

Converting

13 + End to Eng projects.

① Python programming - Basic to Advance

② NLP in Deep Learning →
RNN, LSTM, GRU, RNN, Bidirectional RNN
Encoder, Decoder, Attention Mechanism,
Transformer, BERT.
(Above is Basic building block)

↓
LM models.

③ Langchain Ecosystem → Generative AI framework

→ Pod LM, multimodal, Open Source LM

Agents, Chat with SQL, vector Database,
Retrieval, Text summarization
chat message chat with chat msg history.

Langchain core
Langchain community
Langchain Tools

④ Deployment Techniques

⑤ Fine Tuning Techniques
Quantization, LORA, QDURA

⑥ GenAI in AWS : AWS BedRock, Lambda function.

⑦ Nvidia Nim

Streamlit

Streamlit is an open-source app framework for machine learning and Data science projects. It allows you to create beautiful web applications for your machine learning and data science project with simple python script.

```
st.title('Hello Streamlit') # To write a title for a page.
st.write('') # To write simple text.
st.line_chart(data) # To create line chart.
```

```
name = st.text_input('Enter your name: ')
# for taking input.
```

To run streamlit file we use command:
→ streamlit run filename.py

```
age = st.slider('Select your age:', 0, 100, 25)
# creates a slider with default age 25.
```

```
# we can also create drop down select.
option = ['python', 'c++', 'Java', 'PHP']
choice = st.selectbox('Language', options)
```

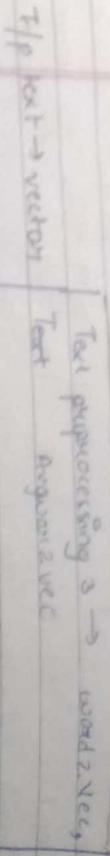
```
file = st.file_uploader('choose a pdf', type='pdf')
# This will create an option to upload your document to the website.
```

Go to streamlit.io for more knowledge and components.

* Machine learning for NLP

Good Luck | Page No. _____ Date _____

Content Page
Date _____



* Cleaning
T/p
Mu T/p

* Tokenization in NLP

Tokenization is a process of converting the paragraph or a single sentence into simple tokens.

- Corpus - paragraphs
 - Documents - sentences
 - Vocabulary - unique words
 - Words - all the words present
- Sentence \rightarrow paragraph
→ from nltk.tokenize import
sent_tokenize
- Paragraph \rightarrow word
sentence \rightarrow word
- from nltk.tokenize import
word_tokenize

Natural Language Toolkit (NLTK)
NLTK is a leading platform program to work with human language. It is open source driven.

→ wordpunct_tokenize
→ TreebankWordTokenizer

* Porter Stemmer
→ from nltk.stem import PorterStemmer

* RegexpStemmer
It is simple and more customizable stemming approach than PorterStemmer.
It uses regular expression to define how to remove suffixes back from words.

→ from nltk.stem import RegexpStemmer
a = RegexpStemmer('ing\\$|s\\$|le\\$|able\\$|eau\\$')

Above RegexpStemmer is configured to remove suffixes like 'ing', 's', 'e' and 'able'.

* Snowball Stemmer
It is a stemming algorithm which is also known as the Porter2 stemming as it is better version of Porter stemmer since some issues of it were fixed in this stemmer

* Lemmatization

Lemmatization technique is like stemming.
The output will get after lemmatization called as 'lemma', which is root word rather than stem word root stem.

In lemmatizer the meaning of the word.

→ from nltk.stem import WordNetLemmatizer
a = WordNetLemmatizer()
a.lemmatize('going', pos='v')

* Stemming
Stemming is the process of reducing a word stem that affixes to suffixes and prefixes or to the root of words known as a lemma.

POS : Part of speech
noun - n
verb - v
adverb - a

* Stopwords

Stopwords are common words in a language that are usually filtered out before processing text data in NLP. These words are typically removed because they occur frequently and contribute little to the meaning or context of the text. e.g.: 'a', 'the', 'is', etc.

→ `import nltk`

```
from nltk.corpus import stopwords
```

`nltk.download('stopwords')`

sentence = 'This is a example'

words = `nltk.word_tokenize(sentence)`

a = [word for word in words if word.lower() not in set(stopwords.words('english'))].
print(a)

Output: ['example']

* Part-of-Speech (POS) tagging

It is a process in NLP that involves labeling each word in a sentence with its corresponding part of speech such as noun, verb, adjective, etc. This is an essential step for understanding the grammatical structure of a sentence.

→ `import nltk`

sentence = 'The brown fox jumps over the

words = `nltk.word_tokenize(sentence)`

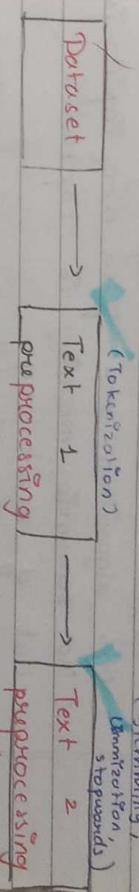
`a = nltk.pos_tag(words)`

→ Output: [('me', 'DT'), ('brown', 'JJ'),

* Name entity recognition (NER)

It is a crucial task in NLP that involves identifying and classifying name entities in text into predefined categories such as person names, organizations, locations, dates and more. e.g.: Akonkha : person name

2003 : date



* One hot encoding

One hot encoding is a fundamental technique in NLP used to represent words to vectors. Each word in a vocabulary is mapped to a unique word vector, where only one element is set to 1 and all other elements to 0.

Example: vocabulary of words: ['cat', 'dog', 'fish']

cat = [1, 0, 0]

dog = [0, 1, 0]

fish = [0, 0, 1]

Advantage :-

- Easy to implement with python.
- It does not make any assumptions about the relationship between words.
- Clearly distinguishes between words.

Disadvantages :-

- sparse matrix → overfitting
- NLP algorithm → fixed size input is required
- No semantic meaning is captured.
- Out of vocabulary
- Large high-dimensional sparse vectors leading to computational inefficiency.

Because of all these disadvantages we don't use one hot encoding in NLP.

★

Bag of words

Bag of words is ~~a~~ a fundamental technique in NLP used for text representation. It converts text into numerical features.

Example:

1. He is a good boy.
2. She is a good girl.
3. Boy and girl are good

Stopwords ↴
Lower all the words
case



Vocabulary frequency

s1 → good boy

good 3

s2 → good girl

boy 2

s3 → boy girl good

girl 2

A bag of words (BOW)

There are 2 types of BOW :-

• Binary BOW = 1 0 0

• BOW = count will get updated based on frequency

Advantages :-

- Simple and intuitive
- Fixed sized IP — machine learning algorithms

Disadvantages :-

- Sparse matrix - overfitting.
- Ordering of the word is getting changed.
- Out of vocabulary
- Semantic meaning is still not captured.

★

TF - IDF

Term Frequency . Inverse Document Frequency
It reflects the importance of a word in a document.

Term frequency = frequency of a word in a document / total no. of words in document.

