



# RNN & Transformer

黃志勝 (Tommy Huang)

義隆電子 人工智慧研發部

國立陽明交通大學 AI 學院 合聘助理教授

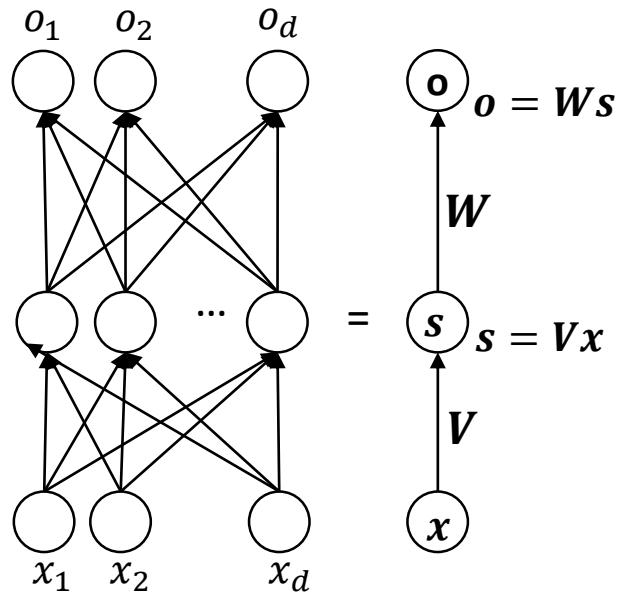
國立台北科技大學 電資學院合聘助理教授



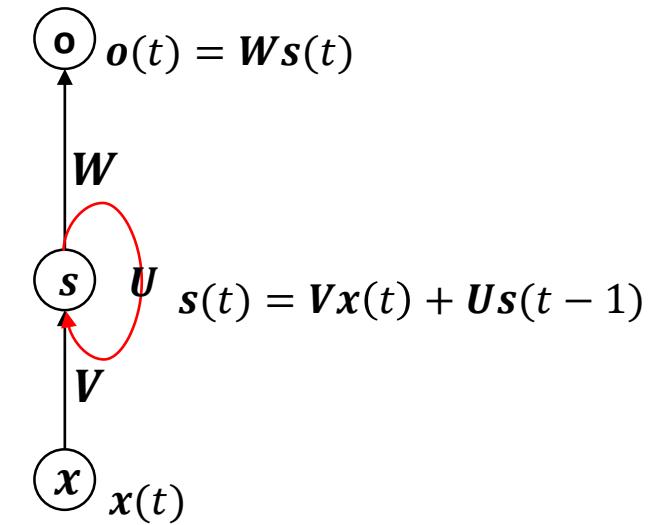
# Recurrent Neural Network (RNN)

- RNN和一般MLP的差異就是引進了具有記憶功能的「狀態(State)」

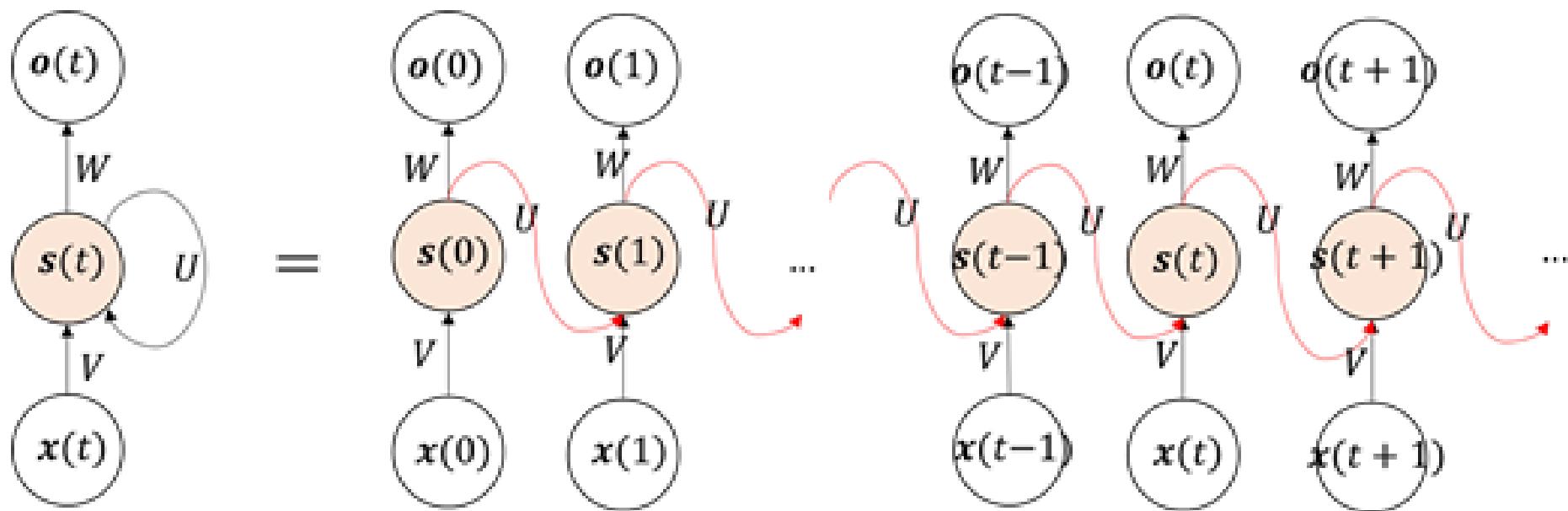
MLP



RNN

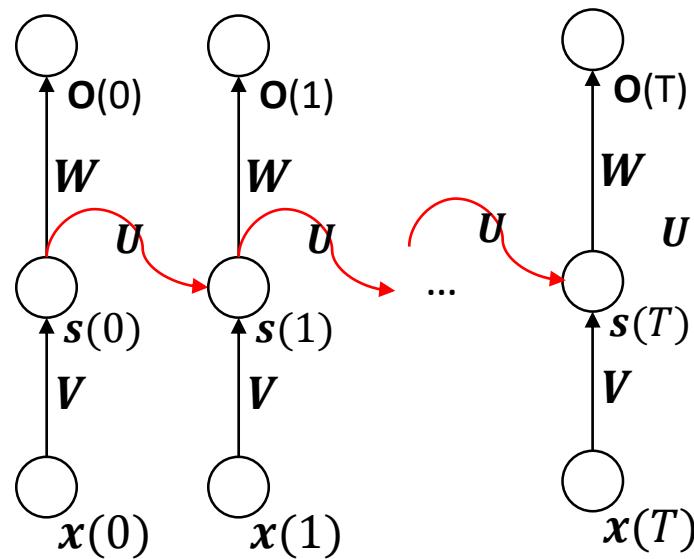


# Recurrent Neural Network (RNN)

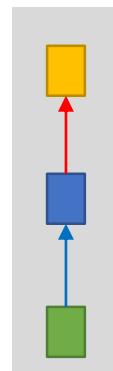


# RNN

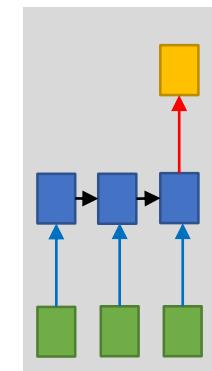
many to many



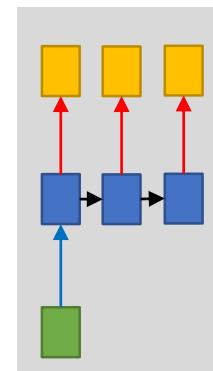
one to one



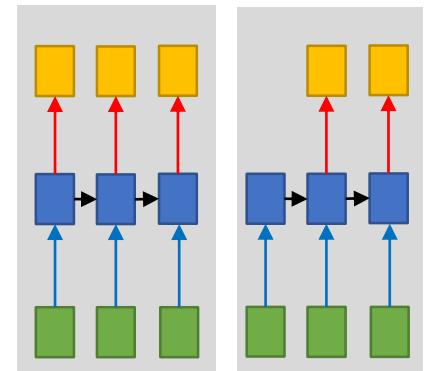
many to one



one to many



many to many

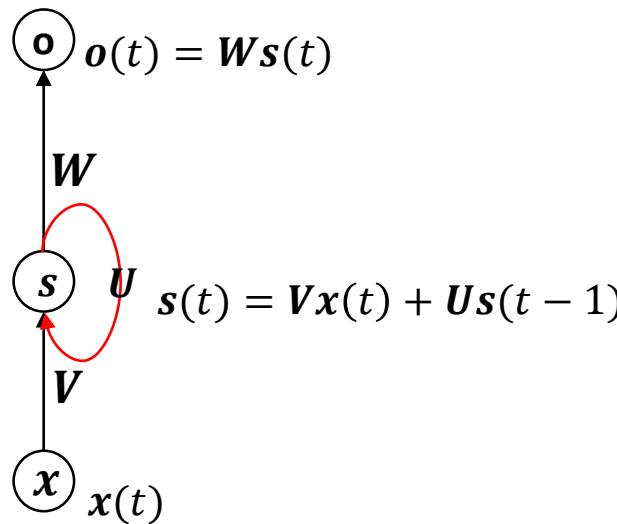


在RNN訓練期間，因為梯度會傳遞到最前一層做乘積，在前面課程提過層數一多可能會發生梯度消失和爆炸的問題，在RNN也是，傳遞不但跟層數相關和時間也先關，可以想像時間等於層數的概念，因此在RNN是很可能發生梯度消失/爆炸。

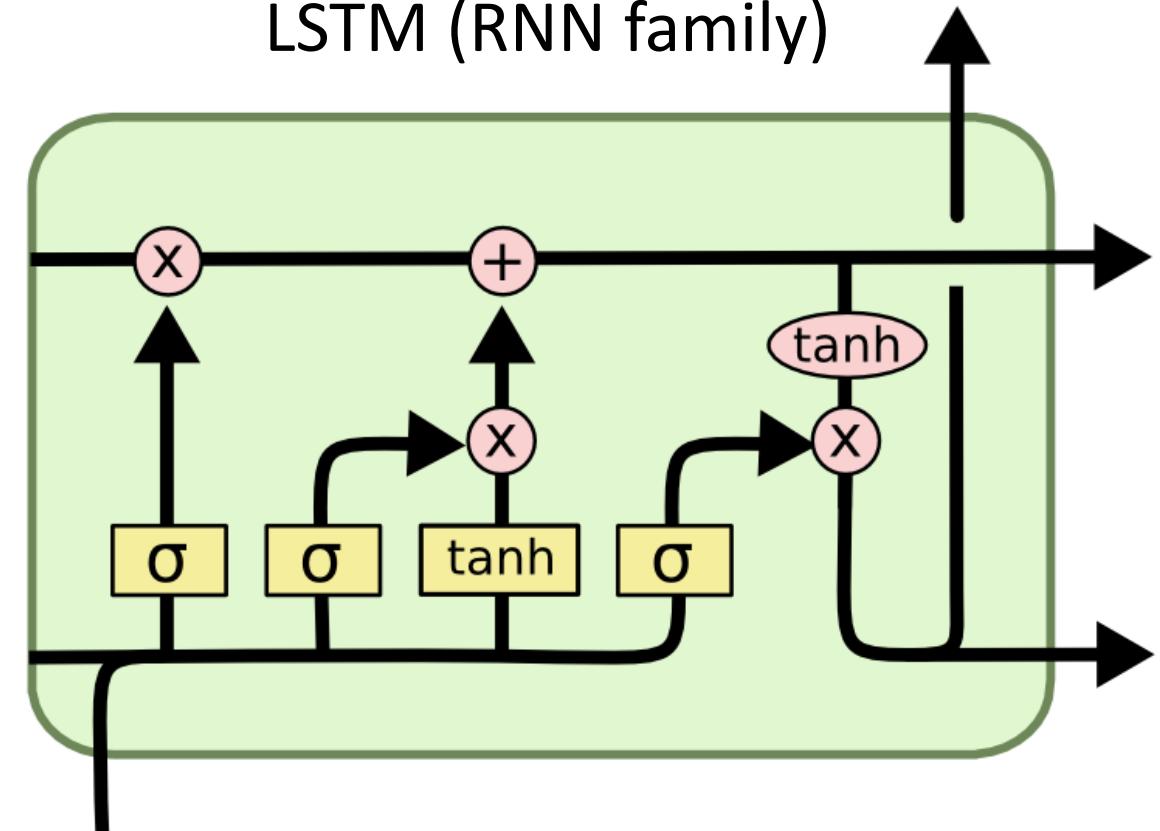


# Long Short Term Memory

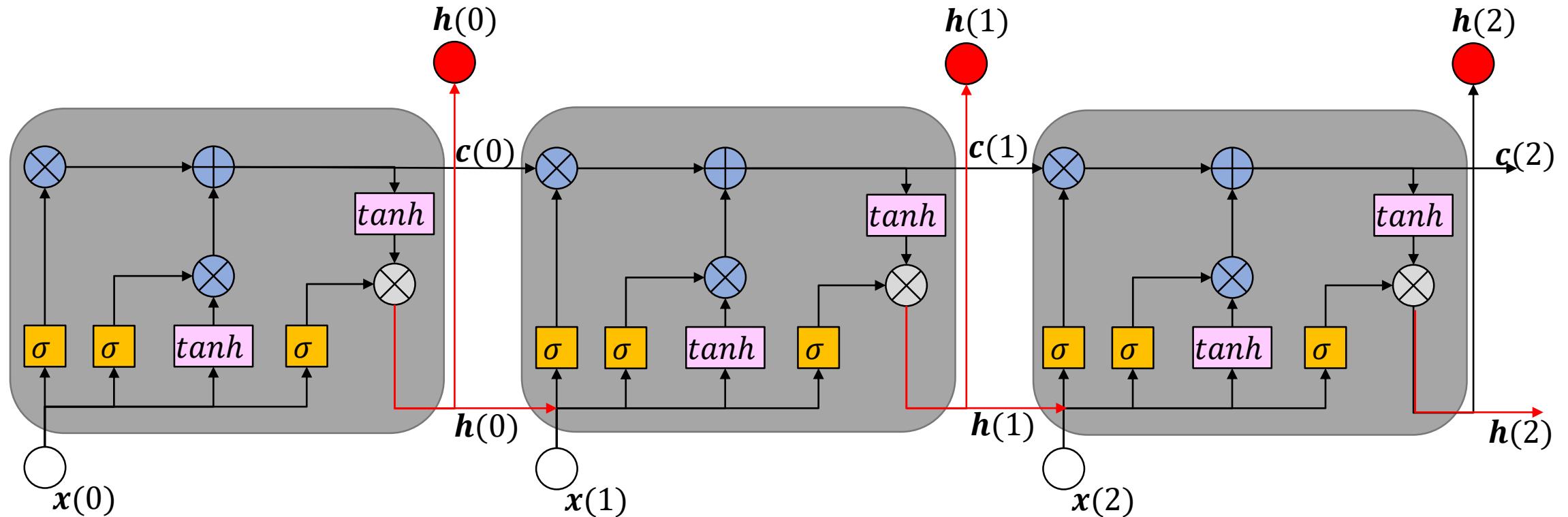
RNN



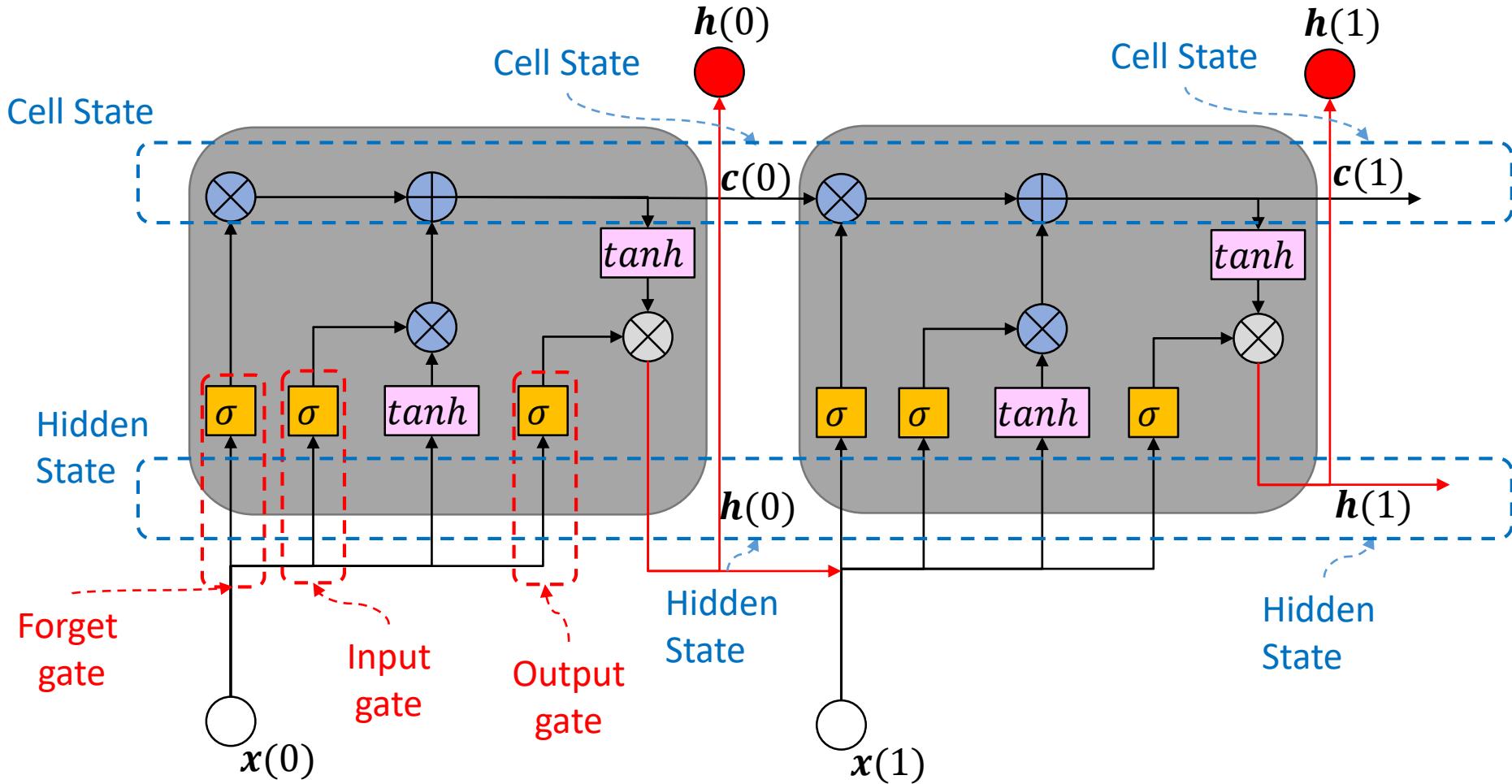
LSTM (RNN family)



