

Encoder Decoder

Seq2Seq Data:

A kind of data in which input & output both will be in sequence.

Example

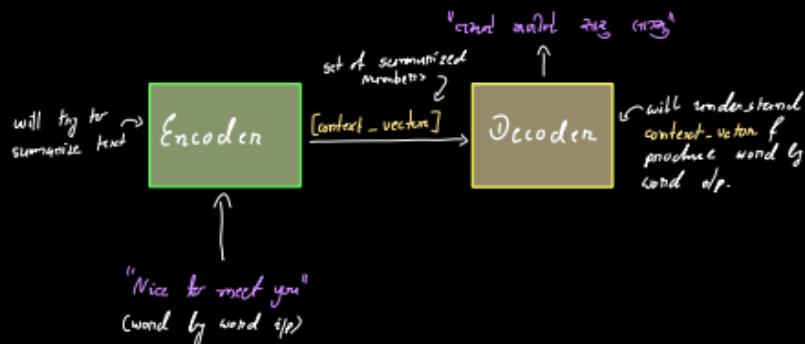
Machine Translation

Challenges

input → Sentence → Variable length

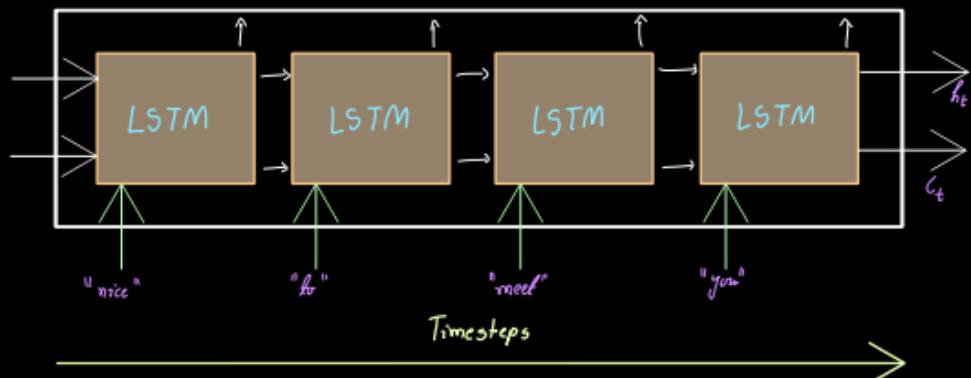
output → Sentence → Variable length

Basic architecture



Encoder Architecture

These np are ignored



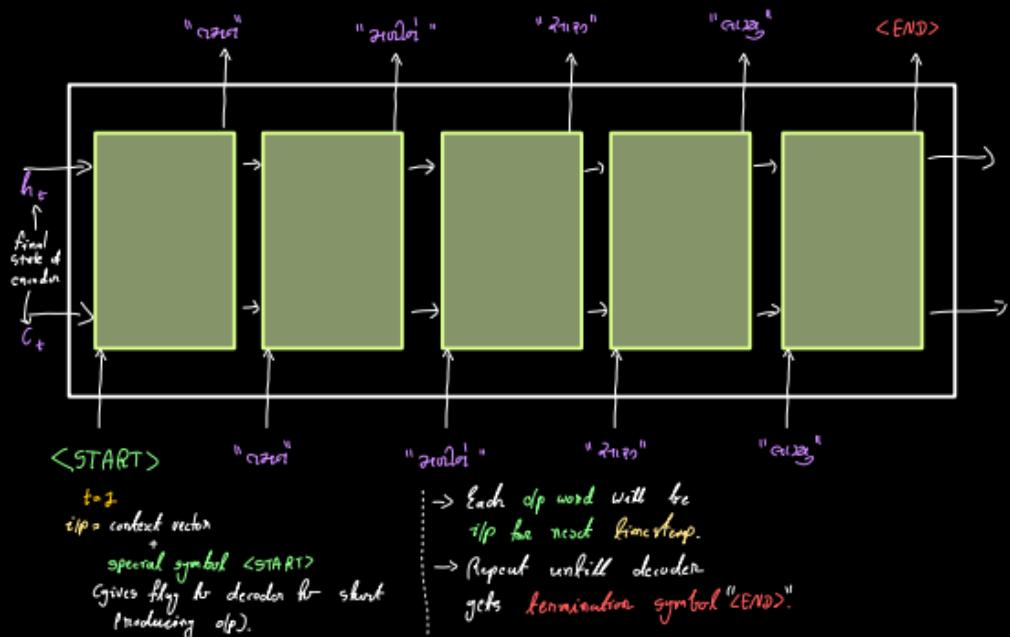
→ Encoder is just a single LSTM cell which we have unfolded in timesteps.

