

# Language Model: Attention and Transformers

- When she tried to print her tickets, she found that the printer was out of toner. She went to the stationery store to buy more toner. It was very overpriced. After installing the toner into the printer, she finally printed her 

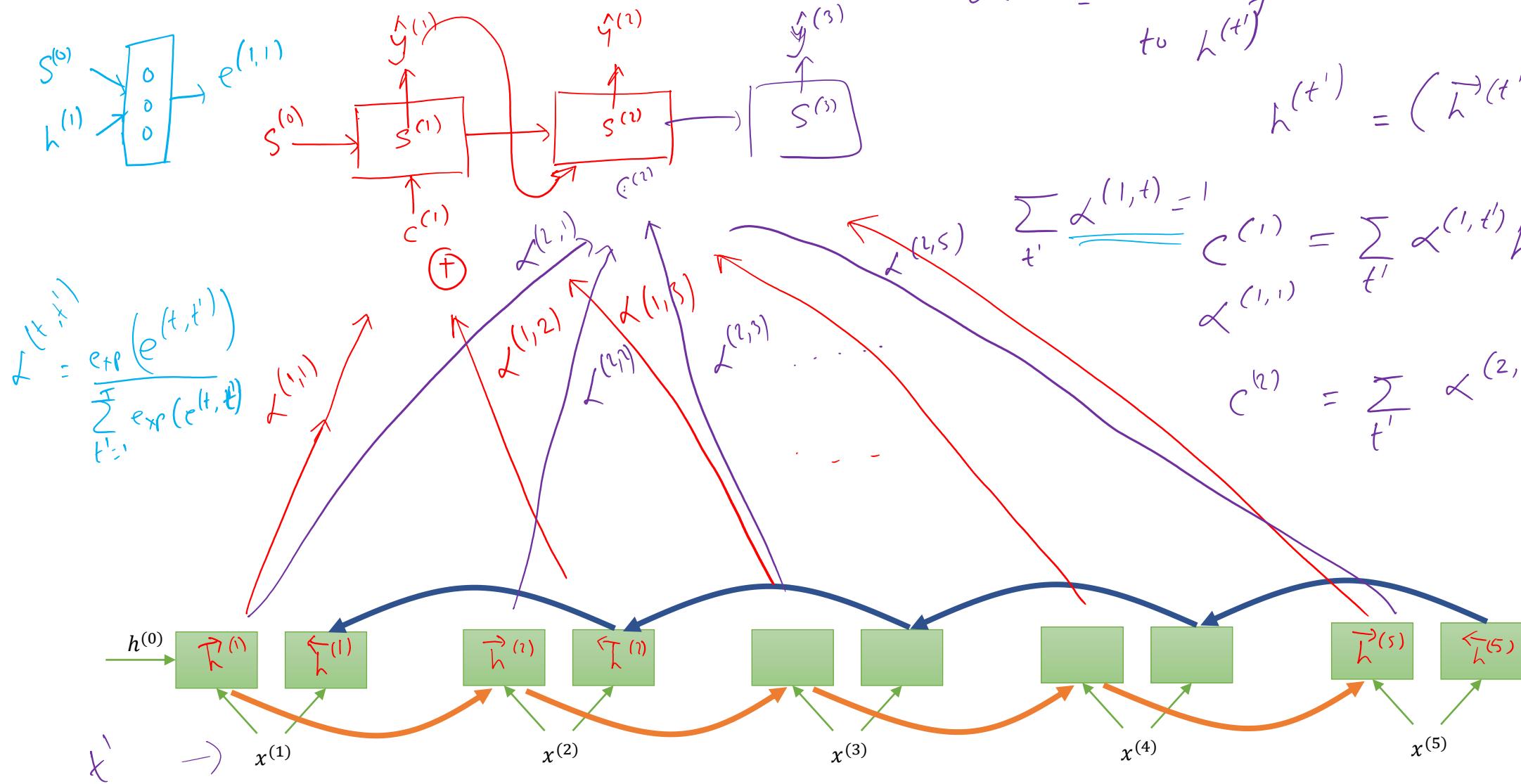
# Sentiment Analysis

- I have bought and returned three of these units now. Each one has been defective, and finally I just gave up on returning the system. The DVD player constantly gives "Bad Disc" errors and skips if there is even the slightest smudge on a disc. The sound quality is very nice for the price, but since the player doesn't work, it's essentially useless. This is a complete rip-off at any price point

# Sentiment Analysis

- I have bought and **returned** three of these units now. Each one has been **defective**, and finally I just gave up on returning the system. The DVD player constantly gives "**Bad Disc**" **errors** and skips if there is even the slightest smudge on a disc. The sound quality is **very nice** for the price, but since the player doesn't work, it's essentially **useless**. This is a **complete rip-off** at any price point

# Attention Model



$\alpha^{(t,t')} = \text{amount of "ATTENTION" } \vec{c}^{(t')} \text{ shall pay to } h^{(t')}$

$$h^{(t')} = (\vec{h}^{(t')}, \vec{l}^{(t')})$$

$$\sum_t \alpha^{(1,t)} = 1$$

$$c^{(1)} = \sum_{t'} \alpha^{(1,t')} h^{(t')}$$

$$c^{(2)} = \sum_{t'} \alpha^{(2,t')} h^{(t')}$$

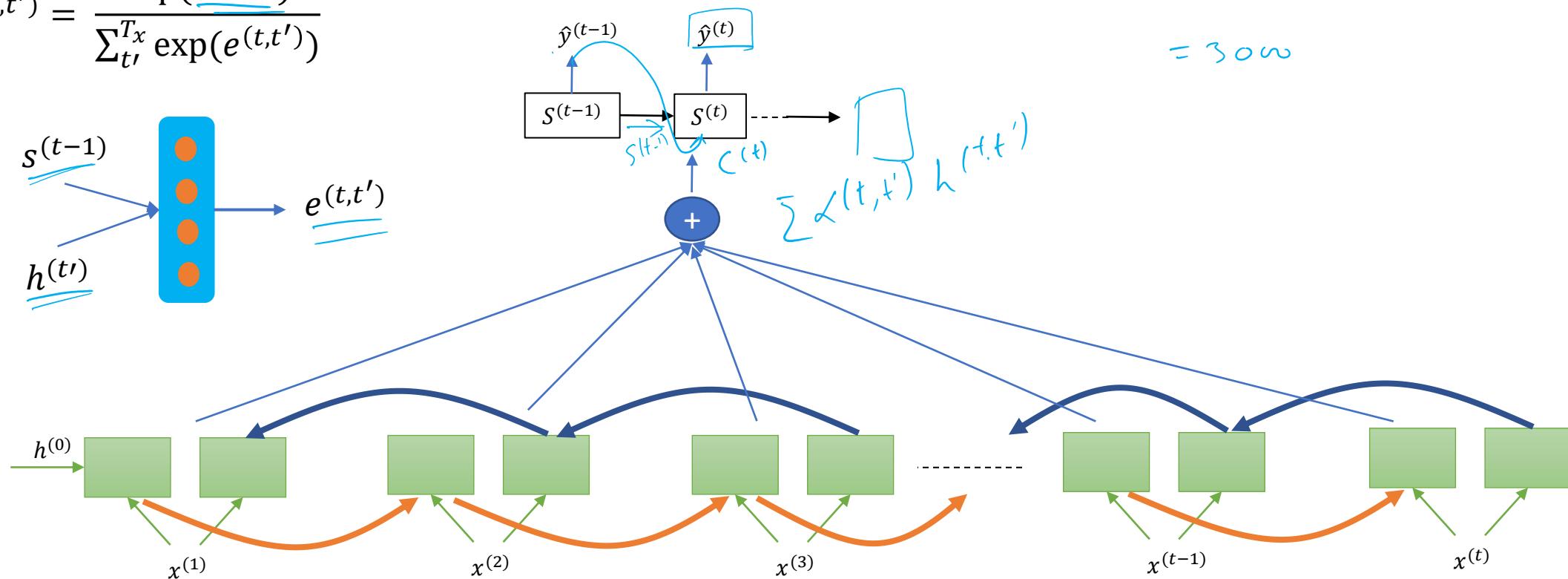
# Attention Model

$\alpha^{(t,t')} = \text{amount of attention } y^{(t)} \text{ should pay to } h^{(t')}$

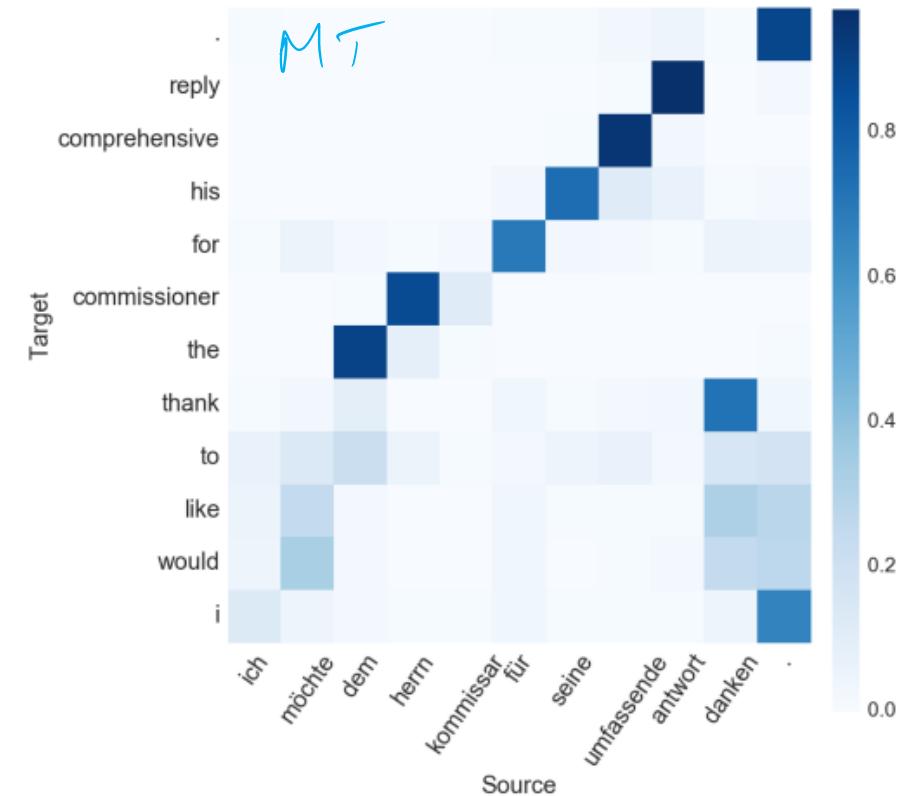
$$\alpha^{(t,t')} = \frac{\exp(e^{(t,t')})}{\sum_{t'}^T \exp(e^{(t,t')})}$$

$n \times m$

$$50 \times 60 \\ = 3000$$



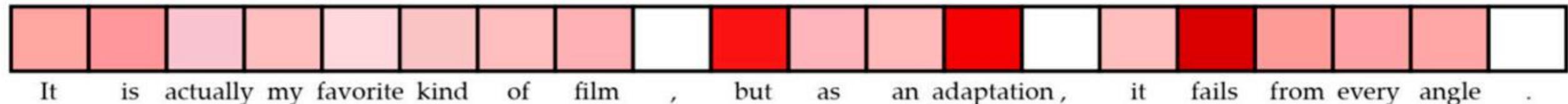
## Machine Translation



## Attention Visualization

### Sentiment Classification

*(Red wavy underline)*



# Transformers

Intuition

---

# Attention Is All You Need

---

**Ashish Vaswani\***  
Google Brain  
avaswani@google.com

**Noam Shazeer\***  
Google Brain  
noam@google.com

**Niki Parmar\***  
Google Research  
nikip@google.com

**Jakob Uszkoreit\***  
Google Research  
usz@google.com

**Llion Jones\***  
Google Research  
llion@google.com

**Aidan N. Gomez\* †**  
University of Toronto  
aidan@cs.toronto.edu

**Lukasz Kaiser\***  
Google Brain  
lukaszkaiser@google.com

**Illia Polosukhin\* ‡**  
illia.polosukhin@gmail.com

## Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

## 1 Introduction

Recurrent neural networks, long short-term memory [13] and gated recurrent [7] neural networks in particular, have been firmly established as state of the art approaches in sequence modeling and

---

\*Equal contribution. Listing order is random. Jakob proposed replacing RNNs with self-attention and started the effort to evaluate this idea. Ashish, with Illia, designed and implemented the first Transformer models and has been crucially involved in every aspect of this work. Noam proposed scaled dot-product attention, multi-head attention and the parameter-free position representation and became the other person involved in nearly every detail. Niki designed, implemented, tuned and evaluated countless model variants in our original codebase and tensor2tensor. Llion also experimented with novel model variants, was responsible for our initial codebase and

# MT and Language Models (Recap)

- Machine Translation: I like her → मैं उसे पसंद करता हूँ
- Language Model: I saw the man walking down the road
- Most of the tasks in NLP can be cast as some form of language modeling
- We can also design some architecture that uses pretty much exclusively some attention blocks and call that architecture a **transformer**

# Attention

- Attention in Computer Vision
  - Used to highlight important parts of an Image that contribute to a desired output
- Attention in NLP
  - Aligned Machine Translation
  - Language Modeling



# RNN and Transformers

Challenges with RNN	Transformer Networks
Long range dependencies	Facilitates long Range dependencies
Gradient Vanishing and Explosion	No Gradient Vanishing and Explosion
Large Number of Training Steps	Relatively Few Training Steps
Recurrence Prevents Parallel Computations	No recurrence → helps parallel processing

# Attention Mechanism

- Mimics retrieval of a value  $v_i$  for a query  $q$  based on a key  $k_i$  in database

- $\text{Attention}(q, k, v) = \sum_i \text{similarity}(q, k_i) \times v_i$

books of RT  
key  
—  
—  
—  
Row

$Q \rightarrow$  Select age from tabl. column Name. Contains(Alt)

key	Name	age	addr	...
n <sub>1</sub>	n <sub>1</sub>	38		
n <sub>2</sub>	n <sub>2</sub>			
n <sub>3</sub>	n <sub>3</sub>			
...	...			
n <sub>k</sub>	n <sub>k</sub>			

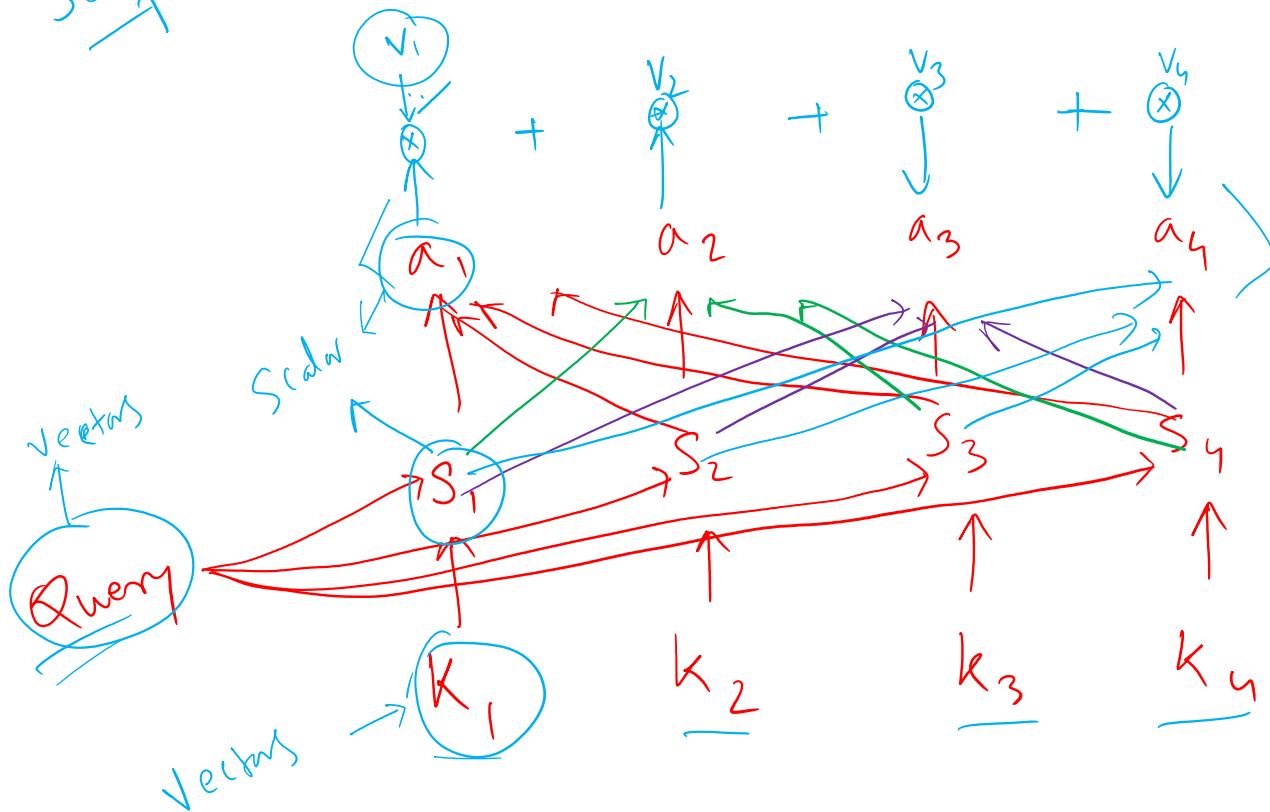
$(0\ 0\ 0\ 0\ 1\ 0\ 0\ 0\ 0)$   
 $(.33, .33, .32, .33, ..) v_i$

