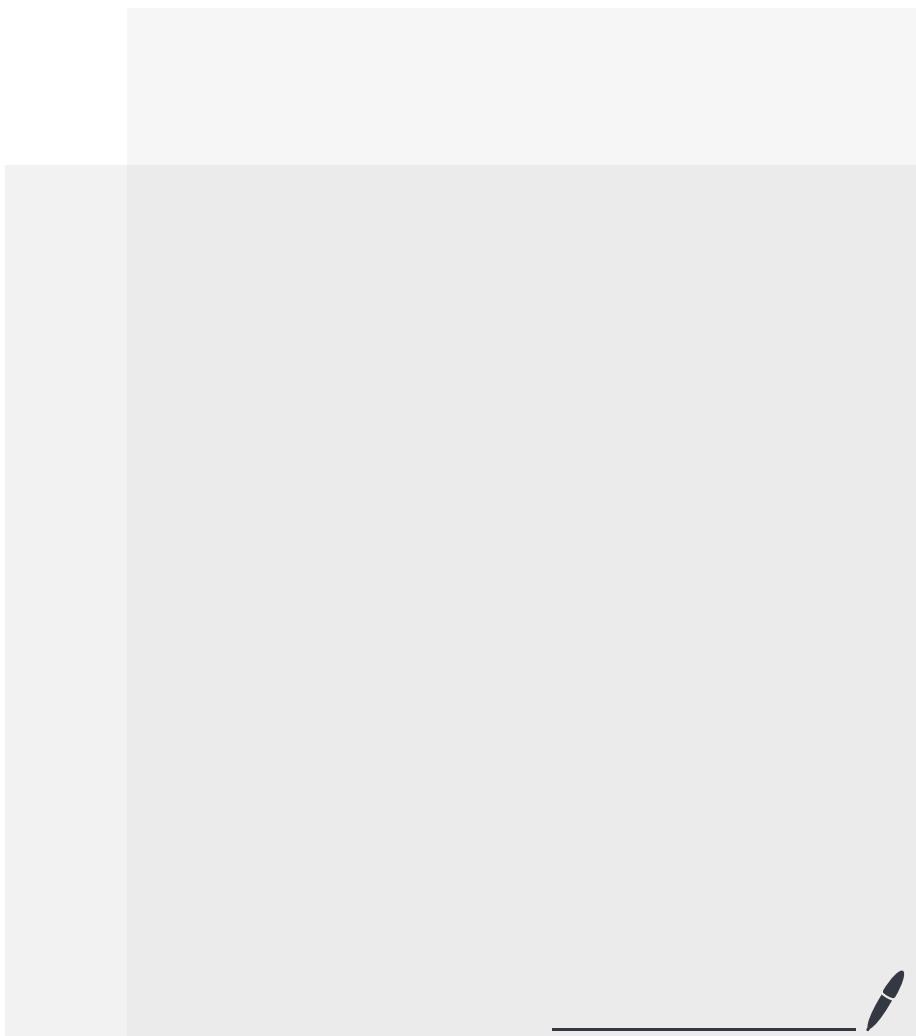


Attention is all you need

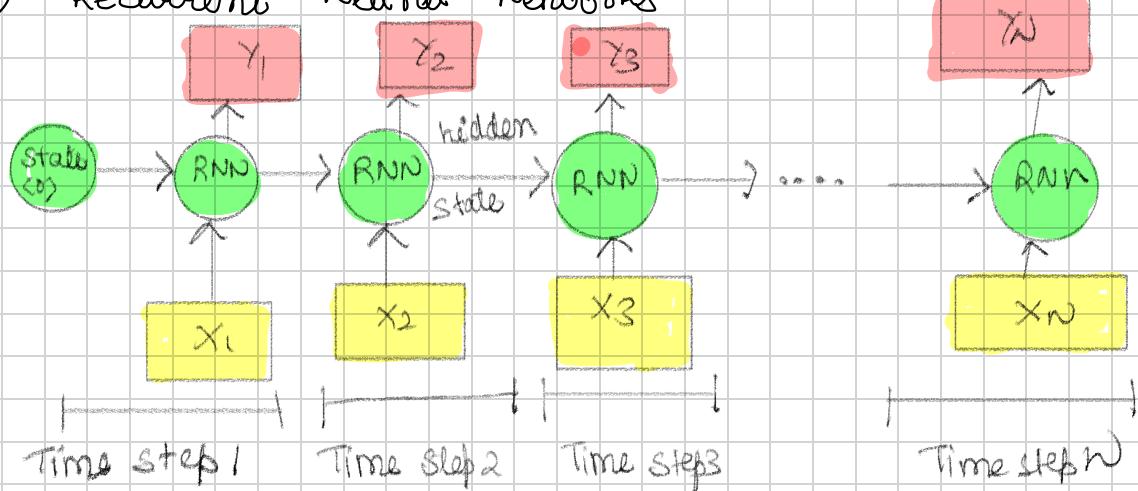


Attention is all you need!

