# Introduction to Deep Learning
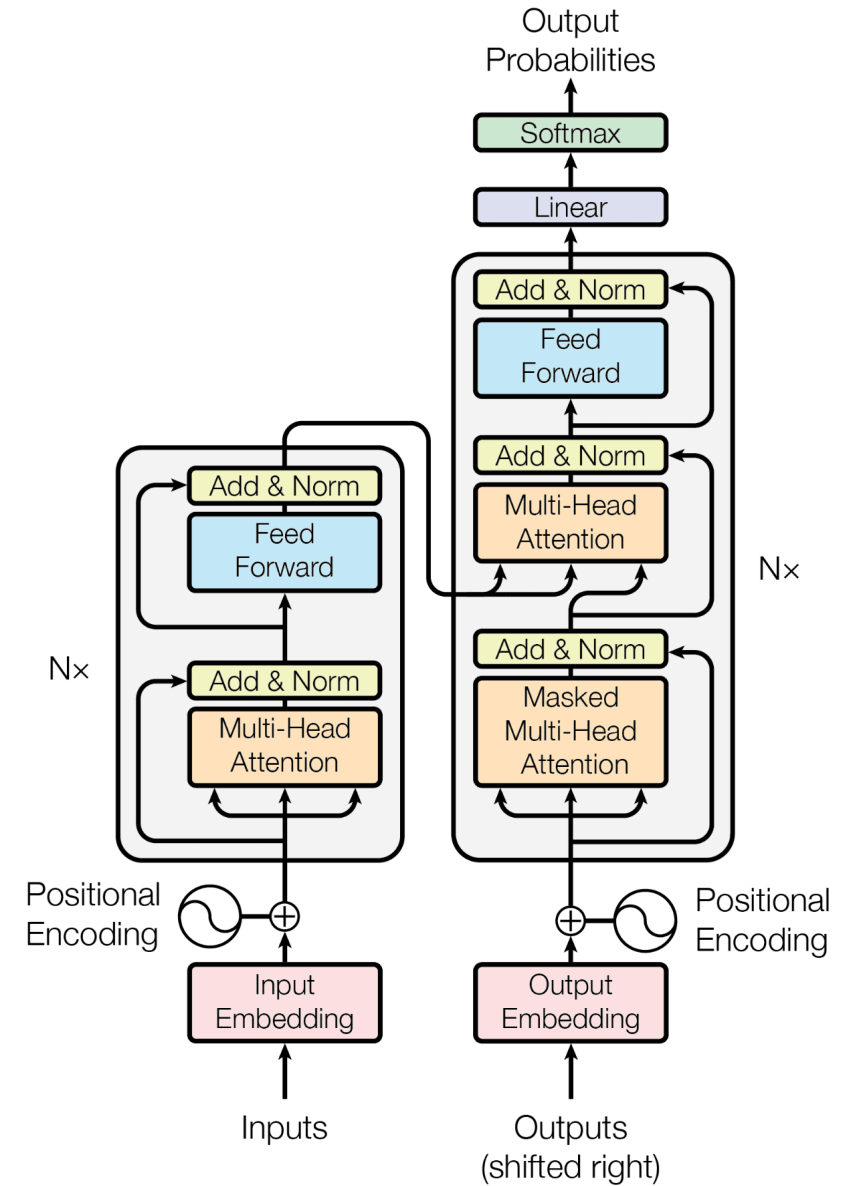## Transformers

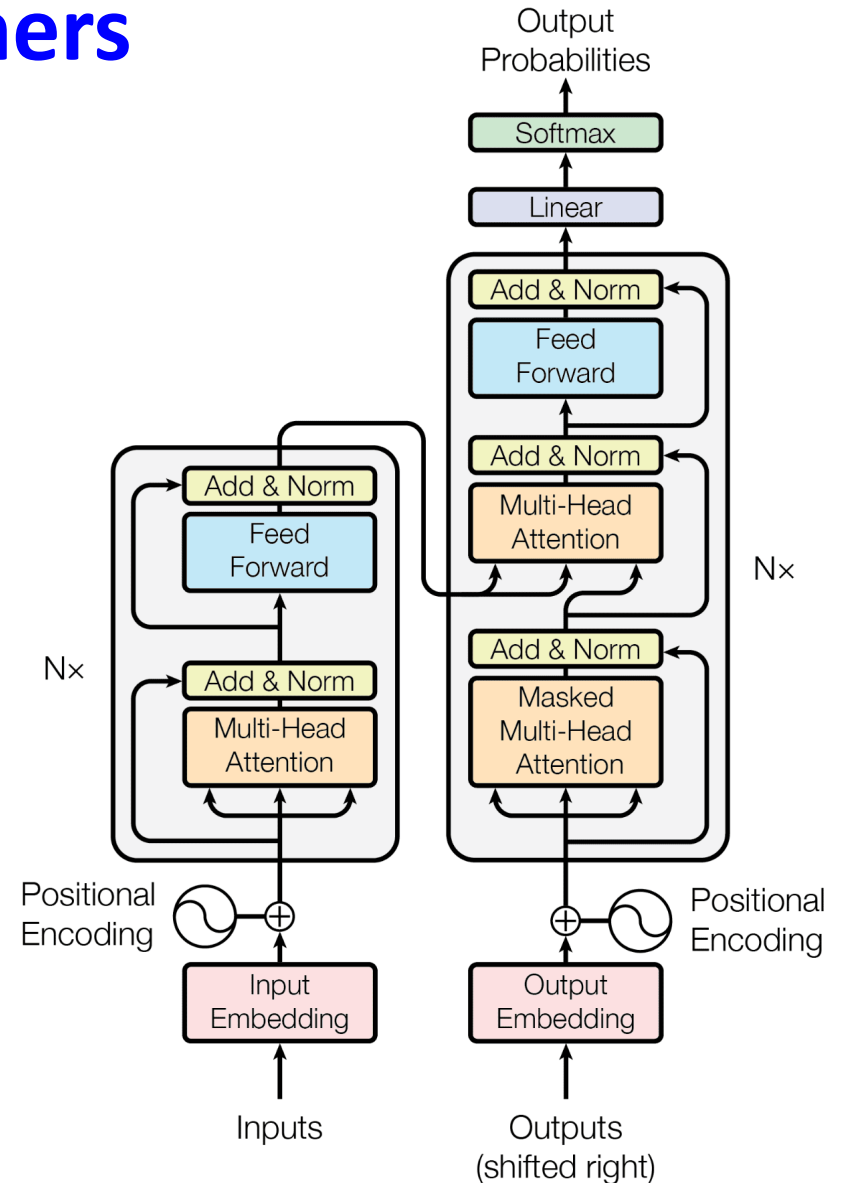**Shikhar Agnihotri**                    **Liangze Li**

**11-785, Fall 2023**

# Transformers

# Transformers

- Tokenizaton
- Input Embeddings
- Position Encodings
- Residuals
- Query
- Key
- Value
- Add & Norm
- Encoder
- Decoder

- Attention
- Self Attention
- Multi Head Attention
- Masked Attention
- Encoder Decoder Attention
- Output Probabilities  / Logits
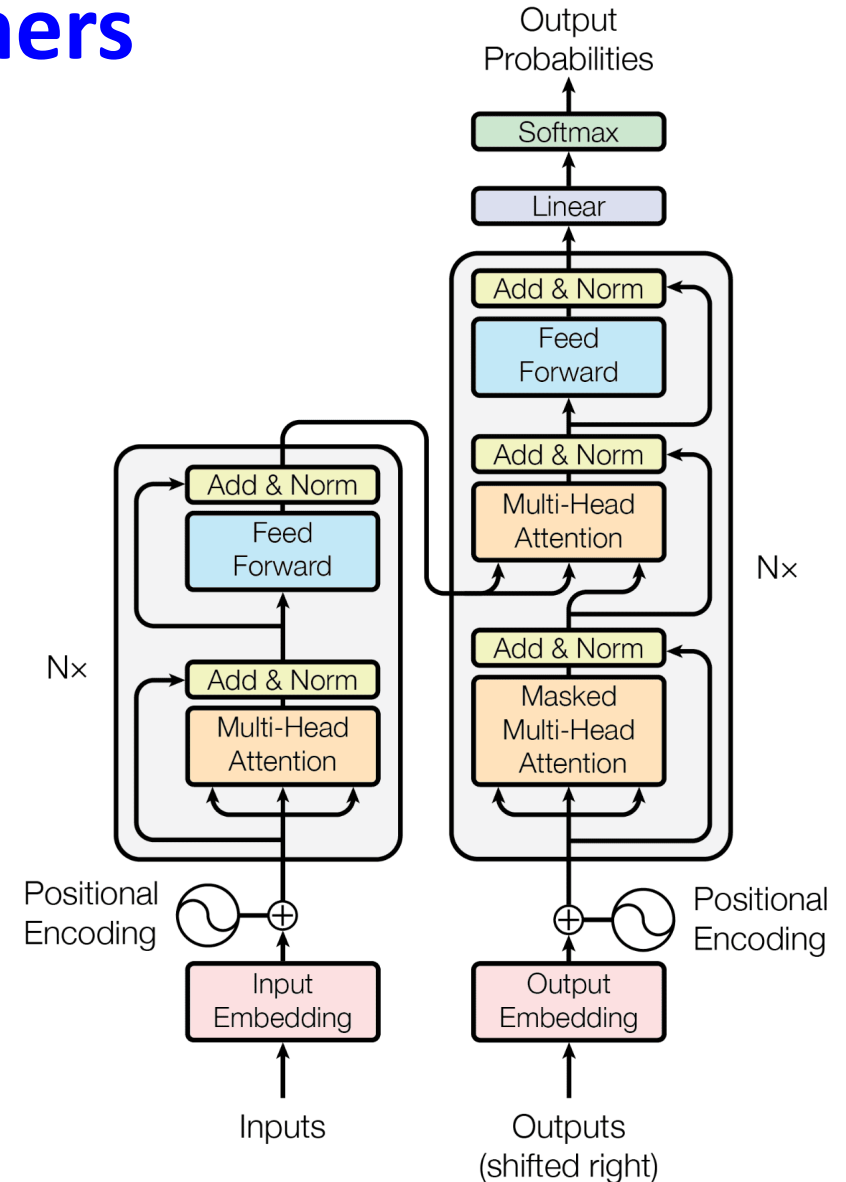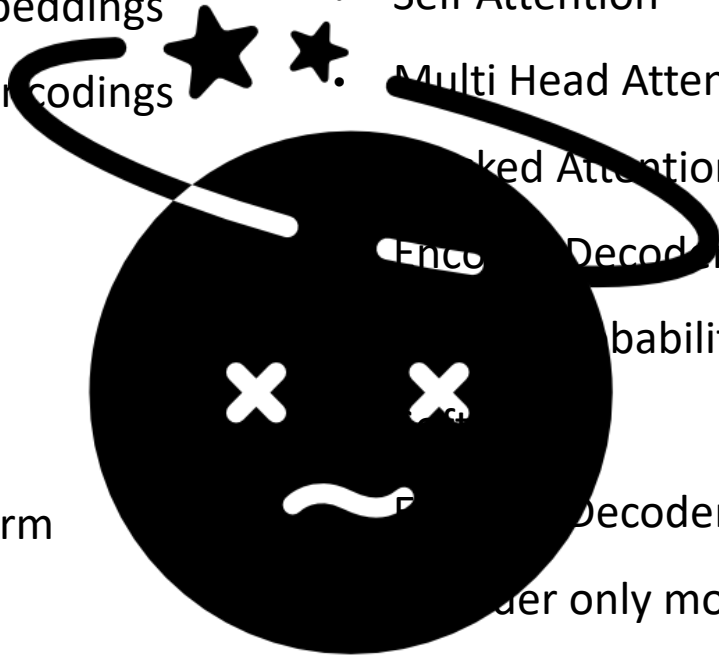- Softmax
- Encoder-Decoder models
- Decoder only models

# Transformers

- Tokenizaton
- Input Embeddings
- Position Encodings
- Residuals
- Query
- Key
- Value
- Add & Norm
- Encoder
- Decoder

- Attention
- Self Attention
- Multi Head Attention
- ~~Masked Attention~~
- ~~Encoder~~ Decoder Attention
- ~~Probabilities / Logits~~
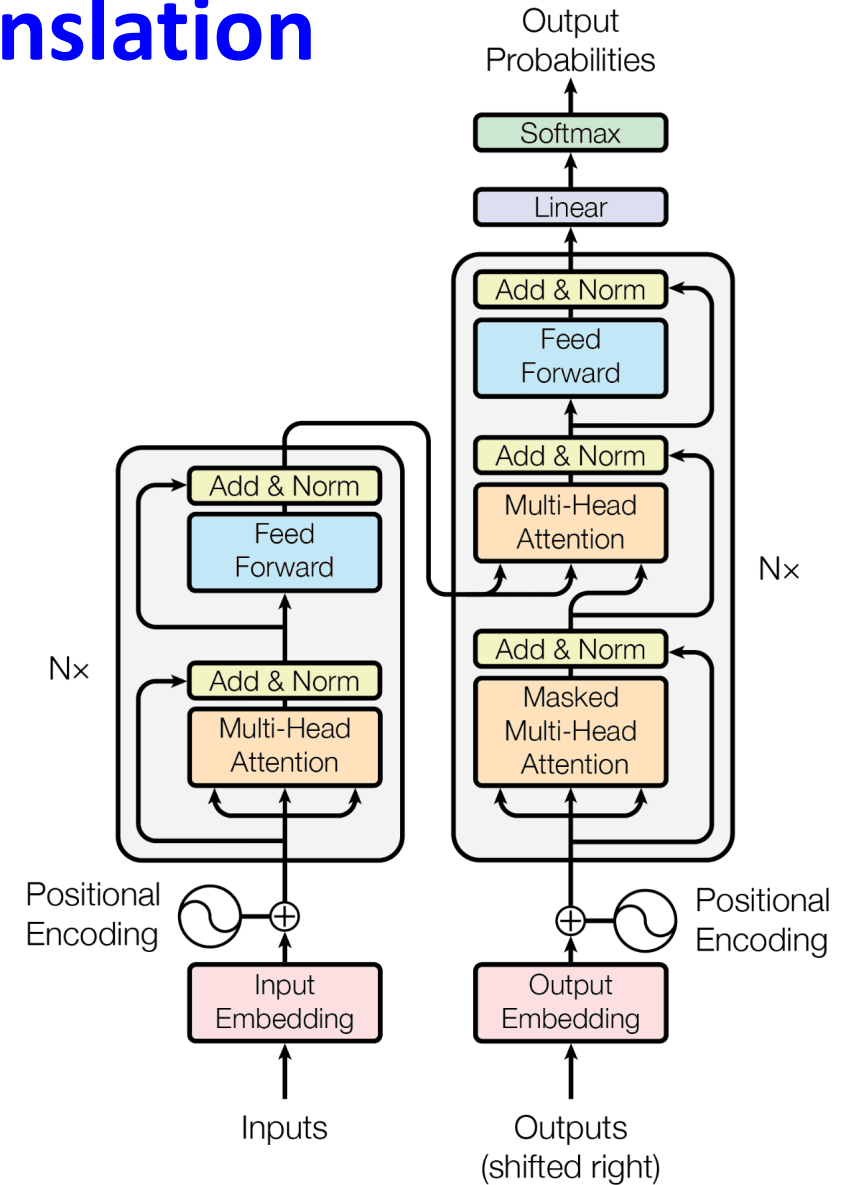- ~~Decoder models~~
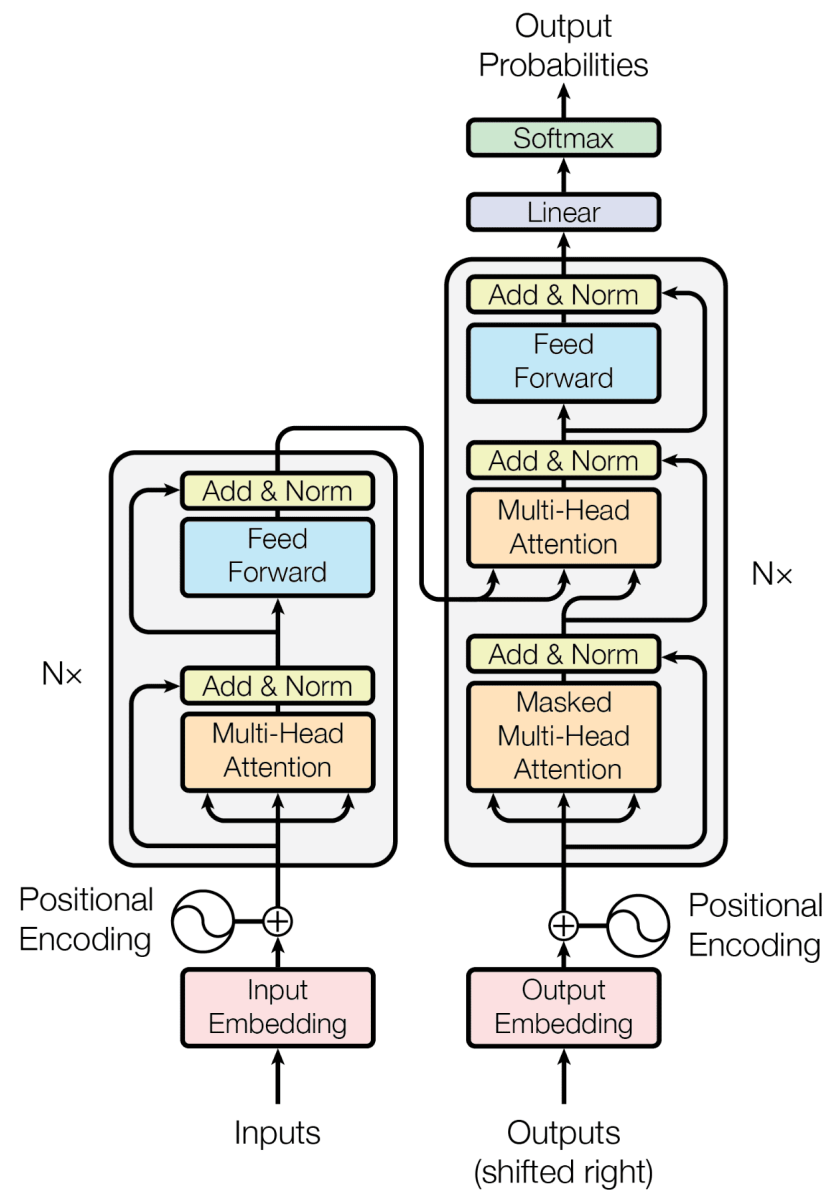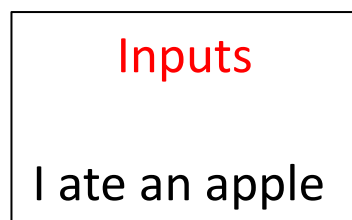- ~~Encoder only models~~

# Machine Translation



Targets

Ich have einen apfel gegessen

Inputs

I ate an apple

# Inputs

## Processing Inputs

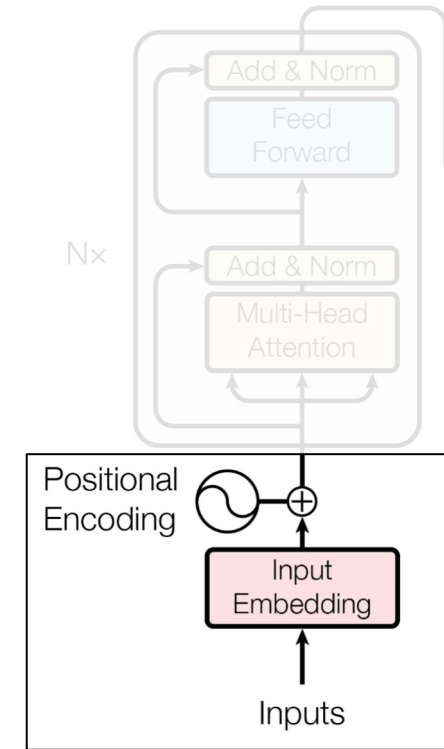| Inputs |
| --- |
| I ate an apple |

# Inputs

| I | ate | an | apple | <eos> |

Tokenizer

I ate an apple

Generate Input Emebeddings

# Inputs

$d_{model}$



Embedding Layer

| I | ate | an | apple | <eos> |

Tokenizer

I ate an apple
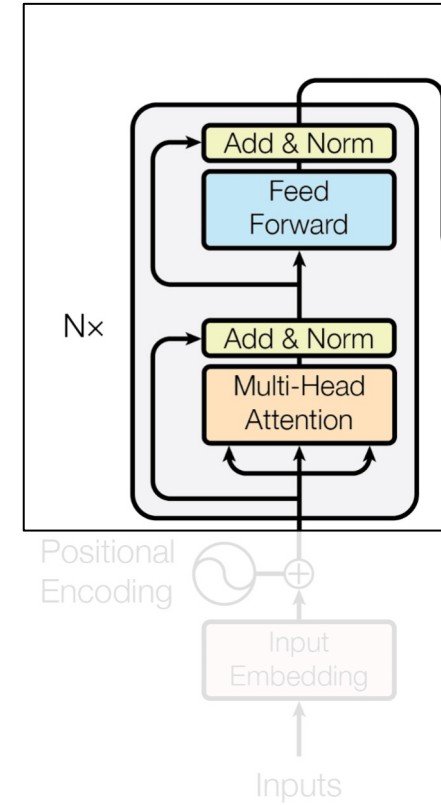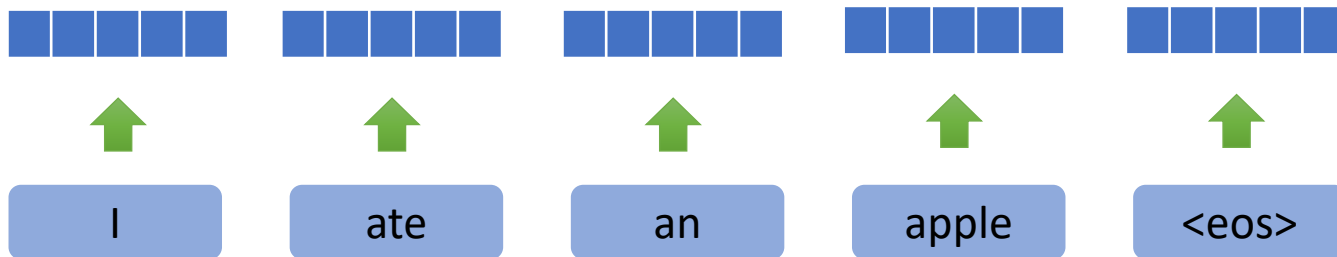
Generate Input Emebeddings

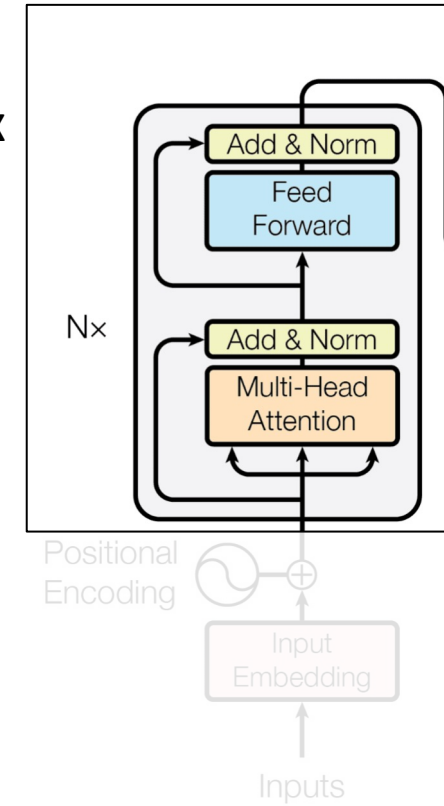Positional Encoding

Input Embedding

Inputs

Add & Norm

Feed Forward

Nx

Add & Norm

Multi-Head Attention

# Encoder

WHERE IS THE CONTEXT ?