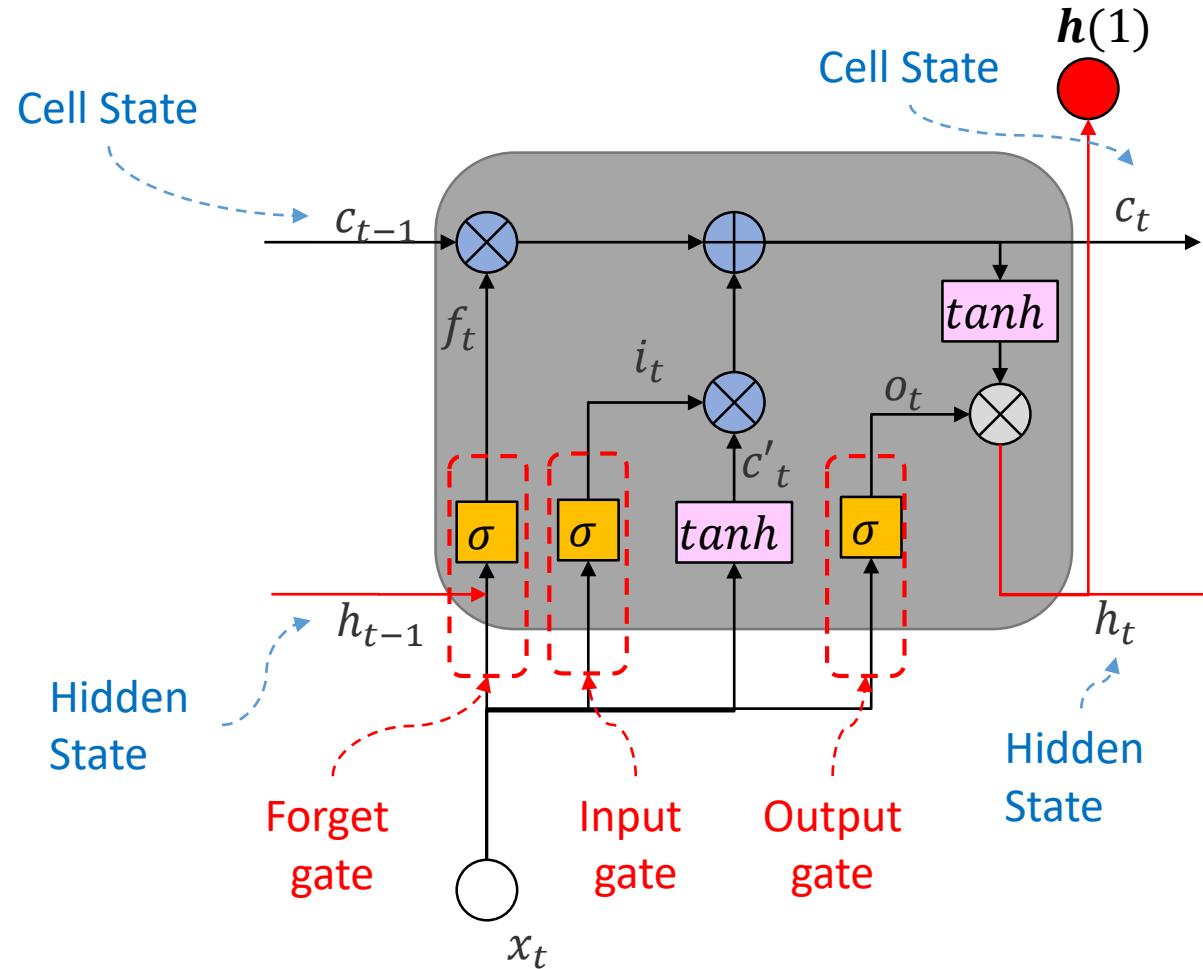
# Long Short Term Memory



# Long Short Term Memory



# Long Short Term Memory



forget gate

$$f_t = \sigma(W_f x_t + U_f h_{t-1} + b_f)$$

Input gate

$$i_t = \sigma(W_i x_t + U_i h_{t-1} + b_i)$$

Cell State

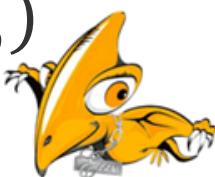
$$c'_t = \tanh(W_c x_t + U_c h_{t-1})$$

$$c_t = (f_t c_{t-1} + i_t c'_t)$$

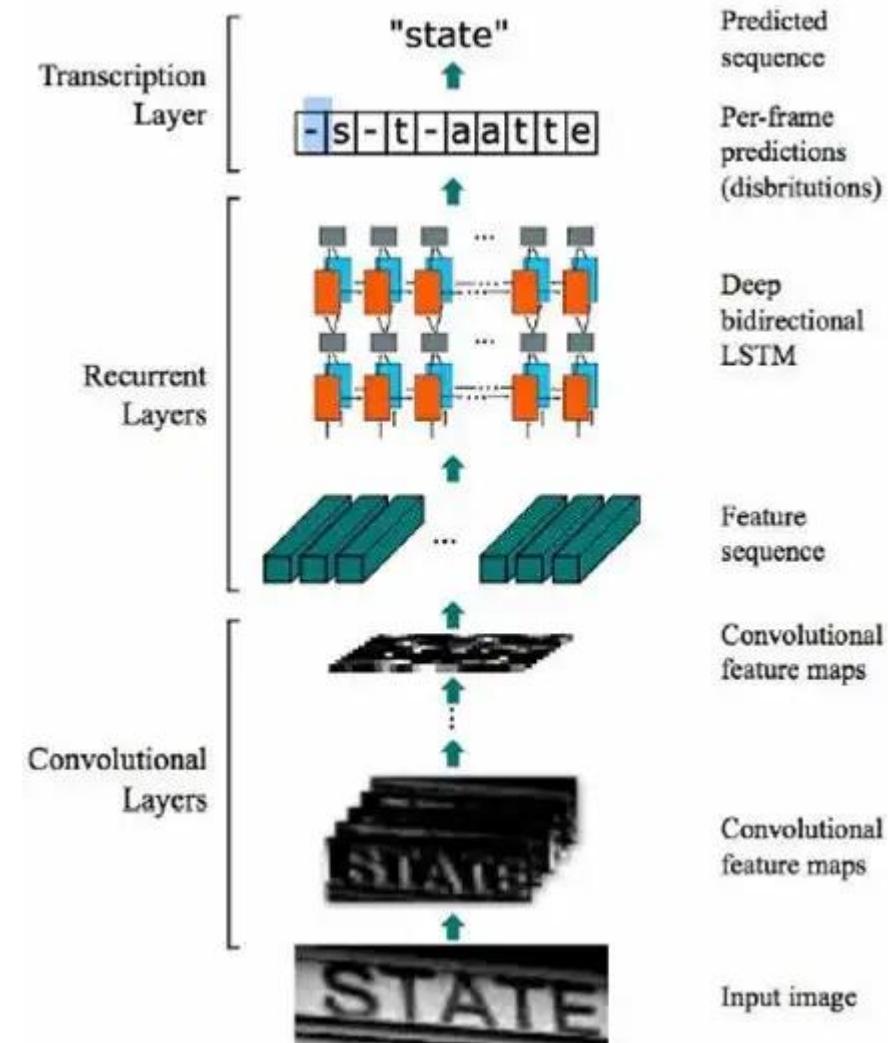
Output gate

$$o_t = \sigma(W_o x_t + U_o h_{t-1} + b_o)$$

$$h_t = o_t \tanh(c_t)$$



# Convolutional Recurrent Neural Network)

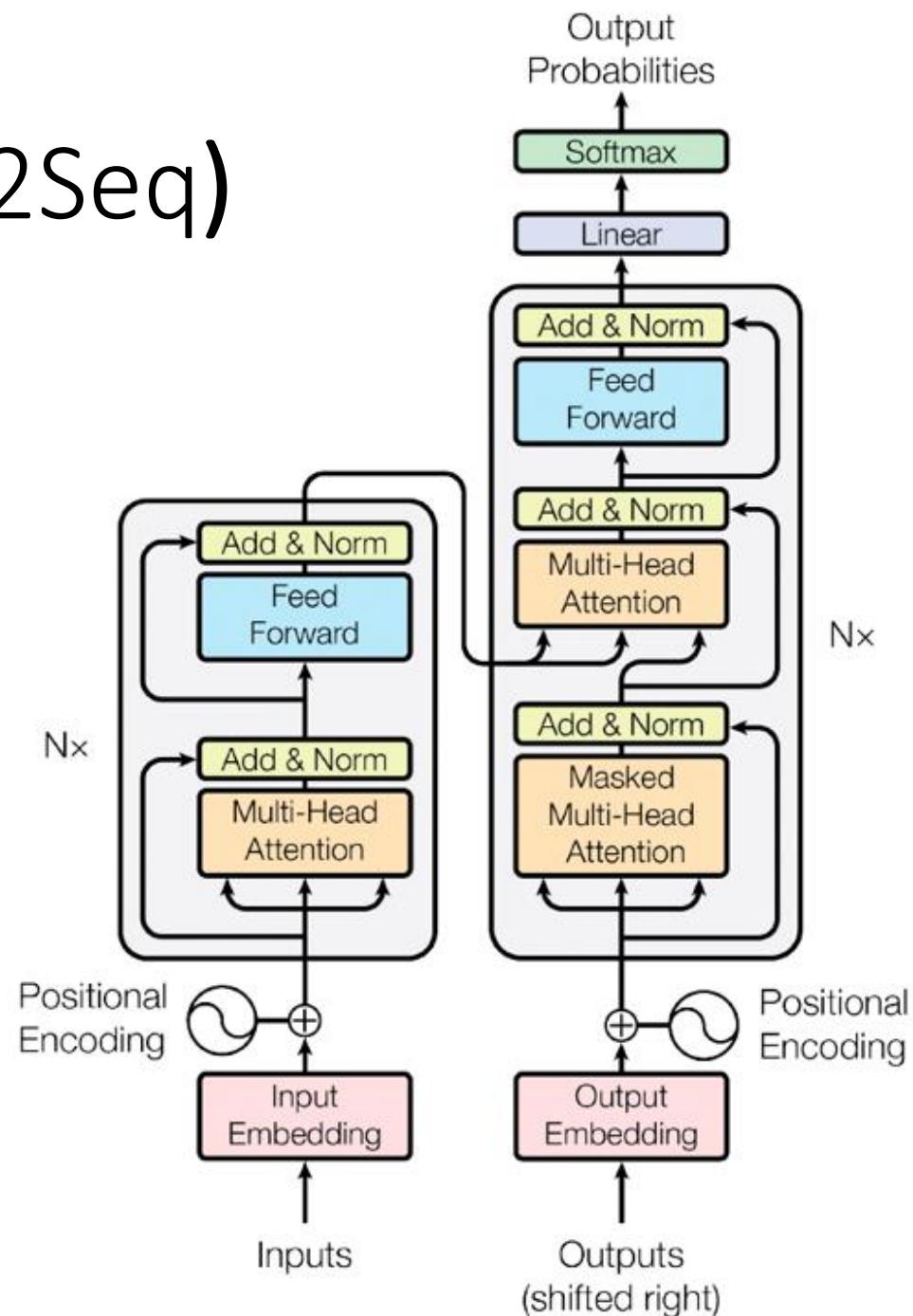
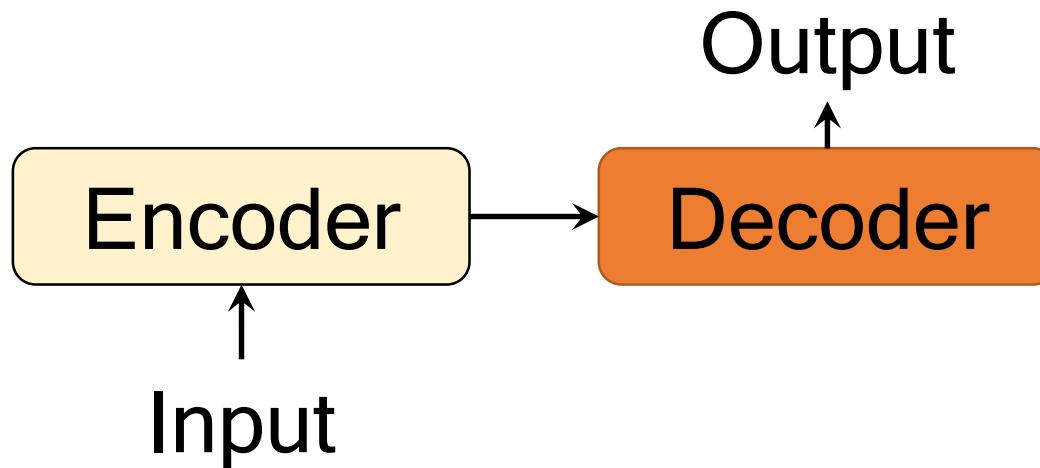


來源: <https://github.com/qjadud1994/CRNN-Keras>

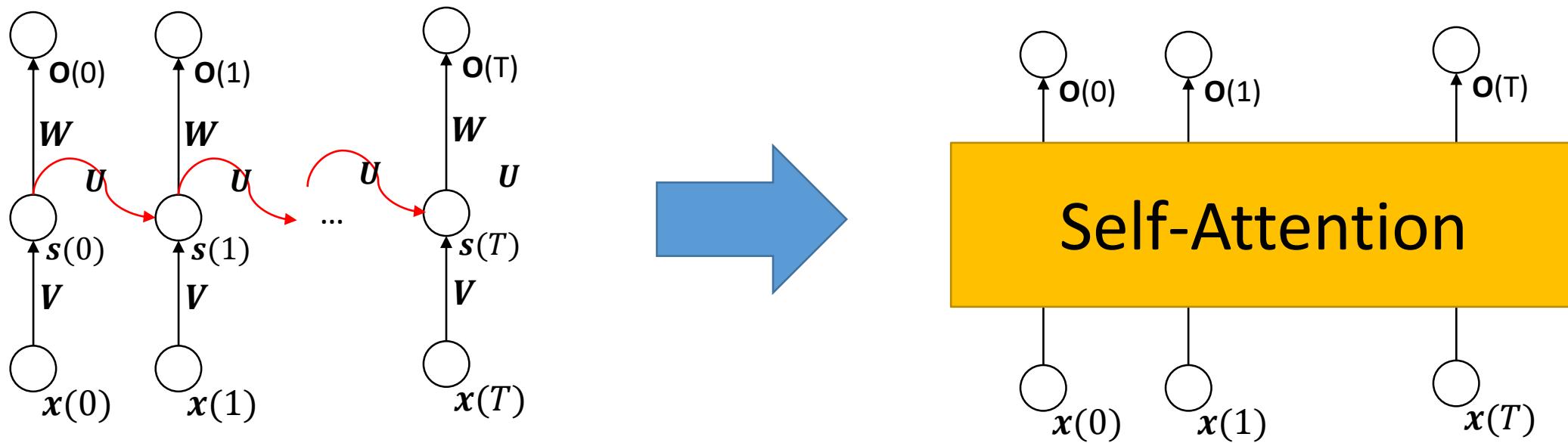


# Transformer (Seq2Seq)

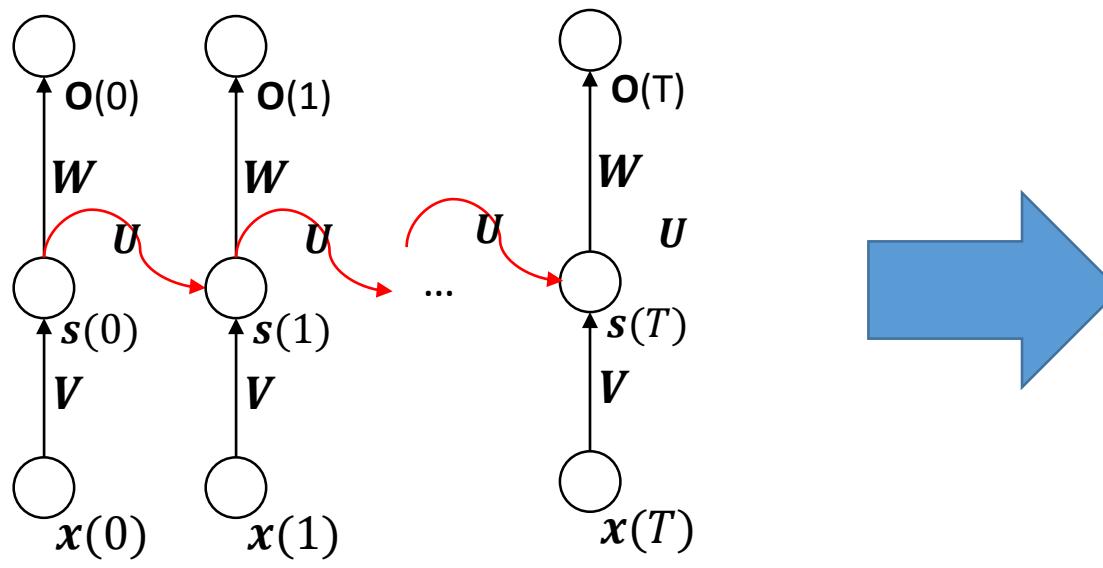
- Sequence-to-sequence: Seq2Seq
- Encoder: takes the input sequence and maps it into a higher dimensional space ( $n$ -dimensional vector).
- That abstract vector is fed into the Decoder which turns it into an output sequence.



