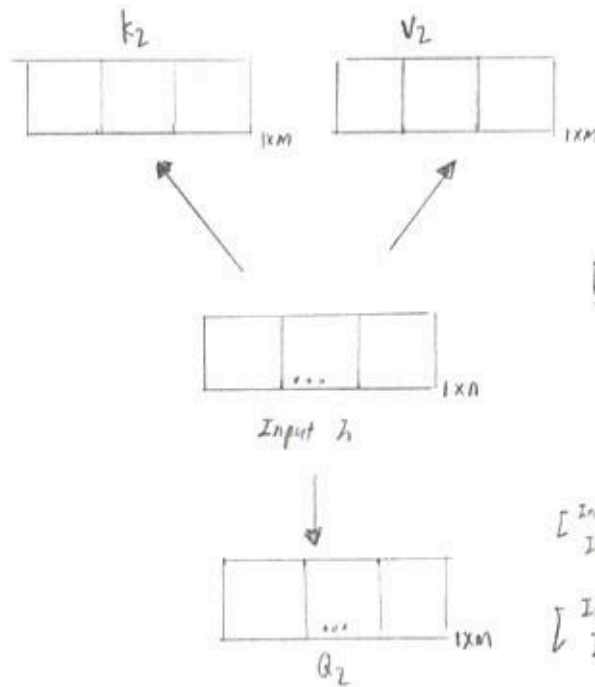
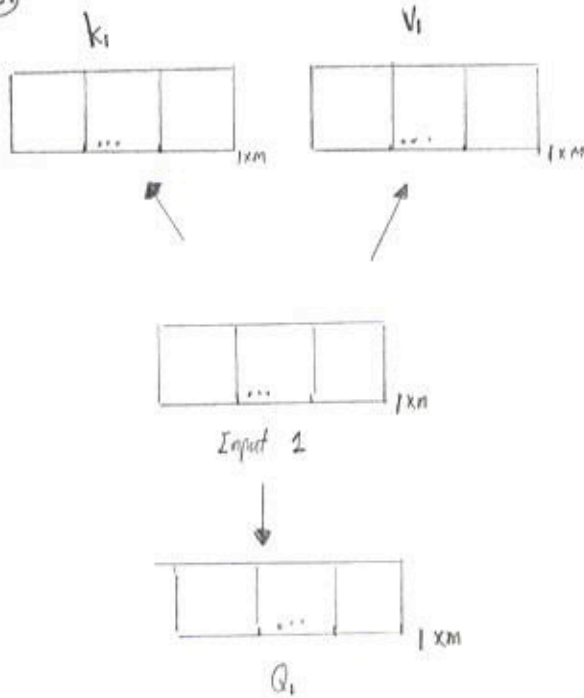


Self-Attention encodes a sequence into a form that accounts for the other respective data in the sequence



* each element in row form is aware of surrounding data

①



$W = [K_{n \times m}, V_{n \times m}, Q_{n \times m}]$
 $n \Rightarrow$ input dim
 $m \Rightarrow$ output dim

Ex: $\begin{bmatrix} \text{Input 1} \\ \text{Input 2} \end{bmatrix}_{2 \times n} K_{n \times m} = \begin{bmatrix} k_1 (1 \times m) \\ k_2 (1 \times m) \end{bmatrix}_{2 \times m}$

$\begin{bmatrix} \text{Input 1} \\ \text{Input 2} \end{bmatrix}_{2 \times n} V_{n \times m} = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}_{2 \times m}$

$\begin{bmatrix} \text{Input 1} \\ \text{Input 2} \end{bmatrix}_{2 \times n} Q_{n \times m} = \begin{bmatrix} q_1 \\ q_2 \end{bmatrix}_{2 \times m}$

② Repeat for each input

For Input 1:

(A)

$$S_{1,(m)} = Q_1 k_1^T$$

$$S_{2,(m)} = Q_1 k_2^T$$

$$S_1 = \begin{bmatrix} s_1 \\ s_2 \end{bmatrix}$$

(B)

$$\text{softmax}(S_1) \rightarrow \hat{S}_1$$

(C)