$f(q, k)$   
 $= q^T \cdot k$   
 $= q^T k$  → dimension  
 $\in q^T w_k i$

key  
 value  
 38

$= w_1^T q^T + w_2^T k^T$

Self Attention



$$\sum (q^T k_i) \times v_i \rightarrow \text{Vector}$$

$$a_i = \frac{\exp(s_i)}{\sum_j \exp(s_j)} \rightarrow \text{Scalar}$$

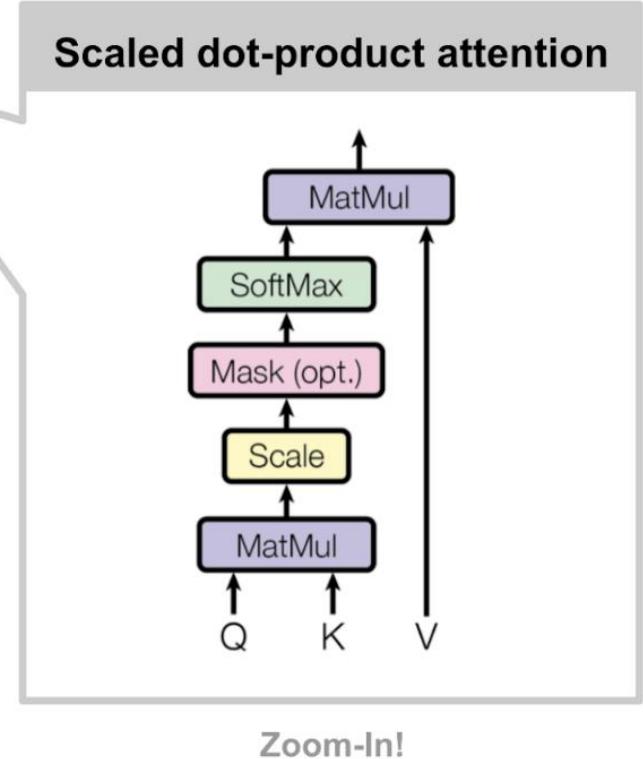
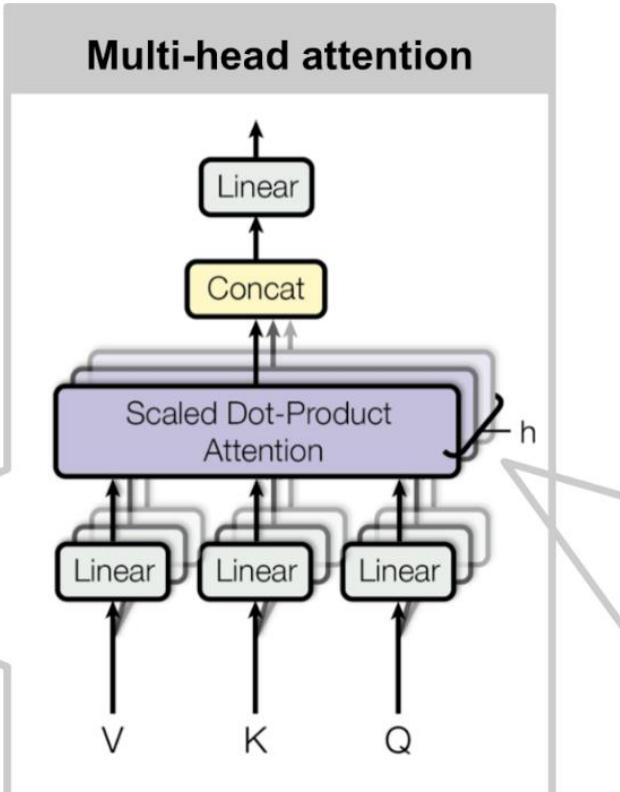
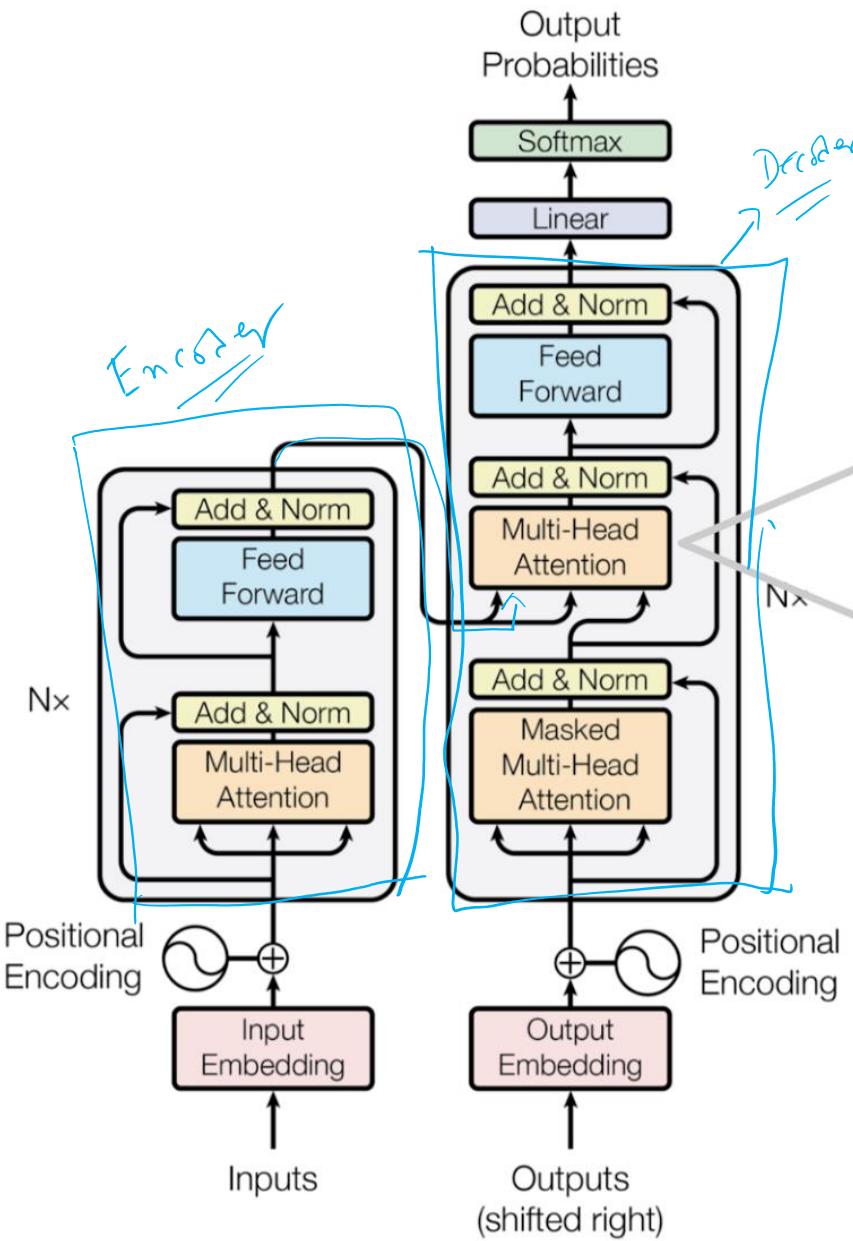
$$s_i = f(q, k_i)$$

$$\left\{
 \begin{array}{l}
 s_i = q^T k_i \rightarrow \text{dot product} \\
 s_i = \frac{q^T k_i}{\sqrt{d}} \rightarrow \text{Scaled dot product} \\
 s_i = [q^T W k_i] \\
 s_i = w_q^T q^T + w_k k_i \rightarrow \text{additive sum}
 \end{array}
 \right.$$

# Attention Mechanism

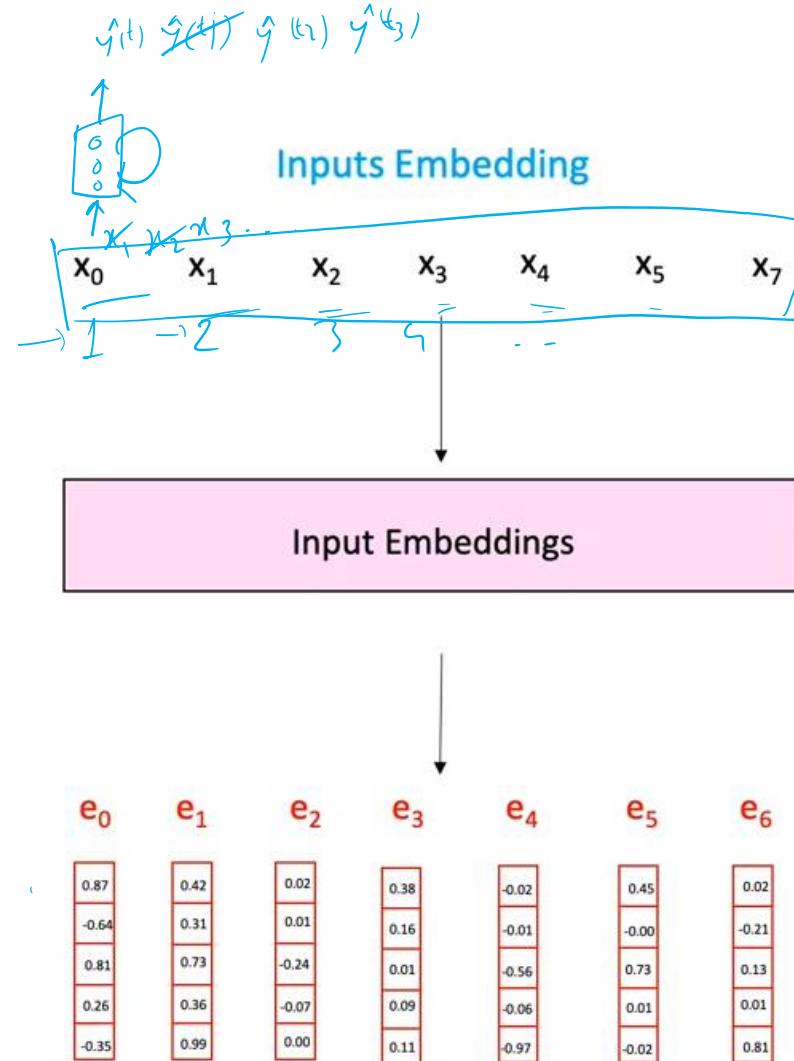
- Neural Architecture
- Example:
  - Machine Translation
    - Query:  $S_i$  hidden vector for the ith output word
    - Key:  $h_j$  hidden vector for the jth input word
    - Value:  $h_j$  hidden vector for the jth input word

# The Transformer Network



# The Transformer Network

- Input Embedding



# The Transformer Network

- Position Embedding

Position Embeddings (Intuition)

$e_0$	$e_1$	$e_2$	$e_3$	$e_4$	$e_5$	$e_6$
0.42	0.87	0.02	0.38	-0.02	0.45	0.02
0.31	-0.64	0.01	0.16	-0.01	-0.00	-0.21
0.73	0.81	-0.24	0.01	-0.56	0.73	0.13
0.36	0.26	-0.07	0.09	-0.06	0.01	0.01
0.99	-0.35	0.00	0.00	-0.97	-0.02	0.81

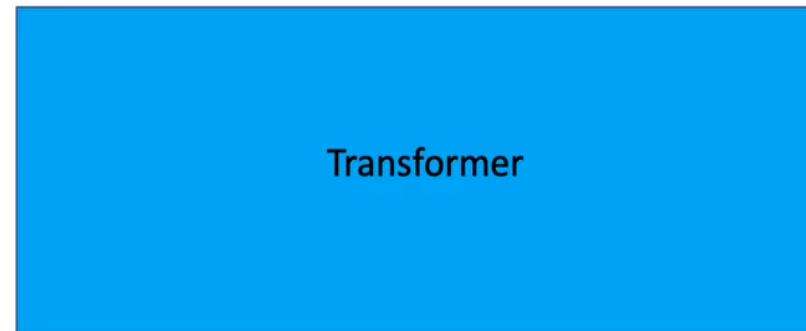
LSTM

# The Transformer Network

- Position Embedding

Position Embeddings (Intuition)

$e_0$	$e_1$	$e_2$	$e_3$	$e_4$	$e_5$	$e_6$
0.42	0.87	0.02	0.38	-0.02	0.45	0.02
0.31	-0.64	0.01	0.16	-0.01	-0.00	-0.21
0.73	0.81	-0.24	0.01	-0.56	0.73	0.13
0.36	0.26	-0.07	0.09	-0.06	0.01	0.01
0.99	-0.35	0.00	0.00	-0.97	-0.02	0.81



# The Transformer Network

- Position Embedding



*Even though she did not win the match, she was satisfied*

Transformer



# The Transformer Network

- Why Position Embedding matters?

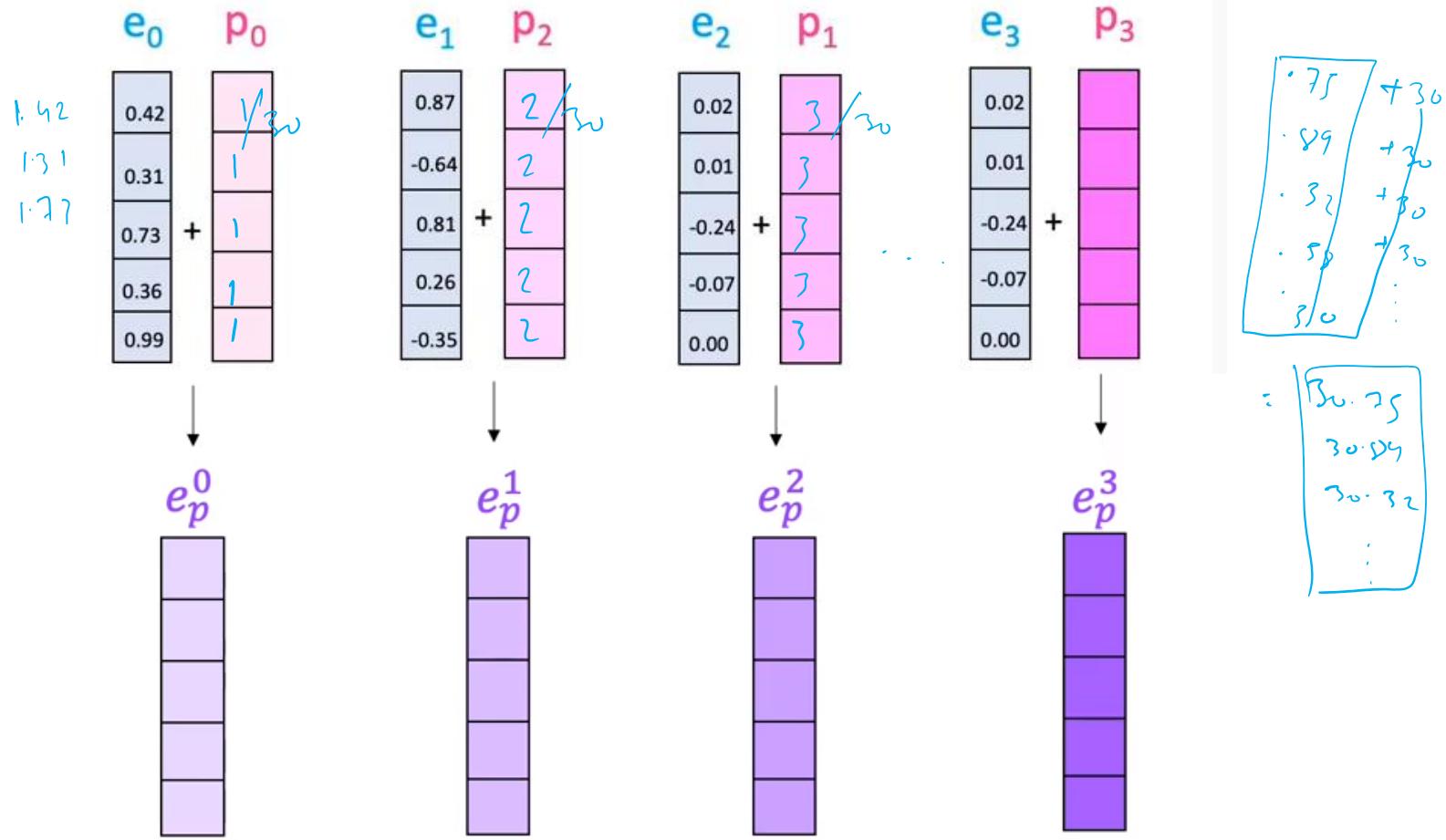
- Even though she did **not** win the match, she was **satisfied**

- Even though she did win the match, she was **not** satisfied

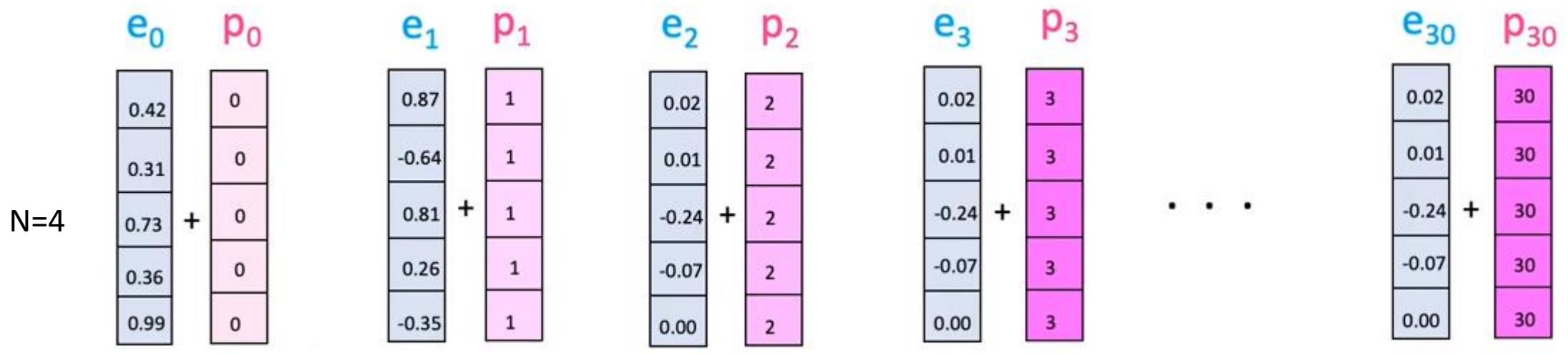
# The Transformer Network

- Position Embedding

$\begin{bmatrix} 0 & 1 \end{bmatrix}$



# The Transformer Network: Position Embeddings



$$0 \times \frac{1}{3}$$

$$1 \times \frac{1}{3}$$

$$2 \times \frac{1}{3}$$

$$3 \times \frac{1}{3}$$

$$\frac{1}{N-1} = \frac{1}{3}$$

# The Transformer Network: Position Embeddings

Sentence 1

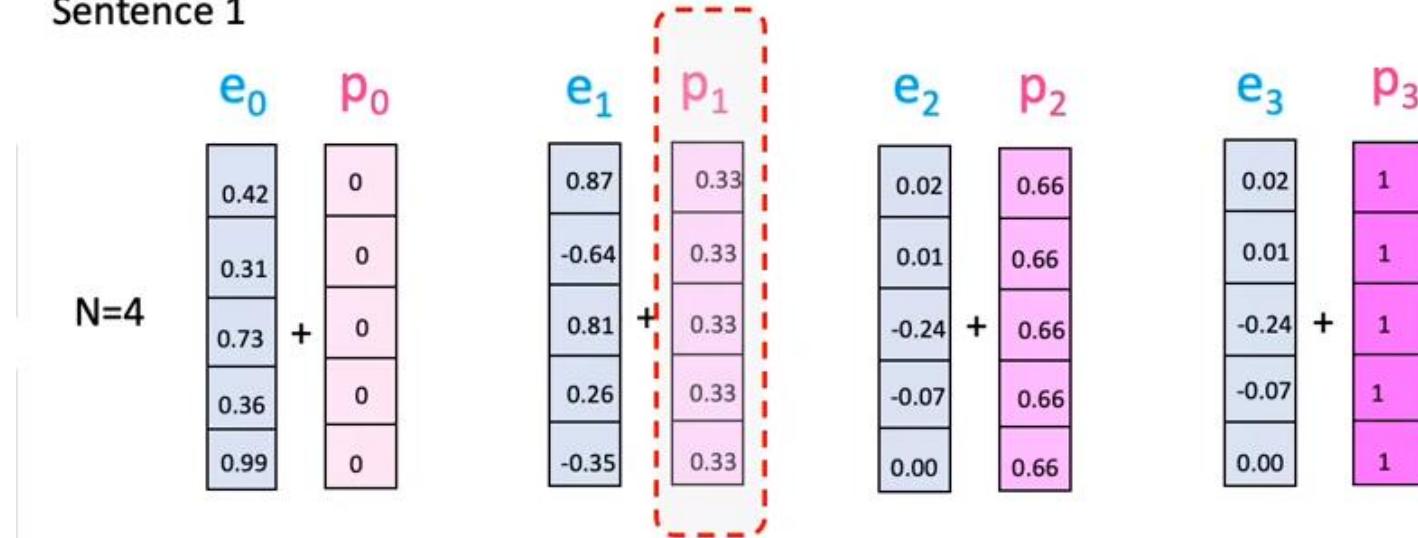
	$e_0$	$p_0$	$e_1$	$p_1$	$e_2$	$p_2$	$e_3$	$p_3$
N=4	0.42	0	0.87	0.33	0.02	0.66	0.02	1
	0.31	0	-0.64	0.33	0.01	0.66	0.01	1
	0.73	0	0.81	0.33	-0.24	0.66	-0.24	1
	0.36	0	0.26	0.33	-0.07	0.66	-0.07	1
	0.99	0	-0.35	0.33	0.00	0.66	0.00	1

Sentence 2

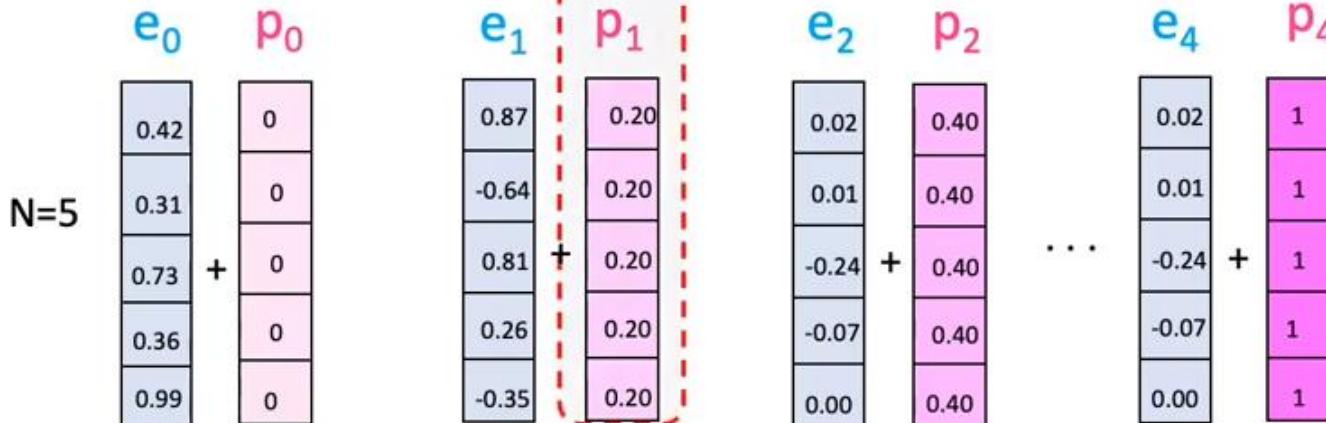
	$e_0$	$p_0$	$e_1$	$p_1$	$e_2$	$p_2$	$e_4$	$p_4$
N=5	0.42	0	0.87	0.20	0.02	0.40	0.02	1
	0.31	0	-0.64	0.20	0.01	0.40	0.01	1
	0.73	0	0.81	0.20	-0.24	0.40	-0.24	1
	0.36	0	0.26	0.20	-0.07	0.40	-0.07	1
	0.99	0	-0.35	0.20	0.00	0.40	0.00	1

# The Transformer Network: Position Embeddings

Sentence 1

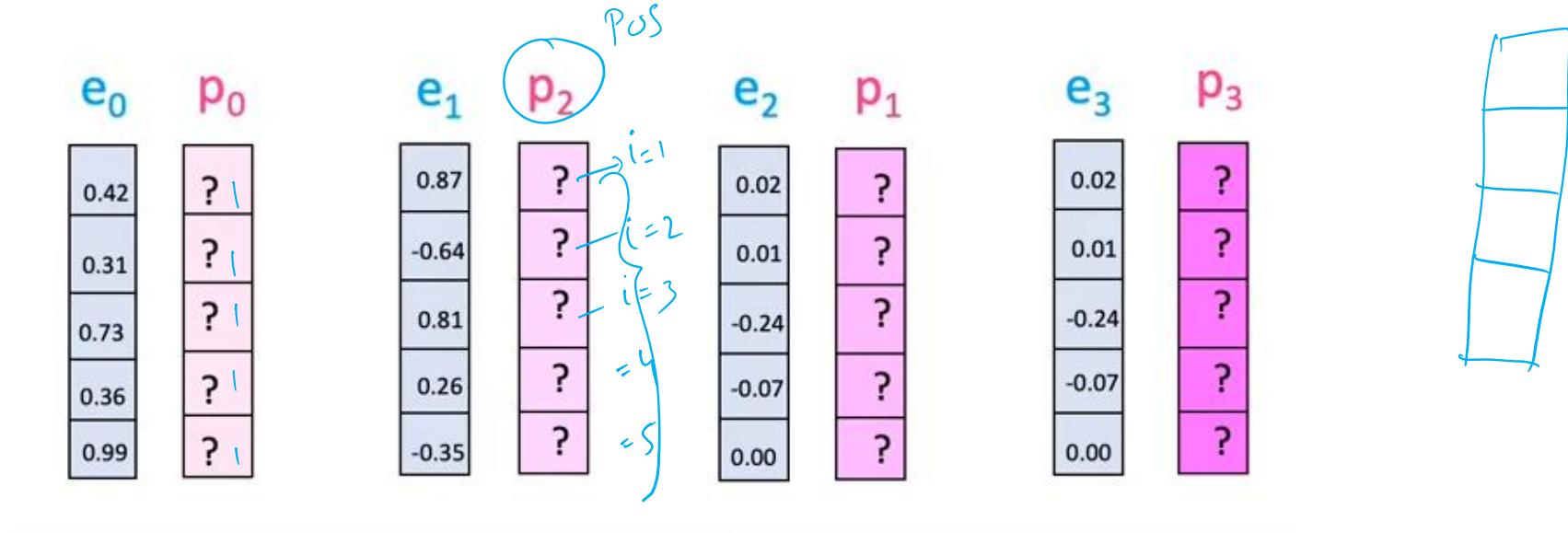


Sentence 2



# The Transformer Network: Position Embeddings

What about frequencies?



