

alignment score

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

alignment

$$\alpha_{11} \rightarrow \text{on} \rightarrow \text{turn}$$

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

given

decoder
→ prev hidden state

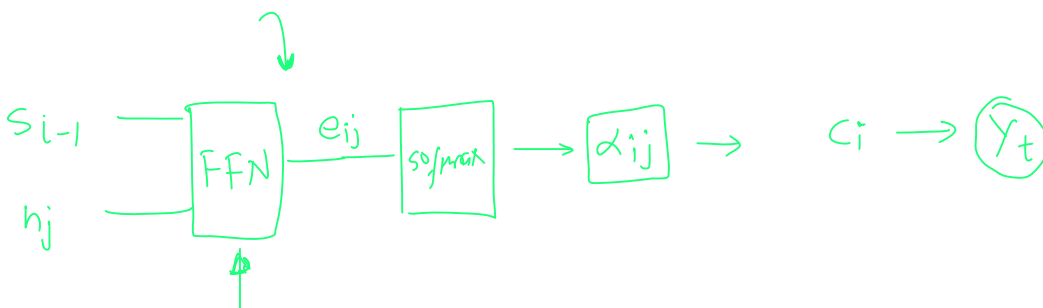
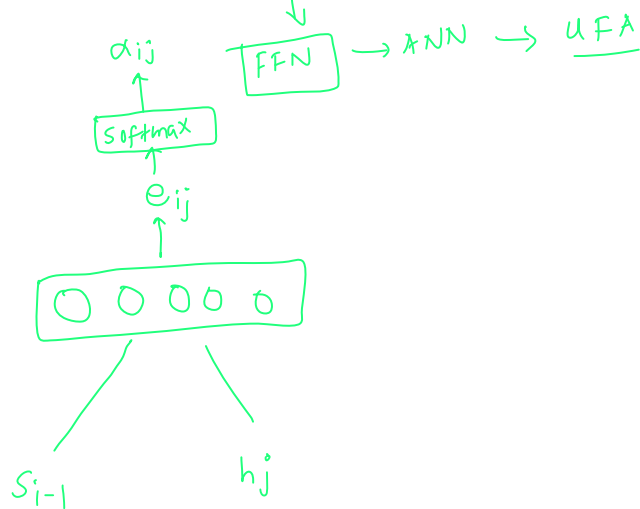
$$\alpha_{11} = f(h_1, s_0) \quad \alpha_{21} = f(h_1, s_1)$$

$$\alpha_{ij} = \text{approximate} \left(\text{matv} \left(h_j, s_{i-1} \right) \right)$$

matv
funct

the
encoder
hidden

deal
prev
hidden
state



$$s_0, y_{t-1}, c_1 \rightarrow \text{LSTM} \rightarrow y_t \text{ (output)} \quad [s_1]$$

encoder hidden
→ h_1, h_2, h_3, h_4

$$[s_{i-1}] h_j$$

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

$$c_1 = \sum_j \alpha_{1j} h_j$$

$$[e_{11}, e_{12}, e_{13}, e_{14}] \rightarrow \text{softmax} \rightarrow [\alpha_{11}, \alpha_{12}, \alpha_{13}, \alpha_{14}]$$

$$c_1 = \frac{\alpha_{11} h_1}{1} + \frac{\alpha_{12} h_2}{1} + \frac{\alpha_{13} h_3}{1} + \frac{\alpha_{14} h_4}{1}$$

unit

3x1

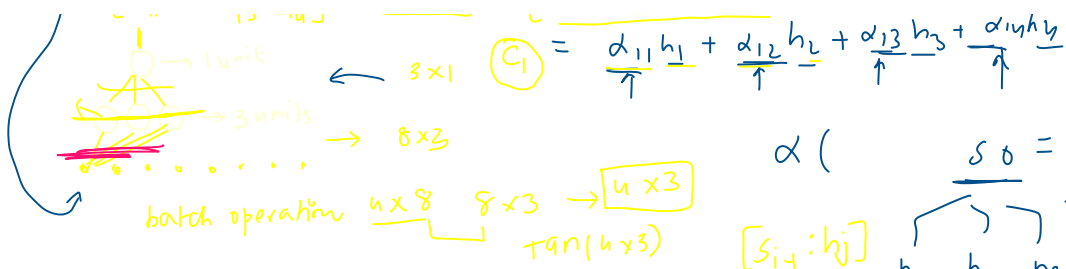


Diagram showing the calculation of α for a sequence $s_0 = [e, f, g, h]$. The sequence is mapped to h_1, h_2, h_3, h_4 . The operation is defined as:

$$\alpha([s_{i-4}:h_j])$$

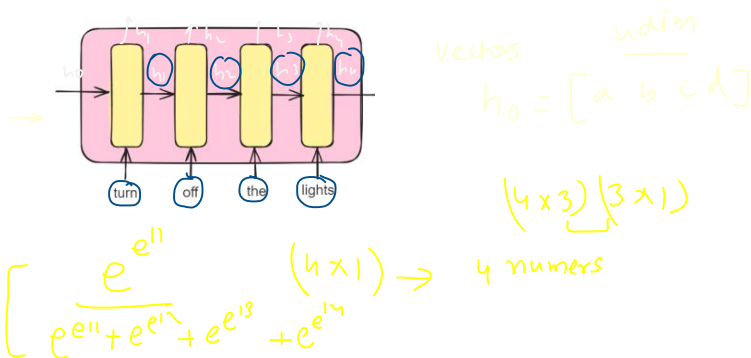
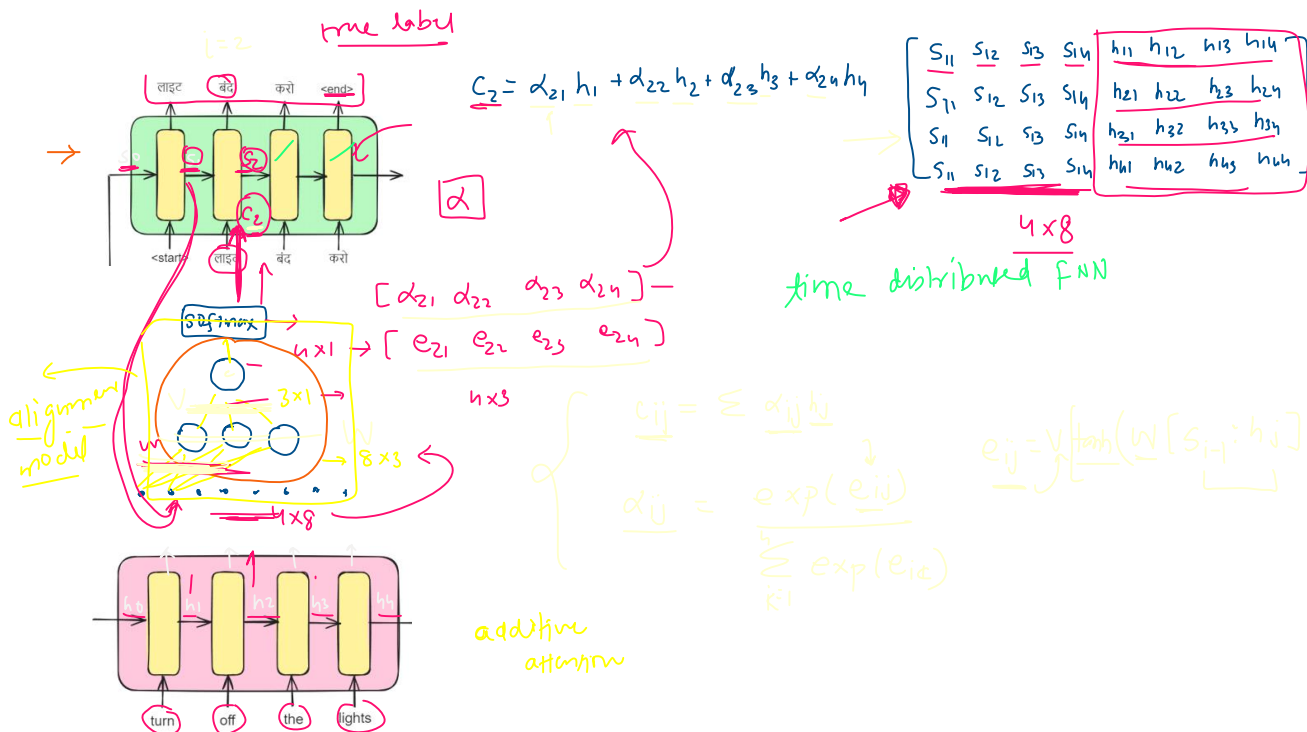


