

Deep Learning



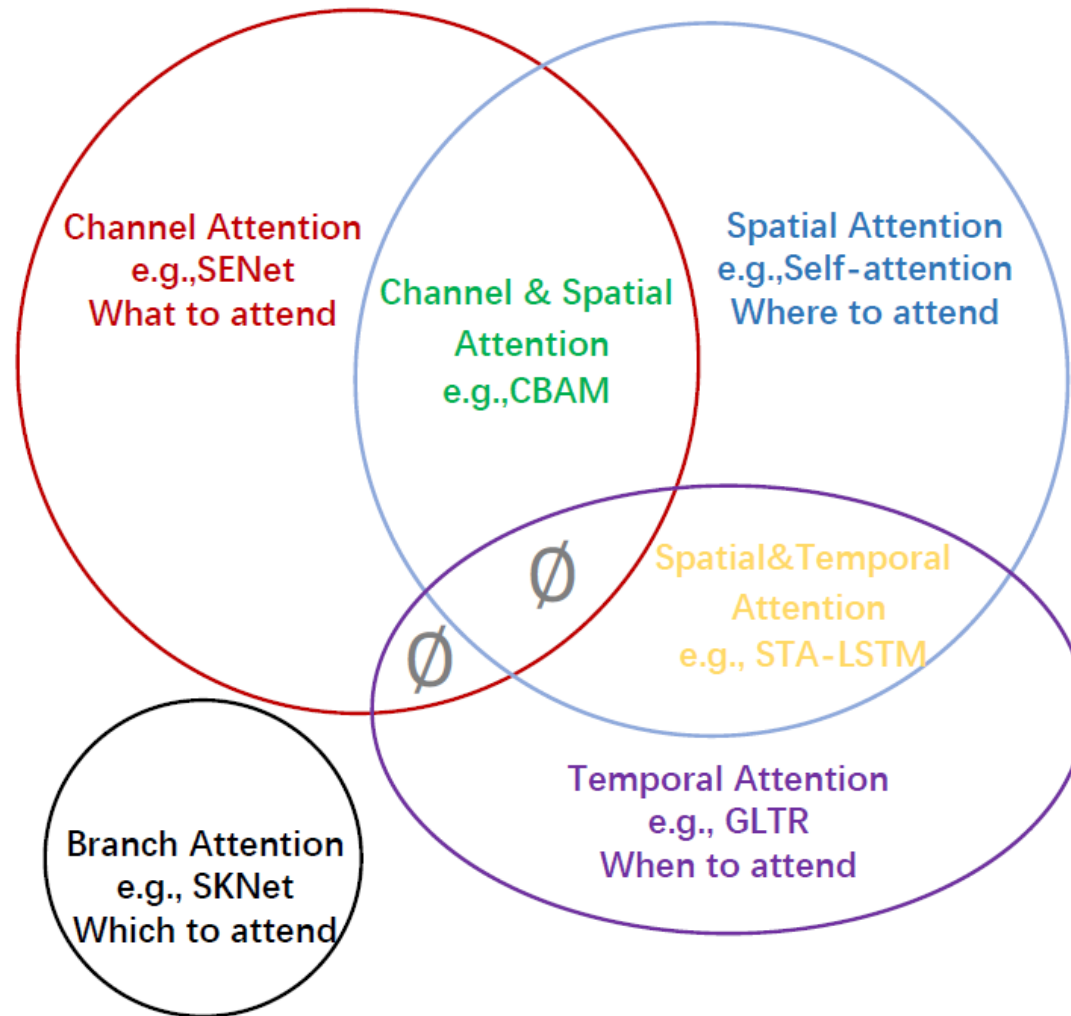
Angshuman Paul

Assistant Professor

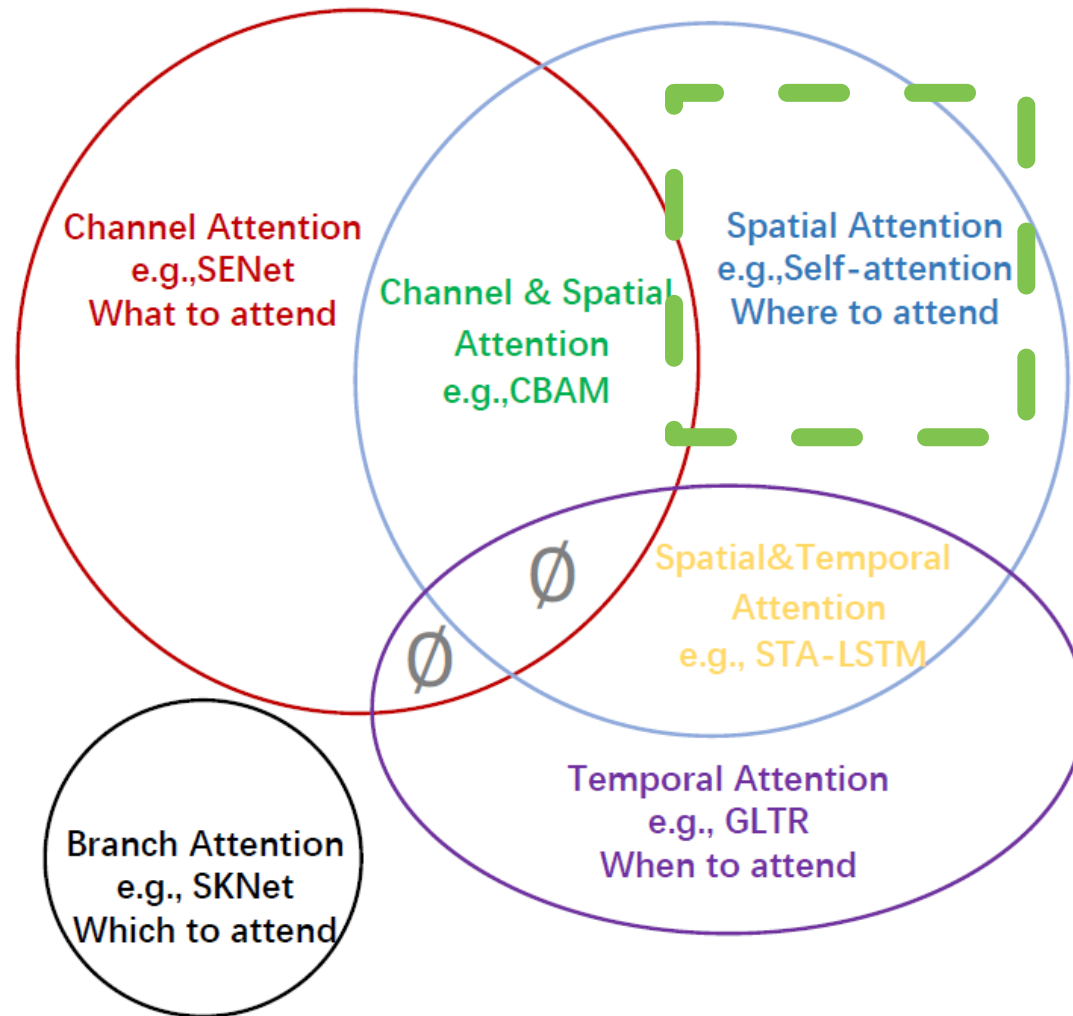
Department of Computer Science & Engineering

Transformers

Attention in Research Papers



Attention in Research Papers



Attention in Language Modelling

- 2017: Attention is All You Need (Vaswani *et al.*)
- Introduction of Transformer

Attention Mechanism

- Mimics the process of retrieval of a value v_i based on key k_i from a database for a query q

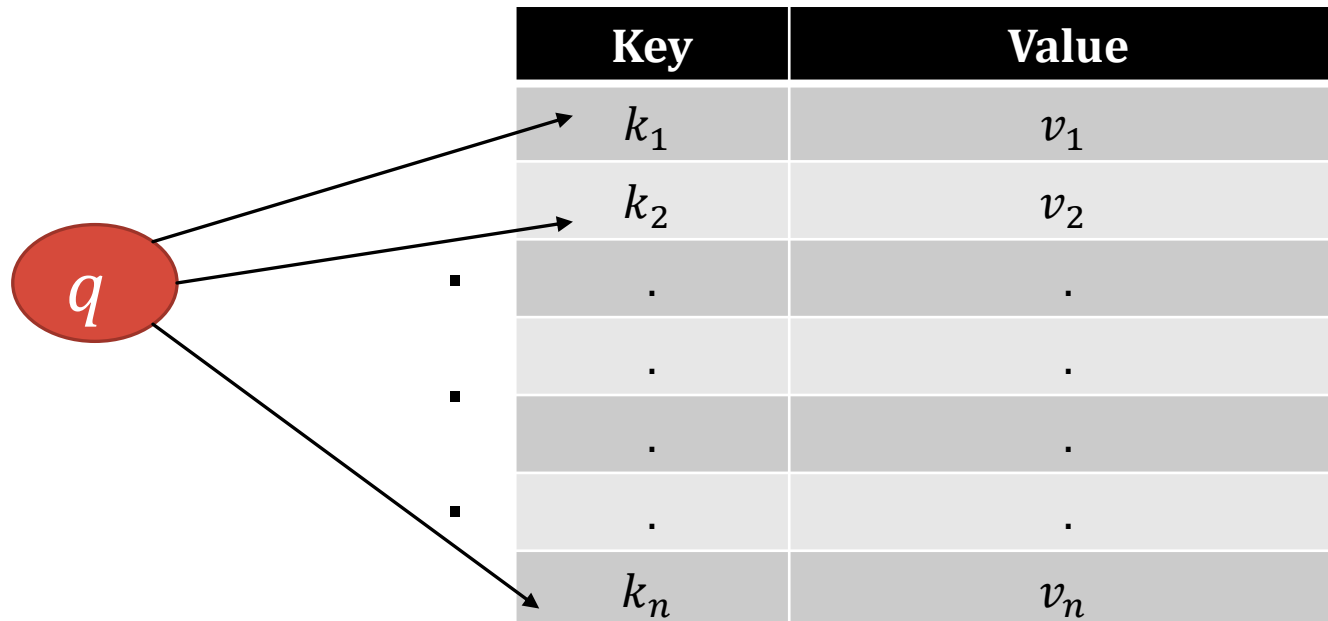
Attention Mechanism

- Mimics the process of retrieval of a value v_i based on key k_i from a database for a query q

Key	Value
k_1	v_1
k_2	v_2
.	.
.	.
.	.
.	.
k_n	v_n

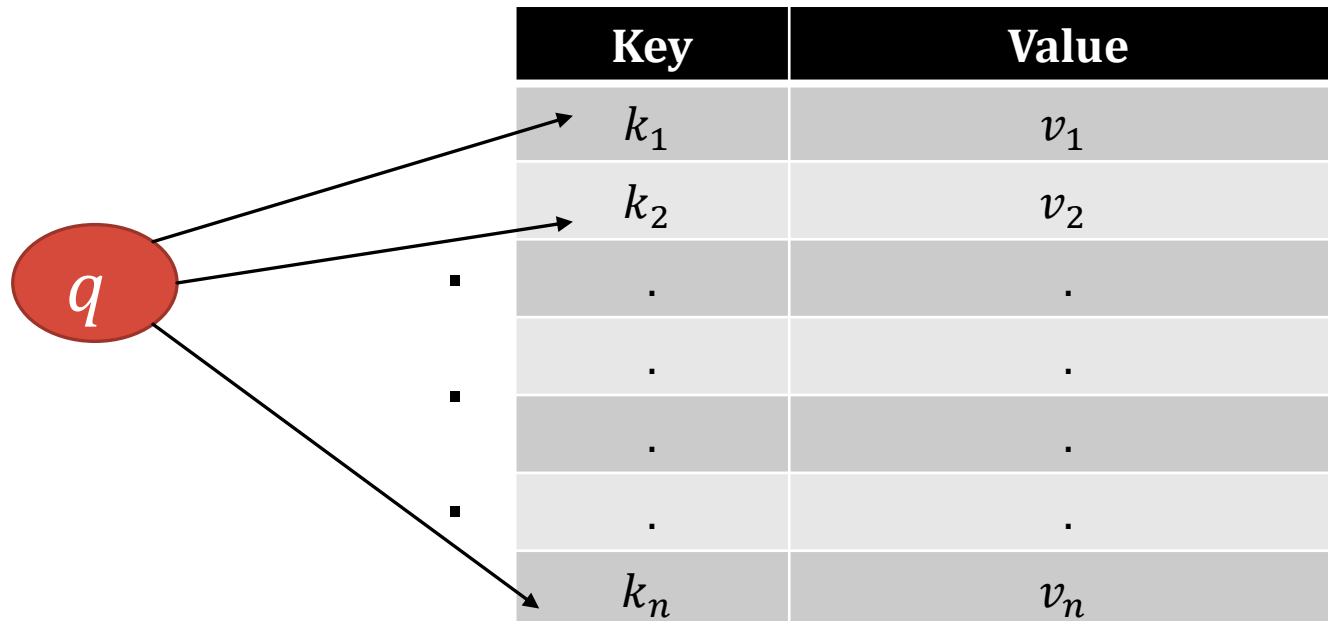
Attention Mechanism

- Mimics the process of retrieval of a value v_i based on key k_i from a database for a query q



Attention Mechanism

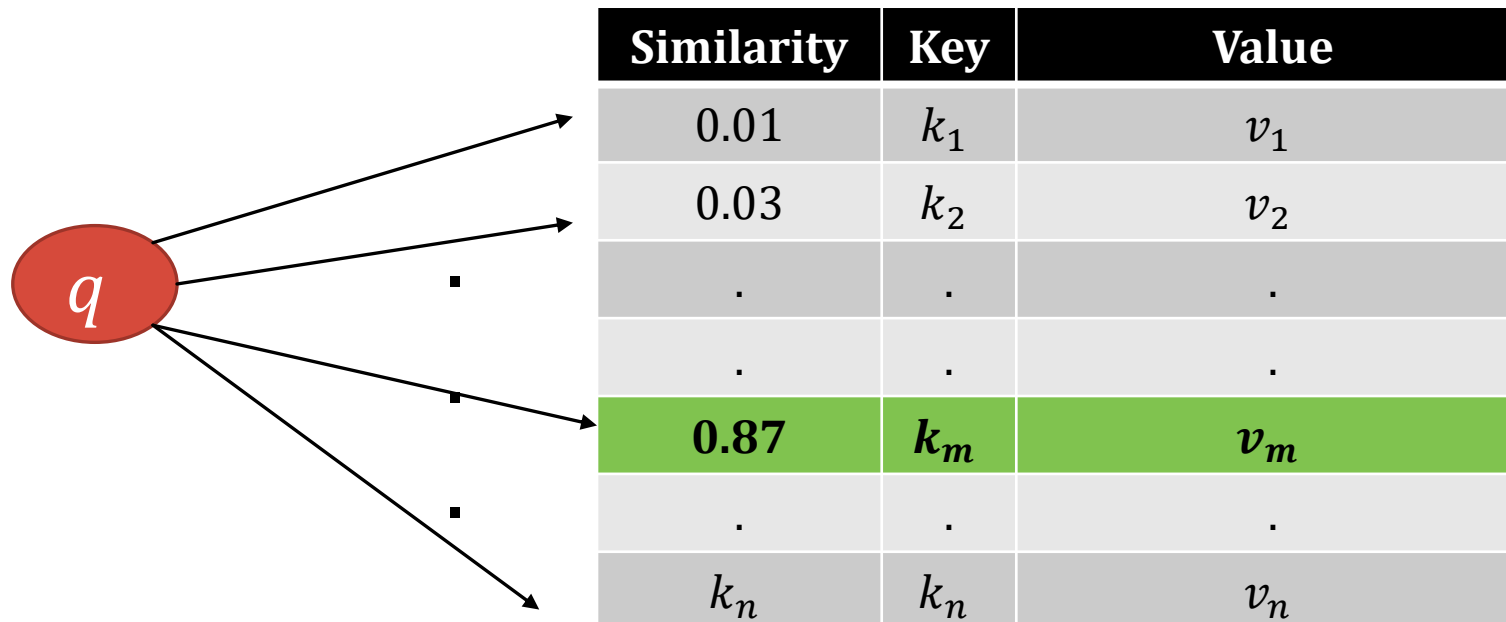
- Mimics the process of retrieval of a value v_i based on key k_i from a database for a query q



Compute similarity
between query and keys

Attention Mechanism

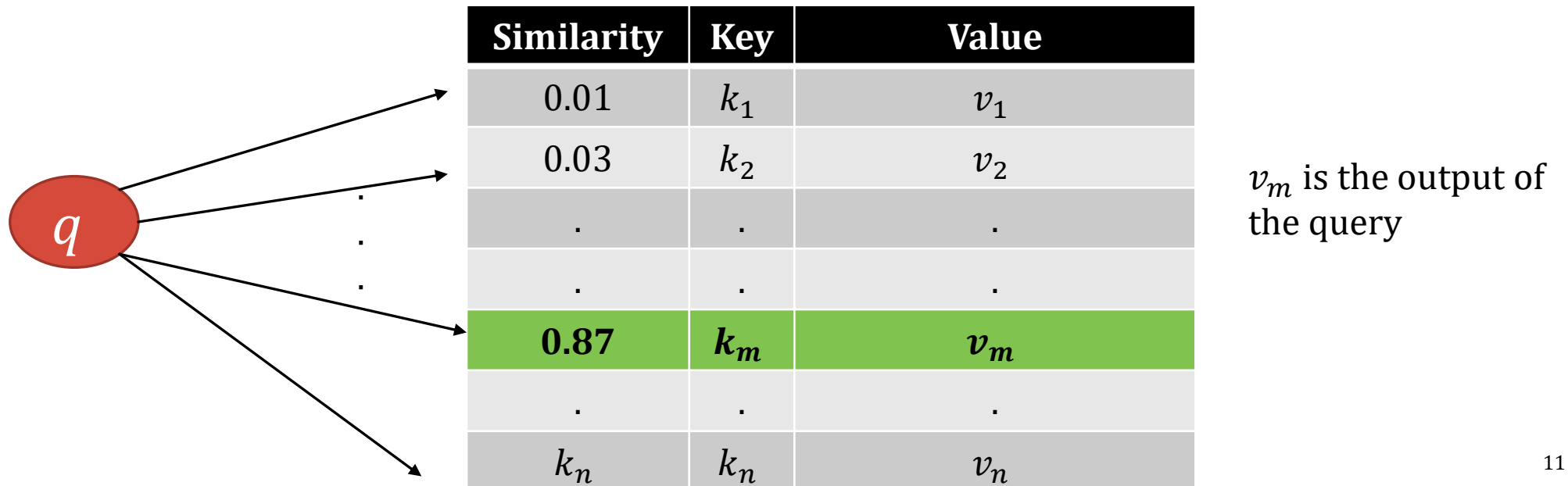
- Mimics the process of retrieval of a value v_i based on key k_i from a database for a query q



Compute similarity
between query and keys

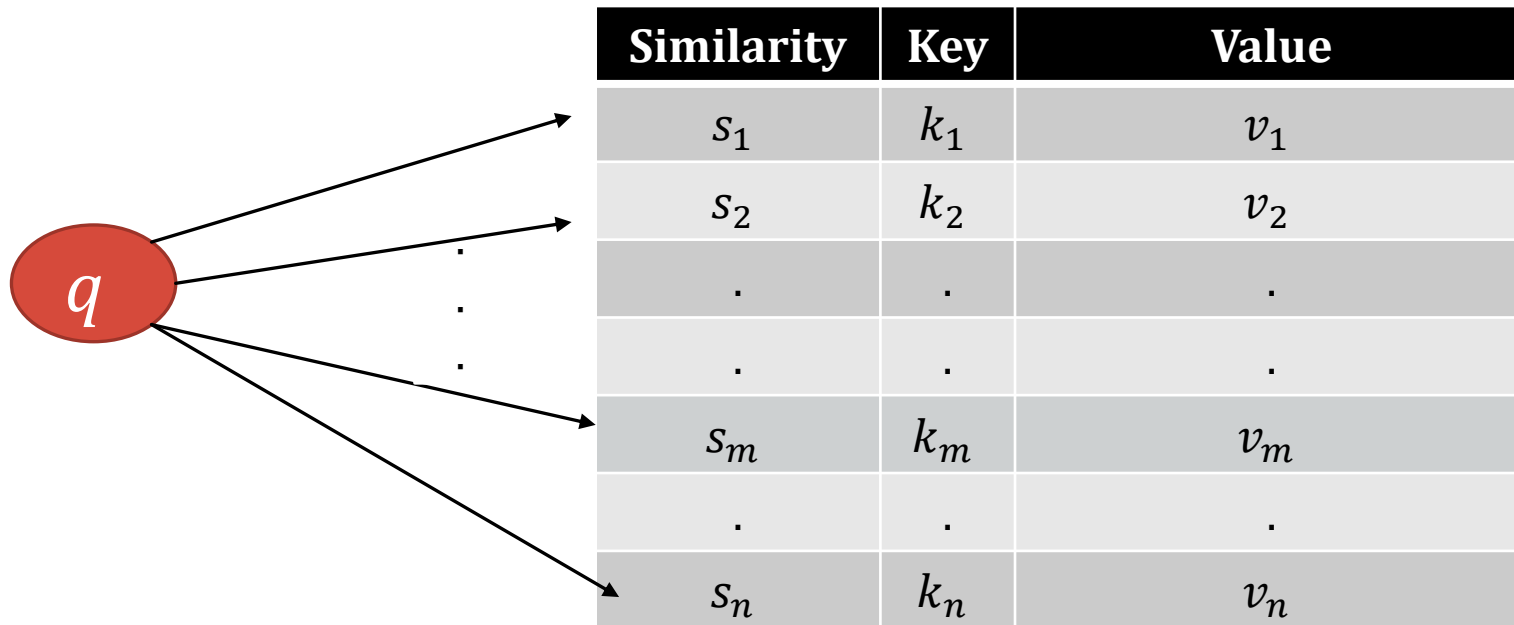
Attention Mechanism

- Mimics the process of retrieval of a value v_i based on key k_i from a database for a query q



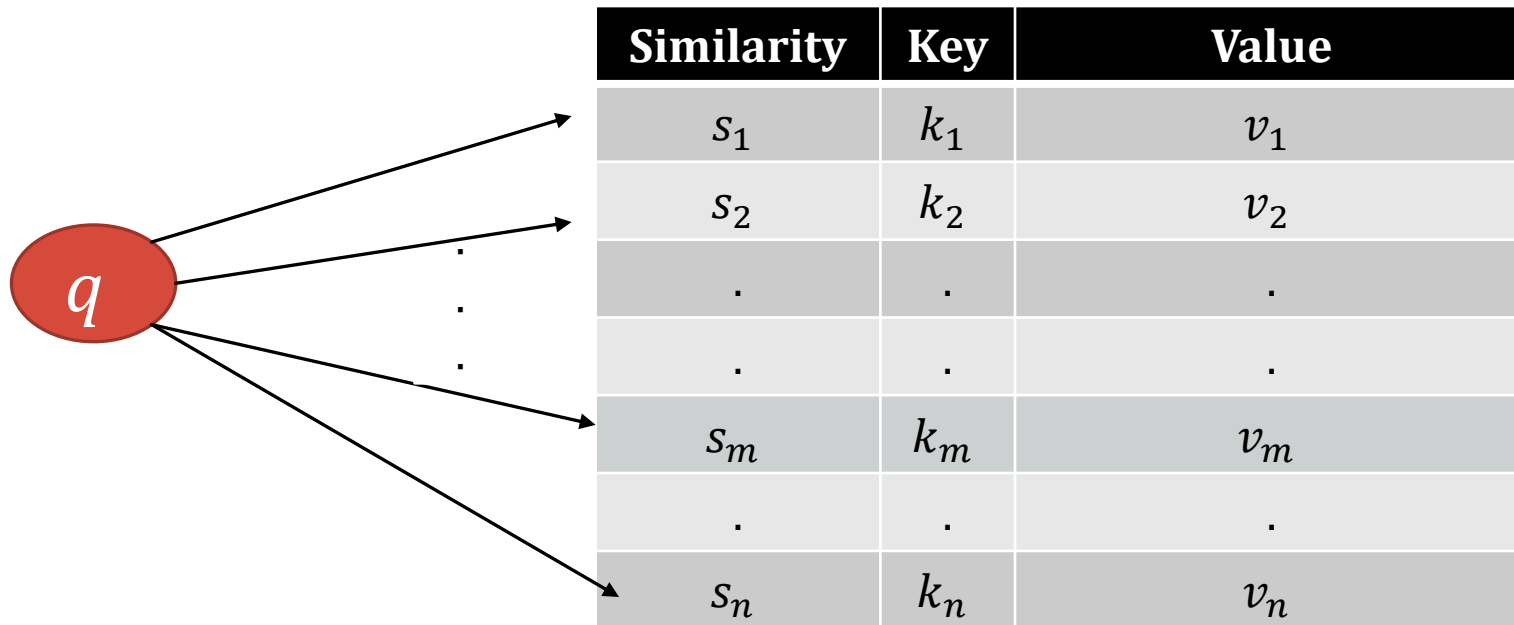
Attention Mechanism

- What else can we do to get the output?



Attention Mechanism

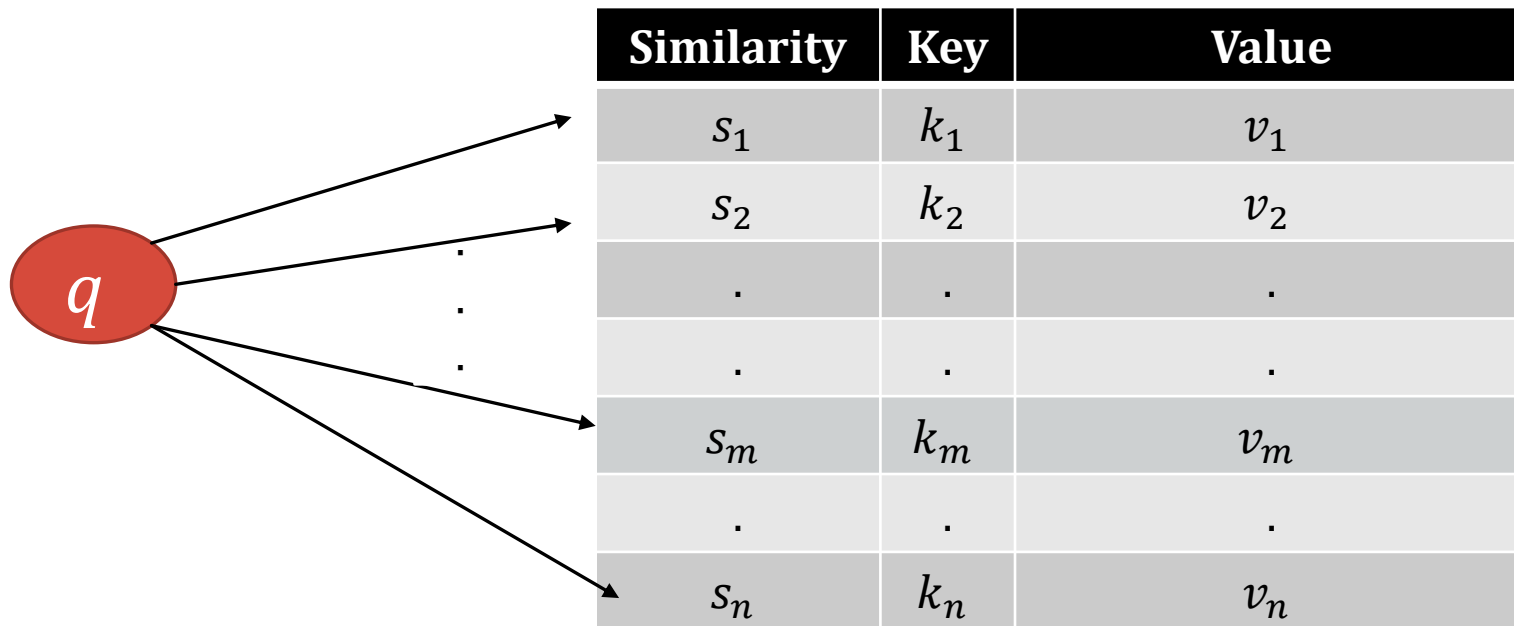
- Instead of a single output, we may take a weighted combination of the outputs



$$Output = \sum_i softmax(s_i) v_i$$

Attention Mechanism

- Instead of a single output, we may take a weighted combination of the outputs



$$\begin{aligned} \text{Output} &= \sum_i \text{softmax}(s_i) v_i \\ &= \sum_i \text{softmax}(\text{similarity}(q, k_i)) \times v_i \end{aligned}$$

How to Calculate Similarity?

- Instead of a single output, we may take a weighted combination of the outputs

$$\text{similarity}(q, k_i) = q^T k_i \quad \textbf{Dot product}$$

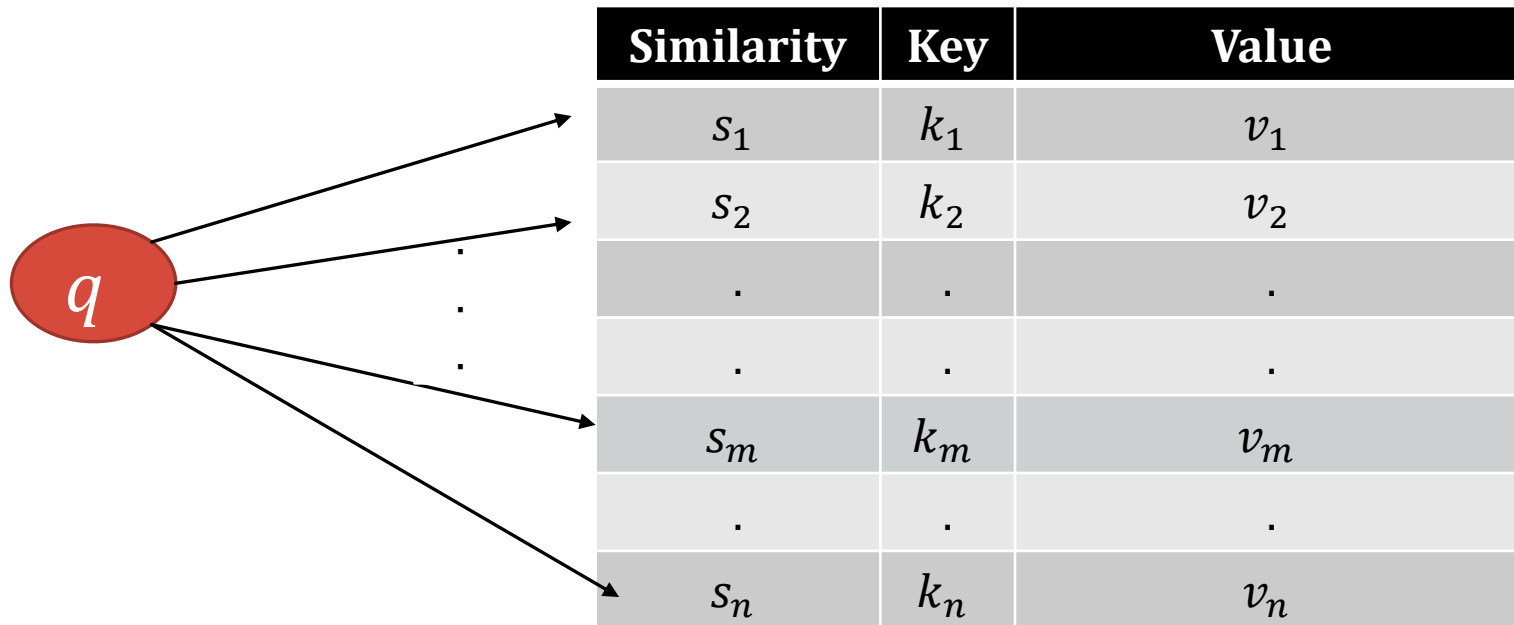
$$\text{similarity}(q, k_i) = \frac{q^T k_i}{\sqrt{d}} \quad \textbf{Scaled Dot product}$$

d is the dimensionality

$$\text{similarity}(q, k_i) = q^T W k_i \quad \textbf{General Dot product}$$

Attention Mechanism

- Attention value can be computed using similarity scores and values

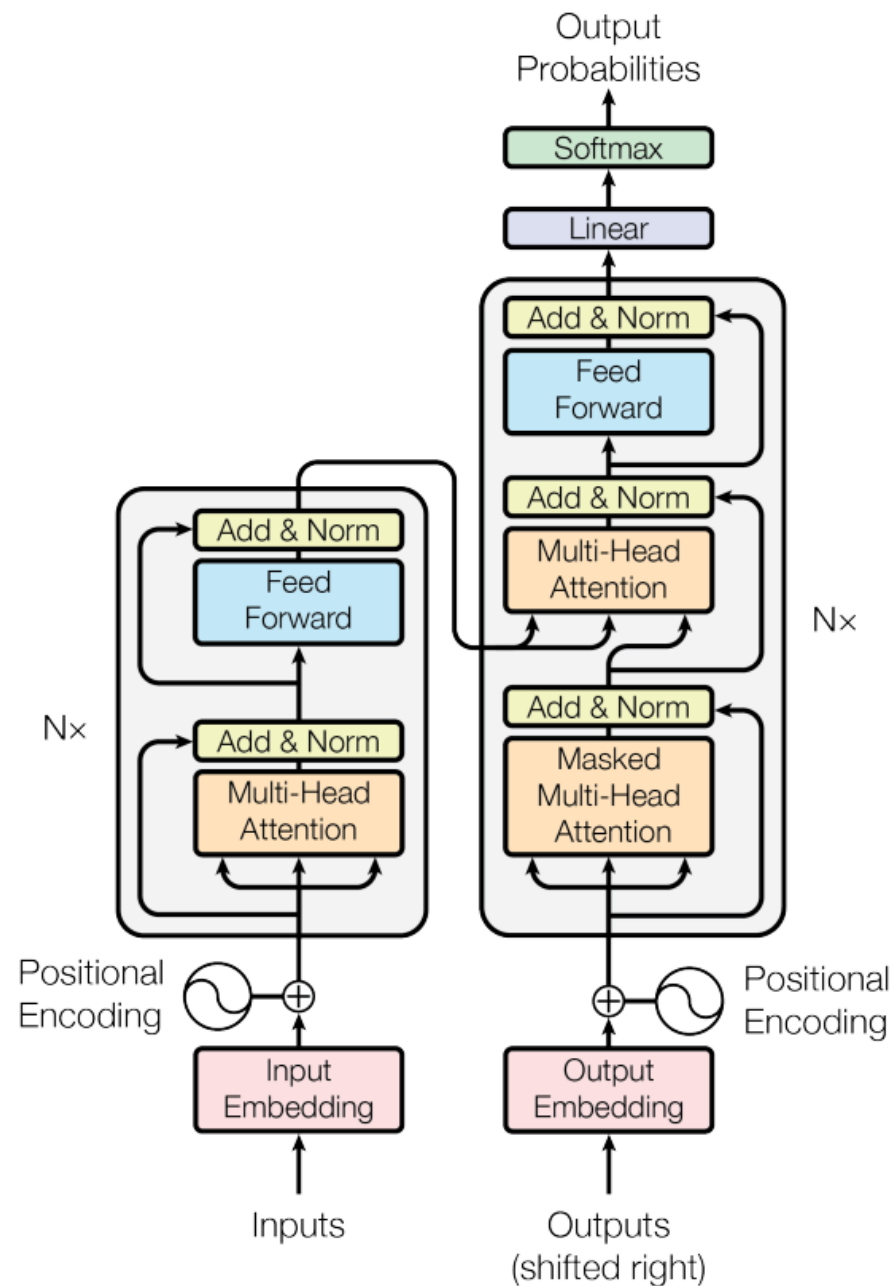


$$\begin{aligned} &\text{Attention value} \\ &= \sum_i \text{softmax}(s_i) v_i \end{aligned}$$

Transformer for NLP

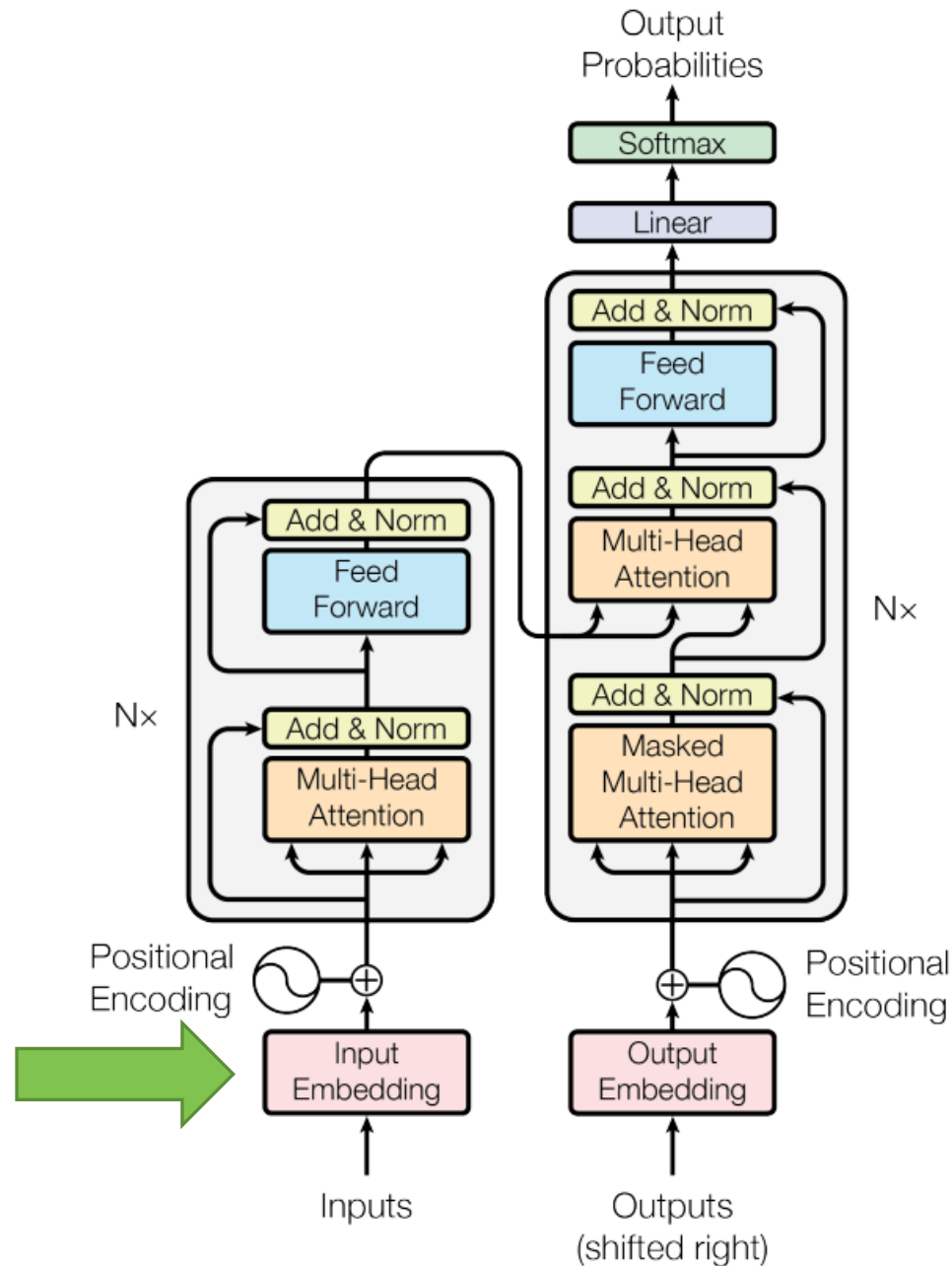
Where there is will, there

**Can an ML model
complete the sentence?**



Transformer for NLP

Word embedding: Vector
representation of words

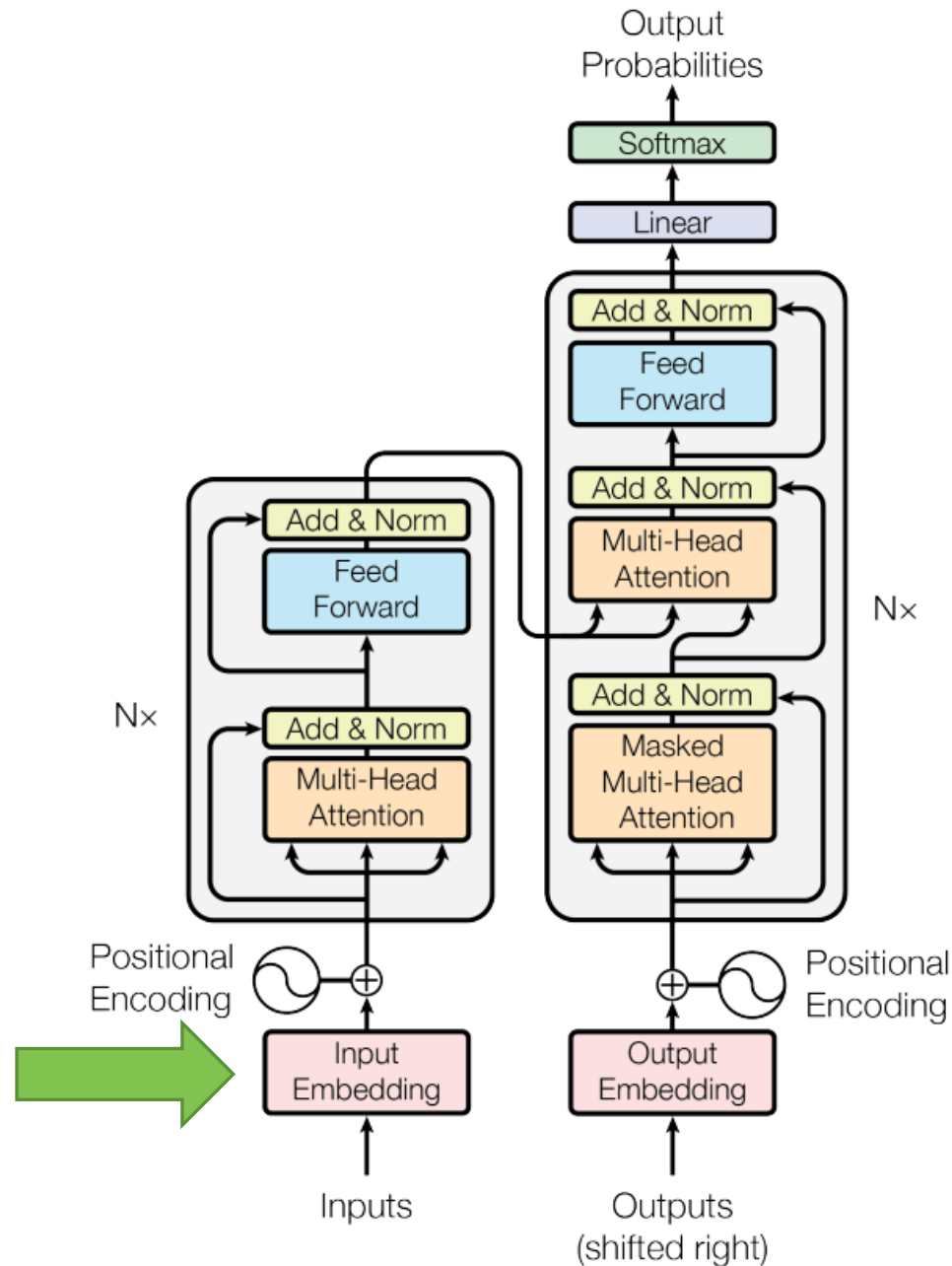


Transformer for NLP

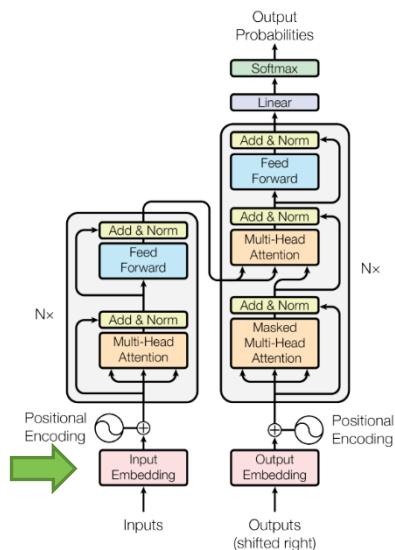
Where there is will, there

Each word is assigned an index based on the dictionary of words

Word	Index
a	0
..	..
is	882
..	..
..	..
there	3582
..	..
where	11442
..	..



Transformer for NLP

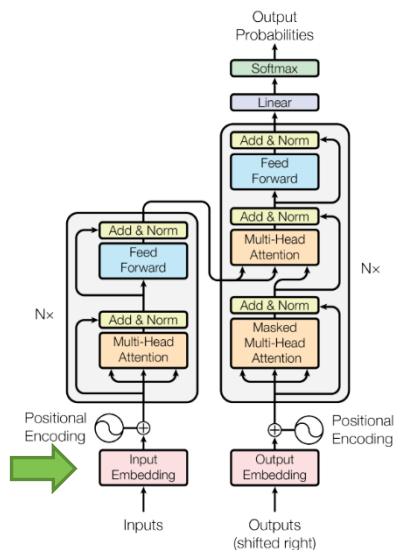


Word	Index
a	0
..	..
is	882
..	..
..	..
there	3582
..	..
where	11442
..	..