GPT: Generatively Pre-Trained Transformers

Recurrent Neural Network \longrightarrow Transformers

i) Recurrent Neural Networks



one series of input to get one series of output

- i) slow computation for long sequences
- ii) vanishing or exploding gradients

$$\begin{array}{c} x \\ y \end{array} \rightarrow f(x,y) = x \cdot y \rightarrow z \rightarrow g(z) = z^2$$

$$\frac{dg}{dz} = \frac{dg}{df} \cdot \frac{df}{dx}$$

longer this multiplication chain
longer computation time.

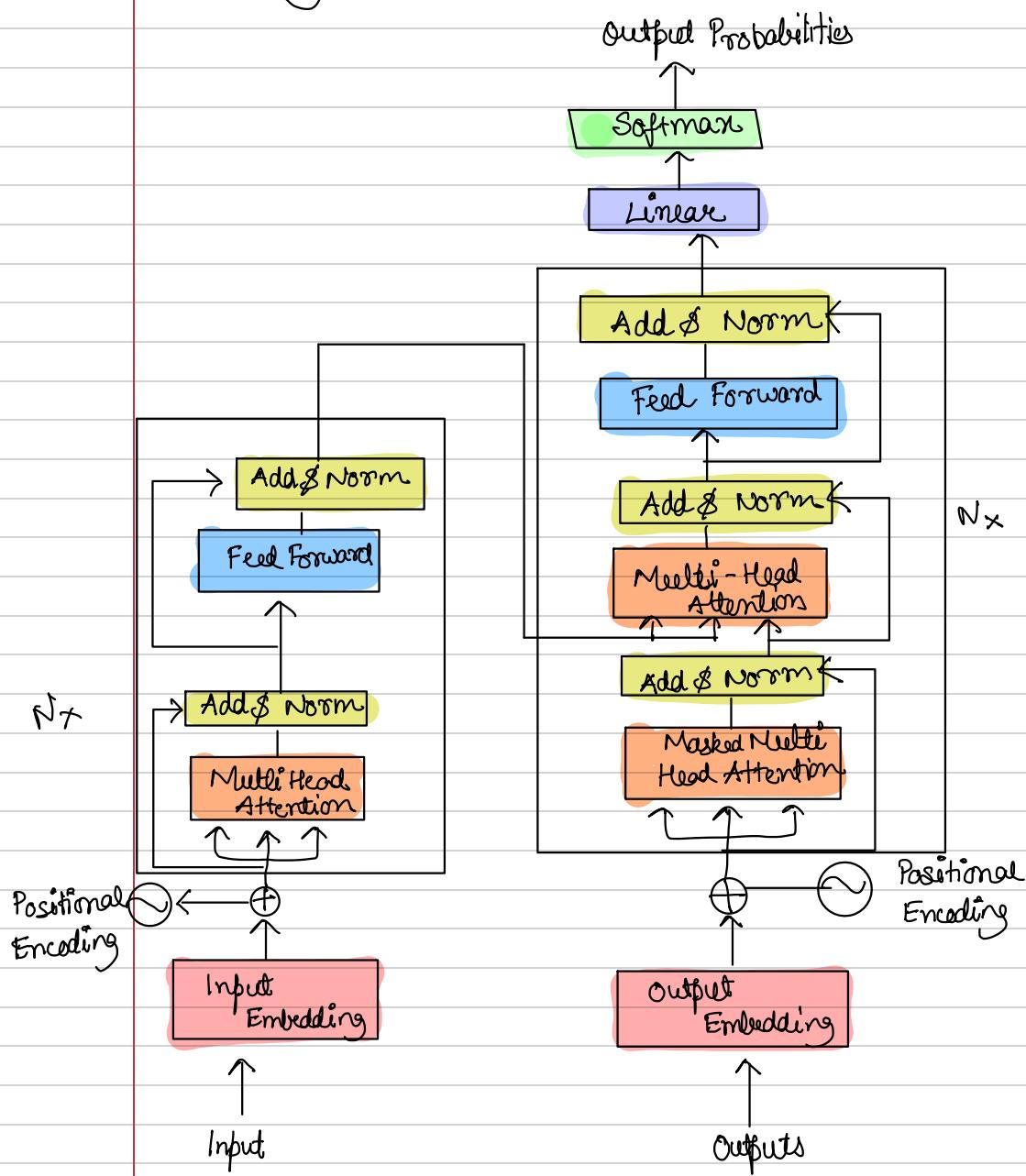
it'll either be a very large number or a very small number.

