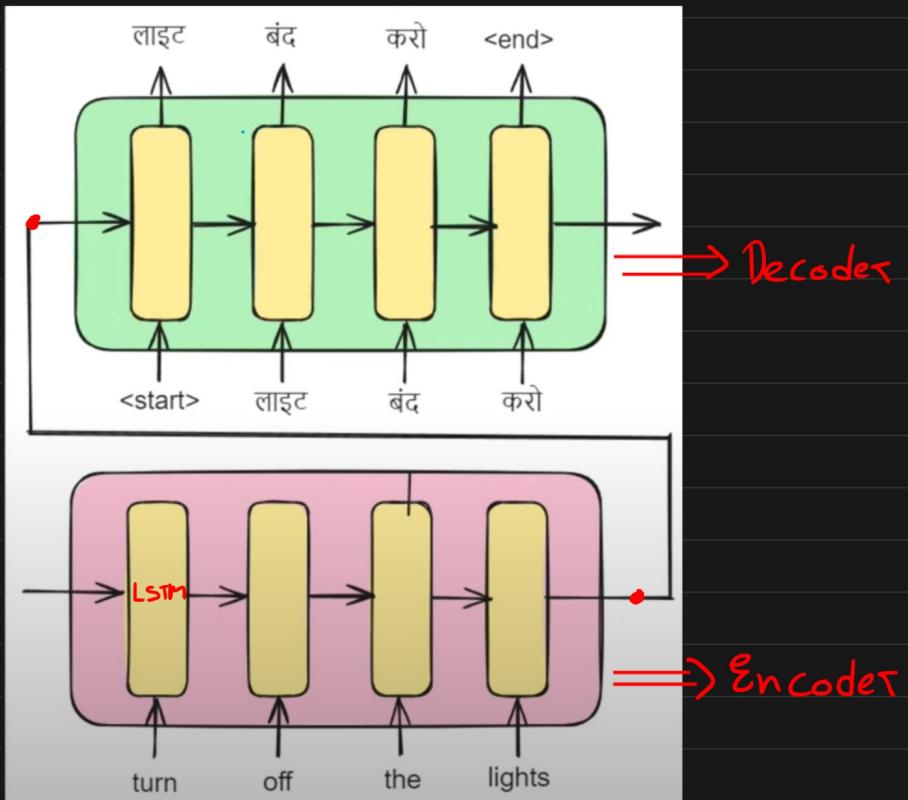


Attention Mechanism

Disadvantage of Encoder - Decoder

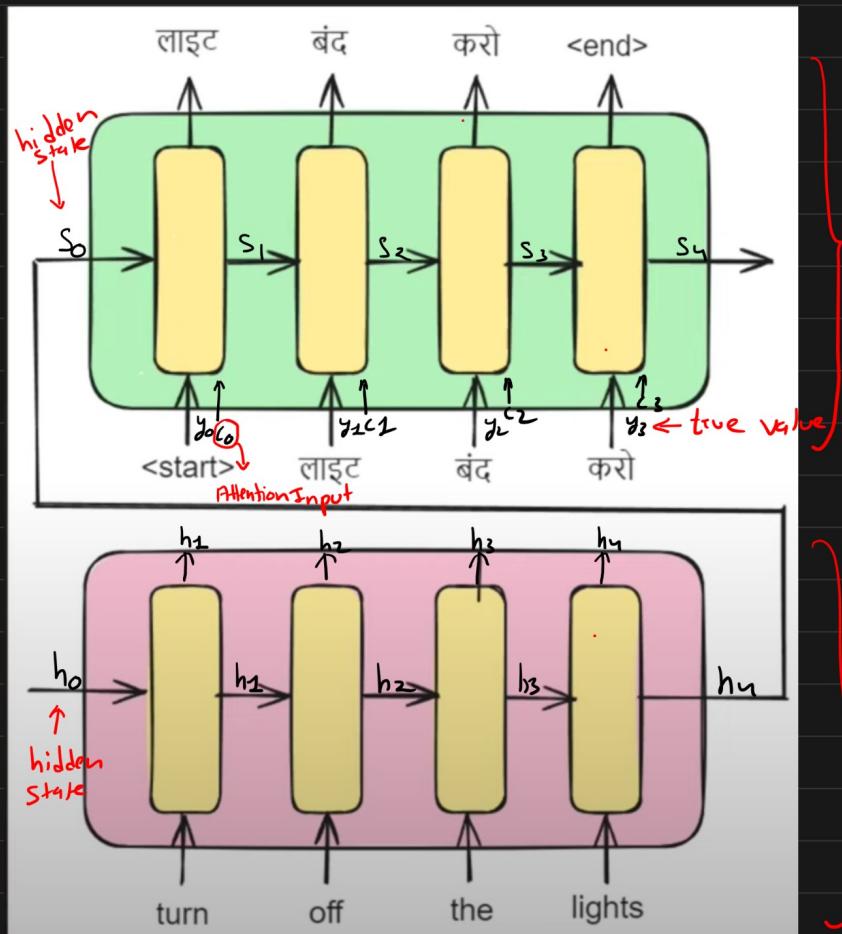


1) In Encoder-Decoder architecture the encoder process the whole sentence and it passes **content vector** to the decoder i.e. unnecessary we are giving huge load to encoder the whole sentence that is only represented by **context vector**.

2) Unnecessary we are giving all the sentences to decoder at every timestep i.e. at every timestep decoder only required set of words for decoding not the whole sentence.

