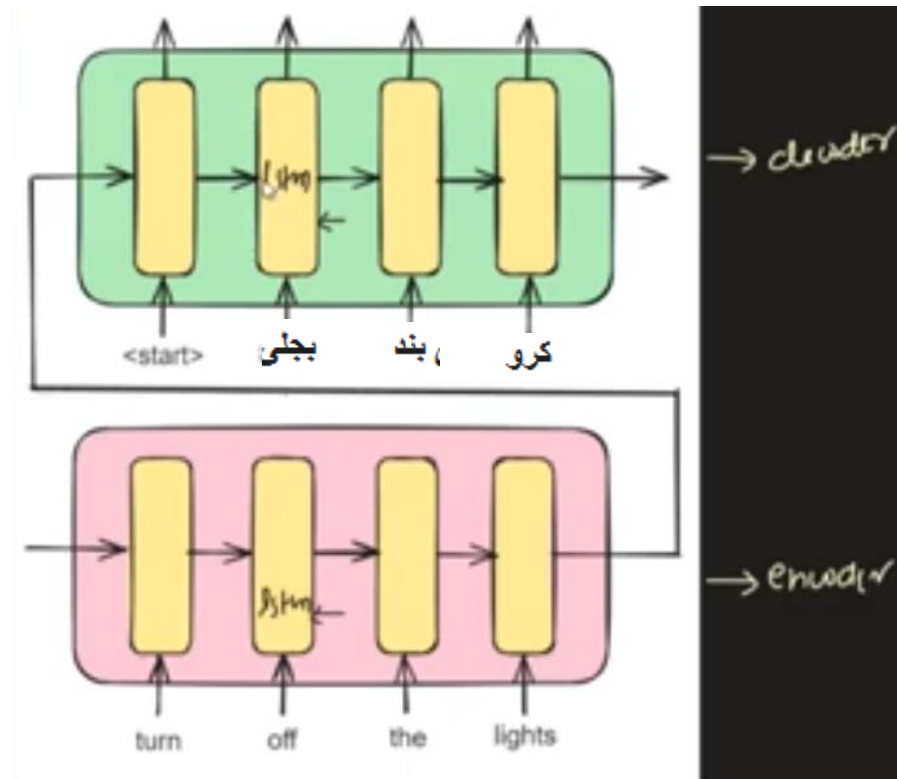


# Attention Models

Dr. Muhammad Safyan

# Attention Mechanism

- A way to improve the Encoder-Decoder Architecture.
- Problems with Encoder/Decoder



# Translate in Urdu

- Translate in Urdu

1<sup>st</sup> Problem:

Once upon a time in a small Pakistani village, a mischievous monkey stole a turban from a sleeping barber, wore it to a wedding, danced with the bewildered guests, accidentally got crowned the 'Banana King' by the local kids, and ended up leading a vibrant, impromptu parade of laughing villagers, cows, and street dogs, all while balancing a stack of mangoes on its head, creating a hilariously unforgettable spectacle and an amusing legend that the village still chuckles about every monsoon season.

Face Difficulty when the sentence is greater than 25 word.

2<sup>nd</sup> problem:

Turn off the light.

Light:

Turn off:

Don't need the complete sentence, only a specific part of that.

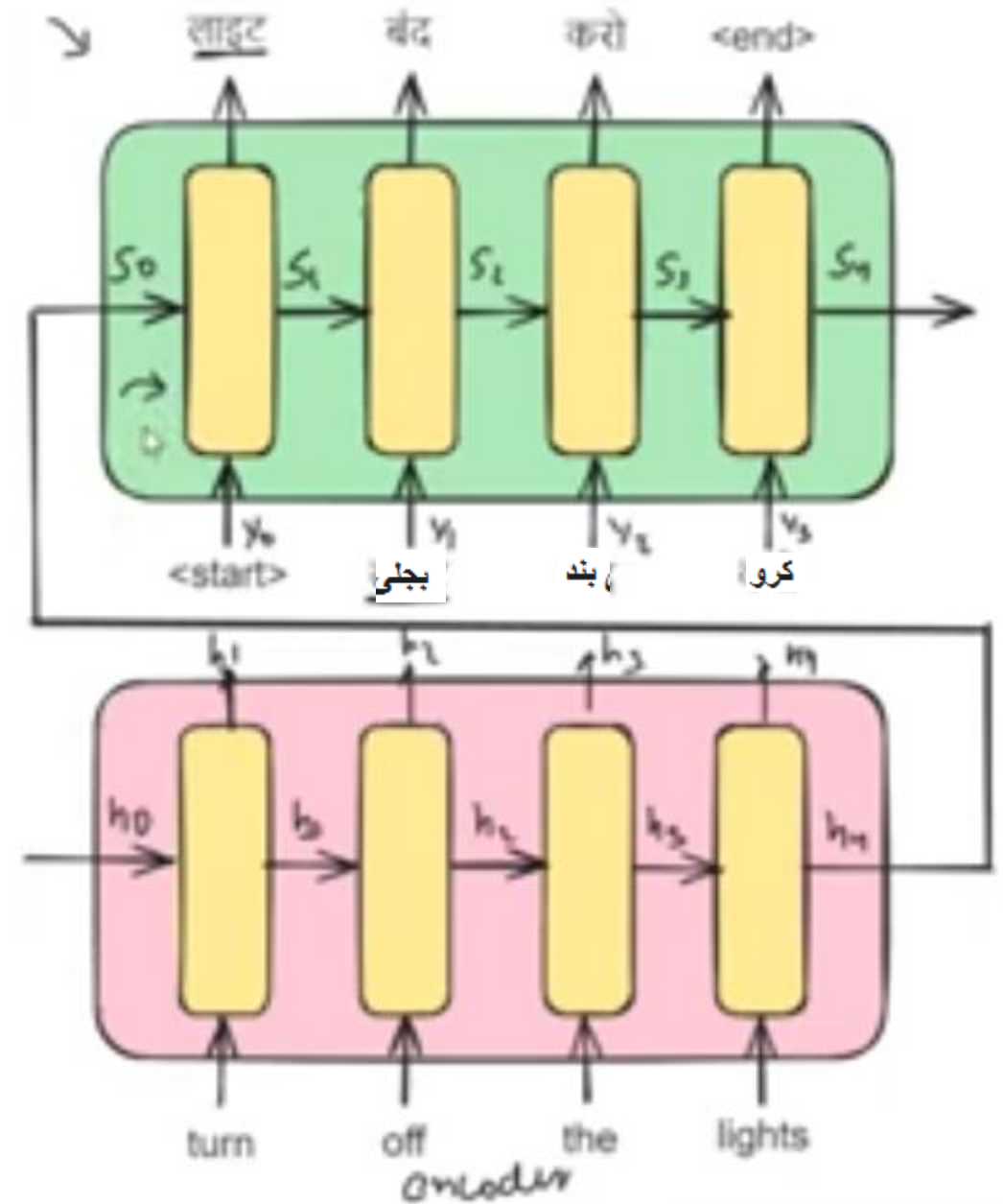
But we give complete sentence to Decoder to translate → Static Representation

It would be good, if we attention a specific part of the sentence at the time of translation.

