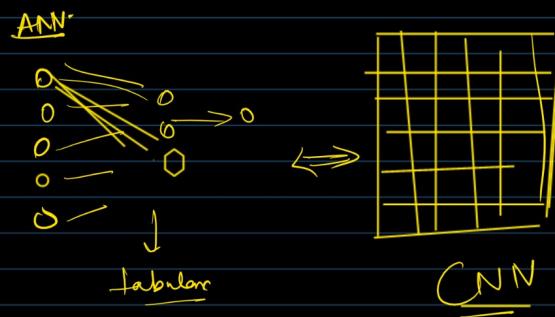
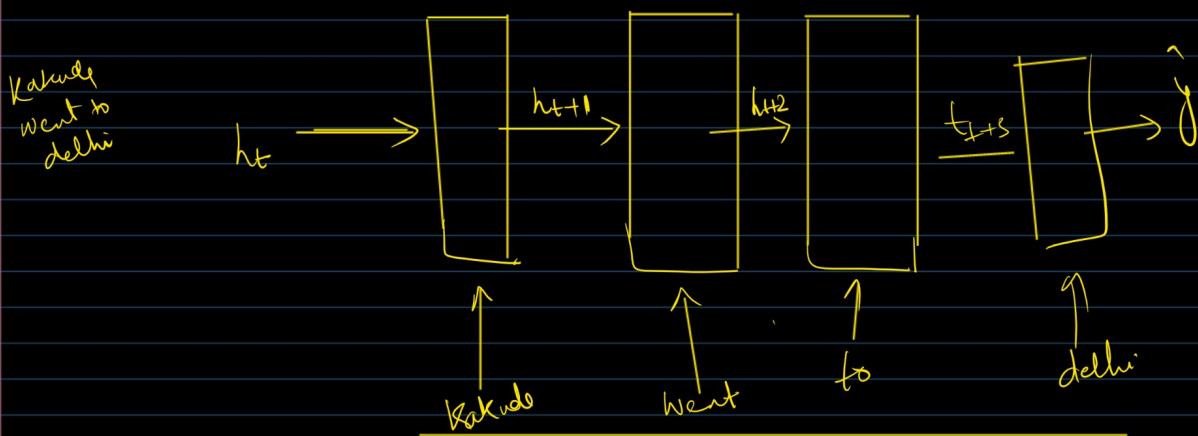


RNN

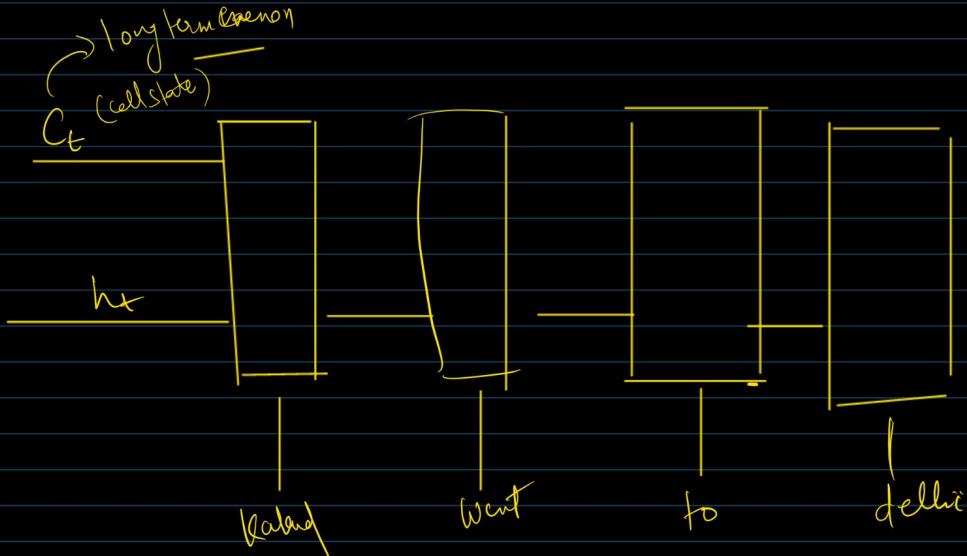


"Kakude went to Delhi. He met Manjo. Manjo loves ice cream. ~~akuda~~
loves ice cream both love each other."

RNN



LSTM



① Forget gate

② Update gate

③ Output gate

GRU $\rightarrow U_t / Z_t \rightarrow$ Update gate

\rightarrow Reset gate $= \sigma_t$

fewer params. / easy to train

sits on the mat.