$d=5$

Vaswani et al 2017

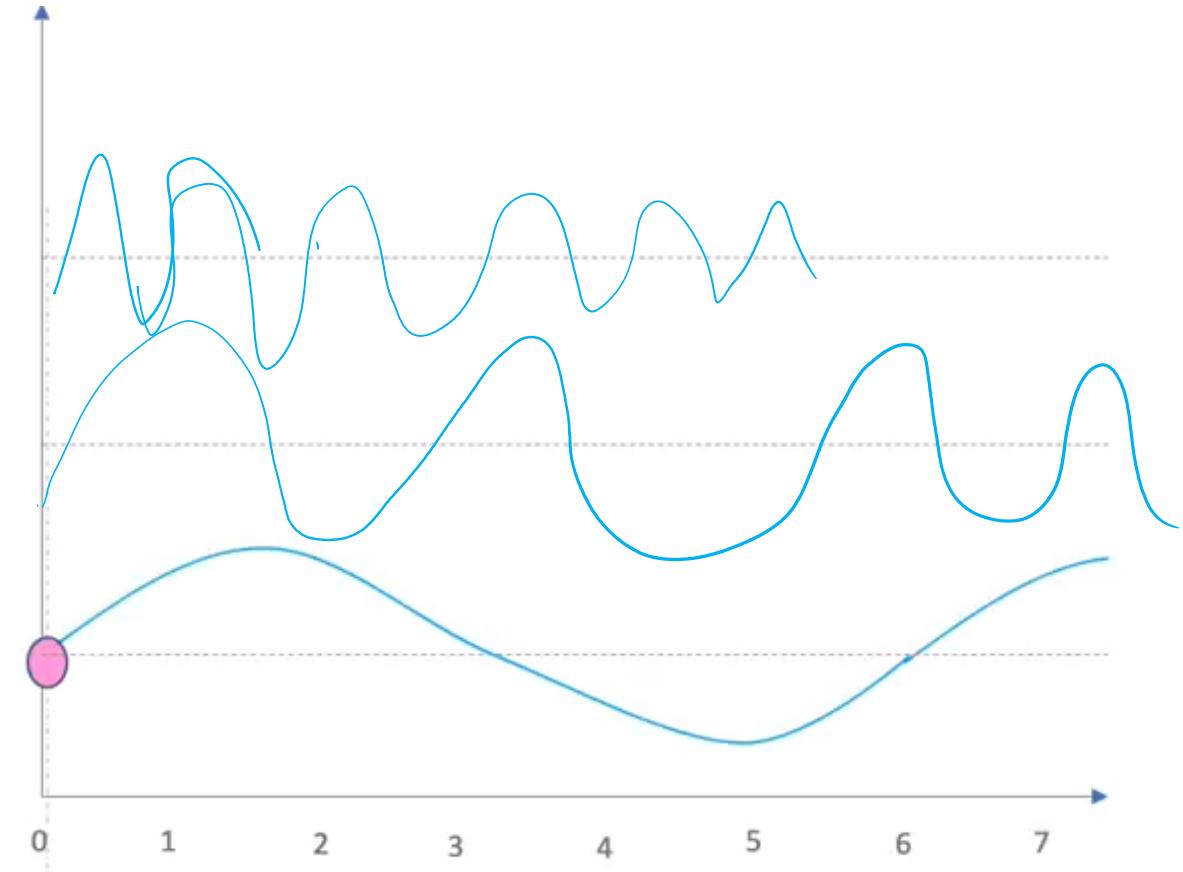
$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000^{\frac{2i}{d}}}\right)$$

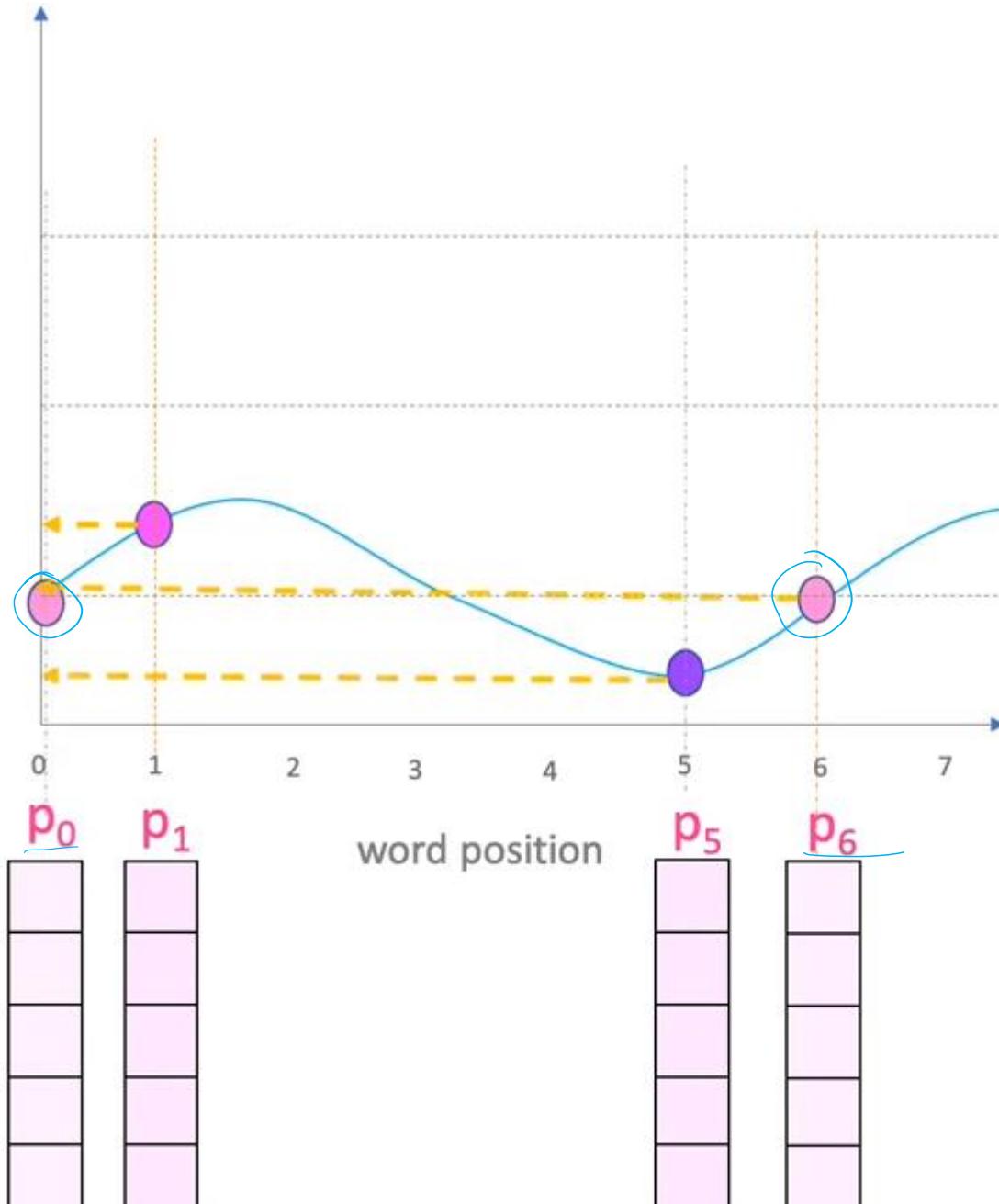
$$PE_{(pos,2i+1)} = \cos\left(\frac{pos}{10000^{\frac{2i}{d}}}\right)$$

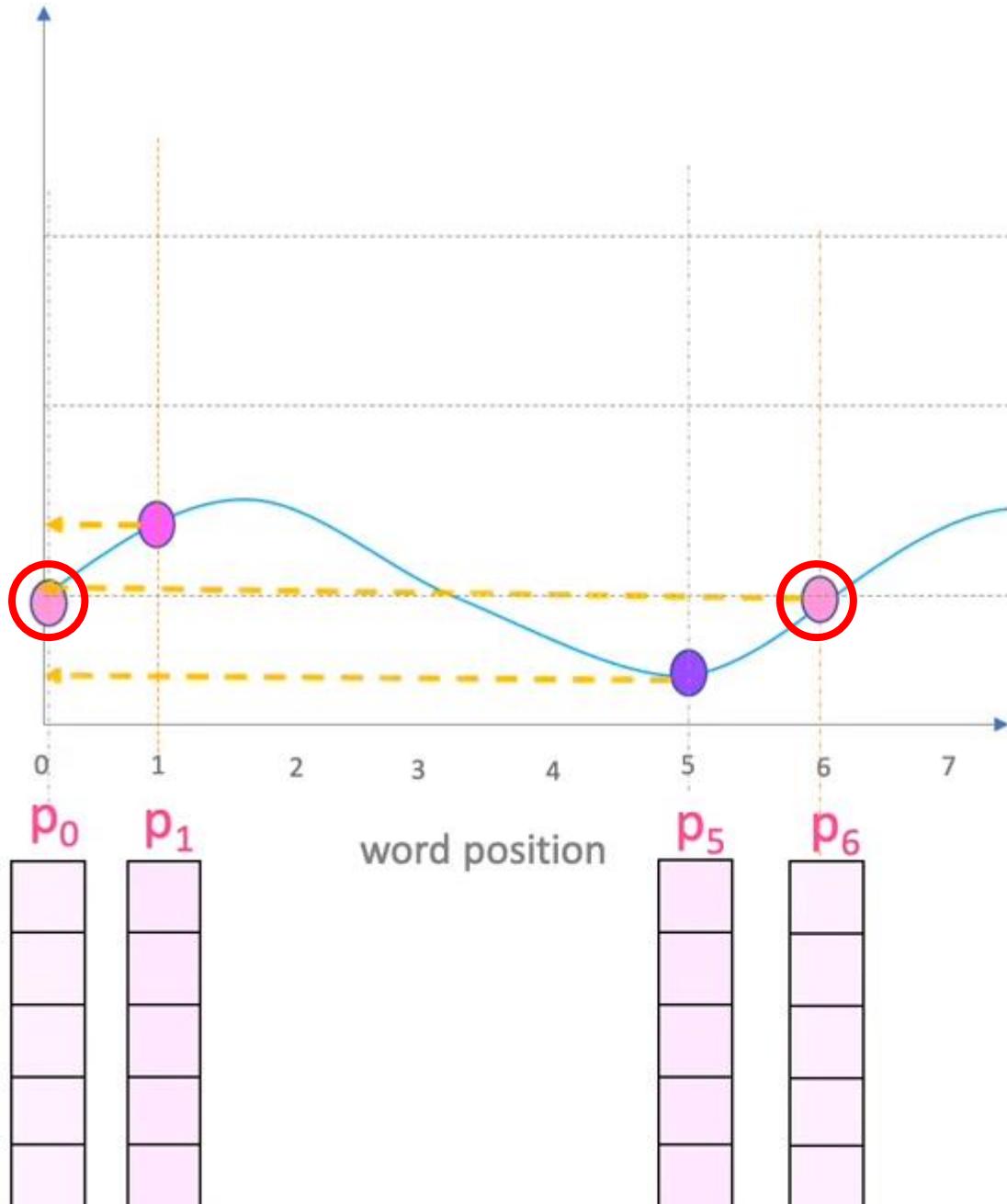


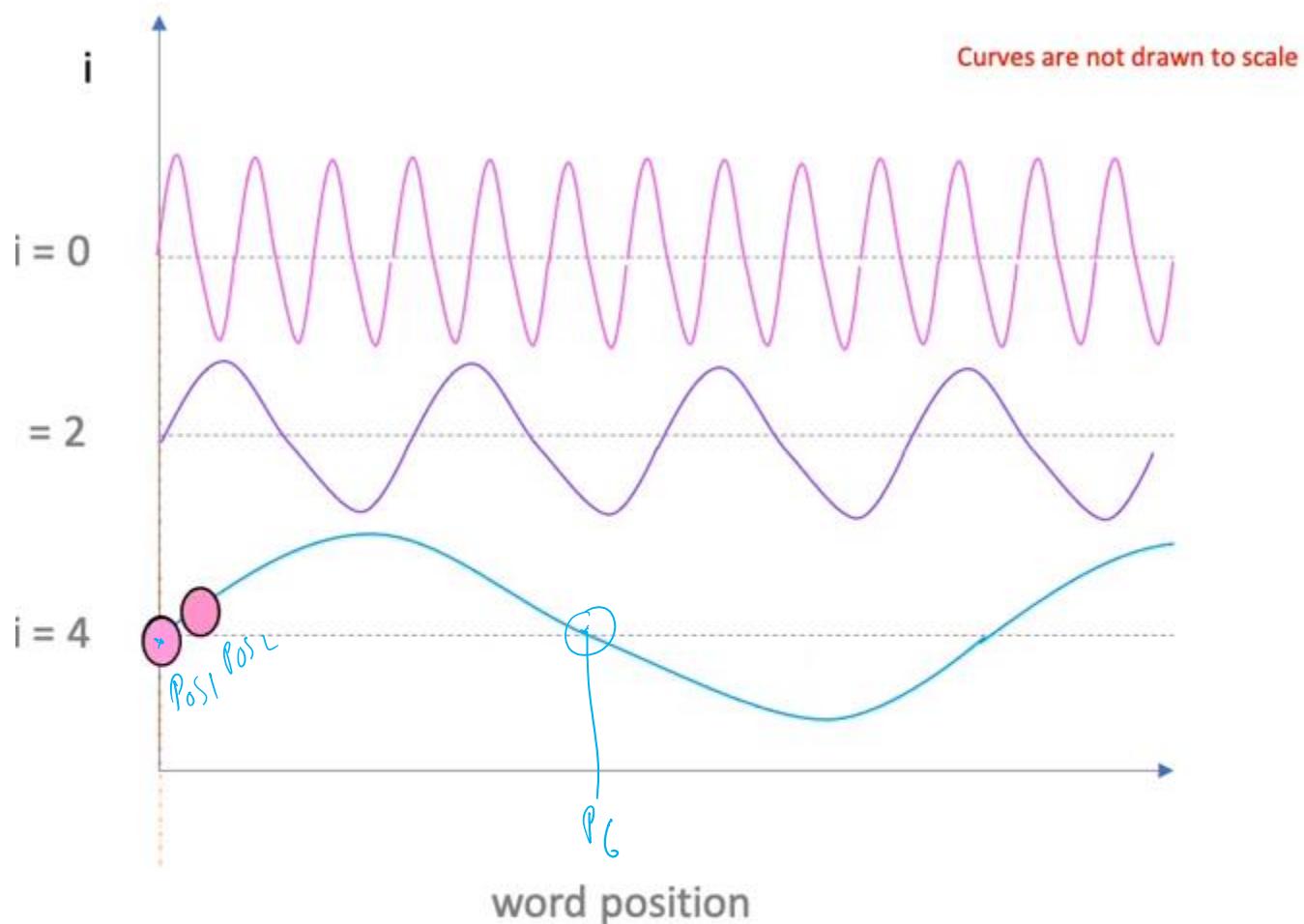
$pos = 0$

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

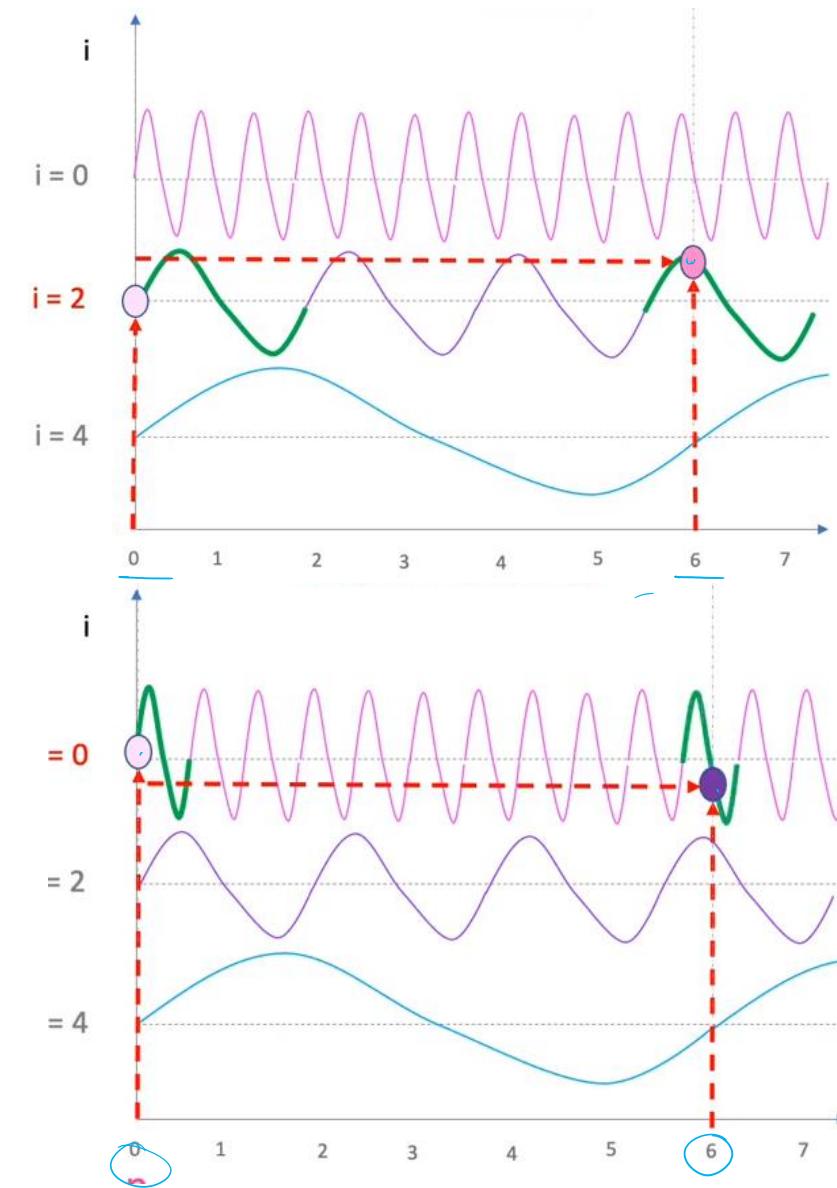
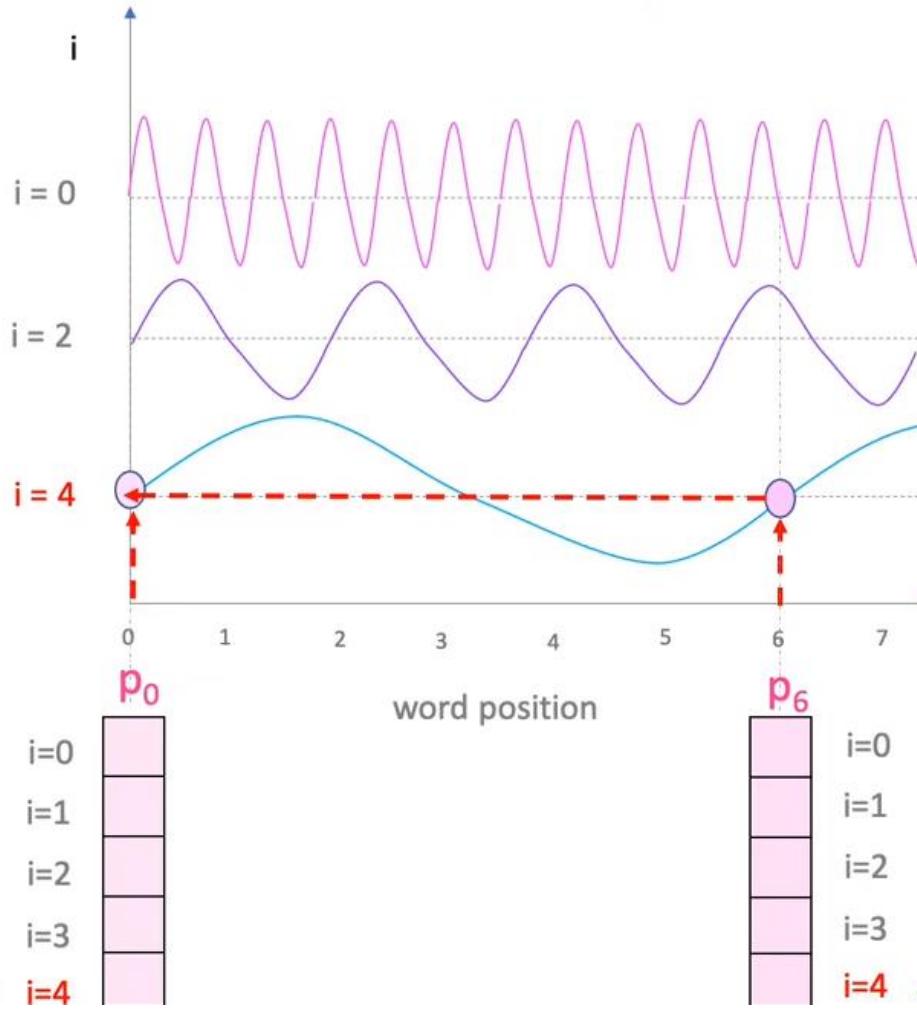


