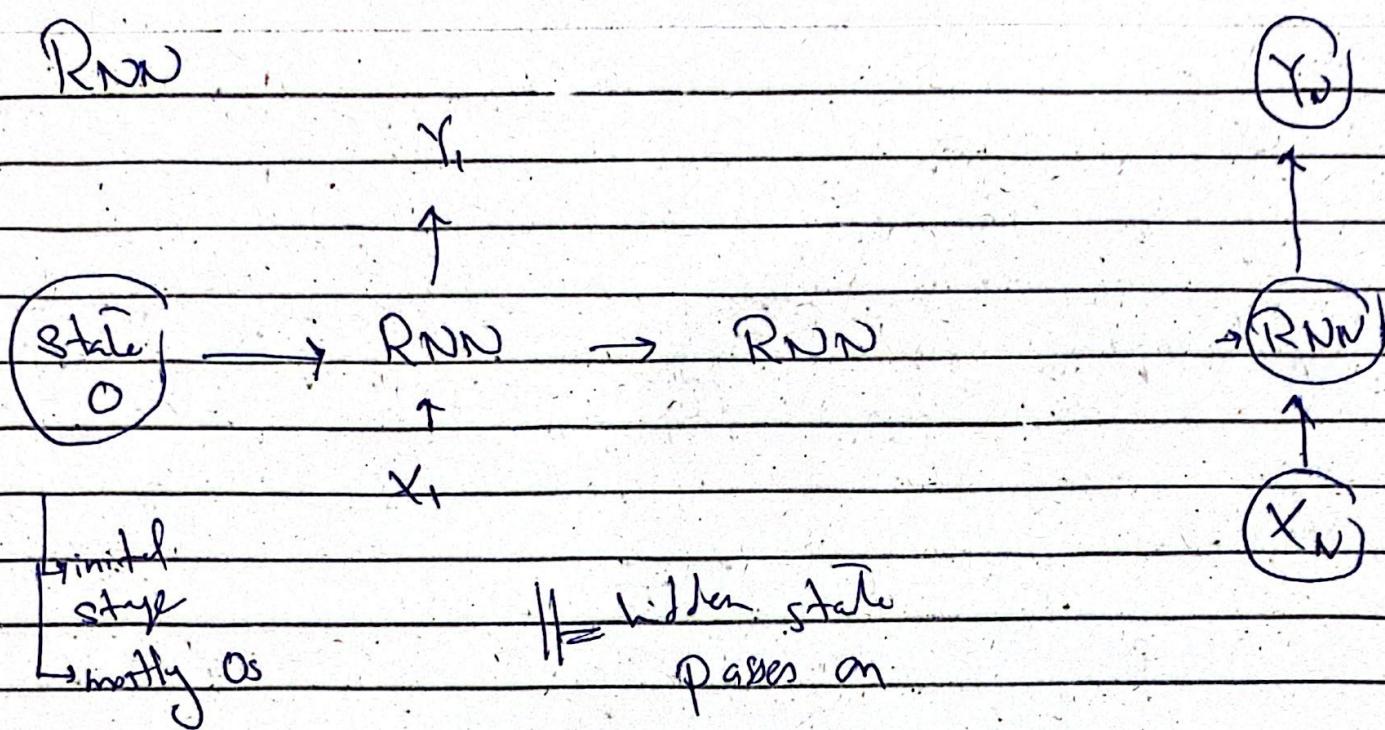


TRANSFORMERS:

RNN



RNN Limitations

- bad at remembering long sequence
- unable to vanish gradient

→ They learn the sequence with time steps

In order to get result for timestamp t , you'll have to wait for full network to finish for $t-1$, that's why it's computationally inefficient. So no parallelization.

Transformer solved this problem by using attention-head.

Why transformer are better

input

chatbot



They use self-attention

output



input → Encoder

Step 1: Tokenize the sentence

2. Feed it to embedding layer (Encoder)

input: 2339 19081 2024 2488

the LSTM
2084 100

after feed to embedding layer, embedding layer will convert each token into vector representation.

It would generate 6, 5×1 vectors.

~~Step 3~~ the output from embedding layer will go through positional encoding

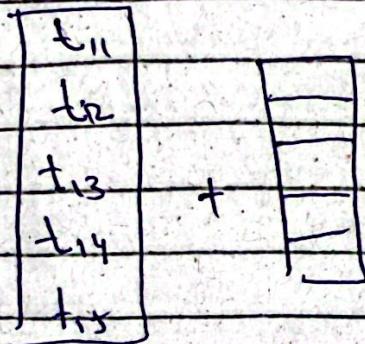
↳ If we replace 'LSTM' and 'transformer'

"Why LSTM are better than Transformers.."

Whole meaning has changed but for model both sentences are equal so it is important to maintain positional information (sentence order).

encoding

Positional ~~index~~ adds positional vector into embedding vector which holds the info about position of each word.

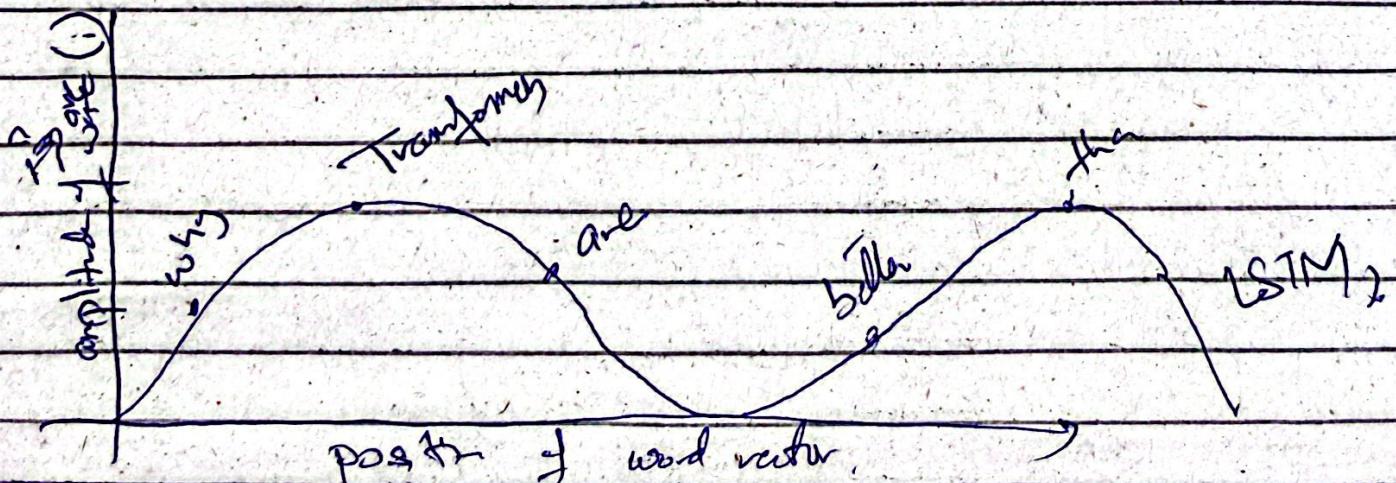


for each token

Embedding + Position vector

Researchers used sine wave to compute position vector by using the

$$PE_{(pos, d)} = \sin\left(\frac{pos}{100^{2/d_{model}}}\right)$$



[wall]

