Attention Models                    20/11/2022

Encoder-Decoder → Attention Models → Transformers → Hugging
                                                     Face
                                                       ↑
                    Conceptual Session          Python
                                                Implementation

Introduction:  Attention: most influential ideas in DL.
                    ↳ NMT (Neural M/c Translation)
                       ↳ Seq 2 Seq Models
                          ↳ Encoder - Decoder architecture

English  ⟶  Marathi

process i/p seq                It is initialized
& encode/compress/             with context vector,
summarize info into            using which it starts
CONTEXT/THOUGHT vector         generating
                               transformed o/p.

I/P English: Rahul is a good boy.

Target Marathi: राहुल चांगला मुलगा आहे.

RNN:     LSTM/GRU → complex sequence → large amt of data.

       Appl.: Speech recognition, NLP, Time forecasting, etc.
                                             series

Seq2Seq

## Machine Language Translation ✓

Les modèles de séquence sont super puissants → Sequence Model → Sequence models are super powerful

## Text Summarization ✓

A strong analyst have 6 main characteristics. One should master all 6 to be successful in the industry :
1. ............
2. ............
→ Sequence Model → 6 characteristics of successful analyst

## Chatbot

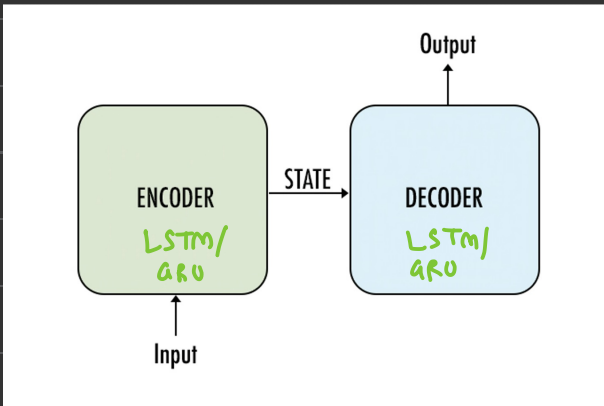eg: Question Answering.

How are you doing today? → Sequence Model → I am doing well. Thank you. How are you doing today?

Prerequisites:   1. Fundamental concepts in ML & NN
             2. Linear Algebra & Prob.
             3. Working knowledge of LSTM.

# Encoder - Decoder Architecture:



**Encoder:** read i/p seq.
↓
Summarize info as Internal State vectors (Hidden state & Cell state).

o/p of encoder are discarded & we only preserve the Internal State.
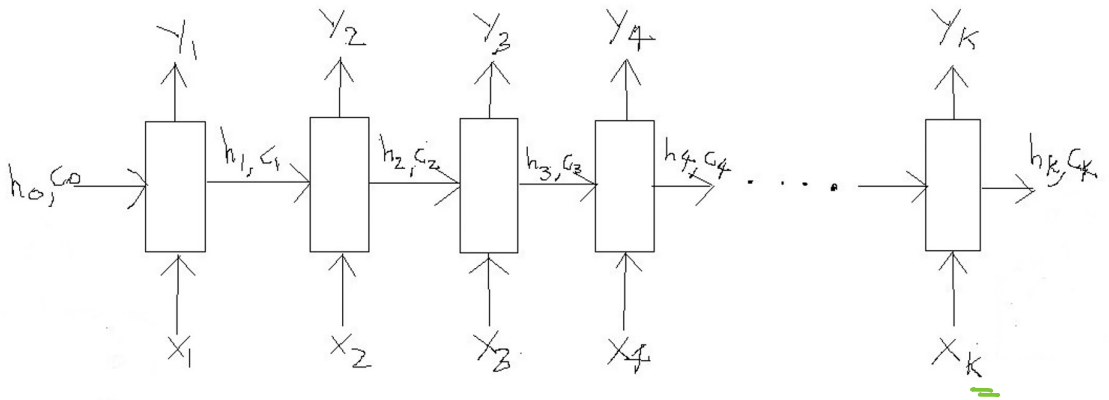
**Decoder:** I/p is final states of encoder
Using this, Decoder starts generating o/p seq.
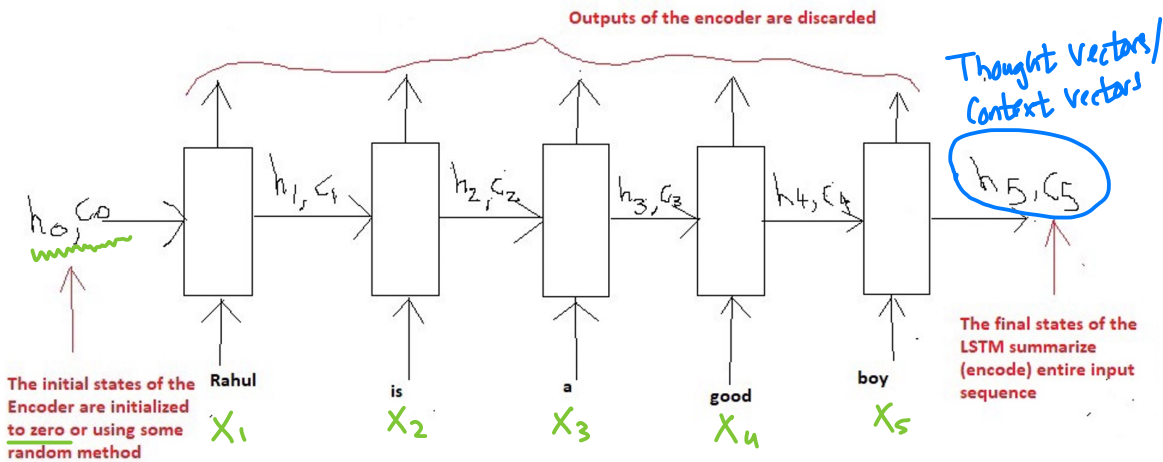
Decoder Behaves differently in ┬ Training (uses Teacher-forcing)
                                └ Inference (i/p to decoder at each time step is o/p from previous time step).

