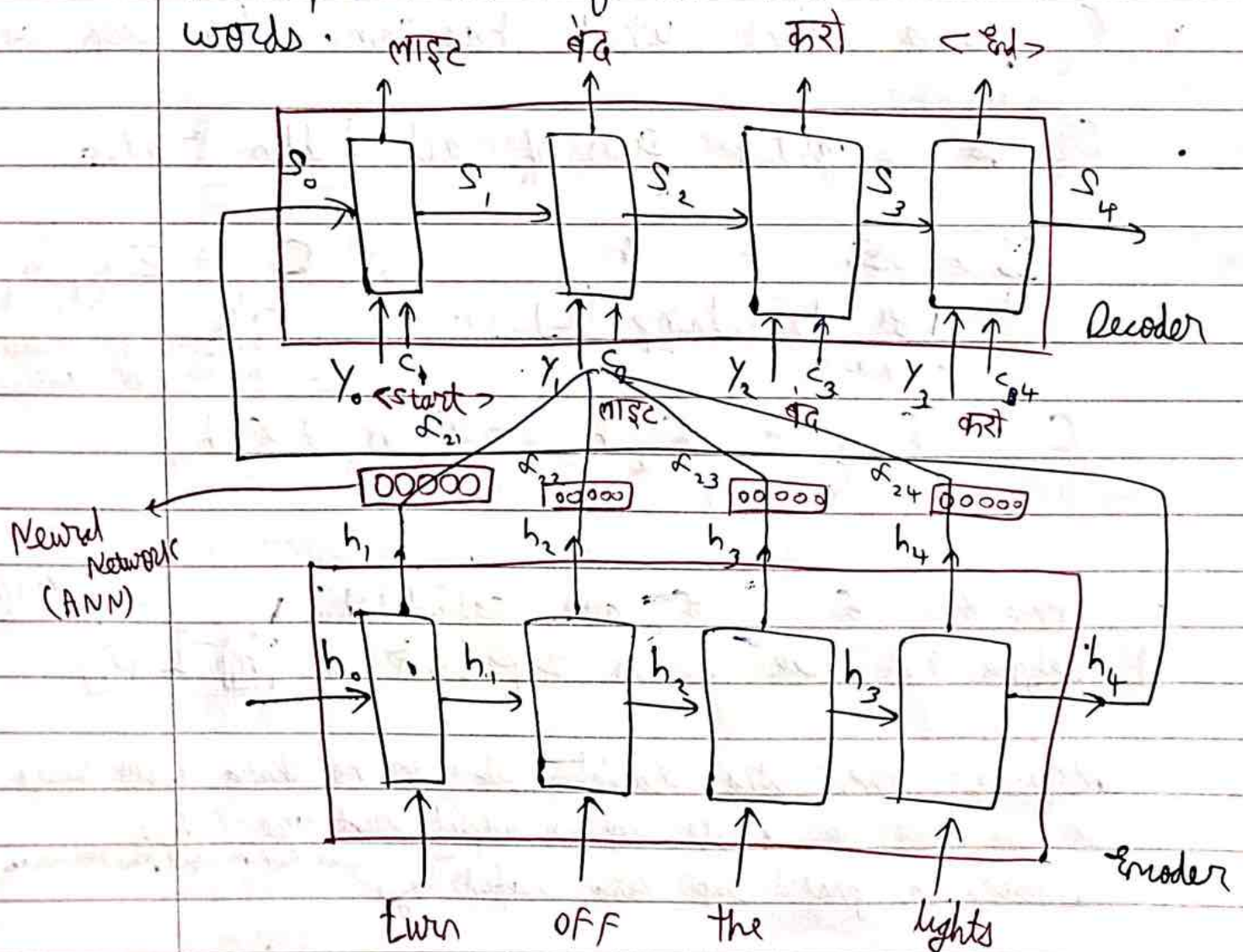
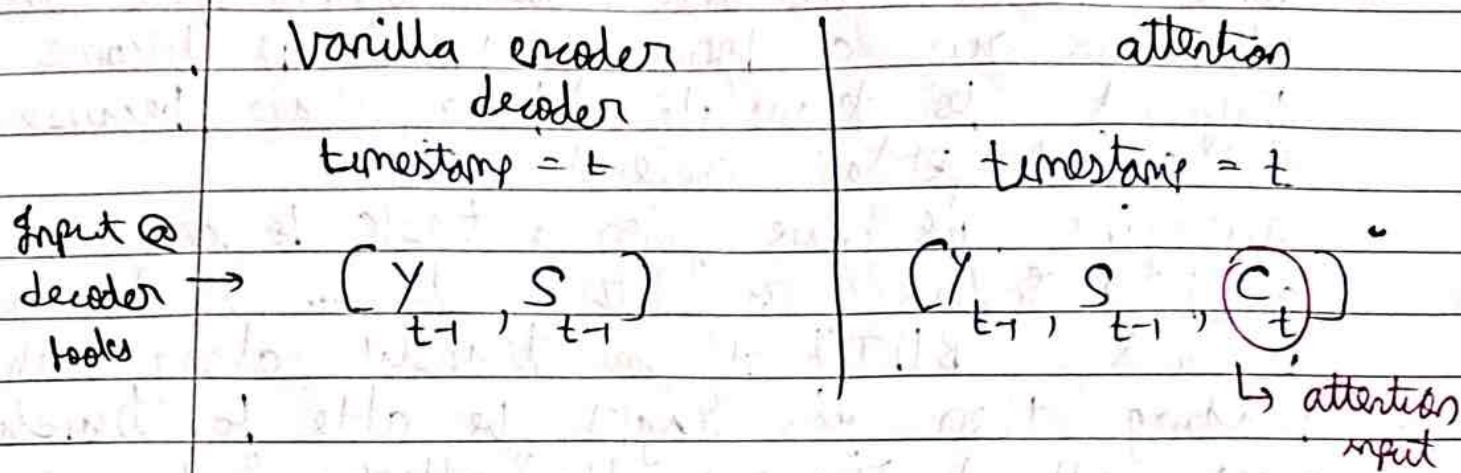