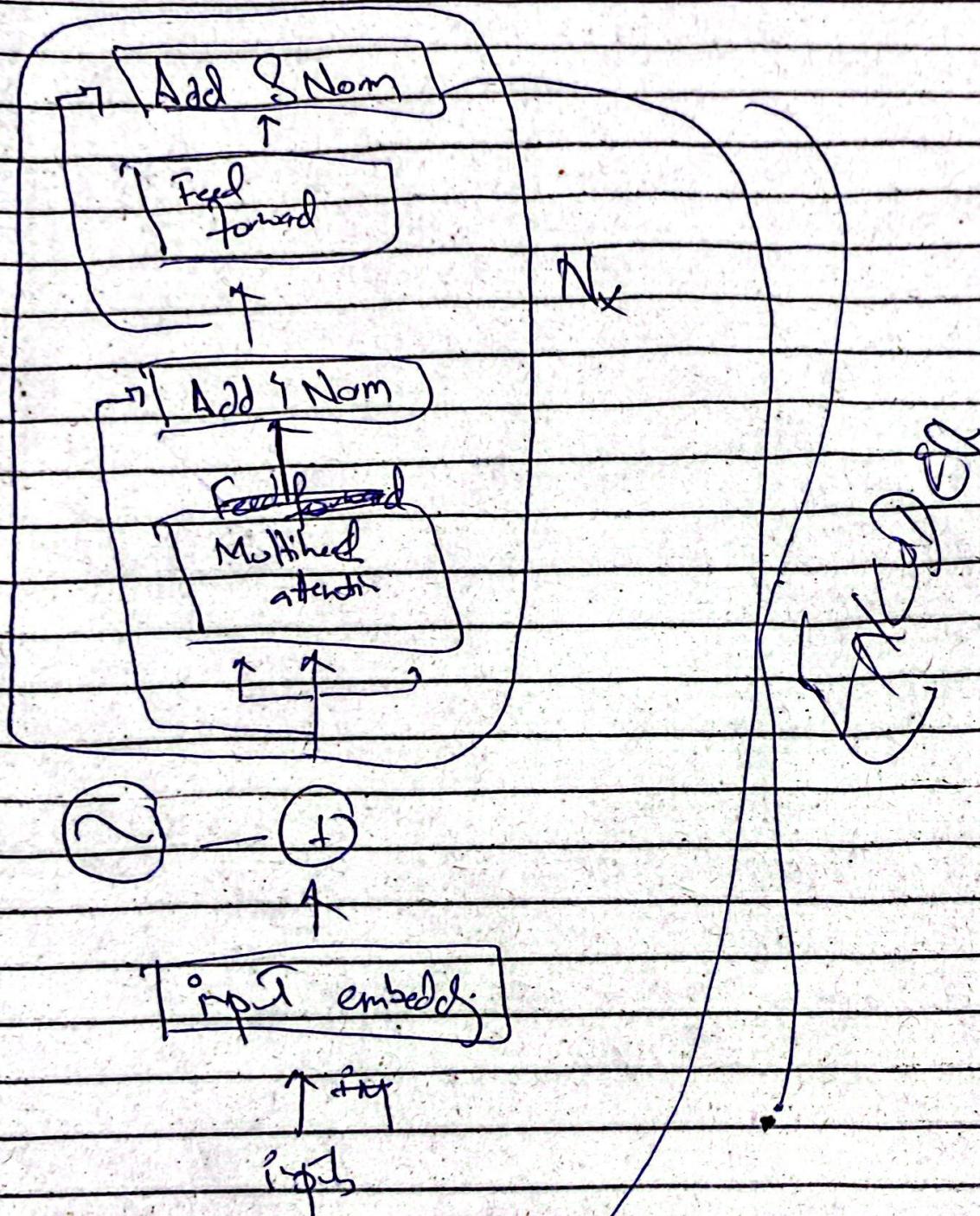
vanishing gradient

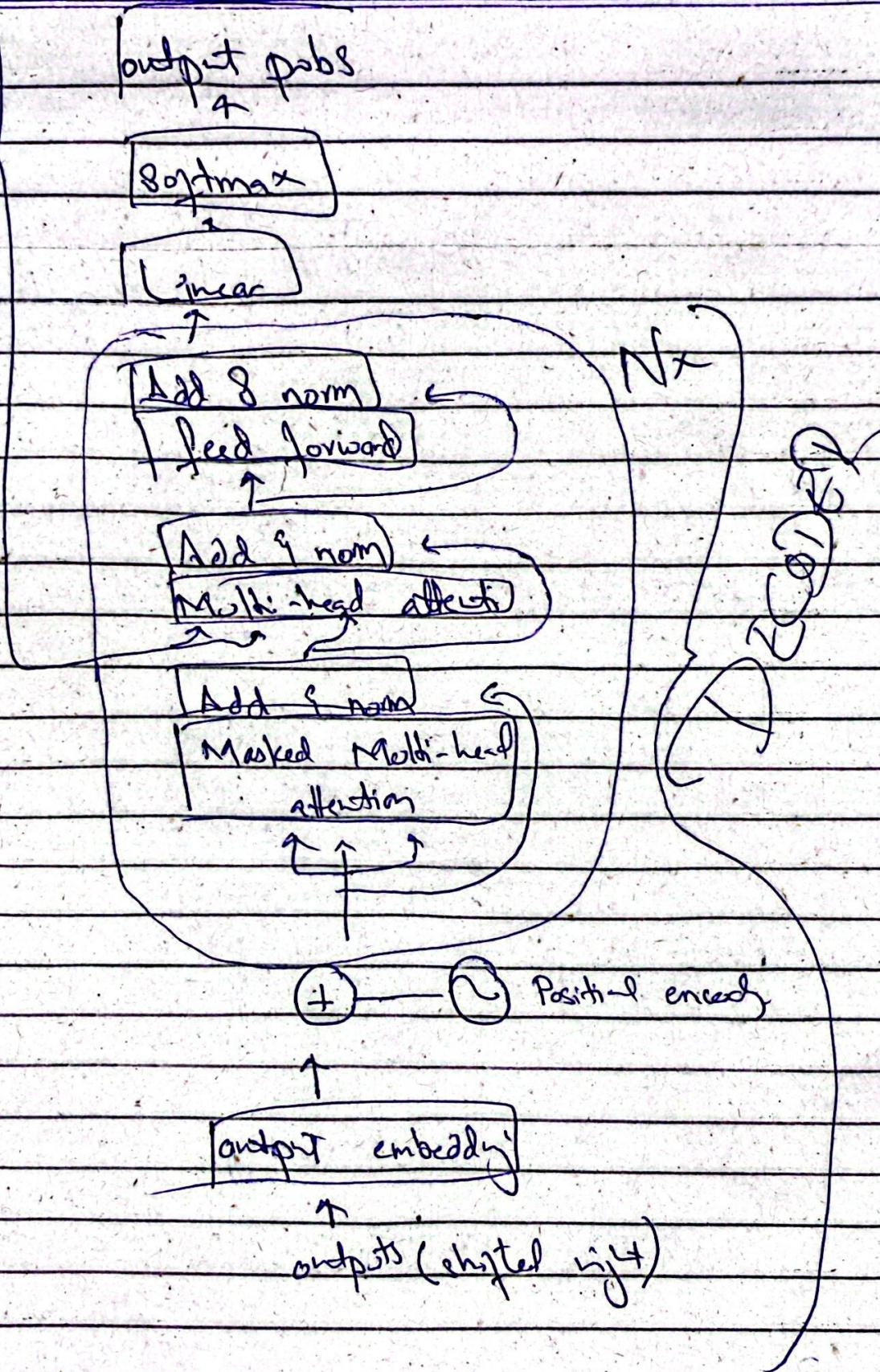
$$\frac{dy}{dx} = \frac{\partial y}{\partial f} \cdot \frac{\partial f}{\partial x}$$

