

Encoder - Decoder

#1) Seq2Seq Data



Example

Nice to meet you → Machine Translation → आप से भी हमें मिला रहा है।
 Input sequence → Output sequence

Q.1) What are the problems faced in Seq2Seq data?

Ans) 1) Variable Length Input

The input can be of any length & it can vary sample to sample

I/p → Sentence → English → variable Length

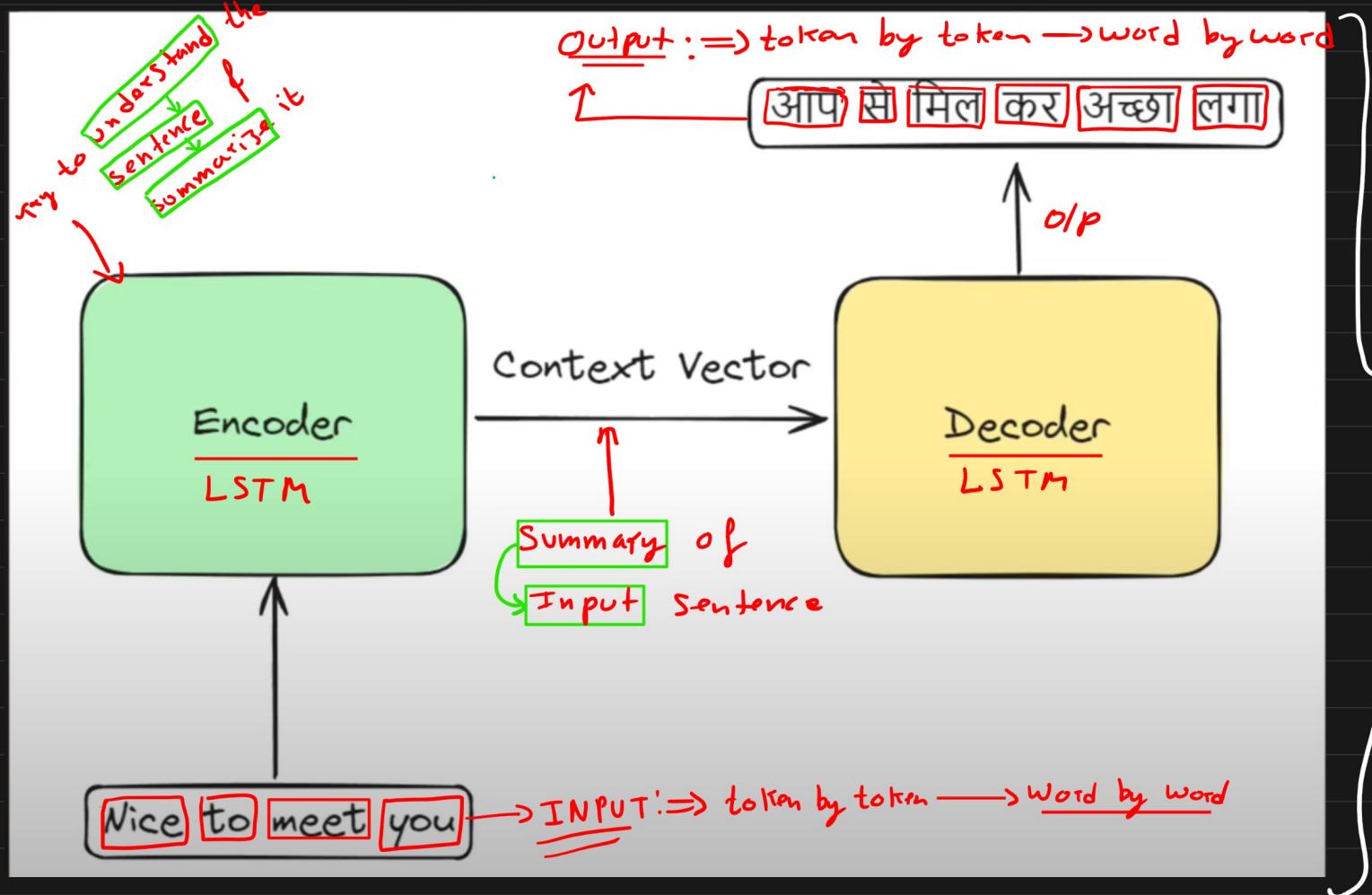
2) Variable Length Output

The output can be of any length & it can vary sample to sample

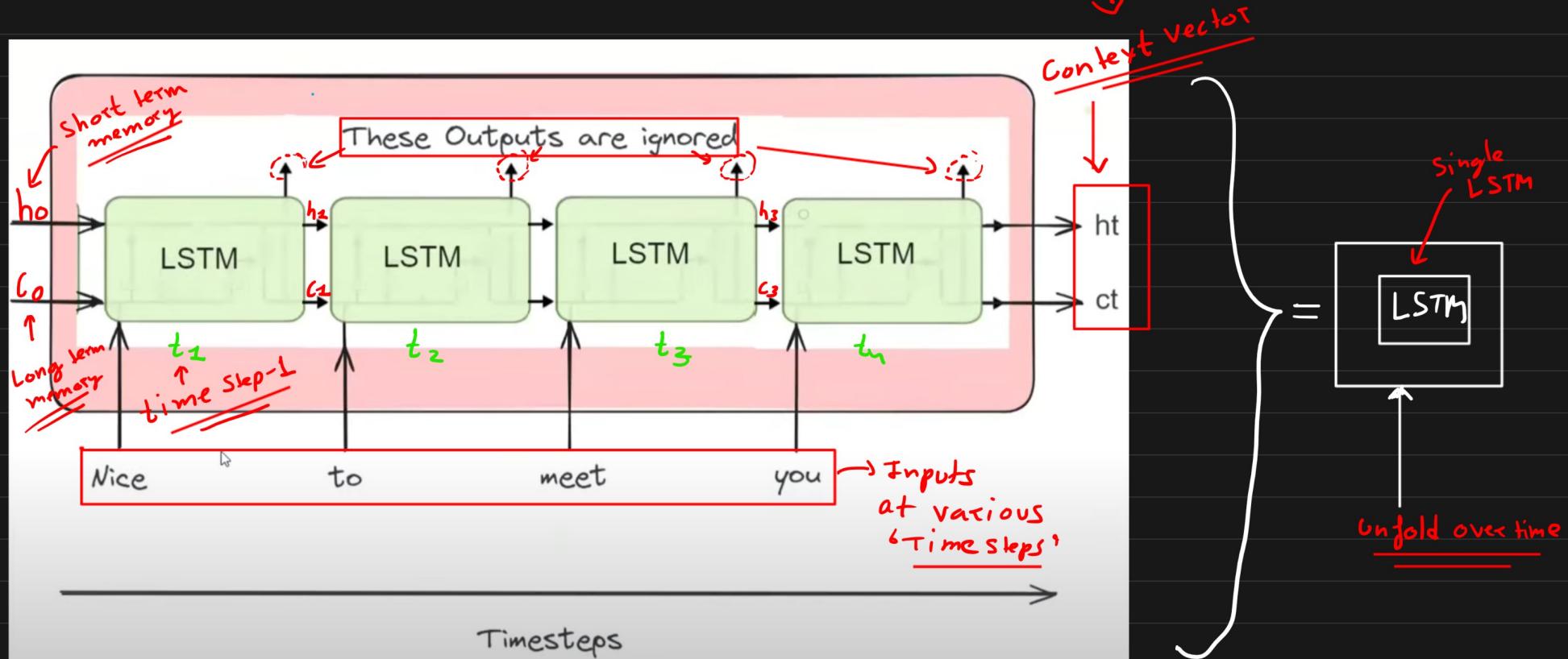
3) No guarantee that input & output for a sample will be of same length

Nice to meet you → Machine Translation → आप से भी हमें मिला रहा है।
 Input Length = 4 → Output Length = 6