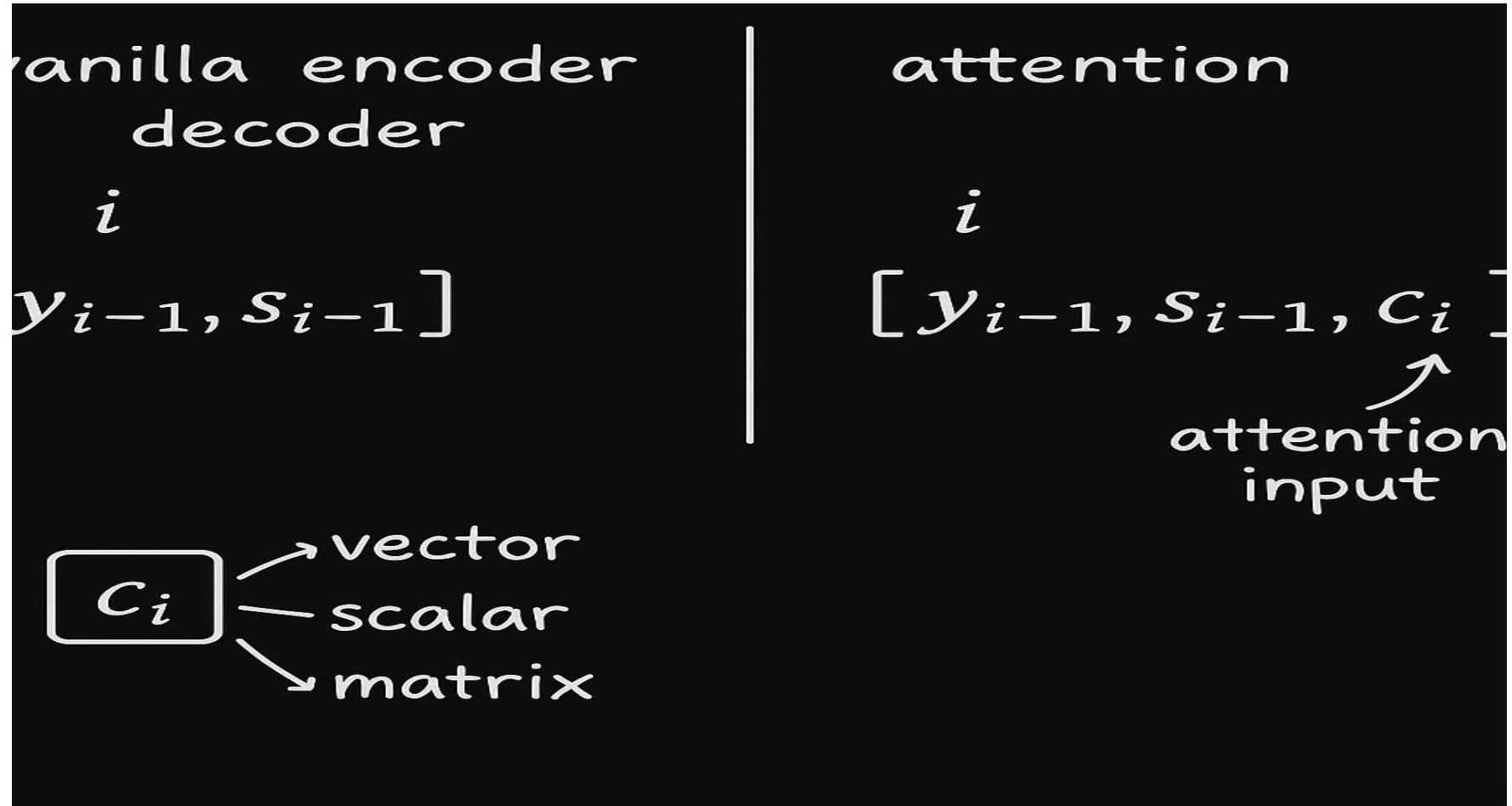


- When we read a larger sentence. Our
- Eyes creates a attention span/region.
- And other things are blurry at that moment.
- We need to introduce this thing in our architecture.
- When translating “light” which portion of the sentence need pass.
- Which time step is important to translate the light.
- This mechanism is called Attention Mechanisms

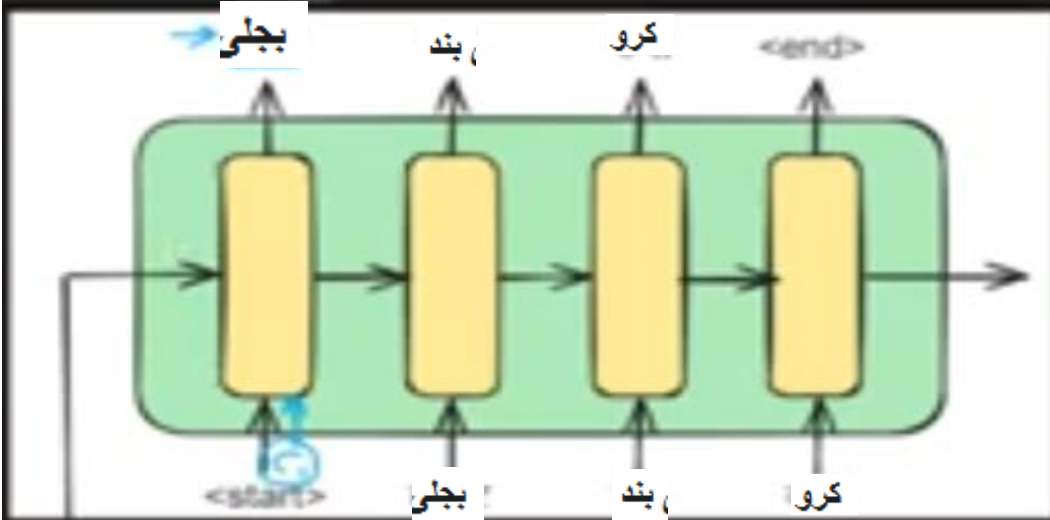
- At time step  $i=2$ , we provided
- $Y_1, s_1$
- In Attention mechanism you provide
- More piece of information ( $h_1, h_2, h_3$ , or  $h_4$ )
- We say it  $c_i \rightarrow$  attention input



