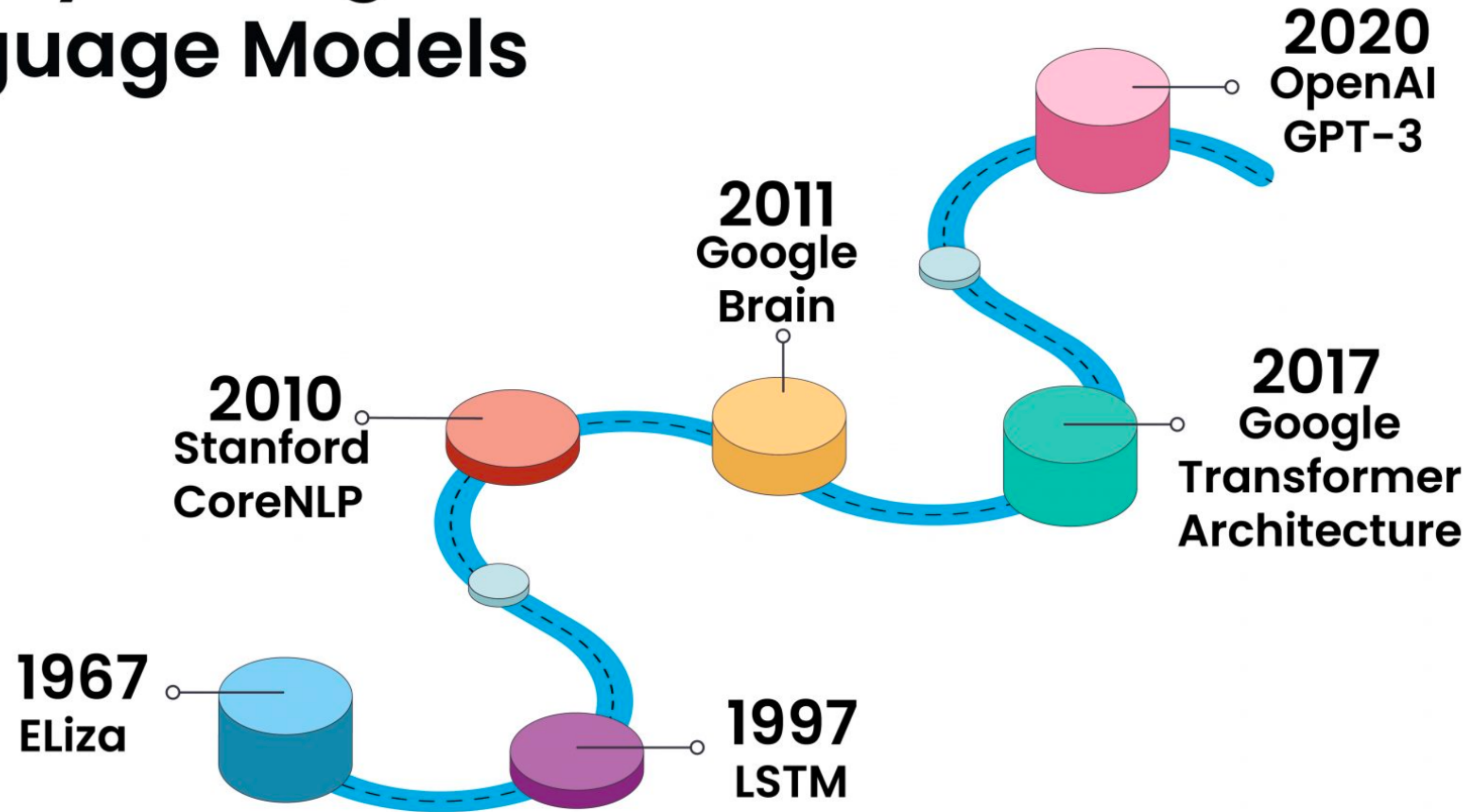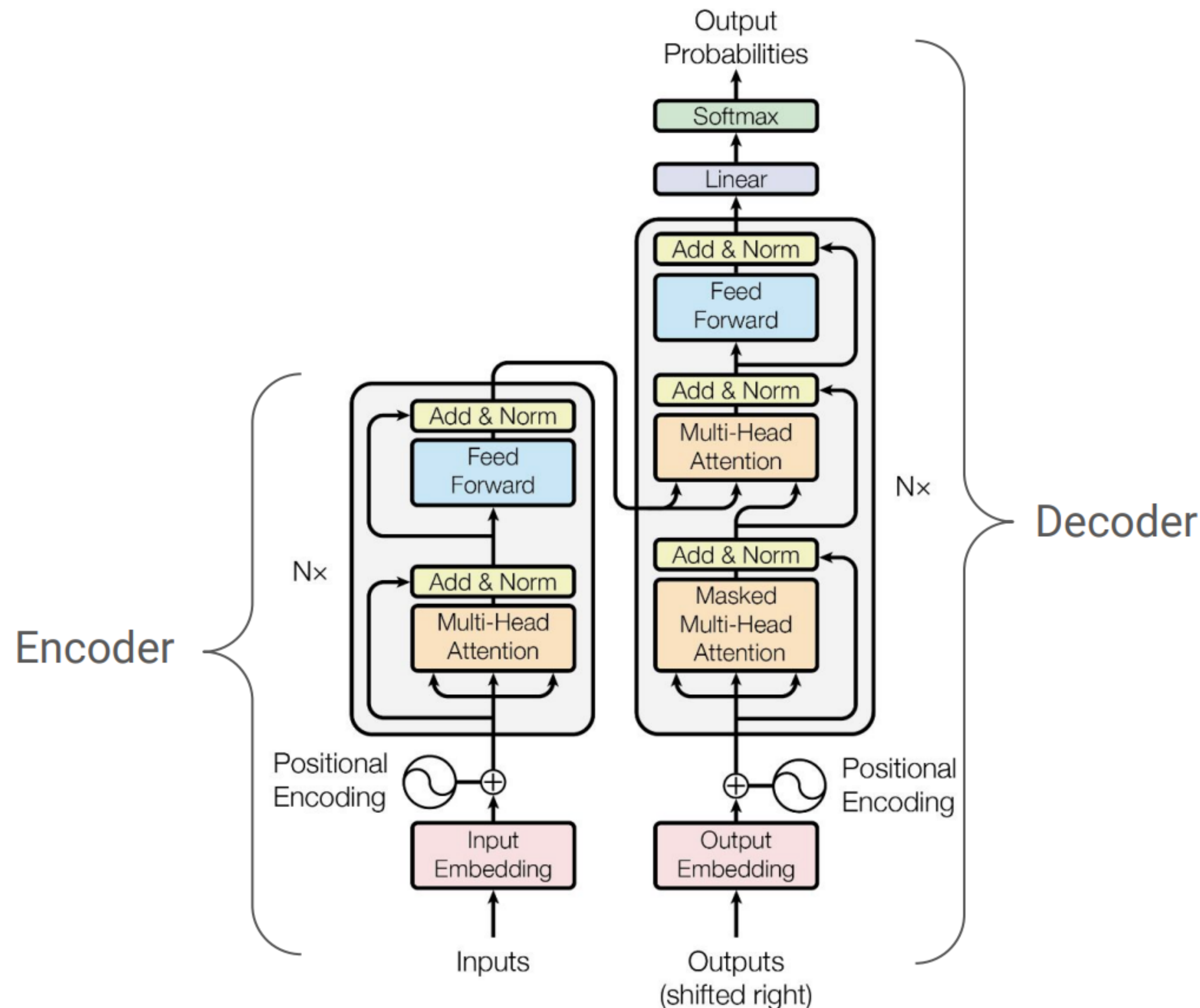
**2020** OpenAI GPT-3

