

What is an input embedding?



What is positional encoding?

- We want each word to carry some information about its position in the sentence.
- We want the model to treat words that appear close to each other as “close” and words that are distant” as “distant”.
- We want the positional encoding to represent a pattern that can be learned by the model.
- With sinusoidal positional encoding, we build one fixed table of vectors for positions 0...max sequence length, and every token simply uses the vector of its position, so all sentences share the same position table independent of vocabulary size.

What is positional encoding?

Original sentence

YOUR

CAT

IS

A

LOVELY

CAT

Embedding

(vector of size 512)

9551.617
3473.660
1372.660
—
1.680
2211.029

773.411
7778.580
9772.379
—
1678.523
6005.173

638.680
1867.330
51.860
—
1678.290
6005.685

778.580
1867.330
51.860
—
7774.116
5772.830

778.878
3778.580
1972.779
—
2002.210
3778.839

773.411
7778.580
9772.379
—
2603.617
5772.613

+

+

+

+

+

+

Position Embedding

(vector of size 512).
Only computed once
and reused for every
sentence during
training and inference.

—
—
—
—
—
—
—

1544.580
6888.830
2108.590
—
9385.603
2128.115

—
—
—
—
—
—
—

—
—
—
—
—
—
—

—
—
—
—
—
—
—

1061.610
7778.790
774.116
2671.580
1686.237
7118.680

+

+

+

+

+

+

Encoder Input

(vector of size 512)