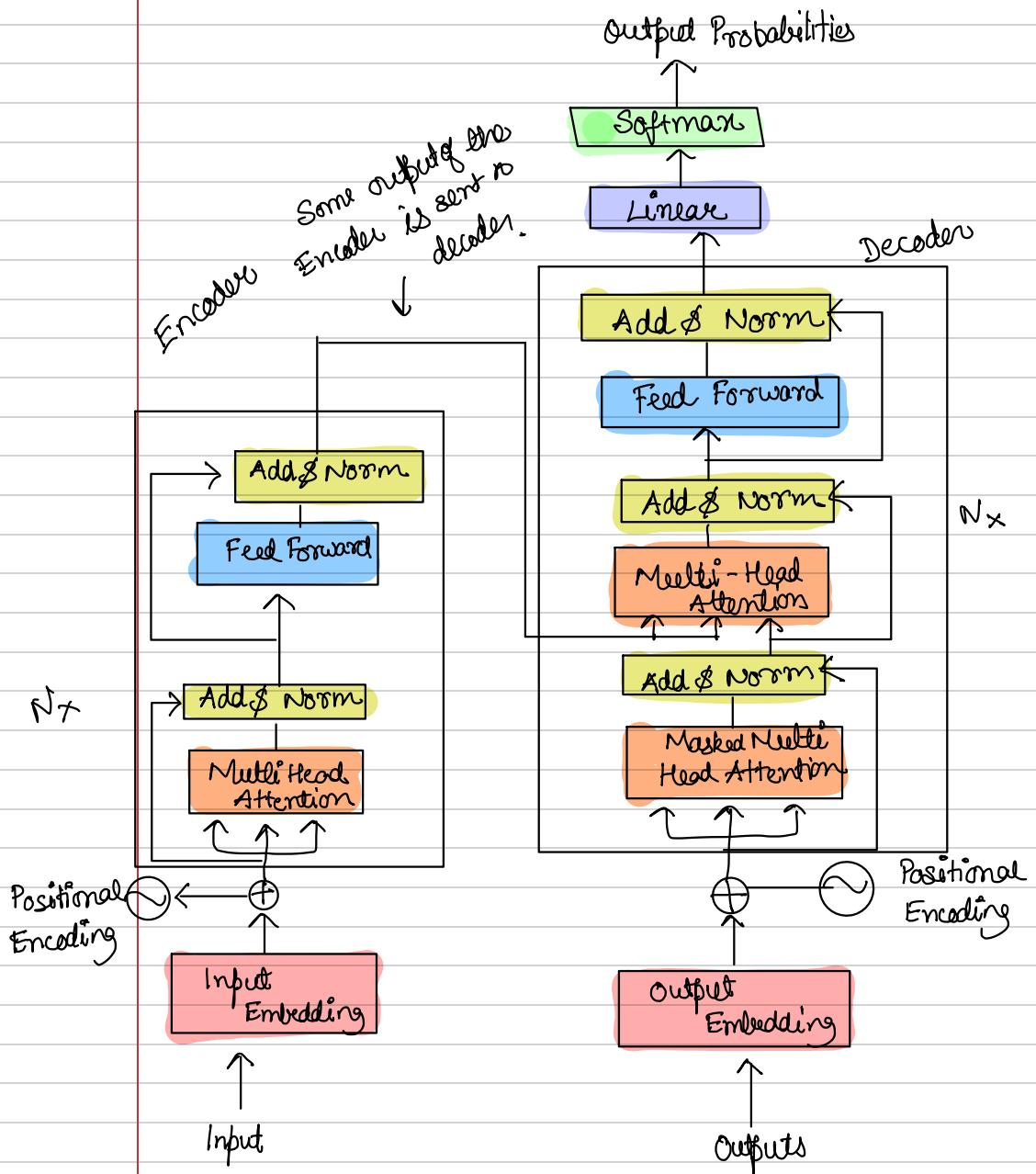
vanishing gradient

exploding gradient

iii) Difficulty in accessing information from long terms ago.

