



الجامعة السورية الخاصة  
SYRIAN PRIVATE UNIVERSITY

المحاضرة الثالثة

كلية الهندسة

الذكاء الصناعي العملي

# Transformers: Exploring Transformers Through BERT

د. رياض سنبل

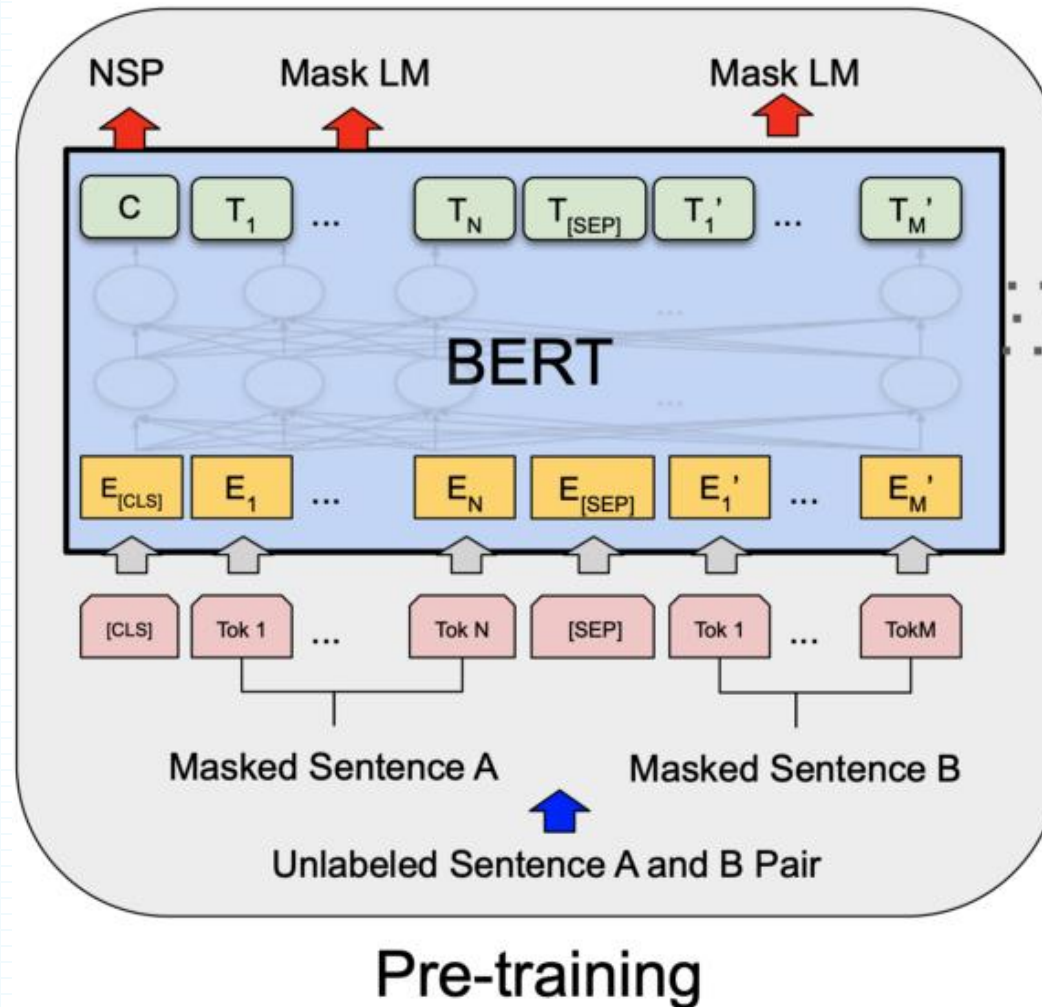
# BERT

Next Sentence  
Prediction

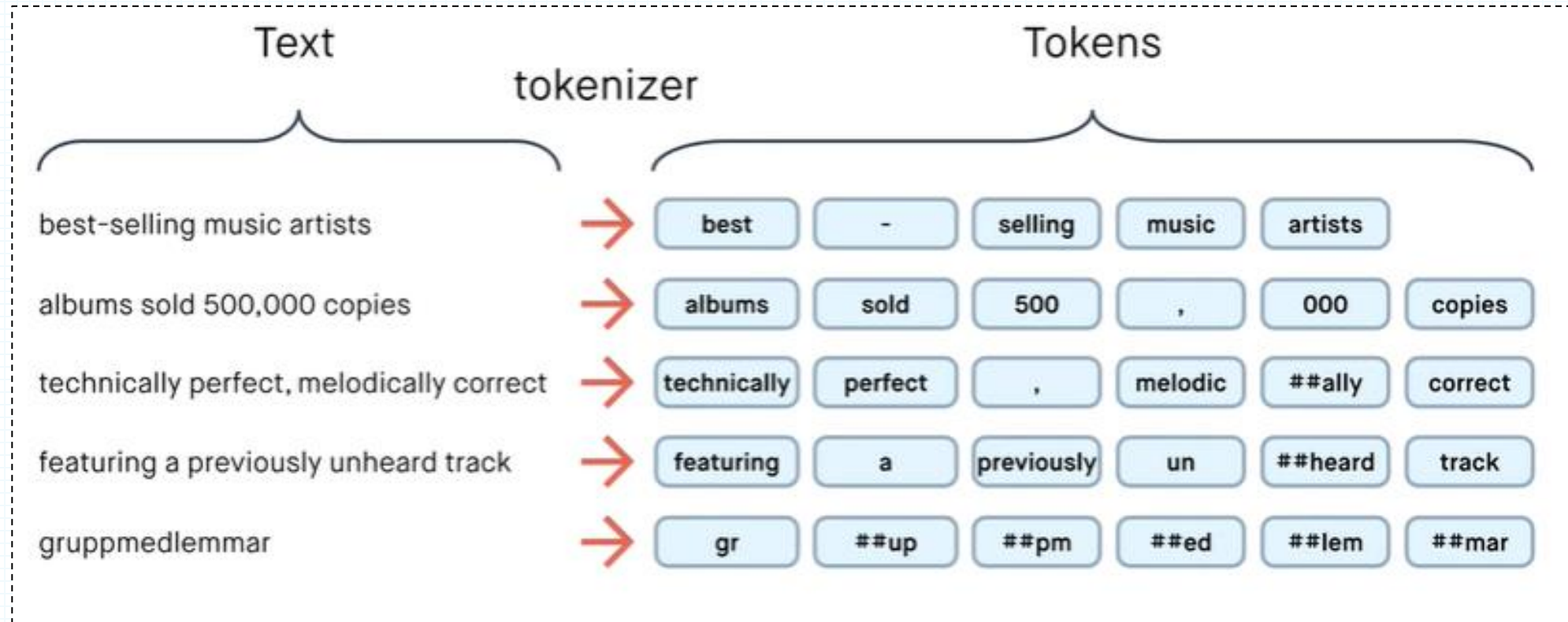
BERT: Pre-training of Deep Bidirectional Transformers for  
Language Understanding

2018

Jacob Devlin Ming-Wei Chang Kenton Lee Kristina Toutanova  
Google AI Language  
{jacobdevlin, mingweichang, kentonl, kristout}@google.com



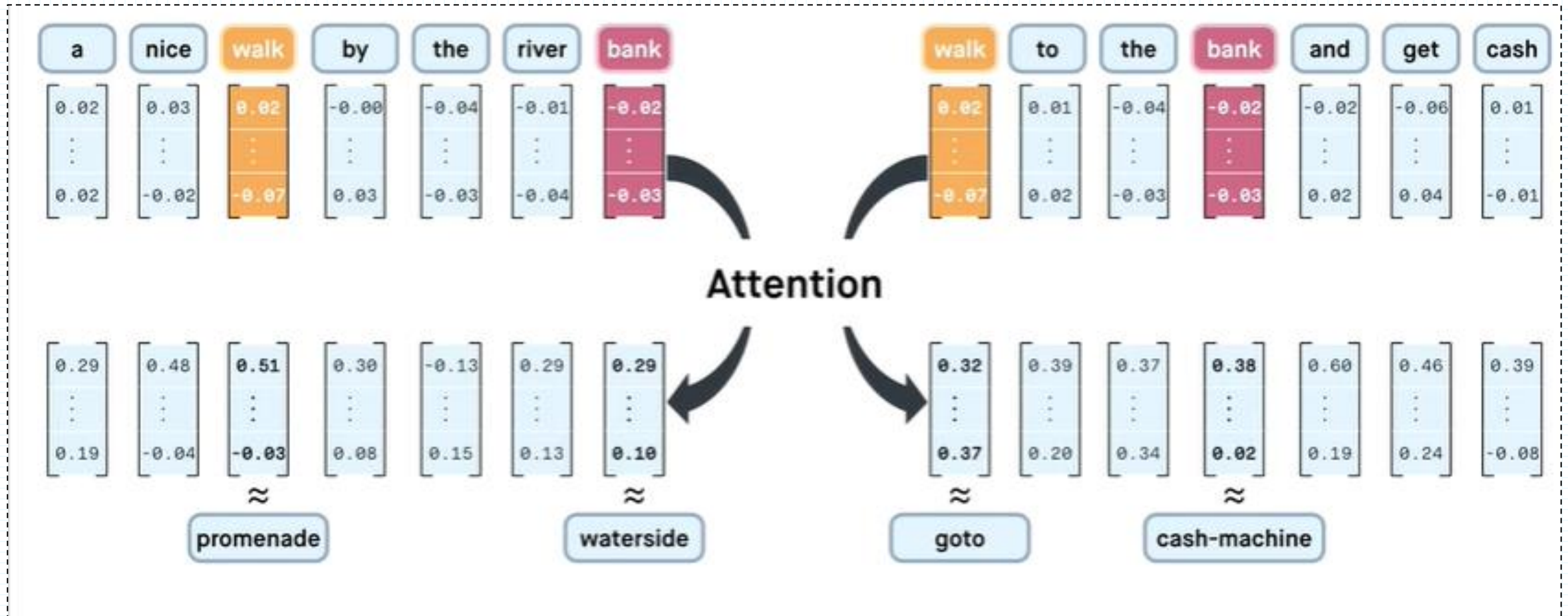
# Tokenization

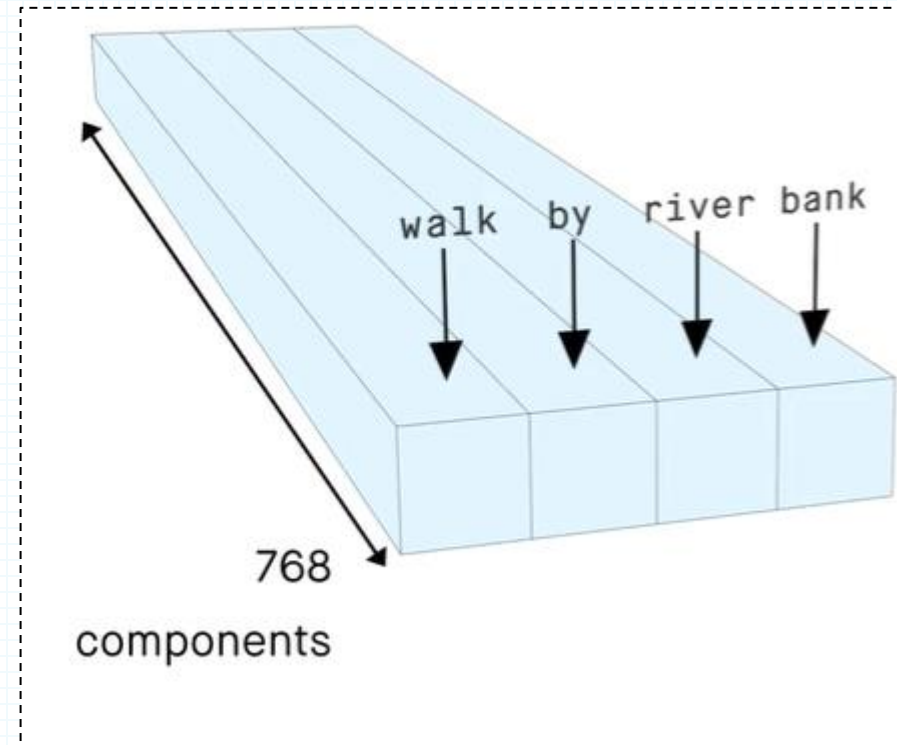
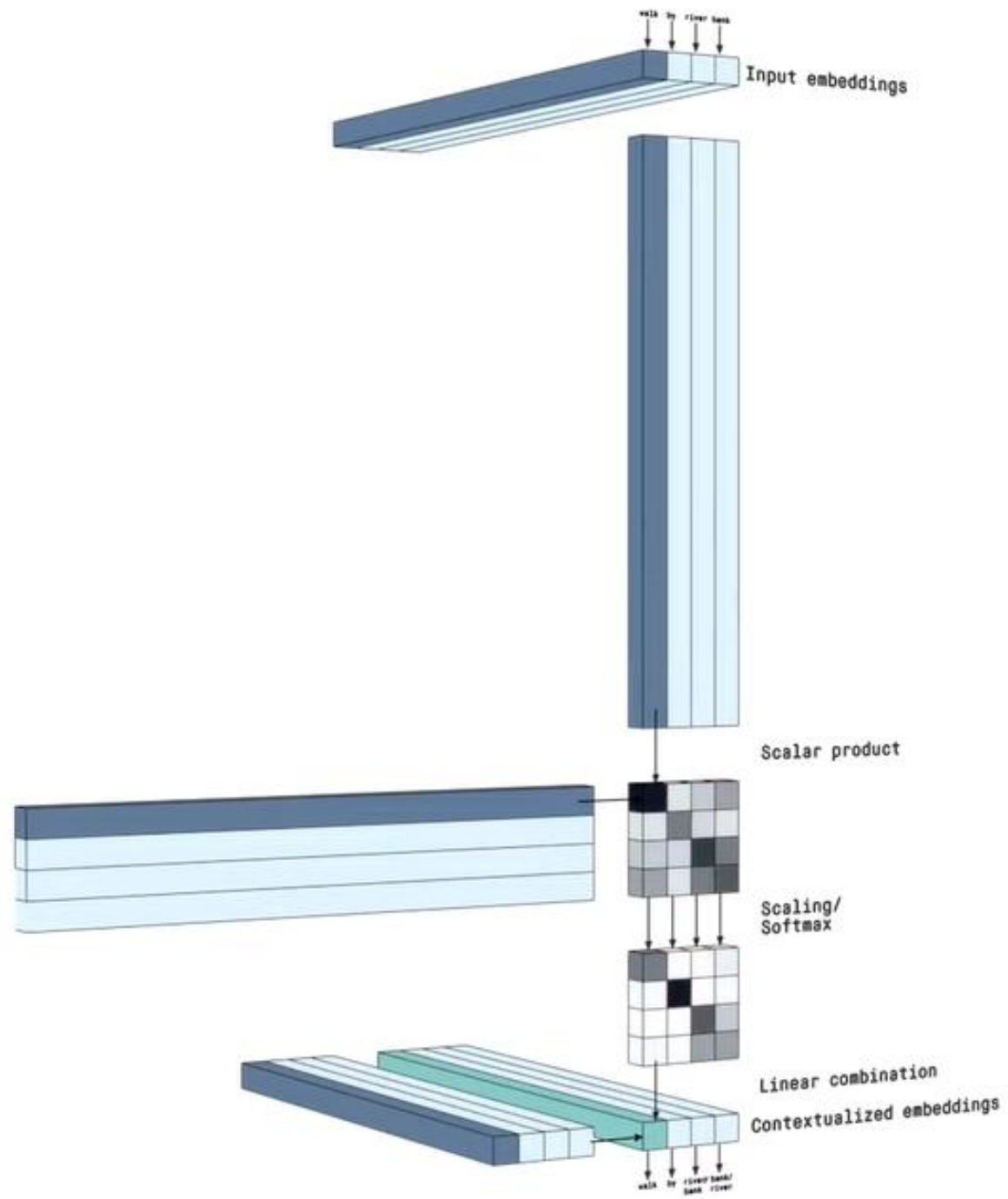