—
—
—
—
—
—
—

7556.679
31611.513
10570.390
—
13079.670
13520.637

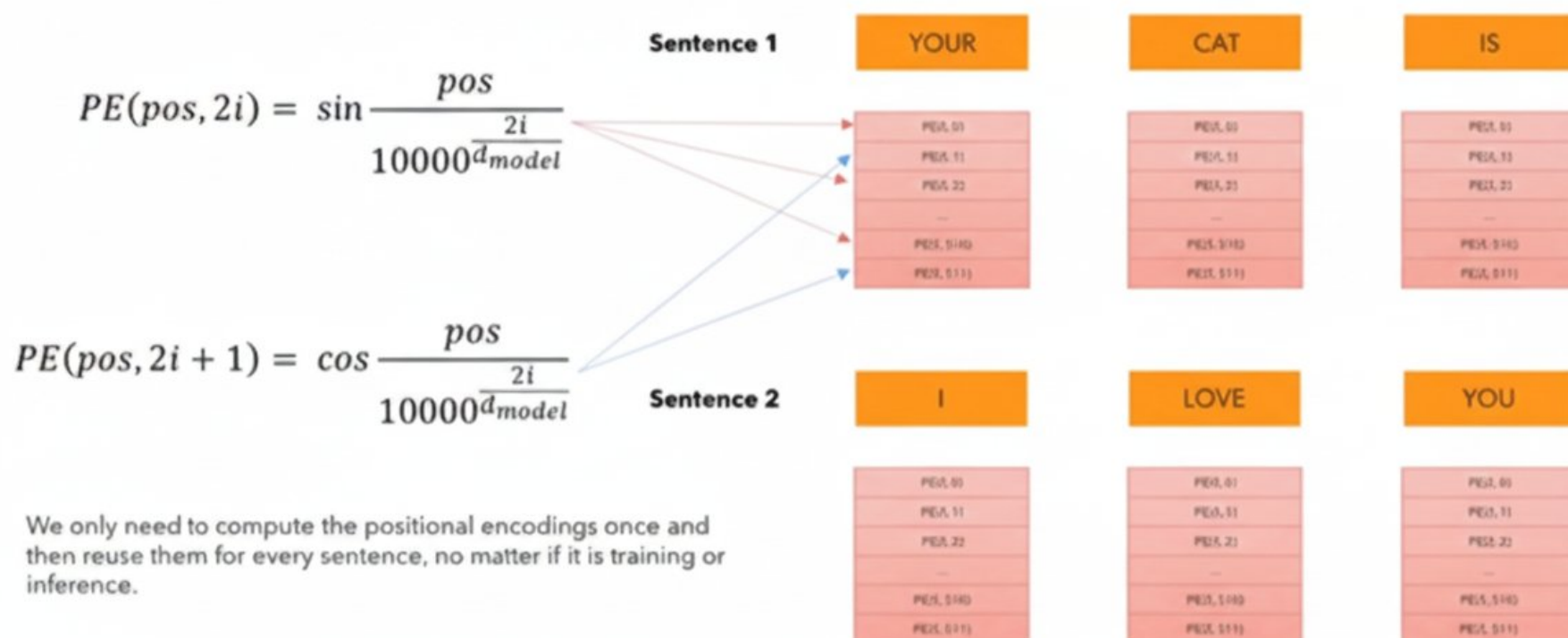
—
—
—
—
—
—
—

—
—
—
—
—
—
—

—
—
—
—
—
—
—

1622.602
17379.30
11510.399
—
1078.617
773.637

What is positional encoding?



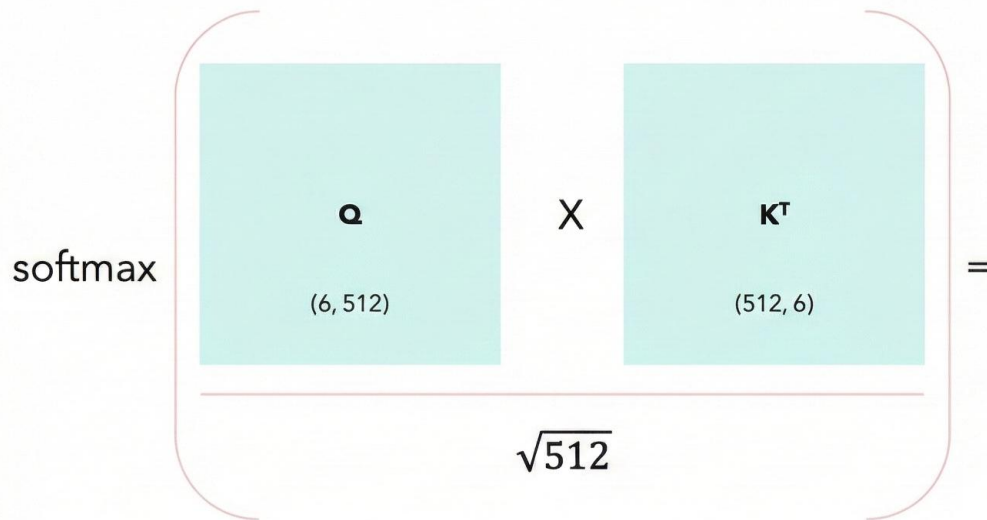
What is Self-Attention?

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	0.266	0.119	0.134	0.148	0.179	0.152	1
CAT	0.124	0.278	0.201	0.128	0.145	0.115	1
IS	0.147	0.132	0.392	0.281	0.218	0.145	1
A	0.210	0.106	0.206	0.193	0.119	0.125	1
LOVELY	0.146	0.146	0.142	0.143	0.227	0.219	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 $d_{\text{model}} = d_v$.

(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	0.266	<u>0.119</u>	<u>0.134</u>	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	0.193	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