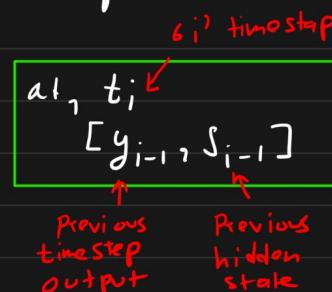
Attention Mechanism

⇒ Dynamically we need to generate the information, at every timestep of decoder^{that} which timestep or set of timestep of encoder, are important CURRENTLY for it.



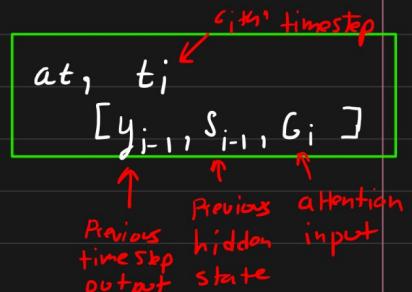
Vanilla Encoders-Decoders

- 1) There are only 3 inputs provided to decoder at every timestep.



Attention Encoders-Decoders

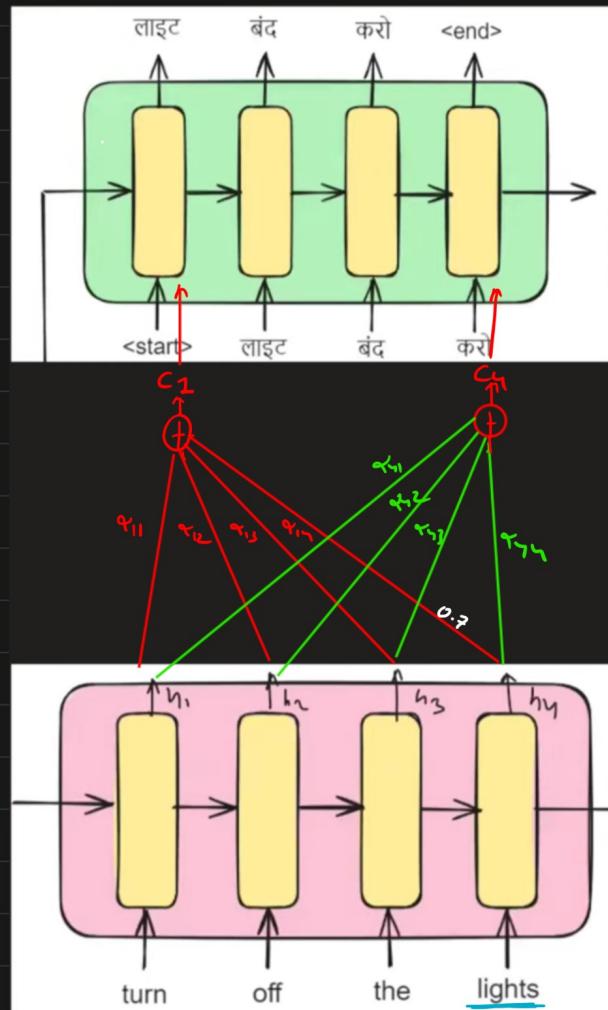
- 1) There are 3 inputs provided to decoder at every timestep.