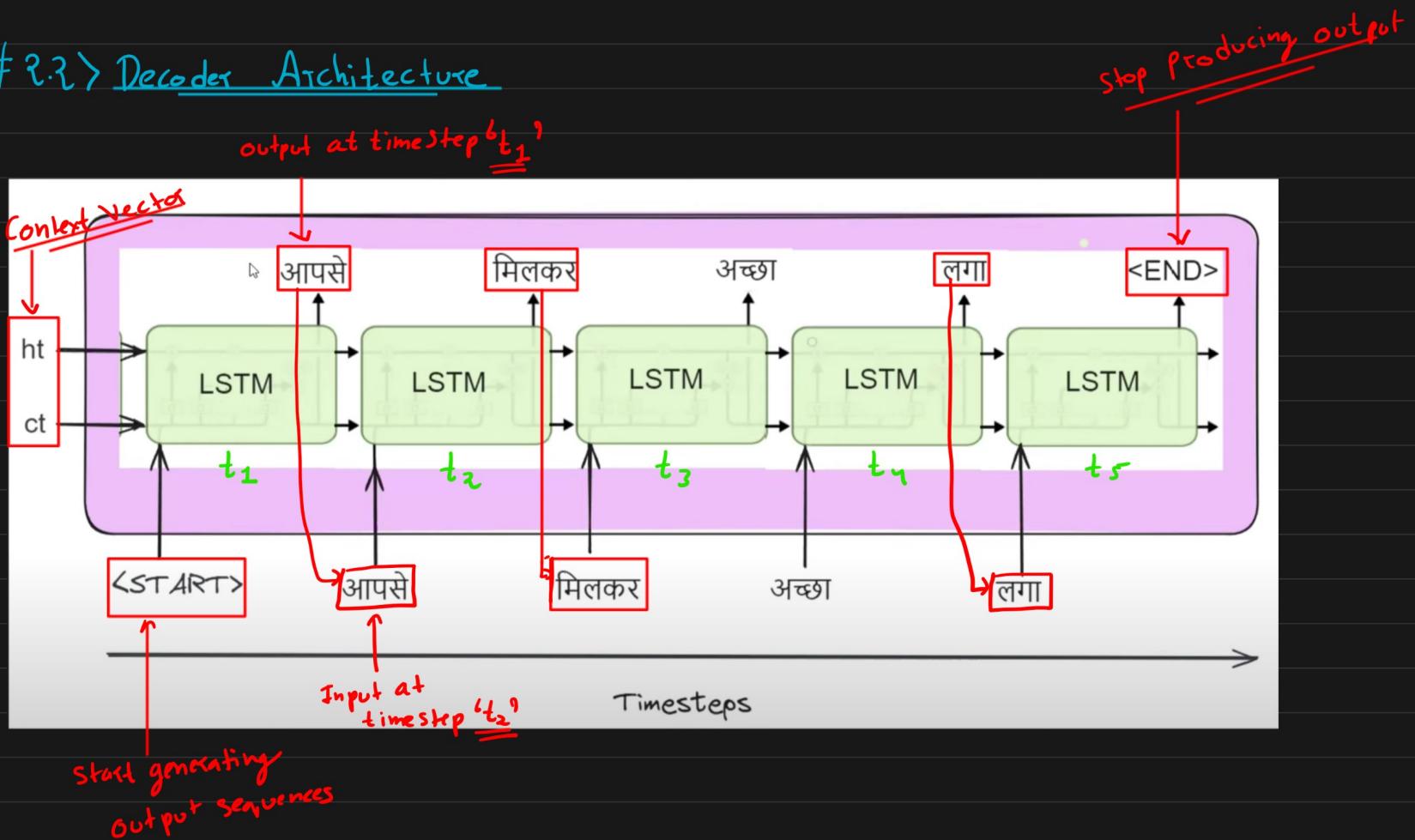
2 > High Level Overview : Encoder Decoder



#2.1 > Encoder Architecture



2.3 > Decoder Architecture



NOTES : i) **<START>** & **<END>** \Rightarrow Represent starting and ending of output sentence.

ii) The context vector i.e. [the final states generated by encoder] is pass as the initial states to decoder.

