

Transformers

Introduced in

2017

- ANN → tabular data
- CNN → Image data
- RNN → Sequential data
- Transformers → Sequence-to-sequence task
(inp & o/p is sequence)

→ Machine translation
→ Q&A system
→ Text summarization

→ In transformers, we use Self-attention. That's why it can handle massive amount of data parallelly.

→ Hence, transformers are a neural network architecture for handling sequence-to-sequence data.

→ Transformer have encoder-decoder, but instead of LSTM it uses self-attention mechanism.

→ Due to parallel data processing, training becomes faster so it is useful for large datasets.



Impact of Transformers

Revolution in NLP.	Demonstrating AI	Multi-modal capability	Acceleration of Gen AI	Unification of Deep-learning
--------------------	------------------	------------------------	------------------------	------------------------------

Advantages

- Scalability
- Transfer learning
- Multi-modal I/O
- flexible Architecture
- Eco-system
- Integrated AI Techniques

- GAN + Transformers = Image generator
- Reinforcement learning + Transformers = Game playing agent.
- CNN + Transformers = Image generator / Visual (Vision Transformer)



Disadvantages

- High computational infrastructure
- Data
- Energy
- Overfitting
- Bias

Future

- Improvement in efficiency
 - pruning
 - Quantization
 - knowledge distillation
- Multi-model capability
 - Text
 - Image / Speech
 - Sensory data
 - Biometric feed back
 - Time Series
- Responsible
 - Bias elimination
 - Focusing on ethical concerns.
- Domain Specific
- Multi-lingual (Multi-language Support)
- Interoperability
 - Difficult to understand external process.
 - Black Box

Self-attention Mechanism

What is self attention....?

→ In NLP applications, the most important part is vectorization (words-to-number conversion), like how you convert words into numbers.

like:

OneHotEncoding	1. cat mat nut	cat	mat	nut	Inefficient
	2. mat mat mat	mat	mat	mat	
		1	0	0	
		0	1	0	
		0	0	1	

Bag of Words	→ Sentence 1	[mat	cat	nut]
	Sentence 1	[1	1	2]
	Sentence 2	[2	0	1]

TF-IDF (Term frequency-Inverse Document Frequency)

Word Embeddings → Capability of Skinning Semantic Meaning.

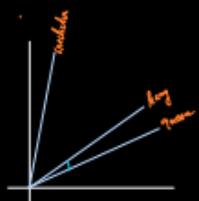
→ Process :

1. Get training data



Hypothetical graph of 2D.

2. Pass it to neural network & it will convert it into n-dimensional vector space



For

king $\rightarrow [0.1, 0.38, 0.46, 0.67, 0.49]$

queen $\rightarrow [0.51, 0.37, 0.47, 0.66, 0.49]$

computer $\rightarrow [0.2, 0.9, 0.77, 0.60, 0.39]$

- Similar words will have very similar vectors.
- Based on data, each number reveals some aspects.
- Word-embedding captures average meaning, like "On an average with which frequency, that particular word used in your data."

\rightarrow Now suppose we're generating embedding of a word "apple" of 2D $[x, y]$ where: x : Taste & y : Technology.

Example:

- An apple a day keeps doctor away $\rightarrow [0.6, 0]$
- Apple is healthy. $\rightarrow [0.7, 0]$
- Apple better than orange $\rightarrow [0.8, 0]$
- Apple makes great phones. $\rightarrow [0.8, 0.2]$

Now suppose, if our dataset is of 10,000 lines. where 7000 rows, Apple is represented as fruit & in 3000 rows, Apple is represented as company. Now, obviously overall vector of Apple will look like this $[0.7, 0.3]$, because the

average meaning of apple word in dataset is more biased towards taste technology

- So word-embeddings are generated once & will be used multiple times.
- It means that they are static, which is an issue.

Sample sentence.

apple launched a new phone, while i was eating an orange.

- Currently, embedding of apple is $[0.7, 0.3]$ which tells that whenever word is used, it always be in the context of a fruit.
- So there should be Context-based Embedding instead of Static Embedding.

Value should get changed based on context.

like

apple launched a new phone, while i was eating an orange.
[0.3, 0.9]

To solve this problem, we need to use Self-attention.

It can give you smartly converted contextual embedding from your word embedding.

Embeddings of each words of a sentence.

Self-attention Mechanism

→ Newly converted smart contextual embeddings of each word for given sentence.



Created with
Notewise

Behind the scene of Self-attention

Sample text:

SL = Money bank grows , SR = River bank flows

Now let's distinguish the word "river" based on their context

for SL : funds = 0.3 money + 0.7 bank + 0.2 grows

for SL : funds = 0.5 river + 0.4 bank + 0.2 flows

SL

money ~ 0.3 money + 0.2 bank + 0.1 grows

bank ~ 0.25 money + 0.7 bank + 0.05 grows

grows ~ 0.2 money + 0.2 bank + 0.7 grows

SR

river = 0.8 river + 0.15 bank + 0.05 flows

bank = 0.2 river + 0.75 bank + 0.02 flows

flows = 0.4 river + 0.01 bank + 0.57 flows

Mathematical notation

$$e_{\text{money}}^{(\text{new})} = 0.3 e_{\text{money}} + 0.2 e_{\text{bank}} + 0.1 e_{\text{grows}}$$

$$e_{\text{bank}}^{(\text{new})} = 0.25 e_{\text{money}} + 0.7 e_{\text{bank}} + 0.05 e_{\text{grows}}$$

$$e_{\text{grows}}^{(\text{new})} = 0.2 e_{\text{money}} + 0.2 e_{\text{bank}} + 0.7 e_{\text{grows}}$$

$$e_{\text{river}}^{(\text{new})} = 0.8 e_{\text{river}} + 0.15 e_{\text{bank}} + 0.05 e_{\text{flows}}$$

$$e_{\text{bank}}^{(\text{new})} = 0.2 e_{\text{river}} + 0.75 e_{\text{bank}} + 0.02 e_{\text{flows}}$$

$$e_{\text{flows}}^{(\text{new})} = 0.4 e_{\text{river}} + 0.01 e_{\text{bank}} + 0.57 e_{\text{flows}}$$

More depth

$$e_{\text{bank}}^{(\text{new})} = [e_{\text{bank}} \cdot e_{\text{money}}^T] e_{\text{money}} + [e_{\text{bank}} \cdot e_{\text{bank}}^T] e_{\text{bank}} + [e_{\text{bank}} \cdot e_{\text{grows}}^T] e_{\text{grows}}$$

$$\left(\begin{array}{|c|c|c|} \hline & & \\ \hline \end{array} \right) \cdot \left(\begin{array}{|c|c|} \hline & \\ \hline \end{array} \right) = s_{21}$$

query vector key vector

$$\left(\begin{array}{|c|c|c|} \hline & & \\ \hline \end{array} \right) \cdot \left(\begin{array}{|c|c|} \hline & \\ \hline \end{array} \right) = s_{22}$$

query vector key vector

$$\left(\begin{array}{|c|c|c|} \hline & & \\ \hline \end{array} \right) \cdot \left(\begin{array}{|c|c|} \hline & \\ \hline \end{array} \right) = s_{23}$$

query vector key vector

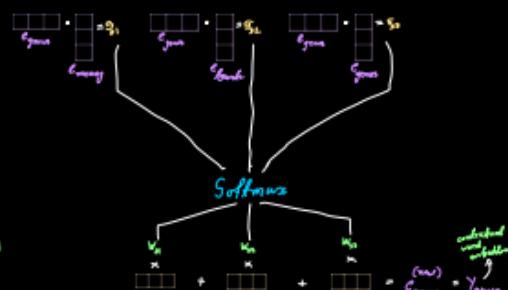
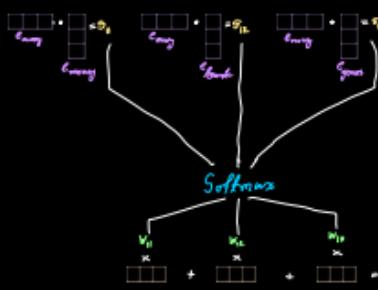
Softmax

$$W_{21} = \frac{e^{s_{21}}}{e^{s_{21}} + e^{s_{22}} + e^{s_{23}}} \quad W_{22} = \frac{e^{s_{22}}}{e^{s_{21}} + e^{s_{22}} + e^{s_{23}}} \quad W_{23} = \frac{e^{s_{23}}}{e^{s_{21}} + e^{s_{22}} + e^{s_{23}}}$$

$$\text{Value Vector} \longrightarrow \left[\begin{array}{|c|c|c|} \hline & & \\ \hline \end{array} \right] + \left[\begin{array}{|c|c|c|} \hline & & \\ \hline \end{array} \right] + \left[\begin{array}{|c|c|c|} \hline & & \\ \hline \end{array} \right] = e_{\text{bank}}^{(\text{new})}$$

contextual word embedding

With this same process we will calculate y_{money} and y_{goats} .



Now, for sentence-2 Money bank goats we have unique contextual embeddings of each words

Consider....!

- Parallel calculation process
- Due to parallel computation, we're losing sequence information.
- No learnable parameters.



Nearly generated embeddings are general contextual embeddings, not task specific contextual embedding, which is a big problem.

So, it is not sufficient to generate only contextual embeddings, we also need task specific contextual embeddings.

(embedding of a word based on its surrounded words).

Piece of cake

meaning-1: देश की गज़िया

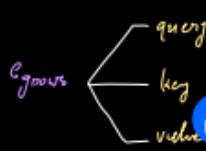
meaning-2: बड़ा बड़ा देश

In short we need to add trainable parameters in calculation process.



Our current process, there are 2 possible locations where we can set-up trainable parameters which is dot-product operation.

Now, here e_{money} or e_{bank} or e_{goats} they're individually performs their different role. one query vector, key vector & value vector



→ Different role as query, key & value vector of a particular word (e_{money} | e_{grows}) is same here which is not recommended. Because those query, key & value vector must be different for their respective task.

→ So we will create different vectors based on their task such as..

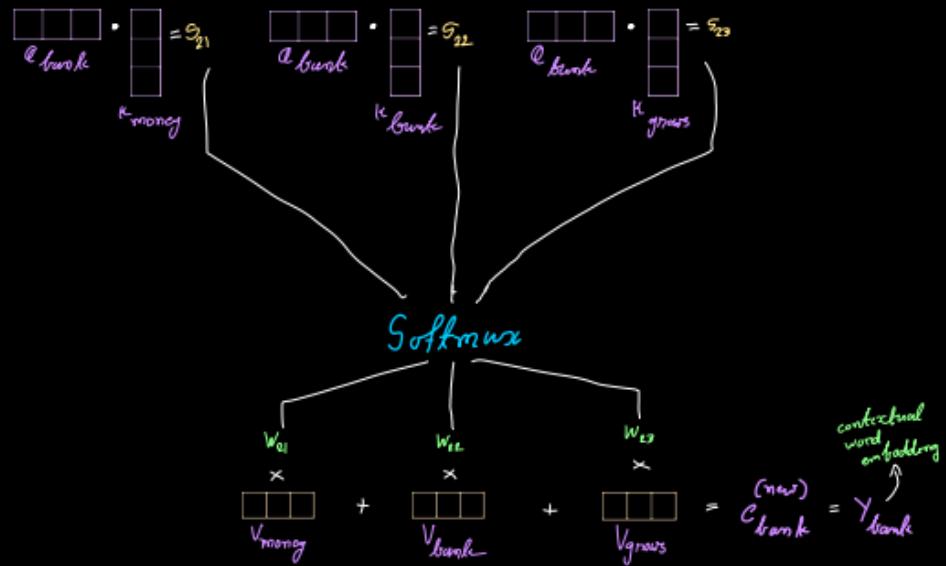
e_{bank} → Question: How much i am similar?

e_{bank} → Key: Solution of query

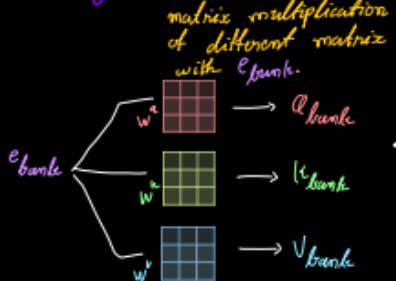
e_{bank} → Value: Answer of query.

These 3 Query, Key & Value Vectors will be made by data.

→ After as we got these data, now we will send Query Vectors of particular word instead of its embedding.