Add & Norm

Feed Forward

N×

Add & Norm

Multi-Head Attention

Positional Encoding

Input Embedding

Inputs

I | ate | an | apple | <eos>

# Encoder

BLACK BOX
OF SORTS

I ate an apple <eos>

# Encoder

**BLACK BOX OF SORTS**

**LEARN TO ADD CONTEXT**

I  ate  an  apple  <eos>

Add & Norm

Feed Forward

N×

Add & Norm

Multi-Head Attention

Positional Encoding

Input Embedding

Inputs

11

# Encoder

**CONTEXTUALLY RICH EMBEDDINGS**

**BLACK BOX OF SORTS**

**LEARN TO ADD CONTEXT**

| I | ate | an | apple | <eos> |

# Encoder

$$\alpha_{[ij]} \ ?$$

**CONTEXTUALLY RICH EMBEDDINGS**

**BLACK BOX OF SORTS**

**LEARN TO ADD CONTEXT**

| I | ate | an | apple | \<eos\> |

Add & Norm

Feed Forward

N×

Add & Norm

Multi-Head Attention

Positional Encoding

Input Embedding

Inputs

# Encoder

$$\alpha_{[ij]} \; ? \; \sum \prod \; ?$$

CONTEXTUALLY RICH EMBEDDINGS

BLACK BOX OF SORTS

LEARN TO ADD CONTEXT

I    ate    an    apple    <eos>

Add & Norm
Feed Forward
Add & Norm
Multi-Head Attention
N×
Positional Encoding
Input Embedding
Inputs

# Attention

$$\alpha_{[ij]} \ ?$$

From lecture 18:

$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

# Attention

$$\alpha_{[ij]} \; ?$$

From lecture 18:

$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$
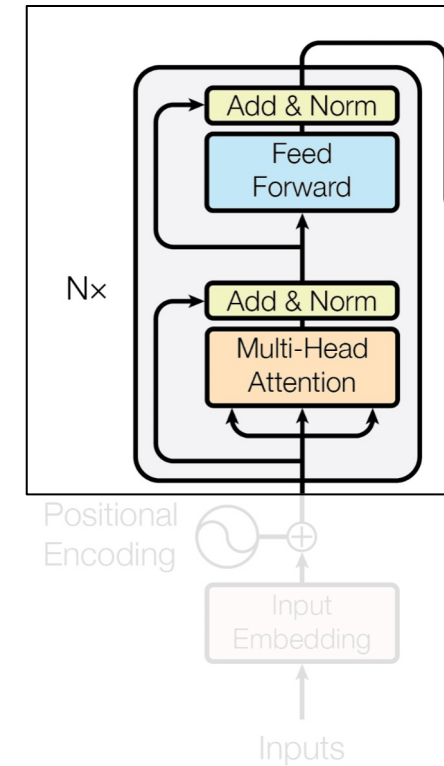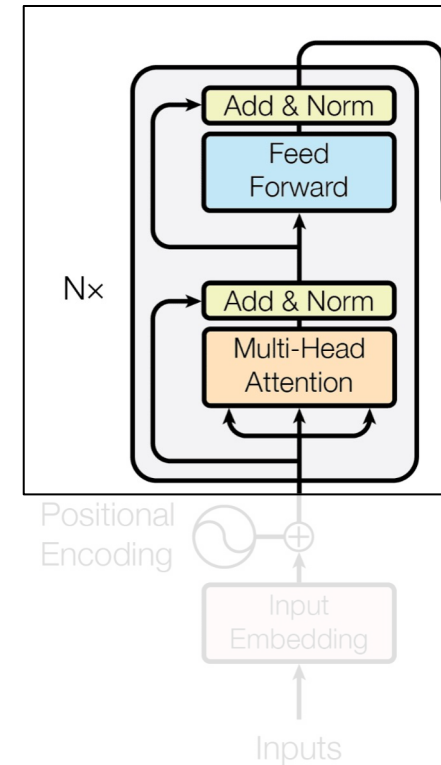
- Query

- Key

- Value

# 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

**Database**

**{Key, Value store}**

{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

{Key, Value store}

{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

{Key, Value store}

{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

{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: "Order details of order_104"}

OR

{Query: "Order details of order_106"}

# Query, Key & Value

**Done at the same time !!**

{Key, Value store}

{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

{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
1. Search for info

## Key
1. Interacts directly with Queries
2. Distinguishes one object from another
3. Identify which object is the most relevant and by how much

## Value
1. Actual details of the object
2. More fine grained

23

# Attention



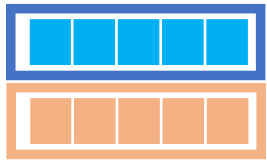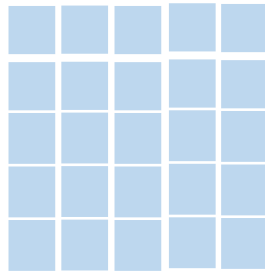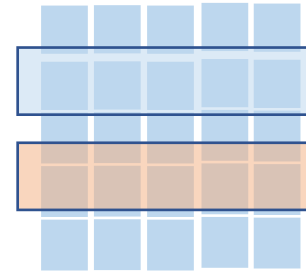Query  Key Value Store  Key  Value

# Attention



Query
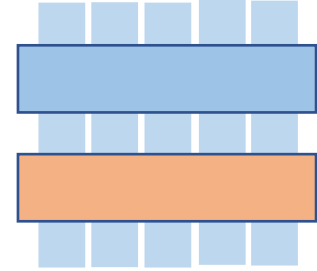
Key Value Store

Key

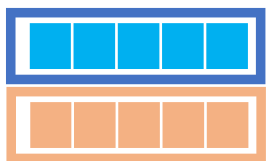Value

# Attention
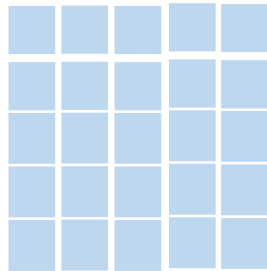
Done at the same time !!

Query

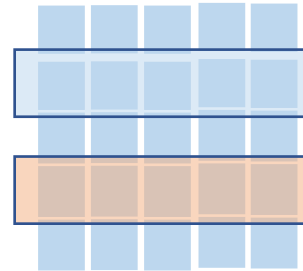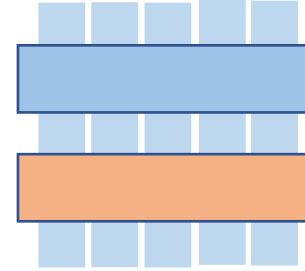Key Value Store

Key

Value

# Attention

**Parallelizable !!!**



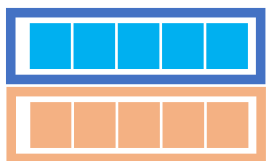| Query | Key Value Store | Key | Value |
|---|---|---|---|

$$Q \qquad QK^T \qquad softmax(\frac{QK^T}{\sqrt{d}}) \qquad softmax(\frac{QK^T}{\sqrt{d}})V$$

# Attention

**Parallelizable !!!**

*Attention Filter*

Query

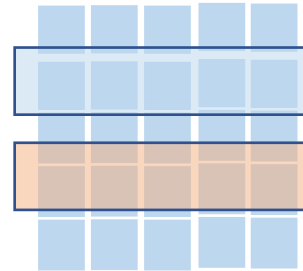Key Value Store

Key

Value

$Q$

$QK^T$

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

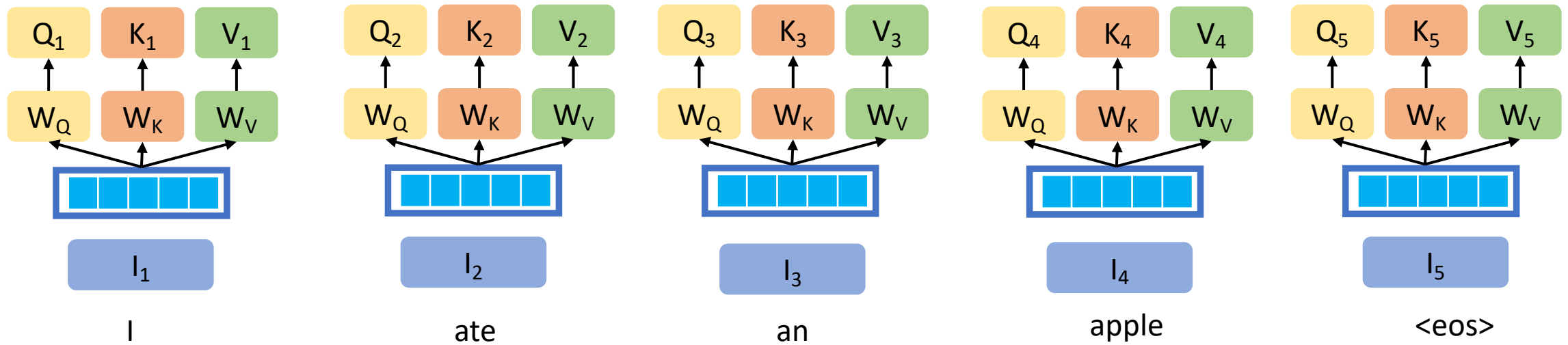$softmax(\frac{QK^T}{\sqrt{d}})V$
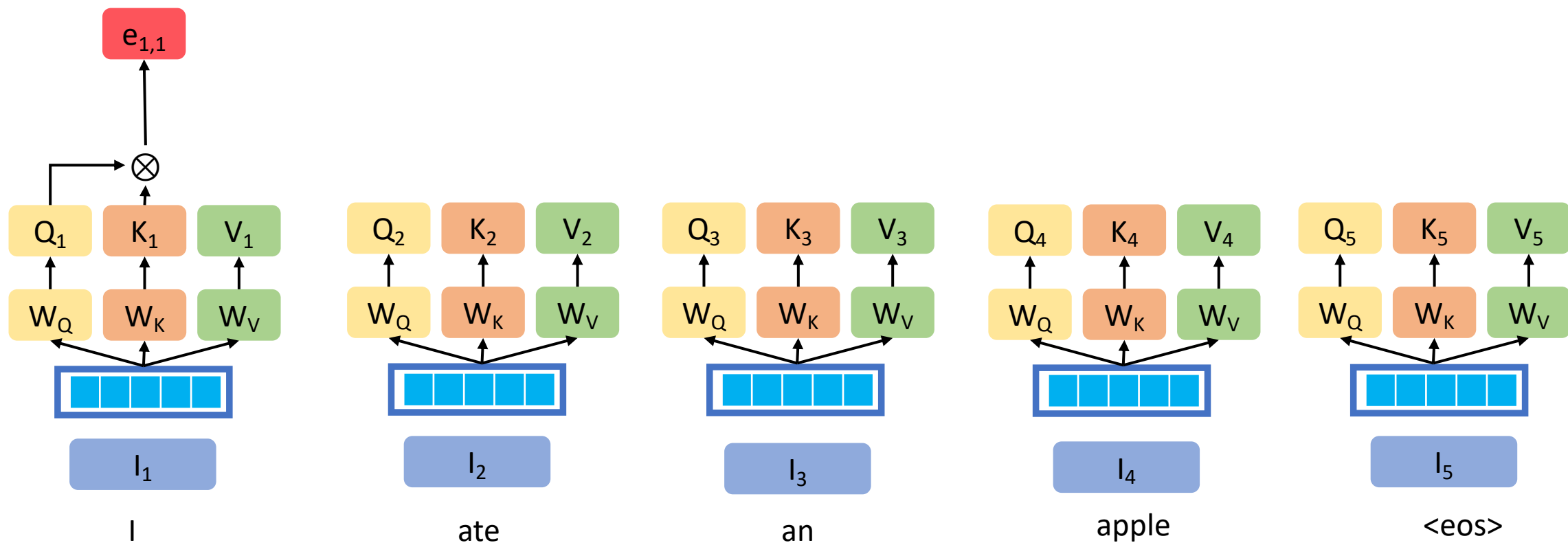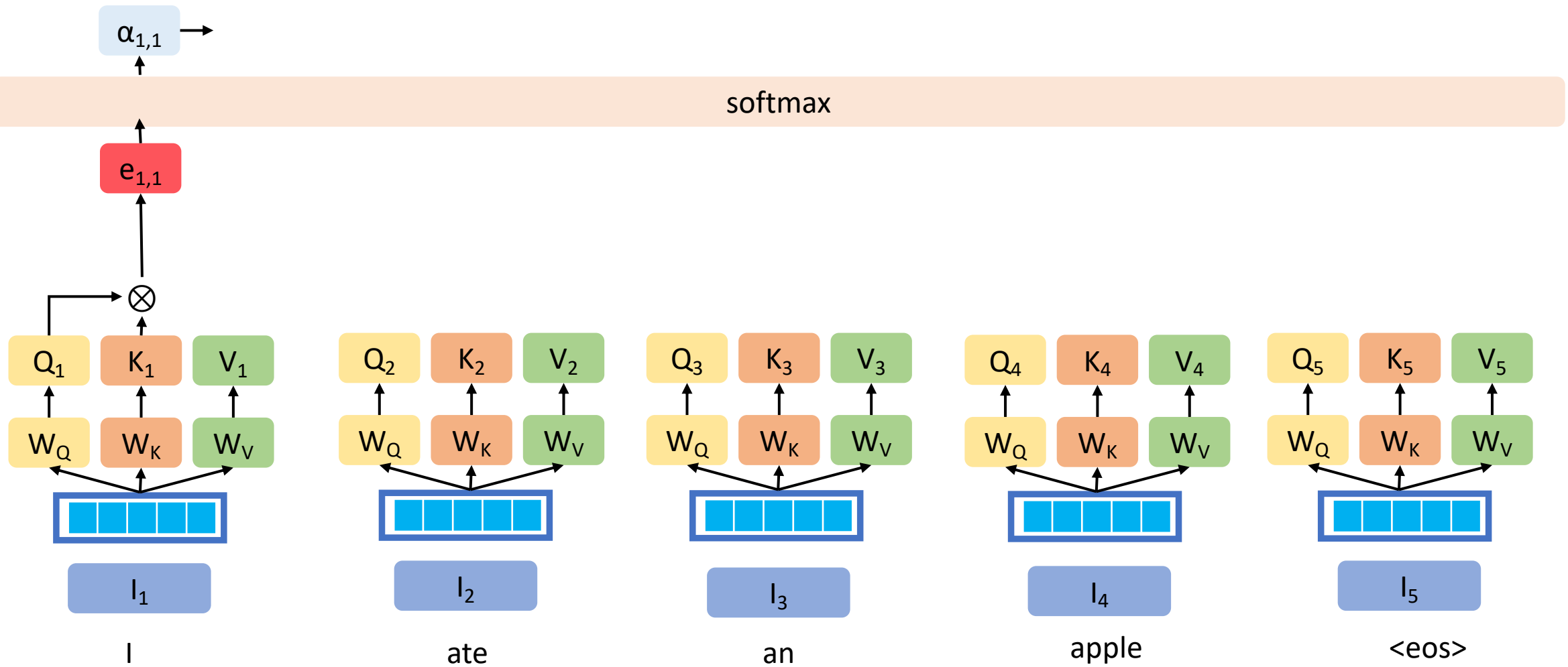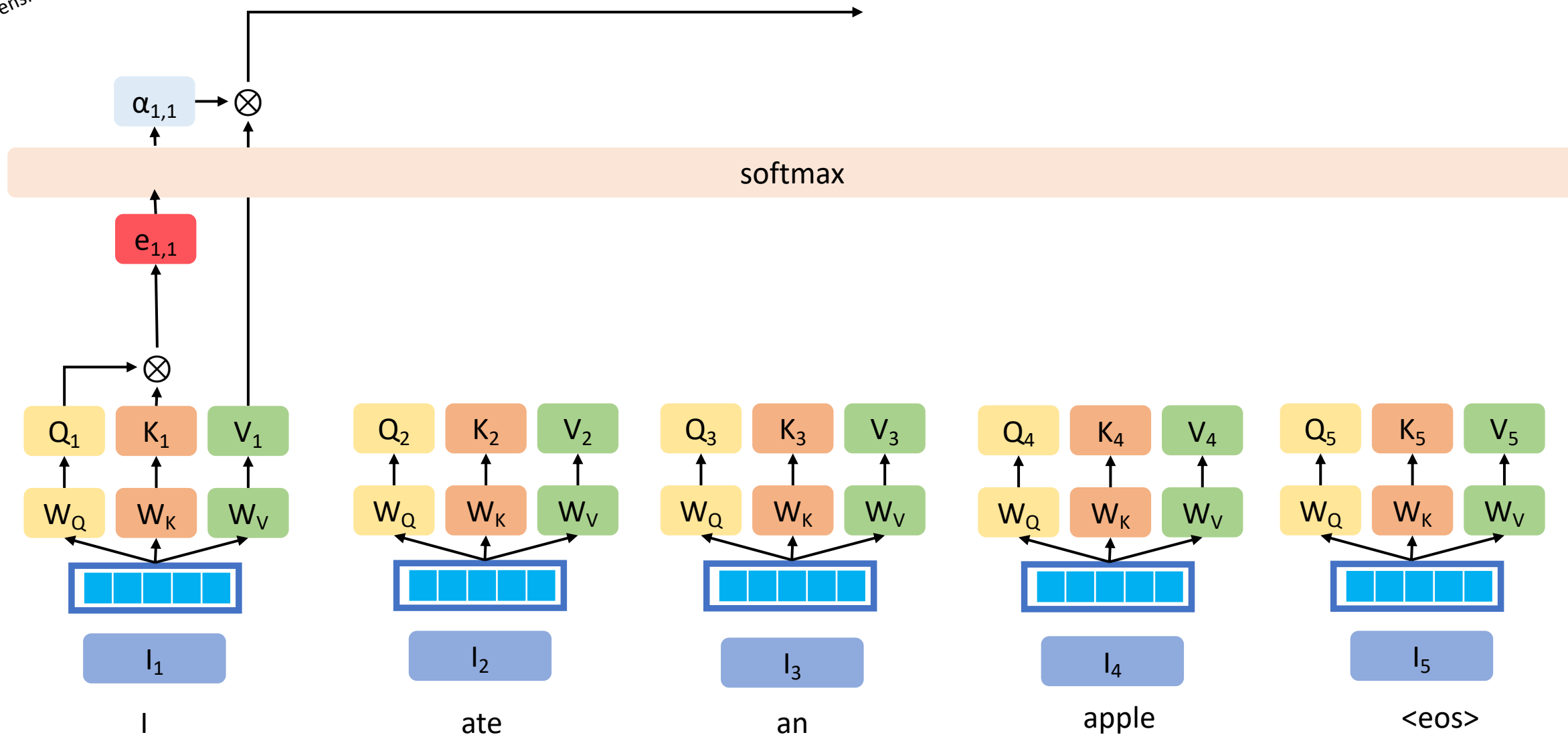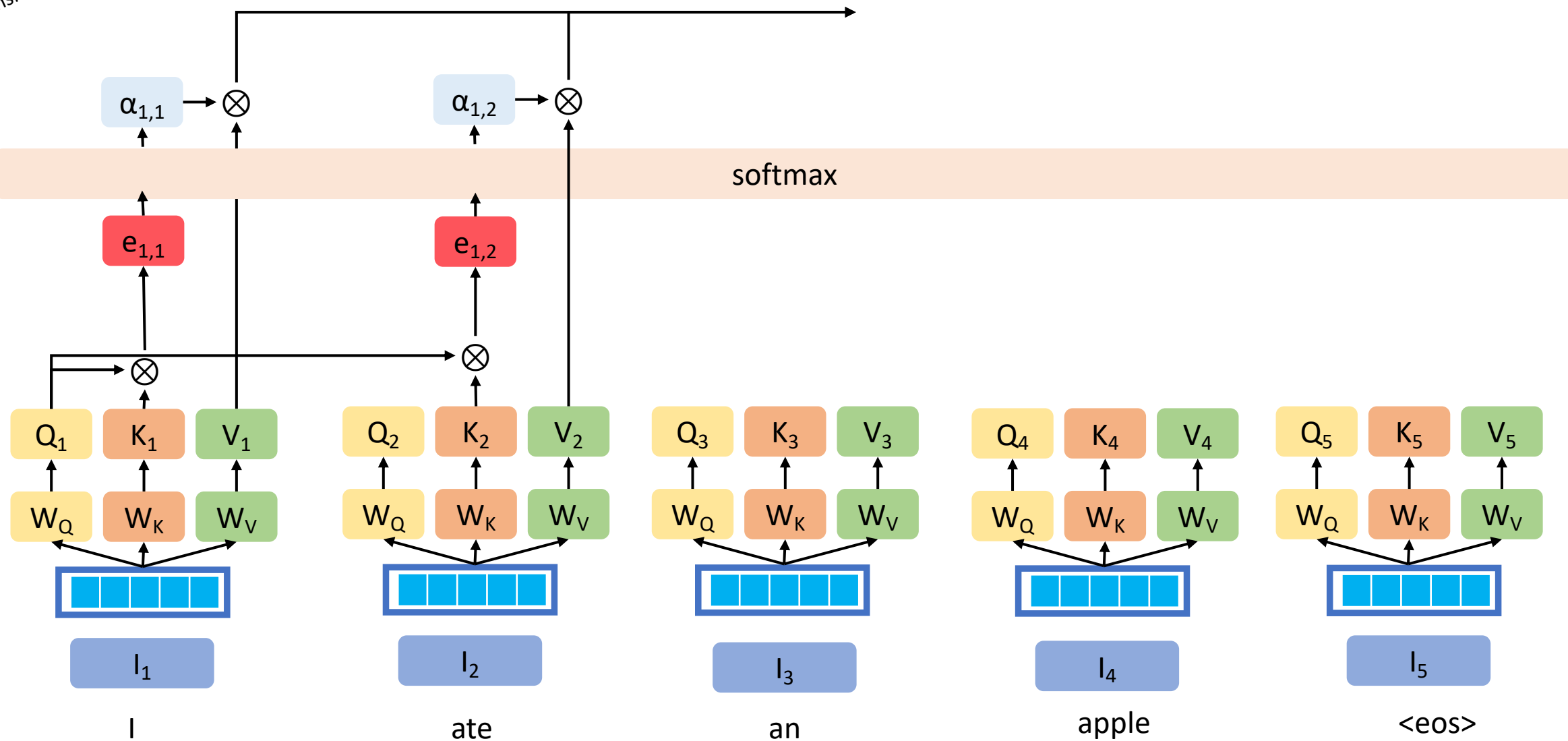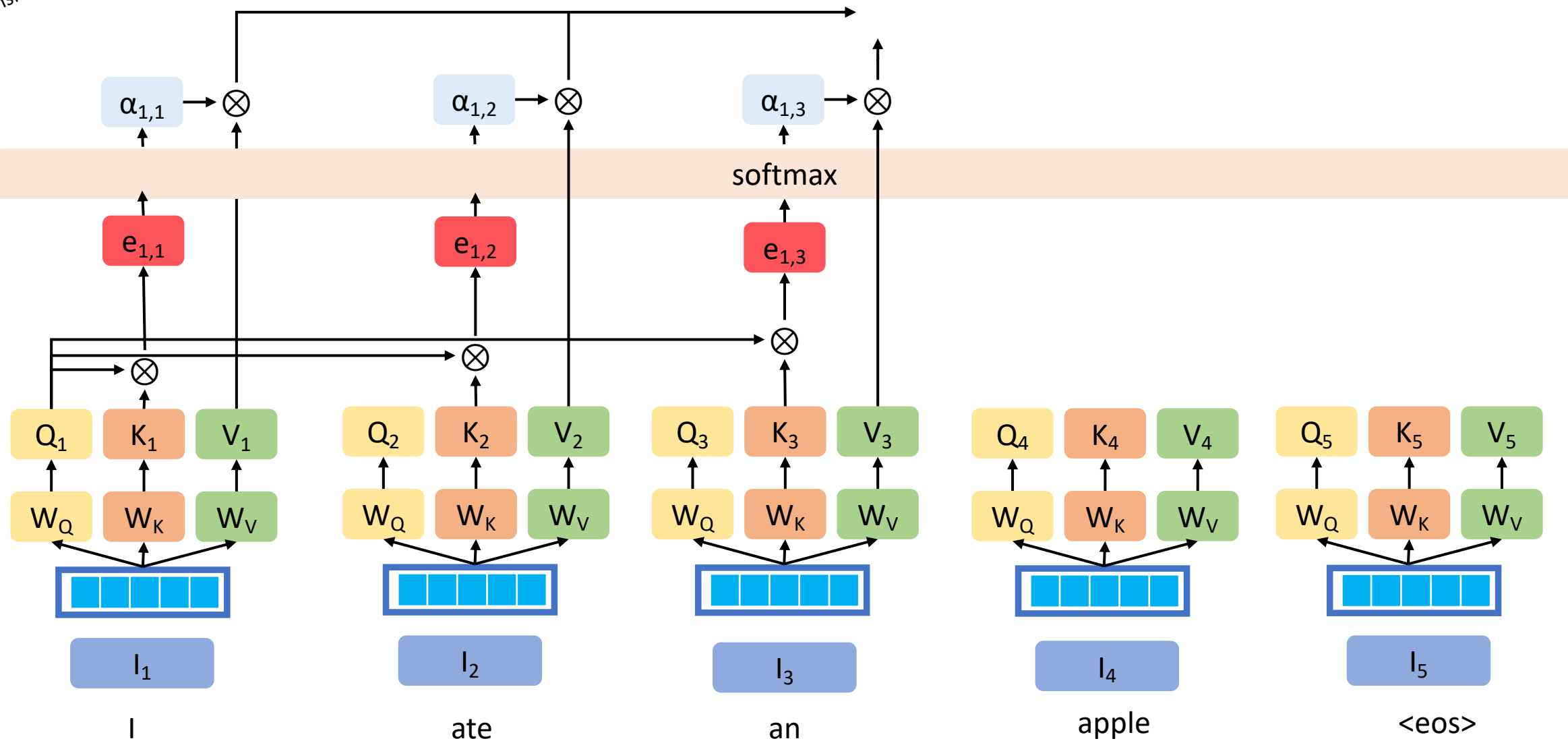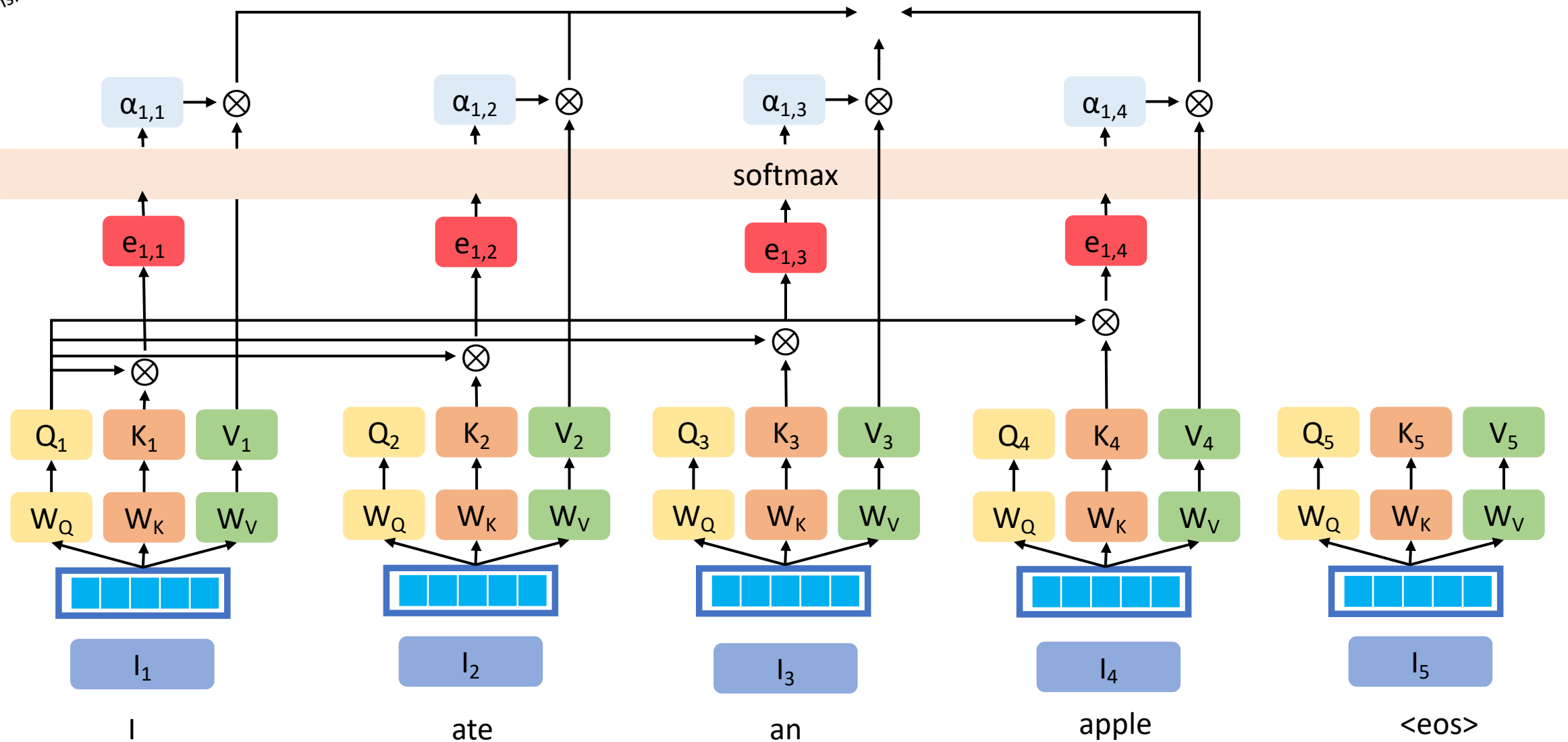
# Attention