<u>Encoder</u>  LSTM: 3 components: 1. $X_i \rightarrow$ i/p seq. at time step i

2. $h_i$ & $c_i \rightarrow$ internal states

3. $y_i \rightarrow$ o/p seq. at time step i



## Word Level NMT:

Outputs of the encoder are discarded

Thought vectors/
Context vectors

$h_5, c_5$

The final states of the LSTM summarize (encode) entire input sequence

The initial states of the Encoder are initialized to zero or using some <u>random</u> method

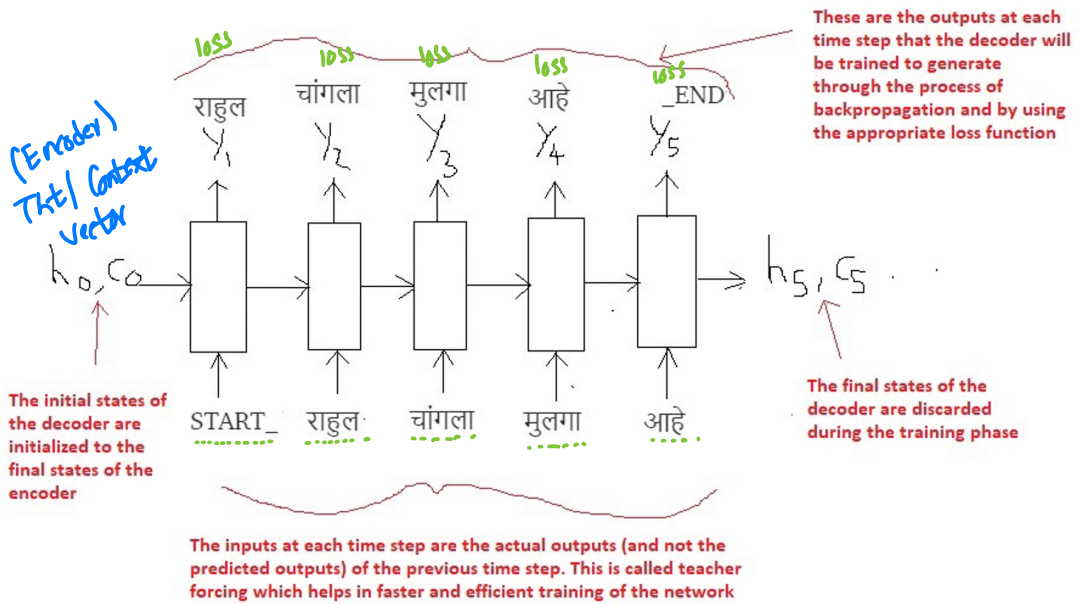Rahul  $X_1$     is  $X_2$     a  $X_3$     good  $X_4$     boy  $X_5$



$h_i$ & $c_i$: They remember what the LSTM has read (learned) till now

$y_i$ : o/p (predictions) of LSTM at each time step)