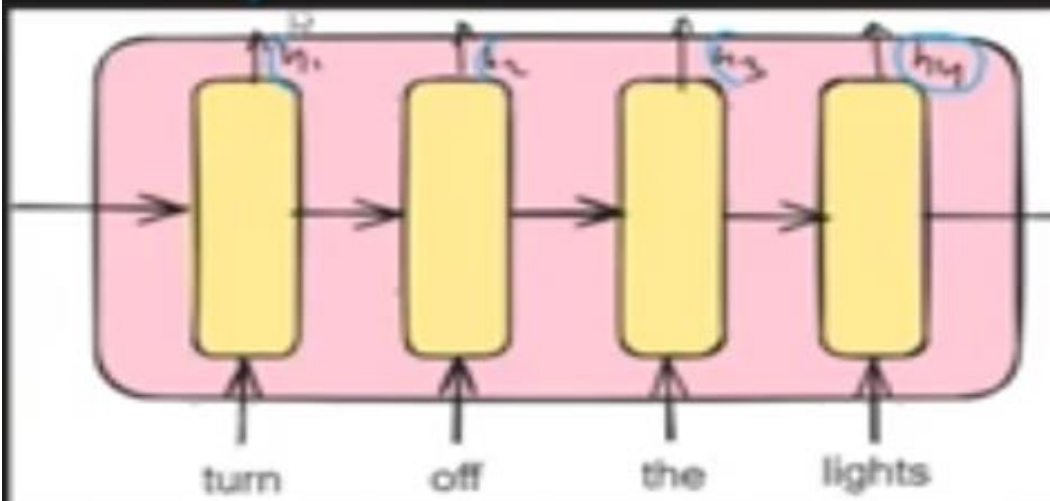
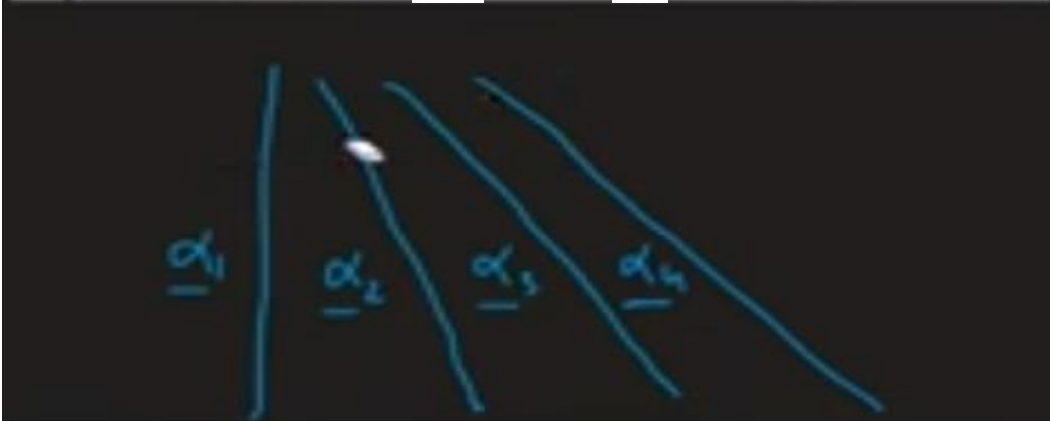
# Vanilla encoder vs. Attention bases



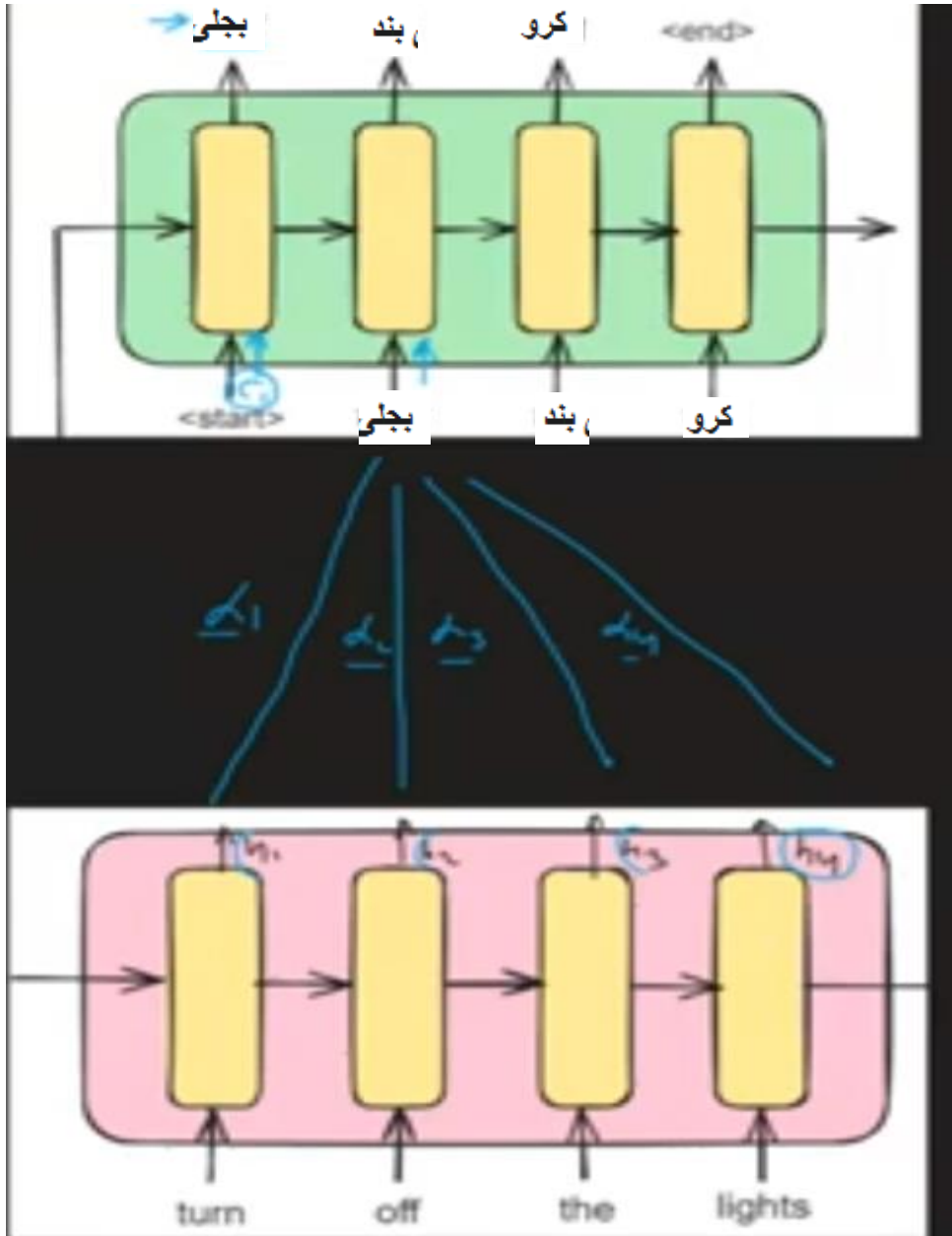
- What is dimension of  $c_i$ ,
- $C_i$  is vector that is may be  $h_1$ ,  $h_2$  or combination of  $h_1$  and  $h_2$
- For both weighted sum Is used.