# Initial Embedding

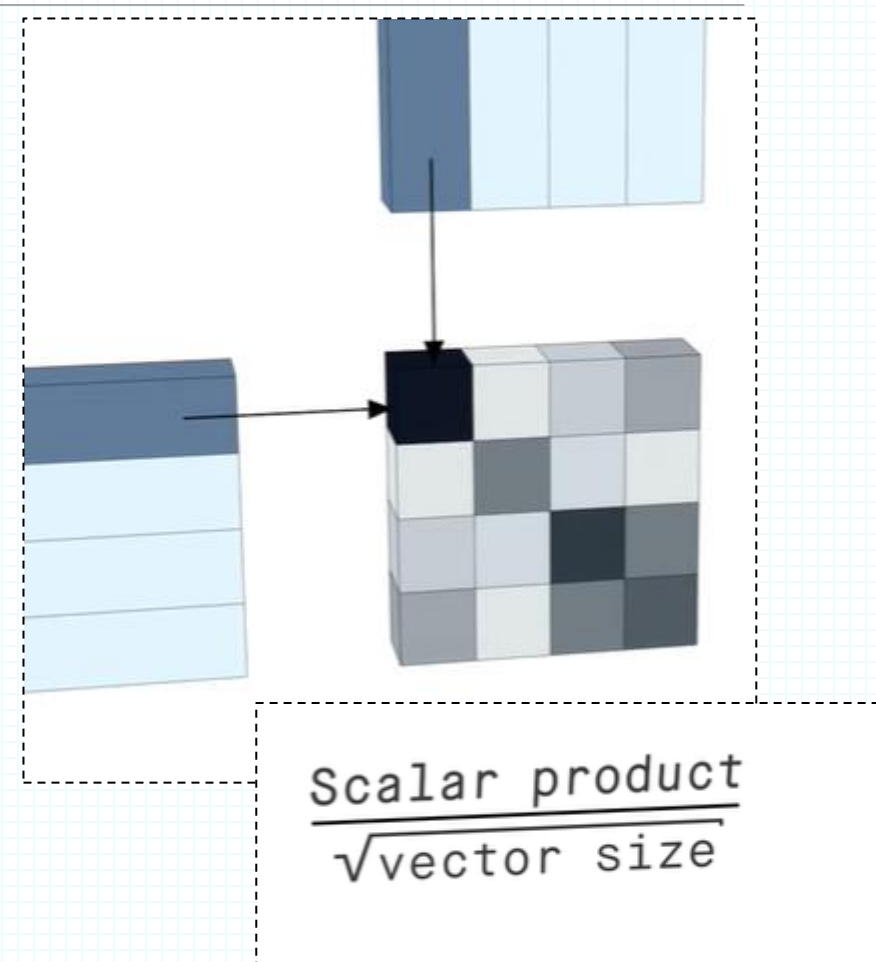
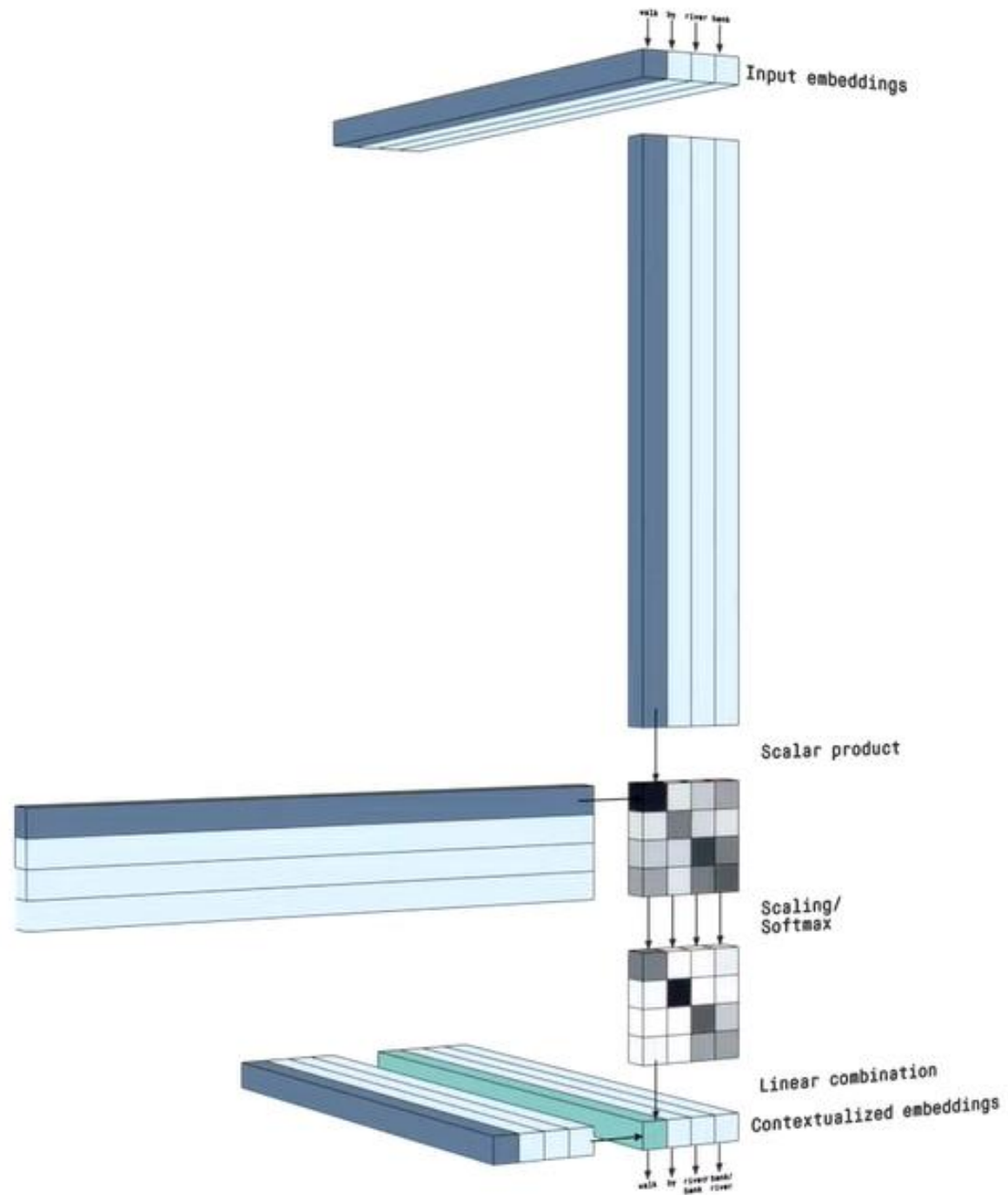