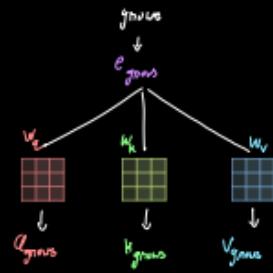
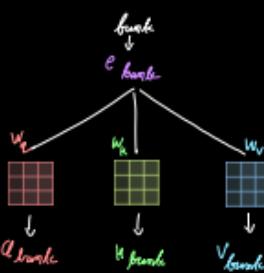
Q. How to generate these new vectors from embeddings...?



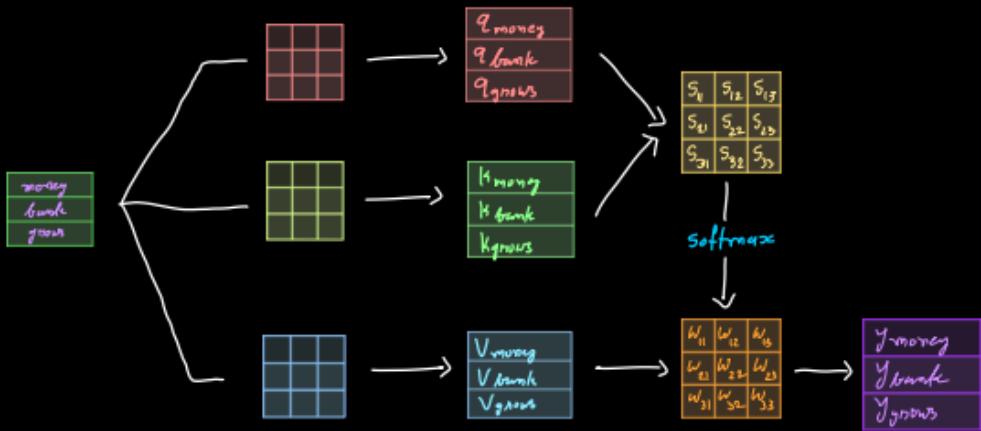
Now the value of all of these individual matrix will be generated by given data during training.

• This process will be repeated for each word in a sentence.





→ Hence all w_q , w_k and w_v will be sums for all words in given sentence.



This is self-attention mechanism.
In short.....

$$\text{Attention}(Q, K, V) = (QK^T)V$$

Formula, mentioned in document....

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



→ This formula is called as Scaled Dot Product.

→ Here $d_n =$ dimensionality of \mathbf{h} (key) vector.

Q. Why scaling...?

A. Due to the nature of dot product.

Nature of dot-product: If we have low dimensional vector space than the dot-product of those vectors will give less variance.

- If we have high dimensional vector space than the dot-product of those vectors will give high variance
- high variance = big problem because softmax gives high probability to higher value & very small probability to small value. if that will create a gap between those numbers & due to this problem backpropagation will focus on larger numbers. if small numbers gets ignored & gradients of those parameters faces vanishing gradient problem & that will not let those parameters update so we need to reduce variance

→ We can divide each value of matrix with a constant one & that will cause low variance.

(Scaling factor)

number of dimensions & variance has linear relationship

in ideal scenario, variance must remain constant even dimension increases.

A mathematical rule.

If you have a random variable X with a variance of $\text{Var}(X)$, & you create a new variable Y by scaling X with a constant c , so that $Y = cX$, the variance of Y ($\text{Var}(Y)$) is related to variance of X by the square of scaling factor c . Mathematically this relationship is expressed as.

$$\text{Var}(Y) = c^2 \text{Var}(X)$$

$X \rightarrow \text{Variance} = \text{Var}(X)$

$Y \rightarrow$ created by scaling X with
constant c .

∴ Hence, variance of $Y = c^2 \text{Var}(X)$

→ Now we need to divide it by such constant number which will always give

$\text{Var}(X)$ in off regardless of number of dimensions.

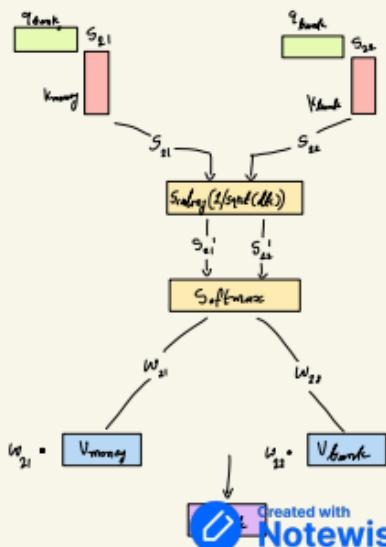
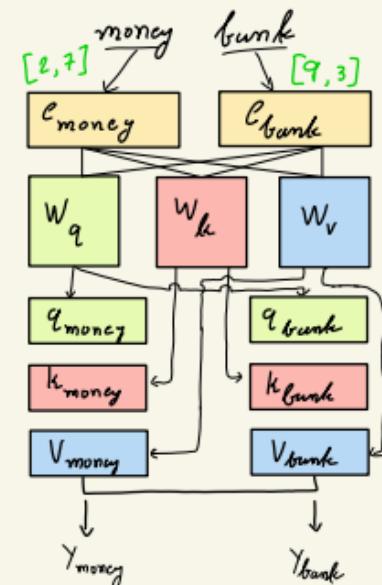
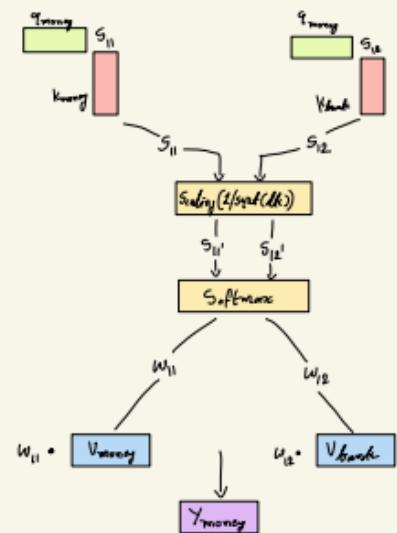
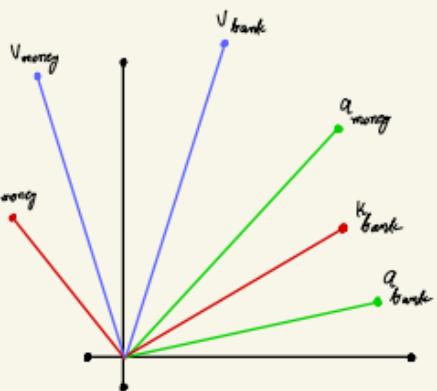
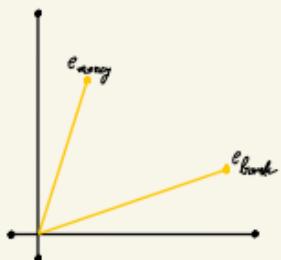
hence we will divide our vector by $\frac{1}{\sqrt{d_n}}$.

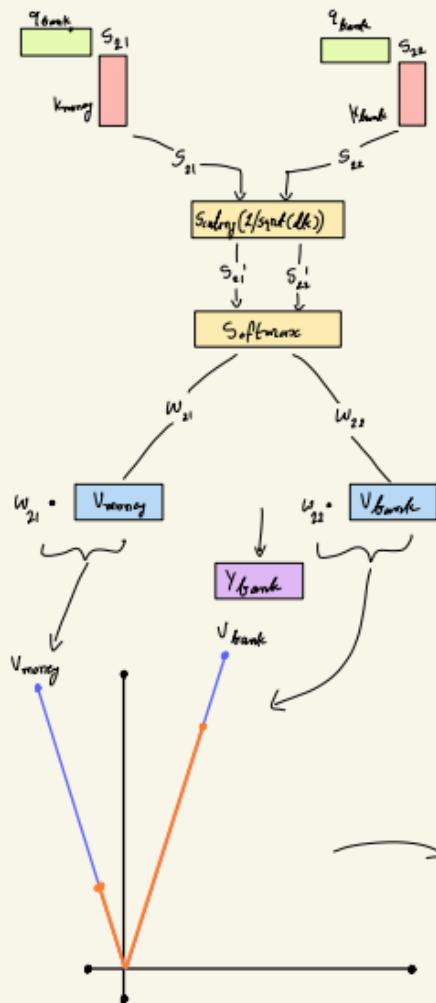


Geometric Intuition

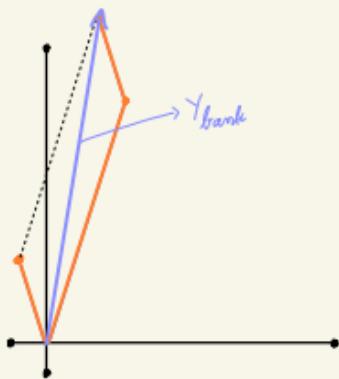
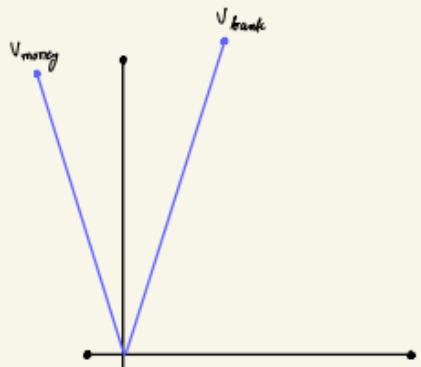
$$W_q = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \quad W_k = \begin{bmatrix} 3 & 4 \\ 5 & 2 \end{bmatrix} \quad W_v = \begin{bmatrix} 4 & 2 \\ 2 & 2 \end{bmatrix}$$

All these values are hypothetical.





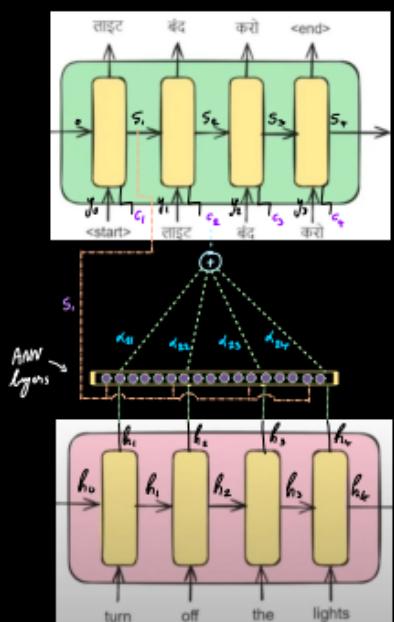
Dot product	$S_{z1} = 20$	$S_{z2} = 38$
Scaling		
	$S_{z1}^1 = \frac{f_0}{\sqrt{1}} = 7.07$	$S_{z2}^1 = \frac{2.2}{\sqrt{1}} = 22.69$
Softmax	$w_{21} = 0.2$	$w_{22} = 0.8$



And that is the brilliance of self-attention.



Q. Now is self-attention similar to Bahdanau & Luong Attention..?



Self-attention



$$\text{Here } Y_{\text{turn}} = W_{11} V_{\text{turn}} + W_{12} V_{\text{off}} + W_{13} V_{\text{the}} + W_{14} V_{\text{light}}$$

Since we're calculating attention between same sequence, hence it is called as "Self-attention".

Disadvantage

→ Self-attention can only figure-out single perspective from given sentence.

Example

- The man saw the astronomer with a telescope.

→ Meaning-1 : A man saw an astronomer on his telescope.

→ Meaning-2 : A man saw an astronomer working with his telescope.

→ We want such mechanism which can figure-out maximum possible perspectives from given sentences/paragraphs.

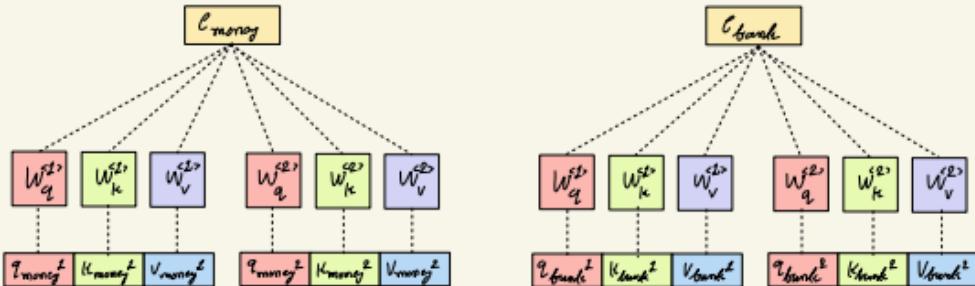
→ Hence, there is a new solution called as Multi-head attention.

→ Multi-head attention gives simple solution use multiple self-attention architecture. It is called as multi-head attention.