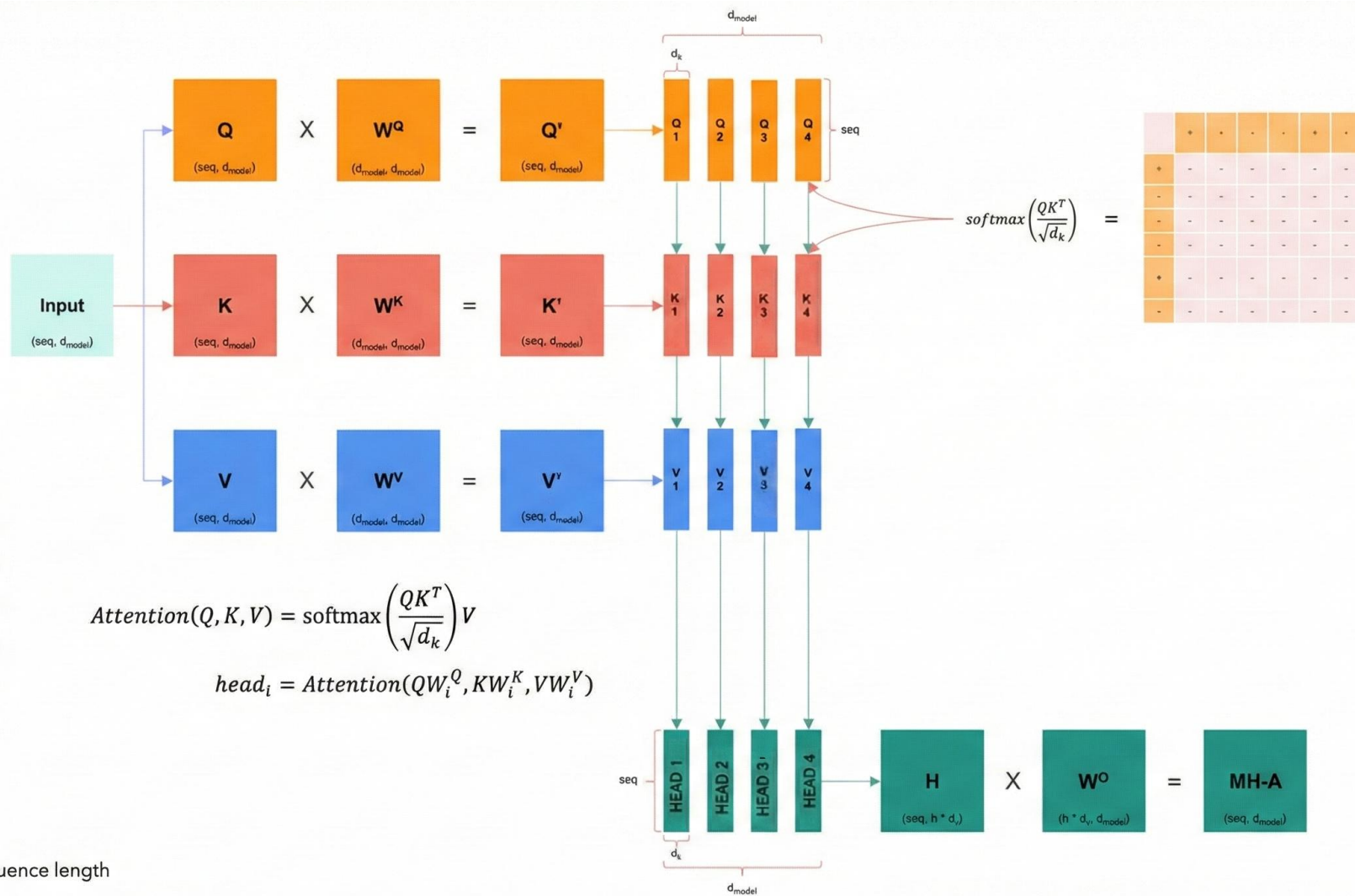
(6, 512)

=

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.



seq = sequence length

d_{model} = size of the embedding vector

h = number of heads

d_k, d_v = d_{model} / h

What is layer normalization?

Batch of 3 items

ITEM 1

58.147
3014.823
...
...
9403.361
8.021

μ_1

σ_1^2

ITEM 2

1242.229
688.123
...
...
423.099
102.442

μ_2

σ_2^2

ITEM 3

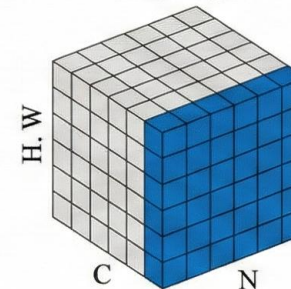
9.310
4664.673
...
...
944.705
21199.444