→ h_t and c_t at np side are final representation of sentence (context-vector).

→ Instead of LSTM cell, we can also use GRU. we don't use RNN due to vanishing gradient problem.

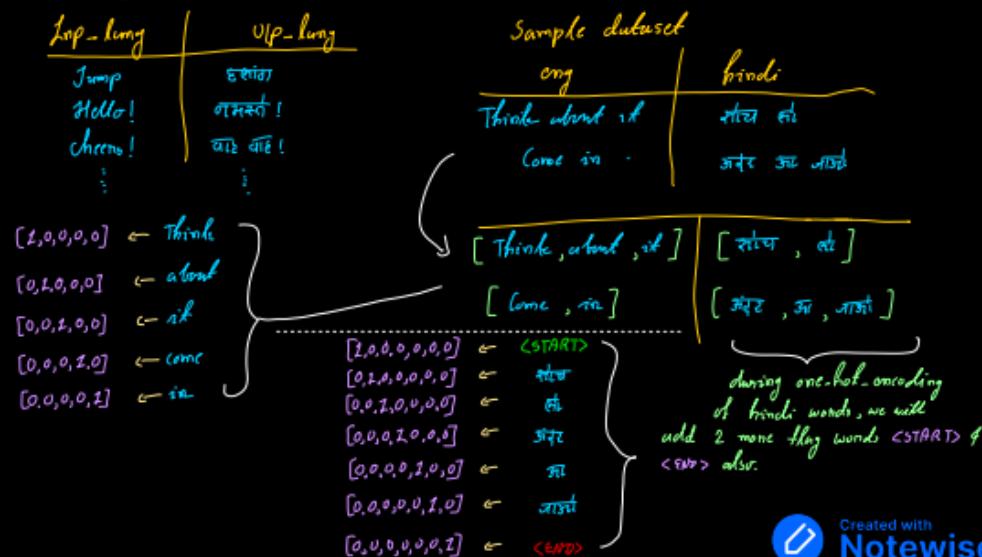


Decoder Architecture



Training the Architecture using backpropagation

→ If we're trying to solve machine translation kind of problems, then we need to use supervised dataset, which looks like this.

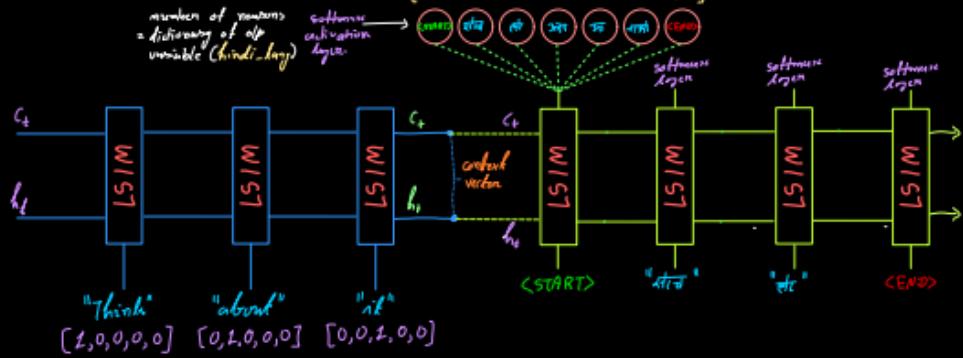


Starting training process

Row-2 : [Think about it] \rightarrow [शब्द, कर]

→ Initially weights & bias will be assigned randomly

$y_tanh[0 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0]$ each neuron will return probability for each word of input
 $y_prod[0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0]$ probability corresponding to word will be final output word.