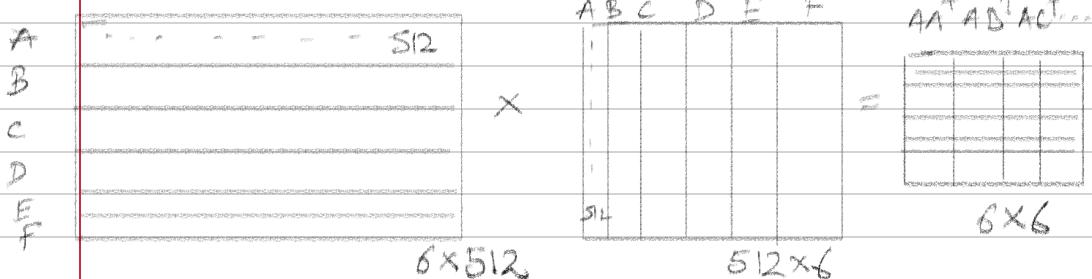
Introducing the Transformer





Notations

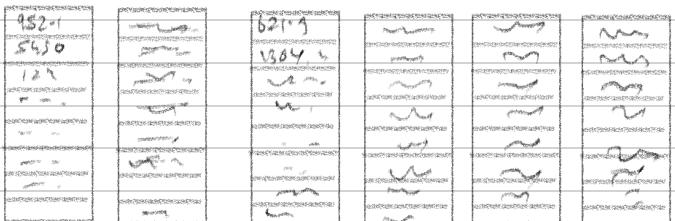
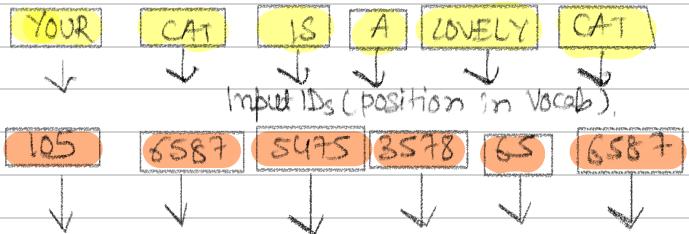
Input Matrix (Sequence, d_{model})



ENCODER

What is an Input Embedding?

Original Sentence (tokens)

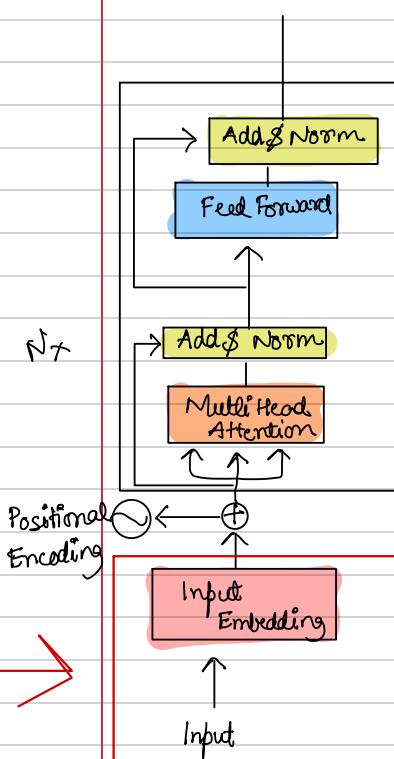


Embedding
(vectors of size d_{model})

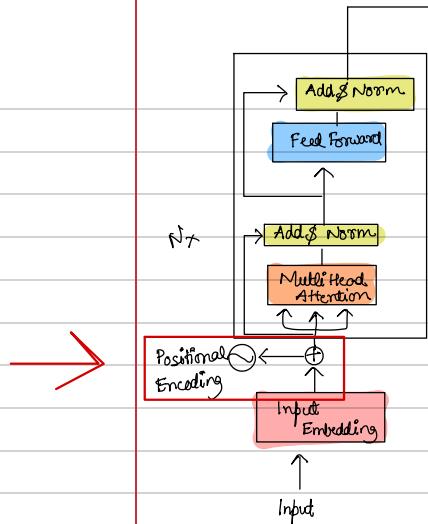
Always map the same word to same embeddings.