$\uparrow \uparrow \uparrow \uparrow \uparrow$

One to many



\uparrow

Cat

eating something

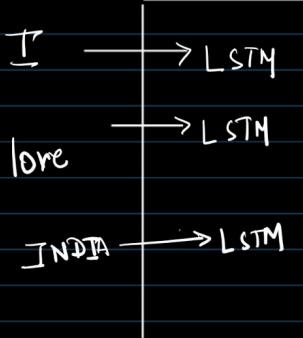
many to many



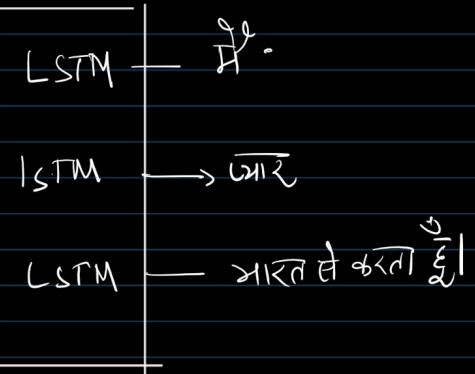
Cat on the mat

Seq 2 Seq model

Encoder



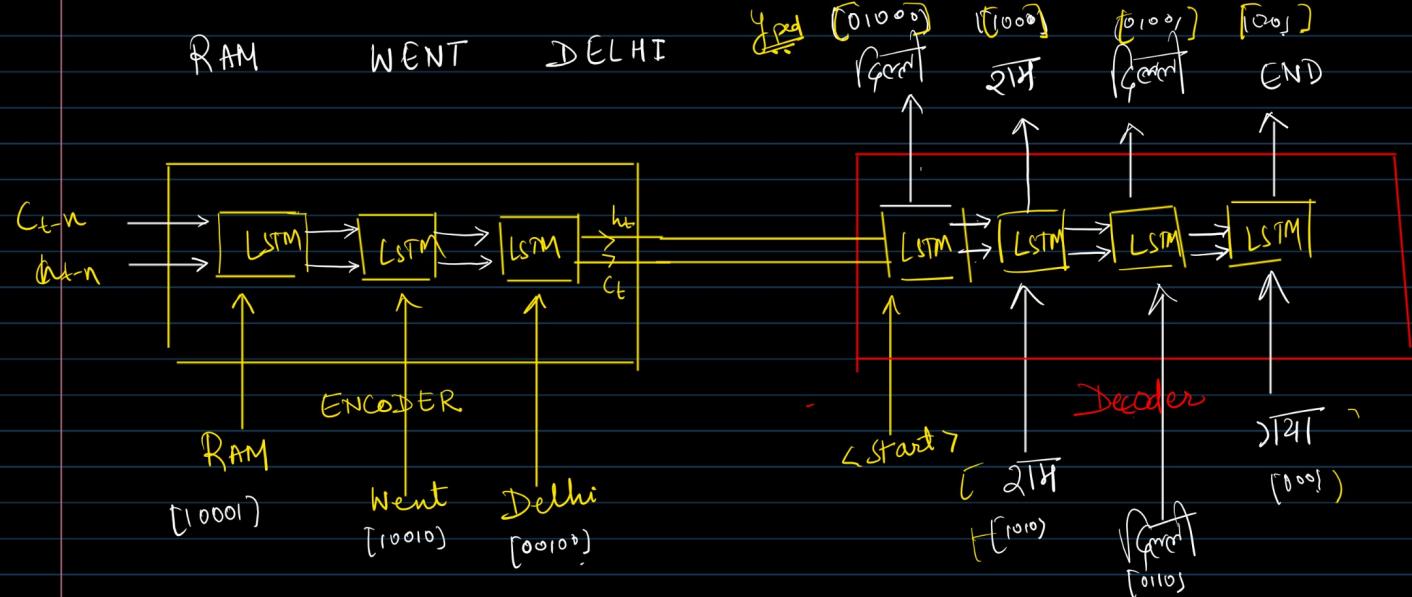
Decoder



sequences. In this paper, we present a general end-to-end approach to sequence learning that makes minimal assumptions on the sequence structure. Our method uses a multilayered Long Short-Term Memory (LSTM) to map the input sequence to a vector of a fixed dimensionality, and then another deep LSTM to decode the target sequence from the vector. Our main result is that on an English to French translation task from the WMT-14 dataset, the translations produced by the LSTM

\rightarrow Ram went delhi \rightarrow राम दिल्ली गया

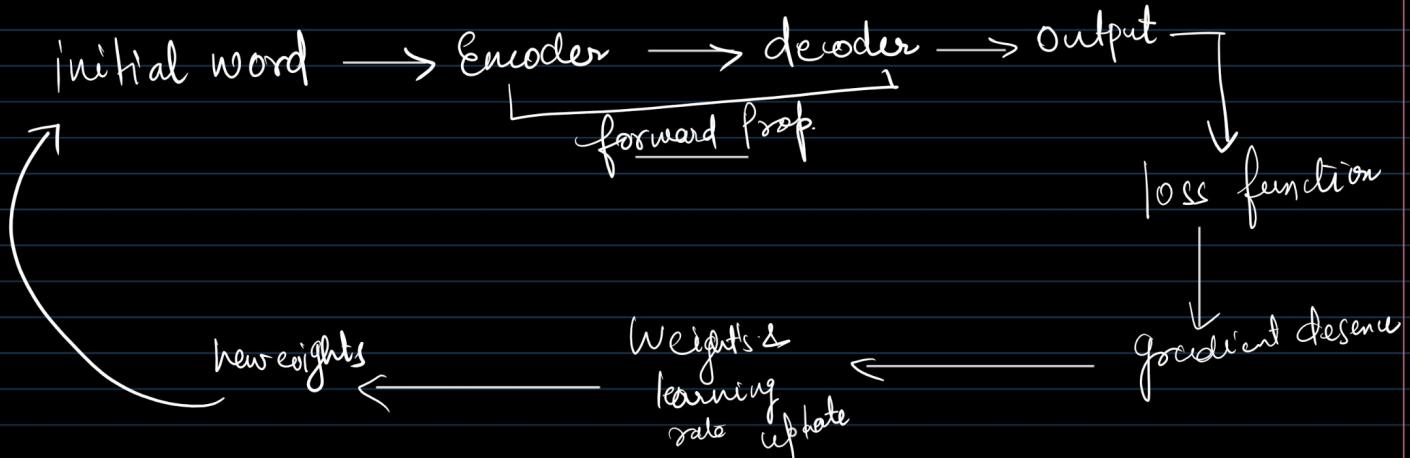
\rightarrow Shyam came home \rightarrow श्याम घर आया