Document → sentence.

| | | | | | |
|----|---|---|---|---|---|
| s1 | [| 1 | 1 | 0 |] |
| s2 | [| 1 | 0 | 1 |] |
| s3 | [| 1 | 1 | 1 |] |

$$IDF = \log \left(\frac{\text{Total no. of documents in corpus}}{\text{no. of documents containing word}} \right)$$

$$TF-IDF = TF * IDF$$

Example:-

$$S_1 = \text{good boy}$$

$$S_2 = \text{good girl}$$

$$S_3 = \text{boy girl good}$$

Term-frequency

$$\begin{matrix} S_1 & S_2 & S_3 \\ \frac{1}{2} & \frac{1}{2} & \frac{1}{2} \\ \text{good} & \text{boy} & \text{girl} \\ \text{boy} & 0 & \frac{1}{3} \\ \text{girl} & 0 & \frac{1}{3} \end{matrix}$$

Inverse-frequencey

$$IDF$$

$$\begin{matrix} \text{good} & \log_e(3/2) = 0 \\ \text{boy} & \log_e(3/2) \\ \text{girl} & \log_e(3/2) \end{matrix}$$

Word embedding

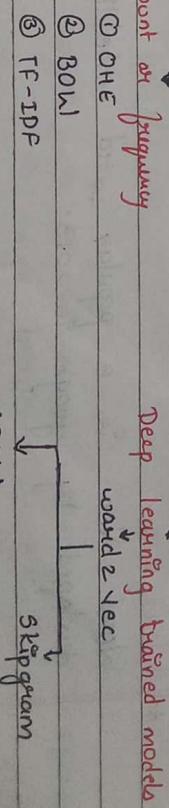
It is technique used in NLP to represent words as vectors in continuous vector space. This allows words with similar meanings to have similar representations. Words that are closer in the vector space are expected to be similar in meaning.

Example:

If we represent the words: angry, happy and excited, in the vector space:-

$$\therefore TF-IDF = S_1 \begin{bmatrix} 0 & \frac{1}{2} \times \log_e \frac{3}{2} \\ 0 & 0 & \frac{1}{2} \log_e \frac{3}{2} \\ 0 & 0 & \frac{1}{3} \log_e \frac{3}{2} & \frac{1}{3} \log_e \frac{3}{2} \end{bmatrix}$$

Word embedding



Advantage

- Intuitive: easily understandable
- Fixed size input
- Word importance is getting captured.

Disadvantage

- Sparse matrix.
- Out of vocabulary

Word2Vec

* N-grams

An n-gram is a sequence of 'n' items from a

given text or speech.

Unigram (1-gram)

Sequence of single word with single item

eg: I love NLP

unigram : ['I' , 'love' , 'NLP']

Bigram (2-gram)

A sequence of two items

eg: I love NLP

bigram : ['I love' , 'love NLP']

Trigram (3-gram)

A sequence of 3 items

eg: I love NLP

bigram : ['I love NLP']

N-gram is essential for capturing the context and structure of language.

sklearn \rightarrow n-gram = (1,1) \rightarrow unigram

= (1,2) \rightarrow unigram, bigram

= (1,3) \rightarrow unigram, bigram, trigram

= (2,3) \rightarrow bigram, trigram.

Word2Vec

Word2Vec is a popular word embedding technique developed by google in 2013.

It is designed to capture transform word into continuous vectors that capture semantic meanings and relationships between words. Word2Vec creates these vectors by using a shallow, two-layer neural network.

Once trained such models can detect synonymous words or suggest words for partial sentence.

Each and every word present in the vocabulary is converted into feature representation.

\downarrow
means converting all the words in vectors based on some features.

Vocabulary

| | boy | girl | king | queen | apple |
|--------|-----|------|-------|-------|-------|
| Gender | -1 | 1 | -0.92 | 0.93 | 0.01 |

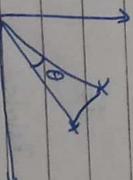
Royal

Age

Food

Features

$$\text{Distance} = 1 - \cosine \theta$$

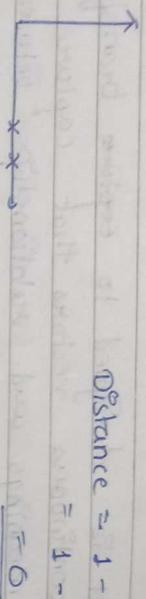


where, $x \rightarrow$ input word vector in space
 $\theta \rightarrow$ angle between word vector

| | |
|-----------|----------|
| Good Luck | Page No. |
| Date | |

| | |
|-----------|----------|
| Good Luck | Page No. |
| Date | 1 |

$$\text{Distance} = 1 - \cos \theta = 1 - 1 = 0$$



In this way recommendation is done.
or same words

when a word is converted to vector, will
we get a size of low dimension
which will our final output
Window size represents basically the feature
representation of a word.

* Skipgram
This model predicts the surrounding context words
for a given target word.

I/P
o/s
O/P

o/s, Related, Data, Science

Related
company, o/s, To, Data

To
o/s, Related, Data, Science

Small dataset \rightarrow CBOW
Huge dataset \rightarrow skipgram

- Advantages

- Sparse Matrix \rightarrow dense matrix
- Semantic information is captured.

- Vocabulary size
- Out of vocabulary solved.

I/P

o/s

iNeuron, Company, Related, To

iS

company, To, Data

Related

i.e., Related, Data, Science

To

and we give window size = 5

window size indicates how many words to be selected initially.

Suppose we have corpus :-

CBOV (Continuous Bag of Word)
CBOV is a fully connected neural network.
This model predicts the current word based on its surrounding words.
Suppose we have corpus :-

iNeuron Company Related To Data Science

* NLP In Deep Learning

This is prerequisite for Generative AI

Encoder - Decoder

Simple RNN \rightarrow LSTM / GRU RNN

Encoder - Decoder \rightarrow Bidirectional RNN

Self attention \rightarrow Transformer

\rightarrow To maintain the sequential information, pass

pass one word at a time instead of passing the entire sentence at once (like in ANN)

ANN : providing all the information then use

as forward propagation and then backward propagation

- \rightarrow We don't use ANN for textual data because:-
- \rightarrow It doesn't maintain the order (sequence) of the data, and because of that information loss takes place.
- \rightarrow It passes the entire input sequence at once to the network instead of word by word.

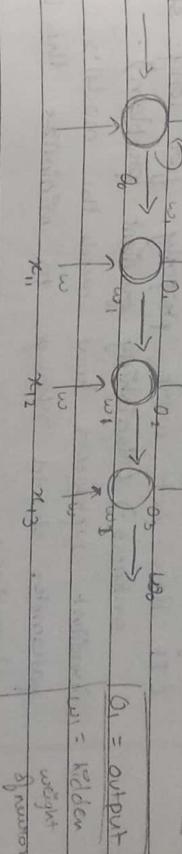
* Simple RNN

Type of neural network designed to process sequential data.

It introduces a feedback loop, allowing information to persist across time steps.

All the 3 weights are updated using this above formula.

* Forward propagation



$x_{11} \rightarrow \text{word 1}$

$x_{12} \rightarrow \text{sentence 1}$

$w_{11} = \text{hidden weight}$

$w_{12} = \text{hidden weight}$

$w_{13} = \text{hidden weight}$

$t=1 \quad t=2 \quad t=3$

Forward propagation

$$o_1 = f(x_{11} \cdot w + b_1)$$

$$o_2 = f(x_{12} \cdot w + o_1 \cdot w_1 + b_2)$$

$$o_3 = f(x_{13} \cdot w + o_2 \cdot w_1 + b_3) \xrightarrow{\text{activation function}} \text{forward output}$$

* Backward propagation

After this calculation of final output there will be some loss

$$\text{loss} = \frac{1}{2} (\hat{y} - y)^2$$

So main goal is to reduce this loss so we can do this by updating the weights

$$[w, w_1, w_2]$$

So during backward propagation these weights will be updated and then again forward propagation is performed.

Backward propagation with time

Update $[w, w_1, w_2]$

weight update formula,

$$w_{new} = w_{old} - \eta \frac{\partial \text{loss}}{\partial w_{old}}$$

η slope of gradient descent

* Problems with RNN

- Long term dependency cannot be captured by RNN
(It does not work properly for long sentences) \rightarrow (Vanishing gradient problem).

Gradient represents how much the model's parameters should change to minimize the loss function.

In backpropagation, the gradients are calculated layer by layer, starting from the output layer and moving backward to input.

During this, the weight will \downarrow so input becomes equivalent to ≈ 0 . thus not having much effect on the output.