Where there is will, there

(11442 3582 882 ... 3582)

Transformer for NLP



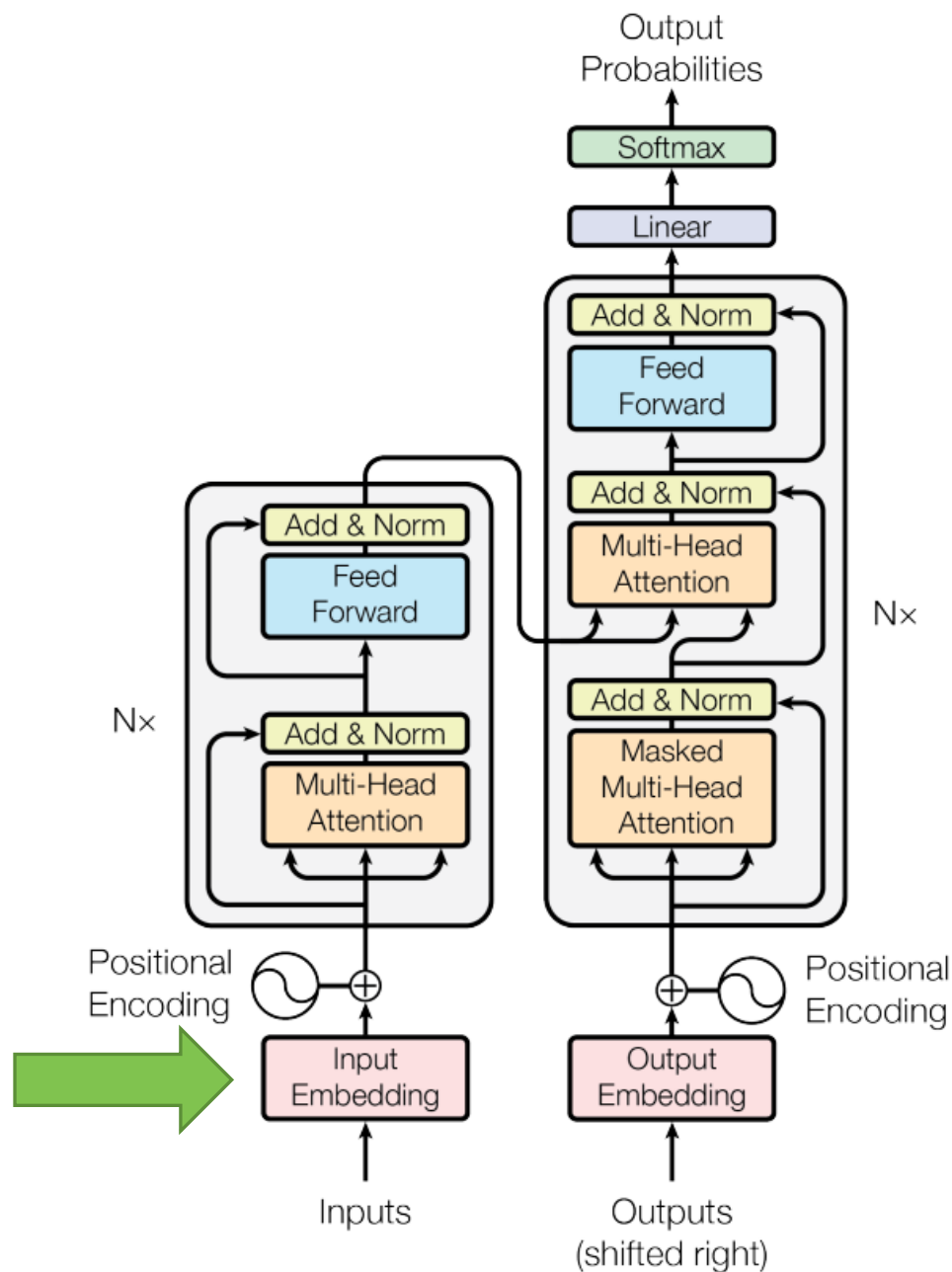
Word	Index
a	0
..	..
is	882
..	..
..	..
there	3582
..	..
where	11442
..	..

Where there is will, there

(11442 3582 882 ... 3582)

Input

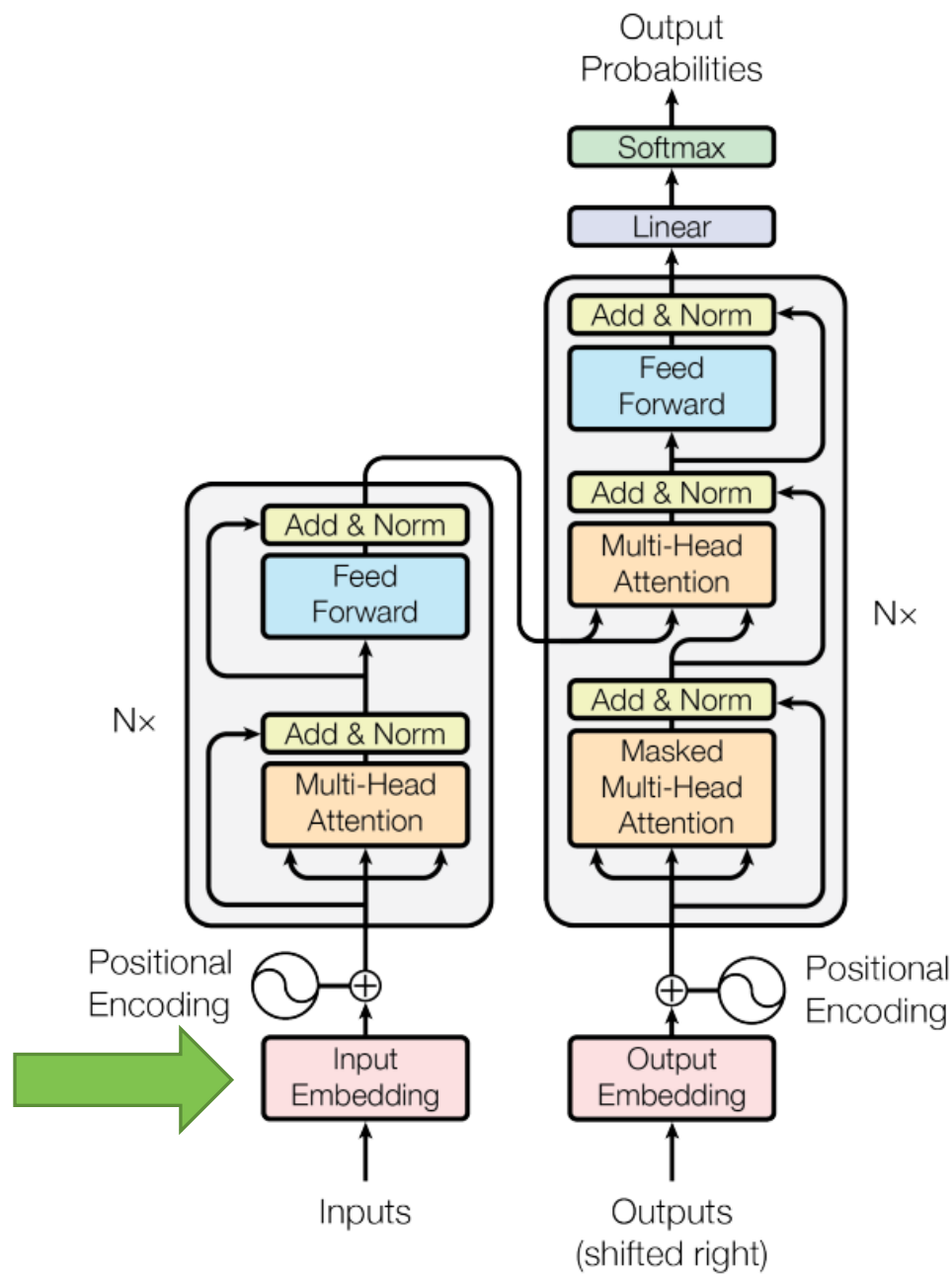
Word Embedding



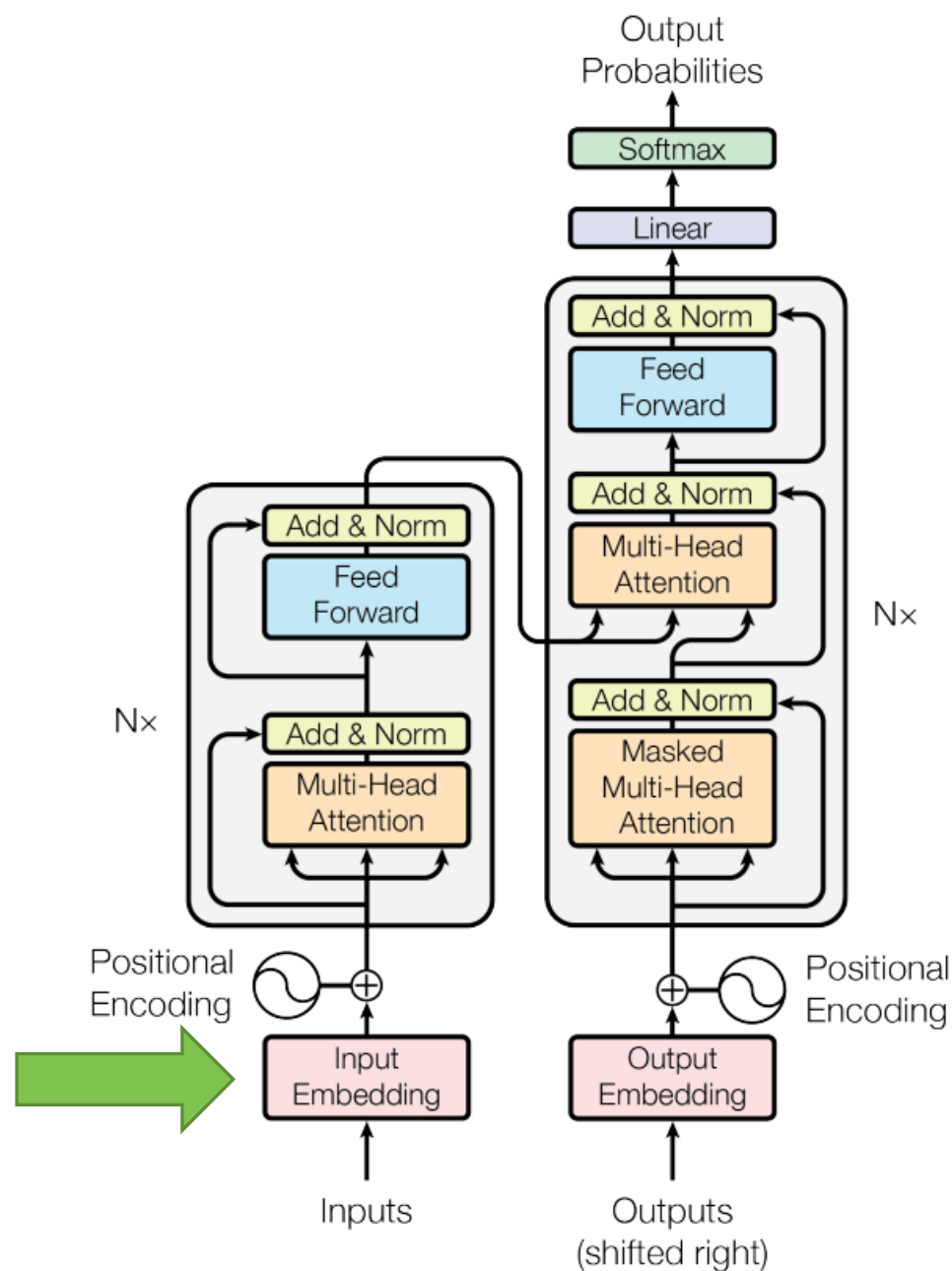
	living being	feline	human	gender	royalty	verb	plural
cat →	0.6	0.9	0.1	0.4	-0.7	-0.3	-0.2
kitten →	0.5	0.8	-0.1	0.2	-0.6	-0.5	-0.1
dog →	0.7	-0.1	0.4	0.3	-0.4	-0.1	-0.3
houses →	-0.8	-0.4	-0.5	0.1	-0.9	0.3	0.8
man →	0.6	-0.2	0.8	0.9	-0.1	-0.9	-0.7
woman →	0.7	0.3	0.9	-0.7	0.1	-0.5	-0.4
king →	0.5	-0.4	0.7	0.8	0.9	-0.7	-0.6
queen →	0.8	-0.1	0.8	-0.9	0.8	-0.5	-0.9

Word Embedding

Word embedding: Learnt during training for transformer

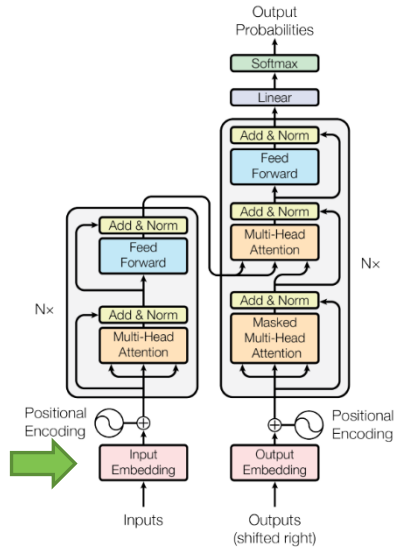


Word Embedding



	Feature 1	Feature 2	Feature 3	Feature 4	Feature 5	Feature 6	Feature 7
<i>cat</i> →	0.6	0.9	0.1	0.4	-0.7	-0.3	-0.2
<i>kitten</i> →	0.5	0.8	-0.1	0.2	-0.6	-0.5	-0.1
<i>dog</i> →	0.7	-0.1	0.4	0.3	-0.4	-0.1	-0.3
<i>houses</i> →	-0.8	-0.4	-0.5	0.1	-0.9	0.3	0.8
<i>man</i> →	0.6	-0.2	0.8	0.9	-0.1	-0.9	-0.7
<i>woman</i> →	0.7	0.3	0.9	-0.7	0.1	-0.5	-0.4
<i>king</i> →	0.5	-0.4	0.7	0.8	0.9	-0.7	-0.6
<i>queen</i> →	0.8	-0.1	0.8	-0.9	0.8	-0.5	-0.9

After Word Embedding



Word	Index
a	0
..	..
is	882
..	..
..	..
there	3582
..	..
where	11442
..	..

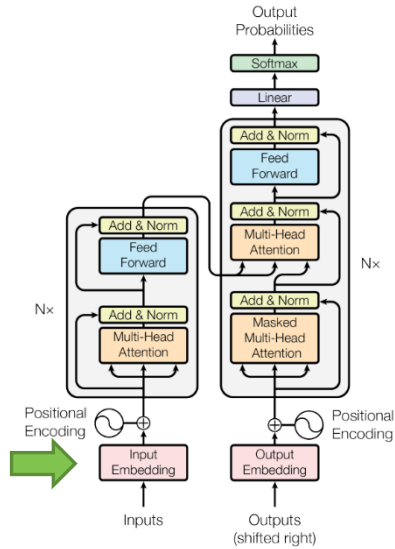
	Feature 1	Feature 2	Feature 3	Feature 4	Feature 5	Feature 6	Feature 7
<i>cat</i> →	0.6	0.9	0.1	0.4	-0.7	-0.3	-0.2
<i>kitten</i> →	0.5	0.8	-0.1	0.2	-0.6	-0.5	-0.1
<i>dog</i> →	0.7	-0.1	0.4	0.3	-0.4	-0.1	-0.3
<i>houses</i> →	-0.8	-0.4	-0.5	0.1	-0.9	0.3	0.8
<i>man</i> →	0.6	-0.2	0.8	0.9	-0.1	-0.9	-0.7
<i>woman</i> →	0.7	0.3	0.9	-0.7	0.1	-0.5	-0.4
<i>king</i> →	0.5	-0.4	0.7	0.8	0.9	-0.7	-0.6
<i>queen</i> →	0.8	-0.1	0.8	-0.9	0.8	-0.5	-0.9

Where there is will, there

(11442 3582 882 ... 3582)

0.30	0.17	...	0.20
0.22	0.29		0.52
0.17	0.10		0.18
0.65	0.32		0.69
0.12	0.11		0.72

After Word Embedding



Word	Index
a	0
..	..
is	882
..	..
..	..
there	3582
..	..
where	11442
..	..

	Feature 1	Feature 2	Feature 3	Feature 4	Feature 5	Feature 6	Feature 7
<i>cat</i> →	0.6	0.9	0.1	0.4	-0.7	-0.3	-0.2
<i>kitten</i> →	0.5	0.8	-0.1	0.2	-0.6	-0.5	-0.1
<i>dog</i> →	0.7	-0.1	0.4	0.3	-0.4	-0.1	-0.3
<i>houses</i> →	-0.8	-0.4	-0.5	0.1	-0.9	0.3	0.8
<i>man</i> →	0.6	-0.2	0.8	0.9	-0.1	-0.9	-0.7
<i>woman</i> →	0.7	0.3	0.9	-0.7	0.1	-0.5	-0.4
<i>king</i> →	0.5	-0.4	0.7	0.8	0.9	-0.7	-0.6
<i>queen</i> →	0.8	-0.1	0.8	-0.9	0.8	-0.5	-0.9

Where there is will, there

(11442 3582 882 ... 3582)

0.30	0.17	0.20
0.22	0.29	0.52
0.17	0.10	0.18
0.65	0.32	0.69
0.12	0.11	0.72

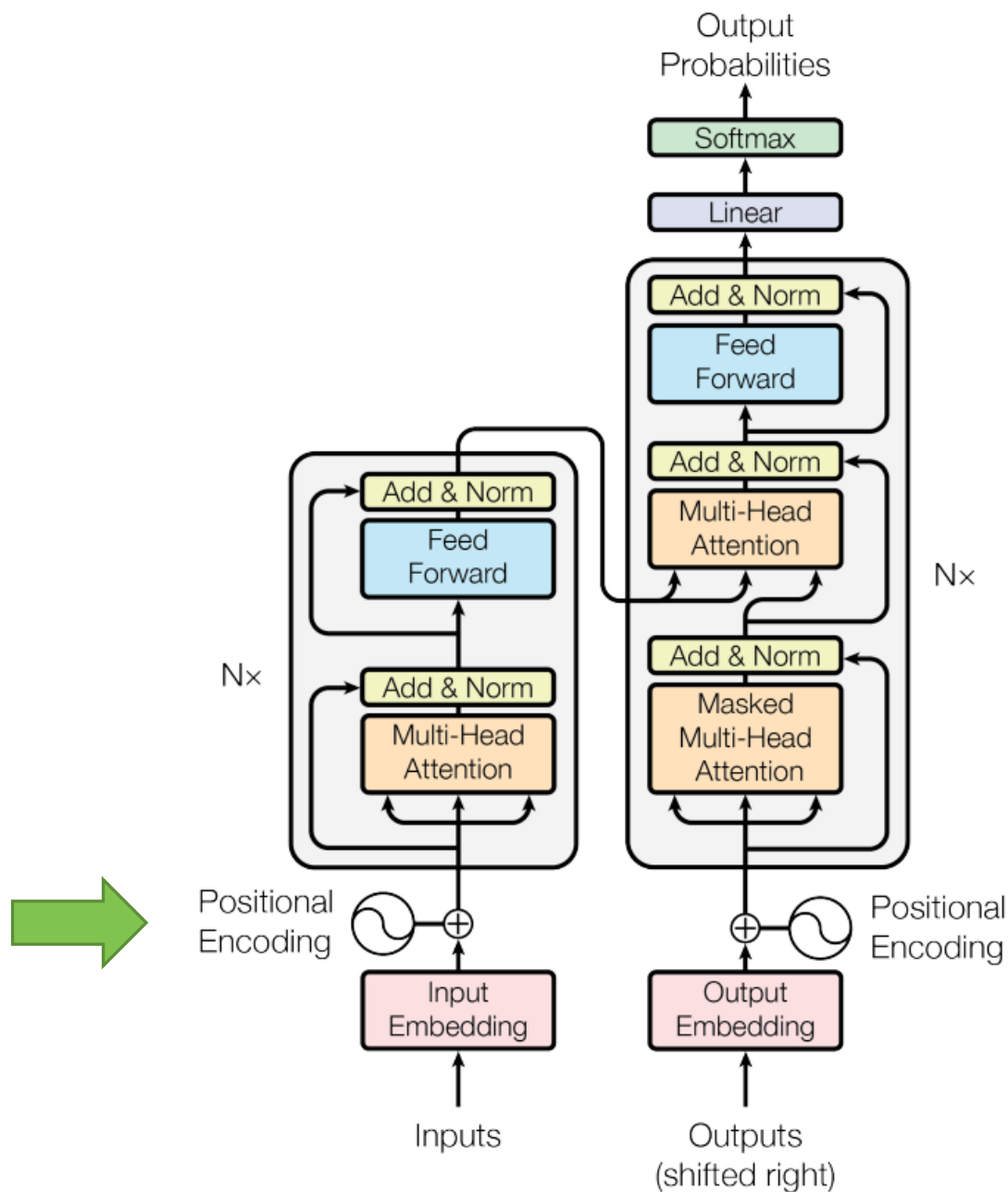
e_0

e_1

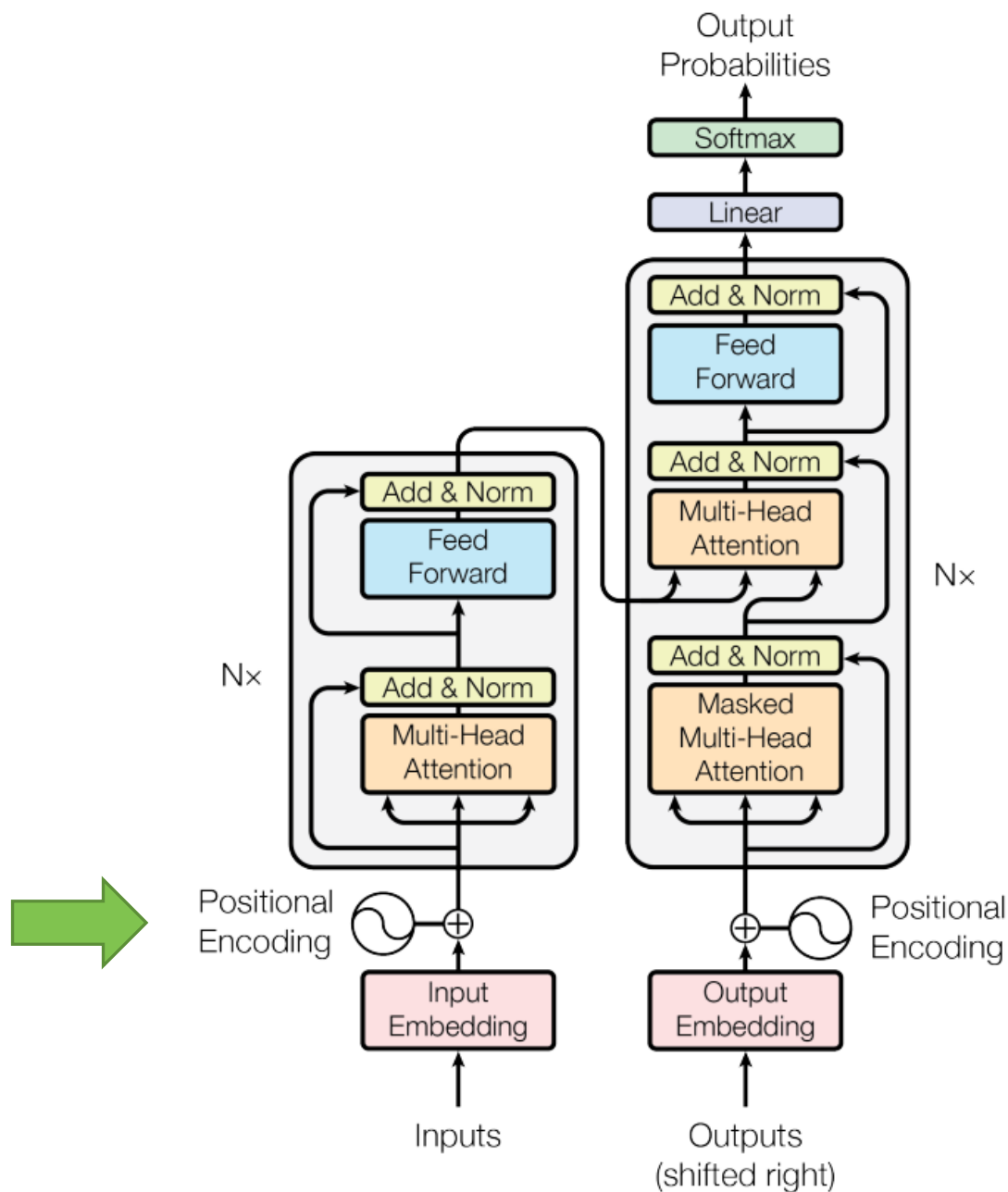
e_n

Positional Embedding

Although the sky is cloudy, it may **not** rain



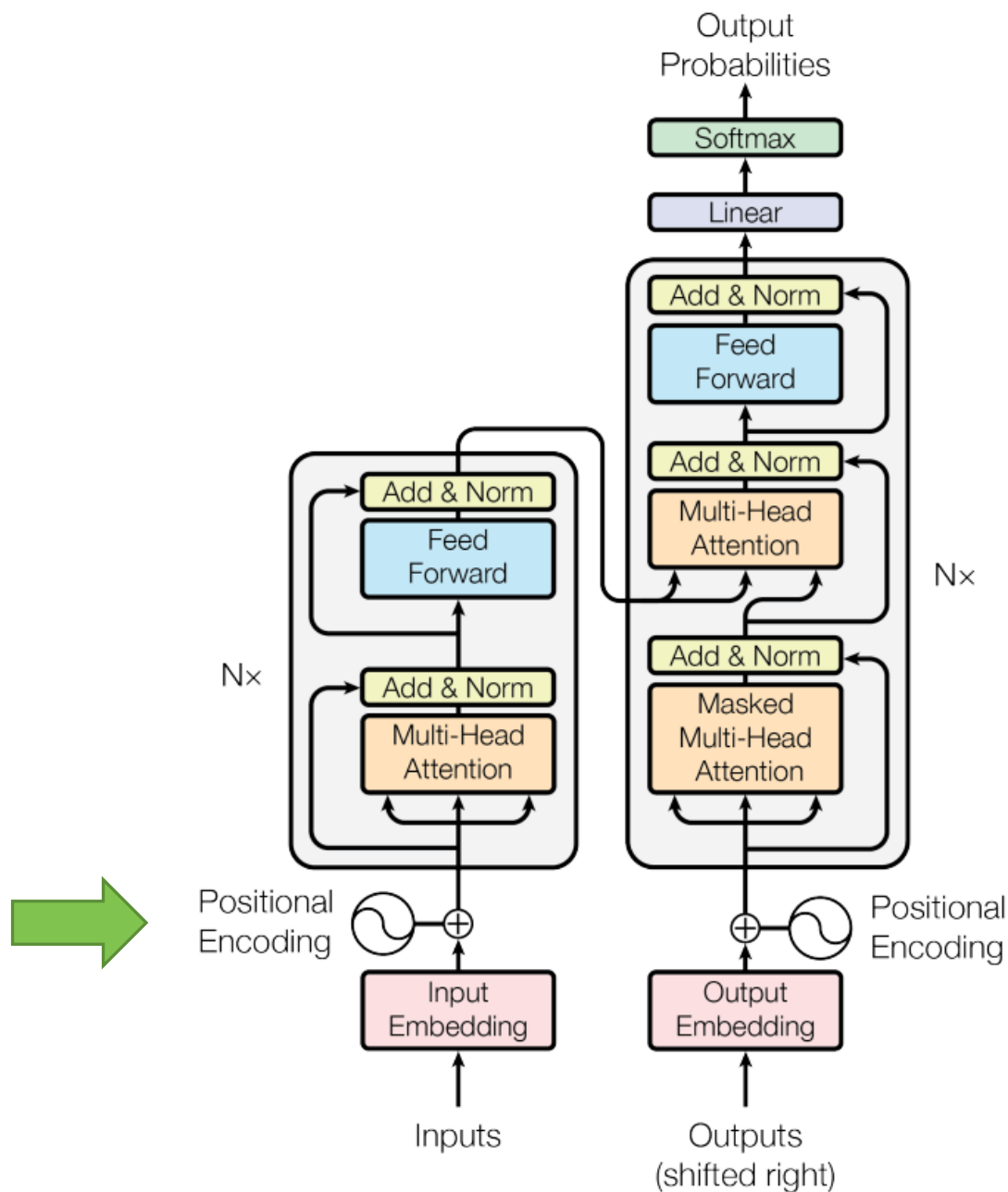
Positional Embedding



Although the sky is cloudy, it may **not** rain

Although the sky is **not** cloudy, it may rain

Positional Embedding

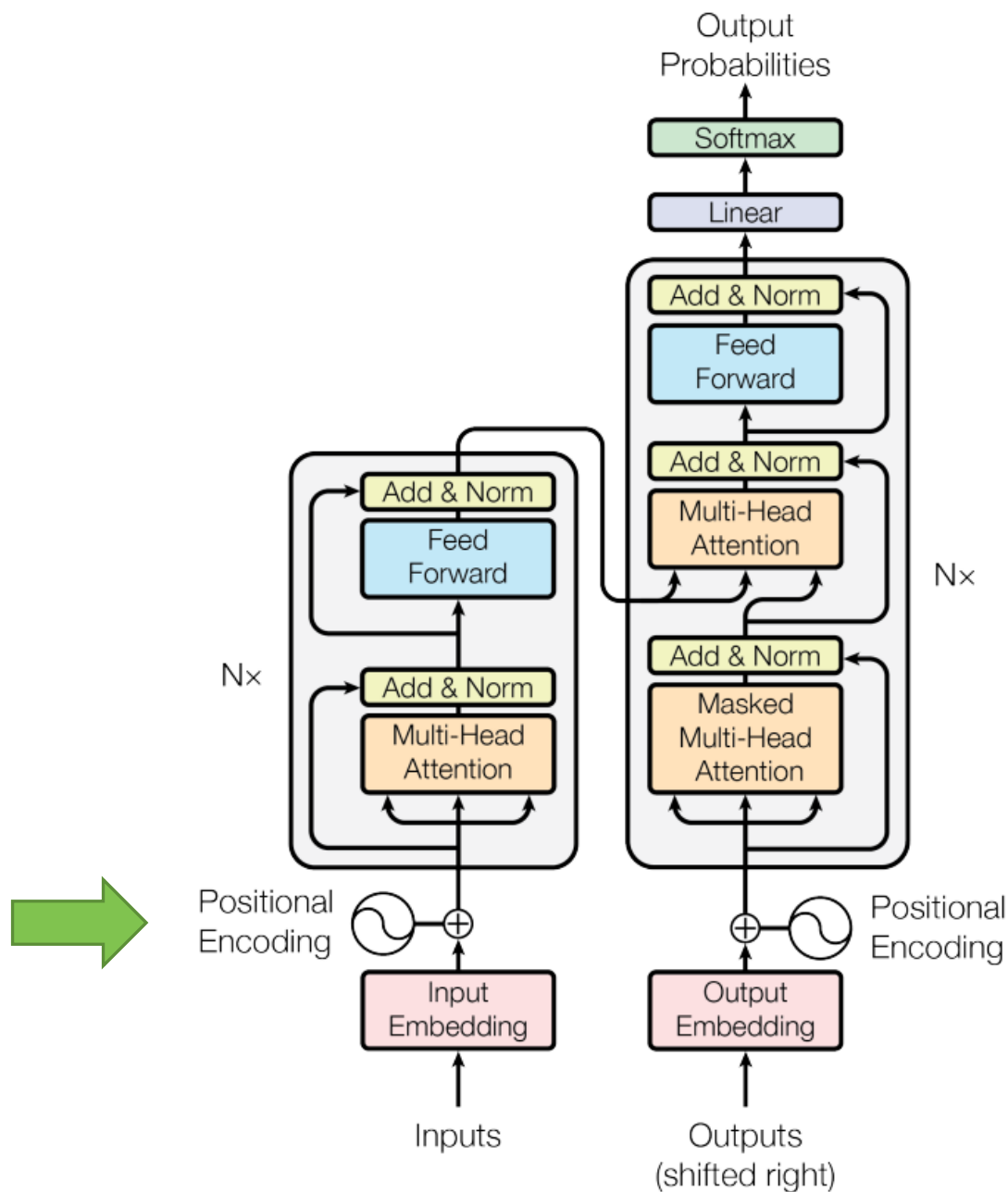


Although the sky is cloudy, it may **not** rain

Although the sky is **not** cloudy, it may rain

The position of the word 'not' changes the meaning of the sentence

Positional Embedding



Although the sky is cloudy, it may **not** rain

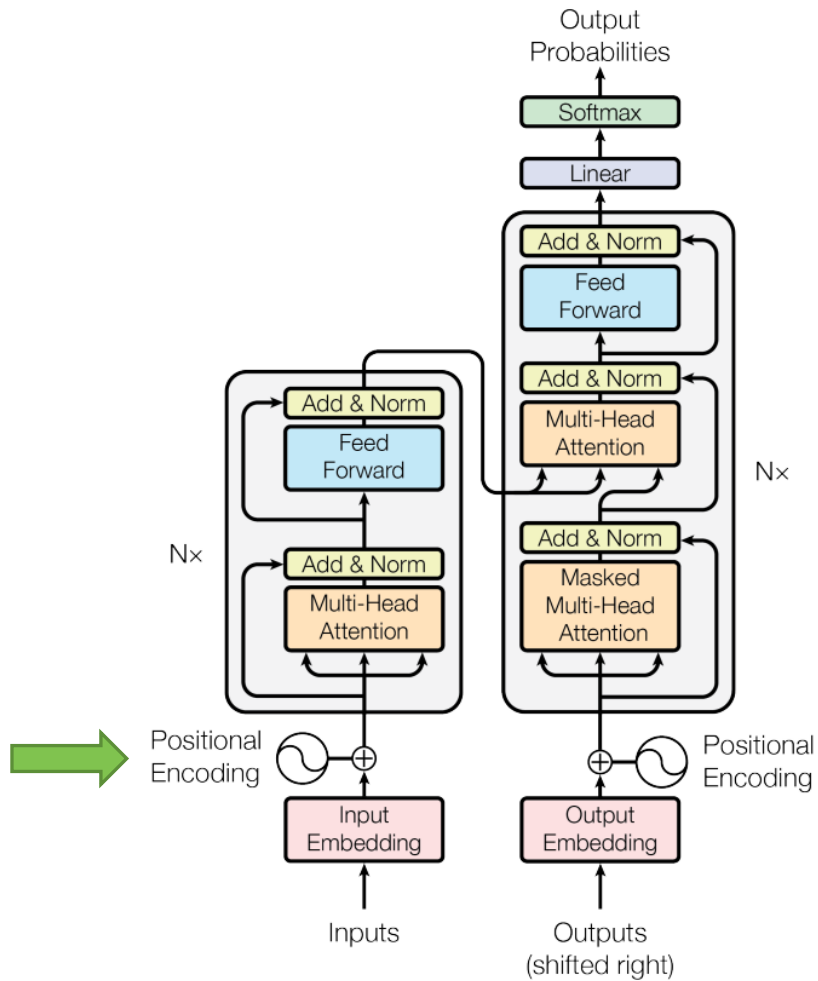
Although the sky is **not** cloudy, it may rain

The position of the word 'not' changes the meaning of the sentence

So, position information should be incorporated with word embedding

Positional Embedding

Where there is will, there



(11442

0.30
0.22
0.17
0.65
0.12

e_0

3582

0.17
0.29
0.10
0.32
0.11

e_1

882 ...

...

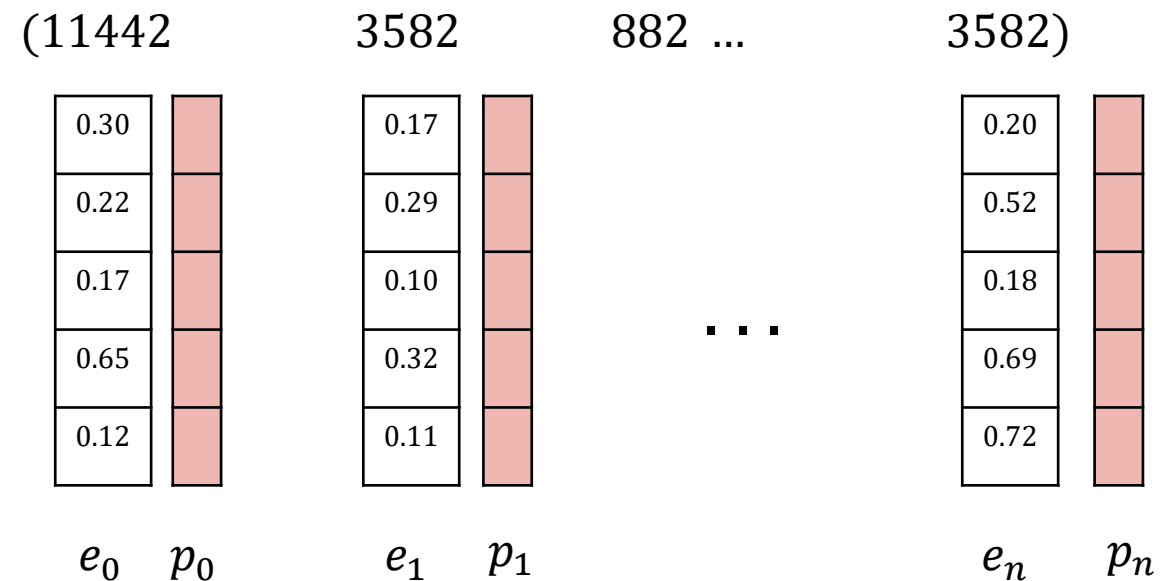
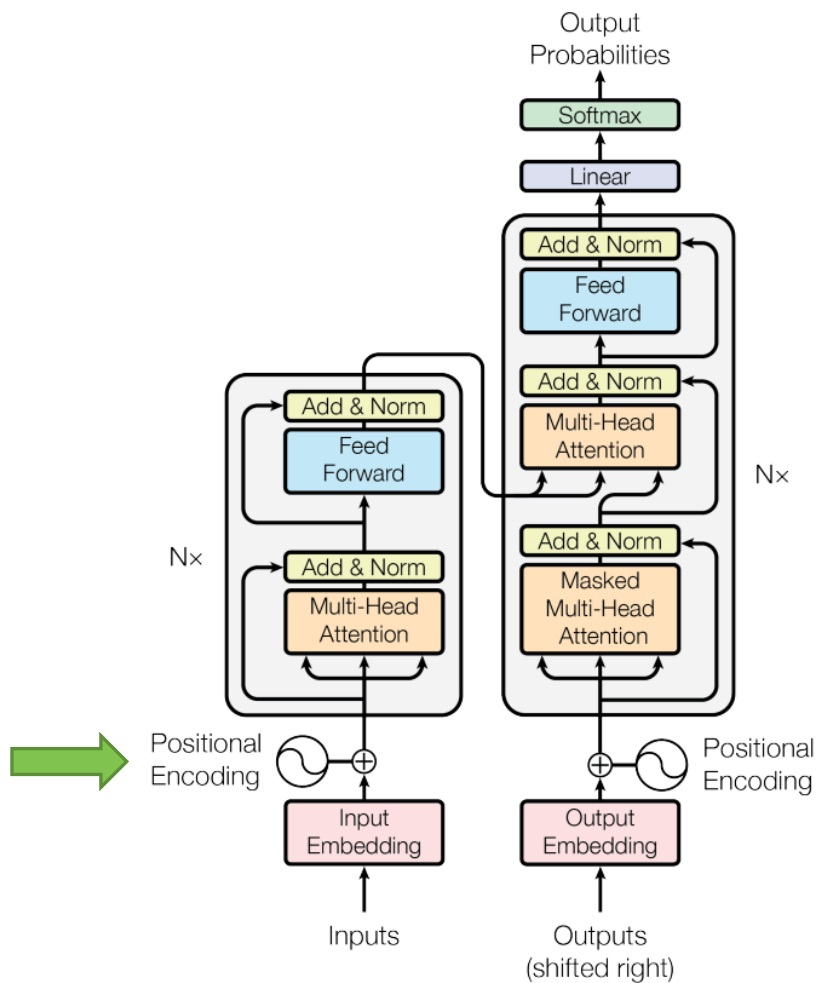
3582)

0.20
0.52
0.18
0.69
0.72

e_n

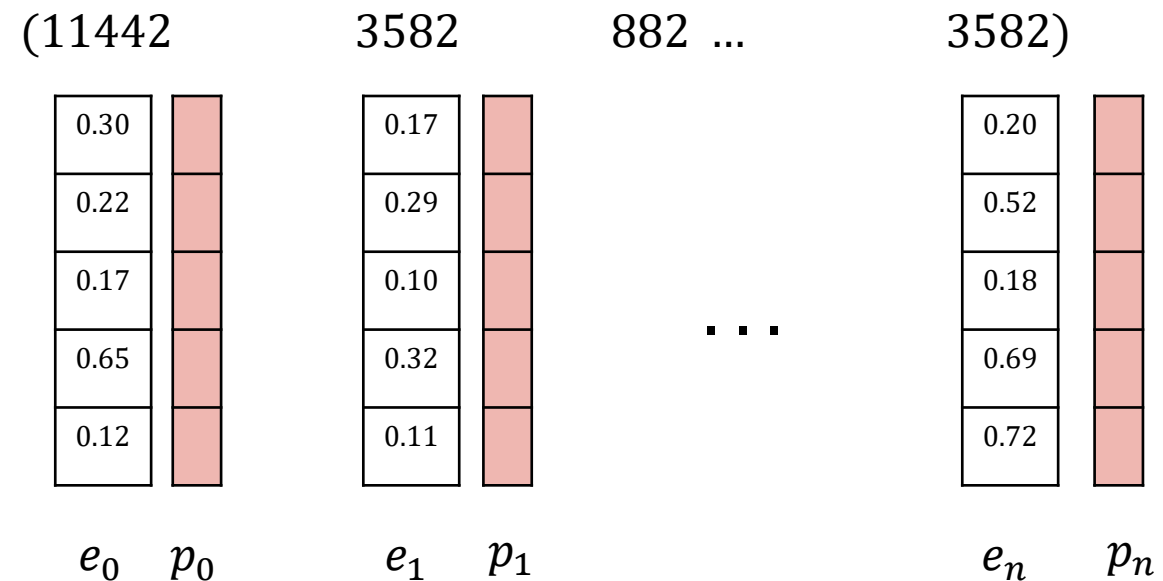
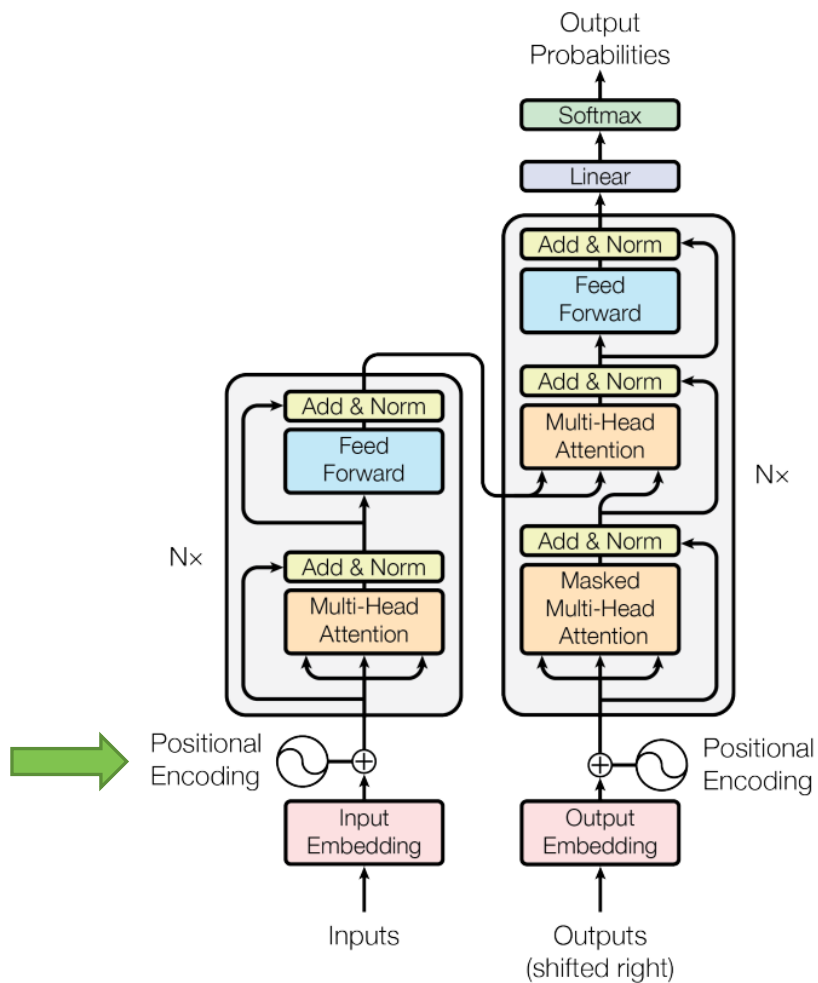
Positional Embedding

Where there is will, there



Positional Embedding

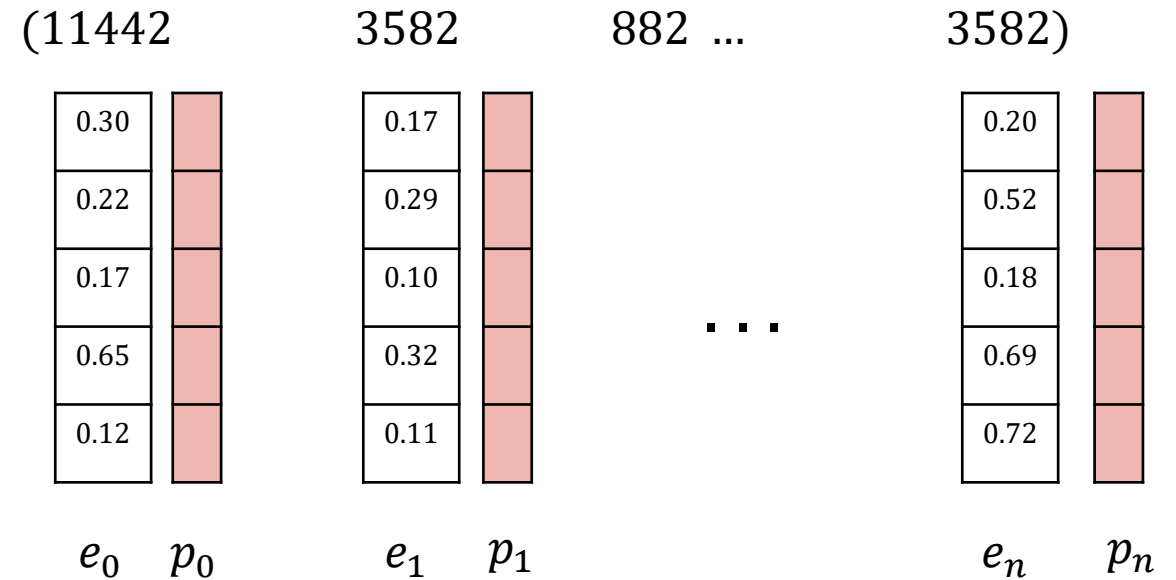
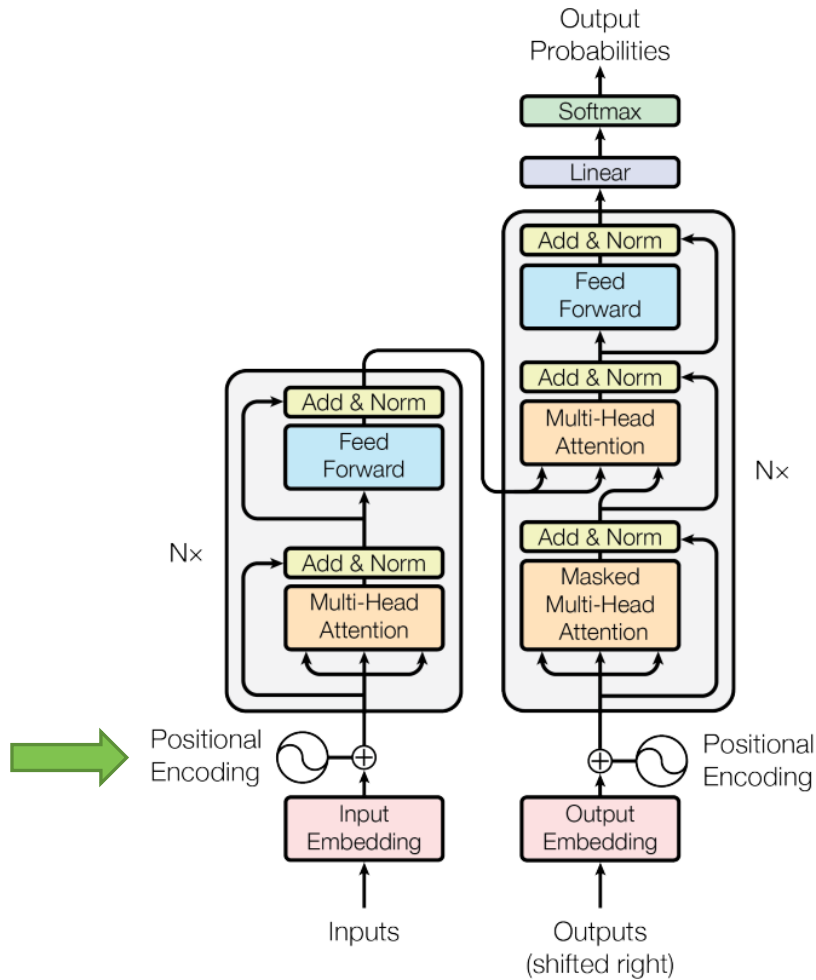
Where there is will, there



How to get p_i ?