# Attention

Dimensions across QKV have been dropped for brevity

# Attention

Dimensions across QKV have been dropped for brevity

# Attention

Dimensions across QKV have been dropped for brevity

# Attention

# Attention

Dimensions across QKV have been dropped for brevity

softmax

$\alpha_{1,1}$  $\otimes$  $\alpha_{1,2}$  $\otimes$

$e_{1,1}$  $e_{1,2}$

$Q_1$  $K_1$  $V_1$  $Q_2$  $K_2$  $V_2$  $Q_3$  $K_3$  $V_3$  $Q_4$  $K_4$  $V_4$  $Q_5$  $K_5$  $V_5$

$W_Q$  $W_K$  $W_V$  $W_Q$  $W_K$  $W_V$  $W_Q$  $W_K$  $W_V$  $W_Q$  $W_K$  $W_V$  $W_Q$  $W_K$  $W_V$

$I_1$  $I_2$  $I_3$  $I_4$  $I_5$

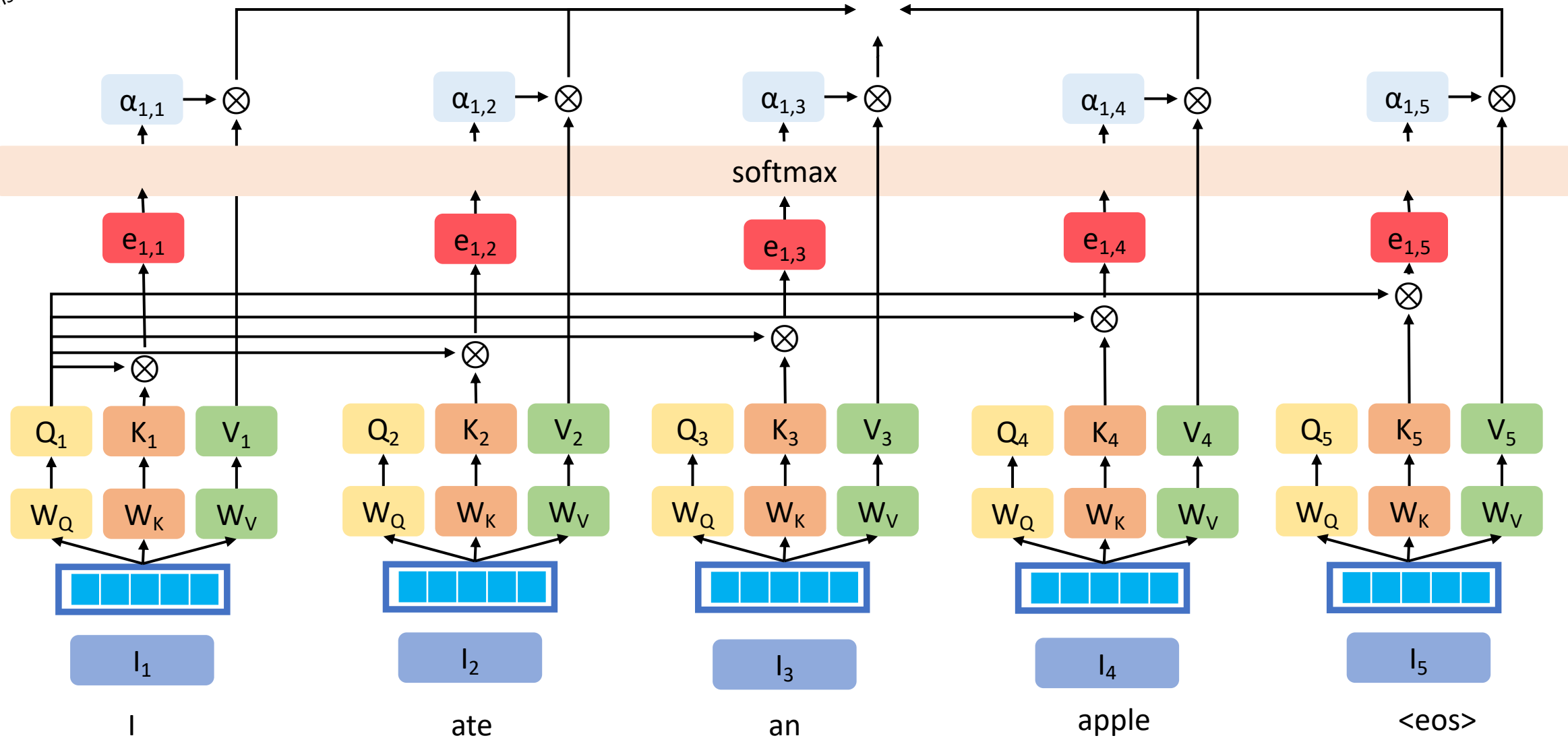I  ate  an  apple  <eos>
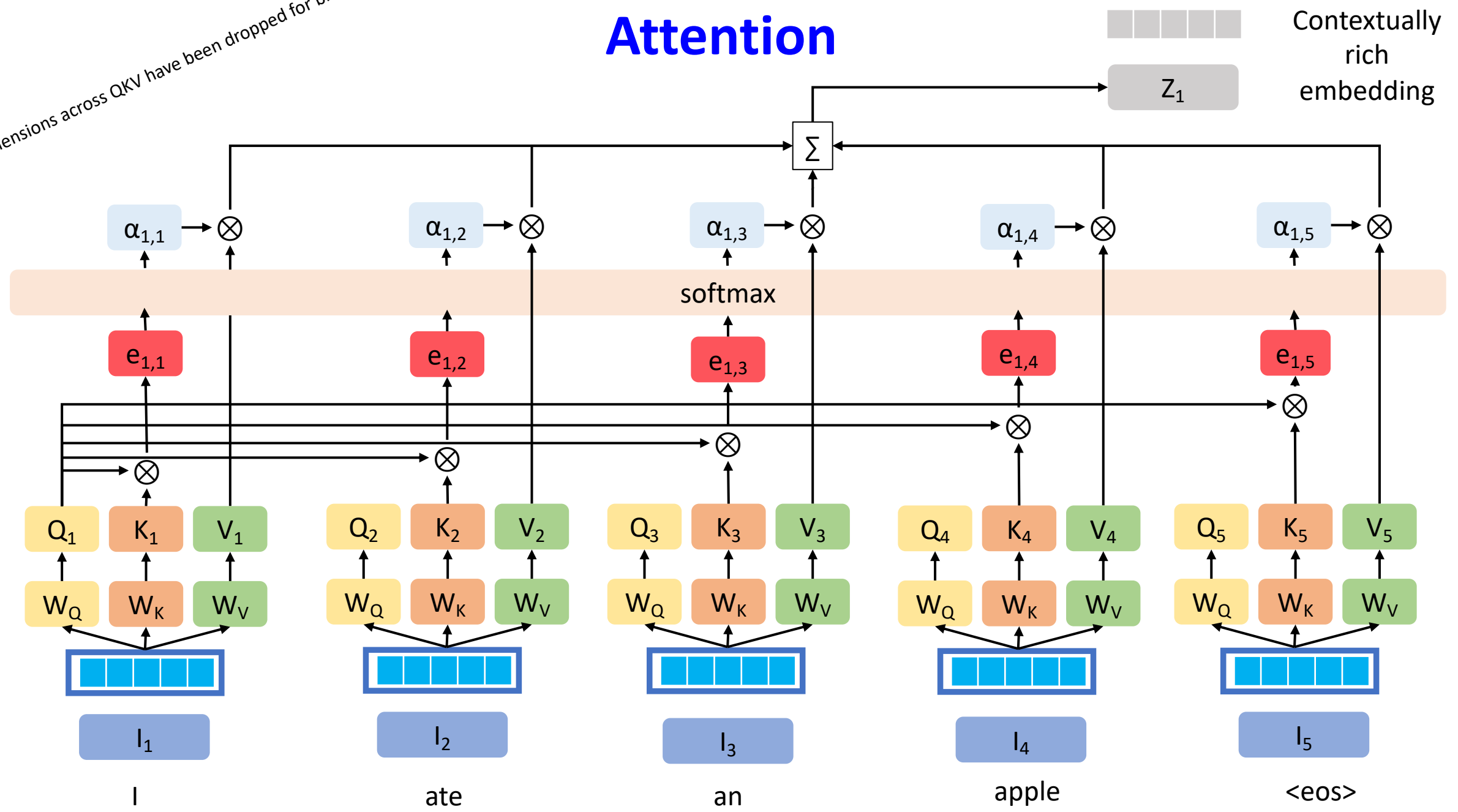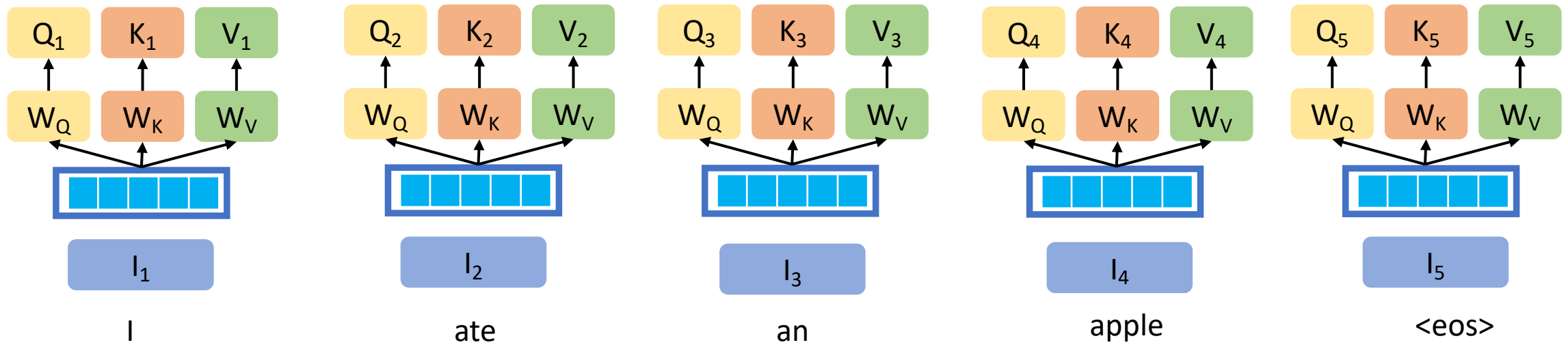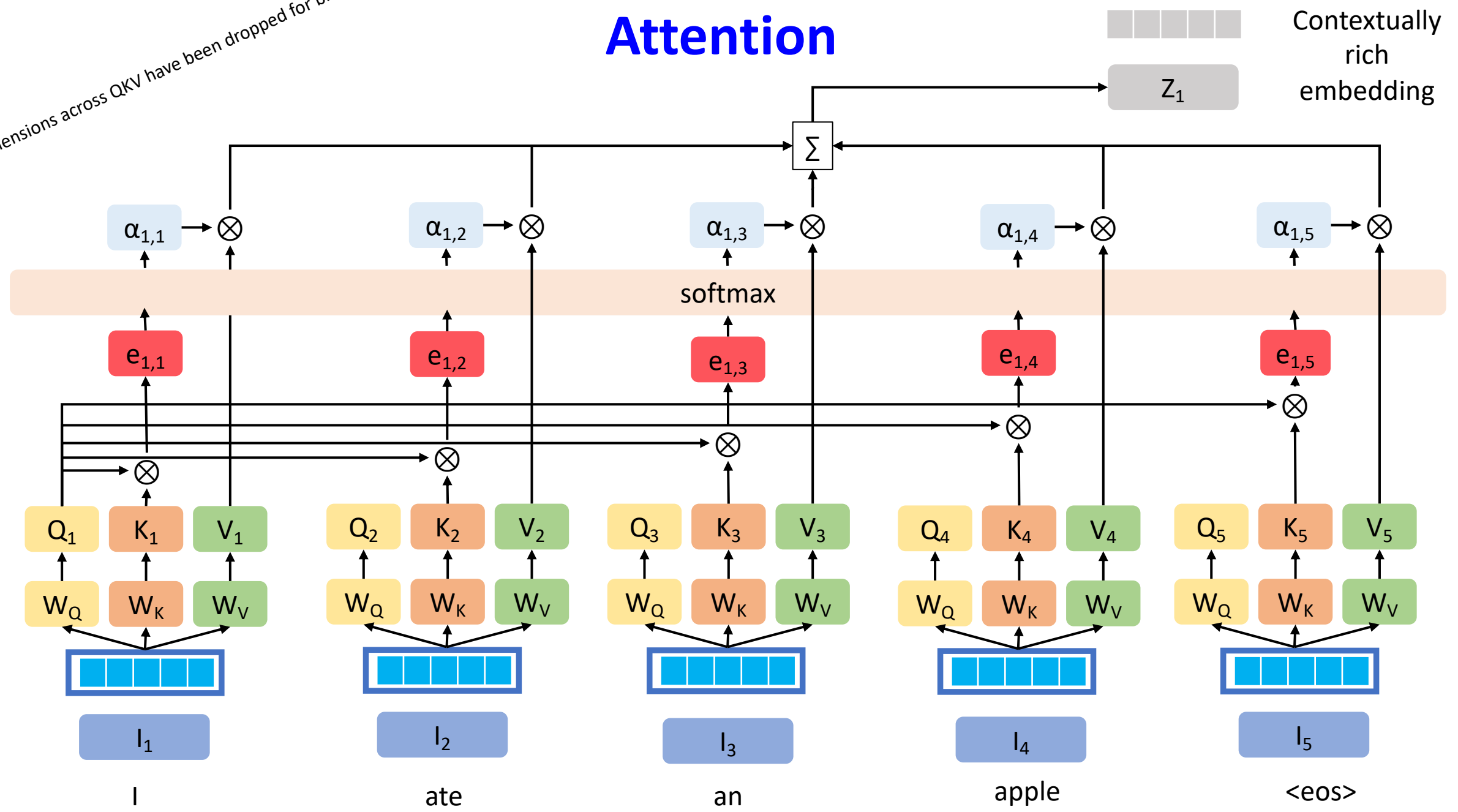
34

# Attention

Dimensions across QKV have been dropped for brevity

35

# Attention

Dimensions across QKV have been dropped for brevity



36

# Attention

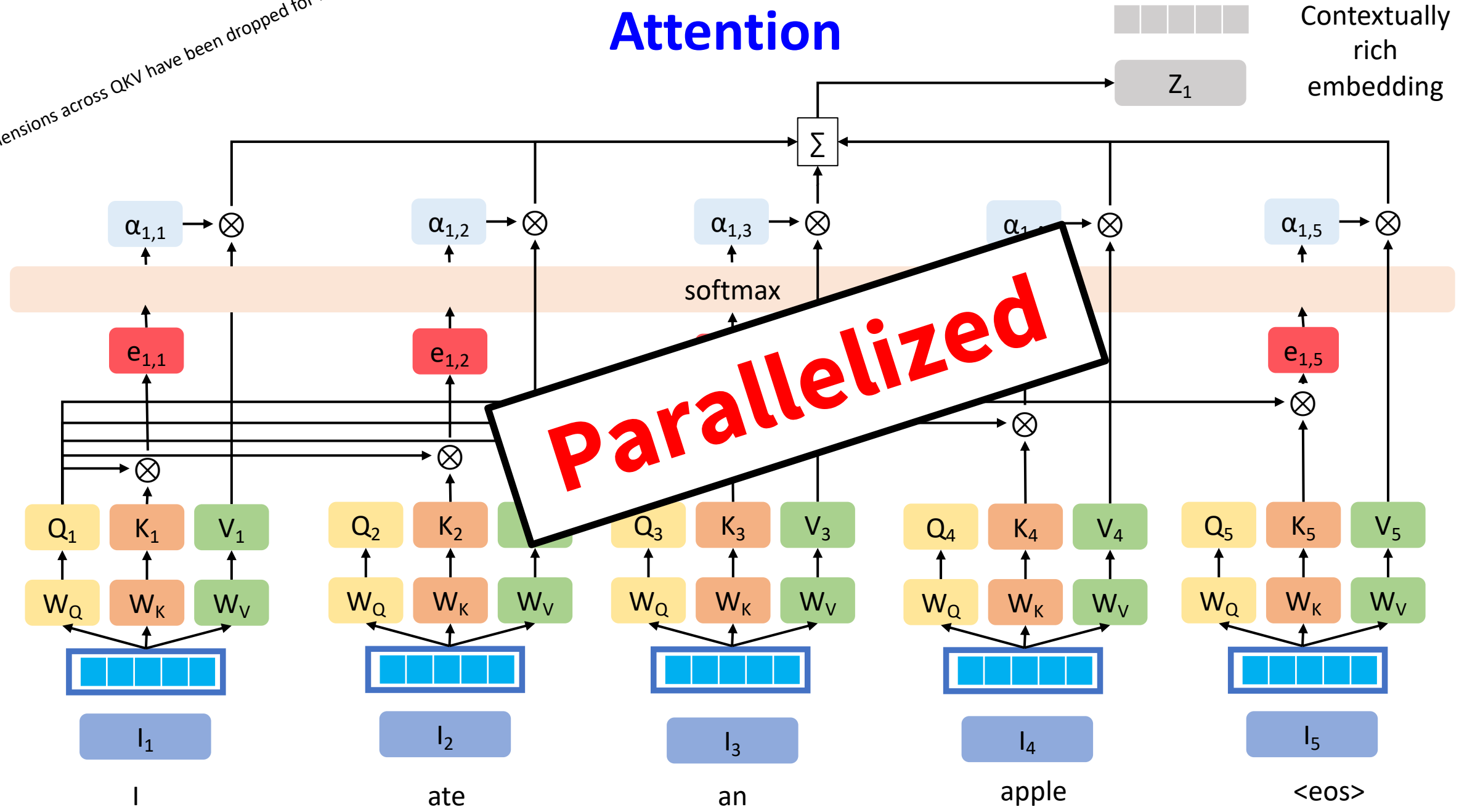Dimensions across QKV have been dropped for brevity



37

# Attention

Dimensions across QKV have been dropped for brevity



39

Attention

# Poll 1 @1296

**Which of the following are true about attention?** (Select all that apply)

a. To calculate attention weights for input $I_2$, you would use key $k_2$, and all queries

b. To calculate attention weights for input $I_2$, you would use query $q_2$, and all keys

c. We scale the $QK^T$ product to bring attention weights in the range of [0,1]
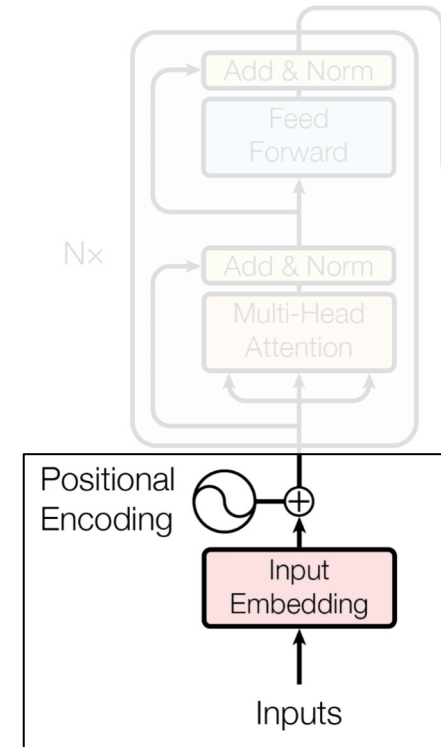
d. We scale the $QK^T$ product to allow for numerical stability

# Poll 1 @1296
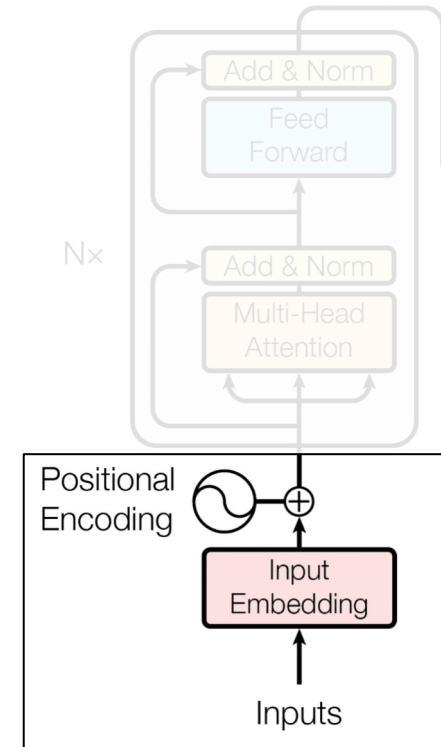
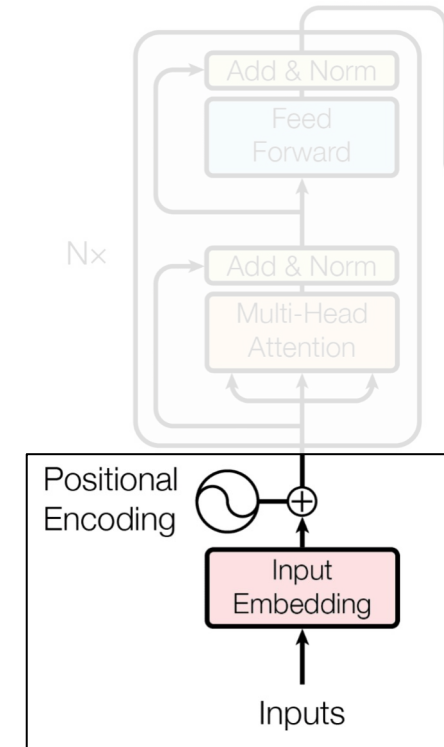**Which of the following are true about attention?** (Select all that apply)

a. To calculate attention weights for input $I_2$, you would use key $k_2$, and all queries

**b. To calculate attention weights for input $I_2$, you would use query $q_2$, and all keys**

c. We scale the $QK^T$ product to bring attention weights in the range of [0,1]

**d. We scale the $QK^T$ product to allow for numerical stability**

# Positional Encoding

| I | ate | an | apple | <eos> |

# Positional Encoding

| I | ate | an | apple | <eos> |

❌

| apple | ate | an | I | <eos> |

Positional Encoding

# Positional Encoding

**Requirements for Positional Encodings**

- Some representation of time ? (like **seq2seq** ?)