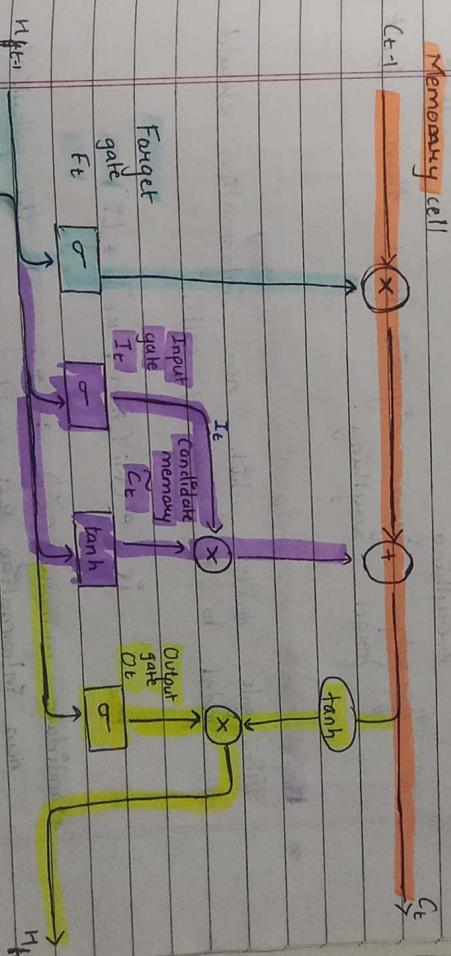
* LSTM RNN

Long - short term memory RNN solves the problems regarding simple RNN.
It is designed to remember information for long periods. (It's like a smart brain cell that can decide what to remember or forget).

LSTM achieves long term memory by using a memory cell and gates to control information flow.

* Basic Architecture

(Basic Components):



This is the way to represent Neural Network larger in our architecture.

meaning :-

a. Neural network

*

Memory cell

and a sigmoid function
 $\begin{bmatrix} 1 \\ 0 \end{bmatrix}$ is applied to it.

Long - short term memory
Long - whatever info. we want will be stored and whatever is not required is removed.

- O : Pointwise operation
Example : $y_1 = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$
 $y_2 = \begin{bmatrix} 4 \\ 5 \\ 6 \end{bmatrix}$

- \rightarrow Vector transfer
- \rightarrow Concatenate Basically combining 2 vectors
- Copy

- h_{t-1} : hidden state of previous timestamp.
- x_t : word passed as input in current timestamp.

The input gate decides how much of the candidate memory should be stored in cell state.

* **Forget gate**: (f_t) =
Here in input (x_t) is taken and concatenated with previous hidden layer (h_{t-1}), then passed through a sigmoid activation function (σ) and then lastly cross multiplied (\otimes) with previous long-term memory (c_{t-1}).

* It is used to forget information from memory cell based on the context.

- sigmoid value ranges between 0 - 1
- if output from sigmoid is 0 then forget everything.
- if output from sigmoid is 1 then don't forget anything.

* **Input Gate and Candidate Memory**
This gate decide which new information should be added to the memory cell.

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

- Candidate memory (c_t^{\sim}) : represents the potential new information that could be added to the c_t .

$$c_t^{\sim} = \tanh(W_c \cdot [h_{t-1}, x_t] + b_c)$$

(Candidate memory : new notes to be added).

Input gate : how much notes should be added

$$\begin{aligned} h_t &= \tanh(W_h \cdot [c_t^{\sim}, x_t] + b_h) \\ c_t &= i_t \cdot c_t^{\sim} + f_t \cdot c_{t-1} \end{aligned}$$

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

long term

$$h_t = O_t + \tanh(C_t)$$

forget gate

Input gate

* **Output gate**: This gate decides what the next hidden state should be;
Calculation of h_t .

$$O_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

∴ Hidden state

$$h_t = O_t + \tanh(C_t)$$

↑
short-term

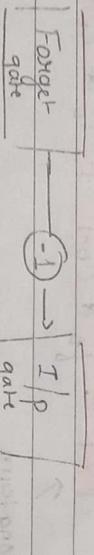
- Variants of LSTM
- 1. Peephole connections used in LSTM

Peephole connection \leftrightarrow
Connection from memory cell to Forget gate
Input gate and output gate.

$$\begin{aligned} f_t &= \sigma(W_f \cdot [c_{t-1}, h_{t-1}, x_t] + b_f) \\ i_t &= \sigma(W_i \cdot [c_{t-1}, h_{t-1}, x_t] + b_i) \\ O_t &= \sigma(W_o \cdot [c_t, h_{t-1}, x_t] + b_o) \end{aligned}$$

Coupling forget and I/P gate.

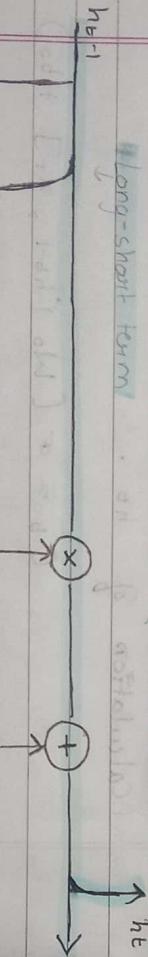
Instead of separately deciding what to forget and what we should add new input, we make this decision together.



- We only forget when we're going to input something in its place and vice versa.

$$c_t = f_t * c_{t-1} + (1 - f_t) * \tilde{c}_t$$

3. GRU (Gated Recurrent Unit) RNN (2014)



(Temporary) Candidate hidden state [current context]

Controls how much of the previous hidden state to keep and how much of new state to update with.

* Bidirectional RNN

When we need context of past words as well as further words there we use bi-RNN

example: I eat — in Bangla
so here I need the context of both previous as well as further words to predict the word.

n_t

Update gate

Reset gate

Temporary hidden state.

$$z_t = \sigma(W_z \cdot [h_{t-1}, n_t])$$

$$n_t = \sigma(W_n \cdot [h_{t-1}, n_t])$$

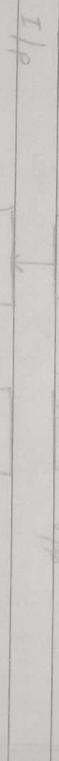
$$\tilde{h}_t = \tanh(W \cdot [n_t * h_{t-1}, n_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

* Types of RNN:

- One to One
- One to many
- Many to one
- Many to many

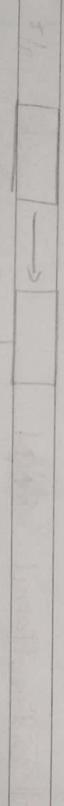
Example: Image captioning : (Google image)



* example : Image captioning

google image : we provide a image and tell us what that image is.

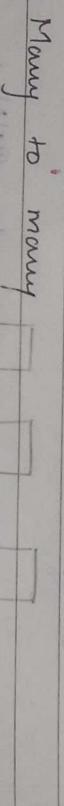
* Many to one



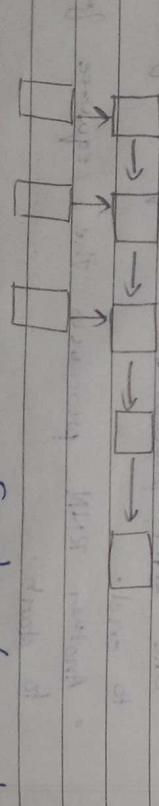
* Many to many



Example : Image search .

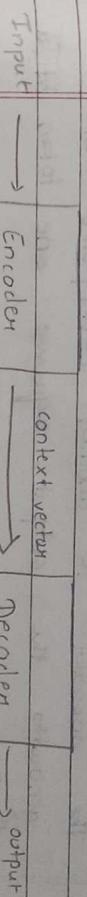


Example : Language translation.



* Encoder and Decoder

- The encoder processes the input sequence one element at a time .
- It converts the sequence into a fixed length context length
- It is used to solve seqⁿ to seqⁿ problems.