Positional Embedding

Where there is will, there



How to get p_i ?

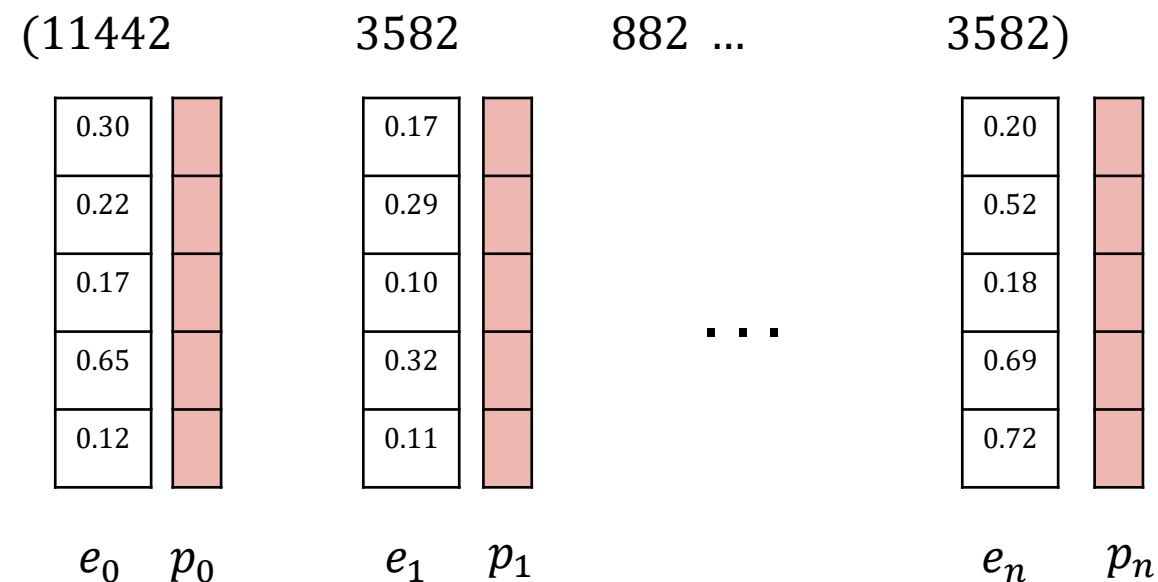
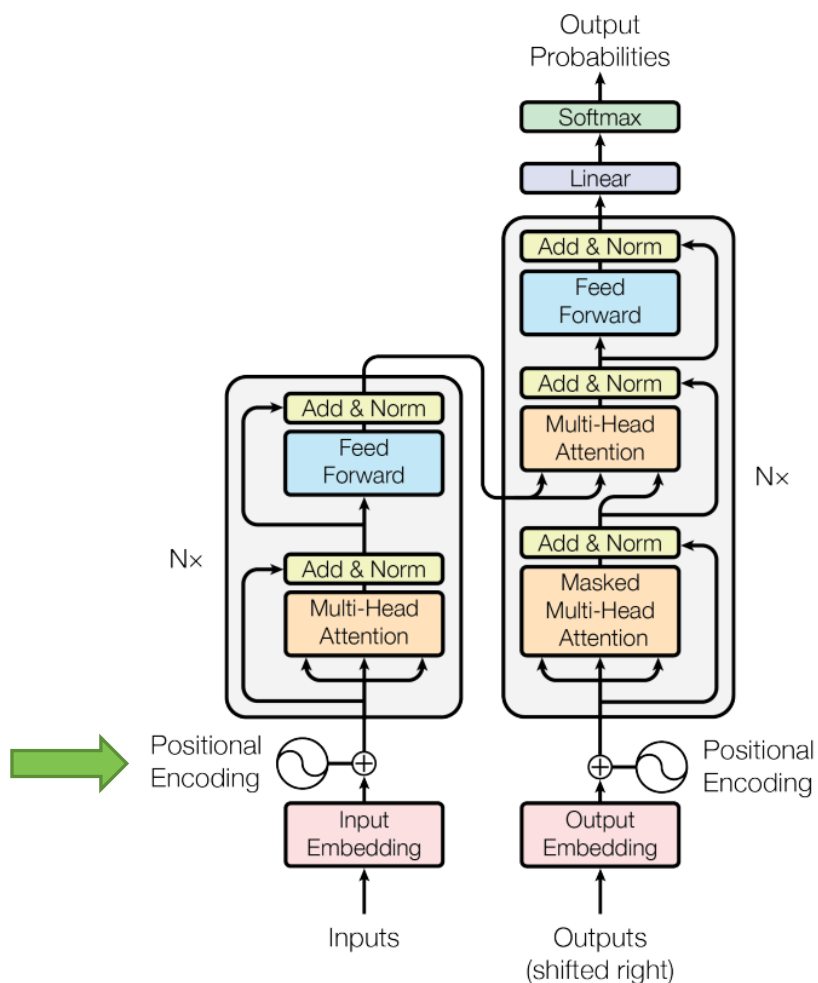
Vaswani *et al.*

$$PE(pos, 2i) = \sin\left(\frac{pos}{1000^{\frac{2i}{d}}}\right) \text{ Dimensions at even positions}$$

$$PE(pos, 2i + 1) = \cos\left(\frac{pos}{1000^{\frac{2i}{d}}}\right) \text{ Dimensions at odd positions}$$

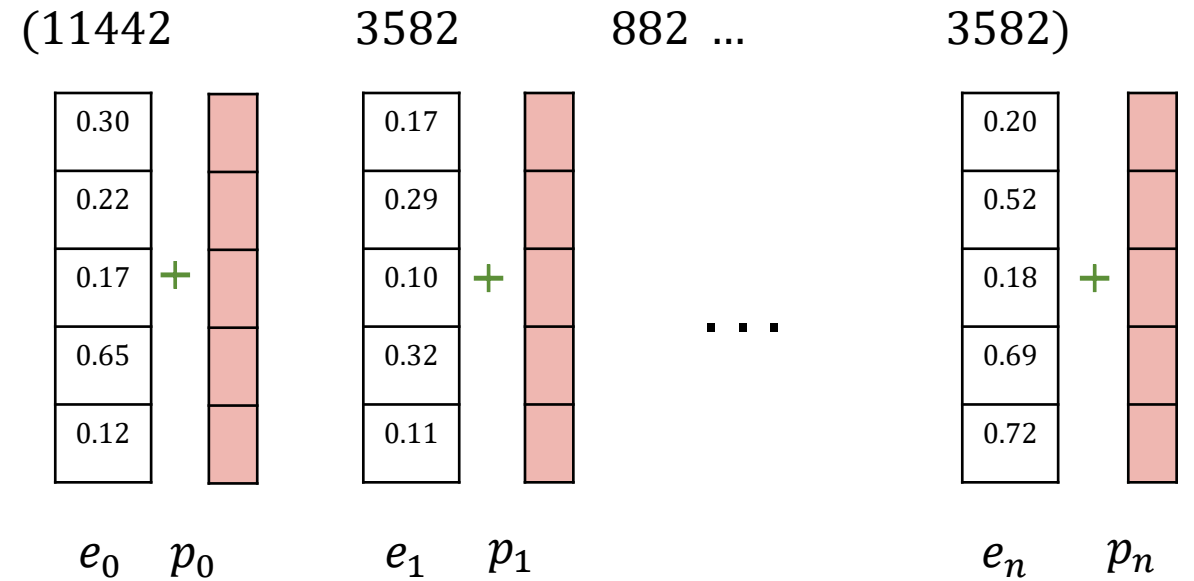
Positional Embedding

Where there is will, there



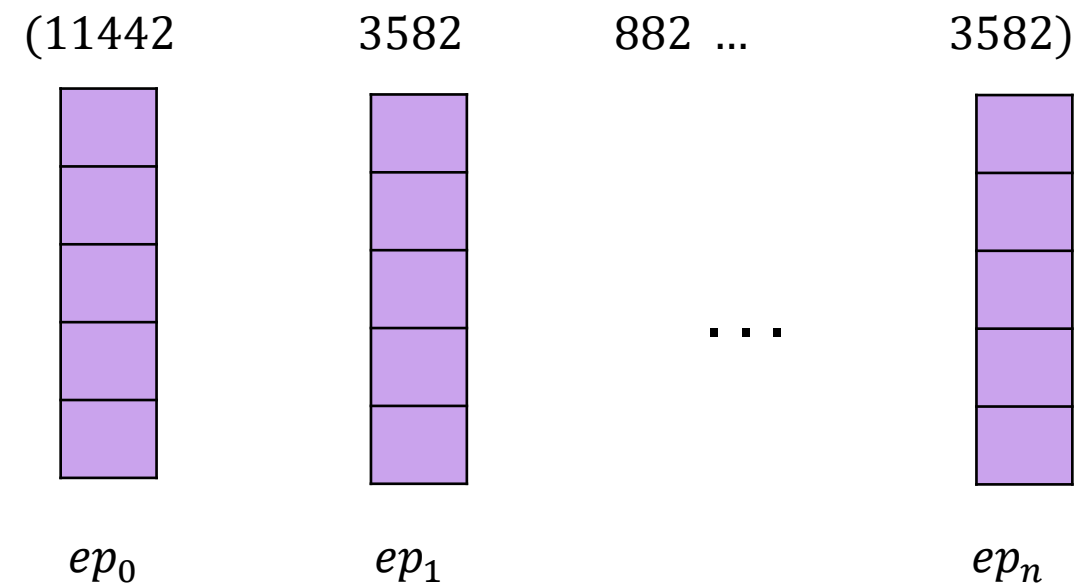
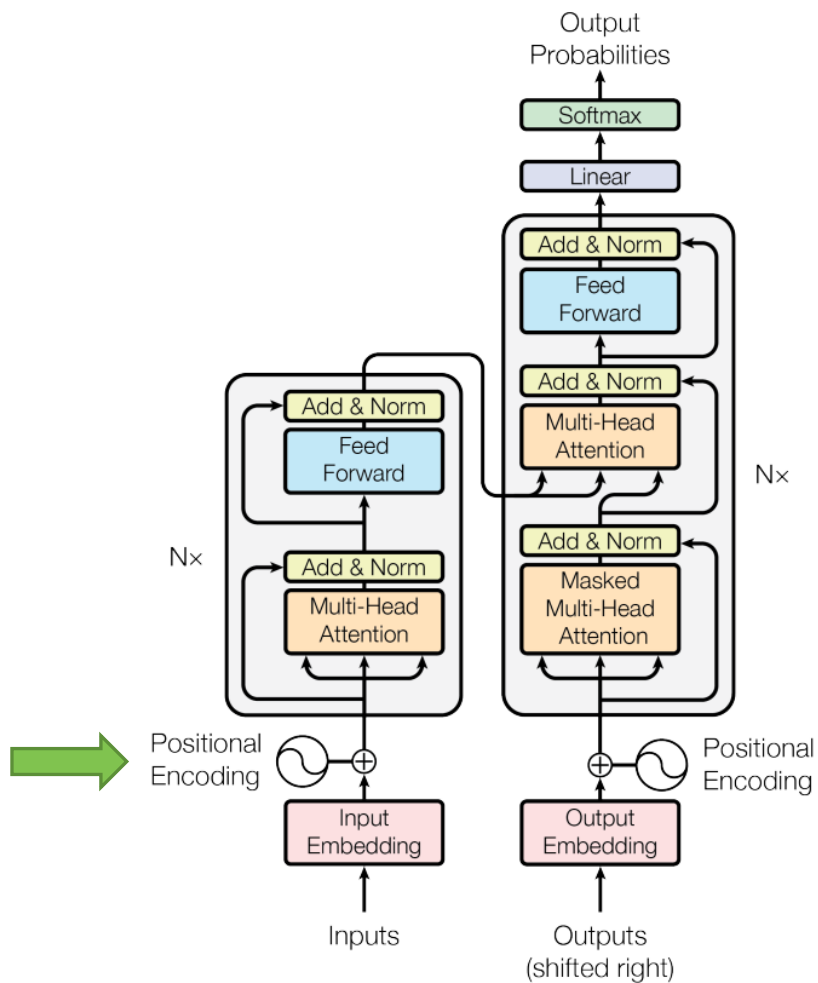
How to combine e_i and p_i ?

Where there is will, there



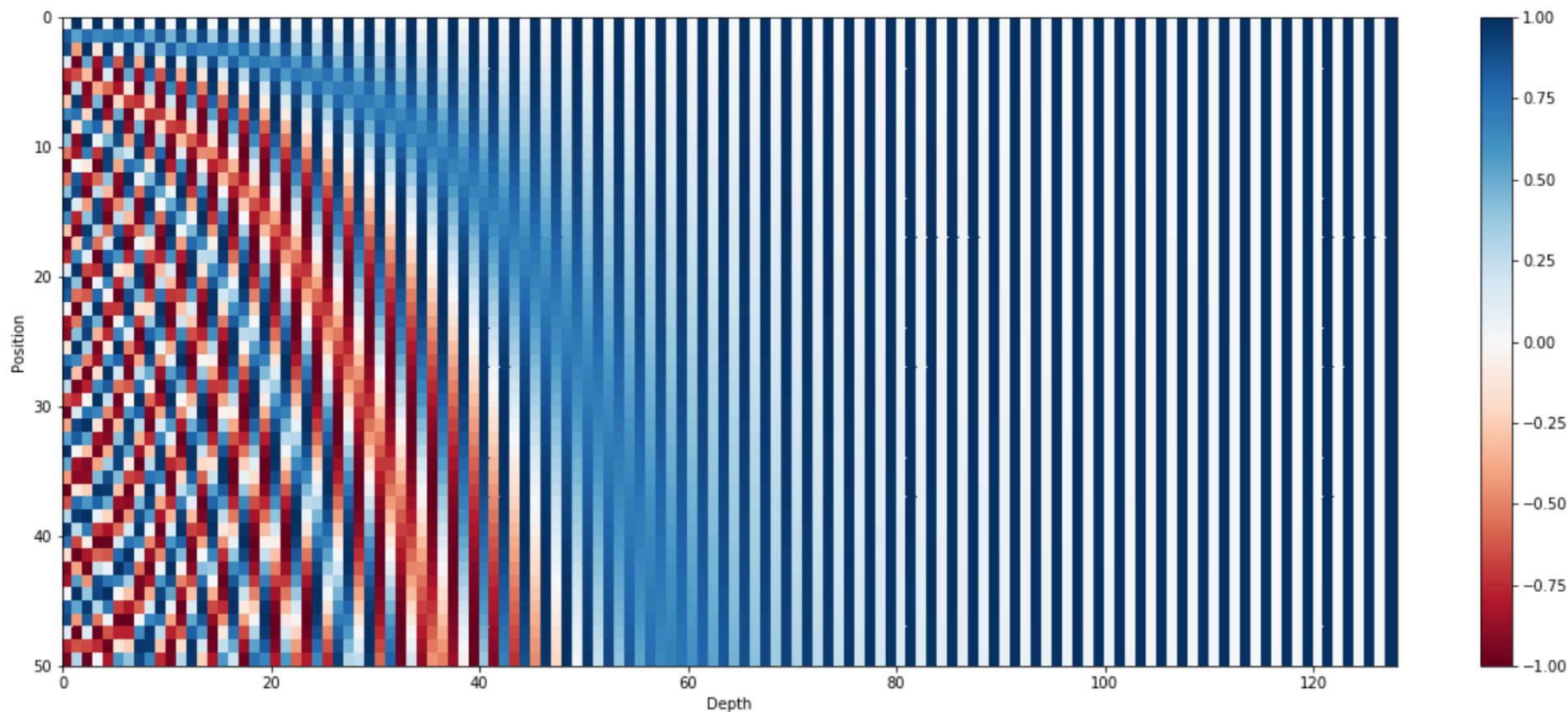
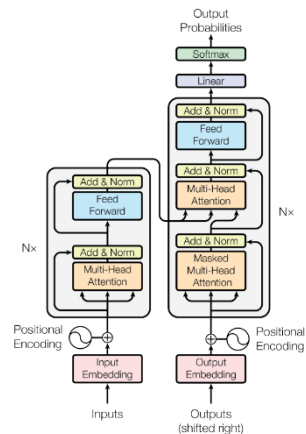
Positional Embedding

Where there is will, there



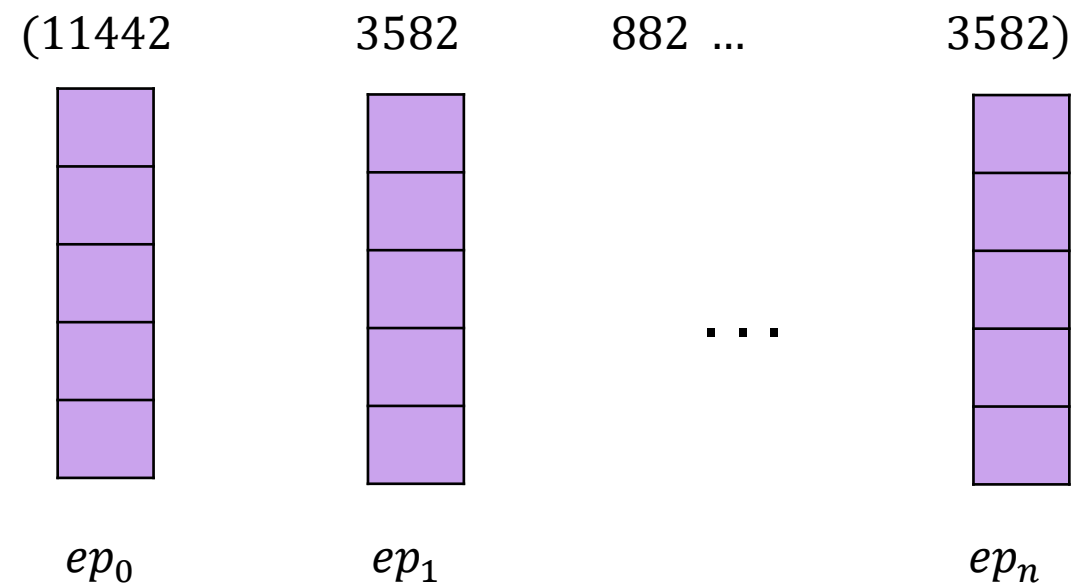
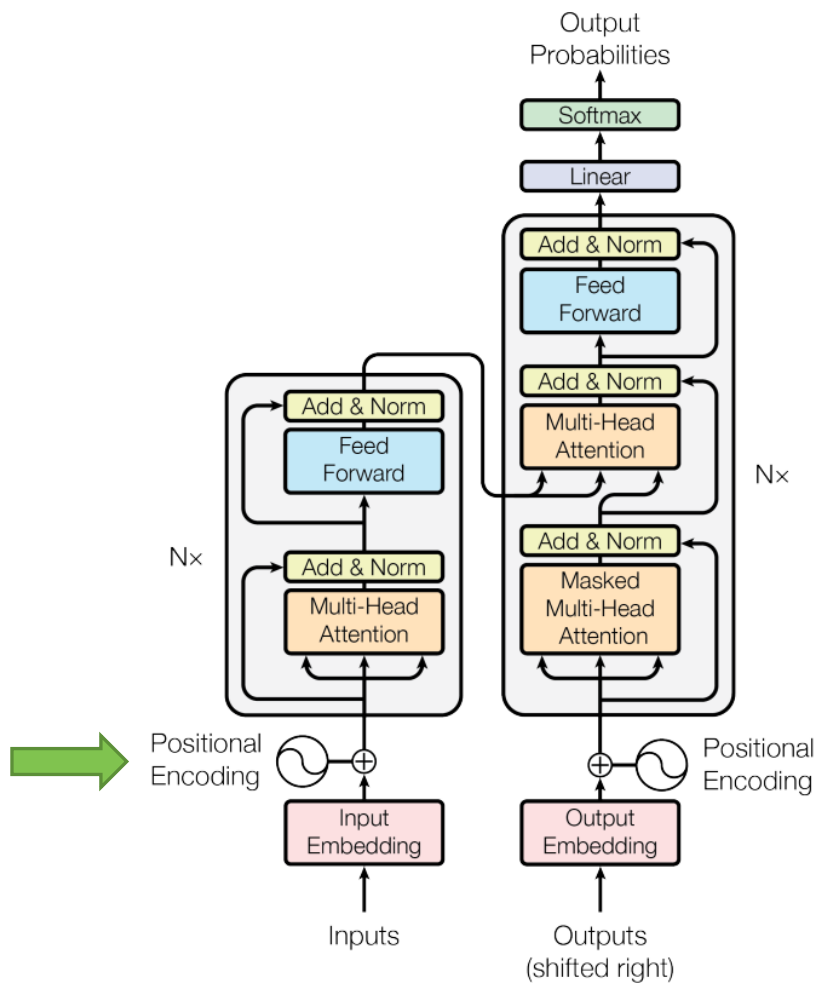
Positional Embedding

Where there is will, there

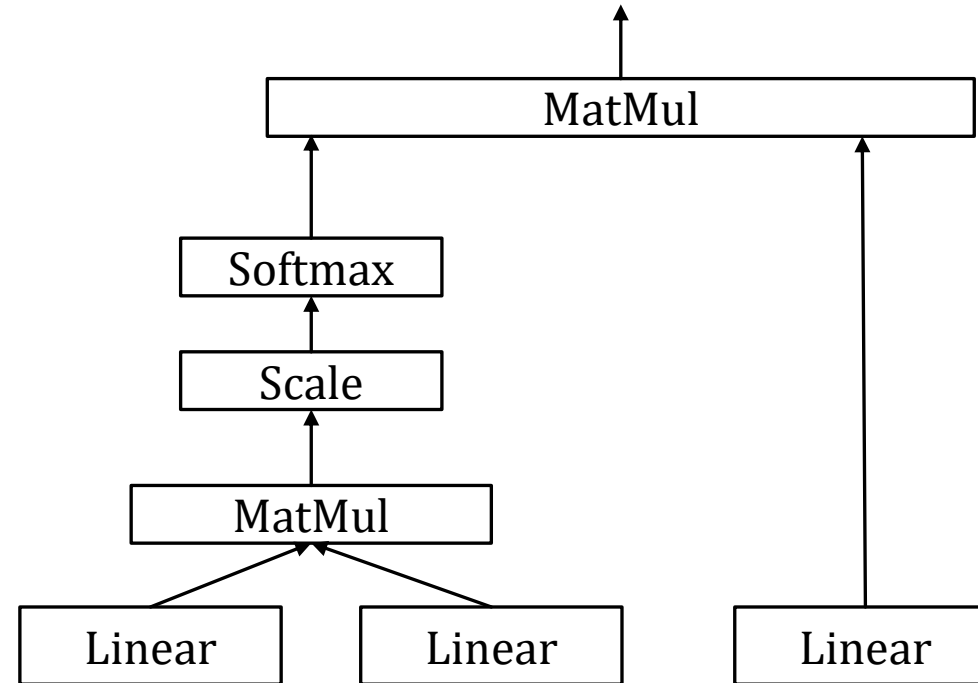
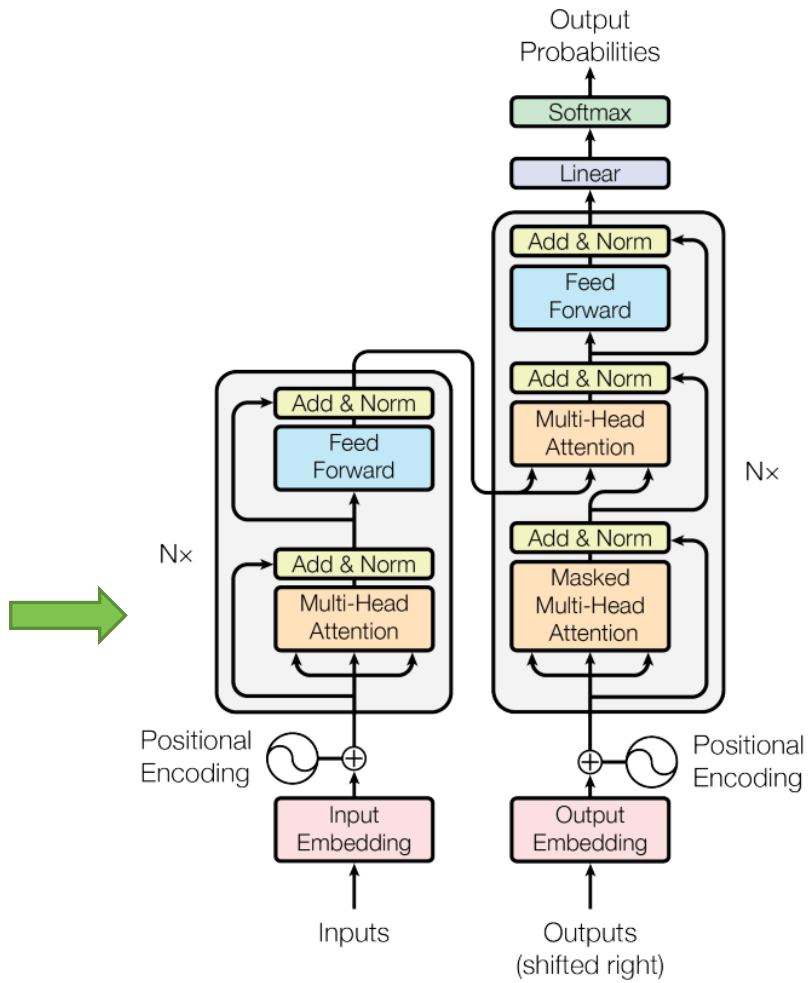


Positional Embedding

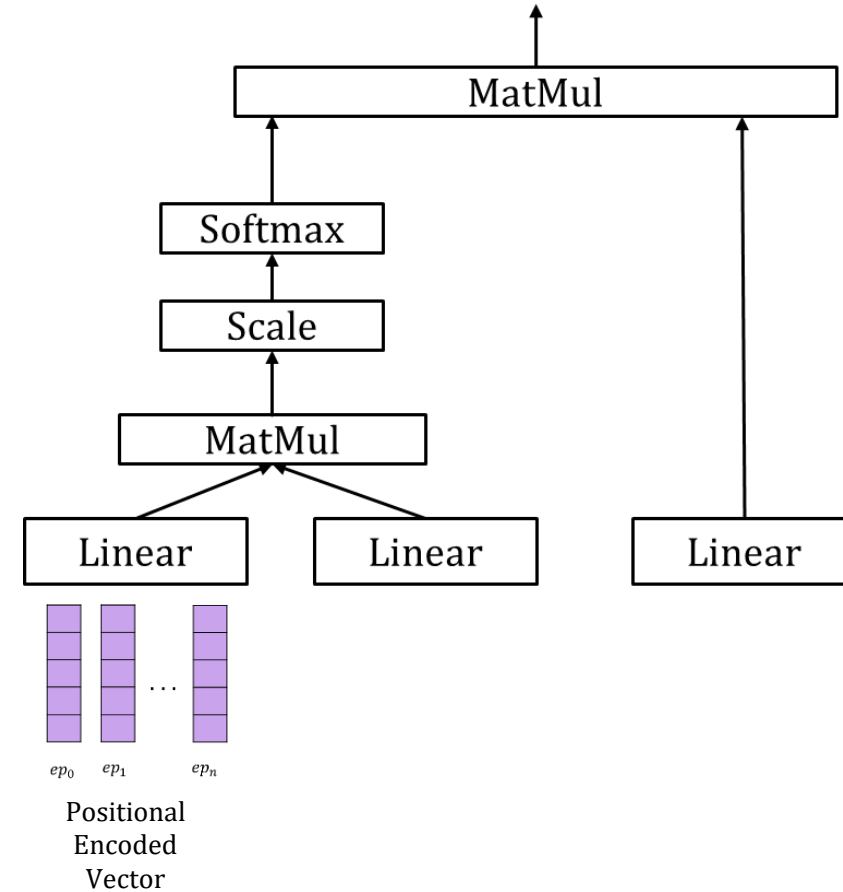
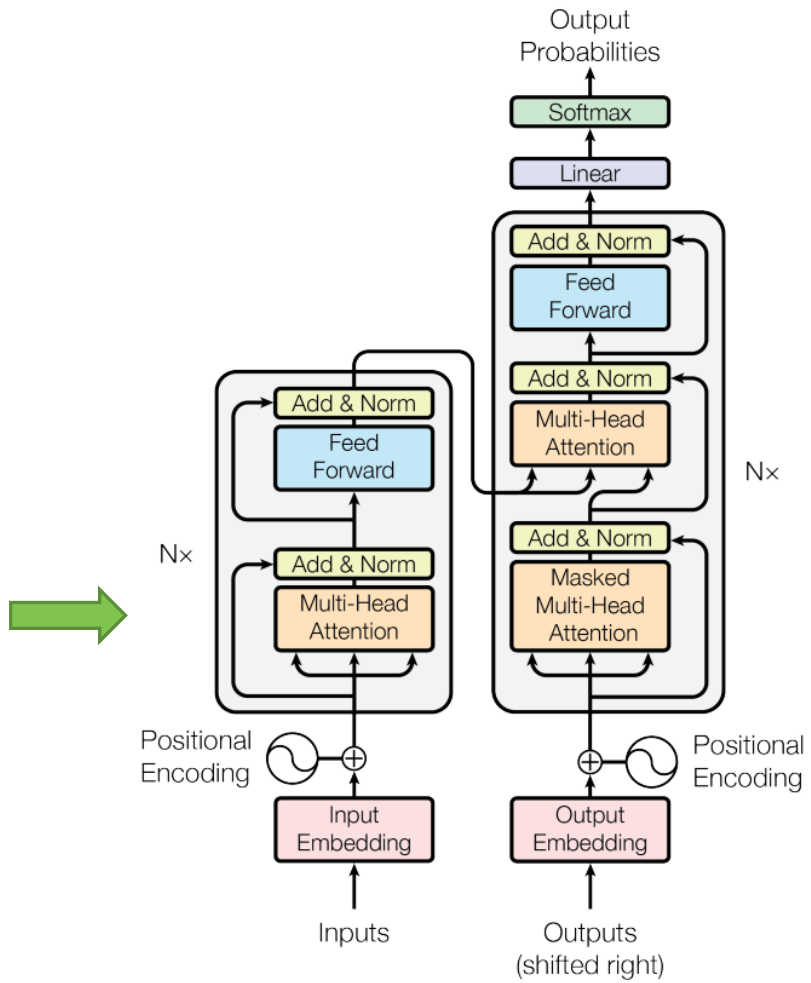
Where there is will, there



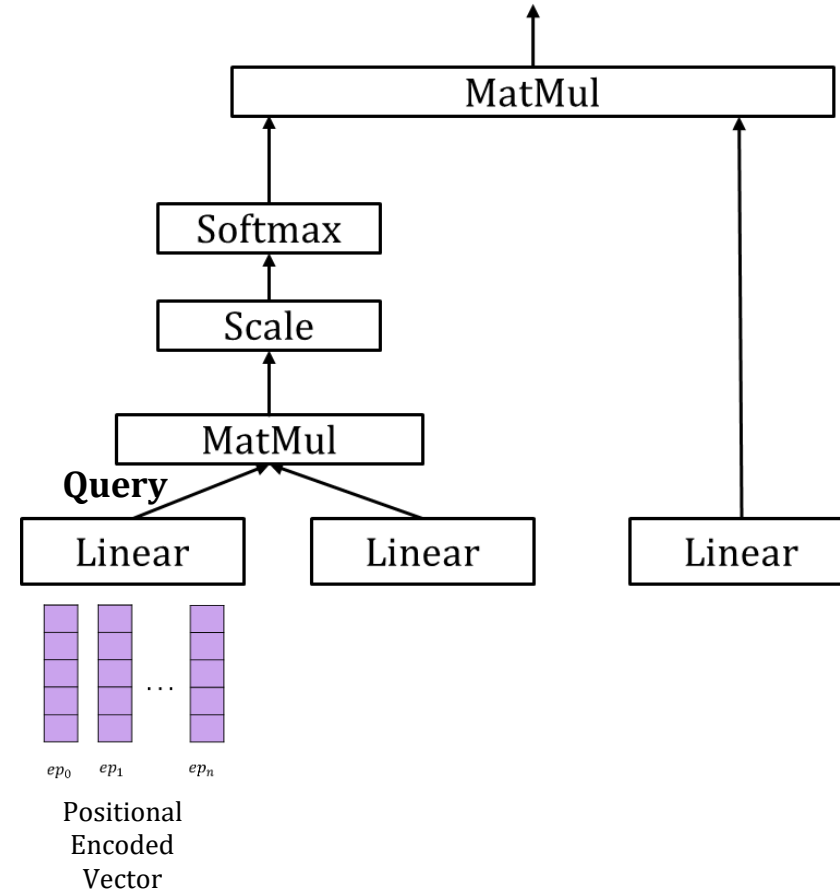
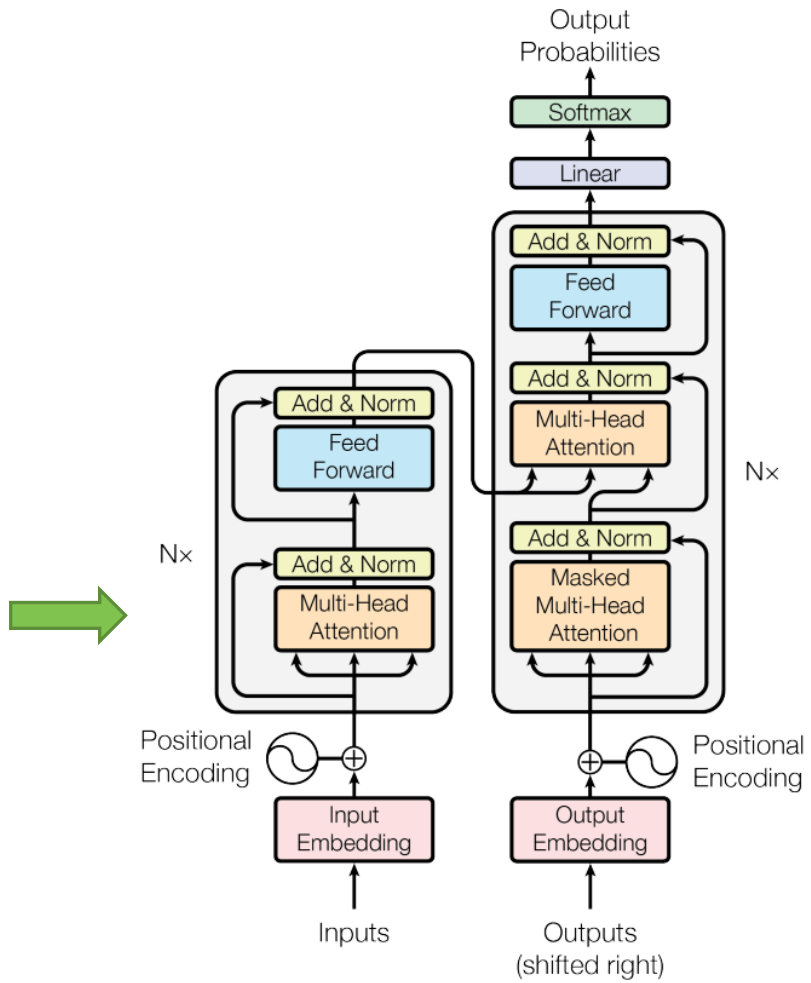
Multi-head Attention



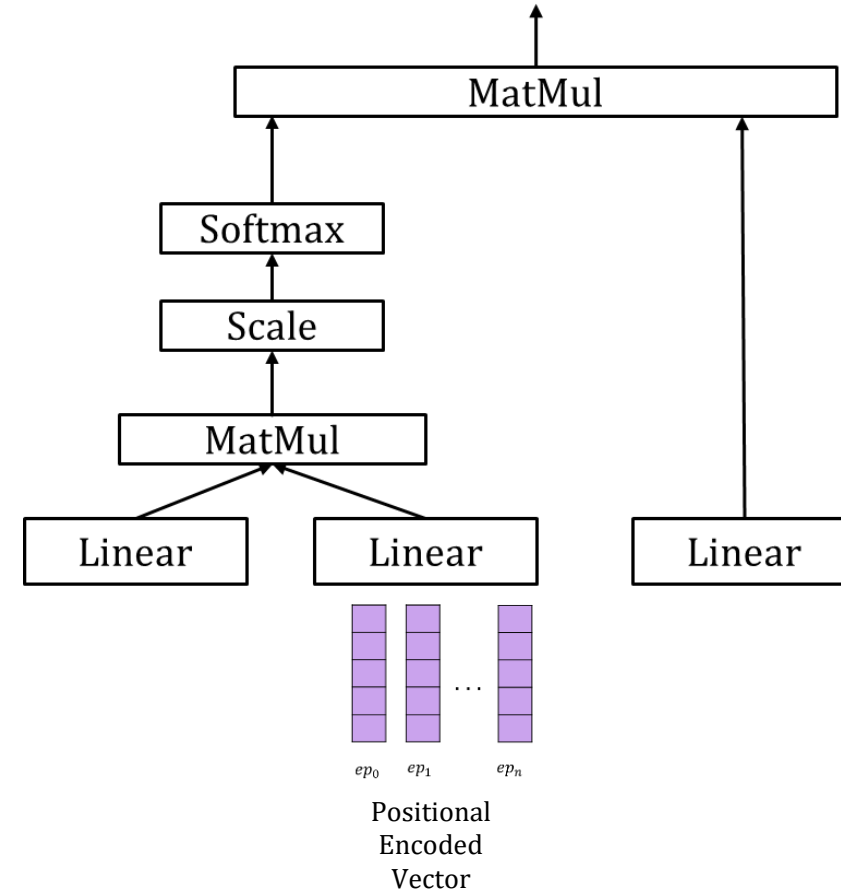
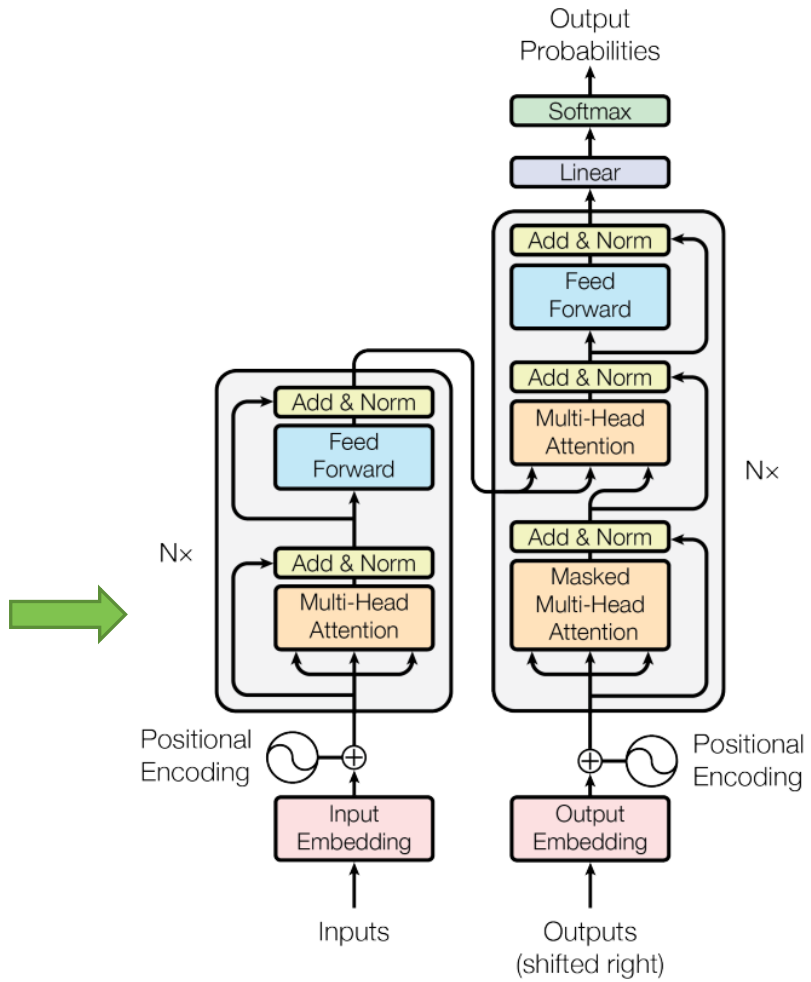
Multi-head Attention



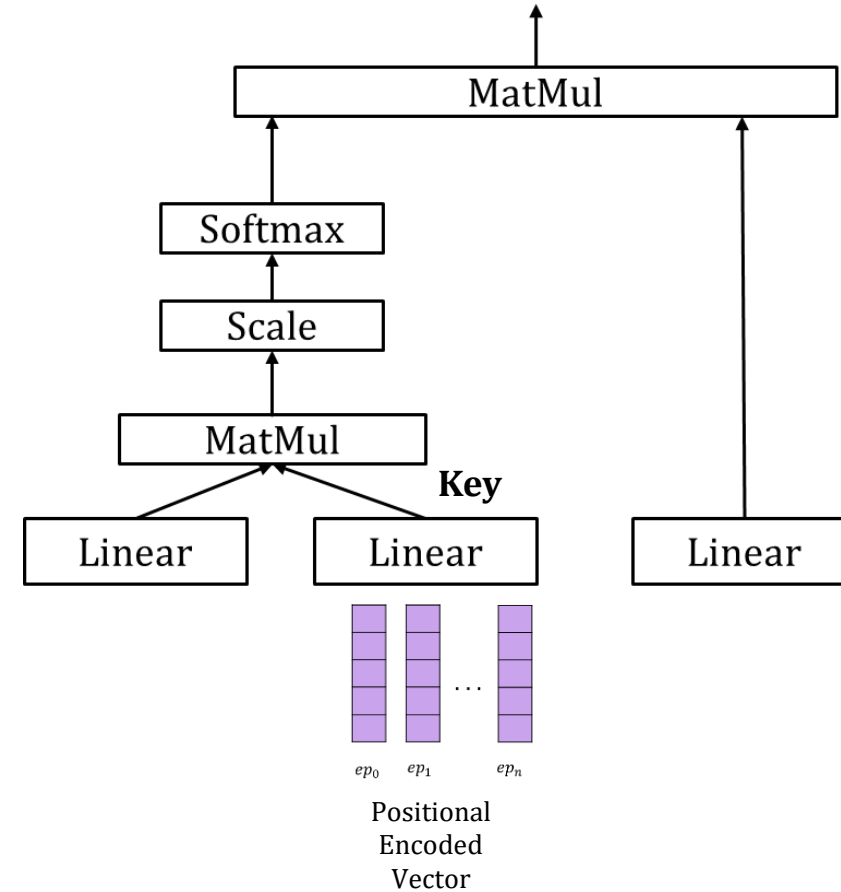
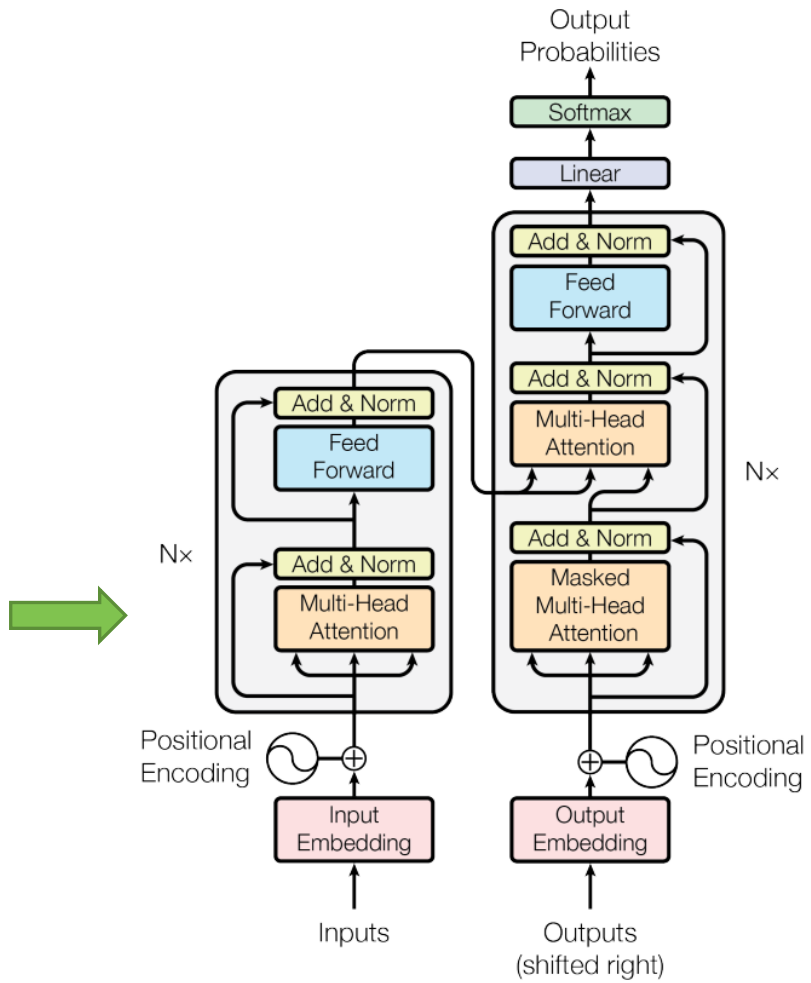
Multi-head Attention



