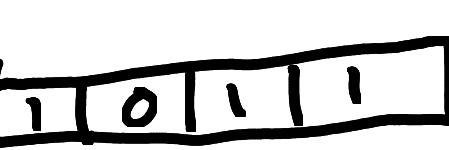


sequences

generate global dependences between inputs & outputs. A self attention module takes  $n$ -inputs & returns  $n$  outputs. It allows the inputs to interact with each other & find out where they should pay most attention.

Example:

Start with some embedding vector (feature vector) of dimension 4

(1,4)  , say we have 3 of those inputs.

$$X = \begin{bmatrix} [Input_1] \\ [Input_2] \\ [Input_3] \end{bmatrix} = \begin{bmatrix} [X_1] \\ [X_2] \\ [X_3] \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} \\ x_{21} & x_{22} & x_{23} & x_{24} \\ x_{31} & x_{32} & x_{33} & x_{34} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 2 & 0 & 2 \\ 1 & 1 & 1 & 1 \end{bmatrix}$$

Every one of the 3 inputs need 3 representations, key, query & value. let's say these representations have dimension 3 (1,3), our weight matrices must have shape (1,3) as  $\begin{array}{c} \text{input} \\ (1,4) \end{array} \times \begin{array}{c} \text{weight} \\ (4,3) \end{array} \rightarrow \begin{array}{c} \text{output} \\ (1,3) \end{array}$

That is to say:

$$Q = XW^q \quad K = XW^k \quad V = XW^v$$

For example shape of  $W^q$ ,  $W^k$ ,  $W^v \rightarrow (4, 3)$ , given  $X \rightarrow (1, 4)$ ,

I vs wanting a  $(1, 3)$  representation,

$$\text{I given: } W^q = \begin{bmatrix} 1 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 1 \end{bmatrix} \quad W^k = \begin{bmatrix} 0 & 0 & 1 \\ 1 & 1 & 0 \\ 0 & 1 & 0 \\ 1 & 1 & 0 \end{bmatrix} \quad W^v = \begin{bmatrix} 0 & 2 & 0 \\ 0 & 3 & 0 \\ 1 & 0 & 3 \\ 1 & 1 & 0 \end{bmatrix}$$

$(4, 3) \quad (4, 3) \quad (4, 3)$

$$\underset{(3, 4)}{X} \underset{(4, 3)}{W^q} = \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 2 & 0 & 2 \\ 1 & 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 1 \end{bmatrix} = \begin{bmatrix} [1 & 0 & 2] \\ [2 & 2 & 2] \\ [2 & 1 & 3] \end{bmatrix} = \begin{bmatrix} [Q_1] \\ [Q_2] \\ [Q_3] \end{bmatrix} = Q$$

$(3, 3)$

$$\underset{(3, 4)}{X} \underset{(4, 3)}{W^k} = \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 2 & 0 & 2 \\ 1 & 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 0 & 0 & 1 \\ 1 & 1 & 0 \\ 0 & 1 & 0 \\ 1 & 1 & 0 \end{bmatrix} = \begin{bmatrix} [0 & 1 & 1] \\ [4 & 4 & 0] \\ [2 & 3 & 1] \end{bmatrix} = \begin{bmatrix} [K_1] \\ [K_2] \\ [K_3] \end{bmatrix} = K$$

$(3, 3)$

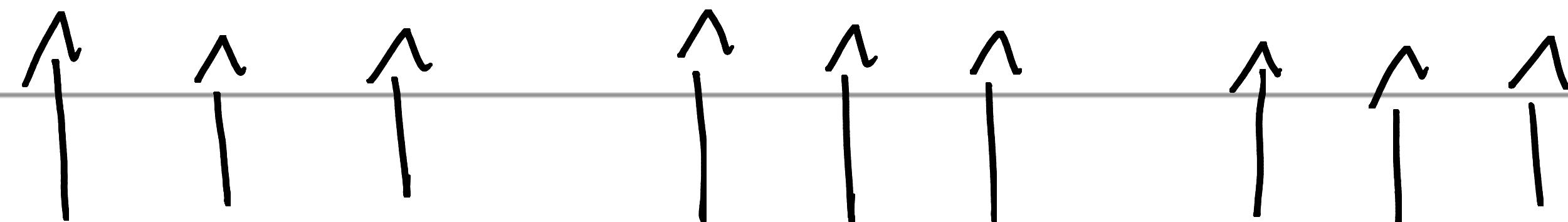
$$\underset{(3, 4)}{X} \underset{(4, 3)}{W^v} = \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 2 & 0 & 2 \\ 1 & 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 0 & 2 & 0 \\ 0 & 3 & 0 \\ 1 & 0 & 3 \\ 1 & 1 & 0 \end{bmatrix} = \begin{bmatrix} [1 & 2 & 3] \\ [2 & 8 & 0] \\ [2 & 6 & 3] \end{bmatrix} = \begin{bmatrix} [V_1] \\ [V_2] \\ [V_3] \end{bmatrix} = V$$

