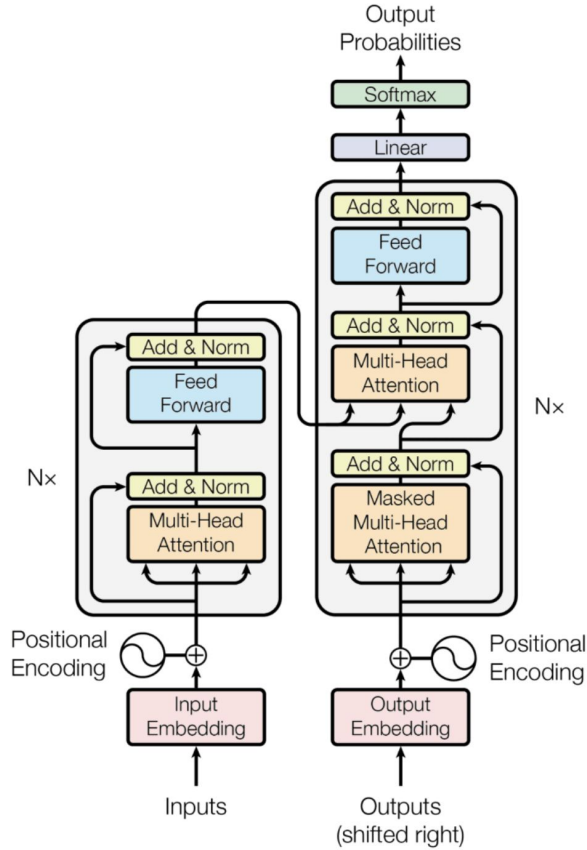
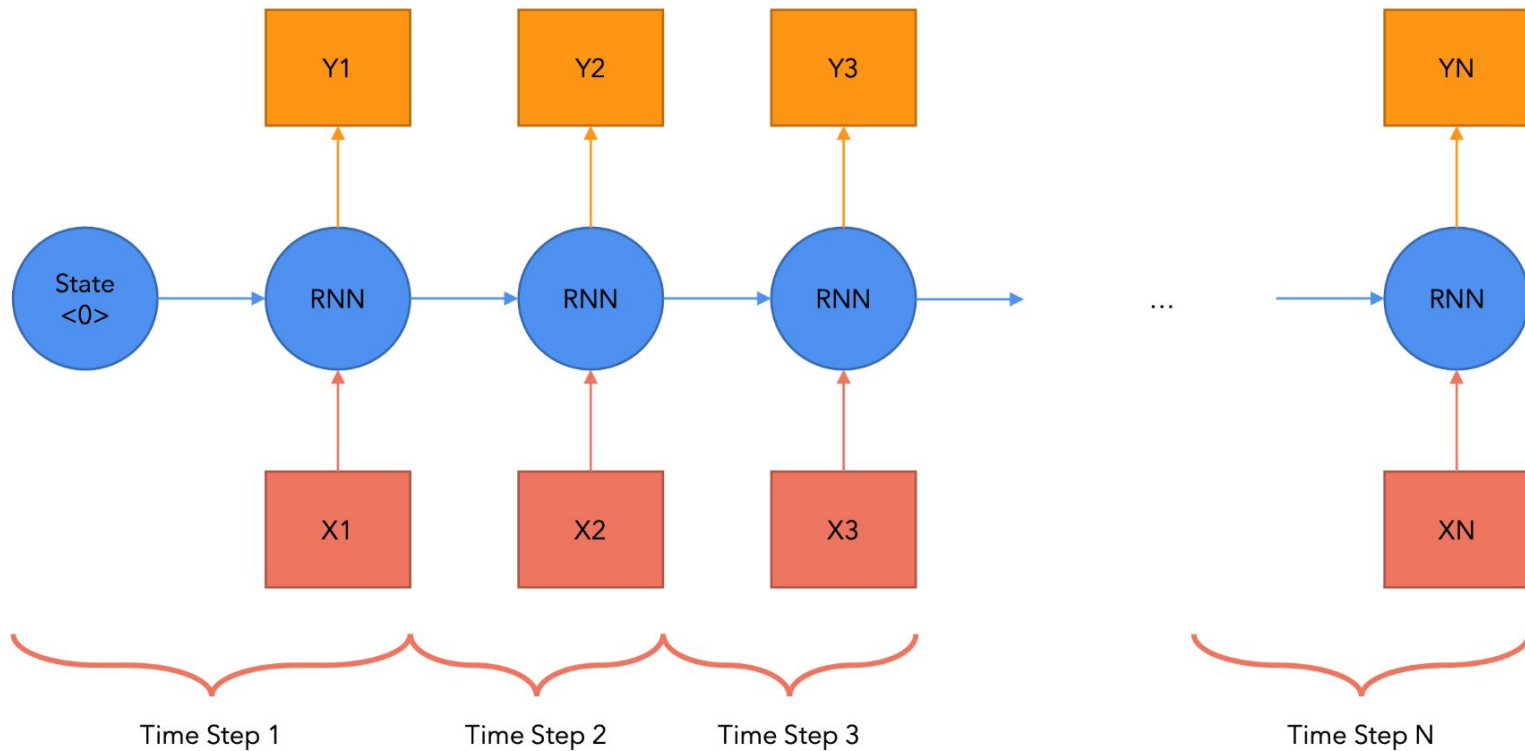


Transformers

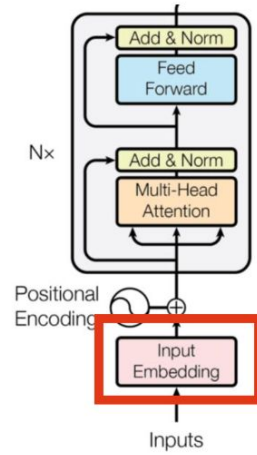


This is a reminder for the presenter of this presentation to not forget to start a recording of this presentation.

Recurrent Neural Networks (RNN)



Encoder



What is an input embedding?

Original sentence
(tokens)

YOUR

CAT

IS

A

LOVELY

CAT

Input IDs (position in
the vocabulary)

105

6587

5475

3578

65

6587

Embedding
(vector of size 512)

952.207

5450.840

1853.448

...

1.658

2671.529

171.411

3276.350

9192.819

...

3633.421

8390.473

621.659

1304.051

0.565

...

7679.805

4506.025

776.562

5567.288

58.942

...

2716.194

5119.949

6422.693

6315.080

9358.778

...

2141.081

735.147

171.411

3276.350

9192.819

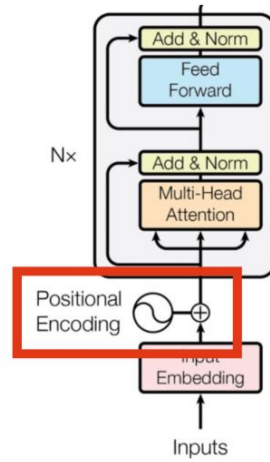
...

3633.421

8390.473

We define $d_{\text{model}} = 512$, which represents the size of the embedding vector of each word

Encoder



What is positional encoding?

Original sentence

YOUR

CAT

IS

A

LOVELY

CAT

Embedding

(vector of size 512)

952.207

5450.840

1853.448

...

1.658

2671.529

171.411

3276.350

9192.819

...

3633.421

8390.473

621.659

1304.051

0.565

...

7679.805

4506.025

776.562

5567.288

58.942

...

2716.194

5119.949

6422.693

6315.080

9358.778

...

2141.081

735.147

171.411

3276.350

9192.819

...

3633.421

8390.473

+

+

+

+

+

+

Position Embedding

(vector of size 512).

Only computed once
and reused for every
sentence during
training and inference.

...

...

...

...

...

...

1664.068

8080.133

2620.399

...

9386.405

3120.159

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

1281.458

7902.890

912.970

3821.102

1659.217

7018.620

=

=

=

=

=

=

Encoder Input

(vector of size 512)

...

...

...

...

...

...

1835.479

11356.483

11813.218

...

13019.826

11510.632

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

1452.869

11179.24

10105.789

...

5292.638

15409.093

What is positional encoding?