Attention Input [c_i]

$c_i \Rightarrow$ Attention input at i^{th} timestep of decoder. i.e. Weighted sum of the hidden states of encoder.

$$\text{i}^{\text{th}} \text{ timestep attention input} \rightarrow c_i = \sum \alpha_{ij} \times h_j \quad \begin{matrix} & j^{\text{th}} \text{ timestep hidden input (vector)} \\ \text{vector} & \uparrow \\ \text{weight (scalar)} & \end{matrix}$$



Example

$$c_1 = \alpha_{11} h_1 + \alpha_{12} h_2 + \alpha_{13} h_3 + \alpha_{14} h_4$$

$$c_4 = \alpha_{41} h_1 + \alpha_{42} h_2 + \alpha_{43} h_3 + \alpha_{44} h_4$$

$$\boxed{\text{Total weights} = i \times j} \Rightarrow 4 \times 4 = 16 \leftarrow \text{No. of weights}$$

(Q) What are weights " α "?

Ans:

$\alpha_{ij} \rightarrow$ Similarity score

It tells for decoder i^{th} timestep **Output** how much contribution does encoder j^{th} timestep has.

Similarity Score [α_{ij}]

$$\alpha_{ij} = f(h_j, s_{i-1})$$

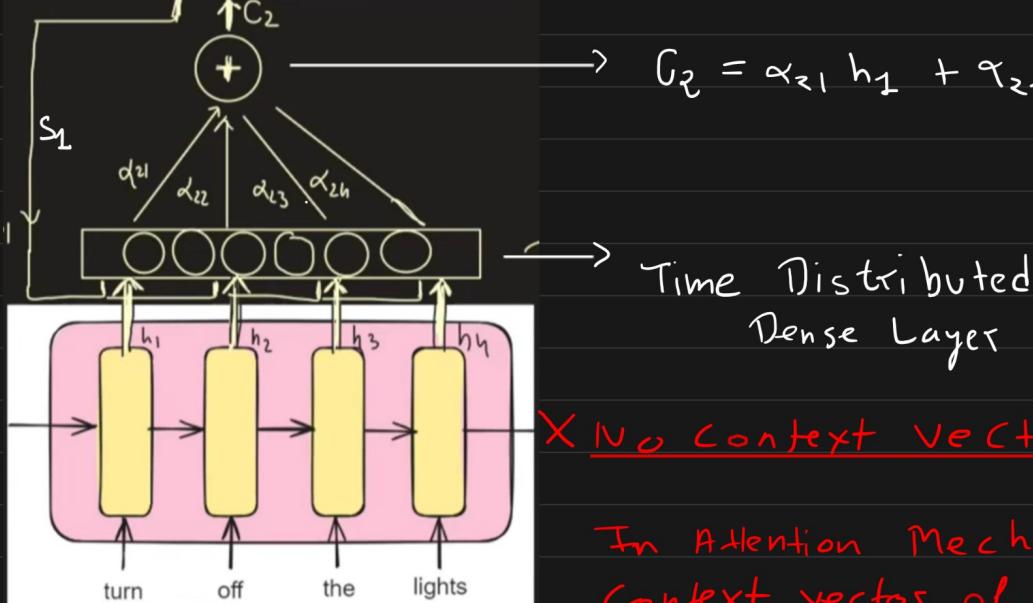
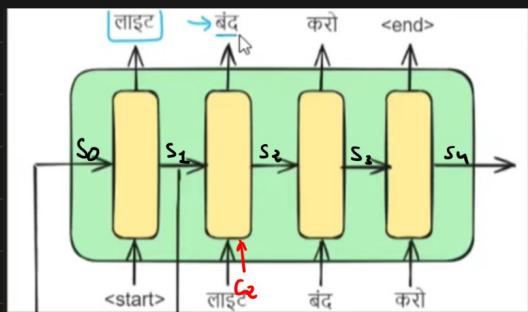
function