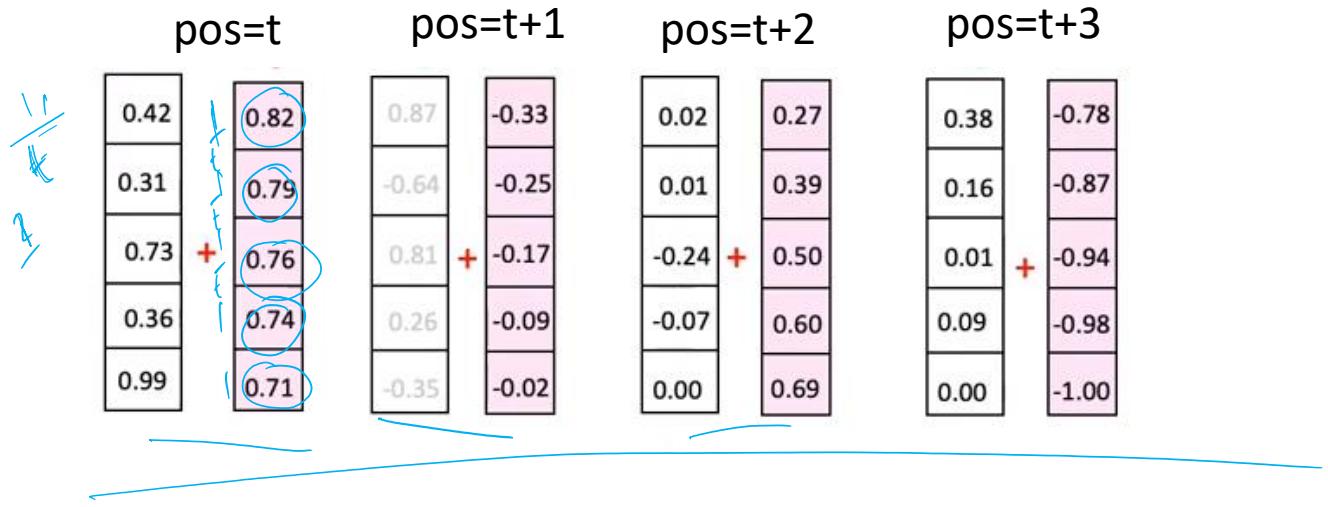






$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000} \frac{2i}{d}\right)$$





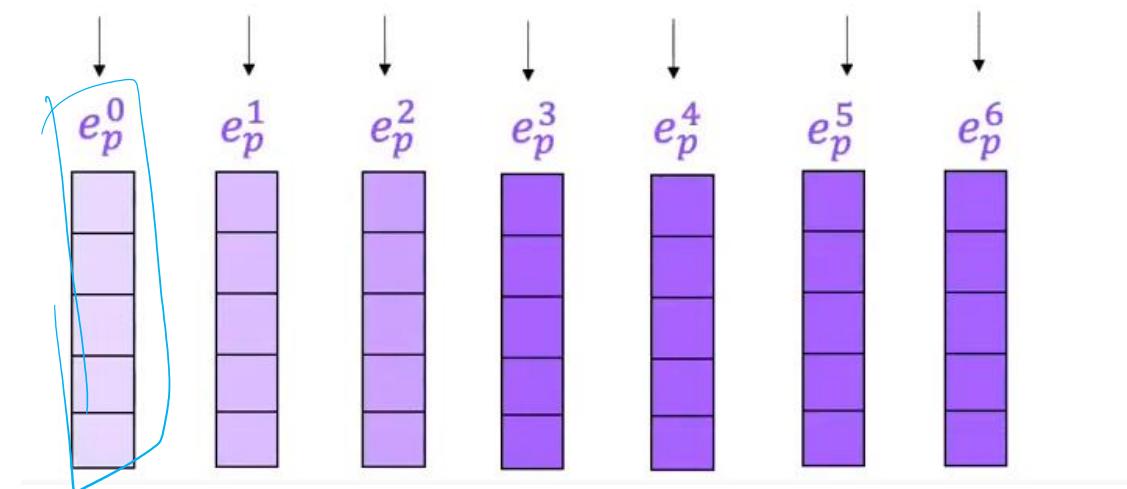
$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000} \frac{2i}{d}\right)$$

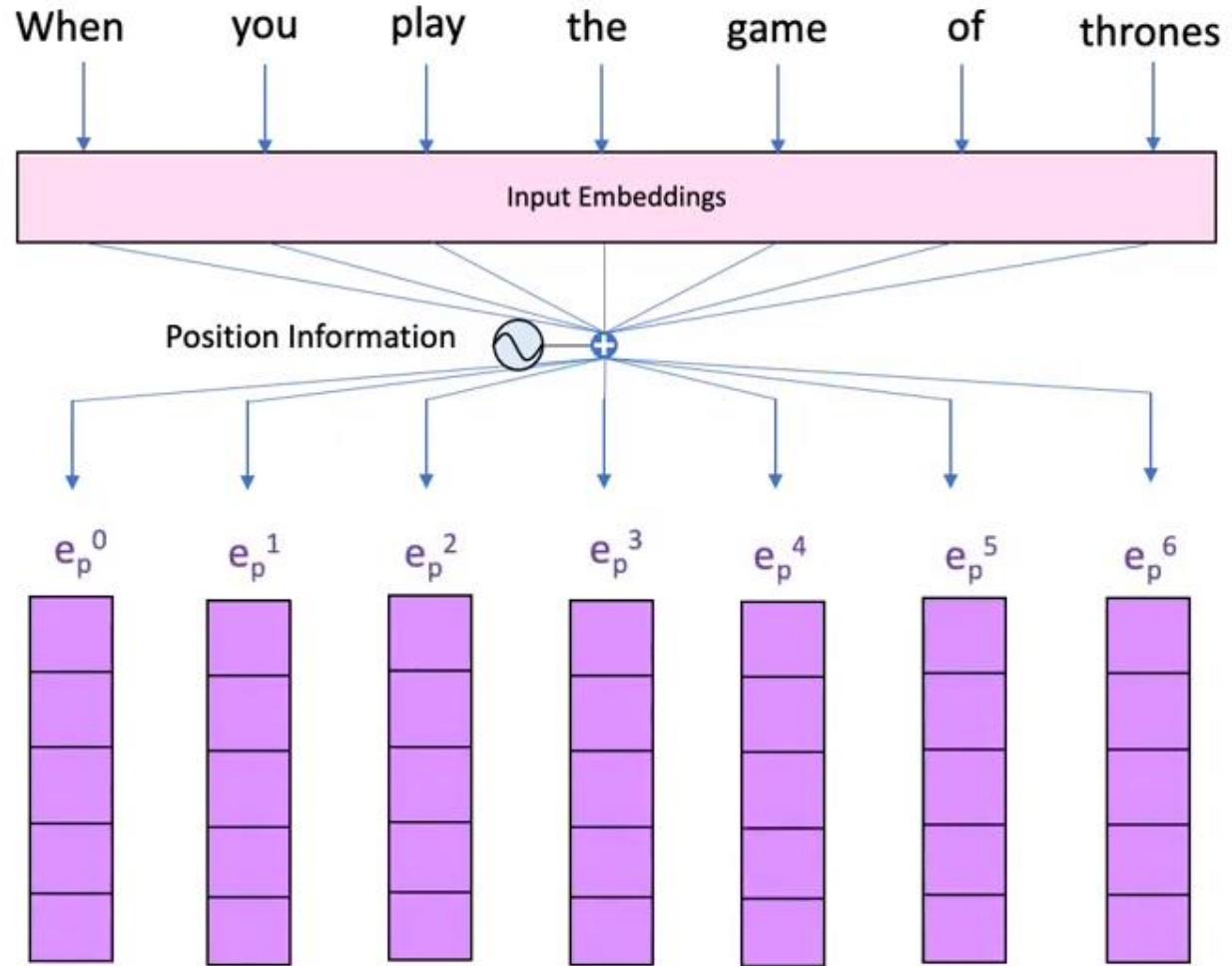
$$PE_{(pos,2i+1)} = \cos\left(\frac{pos}{10000} \frac{2i}{d}\right)$$

$$\vec{p}_t = \begin{bmatrix} \sin(\omega_1 \cdot t) \\ \cos(\omega_1 \cdot t) \\ \vdots \\ \sin(\omega_{d/2} \cdot t) \\ \cos(\omega_{d/2} \cdot t) \end{bmatrix}_{d \times 1}$$

## Position Embeddings

When	you	play	the	game	of	thrones
$e_0 + p_0$	$e_1 + p_1$	$e_2 + p_2$	$e_3 + p_3$	$e_3 + p_3$	$e_3 + p_3$	$e_3 + p_3$
0.47	0.87	0.02	0.38	0.38	0.38	0.38
0.31	-0.64	0.01	0.16	0.16	0.16	0.16
0.73	0.81	-0.24	0.01	0.01	0.01	0.01
0.36	0.26	-0.07	0.09	0.09	0.09	0.09
0.99	-0.35	0.00	0.00	0.00	0.00	0.00



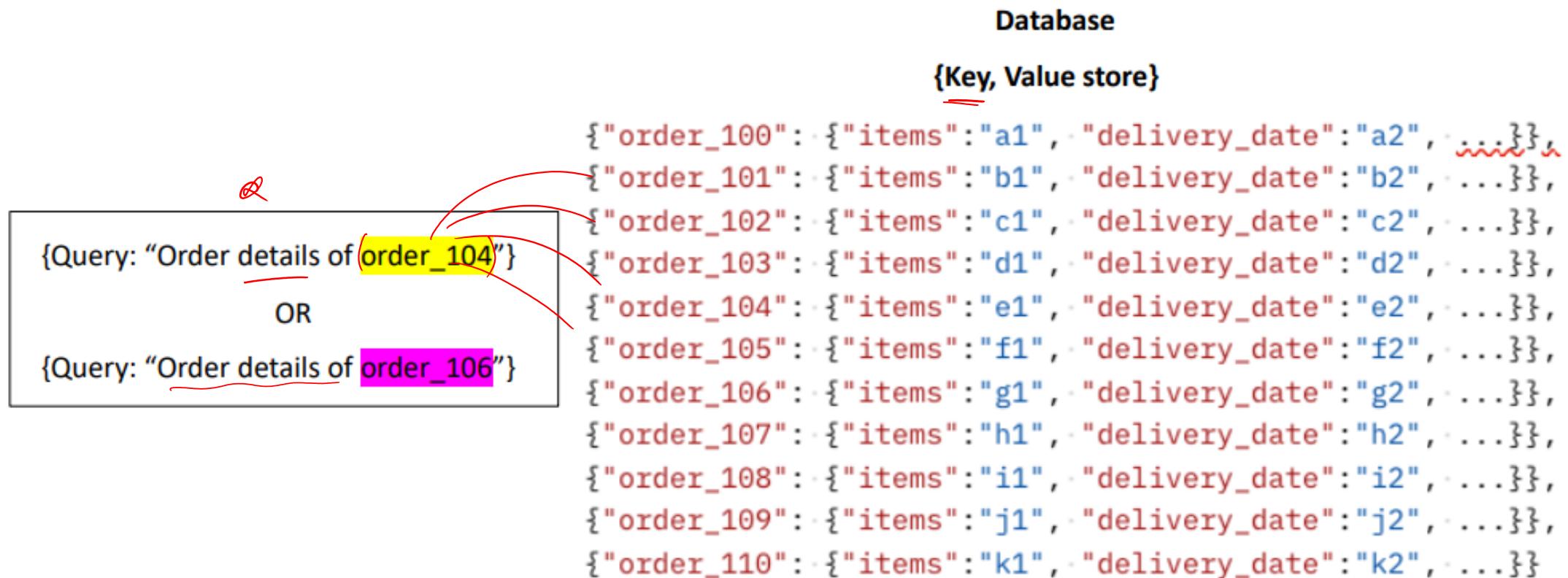


# Query, Key & Value

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

# Query, Key & Value



{Query: "Order details of order\_104"}  
OR  
{Query: "Order details of order\_106"}

{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", ...}}
```

# Query, Key & Value

{Query: "Order details of order\_104"}  
OR  
{Query: "Order details of order\_106"}

{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", ...}, value  
 "order_105": {"items": "f1", "delivery_date": "f2", ...},  
 "order_106": {"items": "g1", "delivery_date": "g2", ...}, value  
 "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, Key & Value

{Query: "Order details of order\_104"}  
OR  
{Query: "Order details of order\_106"}

{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", ...}}
```

# Query, Key & Value

Done at the same time !!

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

OR

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

{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", ...}}
```

# Query, Key & Value

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

OR

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

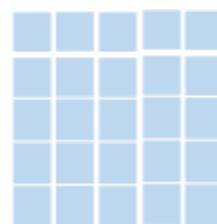
```
{"order_100": {"items": "a1", "delivery_date": "a2", ...}},  
 {"order_101": {"items": "b1", "delivery_date": "b2", ...}},  
 {"order_102": {"items": "c1", "delivery_date": "c2", ...}},  
 {"order_103": {"items": "d1", "delivery_date": "d2", ...}},  
 {"order_104": {"items": "e1", "delivery_date": "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	Key	Value
1. Search for info	1. Interacts directly with Queries 2. Distinguishes one object from another 3. Identify which object is the most relevant and by how much	1. Actual details of the object 2. More fine grained

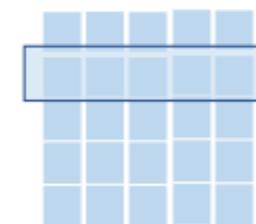
# Attention



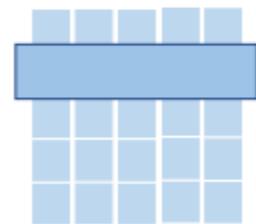
Query



Key Value Store



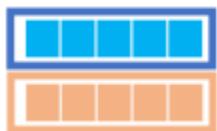
Key



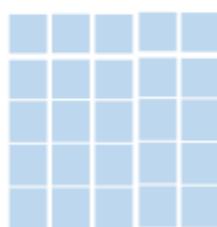
Value

## Attention

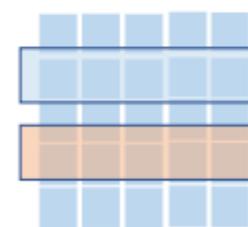
The food is good



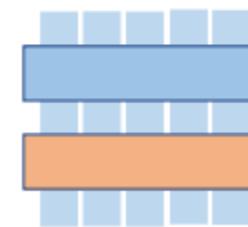
Query



Key Value Store



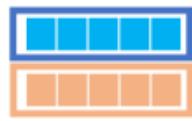
Key



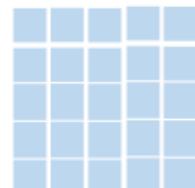
Value

$\sum_{t'} \alpha(t, t') (\vec{Q} \cdot \vec{k}) \vec{V}$   
**Attention**  
 Scaled Attention.

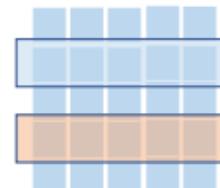
Parallelizable !!!



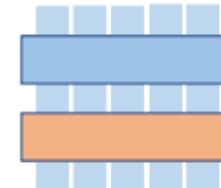
Query



Key Value Store



Key



Value

$Q$

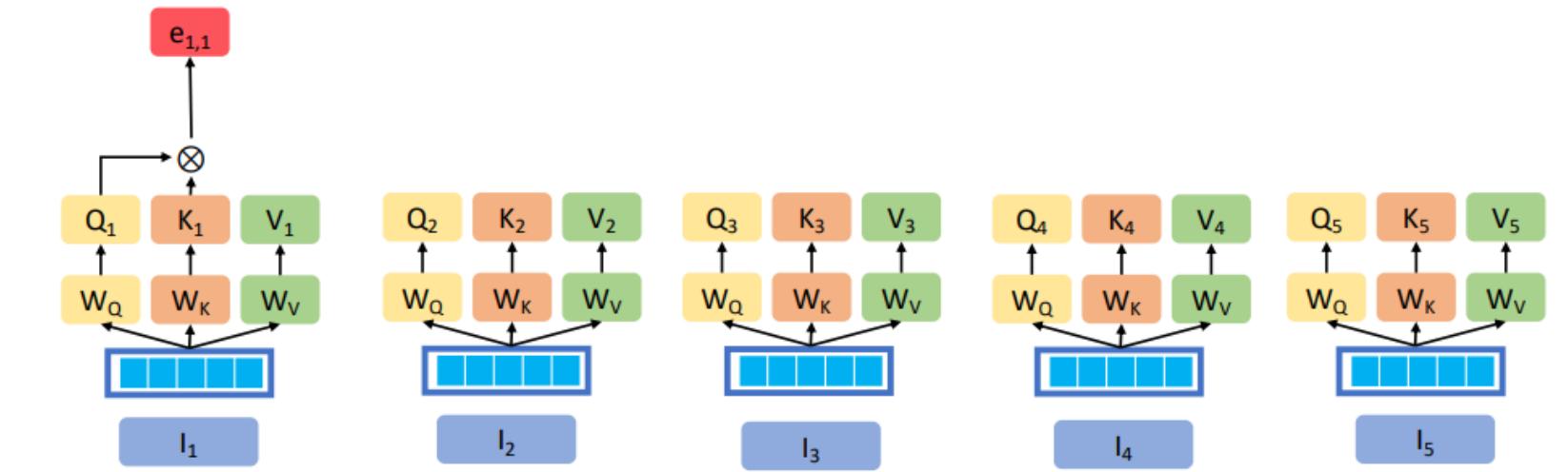
$QK^T$

$softmax\left(\frac{QK^T}{\sqrt{d}}\right)$

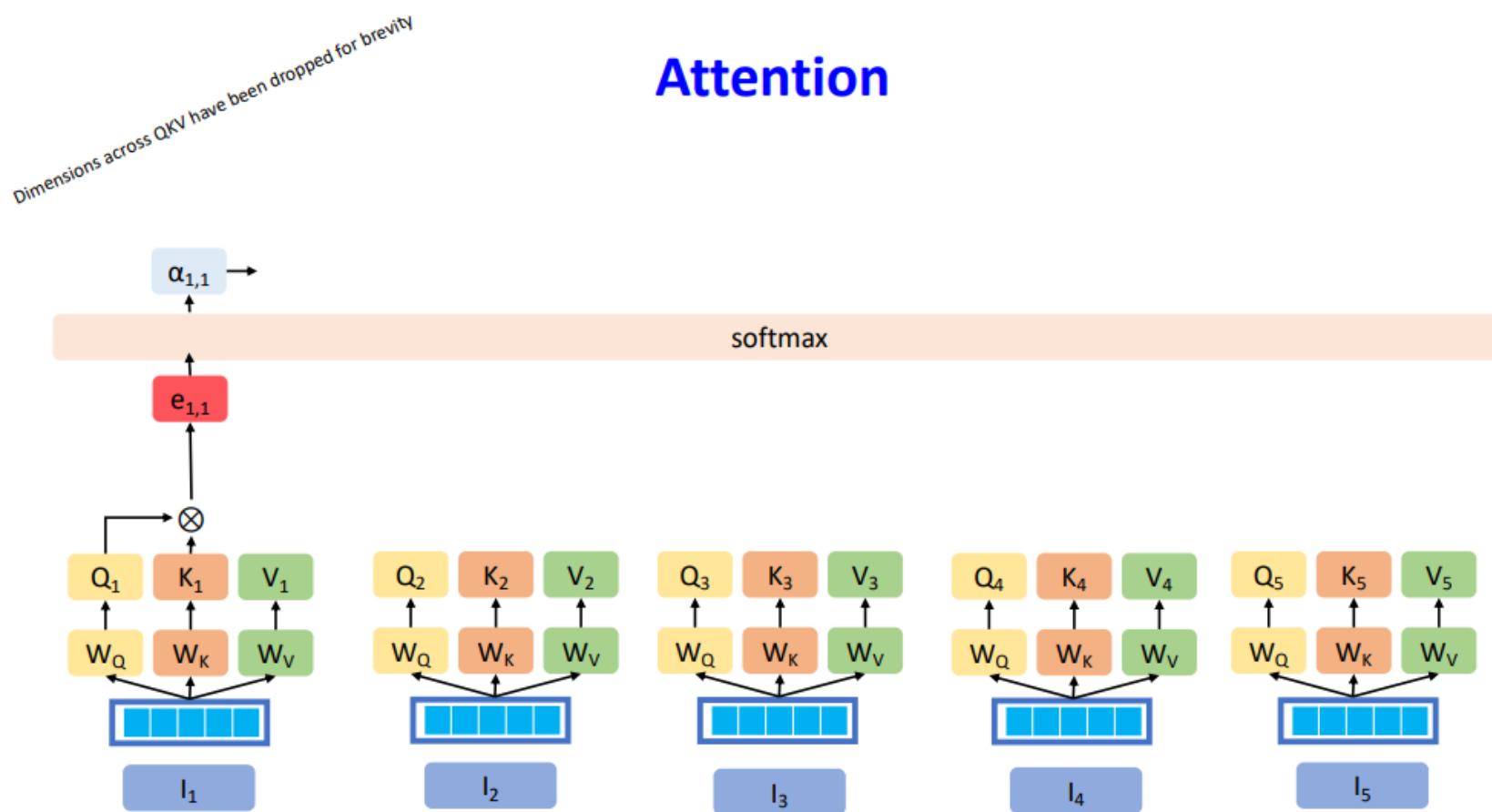
$softmax\left(\frac{QK^T}{\sqrt{d}}\right)V$