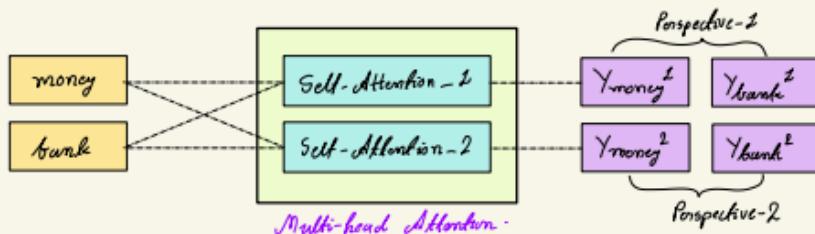
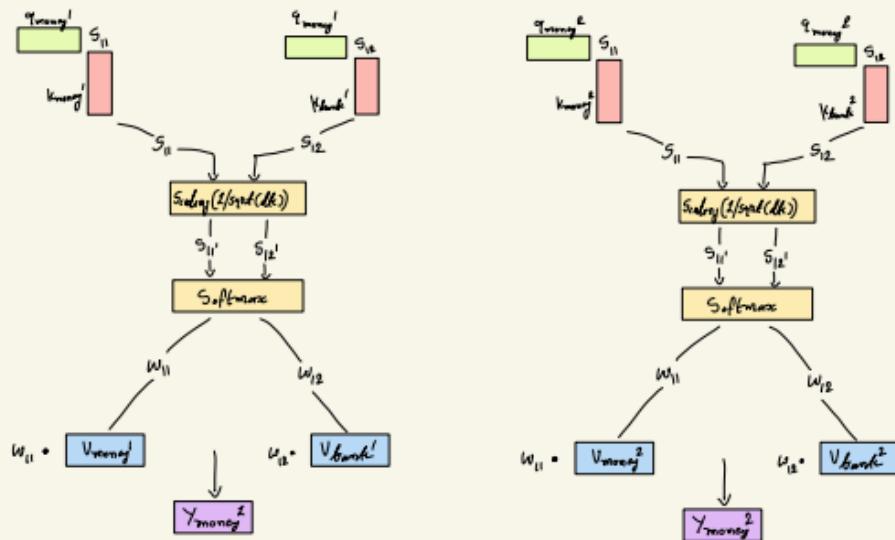
→ For multi-head architecture, we will expand the layer of weights, for each self-attention architecture.

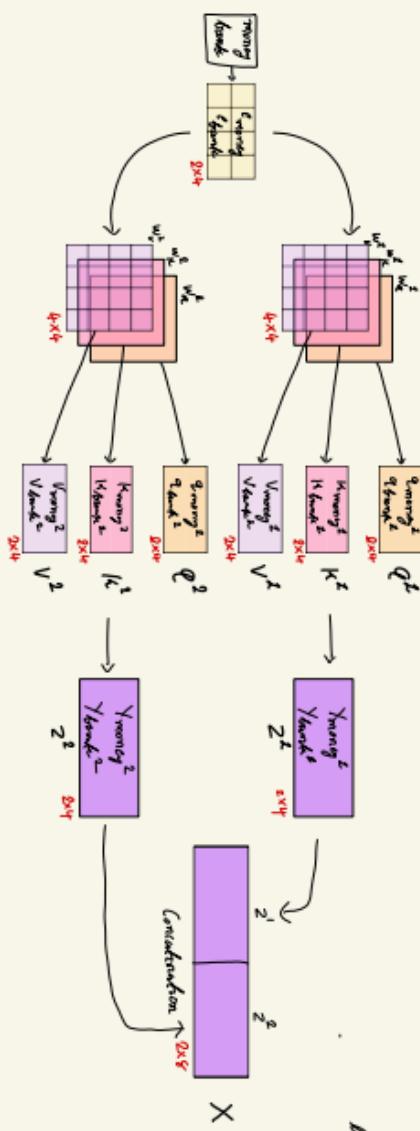


Dual Self-attention

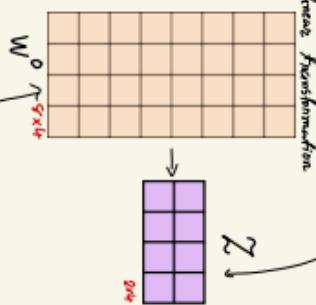


- Now next process will execute twice, due to 2 self-attention architecture.





To maintain its original matrix's shape



Positional Encoding

- Self-attention mechanism can able to handle all words of sentence simultaneously.
- But this feature caused a disadvantage that self-attention mechanism can't capture the order of words in given seq sentence.
- hence, top self-attention mechanism
 - Parth killed lion and
 - Lion killed parth are same. self-attention doesn't know about how words are exist in seq sentence.
- hence positional encoding encodes the position of each word in given sentence.
- The simple way is to append index number at the end of the context vector.
- This solution is unbounded (No upper limit) because, if the index-number of a particular word is huge than it will cause problems like exploding gradient prob.
- We can normalize it by dividing each word index with total-number-of-words.
- But this problem will confuse neural network because due to this approach, there will be different normalized-index of same position of different sentences.
- position number shouldn't change to ensure consistent behaviour across training samples.
- Also we can't catch relative positioning in given sentence.

Parth killed or Lion

1 2 3 4

here we're assigning unique positions to each words it is called as absolute positioning. hence model can't be able to calculate the distance between the and Parth. but we need **Relative Positioning**.

- So, instead of discrete function, it would be better if there is a periodic function.
- This periodic function gives bounded values. also it should continuous & periodic.
- for that we will use **sine** function.

$$Y = \sin(\text{position})$$

$$\sin(1) = 0.84$$

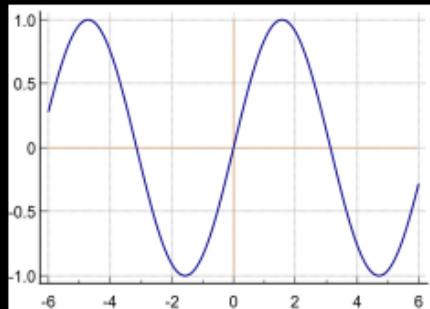
$$\sin(2) = 0.87$$

$$\sin(3) = 0.19$$

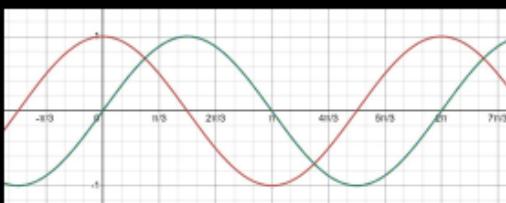
→ But still this method will cause the problem because the positional values are repeating.

→ One solution is to use one more trigonometric function with sin.

$$Y = \sin(\text{position}) \text{ and } Y = \cos(\text{position})$$



→ So we will represent each positional value as a vector instead of single value.



Ex sentence

Parth keetted bin

$$Y_1 = \sin(x) = 0.54$$

$$Y_2 = \cos(x) = 0.50$$

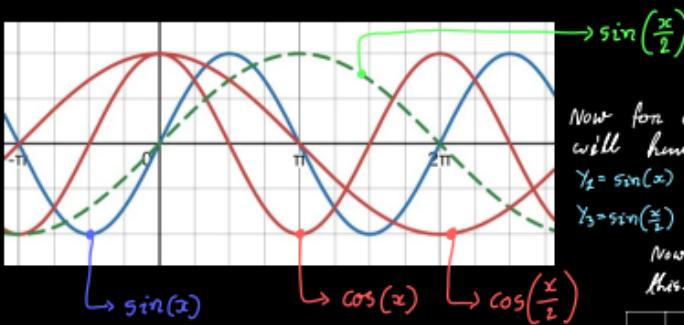
Now the embedding of a word Parth at index position 2 will look like this.

$$\text{Parth} = [0.11 \ 0.23 \ 0.7 \ 0.49 \ 0.27 \ 0.54 \ 0.50]$$

→ Same as we will calculate position of all words & append it at the end of their embeddings.

→ But still it is possible that sin and cos function can also give same value for a particular word.

→ For that we can also add more curves(functions).



Now for each position we will have 4 values.

$$Y_1 = \sin(x) \quad Y_2 = \cos(x)$$

$$Y_3 = \sin\left(\frac{x}{2}\right) \quad Y_4 = \cos\left(\frac{x}{2}\right)$$

Now embedding will look like this.

$$[0.11 \ 0.23 \ 0.7 \ 0.49 \ 0.27 \ 0.54 \ 0.50 \ 0.94 \ 0.74]$$

→ In case of masking we can also add one more layer like $Y = \sin(x/3)$ and $Y = \cos(x/3)$ and continue.....

Points must be remember

- Each positional encoding will be a vector.
- The dimensions of positional encoded vector will be same as dimension of embedding vector.

$$[0.26 \ 0.25 \ 0.36 \ 0.8 \ 0.01 \ 0.21 \ 0.33 \ 0.48]$$

→ But according to original research paper "Attention is all you need" we will perform elementwise addition between positional vector & encoding vector.

$$[0.26 \ 0.25 \ 0.36 \ 0.8] + [0.01 \ 0.21 \ 0.33 \ 0.48] = [0.27 \ 0.36 \ 0.69 \ 0.56]$$



- Now this combined vector will be sent to self-attention mechanism.
- Previous vector concatenation approach is clearly increasing number of dimensions. This will lead to slower training time.
- As we have discussed previously that....

"The dimensions of positional encoded vector will be same as dimension of embedding vector".

Suppose embedding vector is a 6D vector. that means positional encoder will also be a 6D vector. and let's see how we will fill 6D positional vector.

Embedding vector

$\begin{bmatrix} 0.6 & 0.5 & 0.1 & 0.9 & 0.6 & 0.7 \end{bmatrix}$



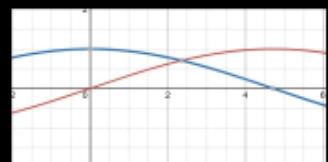
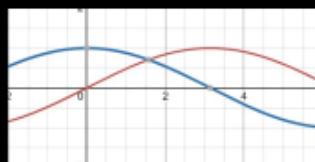
we will work in pairs with sine & cosine curve. for 1st pair $y_1 = \sin(x)$ & $y_2 = \cos(x)$.
 for 2nd pair $y_1 = \sin(2x)$ & $y_2 = \cos(2x)$.
 for 3rd pair $y_1 = \sin(4x)$ & $y_2 = \cos(4x)$.

- In original research paper embedding vector is of 512D, hence positional encoding will also be 512D.
- As we increases pairs, the frequency of next sin, cosine curve frequency gets reduced.

$$y = \sin(x) \quad y = \cos(x)$$

$$y = \sin\left(\frac{x}{2}\right) \quad y = \cos\left(\frac{x}{2}\right)$$

$$y = \sin\left(\frac{x}{3}\right) \quad y = \cos\left(\frac{x}{3}\right)$$



Q. How frequency will be divided...?

- At which rate, we will reduce frequency.
- Till now we have used rational numbers to reduce frequency.

Based on "Attention is all you need"

$$PE_{(pos, 2i)} \overset{\text{position}}{=} \sin\left(\frac{pos \cdot 2i}{20000} \cdot d_{model}\right) \quad PE_{(pos, 2i+1)} = \cos\left(\frac{pos \cdot 2i+1}{20000} \cdot d_{model}\right)$$

dimensionality of embeddings

Let's see an example of a word "River bank".