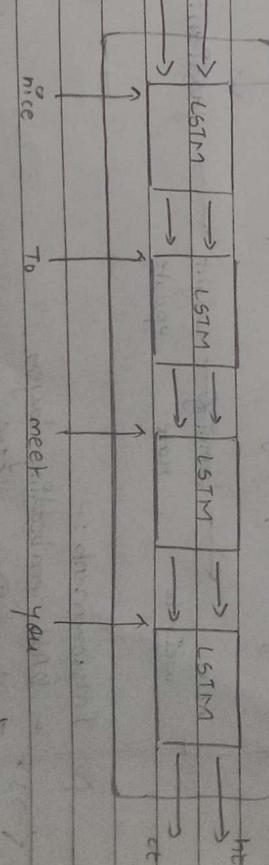


* High Level Overview

The encoder - decoder is a fundamental neural network commonly used in seqⁿ to seqⁿ tasks such as machine translation, text summarization, etc .

Encoder - Decoder architecture consists of two main components:-

- Encoder
- Processes the input sequence and generates a context vector.
- Typically implemented using LSTM or GRU cells.
- Example : English \rightarrow Hindi translation.



Time steps .

$\downarrow \rightarrow$ context return .

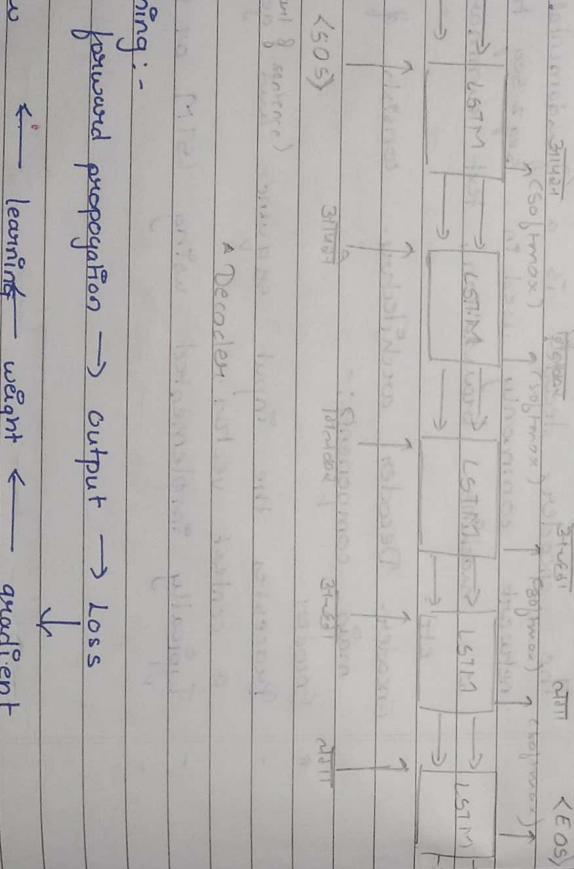
o Context Vector

- The last hidden state from the encoder acts as summary of the entire sequence.
- It is passed to the decoder in its initial state.

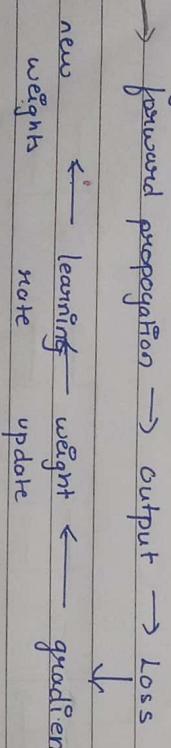
★ Decoder

The decoder is another LSTM / GRU that generates the output sequence one token at a time.

- Decoder is autoregressive, meaning it generates one token at a time and use prior previously generated tokens to predict the next token.



- o Training :-
- forward propagation → output → loss
- new weights ← learning weight ← gradient note update



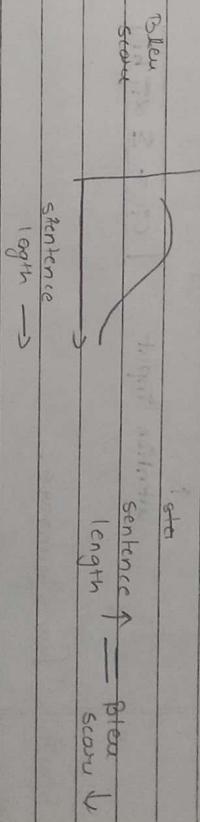
o Improvements :-

- Add embedding layer
- Use Deep LSTMs.

Limitations

- Fixed-length context vector. (Bottleneck Issue)
- If the input sequence is long, this context vector cannot retain all important details, leading to information loss.
- Requires large details for good performance.

Bleu score : performance metric.

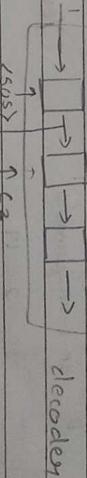


Attention Mechanism

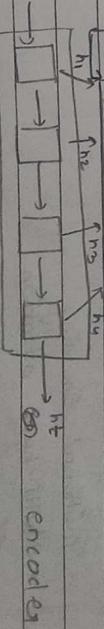
It helps model focus on the most relevant part of the input when generating output.

o It solves the bottleneck issue in encoder-decoder by allowing the decoder to access all input token dynamically, instead of relying on a single fixed-length context vector.

Attention helps decoder 'attend' to specific words



$$\alpha_{21} = f(h_1, s_2)$$



$s_1 \rightarrow$ previous hidden state of decoder.

* Basic Transformer Architecture

$$C_2 = \alpha_{21} h_1 + \alpha_{22} h_2 + \alpha_{23} h_3 + \alpha_{24} h_4$$

In normal encoder-decoder, we usually only take s_i and context vector to calculate

the next word in decoder.

But in attention mechanism new variable is added C_i with s_i and context vector. This C_i is nothing but attention 'attend'.

$$\text{Attention input } C_i = \sum \alpha_i^w h_i$$

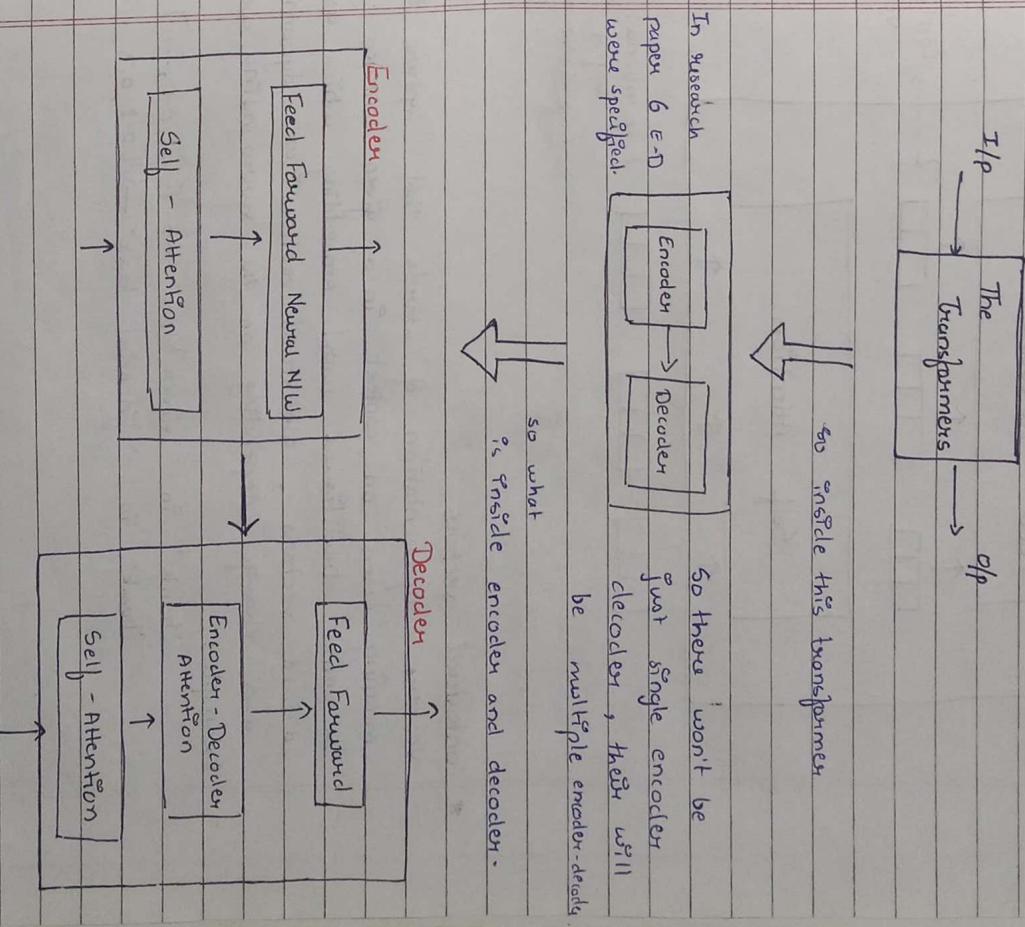
* Transformers.