iii) The input at timestep t_2 is the output generated at previous timestep t_1



3 > Dataset Preprocessing

Task → Machine Translation

Dataset

↓

English	Hindi
Think about it.	सोच ला
Come in	अंदर आ जाओ

Step-1) Tokenize the sentence

↓
word

↓
Dataset

English	Hindi
[Think, about, it]	[सोच, ला]
[Come, in]	[अंदर, आ, जाओ]

word
↓
token

Building block of text

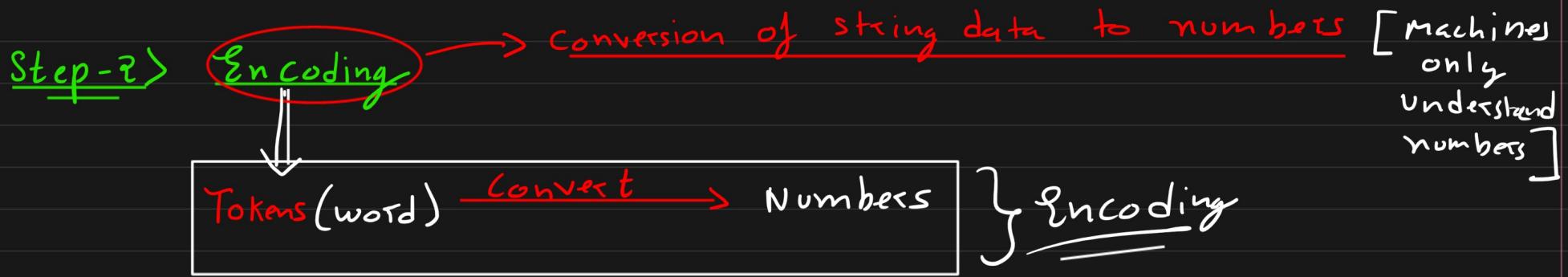
Body of Language



Vocabulary

unique tokens

English	Hindi
Think	सोच
about	ला
it	अंदर
Come	आ
in	जाओ

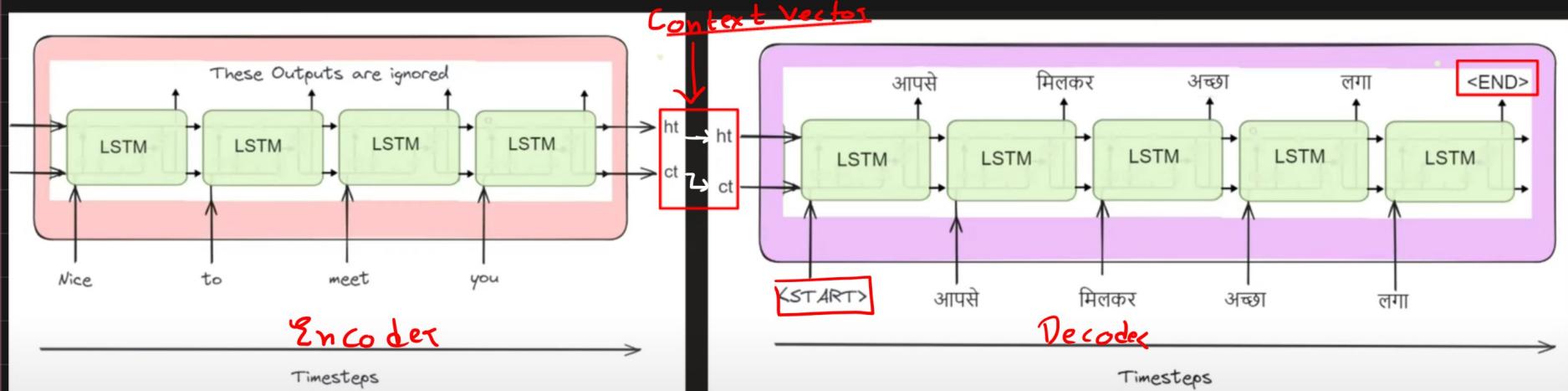


English	One-hot Encoding
① think	[1, 0, 0, 0, 0]
② about	[0, 1, 0, 0, 0]
③ it	[0, 0, 1, 0, 0]
④ come	[0, 0, 0, 1, 0]
⑤ in	[0, 0, 0, 0, 1]
[① ② ③ ④ ⑤]	