Fig. 1. History of Large Language Models from [2]
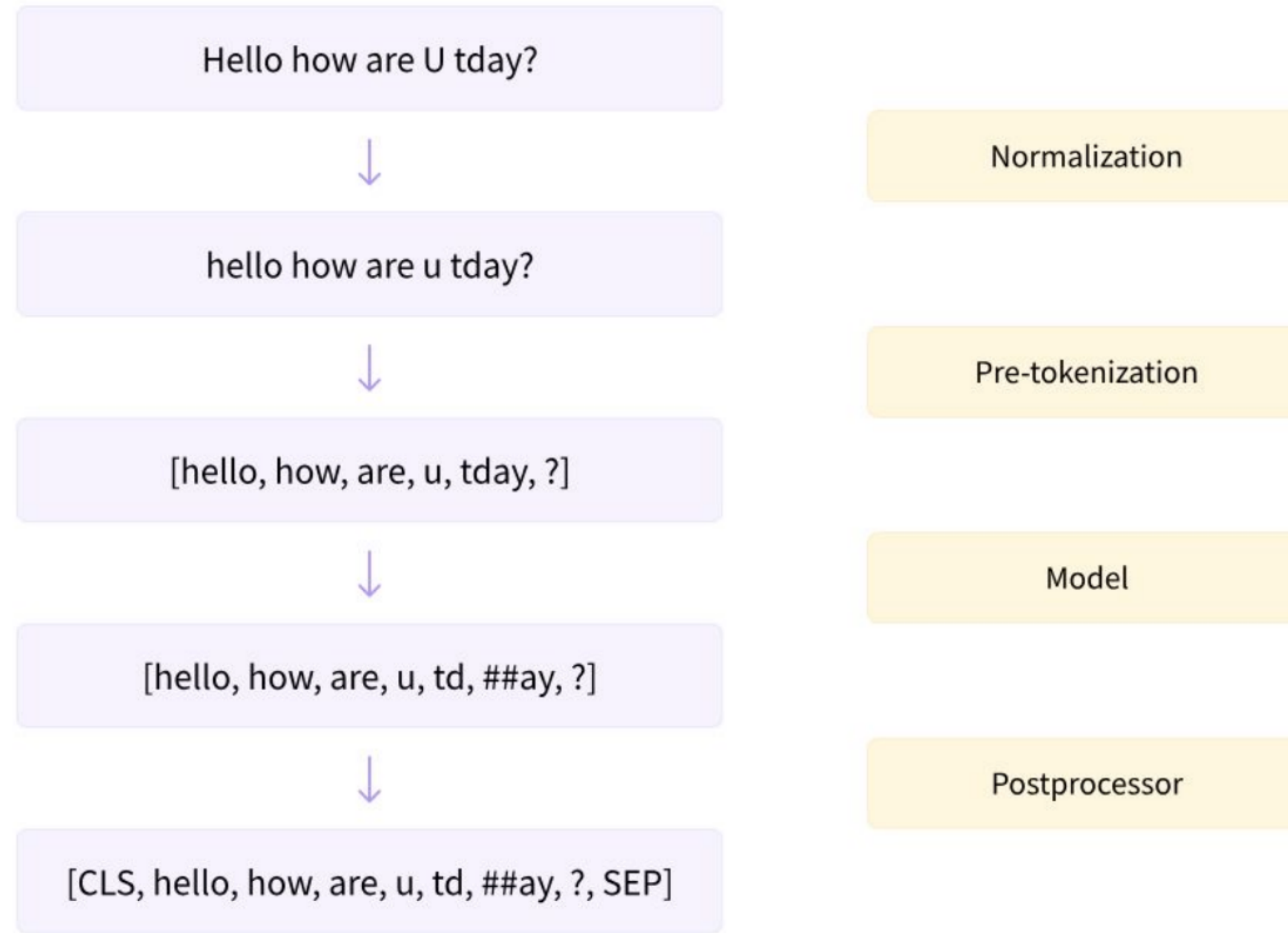
# Transformer Block

