

8) In 'prediction.ipynb':

- Load the ANN trained model, scalar, pickle, onehot
  - Take a sample ip → convert it into suitable format i.e. OHE few columns & scale it.
  - Predict on the ip.
- a) Write 'app.py' → deploy on streamlit
- b) Practice by predicting a value now i.e. practice regression prob. stmt. on this same dataset. (Similar to prev.)
- c) Determining the optimal number of hidden layers & neurons for an ANN.

(Guidelines:

- Start Simple: Begin with a simple architecture & gradually increase complexity if needed.
- Grid Search/Random Search: Use grid search / random search to try different architectures.
- Cross-Validation: Use cross-validation to evaluate the performance of different architectures.

→ Heuristics & Rules of thumb: → A common practice is to start with 1-2 days.

The no. of neurons in the hidden layer should be between the size of the ip layer & the size of op layer

→ Refer 'hyperparameter tuning.ipynb'.

// Refer files for above project.

### NLP IN DEEP LEARNING

Simple RNN → LSTM/GRU RNN → Bidirectional RNN → Encoder Decoder

↓  
Transformers ← Self Attention

### SIMPLE RNN INDEPTH INTUITION

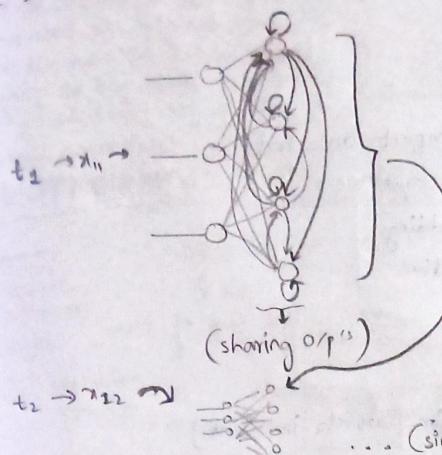
- Limitations of ANN:  
• ANN processes inputs simultaneously, ignoring word order  
• Inputting all words at once hinders context preservation.

- Intro. to RNN (Recurrent Neural Networks):

- Words are passed 1 at a time. (e.g. "The" at t=1, "food" at t=2)
- Feedback loops: Hidden neurons share outputs with themselves & other neurons maintaining context.

- Each timestamp retains info. from previous words, preserving the sequential nature of text.

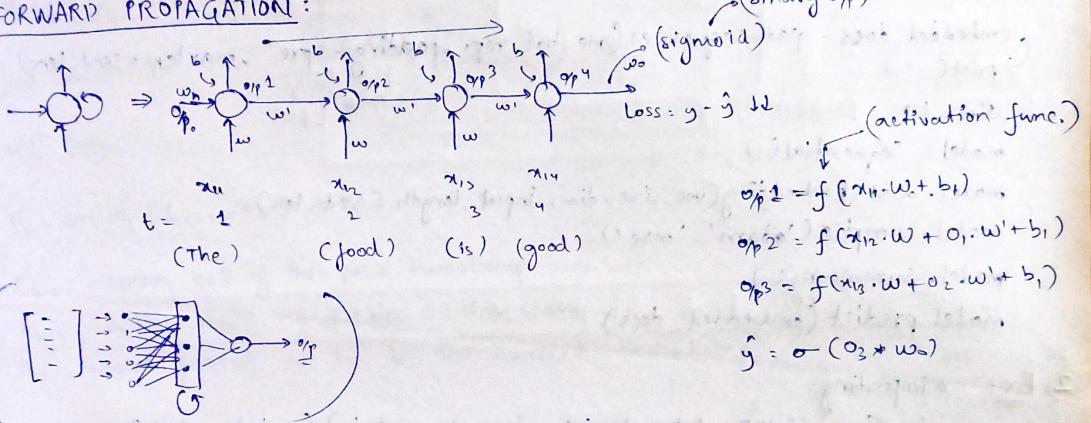
## RNN Architecture:



- eg: The food is good
- The food is bad
- The food is not good

- Isp prep.: → Preprocess the text by removing stop words (e.g. "the food is good" → "food" "good") & tokenize it into words.
  - Create a vocab. of unique words & now convert it into no. is (let  $\rightarrow$  OHF)  $\rightarrow$  food = [1, 0, 0, 0, 0] and = [0, 1, 0, 0, 0]

## • FORWARD PROPAGATION:



## BACKWARD PROPAGATION:

Update  $[w_i, w_n, w_0]$

; Weight updation formula  $\Rightarrow w_{\text{new}} = w_{\text{old}} - \eta \left( \frac{\partial h}{\partial w_{\text{old}}} \right)$

1) Update  $w_0$

2) Update  $w_m$  (Hidden Layer Weights)

3) Update weight  $w_i$

(timestamps)



#### - Problems with RNN:

• Long term dependency cannot be captured by RNN of Vanishing Gradient Prob.)

ex: Sentence length = 50 words.

