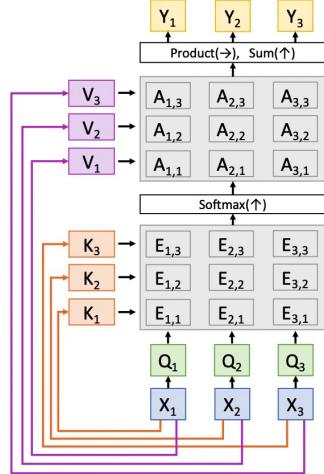
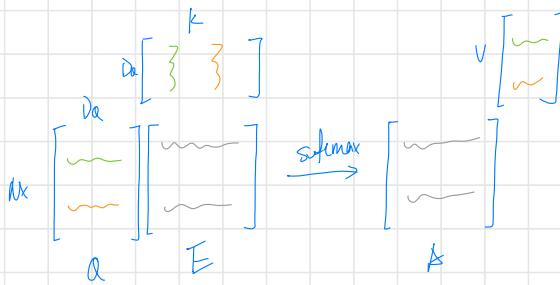


## • Self - Attention

input vectors  $X: N_x \times D_x$   
 key matrix  $W_k: D_x \times D_k$   
 value matrix  $W_v: D_x \times D_v$   
 Query matrix  $W_q: D_x \times D_q$

$N_x$  is the sequence length, not batch size

$$\begin{array}{l} \text{query: } Q = X \cdot W_q \\ \text{key: } K = X \cdot W_k \\ \text{value: } V = X \cdot W_v \end{array} \quad \left. \begin{array}{l} \text{similarities } E = QK^T / \sqrt{D_k} \\ A = \text{softmax}(E, \text{dim}=1) \end{array} \right.$$



self - attention is permutation equivariant

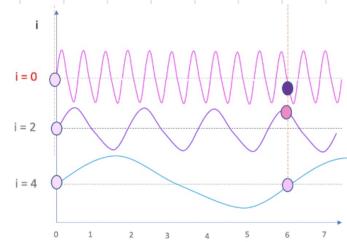
## • Positional Encoding

Requirement: the encoding is deterministic  
 output a unique encoding for each time-step  
 distance between any two timesteps should be consistent across sequences of different lengths

let  $E_t \in \mathbb{R}^d$  be the positional encoding at time-step  $t$

$$E_{2t}[t] = \sin\left(\frac{t}{10000^{2k/d}}\right)$$

$$E_{2t+1}[t] = \cos\left(\frac{t}{10000^{2k/d}}\right)$$



$$\begin{bmatrix} E_{2k}[t+\phi] \\ E_{2k+1}[t+\phi] \end{bmatrix} = \begin{bmatrix} \cos\left(\frac{t+\phi}{10000^{2/\alpha}}\right) \\ \sin\left(\frac{t+\phi}{10000^{2/\alpha}}\right) \end{bmatrix} = \begin{bmatrix} \cos\left(\frac{\phi}{10000^{2/\alpha}}\right) & -\sin\left(\frac{\phi}{10000^{2/\alpha}}\right) \\ \sin\left(\frac{\phi}{10000^{2/\alpha}}\right) & \cos\left(\frac{\phi}{10000^{2/\alpha}}\right) \end{bmatrix} \cdot \begin{bmatrix} \cos\left(\frac{t}{10000^{2/\alpha}}\right) \\ \sin\left(\frac{t}{10000^{2/\alpha}}\right) \end{bmatrix}$$