Vanish

{ if $\frac{\partial y}{\partial f}$ & $\frac{\partial f}{\partial x}$ are < 1 , it
will result in more smaller numbers.

Explode { if $\frac{\partial y}{\partial f} > 1$, result will be a more
greater result.



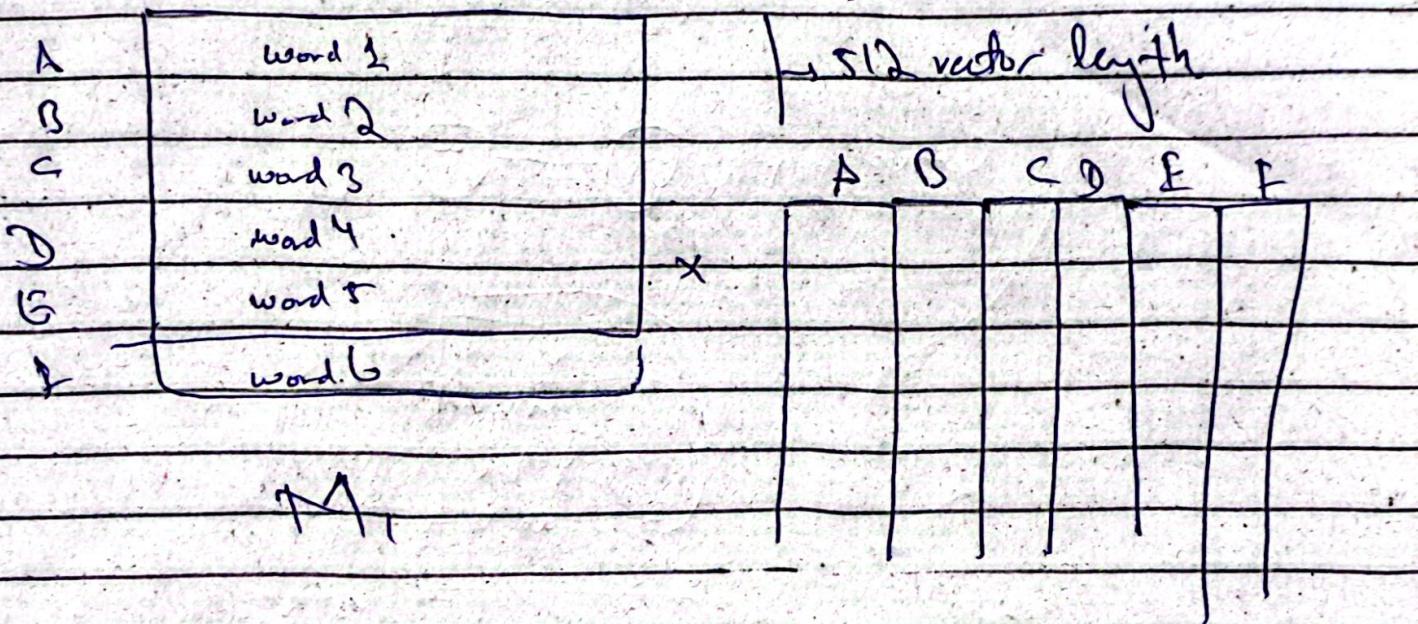


Notations

Input matrix (Sequence, d_{model})

6×512

6 words



512×6

(Transpose)

$\downarrow T_2$

\Rightarrow result of both would be
 (6×6)

here (1,1)

AA	AB	AC	AD	AE	AF
BA	BD				
CA		CC			
DA			DD		
EA				EE	
FA					FF

here $(1,1)$

is dot product

of A in M_1

& A in M_2

6×6

embeddings are

not fixed, they are
parameters to our
model

Encoder

input embedding

sentence (tokens): 100 105 101 101 101 101



input IDs: 100 105 101 101 101 101



Embedding
(vector size 512)

512×1
for each

our model will learn to chose these 512 values for
each token so that it could represent sentence.

Due to 512 size, its 312d model,
 $d_{model} = 512$

Positional Encoding

↳ Each word should carry some info about its position in sentence.

↳ we want model to treat words that are appear close to each other as close and words that are distant as distant.

↳ we want the position encoding to represent a pattern that can be learned by the model.

For each word vector we create some special vector called position encoding of same length i.e. 512 that we add to these embeddings.

Embedding to token

512×1

Position encoding

512×1

= ENCODER INPUT
 (512×1)

It is not learned.

It is computed only once.
it is fixed.

How position encodings are calculated?

1. Your (col) 2. []

$PE(pos, di) = \sin \frac{pos}{10000^{\frac{2i}{d}}} \quad \text{for even } di$,
 $\cos \frac{pos}{10000^{\frac{2i+1}{d}}} \quad \text{for odd } di$

$PE(pos, di+1) = \cos \frac{pos}{10000^{\frac{2i}{d}}} \quad \text{for odd } di+1$,
 $\sin \frac{pos}{10000^{\frac{2i+1}{d}}} \quad \text{for even } di+1$

$PE(pos, di)$

↓
pos \Rightarrow word inside sentence starting from 0.

For each sentence,
Position encoding will be
same

After applying position encoding only on \mathbf{e} , we can reuse the \mathbf{e} for every sentence whether its train or inference.