Attention mechanism (transformation) - Parallelly we cannot send all words in a sentence (They were sent one by one).
- Not scalable : unable to handle large datasets

Transformers \rightarrow LSTM RNN

\downarrow Self attention module
All the words will be parallelly sent to the encoder.

Paper (first research paper) : Attention All You Need

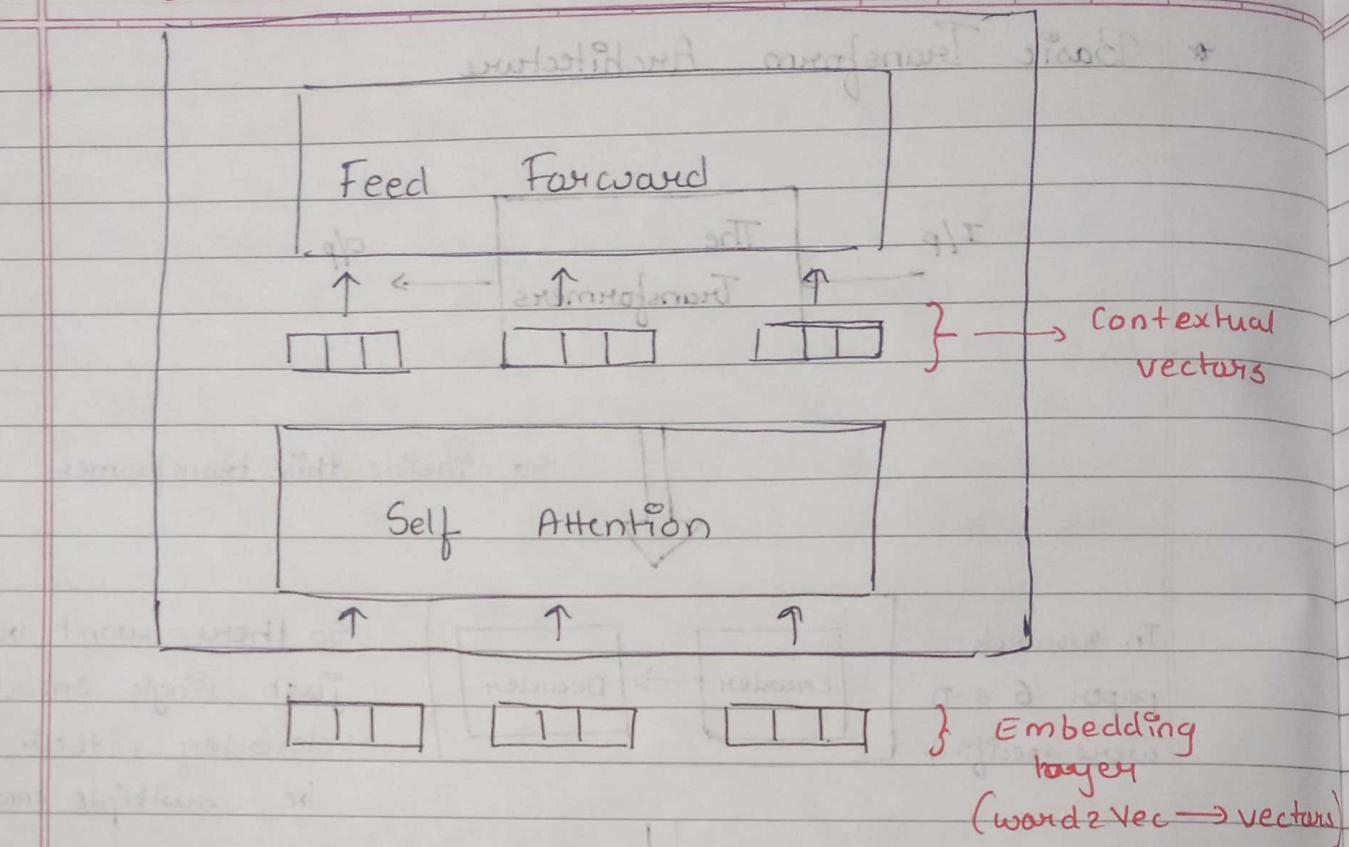


Features :-

- Self-attention mechanism
- Positional Encoding
- Multi-head attention
- Parallel processing.
- Feed forward Neural Network.

Let's look into encoders.

Encoder



* Contextual vectors

- Vectors representation of words that capture their meaning based on context in a given sentence.
- Unlike traditional word embedding which assigns a fixed vector, contextual vectors are dynamic and change depending on the surrounding words.
- Example:

'bank' in 'river bank' $\rightarrow [0.2, 0.5, \dots]$

'bank' in 'financial bank' $\rightarrow [0.8, 0.3, \dots]$

★ 6) Attention at higher and detailed Level
 Self-attention, also known as scaled dot-product attention, is a crucial mechanism in the transformer architecture that allows the model to weigh the importance of different tokens in the input sequence relative to each other.

- Creates contextual vectors.

Scaled Dot-Product-

Steps :-

1) Token Embedding.

Let's take example

Input sequence =

['The', 'cat', 'sat']

Embedding size = 4

$Q, K, V = 4$

$$E_{\text{The}} = [1, 0, 1, 0]$$

$$E_{\text{cat}} = [0, 1, 0, 1]$$

$$E_{\text{SAT}} = [1, 1, 1, 1]$$

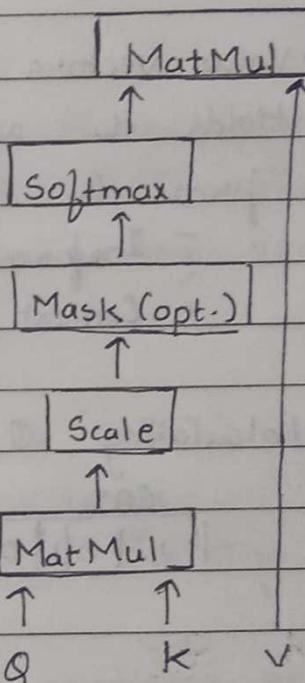
2) Linear Transformation

We create Q, K, V by multiplying the embedding by learned weights matrices W_Q, W_K and W_V .

- Query (Q): represents the token for which we are calculating the attention. They determine the context of importance of ^{other} ~~current~~ tokens in the context of the current token.

- Focus Determination

- Contextual Understanding



• Key Vector ($\mathbf{K} \in \mathbb{R}^k$)

Represents all the tokens in the sequence and are used to compare with query vector to calculate attention score.

- Relevance measurement: compatibility of each token with the current token.
- Information retrieval

• Value vectors (\mathbf{V})

Holds the actual information that will aggregate to form the output of the attention mechanism.