lstm hidden state

Previous hidden state of decoder

Similarity score

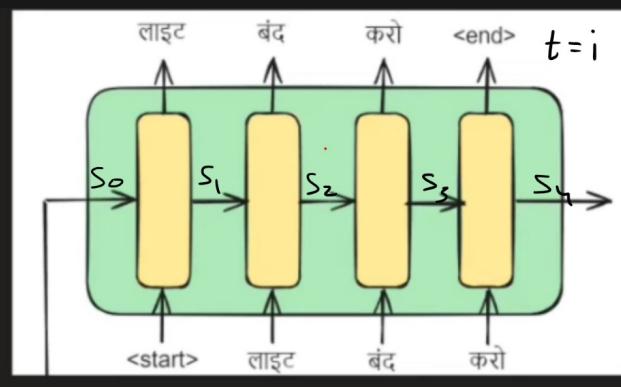
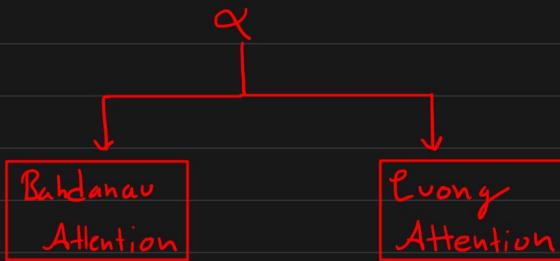
ANN → Universal function Approximation



In Attention Mechanism we do not pass context vector of last time step to decoder

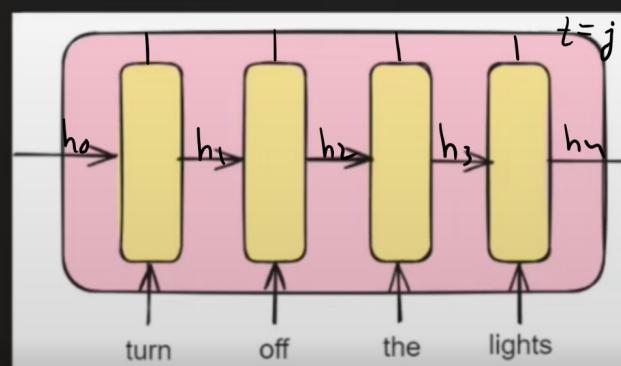
How to find α

⇒ In order to find Similarity score of word to word we can use :



$$C_{ij} = \sum_i \alpha_{ij} h_j$$

Similarity Score



$$\alpha = f(h_j, s_{i-1}) \Rightarrow \text{math function} \begin{pmatrix} \text{Encoder hidden state} & \text{Decoder Previous hidden state} \end{pmatrix}$$

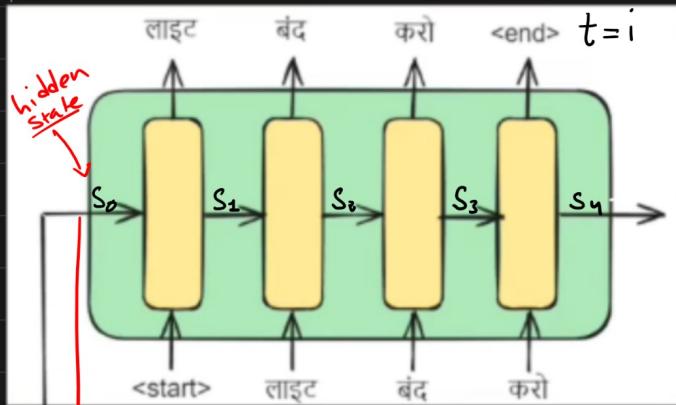
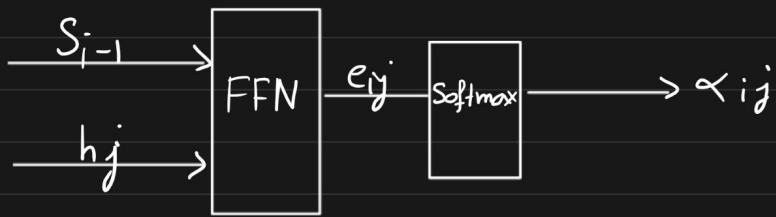
function Approximation