Why \sin & \cos for \mathbf{e} ?

dir from like \cos is naturally represent a pattern that model can recognize as rotated coordinate, so relative position are easier to see for model.

$$p(k, 2i) = \sin\left(\frac{k}{n^{2i/d}}\right)$$

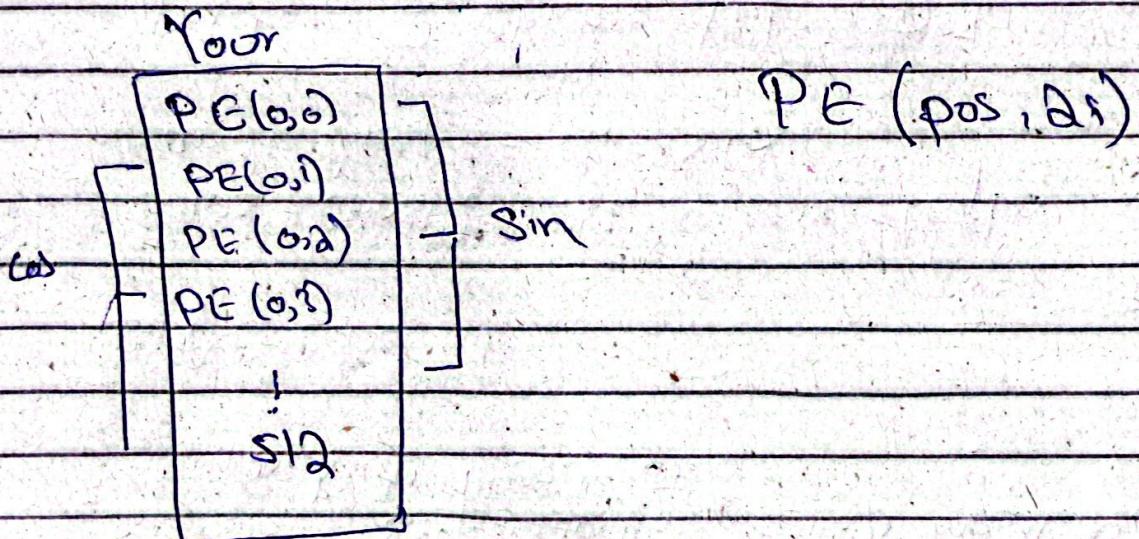
$$p(k, 2i+1) = \cos\left(\frac{k}{n^{2i/d}}\right)$$

k : pos of object in sequence
 d : dimension of output embedding space

$p(k, j)$ - position j , for mapping a position k in the input sequence to index (k, j) of the position matrix

n - user defined scalar, set to 10,000 by author

i - used to for mapping to column indices
 $0 \leq i \leq d_2$ with a size of i maps to
both size and cosine.



We will have same position encoding for each sentence because position encoding are calculated once and are reused for every sentence that our model will see - during inference or training.

Multi-head attention self-attention with single head

→ Single head attention: Self-attention is a mechanism that existed before they introduced transformers. They authors of the transformer changed it to multi-head attention.

→ What is self-attention?

Self-attention matrix for input matrix (Q, K, V) are calculated as

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

To calculate self-attention

Self-attention allows model to relate words to each other

Simple case

We consider the sequence of length 1 length seq = 2 and $d_{\text{model.}} = d_h = 512$.

Matrix Q, K, V are just input sentence.

where Q, K, V are the concatenation of query, key and value vectors.

Self-attention believe we are computing relative b/w each word in sentence so other word in sentence

Q - input sentence

(6, 512)

↳ 6 tokens embedding of these tokens

[Your] [cat] [is] [a] [lonely] [cat]

K^T - same sentence but transposed (512, 6)

$\Rightarrow Q \times K^T$
(6, 512)

$\underline{Q \times K^T}$

$\sqrt{512} \quad \because d_K = 512$

The softmax

softmax $\left(\frac{Q \times K^T}{\sqrt{512}} \right)$

$Q \times K$	Your	cat	is	a	lovely	cat	Σ
$\sqrt{512}$	0.268	row1.col2				row2.col2	1
	0.124						1
	0.1477						1
							1
lovely							1
cat	row2.col1						1

6x6

(1x1) First value represents dot product of first row with the first column

Softmax makes these value in such a way that they sum up to 1.

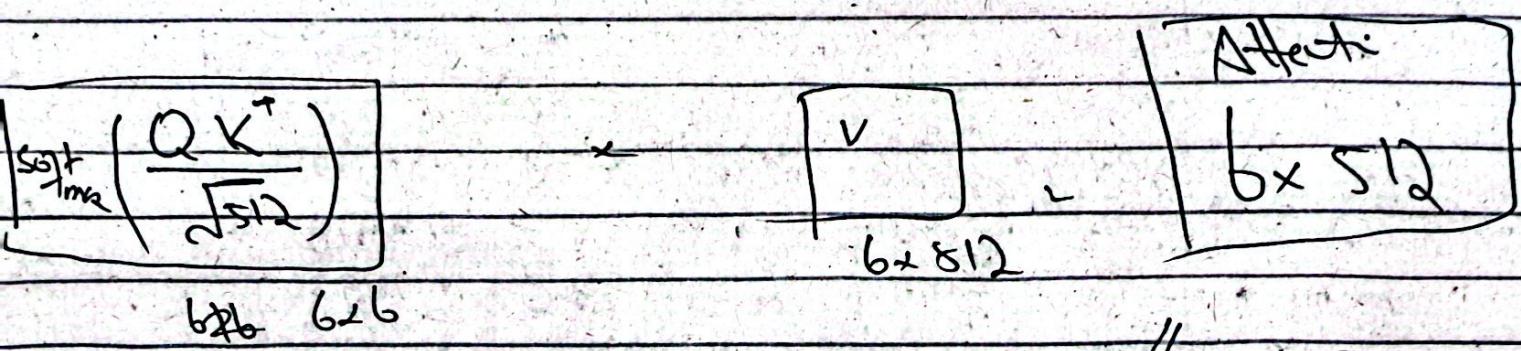
so $[1 \times 1]$ is dot product of first word embeddings of word 'Your'

$[1 \times 2]$ = dot product of embedding of word 'Your' with embedding of word 'Cat' 'Cat'

number shows interest of value

So now multiply with $V_{6 \times 512}$

softmax $\left(\frac{Q \cdot K^T}{\sqrt{512}} \right) V_2$. Attention



Now this new representation represents the words in the sentence with the same shape as the input embeddings.

- meaning of the word - input embedding
- position of word (position embedding)
- relation b/w ^{each} words (attention), with all other words.

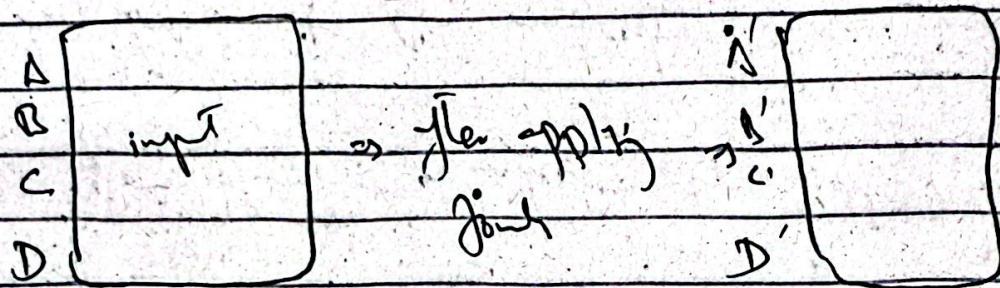
This is self-attention (single head).

V represents the meaning or significance of each token within the context of the entire sentence. This info is crucial for the attention mechanism to effectively weight the importance of each token when generating the output.

very attention captures the position in the sentence (given by partial encoding) but also each word's interaction with other words.

Self-attention in detail

→ Its permutation invariant: (i) we don't consider positional encoding.