Diagram showing a sequence of words: "turn", "off", "the", "lights". The sequence is mapped to a vector $h_0 = [a, b, c, d]$. The operation is defined as:

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

Dimensions: 4×3 (input), 3×1 (output). The operation is $(4 \times 3)(3 \times 1)$.



Luong Attention

17 January 2024 00:09

Luong

parameters \rightarrow slow

$$c_i = \sum \alpha_{ij} h_j \rightarrow \text{FFN} \left\{ \left[V + \tan(W[S_{1:i-1}; h_i] + b) \right] \right\}$$

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

$$\alpha_{ij} = f(s_i, h_j) \rightarrow [s_i^T \cdot h_j] \rightarrow \text{dot product fast}$$

updated info current ① diff

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

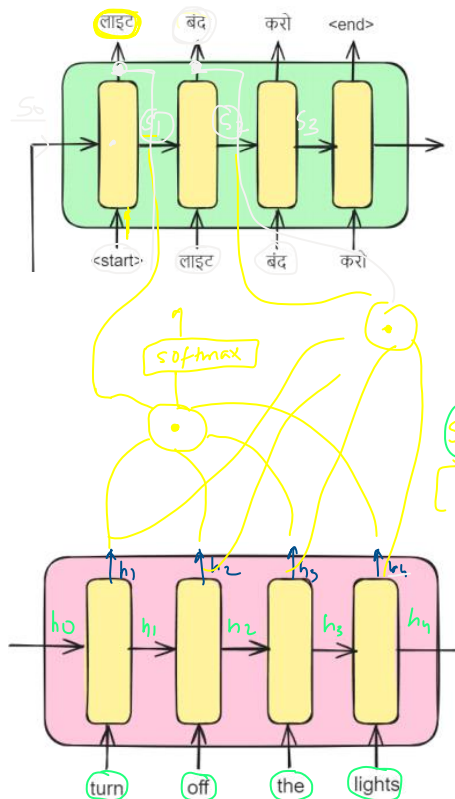
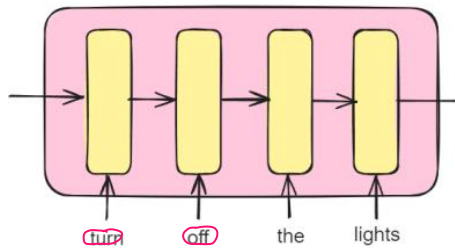
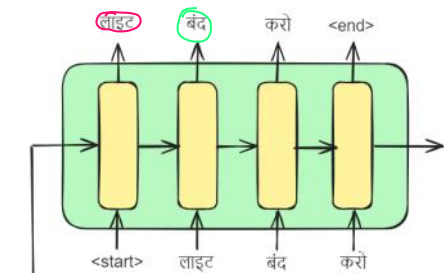
$$h_j = [e \ f \ g \ h]$$

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

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

$$\text{softmax} \leftarrow [e_{ij}]$$

slow \rightarrow attention



output

$$s_1: c_1 \rightarrow \text{softmax} \rightarrow [s_1]$$

$$s_2: c_2 \rightarrow \tilde{s}_2 \rightarrow \text{multiplicative atten}$$

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

$$\rightarrow c_2$$

$$\begin{bmatrix} s_1 h_1 & s_1 h_2 & s_1 h_3 & s_1 h_4 \\ e_{11} & e_{12} & e_{13} & e_{14} \end{bmatrix} \rightarrow \begin{bmatrix} \alpha_{11} & \alpha_{12} & \alpha_{13} & \alpha_{14} \end{bmatrix} \rightarrow c_1$$

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