→ As we know that this architecture can't give any blank word as op. so it will return a vector as op. → So that we can't pass op a current timestep on next timestep.

→ As we have seen example above that we can't pass up a timestamp because wrong prediction at first time step will lead wrong step at next timestamp also.

→ Due to this problem instead of sending step of current timestep on next, we will panic connect step at related timestep intentionally. this process is called as Teacher forcing.

→ Teacher forcing will be apply during training, on testing current predicted word will be zip for next timestamp.

→ Due to teacher forcing method, training will be faster

→ Till now we have completed forward propagation, now let's calculate loss & move for backward propagation

y -true	$[0, 1, 0, 0, 0, 0]$	$[0, 0, 1, 0, 0, 0]$	$[0, 0, 0, 0, 0, 1]$
y -pred	$[0.1, 0.2, 0.3, 0.2, 0.1, 0.3, 0.2]$	$[0.1, 0.2, 0.1, 0.15, 0.3, 0.5, 0.05]$	$[0.1, 0.2, 0.3, 0.1, 0.2, 0.1, 0.4]$

loss function: categorical-crossentropy
(multiclass classification)

$$L = - \sum_{i=1}^n y_i^{text} \log(y_i^{pred})$$

$$\text{average loss} = 0.7981$$

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

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

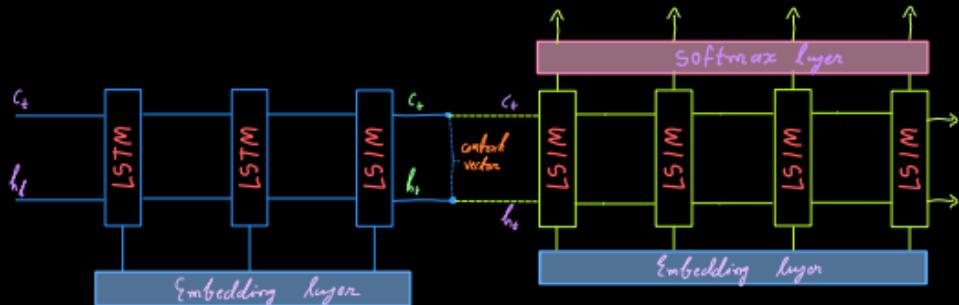
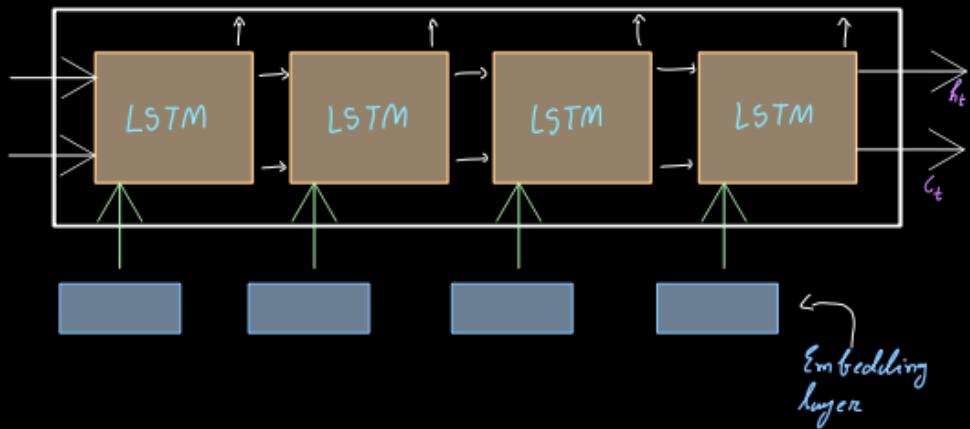
Created with Notev