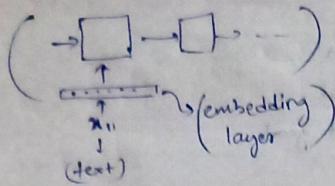
$$\therefore \frac{\partial h}{\partial w_{old}} = \left[ \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial z_0} \cdot \frac{\partial z_0}{\partial w_{old}} - \frac{\partial w_{old}}{\partial w_{old}} \right] + [ \dots ] + t_{50} \left[ \frac{\partial h}{\partial y} \cdot \frac{\partial y}{\partial z_0} \cdot \frac{\partial z_0}{\partial w_{old}} \right]$$

(t=1) ~~~~~ t=2

( $\therefore t=1 \rightarrow$  word not participating in updating the weights)

• Embedding layers = responsible for converting words into vectors.  
↳ uses "Word Embeddings"

# E2E DEEP LEARNING PROJECT WITH SIMPLE RNN



## ① Embedding example

We need all sentences to be of same length to train our RNN model because all the sentences that will be going it will go for fixed no. of timesteps based on sentence length.

① `from tensorflow.keras.preprocessing.text import one_hot  
 .utils import pad_sequences  
 .layers . Embedding  
 .models . Sequential`

`sent = ['...', '...', ...] → sentences`  
`voc_size = 1000 → vocab. size`  
`one_hot_repr = [one_hot(words, voc_size) for words in sent]` → ex. [6186, 7872, 532, 7113].  
`sent_len = 8`

`embedded_docs = pad_sequences(one_hot_repr, padding="pre", maxlen=sent_len)`  
`print( )` ↓  
dim = 10 → Feature representation  
model = Sequential()  
model.add(Embedding(voc\_size, dim, input\_length=sent\_len))  
model.compile('adam', 'mse')  
model.summary()  
model.predict(embedded\_docs)

② Ex:

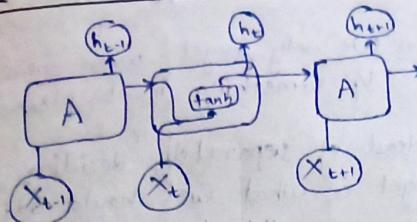
- importing
- loading IMDB dataset into (x-train, y-train), (x-test, y-test)
- mapping of words index back to words (for understanding)
- making sentences to equal length
- training Simple RNN

total 3 layers  
model = Sequential()  
total features → 1000 (vocab. size)  
total features → 500 (total features)  
model.add(Embedding(max\_features, 128, input\_length=max\_len))  
model.add(SimpleRNN(128, activation='relu'))  
model.add(Dense(1, activation='sigmoid'))  
model.compile(optimizer='adam', loss='binary\_crossentropy',  
metrics=['accuracy'])