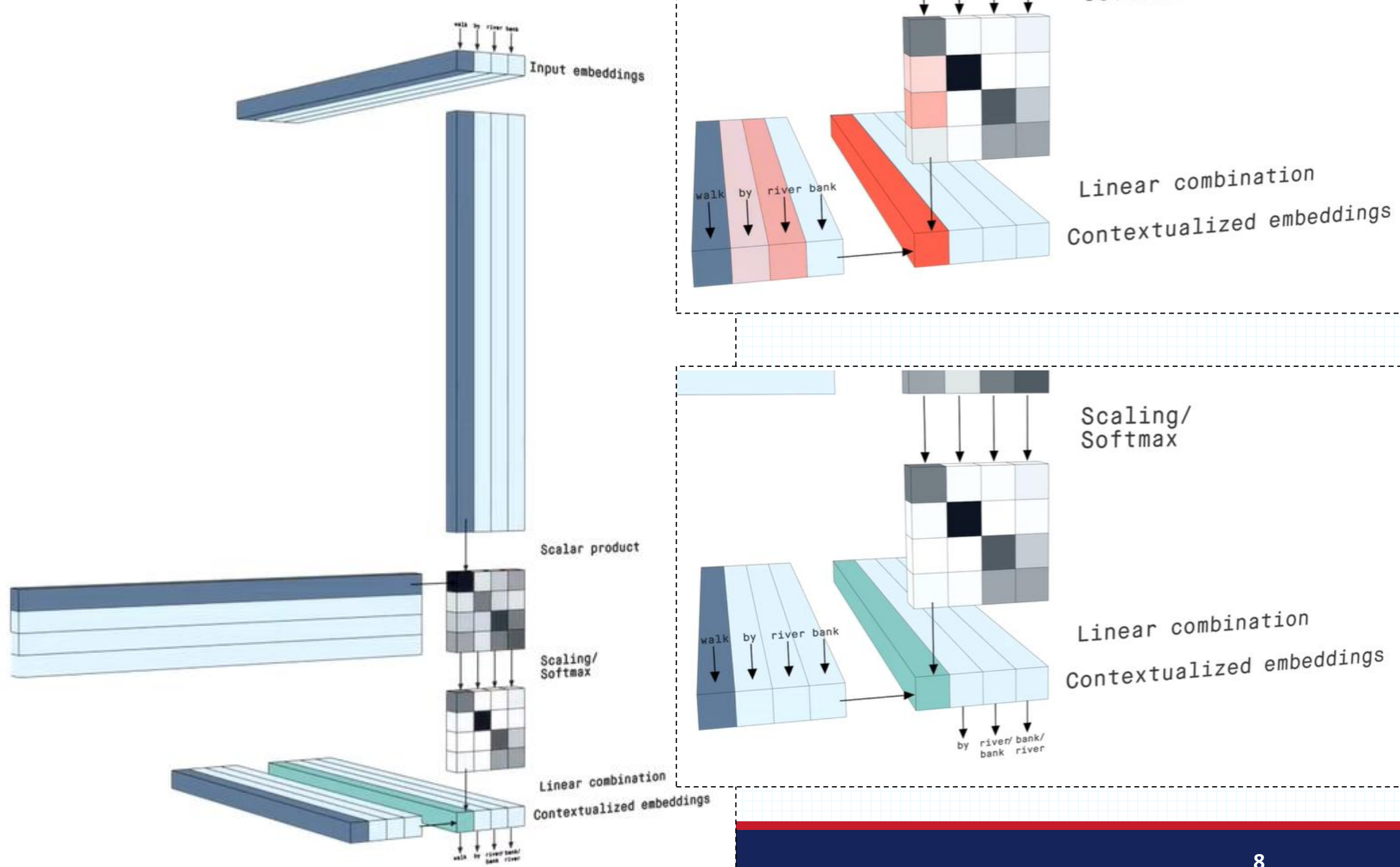
# Contextual Embedding





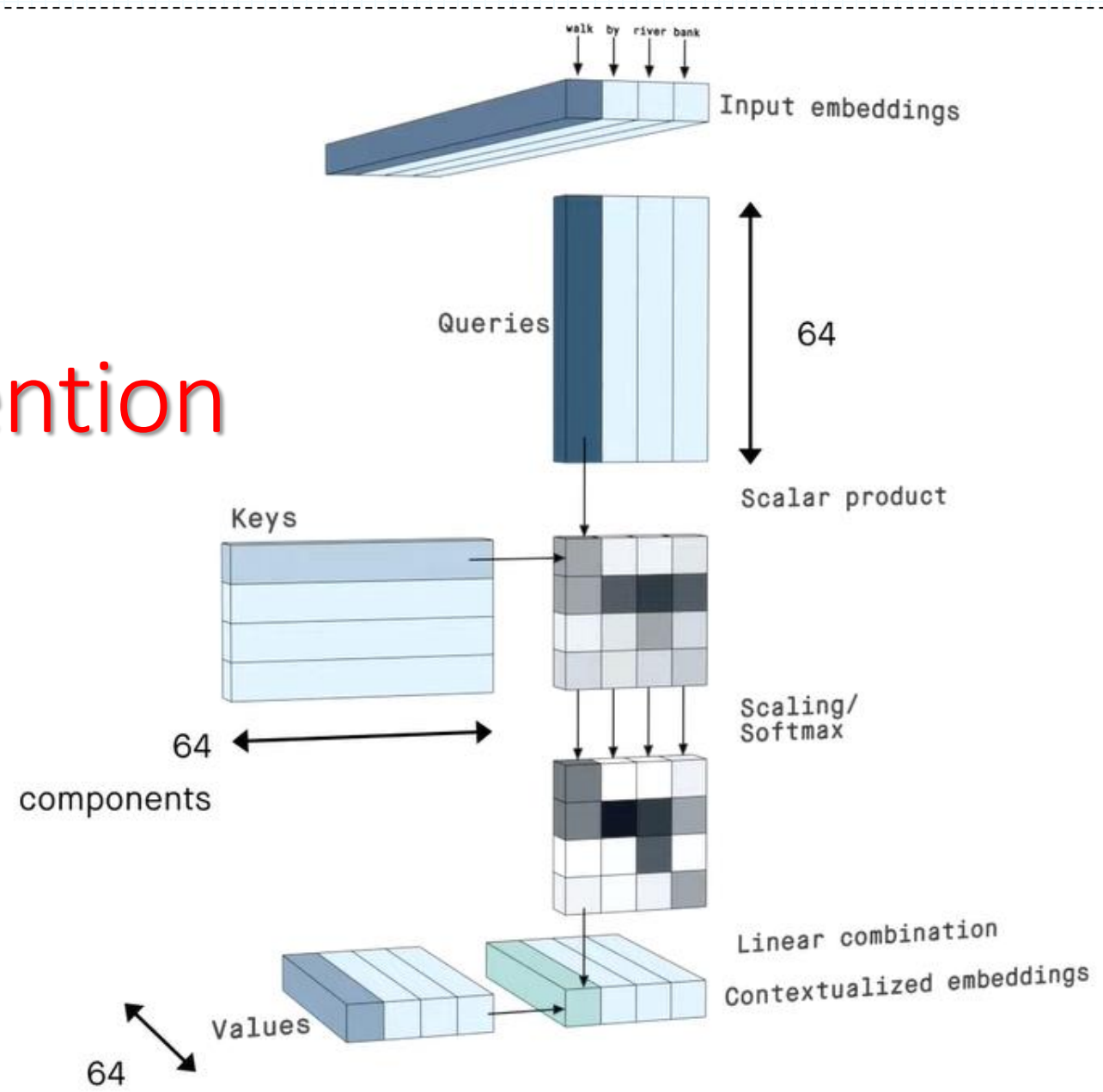
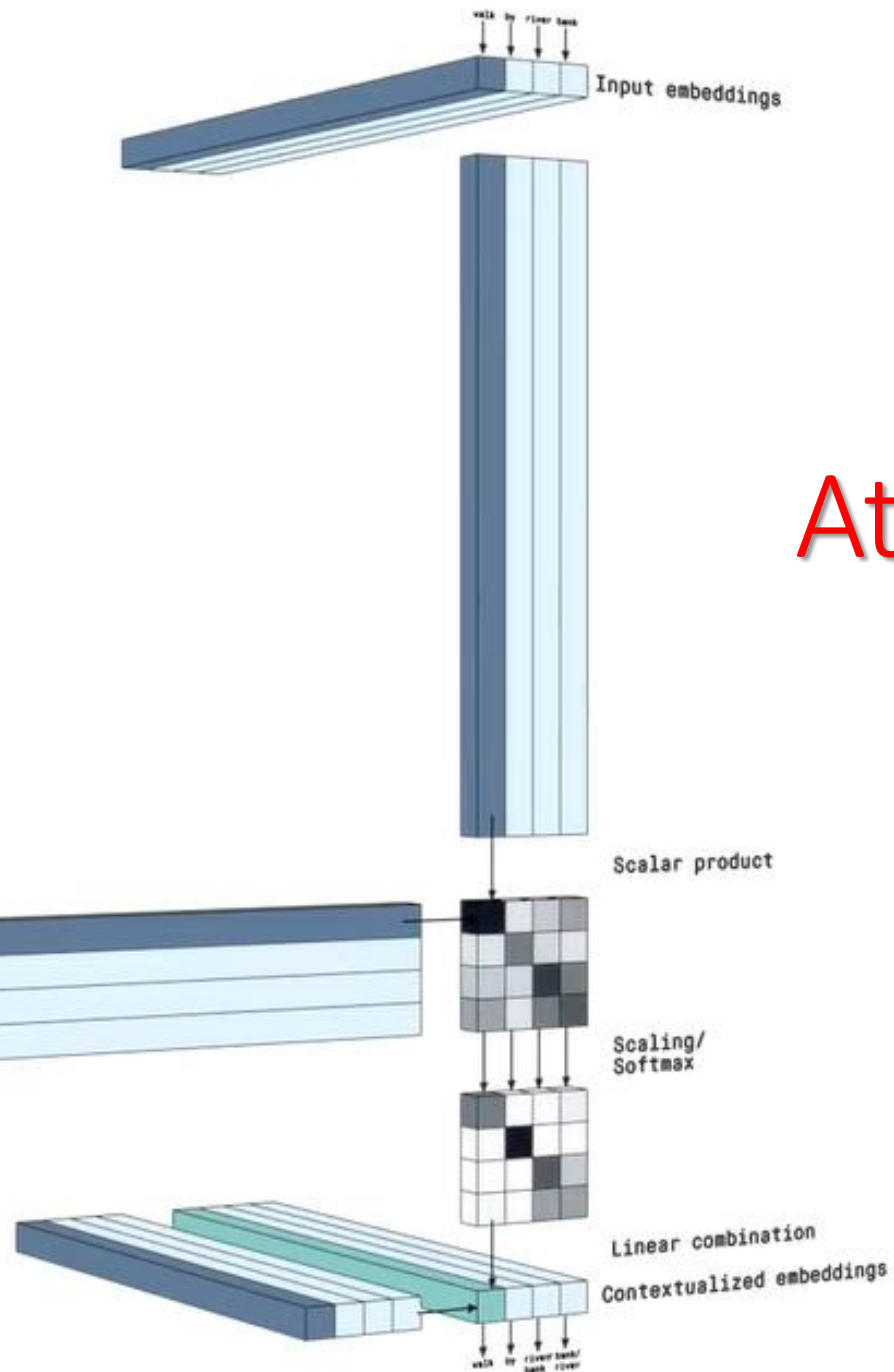