LEARN THOSE PARAMETERS.

$$d_{\text{model}} = 512$$



What is Positional Encoding?



- we want each word to carry some information about its position in the sentence.
- we want models 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.

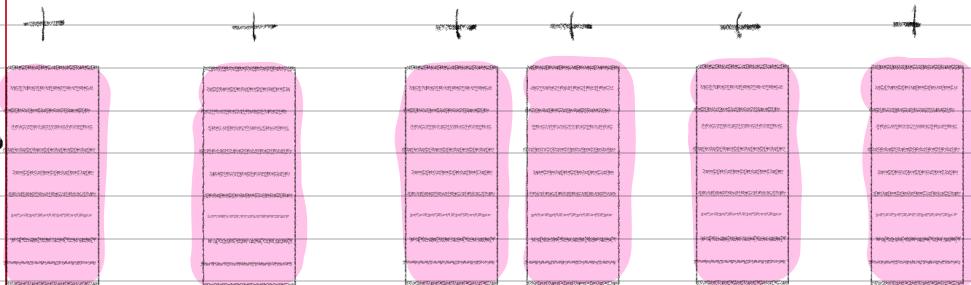
Original Sentence

YOUR CAT IS A LOVELY CAF

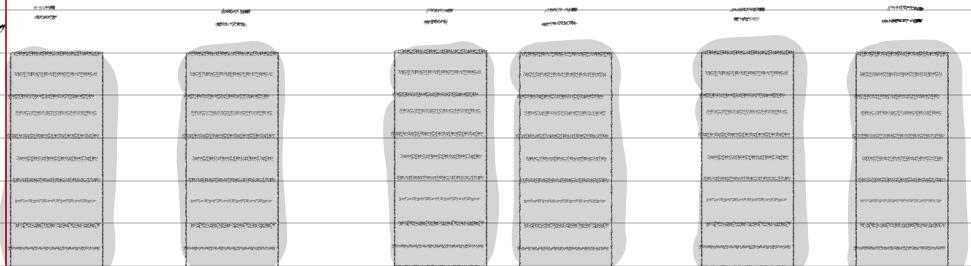
LR →



Embedding vector
($S1 \times 1$)



Encoder input
(Vectors of $S1 \times S1$)

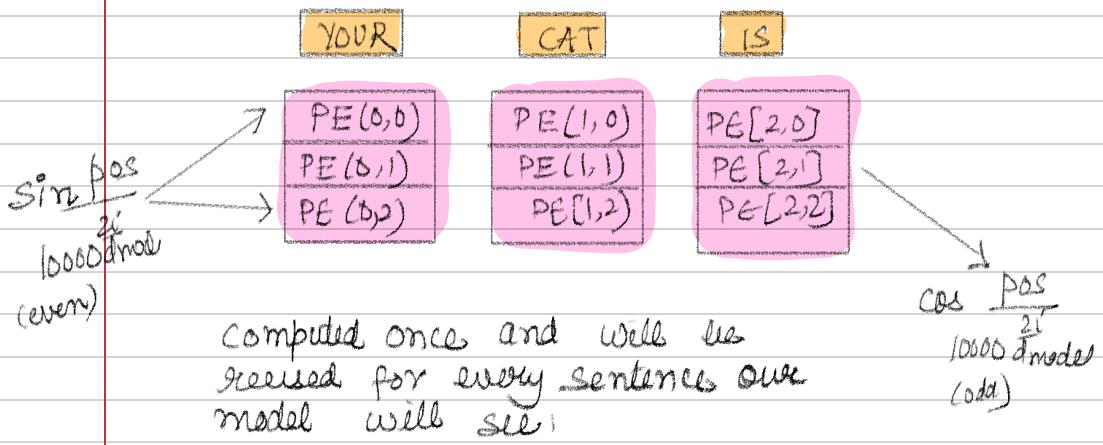


How are Positional Embeddings calculated?

Embedding vector is something that we learn, while the positional Embedding remain fixed.

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

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



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

Why trigonometric functions?

