

10/4/24

Attention

$$\begin{bmatrix} 15 \\ 2 \\ 10 \\ 3 \end{bmatrix} \Rightarrow v \Rightarrow \begin{bmatrix} x_1' \\ x_2' \\ x_3' \\ x_4' \end{bmatrix} \Rightarrow \frac{v}{\sqrt{(x_1')^2 + (x_2')^2 + (x_3')^2 + (x_4')^2}}$$

DL:

Attention is focusing on relevant information/features and not focusing on others.

(Max pool etc.)

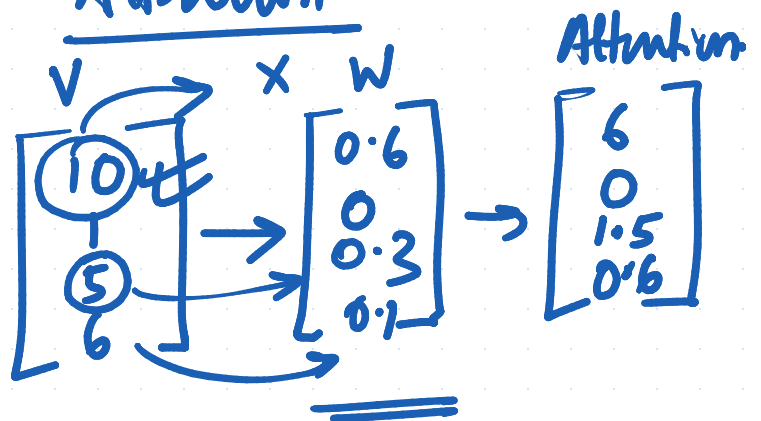
VS

Attention

$$\begin{bmatrix} 10 \\ 1 \\ 5 \\ 6 \end{bmatrix} \rightarrow \begin{bmatrix} 10 \\ 6 \end{bmatrix}$$

dim
(2x1)

↓
compression.



1
dim = same as input

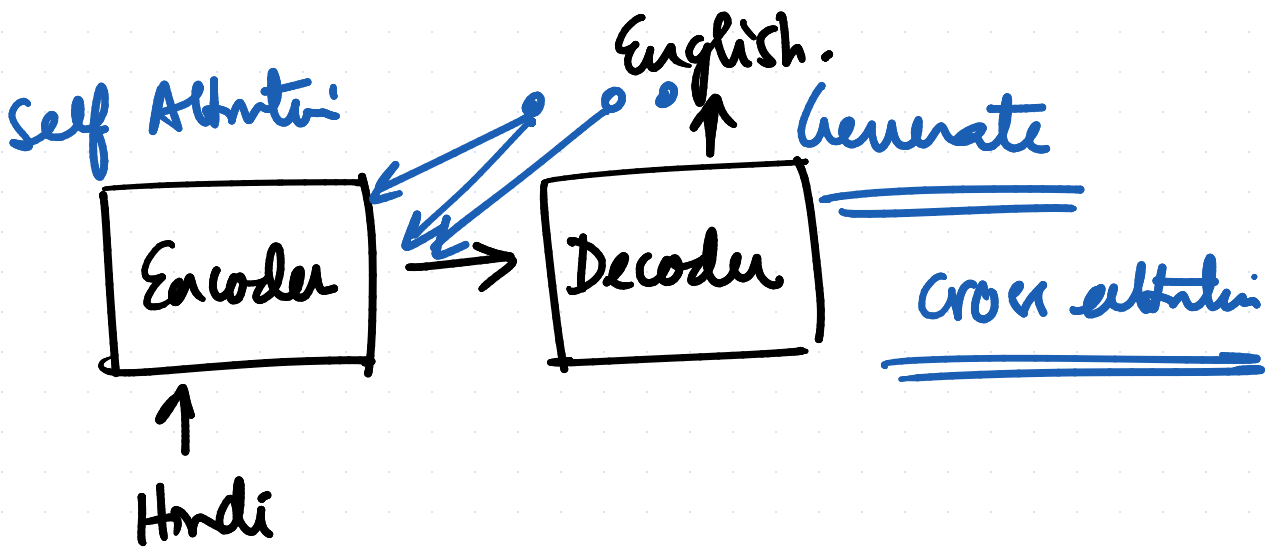
$$\begin{matrix} \text{ } \\ \text{ } \\ \text{ } \\ \text{ } \\ \text{ } \end{matrix} \in \mathbb{R}^{n \times 1} \quad W = \mathbb{R}^{n \times 1} \quad R \in \mathbb{R}^{n \times 1}$$