Hindi	One-hot Encoding
① <START>	[1, 0, 0, 0, 0, 0, 0]
② संग्रह	[0, 1, 0, 0, 0, 0, 0]
③ लै	[0, 0, 1, 0, 0, 0, 0]
④ अंदर	[0, 0, 0, 1, 0, 0, 0]
⑤ आ	[0, 0, 0, 0, 1, 0, 0]
⑥ बाइसी	[0, 0, 0, 0, 0, 1, 0]
⑦ <END>	[0, 0, 0, 0, 0, 0, 1]
[① ② ③ ④ ⑤ ⑥ ⑦]	

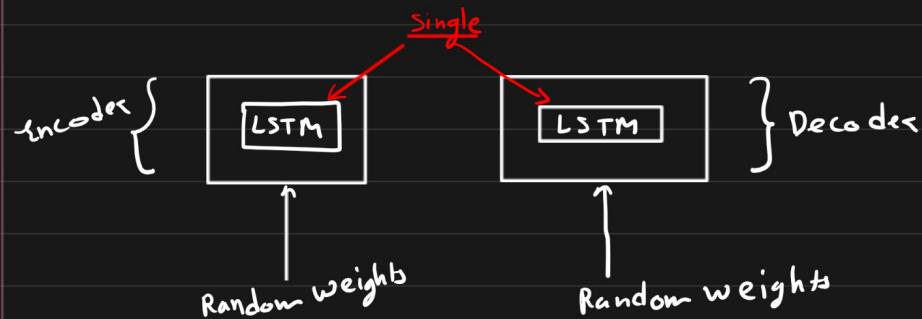
Model Training

→ The training of Encoder & decoder is done **simultaneously**.



Training Process

1> At initial stage the learning parameters (weights + Bias) of encoder & decoder are initialized **randomly**.