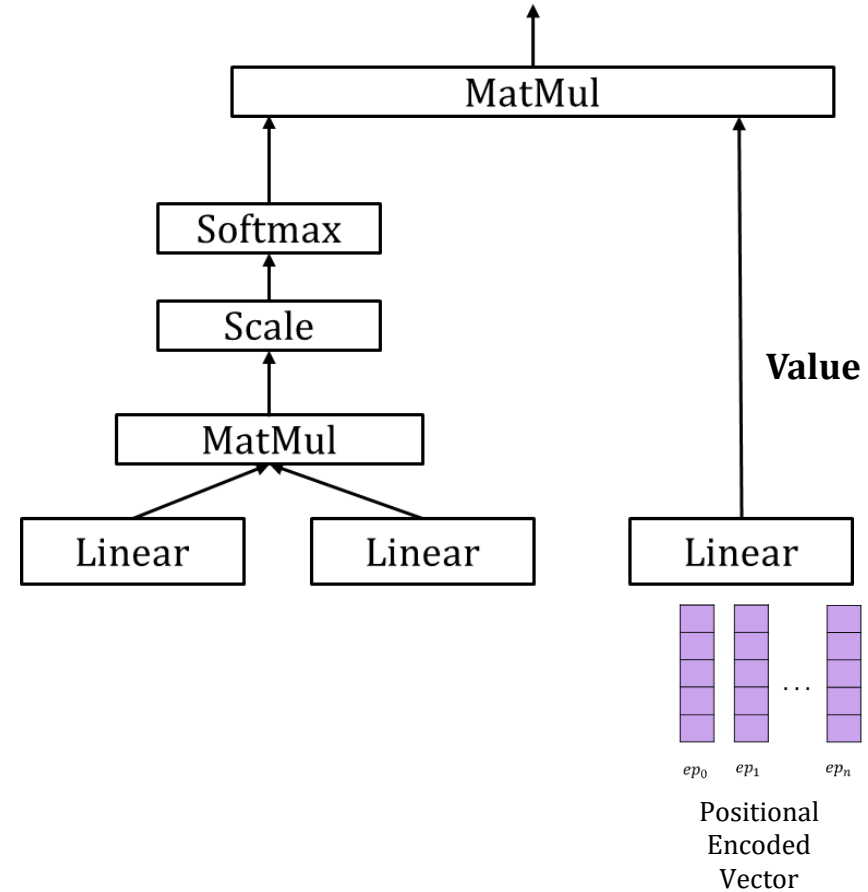
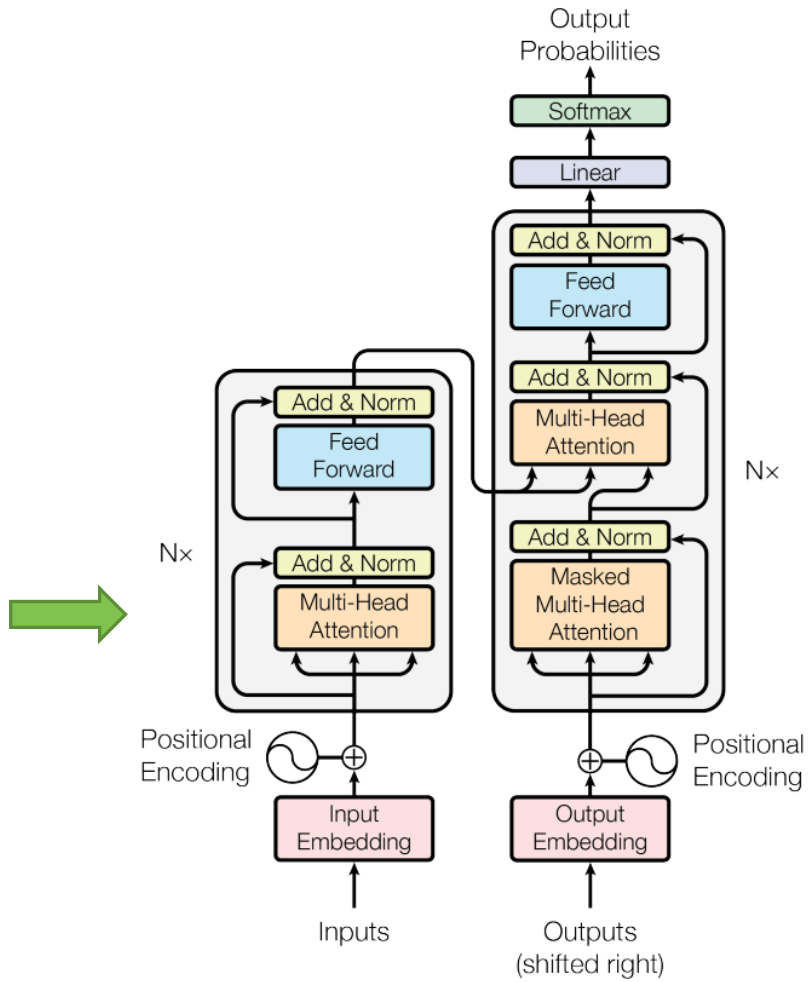
Multi-head Attention



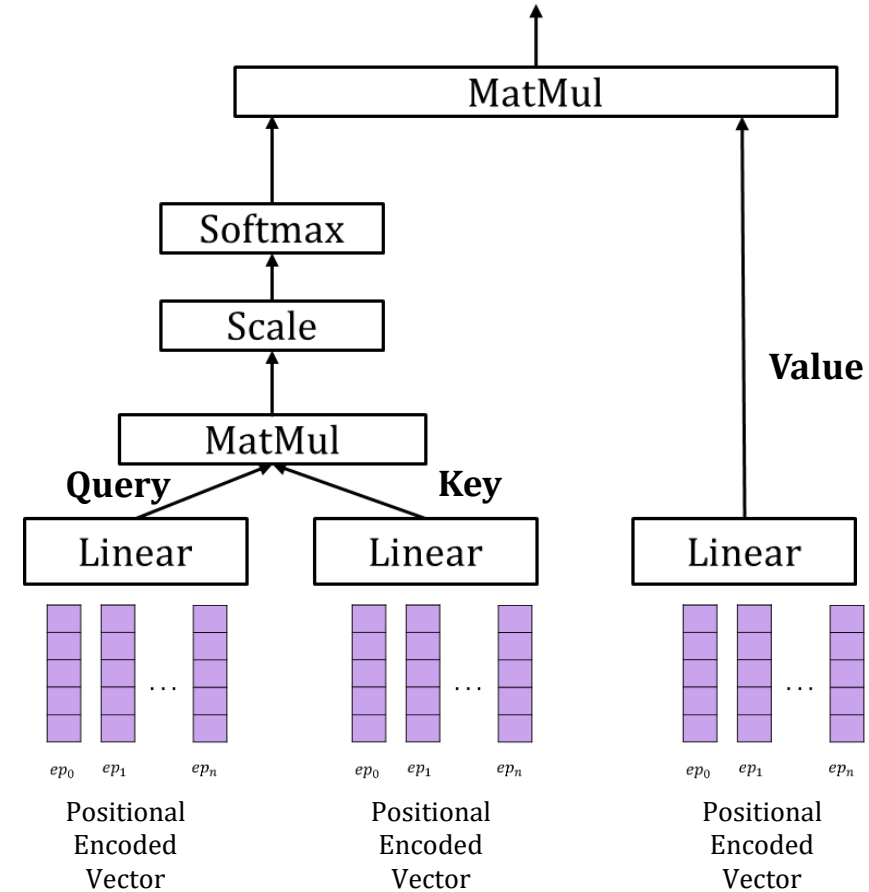
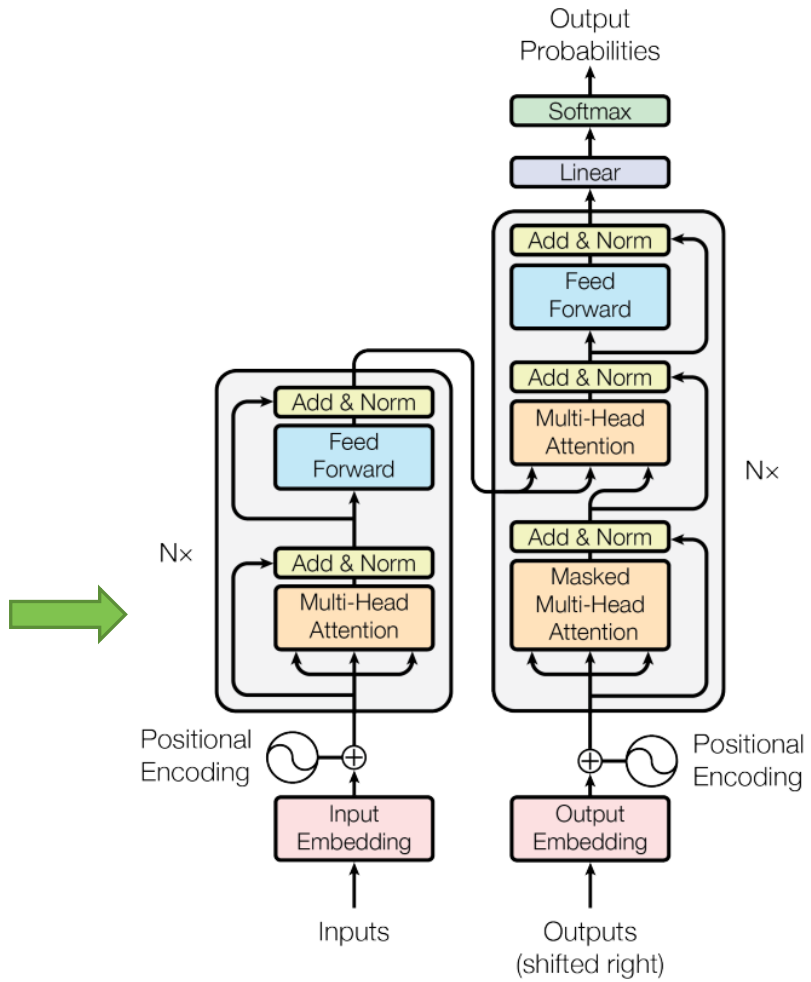
Multi-head Attention



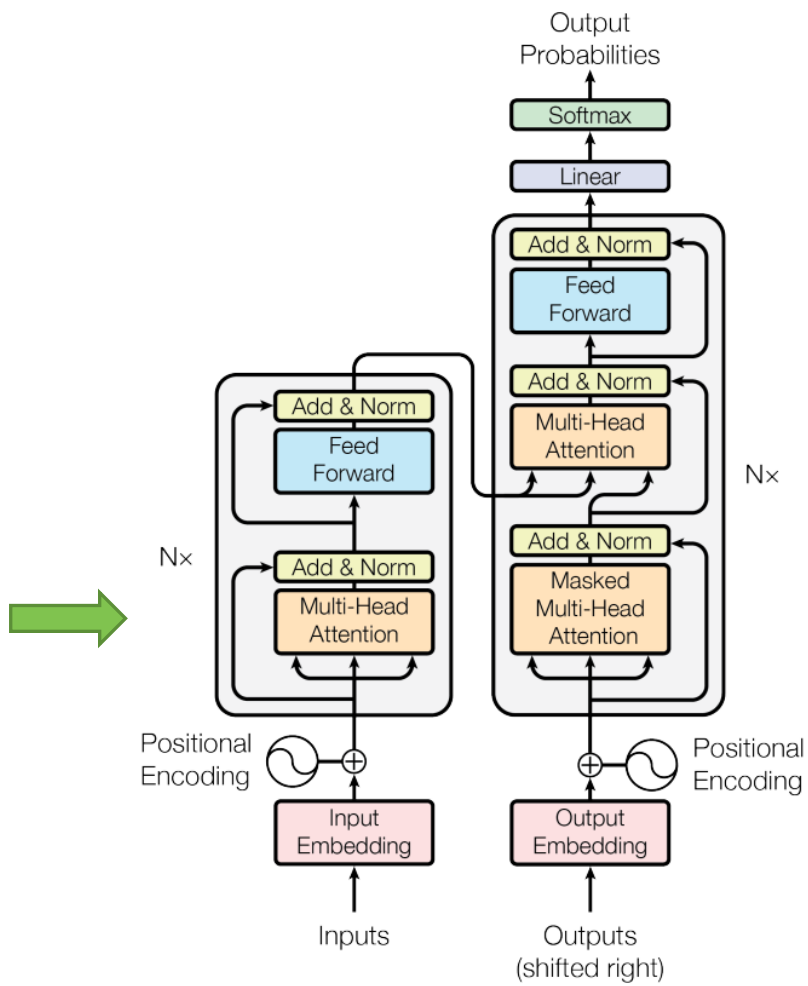
Multi-head Attention



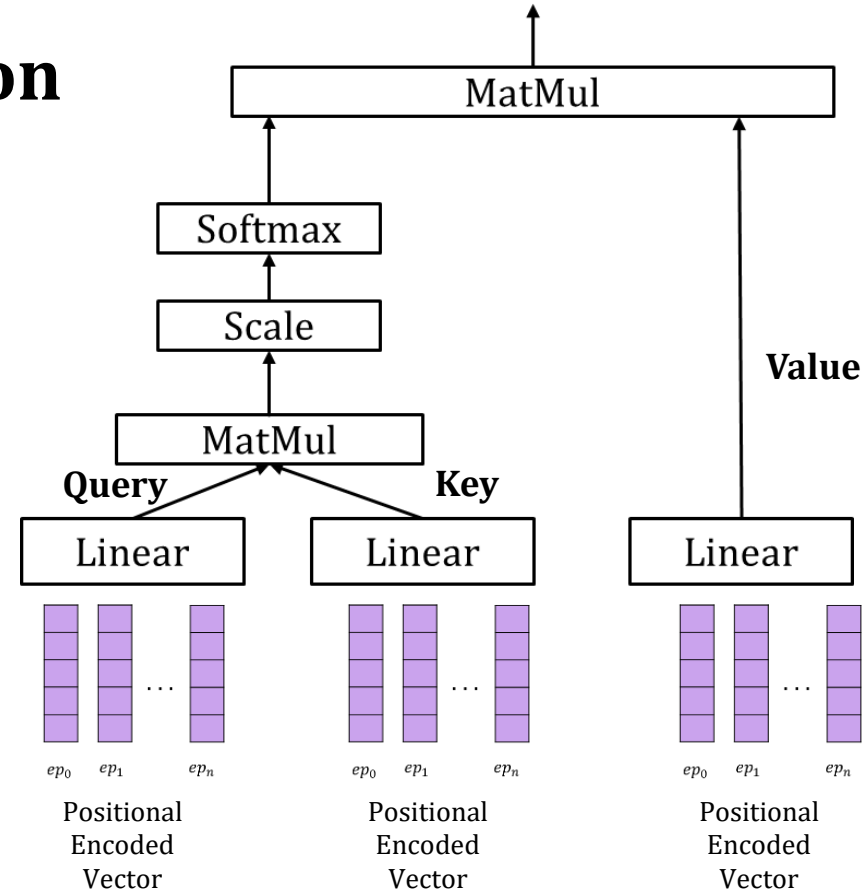
Multi-head Attention



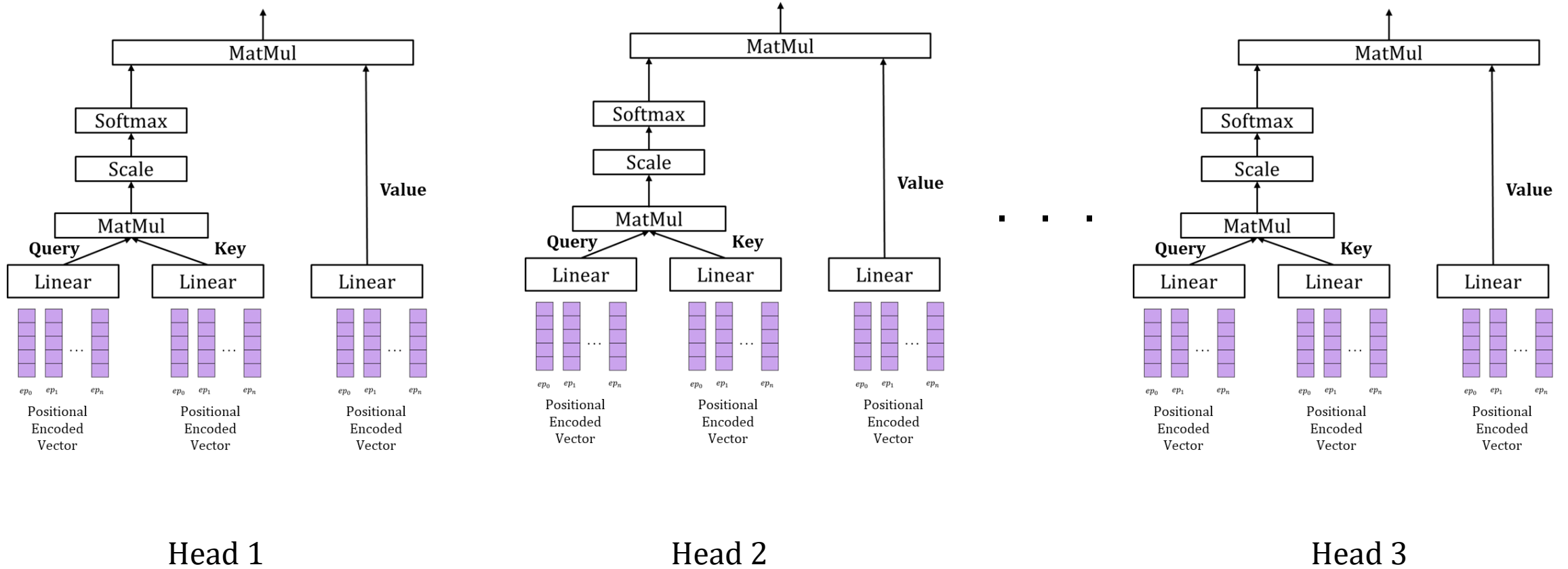
Multi-head Attention



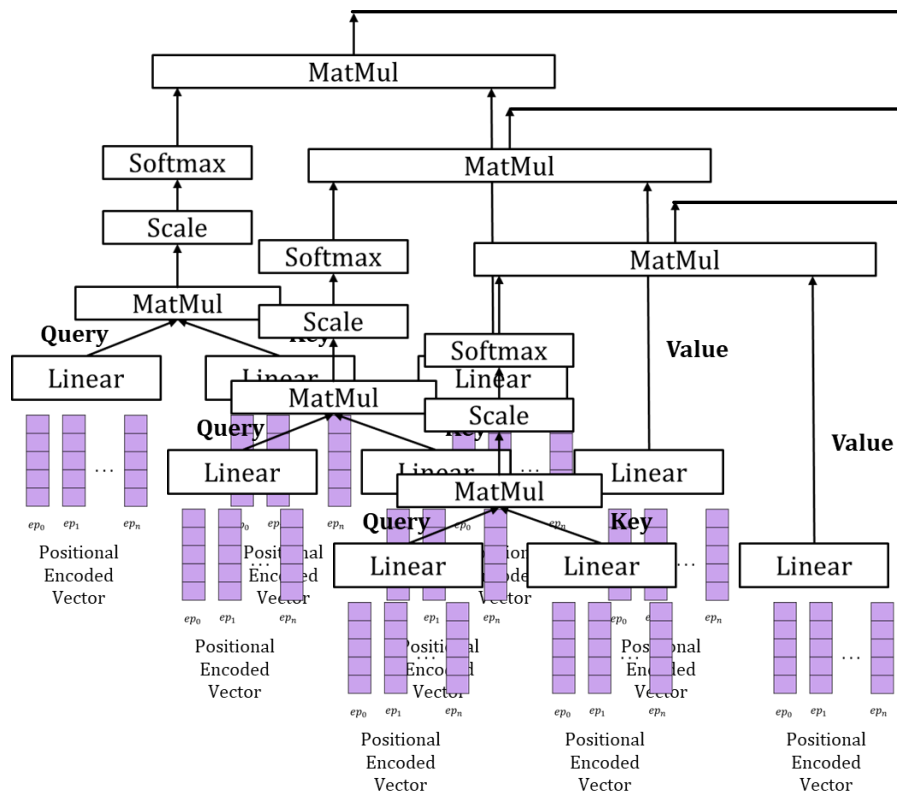
Self-Attention



Multi-head Self-Attention

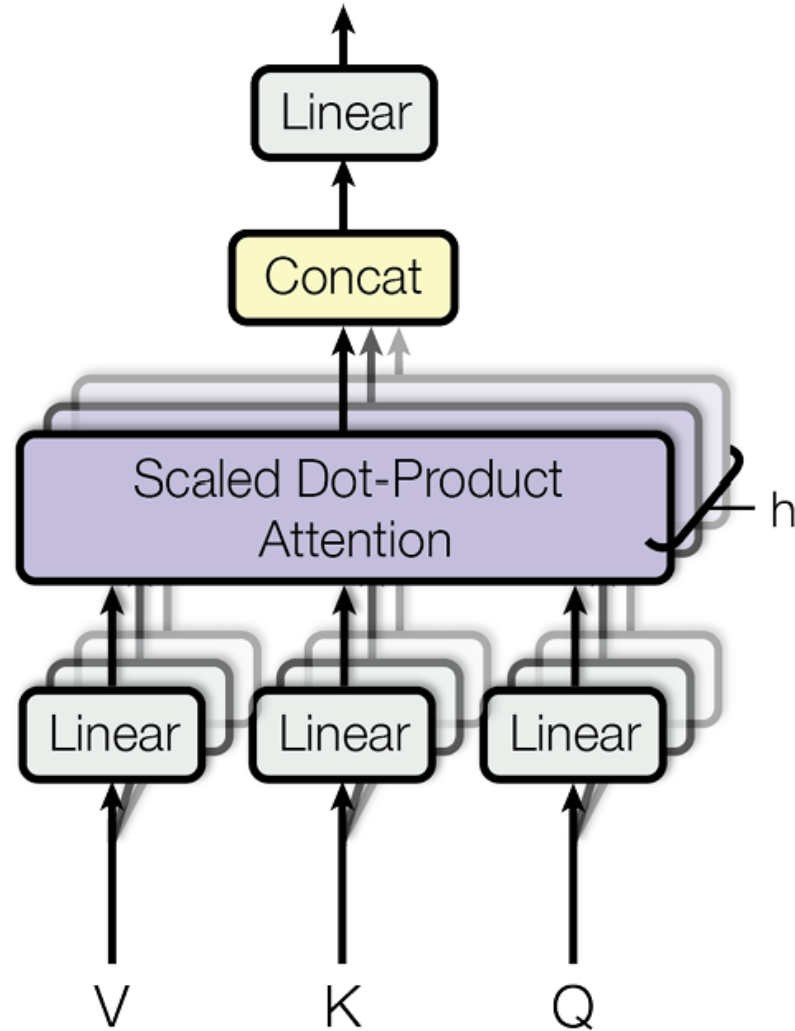


Multi-head Self-Attention



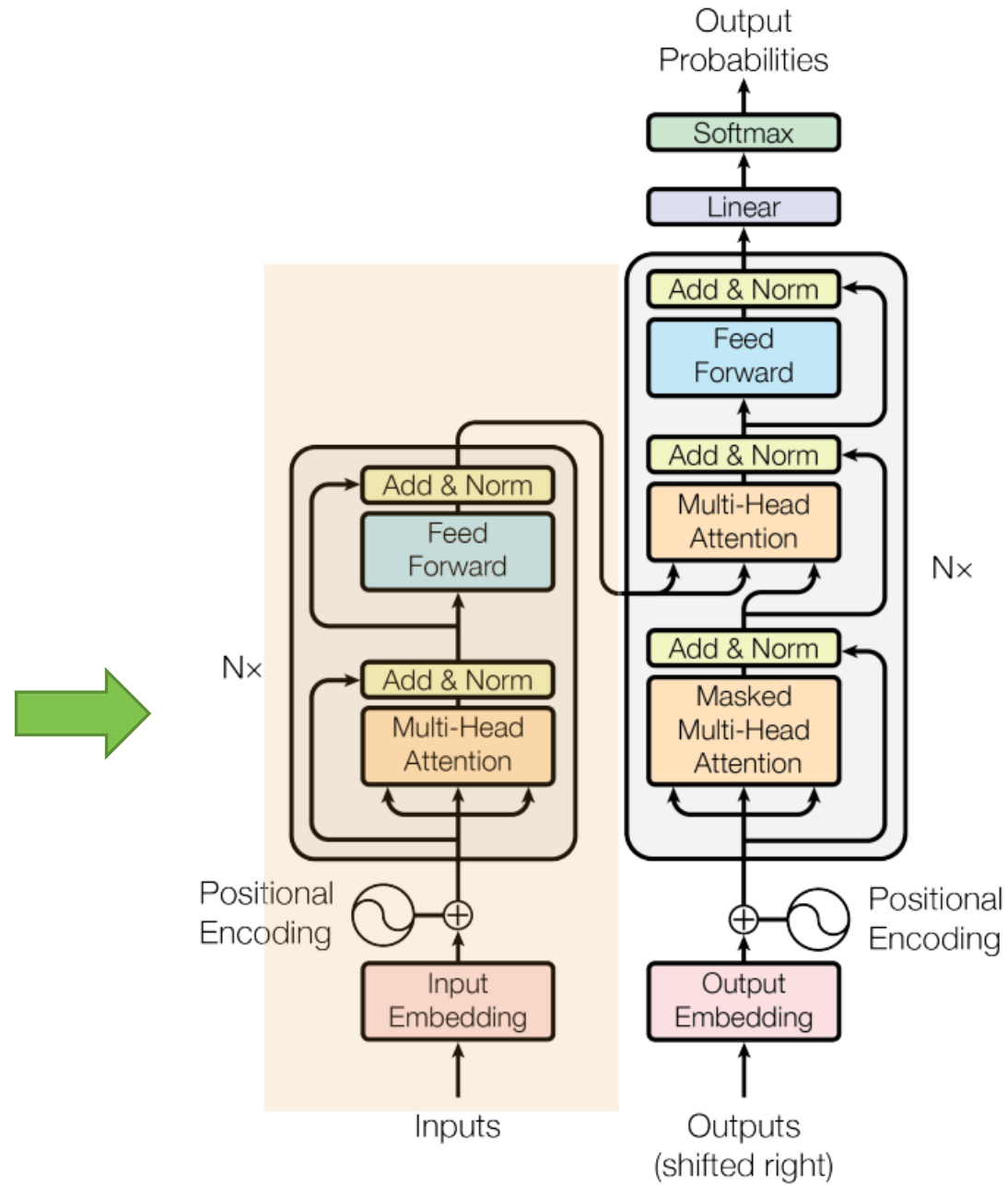
Concatenation

Multi-Head Attention

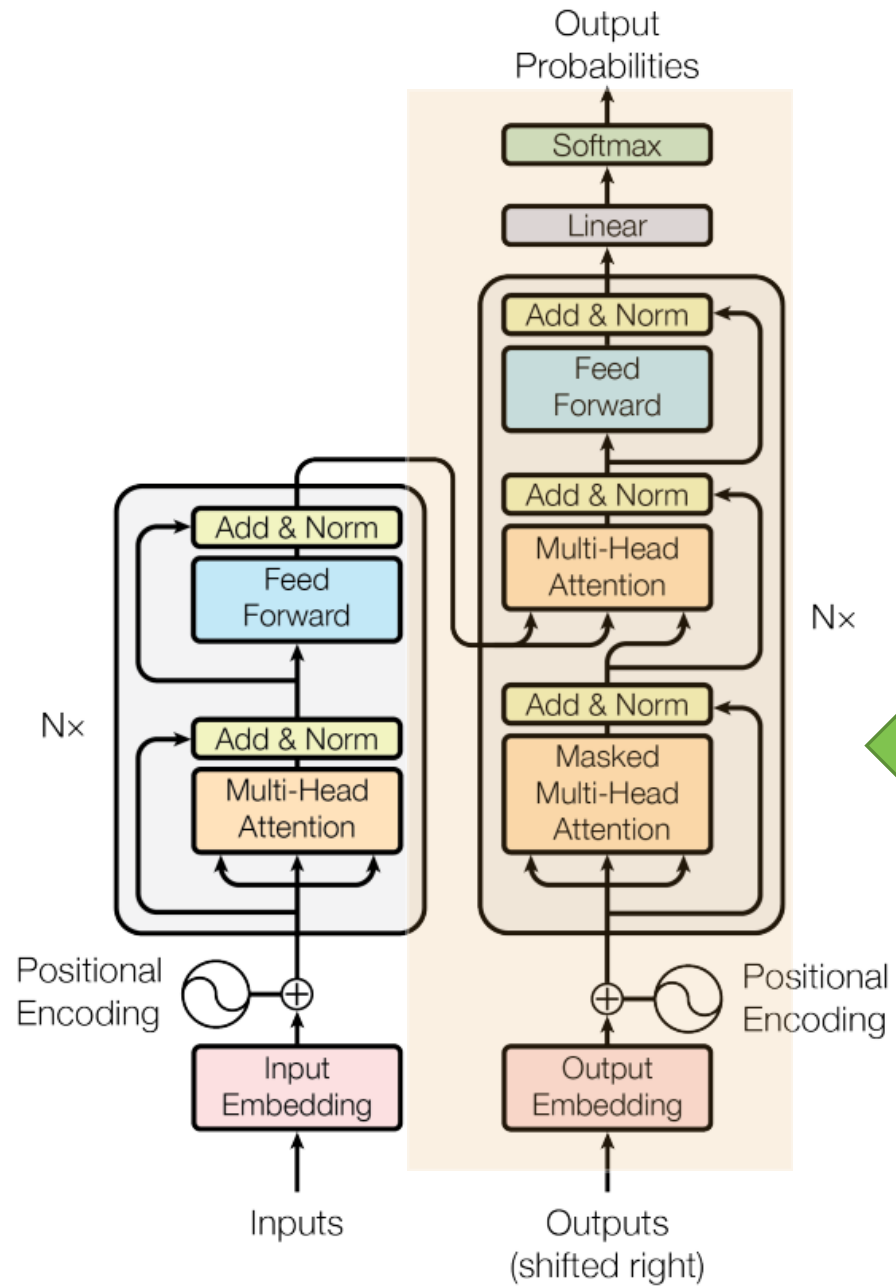


Multi-head Attention

Encoder



Decoder



Transformer: Step by Step

A ... Am

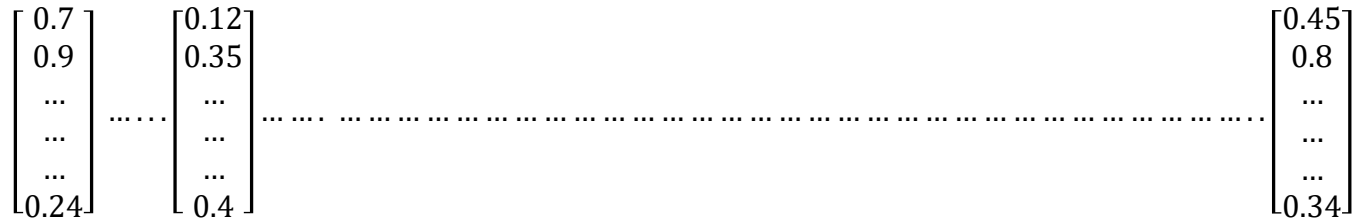
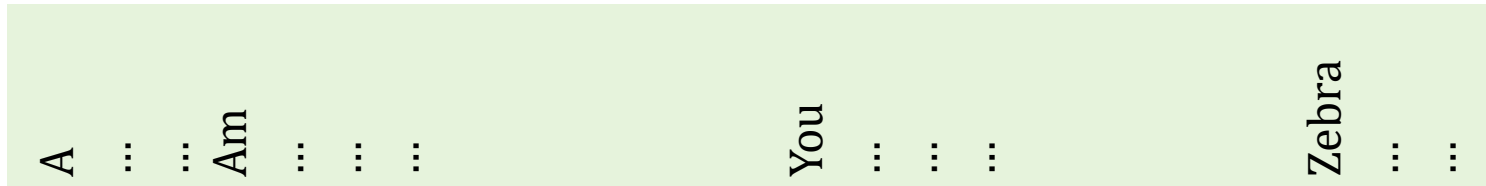
You

Zebra ...

The English vocabulary

(Let's say
~50K words)

Transformer: Step by Step

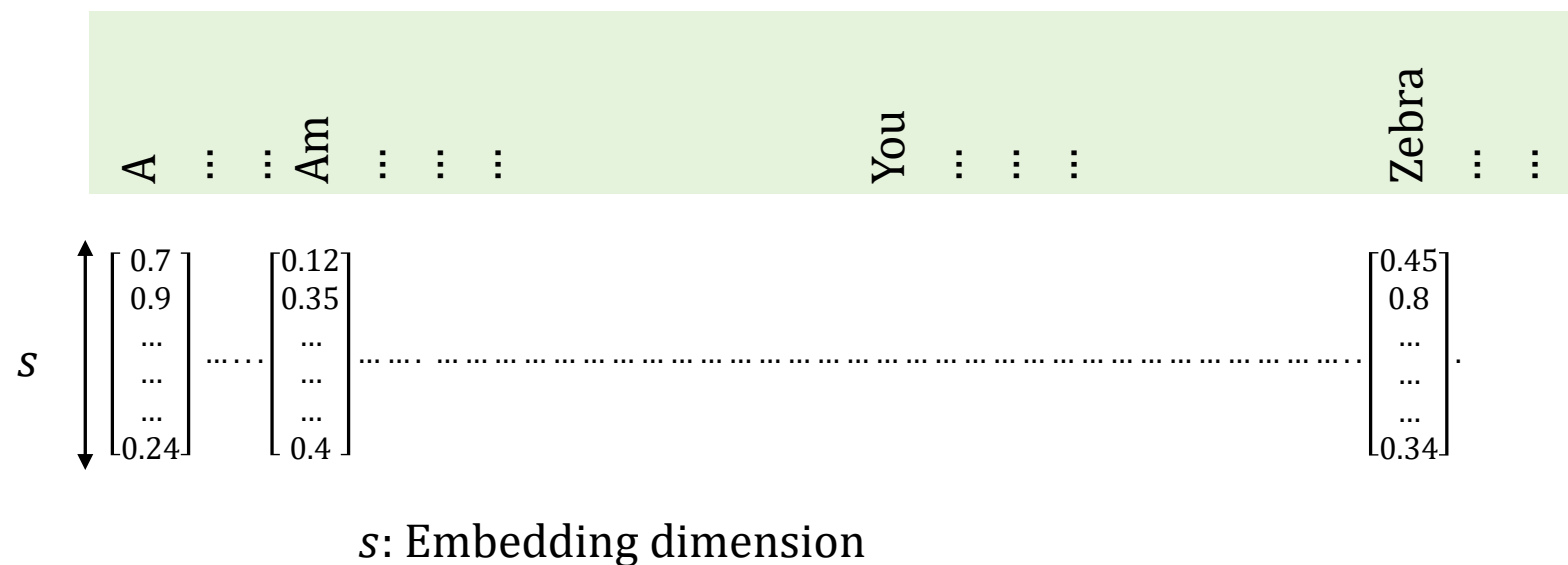


The English vocabulary
(Let's say
~50K words)

Each word
converted to
vector (Word
embedding)

Learnt from
data

Transformer: Step by Step



The English vocabulary
(Let's say
~50K words)

Each word
converted to
vector (Word
embedding)

Learnt from
data

In GPT 3, the vocabulary has 50257 tokens

Embedding dimension: 12288

Transformer: Step by Step

Consider the sentence translation problem. We want to translate the following sentence to Hindi

Where there is a will, there is a way

Transformer: Step by Step

Consider the sentence translation problem. We want to translate the following sentence to Hindi

Where there is a will, there is a way

Assume that the Transformer can take at most 20 token-long sentences as input.

Transformer: Step by Step

Where there is a will, there is a way

Assume that the Transformer can take at most 20 token-long sentences as input.

We first divide our input sentence into tokens and pad suitable number of additional tokens so that the total number of tokens become 20

Transformer: Step by Step

Where there is a will, there is a way

Assume that the Transformer can take at most 20 token-long sentences as input.

We first divide our input sentence into tokens and pad suitable number of additional tokens so that the total number of tokens become 20

Where there is a will , there is a way

Transformer: Step by Step

Where there is a will, there is a way

Assume that the Transformer can take at most 20 token-long sentences as input.

We first divide our input sentence into tokens and pad suitable number of additional tokens so that the total number of tokens become 20

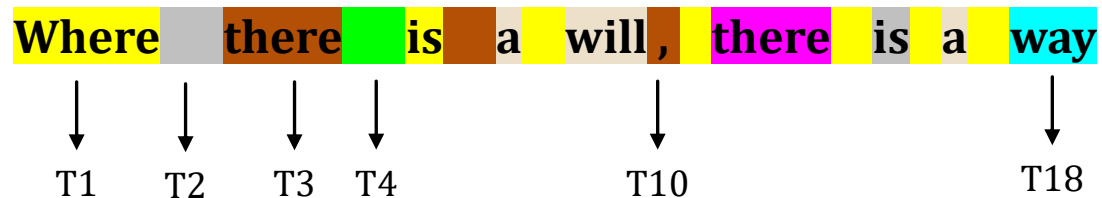
Where [padding] there is a will, there is a way

Transformer: Step by Step

Where there is a will, there is a way

Assume that the Transformer can take at most 20 token-long sentences as input.

We first divide our input sentence into tokens and pad suitable number of additional tokens so that the total number of tokens become 20

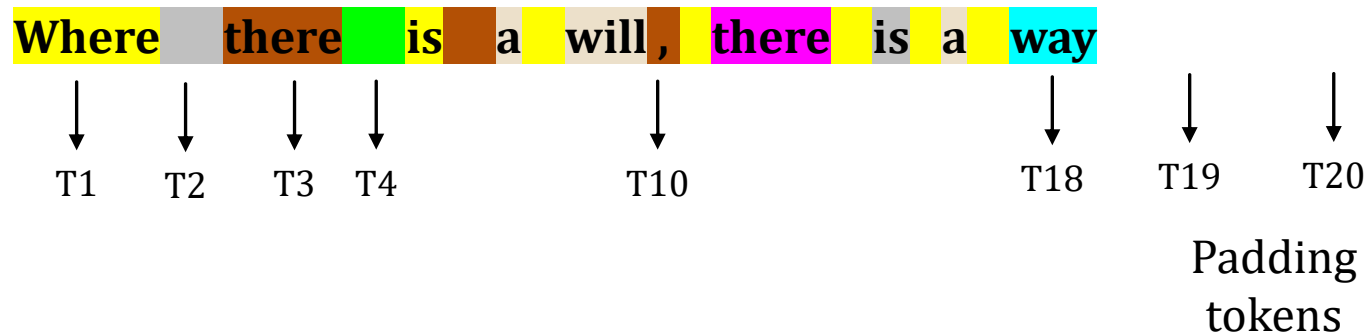


Transformer: Step by Step

Where there is a will, there is a way

Assume that the Transformer can take at most 20 token-long sentences as input.