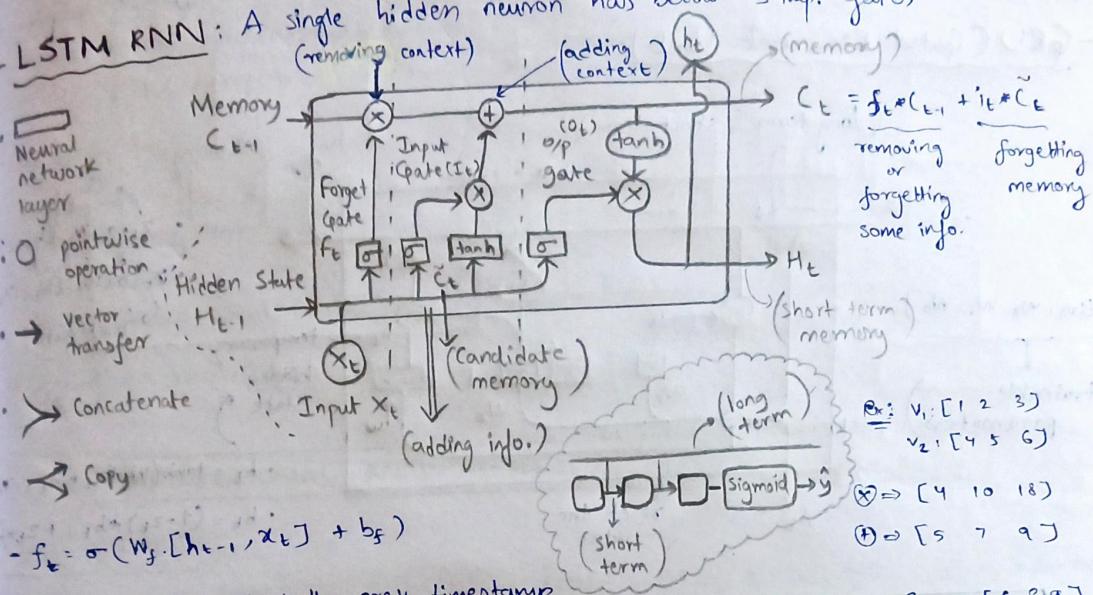
- create instance of Early stopping callback
- train the model with early stopping.
- save the model.
- Refer "app.py" → how to get movie review from saved model  
→ display it using streamlit.

(Refer Notebook)

# LSTM & GRU RNN INDEPTH INTUITION



LSTM RNN: A single hidden neuron has below 3 imp. gates  
(removing context) , (adding context) 



$$f_t = \sigma(w_f \cdot [h_{t-1}, x_t] + b_f)$$

-  $c_{t-1}$  = memory cell of the prev. timestamp

- $c_{t-1}$ : memory cell of previous timestamp
- $h_{t-1}$ : hidden state of previous timestamp

- $h_{t-1}$  = hidden state of
- $x_t$  = word passed as input in the current timestamp

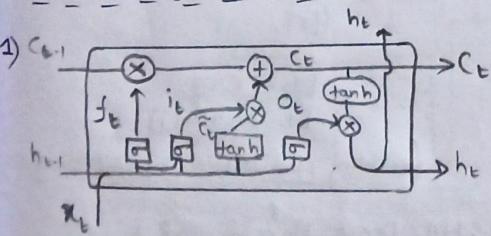
$\therefore$  Forget Gate ( $f_t$ ) = based on the context, forget gate will let go some information or will not let go info. (forgetting)

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

$$-\tilde{c}_t = \tanh(w_c \cdot [h_{t-1}, x_t] + b_c)$$

-  $[w_i, w_c, w_o] \rightarrow$  updating  $\leftarrow$  back propagation

## -Variants of LSTM RNN:



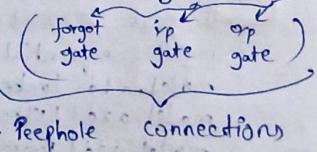
( we let the gate layers look at the cell state )

$$f_t = \sigma(W_f \cdot [C_{t-1}, h_{t-1}, x_t] + b_f)$$

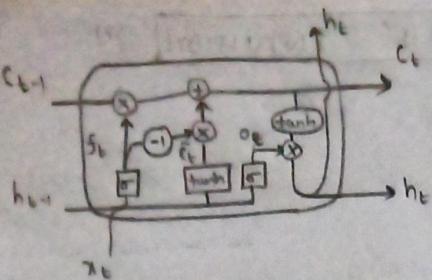
$$i_t = \sigma(W_i \cdot [C_{t-1}, h_{t-1}, x_t] + b_i)$$

$$O_t = \sigma(W_0 \cdot [C_t, h_{t-1}, x_t] + b_0)$$

(connections → from memory cell to )



2)