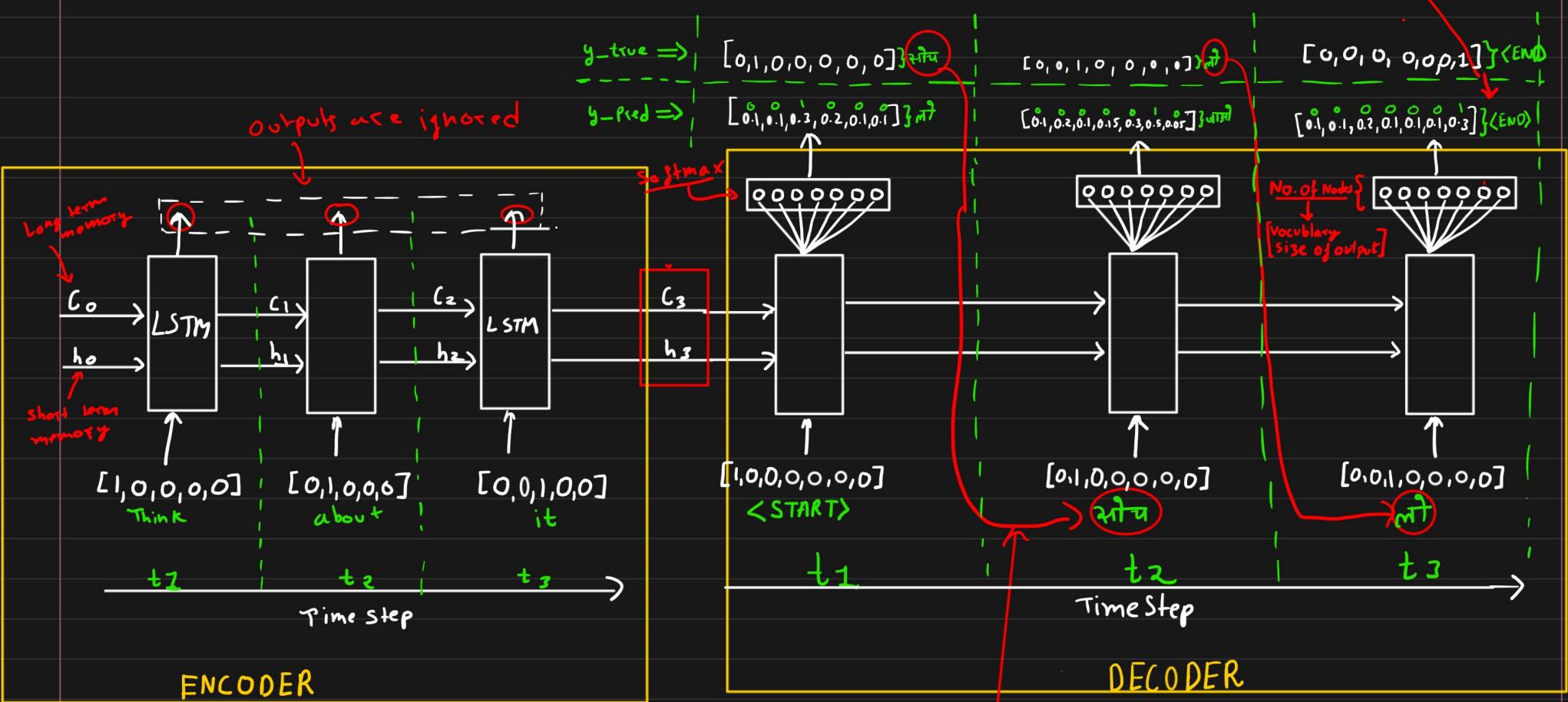


i) Forward Propogation

Sample-1) [Think about it] \rightarrow [सोच मत <END>]

features target

stop the output prediction
if "y-pred" \rightarrow <END>



Teacher Forcing \Rightarrow { Fast Training $\xrightarrow{\text{how}}$ Fast Convergence }

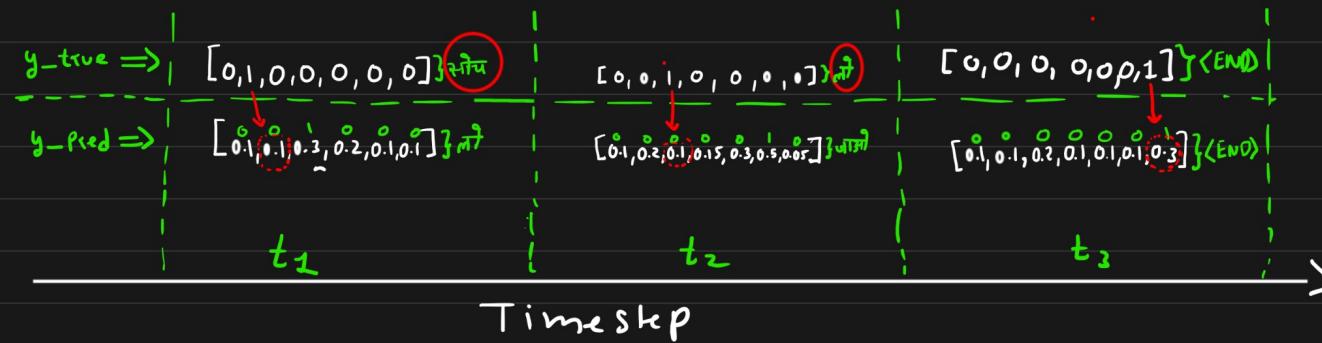
It means we will pass the correct output " y_{true} " at time step t_n as input to next timestep t_{n+1} , irrespective of prediction output " y_{pred} " at timestep t_n .

ii) Calculate Loss

$$\text{Loss function} \Rightarrow \text{Categorical Crossentropy} = - \sum_{i=1}^{\text{no. of classes}} y_i \log(\hat{y}_i)$$

↑ Predicted
↓ True

Since at each timestep we need to output one of the category [vocabulary of output]



Total loss $\rightarrow L = \sum_{i=1}^n L_{ti}$ Loss at each timestep

$$L_{t=1} = -1 \times \log(0.1) = -1 \times (-1) = 1$$

$$L_{t=2} = -1 \times \log(0.1) = -1 \times (-1) = 1$$

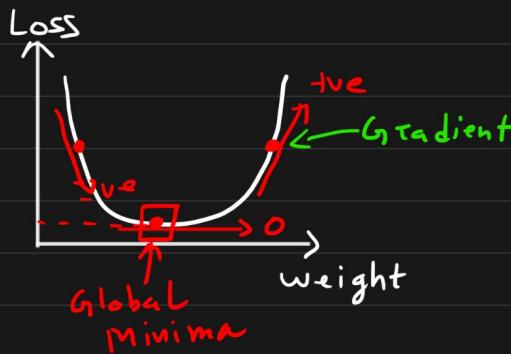
$$L_{t=3} = -1 \times \log(0.3) = -1 \times (-0.52) = 0.52 \leftarrow \text{Minimum Loss}$$

$$L_{\text{total}} = L_{t_1} + L_{t_2} + L_{t_3} = 1 + 1 + 0.52 = 2.52 \leftarrow \text{total Loss}$$

iii) Back Propagation

1) Gradient Calculation : $\Rightarrow \{$ Represent: How much each learning parameter (weight + bias) contribute to loss function $\}$