We first divide our input sentence into tokens and pad suitable number of additional tokens so that the total number of tokens become 20



Transformer: Step by Step

Where there is a will, there is a way

↓
T1

↓
T2

↓
T3

↓
T4

↓
T10

↓
T18

↓
Tc

$$\begin{bmatrix} e_{11} \\ e_{12} \\ \dots \\ \dots \\ \dots \\ e_{1s} \end{bmatrix} \quad \begin{bmatrix} e_{21} \\ e_{22} \\ \dots \\ \dots \\ \dots \\ e_{2s} \end{bmatrix}$$

$$\begin{bmatrix} e_{c1} \\ e_{c2} \\ \dots \\ \dots \\ \dots \\ e_{cs} \end{bmatrix}$$

c is the max
number of
tokens

Transformer: Step by Step

Where there is a will, there is a way

↓
T1

↓
T2

↓
T3

↓
T4

↓
T10

↓
T18

↓
Tc

Add positional
embedding

$$\begin{bmatrix} e_{11} \\ e_{12} \\ \dots \\ \dots \\ \dots \\ e_{1s} \end{bmatrix} \quad \begin{bmatrix} e_{21} \\ e_{22} \\ \dots \\ \dots \\ \dots \\ e_{2s} \end{bmatrix}$$

+

+

$$\begin{bmatrix} p_{11} \\ p_{12} \\ \dots \\ \dots \\ \dots \\ p_{1s} \end{bmatrix} \quad \begin{bmatrix} p_{21} \\ p_{22} \\ \dots \\ \dots \\ \dots \\ p_{2s} \end{bmatrix}$$

$$\begin{bmatrix} e_{c1} \\ e_{c2} \\ \dots \\ \dots \\ \dots \\ e_{cs} \end{bmatrix}$$

+

$$\begin{bmatrix} p_{c1} \\ p_{c2} \\ \dots \\ \dots \\ \dots \\ p_{cs} \end{bmatrix}$$

Where there is a will, there is a way

↓
T1

↓
T2

↓
T3

↓
T4

↓
T10

↓
T18

↓
Tc

$$\begin{bmatrix} e_{11} \\ e_{12} \\ \dots \\ \dots \\ \dots \\ e_{1s} \end{bmatrix}$$

$$\begin{bmatrix} e_{21} \\ e_{22} \\ \dots \\ \dots \\ \dots \\ e_{2s} \end{bmatrix}$$

+

+

$$\begin{bmatrix} p_{11} \\ p_{12} \\ \dots \\ \dots \\ \dots \\ p_{1s} \end{bmatrix}$$

$$\begin{bmatrix} p_{21} \\ p_{22} \\ \dots \\ \dots \\ \dots \\ p_{2s} \end{bmatrix}$$

=

=

$$\begin{bmatrix} ep_{11} \\ ep_{12} \\ \dots \\ \dots \\ \dots \\ ep_{1s} \end{bmatrix}$$

$$\begin{bmatrix} ep_{21} \\ ep_{22} \\ \dots \\ \dots \\ \dots \\ ep_{2s} \end{bmatrix}$$

$$\begin{bmatrix} e_{c1} \\ e_{c2} \\ \dots \\ \dots \\ \dots \\ e_{cs} \end{bmatrix}$$

+

$$\begin{bmatrix} p_{c1} \\ p_{c2} \\ \dots \\ \dots \\ \dots \\ p_{cs} \end{bmatrix}$$

=

$$\begin{bmatrix} ep_{c1} \\ ep_{c2} \\ \dots \\ \dots \\ \dots \\ ep_{cs} \end{bmatrix}$$

Transformer:
Step by Step

Where there is a will, there is a way

↓
T1

↓
T2

↓
T3

↓
T4

↓
T10

↓
T18

↓
Tc

$$\begin{bmatrix} e_{11} \\ e_{12} \\ \dots \\ \dots \\ \dots \\ e_{1s} \end{bmatrix}$$

$$\begin{bmatrix} e_{21} \\ e_{22} \\ \dots \\ \dots \\ \dots \\ e_{2s} \end{bmatrix}$$

+

+

$$\begin{bmatrix} p_{11} \\ p_{12} \\ \dots \\ \dots \\ \dots \\ p_{1s} \end{bmatrix}$$

$$\begin{bmatrix} p_{21} \\ p_{22} \\ \dots \\ \dots \\ \dots \\ p_{2s} \end{bmatrix}$$

=

=

$$\begin{bmatrix} ep_{11} \\ ep_{12} \\ \dots \\ \dots \\ \dots \\ ep_{1s} \end{bmatrix}$$

$$\begin{bmatrix} ep_{21} \\ ep_{22} \\ \dots \\ \dots \\ \dots \\ ep_{2s} \end{bmatrix}$$

ep_1

ep_2

$$\begin{bmatrix} e_{c1} \\ e_{c2} \\ \dots \\ \dots \\ \dots \\ e_{cs} \end{bmatrix}$$

+

$$\begin{bmatrix} p_{c1} \\ p_{c2} \\ \dots \\ \dots \\ \dots \\ p_{cs} \end{bmatrix}$$

=

$$\begin{bmatrix} ep_{c1} \\ ep_{c2} \\ \dots \\ \dots \\ \dots \\ ep_{cs} \end{bmatrix}$$

ep_c

Transformer:
Step by Step

Where there is a will, there is a way



T1

ep_1



T2

ep_2



T3



T4



T10



T18



Tc

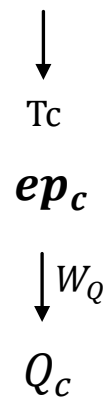
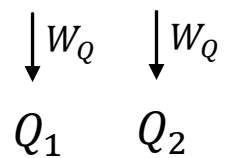
ep_c

Transformer:
Step by Step

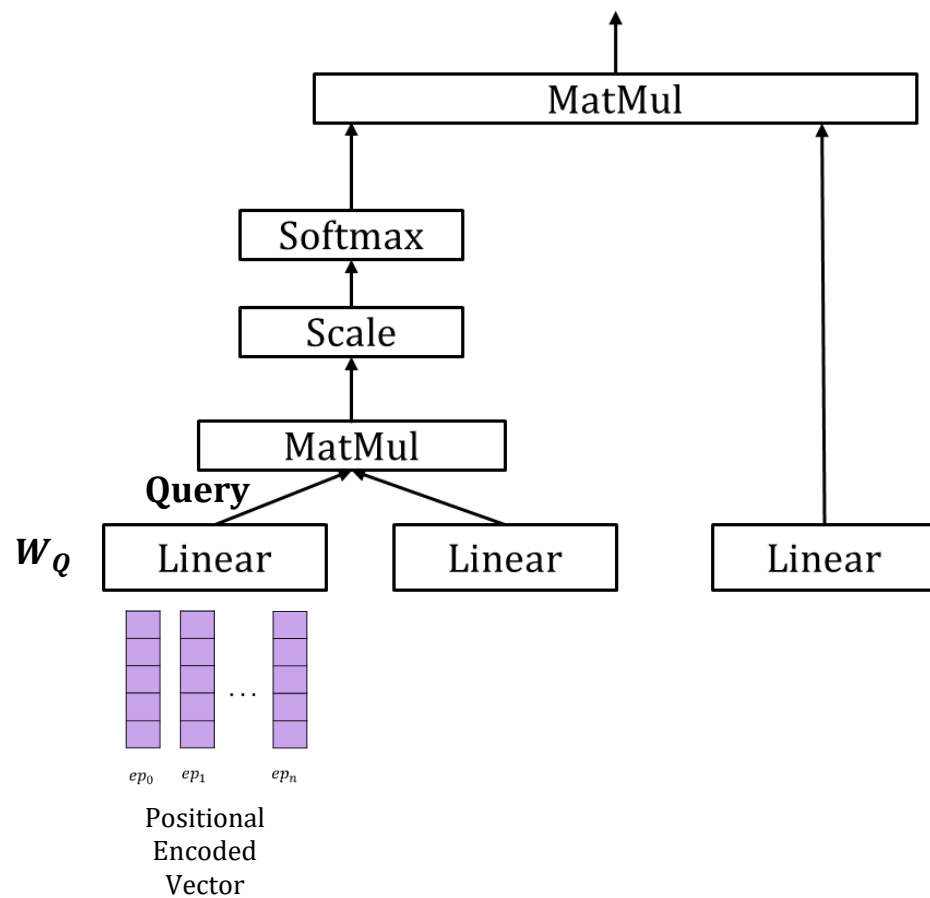
Where there is a will, there is a way



ep_1 ep_2



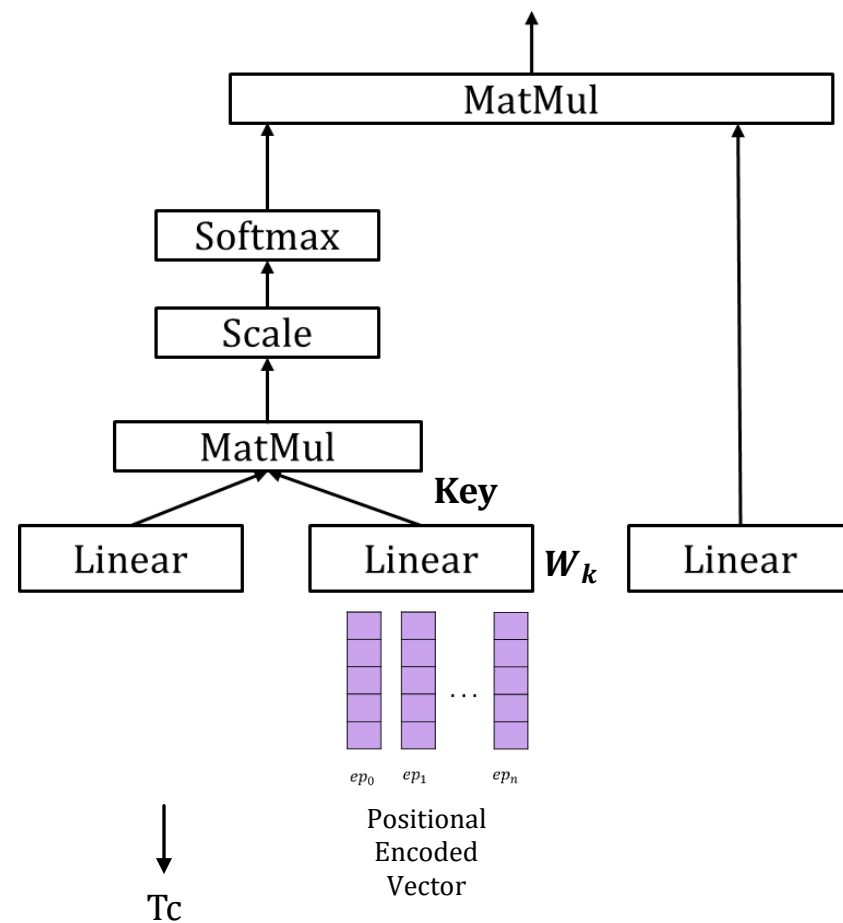
Q_i : Query vector of dimension d for i^{th} token



Where there is a will, there is a way

T1 T2 T3 T4 T10 T18

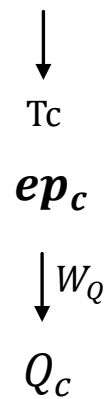
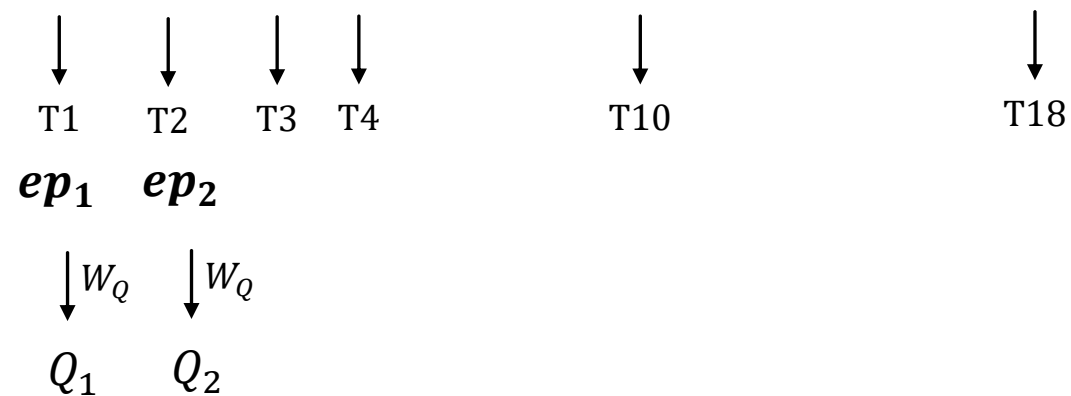
ep_1 ep_2
 $\downarrow W_k$ $\downarrow W_k$
 K_1 K_2



ep_c
 $\downarrow W_k$
 K_c

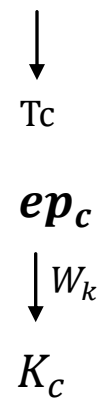
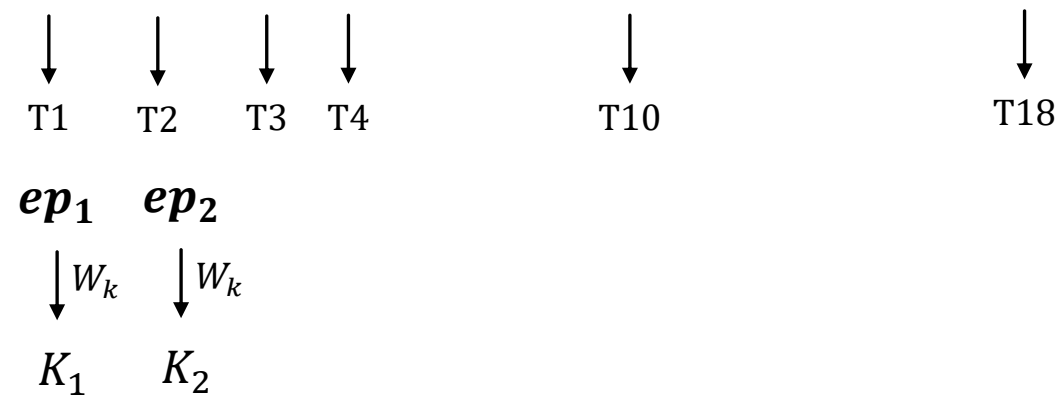
K_i : Key vector of dimension d for i^{th} token

Where there is a will, there is a way



Q_i : Query vector of dimension d for i^{th} token

Where there is a will, there is a way



K_i : Key vector of dimension d for i^{th} token

Where there is a will, there is a way

$$\begin{array}{c} ep_1 \\ \downarrow W_Q \\ Q_1 \end{array}$$
$$\begin{array}{c} ep_2 \\ \downarrow W_Q \\ Q_2 \end{array}$$
$$\begin{array}{c} ep_c \\ \downarrow W_Q \\ Q_c \end{array}$$

Where there is a will, there is a way

$$ep_1 \xrightarrow{W_K} K_1$$
$$ep_2 \xrightarrow{W_K} K_2$$
$$ep_c \xrightarrow{W_K} K_c$$

Where there is a will, there is a way

$$\begin{matrix} ep_1 \\ \downarrow W_Q \\ Q_1 \end{matrix}$$

$$\begin{matrix} ep_2 \\ \downarrow W_Q \\ Q_2 \end{matrix}$$

$$\begin{matrix} ep_c \\ \downarrow W_Q \\ Q_c \end{matrix}$$

Where there is a will, there is a way

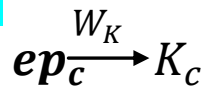
$$ep_1 \xrightarrow{W_K} K_1$$

$$ep_2 \xrightarrow{W_K} K_2$$

$$ep_c \xrightarrow{W_K} K_c$$

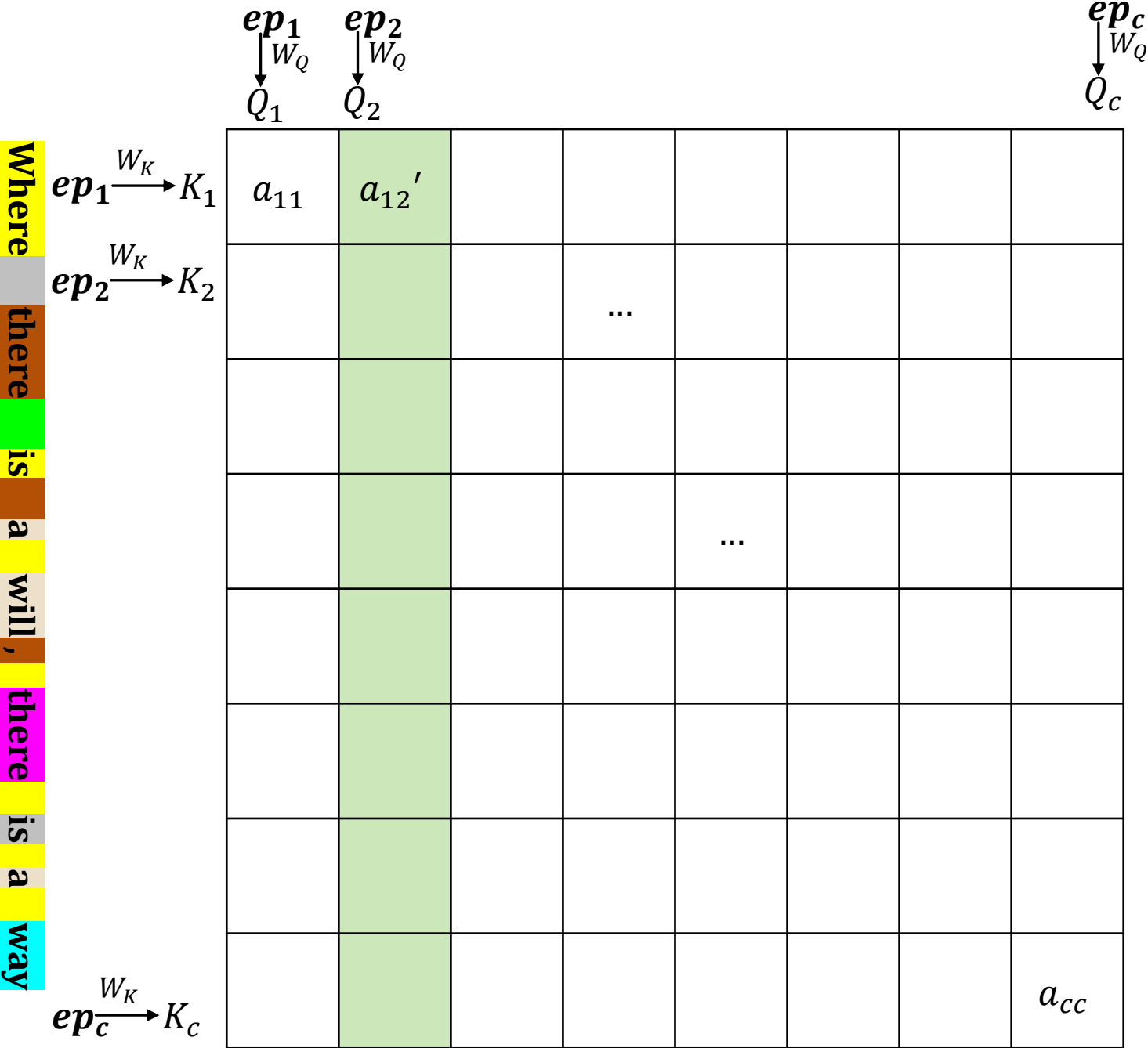
$\frac{Q_1 \cdot K_1}{\sqrt{d}}$	$\frac{Q_2 \cdot K_1}{\sqrt{d}}$						
			...				
				...			
							$\frac{Q_c \cdot K_c}{\sqrt{d}}$

Where there is a will, there is a way

[illegible]

Entry (i, j) indicate some sort of correlation between the token of the i^{th} row and the token of the j^{th} row

Where there is a will, there is a way

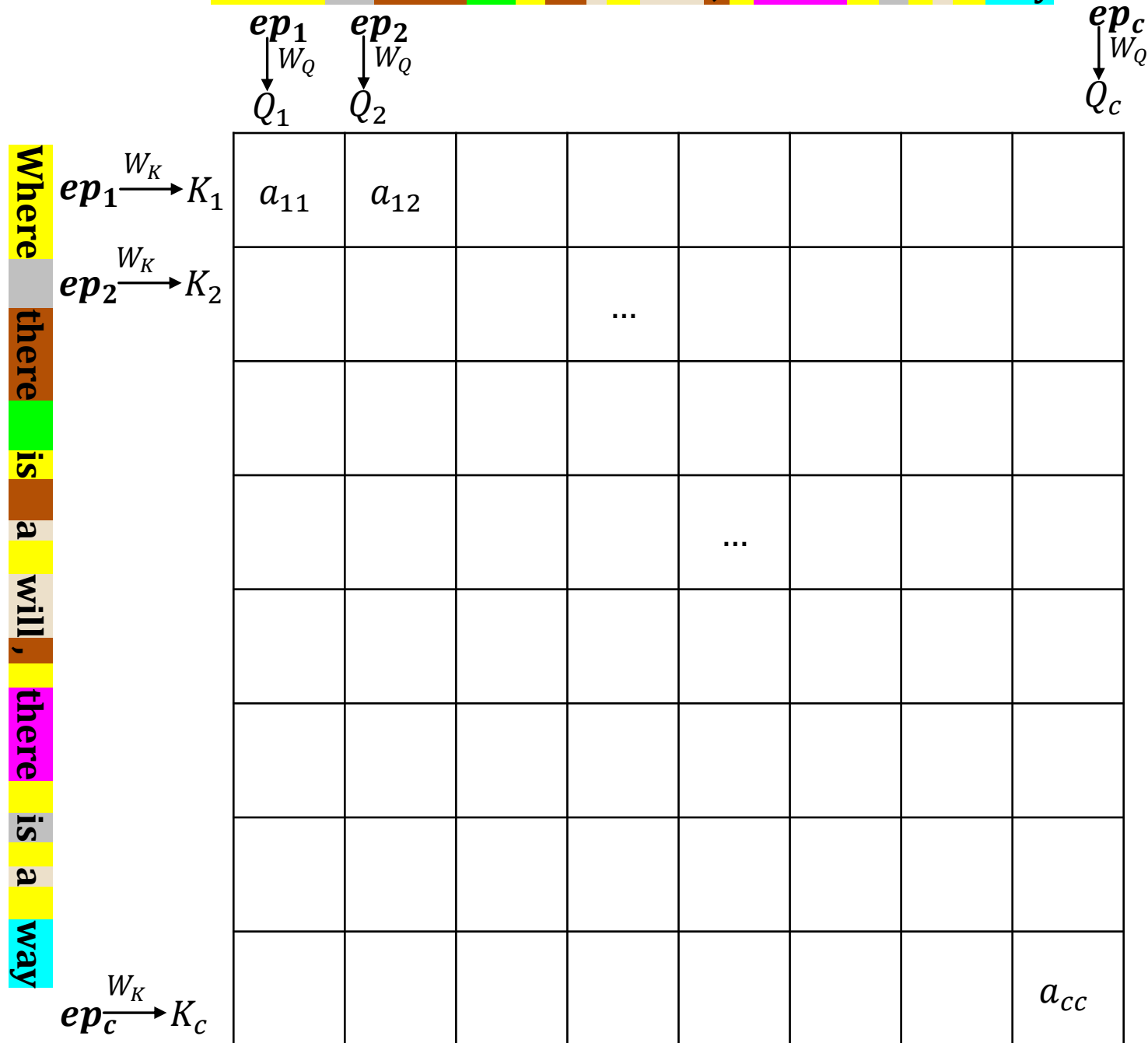


Attention matrix

For each query column, we apply softmax. As a result, each column kind of indicates the probabilities of the corresponding query to match with different keys

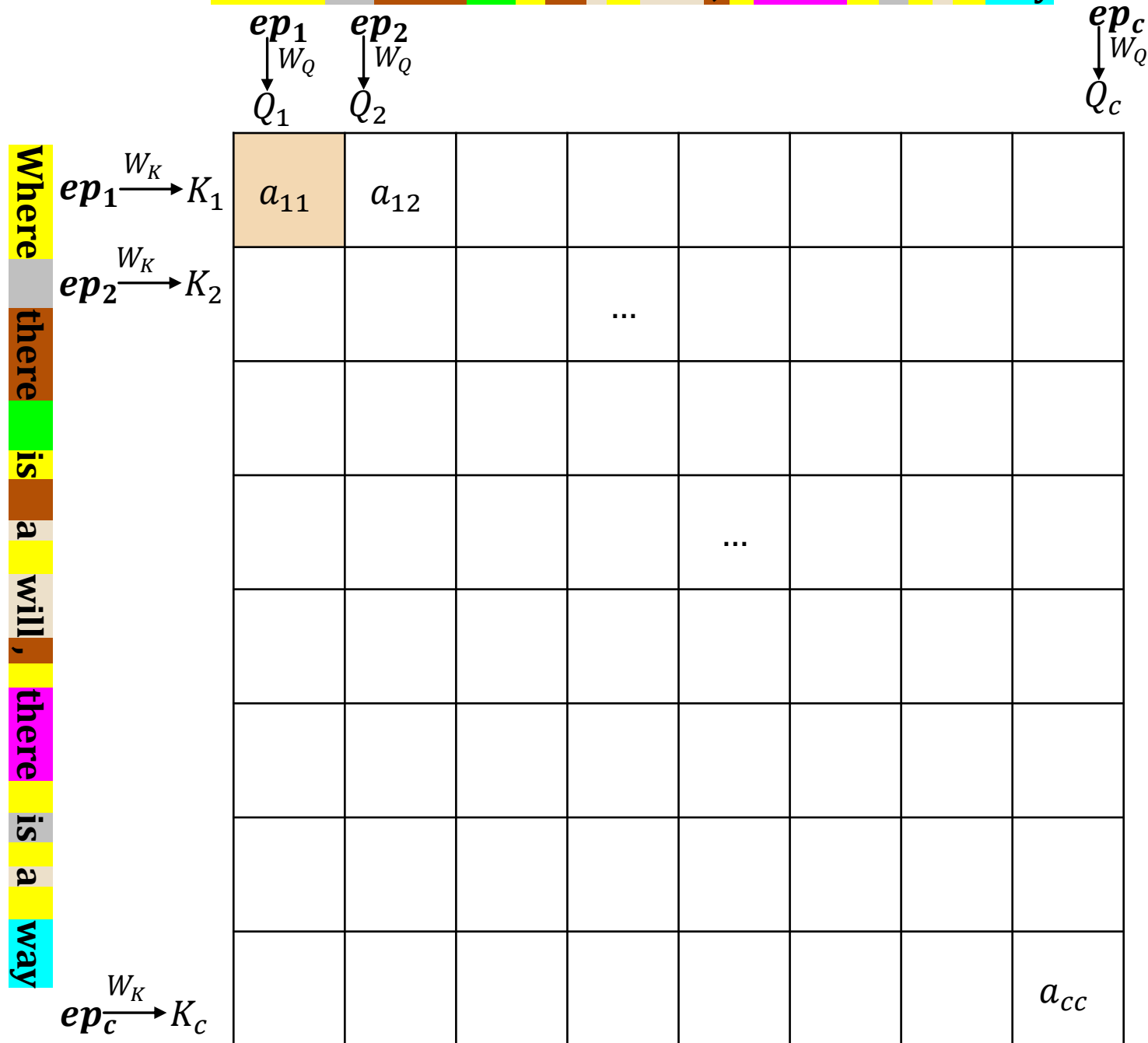
a'_{12} : Value of the entry after softmax

Where there is a will, there is a way



We may want that a query (a token) is influenced by only past keys (past tokens)

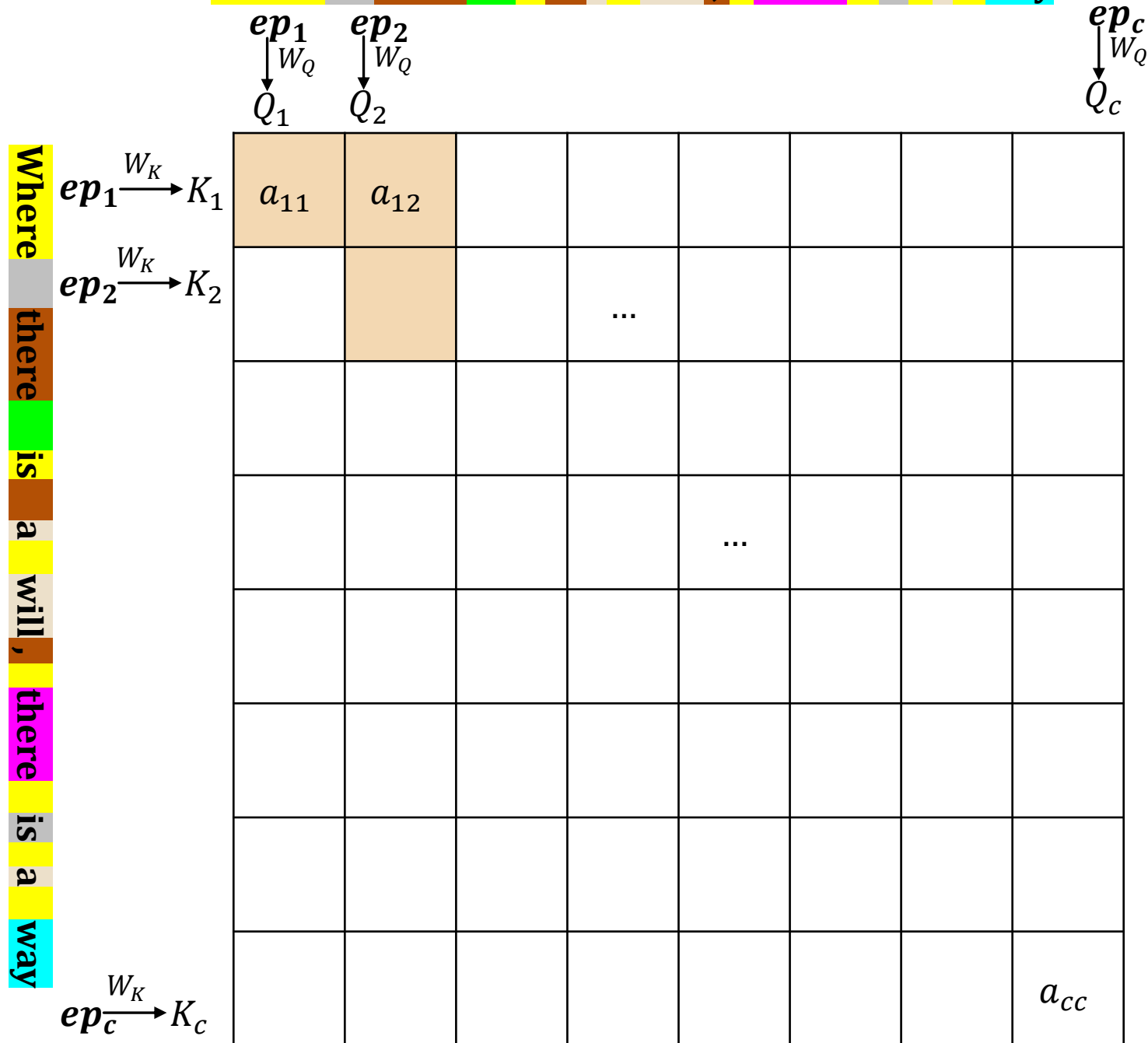
Where there is a will, there is a way



We may want that a query (a token) is influenced by only past keys (past tokens)

Q_1 is influenced by only K_1

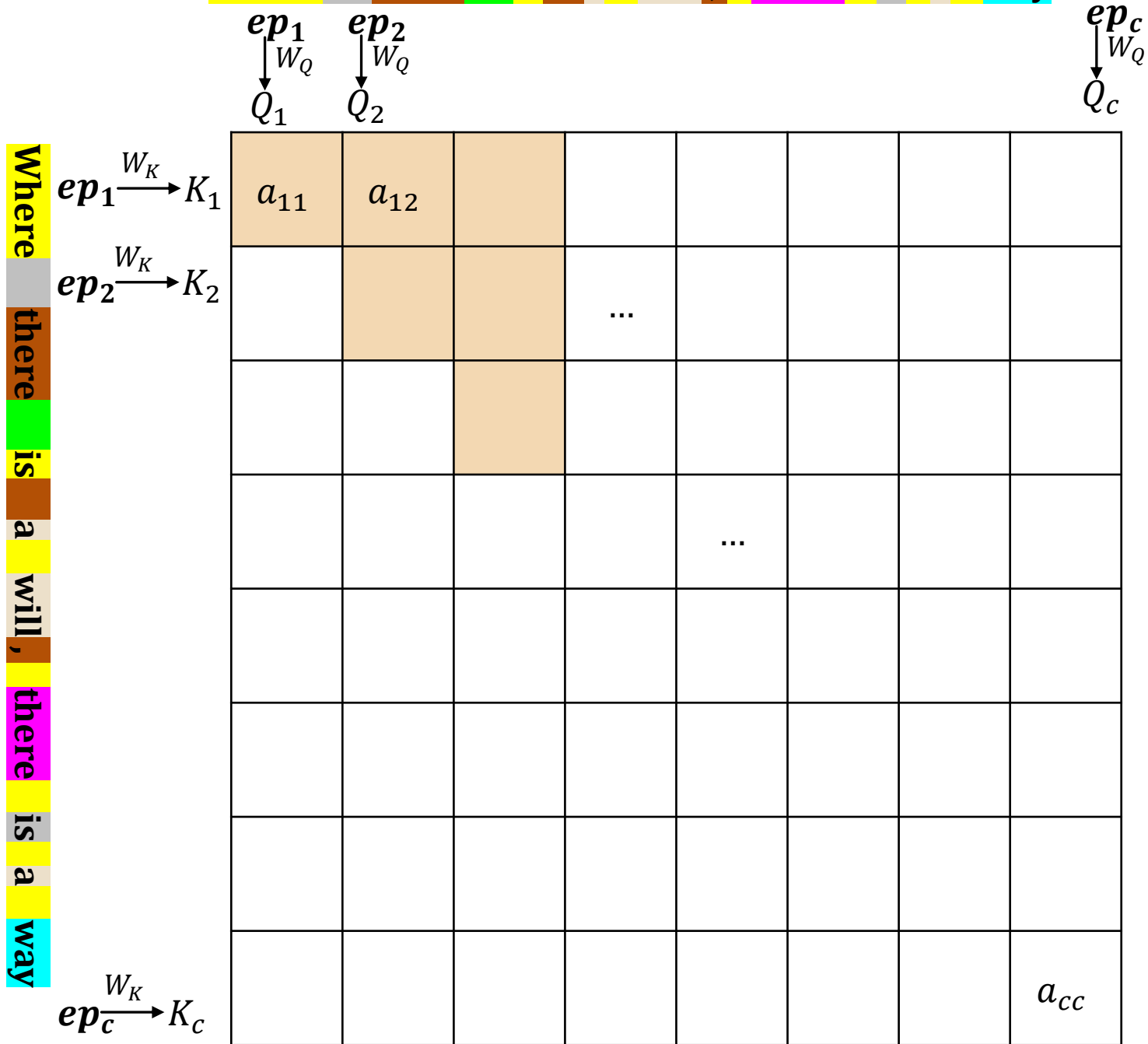
Where there is a will, there is a way



We may want that a query (a token) is influenced by only past keys (past tokens)

Q_1 is influenced by only K_1
 Q_2 is influenced by only K_1, K_2

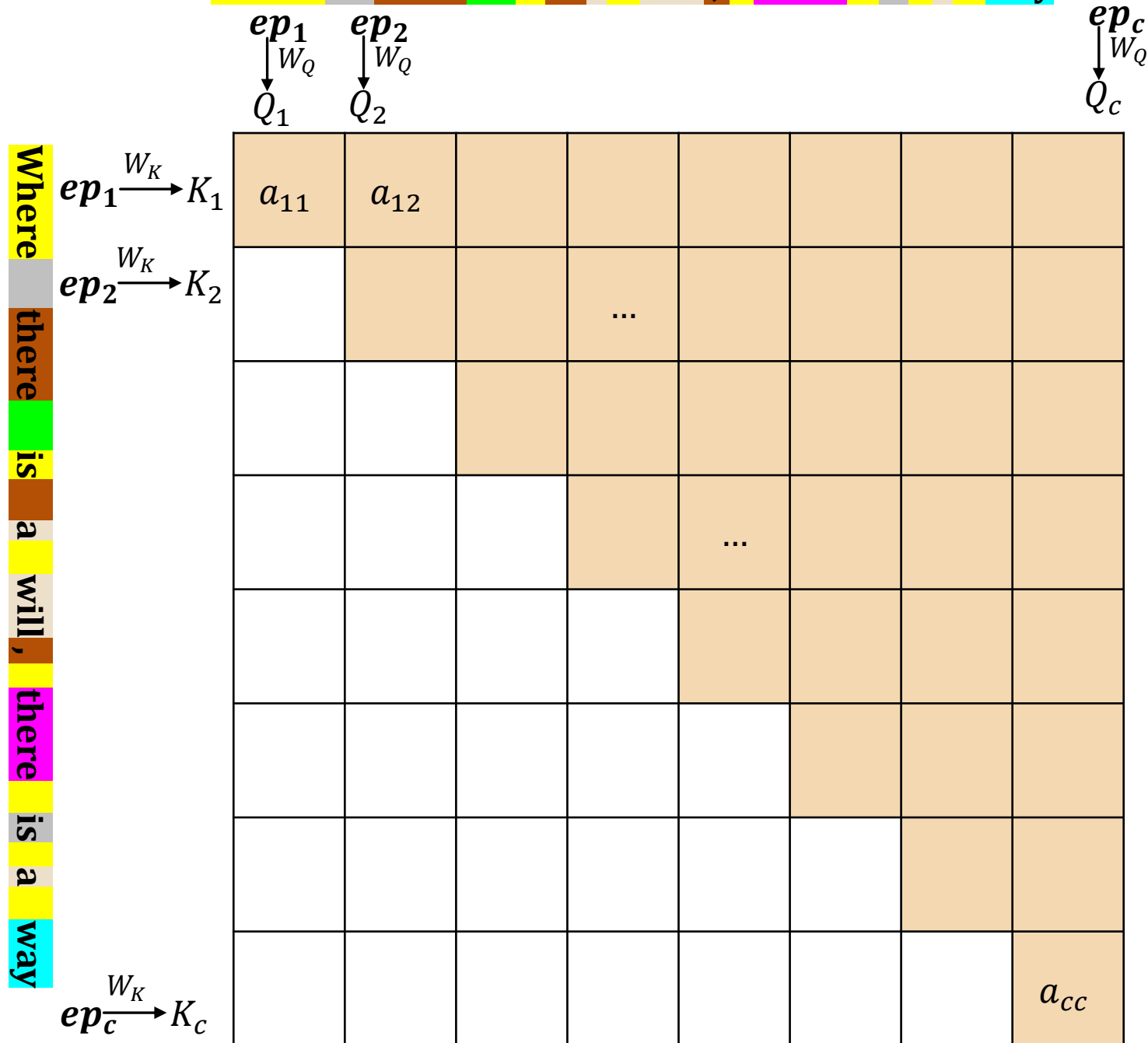
Where there is a will, there is a way



**We may want that a query
(a token) is influenced by
only past keys (past tokens)**

Q_1 is influenced by only K_1
 Q_2 is influenced by only
 K_1, K_2
 Q_3 is influenced by only
 K_1, K_2, K_3 , and so on

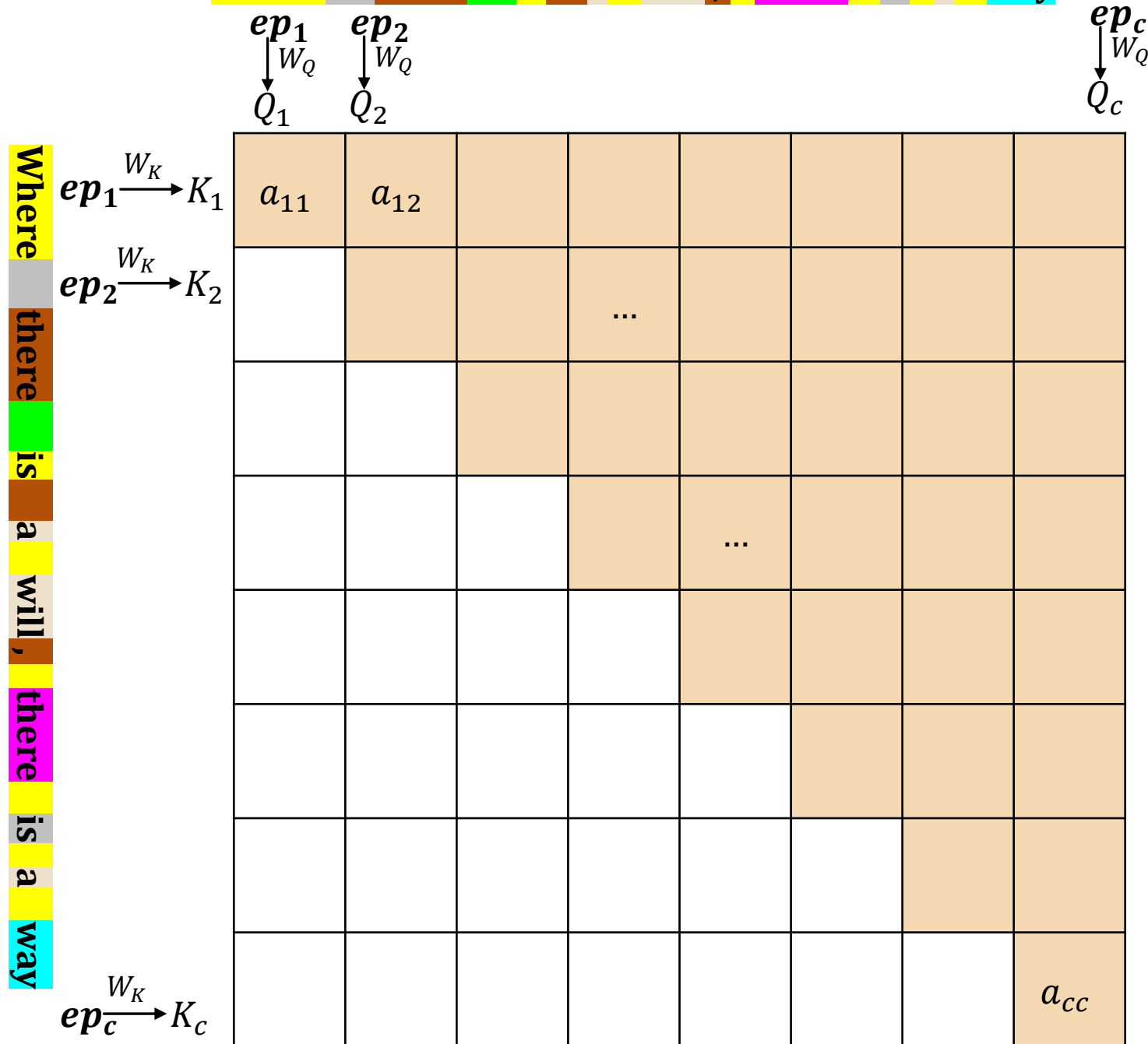
Where there is a will, there is a way



We may want that a query (a token) is influenced by only past keys (past tokens)

Q_1 is influenced by only K_1
 Q_2 is influenced by only K_1, K_2
 Q_3 is influenced by only K_1, K_2, K_3 , and so on

Where there is a will, there is a way

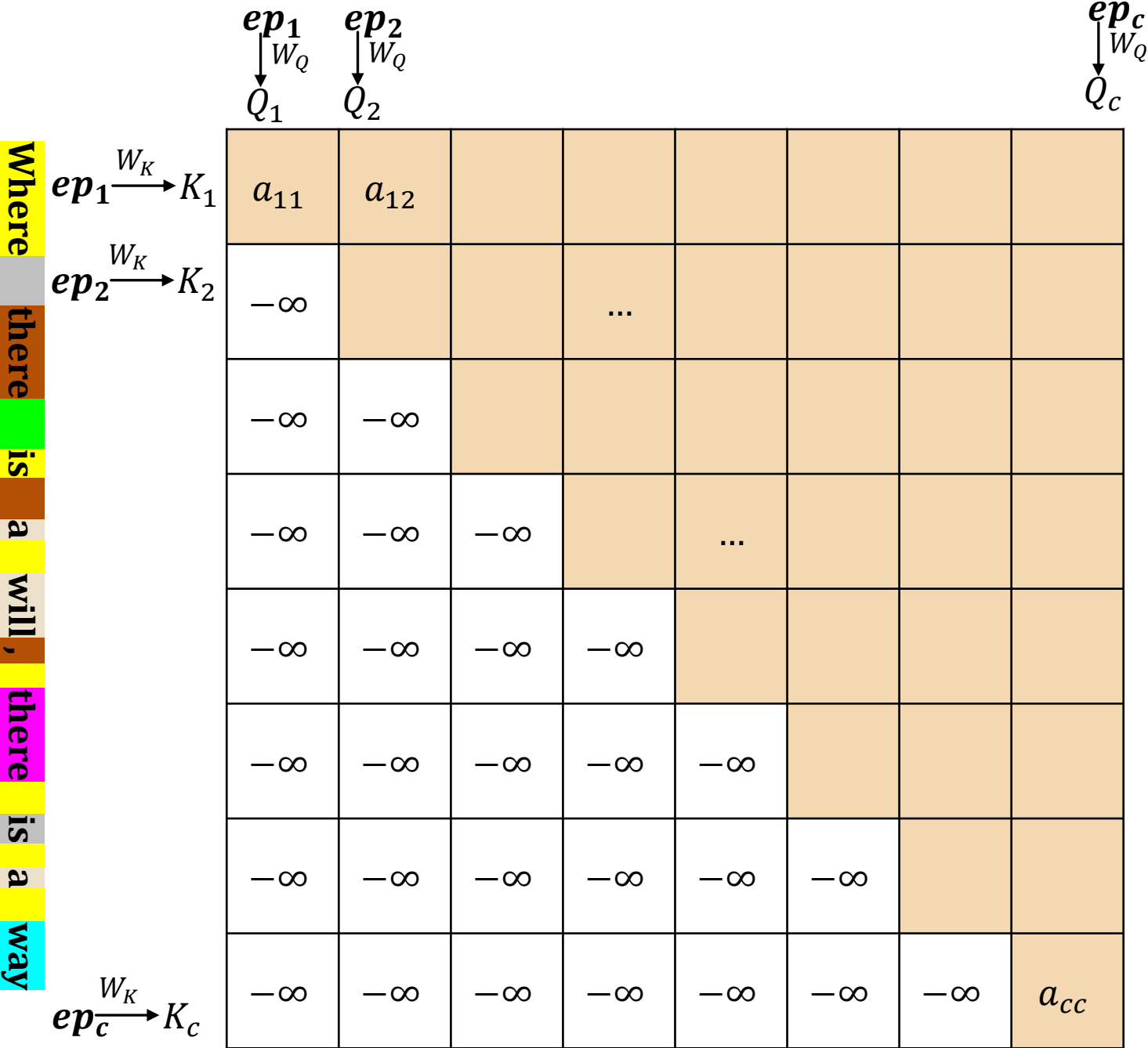


We may want that a query (a token) is influenced by only past keys (past tokens)

Q_1 is influenced by only K_1
 Q_2 is influenced by only K_1, K_2
 Q_3 is influenced by only K_1, K_2, K_3 , and so on

This is equivalent to masking the white region

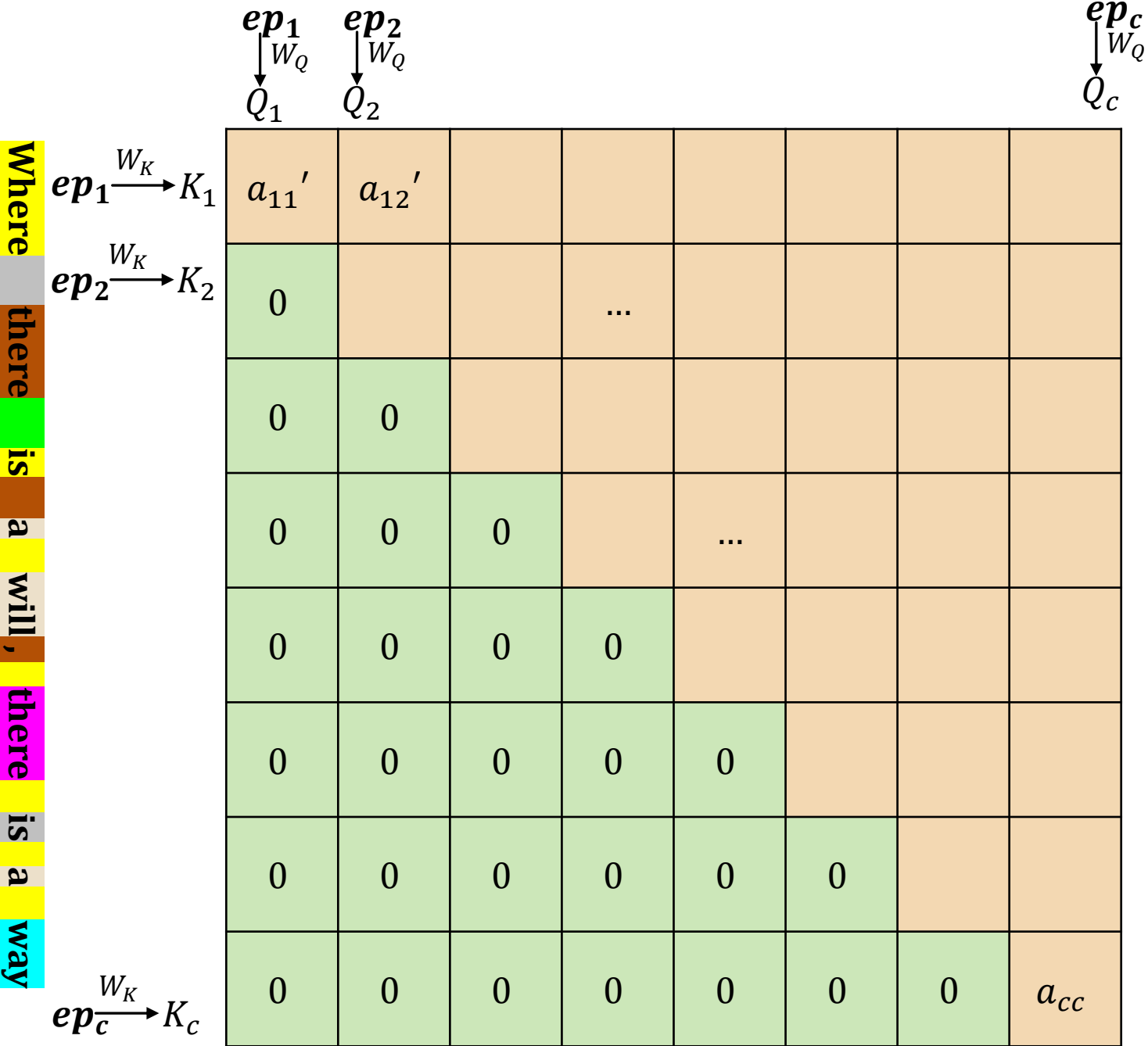
Where there is a will, there is a way



This is equivalent to masking the white region

For that, we make the entries at the white boxes $-\infty$

Where there is a will, there is a way



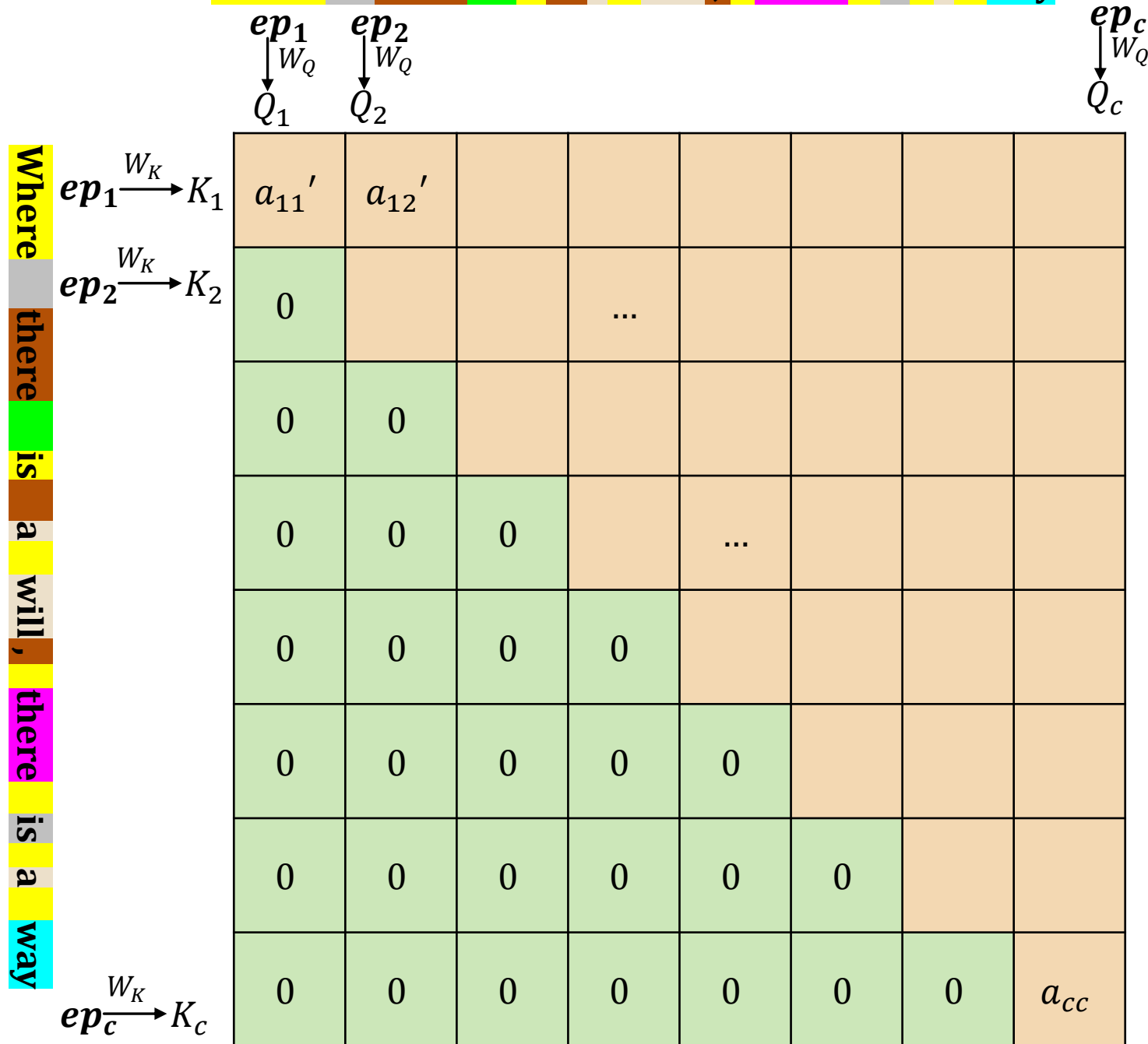
This is equivalent to masking the white region