Back propagation

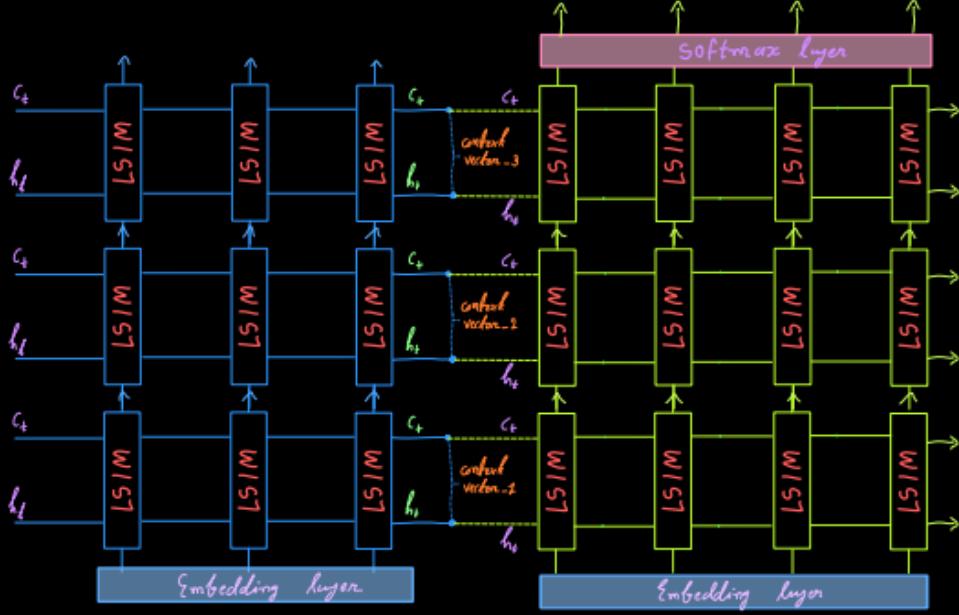
- 1 Gradient calculation
- 2 Update weights



Improvement - 1 (Embedding)



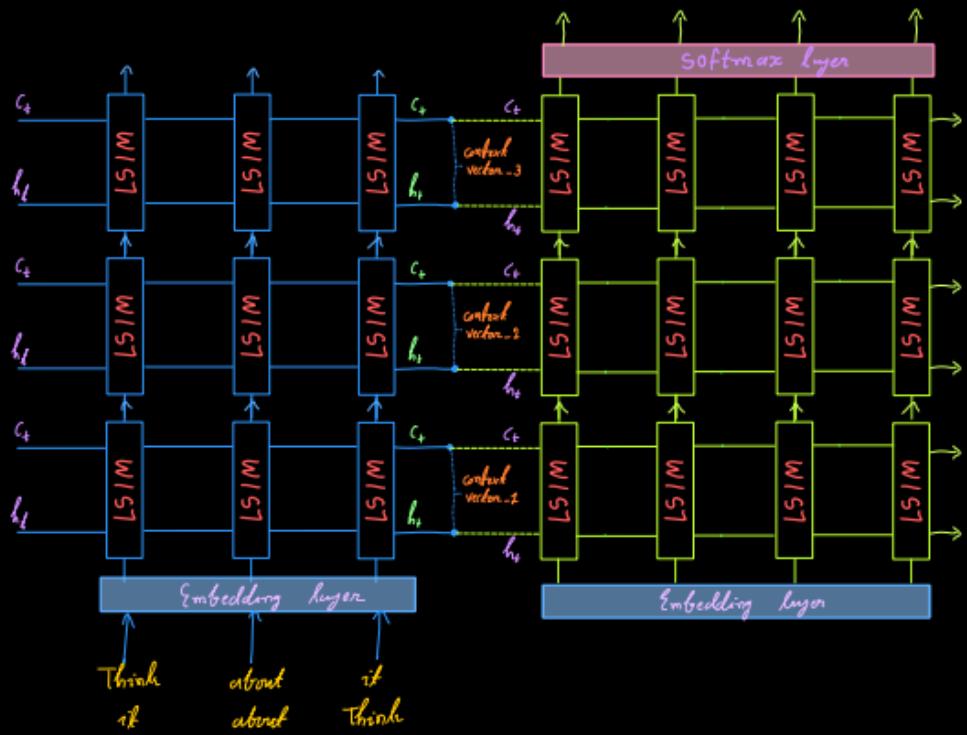
Improvement - 2 (Deep LSTMs)