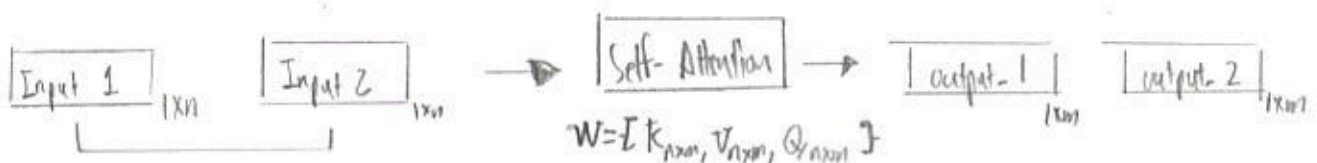
$$\hat{v}_1 = v_1 \odot \hat{S}_1$$

$$\hat{v}_2 = v_2 \odot \hat{S}_2$$

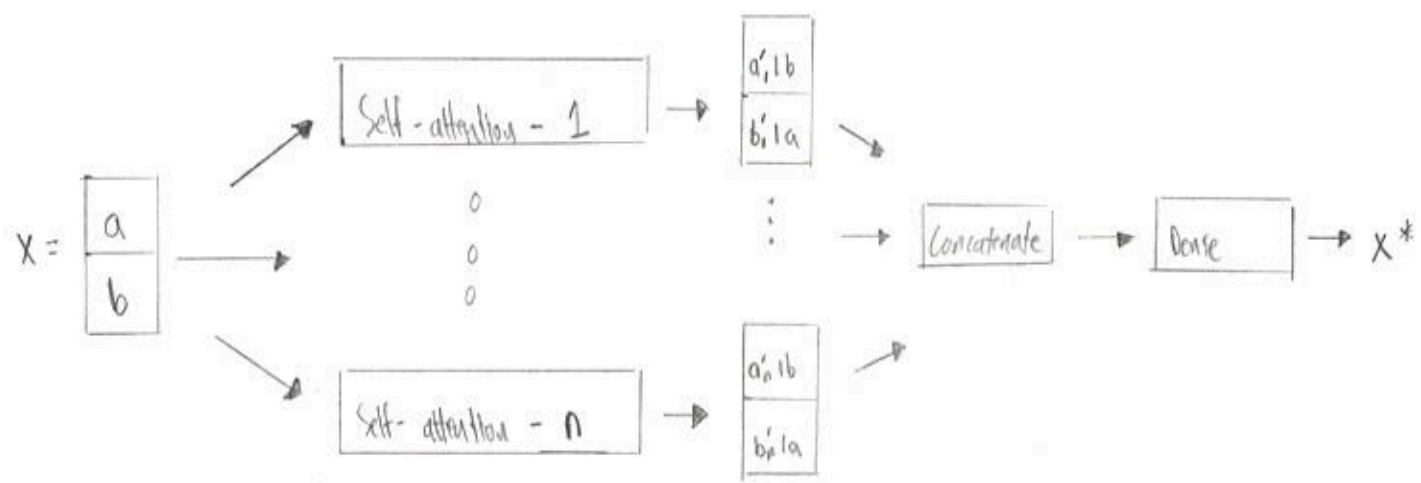
(D)

$$\text{Output 1}_{(1 \times m)} = \sum_{i=1}^2 \hat{v}_i$$

③

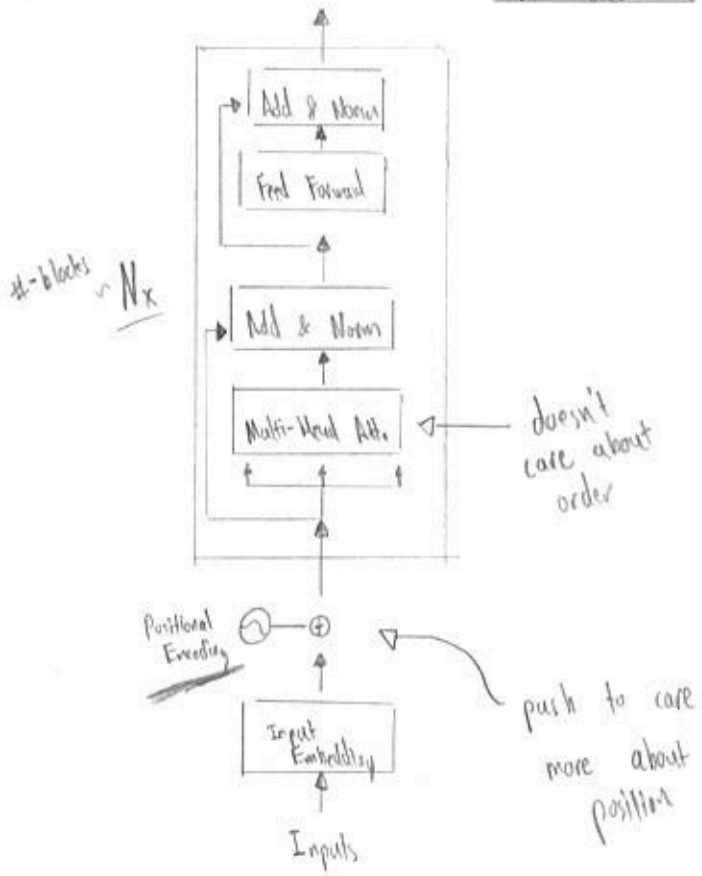


Multi-head attention



Transformers

Encoder



Recurrent Network vs MH-attention

- ~ RN can't be parallelized, MH-att can
- ~ RN not numerically stable