$$C_1 = \alpha_1 h_1 + \alpha_2 h_2 + \alpha_3 h_3 + \alpha_4 h_4$$







$$C_1 = \alpha_{11}h_1 + \alpha_{12}h_2 + \alpha_{13}h_3 + \alpha_{14}h_4$$

At time step = 2

$$C_2 = \alpha_{21}h_1 + \alpha_{22}h_2 + \alpha_{23}h_3 + \alpha_{24}h_4$$

$$C_i = \sum_j \alpha_{ij}h_j$$

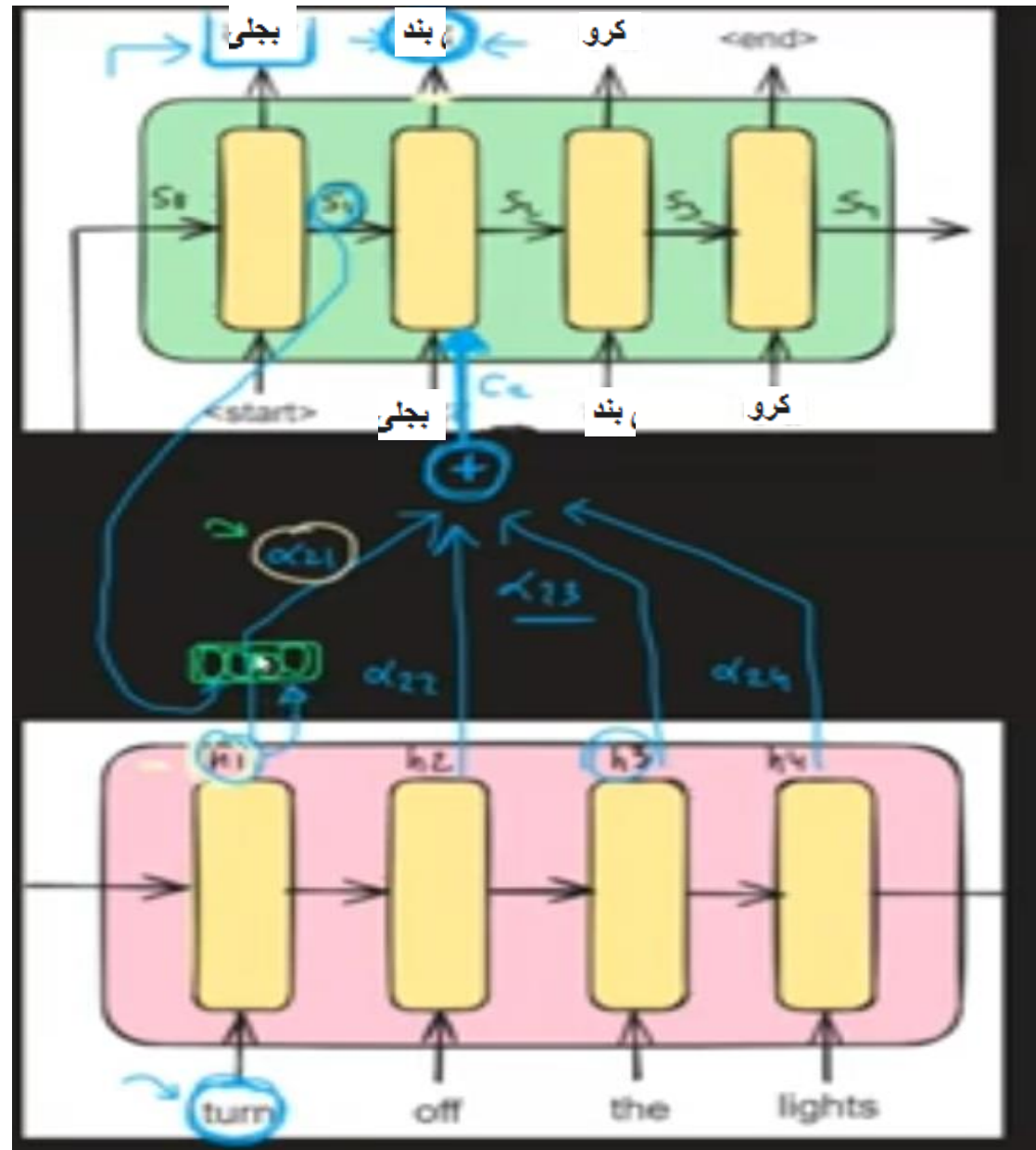
i\*j=16

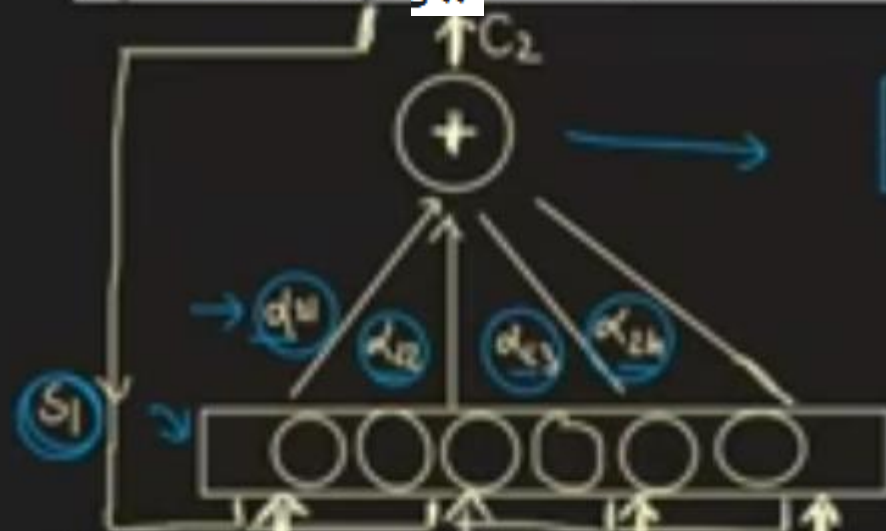
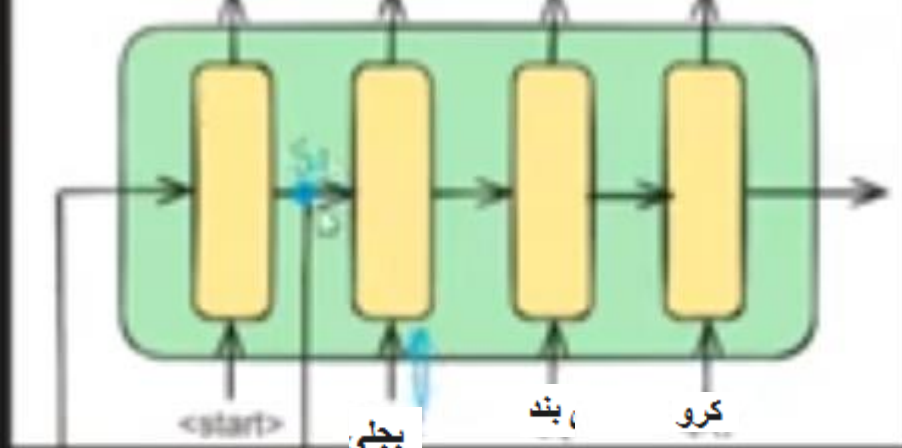
- How the alpha will create
- calculate  $\alpha_{21}$ 
  - Called Alignment Score
  - Similarity Score.
- At timestep,  $i=2$ , the output is printed , what is the role of encoder time step  $j=1$
- i.e. “Turn” or  $h_1$
- Now  $\alpha_{21}$  depends upon on which quantities
- $H_1, s_1$