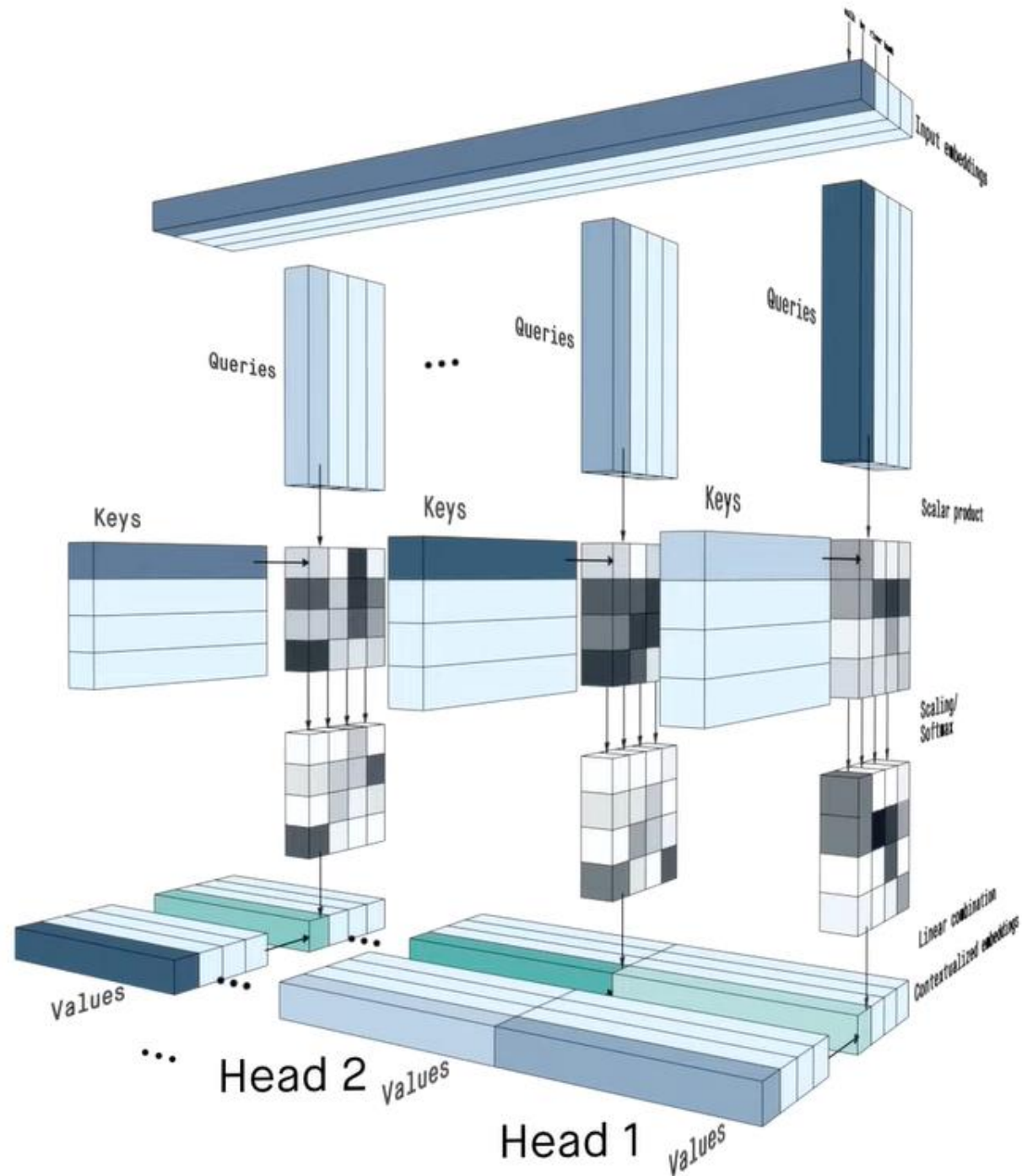




# Attention



# Multi-head attention



# Enrich the input!

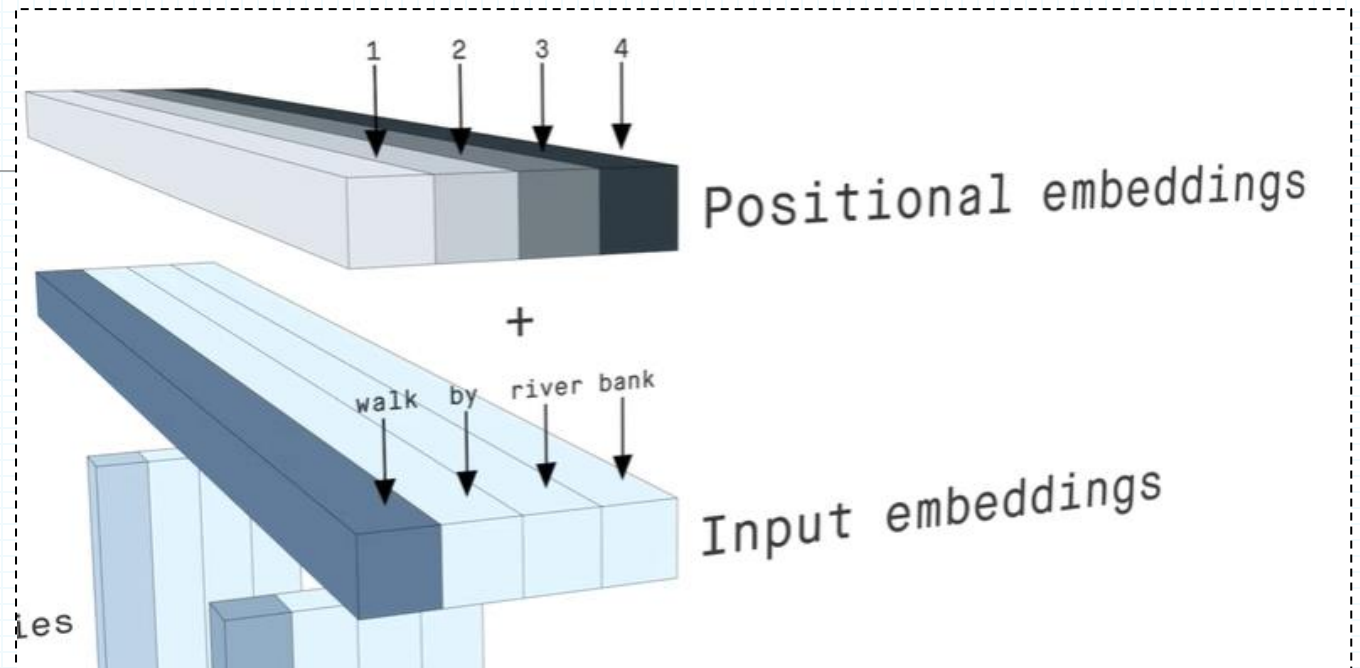
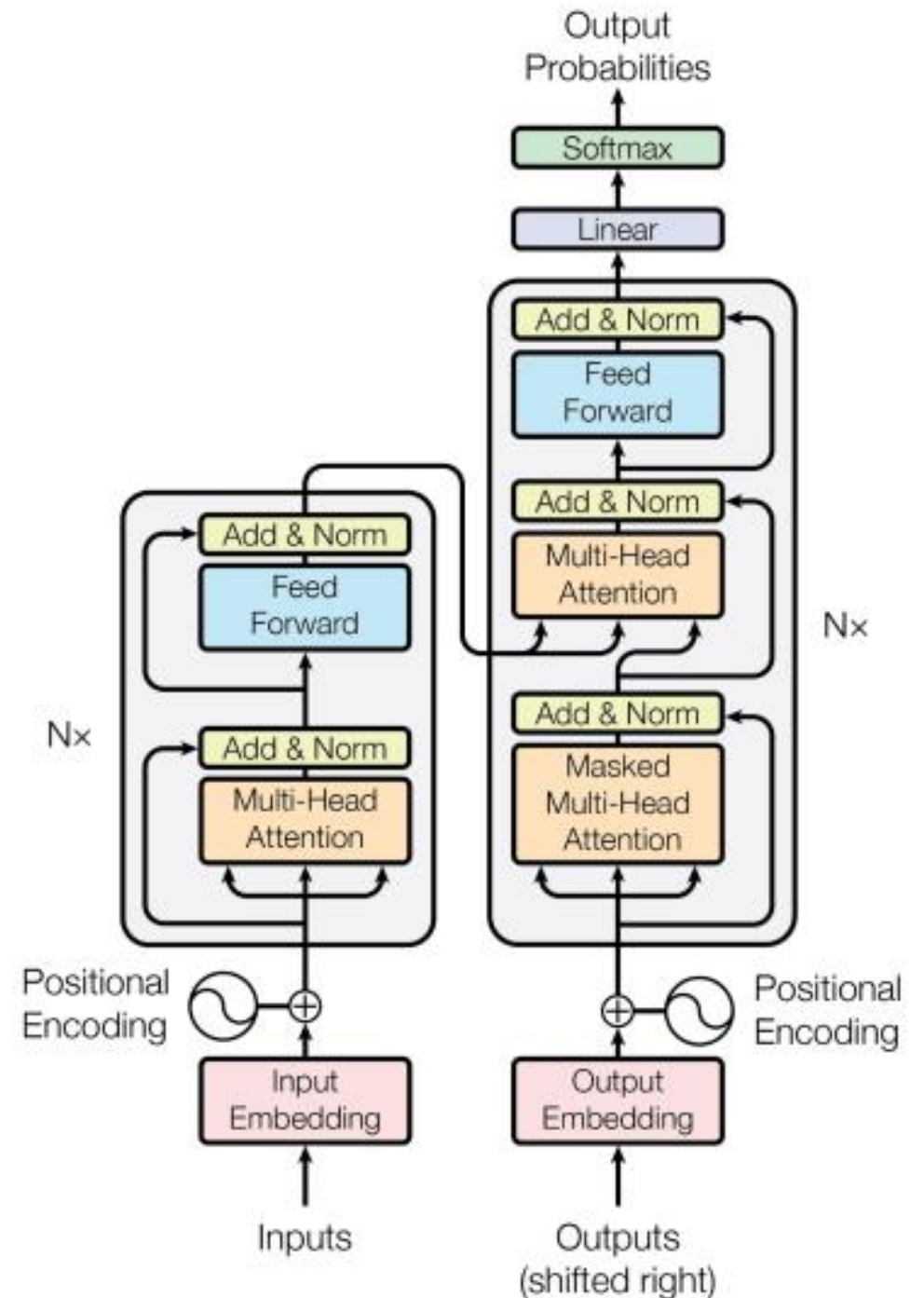


Figure 2: BERT input representation. The input embeddings are the sum of the token embeddings, the segmentation embeddings and the position embeddings.

So.. What are Transformers in general!

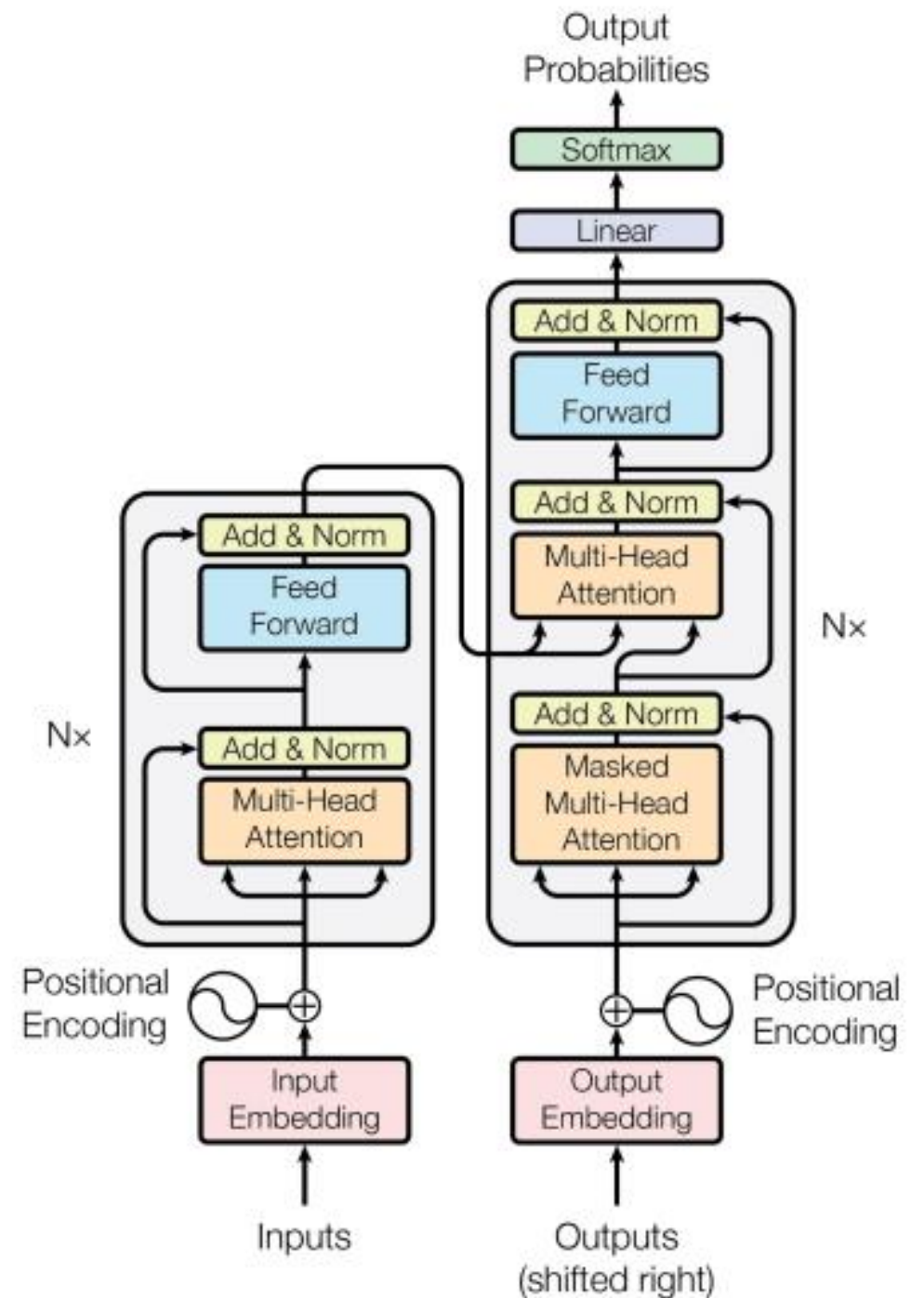
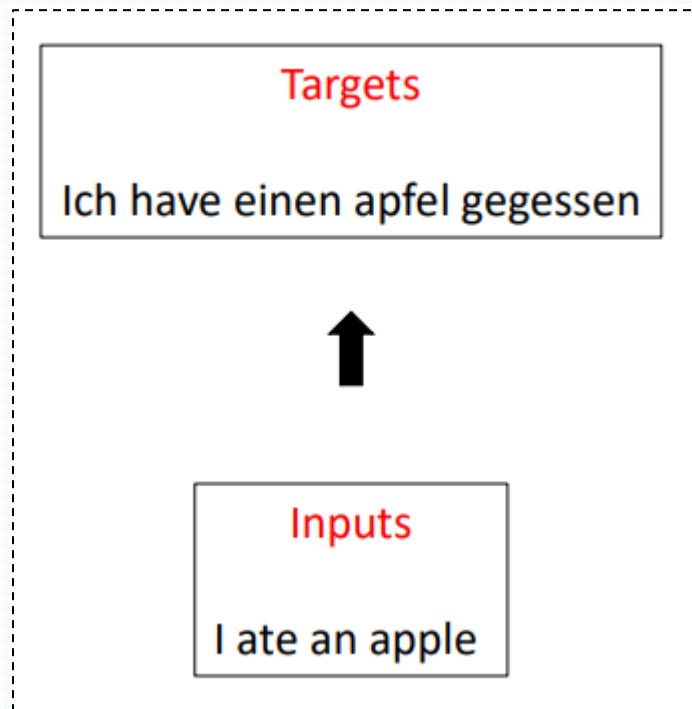
# Transformers

- Transformers are a type of neural network architecture that transforms or changes an input sequence into an output sequence.
- They do this by learning context and tracking relationships between sequence components.
- And break the problem into two parts:
  - An encoder (e.g., Bert)
  - A decoder (e.g., GPT)