1 2

Suppose we're working with 6D embedding vector & positional encoder.

hence, frequency(i) = $\frac{6}{2} = ③$

6: Total dimension

2: Pair



"River"

for $i=0$

$$PE(0,0) = \sin\left(\frac{0}{20000} \cdot 2\pi\right) = 0$$

$$PE(0,1) = \cos\left(\frac{0}{20000} \cdot 2\pi\right) = 1$$

"Bunle"

for $i=0$

$$PE(0,0) = \sin\left(\frac{0}{20000} \cdot 2\pi\right) = 0.44$$

$$PE(0,1) = \cos\left(\frac{0}{20000} \cdot 2\pi\right) = 0.94$$

for $i=1$

$$PE(0,2) = \sin\left(\frac{1}{20000} \cdot 2\pi\right) = 0$$

$$PE(0,3) = \cos\left(\frac{1}{20000} \cdot 2\pi\right) = 1$$

for $i=1$

$$PE(1,0) = \sin\left(\frac{0}{20000} \cdot 2\pi\right) = 0.44$$

$$PE(1,1) = \cos\left(\frac{0}{20000} \cdot 2\pi\right) = 0.94$$

for $i=2$

$$PE(0,4) = \sin\left(\frac{2}{20000} \cdot 2\pi\right) = 0$$

$$PE(0,5) = \cos\left(\frac{2}{20000} \cdot 2\pi\right) = 1$$

for $i=2$

$$PE(1,0) = \sin\left(\frac{0}{20000} \cdot 2\pi\right) = 0.44$$

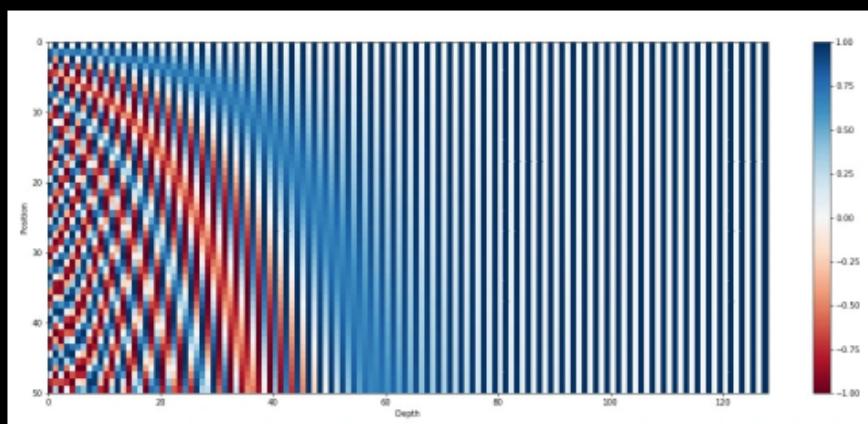
$$PE(1,1) = \cos\left(\frac{0}{20000} \cdot 2\pi\right) = 0.94$$

0	1	0	1	0	1
---	---	---	---	---	---

"River"

0.44	0.54	0.44	0.44	0.44	0.44
------	------	------	------	------	------

"Bunle"



Based on above graph

→ Total words in sentence : 50

→ Embedding dimension of each word : 228 ($d_{model} = 228$)

→ Total positional encodings : 50

→ Dimension of each positional encoding : 228

→ We can see that majority changes are happening on lower dimensions (Depth 0-75)



Created with
Notewise

→ for higher dimensions (Depth > 75) there is no change hence we're not reducing frequency as number of dimensions increases.

Q. How sine, cosine can capture relative positioning?

→ Due to this sine, cosine set-up, our positional encoding token got an interesting property.

→ Let's imagine each 50 words as vectors (V_i). Now if we apply linear transformation on V_{20} then we will reach at V_{20} . some transformation matrix if we apply on V_{20} then we will reach at V_{30} .

→ This particular linear transformation matrix has capability of covering delta of 10 on distance. & this property is available for each delta.

Example

$M_2 = \text{Matrix} \rightarrow$ This matrix M_2 has capability of covering delta/distance of 5.

$V_5 \rightarrow M_2 \rightarrow V_{20}$ This feature helps to understand relative positioning.

$V_{22} \rightarrow M_2 \rightarrow V_{27}$

$V_{44} \rightarrow M_2 \rightarrow V_{49}$

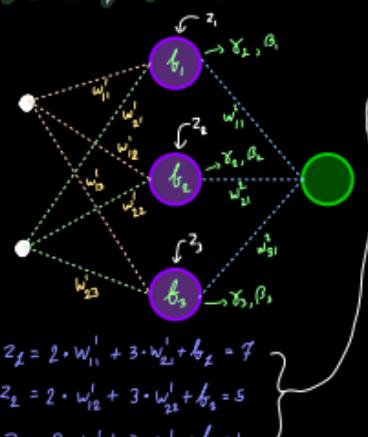
Layer Normalization

Q. What is normalization?

→ In deep learning, normalization refers to the process of transforming data or model op to have specific properties, typically $\mu=0, \sigma=1$.

→ Layer Normalization normalizes activations (Hidden layer).

Quick recap of batch normalization



batch-size = 5

f_1	f_2	z_1	z_2	z_3
2	3	7	5	4
1	1	2	3	4
5	4	1	2	3
6	1	7	5	6
7	1	3	3	4

$$\frac{7 - \bar{u}_2}{\sigma_2} = 0.36 \cdot \bar{y}_2 + \beta_2 = 0.36 \rightarrow \text{activation}$$

$$\frac{2 - \bar{u}_2}{\sigma_2} = 0.71 \cdot \bar{y}_2 + \beta_2 = 0.71$$

$$\frac{5 - \bar{u}_2}{\sigma_2} = -0.21 \cdot \bar{y}_2 + \beta_2 = -0.21$$

$$\frac{6 - \bar{u}_2}{\sigma_2} = 0.11 \cdot \bar{y}_2 + \beta_2 = 0.11$$

- Hence \bar{y}_2 and β_2 are learnable parameters whose critical values will be 1 and 0 respectively.



- We can't use batch normalization on self-attention mechanism.
- We can also send multiple sentences in self-attention mechanism in the form of batches.

Example

Review	Sentiment
Hi pants	1
How are you today	0
I am just	0
You?	1

Embedding - dim = 3

batch-size = 2 (At a time, 2 sentences will be passed in self-attention mechanism)

Step-1

- Embedding vector of first 2 sentences

0.1	0.45	0.71
"Hi"		

0.21	0.3	0.8
"Pants"		

- We need to apply padding on first sentence, because we can clearly see that both sentences have different number of words.

0.1	0.5	0.51
"How"		

0.1	0.0	0.15
"are"		

0.33	0.56	0.9
"you"		

0.11	0.4	0.34
"today"		

0.2	0.45	0.71
"Hi"		

0.21	0.3	0.8
"Pants"		

0	0	0
PADDING		

0	0	0
PADDING		

0.1	0.5	0.51
"How"		

0.1	0.0	0.25
"are"		

0.33	0.56	0.9
"you"		

0.11	0.4	0.34
"today"		

"Hi"	6.5	2.41	3.21
"Pants"	2.21	0.4	3.6
Padding	0	0	0
Padding	0	0	0

Contextual embedding in single tensor of (1, 4, 3) of 2 sentences

7.5	9.2	2.5
2.2	1.1	6.7
2.9	6	9
9.9	2.3	6.5

"How"

"are"

"you"

"today"

↑ 4x8

↑ 4x3

Self-attention

"Hi"	0.2	0.45	0.71
"Pants"	0.21	0.3	0.8
Padding	0	0	0
Padding	0	0	0

Single tensor of (2, 4, 3)

0.1	0.5	0.34
0.1	0.0	0.25
0.33	0.56	0.9
0.11	0.4	0.34

"How"

"are"

"you"

"today"



→ Now we need to normalize o/p matrix. for that we will scale-up both o/p matrix & perform normalization as shown done previously. each column will have it's own

	d_1	d_2	d_3
"Hi"	6.5	2.4	3.2
"Punkt"	2.21	0.4	3.6
"Grußgott"	0	0	0
"Grußgott"	0	0	0
"How"	7.5	9.2	2.5
"are"	2.2	1.1	6.7
"you"	2.9	6	9
"tasty"	9.9	2.3	6.3

$$\begin{matrix} \mu_1 & \mu_2 & \mu_3 \\ \sigma_1 & \sigma_2 & \sigma_3 \\ \gamma_1 & \gamma_2 & \gamma_3 \\ \beta_1 & \beta_2 & \beta_3 \end{matrix}$$

So the solution is Layer Normalization.....

→ in layer Normalization we will apply normalization across features (Horizontal direction)

Example

f_1	f_2	Z_1	Z_2	Z_3
2	3	7	5	4
1	1	2	3	4
5	4	1	2	3
6	1	7	5	6
7	1	3	3	4

only of Z_3

$$\begin{array}{lll} \mu_1, \sigma_1 & \left(\frac{7 - \mu_1}{\sigma_1} \right) Z_1 + \beta_1 & \left(\frac{2 - \mu_2}{\sigma_2} \right) Z_2 + \beta_2 \\ \mu_2, \sigma_2 & \left(\frac{5 - \mu_2}{\sigma_2} \right) Z_1 + \beta_2 & \left(\frac{3 - \mu_3}{\sigma_3} \right) Z_2 + \beta_3 \\ \mu_1, \sigma_3 & \left(\frac{1 - \mu_1}{\sigma_3} \right) Z_1 + \beta_3 & \left(\frac{6 - \mu_3}{\sigma_3} \right) Z_3 + \beta_3 \end{array}$$

and so on.....

Now let's see with embedding example



	d_1	d_2	d_3
"He"	6.5	2.41	3.21
"Punk"	2.21	0.4	3.6
"hitting"	0	0	0
"hitting"	0	0	0
"Now"	7.5	4.2	1.5
"are"	2.2	1.1	6.7
"you"	2.9	6	9
"fully"	9.9	2.3	6.3

$$u_1, \sigma_1 \quad \left(\frac{6.5 - u_1}{\sigma_1} \right) \gamma_1 + \beta_1 \quad \left(\frac{7.5 - u_5}{\sigma_5} \right) \gamma_5 + \beta_5$$

$$u_2, \sigma_2 \quad \left(\frac{2.21 - u_2}{\sigma_2} \right) \gamma_2 + \beta_2 \quad \left(\frac{9 - u_7}{\sigma_7} \right) \gamma_7 + \beta_7$$

$$u_3, \sigma_3 \quad \text{In this process padding values now won't make any impact.}$$