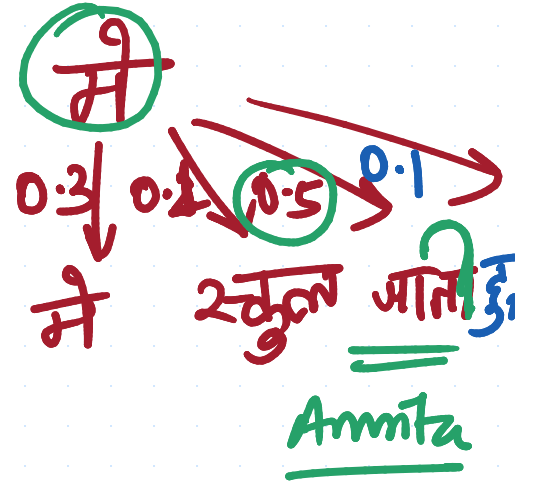
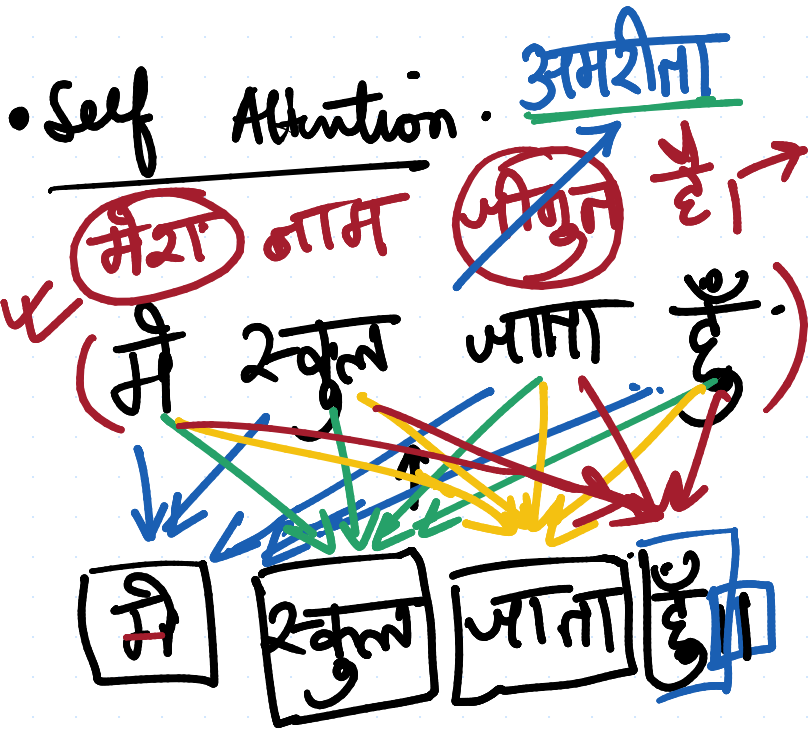
$$\begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix} \otimes \begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

Find the weight matrix which will be used to multiply the initial

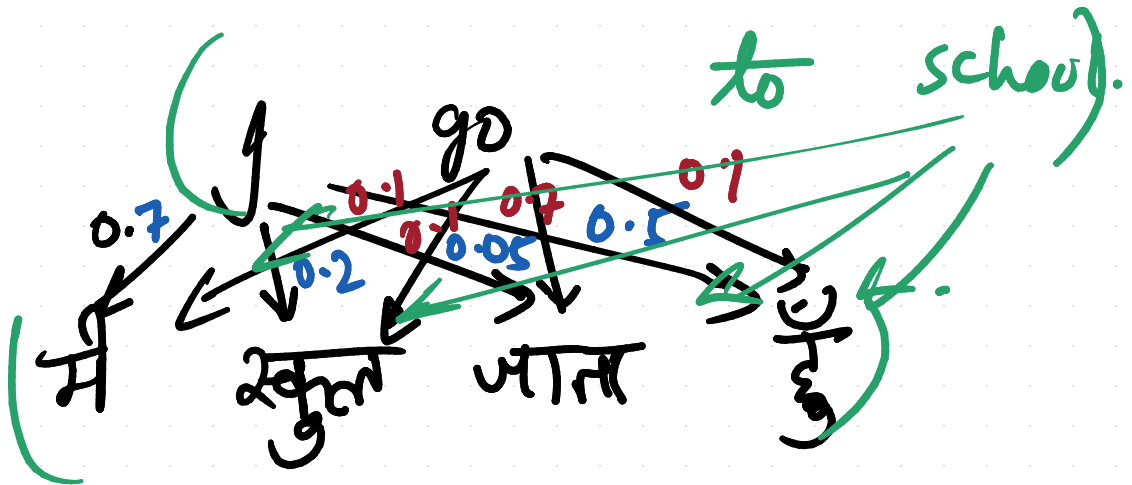
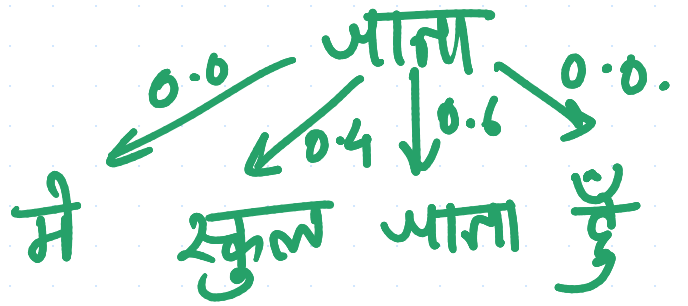
feature volume so that the less important features gets less weight & the more imp feature gets more weight.

- I go to school.
- मैं स्कूल जाता हूँ।





self Attention



;

Self Attention using matrix

$$X_{\text{input}} \boxed{x} \textcircled{W^Q} = \underline{Q}$$

$$X_{\text{input}} \boxed{x} \textcircled{W^K} = K$$

$$X_{\text{input}} \boxed{x} \textcircled{W^V} = V$$

→ Weight matrix.

$$\boxed{\text{Softmax} \left(\frac{Q \times K^T}{\sqrt{d_k}} \right)} \times V = \underline{Z}$$