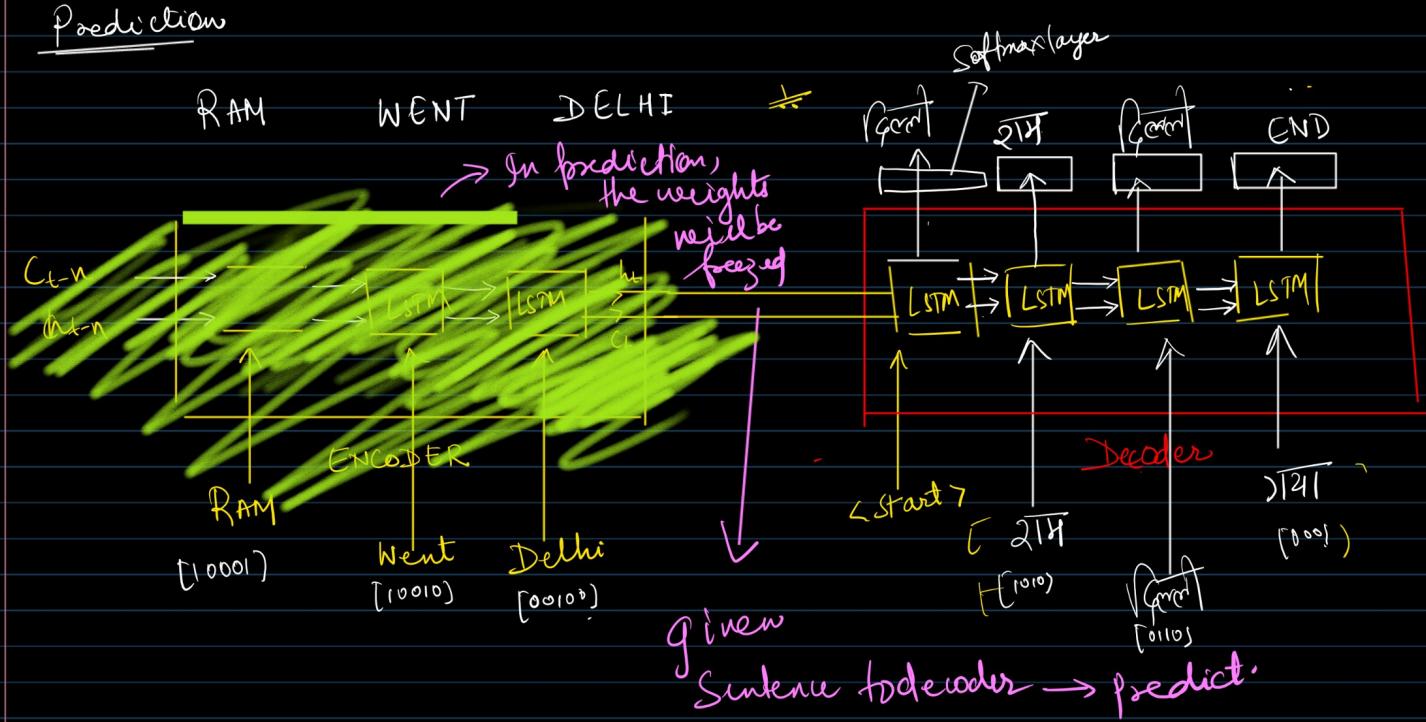
\rightarrow Shyam Came home

there is difference

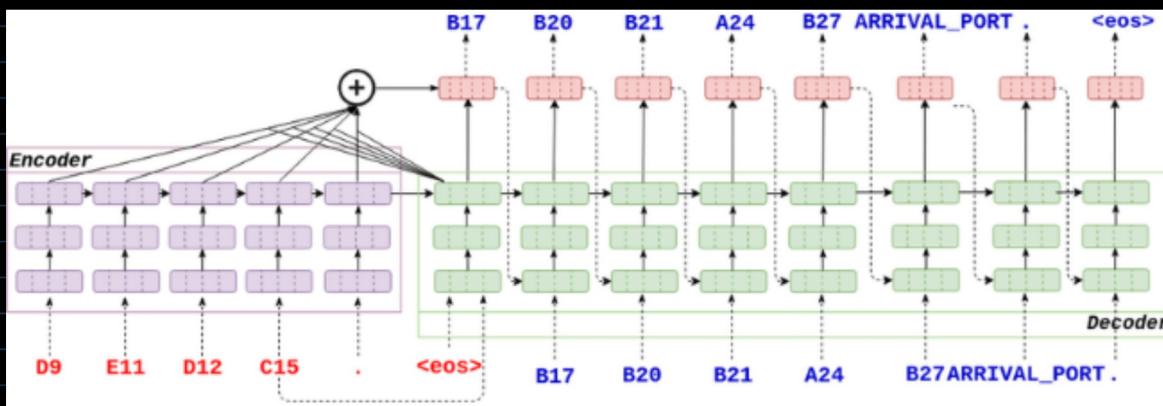
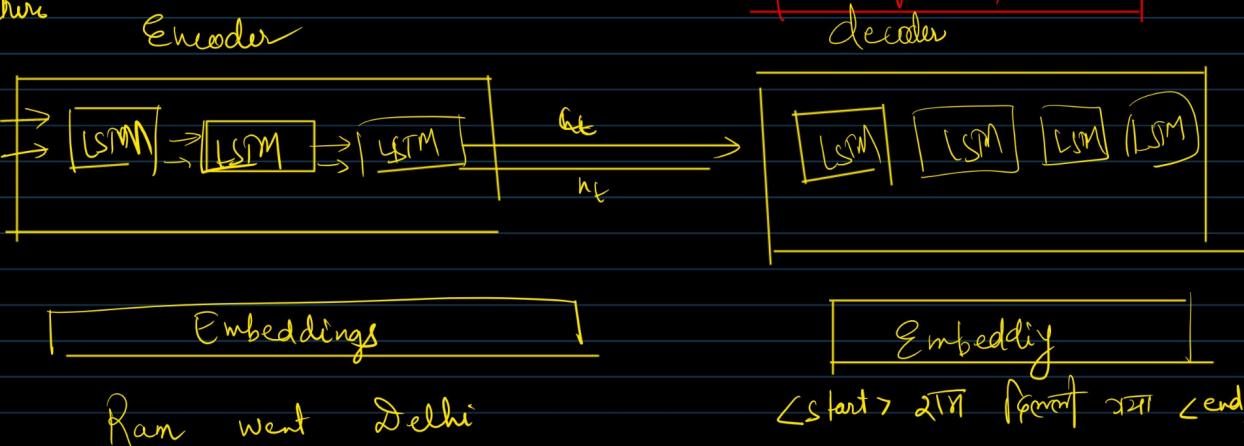
$\left\{ \begin{array}{l} \text{Actual} \rightarrow \text{PREDICTED} \\ \text{Pred} \rightarrow \text{PREDICTED} \end{array} \right.$
 loss function \rightarrow Categorical cross entropy
 $\approx -\sum y_i \log p_i$