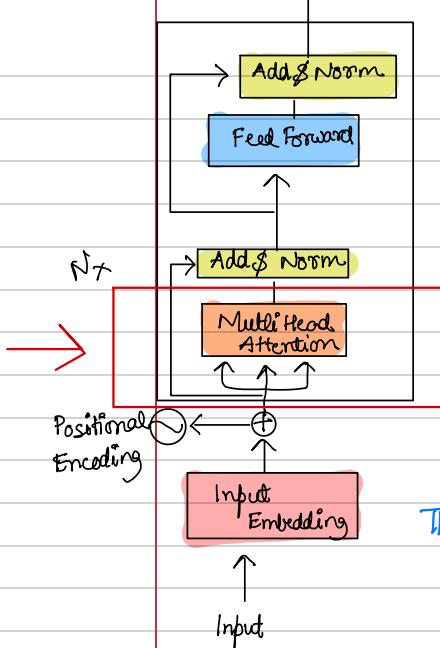
Trigonometric functions like cosine and sine natively represent a patterns that the model can recognize as continuous, so relative positions are easier to see for the model.

Multi-Head Attention

Let's first understand the single head version.

What is Self Attention?

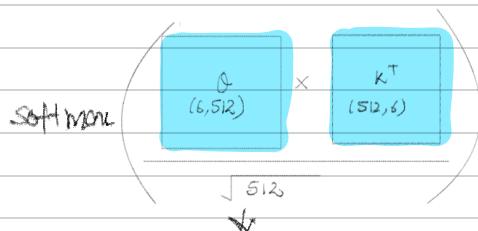
Self Attention allows the model to relate words to each other.



$$\text{Attention } (\theta, k, v) = \text{softmax} \left(\frac{\theta k^T}{\sqrt{d_k}} \right) v$$

Let's assume we have seq=6, $d_{model} = 512$

The matrices θ , k and v are just the input sentence.



	YOUR	CAT	IS	A	LOVELY	CAT	Σ
YOUR	0.268	0.119	0.134	0.145	0.139	0.152	1
CAT	0.124	0.278	0.201	0.125	0.154	0.115	1
IS	0.147	0.132	0.262	0.097	0.218	0.145	1
A	0.210	0.128	0.246	0.212	0.119	0.125	1
LOVELY	0.146	0.158	0.152	0.163	0.223	0.134	1
CAT	0.195	0.114	0.203	0.103	0.153	0.229	1

6x6

	YOUR	CAT	IS	A	LOVELY	CAT
YOUR	0.268	0.119	0.134	0.148	0.179	0.152
CAT	0.124	0.231	0.201	0.128	0.154	0.115
IS	0.147	0.152	0.262	0.097	0.218	0.105
A	0.210	0.128	0.286	0.212	0.119	0.125
LOVELY	0.146	0.158	0.132	0.163	0.227	0.174
CAT	0.195	0.114	0.203	0.103	0.157	0.229

6x6

$$\times \begin{matrix} \checkmark \\ (6, 512) \end{matrix}$$

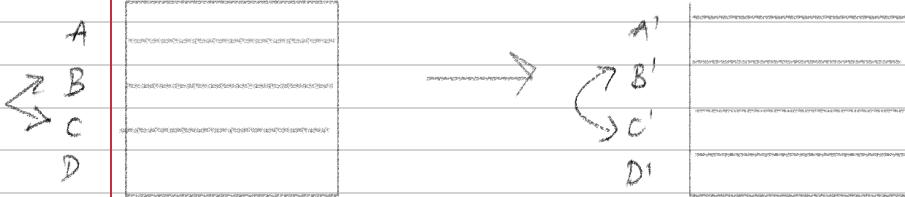
Attention

(6, 512)

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

Self Attention in detail

- Self Attention is permutation invariant, if we don't consider the contribution of the positional encodings -