$$u_4, \sigma_4 \quad \text{and that is very good benefit.}$$

$$u_5, \sigma_5$$

$$u_6, \sigma_6$$

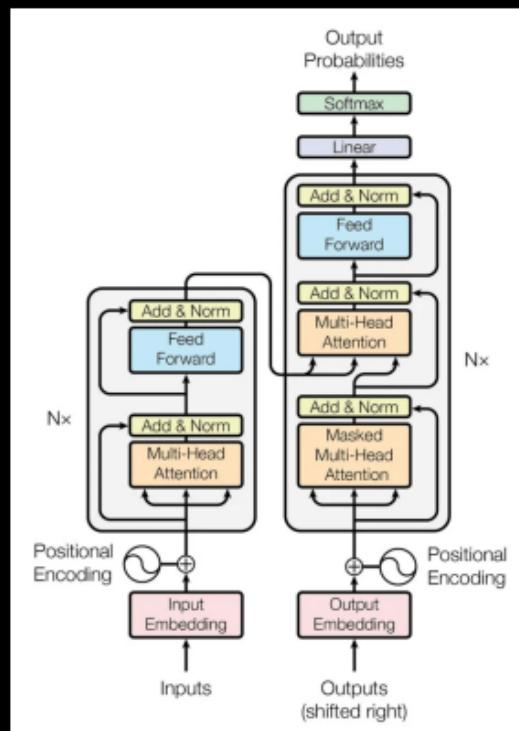
$$u_7, \sigma_7$$

$$u_8, \sigma_8$$

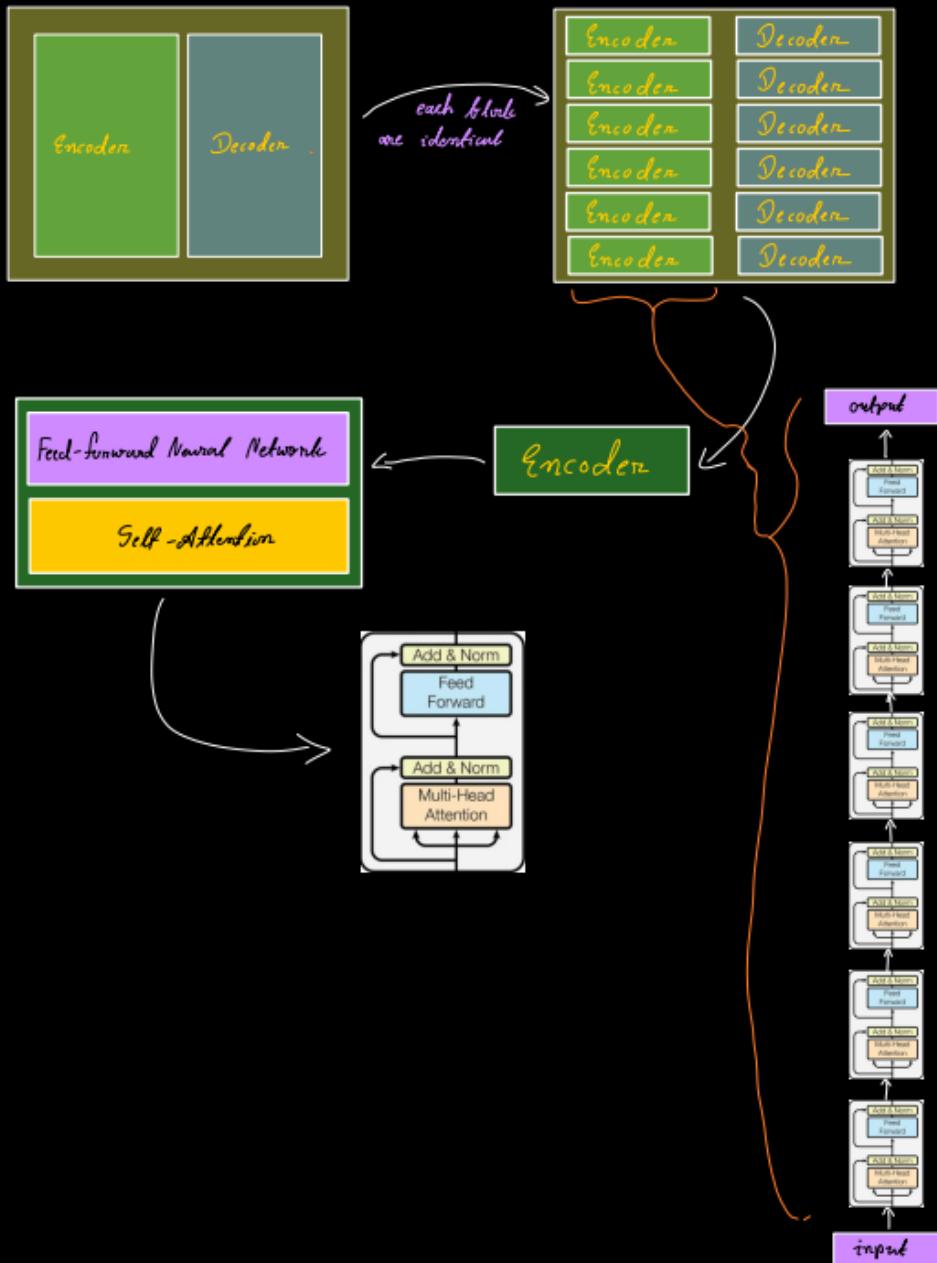
$$\begin{matrix} \gamma_1 \\ \gamma_2 \\ \gamma_3 \end{matrix}$$

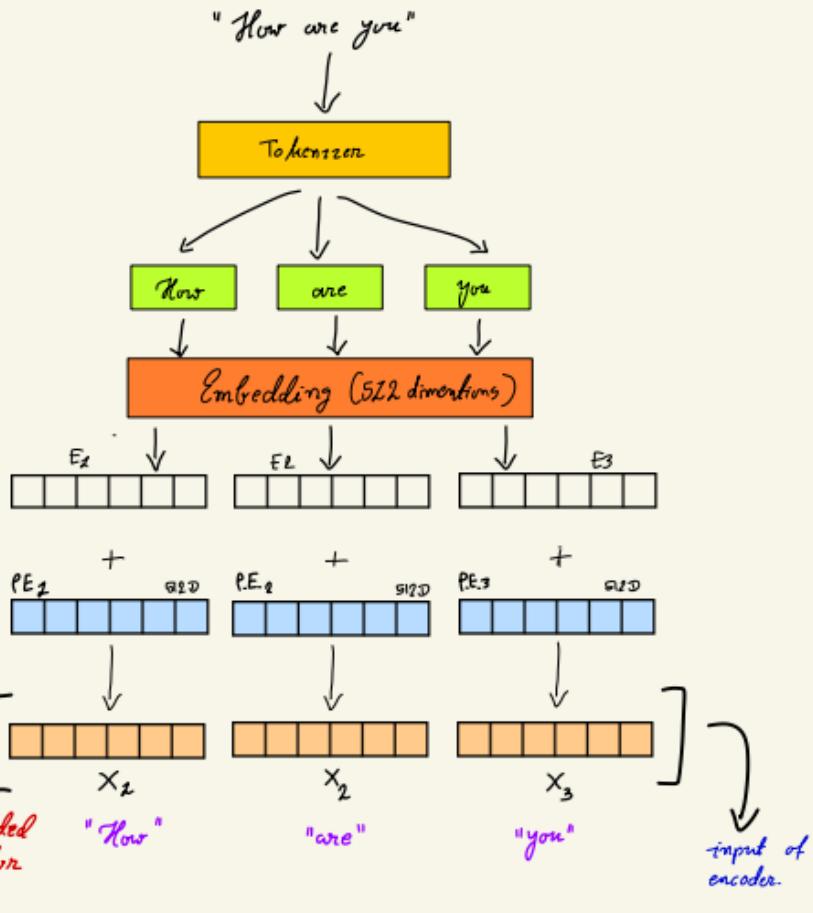
$$\begin{matrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{matrix}$$

Transformer Architecture



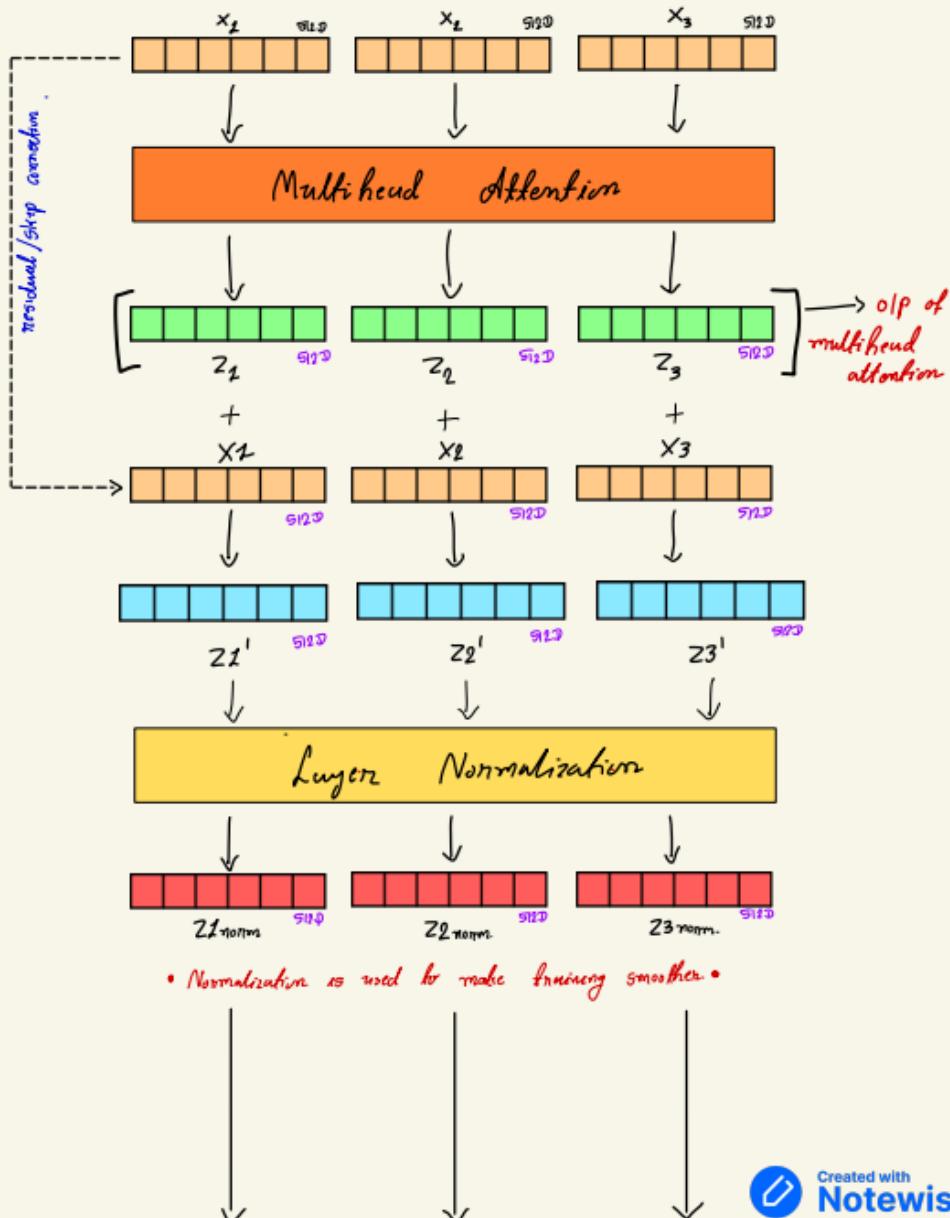
Simplified representation

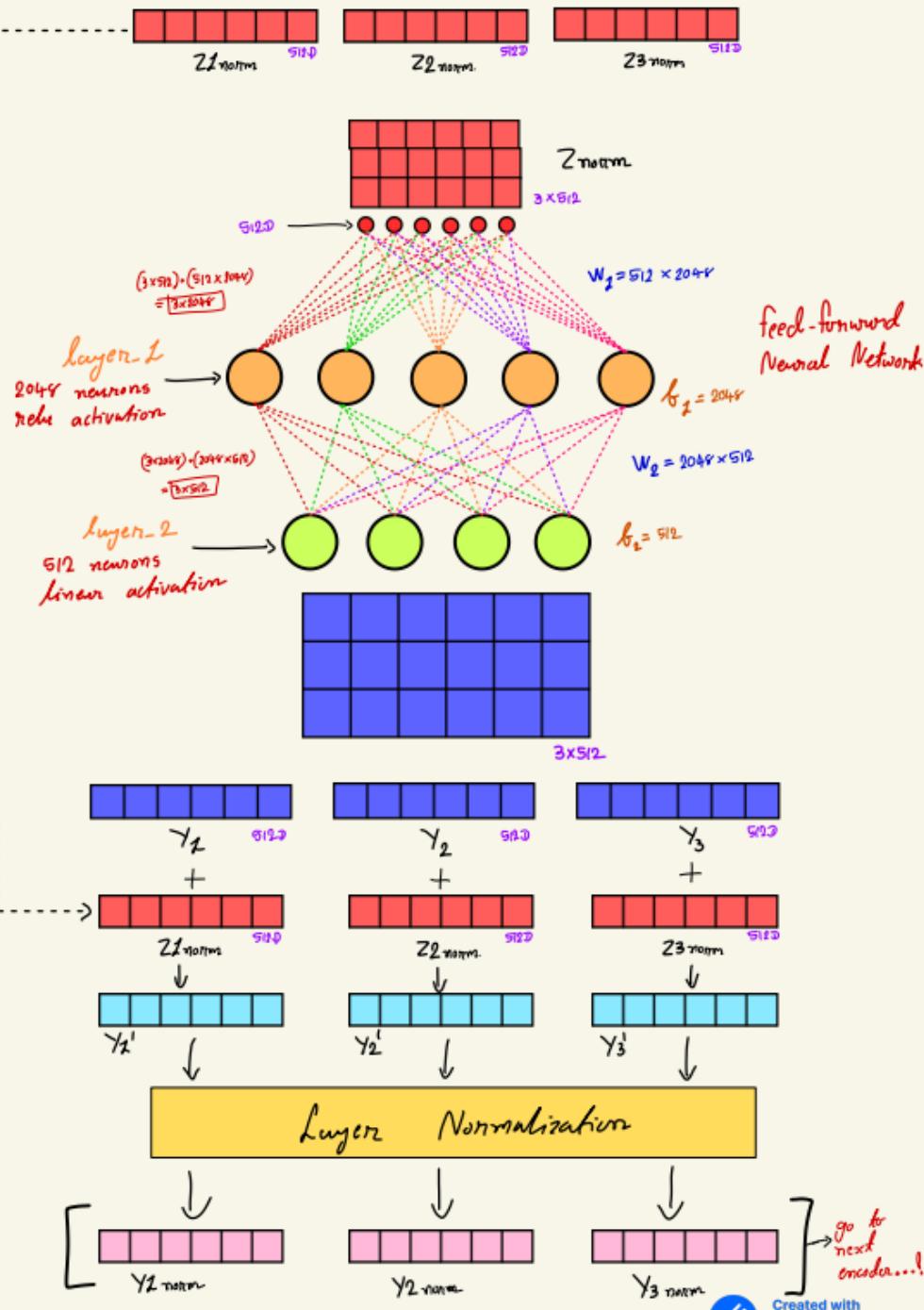




This is what input block does before passing data to encoder architecture for further process.







→ Each encoder block has separate trainable parameters, even though the architecture is same.

Q1. Why use residual connections?

- Stable training.
- Forwards original features as it is.

Q2. Why use a FFNN?

- To capture non-linearity/non-linear relationship.