$(3, 3)$

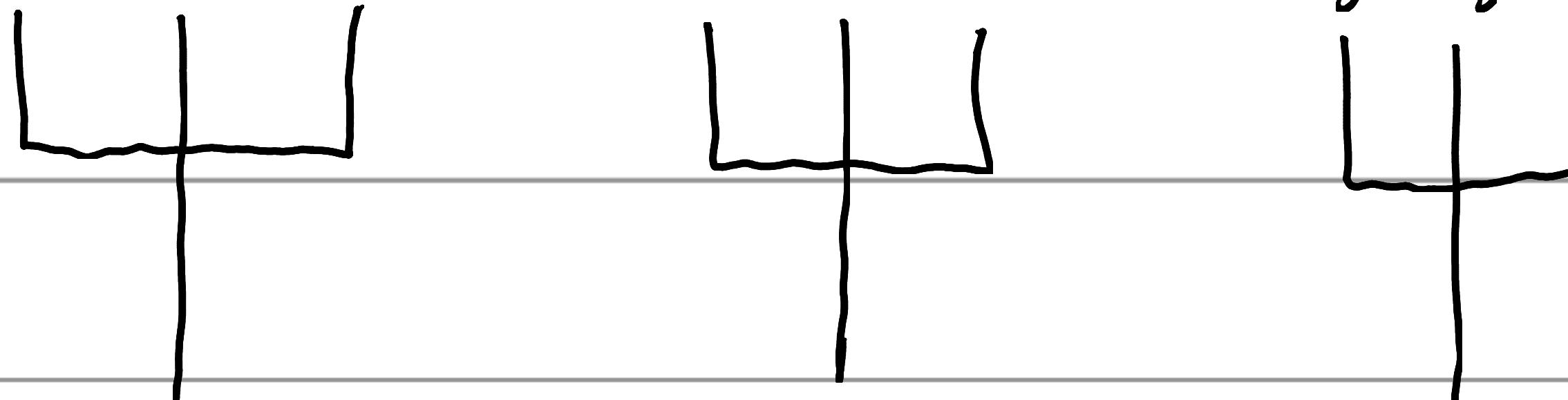
Or diagrammatically

$$[102] \ [011] \ [123] \ [221] \ [440] \ [280] \ [213] \ [231] \ [263] \ [3 \times (1,3)]$$

$$Q_1 \ K_1 \ V_1 \quad Q_2 \ K_2 \ V_2 \quad Q_3 \ K_3 \ V_3 \quad q = xw^T \quad k = xw^k \quad v = xw^v$$



$$W_1^q \ W_1^k \ W_1^v \quad W_2^q \ W_2^k \ W_2^v \quad W_3^q \ W_3^k \ W_3^v \quad [3 \times (4,3)]$$



$$[1010] \ [0202] \ [1111] \cdot [3 \times (1,4)]$$

Input 1

Input 2

Input 3

$x_1$

$x_2$

$x_3$

$$(1,4) \times (4,3) \rightarrow (1,3)$$

| Input                                  | Weights | Representations |
|--|---------|-----------------|
| $(3,4) \times (4,3) \rightarrow (3,3)$ |         |                 |