- Information aggregation
- Context preservation.

Calculating $\mathbf{Q}, \mathbf{K}, \mathbf{V}$:-

$$\begin{bmatrix} 1 & 0 & 1 & 0 \end{bmatrix}_{\text{CAT}} \xrightarrow{\text{dot-product}} \begin{bmatrix} \cdot & \cdot & \cdot & \cdot \end{bmatrix}_{\mathbf{W}_Q} = \mathbf{Q}$$

$$\begin{bmatrix} \cdot & \cdot & \cdot & \cdot \end{bmatrix}_{\mathbf{W}_K} = \mathbf{K}$$

$$\begin{bmatrix} \cdot & \cdot & \cdot & \cdot \end{bmatrix}_{\mathbf{W}_V} = \mathbf{V}$$

so far each token's vectors are calculated.

3) Compute Attention Scores

By calculating the dot product between a query vector and all the key vector, the model assesses how much attention to give to each token relative to the current token.

Example:-

The

$$\text{Score } (\mathbf{Q}_{\text{cur}}, \mathbf{K}_{\text{tar}}) = [1 \ 0 \ 1 0] \cdot [1 \ 0 \ 1 0]^T = 2$$

$$\text{Score } (\mathbf{Q}_{\text{cur}}, \mathbf{K}_{\text{tar}}) = [1 \ 0 \ 1 0] \cdot [1 \ 1 \ 1]^T = 0$$

Similarly we calculate for the rest of the tokens

4) Scaling

We take up the score and scale them down by dividing the score by \sqrt{k} (dimension of key vector), so the here dimension of key vector is $n : T_k = 2$.

Scaling in the attention mechanism is crucial to prevent the dot product from growing too large.
 \Rightarrow Ensure stable gradients during training.

If we don't do scaling following problems occur :-

1. Gradient Exploding
2. Softmax saturation.

\hookrightarrow most of the weight is assigned to a single token and other tokens weight will be assigned near to zero.

• Property of softmax

- $\text{softmax}([6, u]) = \approx (0.88, 0.12)$
- most of the attention weight is assigned to the first key vector and very little to second vector.

$$-\text{softmax}([10, 1]) = \approx (0.99, 0.01)$$

Now with scaling:

$$[6, 4] \Rightarrow \text{scale} \Rightarrow [6/2, 4/2] = [3, 2]$$

$$\text{softmax}([3, 2]) \approx [0.73, 0.27]$$

here, the attention weights are more balanced compared to unscaled case.

- It stabilize Training.
- Preventing saturation.

Scaled-score $(Q_{\text{THE}}, K_{\text{THE}}) = 2/2 = 1$

$$(Q_{\text{THE}}, K_{\text{CAT}}) = 0/2 = 0$$

$$(Q_{\text{THE}}, K_{\text{SAT}}) = 2/2 = 1$$

My scaling will be done to all other tokens.

5) Apply softmax

$$\text{Attention weights} = \text{softmax}([1, 0, 1])$$

$$= [0.4223, 0.1554, 0.4223]$$

$$\text{Attention weights 'cat'} = \text{softmax}([0, 2, 2])$$

$$= [0.1554, 0.422, 0.422]$$

$$\text{Attention weights 'sat'} = \text{softmax}([C_2, 2, u])$$

$$= [0.219, 0.219, 0.472]$$

$$\text{softmax}(n, u) = \begin{bmatrix} e^n & e^u \\ e^{n+u} & e^{n+u} \end{bmatrix}$$

6) Weight sum of values

We multiply attention weight by corresponding value vector.

$$\text{Output 'THE'} = 0.4223 * \sqrt{\text{THE}} + 0.1554 * \sqrt{\text{CAT}} +$$

$$= [1.2669^2, 0.9999, 1.2669, 0.9999]$$

\uparrow
attention head (1)

My for other tokens are calculated.

$$\text{The} \rightarrow [1 | 0 | 1 | 0] \Rightarrow \text{Attention Self} \Rightarrow [1.2669, 0.9999, 1.2669, 0.9999]$$

Ships in short.

↳ Calculated Q, K, V

- Attention score
- Scaled
- Softmax
- weighted sum of values

* Multi-Head Attention

While calculating Q, K and V we change our weight values to different than previous ones we will get the final output (z) as different.

This nothing but multi-head attention.

- It expands model's ability to focus on different position of tokens.

So from our self attention we generate multiple attention heads $\{z_1, z_2, z_3, \dots\}$

So now how do we combine all these things before giving it to the feed forward neural network.

So for that :-

1. Concatenate all the attention heads
2. Multiply with the weight matrix W^b that is trained
3. The result would be the Z matrix that captures information from all the attention heads.

Now we can send this to Feed-Forward Neural Network.

* Positional Encoding

So up till now we are processing the words parallelly \Rightarrow Lack of sequential structure of

The words $\{\text{order}\}$ So how do we know our

encoder knows the order of the words, so here comes the role of positional encoding.

- Representing order of sequence.

We create a positional encoded vector and add it to embedded vectors, so that will give the position information.

Type of positional encoding :-

1) Sinesoidal Positional encoding

It uses \sin and \cos functions of different frequencies to make positional encoding.

Formula :-

$$\text{P.E}(\text{position}, z_i) = \sin \left(\frac{\text{pos} \cdot i}{10000^{\frac{i}{d_{\text{model}}}}} \right)$$

where,

$$i = \text{dimension},$$

d_{model} = dimensionality of the embedding

Example :-

THE CAT SAT

$$i = 0 \quad 1 \quad 2 \quad 3$$

$$\circ \text{ THE} \rightarrow [0.1 \quad 0.2 \quad 0.3 \quad 0.4]$$

$$1 \text{ CAT} \rightarrow [0.5 \quad 0.6 \quad 0.7 \quad 0.8]$$

$$2 \text{ SAT} \rightarrow [0.9 \quad 1.0 \quad 1.1 \quad 1.2]$$

$$\text{so } d_{\text{model}} = 4$$

$$\text{pos of cat} = 1 \quad i = \text{index of over 0 to 3}$$

$$\text{pos of THE} = 0$$

$$\text{pos of SAT} = 2$$

so even numbers position will use \sin formula odd i will use \cos formula

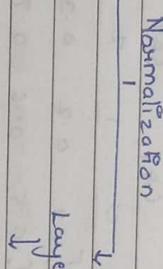
We do not concatenate the positional vector with embedding vector because dimensional will increase that's why we add up.

★ Layer Normalization

Before passing our contextual vector to feedforward or is (Add Norm) Normalized using layer normalization.

Normalization is a technique used to stabilize and standardize the input data and activations at each layer of transformer model.

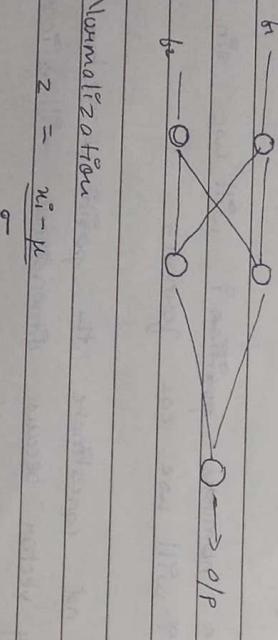
It typically involves normalizing the values so that they fall within a standard range, often between 0 and 1 or -1 and 1.



normalizes across the batch dimension

normalizes across the feature

There are various issues with batch normalization, that's why we use layer normalization.



$$\hat{x}_i = \frac{x_i - \mu}{\sigma}$$

1. Compute mean & variance across the features.
2. Normalize the input

where ϵ = small constant to prevent division by zero.