ANN

Bahdanau
Attention

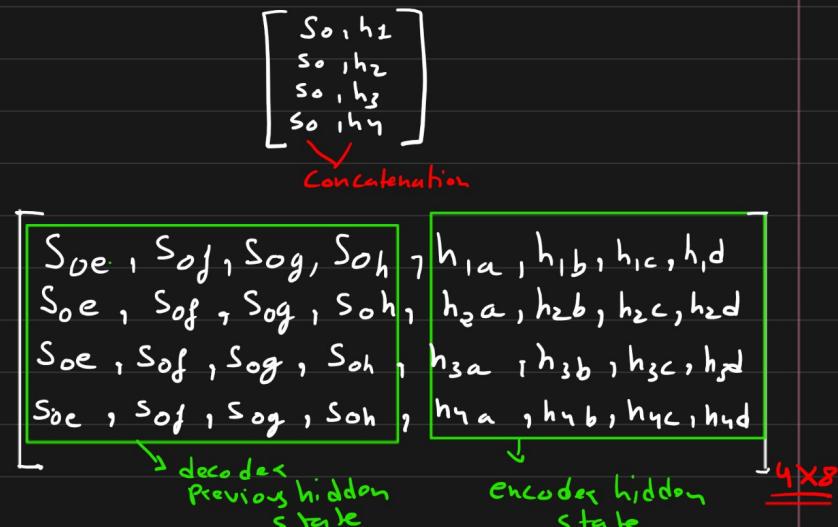
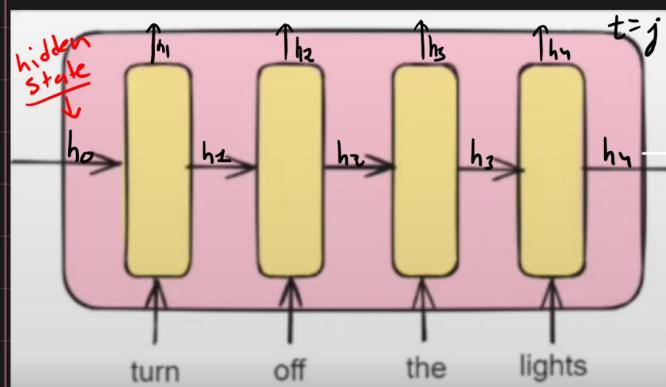
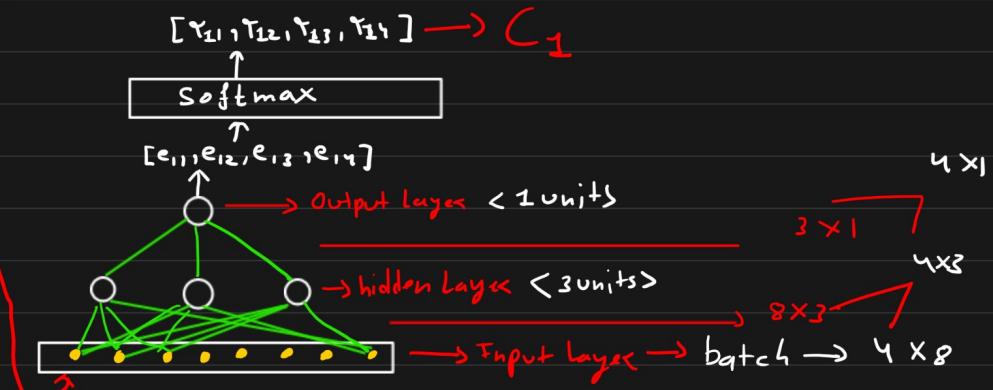
Luong
Attention

Bahdanau Attention



$$\text{hidden state} = S_i = [e, f, g, h]_{4 \times 1}$$

$$C_i = \sum \alpha_{ij} h_j \Rightarrow \alpha_{ij} = f\left(\frac{h_j, S_{i-1}}{\| \cdot \|}\right)$$



$$\text{decoder previous hidden state} \rightarrow h_j = [a, b, c, d]_{4 \times 1}$$