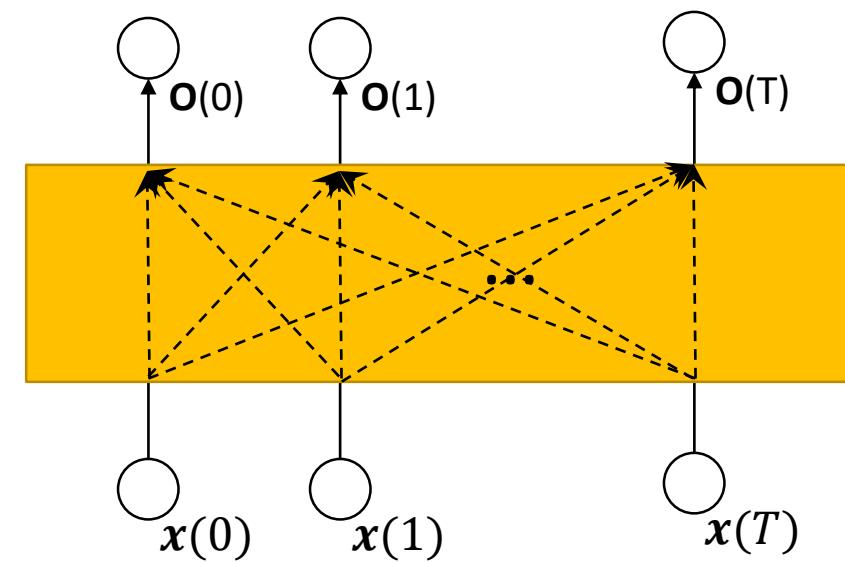
# Transformer: Self-Attention



# Transformer: Self-Attention



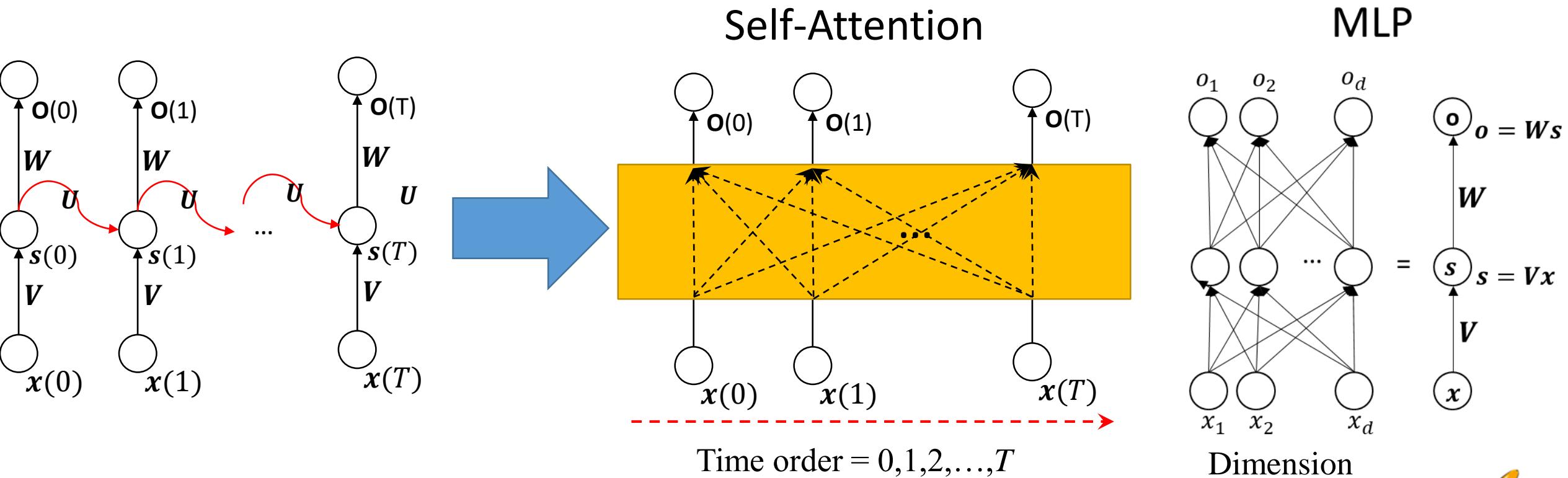
Self-Attention



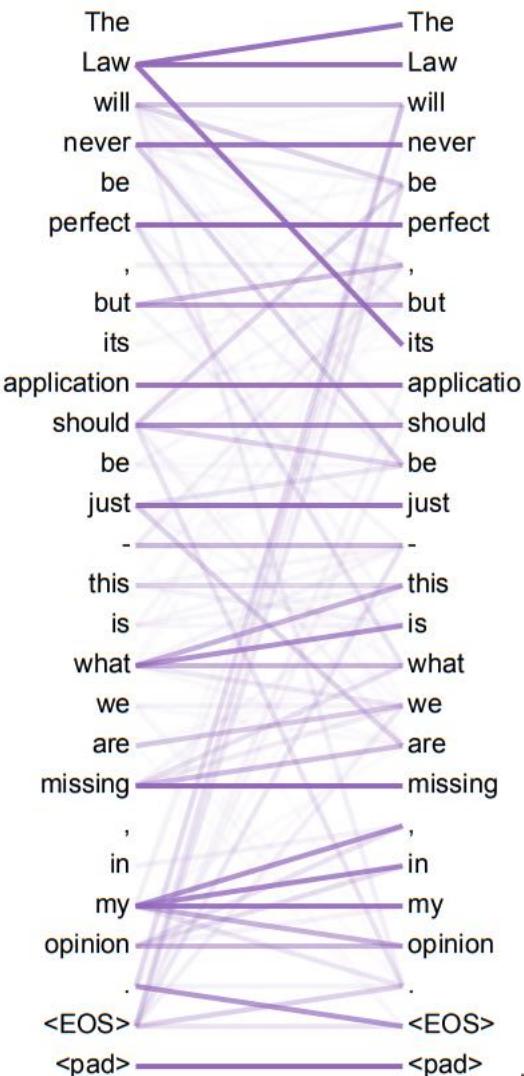
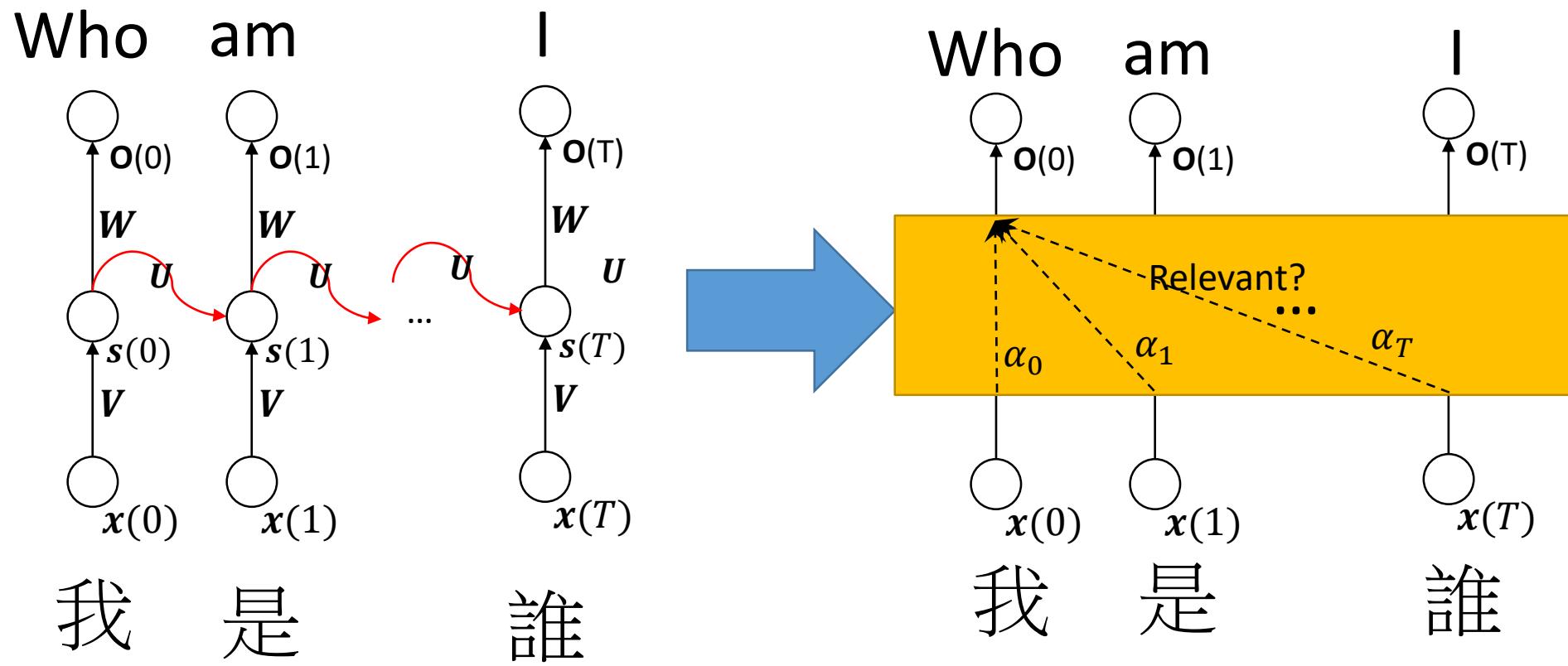
MLP?



# Transformer: Self-Attention



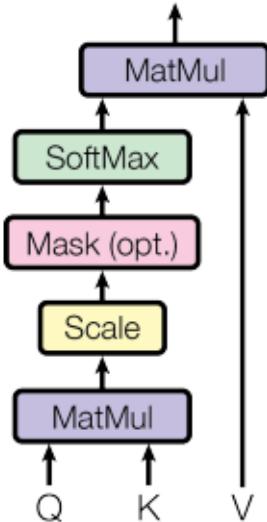
# Transformer: Self-Attention



# Transformer: Self-Attention

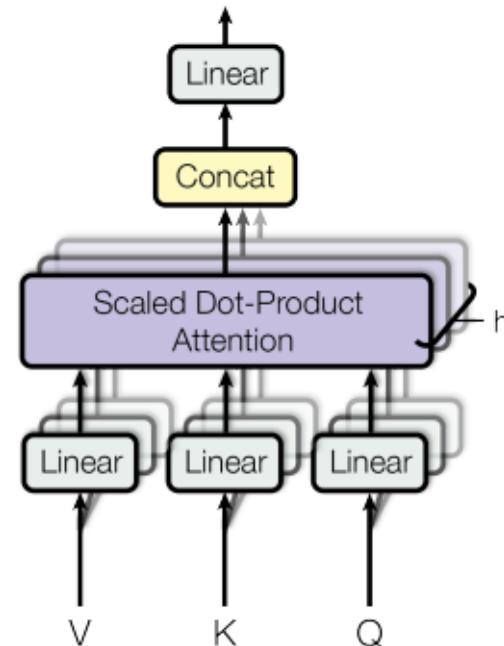
## Single-Head Attention

Scaled Dot-Product Attention

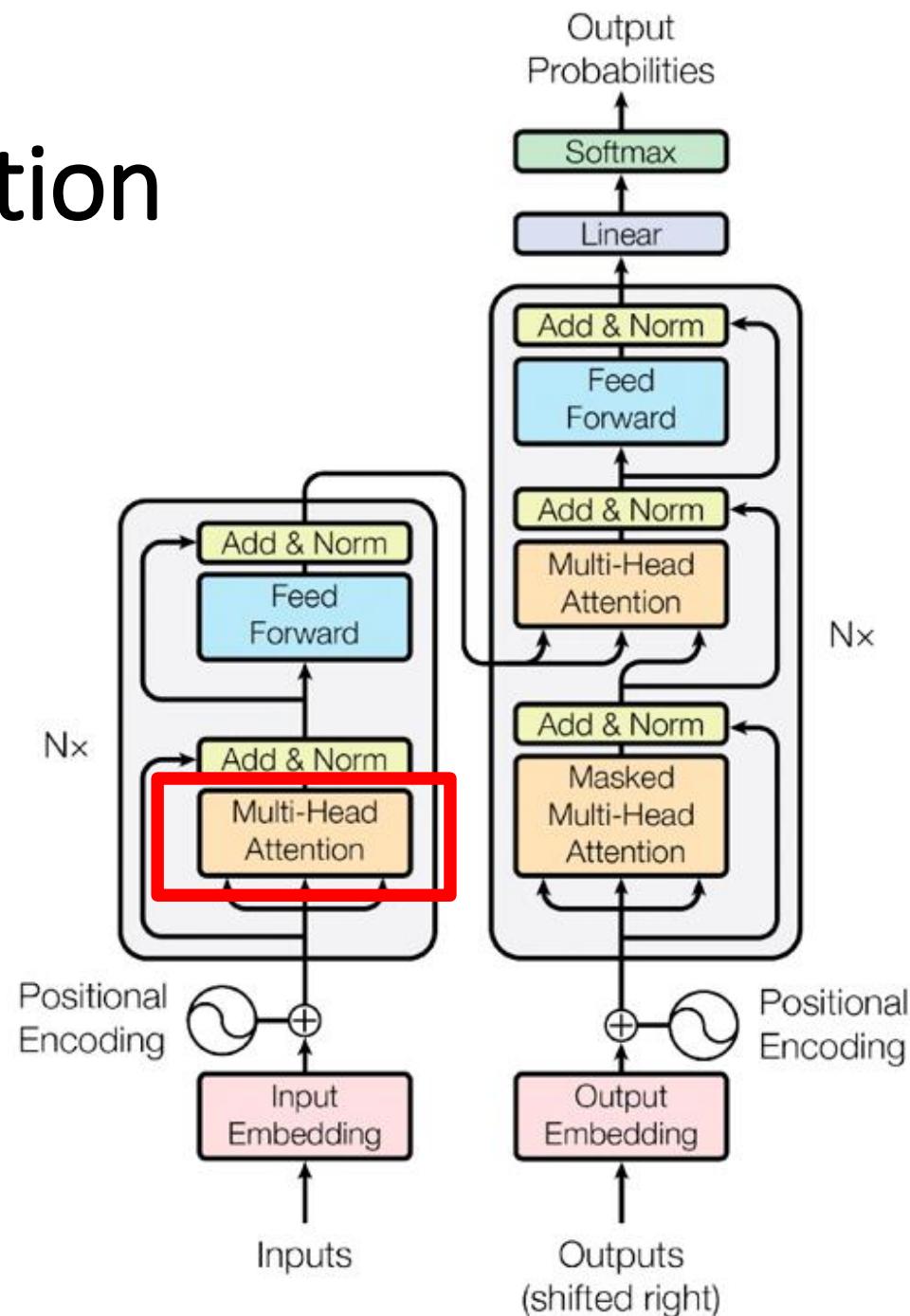


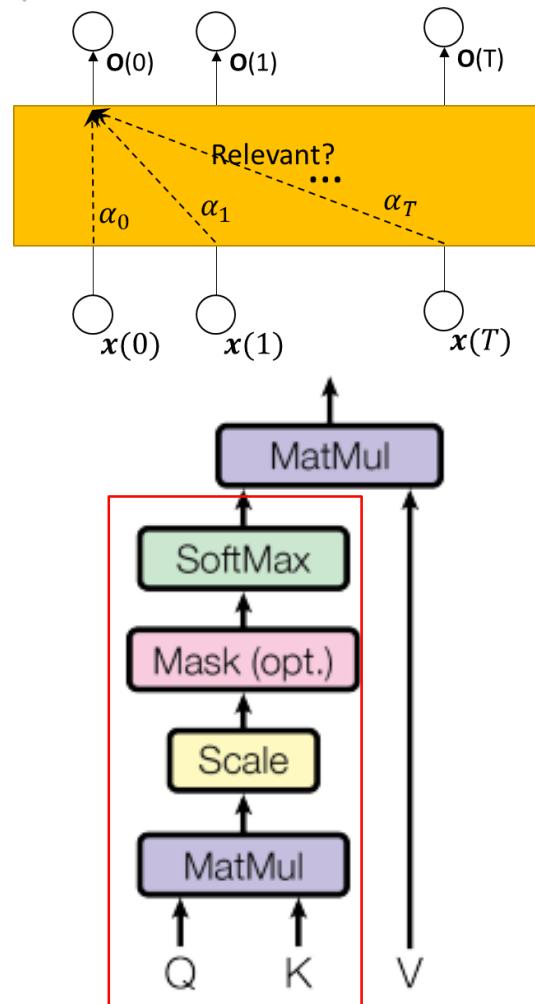
## Multi-Head Attention

Multi-Head Attention



Reference: Attention Is All You Need (2017)