# Attention

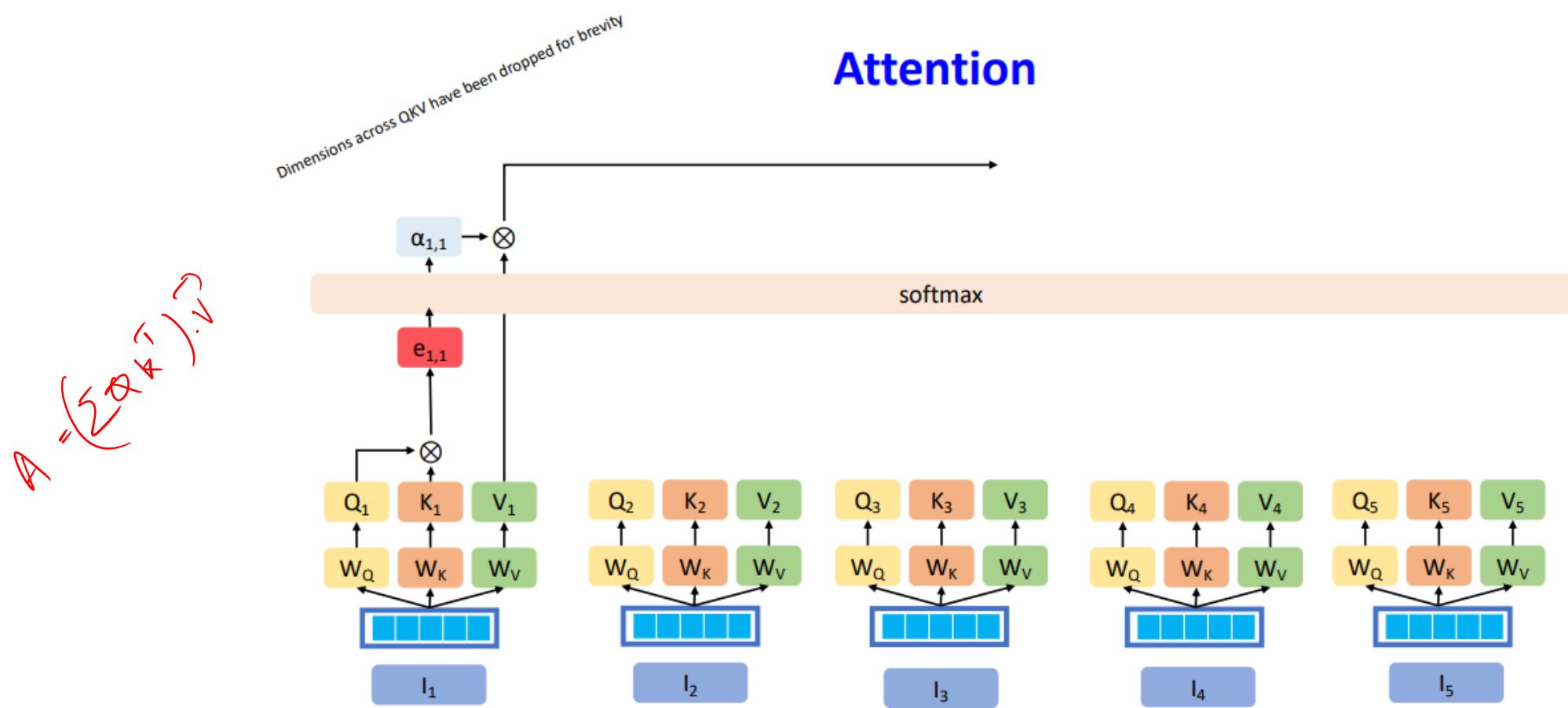
Dimensions across QKV have been dropped for brevity