# \* Attention Mechanism

- The Encoder-Decoder model worked fine, but when we give a long paragraph it became difficult to translate long para's, because it cannot perform context vector.
- Imagine, we have given a task to convert a long English para into hindi, .... it becomes difficult, .... BUT! if you translate along with reading, then you might be able to translate well, all it requires the "attention" to each words.



- at timestamp  $t=2$  (In Decoder block)  
 $\rightarrow [Y, S, ] \sim \text{Input}$



- $C_t$  is a vector which has same dimension as  $h_t$
- $C_t$  is weighted sum of all hidden states

~~For timestamp (t) =~~

For timestamp (t) = )  
 $\uparrow$  scalar

$$C_t = \alpha_1 h_1 + \alpha_2 h_2 + \alpha_3 h_3 + \alpha_4 h_4$$

$\uparrow$  vectors

$$C = \sum_{t_1, t_2} \alpha_{t_1 t_2} h_{t_1}$$

$t_1$ : timestamp for decoder.  
 $t_2$ : timestamp for encoder.

- How  $\alpha_1, \alpha_2, \dots, \alpha_n$  are calculated?

1) Before that, see  $\alpha$  is dependent on  $[h_1]$  &  $[S_1]$

Why  $S_1$ ?

because:- Given, jitra. translation ab tak ho chuka uske basis pe ye batao ki agle step me output print krne k liye encoder ka first step kitna useful hoga  
 (decoder block me output print krna hoga  $h_n$ )

$$\alpha_2 \rightarrow f(h_1, S_1)$$



$$\text{i.e. } \alpha_{t, t_2} = f(h_{t_2}, s_{t_1-1})$$

2) Alpha ( $\alpha$ ) is a function of  $h$  &  $s$

$$\alpha_{2,3} \rightarrow f(h_3, s_1)$$

What is this function??  $\rightarrow$

**ANN!!**

because ANN's are the best function approximators

3) When you provide enough data to ANN's, it will automatically adjust its weights to approximate the function, hence to determine value for ( $\alpha$ ). Brilliant idea, indeed.

4) Initially ANN gives random weightage & value ( $\alpha$ ) to everyone, but with loss function and backpropagation it will learn and adjust itself to give priority to ( $\alpha$ ) to a specific hidden state's (context vector)

## \* Attention Mechanism

- Calculates relevance bet<sup>n</sup> each input word and the current output step.
- Assigns weights (called <sup>attention score</sup> ~~weighted sum~~) ( $C_i$ ) to input tokens based on how important they are.
- Produces a weighted sum of the ~~decoder~~ encoder outputs (called the context vector) for each decoder block cell.
- Helps decoder focus on relevant words from the input sequence while generating each word in the output.