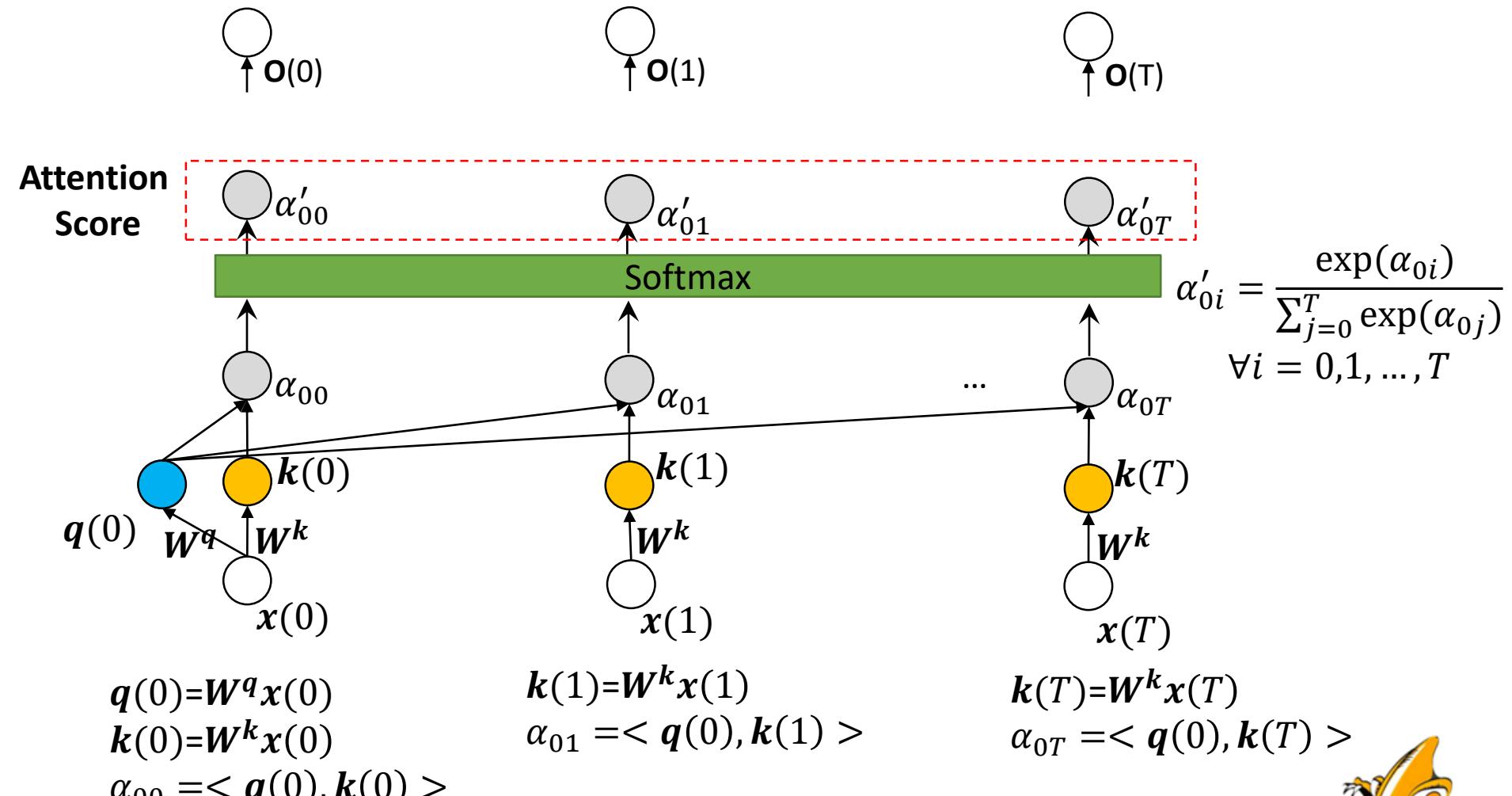


Q: Query

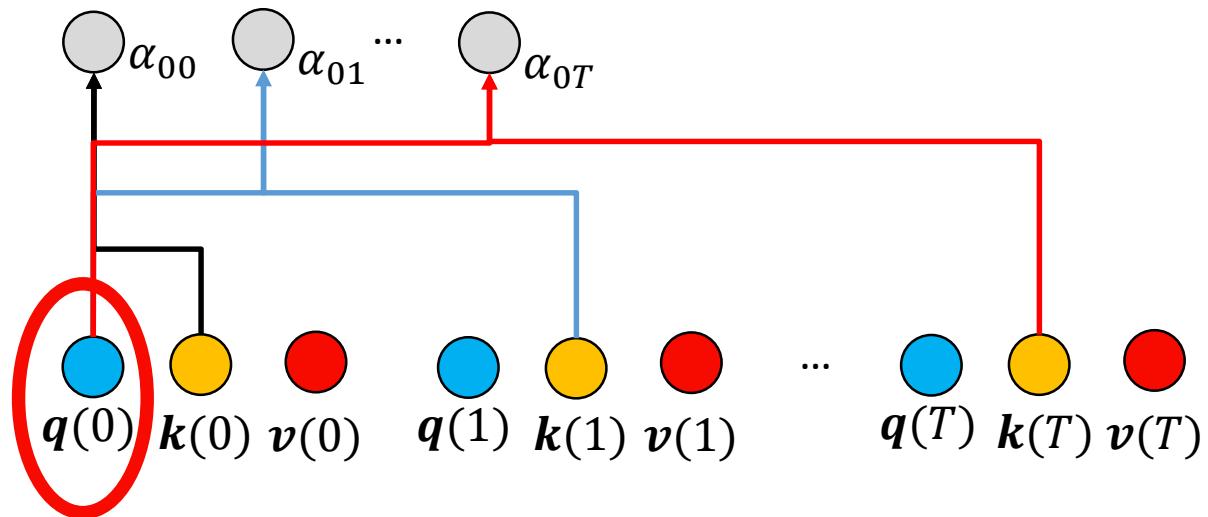
K: Key

V: Value

# Single-Head Self-Attention



# Attention Score



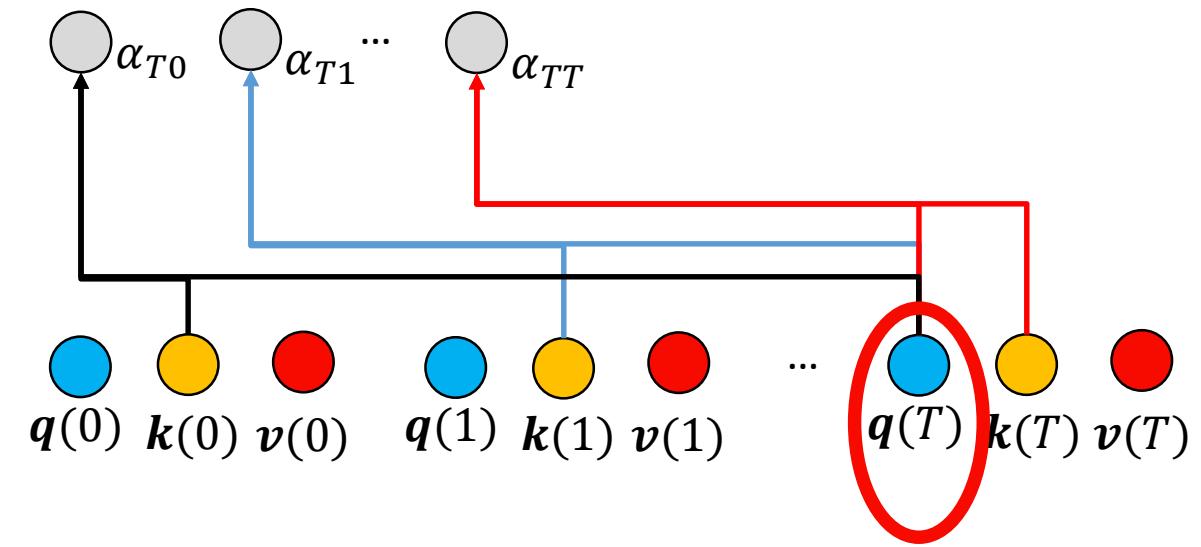
$$\alpha_{00} = \langle \mathbf{q}(0), \mathbf{k}(0) \rangle$$

$$\alpha_{01} = \langle \mathbf{q}(0), \mathbf{k}(1) \rangle$$

...

$$\alpha_{0T} = \langle \mathbf{q}(0), \mathbf{k}(T) \rangle$$

Attention Score at  $q(0)$



$$\alpha_{T0} = \langle \mathbf{q}(T), \mathbf{k}(0) \rangle$$

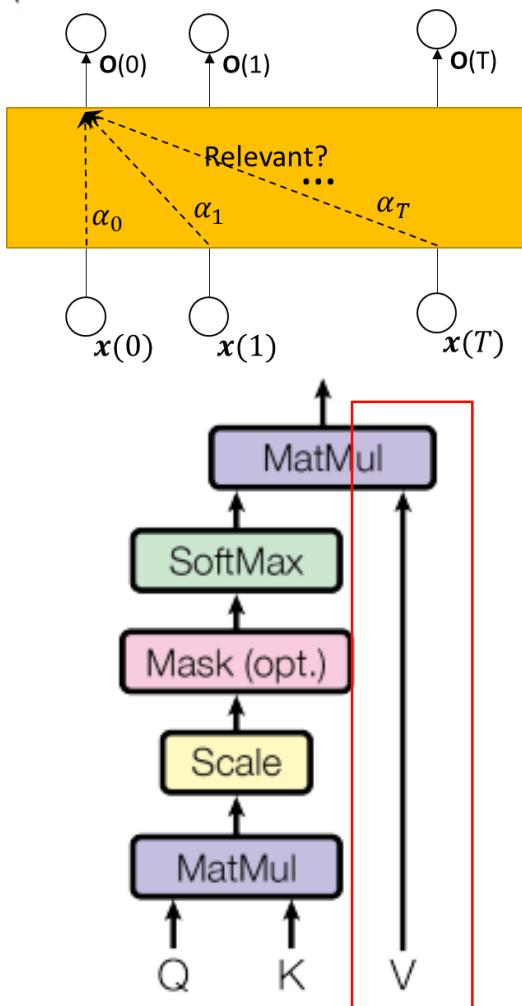
$$\alpha_{T1} = \langle \mathbf{q}(T), \mathbf{k}(1) \rangle$$

...