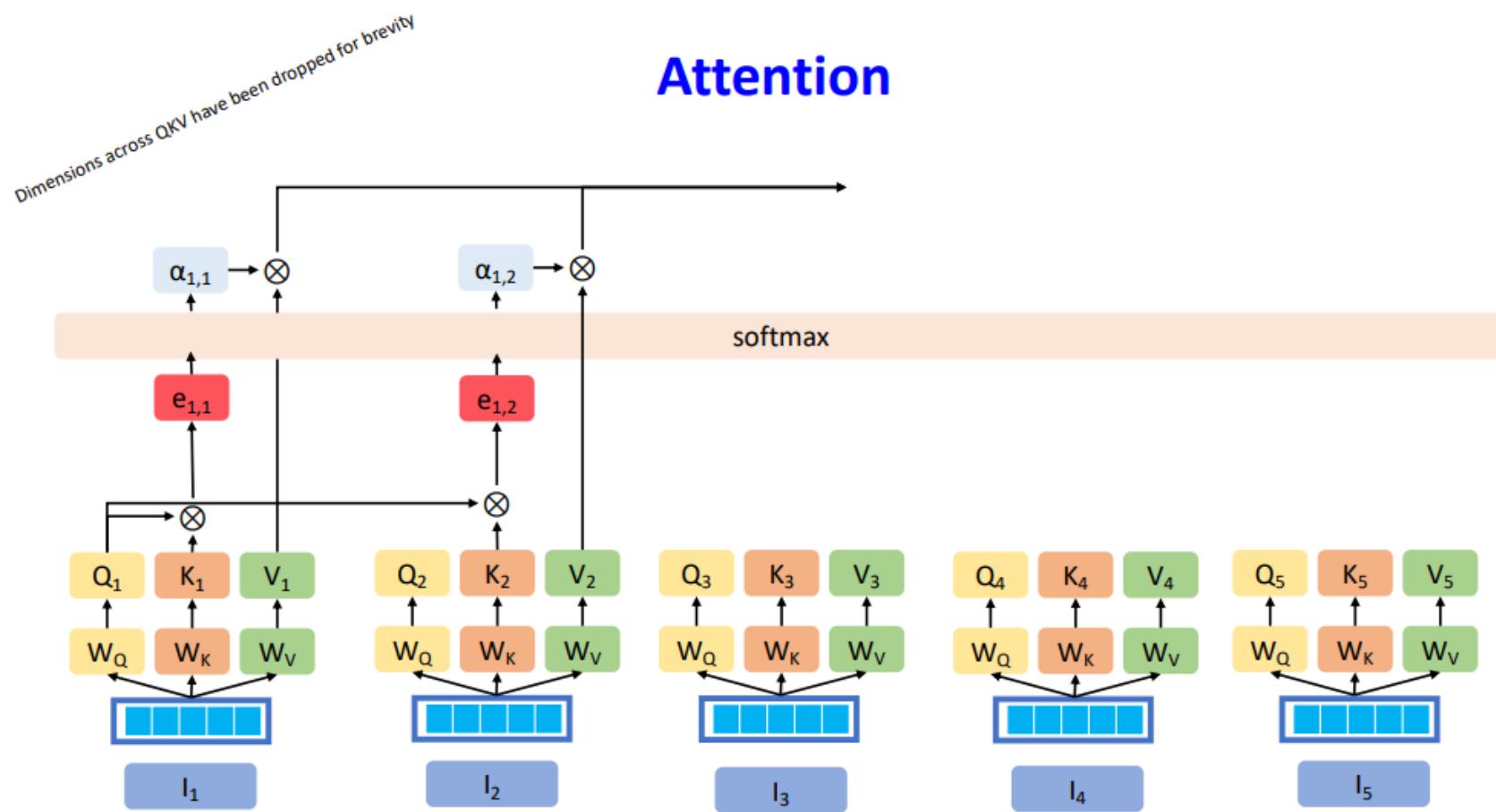
## Attention



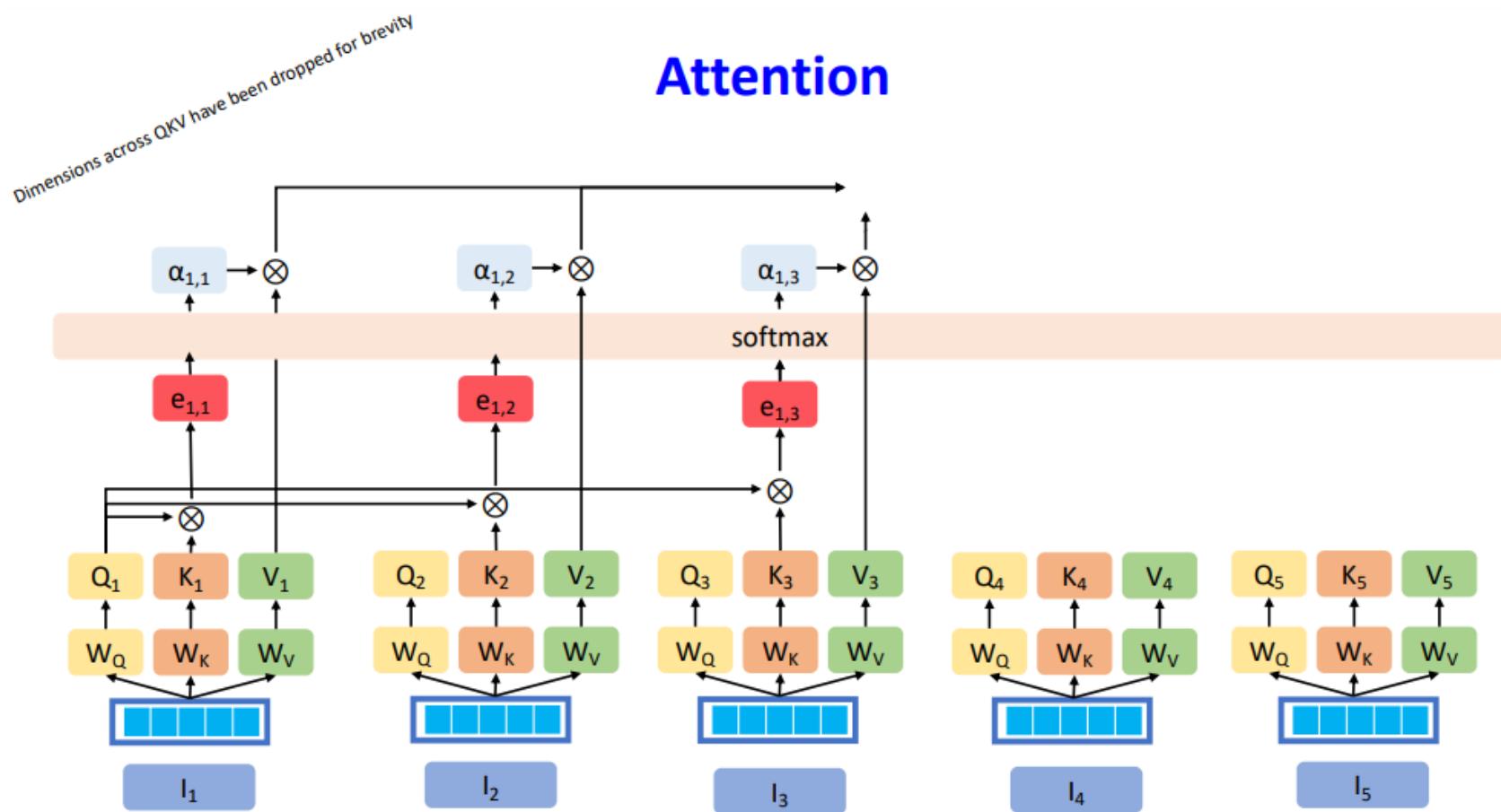
## Attention



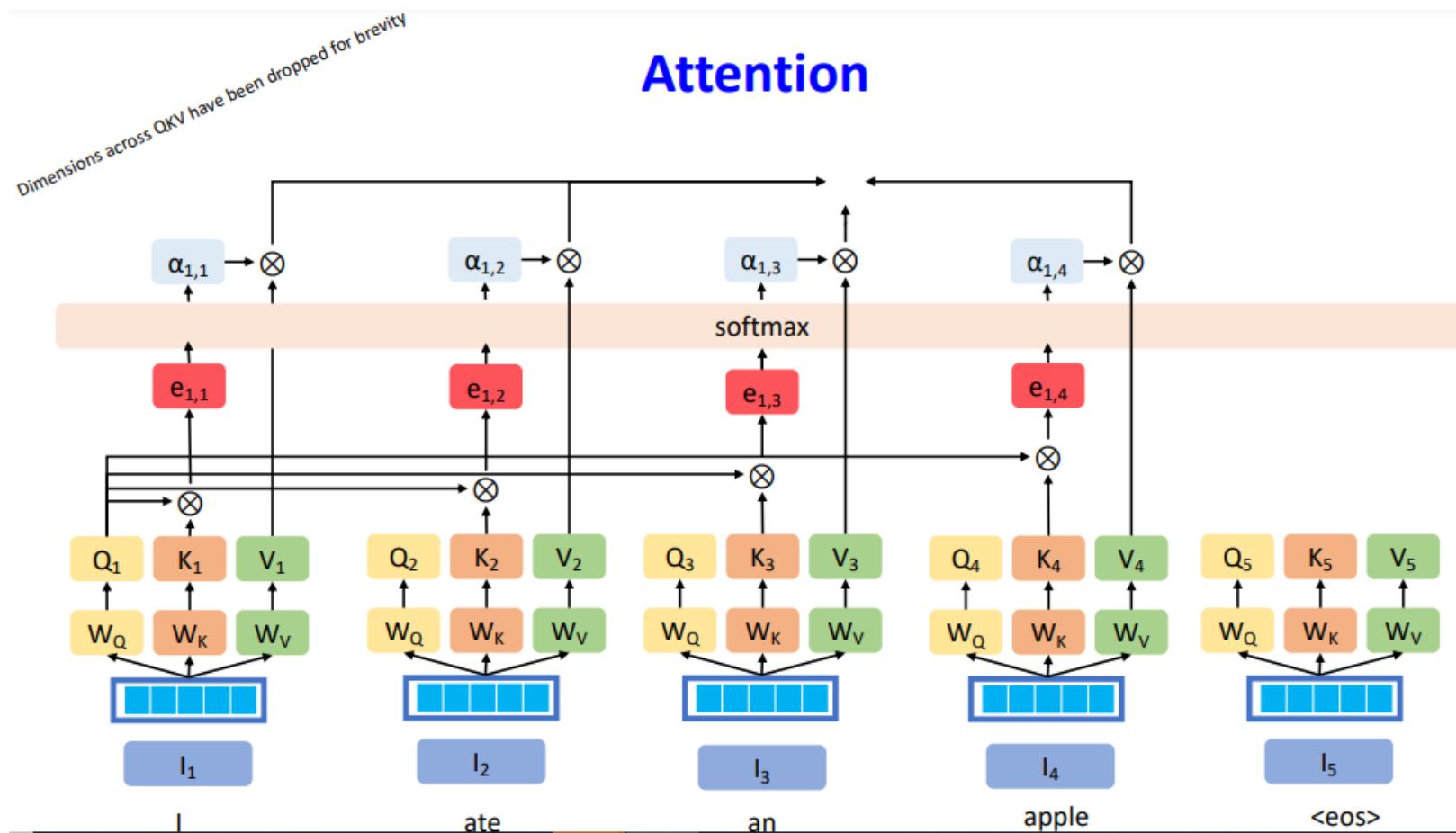
# Attention



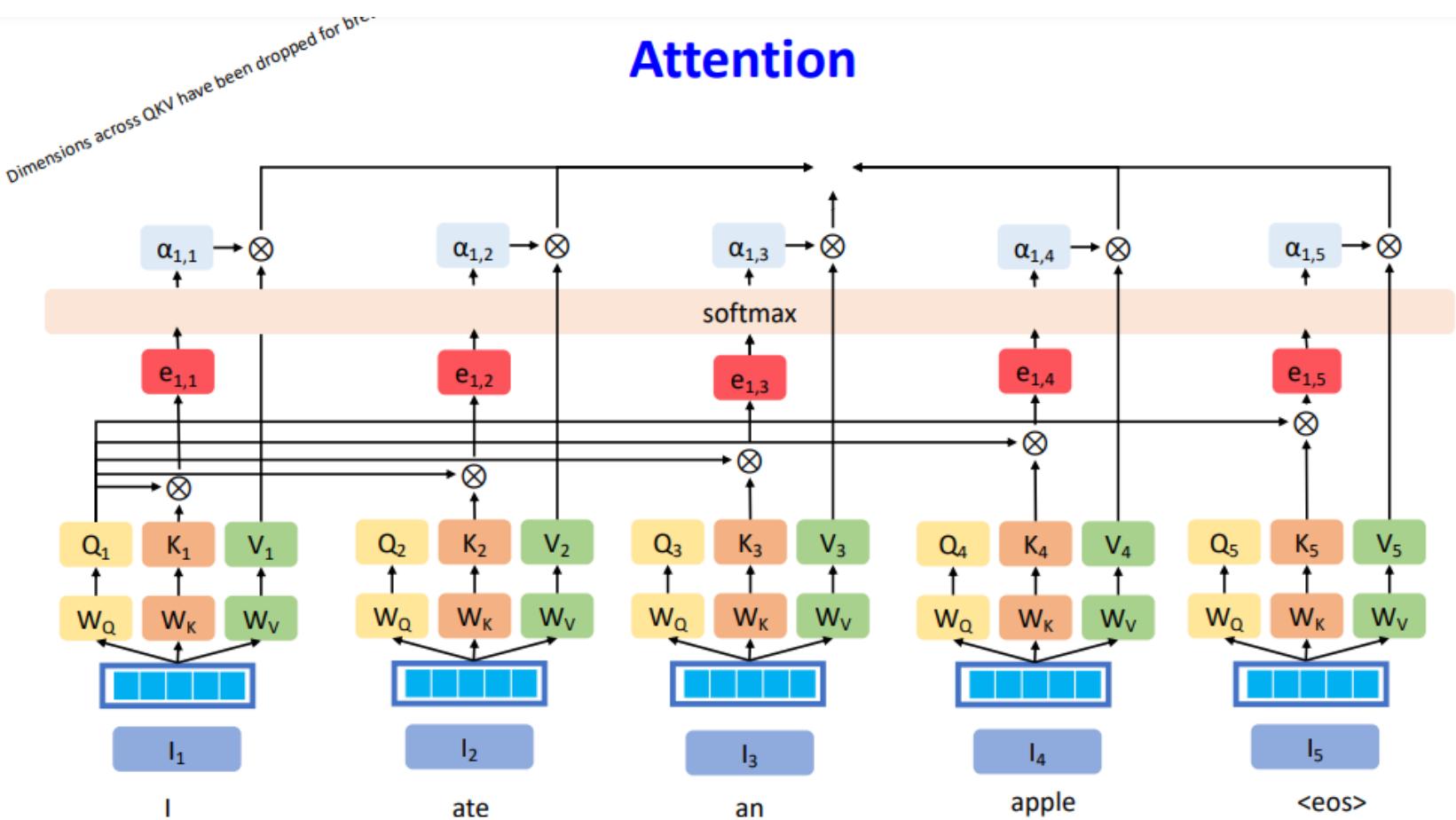
## Attention



## Attention



## Attention

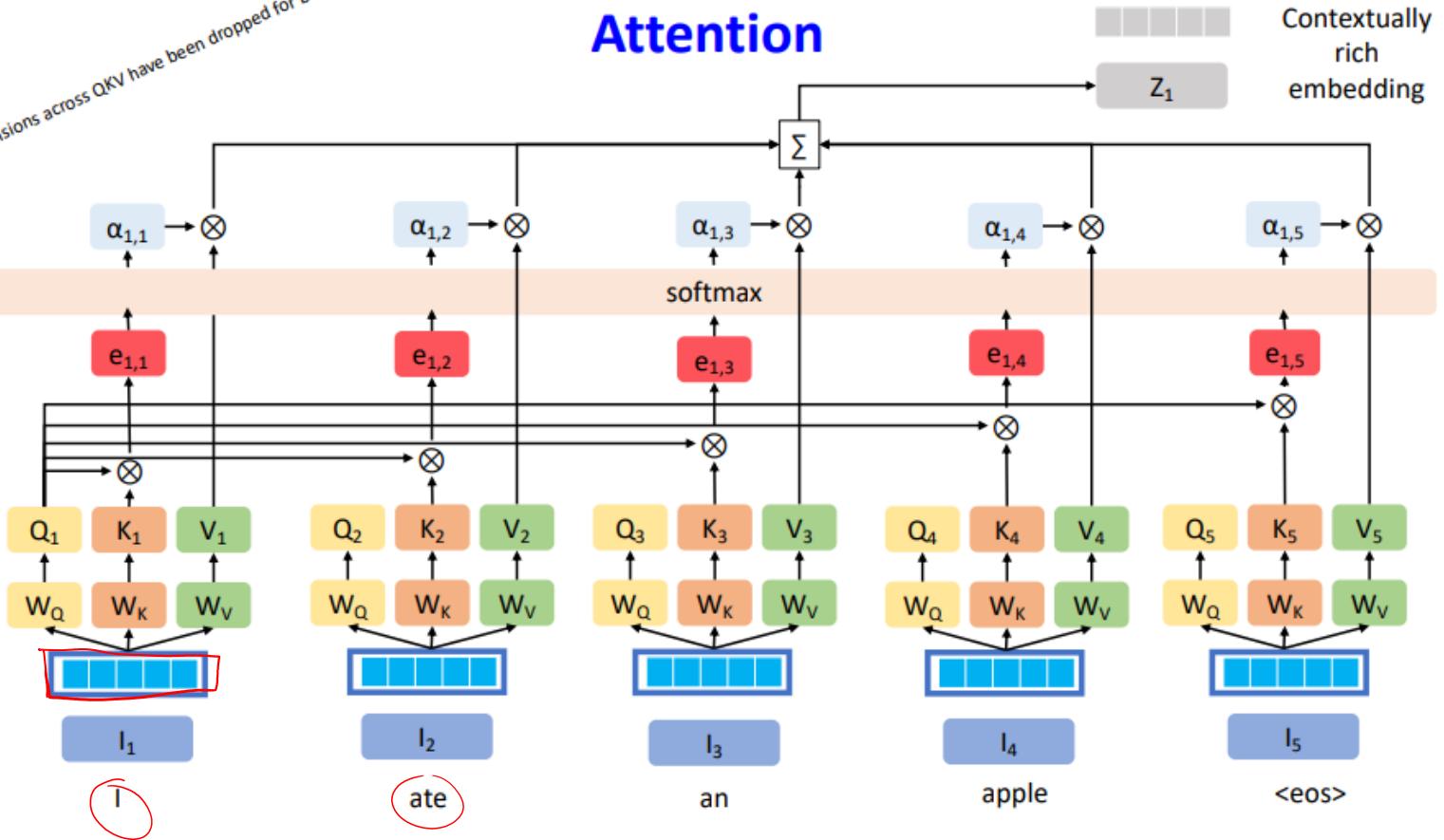


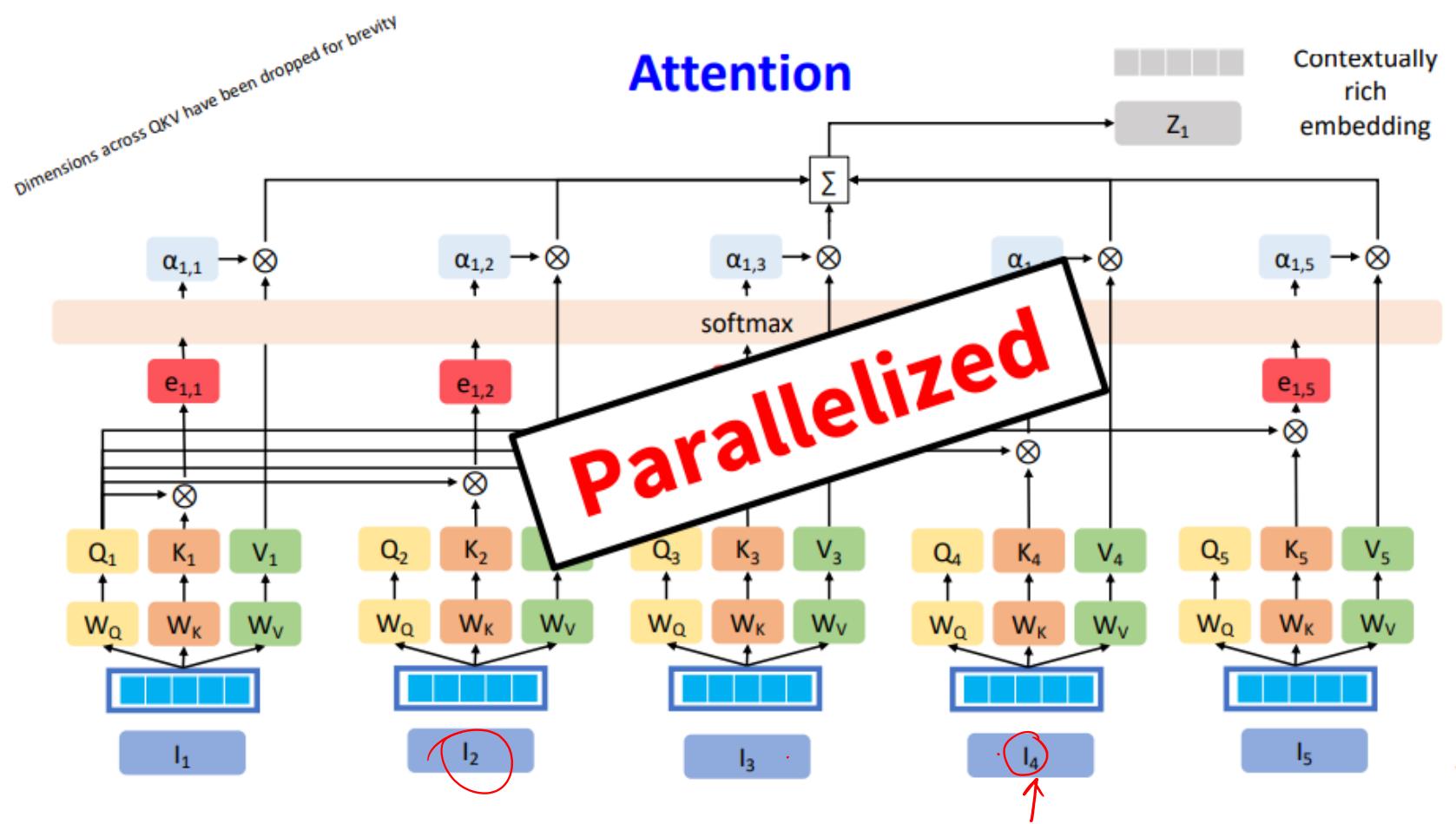
## Self attention

Dimensions across QKV have been dropped for brevity

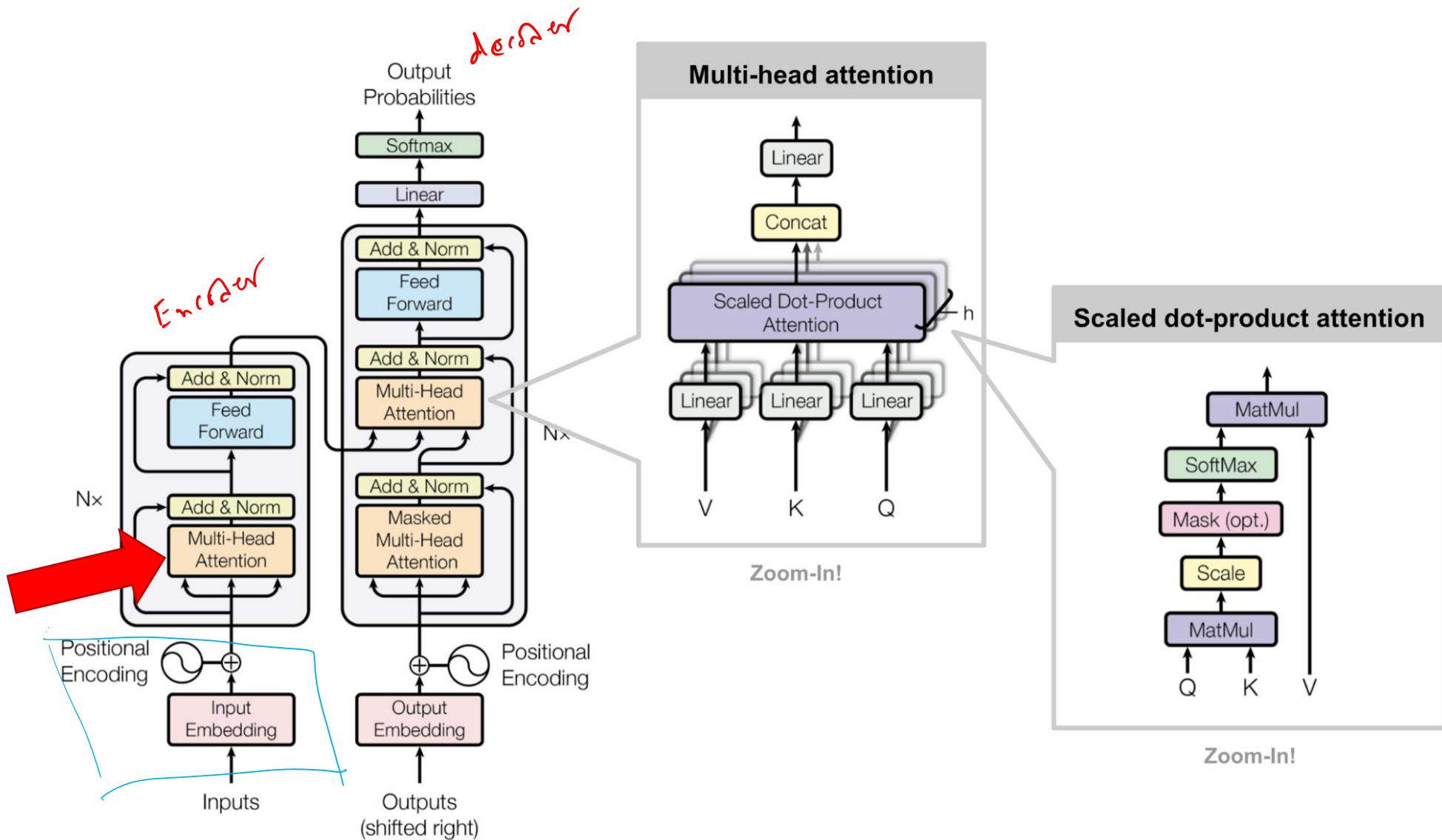
$\Sigma \alpha_k$

## Attention





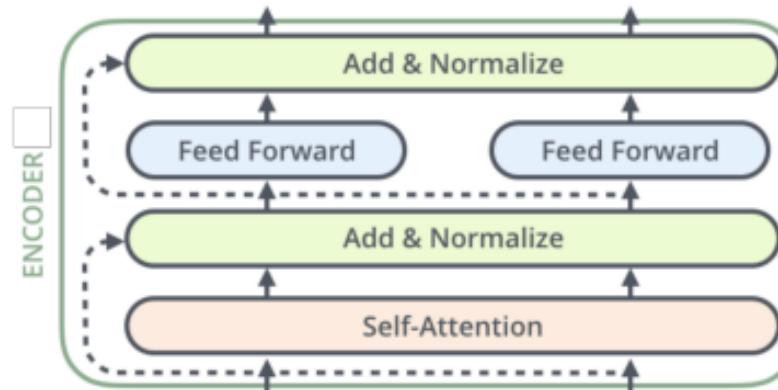
# Multi-Head Attention



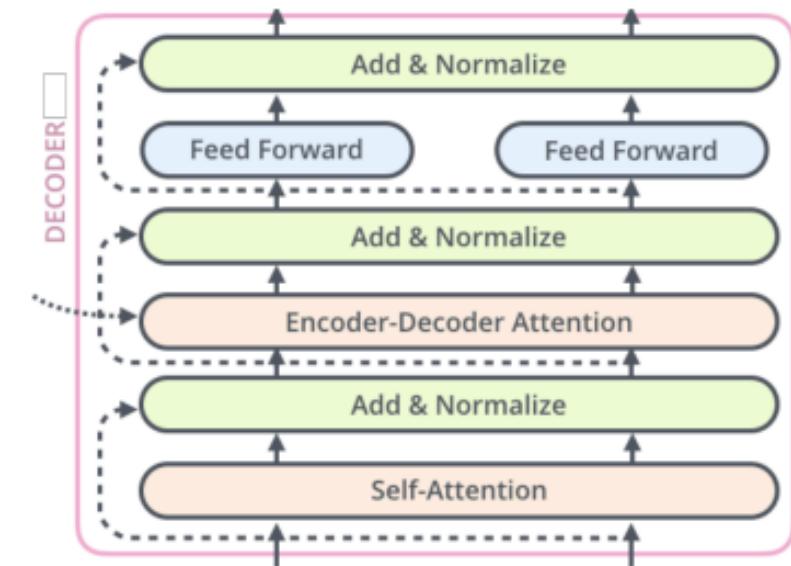
# Transformers

- An attention model with DL best practices!
- Originally introduced for machine translation, and now widely adopted for non-recurrent sequence encoding and decoding

**Transformer encoder**

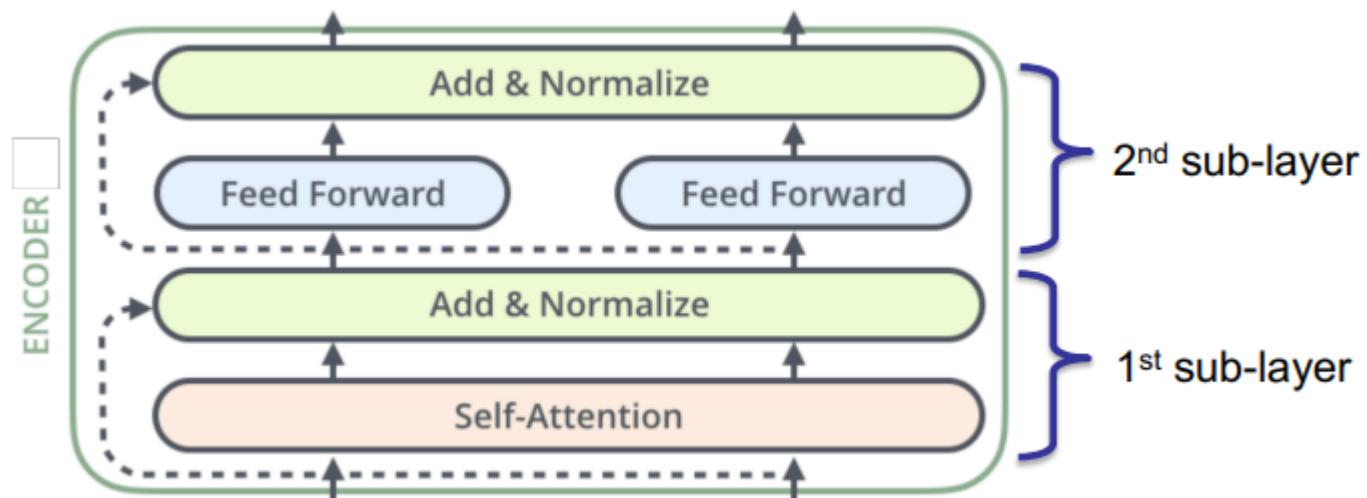


**Transformer decoder**



# Transformer encoder

- Transformer encoder consists of two sub-layers:
- 1st : Multi-head scaled dot-product self-attention
- 2nd : Position-wise feed forward
- Each sub-layer is followed by layer normalization and residual networks ... and drop-outs are applied after each computation



# Transformer encoder

- Let's start from multi-head scaled dot-product self-attention:
- Scaled dot-product attention
- Multi-head attention
- self-attention (recap)

- First, non-normalized attention scores:

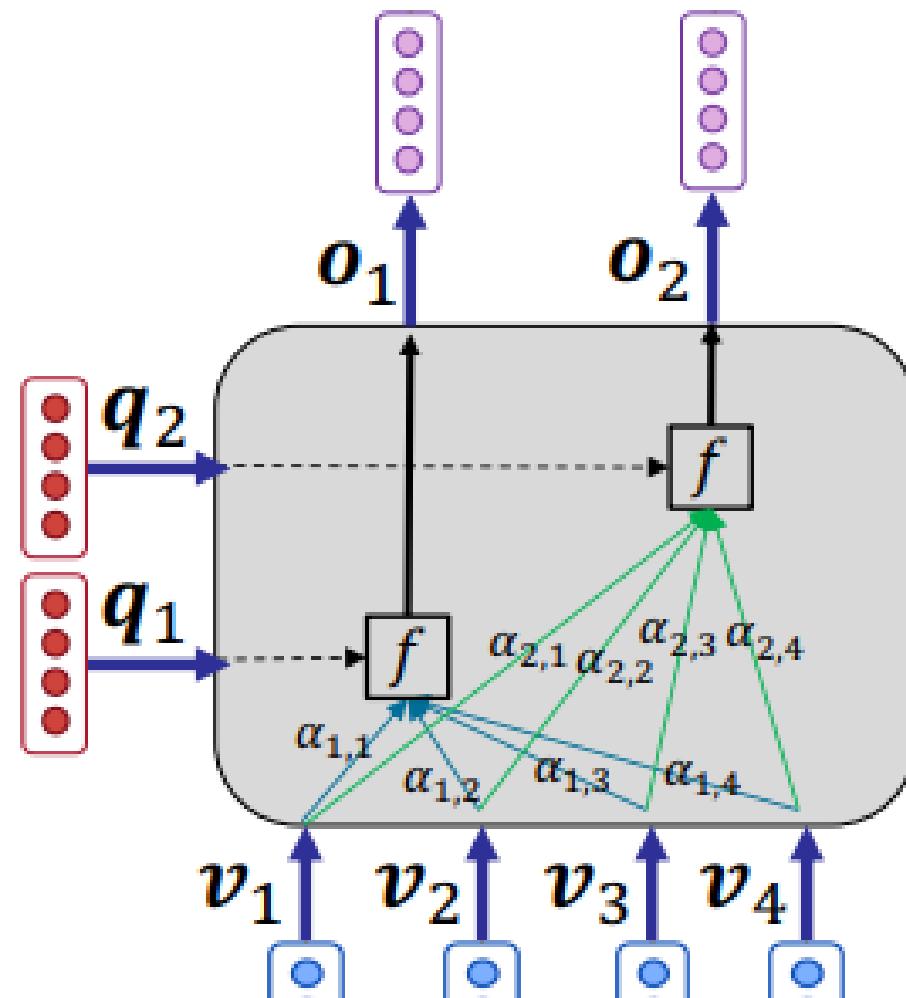
$$\tilde{\alpha}_{i,j} = \underline{q_i v_j^T}$$

- $d = d_q = d_v$  dimension of vectors
- has no parameter!