For that, we make the entries at the white boxes $-\infty$

After column-wise softmax, these entries become 0

$$\sigma(\vec{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

Where there is a will, there is a way



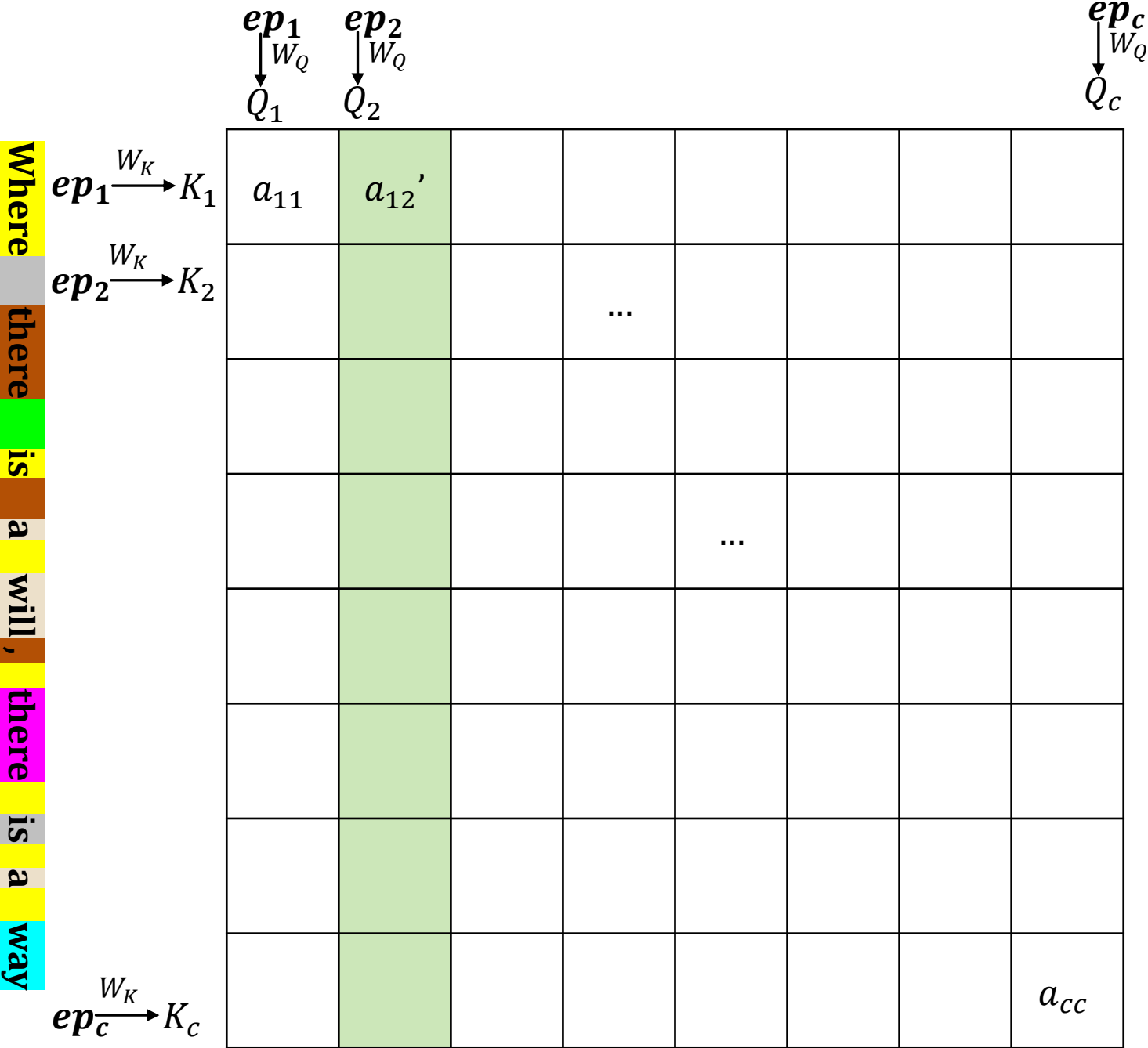
This is equivalent to masking the white region

For that, we make the entries at the white boxes $-\infty$

After column-wise softmax, these entries become 0

This is called **masked attention**

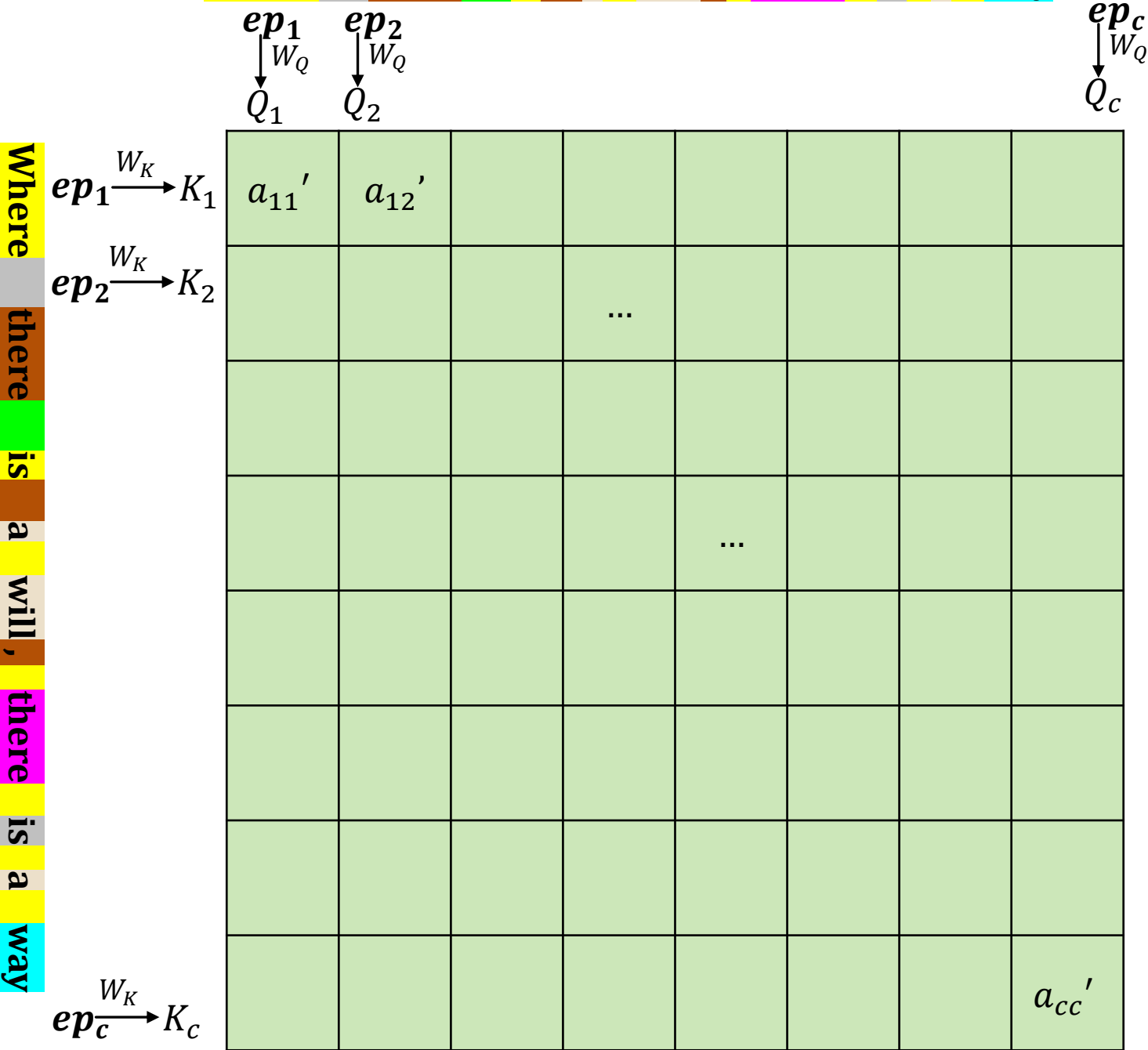
Where there is a will, there is a way



Attention matrix

For each query column, we apply softmax. As a result, each column kind of indicates the probabilities of the corresponding query to match with different keys

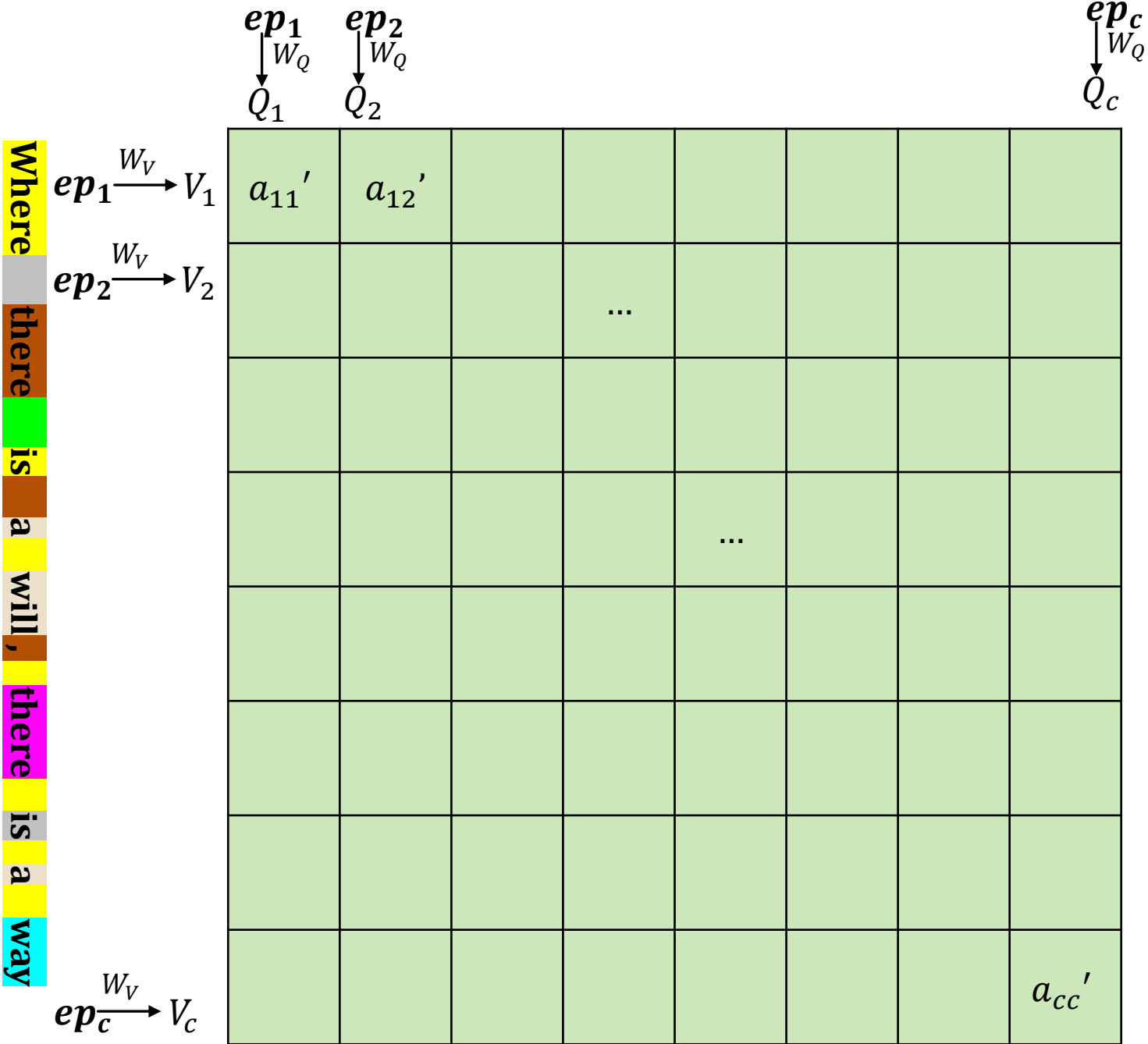
Where there is a will, there is a way



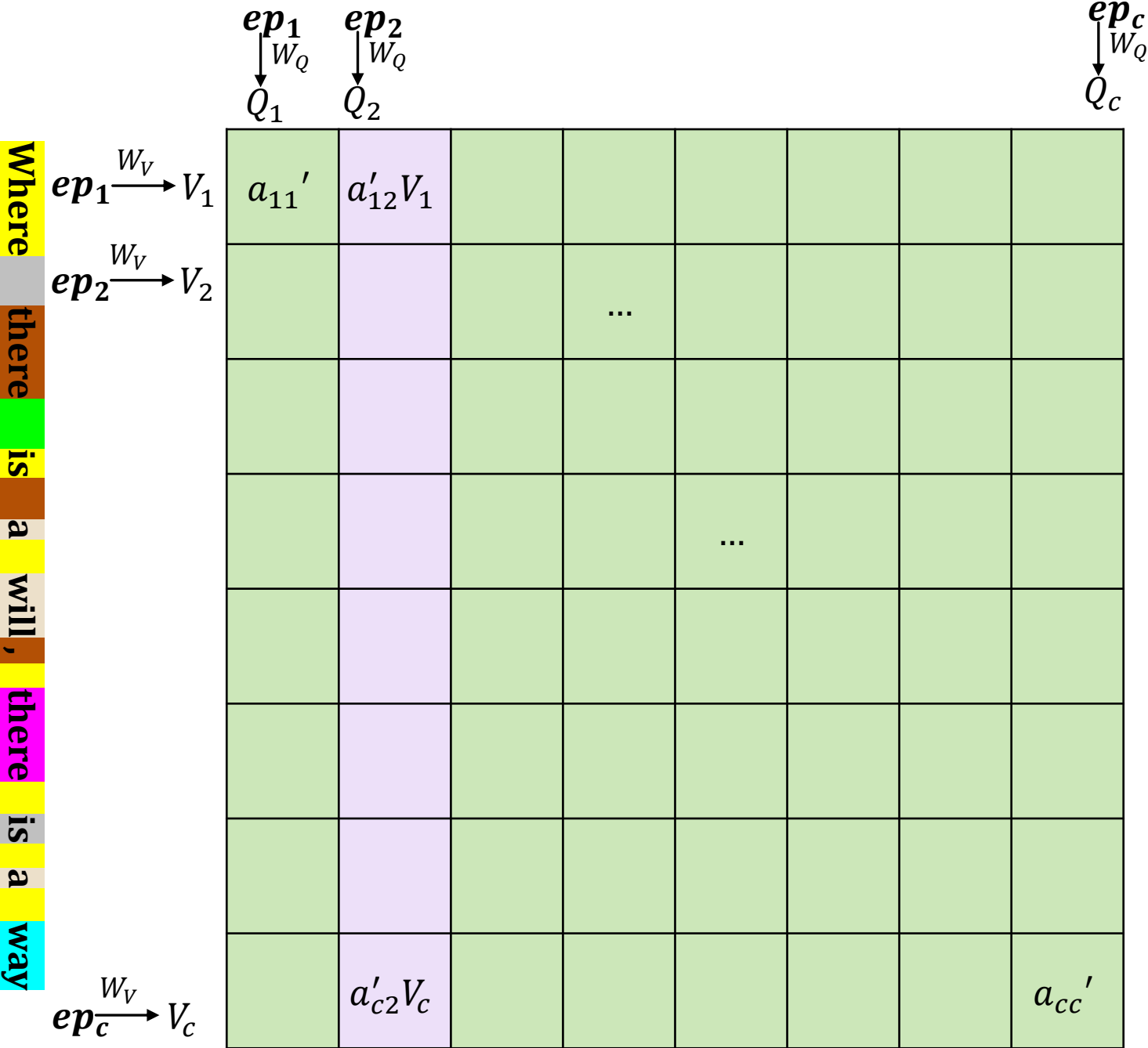
Attention matrix

After applying softmax to every column

Where there is a will, there is a way



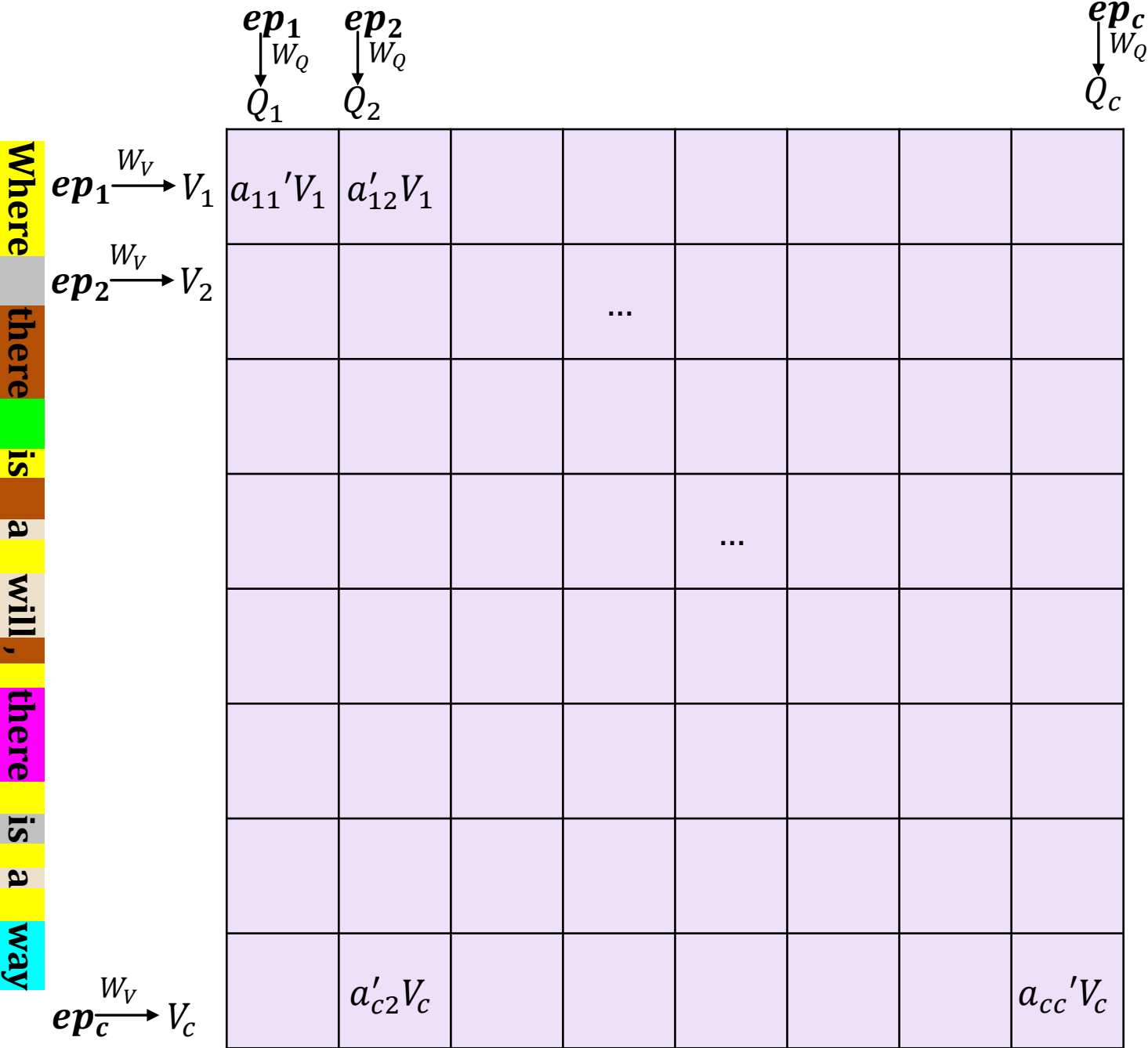
Where there is a will, there is a way



Multiply Value vectors
with every column

Consider column 2

Where there is a will, there is a way



Multiply Value vectors
with every column

Consider column 2

Similarly, for other
columns

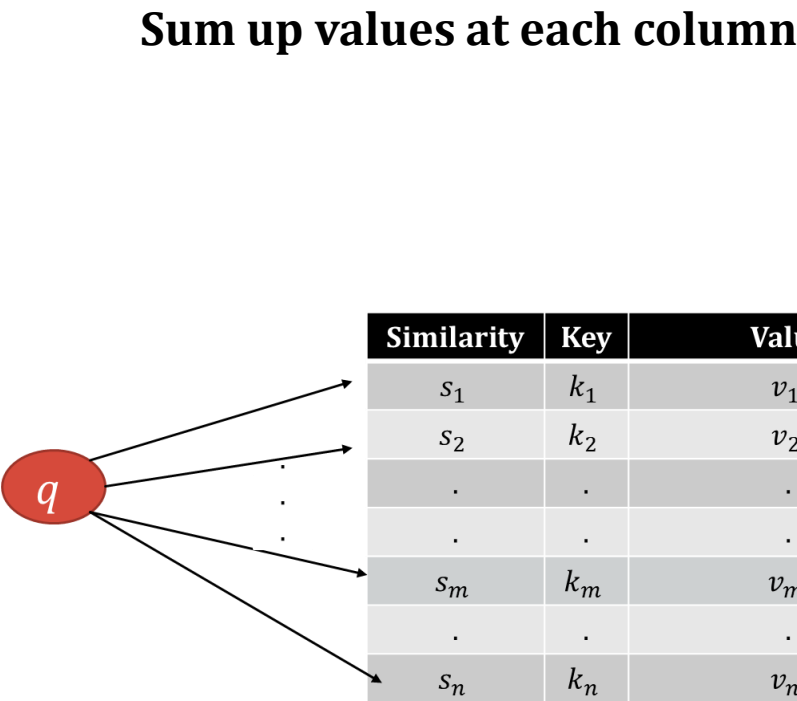
	ep_1	ep_2						ep_c
ep_1	$a_{11}'V_1$	$a_{12}'V_1$						
ep_2				...				
					...			
ep_c		$a_{c2}'V_c$						$a_{cc}'V_c$

**Multiply Value vectors
with every column**

Consider column 2

**Similarly, for other
columns**

	ep_1	ep_2						ep_c
ep_1	$a_{11}'V_1$	$a_{12}'V_1$						
ep_2				...				
					...			
ep_c		$a_{c2}'V_c$						$a_{cc}'V_c$
		Δep_2						



Attention value

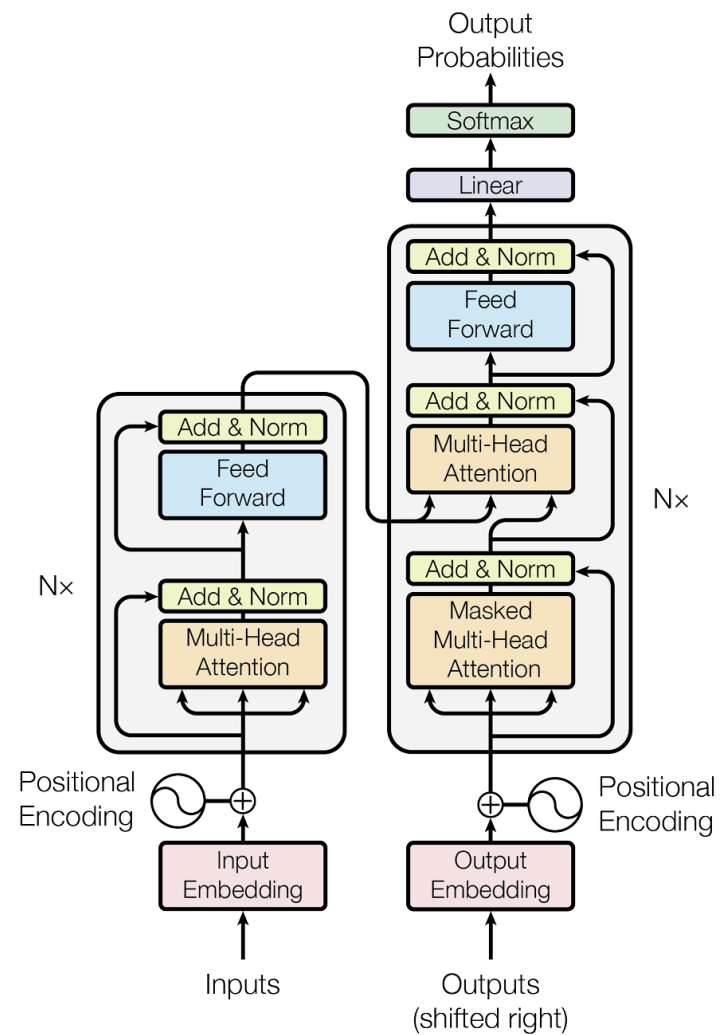
$$= \sum_i softmax(s_i)v_i$$

	ep_1	ep_2						ep_c
ep_1	$a_{11}'V_1$	$a_{12}'V_1$						
ep_2				...				
					...			
ep_c		$a_{c2}'V_c$						$a_{cc}'V_c$
		Δep_2						

Sum up values at each column

Get $ep'_i = ep_i + \Delta ep_i$

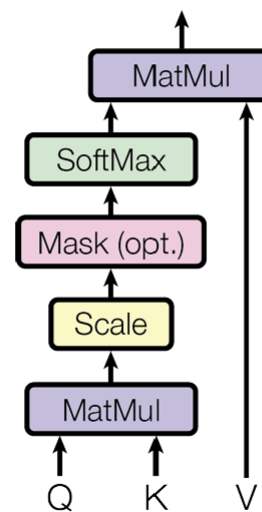
Through residual connection



	ep_1	ep_2						ep_c
ep_1	$a_{11}'V_1$	$a_{12}'V_1$						
ep_2				...				
					...			
ep_c		$a_{c2}'V_c$						$a_{cc}'V_c$
		Δep_2						

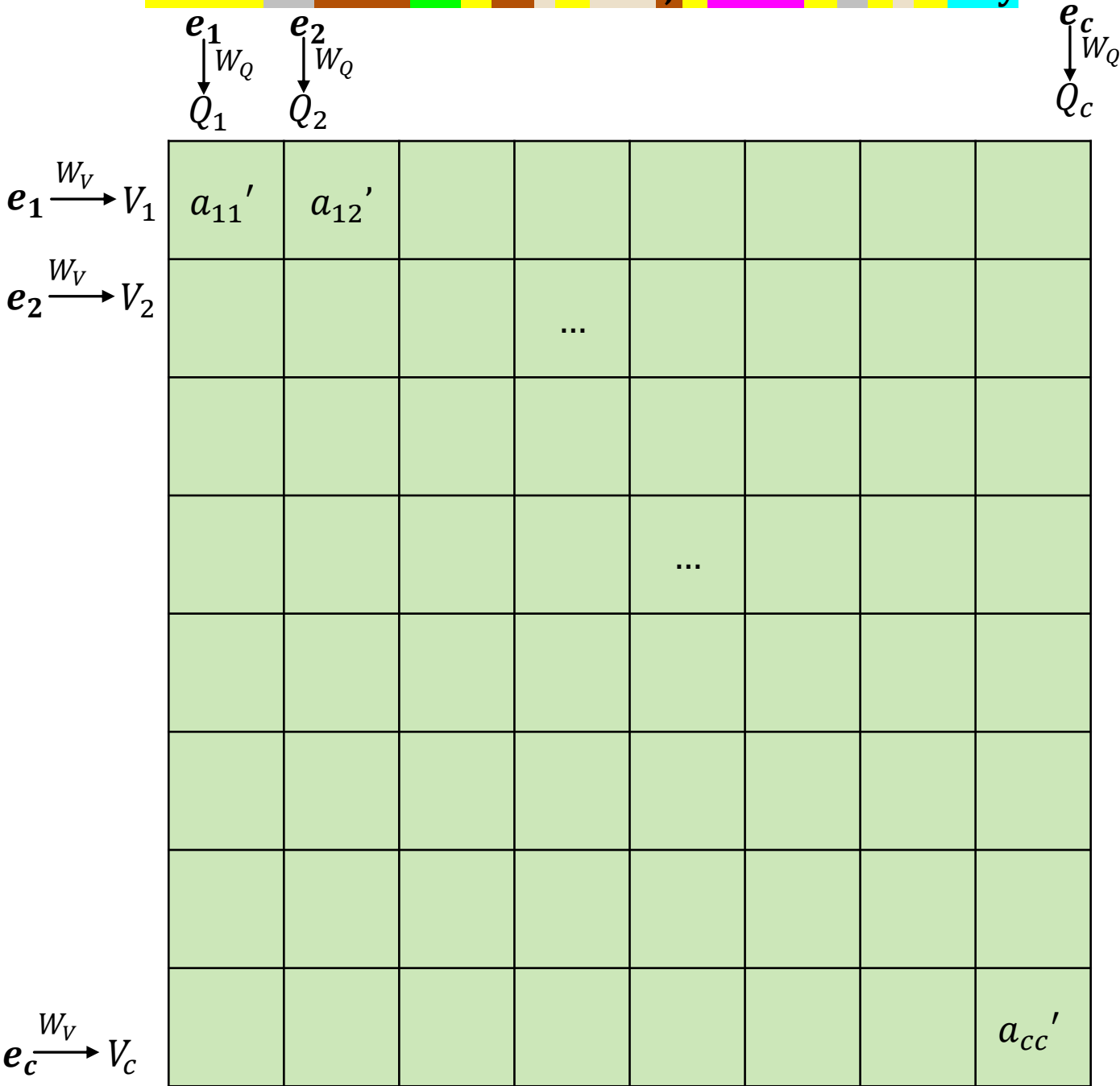
This entire thing is one head of attention

Scaled Dot-Product Attention



Where there is a will, there is a way

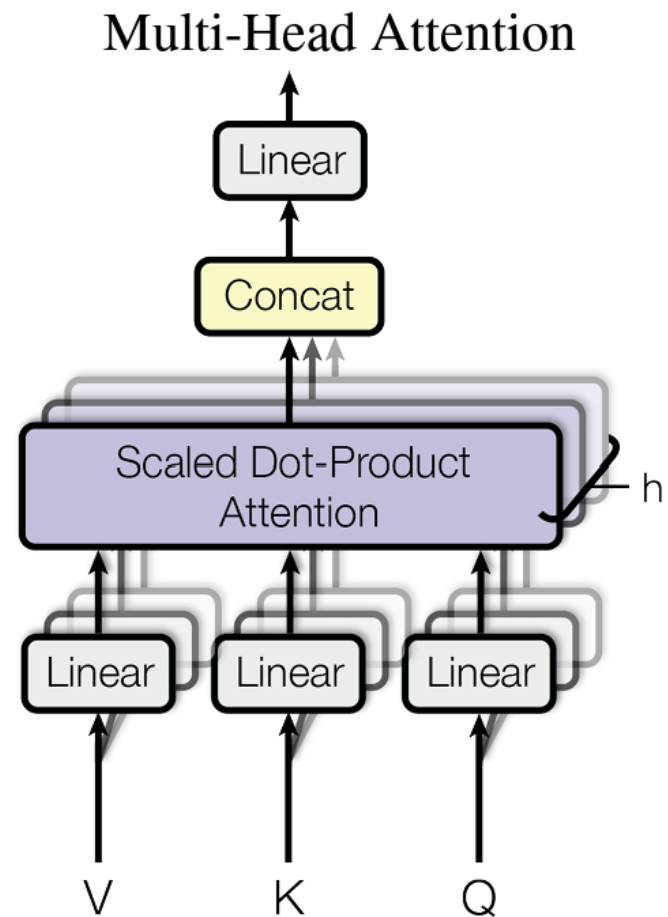
Where there is a will, there is a way



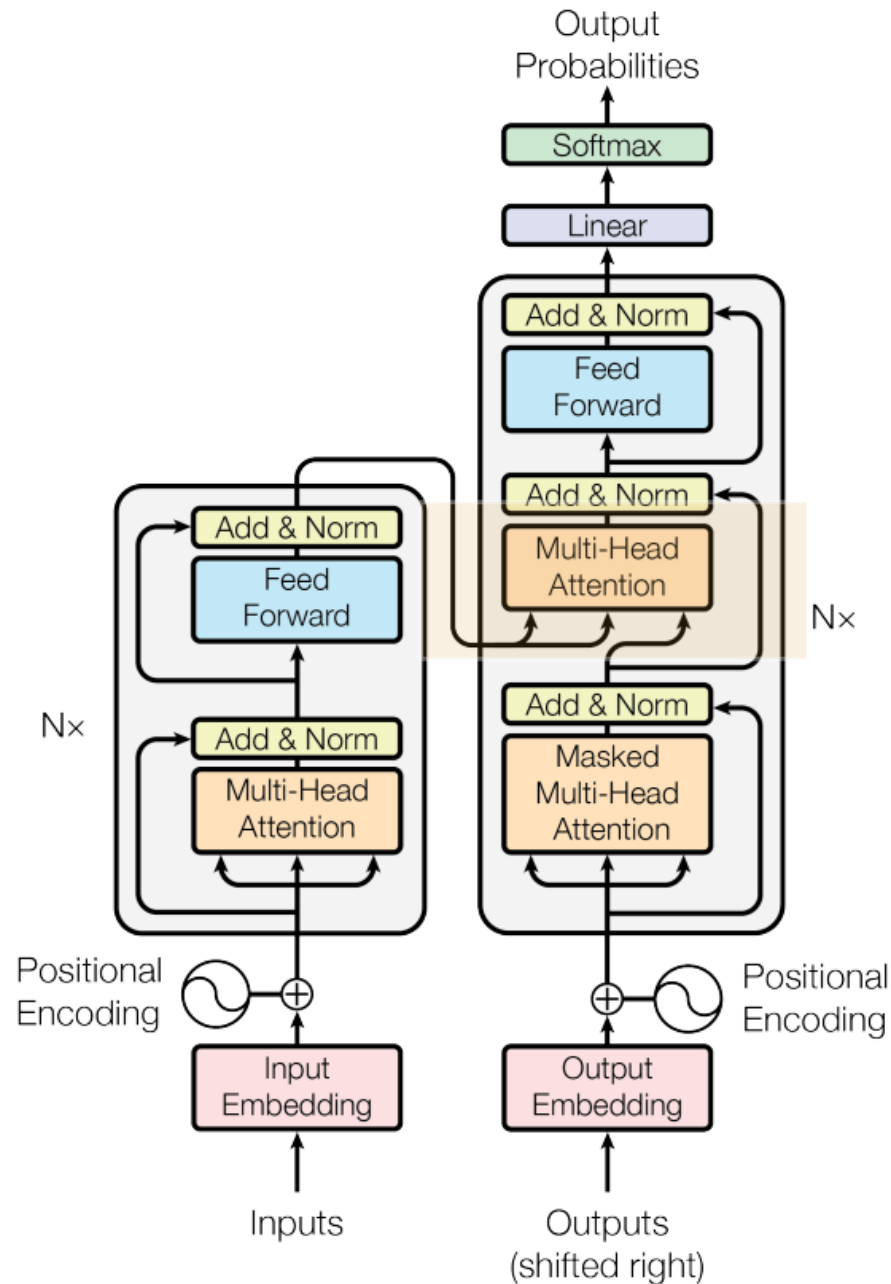
We will have many such heads with different W_V , W_Q , and W_K

	ep_1	ep_2						ep_c
ep_1	$a_{11}'V_1$	$a_{12}'V_1$						
ep_2				...				
					...			
ep_c		$a_{c2}'V_c$						$a_{cc}'V_c$
		Δep_2						

This entire thing is one head of attention



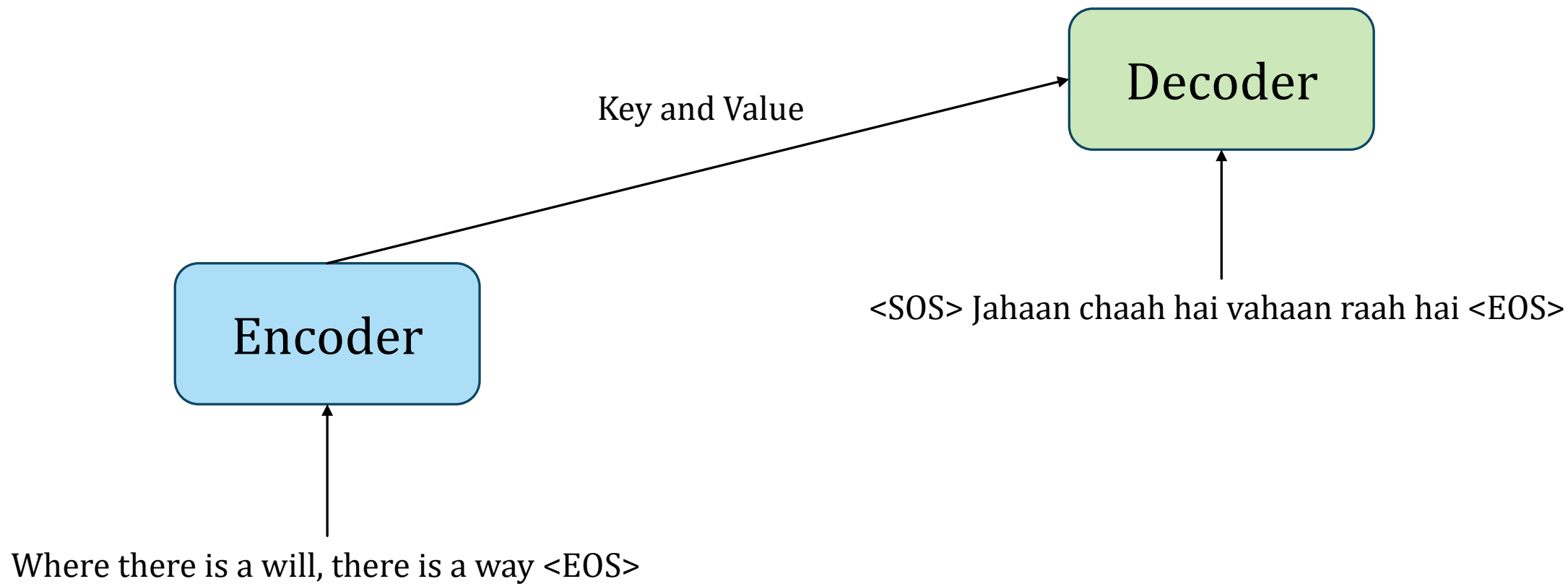
Encoder-decoder Attention



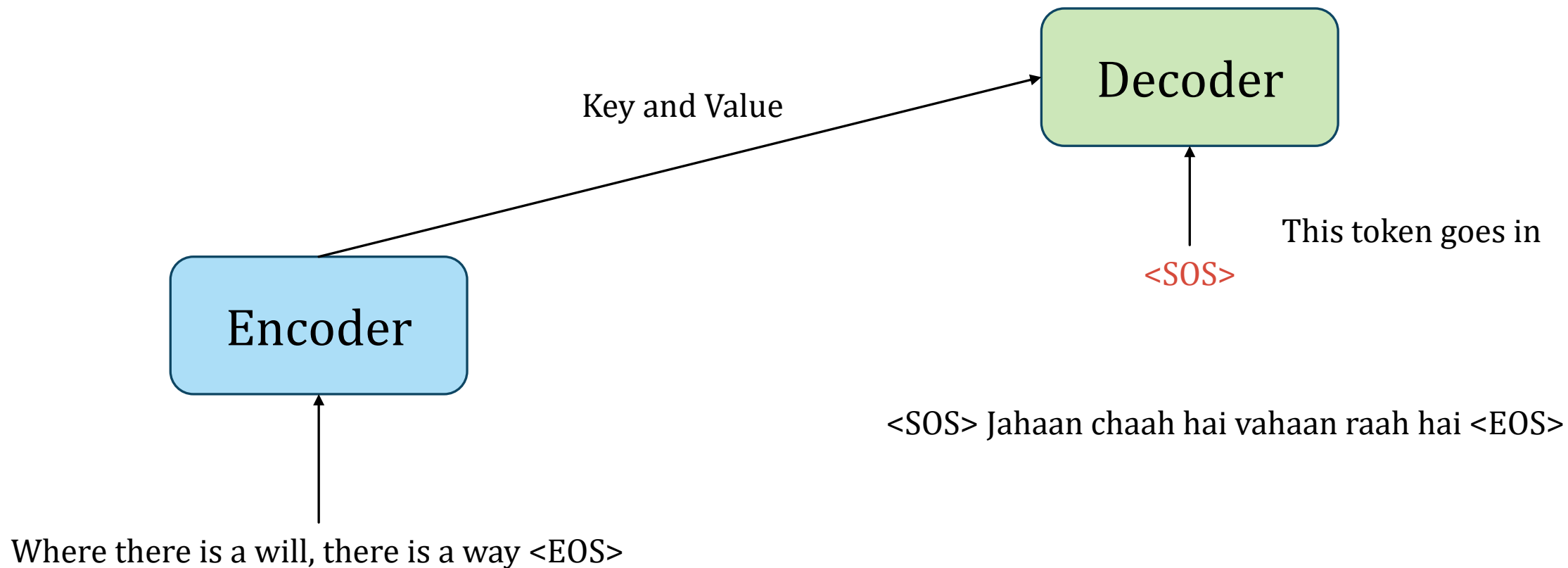
Key and value comes from the encoder through two paths