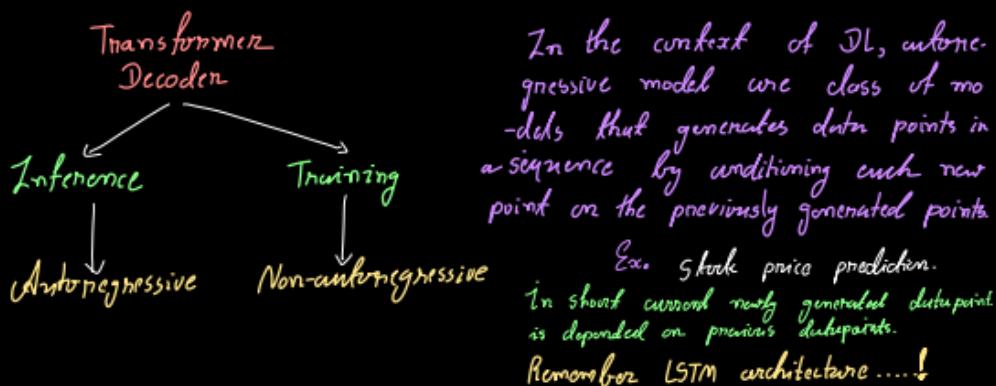
Q3. Why use 6 encoder blocks?

- To capture the context of complex language.

Autoregressive Model

The transformer decoder is autoregressive at inference time and non-autoregressive at training time.

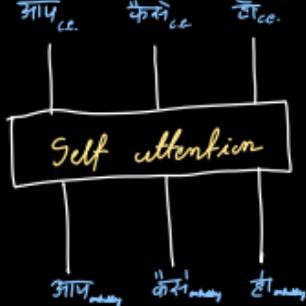
prediction



Transformer's decoders are autoregressive at inference time not at training time is due to Masked Self-attention.

Sentence Number	English Sentence	Hindi Sentence
1	How are you?	आप कैसे हैं?
2	Congratulations	तुमसे हाथ मिलाया जाएगा।
3	Thank you	धन्यवाद।





$$\text{आप}_\text{ce} = 0.8 \text{आप}_\text{में} + 0.1 \text{कैसे}_\text{में} + 0.1 \text{हो}_\text{में}$$

Hypothetical numbers

$$\text{कैसे}_\text{ce} = 0.15 \text{आप}_\text{में} + 0.75 \text{कैसे}_\text{में} + 0.1 \text{हो}_\text{में}$$

$$\text{हो}_\text{ce} = 0.1 \text{आप}_\text{में} + 0.2 \text{कैसे}_\text{में} + 0.7 \text{हो}_\text{में}$$

Big Problem: We're using future embedding (which is not even exist yet) for calculate current o/p token. This isn't possible at inference time.

Hence, we can't solve $\text{आप}_\text{ce} = 0.8 \text{आप}_\text{में} + 0.1 \text{कैसे}_\text{में} + 0.1 \text{हो}_\text{में}$ this equation, because at current stage we're calculating आप_ce and we've not even idea about next words, then we can't calculate over this formula.

This is data leakage problem in parallel approach.

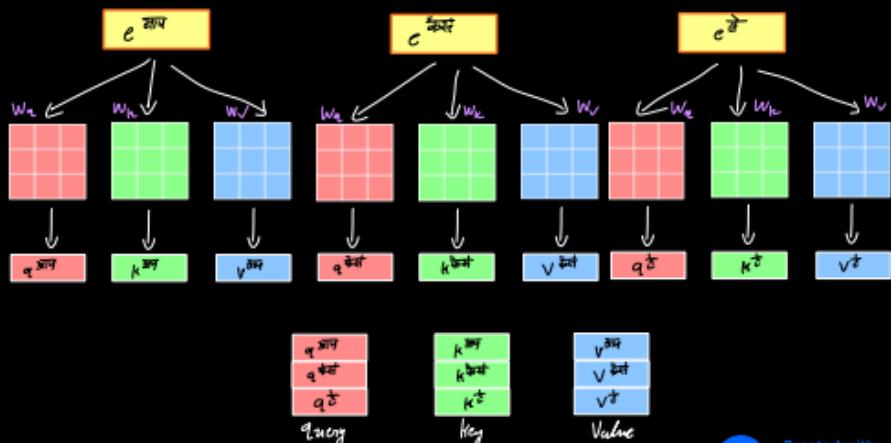
Auto-regressive approach

- ⇒ No data leakage
- ⇒ Slow training

Non-auto-regressive approach

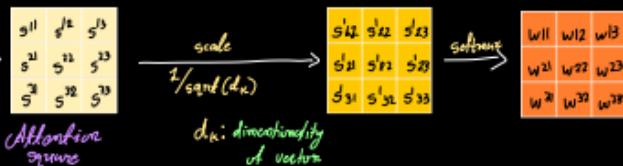
- ⇒ fast training
- ⇒ Data leakage (Your current token can see your future token which is zunkin if it is not possible at prediction time)

finding the answer



Query

Key

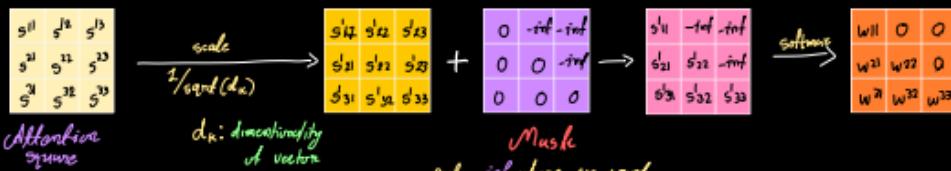


$$\overline{w_{11}y_{ce}} = W_{11} \cdot \boxed{v^{\frac{1}{\sqrt{d_k}}}} + W_{12} \cdot \boxed{v^{\frac{2}{\sqrt{d_k}}}} + W_{13} \cdot \boxed{v^{\frac{3}{\sqrt{d_k}}}}$$

$$\overline{w_{21}y_{ce}} = W_{21} \cdot \boxed{v^{\frac{1}{\sqrt{d_k}}}} + W_{22} \cdot \boxed{v^{\frac{2}{\sqrt{d_k}}}} + W_{23} \cdot \boxed{v^{\frac{3}{\sqrt{d_k}}}}$$

$$\overline{w_{31}y_{ce}} = W_{31} \cdot \boxed{v^{\frac{1}{\sqrt{d_k}}}} + W_{32} \cdot \boxed{v^{\frac{2}{\sqrt{d_k}}}} + W_{33} \cdot \boxed{v^{\frac{3}{\sqrt{d_k}}}}$$

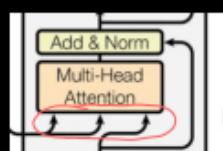
- Here, while calculating $\overline{w_{11}y_{ce}}$ we don't need to include $W_{12} \cdot \boxed{v^{\frac{2}{\sqrt{d_k}}}}$ and $W_{13} \cdot \boxed{v^{\frac{3}{\sqrt{d_k}}}}$ because they aren't in the picture yet so we will make W_{12} & $W_{13} = 0$.
- Similarly while calculating $\overline{w_{21}y_{ce}}$ we don't need $W_{23} \cdot \boxed{v^{\frac{3}{\sqrt{d_k}}}}$ because it isn't picture yet so W_{23} will be assigned 0.
- Here is the masking concept comes in picture to make $W_{12} = W_{13} = W_{23} = 0$.
- For that we will add one more step in above matrix operations.



That's how we achieved non-autoregressive modeling in training phase.

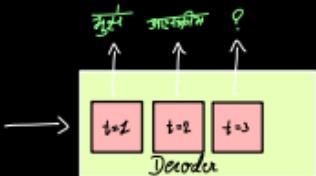


Gross Attention / Encoder-decoder Attention



- A mechanism used in transformer architecture, particularly in tasks involving seq.-to-seq. like translation or summarization. It allows a model to focus on different parts of an input seq. when generating an output seq.

Example of Machine translation



- at t=3 upcoming prediction depends on 2 things.

 1. Generated content till t=3.
 2. Incoming information from encoder.

- 1st problem will be solved by self-attention.
 - 2nd problem to get similarity between two sequences will be done by cross attention.

	कृष्ण	द्वादशवर्षीय	नववर्षा	प्रसंग	प्र.
I	●	●	●	●	●
like	●	●	●	●	●
eating	●	●	●	●	●
icecream	●	●	●	●	●

some sequence \rightarrow similarity
self-attention

Different sequences → similarity

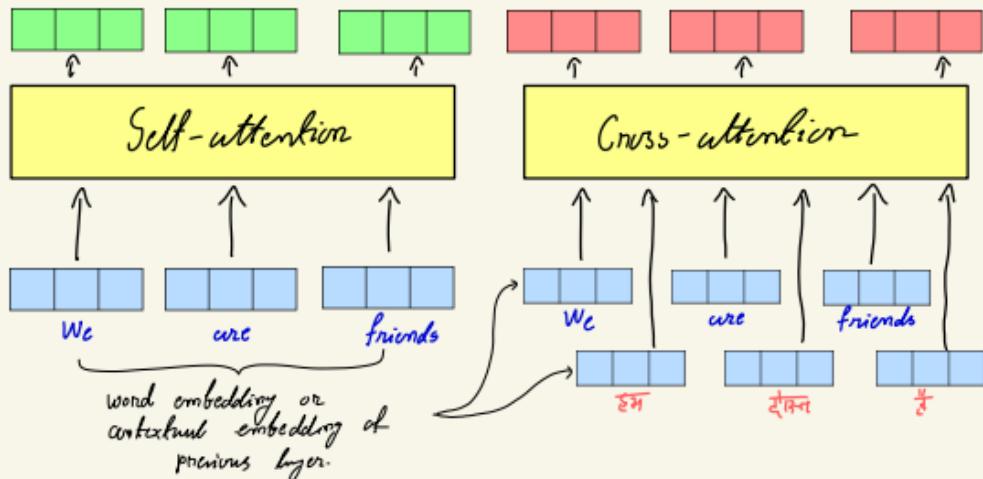
(cross-attention)
(calculating similarity
between different sequences)

Cross-attention is conceptually similar to self-attention.

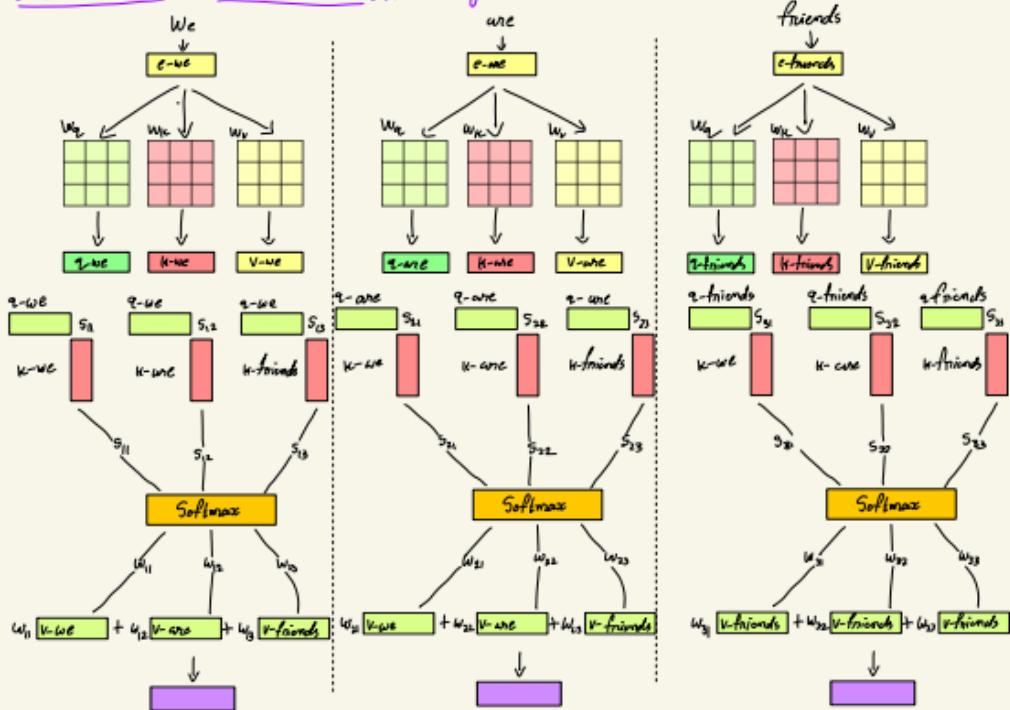
Self attention VS Cross - attention

1. The i/p.
 2. The processing.
 3. The o/p.

Self-attention VS Cross-attention (input)