- Then, softmax over values:

$$\alpha_{i,j} = \text{softmax}(\tilde{\alpha}_i)_j$$

- Output (weighted sum):  $\mathbf{o}_i = \sum_{j=1}^{|V|} \alpha_{i,j} \mathbf{v}_j$



# Scaled dot-product attention

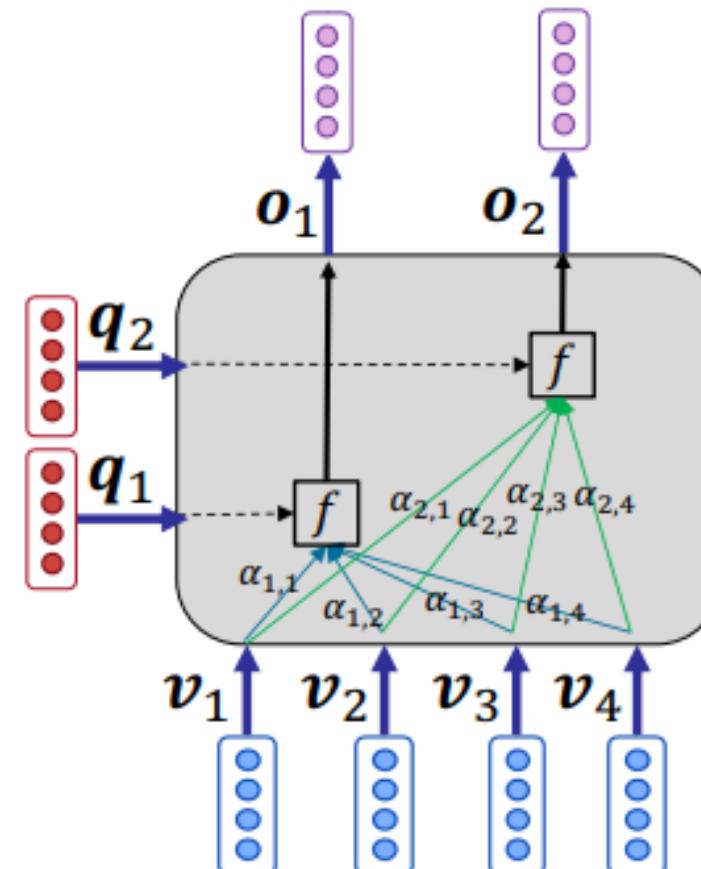
- Problem with basic doc-product attention:
  - As  $d$  gets large, the variance of  $\tilde{\alpha}_{i,j}$  increases ...
  - ... this makes softmax very peaked for some values  $\tilde{\alpha}_i$  ...
  - ... and hence its gradient gets smaller
- Solution: normalize/scale  $\tilde{\alpha}_{i,j}$  by size of  $d$

## Scaled dot-product attention

- Non-normalized attention scores:

$$\tilde{\alpha}_{i,j} = \frac{\mathbf{q}_i \mathbf{v}_j^T}{\sqrt{d}}$$

- Softmax over values:  $\alpha_{i,j} = \text{softmax}(\tilde{\alpha}_i)_j$
- Output (weighted sum):  $\mathbf{o}_i = \sum_{j=1}^{|V|} \alpha_{i,j} \mathbf{v}_j$



# Problem with (single-head) attention

- In all attention networks so far, the final attention of query  $q$  on value vectors  $V$  are normalized with softmax

- Recall that softmax makes the maximum value much higher than the other

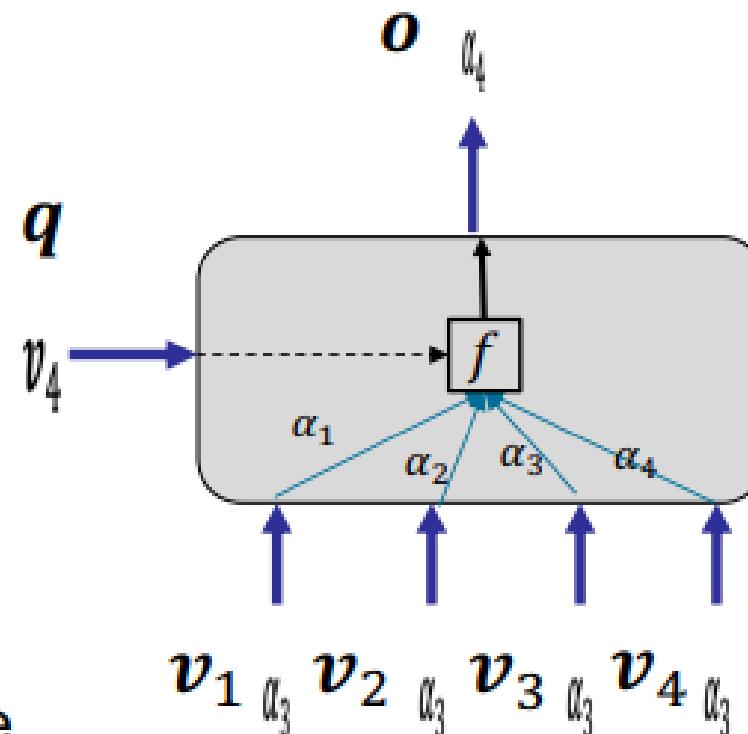
$$z = [1 \ 2 \ 5 \ 6] \rightarrow \text{softmax}(z) = [0.004 \ 0.013 \ 0.264 \ 0.717]$$

- Common in language, a word may be related to several other words in sequence, each through a specific concept

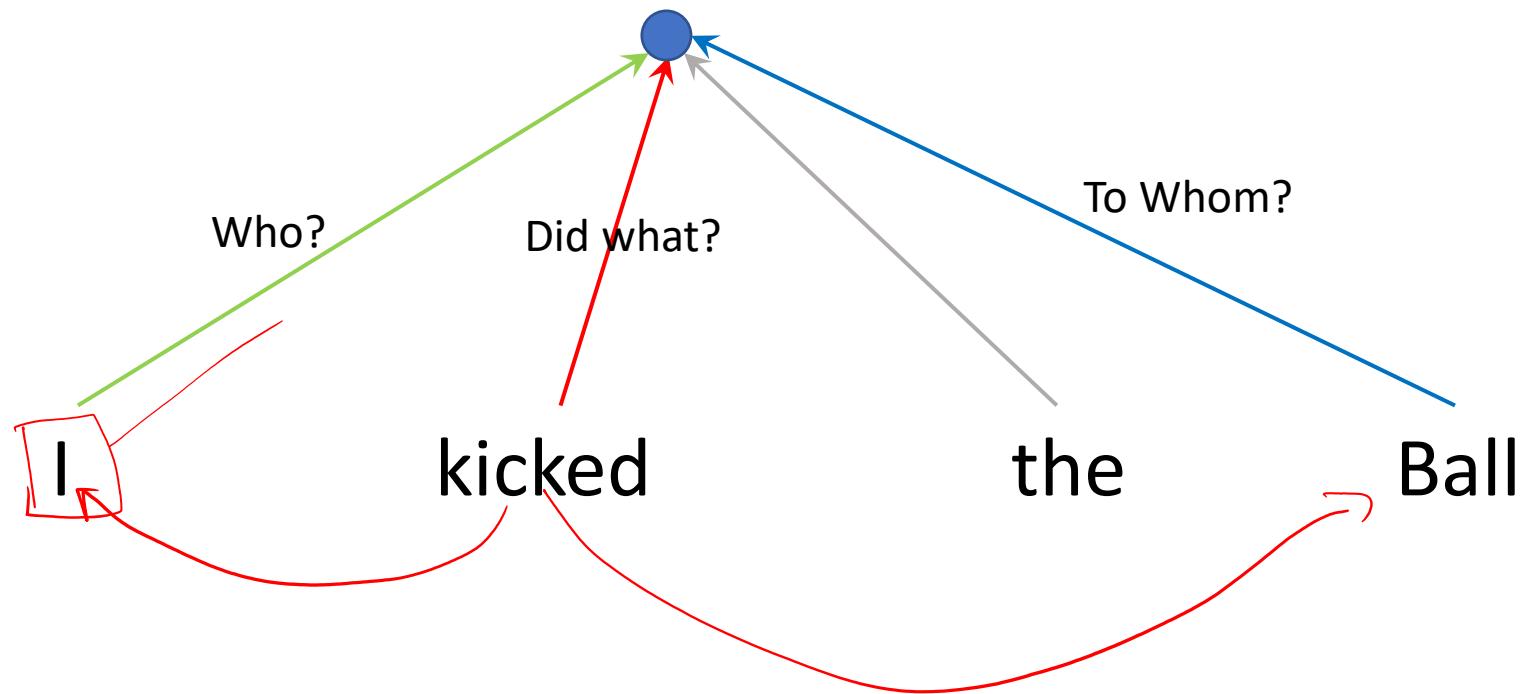
- Like the relations of a verb to its subject and to its object

- However in a (single-head) attention network, all concepts are aggregated in one attention set

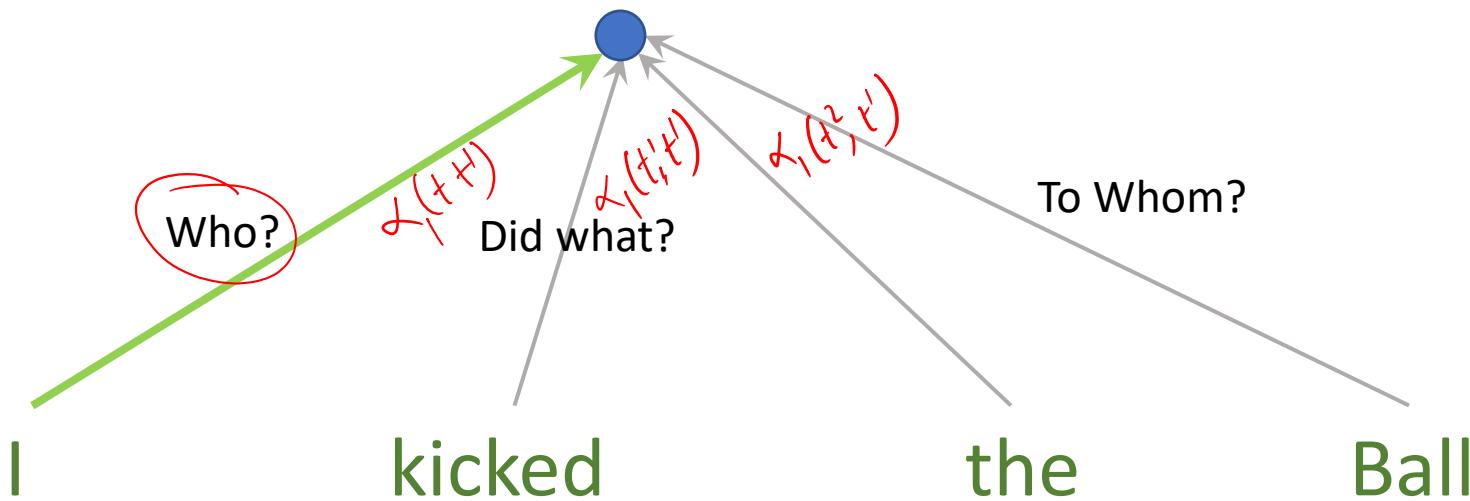
- Due to softmax, value vectors must compete for the attention of query vector  $\rightarrow$  softmax bottleneck



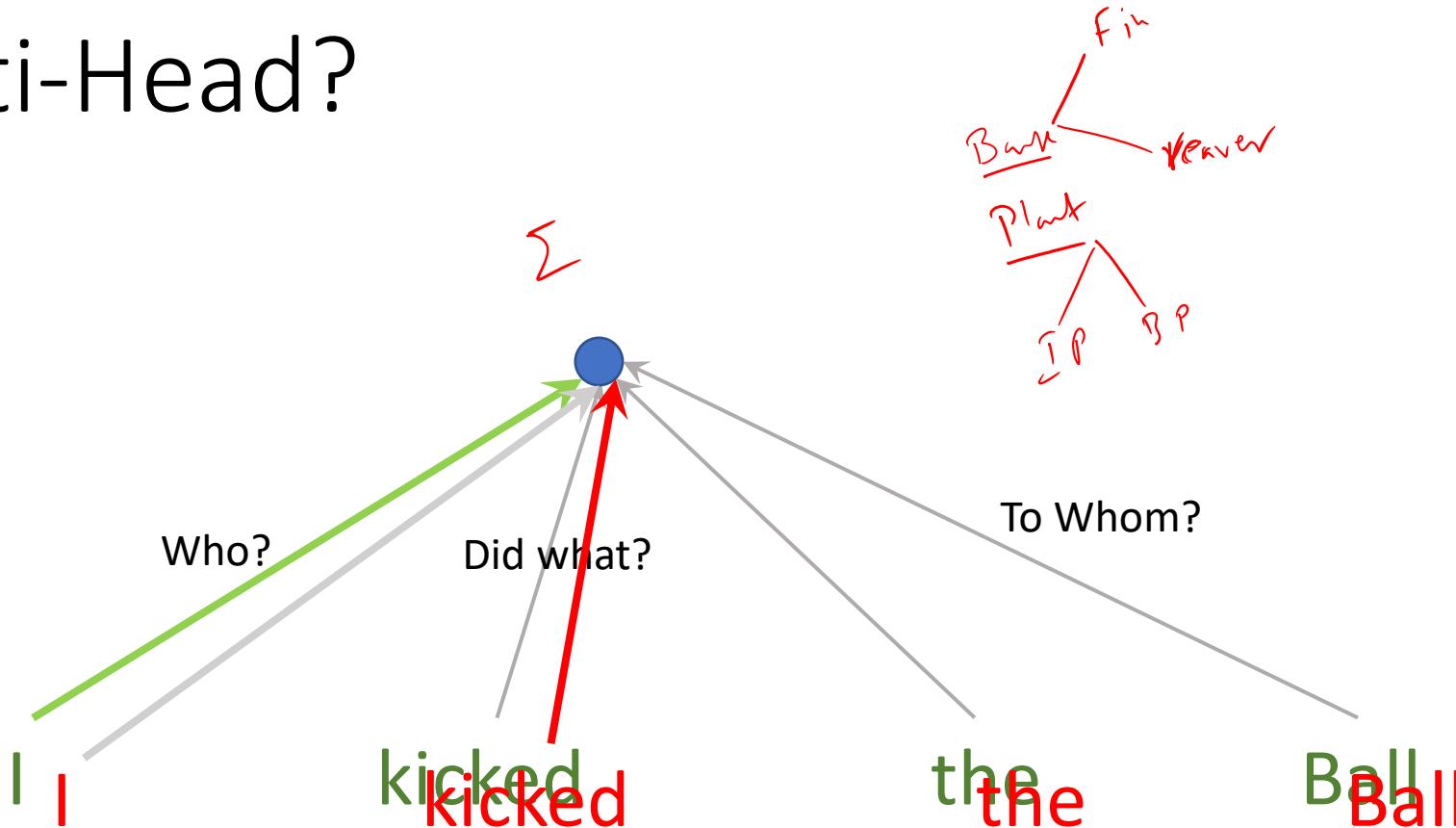
# Multi-Head?



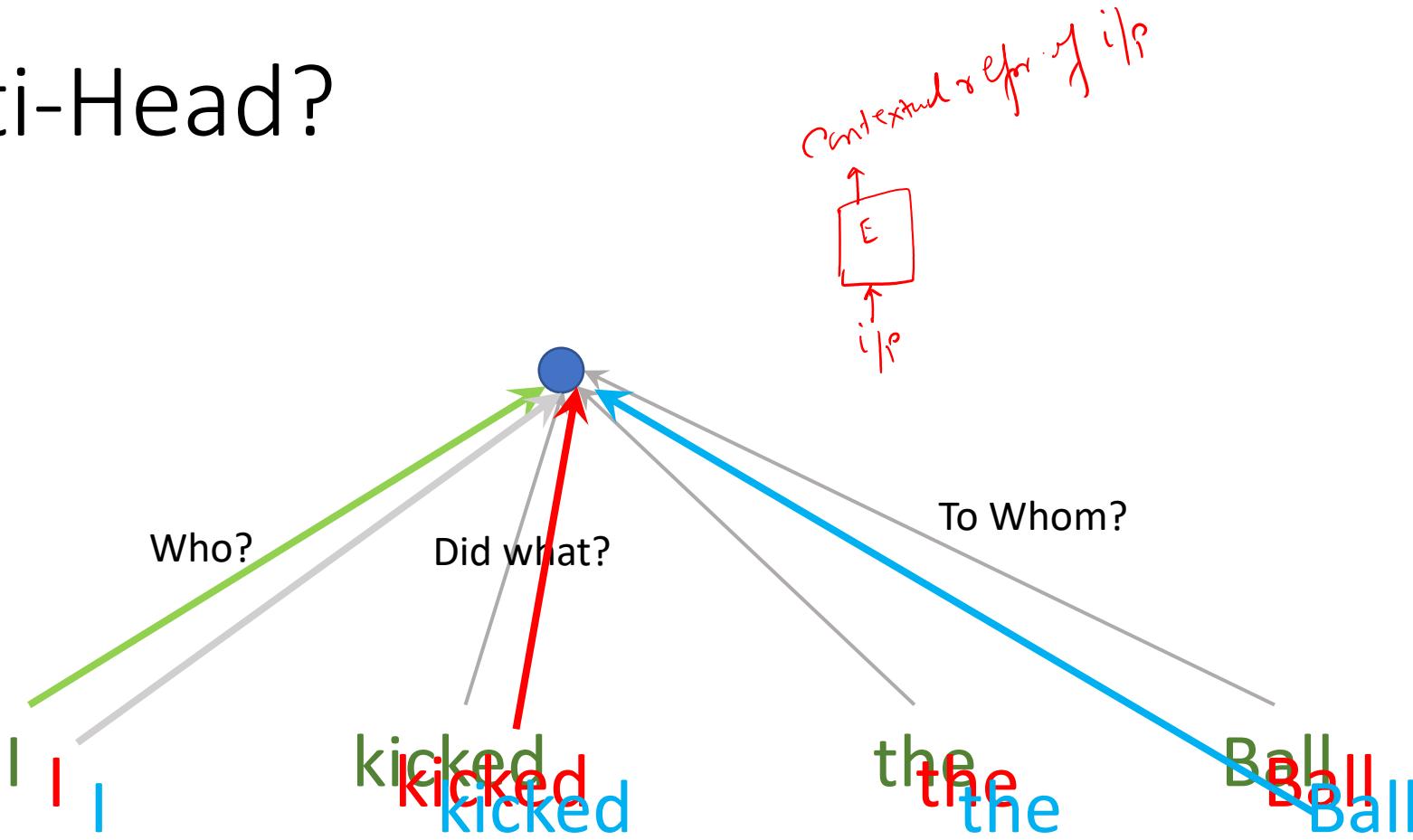
# Multi-Head?