- Should be unique for each position – not cyclic



Positional Encoding

# Positional Encoding

**Requirements for Positional Encodings**

- Some representation of time ? (like **seq2seq** ?)

- Should be unique for each position – not cyclic

Possible Candidates :

$$P_{t+1} = P_t + \Delta c$$

$$P_{t+1} = e^{P_{t_\Delta} c}$$

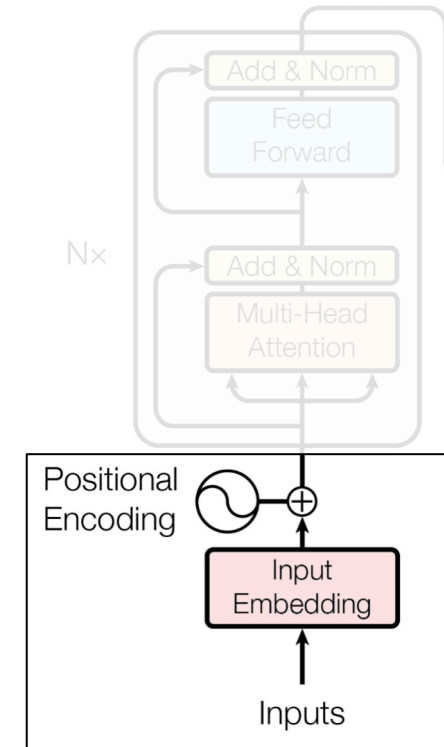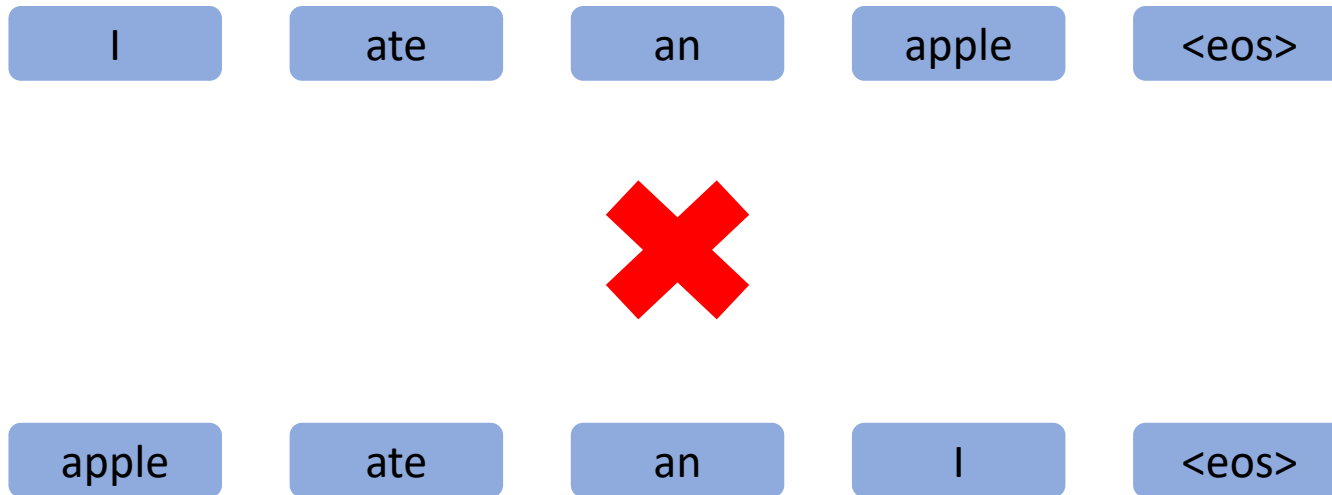$$P_{t+1} = P_t^{\cdot \; t\Delta c}$$



Positional Encoding
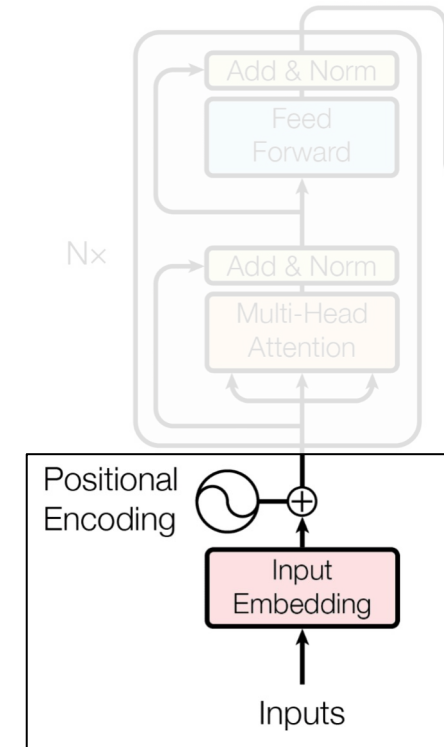Input Embedding
Inputs

<span style="color:red">Positional Encoding</span>

# Positional Encoding
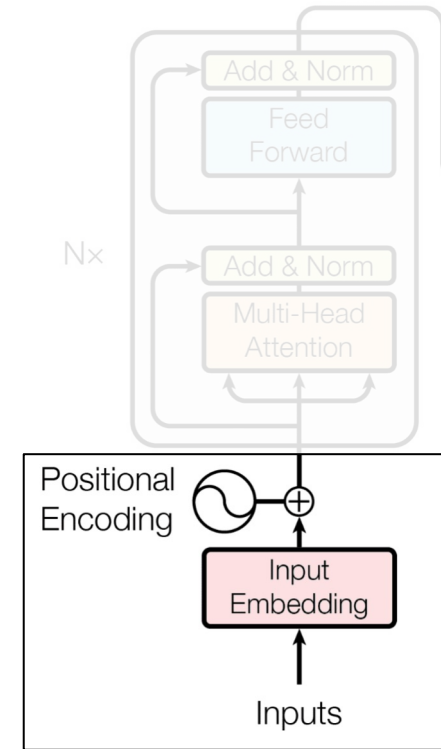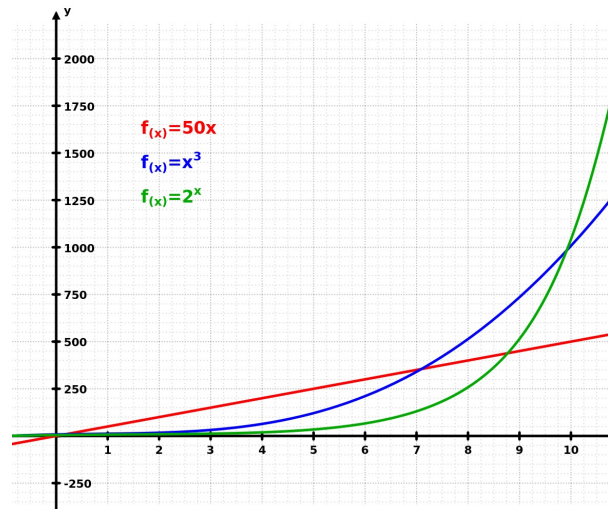
## Requirements for Positional Encodings

- Some representation of time ? (like **seq2seq** ?)

- Should be unique for each position – not cyclic
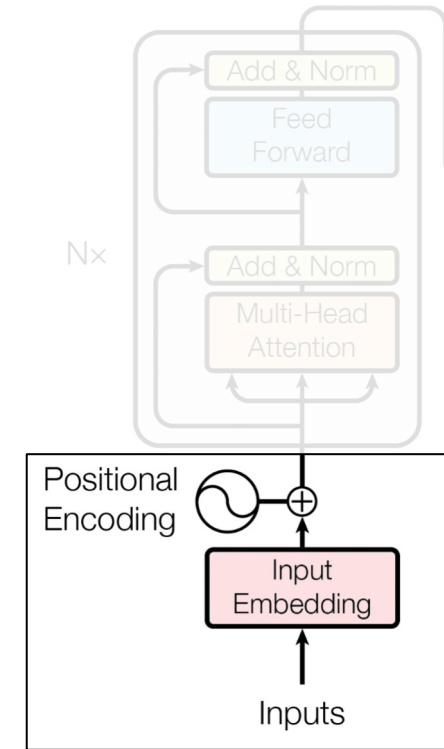
Possible Candidates :

$$P_{t+1} = P_t + \Delta c$$

$$P_{t+1} = e^{P_t \Delta} c$$

$$P_{t+1} = P_t^{\cdot \ t\Delta c}$$



$f_{(x)} = 50x$

$f_{(x)} = x^3$

$f_{(x)} = 2^x$

Positional Encoding

# Positional Encoding

**Requirements for Positional Encodings**

- Some representation of time ? (like **seq2seq** ?)

- Should be unique for each position – not cyclic

- **Bounded**

Possible Candidates :

$$P_{t+1} = P_t + \Delta c$$

$$P_{t+1} = e^{P_t \Delta c}$$

$$P_{t+1} = P_t^{t \Delta c}$$

f<sub>(x)</sub>=50x
f<sub>(x)</sub>=x³
f<sub>(x)</sub>=2ˣ

Positional Encoding

Positional Encoding
Input Embedding
Inputs

Add & Norm
Feed Forward
Nx
Add & Norm
Multi-Head Attention

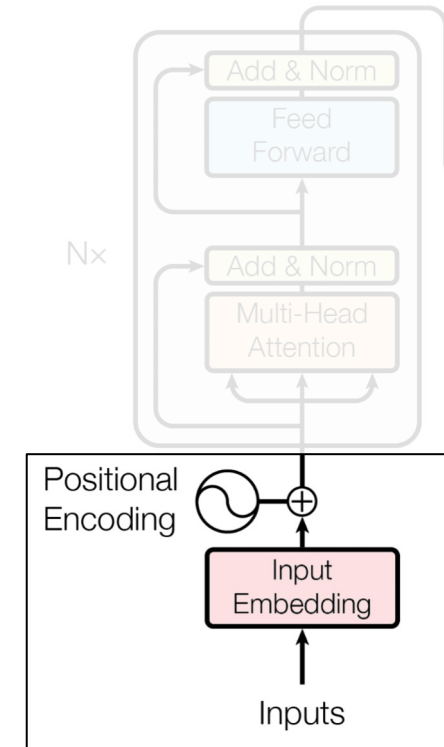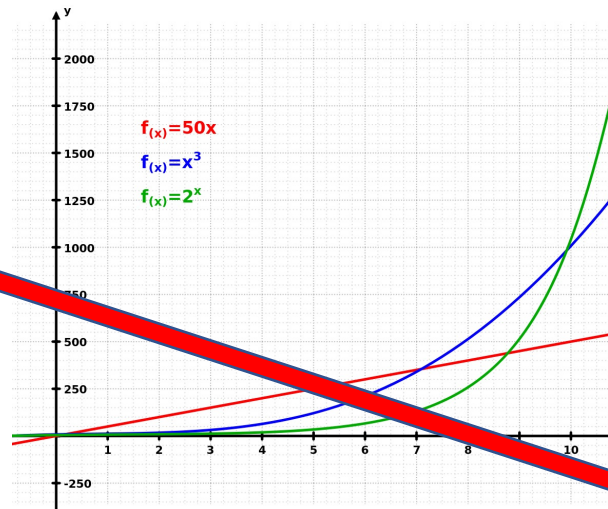# Positional Encoding

**Requirements for Positional Encodings**

- Some representation of time ? (like **seq2seq** ?)

- Should be unique for each position – not cyclic

- **Bounded**

Possible Candidates :

$$P(t + t') = M^{t'} \times P(t)$$
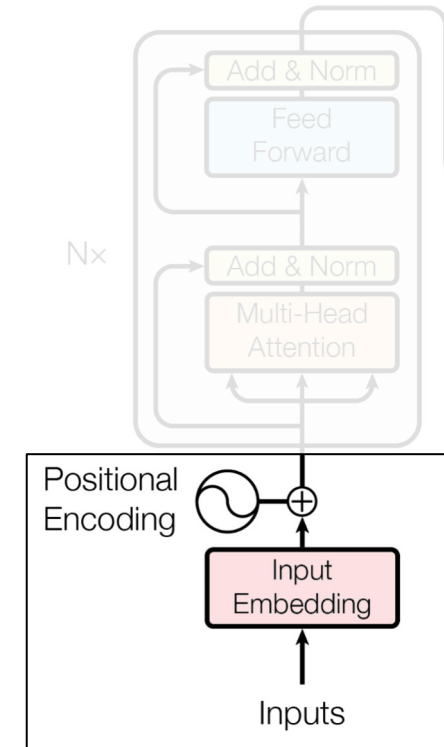


Positional Encoding

# Positional Encoding

**Requirements for Positional Encodings**

• Some representation of time ? (like **seq2seq** ?)

• Should be unique for each position – not cyclic

• **Bounded**

Possible Candidates :

## $P(t + t') = M^{t'} \times P(t)$

## M ?

1. **Should be a unitary matrix**

2. **Magnitudes of eigen value should be 1 -> norm preserving**

<span style="color:red">Positional Encoding</span>

# Positional Encoding

**Requirements for Positional Encodings**

- Some representation of time ? (like **seq2seq** ?)

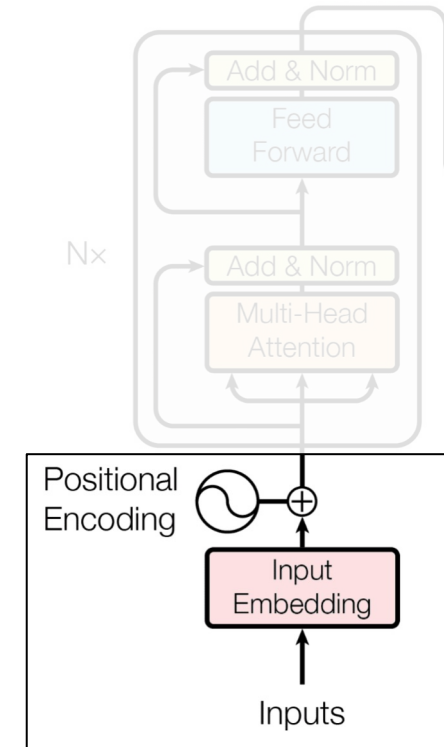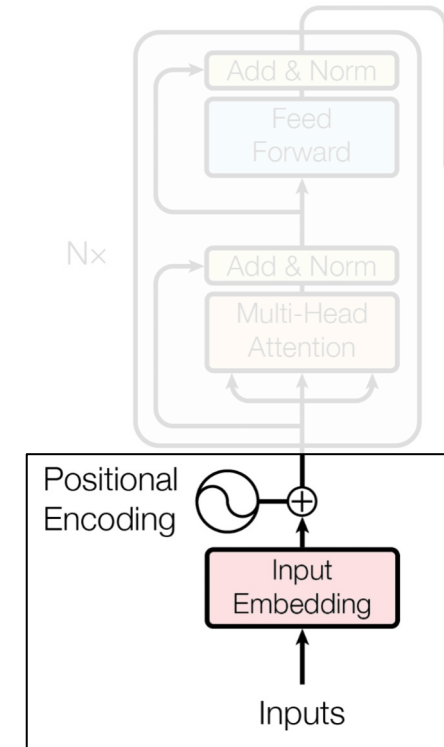- Should be unique for each position – not cyclic

- **Bounded**

Possible Candidates :

$$P(t + t') = M^{t'} \times P(t)$$

$$M$$

1. **The matrix can be learnt**

2. **Produces unique rotated embeddings each time**
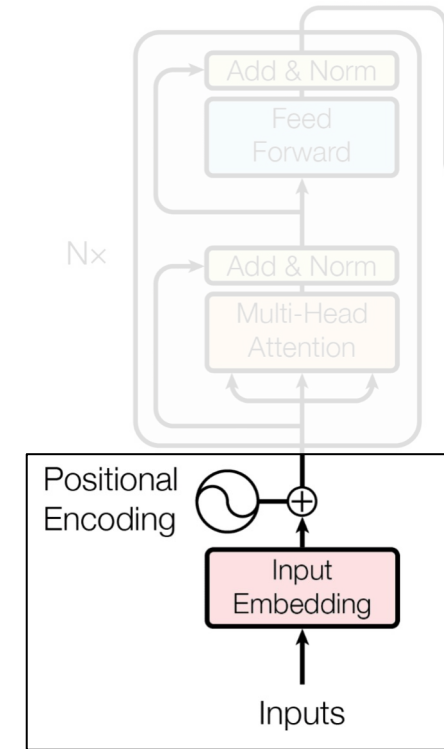
Positional Encoding

# Rotary Positional Embedding

RoFORMER: ENHANCED TRANSFORMER WITH ROTARY POSITION EMBEDDING

$$f_{\{q,k\}}(\boldsymbol{x}_m, m) = \begin{pmatrix} \cos m\theta & -\sin m\theta \\ \sin m\theta & \cos m\theta \end{pmatrix} \begin{pmatrix} W_{\{q,k\}}^{(11)} & W_{\{q,k\}}^{(12)} \\ W_{\{q,k\}}^{(21)} & W_{\{q,k\}}^{(22)} \end{pmatrix} \begin{pmatrix} x_m^{(1)} \\ x_m^{(2)} \end{pmatrix}$$

Table 2: Comparing RoFormer and BERT by fine tuning on downstream GLEU tasks.

| Model | MRPC | SST-2 | QNLI | STS-B | QQP | MNLI(m/mm) |
|---|---|---|---|---|---|---|
| BERTDevlin et al. [2019] | 88.9 | 93.5 | 90.5 | 85.8 | 71.2 | 84.6/83.4 |
| RoFormer | **89.5** | 90.7 | 88.0 | **87.0** | **86.4** | 80.2/79.8 |

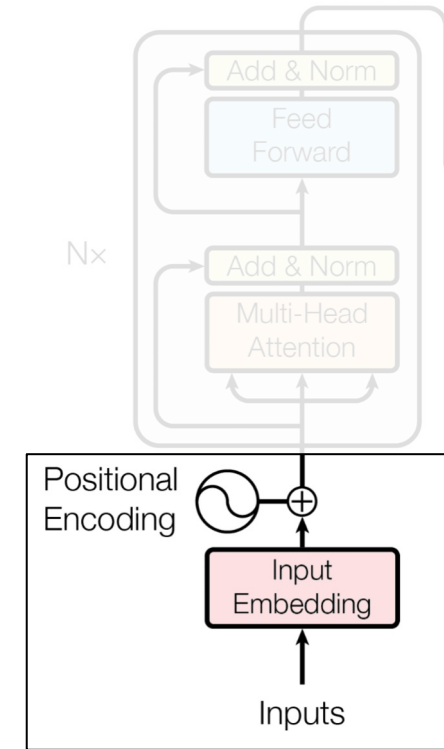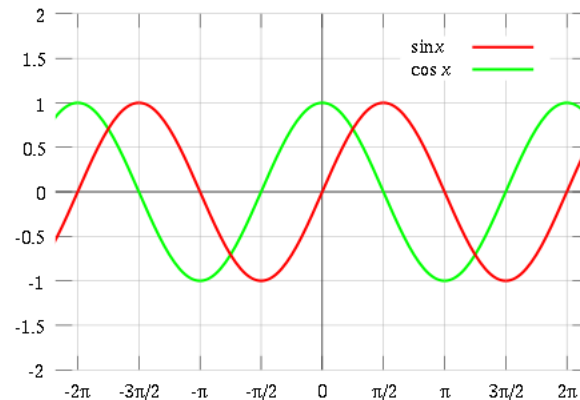REF: Rotary Positional Embeddings 🔗

53

# Positional Encoding

**Requirements for Positional Encodings**

- Some representation of time ? (like **seq2seq** ?)