Self-attention VS Cross-attention (Processing)



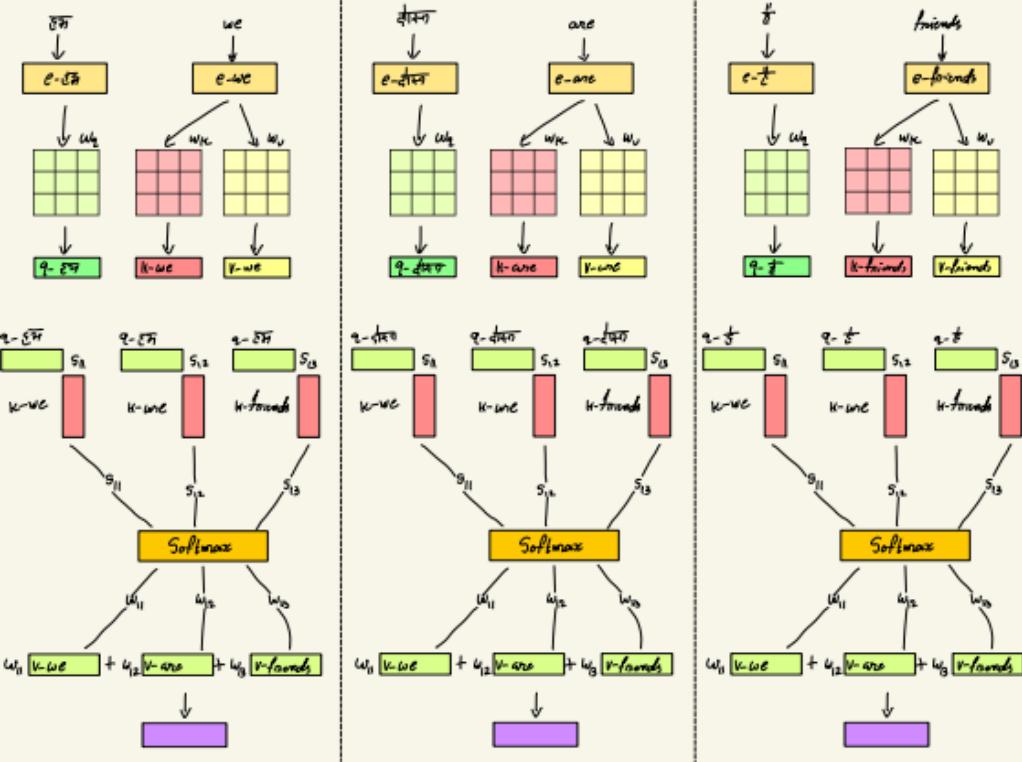
	we	are	friends
we	●	●	●
are	●	●	●
friends	●	●	●

Imaginary similarity



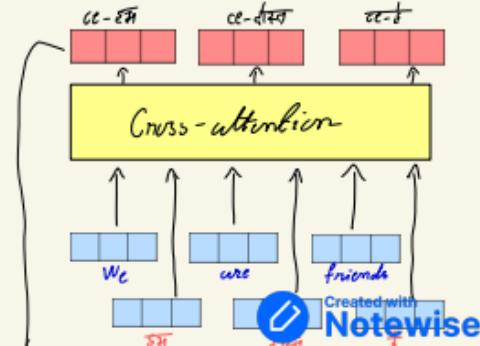
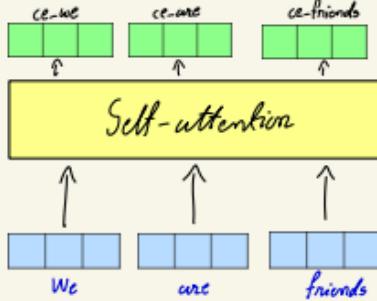
query → Q/p vector
key / Value → K/p vector

Cross-attention



	we	are	friends
we	●	•	•
are	•	-	●
friends	•	●	•

Self-attention Vs Cross-attention (Output)



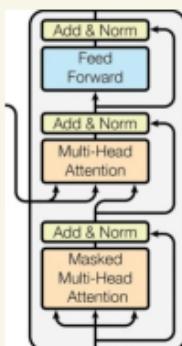
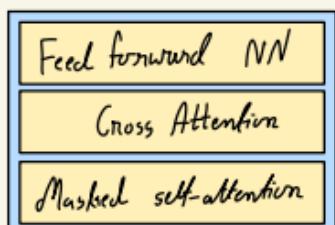
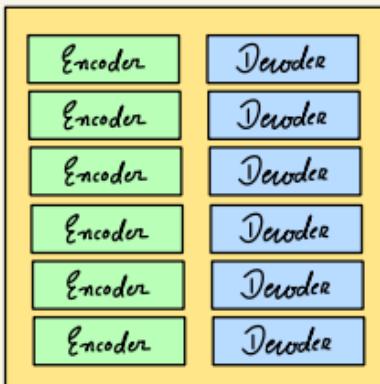
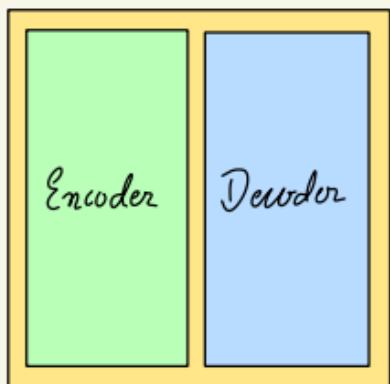
$$\begin{aligned}
 ce_wc &= 0.8 * c_wc + 0.1 * e_wc + 0.1 * c_friends \\
 ce_one &= 0.4 * c_wc + 0.3 * e_wc + 0.3 * c_friends \\
 ce_friends &= 0.2 * c_wc + 0.2 * e_wc + 0.2 * c_friends
 \end{aligned}$$

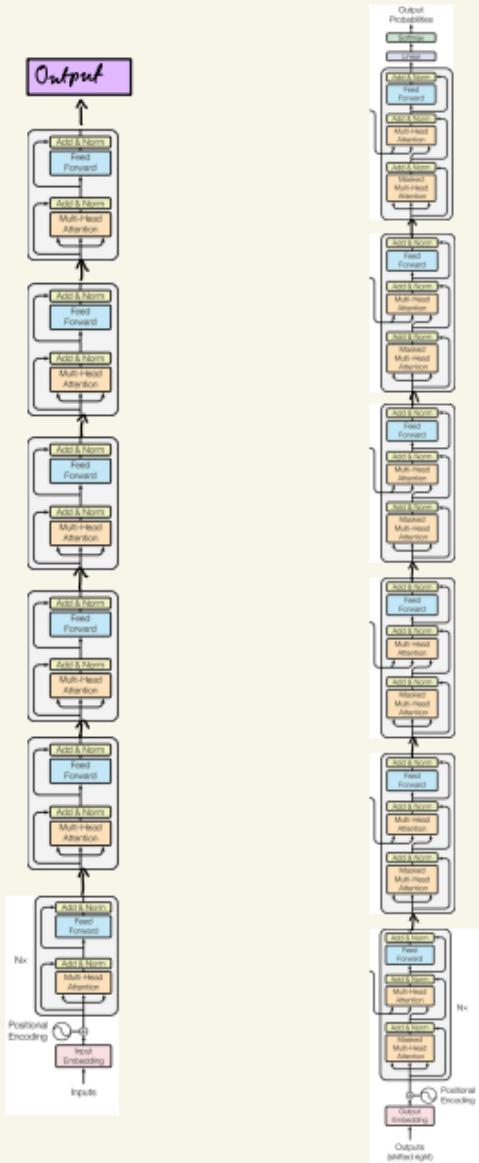
Number of o/p will be same as number of incoming o/p label

$$\begin{aligned}
 ce_st &= 0.5 * c_wc + 0.3 * e_wc + 0.2 * c_friends \\
 ce_two &= 0.2 * c_wc + 0.2 * e_wc + 0.6 * c_friends \\
 ce_t &= 0.3 * c_wc + 0.4 * e_wc + 0.3 * c_friends
 \end{aligned}$$

Use cases

- ⇒ Image - to - text
- ⇒ Text - to - image
- ⇒ Text - to - speech



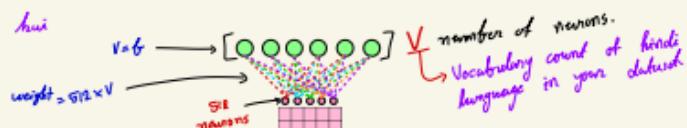


Machine translation task

English - We are friends

Hindi - हम दोस्त हैं

During Training Phase

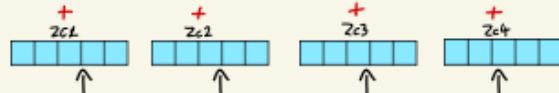
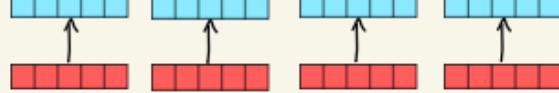
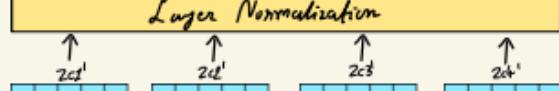
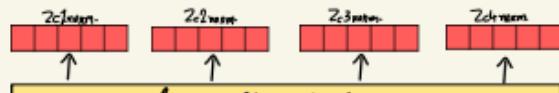
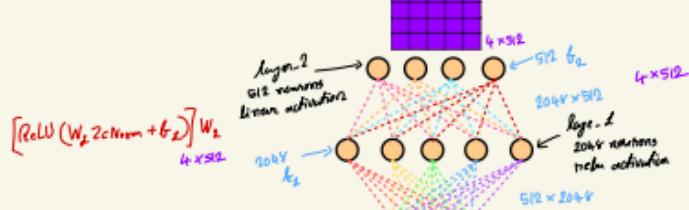
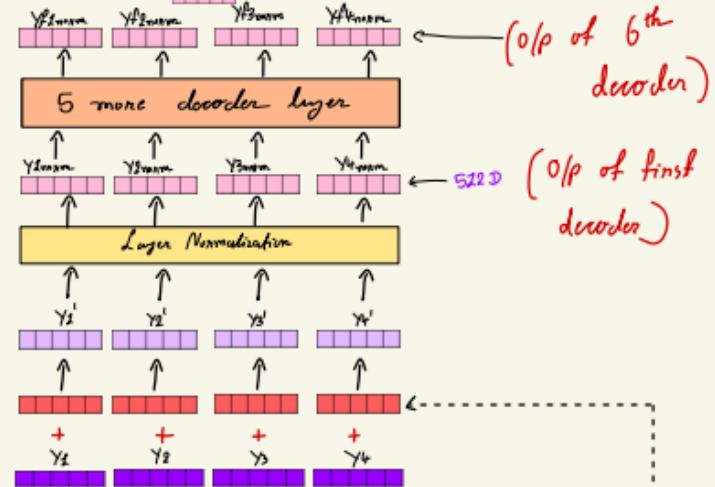