Prediction



Architecture



Use Case

- Text summarization
- Captioning
- QnA

Disadvantage of Seq2Seq models

- * 30-50 words → performs well → the performance starts decreasing
- * Doesn't understand long term context

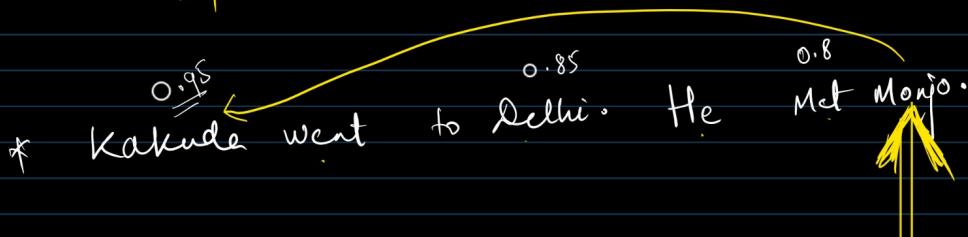
Test Case

Kakuda went to Delhi. He met Monjo. Monjo loves icecream. He loves ice cream both love each other! Both got married. After 5 years they had a child, the child had a child after 30 years.



- * At each state of decoder during time we are accessing all the information.

↓
Weight that we are seeing is distributed and previously occurring words might become less important.

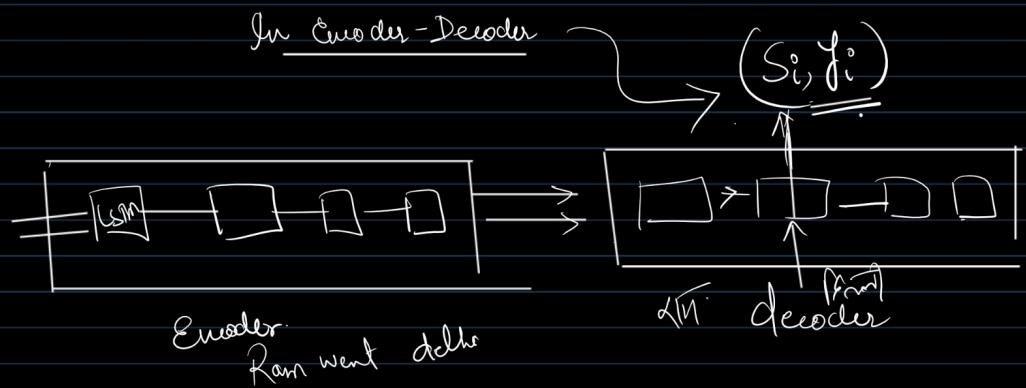


- * Since we are human being, Our mind understands this that Kakuda is very very important

↓
Solu ⇒ Context word Understanding.

↓
By giving more Attention to
imperatant keywords.