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

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

Sentence 1

YOUR

CAT

IS

PE(0, 0)
PE(0, 1)
PE(0, 2)
...
PE(0, 510)
PE(0, 511)

PE(1, 0)
PE(1, 1)
PE(1, 2)
...
PE(1, 510)
PE(1, 511)

PE(2, 0)
PE(2, 1)
PE(2, 2)
...
PE(2, 510)
PE(2, 511)

Sentence 2

I

LOVE

YOU

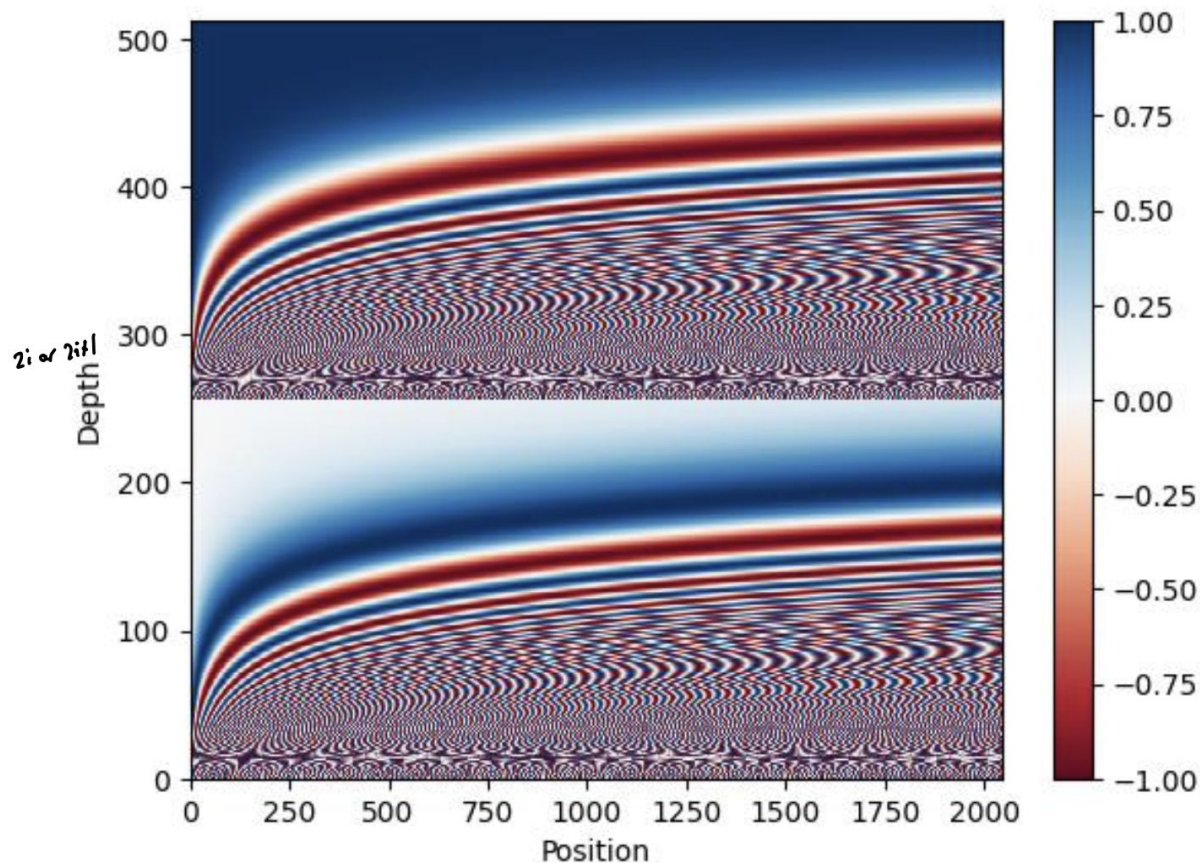
PE(0, 0)
PE(0, 1)
PE(0, 2)
...
PE(0, 510)
PE(0, 511)

PE(1, 0)
PE(1, 1)
PE(1, 2)
...
PE(1, 510)
PE(1, 511)

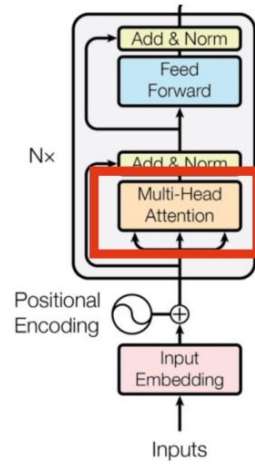
PE(2, 0)
PE(2, 1)
PE(2, 2)
...
PE(2, 510)
PE(2, 511)

We only need to compute the positional encodings once and then reuse them for every sentence, no matter if it is training or inference.

Why trigonometric functions?



Encoder



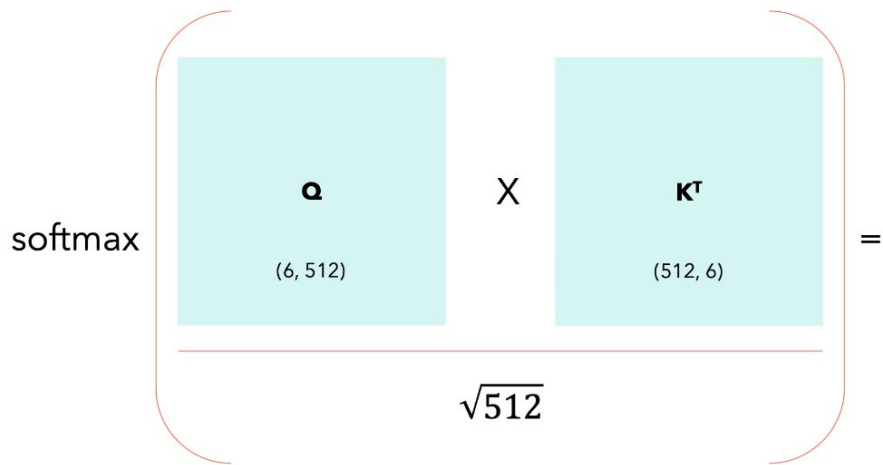
What is Self-Attention?

Self-Attention allows the model to relate words to each other.

In this simple case we consider the sequence length **seq** = 6 and **d_{model}** = **d_k** = 512.

The matrices **Q**, **K** and **V** are just the input sentence.

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



	YOUR	CAT	IS	A	LOVELY	CAT	Σ
YOUR	0.268	0.119	0.134	0.148	0.179	0.152	1
CAT	0.124	0.278	0.201	0.128	0.154	0.115	1
IS	0.147	0.132	0.262	0.097	0.218	0.145	1
A	0.210	0.128	0.206	0.212	0.119	0.125	1
LOVELY	0.146	0.158	0.152	0.143	0.227	0.174	1
CAT	0.195	0.114	0.203	0.103	0.157	0.229	1

* all values are random.

* for simplicity I considered only one head, which makes **d_{model}** = **d_k**.

(6, 6)

How to compute Self-Attention?

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

	YOUR	CAT	IS	A	LOVELY	CAT
YOUR	0.268	0.119	0.134	0.148	0.179	0.152
CAT	0.124	0.278	0.201	0.128	0.154	0.115
IS	0.147	0.132	0.262	0.097	0.218	0.145
A	0.210	0.128	0.206	0.212	0.119	0.125
LOVELY	0.146	0.158	0.152	0.143	0.227	0.174
CAT	0.195	0.114	0.203	0.103	0.157	0.229

(6, 6)

X

V

(6, 512)

=

Attention

(6, 512)

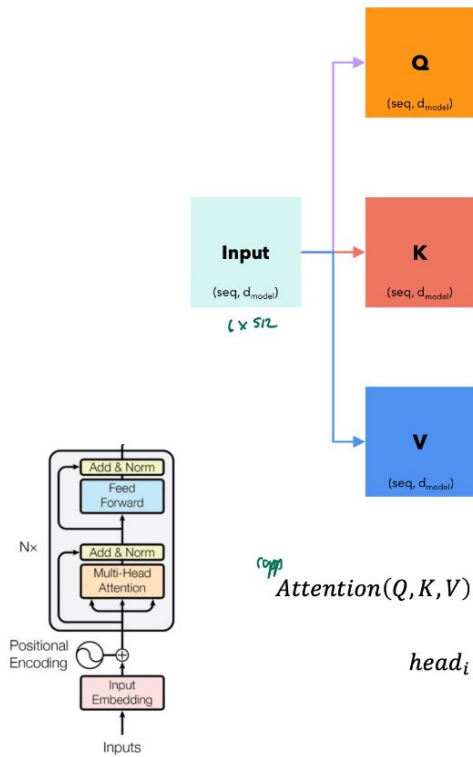
Each row in this matrix captures not only the meaning (given by the embedding) or the position in the sentence (represented by the positional encodings) but also each word's interaction with other words.