- ANN Universal approximate function

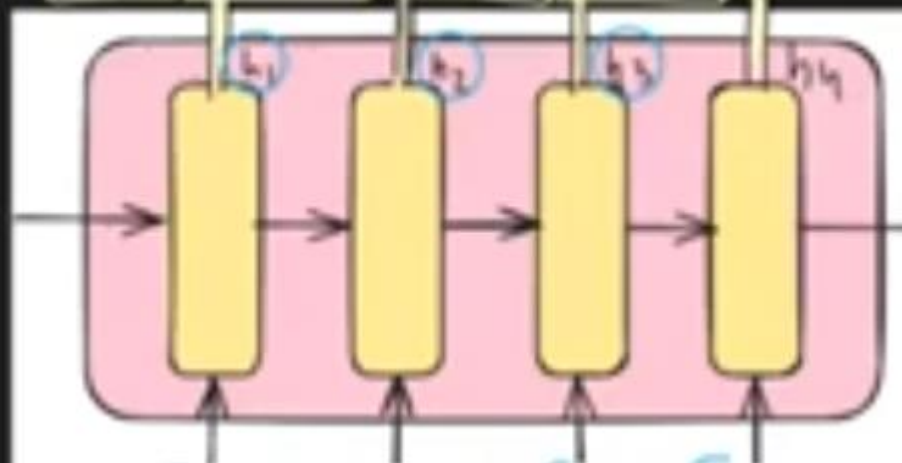
$$\alpha_{21} \rightarrow f \longrightarrow \left[ \alpha_2 = f(h_3, s_1) \right]$$

$$\alpha_{ij} = f h_j = f(h_3, s_1)$$

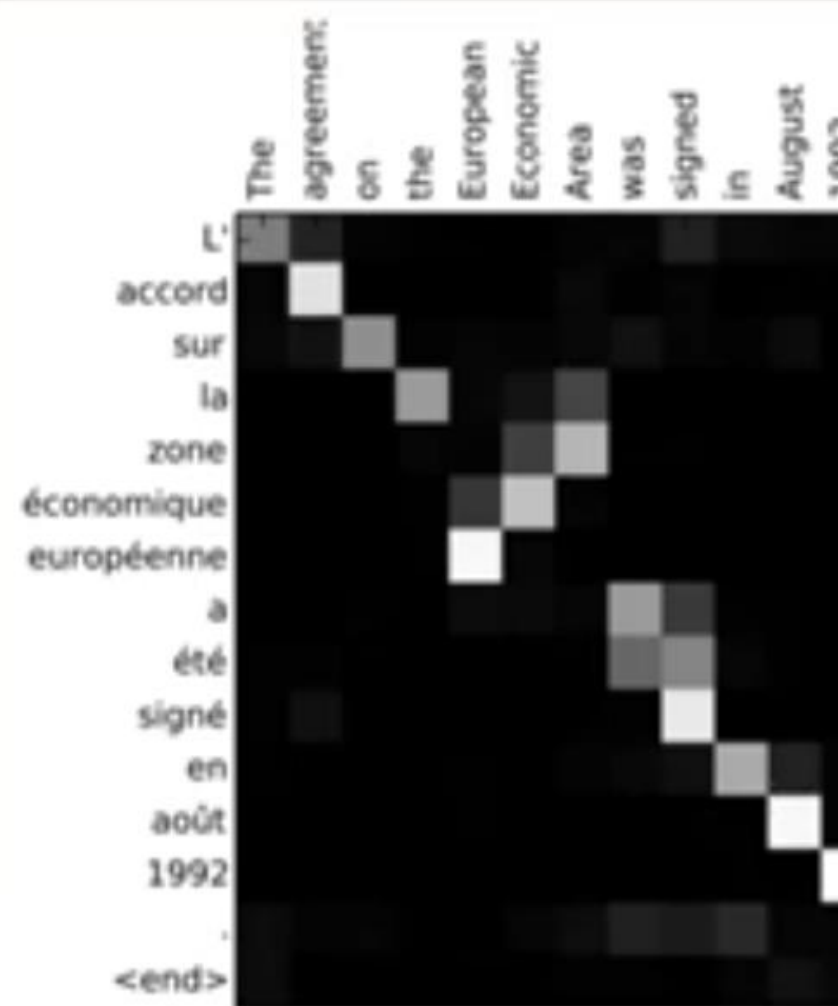
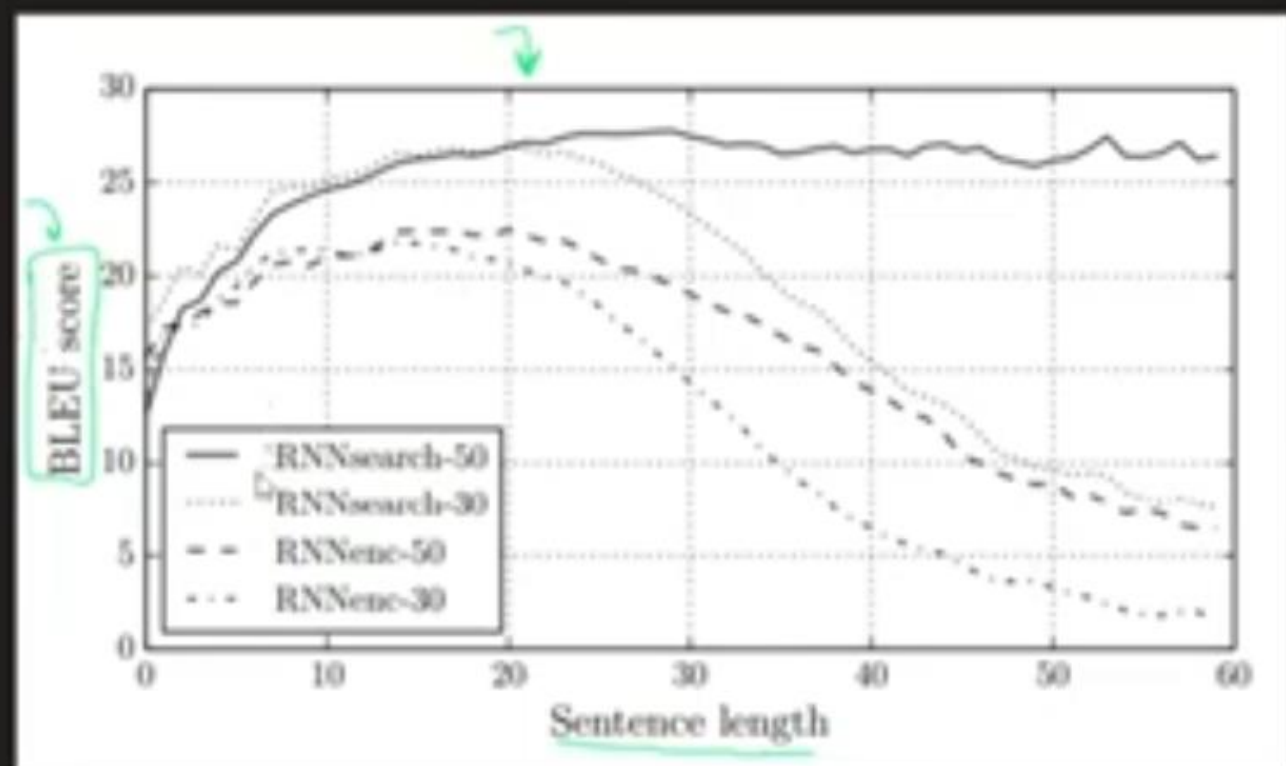