Attention Model

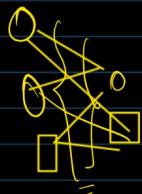
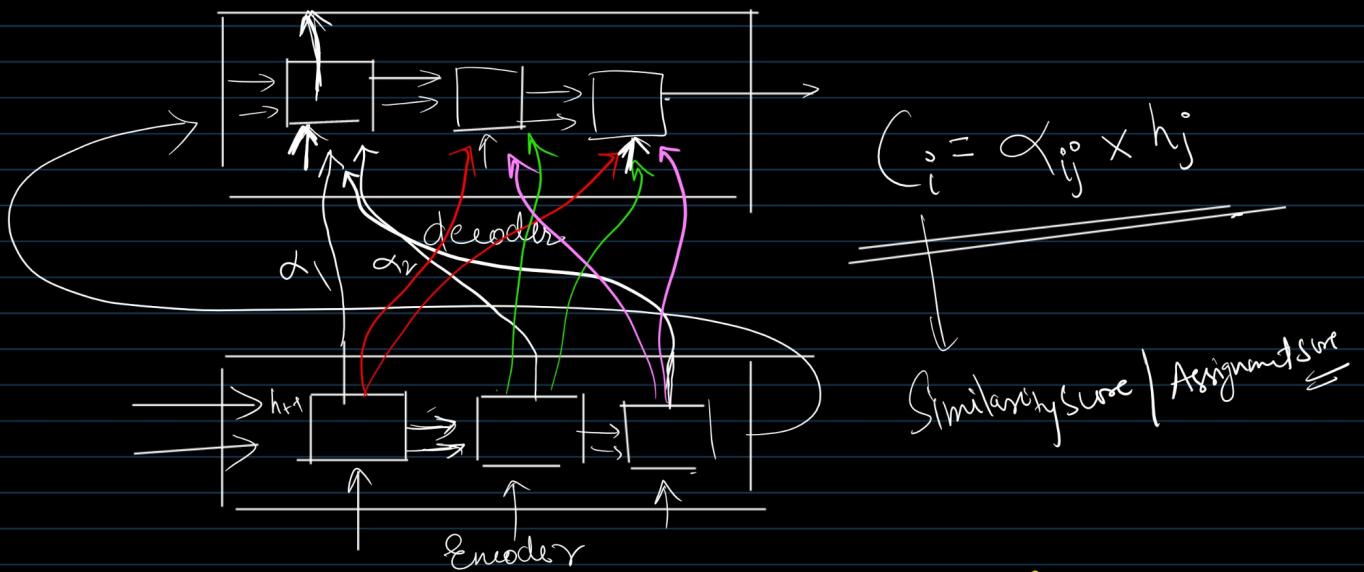


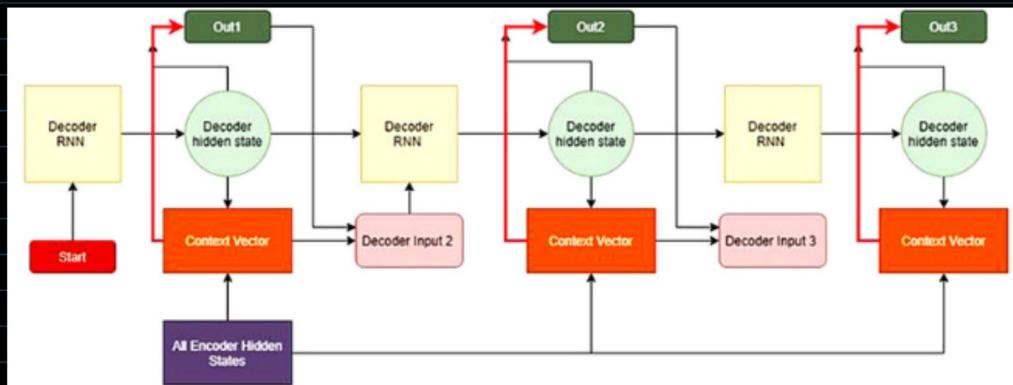
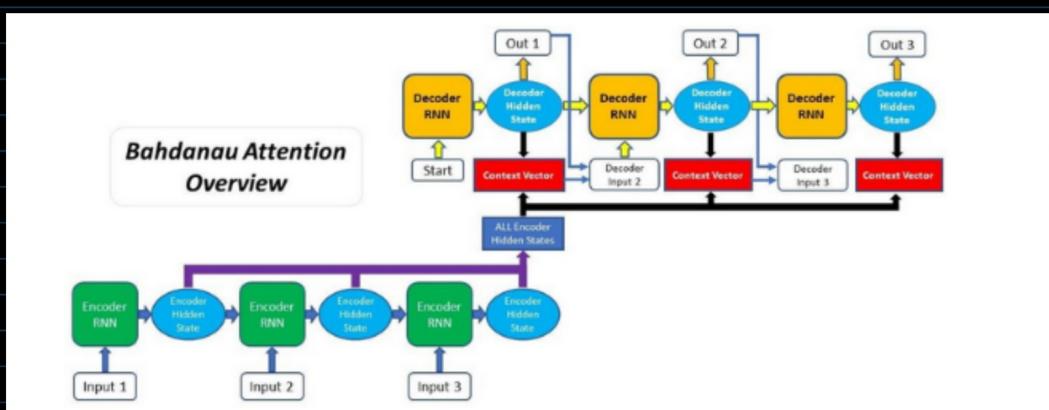
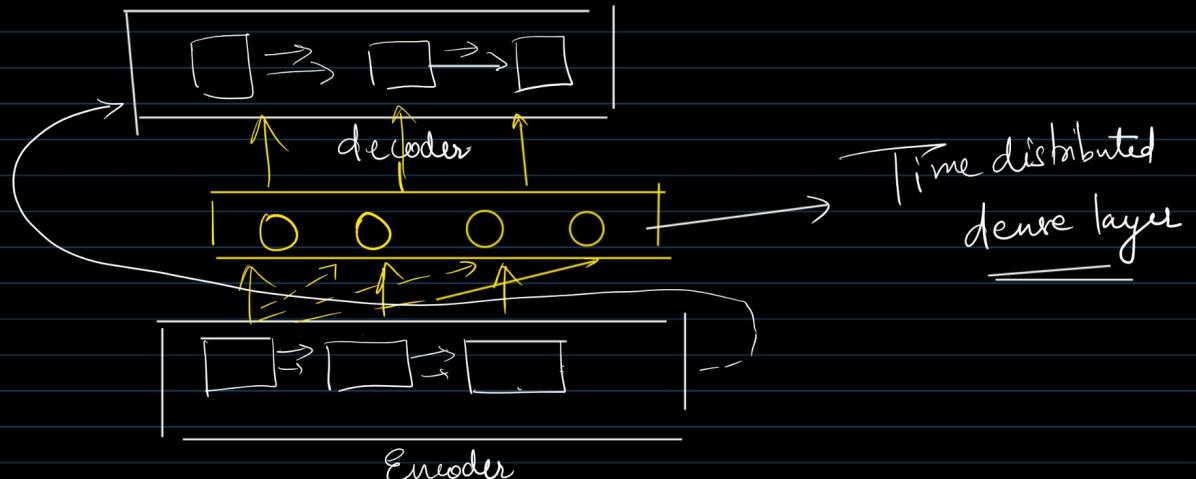
Attention Mechanism