Benefits

- Can easily handle long-term dependencies.
- layered representation easy to understand hierarchical representation.
- Improves learning capability

Improvement -3 (Reversing the input)

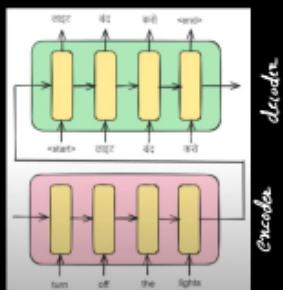


- The input processing word technique won't work in all cases.
- This technique works in such case where initial words contains more context.



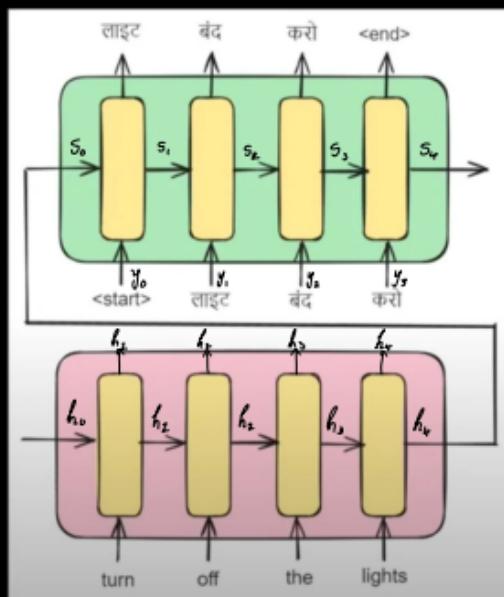
Attention Mechanism

why...?



- Suppose we have a paragraph with almost 50 words. which adds more pressure on encoder for summarize entire paragraph in few numbers which is very hard to manage context.
 - At any timestamp of decoder, we don't need entire sentence, we just need only few words/numbers from encoder to translate at given timestamp.
 - look in this encoder-decoder architecture, entire state in the form of context vector will be passed at each time step and creates a pressure for decoder to generate accurate prediction.

- If we talk about how humans memorize & translate sentences. we won't translate entire paragraph at once. instead we will translate it sentence by sentence by creating an attention window to give more/less attention on each part of sentence rather than giving equal importance to entire paragraph.
 - As we move forward, region around region window gets visible & the translation process will occur based on visibility.