Permutation invariance means that model's output remains the same even if the order of the input elements is changed. In context of self-attention, this means that the attention mechanism does not depend on absolute positions of elements in the input sequence. It only depends on the relationship b/w them the elements.

~~It is good~~

rather than index/post of word in sequence, it depends on relationship b/w words.

mean idea of

If we change B's position, the values of \vec{B}' and \vec{C}' won't change but other positions will change.

→ Self attention requires no parameters.

→ until now, interaction \rightarrow which has been driven by their embedding + the positional encoding. This will change later.

→ No quantity being learnt except embedding ratio

→ we expect values along diagonal to be higher because it's dot product with embedding with itself.

→ If we don't want some position to interact i.e. YOUR aid. CDT, put value = $-\infty$, before softmax.

$$e^{-\infty} \approx 0.$$

To be used in
judder

~~Multi-head attention~~ \rightarrow ~~Sequence input~~ (embedding)

Q - Query vector represents the current state of model or current position in input sequence. It captures the information that the model is currently interested in and helps determine which parts of the input sequence are most relevant to task.

K - Key: They key vector represents the information from all input tokens. It provides a way for the model to assess the relevance of each input token to the current query. Key vector is usually calculated using a learned embedding. It can also be transpose of Q i.e. $K = Q^T$

(V) - Value vector represents the learnt representation of the input tokens. It captures contextual info associated with each token, allowing the attention mechanism to weigh the importance of each token when generating output. The value vector is typically calculated using a learned embedding.

Seq - Sequence length, d_{model} - Size of embedding vector

$$h \approx \# \text{ of heads}$$

$$d_h = d_v \quad d_v = d_{model} / h$$

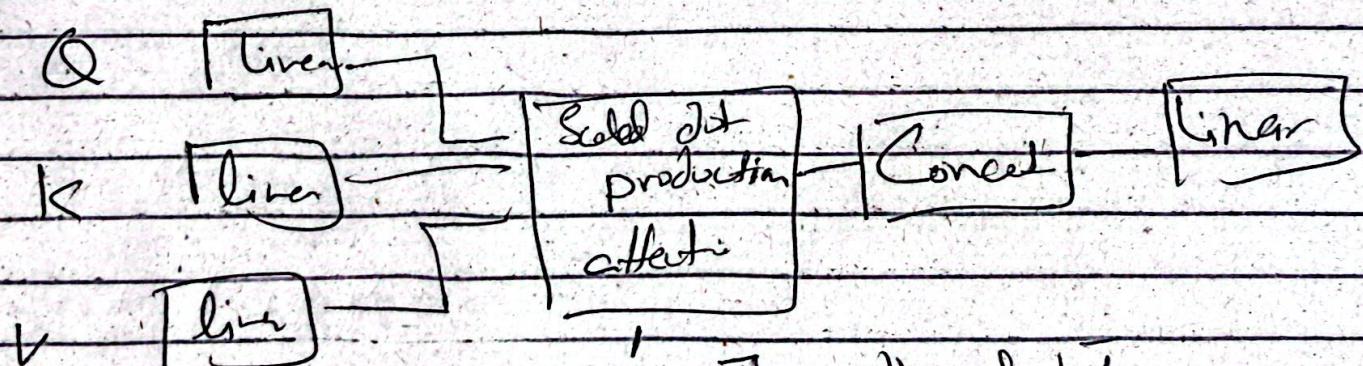
In positional encod - we used precompute segment vector

Query - raw input

Key - After parsing preposing Query

* There are 4th key to model best results & opt value

Input to Q, K IV are see ~~last~~ (Sequence)



i. - This will calculate
i. concatenation by word

Scaled dot Product Attention

↳ let's see how this layer builds correlation

[wait]

Calculating attention comes primarily in three steps:

→ First we take O and each key and compute the similarity b/w two (to obtain a weight).

→ Second step is to use a softmax P_{ij} to normalize these weights and finally to weight these weights in conjunction with corresponding values and obtain the final attention.

In current NLP, key-value mostly

→ Multi-Head Attention

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

$$d_k = \sqrt{512}$$

Multi-head attention (Q, K, V)

$$= \text{Concat}(\text{head}_1, \text{head}_2, \dots, \text{head}_h)W^O$$

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

$$h = \# \text{ of heads} \quad d_k = d_v = d_{\text{model}}$$

$Q, K, V = \text{input sequence}$

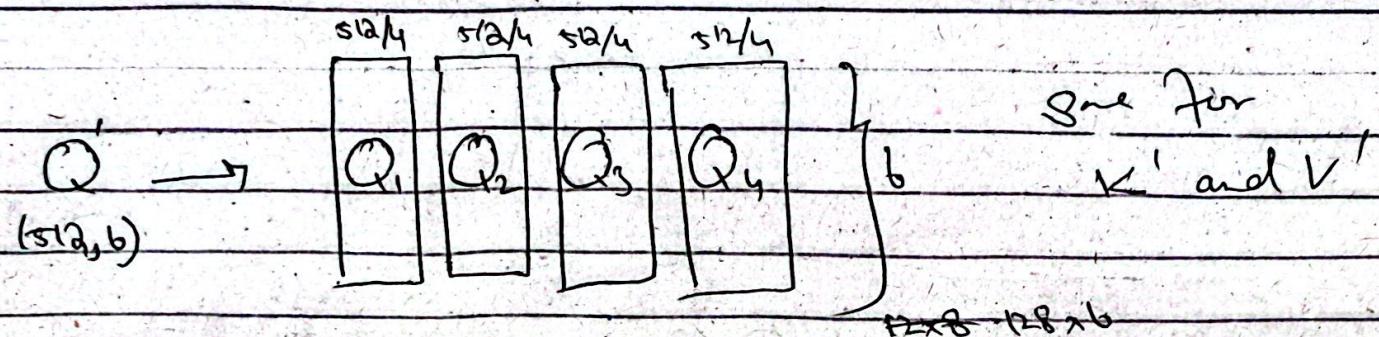
Input

$$(seq, d_{\text{model}}) \rightarrow Q \in \mathbb{R}^{(6, 512)} \times W^Q \in \mathbb{R}^{(512, 512)} = Q' \in \mathbb{R}^{(6, 512)}$$

$$K \in \mathbb{R}^{(6, 512)} \times W^K \in \mathbb{R}^{(512, 512)} = K' \in \mathbb{R}^{(6, 512)}$$

$$V \in \mathbb{R}^{(6, 512)} \times W^V \in \mathbb{R}^{(512, 512)} = V' \in \mathbb{R}^{(6, 512)}$$

Now we have to split Q', K', V' into smaller matrices by sequence model dimension



$$d_K = d_{\text{model}}/h = 512/4 = 128$$

Now calculate attention for b/w these smaller matrices over (Q_1, K_1, V_1) using expression