# Transformers

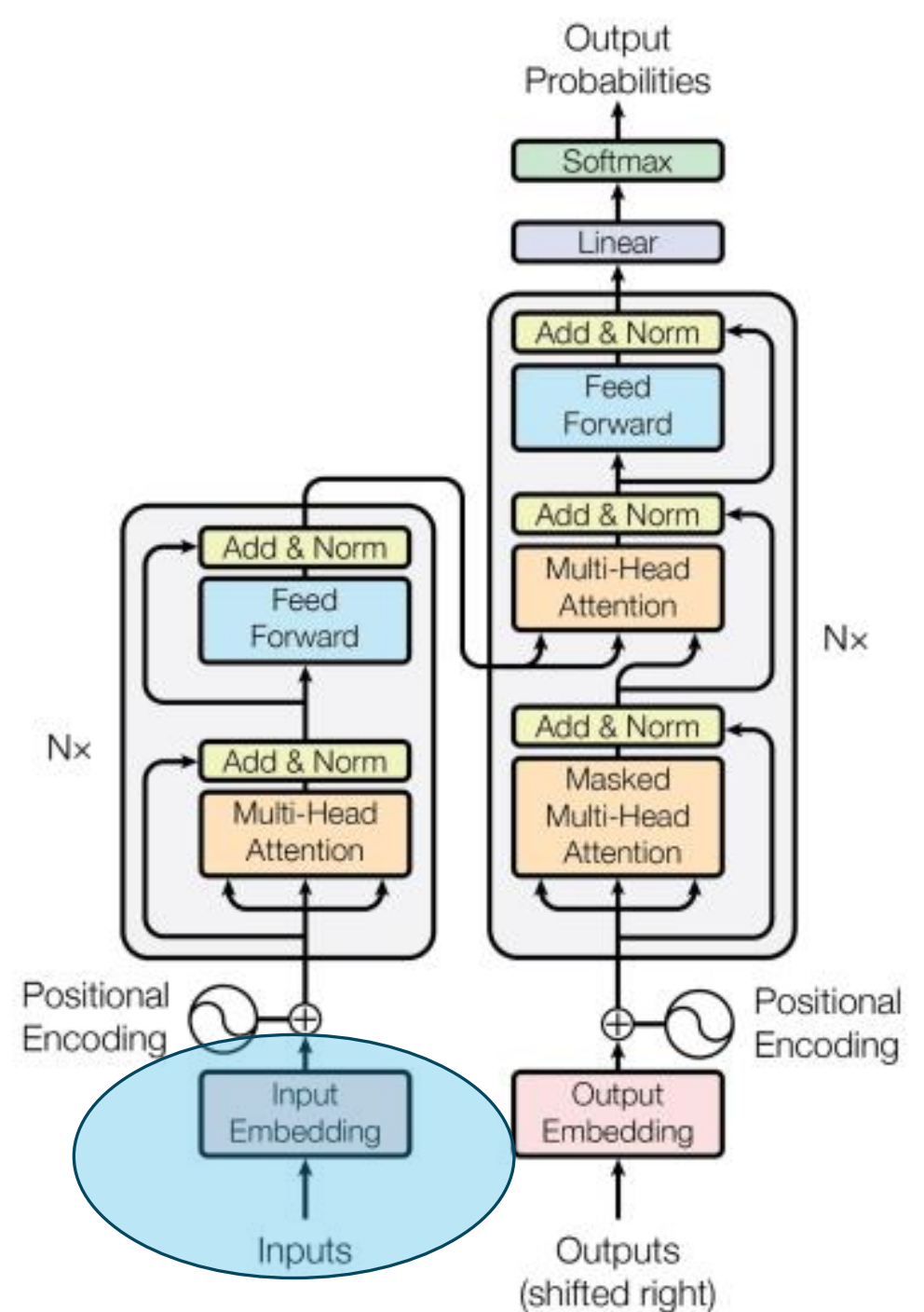
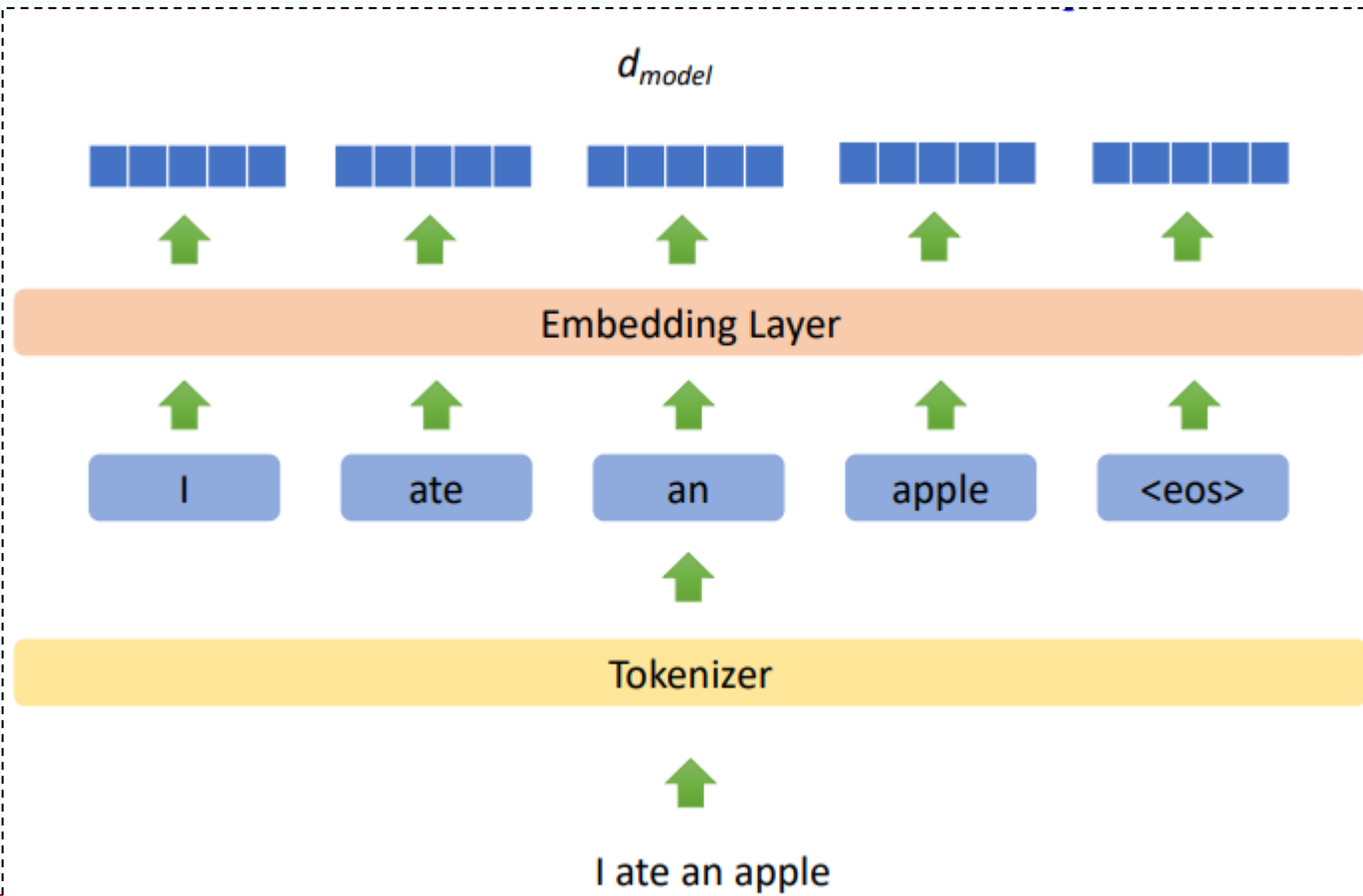
- Example: Machine Translation





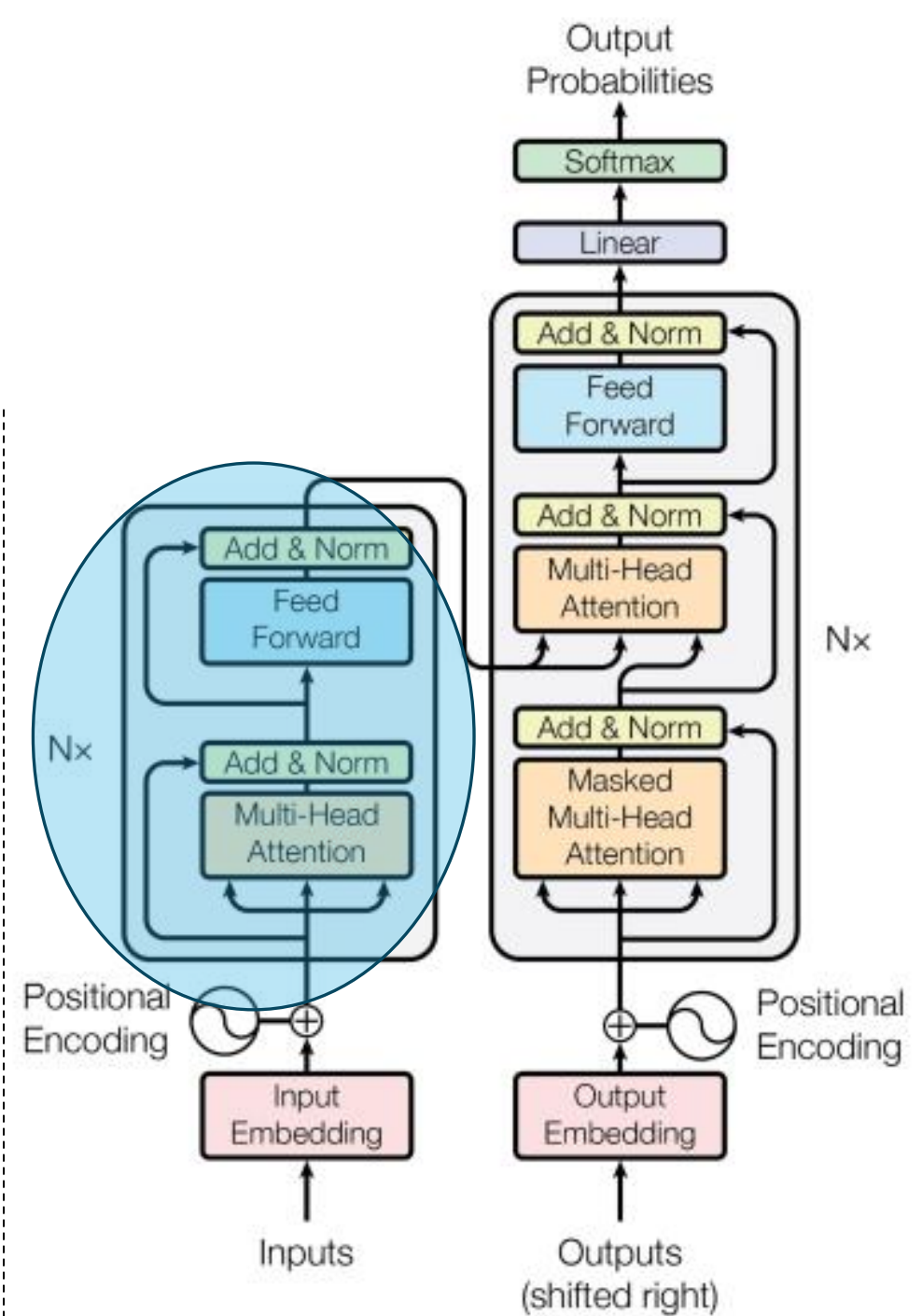
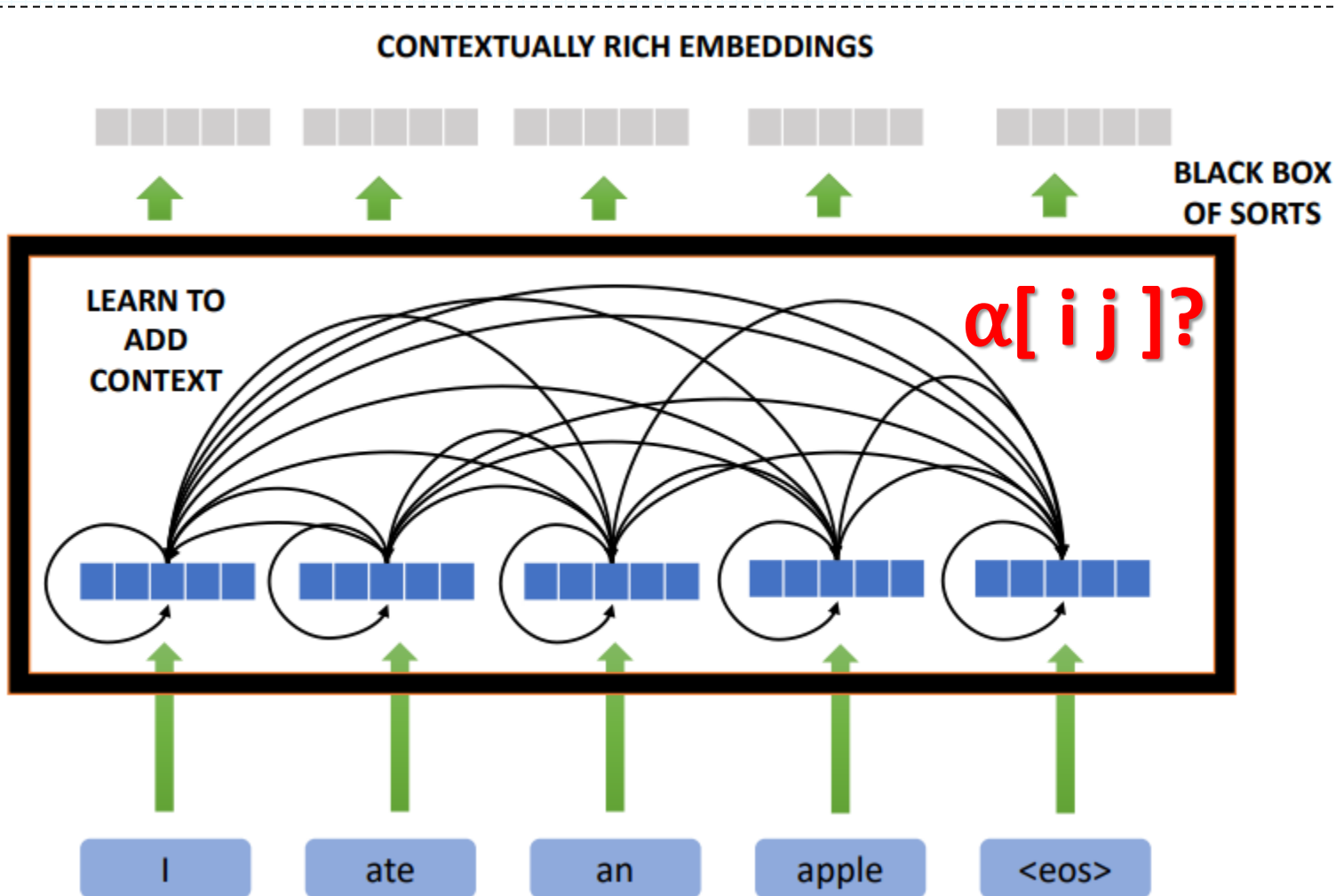
# Transformers