Multi-head Attention

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

$$MultiHead(Q, K, V) = \text{Concat}(head_1 \dots head_h)W^O$$

$$head_i = Attention(QW_i^Q, KW_i^K, VW_i^V)$$



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

$$\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

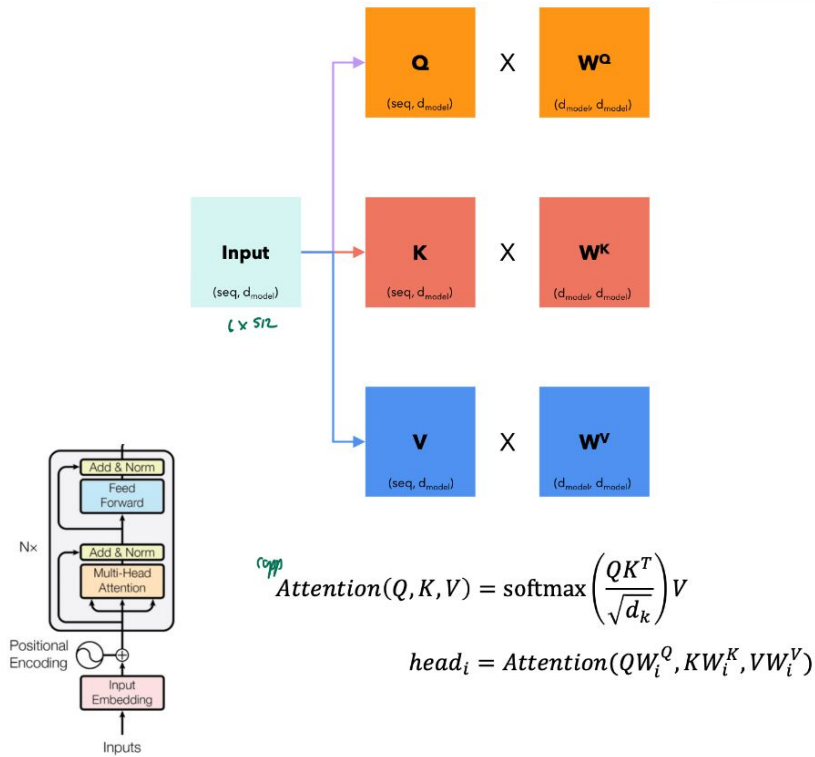
5

seq = sequence length

d_{model} = size of the embedding vector

h = number of heads

$d_k = d_v$ = d_{model} / h



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

$$head_i = Attention(QW_i^Q, KW_i^K, VW_i^V)$$

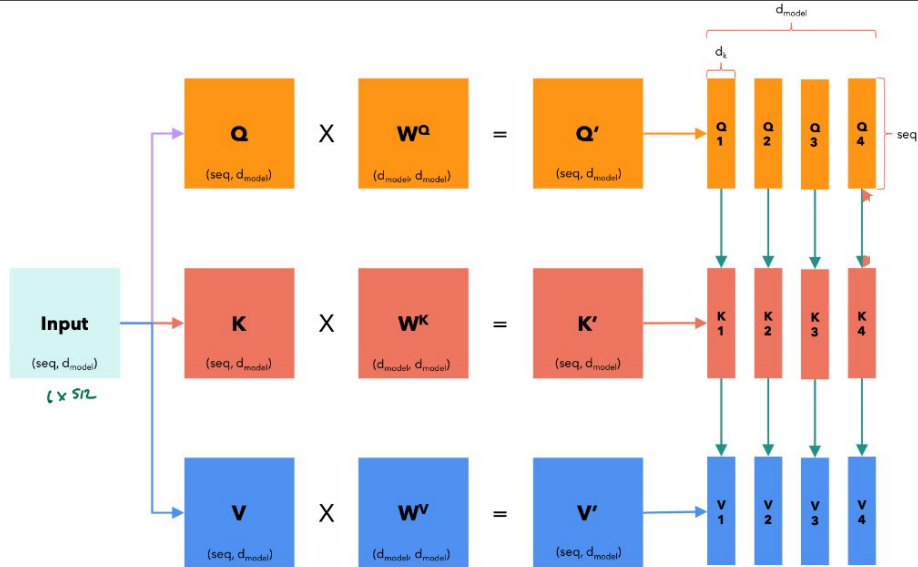
56

seq = sequence length

d_{model} = size of the embedding vector

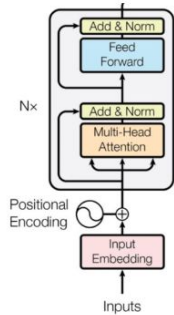
h = number of heads

$d_k = d_v$ = d_{model} / h



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

$$head_i = Attention(QW_i^Q, KW_i^K, VW_i^V)$$

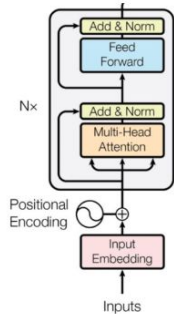
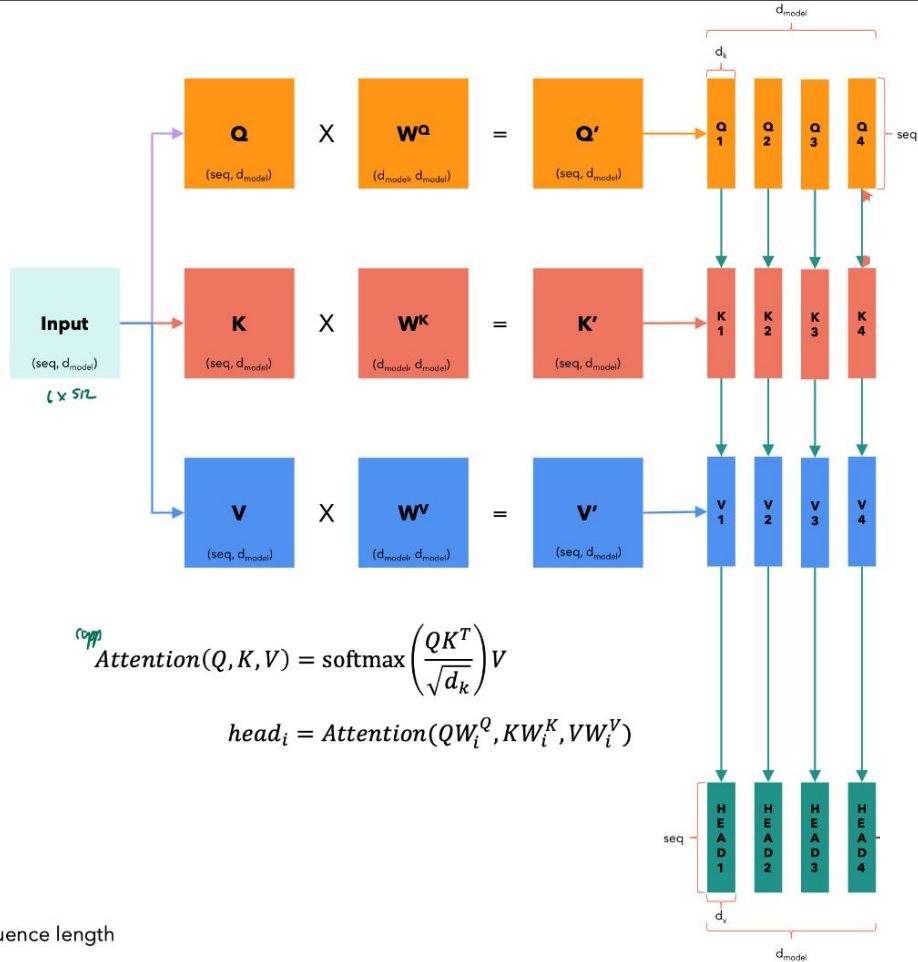