(S_i^*, y_i, c_i)

→ attention input

→ this is a vector with each word importance wrt





Or ANN - tabular data, CNN - Image

Or Seq2seq - RNN \rightarrow LSTM \rightarrow GRU

Seq2seq \rightarrow Encoder-decoder \rightarrow Attention mechanism

Summary

Encoder - RNN | LSTM | GRU.

decoder - RNN.

Attention - extra layer \rightarrow It allowed decoder to look at all encoder hidden states.



disadvantage

\rightarrow slow training. \rightarrow More parameters.

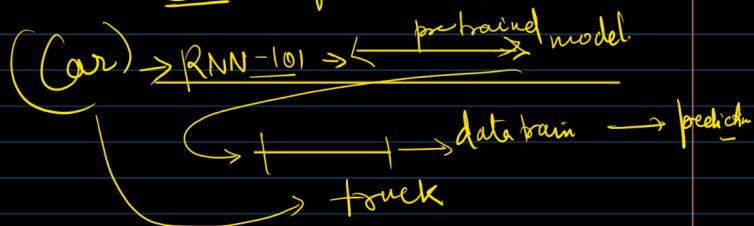
\rightarrow Limited Scalability

2017 \rightarrow Attention is All you need
↓
2018 \rightarrow transformer in NLP

GPT \rightarrow Generative Pre-training

transfer learning

CV \rightarrow famous in CV.



{ 2018 \rightarrow transfer learning was introduced

+ { 2017 \rightarrow Transformer

→ GPT - 2018 \rightarrow Dee.

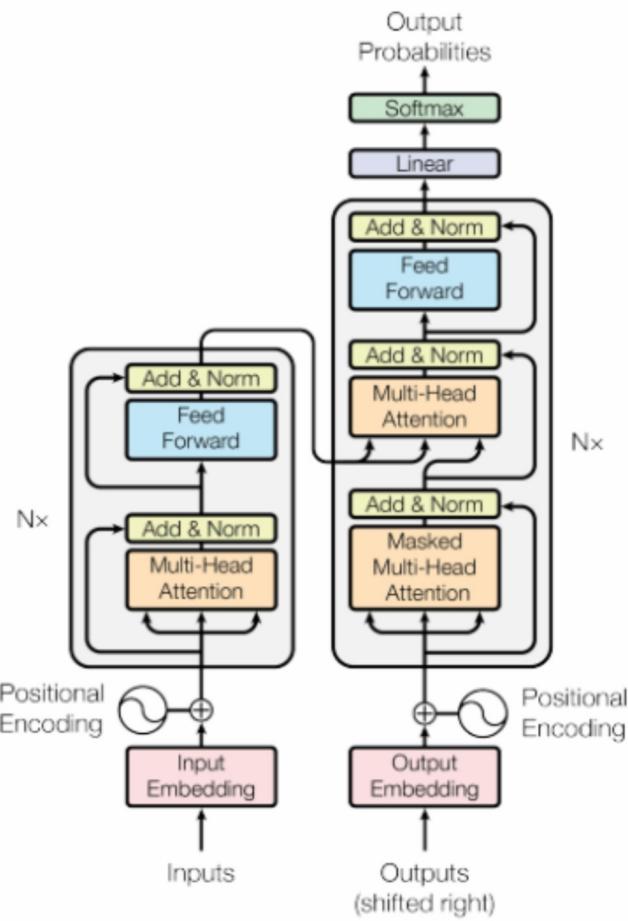


Figure 1: The Transformer - model architecture.

Why Transformers

- RNN / LSTM took lot of training time & not scalable
- removed RNN / LSTM entirely
- kept self attention as the core computation
- parallel computation
- fast (GPU friendly)

Static Embedding → I went to Bank
 $(1000) \quad (0001) \quad (1000) \quad (0001)$

dynamic Embedding

I went to Bank. The river Bank was clean.
 $[0001] \quad [0001]$

So we need contextual embedding.

Self Attention

Ram went Bank. Bank river was clean.

$$\text{bank} = 0.5 \times \text{bank} + 0.2 \times \text{went} + 0.3 \times \text{bank}$$

↓

$$= [0.3 \quad 0.5 \quad 0.2 \quad 0.1]$$

$0.2 \times \text{bank} + 0.5 \times \text{River}$
 $+ 0.3 \times \text{was}$
 $\times 0.1 \text{ clear}$

$[0.6 \quad 0.5 \quad 0.4 \quad 0]$

Both embeddings will be different.