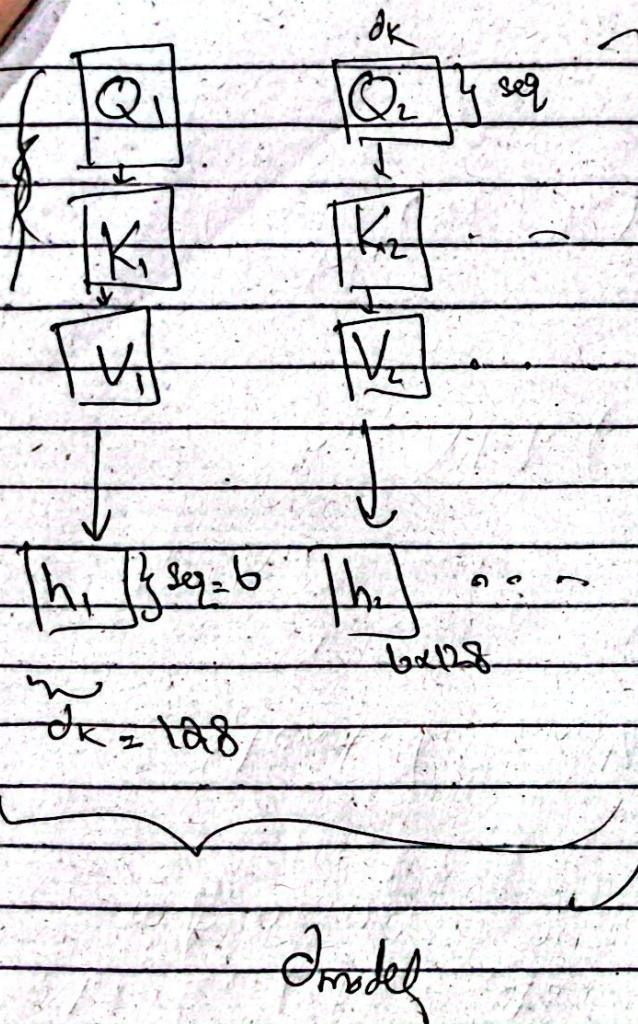
$$\text{Attention}(Q_1, K_1, V_1) = \text{softmax}\left(\frac{Q_1 K_1^T}{\sqrt{128}}\right) \cdot V_1$$

for h_1

$$h_1 = \text{Attention}(Q_1^T, Q_2^T, V_1)$$

$$h_1(\text{seq}, dv)$$

$$dv = d_K$$



the b in last calculation
is done by multiply V .

and if there is 6×128

model

Combine all heads

(6×512)

~~Multi-head (Q, K, V) \rightarrow Concat (h_1, h_2, \dots, h_n)~~
~~along b dimension~~

Combos of heads

$H = \text{Concat}(h_1, h_2, \dots, h_4)$

↳ $(6, 128)$ $(6, 512)$

Multiplying result H by W_0

Multihedral (OK)

$H \times W_0^0$
(512×64) (128×128)

↳ Because we want to each head
to work different aspects of each
word.

↳ because a word can be noun
in one case & ~~marked~~ in other
depending on context -

So, one head may begin to serve as noun
and other as adverb.

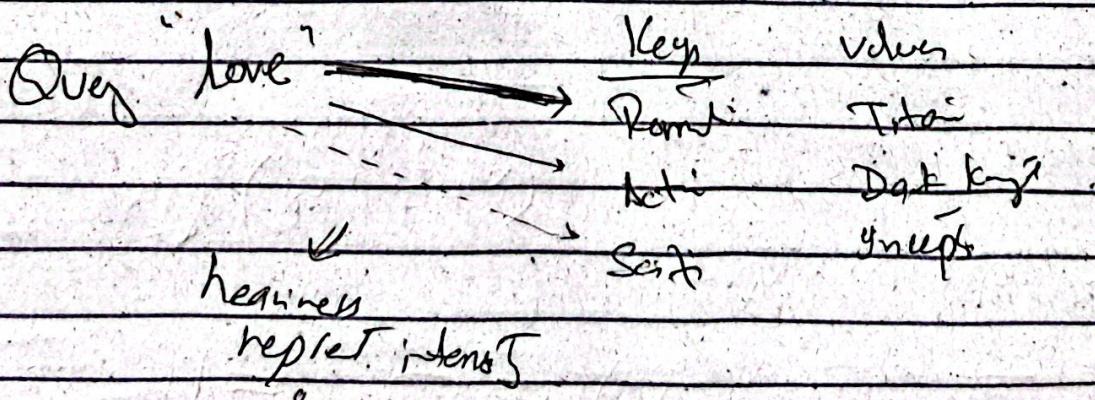
~~n Attest can be calculated.~~

When we calculate att_i^k b/w
 Q_i, K_i is matched.

$$\text{softmax} \left(\frac{Q_i K_i^T}{\sqrt{d_k}} \right)$$

we get a correct match as score
before. This score represents intensity
b/w two words.

why Q, K, V ?



Add. and Norm.

↳ To introduce add. / norm we need layer normalization.

layer Normalization

↳ imagine batch of n -items.

↳ each ~~embed~~ item is an embedding vector

item 1

calculate μ & σ of each of these items
independently from each other and

we replace each value with

another value x_j by

$$x_j' = \frac{x_j - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

Normalizing so, new values are in range
0 to 1

Notes

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

we multiply \hat{x}_j with λ & add β

$$= \lambda(\hat{x}_j) + \beta$$

because for \hat{x}_j might be to
restrict to \hat{x}_j have values 0 or 1.

So the return mean \hat{x}_j , λ to
introduce fluctuation

Batch Norm

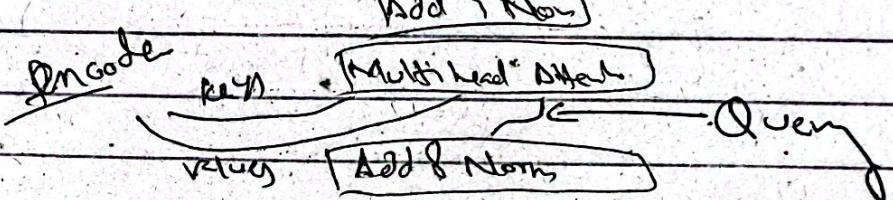
VS Layer Norm

↳ Same feature
for all batch

one token at a time
feed forward

Add f Norm

feed forward



Decoder

Output embeds
are seen in query as
input embeds

Masked multi-head
attents

⊕ ↗ position
encoding

Output embeds

outputs

Multi-head Attn

Take Keys & values from
encoder and Query from decoder

↳ So attention is not self-attende but cross attention

Query & Output of Masked multi-head attend

Marked Multi-head Attn 's effect on

effect of input sentence of decoder

Feed forward is just fully connected layer

Marked Multi-head Attn

Our goal is to make the model causal.

It means the output at a certain position can only depend on the words on previous positions. The model must not be able to see future word.

so,

QK^T

	your	cat	is
your	0.268	$-\infty$	$-\infty$
cat	0.164	0.268	$-\infty$
is	0.147	0.123	0.262



Then softmax will replace $-\infty$ with 0.

we don't want pair to match future words

Cat will only match 'your Cat'

We do this in multi-head attention before

before applying softmax

Inference & Training of a Transformer Model

Training

I love you



Te amo molto

input: < SOS > I love you very much < EOS >

→ will model about what it ~~is~~ ~~st~~

↓
I what is end in ~~is~~ ~~the~~ of sentence

Encoder input:

Decoder input: $\langle \text{SOS} \rangle \text{ Ti amo molto}$

$\langle \text{SOS} \rangle \text{ Ti amo molto}$

↳ we made this sequence of fixed length because when by adding padding.