μ_3

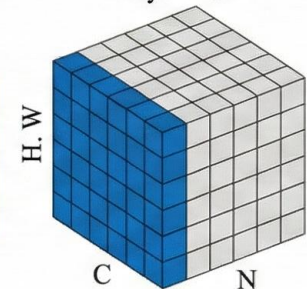
σ_3^2

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

Batch Norm



Layer Norm

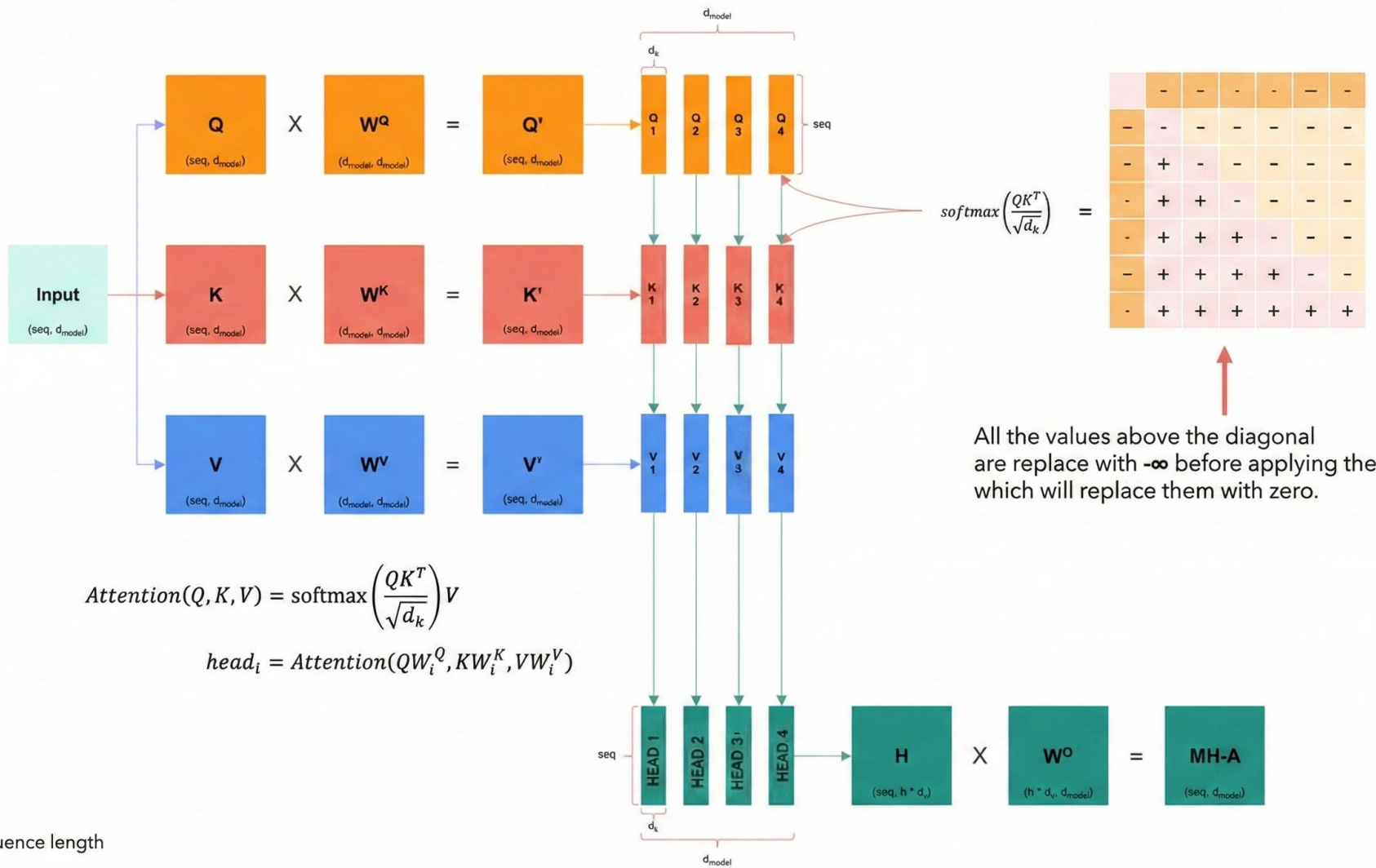


We also introduce two parameters, usually called **gamma** (multiplicative) and **beta** (additive) that introduce some fluctuations in the data, because maybe having all values between 0 and 1 may be too restrictive for the network. The network will learn to tune these two parameters to introduce fluctuations when necessary.

What is Masked Multi-Head Attention?

Our goal is to make the model causal: it means the output at a certain position can only depend on the words on the previous positions. The model **must not** be able to see future words.