- Should be unique for each position – **not cyclic**

- Bounded

Actual Candidates :

$sine(g(t))$

$cosine(g(t))$



Positional Encoding

# Positional Encoding

*Requirements for g(t)*

- *Must have same dimensions as input embeddings*
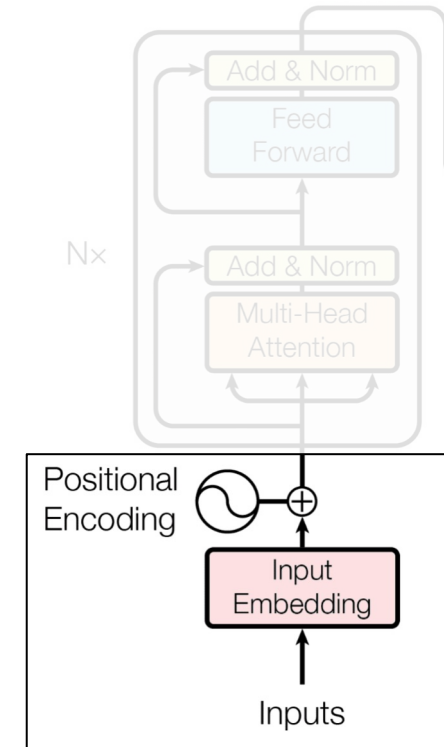
- *Must produce overall unique encodings*

pos -> idx of the token in input sentence

i     -> i[th]   dimension out of d

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

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

<span style="color:red">Positional Encoding</span>

# Positional Encoding

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

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

*Requirements for g(t)*

- *Must have same dimensions as input embeddings*
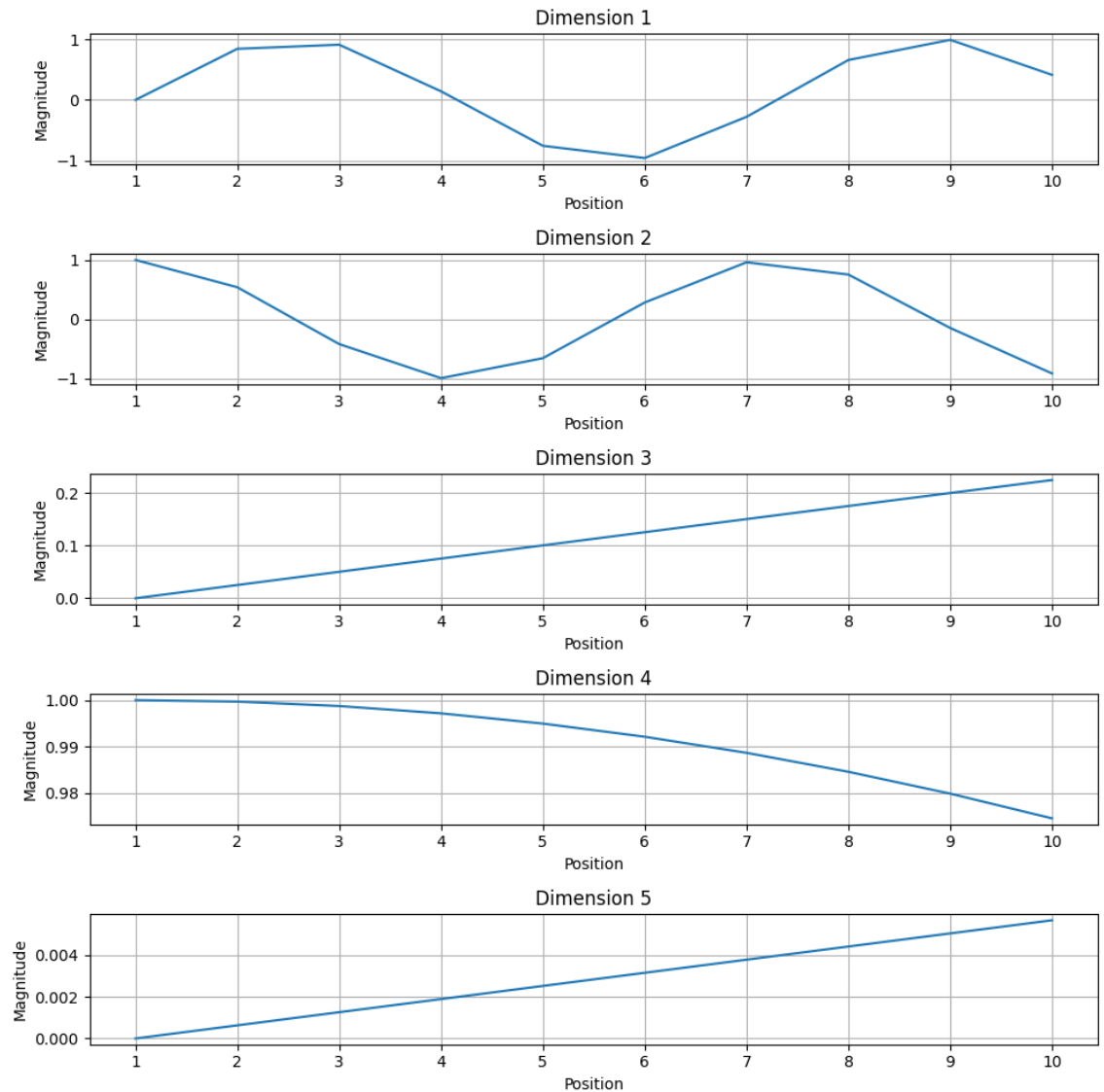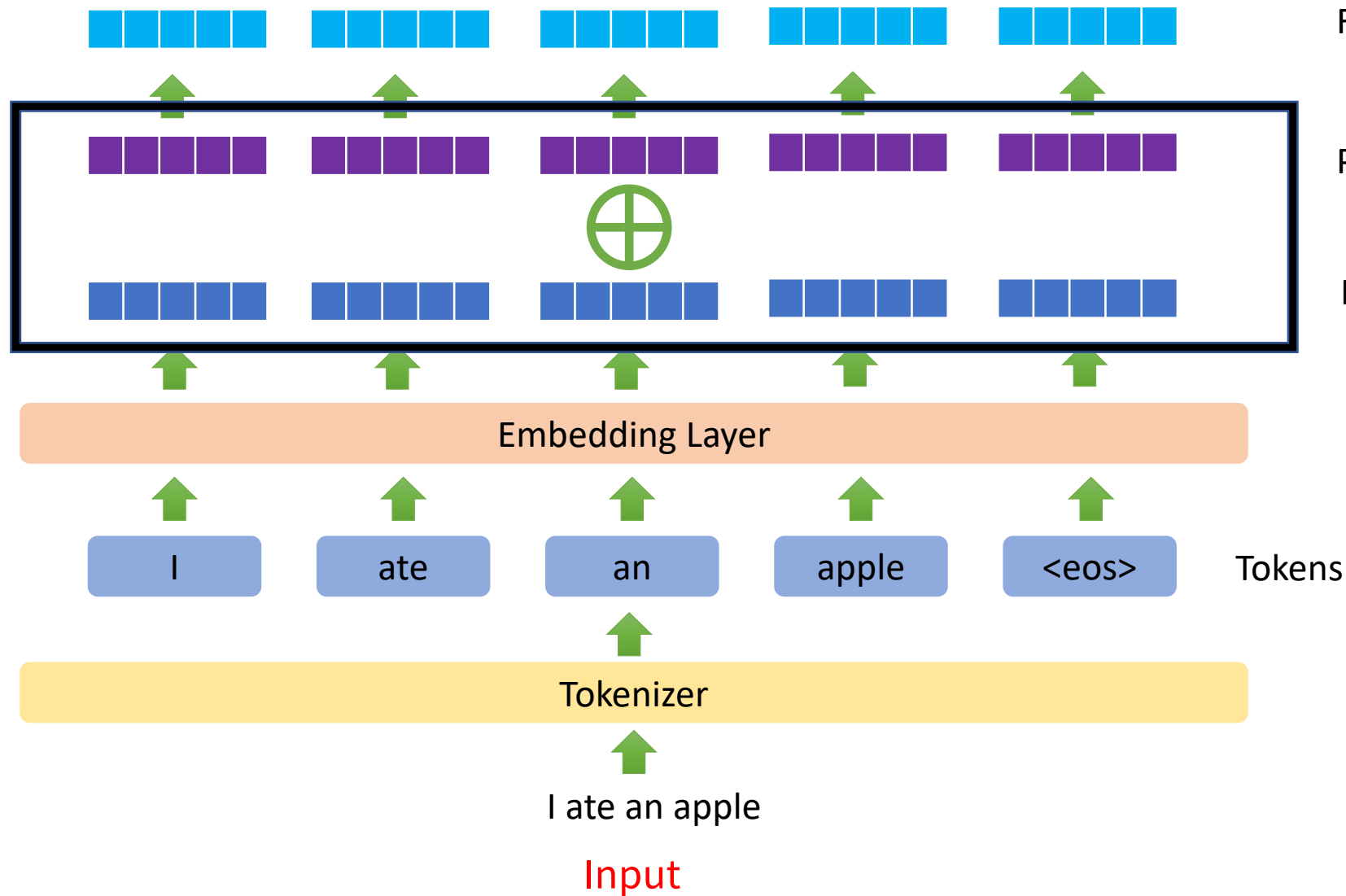
- *Must produce overall unique encodings*

pos -> idx of the token in input sentence

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



```
Positional Encoding:
          0      1       2       3       4
Dim 1 0.000  0.841   0.909   0.141  -0.757
Dim 2 1.000  0.540  -0.416  -0.990  -0.654
Dim 3 0.000  0.025   0.050   0.075   0.100
Dim 4 1.000  1.000   0.999   0.997   0.995
Dim 5 0.000  0.001   0.001   0.002   0.003
```
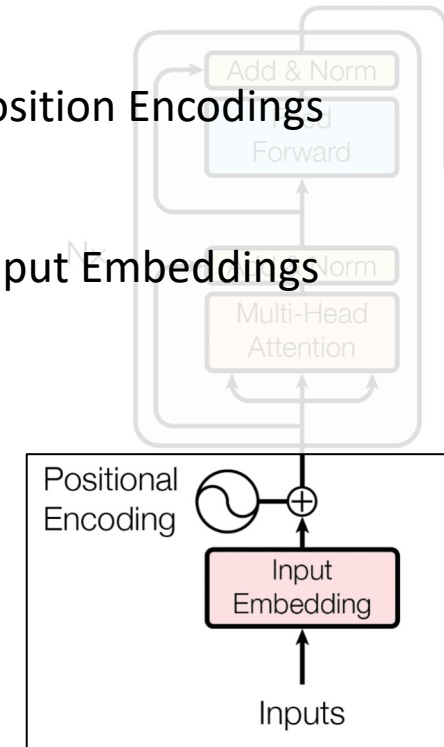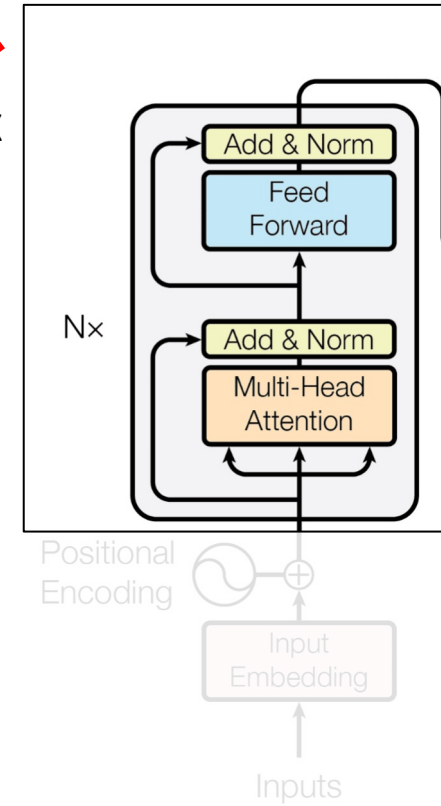
# Positional Encoding

Final Input Embeddings

Position Encodings
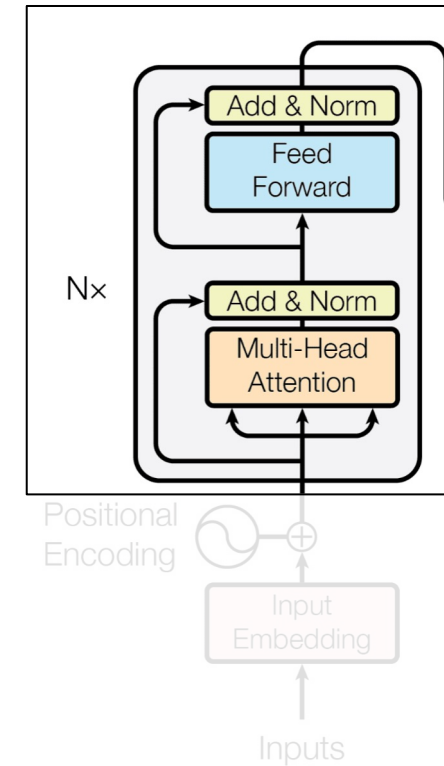
$\oplus$

Input Embeddings

Embedding Layer

| I | ate | an | apple | <eos> | Tokens |

Tokenizer

I ate an apple

Input

Positional Encoding

Input Embedding

Inputs

# Self Attention

From lecture 18:



$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

# Self Attention
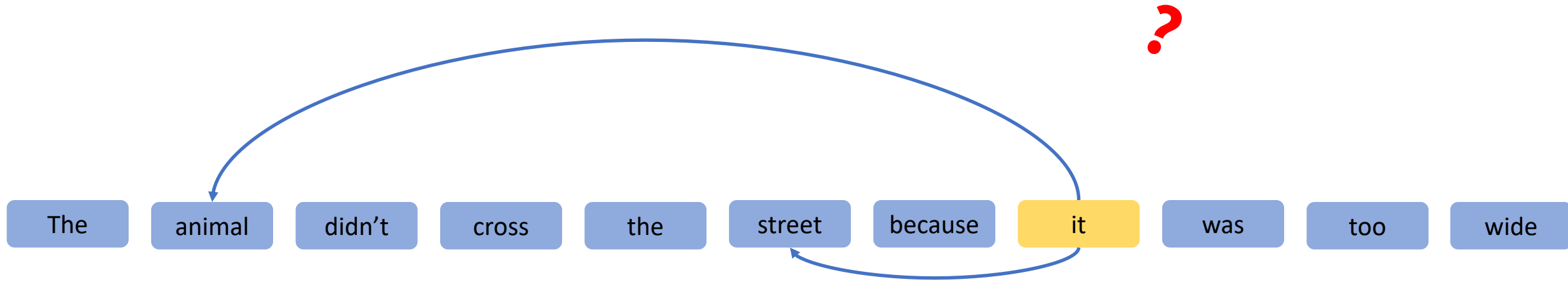
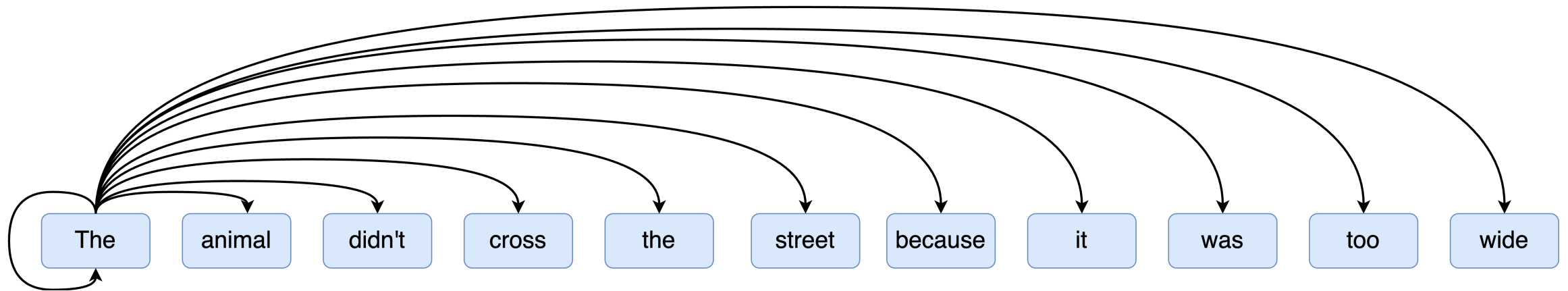The | animal | didn't | cross | the | street | because | it | was | too | wide
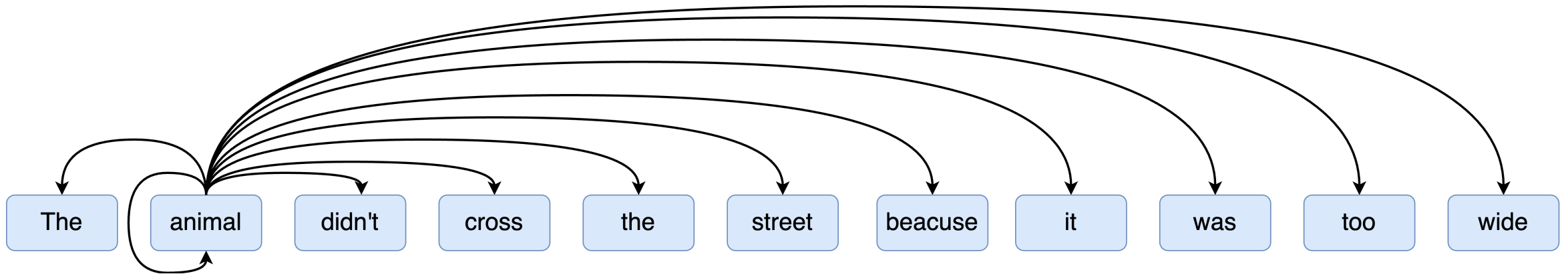
# Self Attention

?

The | animal | didn't | cross | the | street | because | it | was | too | wide

coreference resolution ?

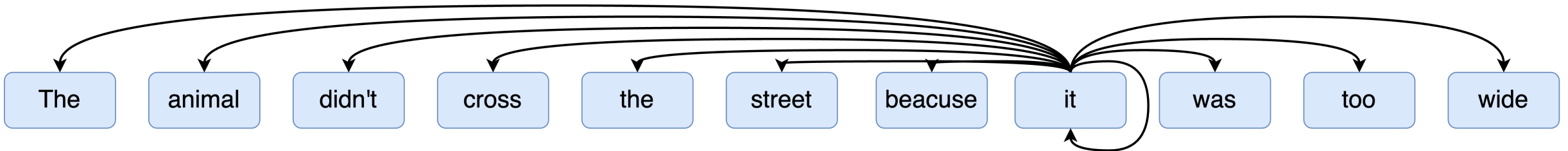# Self Attention

# Self Attention



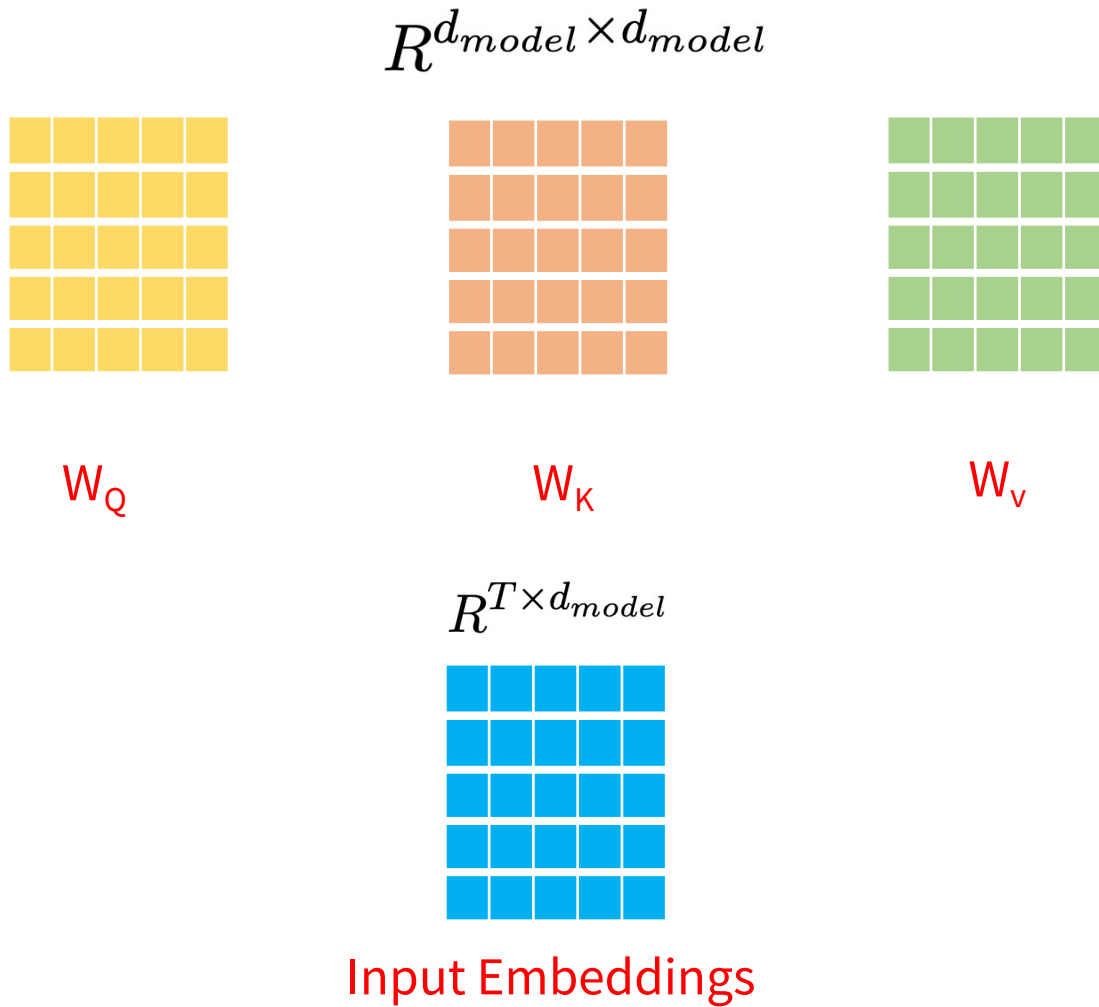The | animal | didn't | cross | the | street | beacuse | it | was | too | wide

# Self Attention



The | animal | didn't | cross | the | street | beacuse | it | was | too | wide

# Self Attention



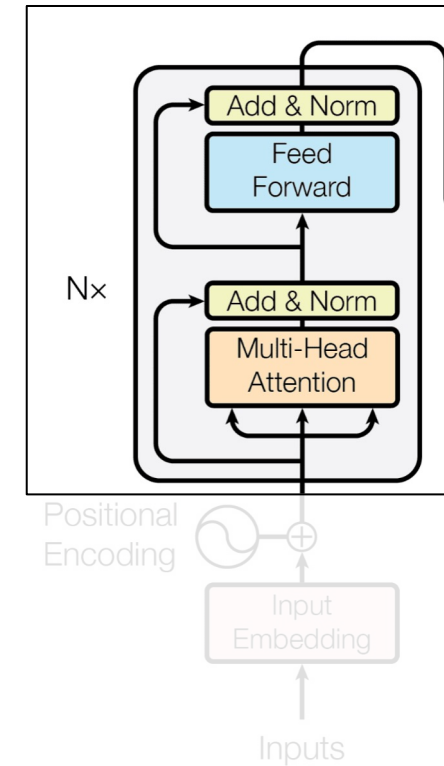The | animal | didn't | cross | the | street | beacuse | it | was | too | wide