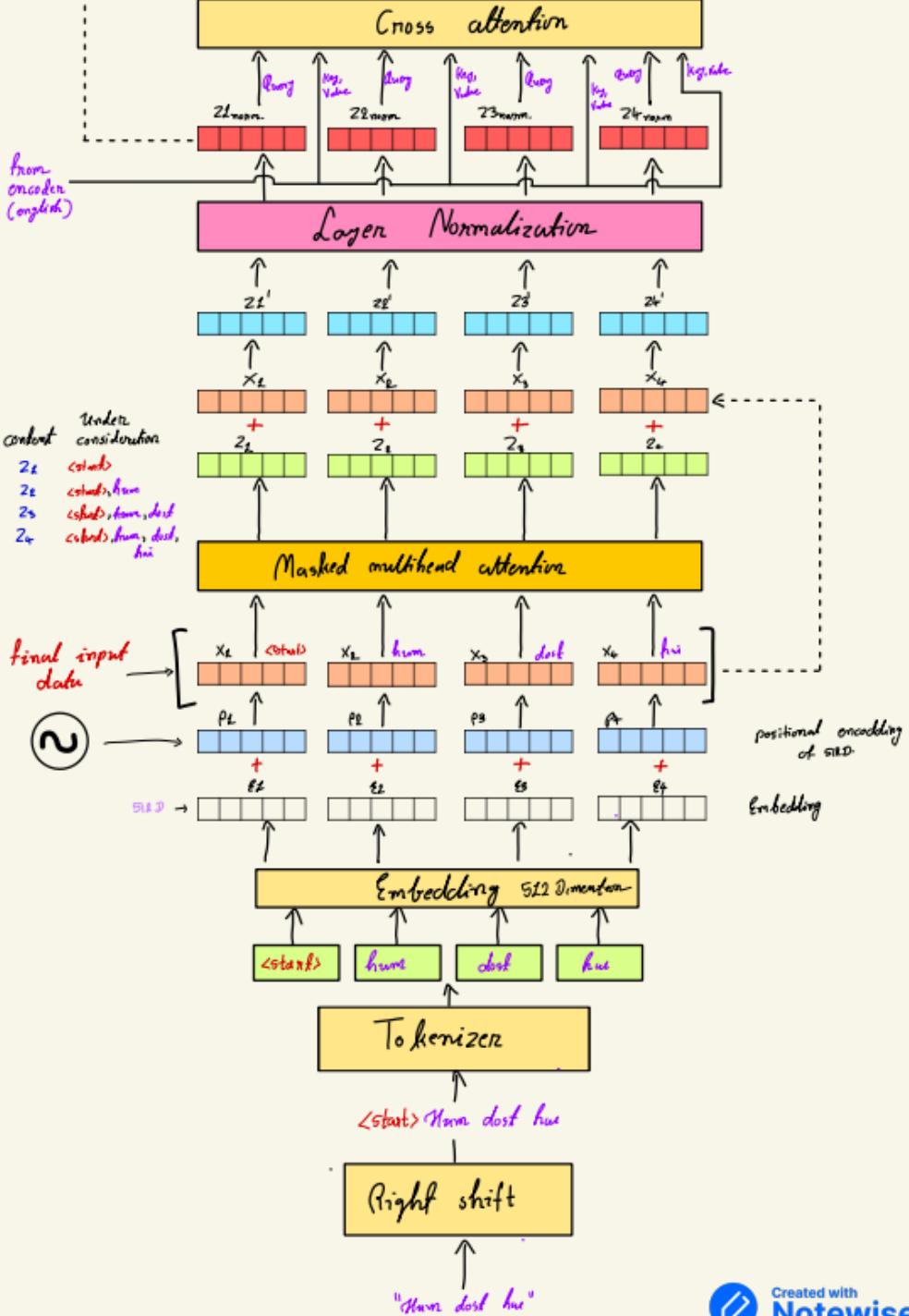
Input operation

1. Shifting
2. Tokenization
3. Embedding
4. Positional encoding

Decoder operations

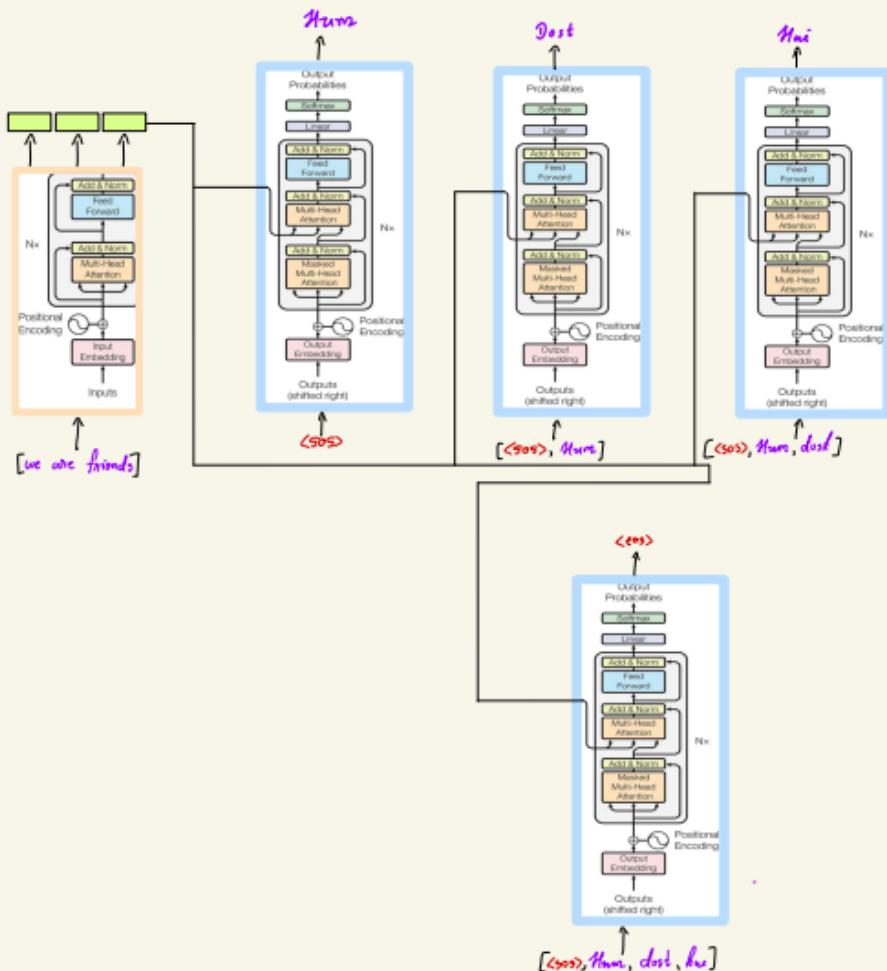
1. Masking
2. Multilevel attention
3. FFNN



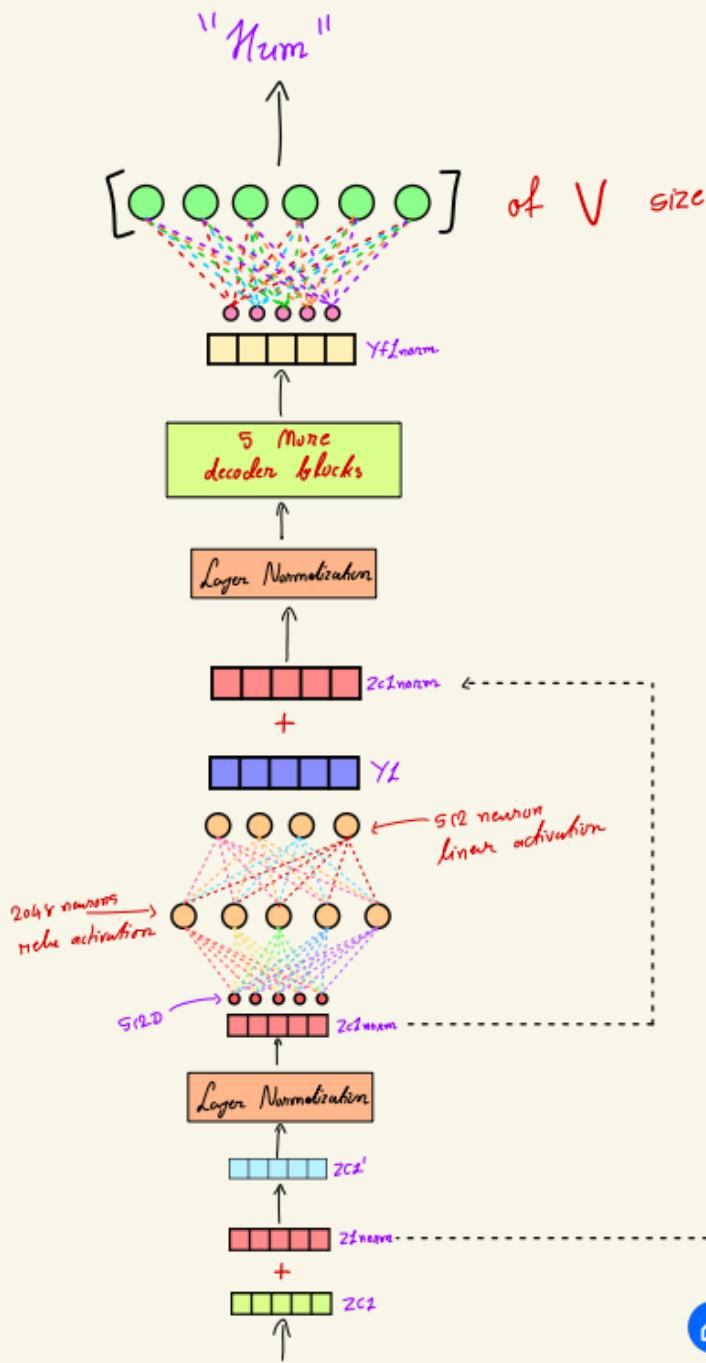


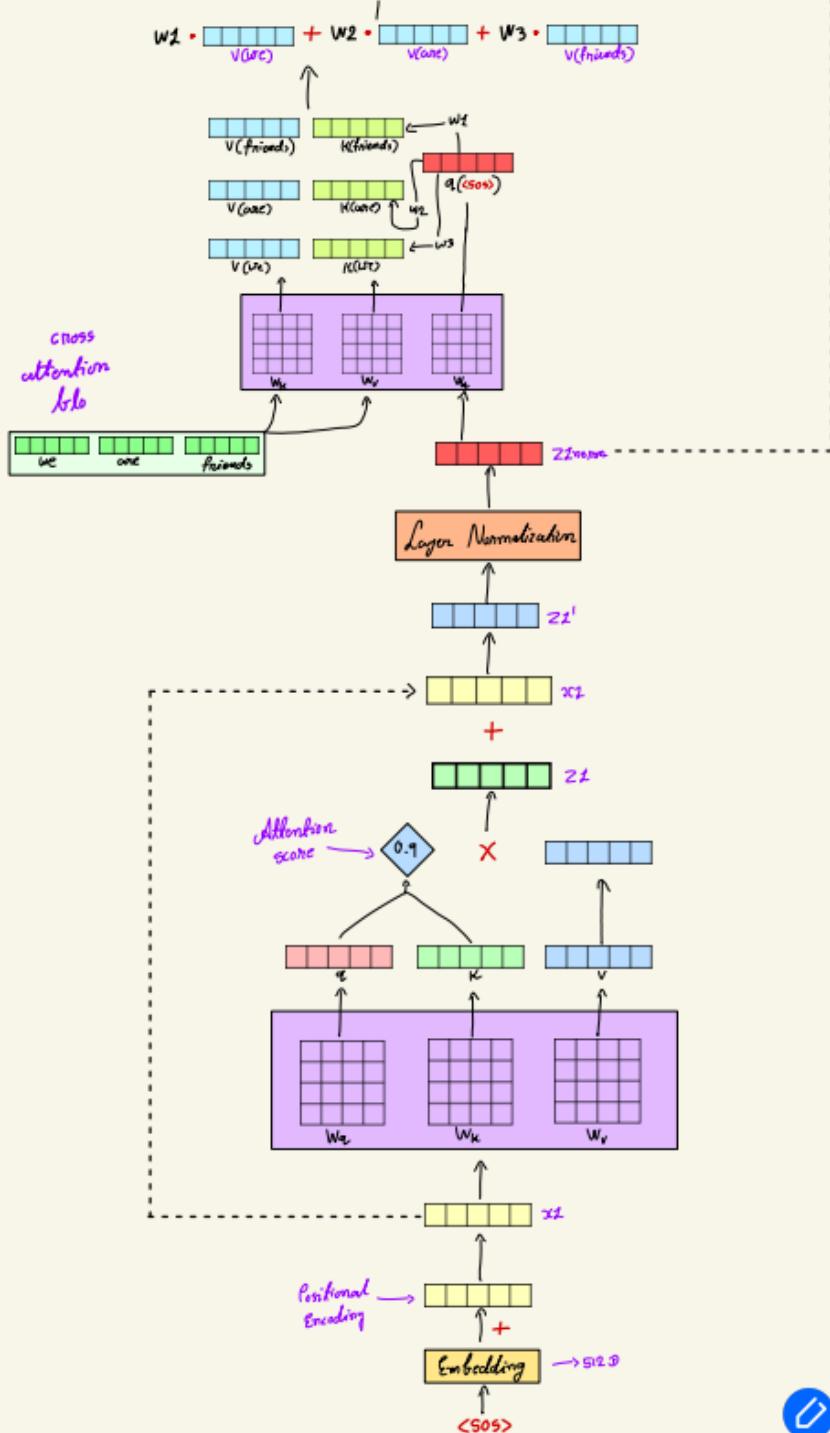
Inferential (Prediction) Phase

- During training phase, decoder works as non-autoregressor & auto-regressor during inferential phase.
- While encoder stays non-autoregressor in both phase.
- Finally, we will send **<SOS>** (Start of sentence) flag to start decoding process.



Phase-1





Phase-2