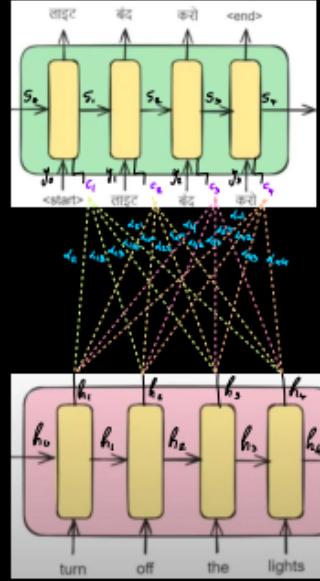


- At each time step of decoder, we also need to pass the most important timestamp of encoder at current timestamp.
 - At $t+2$ of decoder, t_{i+2} of encoder is important.
At t_i of decoder, t_{i+2} and t_i of encoder is important.
 - At any timestamp t_i , it will be $[s_i, y_i]$
Additionally we will also pass one more piece of information, which is most useful timestamp of encoder at current decoder's timestamp.
 - Now the input at current timestamp t_i , it will be $[s_i, y_i, c_n]$
 - So at particular time t_i , it will be $[s_{i-2}, y_{i-2}, c_n]$
 - \hookrightarrow attention step
 - Here c_n will be either vector, scalar or matrix. depends on y_i .

→ The dimension of C_n will equal to R_n . In case of multiple context vectors R_n , we will do weighted sum of them.

Q. How to identify most relevant hidden state of encoder for decoder at any given timestamp....?

→ Suppose at $t=1$ of decoder we want most relevant hidden state that we will assign weights α to the o/p at each timestep of encoder & will perform element-wise operation like---



$$c_2 = \alpha_{02} h_2 + \alpha_{12} h_2 + \alpha_{22} h_2 + \alpha_{32} h_2$$

where:
 α = scalar weights
 h = vectors

→ This process will be repeated for c_3 , c_4 and c_5 also.

→ While calculating c_2 , whichever weight has higher value, will be transferred for further use.

→ Now suppose if decoder timestep is denoted by i and encoder timestep is denoted by j .

$$c_i = \sum_{j=1}^n \alpha_{ij} h_j$$

→ for each h_j we need to create different α_{ij} so the total number of α s will be $i \times j = 26$.

↳ calculation

→ α is called as alignment/similarity score.

→ Example

for g_0 calculation via c_0 , hence let's suppose we pick α_{21} , which depends on h_2 and s_2 .
(previous hidden-state of decoder)

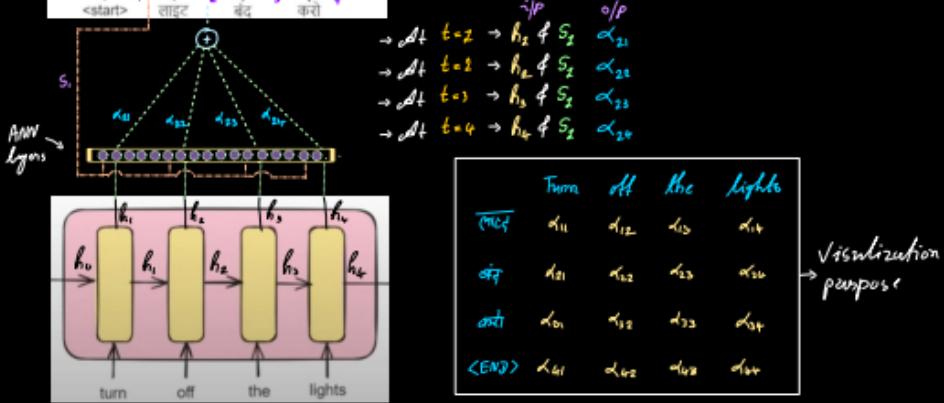
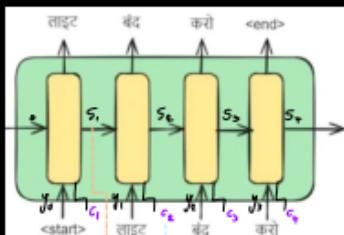
→ S_0 , using attention mechanism, we wants to find that, "based on previous completed state S_1 , tell me now, how much h_2 is important to calculate current o/p at $t=2$ timestep."

$$\alpha_{21} \rightarrow f(h_2, s_2) \quad \alpha_{ij} \rightarrow f(h_j, s_{i-1})$$

$\alpha_{22} \rightarrow f(h_2, s_2)$

→ this mathematical function f is approximation of ANN.