\Rightarrow Provide the direction of gradient i.e. [-ve, 0, +ve]



?> Weight Updation

for K iteration:

$$w_{i\text{new}} = w_{i\text{old}} - \eta \frac{\delta L}{\delta w_{i\text{old}}}$$

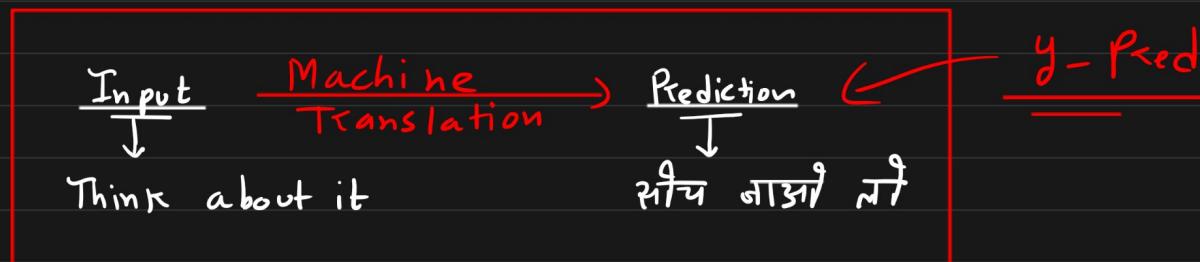
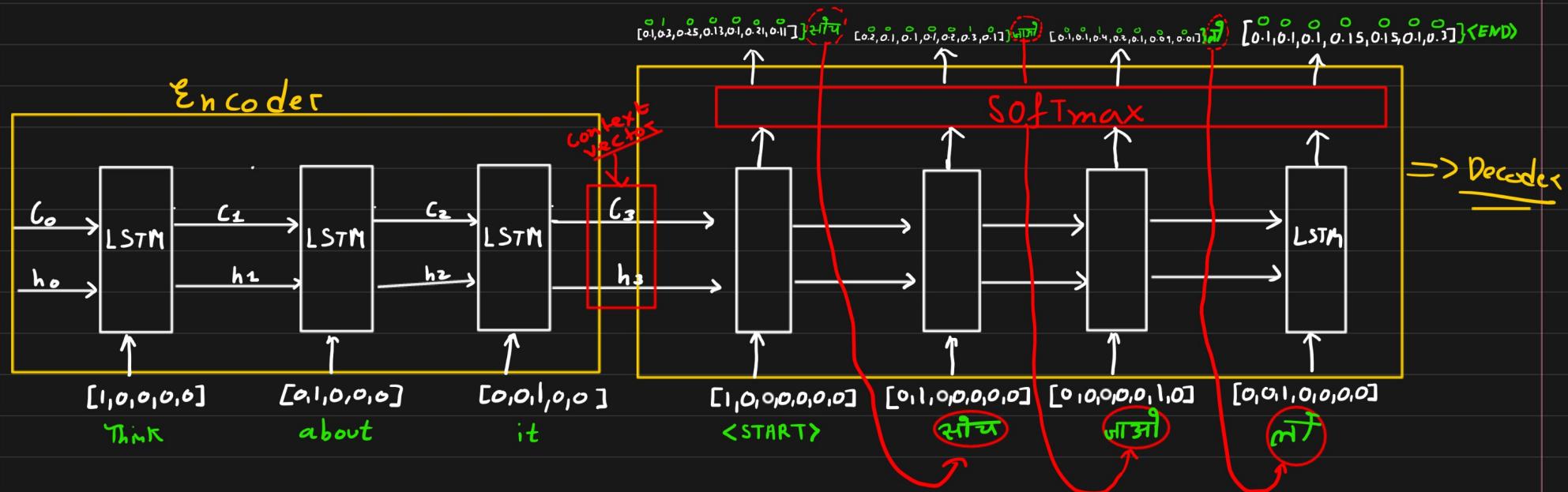
↑
Learning Rate