Query comes from the decoder

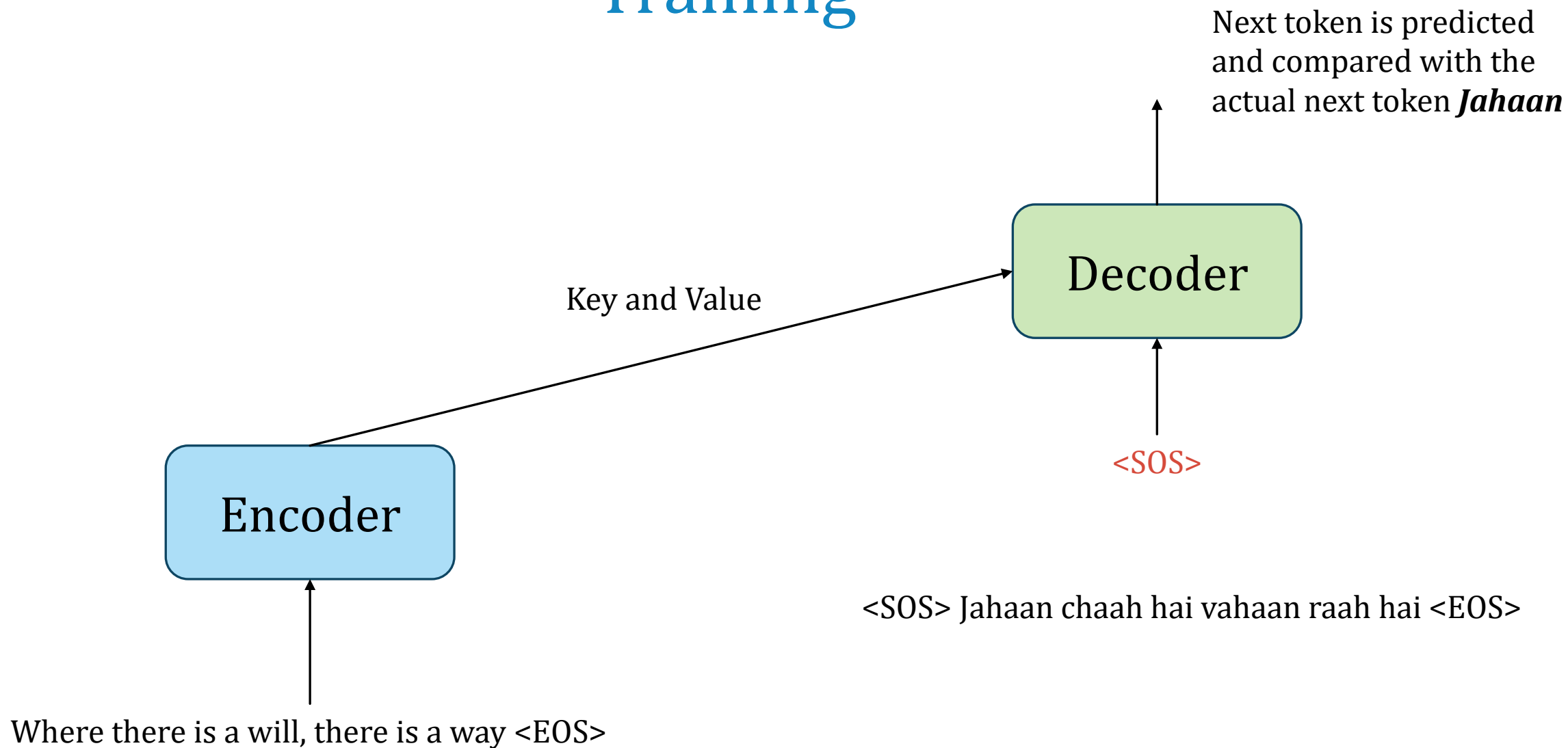
Training



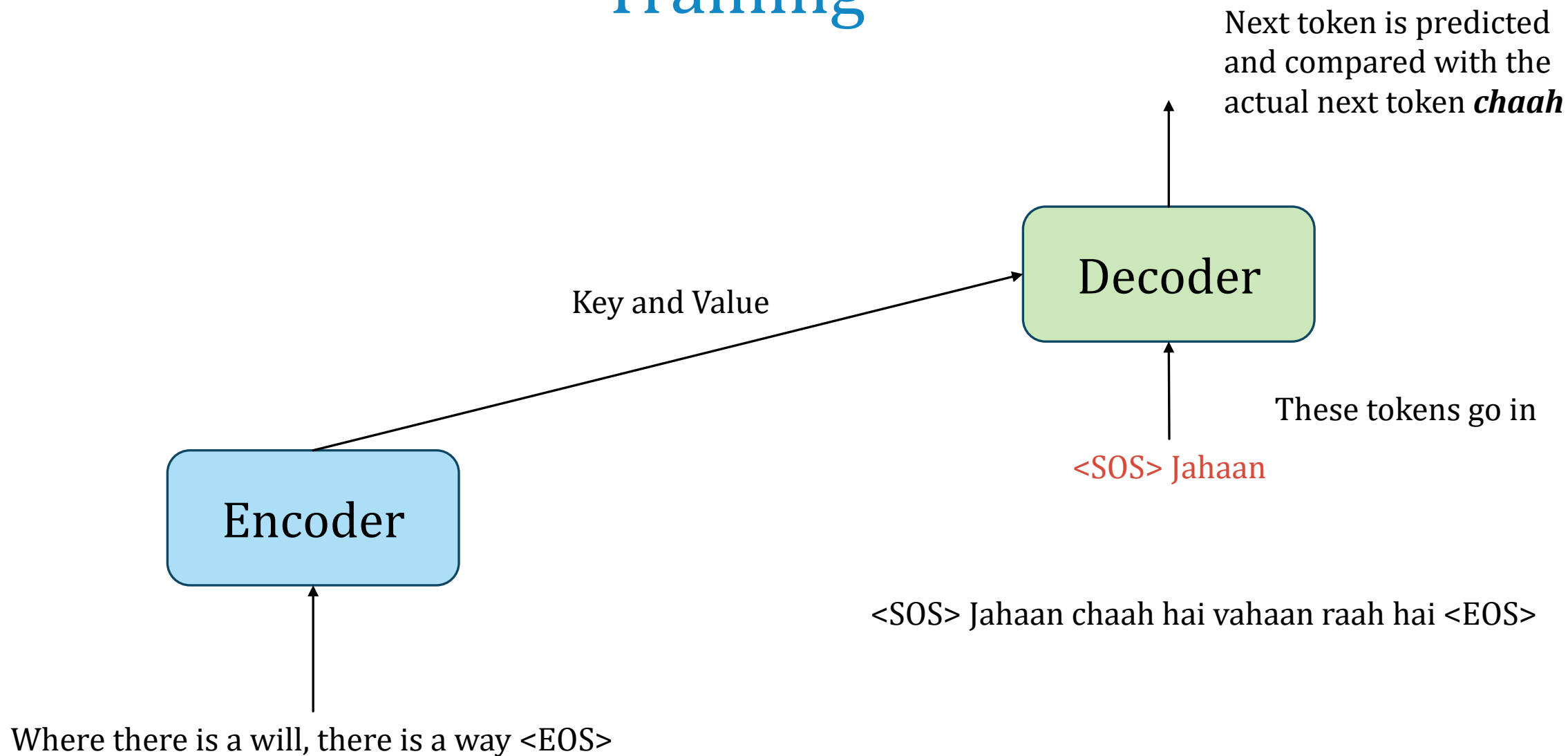
Training



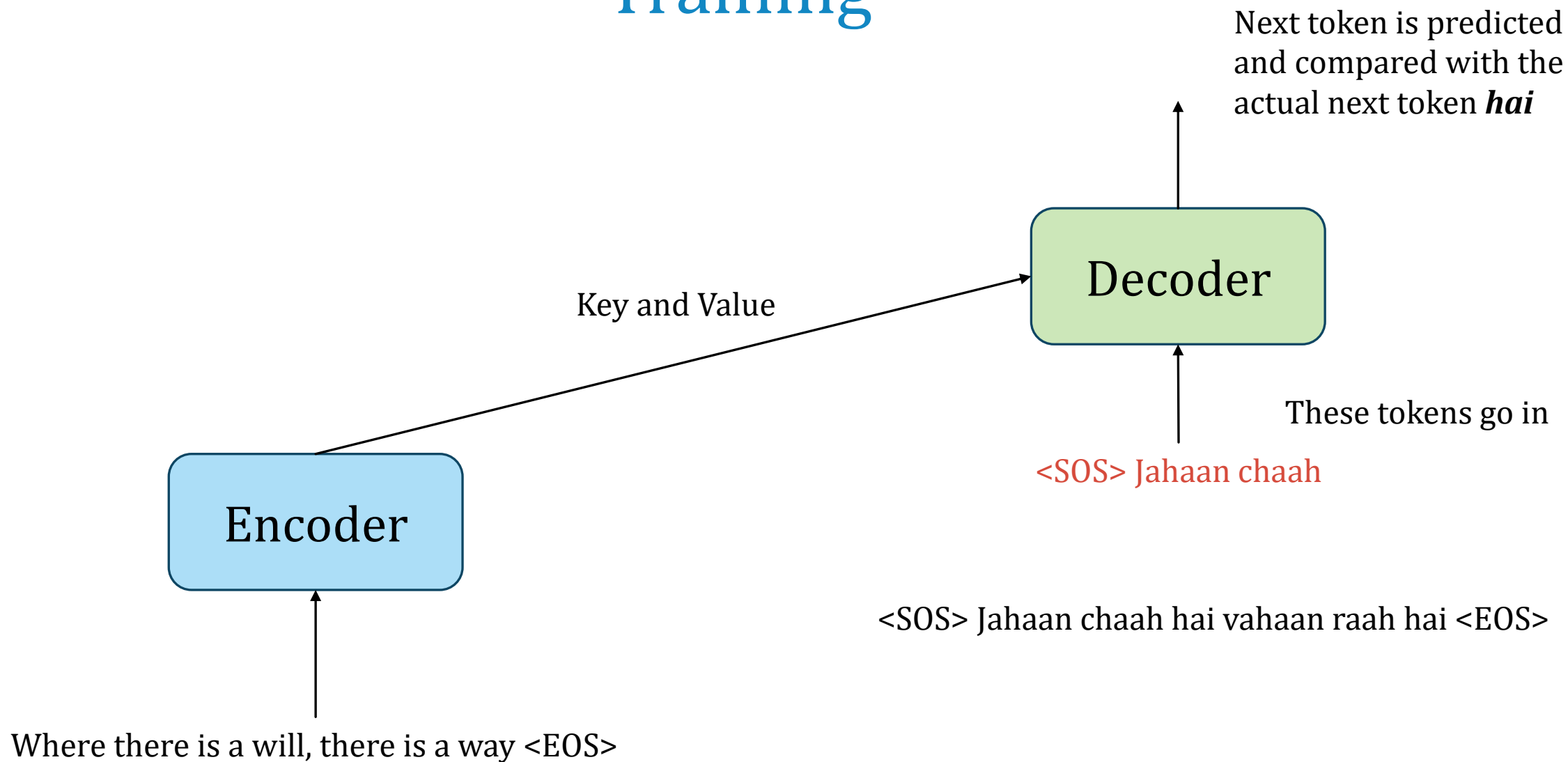
Training



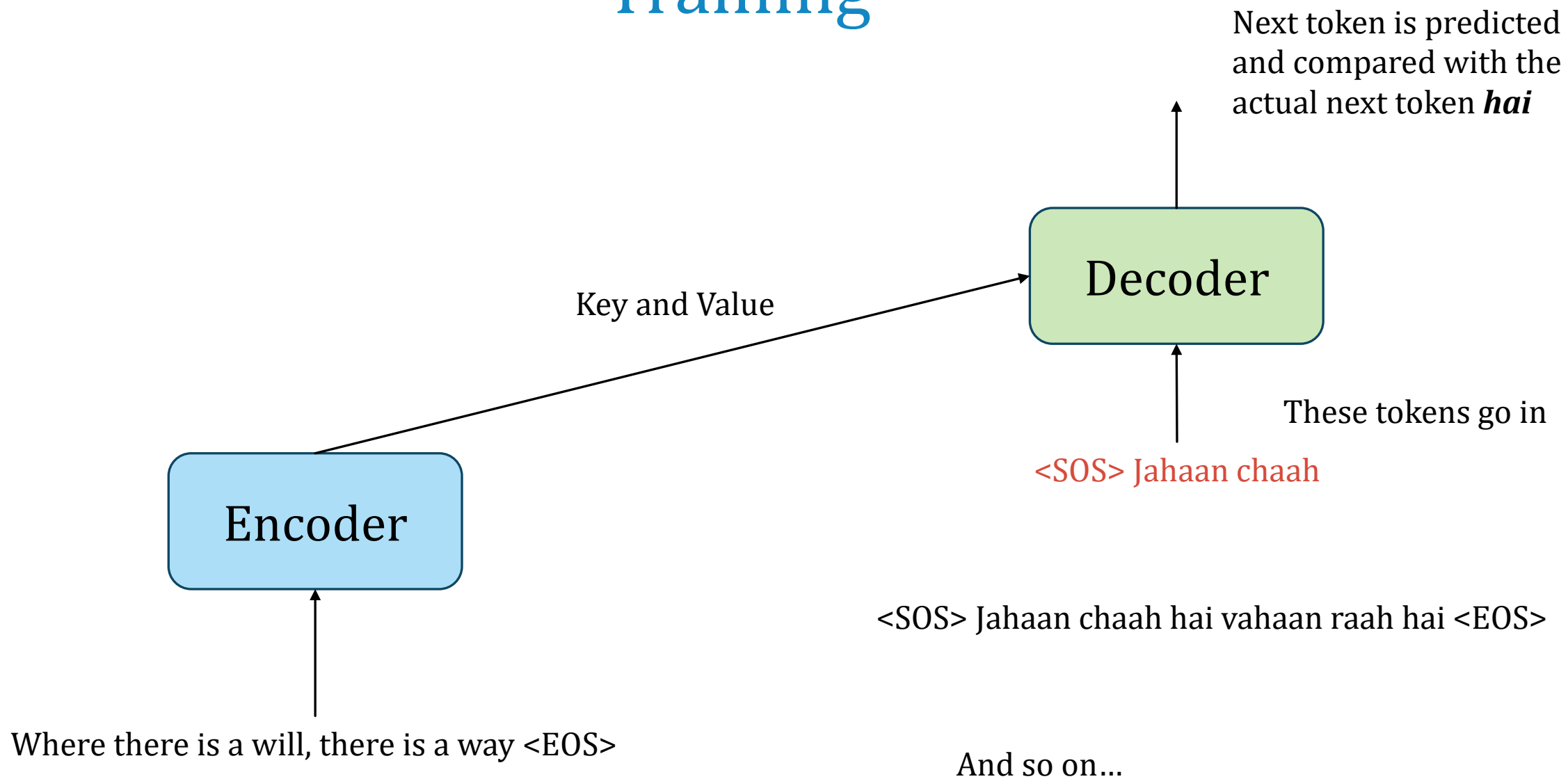
Training



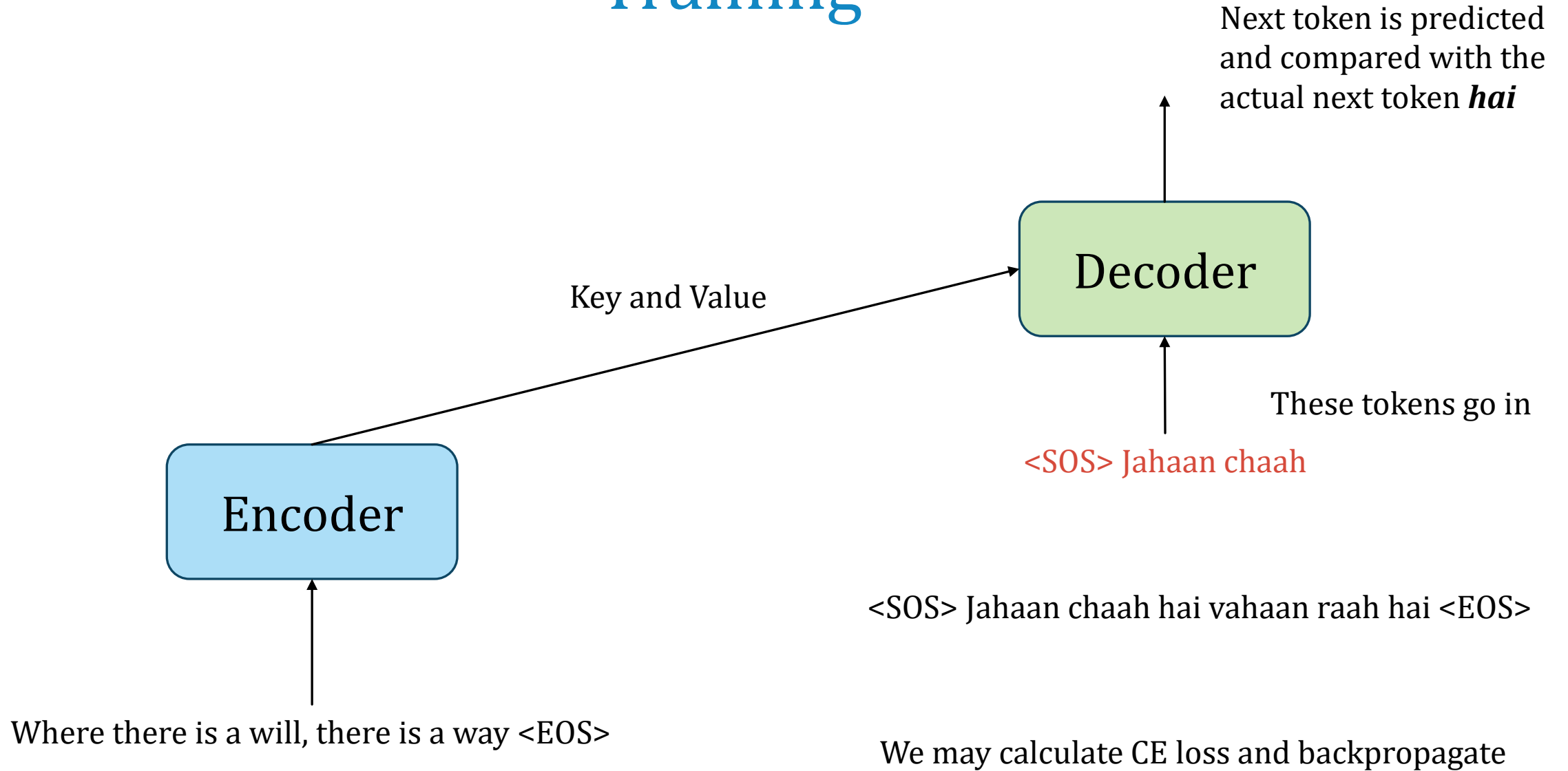
Training



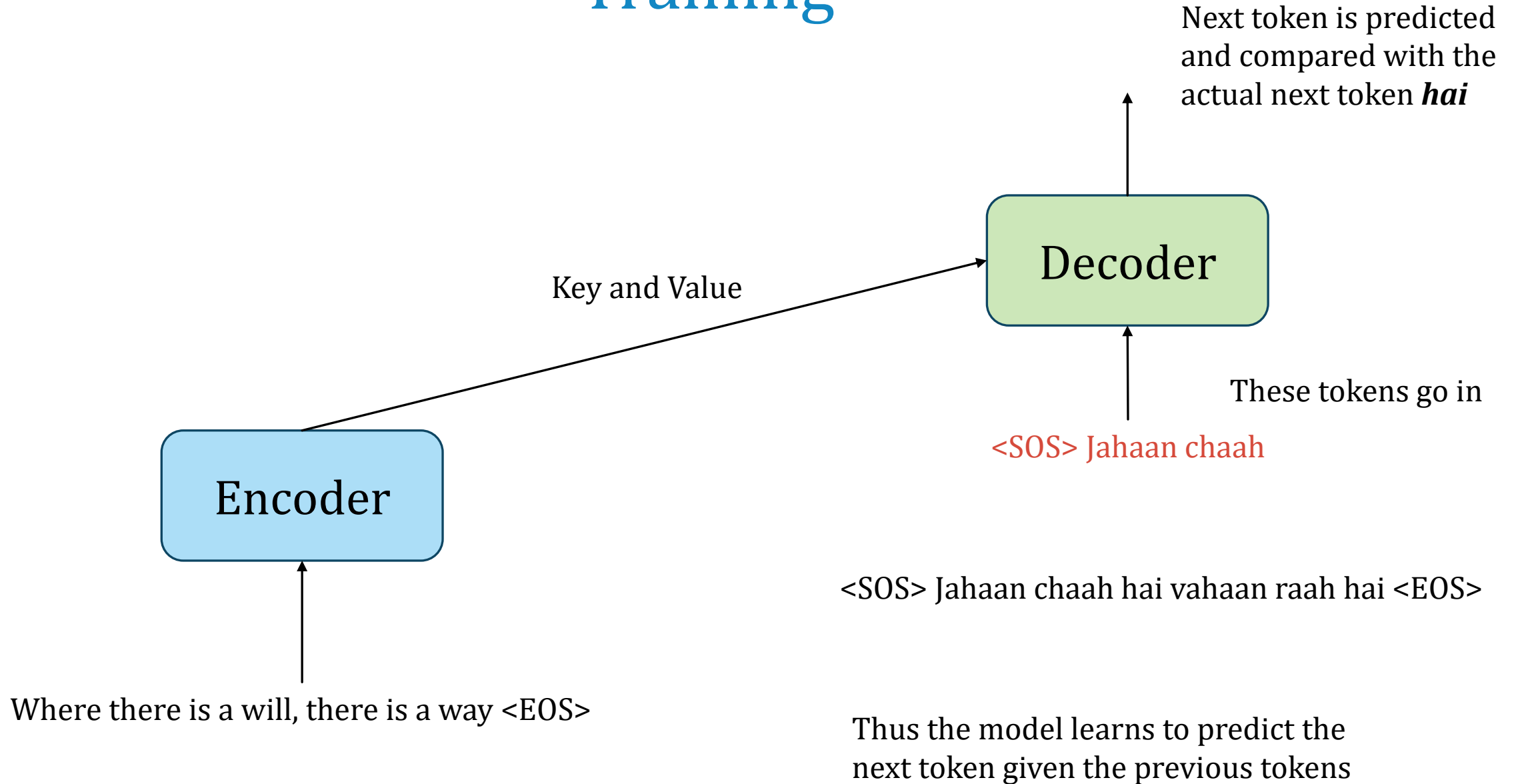
Training



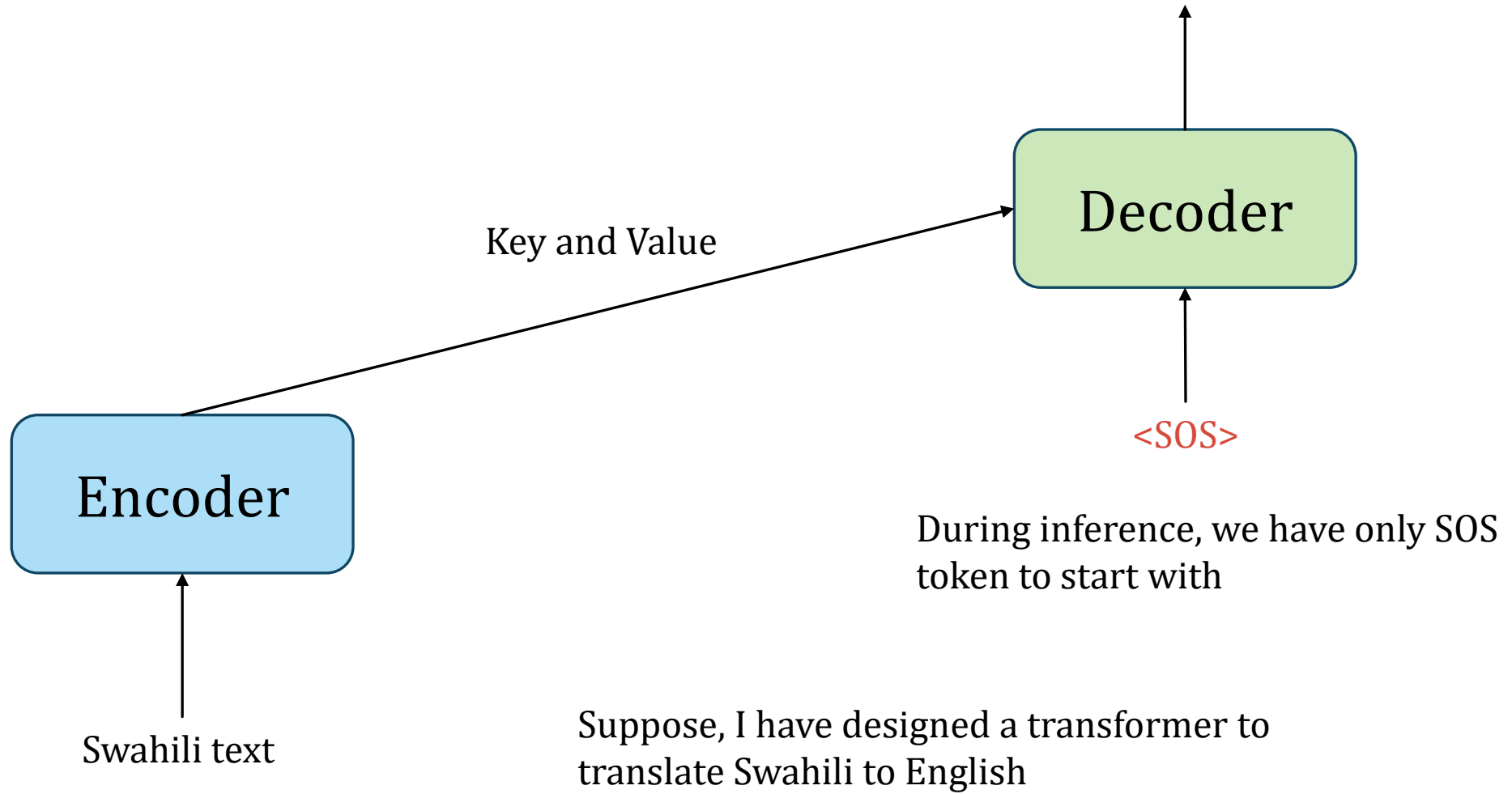
Training



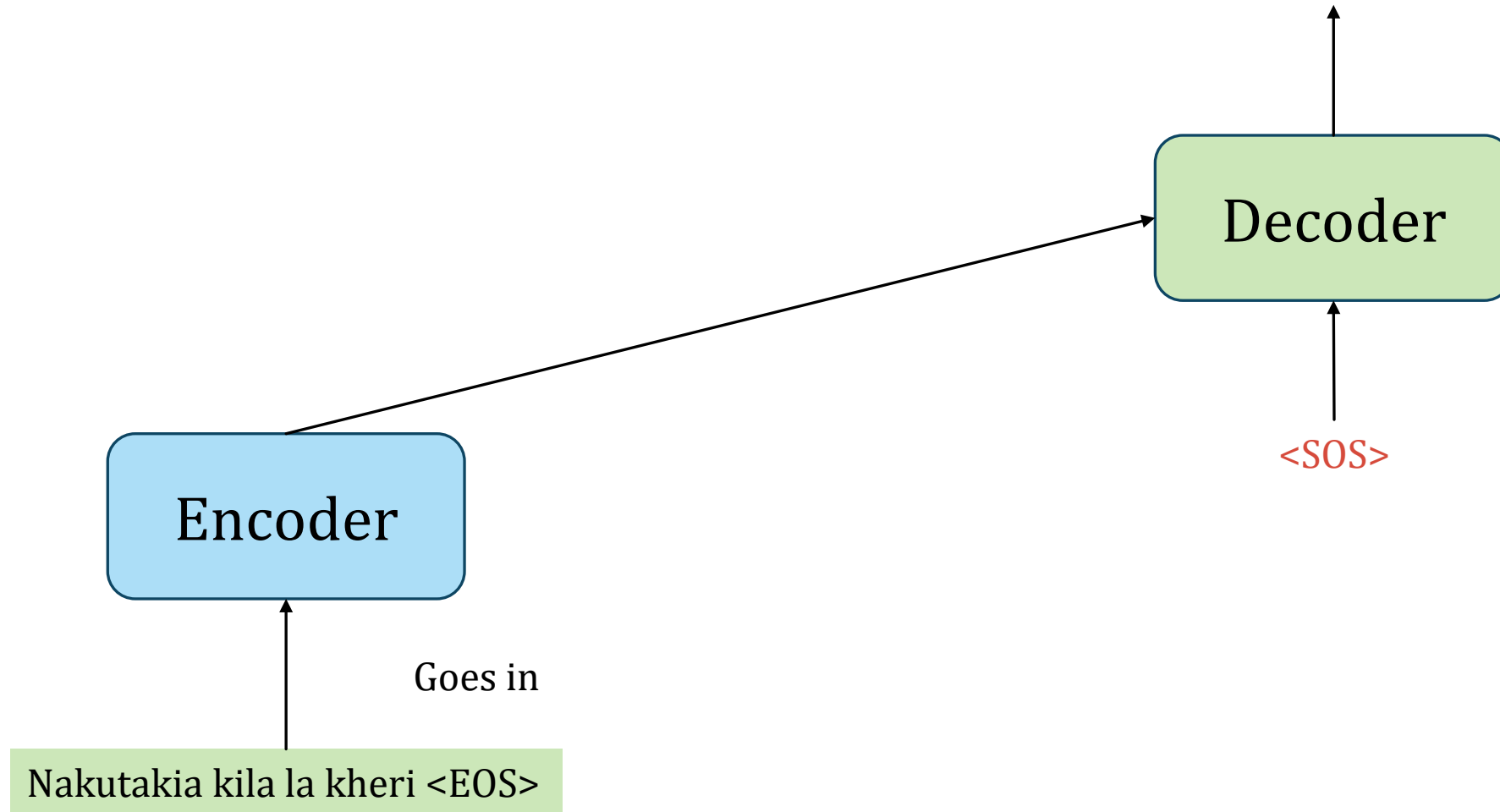
Training



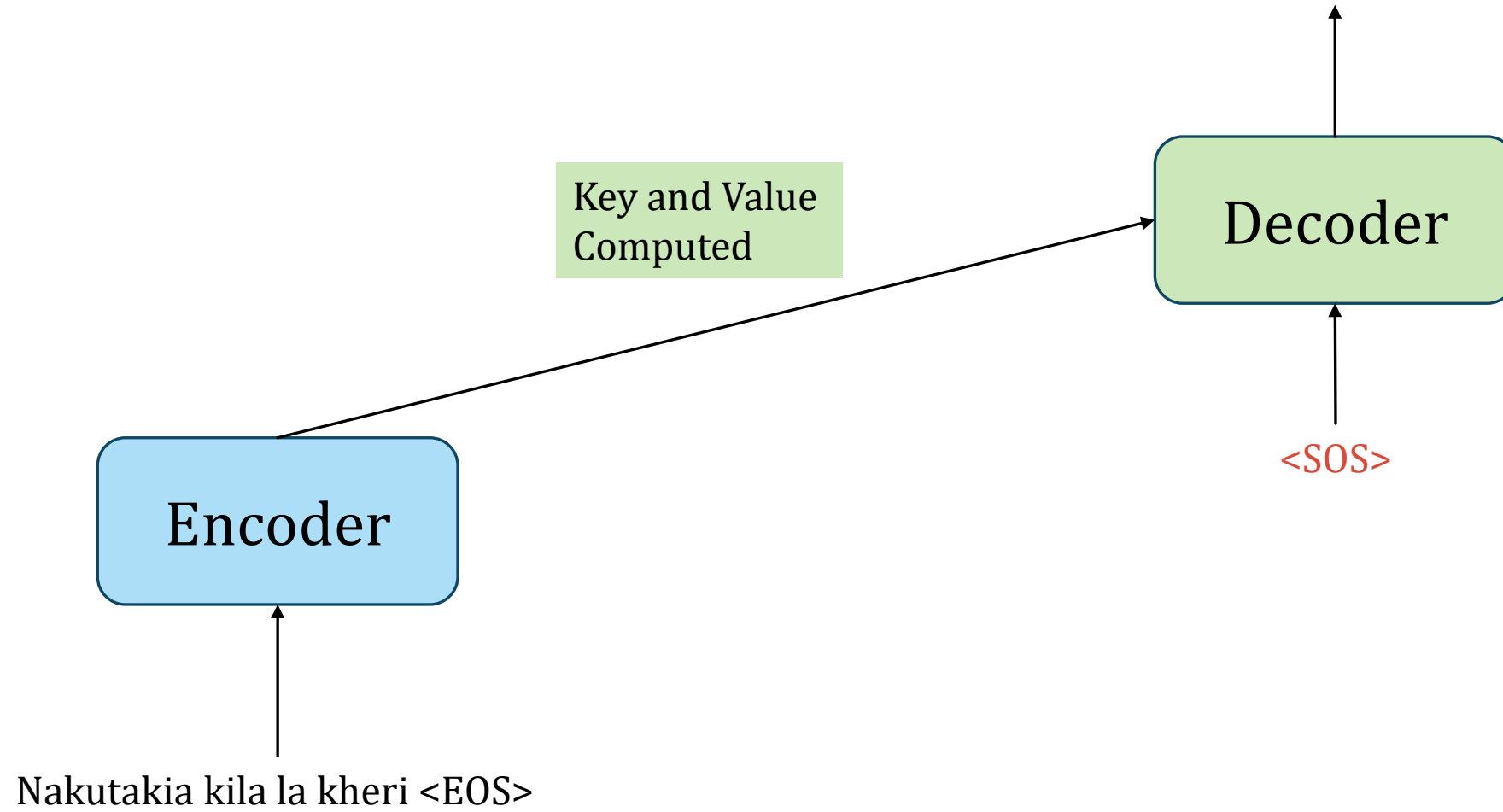
Inference



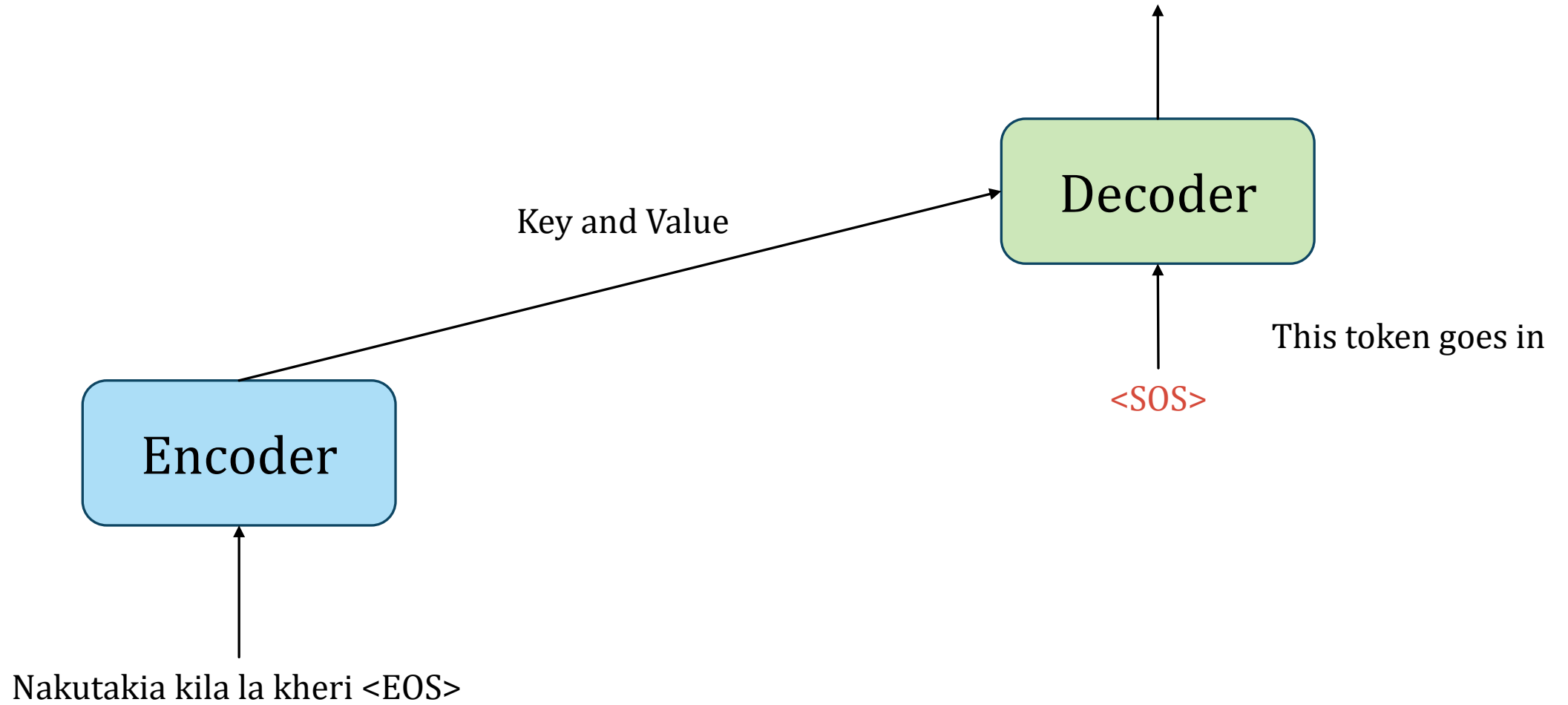
Inference



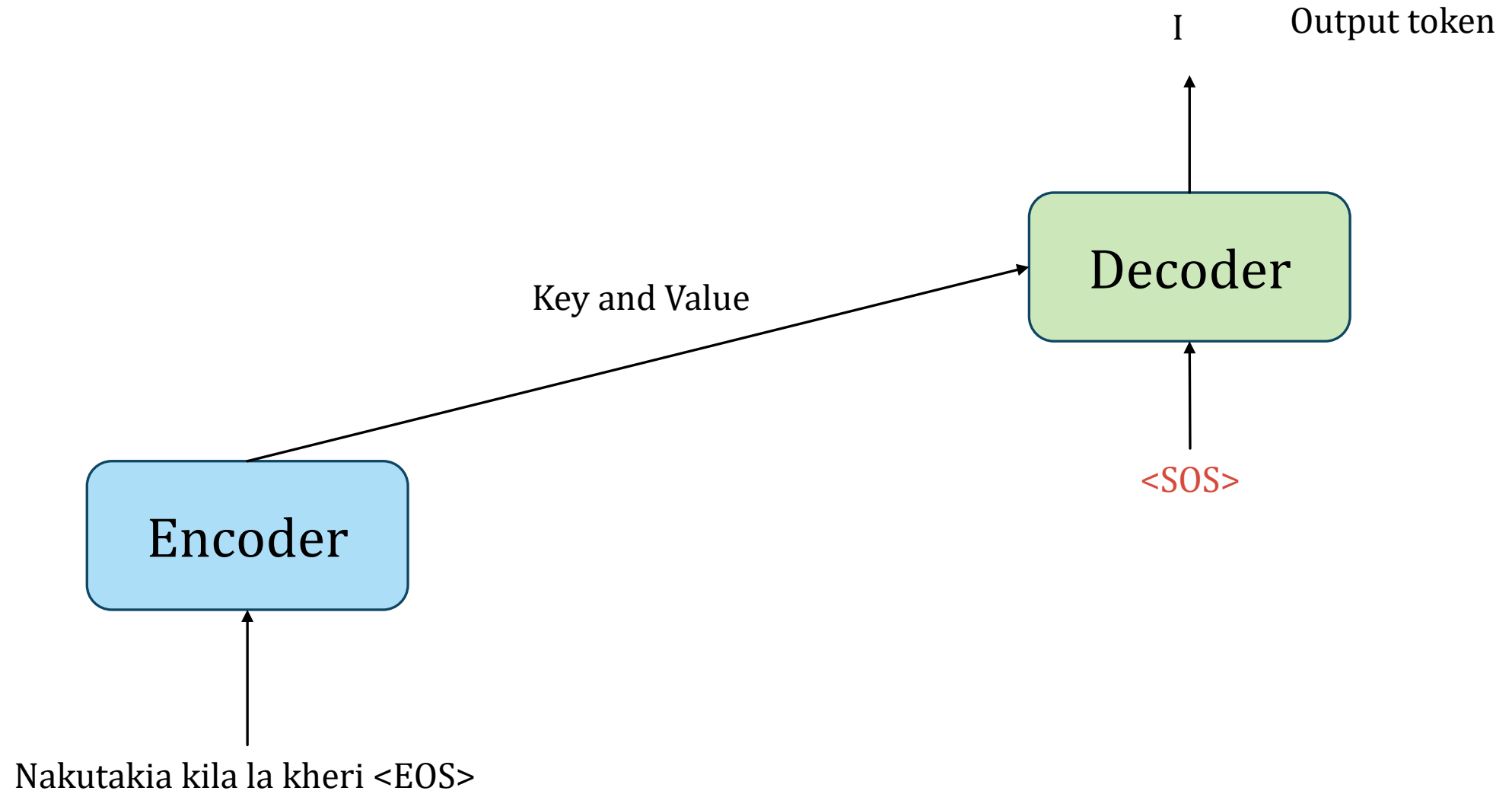
Inference



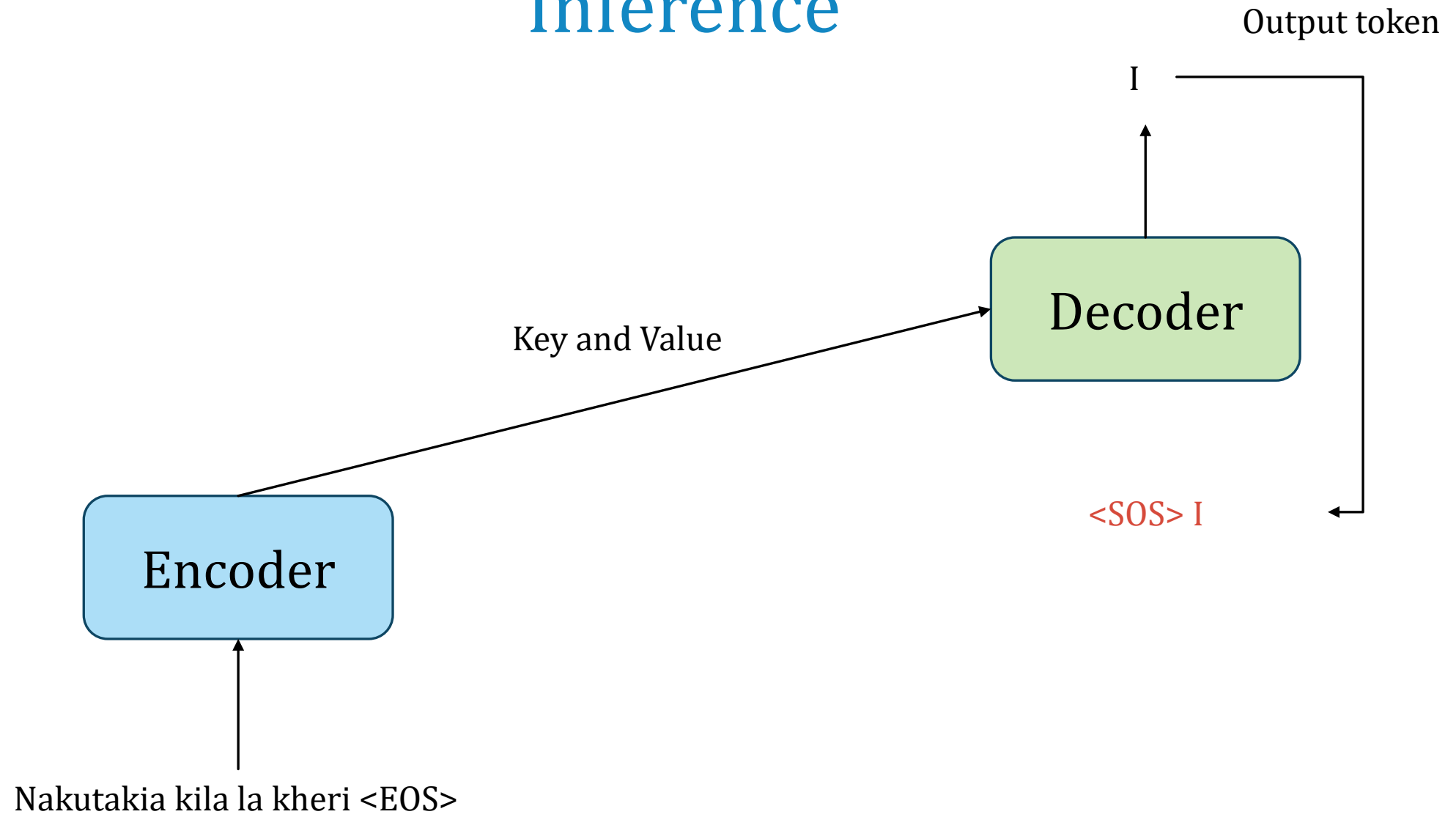
Inference



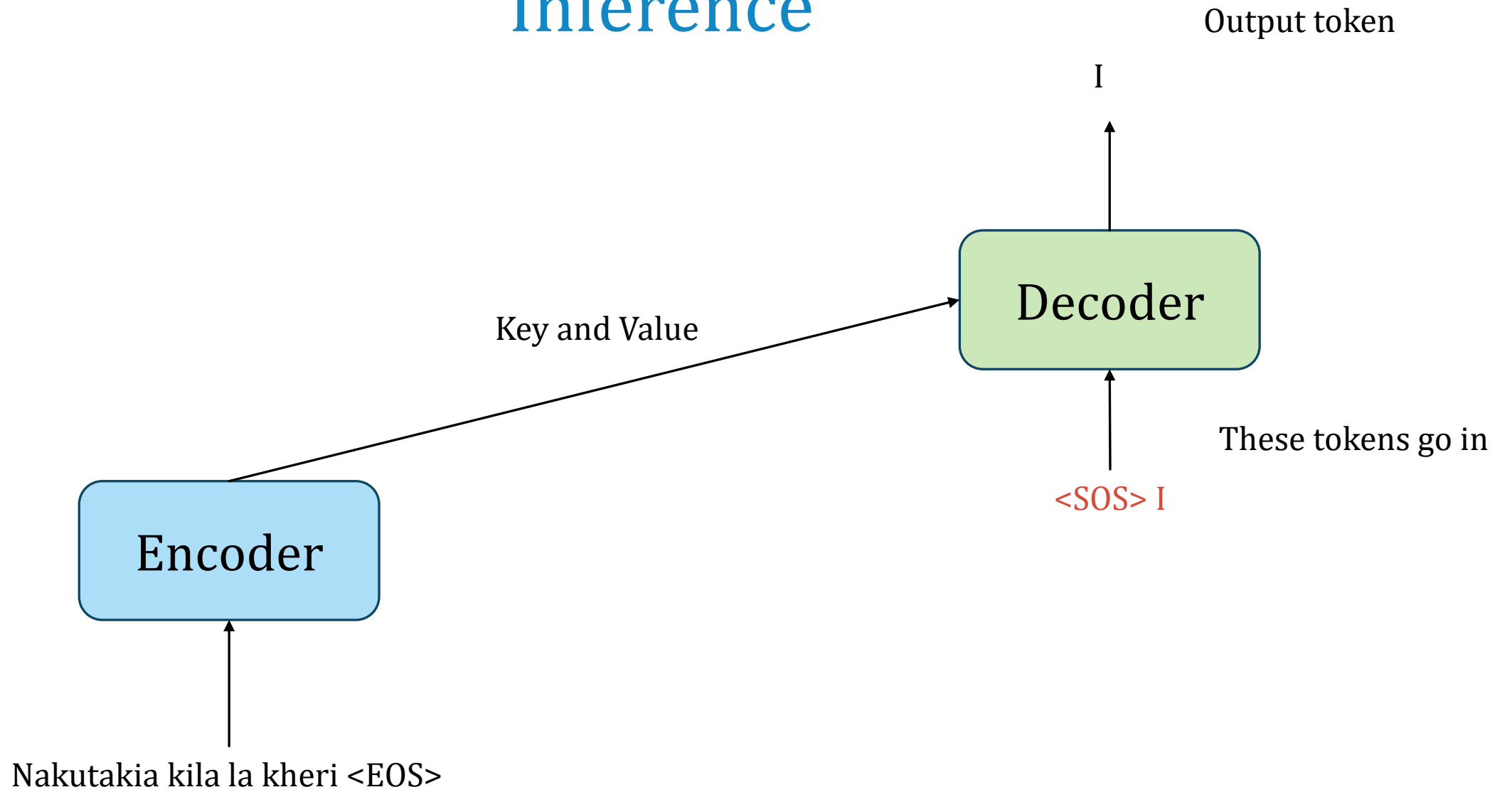
Inference



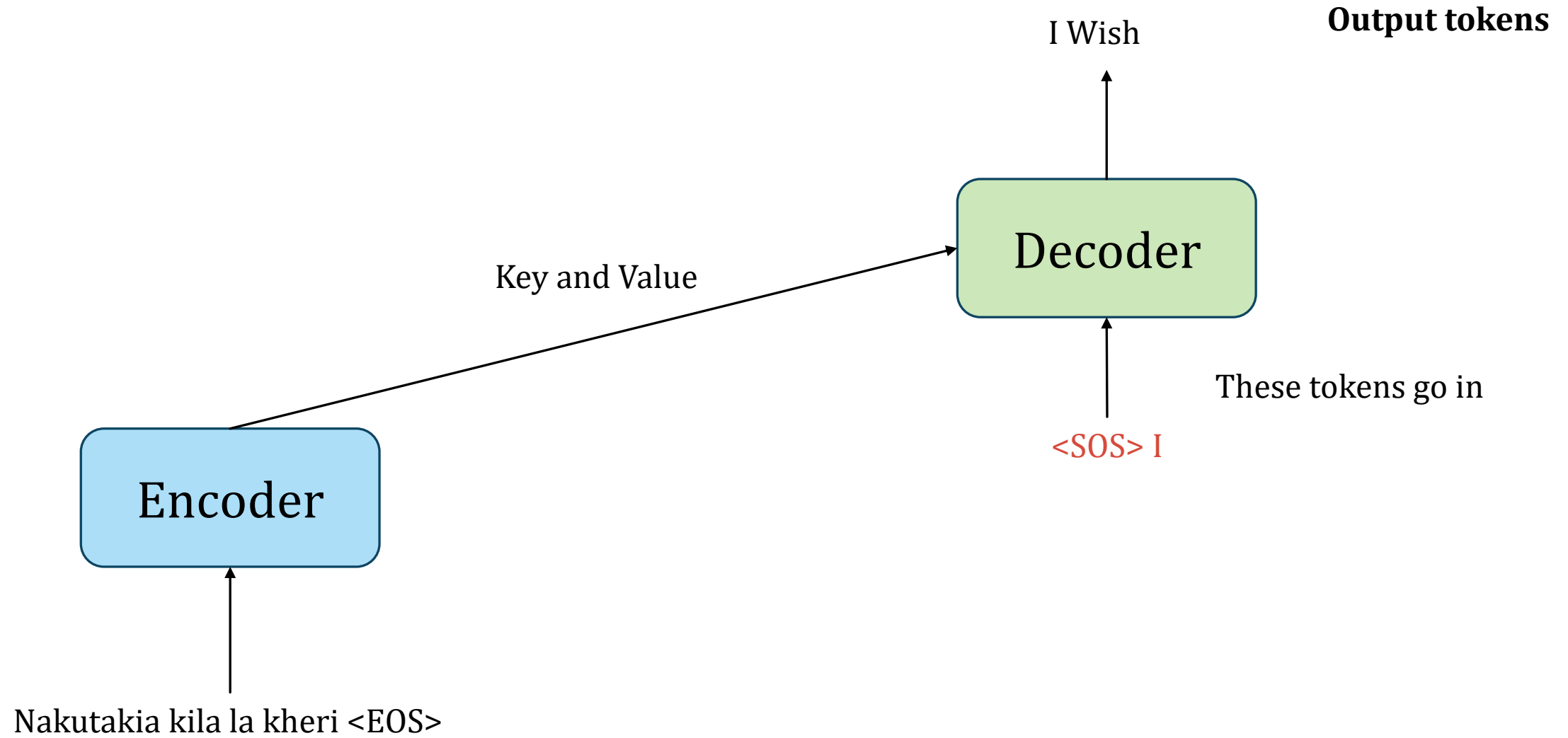