Input      weight      representation

$x_3$  for  $Q, K, T$

To obtain attention scores we take the dot product between each input query & all keys (including itself), we obtain 3 attention scores for each input (for this example)

$$Q = \begin{bmatrix} Q_1 \\ Q_2 \\ Q_3 \end{bmatrix} = \begin{bmatrix} q_{11} & q_{12} & q_{13} \\ q_{21} & q_{22} & q_{23} \\ q_{31} & q_{32} & q_{33} \end{bmatrix}$$

$$K = \begin{bmatrix} K_1 \\ K_2 \\ K_3 \end{bmatrix} = \begin{bmatrix} k_{11} & k_{12} & k_{13} \\ k_{21} & k_{22} & k_{23} \\ k_{31} & k_{32} & k_{33} \end{bmatrix}$$

$$K^T = \begin{bmatrix} [k_1^T] & [k_2^T] & [k_3^T] \end{bmatrix} = \begin{bmatrix} k_{11} & k_{21} & k_{31} \\ k_{12} & k_{22} & k_{32} \\ k_{13} & k_{23} & k_{33} \end{bmatrix}$$

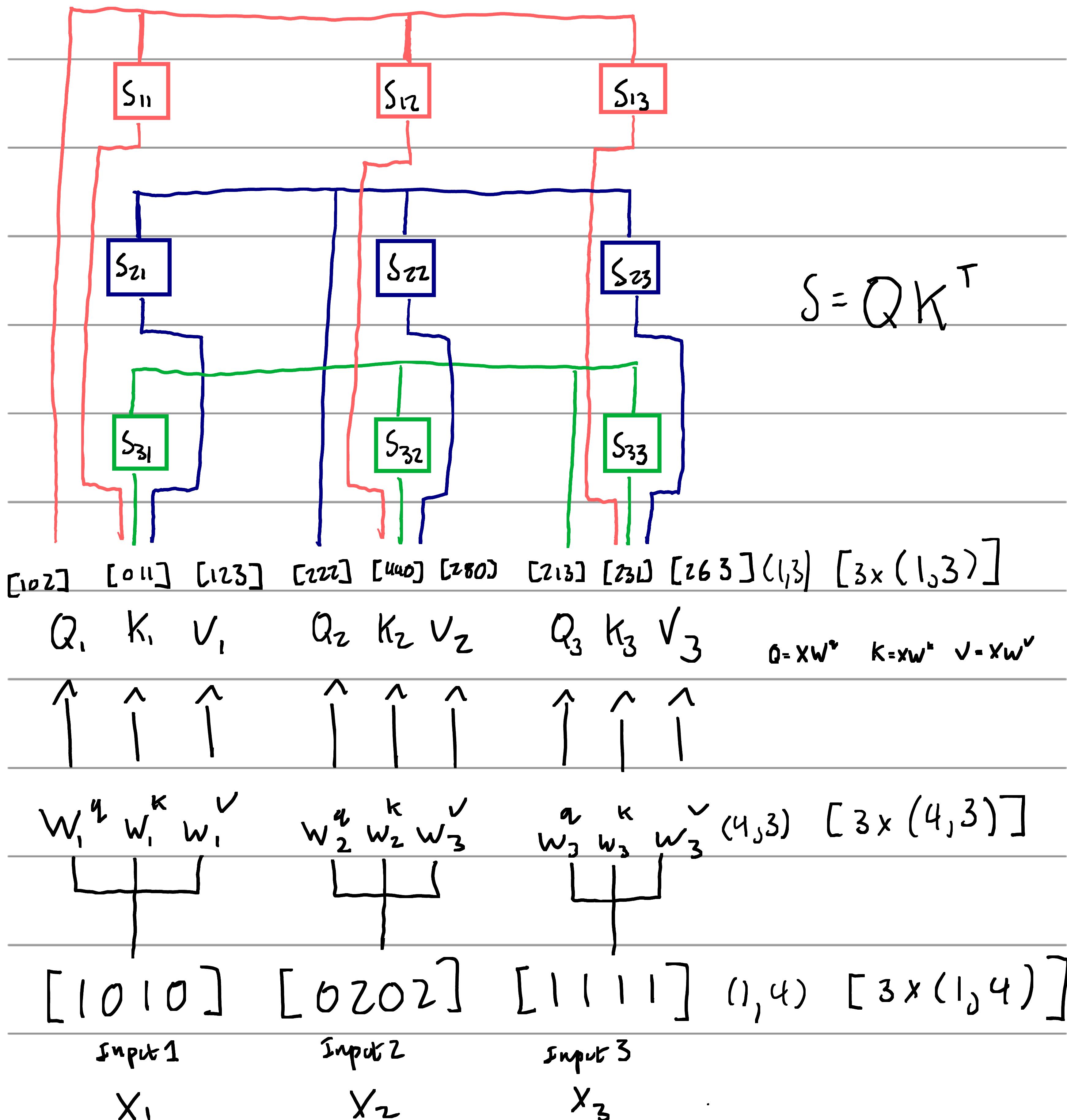
Similarity score  $S = Q K^T$

$$\begin{matrix} (3, 3) & (3, 3) & (3, 3) \end{matrix}$$

$$S = \begin{bmatrix} Q_1 \\ Q_2 \\ Q_3 \end{bmatrix} \begin{bmatrix} [k_1^T] & [k_2^T] & [k_3^T] \end{bmatrix}$$

$$S = \begin{bmatrix} Q_1 K_1^T & Q_1 K_2^T & Q_1 K_3^T \\ Q_2 K_1^T & Q_2 K_2^T & Q_2 K_3^T \\ Q_3 K_1^T & Q_3 K_2^T & Q_3 K_3^T \end{bmatrix} = \begin{bmatrix} S_{11} & S_{12} & S_{13} \\ S_{21} & S_{22} & S_{23} \\ S_{31} & S_{32} & S_{33} \end{bmatrix}$$