↑ Partial derivative
of 'Loss'
w.r.t w_i



iv) Prediction

F/P: Think about it



Improvement → Embedding

Embedding

→ Low Dimensional dense Representation.
Similar words have similar vectors

Example:

Consider you have English vocabulary size is 13K.

One-hot Encoding 3-dimensional Embedding

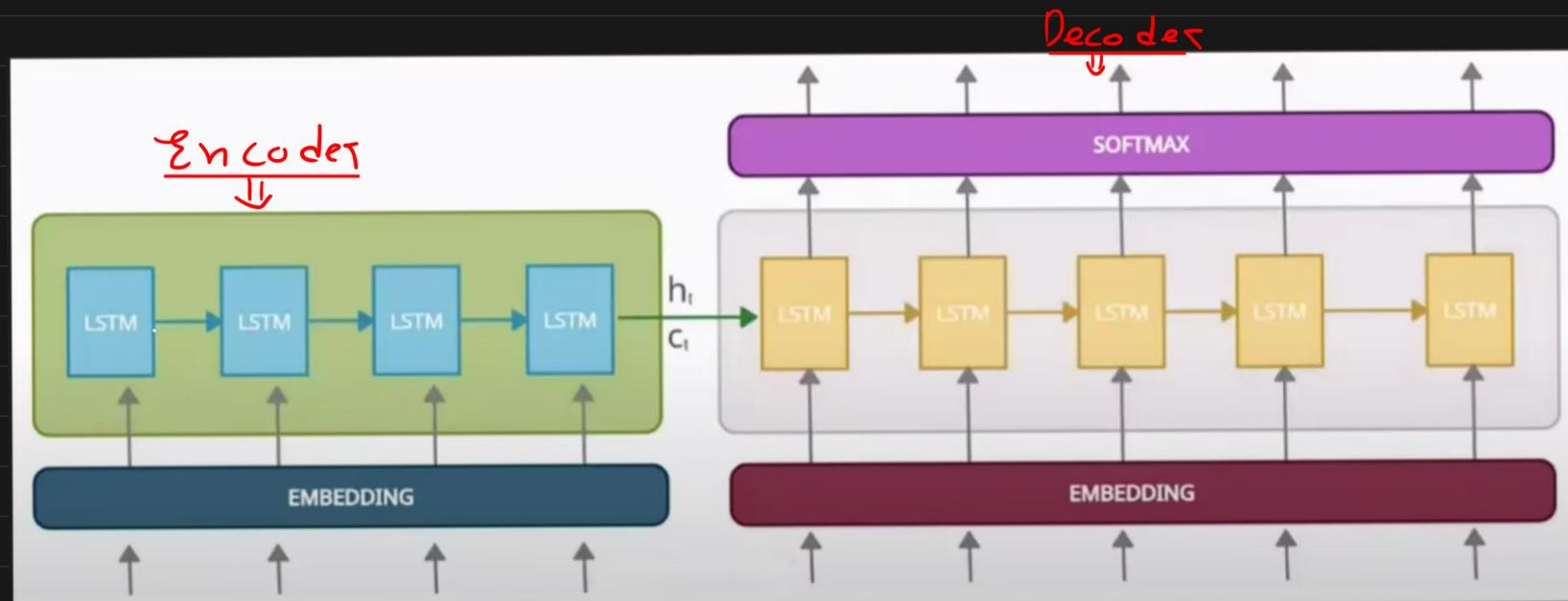
↓ ↓

Think → [1, 0, 0, ..., 0] [2.3, 0.4, 1.8]
about → [0, 1, 0, ..., 0] [2.1, 3.6, 1.2]
it → [0, 0, 1, ..., 0] [0.3, 0.9, 1.1]

1x13K 1x3

Sparse vector Dense vector

Shape Shape



2) Improvement → Deep LSTM