Bahdanau Attention Mathematics

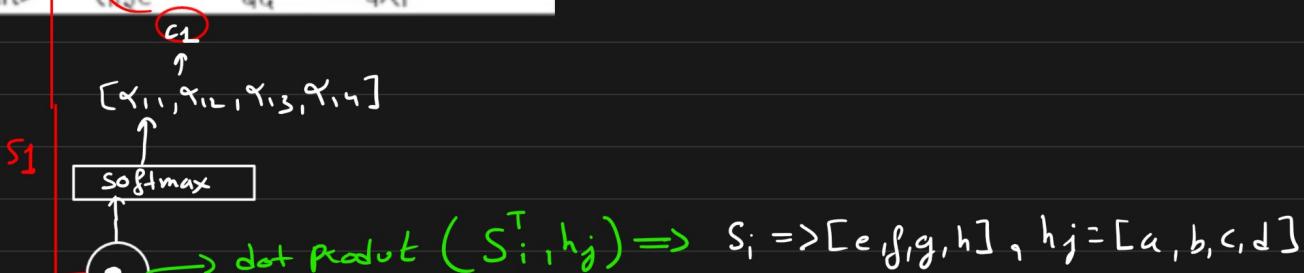
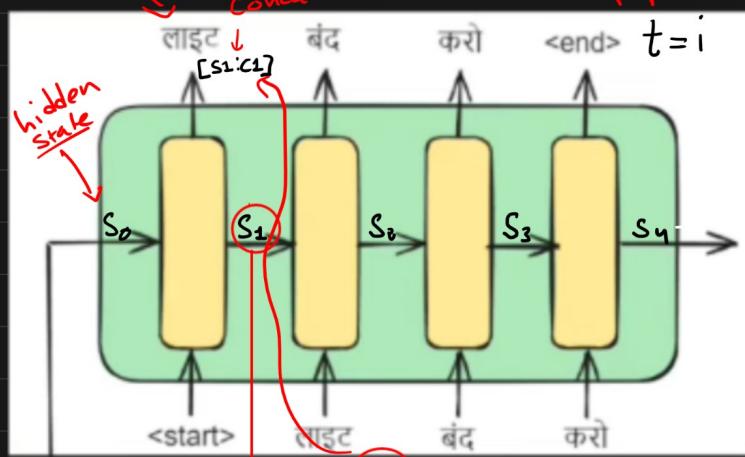
$$C_{ij} = \sum \alpha_{ij} h_j \quad \text{①} \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{Attention Input}$$
$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^n \exp(e_{ik})} \quad \text{②} \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{Softmax}$$
$$e_{ij} = \text{softmax}\left(V \cdot \sigma\left(W \cdot [S_{i-1}; h_j]\right)\right) \quad \text{③}$$

Activation function
weight weight Concatenation

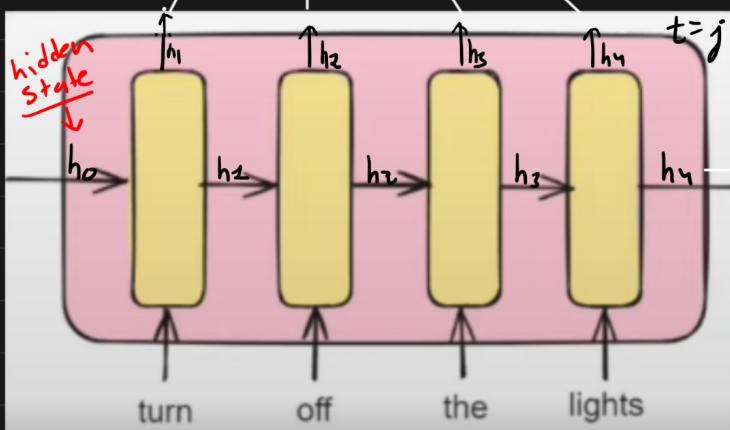
Back Tracking

Long Attention

② difference



$$s_i^T \cdot h_j = \begin{bmatrix} e \\ f \\ g \\ h \end{bmatrix}_{1 \times 4} \cdot [a, b, c, d]_{4 \times 1} = a_{ij} \rightarrow \text{scalar}$$



Mathematical function

$$c_{ij} = \sum a_{ij} h_i$$

$$a_{ij} = f(s_i, h_j) \rightarrow (\cdot) \Rightarrow s_i^T \cdot h_j$$

Latest hidden state

dot product

① difference

Advantage

- i) It is fast as we use (\cdot) → dot product for mathematical function.
- ii) We use latest information (context vector) i.e. $\underline{s_i}$ instead of previous hidden state $\underline{h_{i-1}}$.