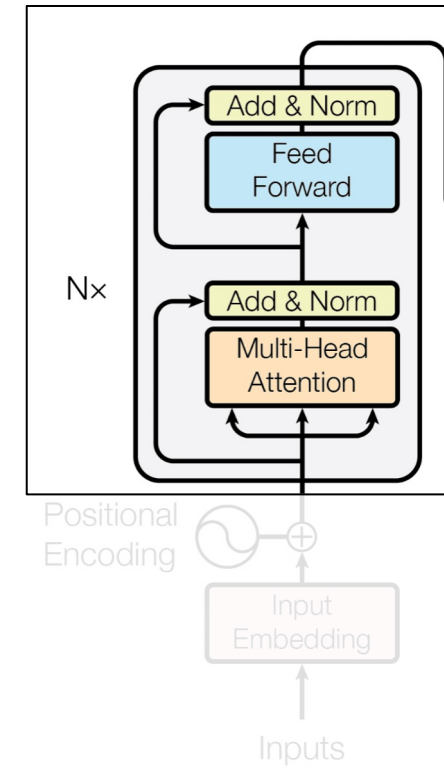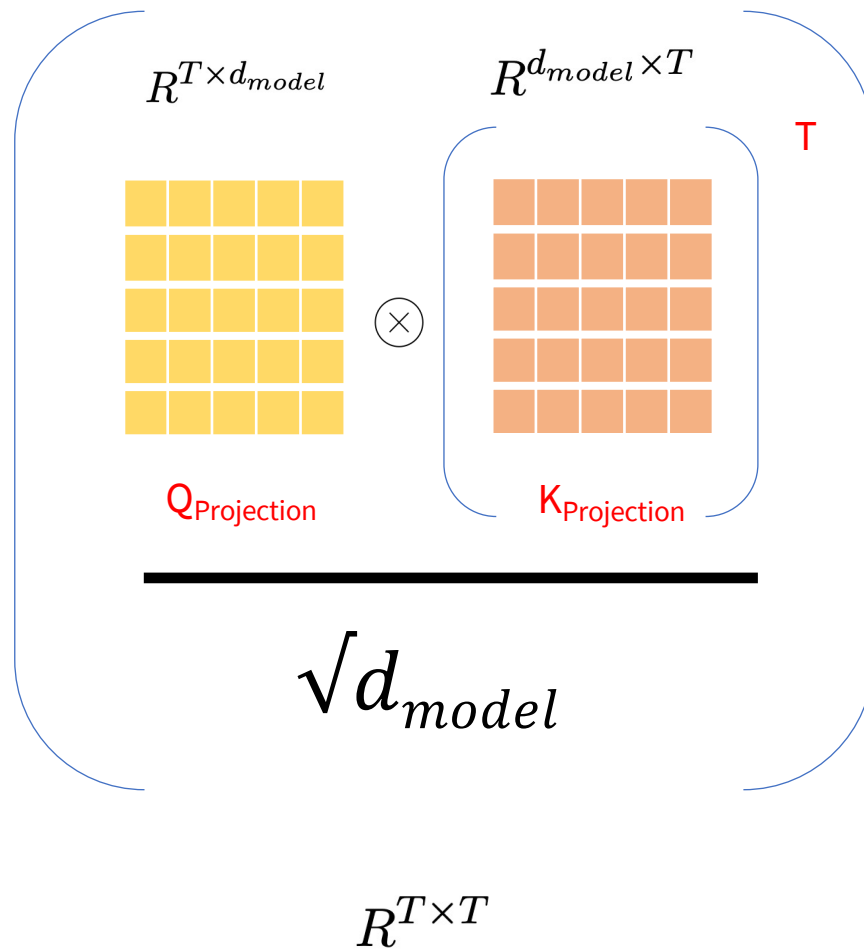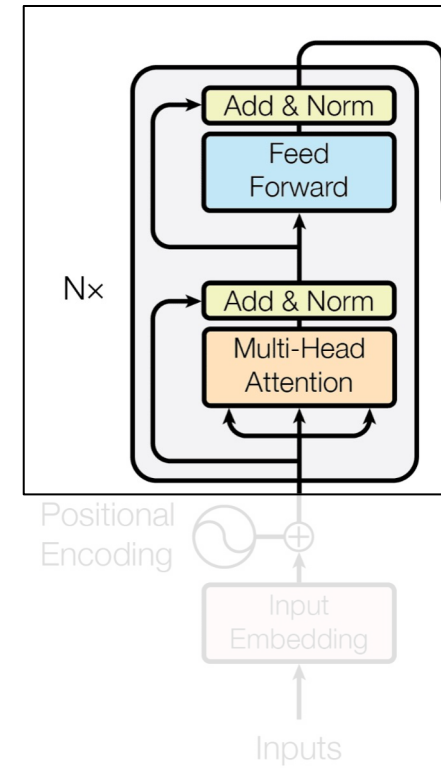**SELF**

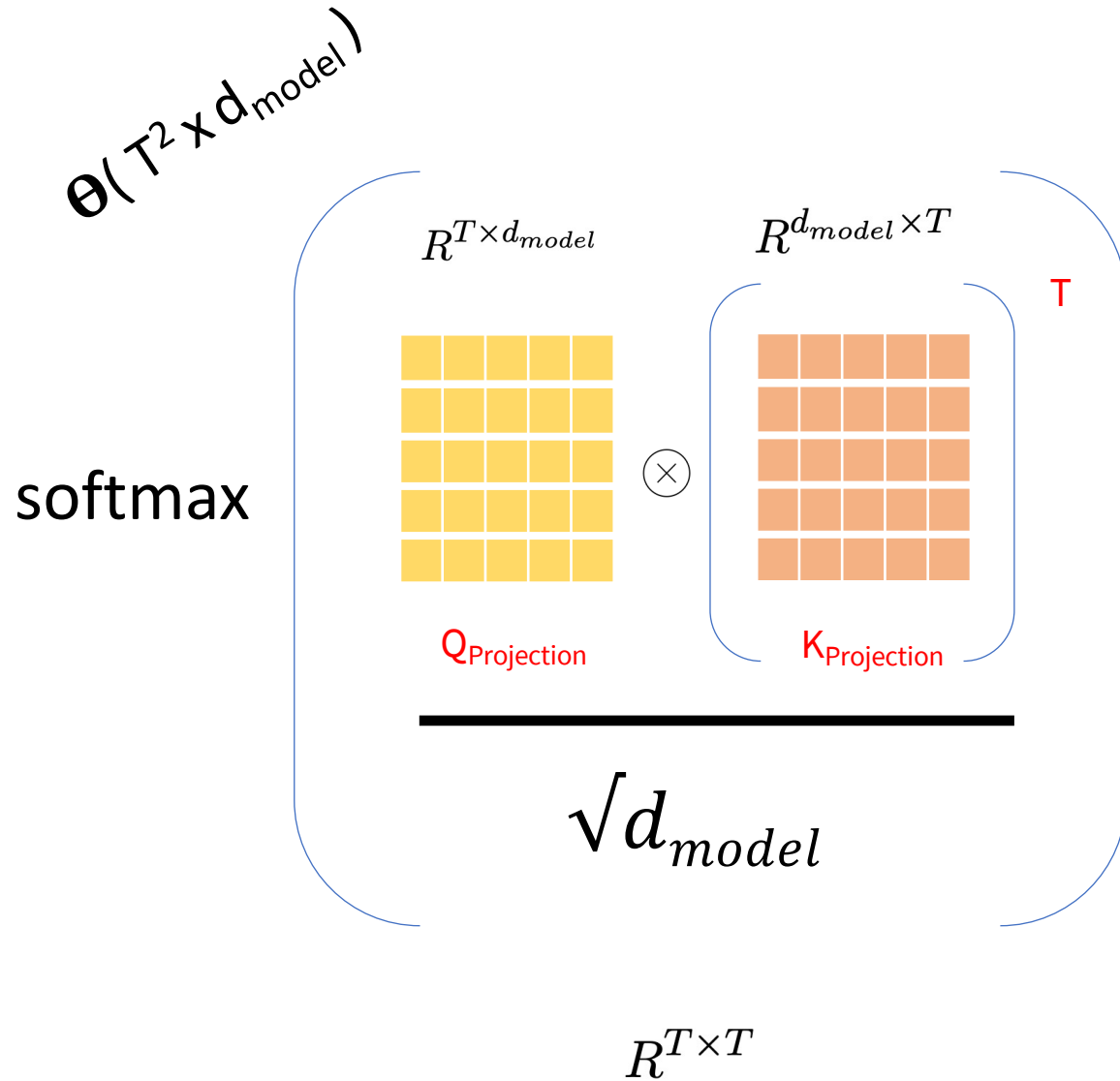Query Inputs = Key Inputs = Value Inputs

# Self Attention

$$R^{d_{model} \times d_{model}}$$



$W_Q$

$W_K$

$W_v$

$$R^{T \times d_{model}}$$



Input Embeddings

# Self Attention

$R^{T \times d_{model}}$    $R^{d_{model} \times d_{model}}$



Input Embeddings    $W_Q$    Q Projections    $R^{T \times d_{model}}$

Input Embeddings    $W_K$    K Projections    $R^{T \times d_{model}}$

Input Embeddings    $W_V$    V Projections    $R^{T \times d_{model}}$

Add & Norm
Feed Forward
N×    Add & Norm
Multi-Head Attention

Positional Encoding

Input Embedding

Inputs

67

# Self Attention

$$\text{softmax} \frac{R^{T \times d_{model}} \otimes R^{d_{model} \times T}}{\sqrt{d_{model}}}$$

$R^{T \times d_{model}}$

$R^{d_{model} \times T}$

T

Q$_{Projection}$

K$_{Projection}$

$R^{T \times T}$



Add & Norm

Feed
Forward

N×

Add & Norm

Multi-Head
Attention

Positional
Encoding

Input
Embedding

Inputs

# Self Attention

$$\theta(T^2 \times d_{model})$$

$$R^{T \times d_{model}}$$

$$R^{d_{model} \times T}$$

T

softmax

$\otimes$

$Q_{Projection}$

$K_{Projection}$

$$\sqrt{d_{model}}$$

$$R^{T \times T}$$

Add & Norm

Feed Forward

N×

Add & Norm

Multi-Head Attention

Positional Encoding

Input Embedding

Inputs

# Self Attention

$\theta(T^2 \times d_{model})$

softmax

$R^{T \times d_{model}}$

$R^{d_{model} \times T}$

T

$R^{T \times d_{model}}$



$Q_{Projection}$

$K_{Projection}$

$V_{Projection}$

$$\frac{}{\sqrt{d_{model}}}$$

$R^{T \times T}$

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

N×

Positional Encoding

Input Embedding

Inputs

# Self Attention

$$R^{T \times d_{model}}$$



Attention: Z

# Self Attention



The | animal | didn't | cross | the | street | because | it | was | too | wide

coreference resolution ✓

# Self Attention

The animal didn't cross the street because **it** was too wide

**?**

Sentence boundaries ?

coreference resolution ✔

Context ?

Semantic relationships ?
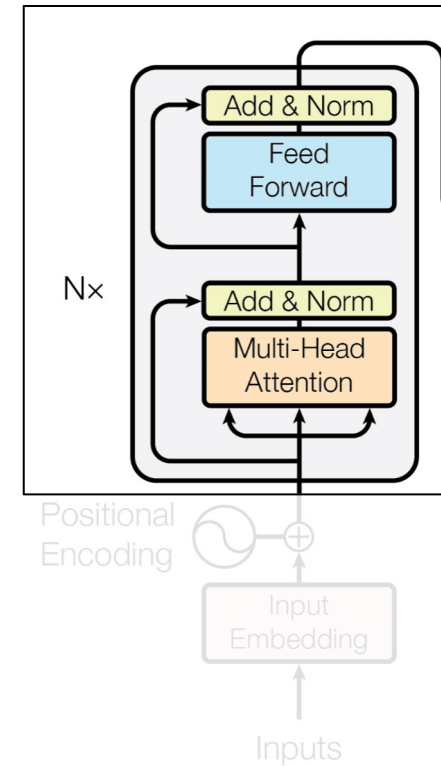
Part of Speech ?

Comparisons ?

# Self Attention

$$R^{d_{model} \times d_{model}}$$



$W_Q$          $W_K$          $W_V$

$$R^{T \times d_{model}}$$



Input Embeddings

# Multi-Head Attention

$$R^{d_{model} \times d_h}$$



$W_{Q1}, W_{Q2}, \dots W_{QH,}$

$W_{K1}, W_{K2}, \dots W_{KH,}$

$W_{V1}, W_{V2}, \dots W_{VH,}$

Input Embeddings

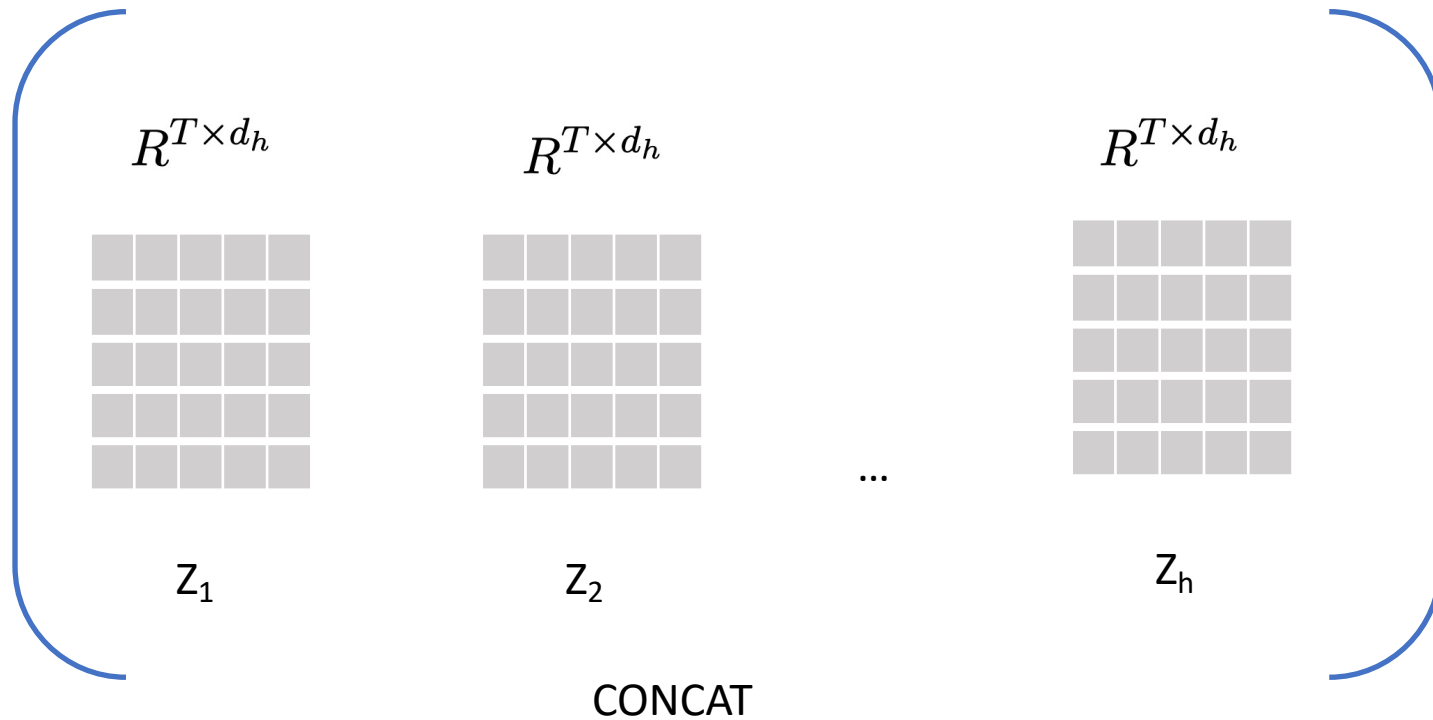$$R^{T \times d_{model}}$$

$$d_h = \frac{d_{model}}{h}$$

# Multi-Head Attention

$R^{T \times d_{model}}$

$R^{d_{model} \times d_h}$

Inputs

$\otimes$

$W_{Qi}$    1    2    ..    H

$\longrightarrow$

$Q_i$    1    2    ..    H

$R^{T \times d_h}$

Inputs

$\otimes$

$W_{Ki}$    1    2    ..    H

$\longrightarrow$

$K_i$    1    2    ..    H

$R^{T \times d_h}$

Inputs

$\otimes$

$W_{Vi}$    1    2    ..    H

$\longrightarrow$

$V_i$    1    2    ..    H

$R^{T \times d_h}$

N×

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Positional Encoding

Input Embedding

Inputs

76

# Multi-Head Attention



softmax

$R^{T \times d_h}$     $R^{d_h \times T}$     $R^{T \times d_h}$

T

$\otimes$     $\otimes$

$Q_i$     $K_i$     $V_i$

$$\sqrt{d_{model}}$$

$R^{T \times T}$

for all i ∈ [1, h]

77

# Multi-Head Attention

$R^{T \times d_h}$  $R^{T \times d_h}$  $R^{T \times d_h}$

$Z_1$  $Z_2$  ...  $Z_h$

CONCAT

Multi Head Attention : Z

$d_h = \dfrac{d_{model}}{h}$

$R^{T \times d_{model}}$

Add & Norm

Feed Forward

N×

Add & Norm

Multi-Head Attention

Positional Encoding

Input Embedding

Inputs

# Multi-Head Attention

| The | animal | didn't | cross | the | street | because | it | was | too | wide |

Sentence boundaries ?  ✓   coreference resolution  ✓   Context ?  ✓

Semantic relationships ?  ✓   Part of Speech ?  ✓   Comparisons ?  ✓

# Add & Norm

Input

Norm(Z)

**Normalization(Z)**

- Mean 0, Std dev 1

- Stabilizes training

- Regularization effect

**Add -> Residuals**

- Avoid vanishing gradients

- Train deeper networks

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

Nx

Positional
Encoding

Input
Embedding

Inputs

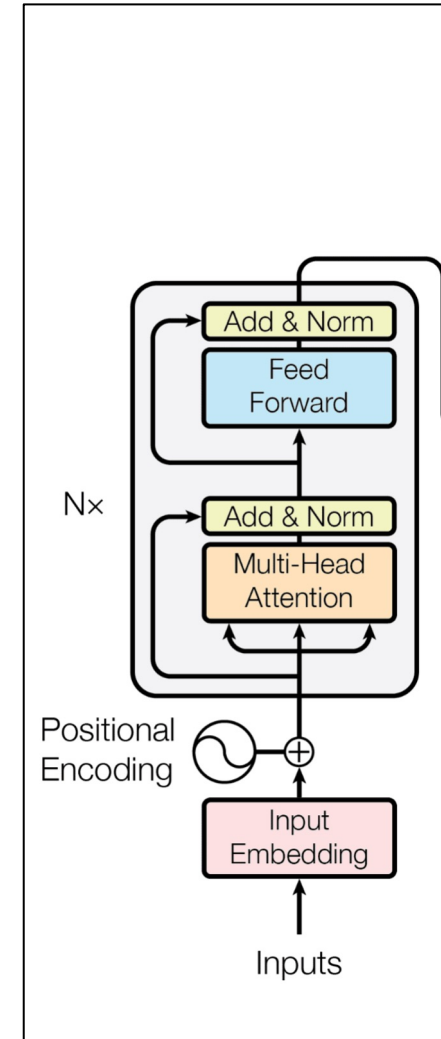# Feed Forward

**Feed Forward**

- Non Linearity

- Complex Relationships

- Learn from each other

# Add & Norm

# Encoders

## Encoder

ENCODER

# Encoders

**Encoder**

ENCODER

.
.
.

ENCODER

Input to Encoder$_{i+1}$

↑

Output from Encoder$_i$

ENCODER



Add & Norm

Feed Forward

N×

Add & Norm

Multi-Head Attention

Positional Encoding

Input Embedding

Inputs

# Transformers

- ✓ **Tokenizaton**
- ✓ **Input Embeddings**
- ✓ **Position Encodings**
- ✓ **Residuals**
- ✓ **Query**
- ✓ **Key**
- ✓ **Value**
- ✓ **Add & Norm**
- ✓ **Encoder**
- **Decoder**

- ✓ **Attention**
- ✓ **Self Attention**
- ✓ **Multi Head Attention**
- **Masked Attention**
- **Encoder Decoder Attention**
- **Output Probabilities  / Logits**
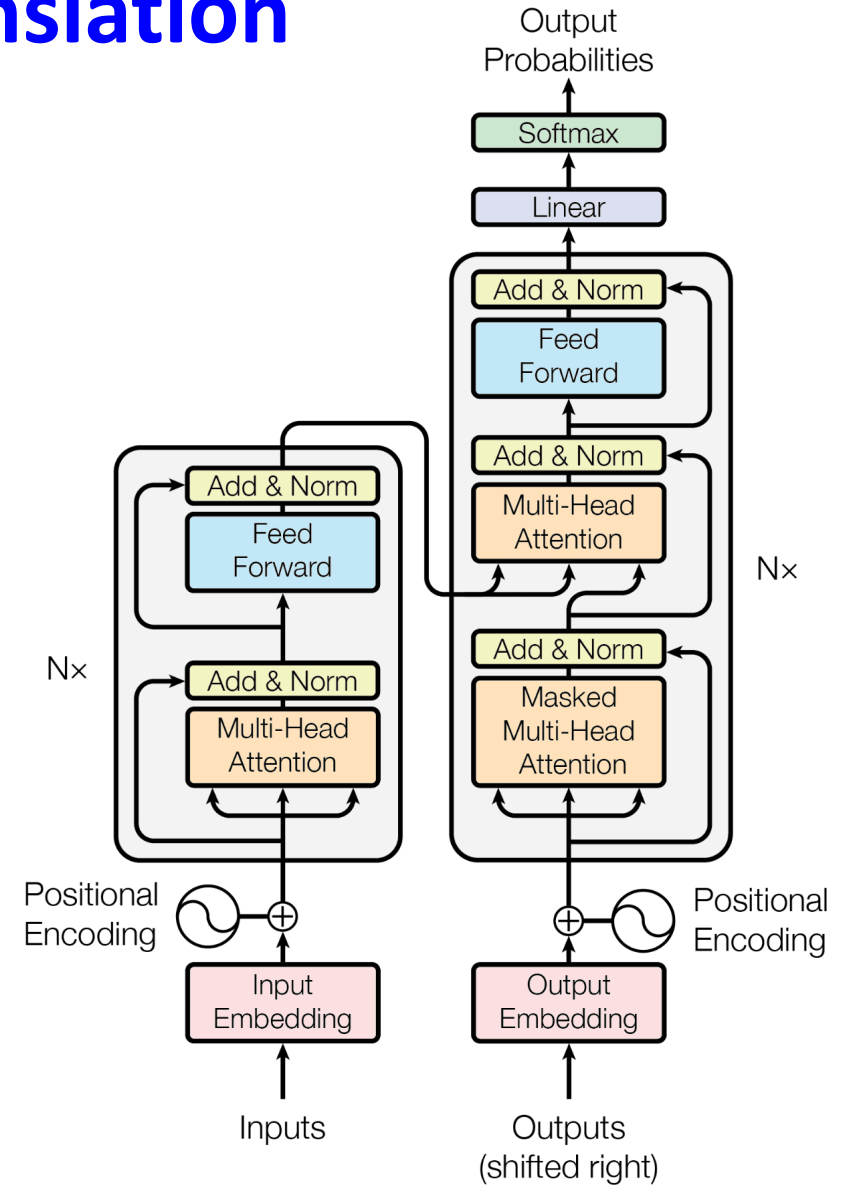- **Softmax**
- **Encoder-Decoder models**
- **Decoder only models**