$y_i$ is prob. dist over entire vocabulary
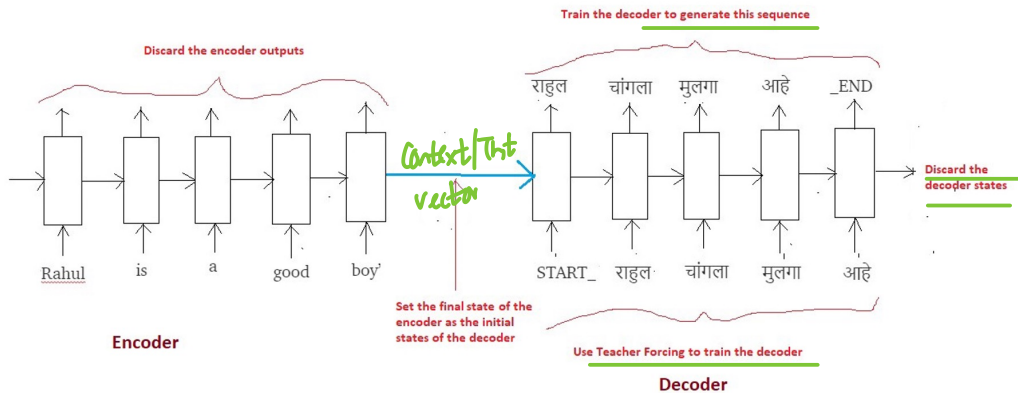
# Decoder LSTM - Training Mode.



These are the outputs at each time step that the decoder will be trained to generate through the process of backpropagation and by using the appropriate loss function

(Encoder) Tht/ Context vector

The initial states of the decoder are initialized to the final states of the encoder

The final states of the decoder are discarded during the training phase

The inputs at each time step are the actual outputs (and not the predicted outputs) of the previous time step. This is called teacher forcing which helps in faster and efficient training of the network

(TEACH)

Goal: Train the decoder to o/p राहुल चांगला मुलगा आहे.
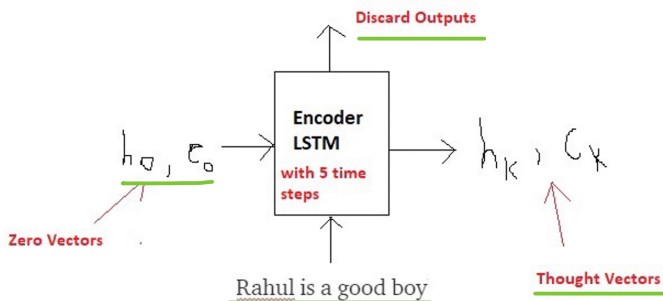
"START_  राहुल  चांगला  मुलगा  आहे  _END"

Teacher Forcing: Here i/p at each time step is given as actual o/p (& not the predicted output) from the previous time step
  - helps to train more fast.

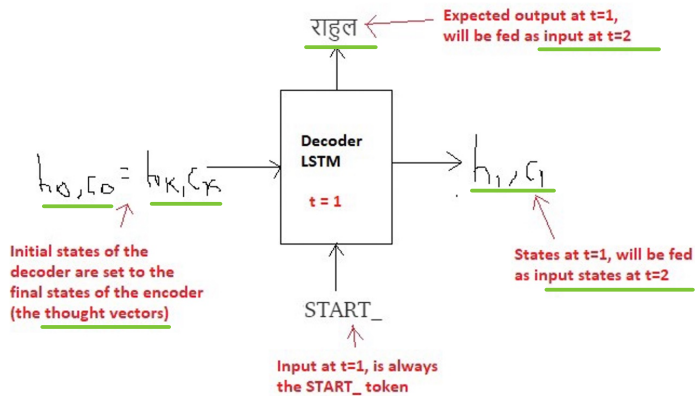https://machinelearningmastery.com/teacher-forcing-for-recurrent-neural-networks/

Encoder-Decoder diagram. Encoder inputs: Rahul, is, a, good, boy'. "Discard the encoder outputs". Context/Tht vector. "Set the final state of the encoder as the initial states of the decoder". Decoder inputs: START_ राहुल चांगला मुलगा आहे, outputs: राहुल चांगला मुलगा आहे _END. "Train the decoder to generate this sequence". "Use Teacher Forcing to train the decoder". "Discard the decoder states".