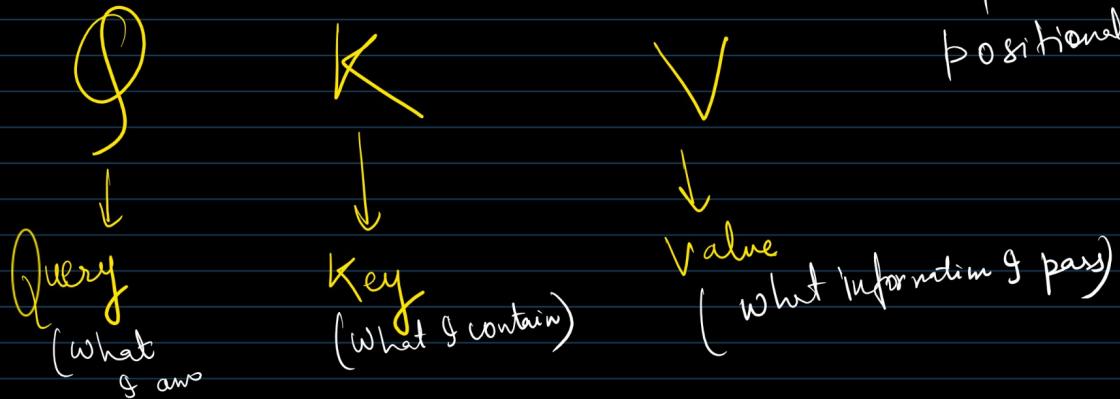


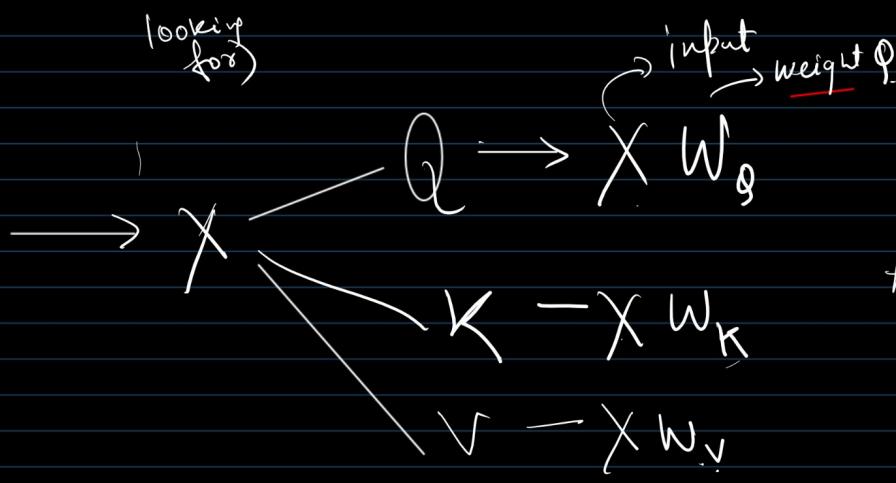
→ operation is parallel.

So training it will
be very cost but
on a cost of loops
Sequence.

* Contextual embeddings → became very powerful.

+ positional encoder
(add order)





$$\text{Attention Score} = Q K^T$$

\sum_{d_k}
 gradient
 stability

Multi head mechanism

$W_Q, W_K, W_V \rightarrow$ learns different relation =

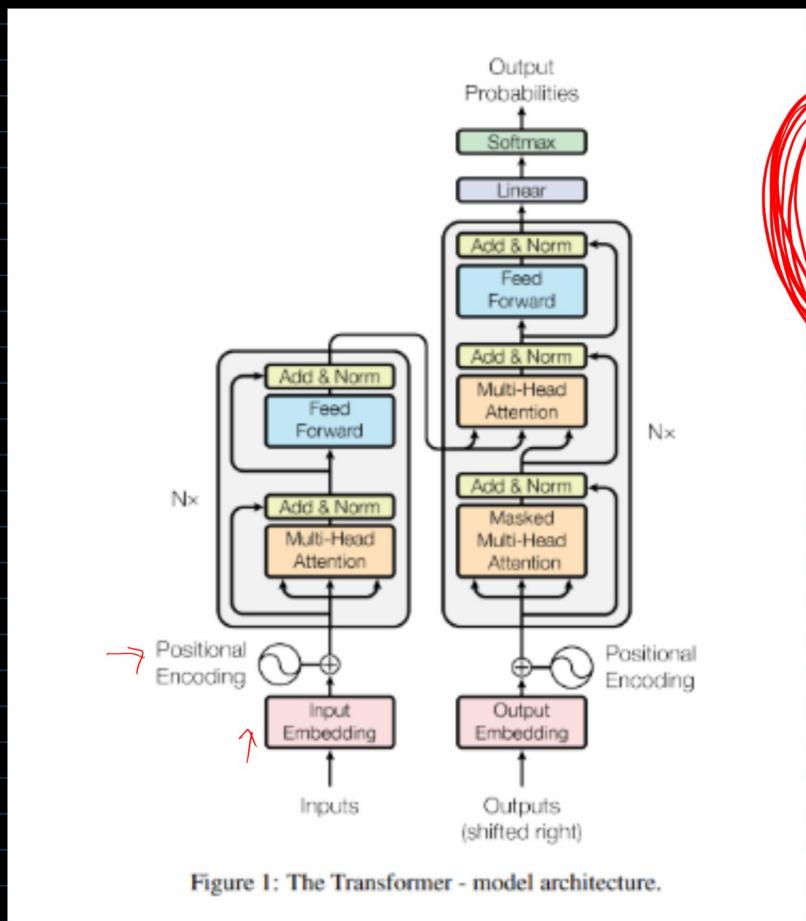


Figure 1: The Transformer - model architecture.