# Machine Translation
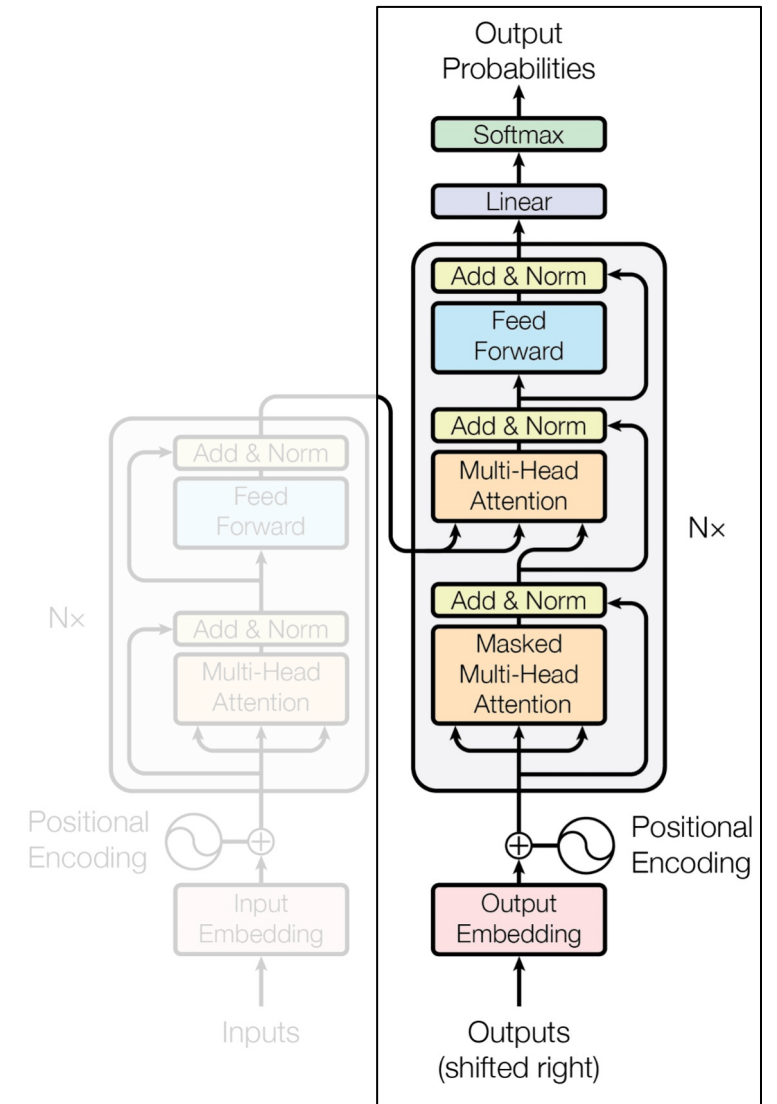
Targets

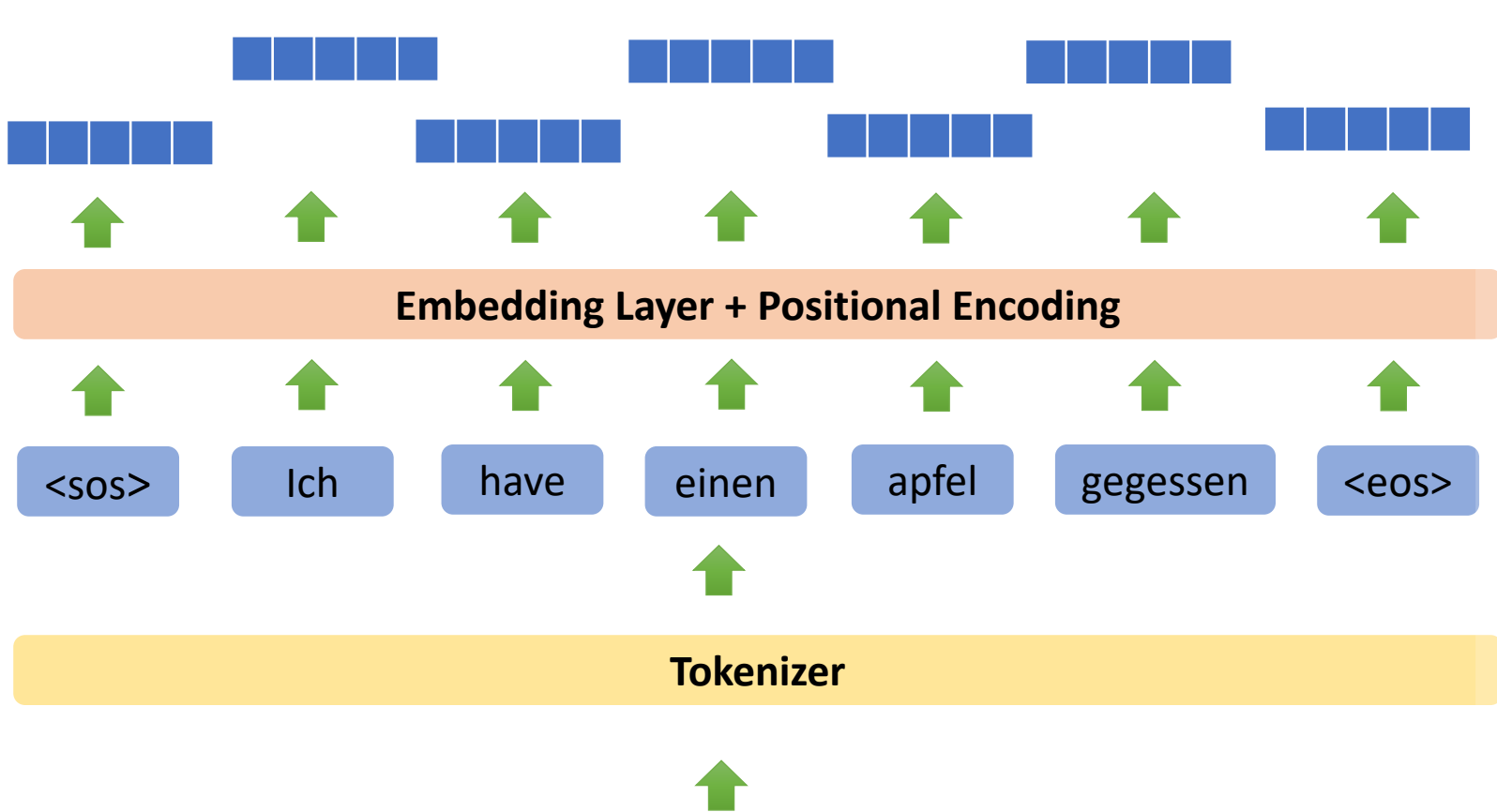Ich have einen apfel gegessen

Inputs

I ate an apple



Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Masked Multi-Head Attention

N×

N×

Positional Encoding

Positional Encoding

Input Embedding

Output Embedding

Inputs

Outputs (shifted right)

# Targets

Targets

Ich have einen apfel gegessen

# Targets



Embedding Layer + Positional Encoding
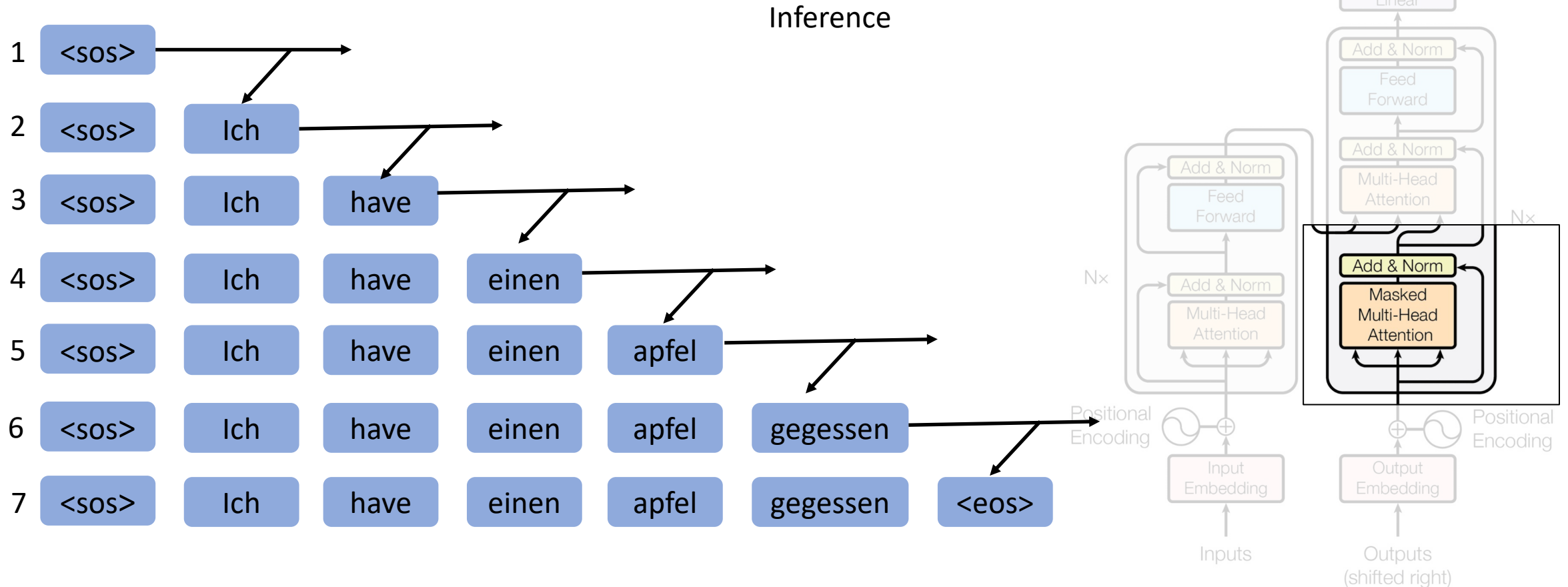
| <sos> | Ich | have | einen | apfel | gegessen | <eos> |

Tokenizer

Ich have einen apfel gegessen

Generate Target Emebeddings

# Masked Multi Head Attention



<sos>     Ich     have     einen     apfel     gegessen     <eos>

# Masked Multi Head Attention

Decoding step by step (using Teacher Forcing)
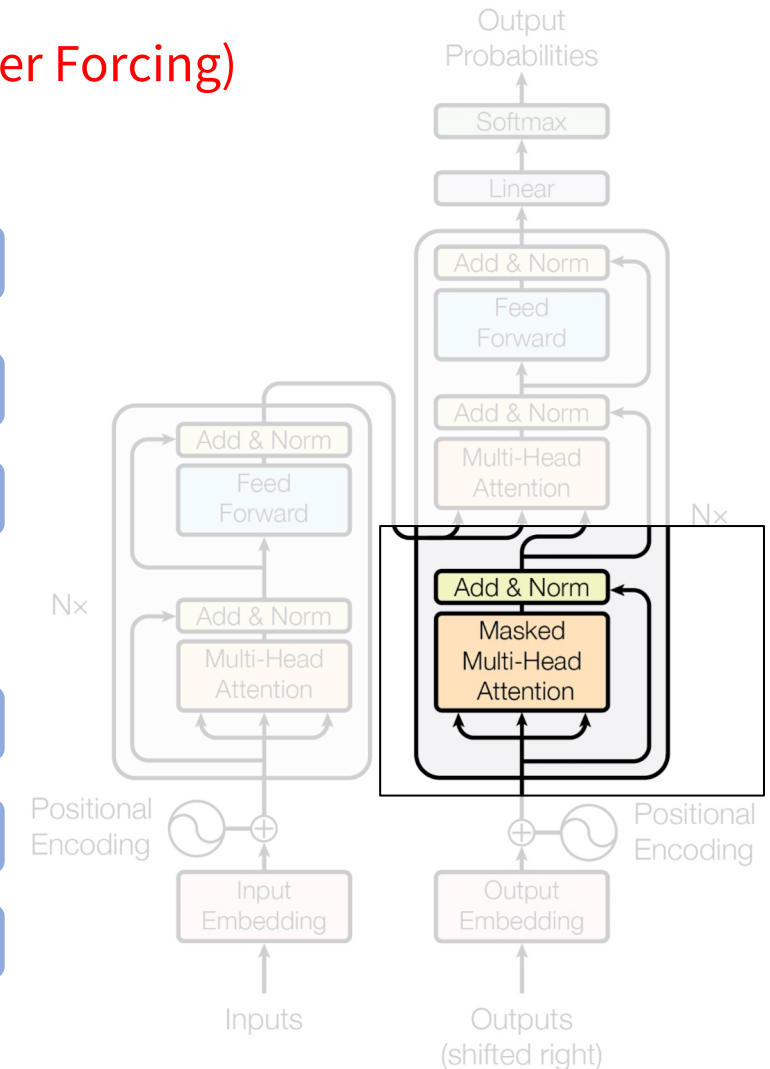
Inference

1  | <sos> |

2  | <sos> | Ich |

3  | <sos> | Ich | have |

4  | <sos> | Ich | have | einen |

5  | <sos> | Ich | have | einen | apfel |

6  | <sos> | Ich | have | einen | apfel | gegessen |

7  | <sos> | Ich | have | einen | apfel | gegessen | <eos> |

# Masked Multi Head Attention

Decoding step by step (using Teacher Forcing)

# Masked Multi Head Attention

Decoding step by step (using Teacher Forcing)

Training

| <sos> | Ich | have | einen | apfel | gegessen | <eos> |

# Masked Multi Head Attention

Decoding step by step (using Teacher Forcing)

Training

| <sos> | Ich | have | einen | apfel | gegessen | <eos> |

*Outputs at time T should only pay attention to outputs*

*until time T-1*

# Masked Multi Head Attention

Decoding step by step (using Teacher Forcing)

# Masked Multi Head Attention

Decoding step by step (using Teacher Forcing)



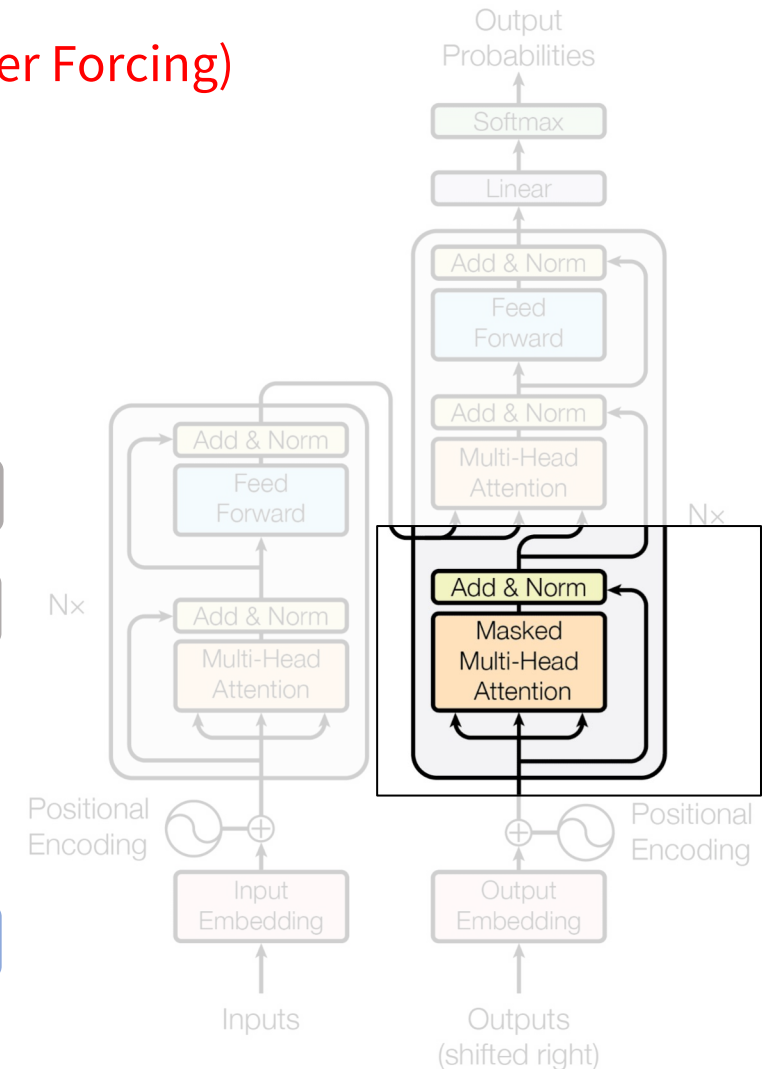Mask the available attention values ?

# Masked Multi Head Attention

Decoding step by step (using Teacher Forcing)

# Masked Multi Head Attention

Masked Multi Head Attention



$R^{T \times T}$      $R^{T \times T}$

QK$^T$      Attention Mask: M      Masked Attention
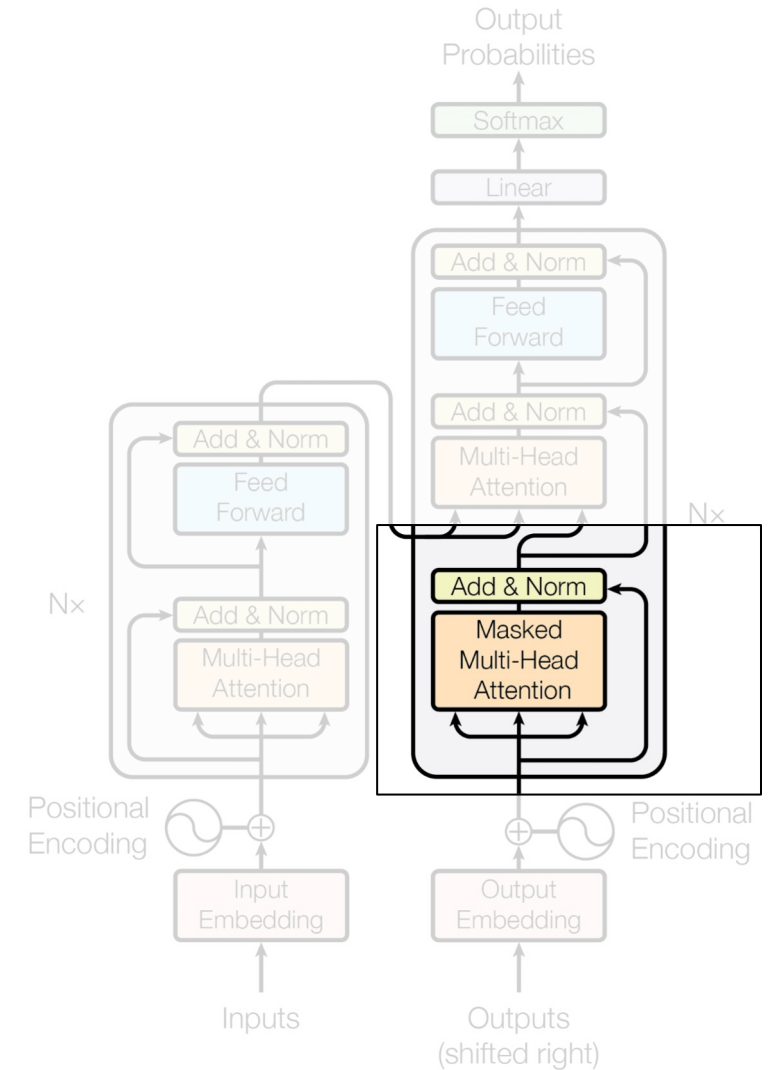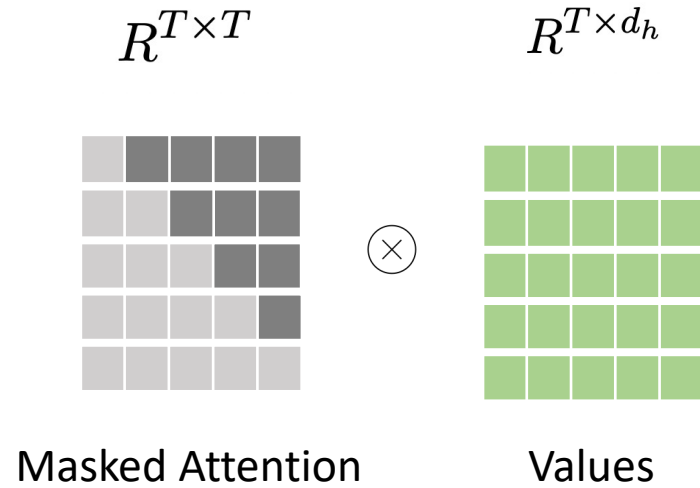
## Masked Multi Head Attention : Z'

# Masked Multi Head Attention

Masked Multi Head Attention

$R^{T \times T}$          $R^{T \times d_h}$

Masked Attention          Values

Masked Multi Head Attention : Z'

# Encoder Decoder Attention



Encoder Decoder Attention ?

Add & Norm

Input
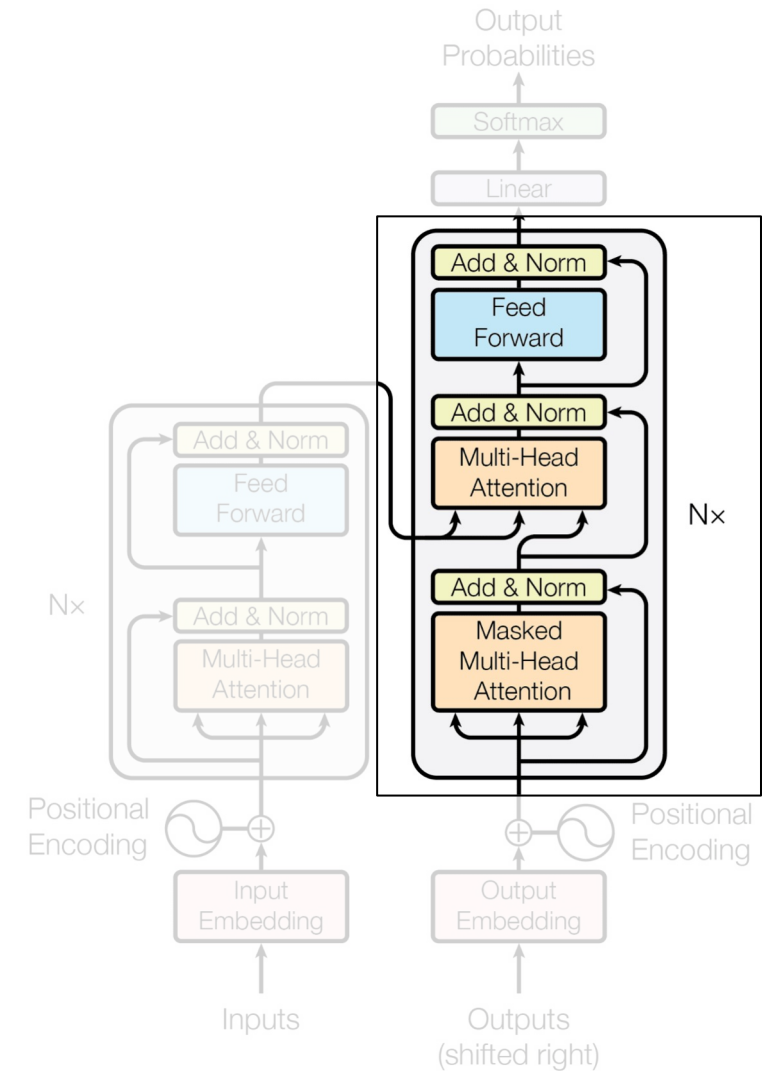
Norm(Z')

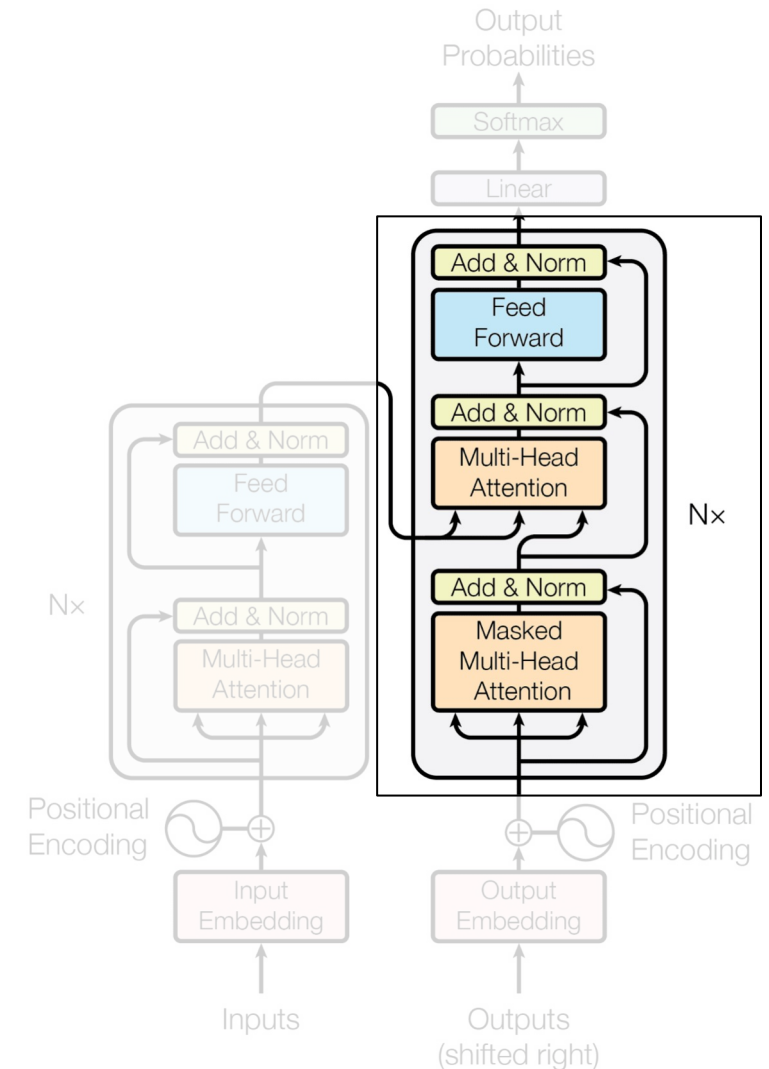# Encoder Decoder Attention

**Encoder Decoder Attention ?**

# Encoder Decoder Attention

Encoder **Self** Attention

1. Queries from Encoder Inputs
2. Keys from Encoder Inputs
3. Values from Encoder Inputs

Decoder **Masked Self** Attention

1. Queries from Decoder Inputs
2. Keys from Decoder Inputs
3. Values from Decoder Inputs

# Attention

{Key, Value store}

{Query: "Order details of order_104"}

{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", ...}}
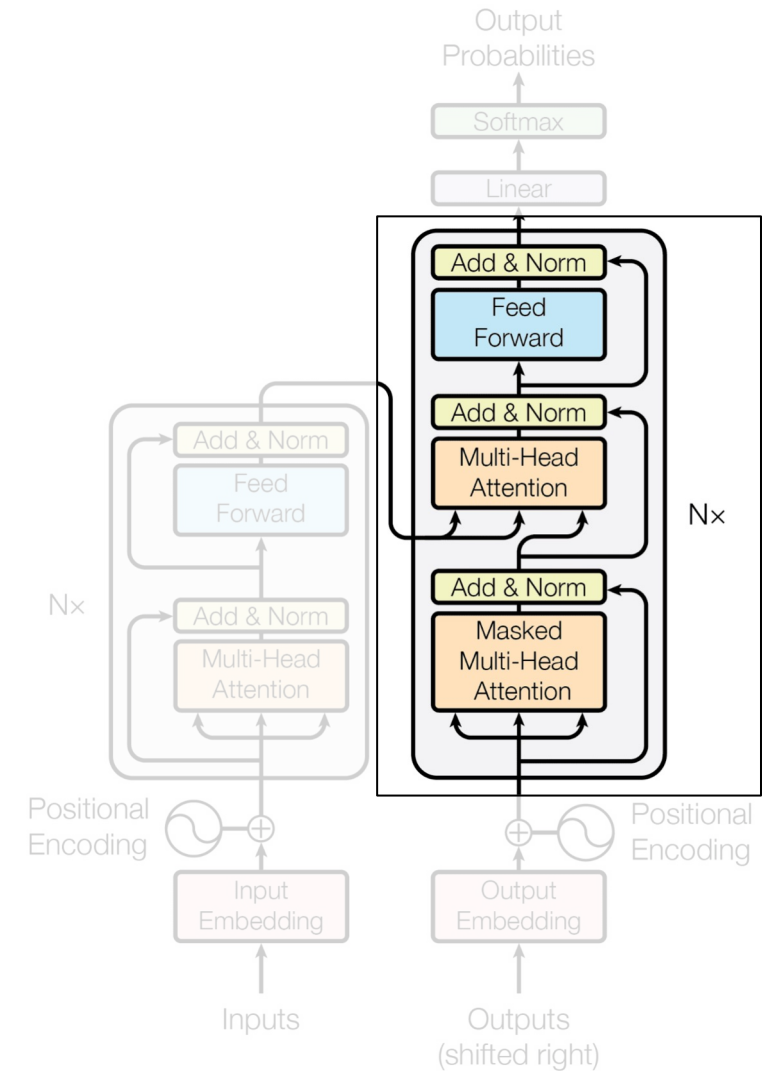
# Encoder Decoder Attention