$$\alpha_{TT} = \langle \mathbf{q}(T), \mathbf{k}(T) \rangle$$

Attention Score at  $q(T)$

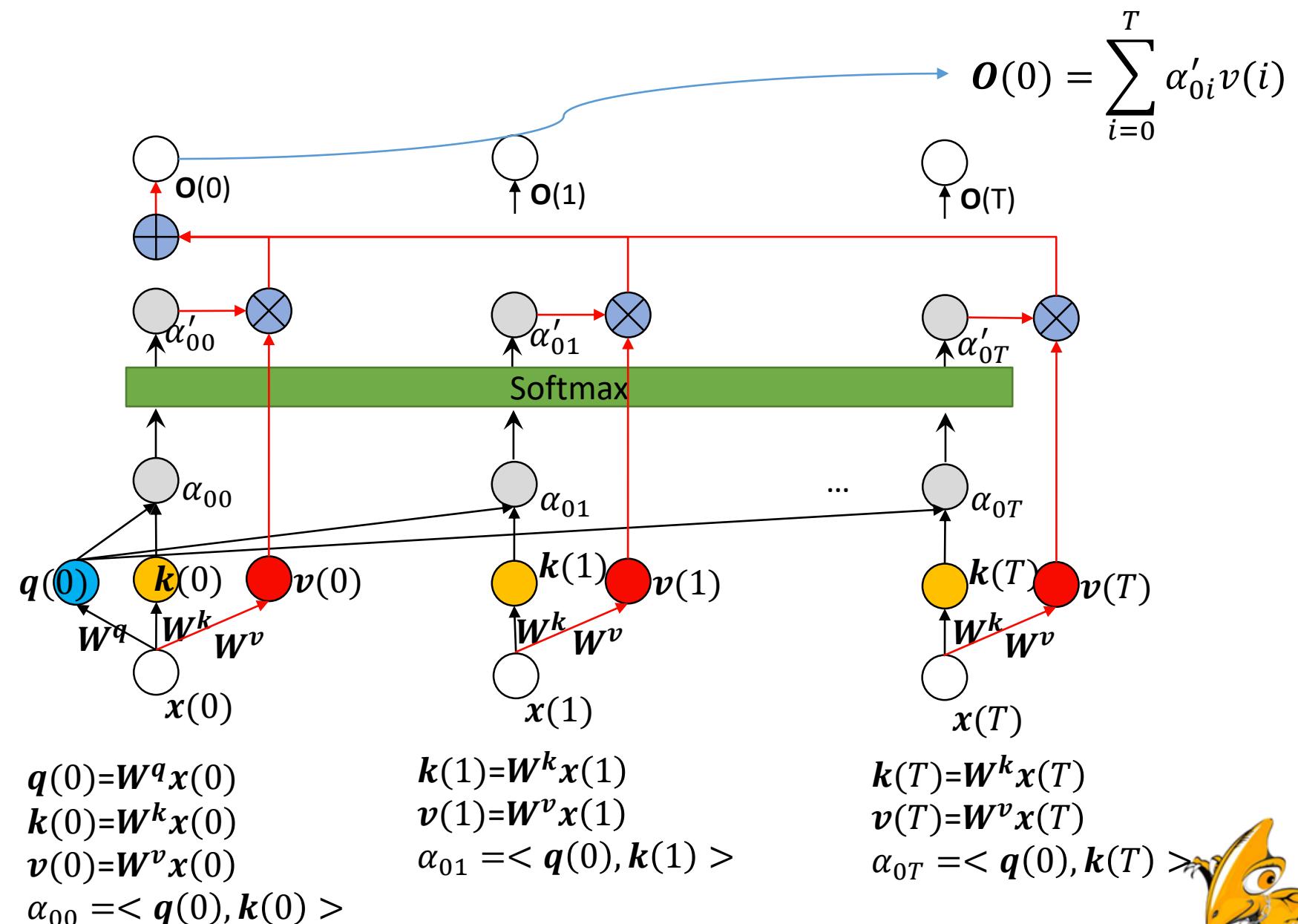


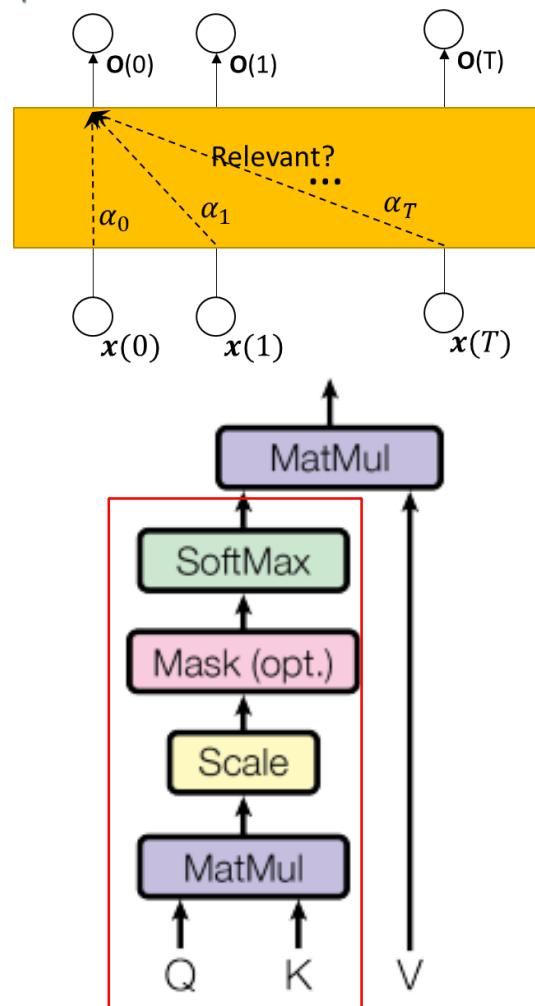


Q: Query

K: Key

V: Value

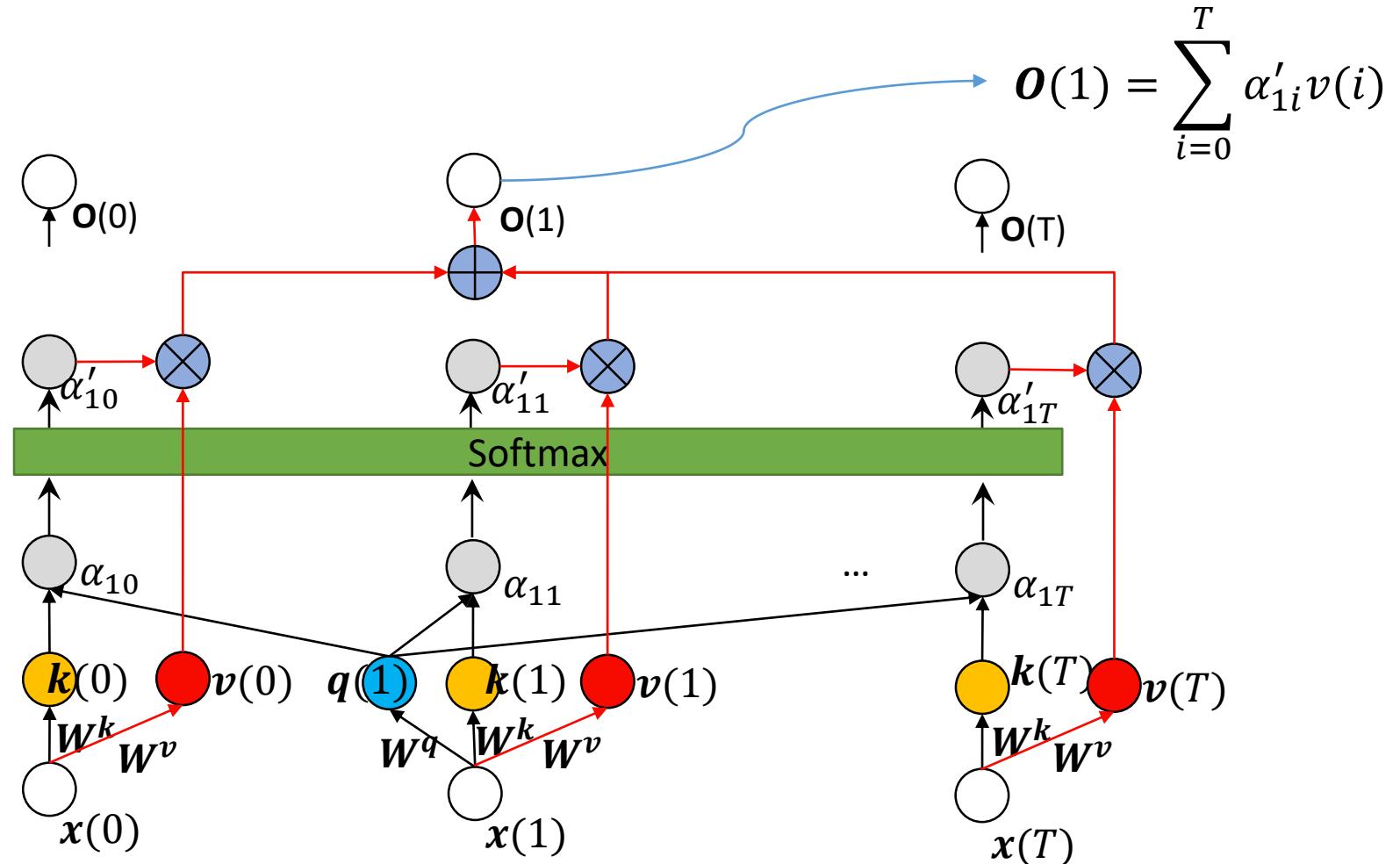




Q: Query

K: Key

V: Value



$$k(0) = W^k x(0)$$

$$v(0) = W^v x(0)$$

$$\alpha_{10} = < q(1), k(0) >$$

$$q(1) = W^q x(1)$$

$$k(1) = W^k x(1)$$

$$v(1) = W^v x(1)$$

$$\alpha_{11} = < q(1), k(1) >$$

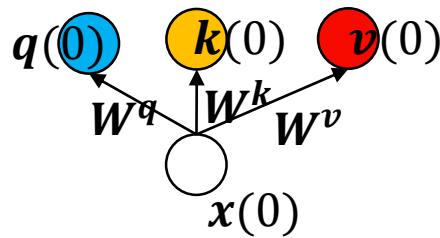
$$k(T) = W^k x(T)$$

$$v(T) = W^v x(T)$$

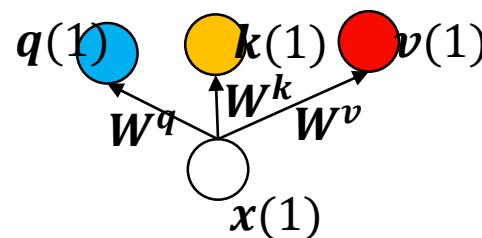
$$\alpha_{1T} = < q(1), k(T) >$$



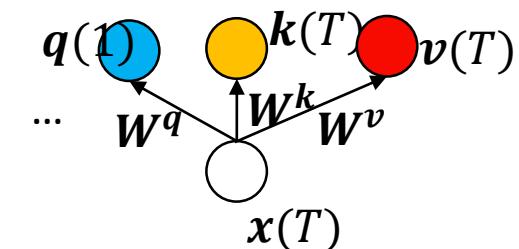
# 用大矩陣直接運算QKV



$$\begin{aligned} q(0) &= W^q x(0) \\ k(0) &= W^k x(0) \\ v(0) &= W^v x(0) \end{aligned}$$



$$\begin{aligned} q(1) &= W^q x(1) \\ k(1) &= W^k x(1) \\ v(1) &= W^v x(1) \end{aligned}$$



$$\begin{aligned} q(T) &= W^q x(T) \\ k(T) &= W^k x(T) \\ v(T) &= W^v x(T) \end{aligned}$$



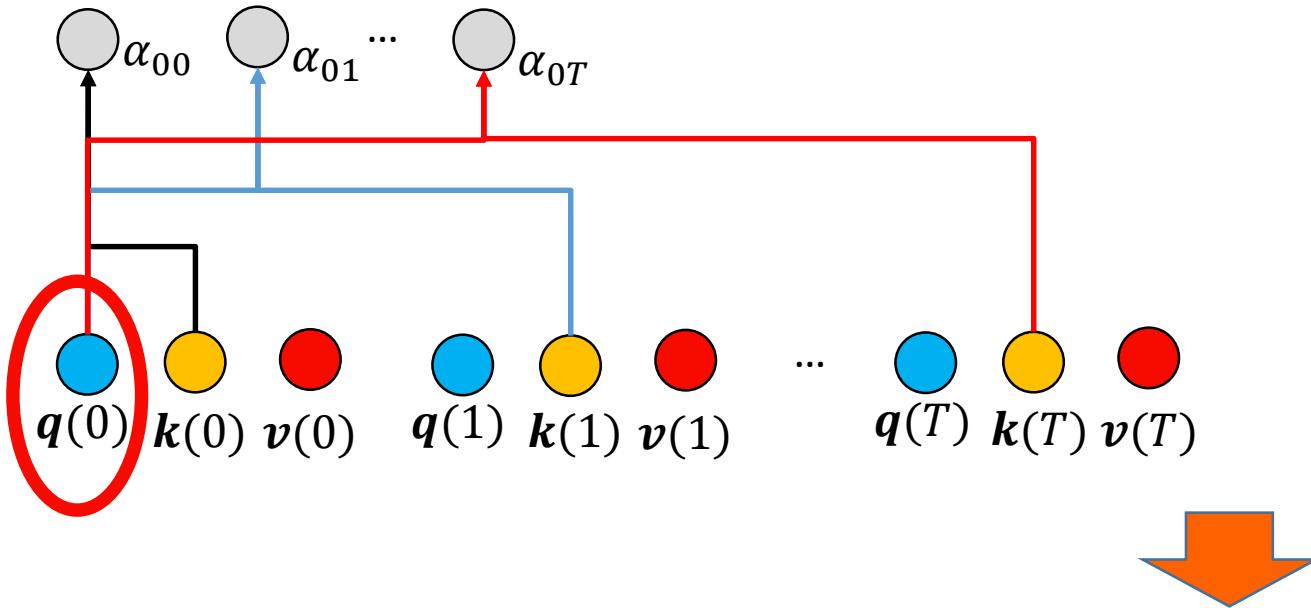
$$\mathbf{Q}_{d \times T} = \left[ \begin{array}{c|c|c} q(0) & q(1) & \cdots & q(T) \\ \hline d \times 1 & d \times 1 & & d \times 1 \end{array} \right]_{d \times T} = \boxed{W^q} \left[ \begin{array}{c|c|c} x(0) & x(1) & \cdots & x(T) \\ \hline d \times 1 & d \times 1 & & d \times 1 \end{array} \right]_{d \times T}$$

$$\mathbf{K}_{d \times T} = \left[ \begin{array}{c|c|c} k(0) & k(1) & \cdots & k(T) \\ \hline d \times 1 & d \times 1 & & d \times 1 \end{array} \right]_{d \times T} = \boxed{W^k} \left[ \begin{array}{c|c|c} x(0) & x(1) & \cdots & x(T) \\ \hline d \times 1 & d \times 1 & & d \times 1 \end{array} \right]_{d \times T}$$

$$\mathbf{V}_{d \times T} = \left[ \begin{array}{c|c|c} v(0) & v(1) & \cdots & v(T) \\ \hline d \times 1 & d \times 1 & & d \times 1 \end{array} \right]_{d \times T} = \boxed{W^v} \left[ \begin{array}{c|c|c} x(0) & x(1) & \cdots & x(T) \\ \hline d \times 1 & d \times 1 & & d \times 1 \end{array} \right]_{d \times T}$$



# 用大矩陣直接運算A



$$\begin{bmatrix} \alpha_{00} & \alpha_{10} & \dots & \alpha_{T0} \\ \alpha_{01} & \alpha_{11} & \dots & \alpha_{T1} \\ \vdots & \ddots & \ddots & \vdots \\ \alpha_{0T} & \alpha_{1T} & \dots & \alpha_{TT} \end{bmatrix}_{T \times T} = \mathbf{A} = \langle \mathbf{Q}, \mathbf{K} \rangle = \mathbf{K}^T \mathbf{Q} =$$

Attention Score at  $q(0)$

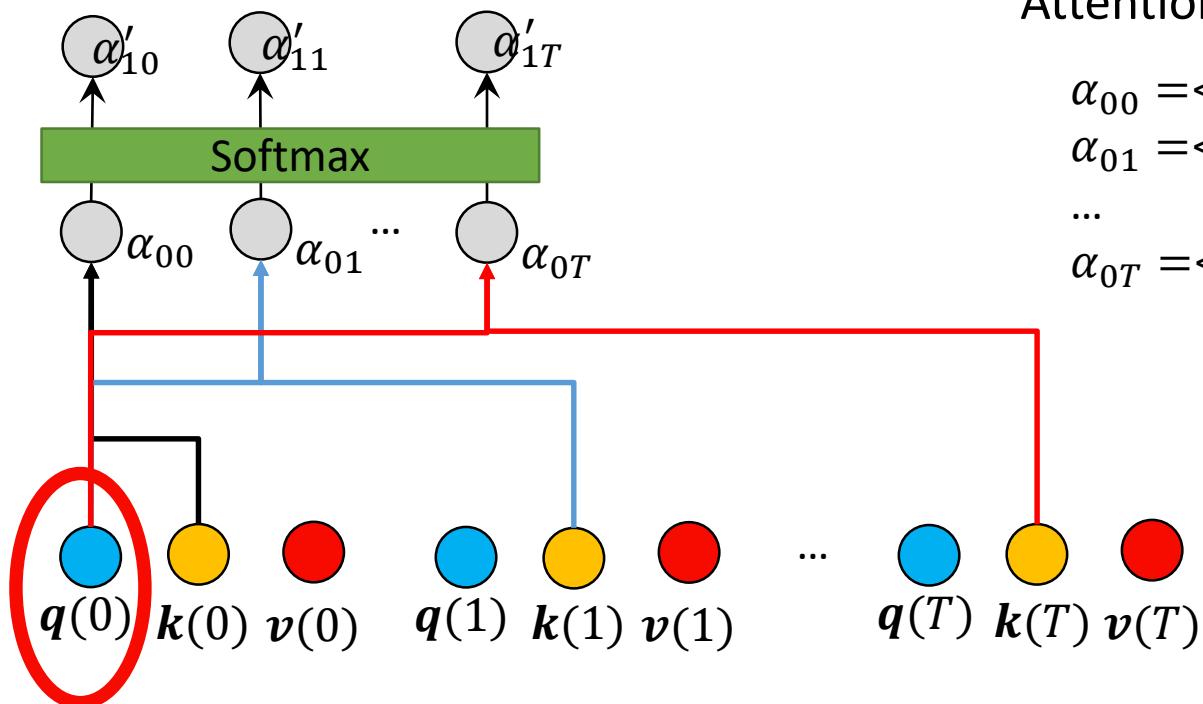
$\alpha_{00} = \langle \mathbf{q}(0), \mathbf{k}(0) \rangle$   
 $\alpha_{01} = \langle \mathbf{q}(0), \mathbf{k}(1) \rangle$   
 ...  
 $\alpha_{0T} = \langle \mathbf{q}(0), \mathbf{k}(T) \rangle$

Attention Score at  $q(0)$

$$\mathbf{K}^T_{T \times d} \quad \mathbf{Q}_{d \times T}$$

$\mathbf{K}^T$  is a  $T \times d$  matrix where each row is a key vector  $\mathbf{k}(0), \mathbf{k}(1), \dots, \mathbf{k}(T)$ , each of size  $d \times 1$ .  
 $\mathbf{Q}$  is a  $d \times T$  matrix where each column is a query vector  $\mathbf{q}(0), \mathbf{q}(1), \dots, \mathbf{q}(T)$ , each of size  $d \times 1$ .