	YOUR	CAT	IS	A	LOVELY	CAT
YOUR	0.268	0.316 $-\infty$	0.119	0.268	0.209	0.514
CAT	0.134	0.279	0.247	0.246	0.246	0.265
IS	0.147	0.132	0.132	0.282	0.272	0.111
A	0.210	0.128	0.206	0.212	0.279	0.264
LOVELY	0.146	0.152	0.132	0.140	0.227	0.212
CAT	0.195	0.114	0.203	0.103	0.157	0.229



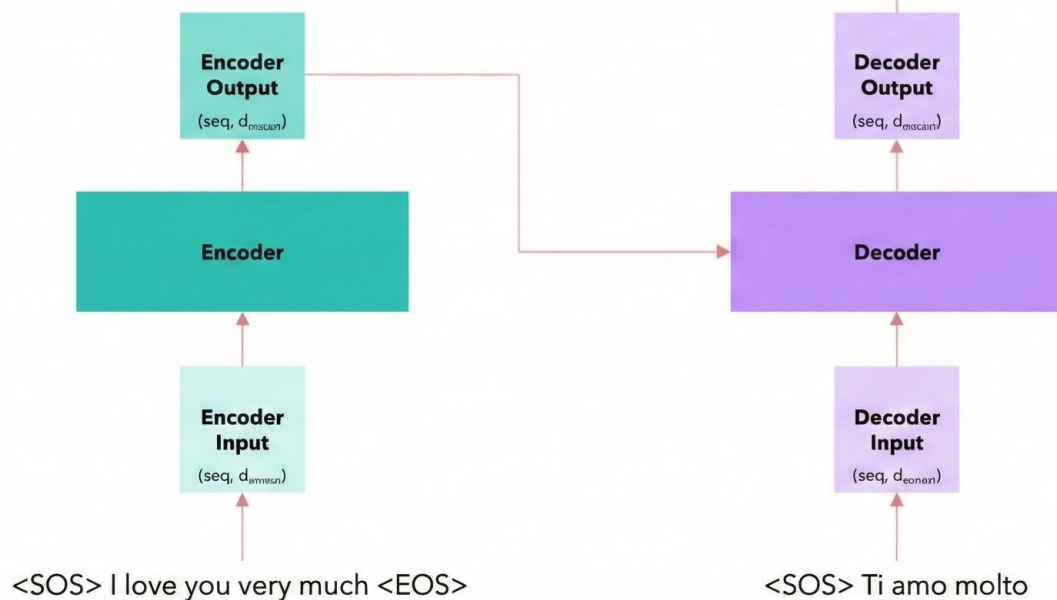
- seq = sequence length
- d_{model} = size of the embedding vector
- h = number of heads
- d_k, d_v = d_{model} / h

Training

Time Step = 1

It all happens in one time step!

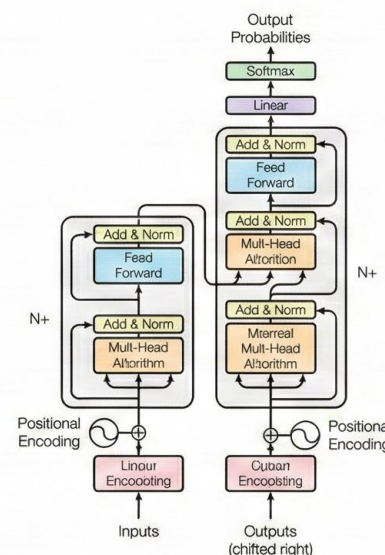
The encoder outputs, for each word a vector that not only captures its meaning (the embedding) or the position, but also its interaction with other words by means of the multi-head attention.