Next, we will learn exactly how the Transformer architecture works:

First, we will talk about the Encoder!

Next, we will go through the Decoder (which is quite similar)!

# Tokenization: decomposing a sentence into a sequence of tokens

Hello how are U tday?

↓

Normalization

hello how are u tday?

↓

Pre-tokenization

[hello, how, are, u, tday, ?]

↓

Model

[hello, how, are, u, td, ##ay, ?]

↓

Postprocessor

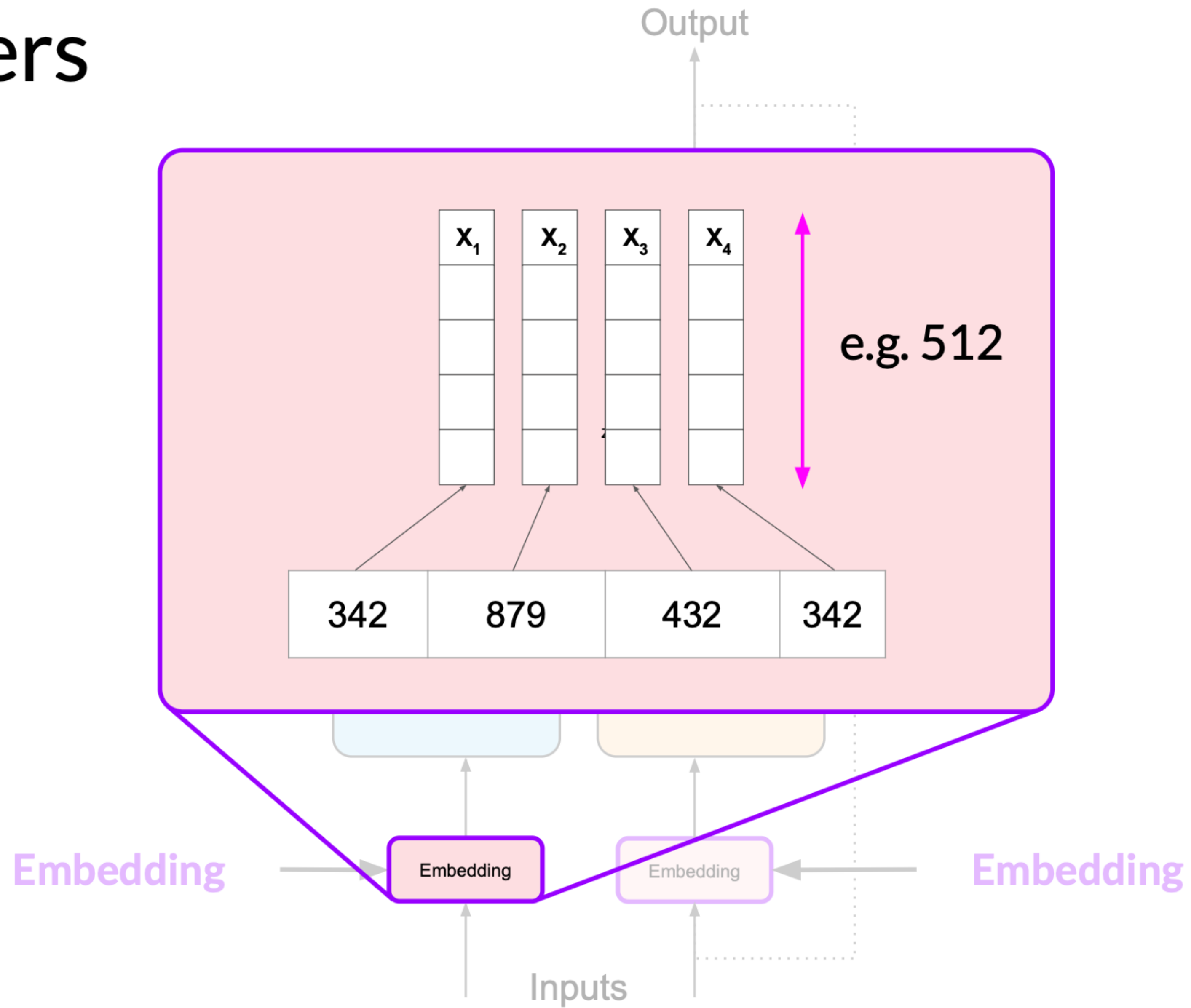[CLS, hello, how, are, u, td, ##ay, ?, SEP]