Attention Score at  $q(0)$

$$\begin{aligned}\alpha_{00} &= \langle \mathbf{q}(0), \mathbf{k}(0) \rangle \\ \alpha_{01} &= \langle \mathbf{q}(0), \mathbf{k}(1) \rangle \\ &\dots \\ \alpha_{0T} &= \langle \mathbf{q}(0), \mathbf{k}(T) \rangle\end{aligned}$$

Softmax

$$\alpha'_{0i} = \frac{\exp(\alpha_{0i})}{\sum_{j=0}^T \exp(\alpha_{0j})} \quad \forall i = 0, 1, \dots, T$$



Attention Score Matrix

$$\begin{bmatrix} \alpha_{00} & \alpha_{10} & \dots & \alpha_{T0} \\ \alpha_{01} & \alpha_{11} & \dots & \alpha_{T1} \\ \vdots & \vdots & \ddots & \vdots \\ \alpha_{0T} & \alpha_{1T} & \dots & \alpha_{TT} \end{bmatrix}_{T \times T} = \mathbf{A}$$

Softmax  
For each column

$\hat{\mathbf{A}} =$

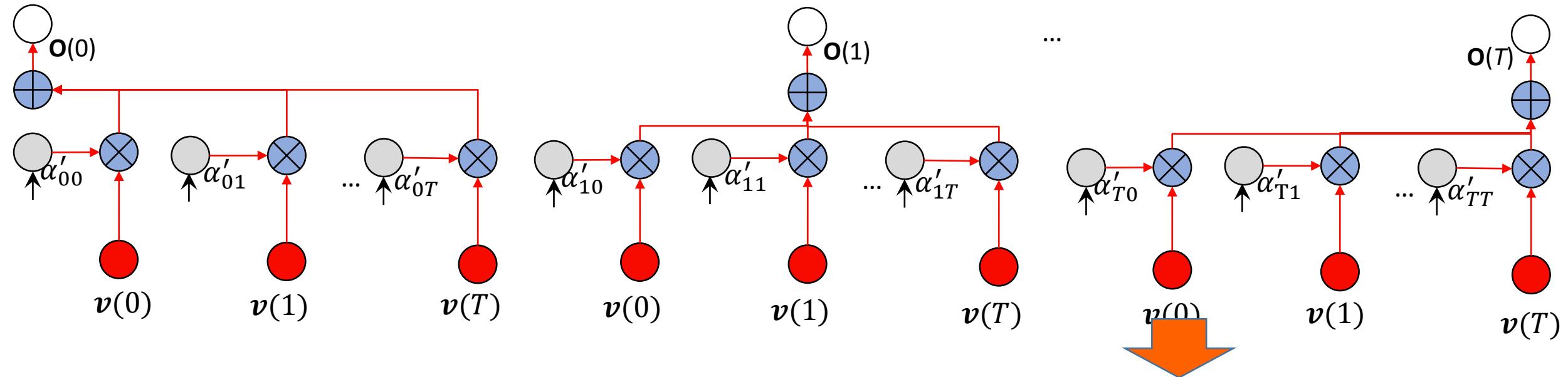
$$\begin{bmatrix} \alpha'_{00} & \alpha'_{10} & \dots & \alpha'_{T0} \\ \alpha'_{01} & \alpha'_{11} & \dots & \alpha'_{T1} \\ \vdots & \vdots & \ddots & \vdots \\ \alpha'_{0T} & \alpha'_{1T} & \dots & \alpha'_{TT} \end{bmatrix}_{T \times T}$$

$$\hat{\mathbf{A}} = \text{softmax}(\mathbf{A})$$

Attention Score at  $q(0)$



# 用大矩陣直接運算 O



$$\mathbf{O} = \begin{bmatrix} \boxed{\mathbf{o}(0)} & \boxed{\mathbf{o}(1)} & \cdots & \boxed{\mathbf{o}(T)} \end{bmatrix}_{d \times T} = \mathbf{V} \widehat{\mathbf{A}} = \begin{bmatrix} \boxed{\mathbf{v}(0)} & \boxed{\mathbf{v}(1)} & \cdots & \boxed{\mathbf{v}(T)} \end{bmatrix}_{d \times T} \begin{bmatrix} \alpha'_{00} & \alpha'_{10} & \cdots & \alpha'_{T0} \\ \alpha'_{01} & \alpha'_{11} & \cdots & \alpha'_{T1} \\ \vdots & \ddots & \ddots & \vdots \\ \alpha'_{0T} & \alpha'_{1T} & \cdots & \alpha'_{TT} \end{bmatrix}_{T \times T}$$



# 矩陣直接運算Self-Attention

$$\mathbf{Q}_{d \times T} = \begin{bmatrix} \mathbf{q}(0) & \mathbf{q}(1) & \cdots & \mathbf{q}(T) \\ d \times 1 & d \times 1 & \cdots & d \times 1 \end{bmatrix}_{d \times T} = \mathbf{W}^q \begin{bmatrix} \mathbf{x}(0) & \mathbf{x}(1) & \cdots & \mathbf{x}(T) \\ d \times 1 & d \times 1 & \cdots & d \times 1 \end{bmatrix}_{d \times T}$$

$$\mathbf{K}_{d \times T} = \begin{bmatrix} \mathbf{k}(0) & \mathbf{k}(1) & \cdots & \mathbf{k}(T) \\ d \times 1 & d \times 1 & \cdots & d \times 1 \end{bmatrix}_{d \times T} = \mathbf{W}^k \begin{bmatrix} \mathbf{x}(0) & \mathbf{x}(1) & \cdots & \mathbf{x}(T) \\ d \times 1 & d \times 1 & \cdots & d \times 1 \end{bmatrix}_{d \times T}$$

$$\mathbf{V}_{d \times T} = \begin{bmatrix} \mathbf{v}(0) & \mathbf{v}(1) & \cdots & \mathbf{v}(T) \\ d \times 1 & d \times 1 & \cdots & d \times 1 \end{bmatrix}_{d \times T} = \mathbf{W}^v \begin{bmatrix} \mathbf{x}(0) & \mathbf{x}(1) & \cdots & \mathbf{x}(T) \\ d \times 1 & d \times 1 & \cdots & d \times 1 \end{bmatrix}_{d \times T}$$

$$A = \langle Q, K \rangle = K^T Q$$

$T \times T$        $T \times d$        $d \times T$

論文寫法

$$\hat{A} = softmax(A)$$

$$O = V\hat{A}$$

$d \times T$

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

$$O = V \times softmax(K^T Q)$$

差異在矩陣擺法，和一個 $\sqrt{d_k}$



# 矩陣直接運算Self-Attention

論文寫法

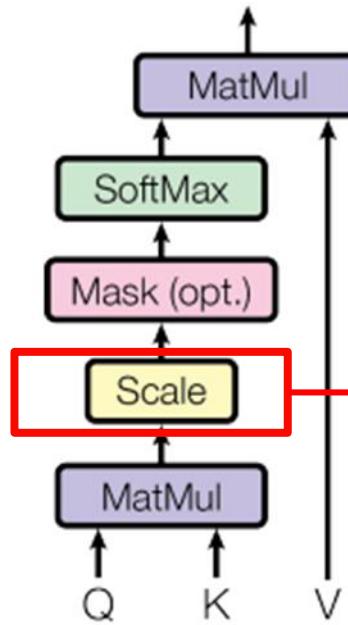
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

$$O = V \times \text{softmax}(K^T Q)$$

差異在矩陣擺法，和一個  $\sqrt{d_k}$

$$\begin{bmatrix} \mathbf{k}(0)_{d \times 1} \\ \mathbf{k}(1)_{d \times 1} \\ \vdots \\ \mathbf{k}(T)_{d \times 1} \end{bmatrix}_{T \times d} \quad \begin{bmatrix} \mathbf{q}(0)_{d \times 1} & \mathbf{q}(1)_{d \times 1} & \cdots & \mathbf{q}(T)_{d \times 1} \end{bmatrix}_{d \times T}$$

**Q**  
 $d \times T$



K和Q在不同時間點上  
d個維度相乘後的總和。

所以維度越大值越大，因此需要正規化(在右圖的Scale)，將每一個值域拉到一致。

最簡單方式有d個維度就除上d  
但論文第四頁內提到當維度很大的時候，預防梯度過小的問題，  
因此開根號。



# Self-Attention: 讓公式和論文一樣

$$O = V \times \text{softmax}(K^T Q)$$

$$\begin{matrix} Q & K & V & O \\ d \times T & d \times T & d \times T & d \times T \end{matrix}$$

差異在矩陣擺法

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

$$\begin{matrix} Q & K & V & O \\ T \times d & T \times d & T \times d & T \times d \end{matrix}$$

我的擺法是讓大家比較能看得懂時序的資料是橫向的。

為了讓公式和論文一樣，我們轉成論文的公式

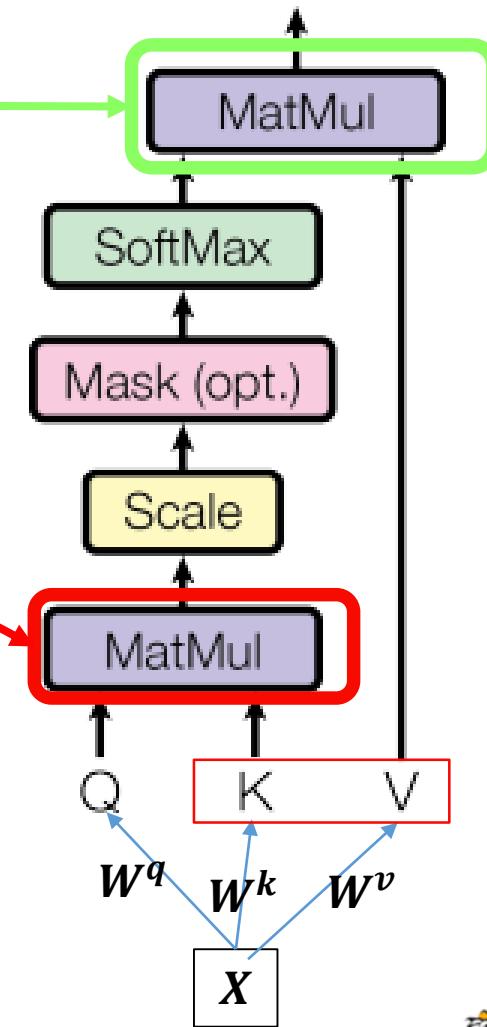
*Linear Projection*  $W^q_{d_x \times d}$

$$X_{T \times d_x} \xrightarrow{\quad} W^k_{d_x \times d} \quad W^v_{d_x \times d}$$

$$Q_{T \times d} = X_{T \times d_x} W^q_{d_x \times d}$$

$$K_{T \times d} = X_{T \times d_x} W^k_{d_x \times d}$$

$$V_{T \times d} = X_{T \times d_x} W^v_{d_x \times d}$$



$X_{T \times d_x}, Q_{T \times d}, K_{T \times d}, V_{T \times d}$

$$T = 3, d_x = 4, d = 2$$

**Linear Projection**  $W^q_{d_x \times d}$

$X_{T \times d_x}$



$W^k_{d_x \times d}$

$W^v_{d_x \times d}$

$$Q_{T \times d} = X_{T \times d_x} W^q_{d_x \times d}$$

$$K_{T \times d} = X_{T \times d_x} W^k_{d_x \times d}$$

$$V_{T \times d} = X_{T \times d_x} W^v_{d_x \times d}$$

|          |          |          |          |
|----------|----------|----------|----------|
| $x_{11}$ | $x_{12}$ | $x_{13}$ | $x_{14}$ |
| $x_{21}$ | $x_{22}$ | $x_{23}$ | $x_{24}$ |
| $x_{31}$ | $x_{32}$ | $x_{33}$ | $x_{34}$ |

$X_{3 \times 4}$

$$W^q_{d_x \times d} = W^q_{4 \times 2}$$

|  |  |
|--|--|
|  |  |
|  |  |
|  |  |
|  |  |
|  |  |

$$W^k_{d_x \times d} = W^k_{4 \times 2}$$

|  |  |
|--|--|
|  |  |
|  |  |
|  |  |
|  |  |
|  |  |

$$W^v_{d_x \times d} = W^v_{4 \times 2}$$

|  |  |
|--|--|
|  |  |
|  |  |
|  |  |
|  |  |
|  |  |

$Q_{3 \times 2}$

$K_{3 \times 2}$

$V_{3 \times 2}$



# 解釋Attention Score

$$O_{T \times d}, Q_{T \times d}, K_{T \times d}, V_{T \times d}$$

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

$$O \approx QK^T V$$

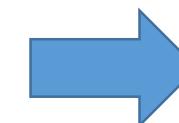
$$A_{T_Q \times T_K} = \begin{bmatrix} a_{q_1 k_1} & a_{q_1 k_2} & \cdots & a_{q_1 k_T} \\ a_{q_2 k_1} & a_{q_2 k_2} & \cdots & a_{q_2 k_T} \\ \vdots & \vdots & \ddots & \vdots \\ a_{q_T k_1} & a_{q_T k_2} & \cdots & a_{q_T k_T} \end{bmatrix} = \begin{bmatrix} q_{11} & q_{12} & \cdots & q_{1d} \\ q_{21} & q_{22} & \cdots & q_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ q_{T1} & q_{d2} & \cdots & q_{Td} \end{bmatrix} \begin{bmatrix} k_{11} & k_{12} & \cdots & k_{1T} \\ k_{21} & k_{22} & \cdots & k_{2T} \\ \vdots & \vdots & \ddots & \vdots \\ k_{d1} & k_{d2} & \cdots & k_{dT} \end{bmatrix}$$

$a_{q_1 k_1}$ : Q在第1個時間資料和K第1個時間資料的內積

⋮

$a_{q_T k_1}$ : Q在第T個時間資料和K第1個時間資料的內積

$a_{q_T k_T}$ : Q在第T個時間資料和K第T個時間資料的內積



Q和K 時間資料之間的相關矩陣



# 解釋Attention Score

$$O_{T \times d}, Q_{T \times d}, K_{T \times d}, V_{T \times d}$$

$$T = 3, d = 2$$

$$O_{3 \times 2}$$

|          |          |
|----------|----------|
| $o_{11}$ | $o_{12}$ |
| $o_{21}$ | $o_{22}$ |
| $o_{31}$ | $o_{32}$ |

$$Q_{3 \times 2}$$

|          |          |
|----------|----------|
| $q_{11}$ | $q_{12}$ |
| $q_{21}$ | $q_{22}$ |
| $q_{31}$ | $q_{32}$ |

$$K_{3 \times 2}$$

|          |          |
|----------|----------|
| $k_{11}$ | $k_{12}$ |
| $k_{21}$ | $k_{22}$ |
| $k_{31}$ | $k_{32}$ |

$$V_{3 \times 2}$$

|          |          |
|----------|----------|
| $v_{11}$ | $v_{12}$ |
| $v_{21}$ | $v_{22}$ |
| $v_{31}$ | $v_{32}$ |

$$QK^T$$

|          |          |          |
|----------|----------|----------|
| $a_{11}$ | $a_{12}$ | $a_{13}$ |
| $a_{21}$ | $a_{22}$ | $a_{23}$ |
| $a_{31}$ | $a_{32}$ | $a_{33}$ |

$$= \begin{matrix} t = 1 \\ t = 2 \\ t = 3 \end{matrix}$$

$$Q$$

|          |          |
|----------|----------|
| $q_{11}$ | $q_{12}$ |
| $q_{21}$ | $q_{22}$ |
| $q_{31}$ | $q_{32}$ |

$$K^T$$

|          |          |          |
|----------|----------|----------|
| $k_{11}$ | $k_{21}$ | $k_{31}$ |
| $k_{12}$ | $k_{22}$ | $k_{32}$ |

Q和K 時間資料之間的相關矩陣

$$\hat{A} = softmax(QK^T)$$

|           |           |           |
|-----------|-----------|-----------|
| $a'_{11}$ | $a'_{12}$ | $a'_{13}$ |
| $a'_{21}$ | $a'_{22}$ | $a'_{23}$ |
| $a'_{31}$ | $a'_{32}$ | $a'_{33}$ |

每個row做softmax

$$a'_{11} = \frac{a_{11}}{a_{11} + a_{12} + a_{13}}$$

$$a'_{12} = \frac{a_{12}}{a_{11} + a_{12} + a_{13}}$$

$$a'_{13} = \frac{a_{13}}{a_{11} + a_{12} + a_{13}}$$

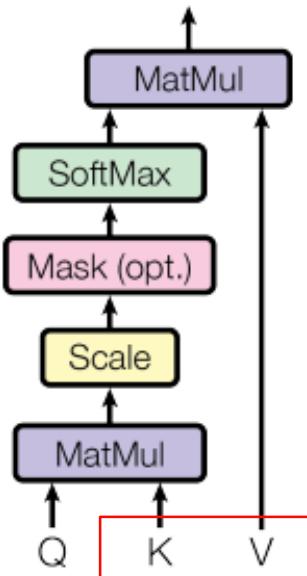
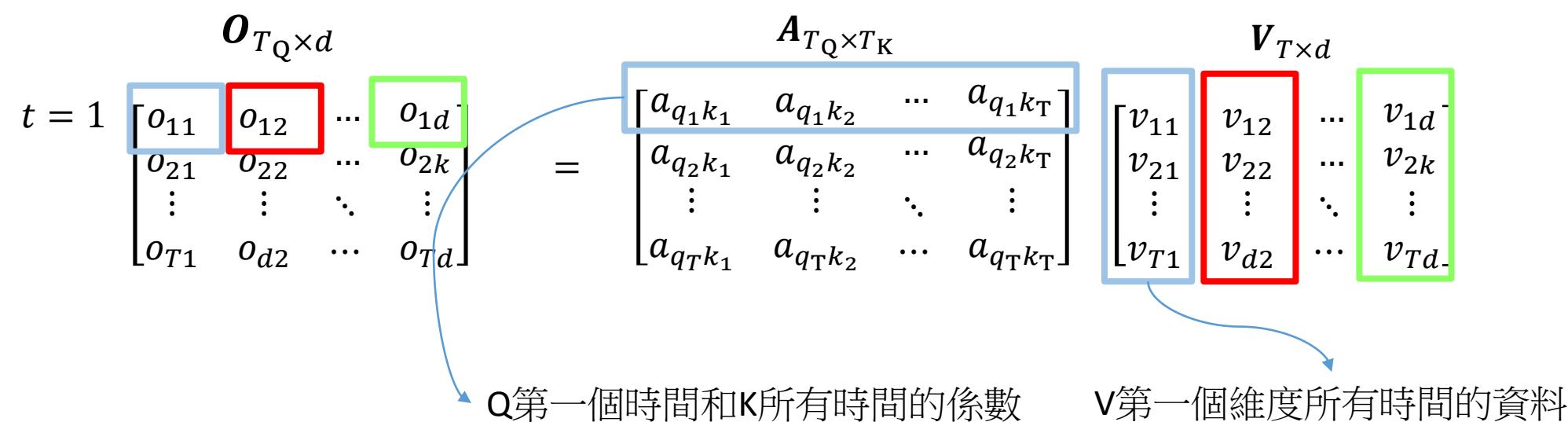