**Encoder**

Keys from **Encoder Outputs**
Values from **Encoder Outputs**

**Decoder**

Queries from **Decoder Inputs**

NOTE: Every decoder block receives the same FINAL encoder output

# Encoder Decoder Attention

$R^{T_d \times d_{model}}$        Z''

$R^{T_d \times T_e}$      $\text{softmax}(\frac{Q_d K_e^{\ T}}{\sqrt{d_{model}}})$.    $V_e$   $R^{T_e \times d_{model}}$
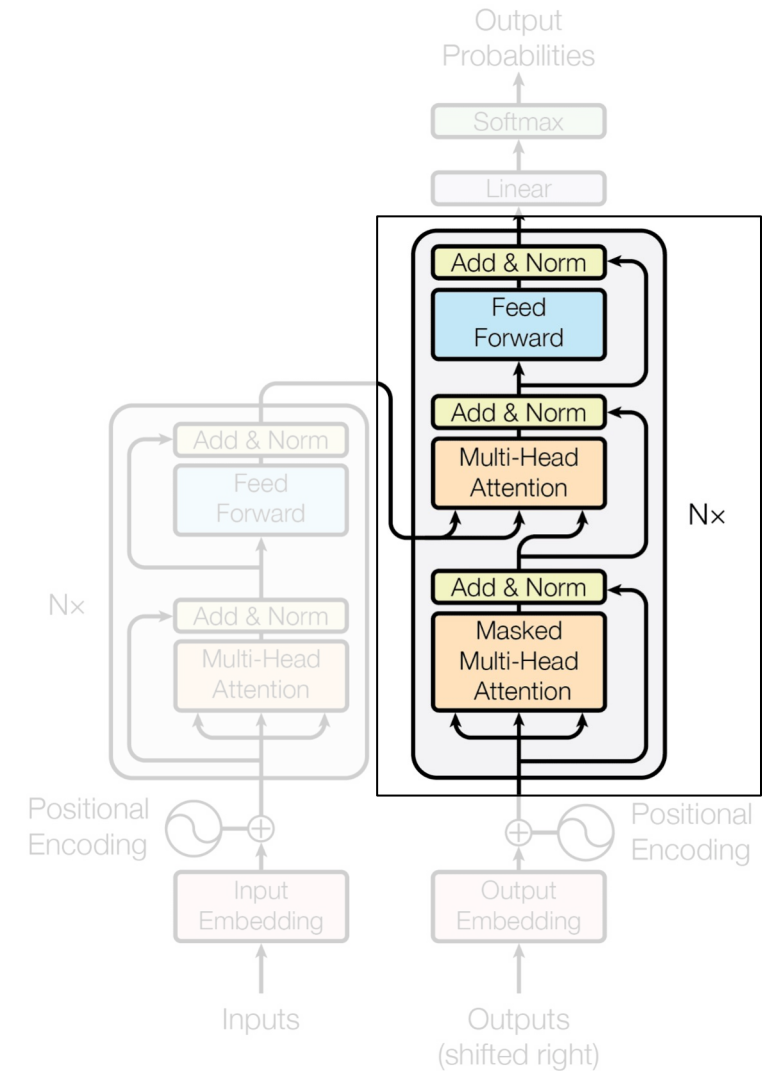
$R^{T_d \times T_e}$      $\text{softmax}(\frac{Q_d K_e^{\ T}}{\sqrt{d_{model}}})$
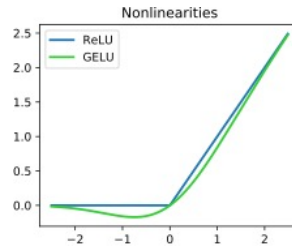
$R^{T_d \times d_{model}}$      $Q_d$   $K_e$      $R^{T_e \times d_{model}}$
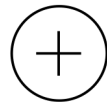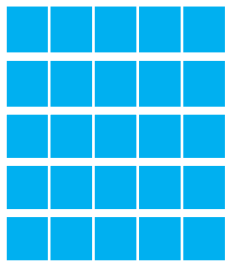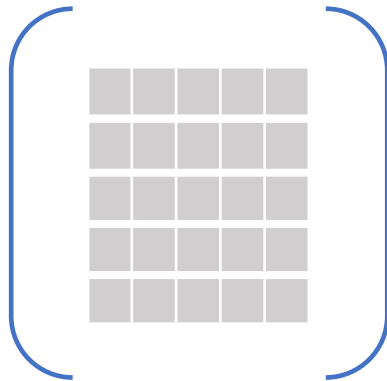
# Encoder Decoder Attention



- Non Linearity
- Complex Relationships
- Learn from each other

Feed Forward

Residuals

Norm(Z'')

Add n Norm Decoder Self Attn
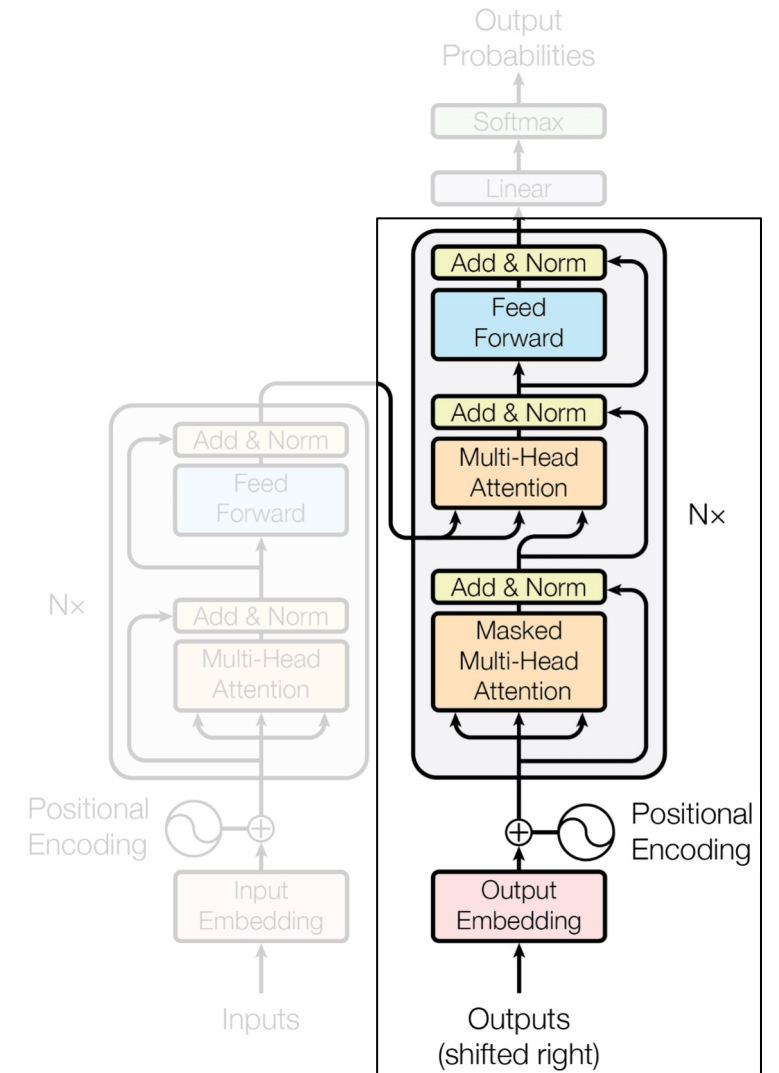
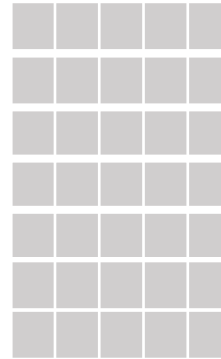# Decoder

DECODER

# Decoder

**DECODER**

.
.
.

**DECODER**

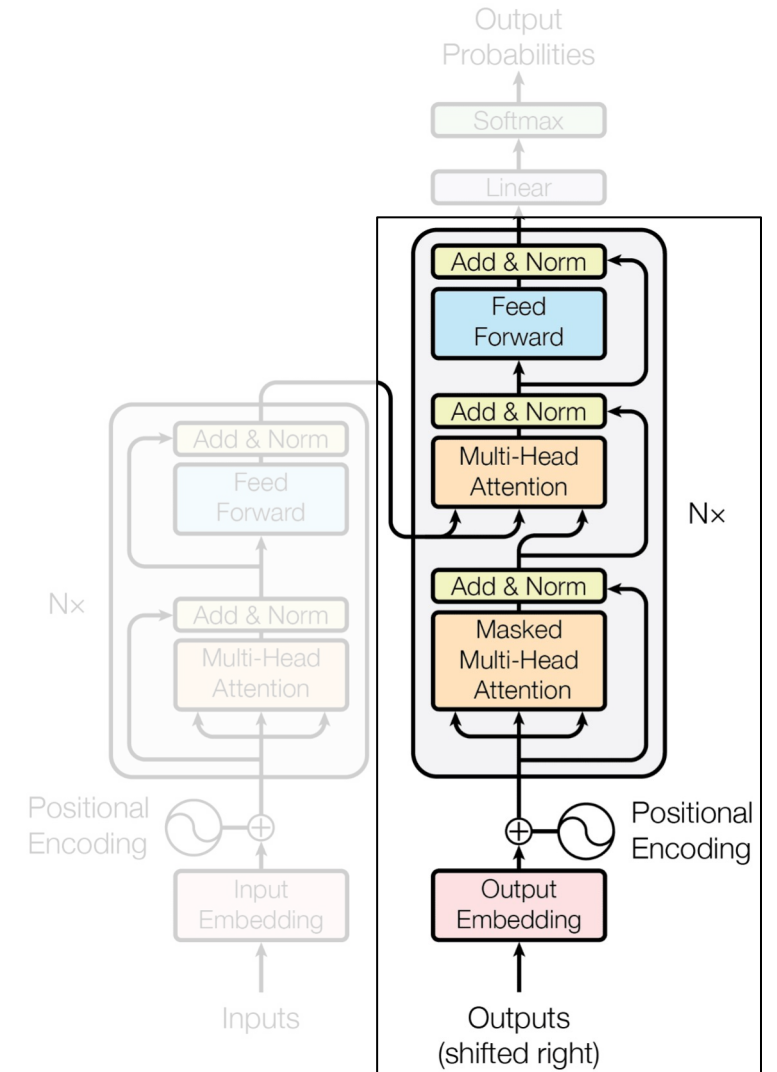**DECODER**

$R^{T_d \times d_{model}}$



Decoder output

# Linear

$R^{V \times d_{model}}$

**Linear weights are often tied with input embedding matrix**

$R^{T_d \times d_{model}}$

softmax

$\bigotimes$

...

Final Decoder Output

$R^{T_d \times V}$

Linear



109

# Softmax

Output Probabilities

V →

$T_d$

$R^{T_d \times V}$

# Poll 2 (@1297)

**Which of the following are true about transformers?**

a.   Transformers can always be run in parallel

b.   Transformer decoders can only be parallelized during training

c.   Positional encodings help parallelize the transformer encoder

d.   Queries, keys, and values are obtained by splitting the input into 3 equal segments

e.   Multiheaded attention helps transformers find different kinds of relations between the tokens

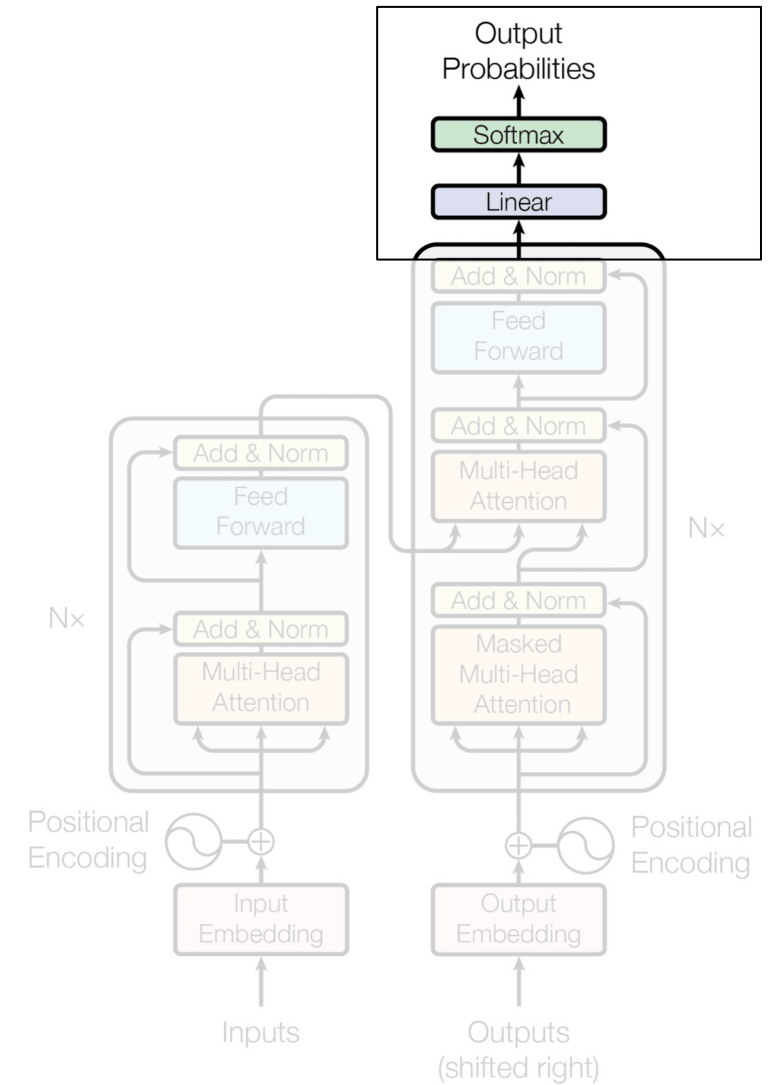f.   During decoding, decoder outputs function as queries and keys while the values come from the encoder
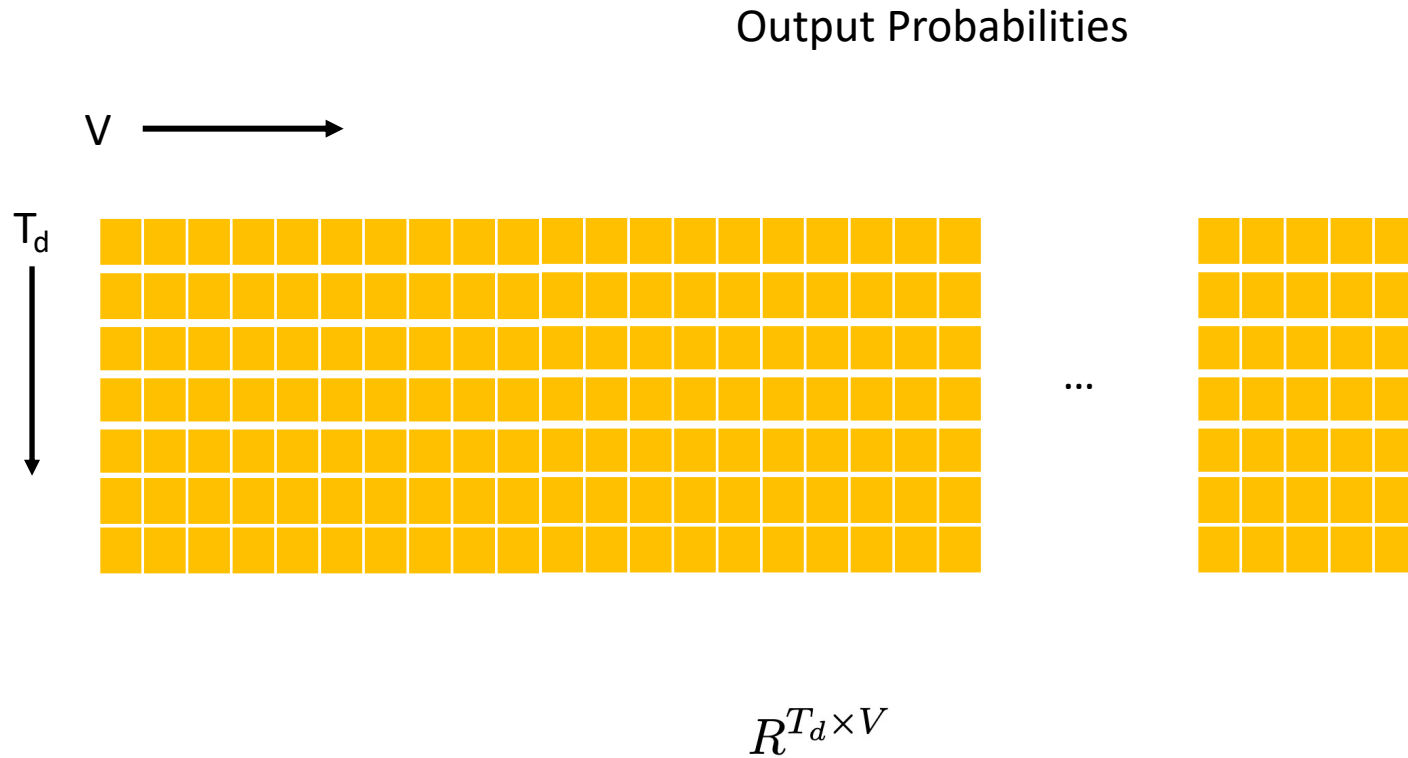
# Poll 2 (@1126)

**Which of the following are true about transformers?**

    a.    Transformers can always be run in parallel

    b.    **Transformer decoders can only be parallelized during training**

    c.    Positional encodings help parallelize the transformer encoder

    d.    Queries, keys, and values are obtained by splitting the input into 3 equal segments

    e.    **Multiheaded attention helps transformers find different kinds of relations between the tokens**

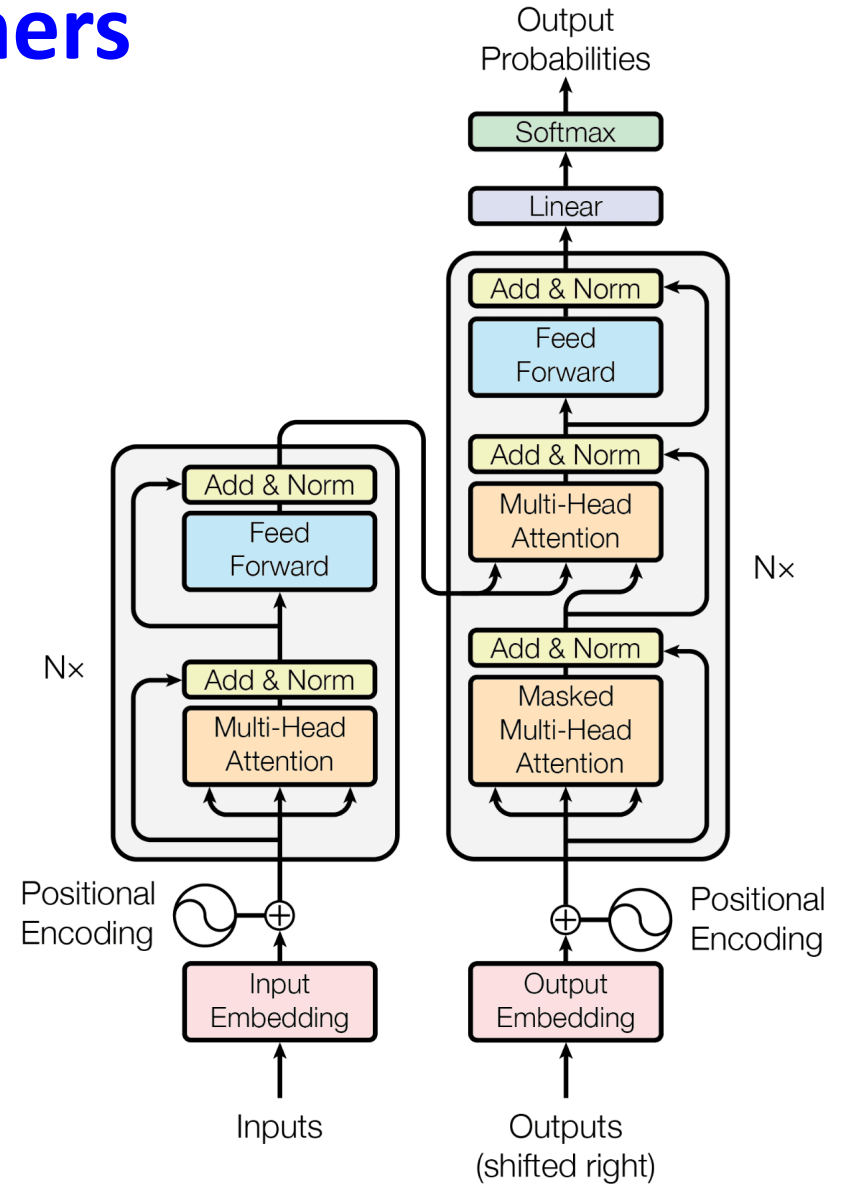    f.    During decoding, decoder outputs function as queries and keys while the values come from the encoder

# Transformers

Targets

Ich have einen apfel gegessen

Inputs

I ate an apple

**Machine Translation**

# Transformers

- ✓ Tokenizaton
- ✓ Input Embeddings
- ✓ Position Encodings
- ✓ Residuals
- ✓ Query
- ✓ Key
- ✓ Value
- ✓ Add & Norm
- ✓ Encoder
- ✓ Decoder

- ✓ Attention
- ✓ Self Attention
- ✓ Multi Head Attention
- ✓ Masked Attention
- ✓ Encoder Decoder Attention
- ✓ Output Probabilities / Logits
- ✓ Softmax
- • Encoder-Decoder models
- • Decoder only models