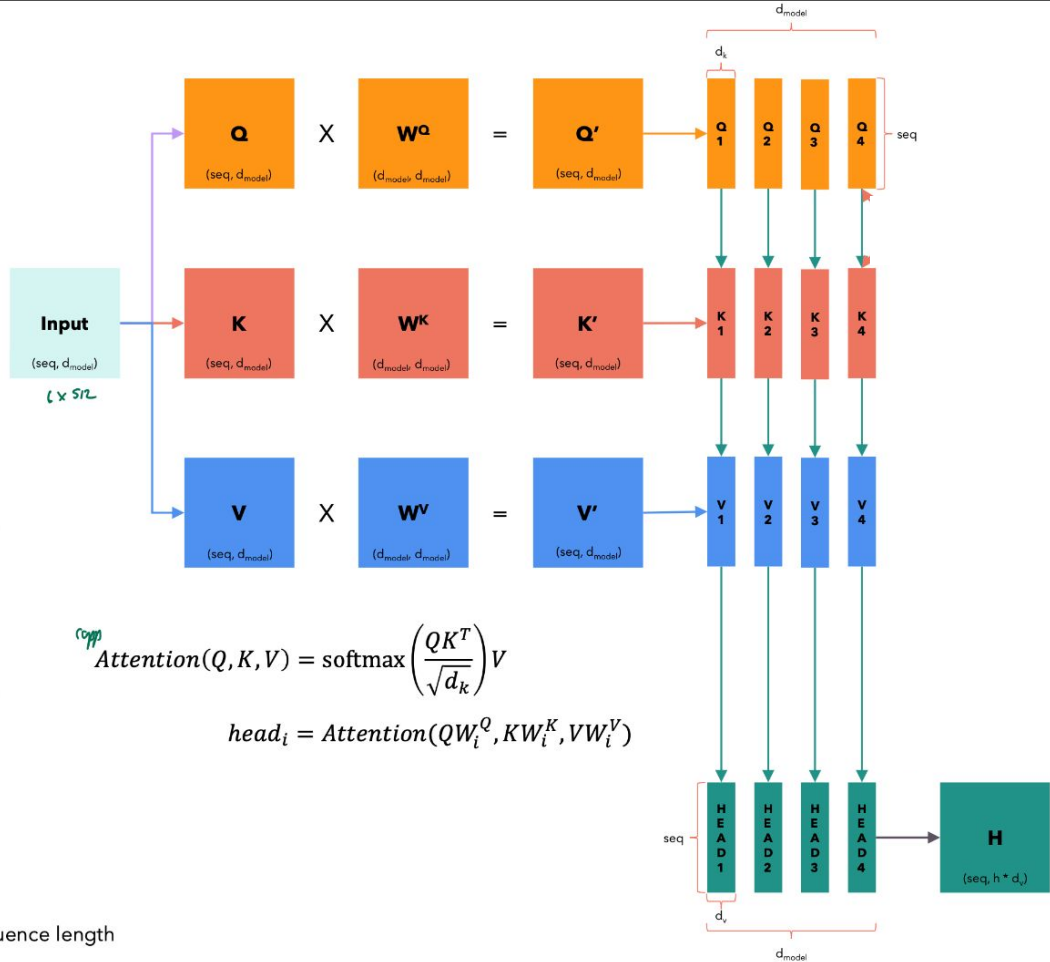
seq = sequence length

d_{model} = size of the embedding vector

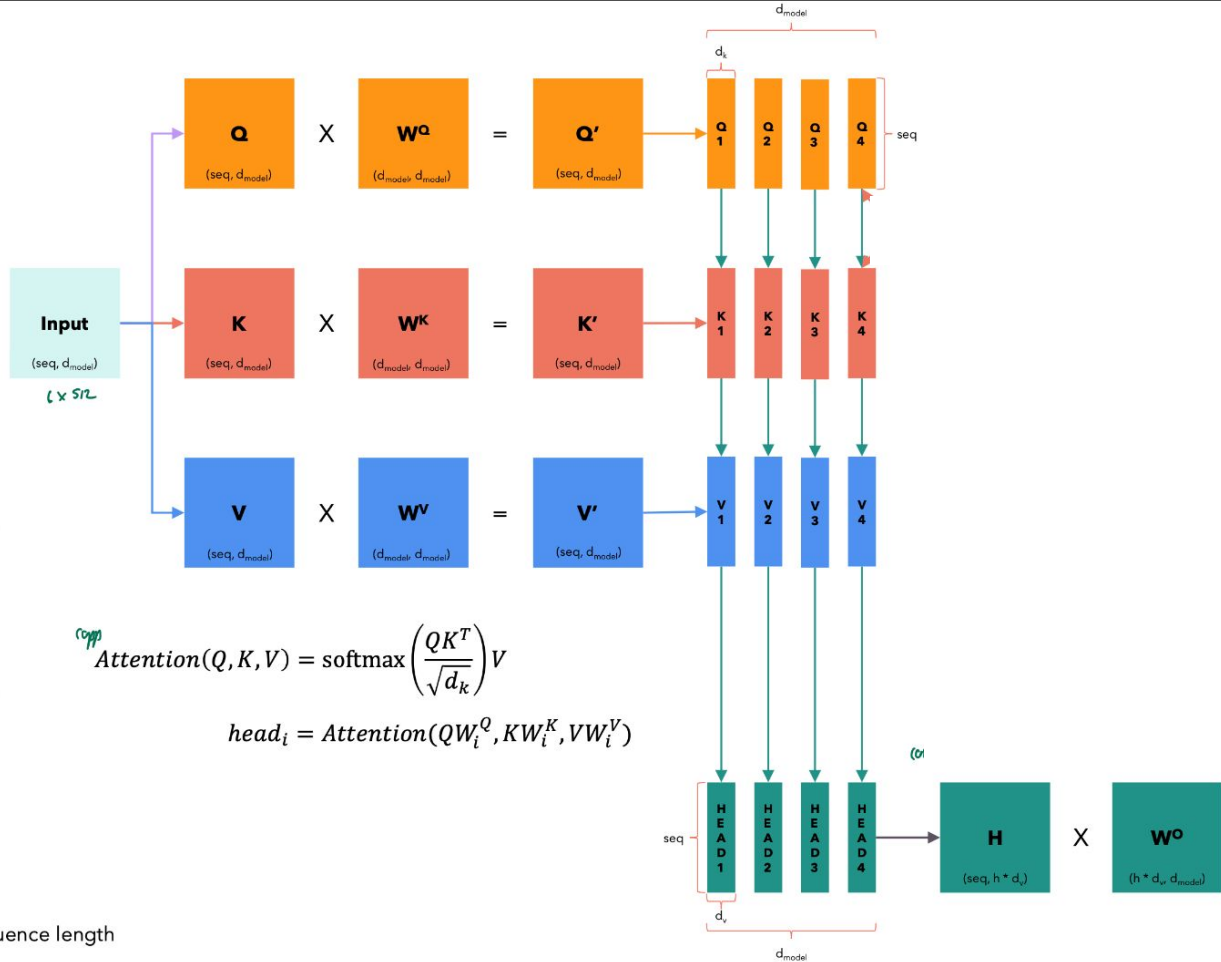
h = number of heads

$d_k = d_v$ = d_{model} / h

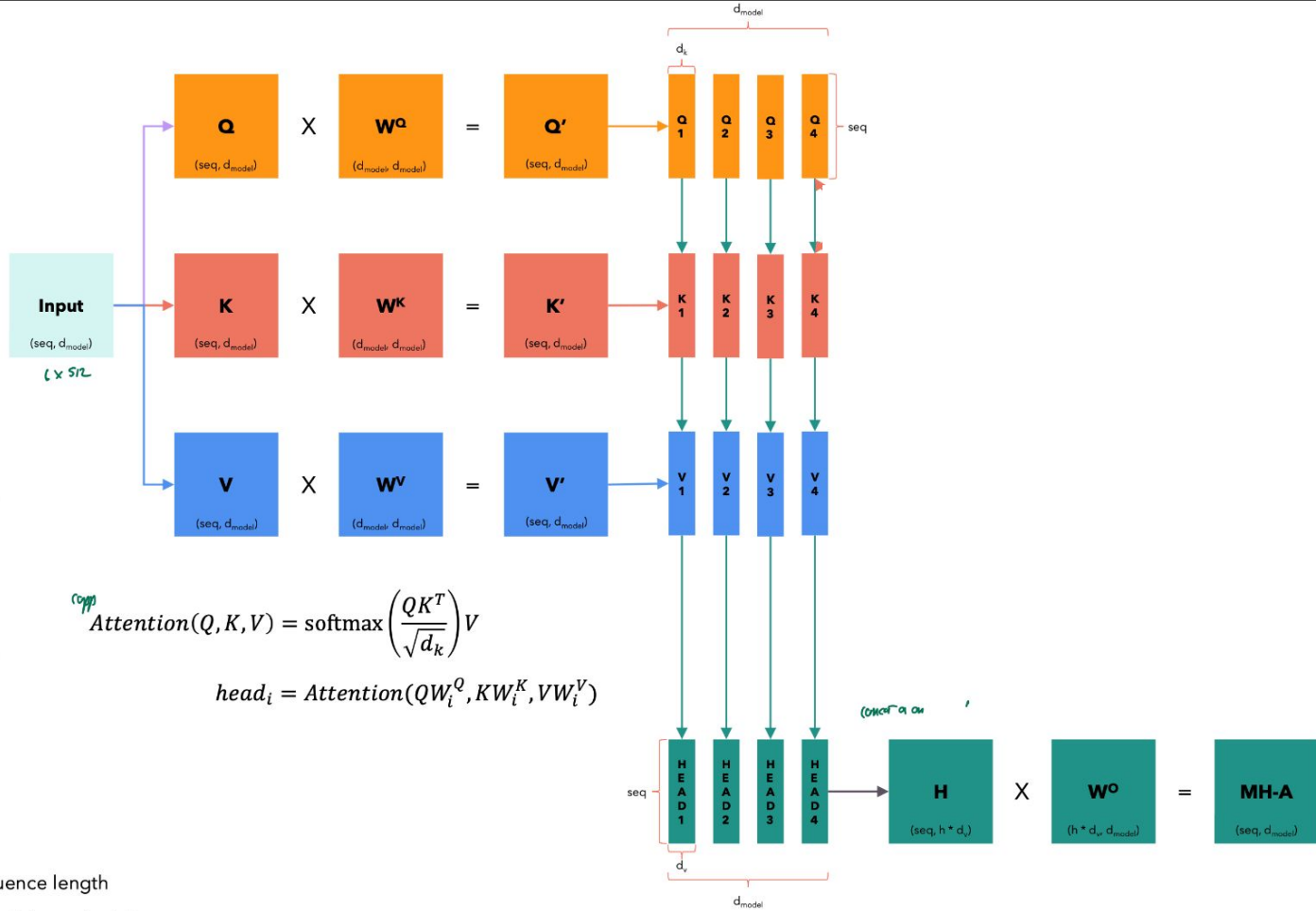




seq = sequence length
 d_{model} = size of the embedding vector
 h = number of heads
 $d_k = d_v$ = d_{model} / h



seq = sequence length
 d_{model} = size of the embedding vector
 h = number of heads
 $d_k = d_v = d_{model} / h$

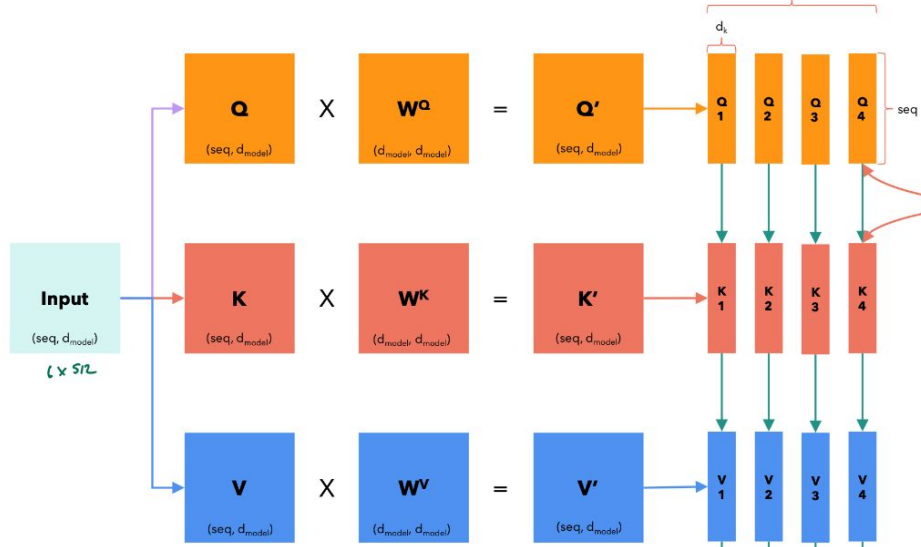
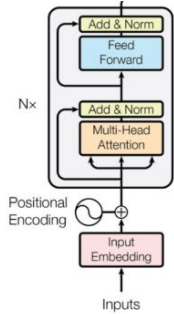


seq = sequence length

d_{model} = size of the embedding vector

h = number of heads

$d_k = d_v = d_{model} / h$



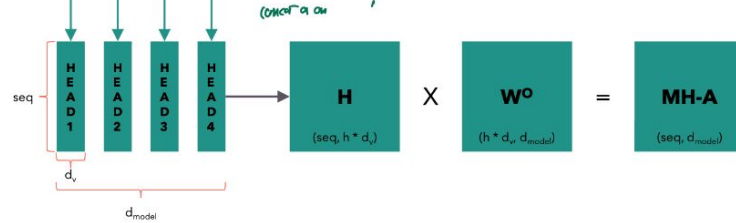
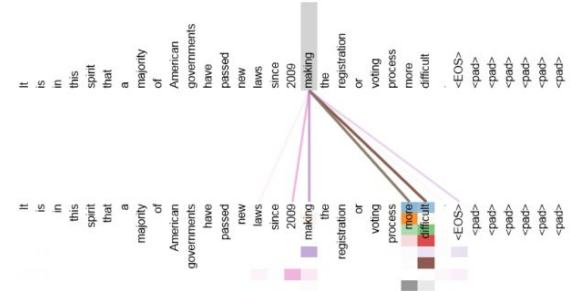
$$Attention(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

$$head_i = Attention(QW_i^Q, KW_i^K, VW_i^V)$$

$$\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) =$$

		1	2	3	4	5	6	7	8
1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
2	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
3	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
4	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
5	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
6	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