# Value和Attention Score

$$O_{T \times d}, Q_{T \times d}, K_{T \times d}, V_{T \times d}$$

$$O \approx QK^T V$$

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



$O_{T_Q \times d}$  : 將V從K空間投影到Q空間





# Value和Attention Score

$$\mathbf{O}_{T \times d}, \mathbf{Q}_{T \times d}, \mathbf{K}_{T \times d}, \mathbf{V}_{T \times d}$$

$$T = 3, d = 2$$

$$\mathbf{O} = \text{softmax}(\mathbf{Q}\mathbf{K}^T) \mathbf{V} \\ = \widehat{\mathbf{A}}\mathbf{V}$$

$$\mathbf{O}_{3 \times 2} \quad \mathbf{Q}_{3 \times 2}$$

|          |          |
|----------|----------|
| $o_{11}$ | $o_{12}$ |
| $o_{21}$ | $o_{22}$ |
| $o_{31}$ | $o_{32}$ |

|          |          |
|----------|----------|
| $q_{11}$ | $q_{12}$ |
| $q_{21}$ | $q_{22}$ |
| $q_{31}$ | $q_{32}$ |

$$\mathbf{O}_{3 \times 2} \quad \widehat{\mathbf{A}} = \text{softmax}(\mathbf{Q}\mathbf{K}^T) \quad \mathbf{V}_{3 \times 2}$$

|          |          |
|----------|----------|
| $o_{11}$ | $o_{12}$ |
| $o_{21}$ | $o_{22}$ |
| $o_{31}$ | $o_{32}$ |

 $=$ 

|           |           |           |
|-----------|-----------|-----------|
| $a'_{11}$ | $a'_{12}$ | $a'_{13}$ |
| $a'_{21}$ | $a'_{22}$ | $a'_{23}$ |
| $a'_{31}$ | $a'_{32}$ | $a'_{33}$ |

 $\times$ 

|          |          |
|----------|----------|
| $v_{11}$ | $v_{12}$ |
| $v_{21}$ | $v_{22}$ |
| $v_{31}$ | $v_{32}$ |

$$\mathbf{K}_{3 \times 2} \quad \mathbf{V}_{3 \times 2}$$

|          |          |
|----------|----------|
| $k_{11}$ | $k_{12}$ |
| $k_{21}$ | $k_{22}$ |
| $k_{31}$ | $k_{32}$ |

|          |          |
|----------|----------|
| $v_{11}$ | $v_{12}$ |
| $v_{21}$ | $v_{22}$ |
| $v_{31}$ | $v_{32}$ |

$$t = 1 \quad \begin{matrix} o_{11} & o_{12} \end{matrix} \quad = \quad \begin{matrix} a'_{11} & a'_{12} & a'_{13} \end{matrix} \quad \times \quad \begin{matrix} v_{11} & v_{12} \\ v_{21} & v_{22} \\ v_{31} & v_{32} \end{matrix}$$

$$= \quad \begin{matrix} a'_{11} \\ =0.6 \end{matrix} \quad \times \quad \begin{matrix} v_{11} & v_{12} \end{matrix} \\ + \quad \begin{matrix} a'_{12} \\ =0.3 \end{matrix} \quad \times \quad \begin{matrix} v_{21} & v_{22} \end{matrix} \\ + \quad \begin{matrix} a'_{13} \\ =0.1 \end{matrix} \quad \times \quad \begin{matrix} v_{31} & v_{32} \end{matrix}$$

第一個時間的輸出 =  $a'_{11} * \text{第一個時間的value} + a'_{12} * \text{第2個時間的value} + a'_{13} * \text{第3個時間的value}$



# 當Attention Score為單位矩陣

$$O_{T \times d}, Q_{T \times d}, K_{T \times d}, V_{T \times d}$$

$$O \approx QK^T V$$

$$\begin{matrix} O_{T_Q \times d} \\ \left[ \begin{matrix} o_{11} & o_{12} & \cdots & o_{1d} \\ o_{21} & o_{22} & \cdots & o_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ o_{T1} & o_{d2} & \cdots & o_{Td} \end{matrix} \right] \end{matrix} = \begin{matrix} A_{T_Q \times T_K} \\ \left[ \begin{matrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{matrix} \right] \end{matrix} \begin{matrix} V_{T \times d} \\ \left[ \begin{matrix} v_{11} & v_{12} & \cdots & v_{1d} \\ v_{21} & v_{22} & \cdots & v_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ v_{T1} & v_{d2} & \cdots & v_{Td} \end{matrix} \right] \end{matrix}$$

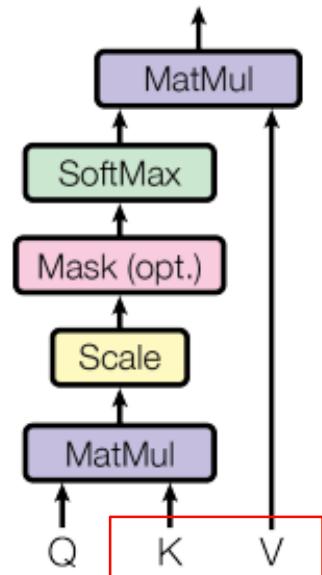
$$O = V$$

第一個時間的輸出 =  $a'_{11} * 第一個時間的value + a'_{12} * 第2個時間的value + a'_{13} * 第3個時間的value$

= 1 \* 第一個時間的value + 0 \* 第2個時間的value + 0 \* 第3個時間的value

第一個時間的輸出 = 第一個時間的value

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

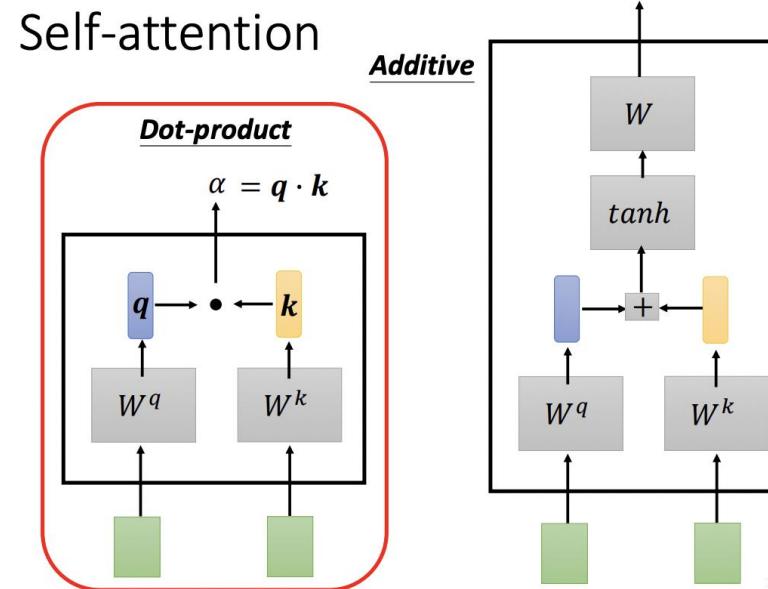


每個時間之間沒有關聯性 = A單位矩陣



# 論文提到的Self-Attention

The two most commonly used attention functions are additive attention [2], and dot-product (multiplicative) attention. Dot-product attention is identical to our algorithm, except for the scaling factor of  $\frac{1}{\sqrt{d_k}}$ . Additive attention computes the compatibility function using a feed-forward network with a single hidden layer. While the two are similar in theoretical complexity, dot-product attention is much faster and more space-efficient in practice, since it can be implemented using highly optimized matrix multiplication code.



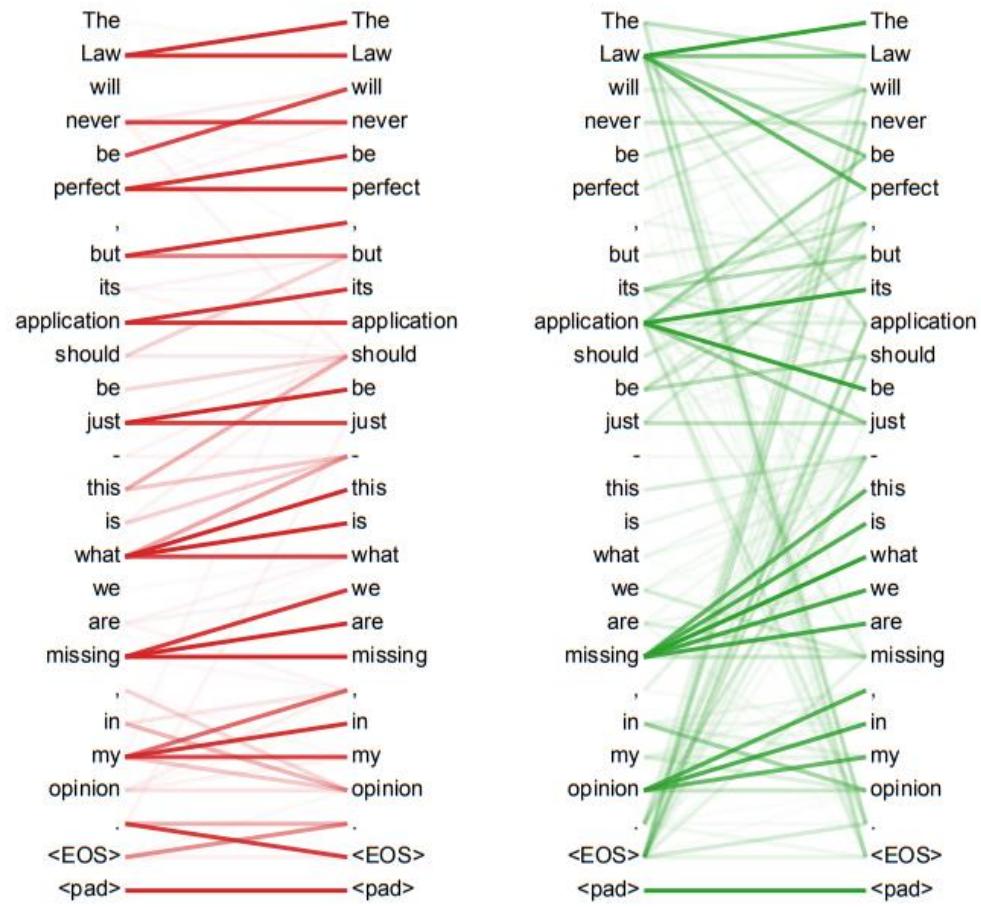
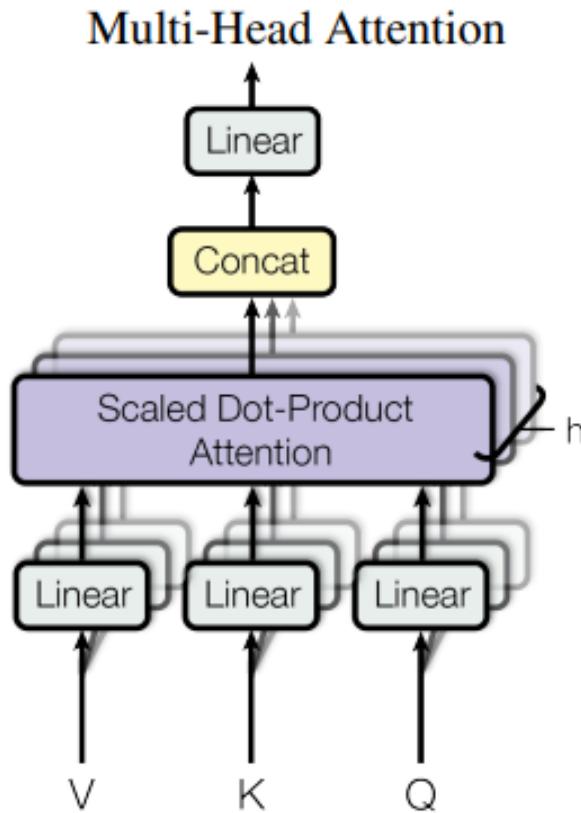
15

Reference: 李宏毅老師



# Multi-head Self-attention

- 相關性在不同的位置可能有不同類別，所以需要多Multi-head來讓模型可以得到更多不同位置的表示。

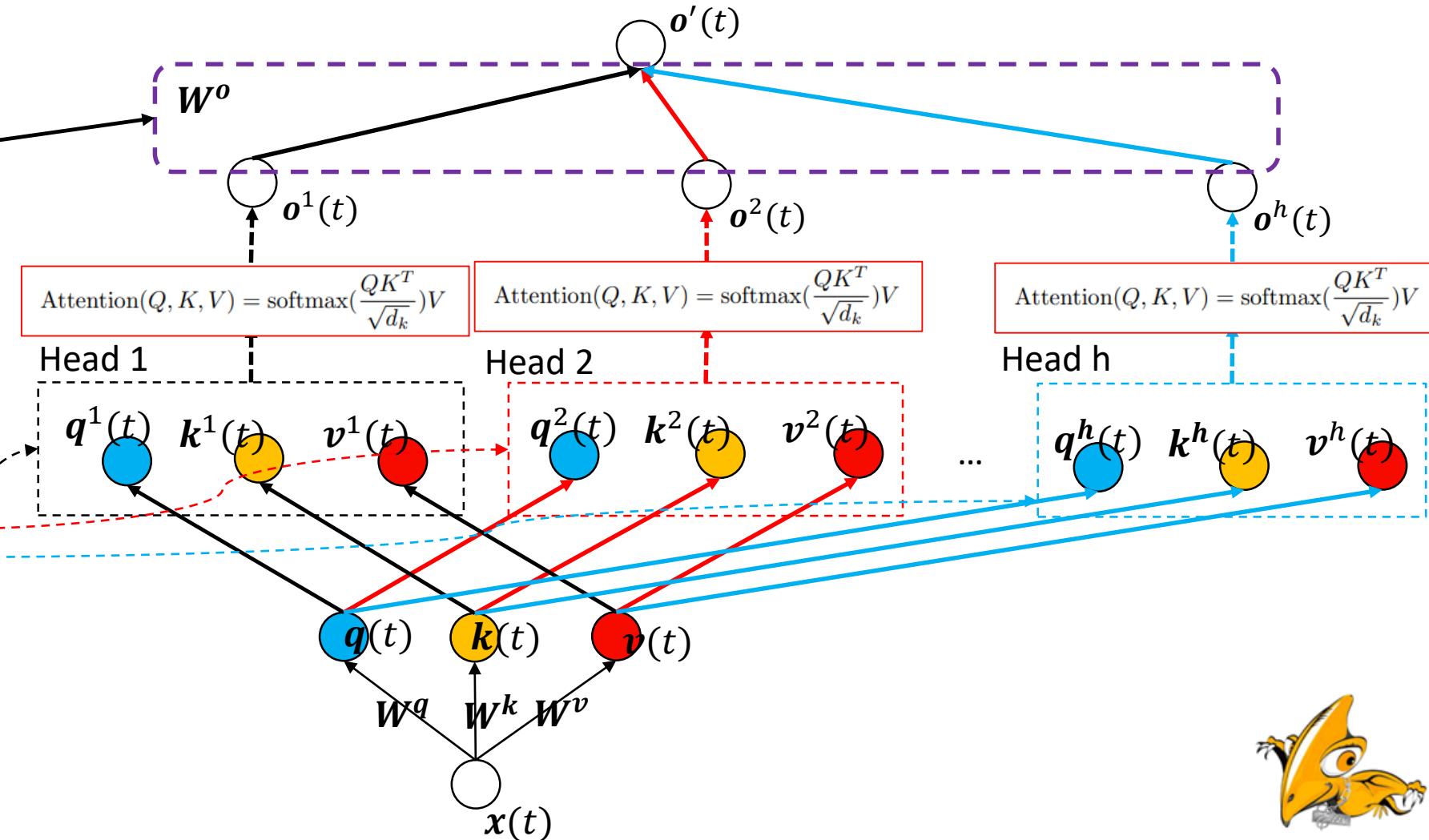
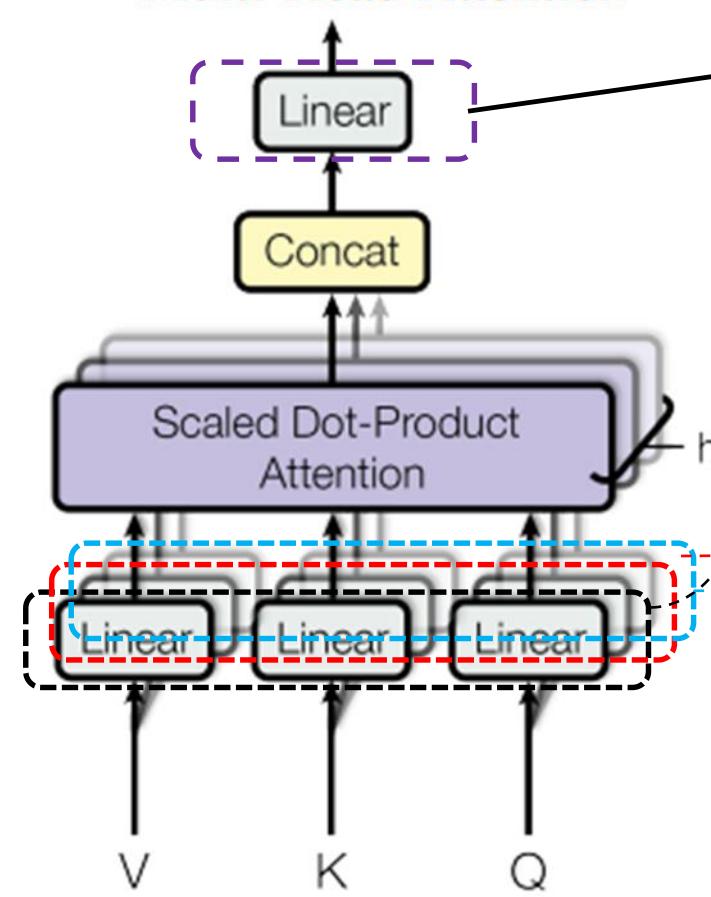


Many of the attention heads exhibit behaviour that seems related to the structure of the sentence.