Ti amo molto <EOS>

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

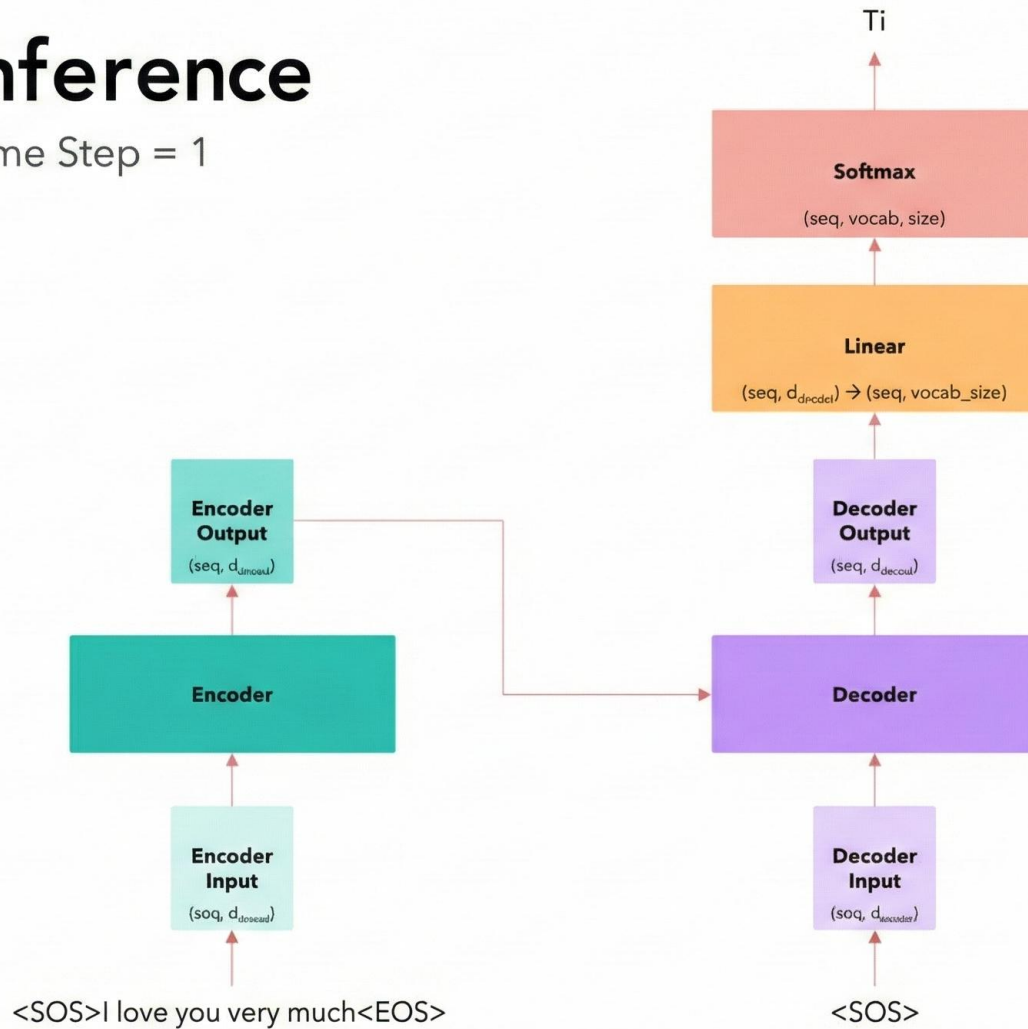
Cross Entropy Loss



We prepend the <SOS> token at the beginning. That's why the paper says that the decoder input is shifted right.

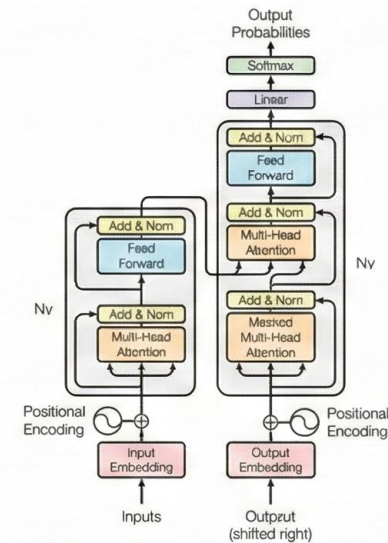
Inference

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**



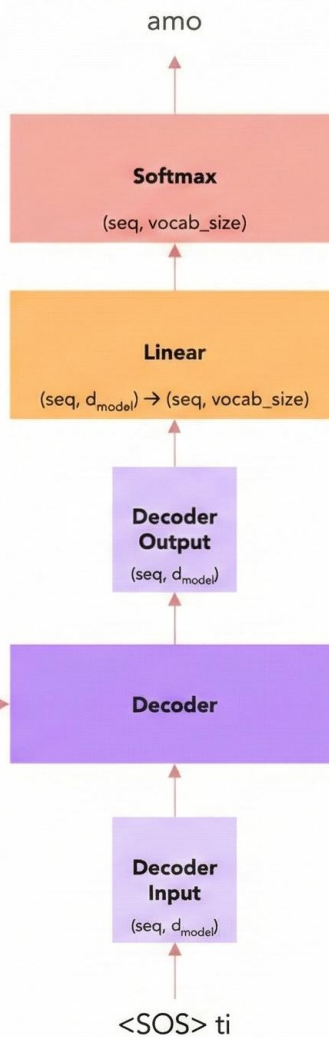
* Both sequences will have same length thanks to padding