7	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
8	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0

Attention Visualizations



$$MultiHead(Q, K, V) = \text{Concat}(head_1 \dots head_h)W^O$$

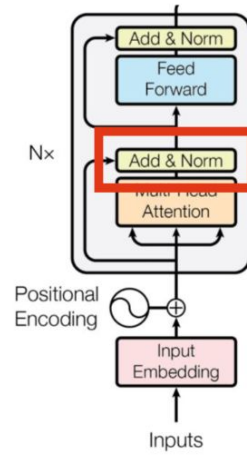
seq = sequence length

d_{model} = size of the embedding vector

h = number of heads

$d_k = d_v$ = d_{model} / h

Encoder



What is layer normalization?

Batch of 3 items

ITEM 1

50.147
3314.825
...
...
8463.361
8.021

μ_1

σ_1^2

ITEM 2

1242.223
688.123
...
...
434.944
149.442

μ_2

σ_2^2

ITEM 3

9.370
4606.674
...
...
944.705
21189.444

μ_3

σ_3^2

$$\hat{x}_j = \frac{x_j - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

What is layer normalization?

Batch of 3 items

ITEM 1

50.147
3314.825
...
...
8463.361
8.021

μ_1

σ_1^2

ITEM 2

1242.223
688.123
...
...
434.944
149.442

μ_2

σ_2^2

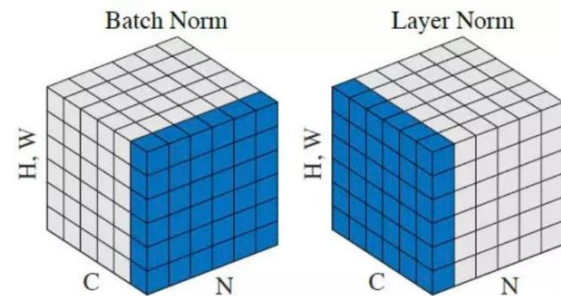
ITEM 3

9.370
4606.674
...
...
944.705
21189.444

μ_3

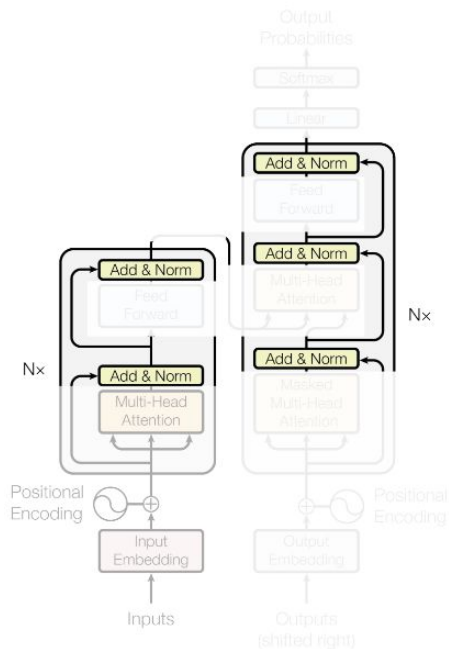
σ_3^2

$$\hat{x}_j = \frac{x_j - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$



Attention Is All You Need - The Transformer architecture

2017, Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, Illia Polosukhin



Residual connections

Each module's output has the exact same shape as its input.