3. Scale and shift
- Apply learnable parameters γ (scale) and β (shift) to maintain the model's ability to represent complex functions:
- $$y_i = \gamma \hat{x}_i + \beta$$

$$\begin{matrix} b_1 & b_2 & z_1 & z_2 \\ - & - & - & - \\ \xrightarrow{\mu_1, \sigma_1} & \xrightarrow{\mu_2, \sigma_2} & \xrightarrow{\mu_3, \sigma_3} & \end{matrix}$$

Layer

$$z_{score} = \frac{x_i - \mu_i}{\sigma_i}$$

$$z_{norm} = \frac{x_i - \mu_i}{\sigma_i}$$

$$z_1 = \sigma [w_i^T x + b_i]$$

$$y = \gamma \left[\frac{z_1 - \mu_1}{\sigma_1} \right] + \beta$$

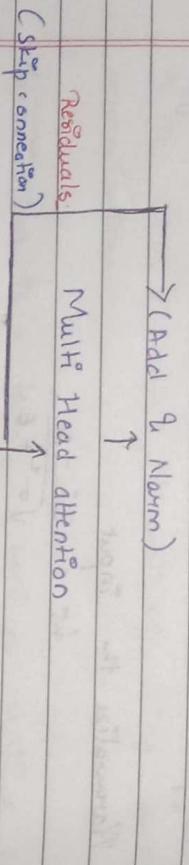
$$z = \frac{x_i - \mu}{\sigma}$$

$$\mu = 0$$

$$\sigma = 1$$

Feed Forward

Decoder



Text embeddings + Positional Encoding

I/P sequence

The transformer decoder is responsible for generating the output sequence one token at a time, using the encoder's output and the previously generated tokens.

Masked Multi-head attention mechanism

1. I/P embedding and position embedding

2. Linear projection for Q,K,V

3. Scaled Dot product attention

4. Mask application

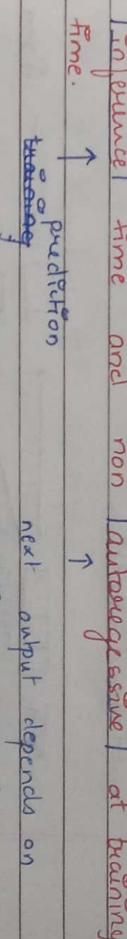
5. Multi-head attention mechanism

6. Concatenation and final linear projection

7. Residual connections and Layer normalization.

Masked Multi-head attention

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



It processes the output from the self-attention mechanism and helps in features transformation and learning complex representations.

- Adding non-linearity.

- Adds depth to the model. → more learning

ENN consists of :-

1. Few fully connected layers.
2. Non-linearity
3. Dropout (Regularization)

Usually ENN is applied after first layer → non-linearity
3. Dropout (Regularization)
Burns overfitting.

wrong

dataset.

In parallel processing we calculate contextual vectors. So there's a problem.

For Example : English \rightarrow Hindi
translation

late How are you \rightarrow तुम कैसे हो

$$(\text{तुम} \rightarrow 0.9 * \sqrt{\text{तुम}} + 0.06 * \sqrt{\text{कैसे}} + 0.1 * \sqrt{\text{हो}})$$

But at this point, this our first word prediction we don't know $\sqrt{\text{तुम}}$ and $\sqrt{\text{हो}}$ is going to come we still $\sqrt{\text{कैसे}}$ have given the value \rightarrow this is a problem of data leakage.

In this type it will work good on training data but not on testing data.

So to solve the problem of data leakage and process data parallelly we use ~~matrix~~ masked multi-head attention.

$$\text{तुम} \begin{bmatrix} 1 & 1 \end{bmatrix} \xrightarrow{\text{embedding vector}} \begin{bmatrix} w_a \\ w_k \\ w_v \end{bmatrix} = \begin{bmatrix} w_{11} \\ w_{21} \\ w_{31} \end{bmatrix}$$

$$\text{तुम} = w_{11} * \sqrt{w_{11}} + w_{12} * \sqrt{w_{12}} + w_{13} * \sqrt{w_{13}}$$

so whole calculation we do not need w_{12} & w_{13} contribution and w_{11} ka contribution

likewise for $\sqrt{\text{हो}}$ we do not need contribution of w_{11} .

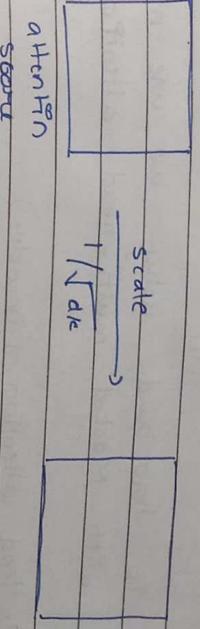
so far that we don't need $w_{12} = w_{13} = w_{14} = 0$
if the weight are zero then the term will be zero and thus no contribution.

so to make it happen we use masked model.

Q matrix
value
metrix
metrix
metrix

| | | |
|-------|-------|-------|
| Q तुम | V तुम | K तुम |
| Q तुम | V तुम | K तुम |
| Q तुम | V तुम | K तुम |

My job $\frac{w_{11}}{w_{11}}$ and $\frac{w_{11}}{w_{11}}$.



Then we calculate attention score.

$$Q * k = \text{attention score}$$

$$\begin{array}{c} \text{Attention} \\ \text{matrix} \\ \text{value} \\ \text{metrix} \\ \text{metrix} \\ \text{metrix} \\ \text{metrix} \end{array} \xrightarrow{\text{scale}} \begin{array}{c} \text{Attention} \\ \text{matrix} \\ \text{value} \\ \text{metrix} \\ \text{metrix} \\ \text{metrix} \\ \text{metrix} \end{array} + \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \xrightarrow{\text{norm}} \begin{bmatrix} -\infty & -\infty \\ -\infty & -\infty \\ -\infty & -\infty \end{bmatrix}$$

some values

softmax($-\infty$) = 0

| | | |
|-----------------|-----------------|-----------------|
| w ₁₁ | 0 | 0 |
| w ₂₁ | w ₂₂ | 0 |
| w ₃₁ | w ₃₂ | w ₃₃ |

This is how and why we use mask.
Done with Masked multi-head attention.

* Multi-Head attention. (Decoder)
also known as cross-attention.

- The information received from encoder (input sequence).
- The decoder generates output - one word at a time, depending on:
 - What was been generated so far.
 - The information received from encoder (input sequence).

Add & Norm

Feed Forward

Add & Norm

Multi-head
Attention

Self Attention

Cross Attention

Input

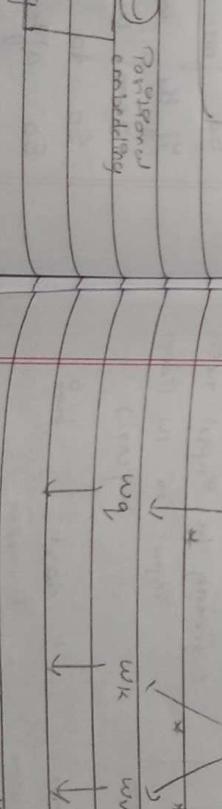
★ Processing
Self attention use a single sequence (input)
Cross attention use two sequences (input and output)

(cross-attention) is a mechanism used in transformer architectures, particularly in tasks involving sequence-to-sequence data like translation or summarization.

It allows a model to

focus on different parts of an input sequence when generating an output sequence.

Decoder



Output

Why for OKV for doing and ?

• Generates query vector from the output sequence (\hat{z}_{text}) and key/value vectors from the input sequence.

Performing dot products to derive attention scores, which indicates the relationship between

input sequence and output sequence

| | | |
|----|-----|---------|
| we | are | friends |
| 0 | 0 | 0 |
| 0 | 0 | 0 |

★ Decoder during training

works with highest probability is predicted among all words.

softmax

converts it into probability distribution (addition of all = 1)

[0.3 0.6 0.1, ...]

(10000 neurons) \rightarrow 0 0 0 0 0 [linear layer]

10000 unique kind words.