- Processing Input



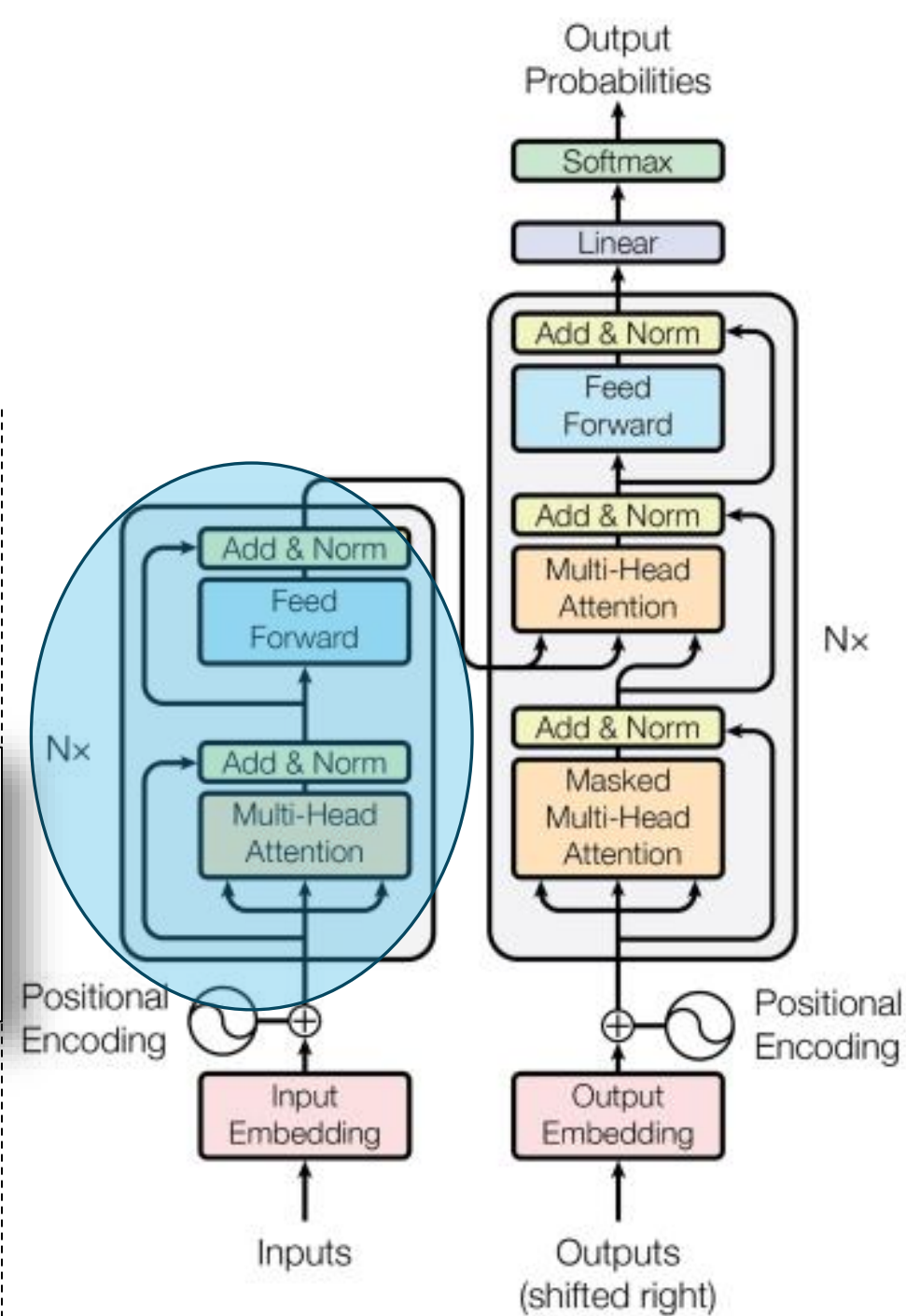
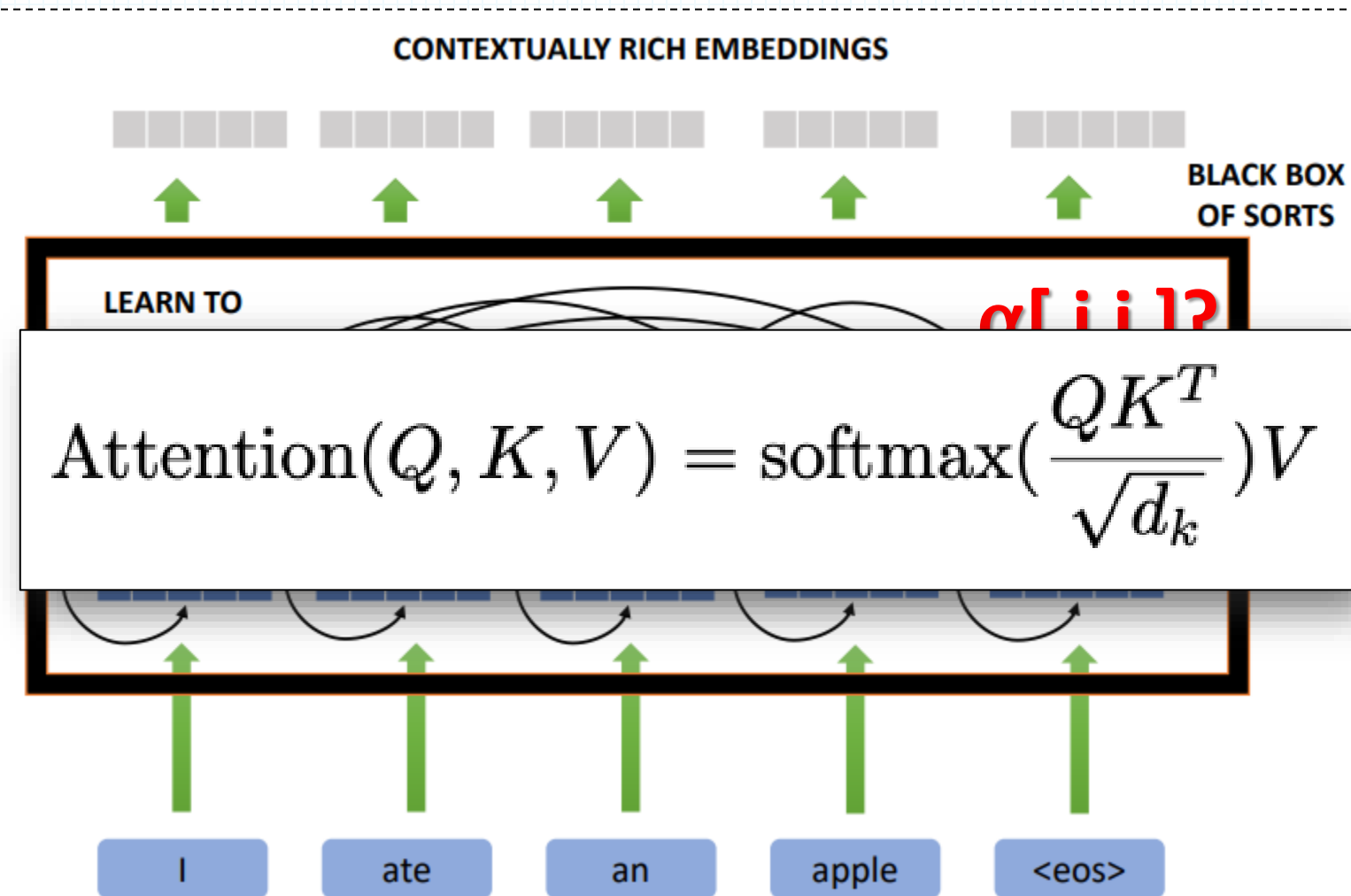
# Transformers

- Learn to add context



# Transformers

- Learn to add context



# Attention Concepts

- In the attention mechanism, we can draw an analogy to Information Retrieval (IR).

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

OR

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

A **query** (an embedding vector representing what we want to find more about — e.g., a specific token or position in a sequence)

A set of **keys** (representing the indexed or stored information — i.e., all other tokens in the context)

{Key, Value store}

```
{
  "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",
    ...
  }
}
```

And a set of **values** (the actual content or information associated with each key).



# Attention Concepts

A set of **keys** (representing the index of the tokens in the input sequence)

- In the attention mechanism, we use an analogy to Information Retrieval

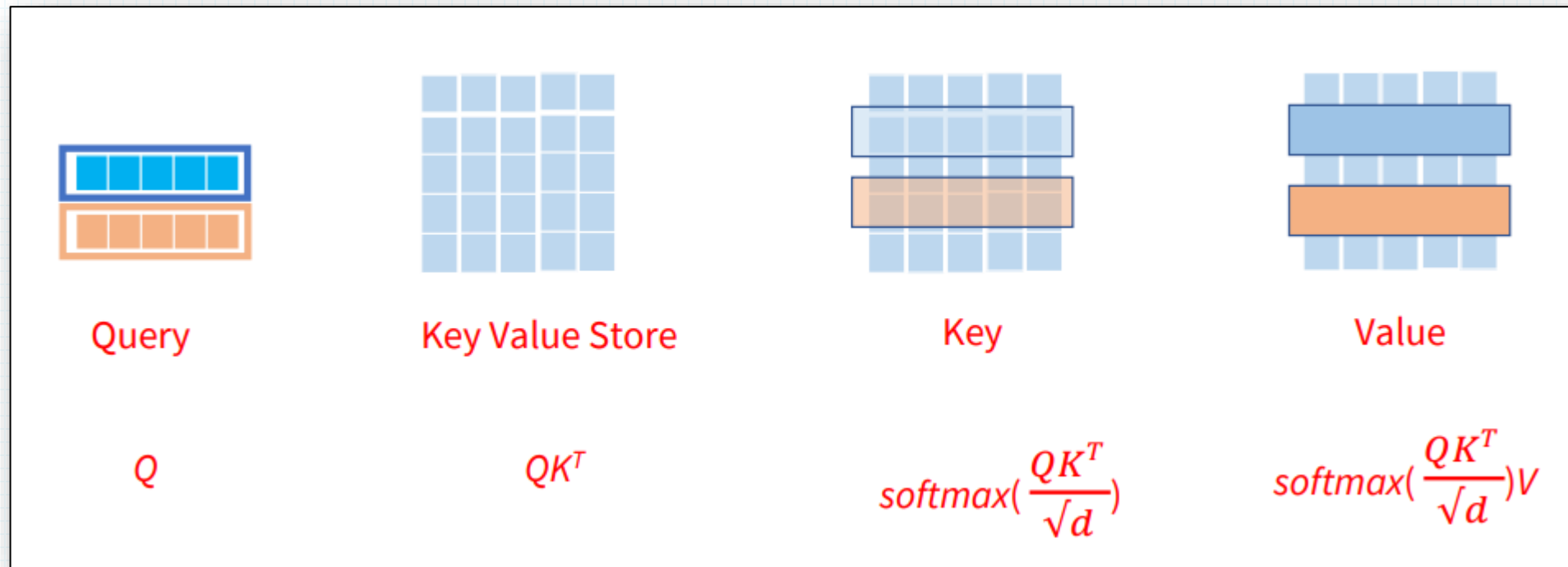
**Attention allows the model to dynamically retrieve relevant information (values) from the context, based on how similar other tokens (keys) are to the current focus (query).**

**The attention process works as follows:**

1. We compute the similarity between the query and each key (usually using a dot product).
2. This gives us a set of **attention scores**, indicating how much focus we should give to each position in the sequence.
3. We then apply a **softmax to normalize the scores** into a probability distribution.
4. Finally, we **compute a weighted sum of the values** based on these scores to produce the attention output.

# Attention

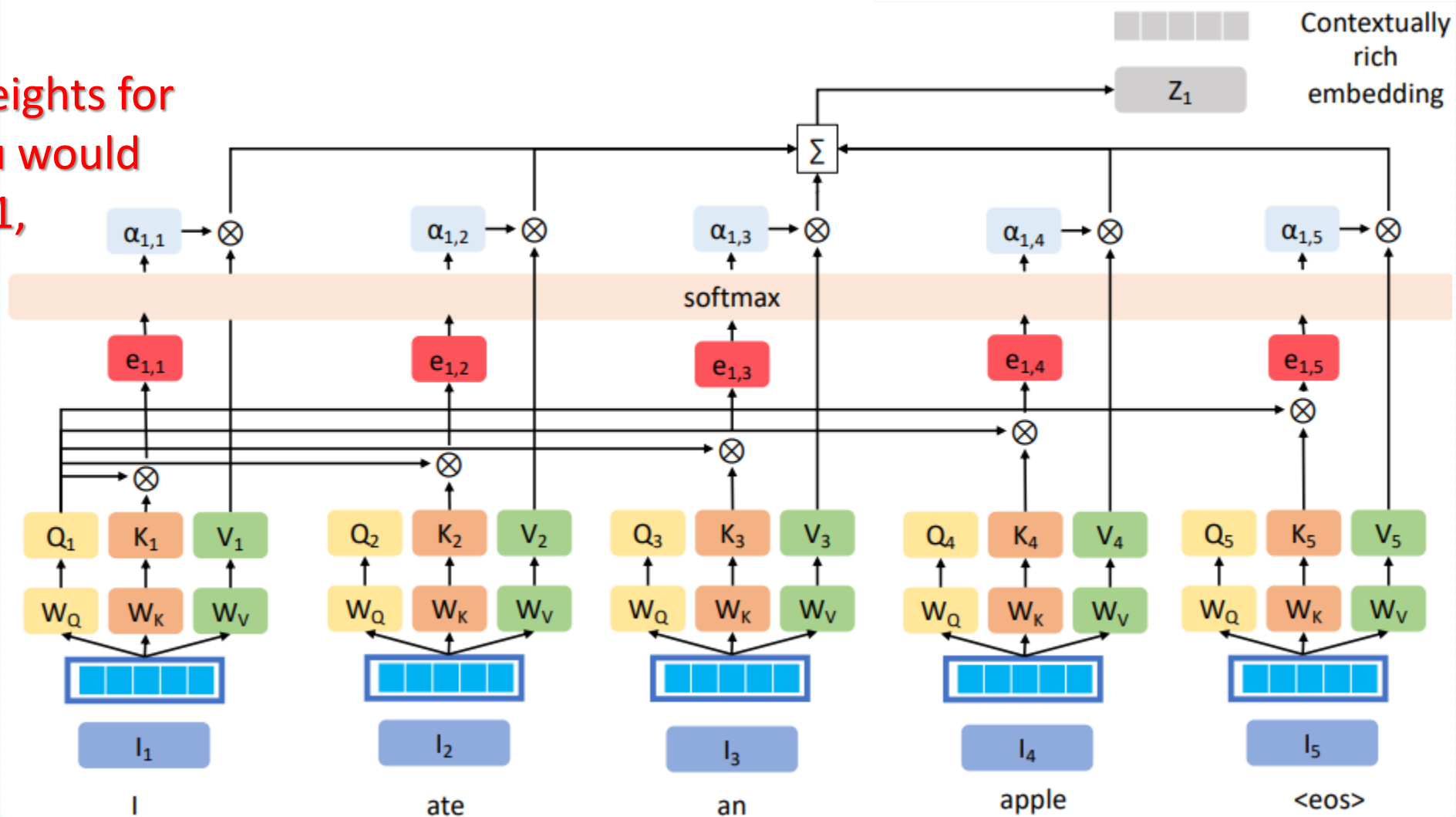
- More formal!