Following ResNets, the module computes a "residual" instead of a new value:

$$z_i = \text{Module}(x_i) + x_i$$

This was shown to dramatically improve trainability.

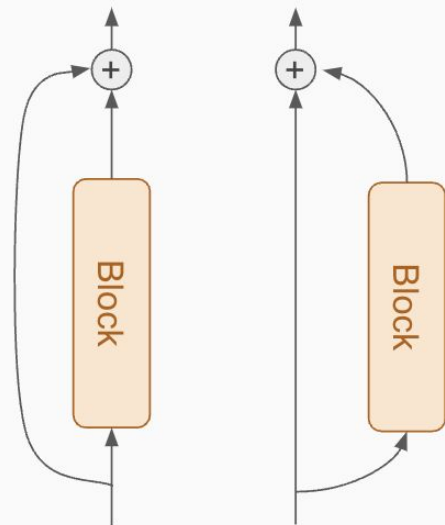
LayerNorm

Normalization also dramatically improves trainability.

There's **post-norm** (original)

$$z_i = \text{LN}(\text{Module}(x_i) + x_i)$$

"Skip connection" == "Residual block"



and **pre-norm** (modern)

$$z_i = \text{Module}(\text{LN}(x_i)) + x_i$$

Attention Is All You Need - The Transformer architecture

2017, Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, Illia Polosukhin

Point-wise MLP

A simple MLP applied to each token individually:

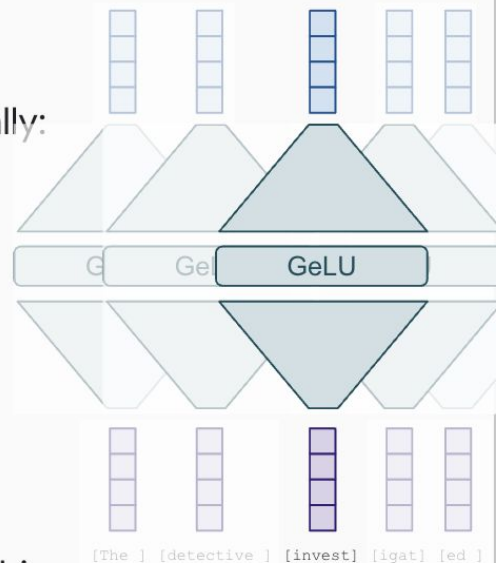
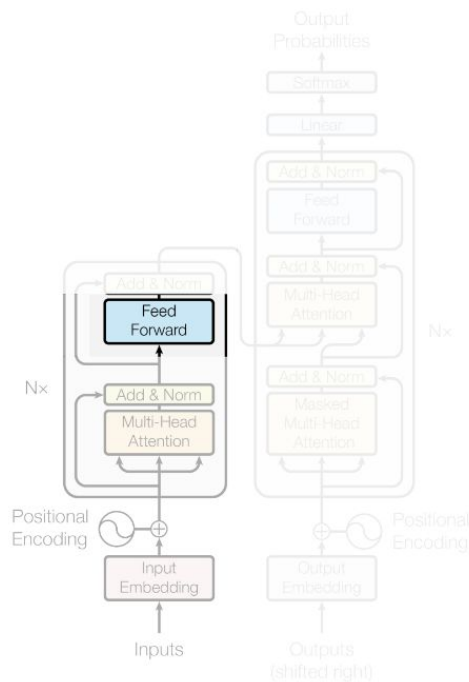
$$z_i = W_2 \text{GeLU}(W_1 x + b_1) + b_2$$

Think of it as each token pondering for itself about what it has observed previously.

There's some weak evidence this is where "world knowledge" is stored, too.

It contains the bulk of the parameters. When people make giant models and sparse/moe, this is what becomes giant.

Some people like to call it 1x1 convolution.



Query = **"love"**

The diagram illustrates a query-based search process. A central query, "love", is shown on the left. Five arrows originate from this query and point to a table on the right. The table is organized into two columns: "Keys" and "Values". The "Keys" column lists movie genres, and the "Values" column lists specific movie titles. The rows represent the results of the query, showing that "love" can be associated with various genres and their corresponding movies.

Keys	Values
ROMANTIC	TITANIC
ACTION	THE DARK KNIGHT
SCIFI	INCEPTION
HORROR	THE SHINING
COMEDY	THE INTOUCHABLES

* this could be a Python dictionary
or a database table.

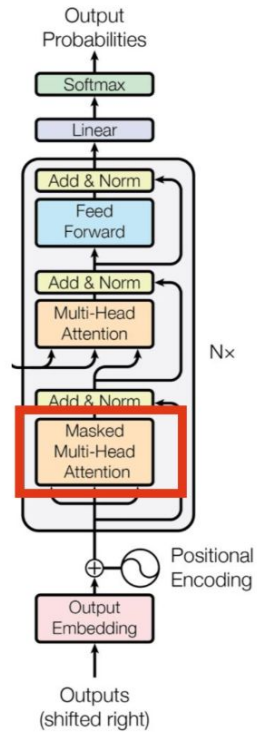
Query = **"love"**

The diagram illustrates a query 'love' branching into five movie genres. Each genre is represented by a colored box (purple for genres, teal for movie titles) in a table. Arrows point from the query to each genre box. Above the table, 'Keys' and 'Values' are labeled with brackets. The genres are ROMANTIC, ACTION, SCIFI, HORROR, and COMEDY. The corresponding movie titles are TITANIC, THE DARK KNIGHT, INCEPTION, THE SHINING, and THE INTOUCHABLES.

Keys	Values
ROMANTIC	TITANIC
ACTION	THE DARK KNIGHT
SCIFI	INCEPTION
HORROR	THE SHINING
COMEDY	THE INTOUCHABLES

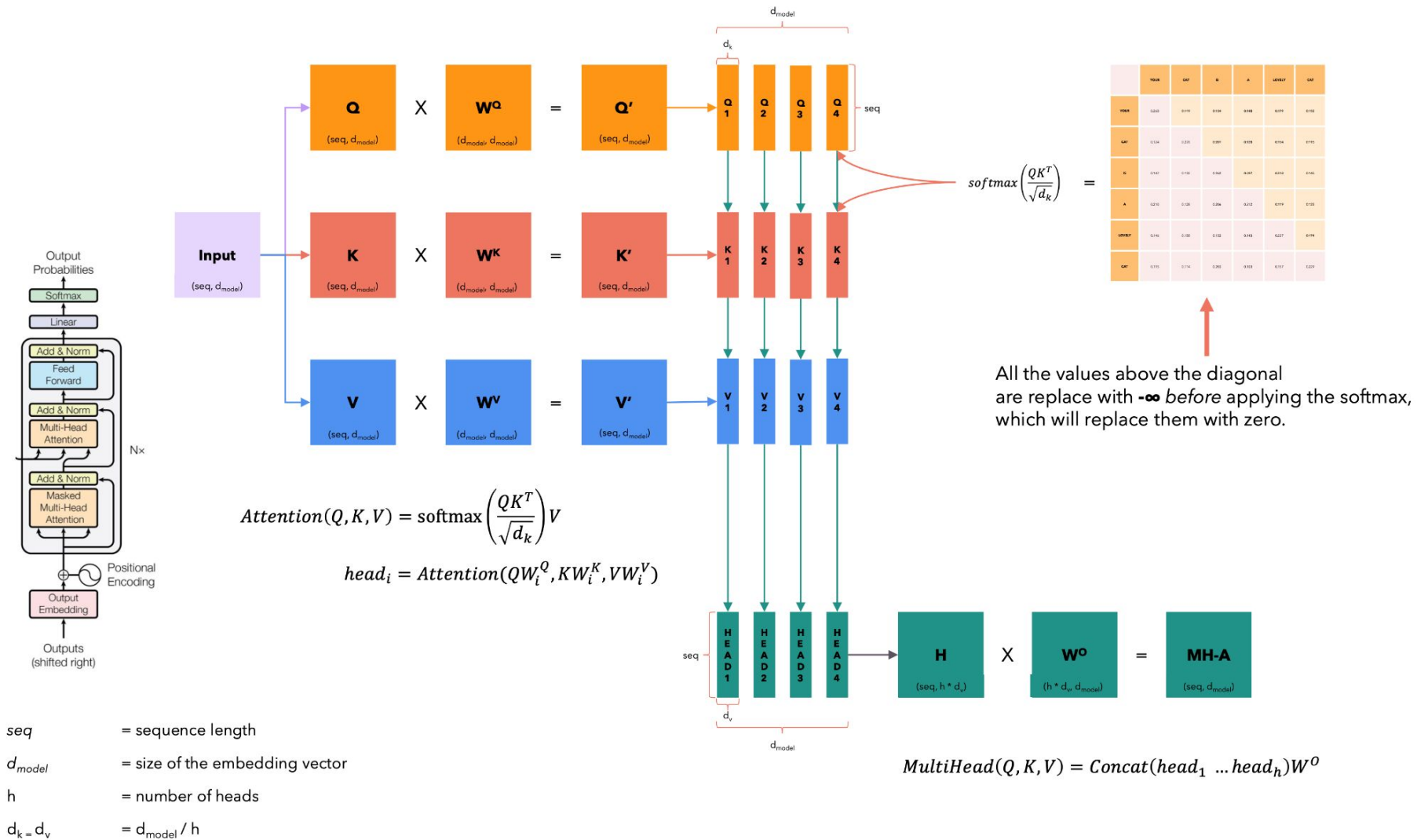
* this could be a Python dictionary
or a database table.