If our model supports seq length of 1000,

here we have 4 tokens ($\langle \text{SOS} \rangle \text{ Ti amo molto}$),

will add 996 of padding to make the sequence long enough.

To get back from embedding to ^{relative} dictation

Decoder output \rightarrow [Linear layer] \rightarrow (b, vocab size)
(b, 512)

8. it will tell, for every embedding it sees, its ^(token) position in our vocab.

label: Ti amo molto <EOS>

This is what if we expect from our model give <EOS> y like you very much.

Cross entropy for

Then calculate ~~loss~~ after softmax.

Cross entropy for $(\hat{y} \text{ vs, label})$

but why <EOS> & <EOS> only

We expect our model to output

Ti amo molto <EOS>

\hookrightarrow this is label/target

Here, in decoder input, our sequence length is 4. " $\langle \text{SOS} \rangle$ Te amo molto" like this are padded here.

and we want

"Te amo molto $\langle \text{EOS} \rangle$ "

So, when our model will see $\langle \text{SOS} \rangle$, it will output first token as output

"Te"

and when our model will see "Te", it will see "amo" and when it will see output

"amo" it will output "muito" and when it will see "muito" it will output $\langle \text{EOS} \rangle$ which indicates translation is done.

This all happens in one time step.

In RNN we have n time steps to map an input sequence into an output sequence, but this problem is solved with transformer.

Line layer to project embeddy to vocab

Different

I love you very much



Ti amo molto

Encoder input - <SOS> I love you very much <EOS>

and

Decoder input: <SOS> with Paddy

logits: output of last layer is known as logits
↳ line ^{layer} to project embeddy back to our vocab. This projector is called logits.

logits will go to system

After system it will produce. T_p

that will be fine step = 2.

Training happen in one pass, while inference
we'll do it token by token and we'll
also see what's the CR

At fine step = 2, we won't calculate
encoder output again because our
English sentence didn't change

for $T=2$,

decoder input = $\langle \text{EOS} \rangle \text{Tr}$
at end we'll get next token "amo"
is appended.

for $T=3$,

decoder input = " $\langle \text{EOS} \rangle \text{Tr amo}$ "
end output = " $\leftarrow \text{so molto}$ "

for $T=4$,

decoder input = " $\langle \text{EOS} \rangle \text{Tr amo molto}$ "
end output = $\langle \text{EOS} \rangle$

It take $\mathcal{O}(n^k)$ time steps.

Inference Strategy

→ We selected, at every step, the word Greedy with the softmax rule. Search

→ A better strategy is to select at each level step the top B words and evaluate all the possible next words for each of them and at each step, keep the top B most probable sequences.

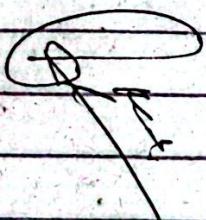
Beam
Search

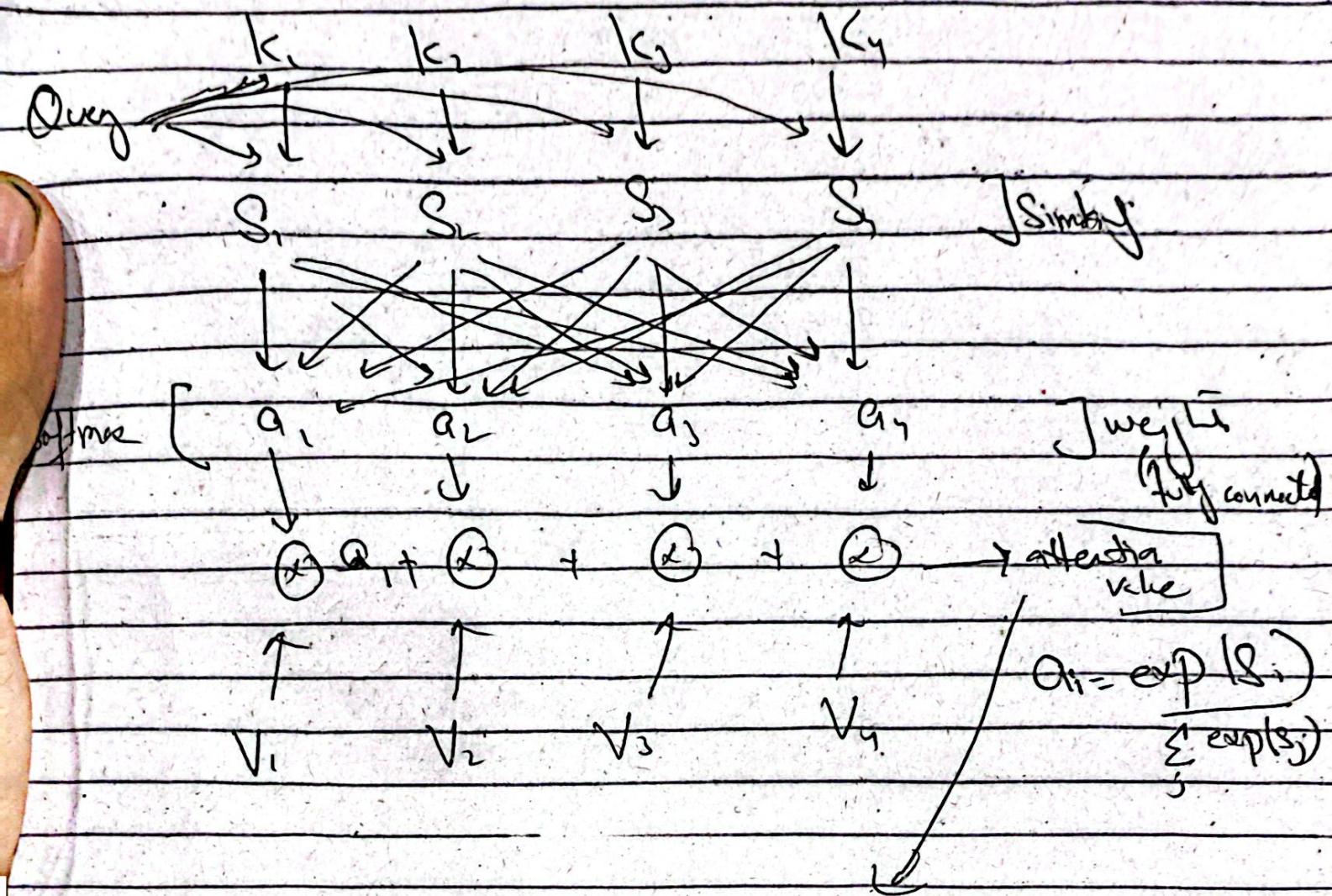
[wait]

$$\text{Similarity } f(Q, K) = \begin{cases} Q^T K - \text{dot product} \\ \frac{Q^T K}{\sqrt{d_K}} - \text{Scaled dot product} \\ Q^T W K - \text{General dot product} \\ W_Q^T Q + W_K^T K - \text{additive similarity} \end{cases}$$