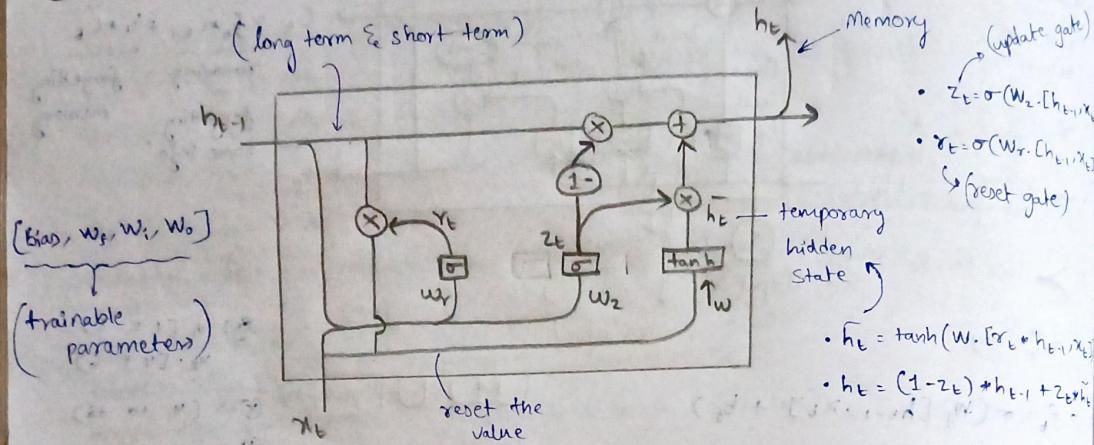


$$c_t = f_t * c_{t-1} + (1-f_t) * \tilde{c}_t$$

Goal: we only forget when we're going to i/p. something in its place.

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

### - GRU (Gated Recurrent Unit):



### LSTM & GRU E2E DL PROJECT - PREDICTING NEXT WORD

Ex: (Refer Notebook)

# Data Collection

```
import nltk
data = gutenberg.raw('shakespeare-hamlet.txt')
```

# Import → Tokenizer, pad\_sequences, train-test-split

tokenizer = Tokenizer()

tokenizer.fit\_on\_texts([text]) → assigns unique integer index to each word in a dataset

tokenizer.word\_index

# Create i/p sequence → refer notebook

↳ created in this way:

# Pad sequence

# Create predictors & labels

# Split into training & testing

# Define early stopping

# Define & train LSTM/GRU model

↳ Import Sequential, Embedding, LSTM, Dense, Dropout, GRU

model = Sequential()

model.add(Embedding(total\_words, 100, input\_length=max\_sequence\_len-1))

model.add(LSTM(150, return\_sequences=True))

model.add(Dropout(0.2))

model.add(LSTM(100))

model.add(Dense(total\_words, activation="softmax"))

model.compile(loss="categorical\_crossentropy", optimizer="adam", metrics=[accuracy])

[1,

[1, 687],

[1, 687, 9],

[1, 687, 9, 45]]

(disabling some neurons to avoid overfitting)

model.summary() # Replace 'LSTM' with 'GRU' to use 'GRU'  
 history = model.fit(x-train, y-train, epochs=50, validation\_data=(x-test, y-test),  
 verbose=1, callbacks=[EarlyStopping()])

# func. to predict next word

# Save the model & tokenizer  
 model.save('next\_word\_lstm.h5')

import pickle

pickle.open('tokenizer.pickle', 'wb') as handle:

pickle.dump(tokenizer, handle, protocol=pickle.HIGHEST\_PROTOCOL)

(ensures that we use the most eff. serialization format available)

# Custom text

input\_text:

max\_sequence\_len = model.input\_shape[1]

next\_word =

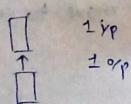
print(f"Next Word Prediction : {next\_word}")

# Deploy on Streamlit

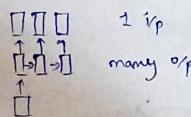
## BIDIRECTIONAL RNN ARCHI. & IN DEPTH INTUITION

Types of RNN:

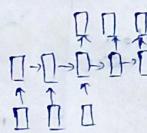
1) One to One RNN:



2) One to Many RNN:



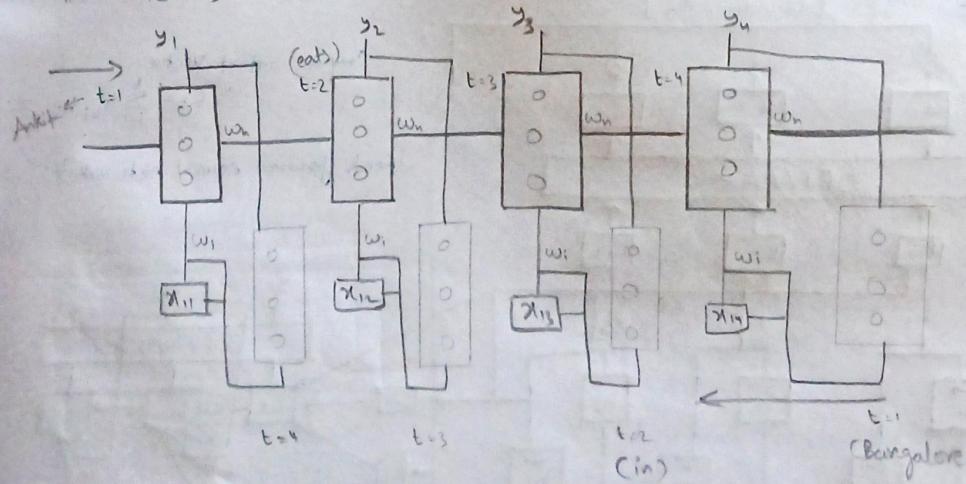
3) Many to Many RNN:



e.g. image captioning

Bidirectional RNN:

e.g. Ankit eats in Bangalore



Now only few prop. couldn't predict ' $x_{13}$ ' accurately, so we run another RNN backwards, so now the model can accurately predict ' $x_{13}$ ' since it has context of its upcoming sentences as well.

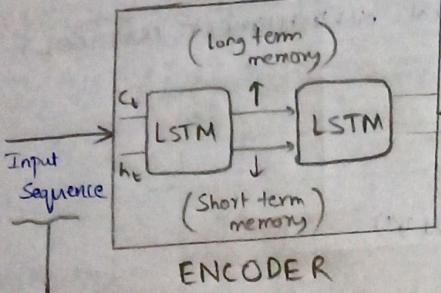
# ENCODER DECODER + ATTENTION MECHANISM - SEQUENCE2SEQUENCE

↓  
ex: One language to other!

## Encoder-Decoder:

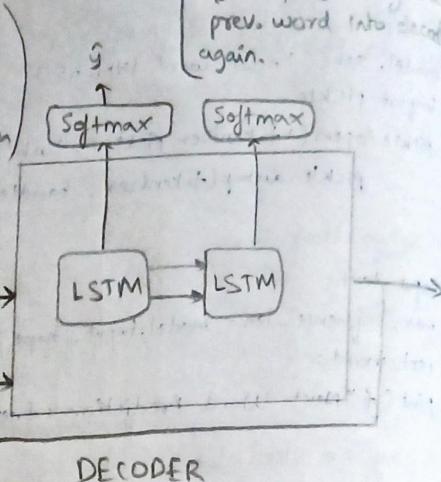
Architecture

(converts words to array of numbers)



Hidden states get passed through

0.1  
0.8  
-0.3  
0.6  
0.1



DECODER

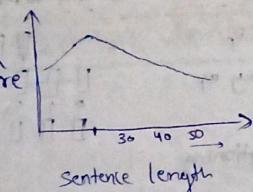
represents the entire sentence

Embedded Layer

$x_1, x_2, x_3 \rightarrow$  passed according to time  
(SOS) ↓  $\langle \text{EOS} \rangle$  → end of sentence  
start of sentence

Drawbacks:

{ performance metrics }

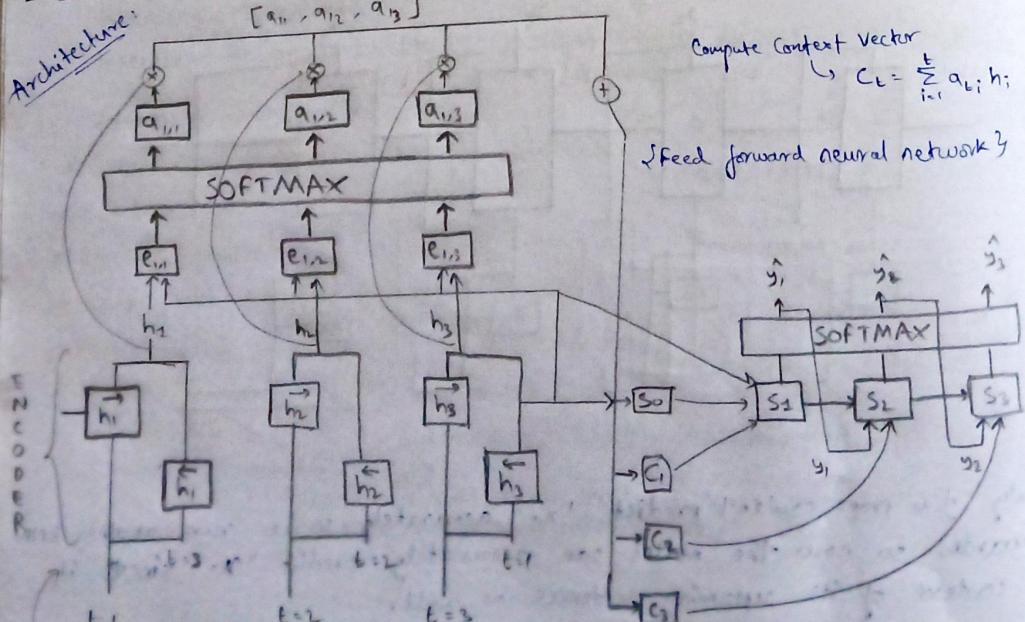


RNN → had vanishing gradient problem, so we have to move to encoder-decoder.

Attention Mechanism: → To resolve above drawback, we use this! {Refer video for detailed explanation}

Architecture:

$[a_{1,1}, a_{1,2}, a_{1,3}]$



Compute Context Vector

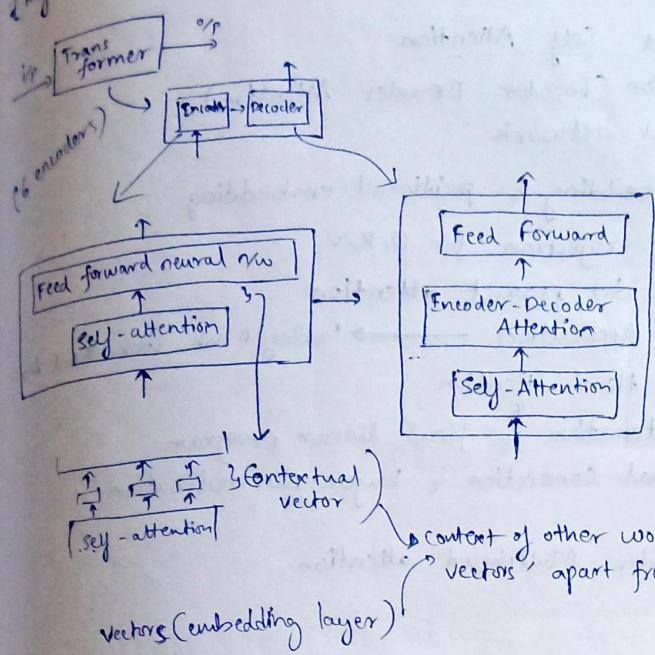
$$C_t = \sum_{i=1}^t a_{t,i} h_i$$

{Feed forward neural network}

(bidirectional RNN)

# TRANSFORMERS

{Refer Architecture Diagram in Pdf 3}



REFER PDF  
FOR INDEPTH  
EXPLANATION

content of other words also feeded in 'contextual vectors' apart from the current word.

Vectors (embedding layer)

Self-attention: aka scaled dot-product attention, is a crucial mechanism in transformer archi. that allows the model to weigh the importance of different tokens in the ip sequence relative to each other.

Inputs = Queries, Keys, Values

Query Vector: represent the token for which we are calculating the attention

Key Vector: represent all the tokens in the sequence & are used to compare with the query vectors to calculate attention scores.

Value Vector: holds the actual info. that will be aggregated to form the op of the attention mechanism.

- 1) Token Embedding
- 2) Linear Transformation
- 3) Compute Attention Scores
- 4) Scaling
- 5) Apply Softmax
- 6) Weight sum of values

Imp. Terms:

• Self Attention with multi heads  
 { it expands model's ability to focus on different position of token }

• Positional Encoding  
 ↘ Sinusoidal PE      ↘ Learned PE  
 (representing order of sequence)

- Layer Normalization & Batch Norm.
- Residual Connection

- The transformer decoder is responsible for generating the output sequence 1 token at a time, using the encoder's output & the previously generated tokens.

1) Masked Multi Head Self Attention

2) Multi Head Attention (Encoder - Decoder Attention)

3) Feed forward neural network

→ a) Input embedding & positional embedding

b) Linear projection for Q, K, V

c) Scaled dot product attention

d) Mask Application → Padding Mask, Look Ahead Mask

e) Multi Head Attention

f) Concatenation & final linear projection

g) Residual connection & Layer Normalization

• Imp Concepts: Encoder Decoder Multihead attention

E2E PROJ. (MLOPS) WITH ETL PIPELINES - BUILDING

NETWORK SECURITY SYSTEM

!! Refer Obsidian !