→ So the architecture may looks like this.



- Researchers have used Bi-directional LSTM only in encoder.
- Now let's see 2 different type of α calculation technique.
 1. Bahdanau Attention
 2. Luong Attention

#2 Bahdanau Attention

Q. What is α ...?

A. α is alignment score.

→ It shows that at any particular timestamp of decoder, what is the importance of every hidden state of encoder.

$$\begin{aligned} d_{11} &\longrightarrow \text{turn} & \alpha_{11} &= f(h_1, s_0) \\ d_{12} &\longrightarrow \text{off} & \alpha_{21} &= f(h_1, s_1) \\ d_{31} &\longrightarrow \text{the} & \alpha_{31} &= f(h_1, s_2) \end{aligned}$$

d says that, based on till now decoded data, α of current timestamp. how much it depends on each hidden state of encoder.

s_i stores context of decoded data till current timestamp.

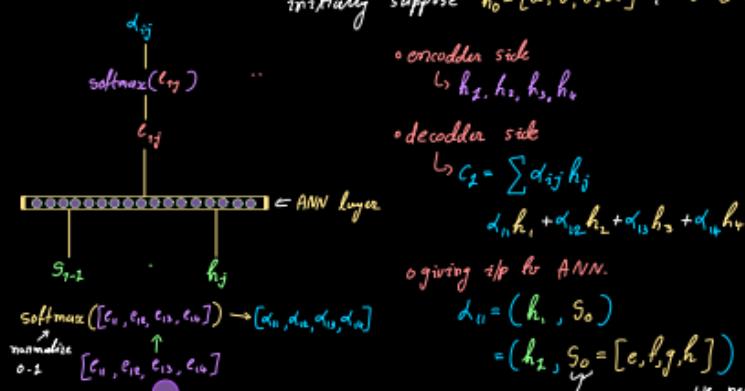
$$\alpha_{ij} = f(h_j, s_{i-1}) \rightarrow \text{alignment score will be a } \left\{ \begin{array}{l} \text{hidden state of} \\ \text{encoder} \\ \text{mathematical function of 2} \\ \text{things. softmax, exp, sigmoid etc...} \end{array} \right\}$$



→ Here, selection of mathematical function $\text{Bidiagonal Attention}$ will help. it says that we can approximate this function using Feed Forward Neural-Network.

→ Because Feed Forward Neural-Networks are Universal Function Approximators. if we provide enough data to ANNs.

initially suppose $h_0 = [a, b, c, d]$ & $s_0 = [e, f, g, h]$



encoder side

$$h_1, h_2, h_3, h_4$$

decoder side

$$c_2 = \sum c_{ij} h_j$$

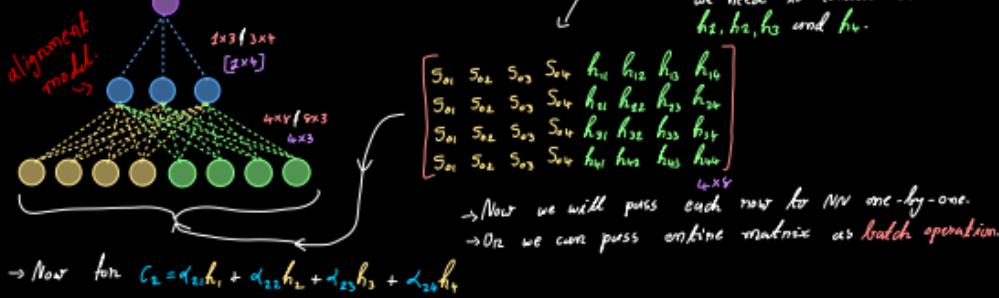
$$d_{11} h_1 + d_{12} h_2 + d_{13} h_3 + d_{14} h_4$$

giving zip for ANN.

$$d_{11} = (h_{11}, s_0)$$

$$= (h_{11}, [e, f, g, h])$$

we need to concat s_0 with h_2, h_3, h_4 and h_{11} .