512 neurons \rightarrow 0 0 0 0 0

(tanh activation function)

Converting the input 4×512

input to the linear layer

(layer)

<start>

num

dot

had

from

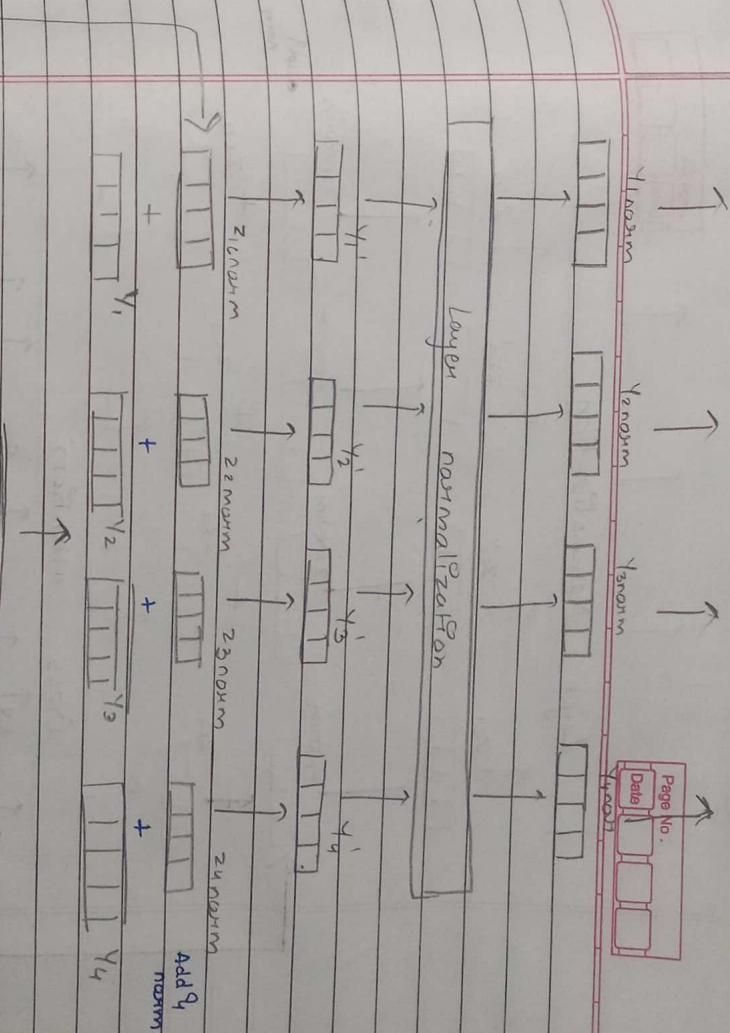
from

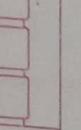
from

to

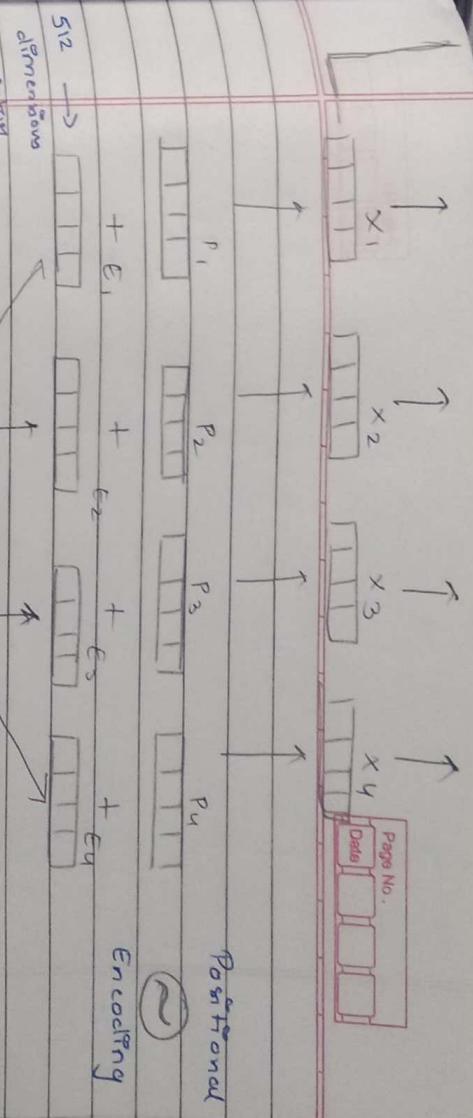
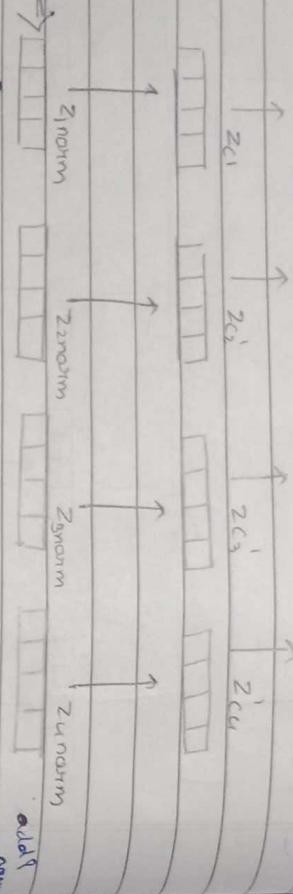
to

5 more decoder layers

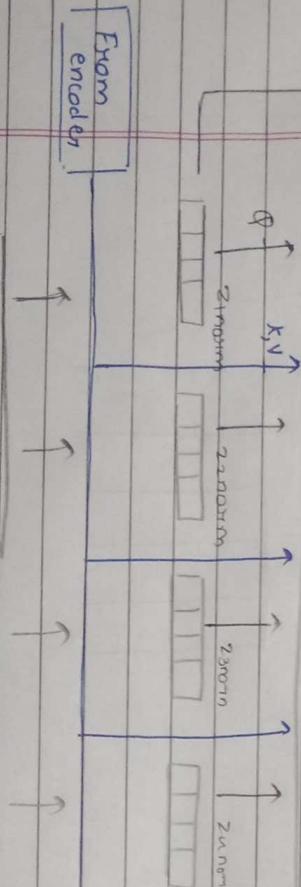




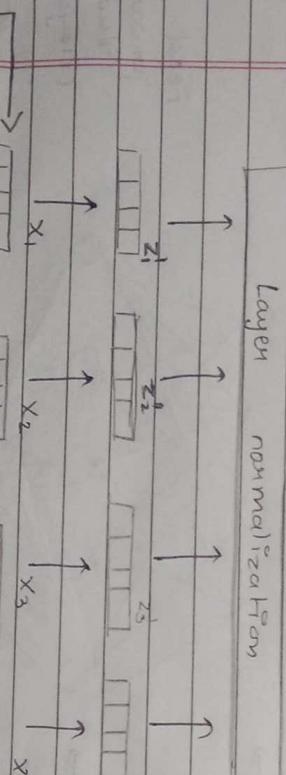
Layer Normalization



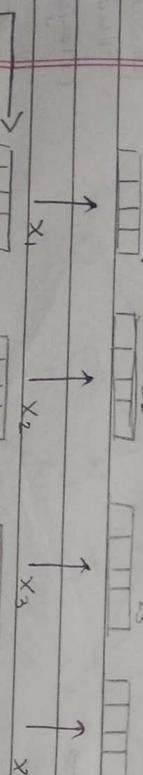
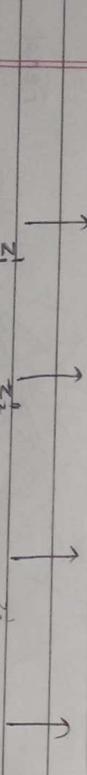
Cross attention



From encoder



Layer normalization



▲ Broad overview of decoder.

num dost ha

Right shift

To tokenization

L start num dost ha

→

