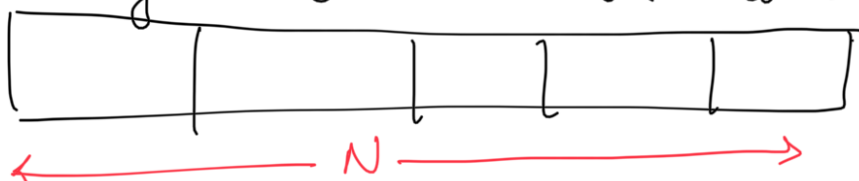


Key, Query in self attention

Sentence \rightarrow Mary had a little lamb

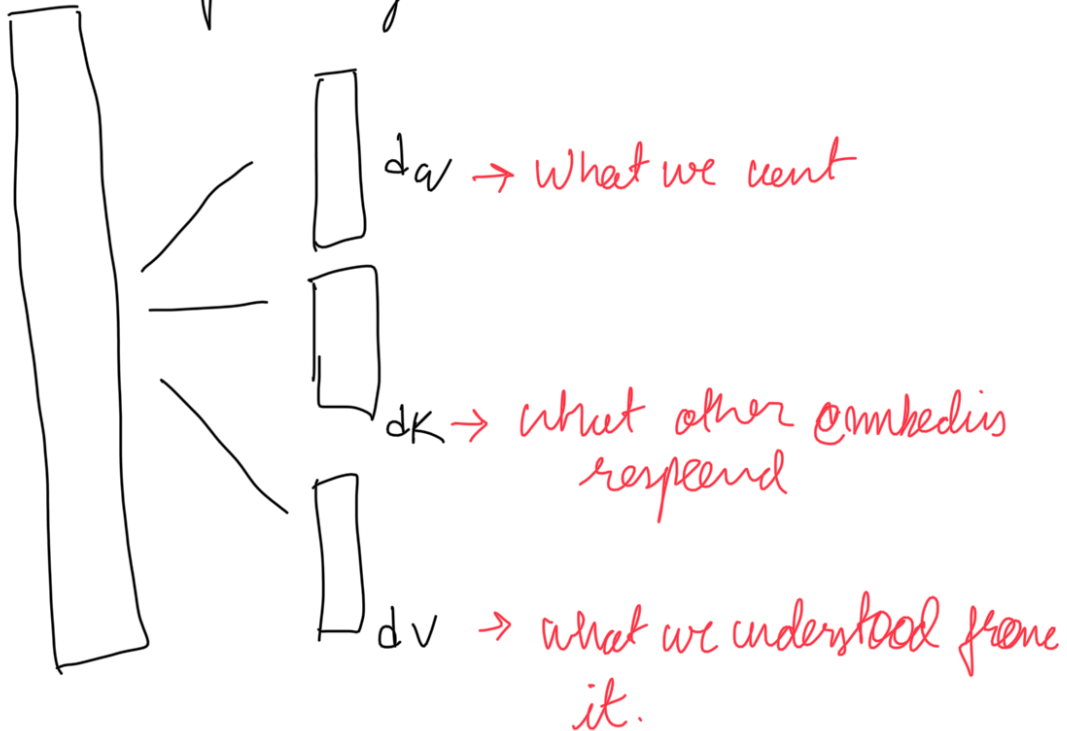


Attention $(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$

\hookrightarrow self attn

Query $\rightarrow Q$
Key $\rightarrow K$
Value $\rightarrow V$

embeddings of mary



$$d_v = d_q = d_k$$

We can pack all values in

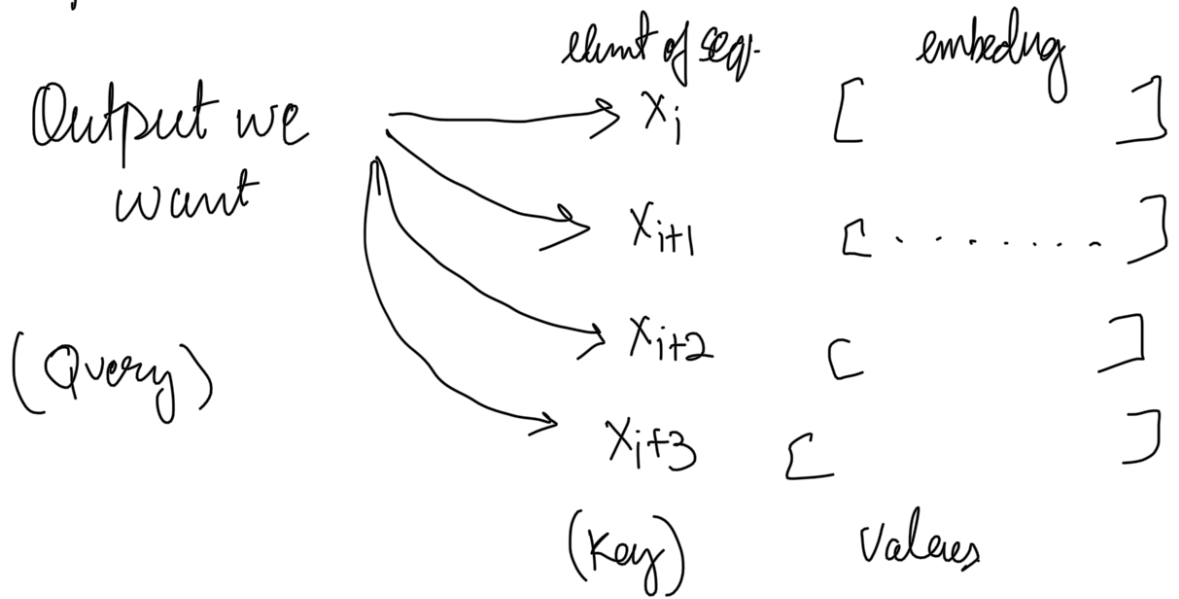
$$\begin{matrix} \uparrow N \\ \downarrow \end{matrix} \begin{bmatrix} \xleftarrow{d_q} \\ N \times d_q \end{bmatrix} \times \begin{bmatrix} d_k \times N \end{bmatrix} = \begin{bmatrix} N \times N \end{bmatrix}$$

so as compute everything at once

row 1 ~

$$\rightarrow [N \times N] \times [N \times d_v] = \text{result}$$

logic



Pick the one we want

masked self attn continued

$$Z_{13} = \frac{K_1' q_3}{\sqrt{d}} \quad (\text{how? look in the PAPER})$$

$$\begin{bmatrix} w_{13} \\ w_{23} \\ w_{33} \\ w_{43} \\ w_{53} \end{bmatrix} = \begin{bmatrix} \exp(Z_{13}) \\ \exp(Z_{23}) \\ \exp(Z_{33}) \\ 0 \\ 0 \end{bmatrix} \Rightarrow \begin{bmatrix} w_{13} \\ w_{23} \\ w_{33} \\ 0 \\ 0 \end{bmatrix} \frac{1}{\sum_{i=1}^3 w_{i3}}$$

$$= \begin{bmatrix} \text{softmax}(Z_{13}, Z_{23}, Z_{33}) \\ 0 \\ 0 \end{bmatrix} \quad \text{nothing but softmax}$$

$$\Rightarrow y_3 = \sum_{i=1}^5 v_i w_{i3}$$

Since $w_{43}, w_{53} = 0$ it doesn't matter