Decoder



What is Masked Multi-Head Attention?

	YOUR	CAT	IS	A	LOVELY	CAT
YOUR	0.268	0.119	0.134	0.148	0.179	0.152
CAT	0.124	0.278	0.201	0.128	0.154	0.115
IS	0.147	0.132	0.262	0.097	0.218	0.145
A	0.210	0.128	0.206	0.212	0.119	0.125
LOVELY	0.146	0.158	0.152	0.143	0.227	0.174
CAT	0.195	0.114	0.203	0.103	0.157	0.229



Training



I love you very much

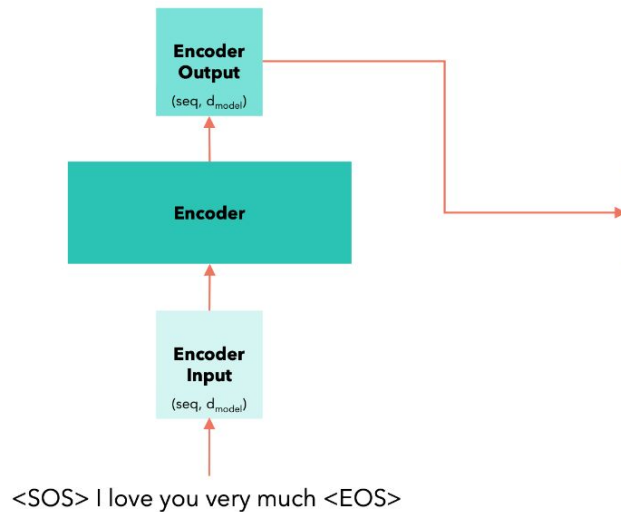


Ti amo molto

Training

Time Step = 1

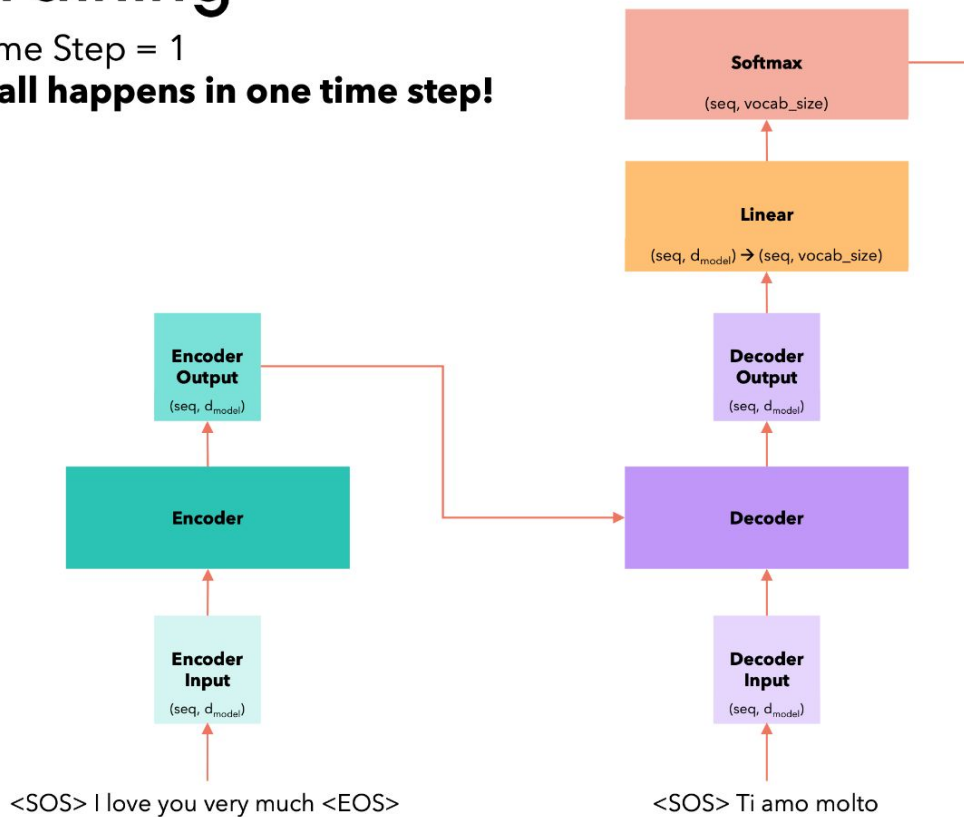
It all happens in one time step!



Training

Time Step = 1

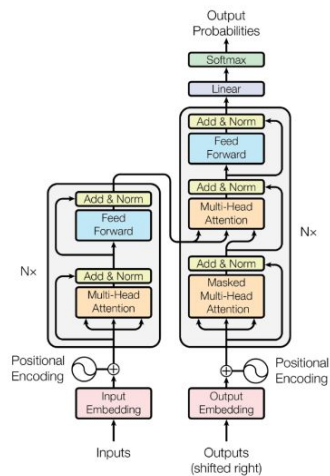
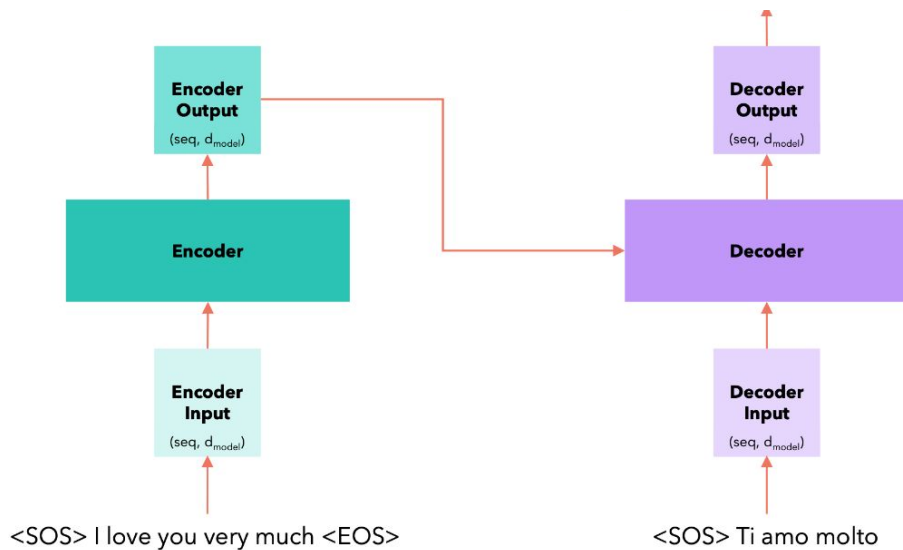
It all happens in one time step!



Training

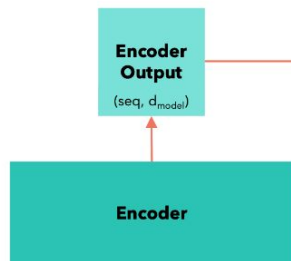
Time Step = 1

It all happens in one time step!