Inference



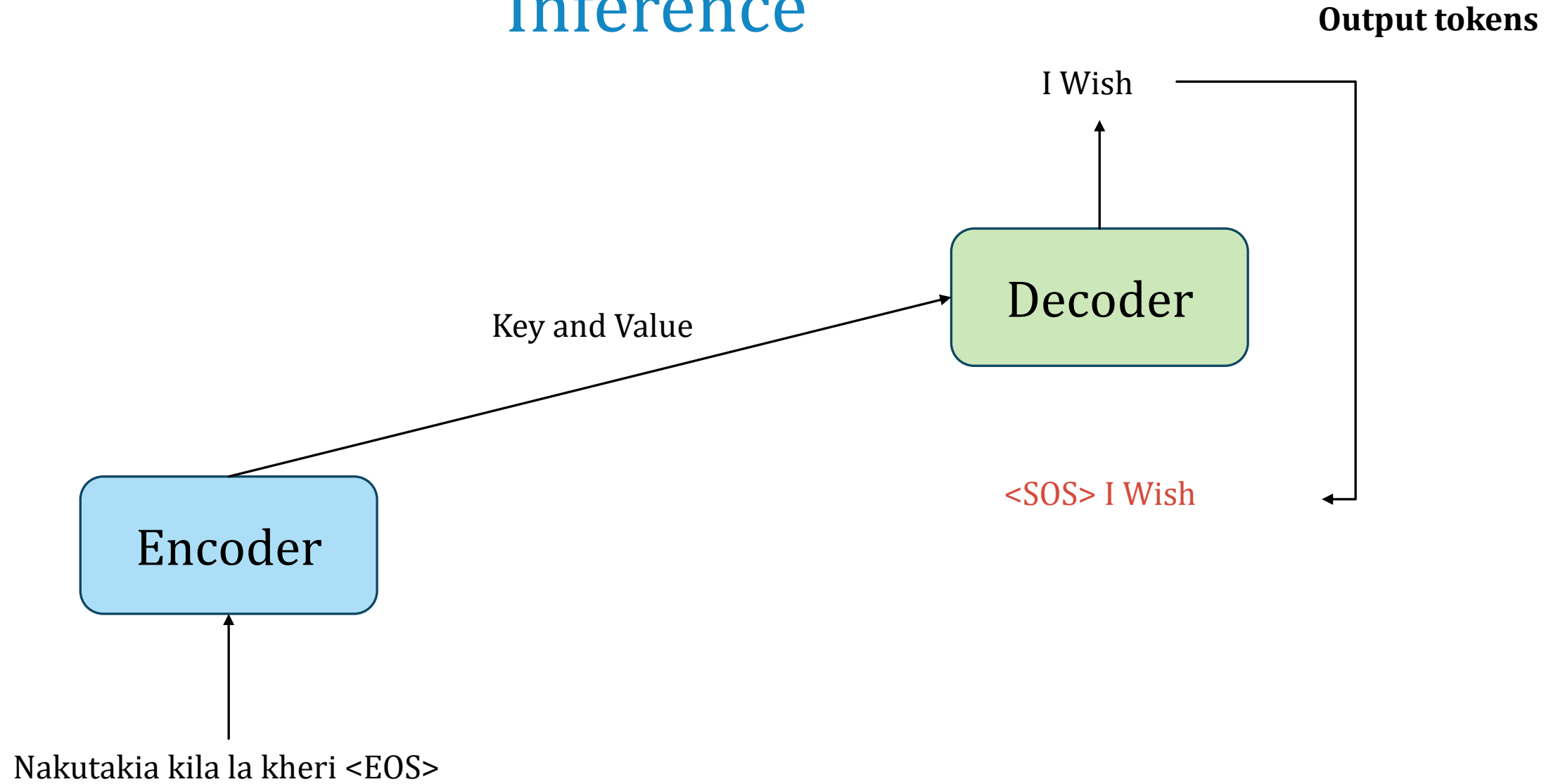
Inference



Inference

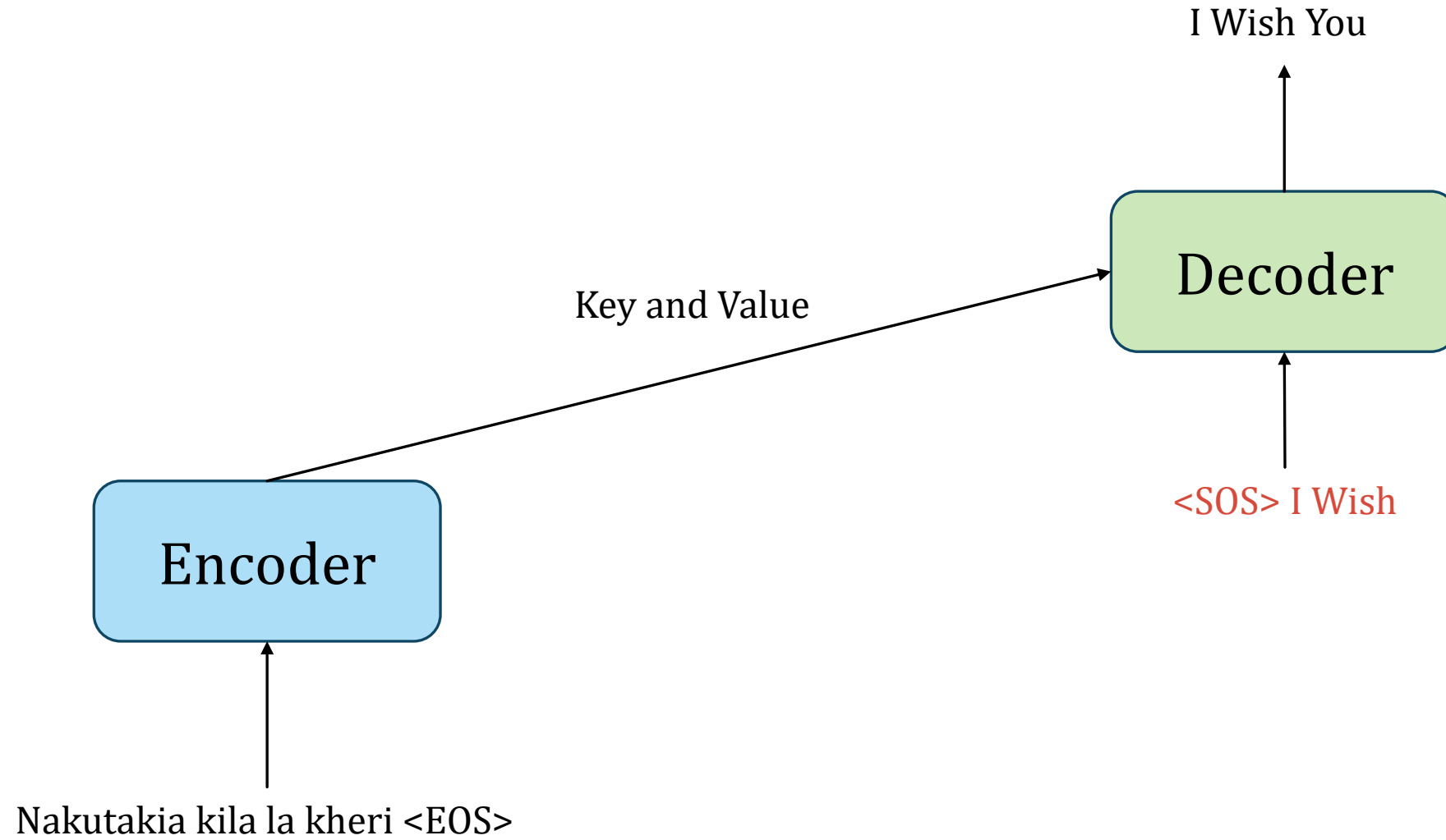


Inference

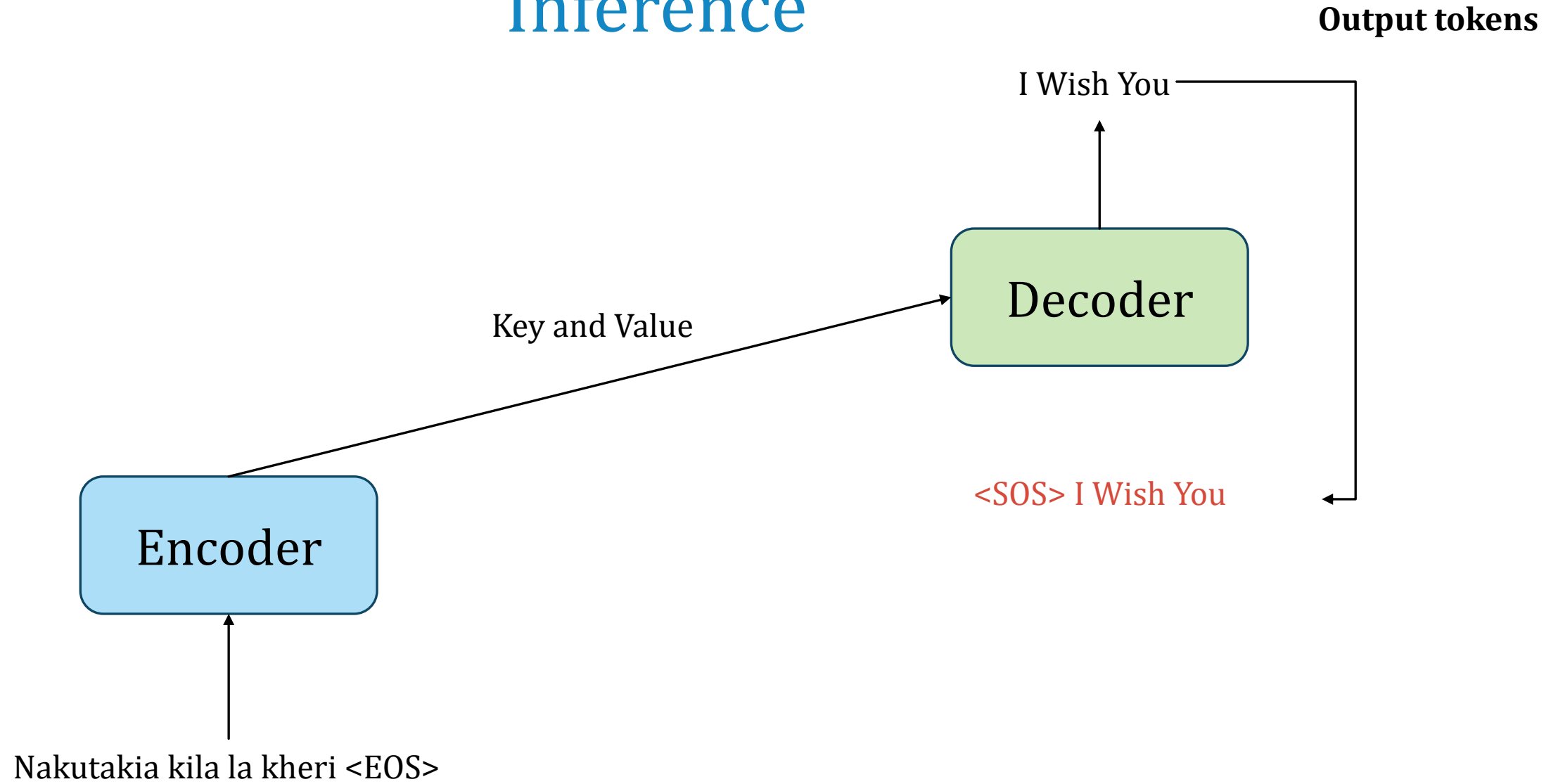


Inference

Output tokens

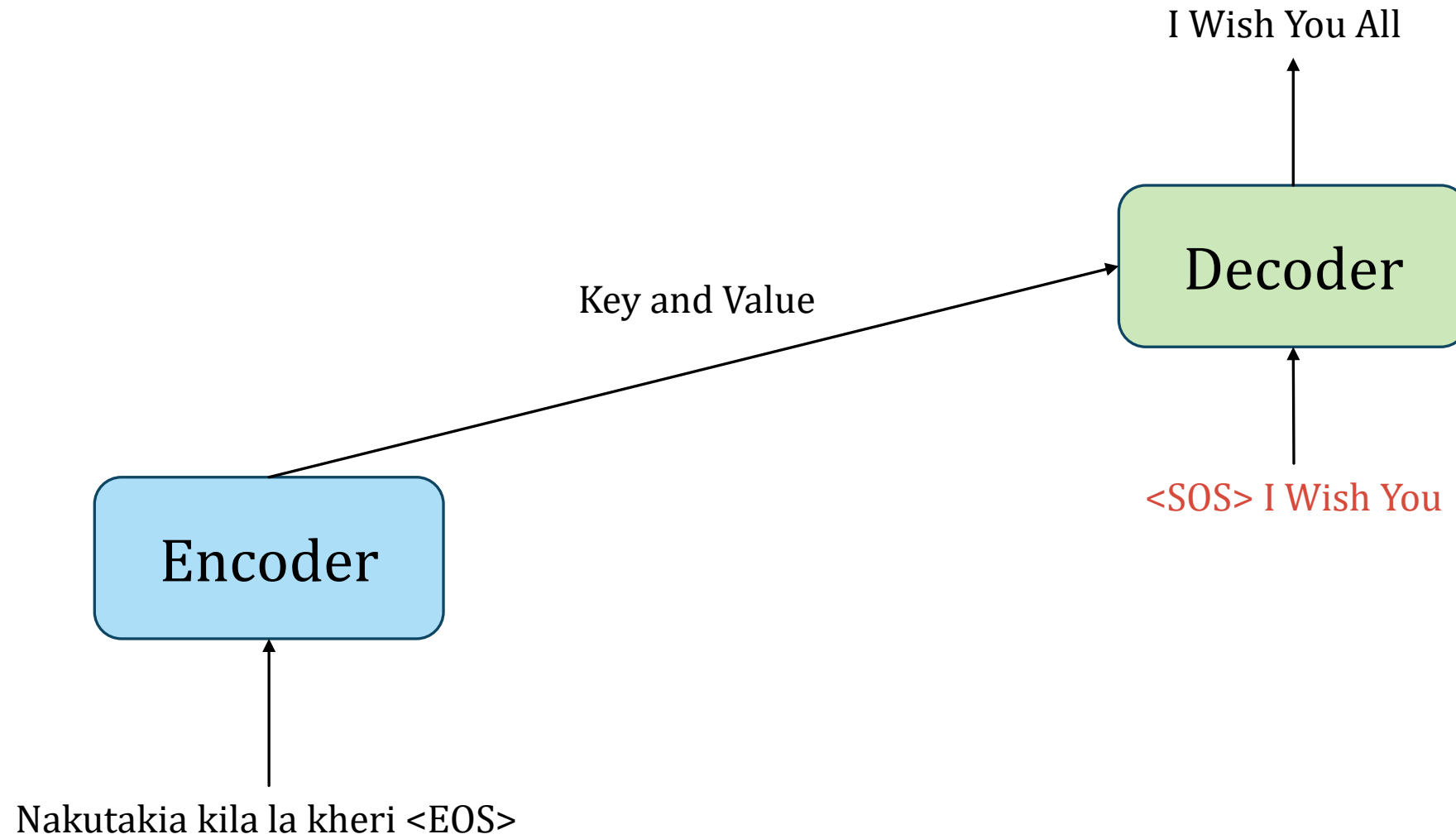


Inference

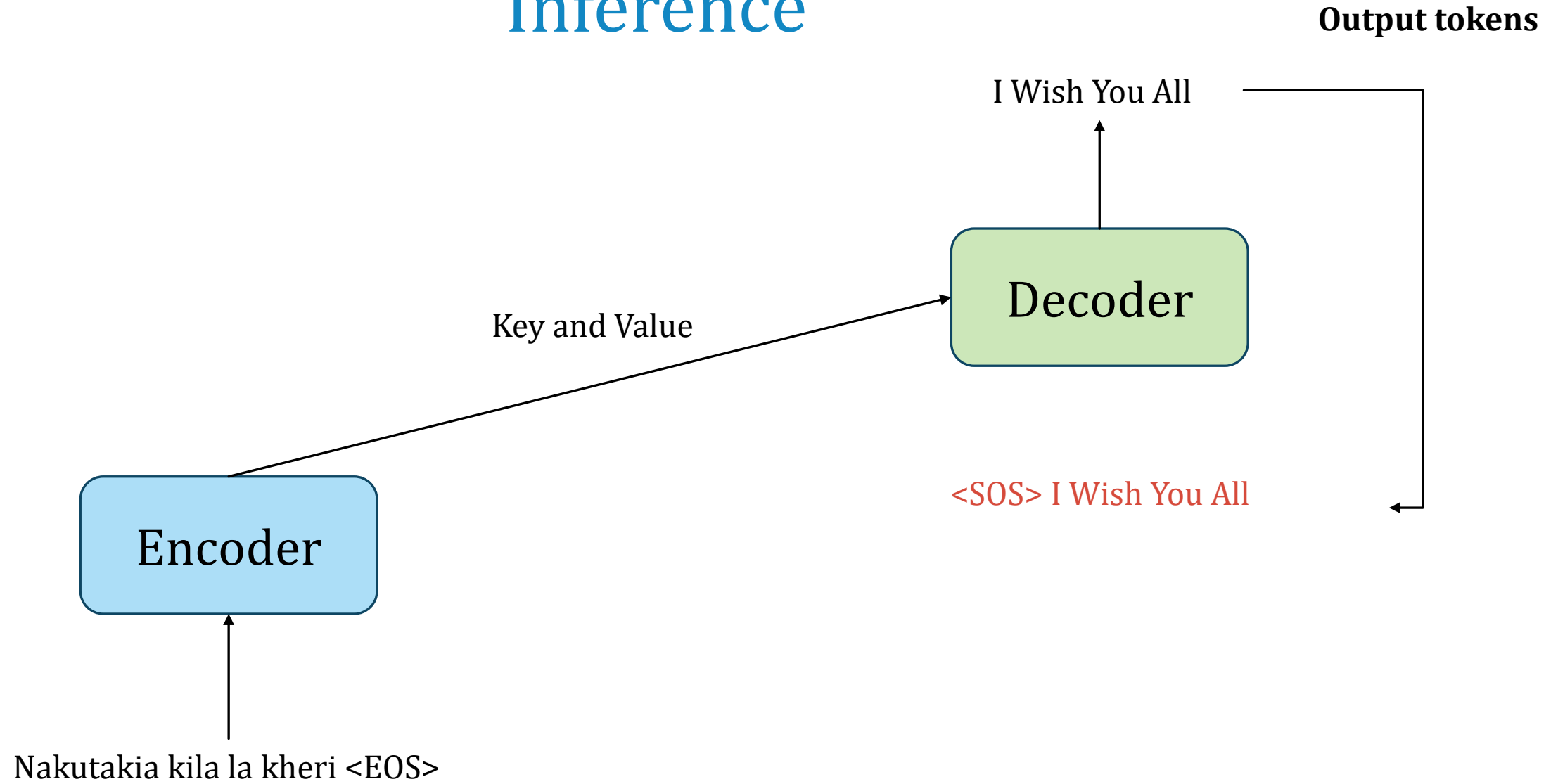


Inference

Output tokens

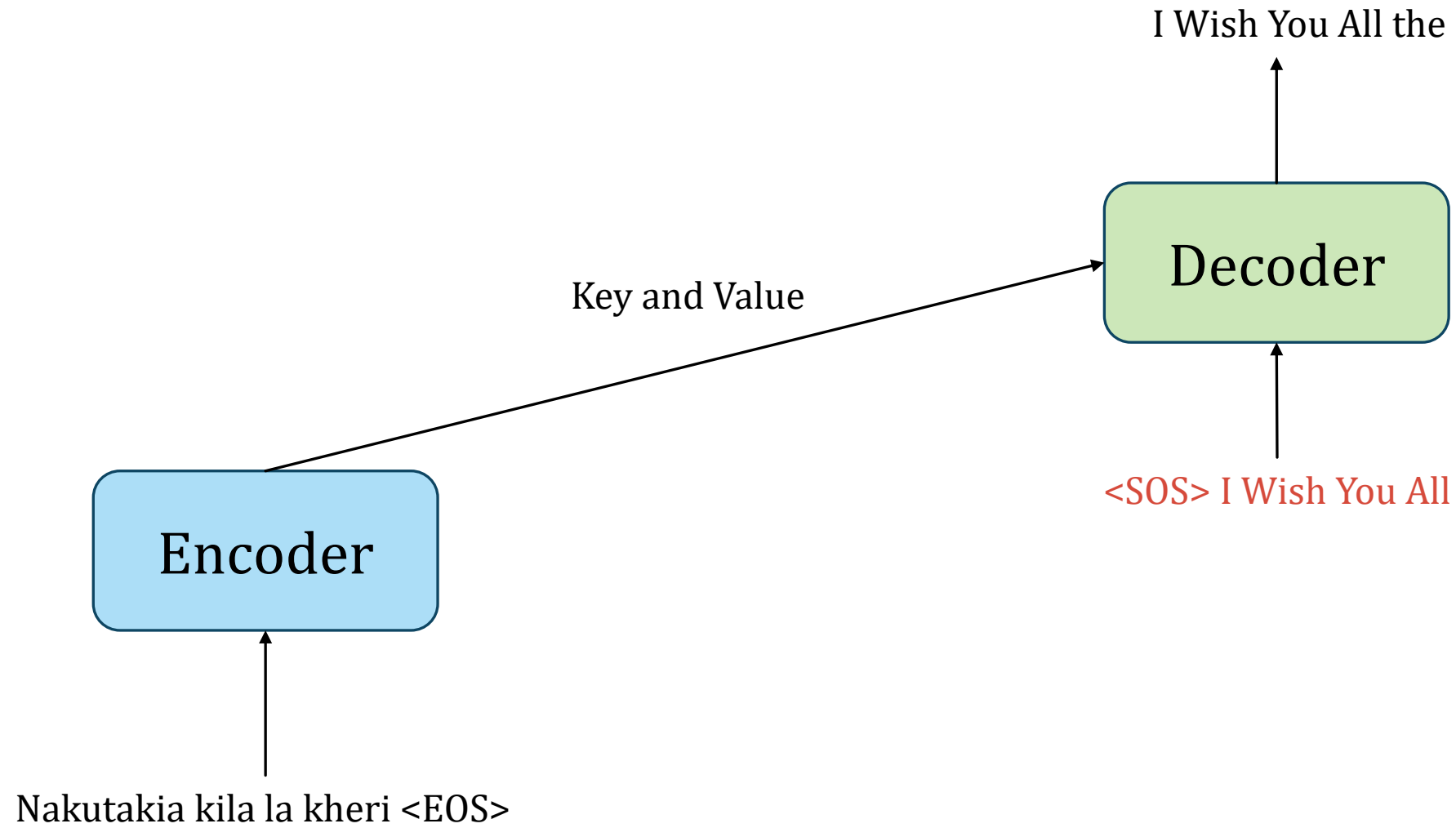


Inference

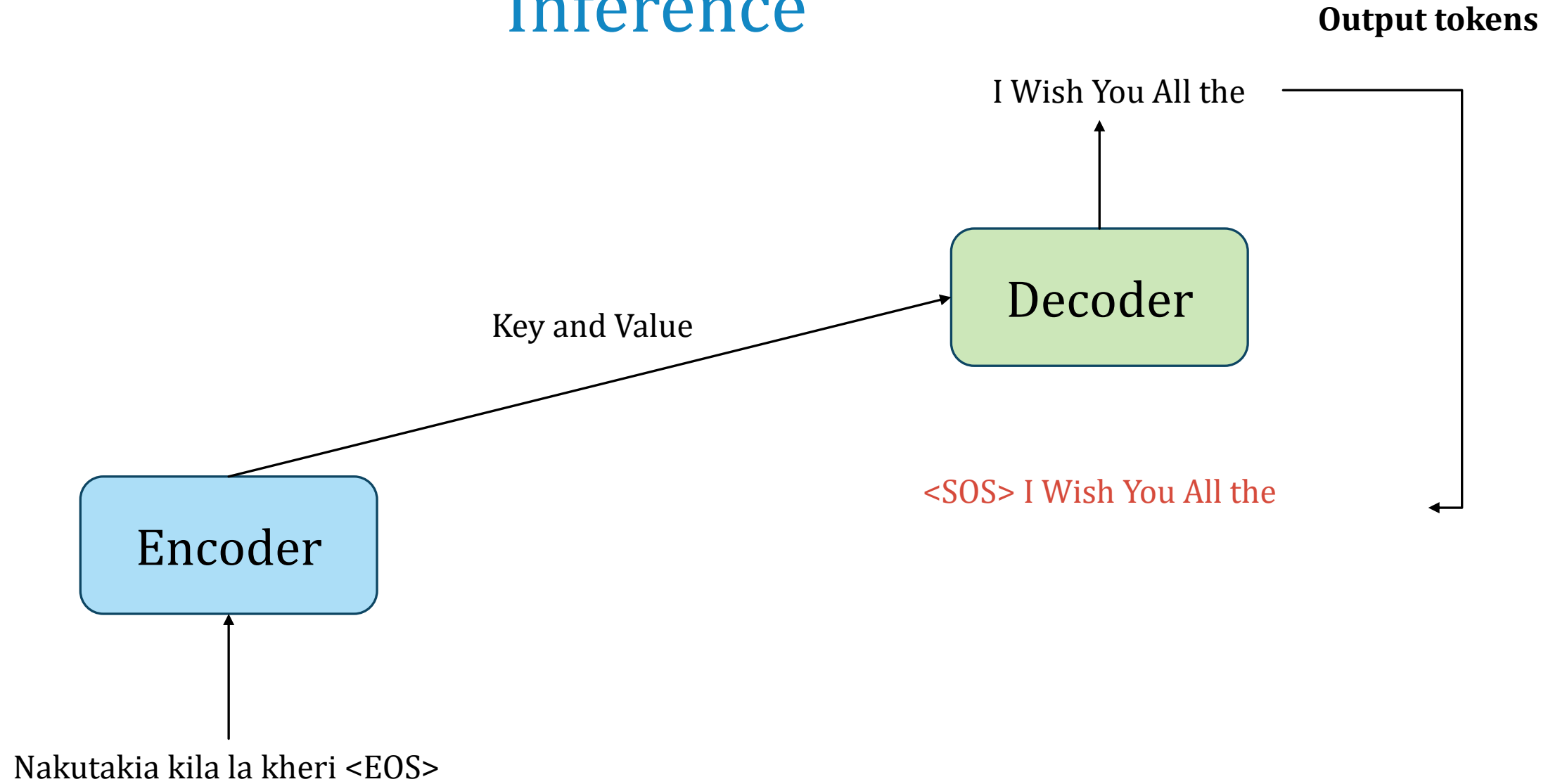


Inference

Output tokens

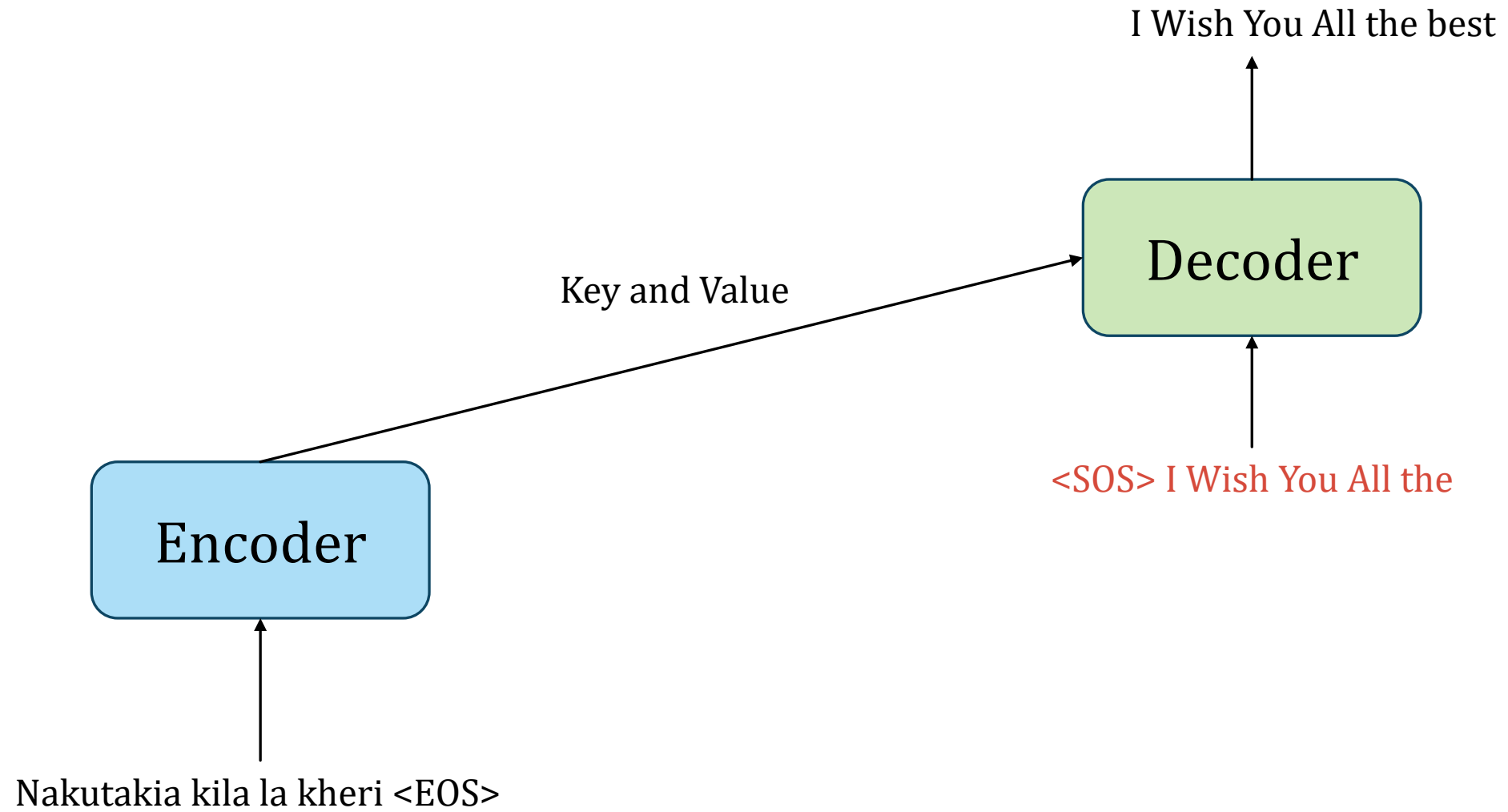


Inference

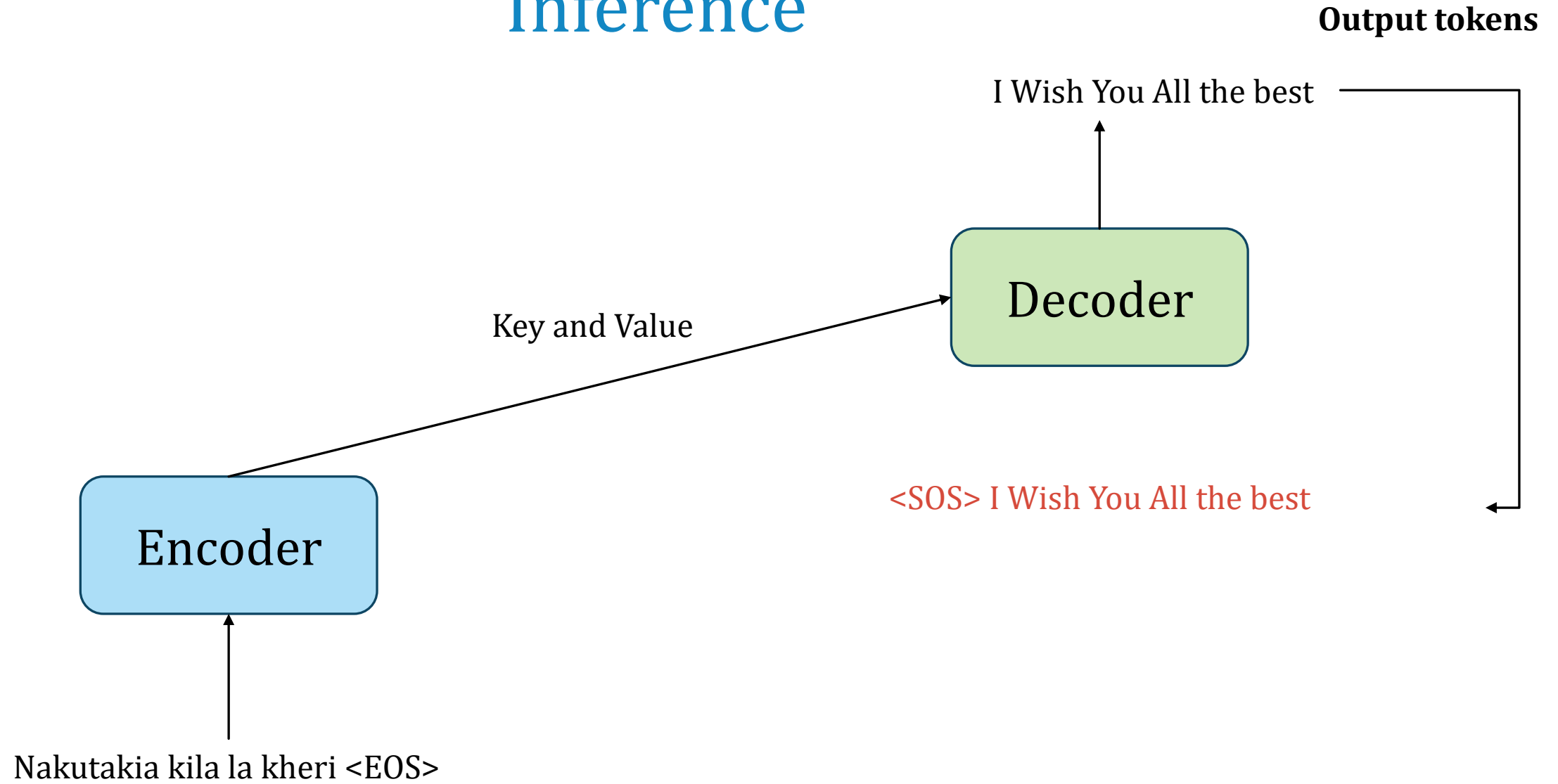


Inference

Output tokens

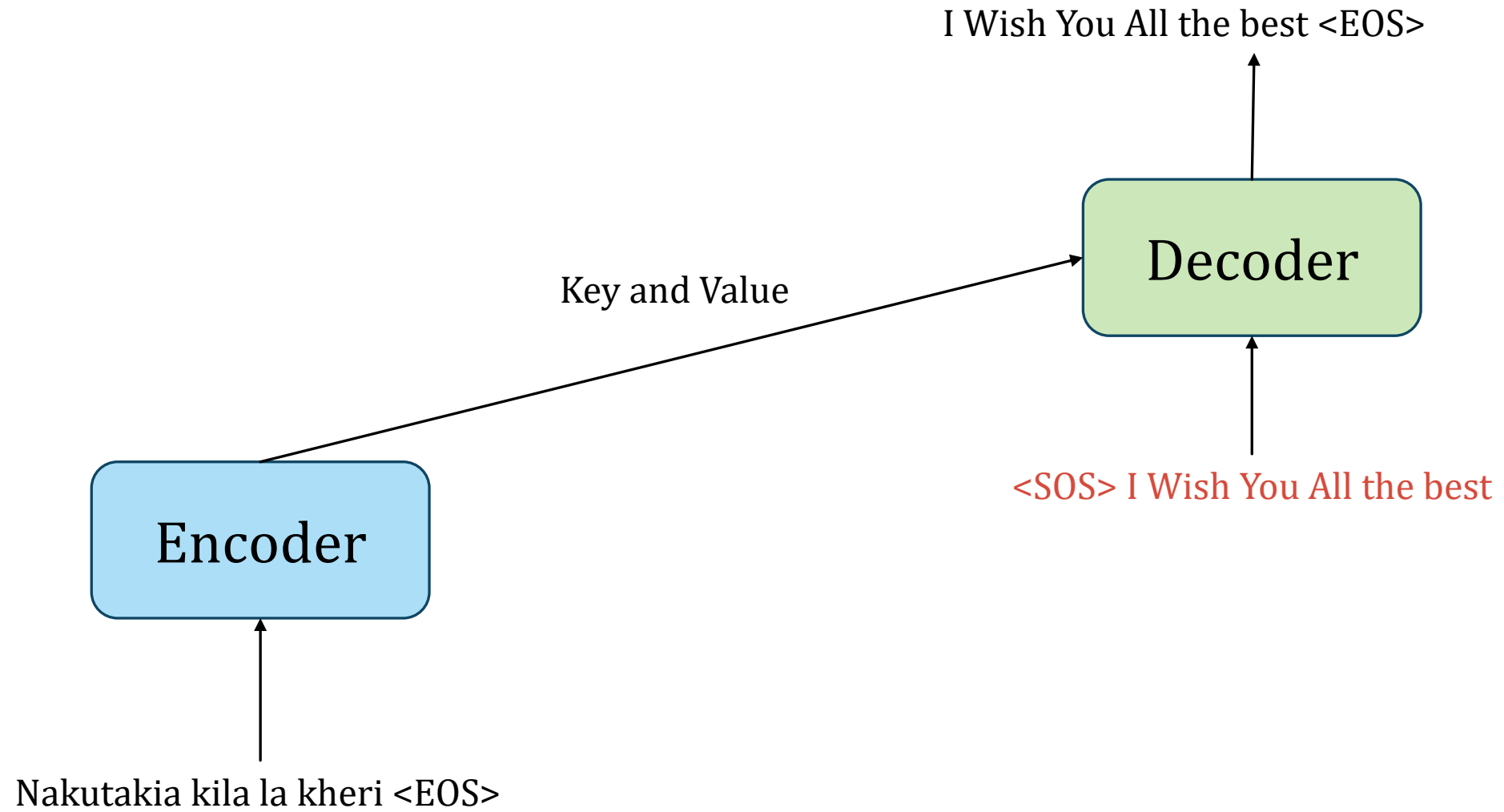


Inference

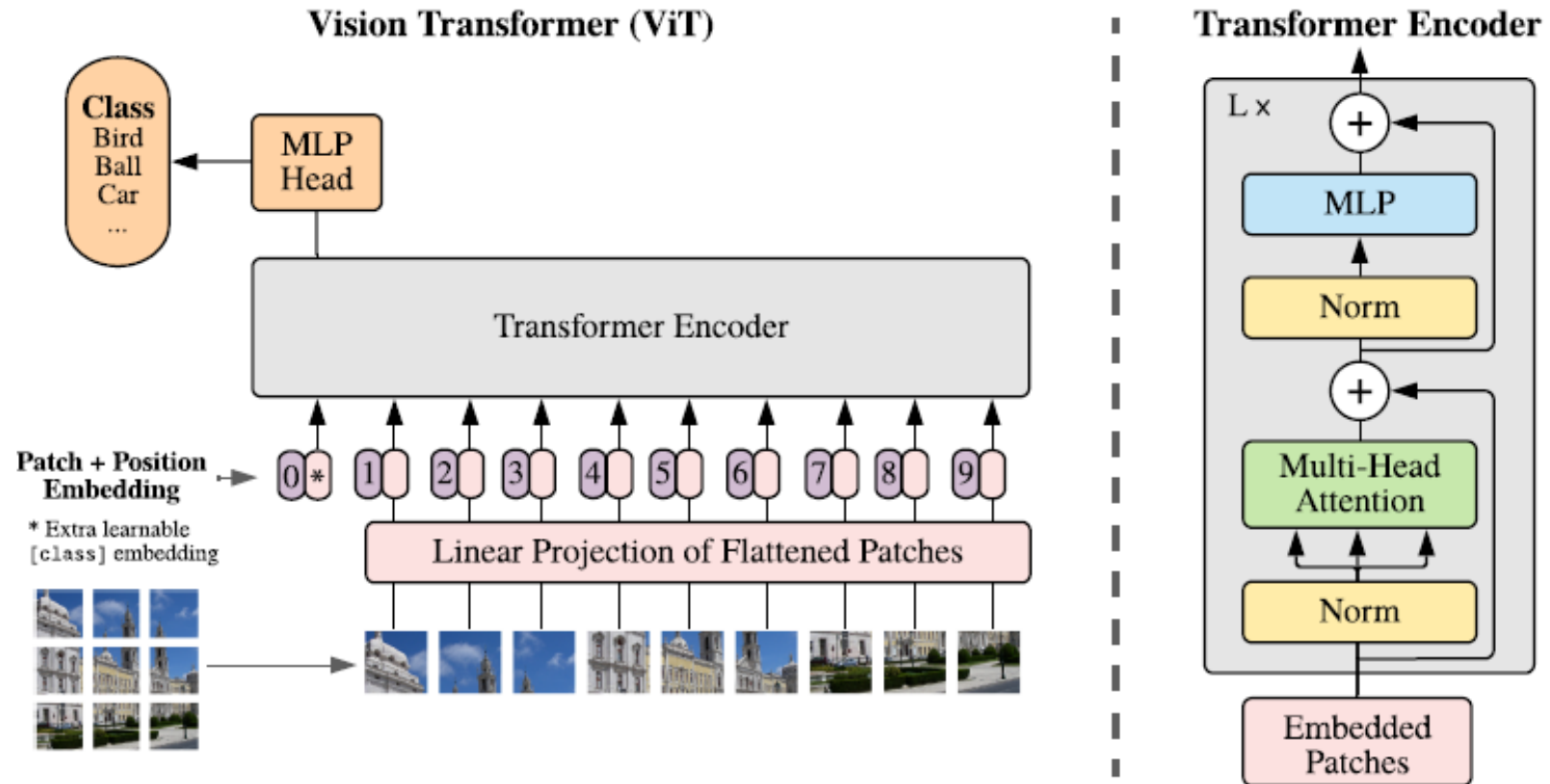


Inference

Output tokens



Vision Transformer



Detection Transformer