Or diagrammatically:



$$Y = \text{attention}(Q, K, V) \quad d_x \text{ is the embedding vector dimension}$$

$$= \text{softmax} \left( \frac{QK^T}{\sqrt{d_x}} \right) \quad \text{for this example } d_x = 4$$

$$= \text{softmax} \left( \frac{S}{\sqrt{d_x}} \right) \quad S = \sum_{i=1}^n \frac{\exp(s_i)}{\sum_j \exp(s_j)}$$

$$\begin{matrix} S'' \\ \uparrow \\ [3 \times 3] \end{matrix} \quad S'' = \text{softmax} \left( \frac{S}{\sqrt{d_x}} \right)$$

$$\text{Softmax} \quad [3 \times 3]$$

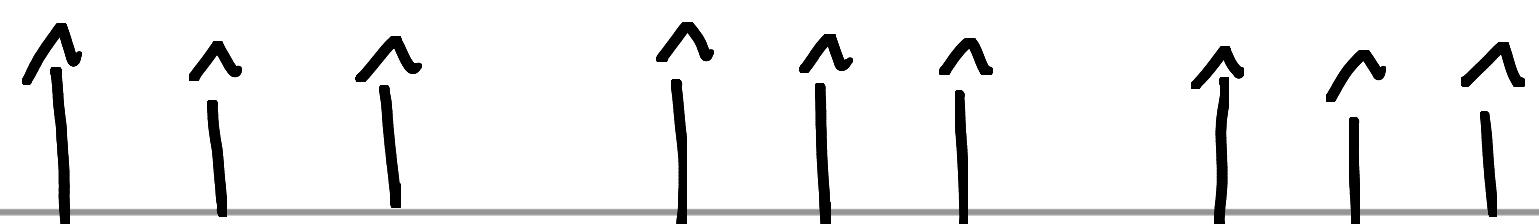
$$\begin{matrix} S' \\ \uparrow \\ [3 \times 3] \end{matrix} \quad S' = \frac{S}{\sqrt{d_x}}$$

$$\text{Normalize by } \sqrt{d_x}$$

$$S = QK^T$$

$$[102] [011] [123] [222] [440] [280] [213] [231] [263] (1,3) \quad [3 \times [3 \times (1,3)]] = [3 \times (3 \times 3)]$$

$$Q_1 \quad K_1 \quad V_1 \quad Q_2 \quad K_2 \quad V_2 \quad Q_3 \quad K_3 \quad V_3 \quad Q = xw^T \quad K = xw^T \quad V = xw^T$$



$$w_1^q \quad w_1^k \quad w_1^v$$

$$w_2^q \quad w_2^k \quad w_2^v$$

$$w_3^q \quad w_3^k \quad w_3^v \quad (4,3) \quad [3 \times (4,3)]$$

$$[1010]$$

$$[0202]$$

$$[1111]$$

$$(1,4) \quad [3 \times (1,4)]$$

$$x_1$$

$$x_2$$

$$x_3$$