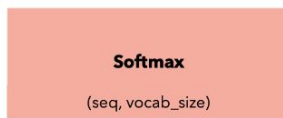
Training

<SOS> I love you very much <EOS>

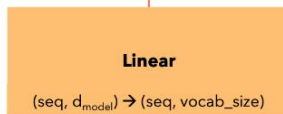


Encoder Input
(seq, d_{model})

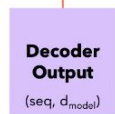
Encoder Output
(seq, d_{model})



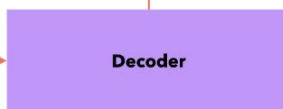
Softmax
(seq, vocab_size)



Linear
(seq, d_{model}) \rightarrow (seq, vocab_size)



Decoder Output
(seq, d_{model})



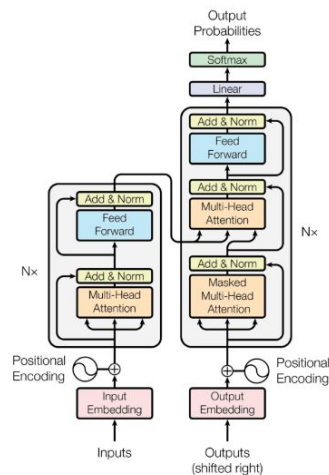
Decoder Input
(seq, d_{model})

<SOS> Ti amo molto

Ti amo molto <EOS>

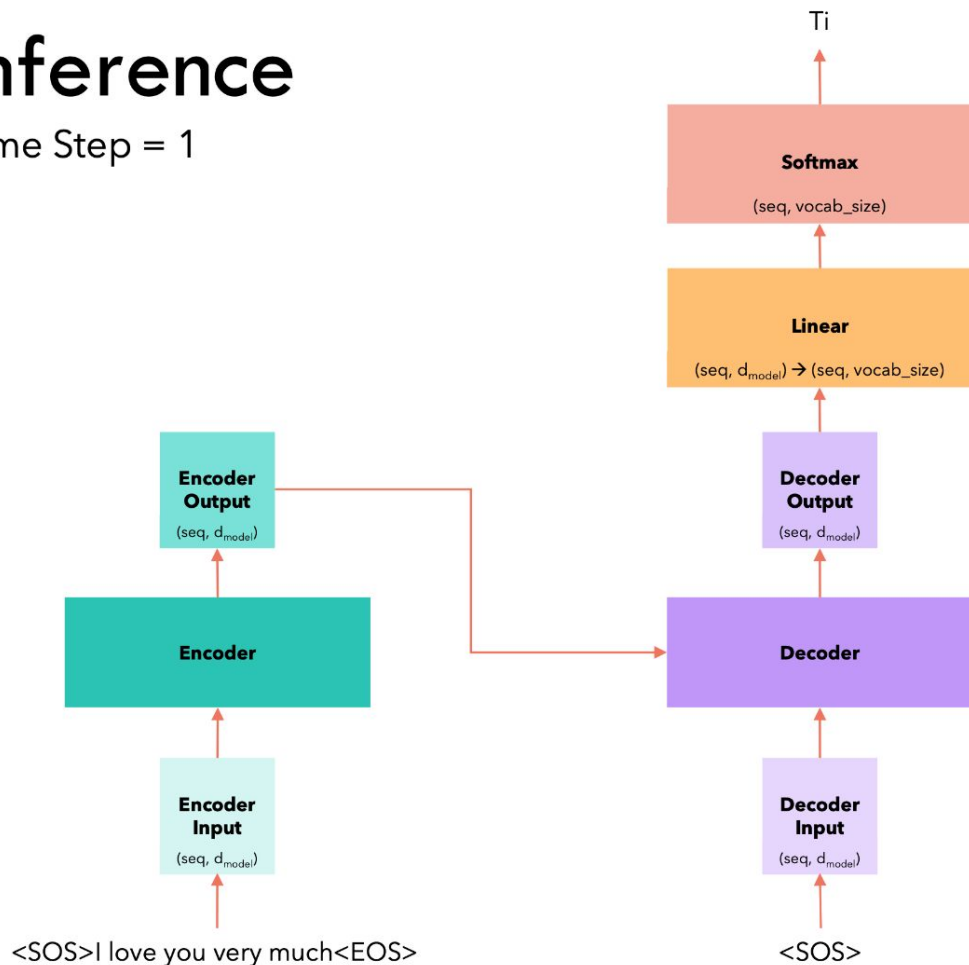
* This is called the "label" or the "target"

Cross Entropy Loss



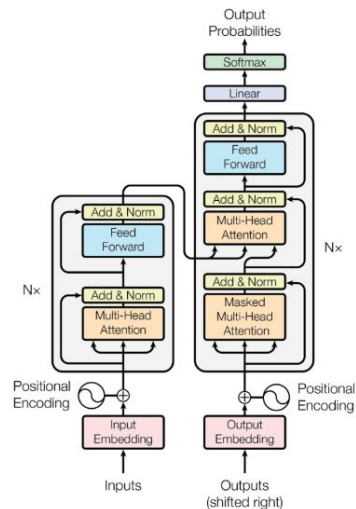
Inference

Time Step = 1



We select a token from the vocabulary corresponding to the position of the token with the maximum value.

The output of the last layer is commonly known as **logits**



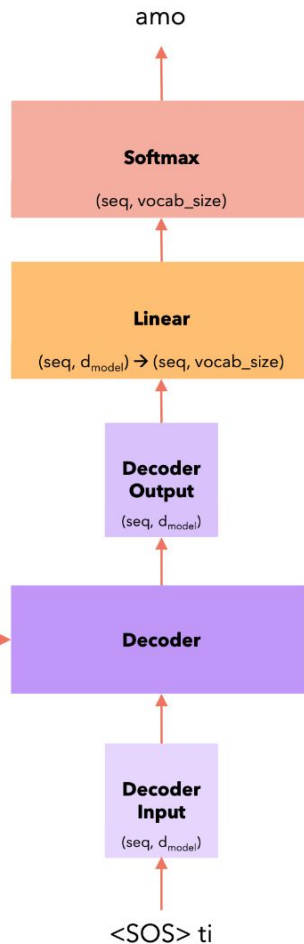
* Both sequences will have same length thanks to padding

Inference

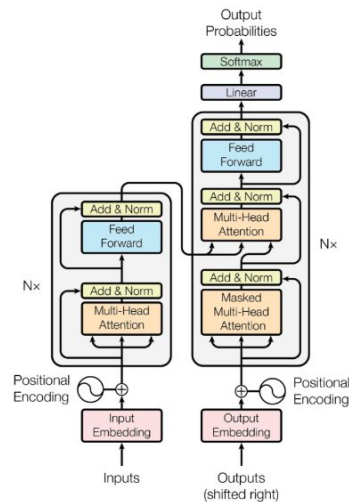
Time Step = 2

Use the encoder output from the first time step

<SOS>I love you very much<EOS>

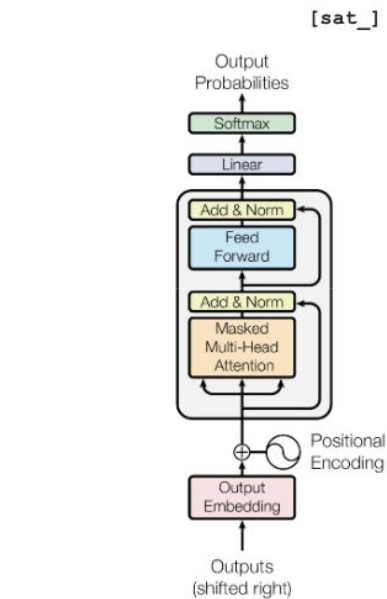


Since decoder input now contains **two** tokens, we select the softmax corresponding to the second token.



Append the previously output word to the decoder input

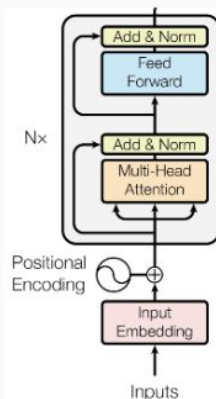
Decoder-only GPT



[START] [The_] [cat_]

Encoder-only BERT

[*] [*] [sat_] [*] [the_] [*]



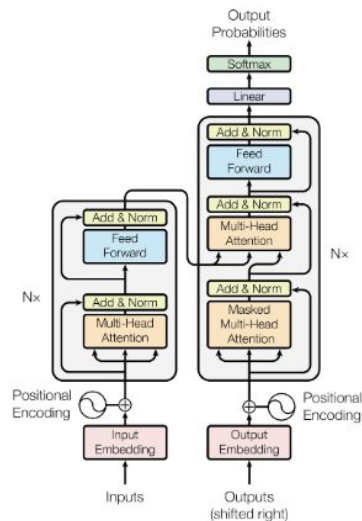
[The_] [cat_] [MASK] [on_] [MASK] [mat_]

Enc-Dec T5

Das ist gut.

A storm in Attala caused 6 victims.

This is not toxic.



Translate EN-DE: This is good.

Summarize: state authorities dispatched...

Is this toxic: You look beautiful today!