Every single explanation you will ever see about Language Models use **words**, **BUT** in reality the unit object is **tokens**
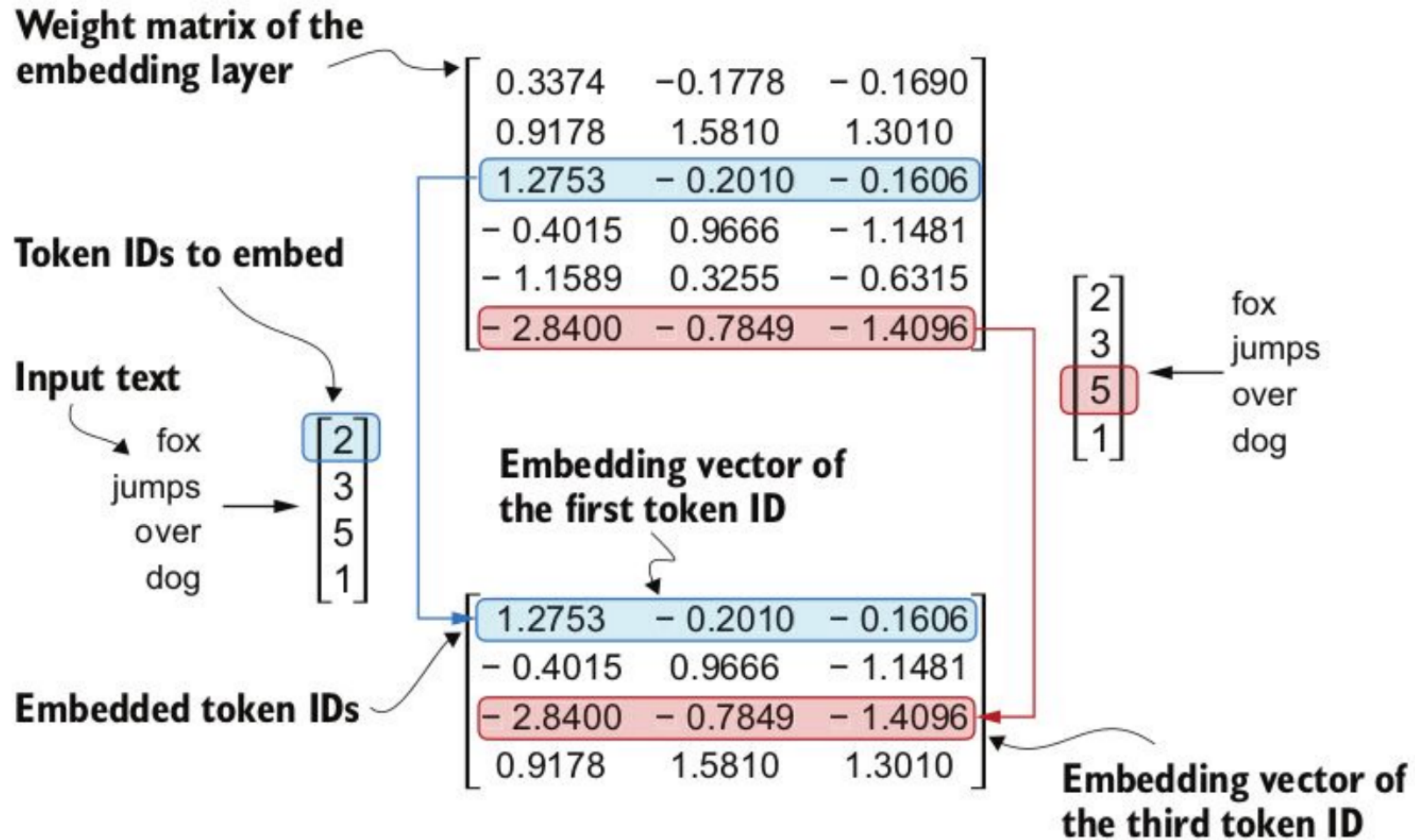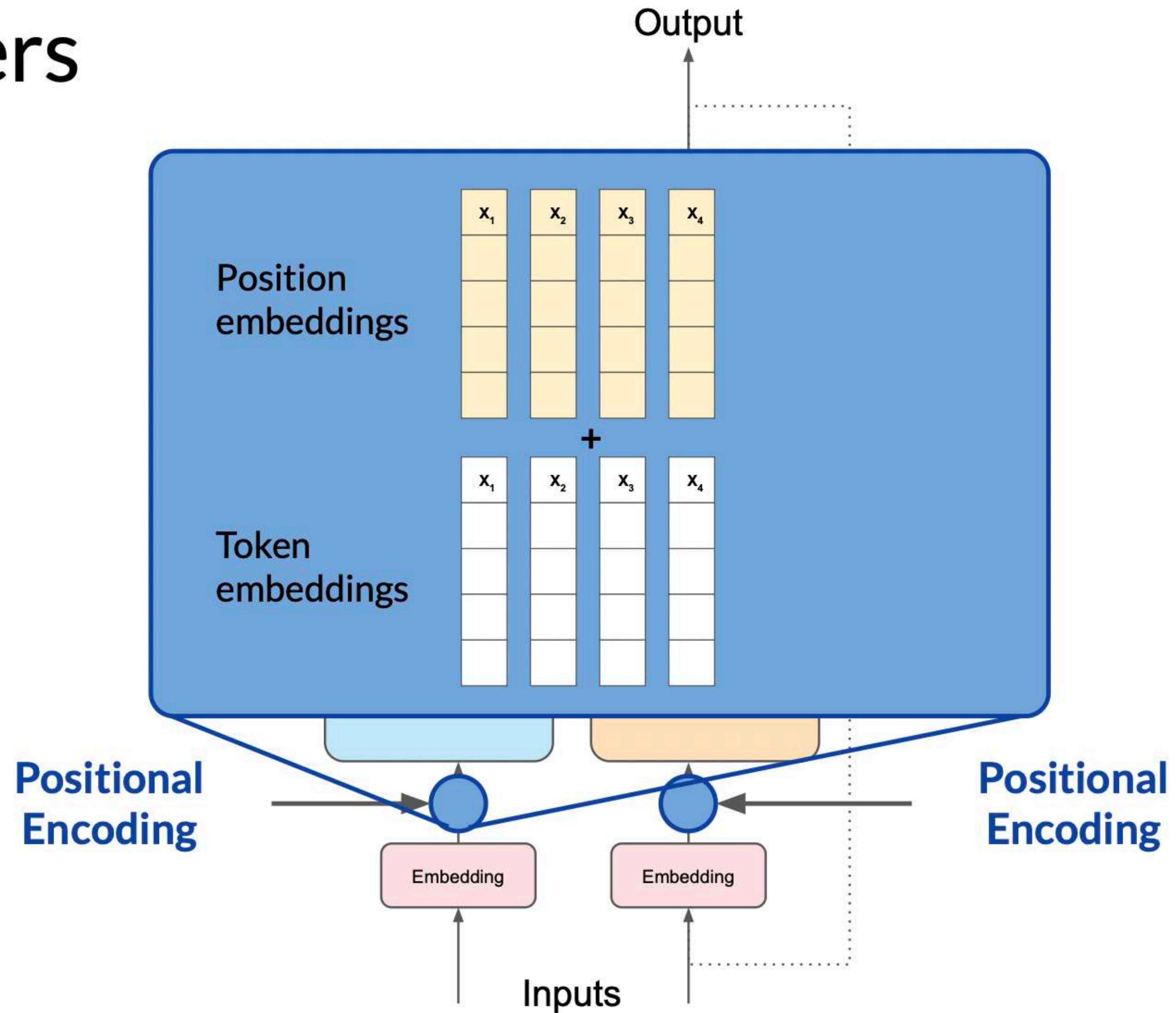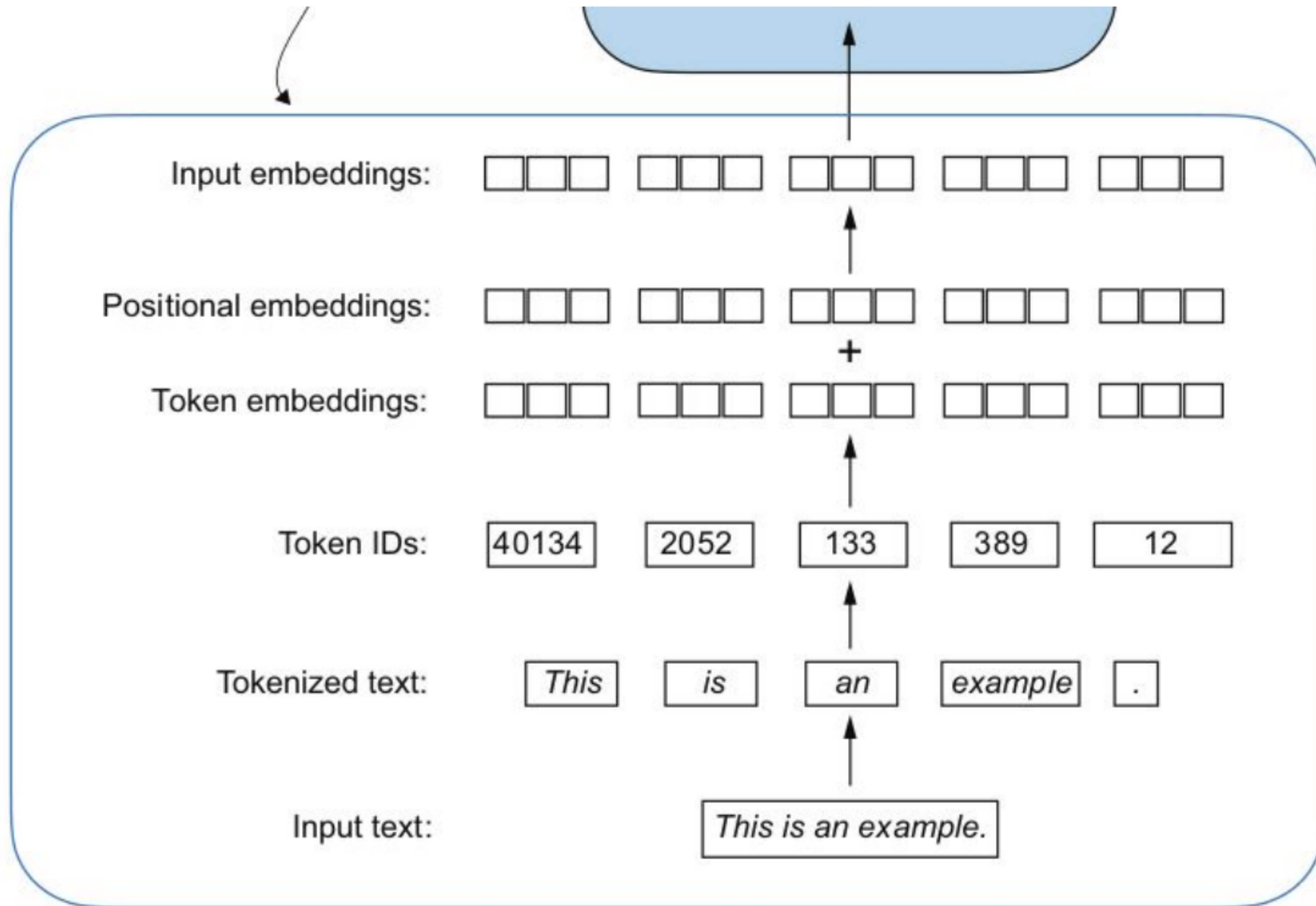
# WORDS != TOKENS