$$c_2 = \alpha_{21}h_1 + \alpha_{22}h_2 + \alpha_{23}h_3 + \alpha_{24}h_4$$



Time  
Distributed  
Dense Layer



(a)

$\alpha_{11}$	$\frac{\partial \Pi}{\partial T_c}$	turn	off	the	light
$\alpha_{12}$	$\frac{\partial \Pi}{\partial k}$	$\alpha_{11}$	$\alpha_{12}$	$\alpha_{13}$	$\alpha_{14}$
	$\frac{\partial \Pi}{\partial \lambda_1}$				
	end				

- Bahduana and Ioung Attention Model



- Two types to find the Alphas
  - Bahduana
  - Lounɡ
- Alignment Score

$$C_i = \sum_j \alpha_{ij} h_j$$

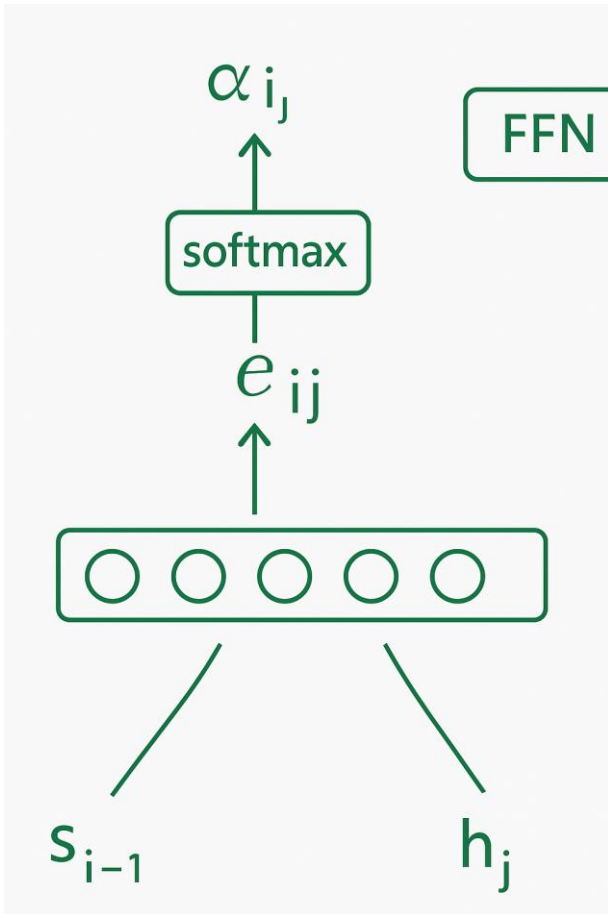
alignment

$\alpha_{11} \rightarrow \overline{01152} \rightarrow \underline{\text{turn}}$

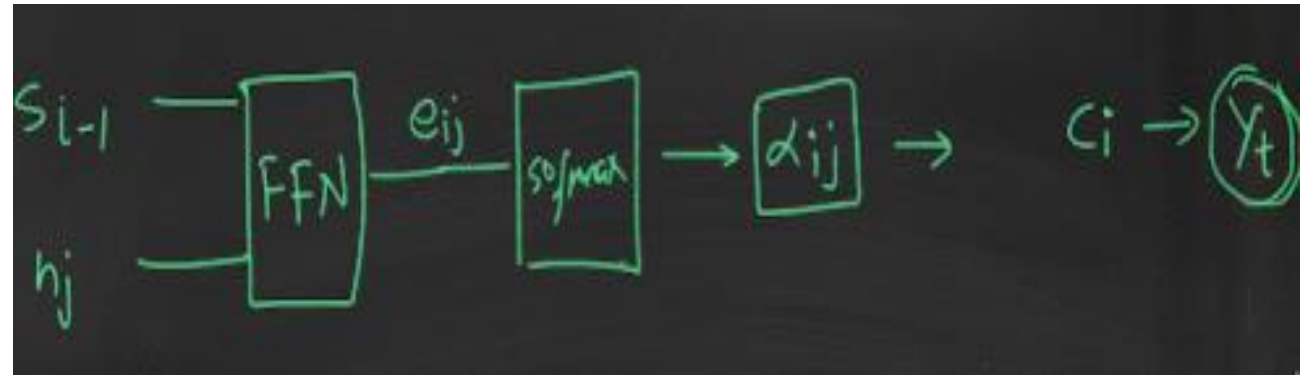
$\alpha_{12} \rightarrow \overline{01152} \rightarrow \text{off}$

$$\underline{\alpha_{11}} = f(h_1, \underline{s_0}) \quad \underline{\alpha_{21}} = f(h_1, s_1)$$