$Q^T W K$:

Projecting query Q into new space using weight matrix W and then taking dot product with K





$$\text{attention value} = \sum a_i \cdot v_i$$

So, in Qwik, we don't specify W . There are values that get optimized by NN

itself. S, W will be adjusted by back propagation so that NN can learn on its own who might be a good source to pass into

K_i, V_i are vectors

S_i, A_i are scalars
(~~handwritten~~)

without position encoding, it will treat words like separated by
j
and

Multi-head attention

↳ we feed in a vector (input) other
multi-head will compute attention by
every position. and every other position.
So, we have vectors that embed words
in each one of these positions and now,
we simply carry out an attention, competition
where we treat each word as a query,
and find some keys that correspond
to other word in the sentence and then
take a convex combination of the
corresponding vector. Their values V are
going to be same as key, then take
dot product of that to produce settle
embeddy

So idea is that this attention will
pick each word, combine it with one
of other words through the attention
mechanism to essentially produce a better
embeddy that ways together info from pairs

So, this encoding block will repeat N times
meaning N blocks of block.

In first block, it will look at
pairs and in next pair of pairs and so on.

So, we are combining N pairs of words.

→ ~~Multiple Attn~~ ~~Stack Attn~~ repeat.

Add 1 Nam

encoder output

will be one embedding per
position. It captures the
original word but also info
from other words that it attends
to in the sentence. So, it
can be thought of a large
embedding of all the words corresponding
to each position.

Decoder

Masked Multihead

b. to combine output words with
precision output words.

Multi-head Attention

compute multiple attentions
per query with different weights.

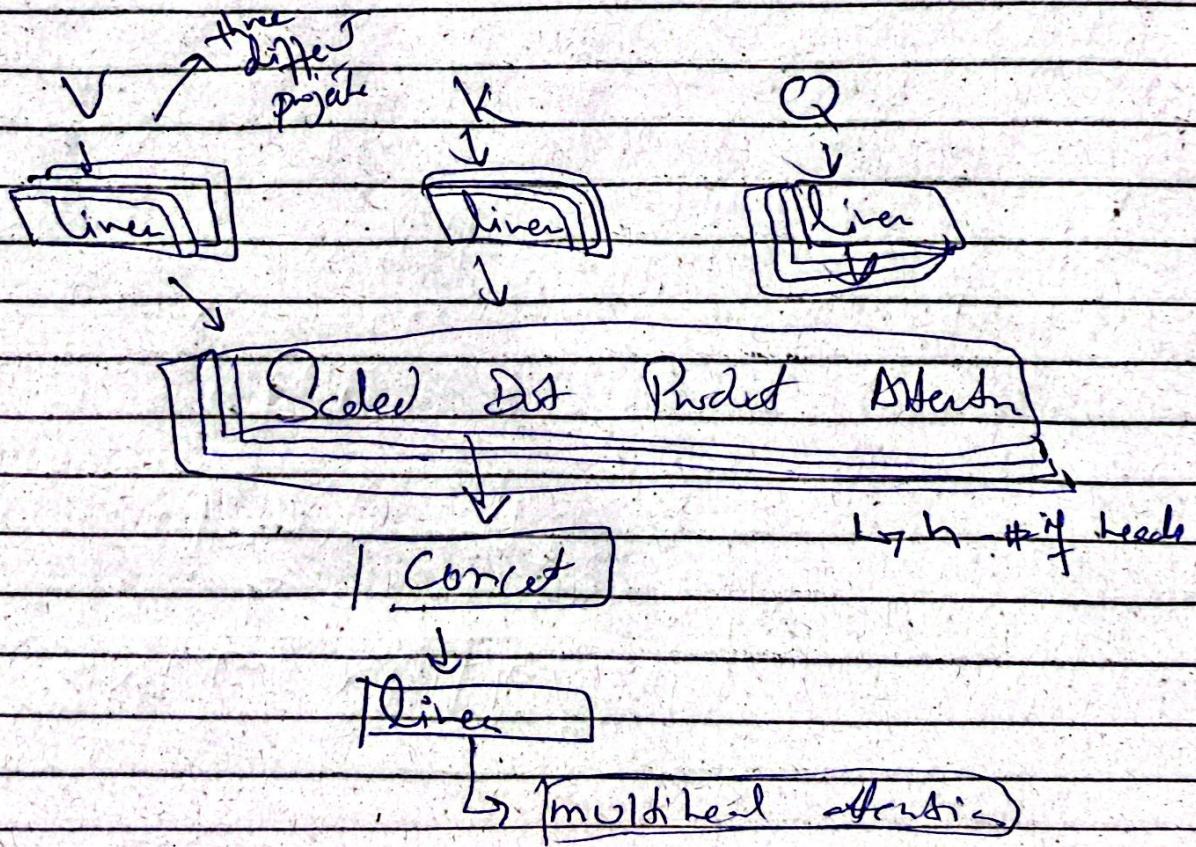
$\text{multihead}(Q, K, V) = W^0 \text{ concat}(h_1, h_2, \dots, h_n)$

$\text{head}_i = \text{attention}(W^0(Q, W^i; K, W^i; V))$

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

The general way of thinking about attention is that we have some key-value pairs, then there's query that going to compare to each key and the keys with greater similarity are going to have higher weights and then, we can take weighted combination of corresponding values to produce the output.

So we are going to feed K, V, Q into some linear layer, then we'll calculate scaled dot product attention.



masked Attention (Q, K, V) - $\text{softmax} \left(\frac{Q^T K + M}{\sqrt{d_k}} \right) V$

↳ output value shouldn't depend on future value.

before softmax, pt value = $-\infty$

Training

Teacher forcing: when you train any NN, you have both input sentence and output sentence/label and then you assume that my outputs are correct everywhere so I'm gonna feed that as decoder input, so I'm gonna feed correct output words are previous positions and the NN simply try to predict next word based on that.

GPT

Jacobian

derivative = slope

↳ can work for $f(z)$ but not
 $\text{for } f(\text{xyz})$

↳ derivative / slope

On a similar note, integrals are thought of as area under the curve

$$\text{Area} = \int_a^b f(z) dz$$

Let's rethink derivative & integral

Linear maps :

Linear maps (from $\mathbb{R}^2 \rightarrow \mathbb{R}^2$) are required to have these properties

- (1) → parallel lines stay parallel
- (2) → Even spacing preserved
- (3) → Origin is fixed

⇒ Determinant is a scaling factor for area (A)

↳ length (l)

Avoid overfitting

- Increase size & quality of your data
- Data augmentation
- Early stopping
- Regularization
- Dropout
- Decrease model complexity
- Feature selection
- Ensemble
- Cross Validation
- Transfer learning

Gradient accumulation

- with gradient grad-zeno

Avoid underfitting

- Increase model complexity
- Increase amount of training data
- Reduce epochs
- Increase # of epochs
- Use regular (reduce regular strength)
- Cross validate
- Ensemble learning