# Attention

To calculate attention weights for input  $I_1$ , you would use query  $q_1$ , and all keys

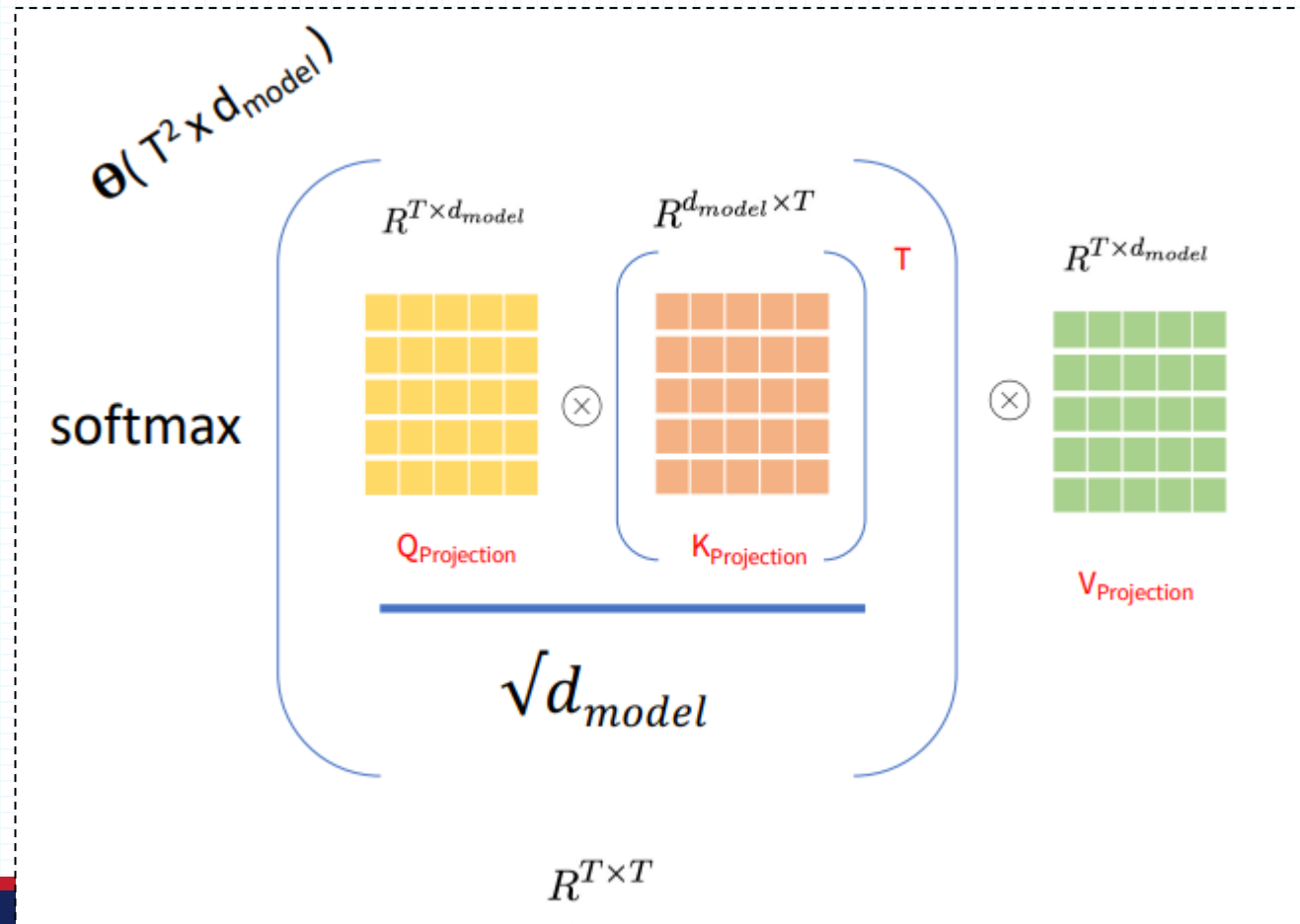


Of  $I_1$

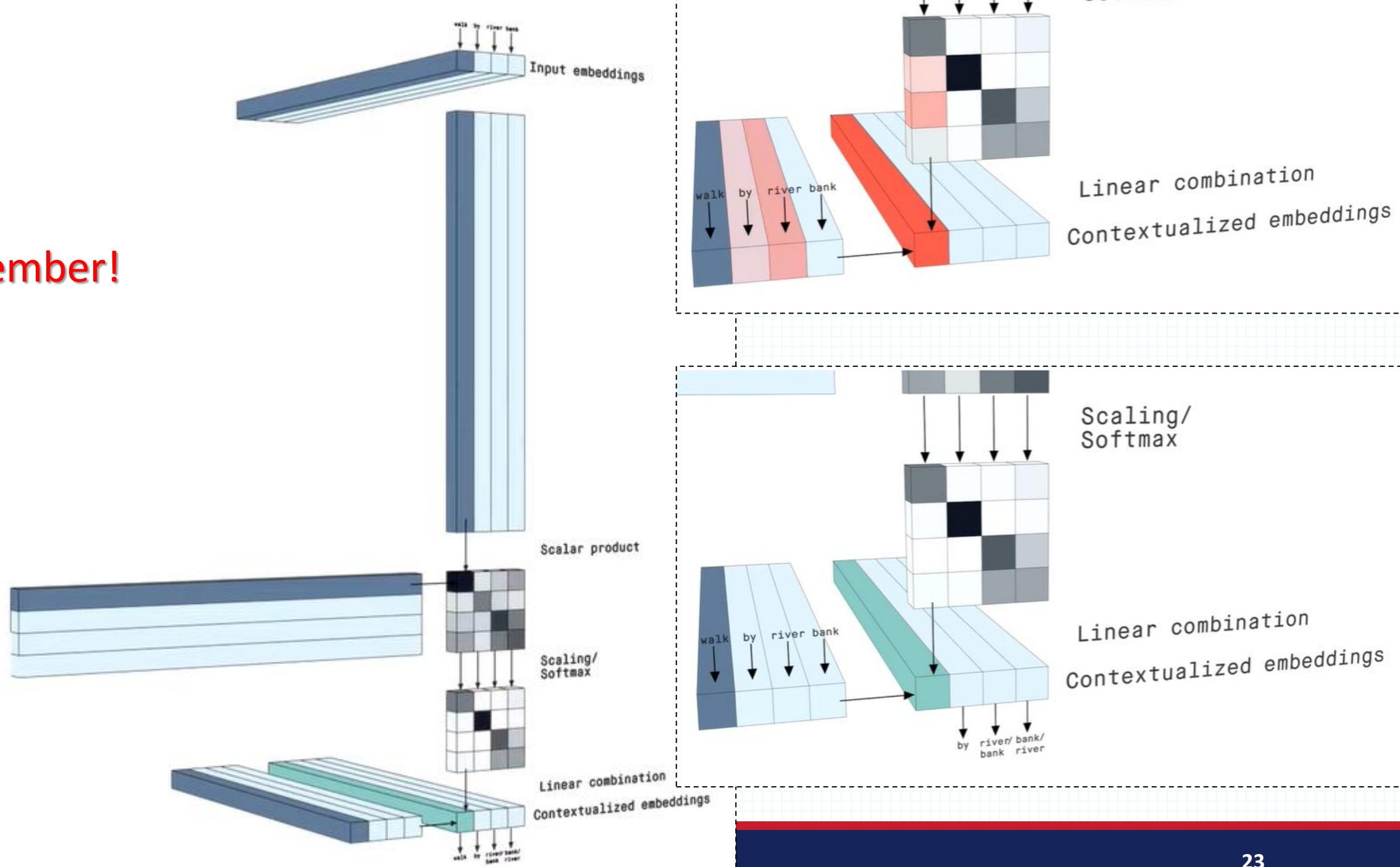
We do the same for  $I_2$ ,  $I_3$ , etc

# Self Attention

Query Inputs = Key Inputs = Value Inputs

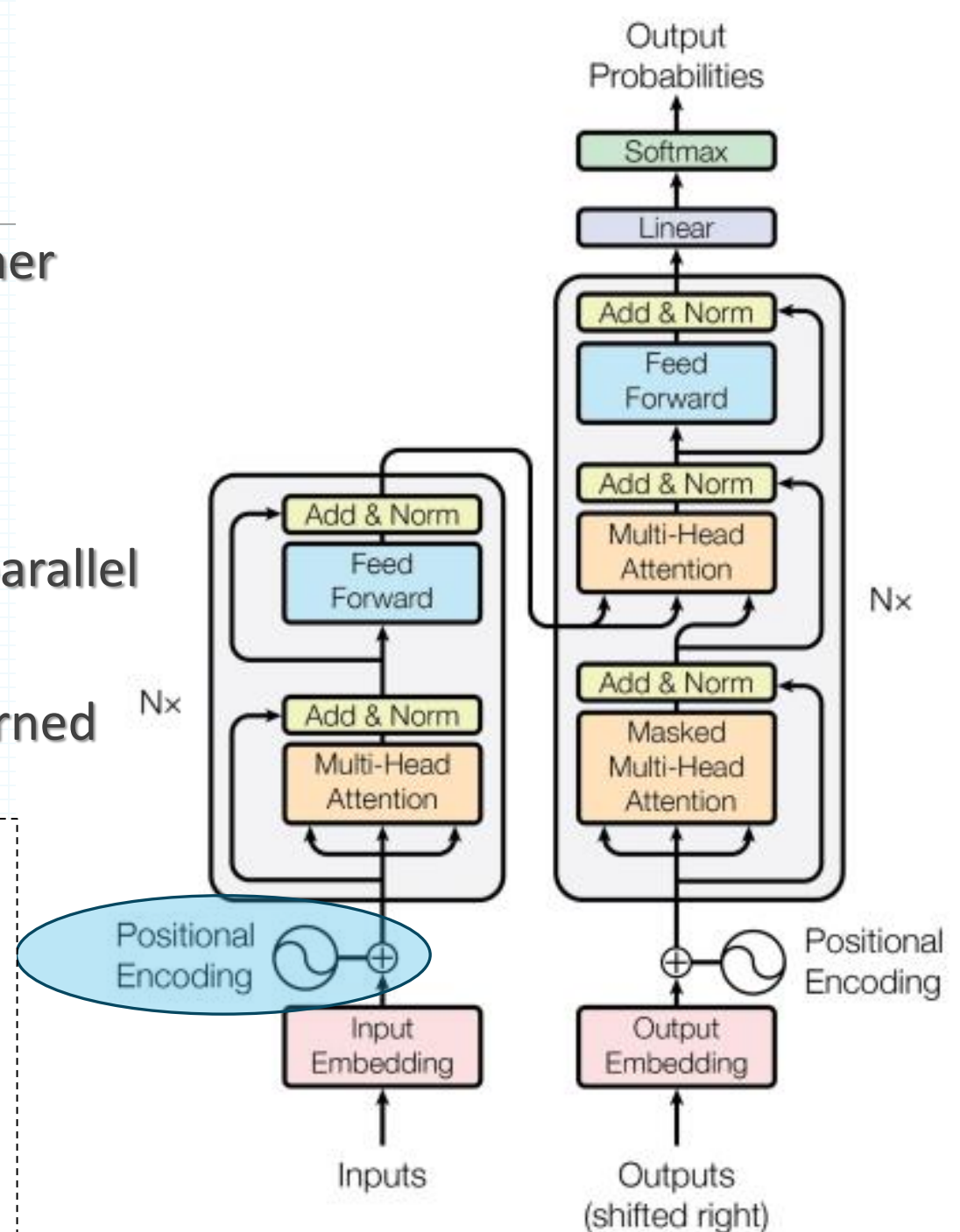
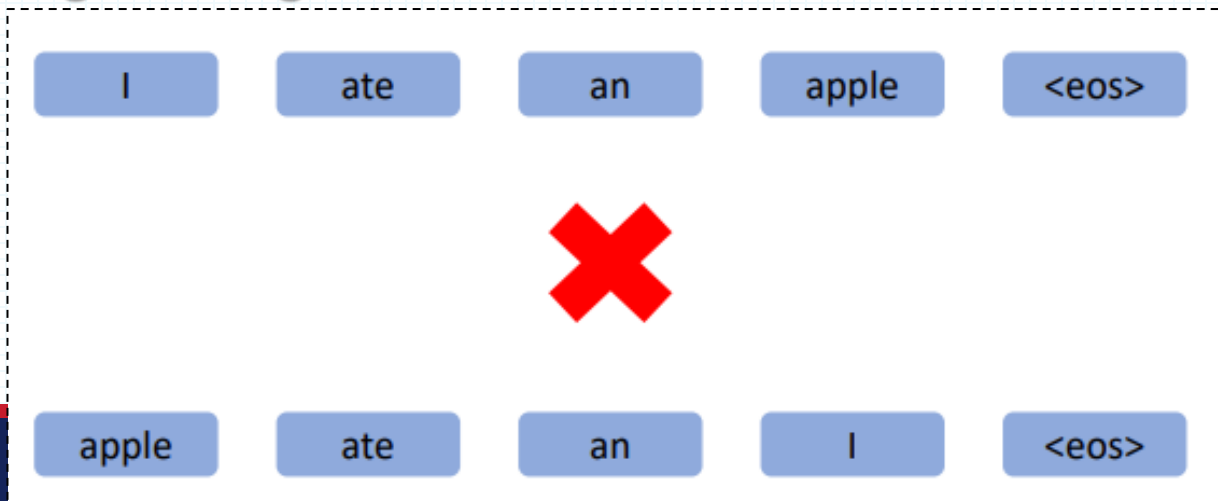


Remember!