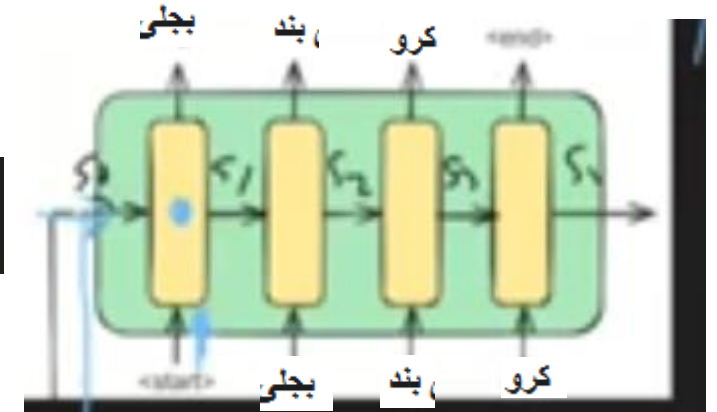
# Bahduana



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



$s_0 = [e \ f \ g \ h]$



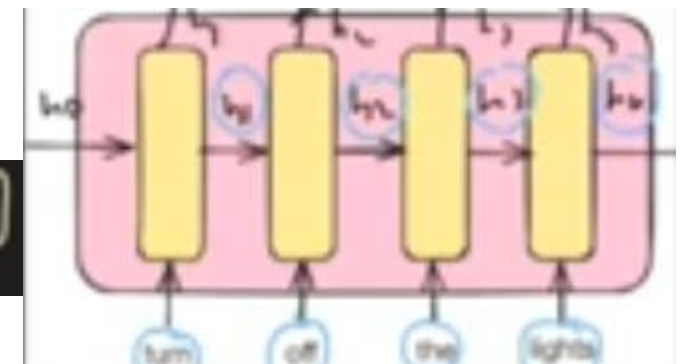
- Enter the encoder sentences.
  - Turn off the light
  - You will be get hidden stats( $h_1, h_2, h_3, h_4$ )
  - Assume
  - At  $i=1$ , decoder will calculate

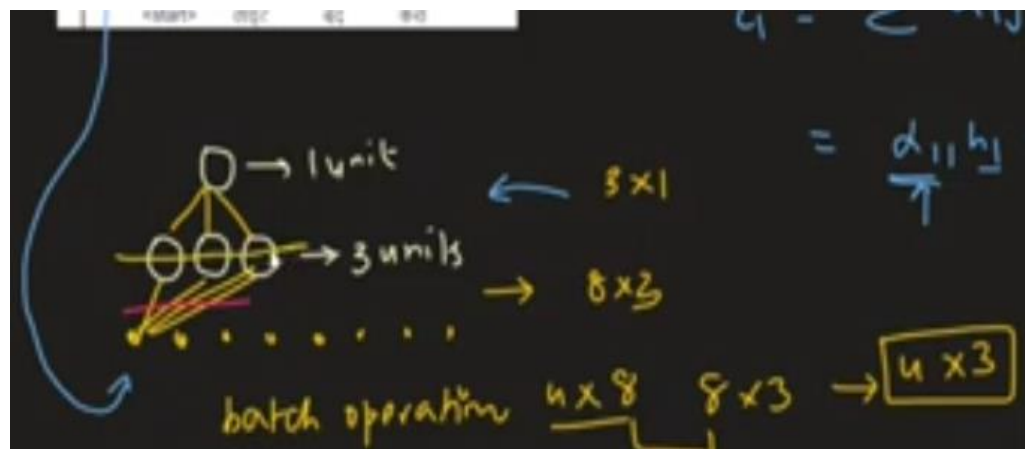
0  
000

$$c_1 = \sum \alpha_{ij} h_j$$

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

$h_0 = [a \ b \ c \ d]$





$(4 \times 1) \rightarrow 4 \text{ numbers}$

$$s_0 = [e \ f \ g \ h]$$

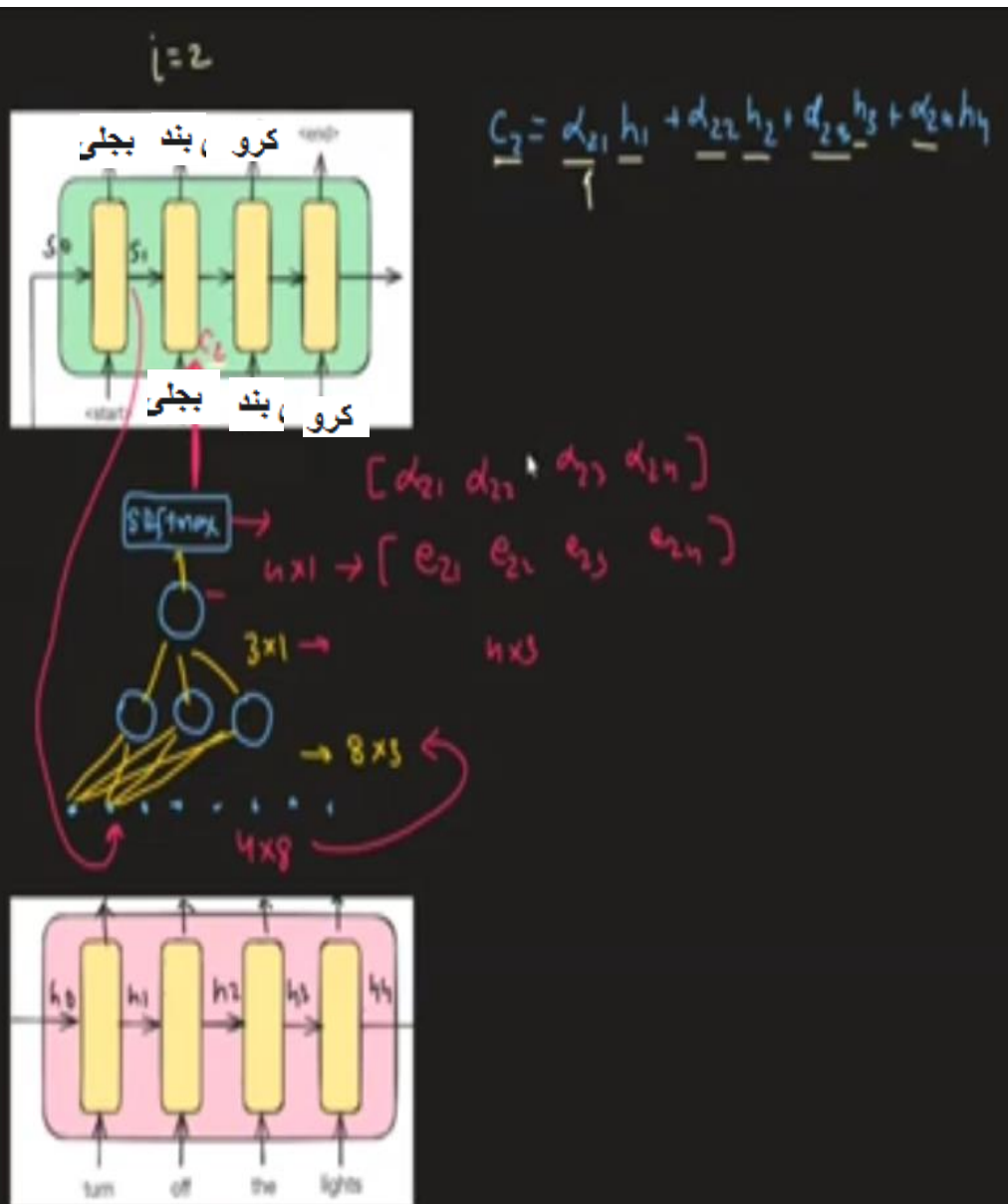
Diagram showing a vector  $s_0$  connected to four hidden units  $h_1, h_2, h_3, h_4$  via a tree structure.