# Transformers



Encoder

Decoder

Embedding

Embedding

Inputs

Tokenizer

| 156 | 790 | 321 | 890 | 156 |

Token IDs

Input:　the　teach　er　taught　the

# Transformers

# From text to vectors



**Weight matrix of the embedding layer**

$$\begin{bmatrix} 0.3374 & -0.1778 & -0.1690 \\ 0.9178 & 1.5810 & 1.3010 \\ 1.2753 & -0.2010 & -0.1606 \\ -0.4015 & 0.9666 & -1.1481 \\ -1.1589 & 0.3255 & -0.6315 \\ -2.8400 & -0.7849 & -1.4096 \end{bmatrix}$$

**Token IDs to embed**

**Input text**

fox
jumps
over
dog

$$\begin{bmatrix} 2 \\ 3 \\ 5 \\ 1 \end{bmatrix}$$

$$\begin{bmatrix} 2 \\ 3 \\ 5 \\ 1 \end{bmatrix}$$ fox jumps over dog

**Embedding vector of the first token ID**

**Embedded token IDs**

$$\begin{bmatrix} 1.2753 & -0.2010 & -0.1606 \\ -0.4015 & 0.9666 & -1.1481 \\ -2.8400 & -0.7849 & -1.4096 \\ 0.9178 & 1.5810 & 1.3010 \end{bmatrix}$$

**Embedding vector of the third token ID**

# Transformers

Input embeddings:

Positional embeddings:

+

Token embeddings:

Token IDs:    40134    2052    133    389    12

Tokenized text:    This    is    an    example    .