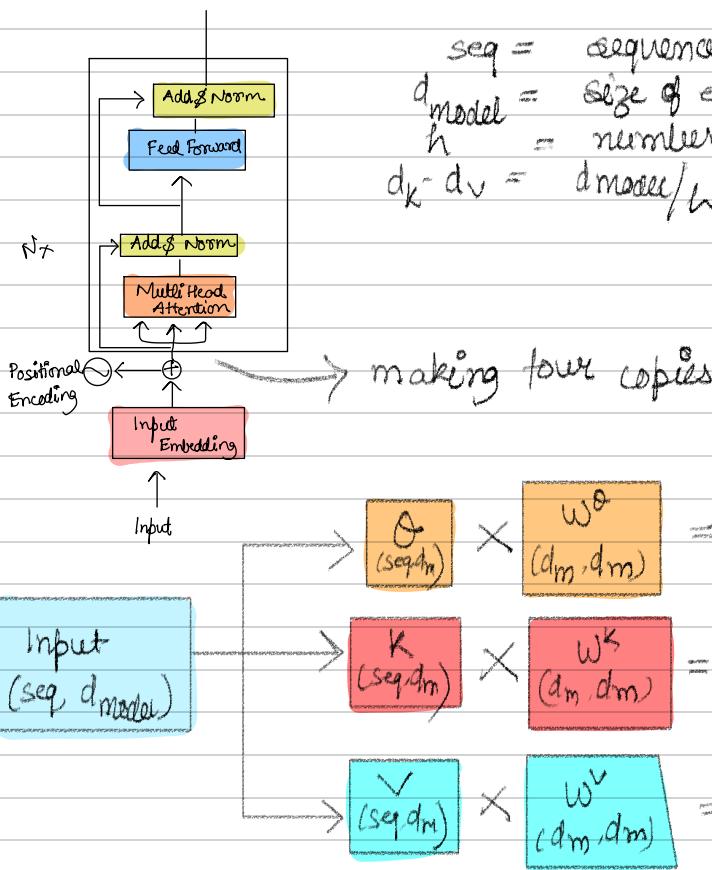
- Self Attention requires no parameters. Up to now the interaction between words has been driven by their embedding and the positional encoding. will change later.
- We expect values along the diagonal to be the highest.
- If we don't want some positions to interact, we can always set their values to $-\infty$ before applying the softmax in this matrix and the model will not learn those interactions.

Multi - Head Attention

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

$$\text{Multi Head}(Q, K, V) = \text{concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

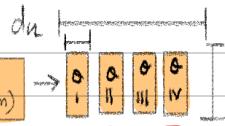
$$\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$



$$d_K = d_{\text{model}}/h$$

$$\begin{aligned} \text{seq} &= \text{sequence length} \\ d_{\text{model}} &= \text{size of embedding vector} \\ h &= \text{number of heads} \\ d_K = d_V &= d_{\text{model}}/h \end{aligned}$$

d_{model}



Input
(seq, d_{model})

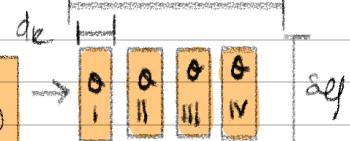
seq

$$\Theta \times W^Q = Q \quad (seq, d_m) \rightarrow Q_1, Q_2, Q_3, Q_4$$

$$K \times W^K = K \quad (seq, d_m) \rightarrow K_1, K_2, K_3, K_4$$

$$V \times W^V = V \quad (seq, d_m) \rightarrow V_1, V_2, V_3, V_4$$

d_{model}



Input
(seq, d_{model})

seq

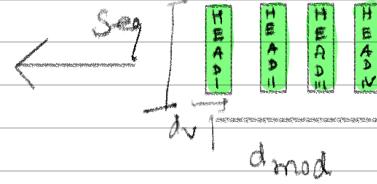
$$\Theta \times W^Q = Q \quad (seq, d_m) \rightarrow Q_1, Q_2, Q_3, Q_4$$

$$K \times W^K = K \quad (seq, d_m) \rightarrow K_1, K_2, K_3, K_4$$

$$V \times W^V = V \quad (seq, d_m) \rightarrow V_1, V_2, V_3, V_4$$

Multi Head
(Θ, K, V)

H
(seq, $h \times d_v$)