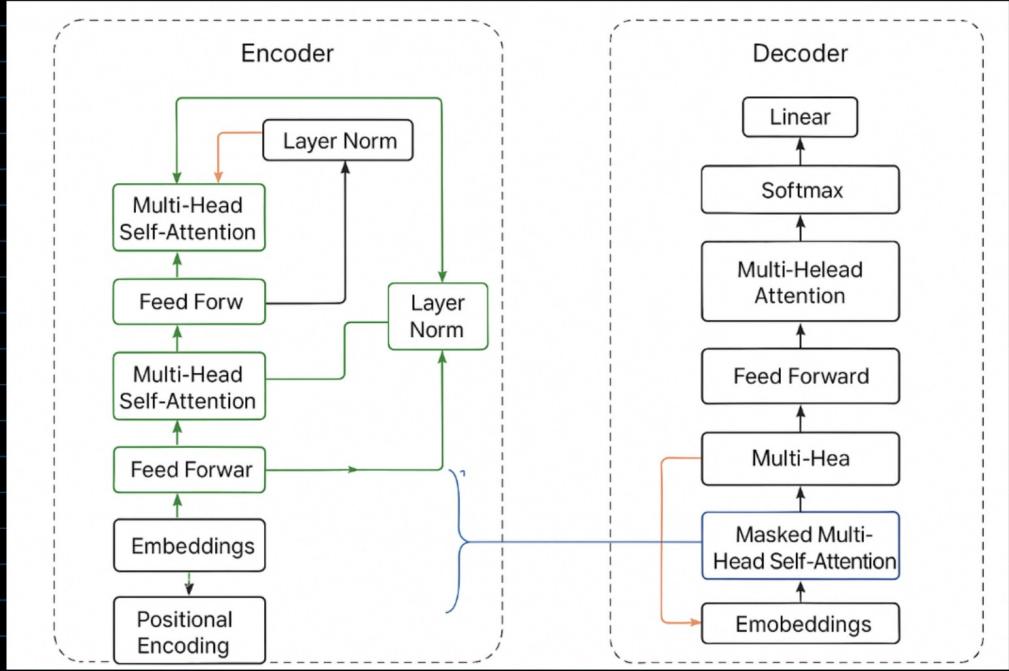
Encoder

- ① Input Embedding + positional Embeddings
- ② Multi-head Self attention \rightarrow Residual + norm
- ③ FFN (dense layer)

↓
 Residual
 ↓
 Yact (y)
 Cost funct.
Contextual Embedding

② Decoder

- \rightarrow Output Embedding + positional Embedding.
- \rightarrow Multi head Attention
- \rightarrow Cross attention
- \rightarrow FFN (dense layer)



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

Bert \rightarrow Bi-directional Encoder Representation from Transformer

Before Transformers \rightarrow Word2Vec (static)

Transformer \rightarrow It processes the complete data in one go.

\rightarrow Contextual Embedding

\rightarrow Processed from left to Right (GPT)
and right to left

Bert

\hookrightarrow bidirectional context problem got solved

River bank Money bank
→ Only Encoder part of Transformer
is Bert.

Bert → RoBERTa.

→ ALBERT

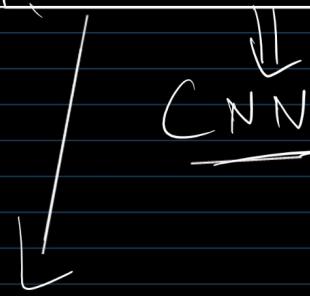
→ DISTILBERT

→ SpanBERT

→ SciBert → Scientific text

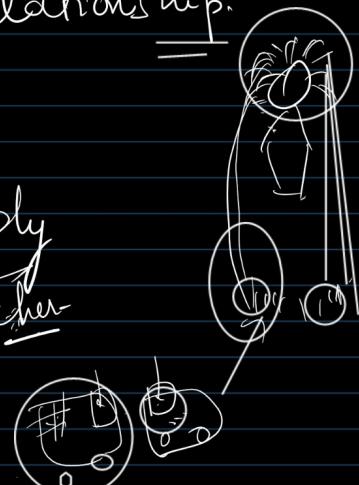
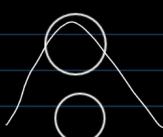
ViT → Vision transformer

CV → Yolov7 or Mask RNN



Struggles with global relationships

Transformer in NLP → apply on image patches



idea

ViT → Vision transformer → treats an image as
a sequence of patches, just like
a sentence is a sequence of words

Analogy

→ Sentence → Words → Embedding → Transformer.

↳ Image → Patches → Embedding → Transformer.

disadvantage

↳ Large dataset for pre-training

DeiT

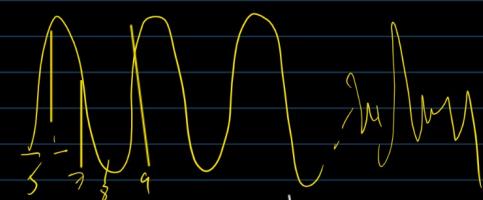
beit

Image ($H \times W \times C$)
↓ (split into patches)

patches

Embedding

↓
Transformers



→ long seg of amplitude

↓
Fourier Transform, Spectrogram

RNN, LSTM

Wav2Vec