→ Now we will pass each row to NN one-by-one.

→ Or we can pass entire matrix as batch operation.

$$\rightarrow \text{Now for } C_2 = d_{21} h_1 + d_{22} h_2 + d_{23} h_3 + d_{24} h_4$$

$$\left[\begin{array}{cccc} S_{11} & S_{12} & S_{13} & S_{14} \\ S_{11} & S_{12} & S_{13} & S_{14} \\ S_{11} & S_{12} & S_{13} & S_{14} \\ S_{11} & S_{12} & S_{13} & S_{14} \end{array} \right] \left[\begin{array}{c} h_1 \\ h_2 \\ h_3 \\ h_4 \end{array} \right] \left[\begin{array}{c} d_{21} \\ d_{22} \\ d_{23} \\ d_{24} \end{array} \right] \rightarrow \text{pass for ANN} \rightarrow [d_{21}, d_{22}, d_{23}, d_{24}] \rightarrow C_2$$

$$C_2, S_2, Y_2 \rightarrow \text{LSTM} \rightarrow Y_2 \rightarrow S_t$$

→ At each timestamp we will get different d values.

→ In decoder, this FNN will be same at each timestamp. hence it is called as

Time distributed Fully-connected Neural Network.

→ The parameters of this network will be updated during backpropagation.

$$C_{ij} = \sum d_{ij} h_j$$

$$\left[\text{softmax}(c_{ij}) \right]$$

$$\left[\tanh(W \cdot [S_{t-1}; h_j] + b) \cdot V \right]$$

weights of zip layer of first hidden layer.

→ So, above shown mathematical process is called as Bahdanau Attention.
also called as additive Attention.

#2 Luong Attention / Multiplicative Attention

→ In this technique, we will calculate attention score.

→ Process will be same as Bahdanau attention. but the α calculation process is slightly different.

$$d_{ij} = f(h_j, s_i)$$

↳ We're not dependent on previous hidden state.

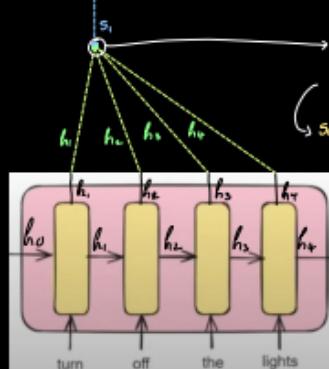
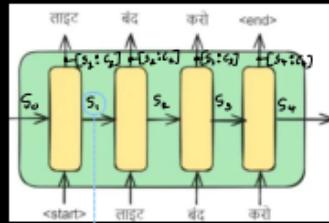
here, function f does the dot product of h_j and s_i .

$$s_i = [a, b, c, d] \quad h_j = [e, f, g, h] \quad s_i^T \cdot h_j = \begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix} \cdot [e, f, g, h] = ae + bf + cg + dh$$

scalar / attention value

$$\underbrace{e_{ij}}_{\text{softmax}(e_{ij})} \quad \underbrace{\alpha_{ij}}$$

→ here, we're using s_i instead of s_{i-1} because, by using s_i we will get more updated hidden state.



→ This process will continue for s_2 , s_3 and s_4 .

$$[s_i h_1, s_i h_2, s_i h_3, s_i h_4]$$

$$[e_{i1}, e_{i2}, e_{i3}, e_{i4}]$$

$$\text{softmax}([e_{i1}, e_{i2}, e_{i3}, e_{i4}])$$

$$[\alpha_{i1}, \alpha_{i2}, \alpha_{i3}, \alpha_{i4}]$$

$$c_L$$

(This c_L will be concatenated with s_2)

$$[s_2 : c_L]$$

$$\tilde{s}_2 \rightarrow \text{FNN} \rightarrow \text{softmax}$$