Input text:    This is an example.

# Statistics

The smallest GPT-2 models (117M and 125M parameters) use an embedding size of 768 dimensions.

The largest GPT-3 model (175B parameters) uses an embedding size of 12,288 dimensions.

# A self-attention head

**Input**: an embedding vector x(i) for each token i

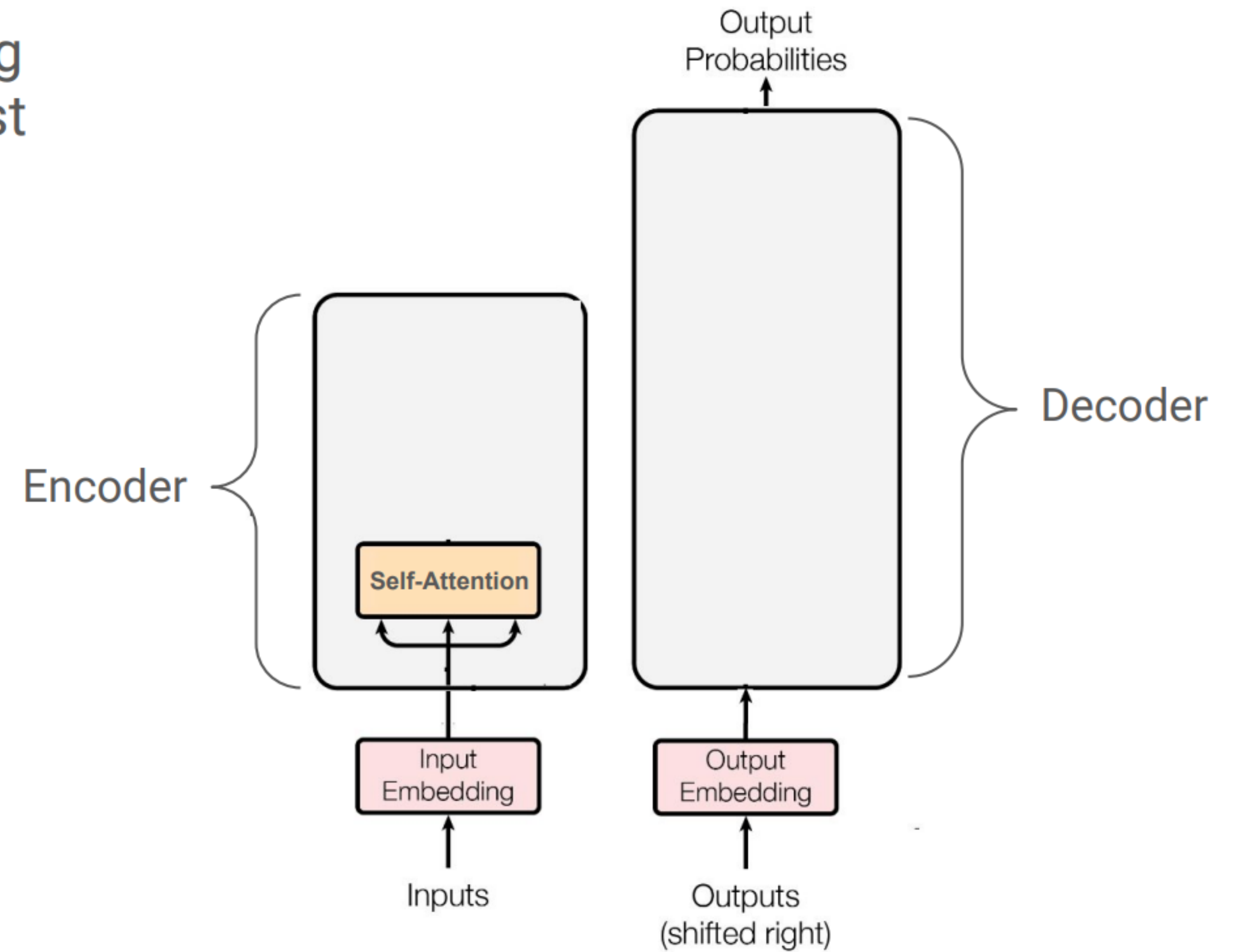**Output**: a context vector z(i) for each token i

*Intuition*: z(i) gathers *contextual* information

# Computing context vectors

Computing context vectors is very easy assuming we have computed **attention weights**: alpha(i,j) describes the importance of token j for token i.

# Encoder: Self-Attention

Self-Attention is the core building block of Transformer, so let's first focus on that!

# Keys, queries, and values

**Input**: an embedding vector x(i) for each token i

**Output**: for each token i:

- A _query_ vector q(i), describing the information token i is interested in,
- A _key_ vector k(i), whose goal is to match the relevant queries for token i,
- A _value_ vector v(i), describing the information contained by token i.