$s_{01}$	$s_{02}$	$s_{03}$	$s_{04}$	$h_{11}$	$h_{12}$	$h_{13}$	$h_{14}$
$s_{01}$	$s_{02}$	$s_{03}$	$s_{04}$	$h_{21}$	$h_{22}$	$h_{23}$	$h_{24}$
$s_{01}$	$s_{02}$	$s_{03}$	$s_{04}$	$h_{31}$	$h_{32}$	$h_{33}$	$h_{34}$
$s_{01}$	$s_{02}$	$s_{03}$	$s_{04}$	$h_{41}$	$h_{42}$	$h_{43}$	$h_{44}$

$$\frac{e^{e^{11}}}{e^{e^{11}} + e^{e^{12}} + e^{e^{13}} + e^{e^{14}}} (4 \times 1) \rightarrow$$

$$\text{So } Y_{t+1} C_1 \rightarrow \text{LSVM} \rightarrow Y_t (\text{ans})$$

At time=2



$$C_2 = \alpha_{21} h_1 + \alpha_{22} h_2 + \alpha_{23} h_3 + \alpha_{24} h_4$$

$$\begin{bmatrix} S_{11} & S_{12} & S_{13} & S_{14} & h_{11} & h_{12} & h_{13} & h_{14} \\ S_{21} & S_{22} & S_{23} & S_{24} & h_{21} & h_{22} & h_{23} & h_{24} \\ S_{31} & S_{32} & S_{33} & S_{34} & h_{31} & h_{32} & h_{33} & h_{34} \\ S_{41} & S_{42} & S_{43} & S_{44} & h_{41} & h_{42} & h_{43} & h_{44} \end{bmatrix}$$

$$c_2, s_1, y_1 \rightarrow \text{let } y_2 \rightarrow \frac{\dot{c}_2}{s_2}$$

# Mathematical Form

$$c_{ij} = \sum_{\underline{j}} \alpha_{ij} \underline{h_j}$$

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{\underline{k}} \exp(e_{ik})}$$

$$e_{ij} = \tanh(W[s_{i-1} : h_j] + b)$$

$$e_{ij} = v \left[ \tanh(W[s_{i-1} : h_j] + \underline{b}) \right]$$



- Also called additive algorithm
- Model is called alignment model

# Loung Attention

$$c_i = \sum \alpha_{ij} h_j$$

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

$$\alpha_{ij} = f(s_i, h_j)$$

↑  
current    ② diff

$$\underline{s_i}^T \cdot \underline{h_j} \rightarrow \text{dot product}$$

① diff

$$s_i = [a \ b \ c \ d]$$

$$h_j = [e \ f \ g \ h]$$

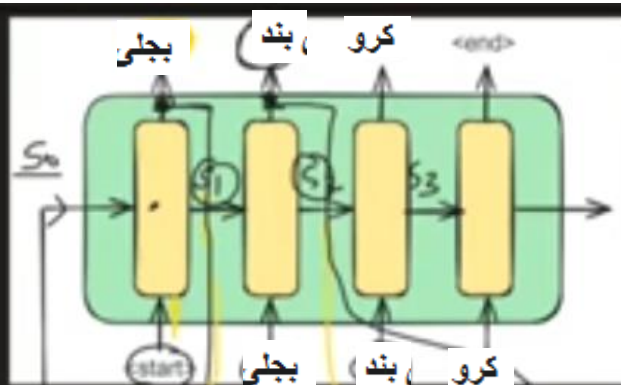
$$\begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix} [e \ f \ g \ h]$$

$$[ae + bf + cg + dh]$$

$$\leftarrow \boxed{e_{ij}}$$

← scales → attention

- Current state has fore info then previous
- Experimentally proved



output  
 $S_1: \underline{c_1} \rightarrow \text{softmax} \rightarrow \tilde{S}_1$

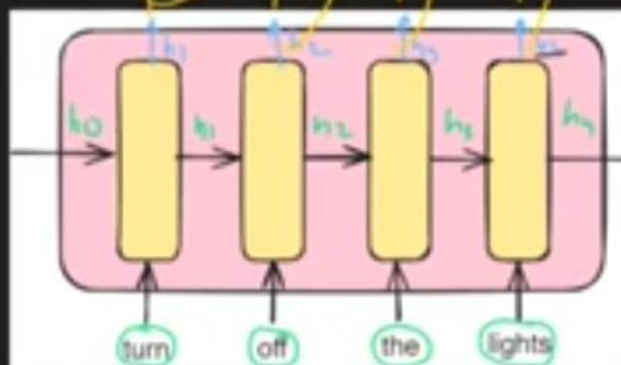
$S_2: c_2 \rightarrow \tilde{S}_2 \rightarrow$

$[e_{21} \ e_{22} \ e_{23} \ e_{24}] \text{ softmax} \rightarrow \alpha_{21} \ \alpha_{22} \ \alpha_{23} \ \alpha_{24} \rightarrow c_2$

$\sum \alpha_{i,j} h_j$

$S_1 \ h_1$	$S_1 \ h_2$	$S_1 \ h_3$	$S_1 \ h_4$
$[e_{11}]$	$e_{12}$	$e_{13}$	$e_{14}]$
$\downarrow$	$\downarrow$	$\checkmark$	$\checkmark$
$\alpha_{11}$	$\alpha_{12}$	$\alpha_{13}$	$\alpha_{14}$

$\rightarrow \boxed{c_1}$



- Called multiplicative
- Necessary for understanding of transformer → self attention.