Inference

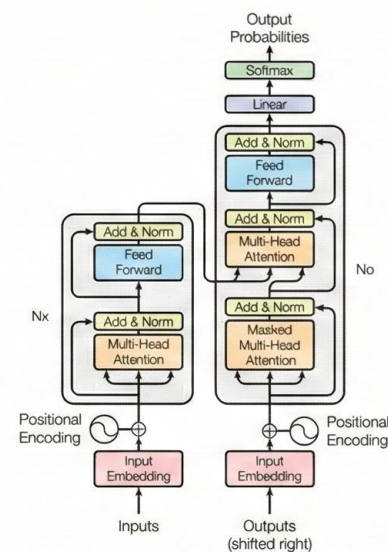
Time Step = 2

Use the encoder output from the first time step

<SOS>I love you very much<EOS>



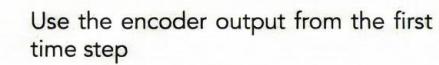
Since decoder input now contains **two** tokens, we select the softmax corresponding to the second token.



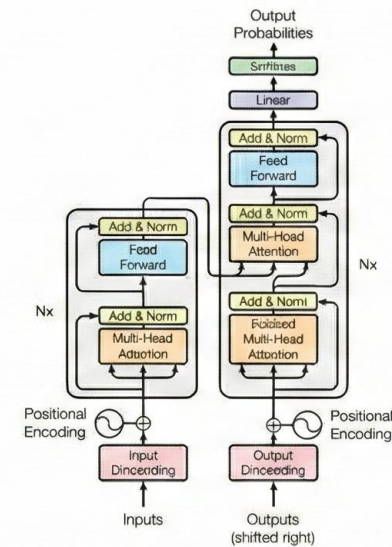
Append the previously output word to the decoder input

Time Step = 3

Time Step = 3



Append the previously output word to the decoder input

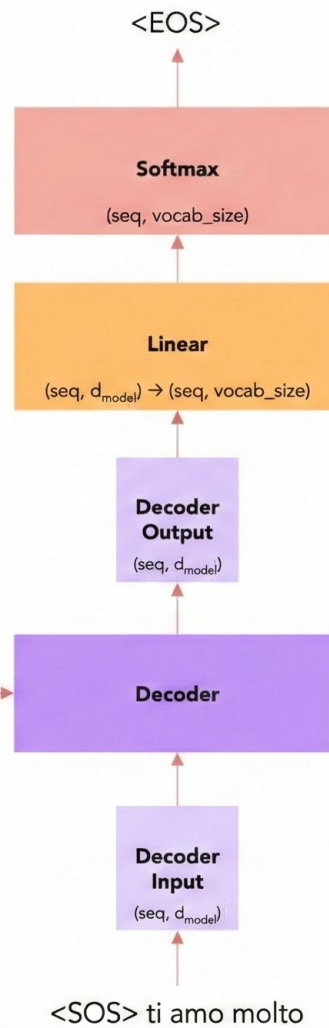


Inference

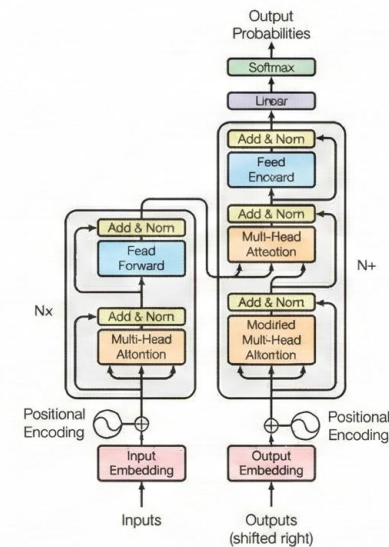
Time Step = 4

Use the encoder output from the first time step

<SOS>I love you very much<EOS>



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



Append the previously output word to the decoder input

Inference strategy

- We selected, at every step, the word with the maximum softmax value. This strategy is called **greedy** and usually does not perform very well.
- A better strategy is to select at each step the top B words and evaluate all the possible next words for each of them and at each step, keeping the top B most probable sequences. This is the **Beam Search** strategy and generally performs better.