## Decoder LSTM — Inference Mode:

1. Encode i/p seq. into Thought Vectors.



$h_0, c_0$ — Zero Vectors → Encoder LSTM with 5 time steps → $h_k, c_k$ — Thought Vectors. Input: Rahul is a good boy. "Discard Outputs".
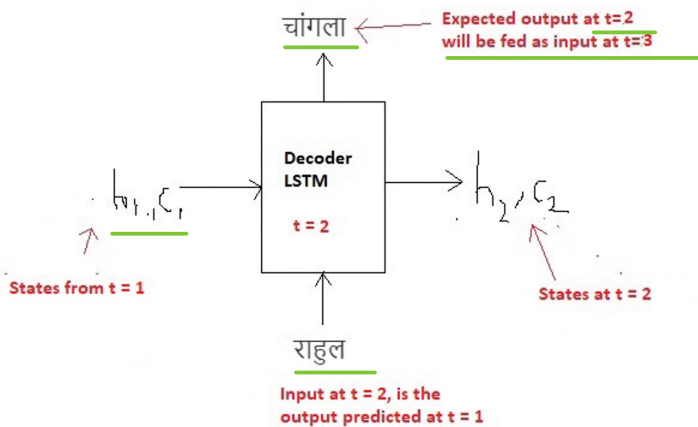
2. Start generating o/p seq. in a loop, word by word.

At $t = 1$,

राहुल ← Expected output at t=1, will be fed as input at t=2

$h_0, c_0 = h_K, c_K$

Decoder LSTM

t = 1

$h_1 / c_1$

Initial states of the decoder are set to the final states of the encoder (the thought vectors)

States at t=1, will be fed as input states at t=2

START_

Input at t=1, is always the START_ token

At t = 2,

चांगला ← Expected output at t=2 will be fed as input at t=3

$h_1, c_1$

Decoder LSTM

t = 2

$h_2 / c_2$

States from t = 1

States at t = 2

राहुल

Input at t = 2, is the output predicted at t = 1

At t = 3,

मुलगा ← Expected output at t=3 will be fed as input at t=4

$h_2, c_2$

Decoder LSTM

t = 3

$h_3 / c_3$

States from t = 2

States at t = 3

चांगला

Input at t = 3, is the output predicted at t = 2

**At t=4,**

आहे ← Expected output at t= 4 will be fed as input at t= 5

**Decoder LSTM**

t = 4

$h_3, C_3$

↑
States from t = 3

→ $h_4, C_4$

States at t = 4

मुलगा

Input at t = 4, is the output predicted at t = 3

**At t=5,**

_END ← Expected output at t=5 This indicates, that we stop

**Decoder LSTM**

t = 5

$h_4, C_4$

↑
States from t = 4

→ $h_5, C_5$

States at t = 5 Since we stop, we discard these states

आहे

Input at t = 5, is the output predicted at t = 4

# Inference Algorithm:



**LOOP**

The states and outputs at each time step, become the states and input respectively for the next time step

Discard the encoder outputs

राहुल  चांगला  मुलगा  आहे  _END

Rahul  is  a  good  boy'

**Encoder**

Set the final state of the encoder as the initial states of the decoder

START_  राहुल  चांगला  मुलगा  आहे

**Decoder**

# Why name Attention?



$h_1$  $h_2$  $h_3$  $h_4$  $h_5$

Rahul  is  a  good  boy

# Encoder GRU:



"Rahul is a good boy" → **ENCODER GRU** → [h1, h2, h3, h4, h5]

(All the intermediate states)

# How does Attention work?

# Decoding at time step1.



"Rahul is a good boy" → Encoder GRU → [h1, h2, h3, h4, h5] → **Feed Forward Neural Network** STEP 1 → [s1, s2, s3, s4, s5] **Softmax** → [e1, e2, e3, e4, e5]

(These are the attention weights)

h5

(this is the internal state of the decoder from the previous time step)

(these are the scores) STEP 2

Compute Context Vector (CV)

STEP 3

CV = e1*h1 + e2*h2 + e3*h3 + e4*h4 + e5*h5

Concatenate context vector with the output from previosu time step

The output of the decoder is the next word in the sequence as well as the internal hidden state.

← Decoder GRU ← concat_vector = [context_vector + "<START>"]

In time step 1:
    Output = "राहुल"
    Hidden state = d1

STEP 5

STEP 4

"राहुल"

Rahul

Compute a score

Compute attention weights : Apply softmax    e1, e2, e3, e4, e5
                            [0-1]

$e_1 + e_2 + e_3 + e_4 + e_5 = 1$

$e_1 = 0.75, e_2 = 0.20, e_3 = 0.02, e_4 = 0.2, e_5 = 0.01$

Compute Context Vector $= e_1 \times h_1 + e_2 \times h_2 + e_3 \times h_3 + e_4 \times h_4 + e_5 \times h_5$

Concaterate context vector with o/p of prev. step

$$\text{_____} \times \text{_____}$$