# Positional Encoding

- Injects sequence order information into transformer models.
- Added to token embeddings to help the model understand word positions.
- Needed because transformers process tokens in parallel and lack inherent order sensitivity.
- Can be fixed (e.g., sinusoidal – sine, cosine) or learned during training.



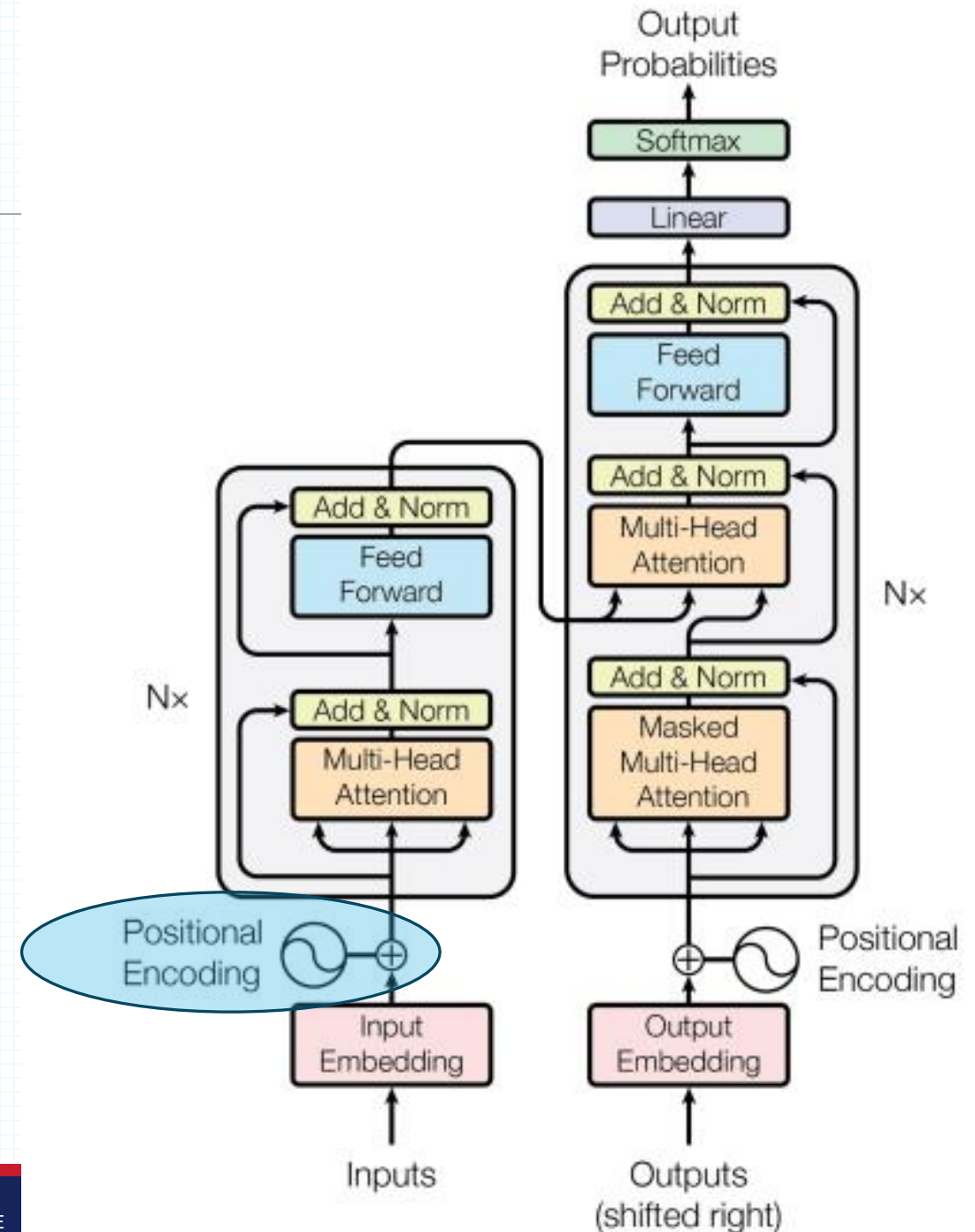
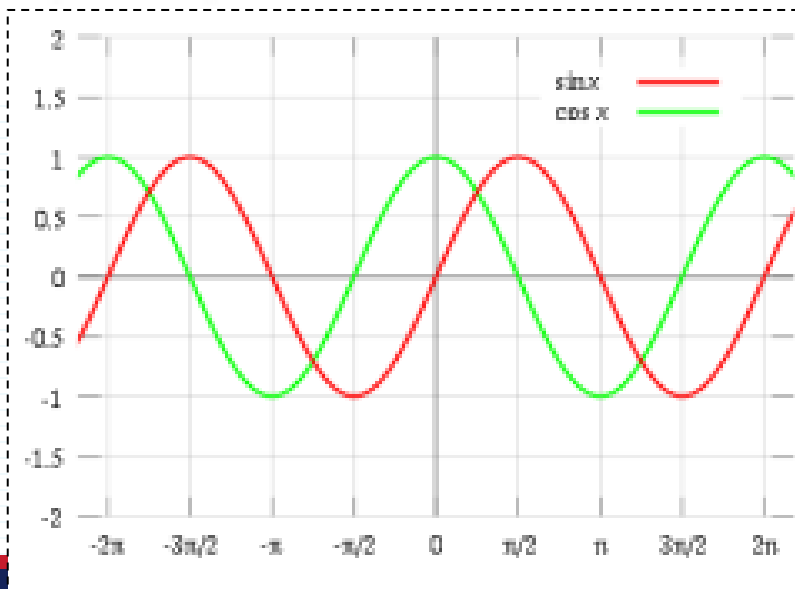
# Positional Encoding

pos -> idx of the token in input sentence

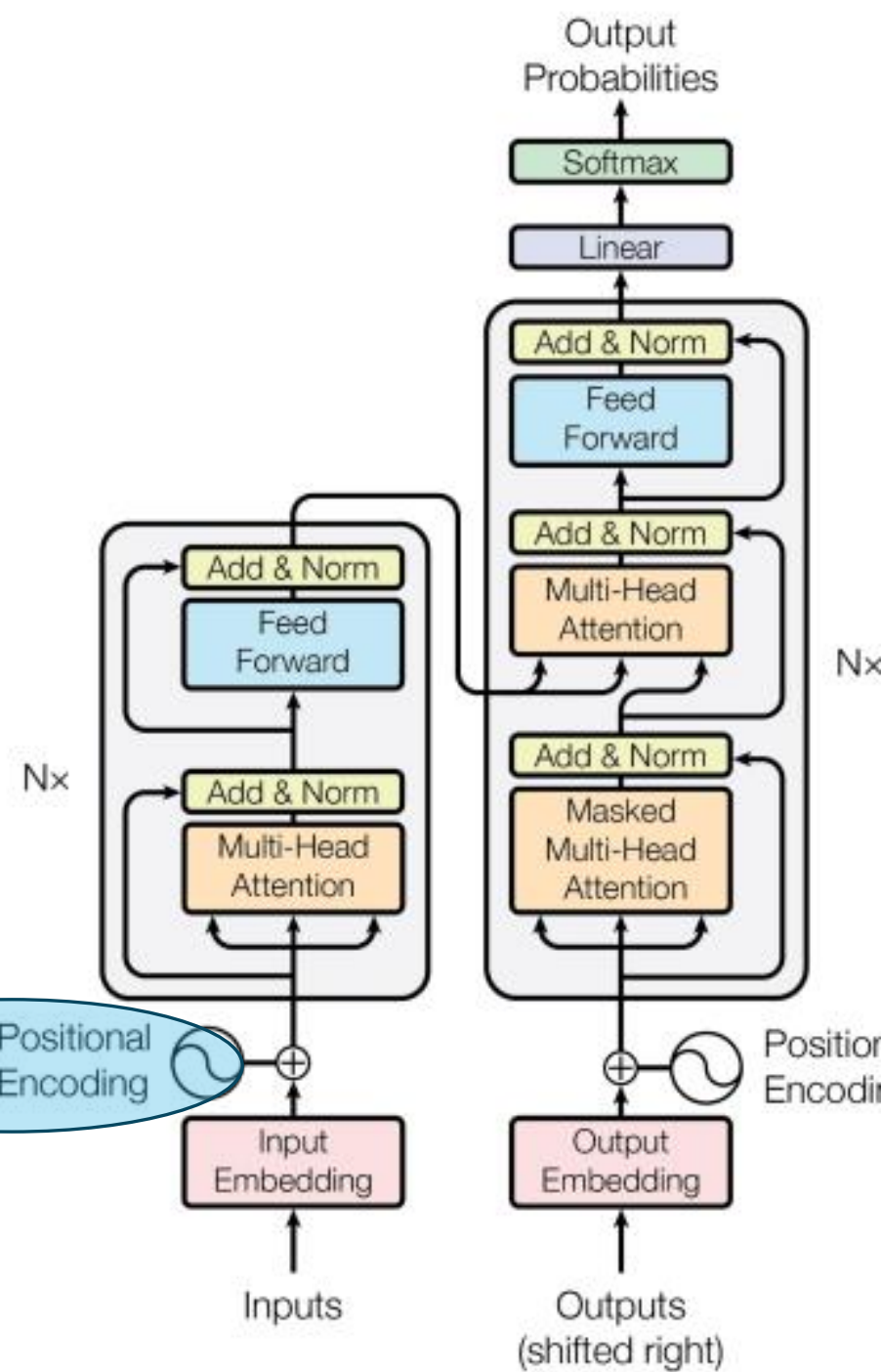
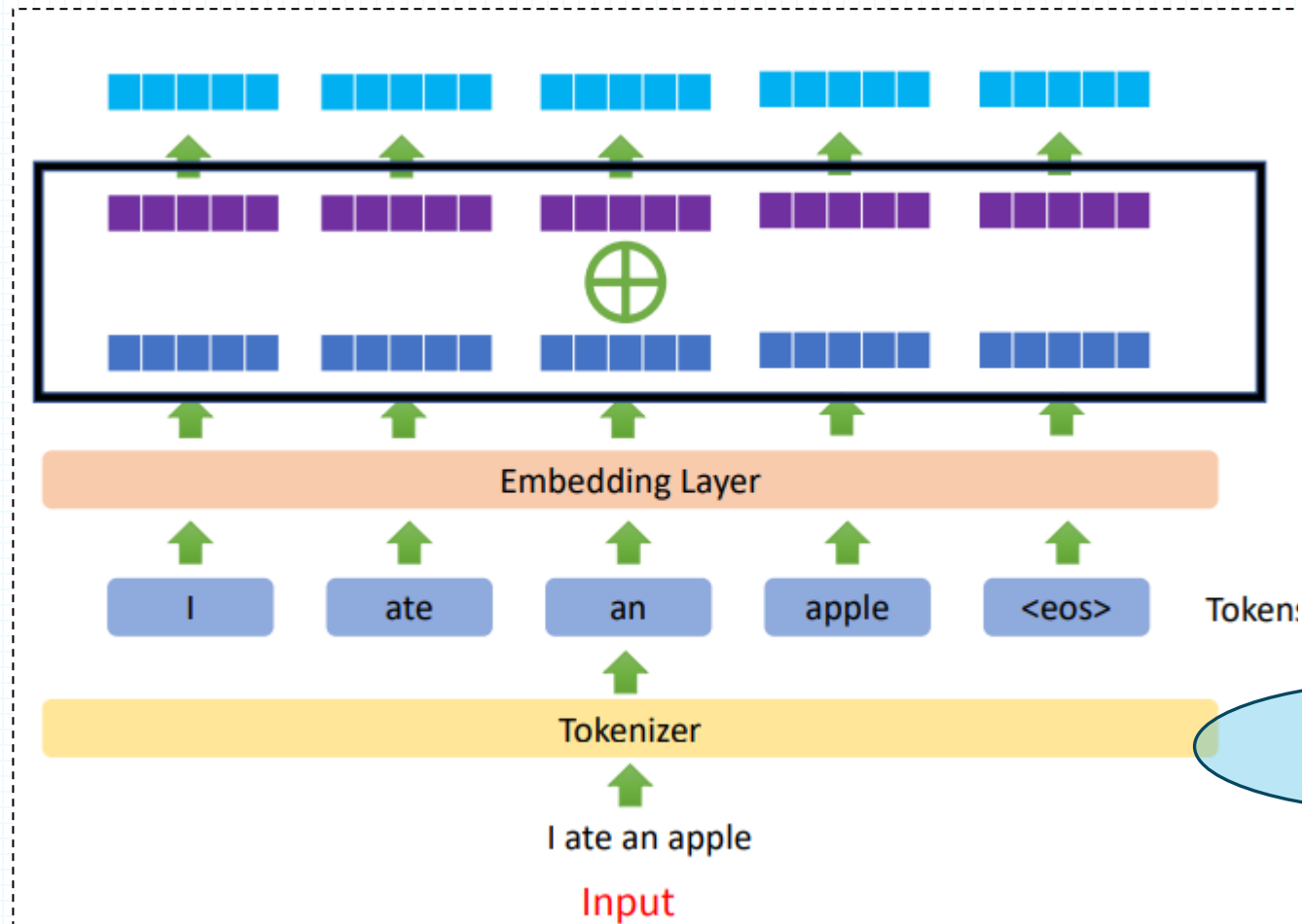
i ->  $i^{\text{th}}$  dimension out of d

$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{model}})$$

$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{model}})$$



# Position Embedding





# Next Lecture