# Recipe for (Vectorized) Self-Attention in the Transformer Encoder

- Step 1: With embeddings stacked in $X$, calculate queries, keys, and values.

$$Q = XW^Q \quad K = XW^K \quad V = XW^V$$

- Step 2: Calculate attention score between query and keys.

$$E = QK^T$$
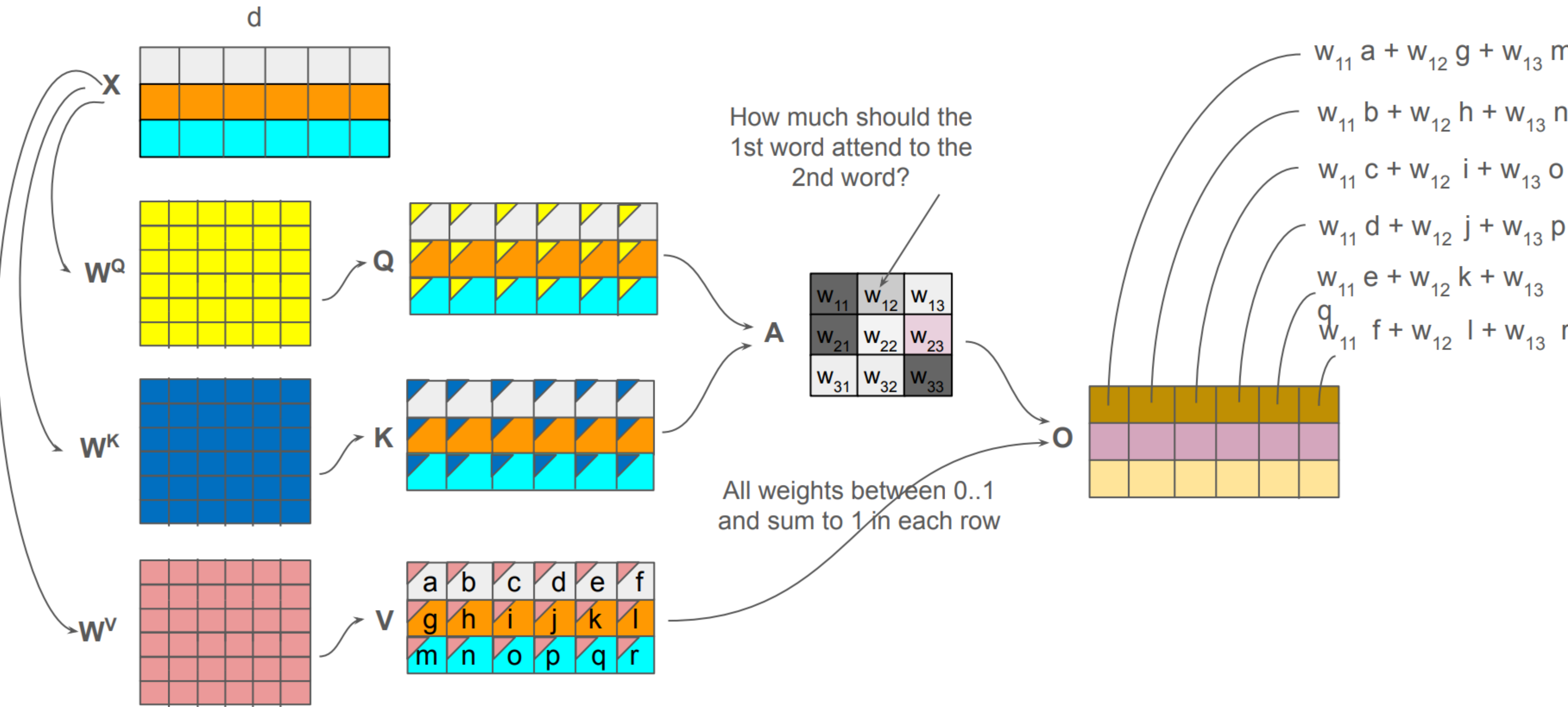
- Step 3: Take the softmax to normalize attention scores.

$$A = softmax(E)$$

- Step 4: Take a weighted sum of values.

$$Output = AV$$

$$\boxed{Output = softmax(QK^T)V}$$

# In Pictures (N = 3, d = 6, h = 1)



$w_{11} a + w_{12} g + w_{13} m$

$w_{11} b + w_{12} h + w_{13} n$

$w_{11} c + w_{12} i + w_{13} o$

$w_{11} d + w_{12} j + w_{13} p$

$w_{11} e + w_{12} k + w_{13} q$

$w_{11} f + w_{12} l + w_{13} r$

How much should the 1st word attend to the 2nd word?

All weights between 0..1 and sum to 1 in each row

# Computing attention scores and weights

Now we focus on the core computation: attention scores and weights.

We first compute **attention scores**, and then normalise them into **attention weights**.