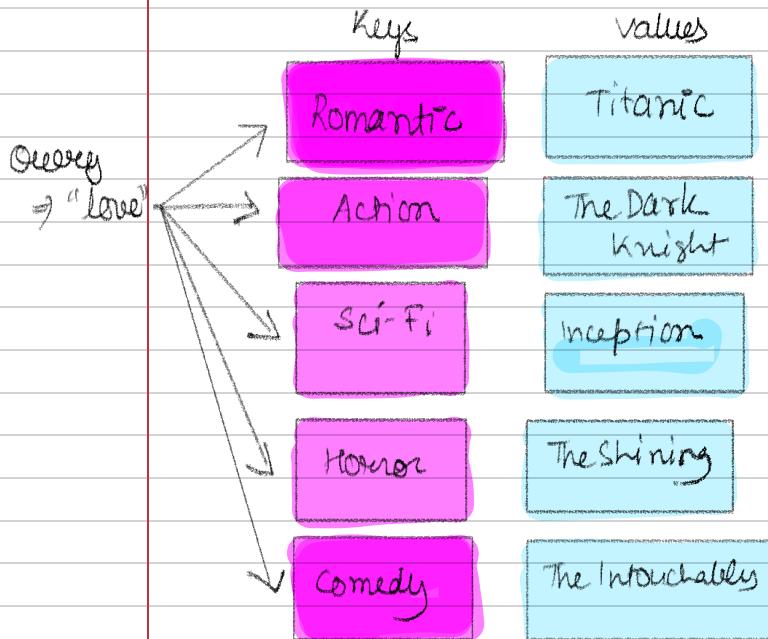
$MH - A$
(seq, d_{model})

W^O
($h \times d_v, d_{model}$)

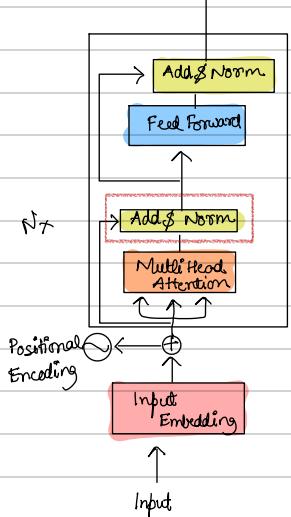
\times

Why query, key and values?

The internet says that these terms come from the database terminology or the Python-like dictionaries.

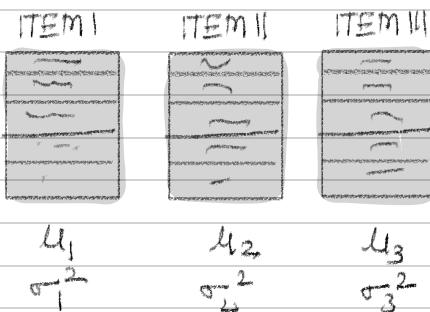


Encoders



what is layer normalization?

Batch of 3 items

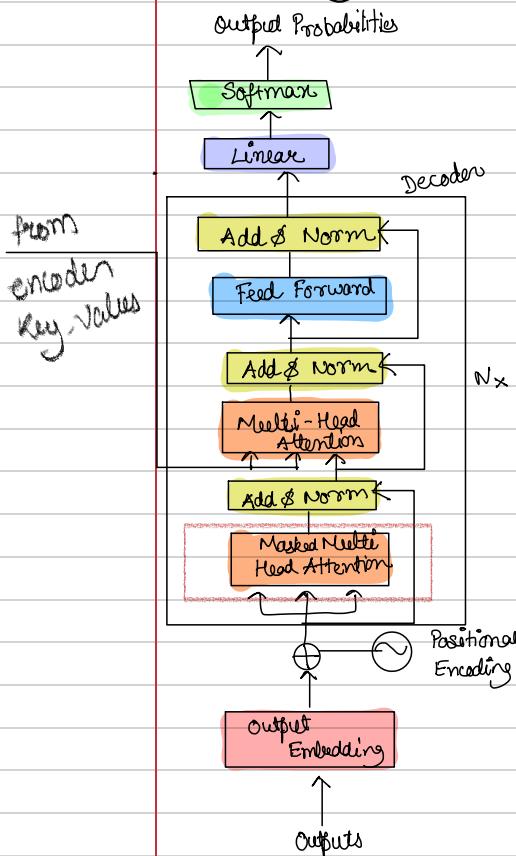


and then replace each value with

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

We also introduce two parameters usually called $\gamma(x)$ and $\beta(t)$ that introduce some fluctuation in the data, because having all the values b/w 0 and 1 may be too restrictive for the network.

DECODER



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	0.124	0.147	0.210	0.146	0.195
CAT	0.124	0.271	0.261	0.218	0.154	0.155
IS	0.147	0.261	0.282	0.218	0.2048	0.145
A	0.210	0.218	0.286	0.212	0.149	0.145
LOVELY	0.146	0.158	0.152	0.163	0.227	0.124
CAT	0.195	0.114	0.203	0.103	0.157	0.229

6x6

I want the words to only see the words before them.

Replace them with $-\infty$

INFERENCE AND TRAINING OF A TRANSFORMER

Training

I love you very much → Ti amo molto



