# Multi-Head?

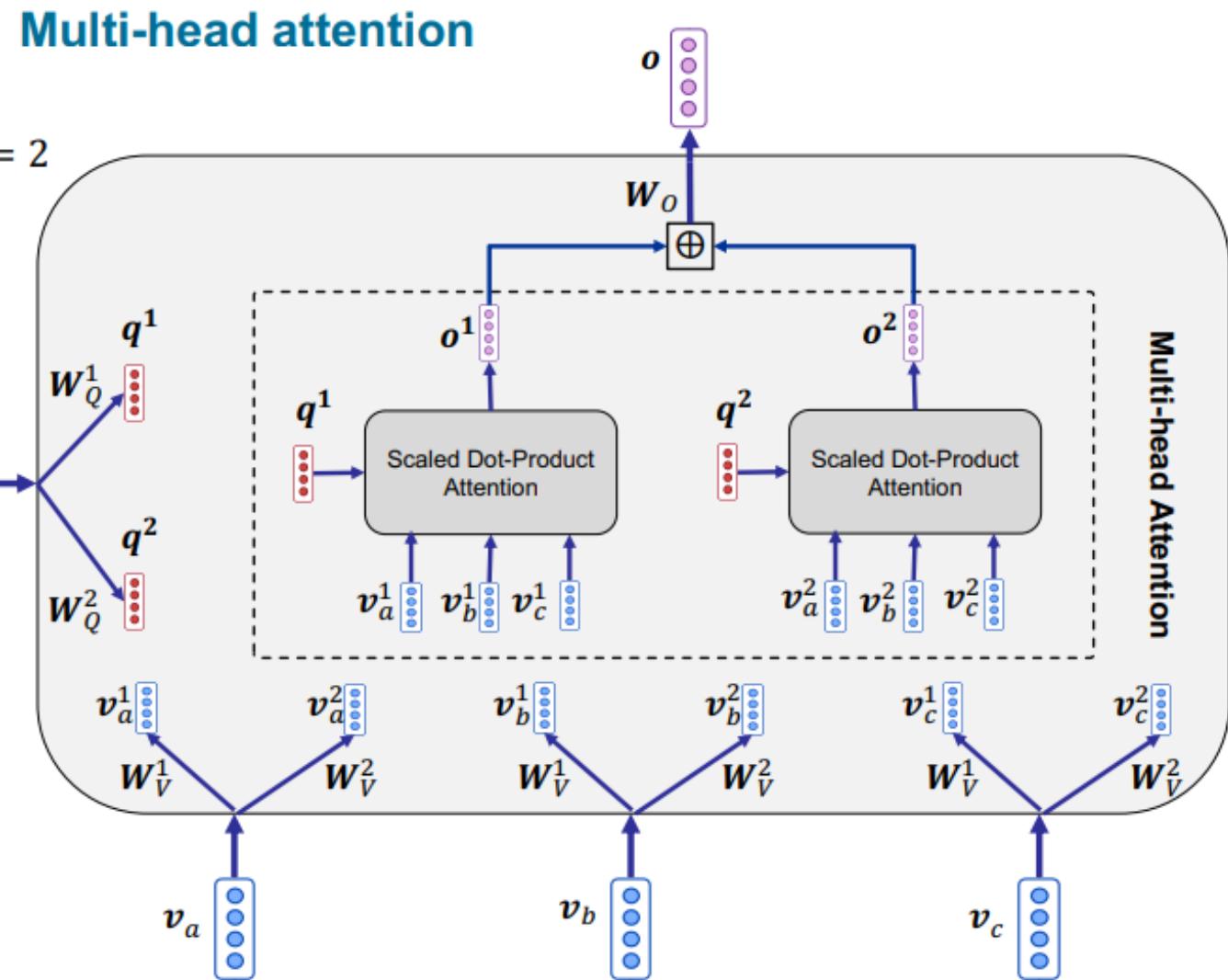


# Multi-Head?



# Multi-head attention

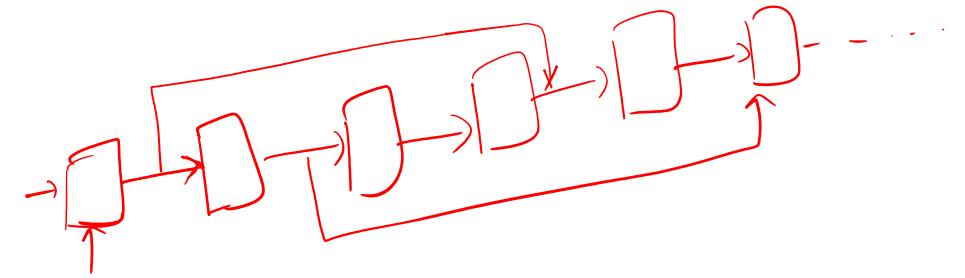
- Multi-head attention approaches this by calculating **multiple sets of attentions** between queries and values
- In multi-head attention, **each head** (and each subspace) can specialize on **capturing a specific kind of relations**



# Residuals

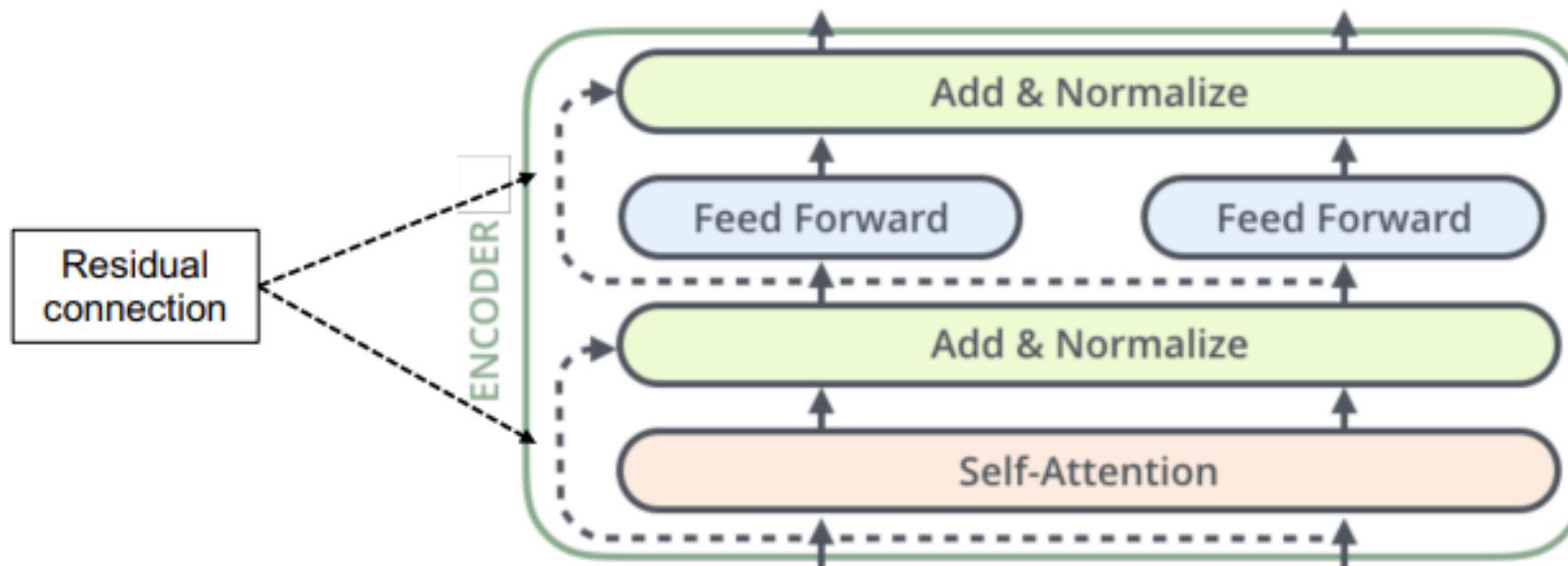
- Residual (short-cut) connection:

$$\text{output} = f(x) + x$$



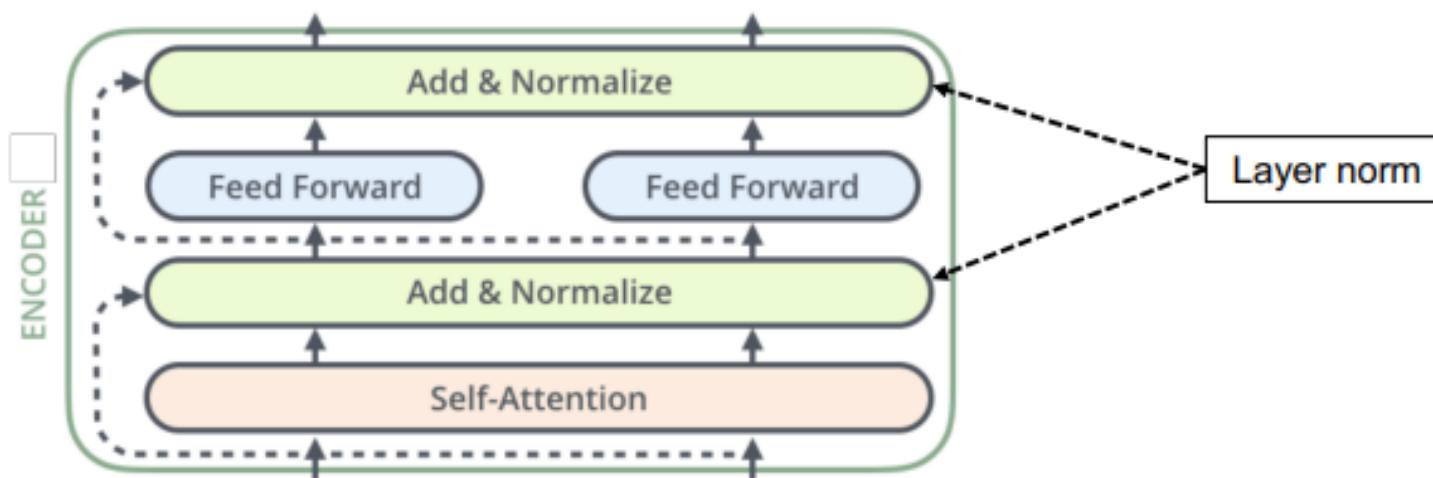
- Learn in detail:

- He, Kaiming; Zhang, Xiangyu; Ren, Shaoqing; Sun, Jian (2016). "Deep Residual Learning for Image Recognition". In proc. of CVPR
- Srivastava, Rupesh Kumar; Greff, Klaus; Schmidhuber, Jürgen (2015). "Highway Networks". <https://arxiv.org/pdf/1505.00387.pdf>



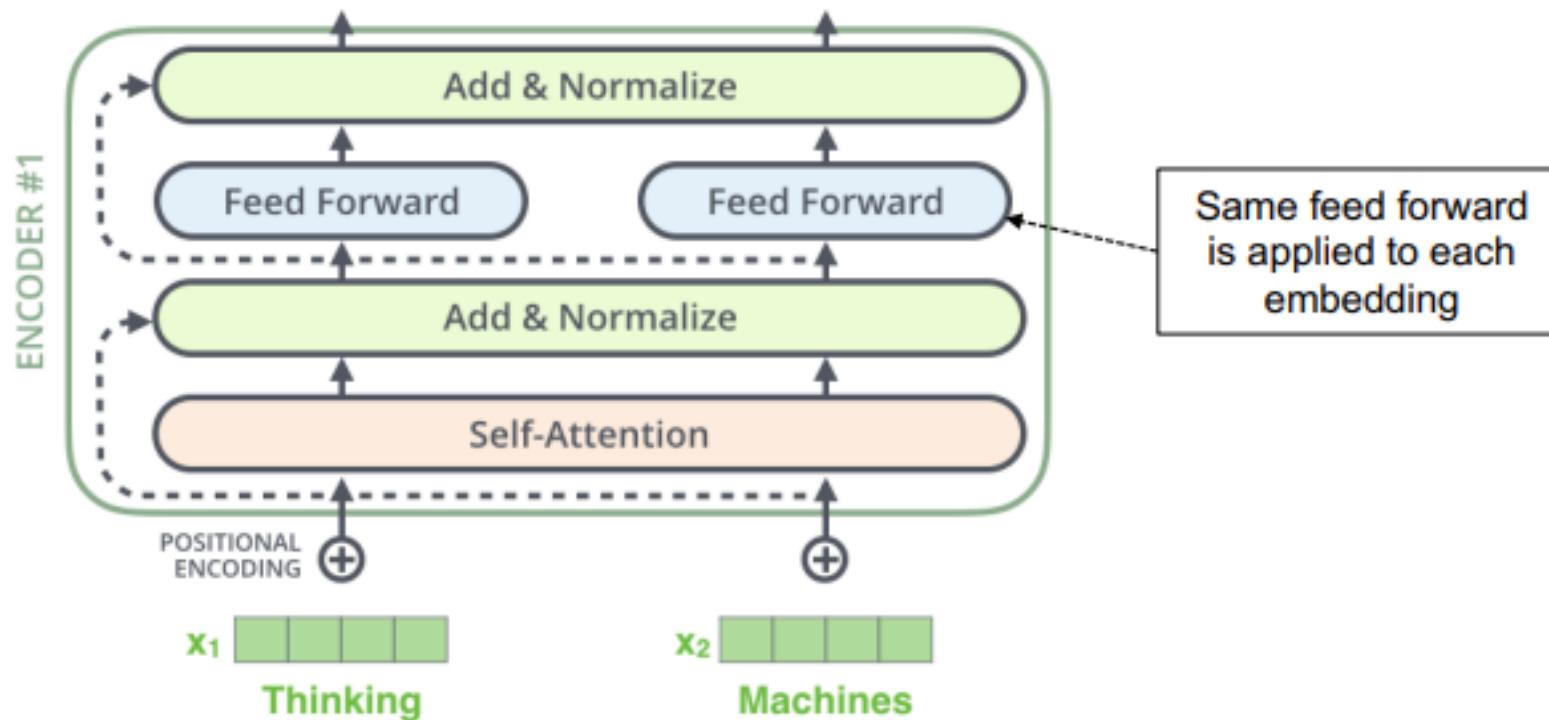
# Layer normalization

- Layer normalization changes each vector to have mean 0 and variance 1 ...
  - ... and learns two more parameter vectors per layer that set new means and variances for each dimension of the vectors
- Learn in detail:
  - Batch Normalization: Deep Learning book section 8.7.1  
<http://www.deeplearningbook.org/contents/optimization.html>
  - Talk by Goodfellow <https://www.youtube.com/watch?v=Xogn6veSxA&feature=youtu.be>
  - Paper: <https://arxiv.org/pdf/1607.06450.pdf>



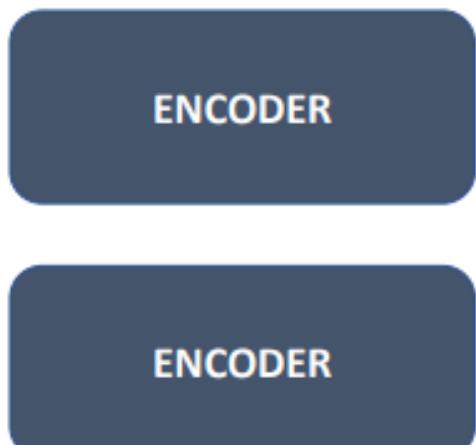
## Feed Forward on embedding

- In Transformers, a two-layer feed forward neural network (with ReLU) is applied to each embedding
  - With the feed forward network, the Transformers gain the capacity to learn non-linear transformations over each (contextualized) embedding



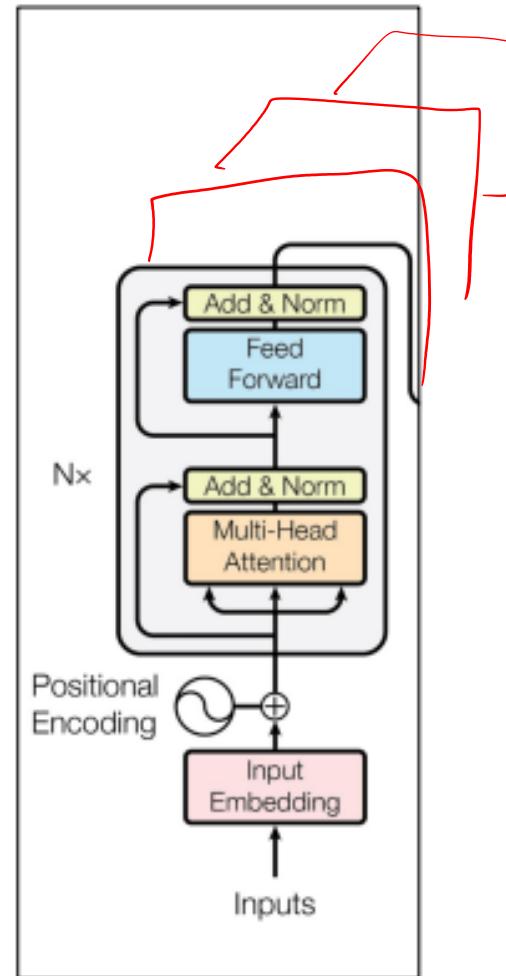
# Encoders

Encoder

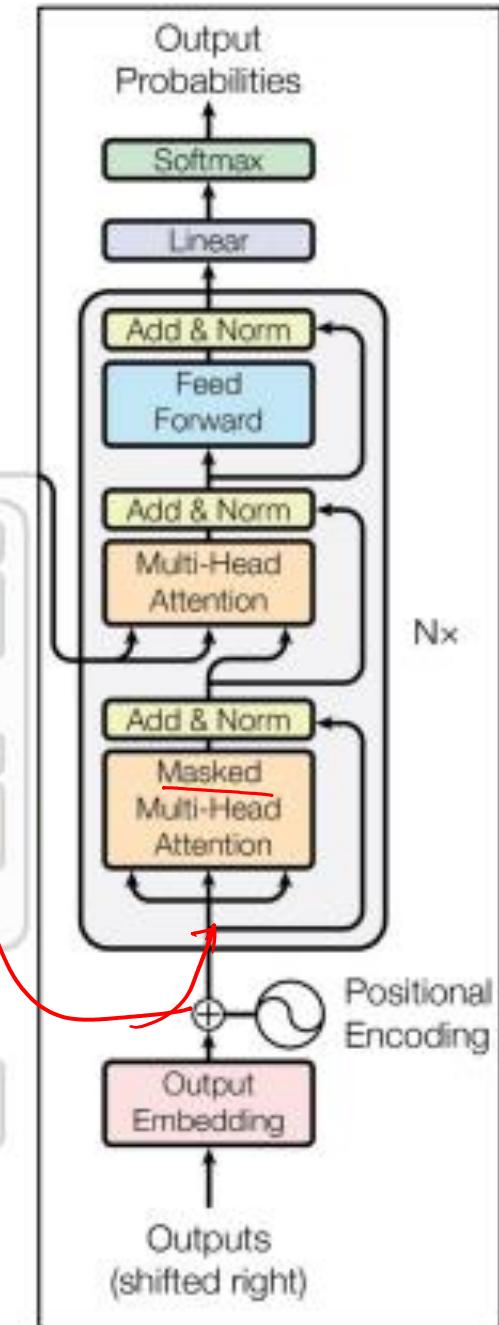
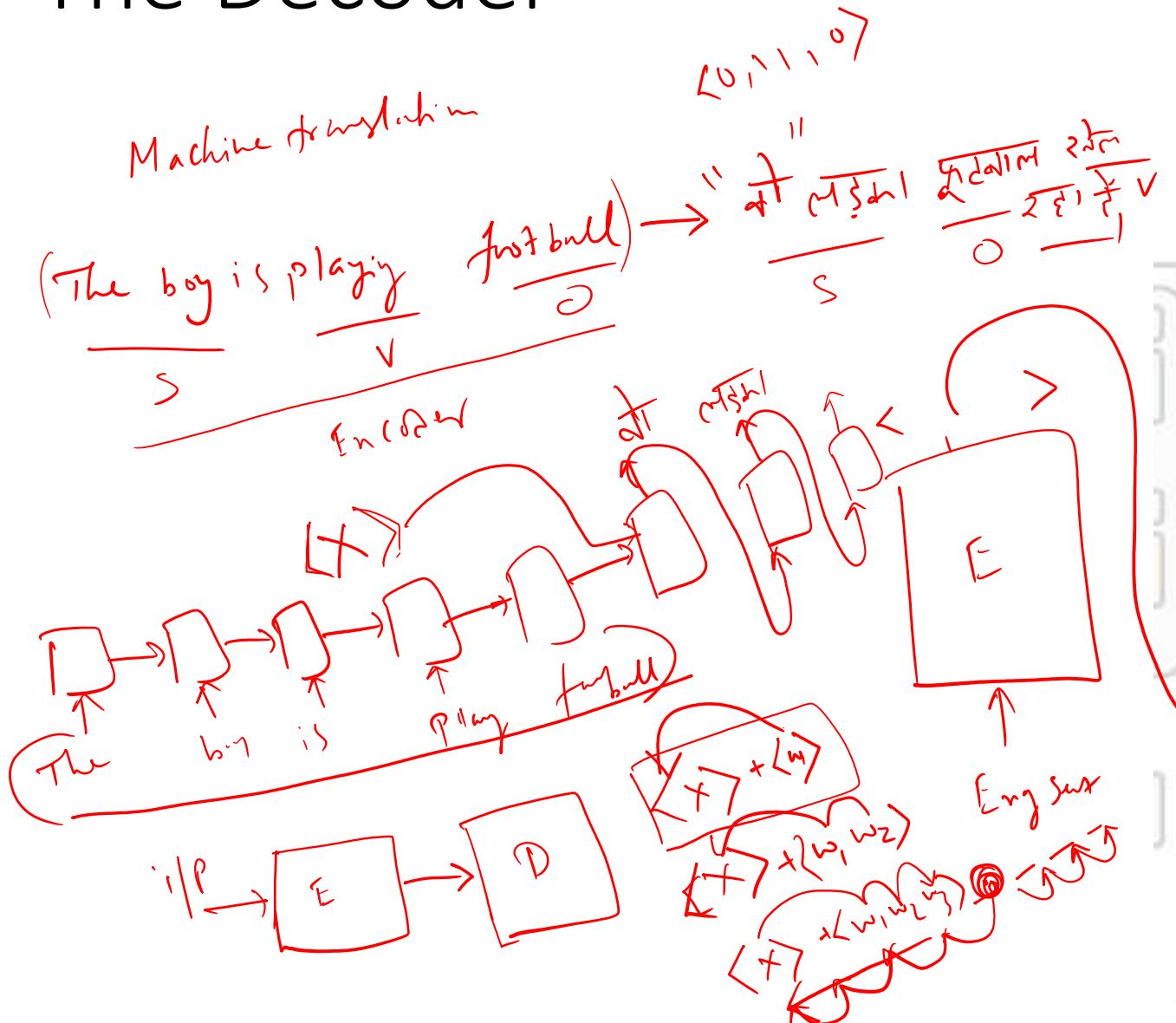


Input to Encoder<sub>i+1</sub>

Output from Encoder<sub>i</sub>



# The Decoder



$w_1, w_2, w_3, \langle w'_4, w'_5, w'_6 \rangle$

Autoregressive LM mostly  
→ BERT

Mixed LM

Permutation LM.  
↪ XLNet

# Transformer encoder – summary

- Multi-head self-attention model followed by a feed-forward layer
- Benefits (as in attentions)
  - No locality bias - A long-distance context has “equal opportunity”
  - Friendly with high parallel computations of GPU
- Look here for self-teaching and the PyTorch implementation:  
<http://nlp.seas.harvard.edu/2018/04/03/attention.html>
- also available on Google Colab