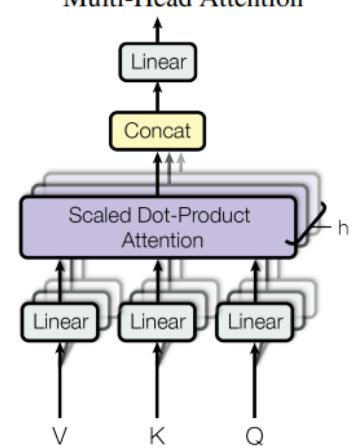


# Multi-head Self-attention

Multi-Head Attention

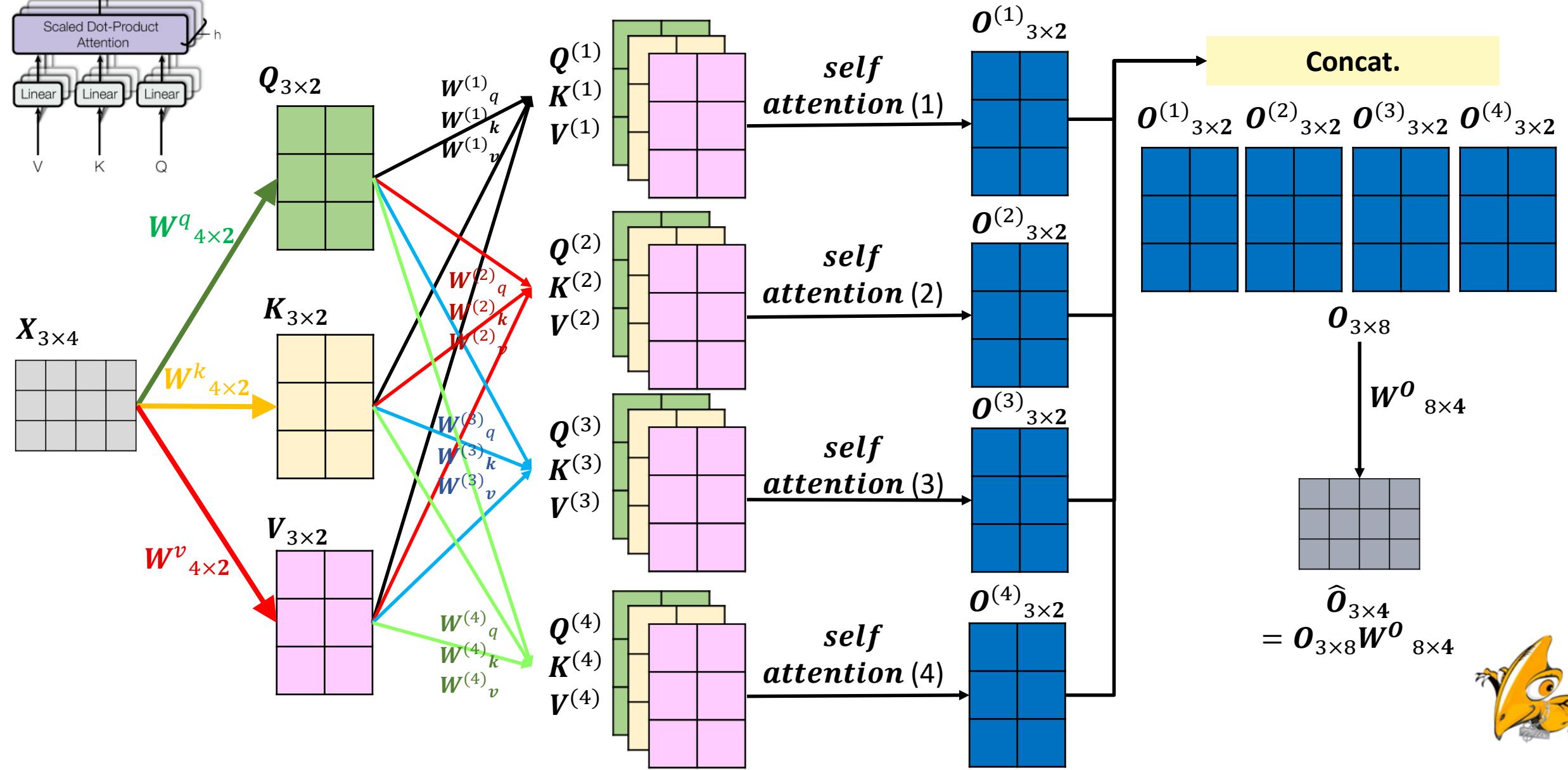


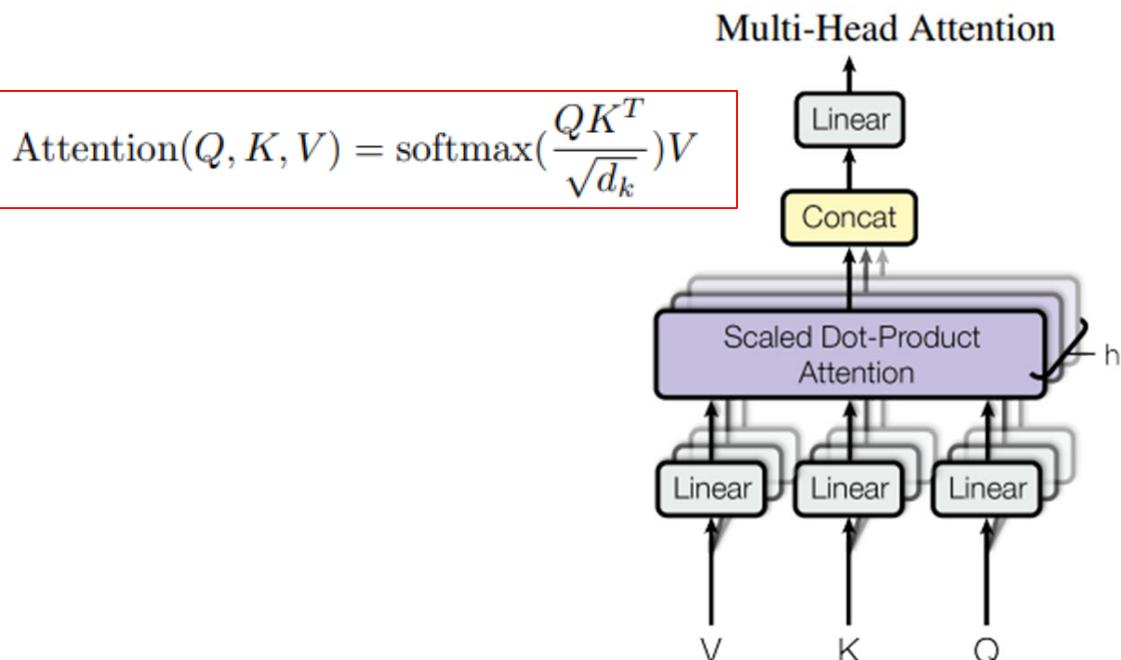
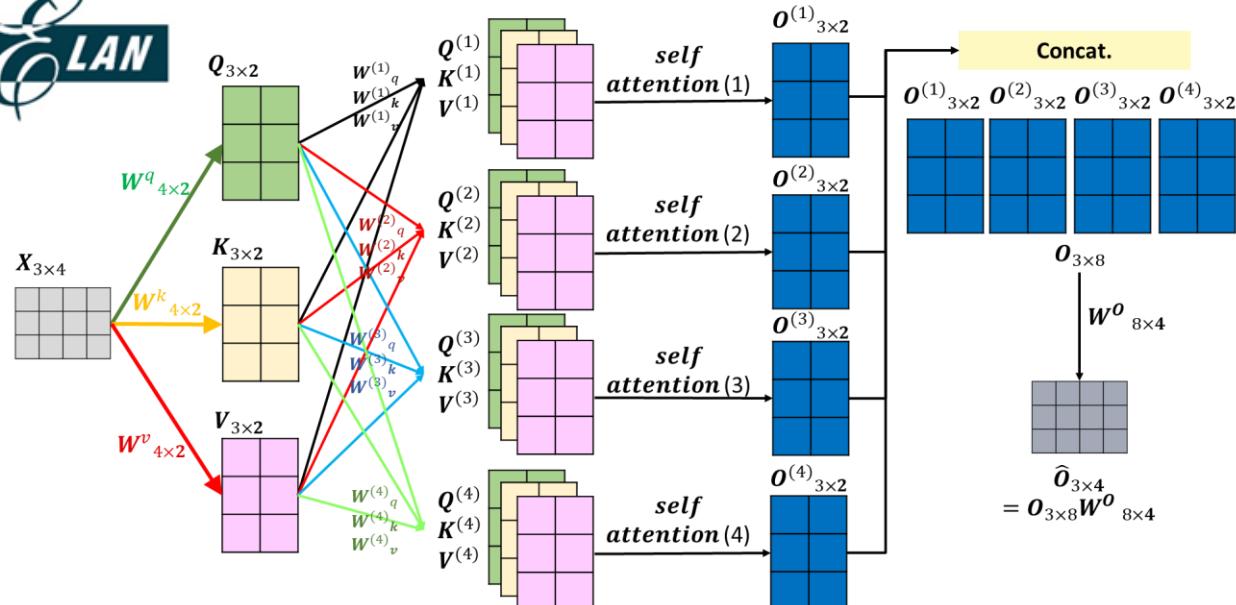
Multi-Head Attention



$$\mathbf{Q}^{(h)}_{T \times d} = \mathbf{W}^{(h)}_q \mathbf{Q}_{T \times d}; \quad \mathbf{K}^{(h)}_{T \times d} = \mathbf{W}^{(h)}_k \mathbf{K}_{T \times d}; \quad \mathbf{V}^{(h)}_{T \times d} = \mathbf{W}^{(h)}_v \mathbf{V}_{T \times d}$$

範例有4個head  $\rightarrow h = 1, 2, 3, 4$





## 參數初估

### Multi-Head Cross Attention:

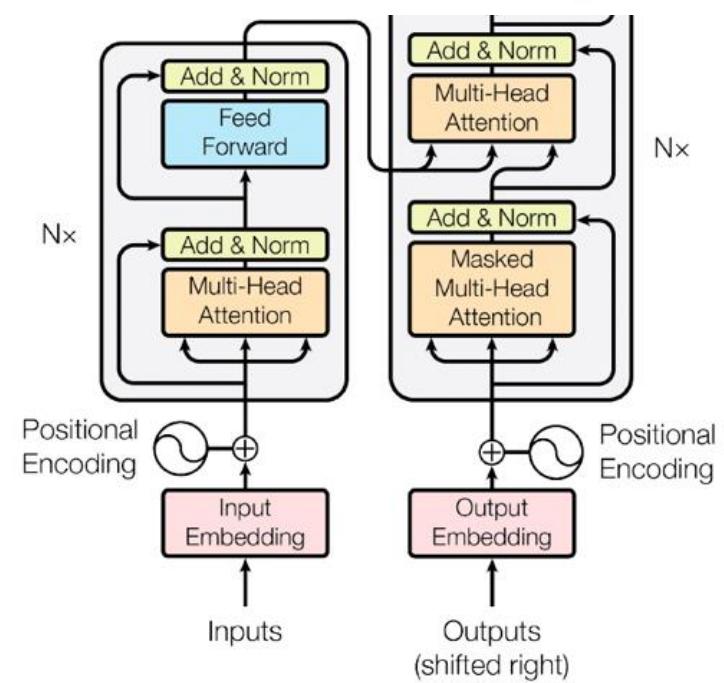
$$X_{T \times dx}, Q_{T \times d}, K_{T \times d}, V_{T \times d}, O_{T \times dx} \\ W^q, W^k, W^v \in \mathbb{R}^{dx \times d}$$

$$W^{(h)}_q, W^{(h)}_k, W^{(h)}_v \in \mathbb{R}^{T \times d}, \quad W^o \in \mathbb{R}^{(n_h \times d) \times dx}$$

$$W^q + W^k + W^v + W^o + \sum_{h=1}^{n_h} W^{(h)}_q + W^{(h)}_k + W^{(h)}_v$$

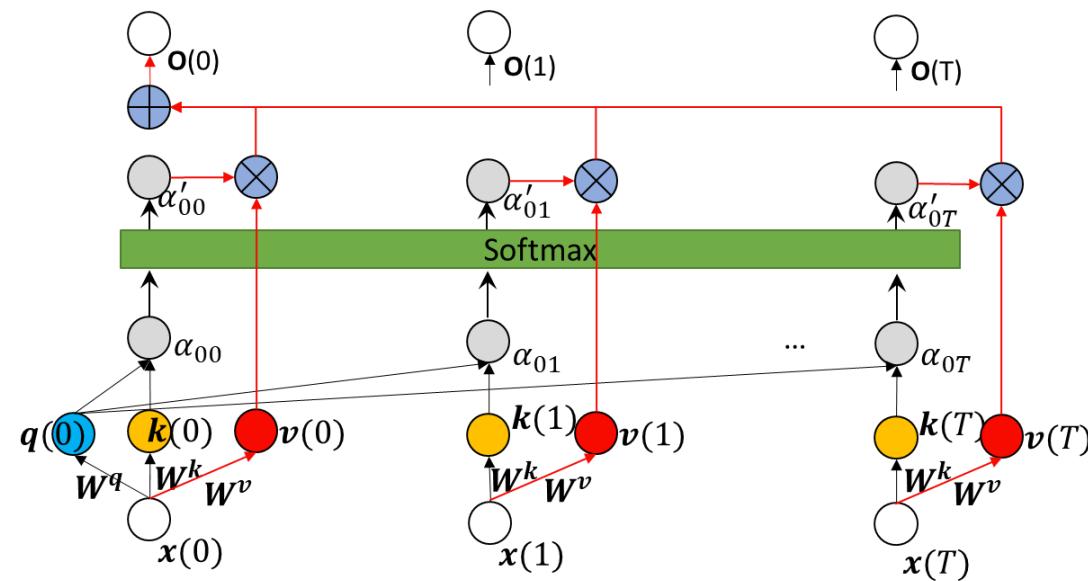
$$T = 10, dx = 300, d = 300, n_h = 8$$

$$dx \times d \times 3 + n_h \times (T \times d \times 3) + (n_h + d) \times dx \\ = 300 \times 300 \times 3 + 8 \times (10 \times 300 \times 3) + (8 + 300) \times 300 \\ = 434,400$$

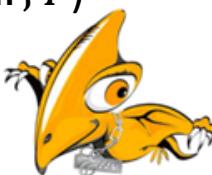


# Positional Encoding

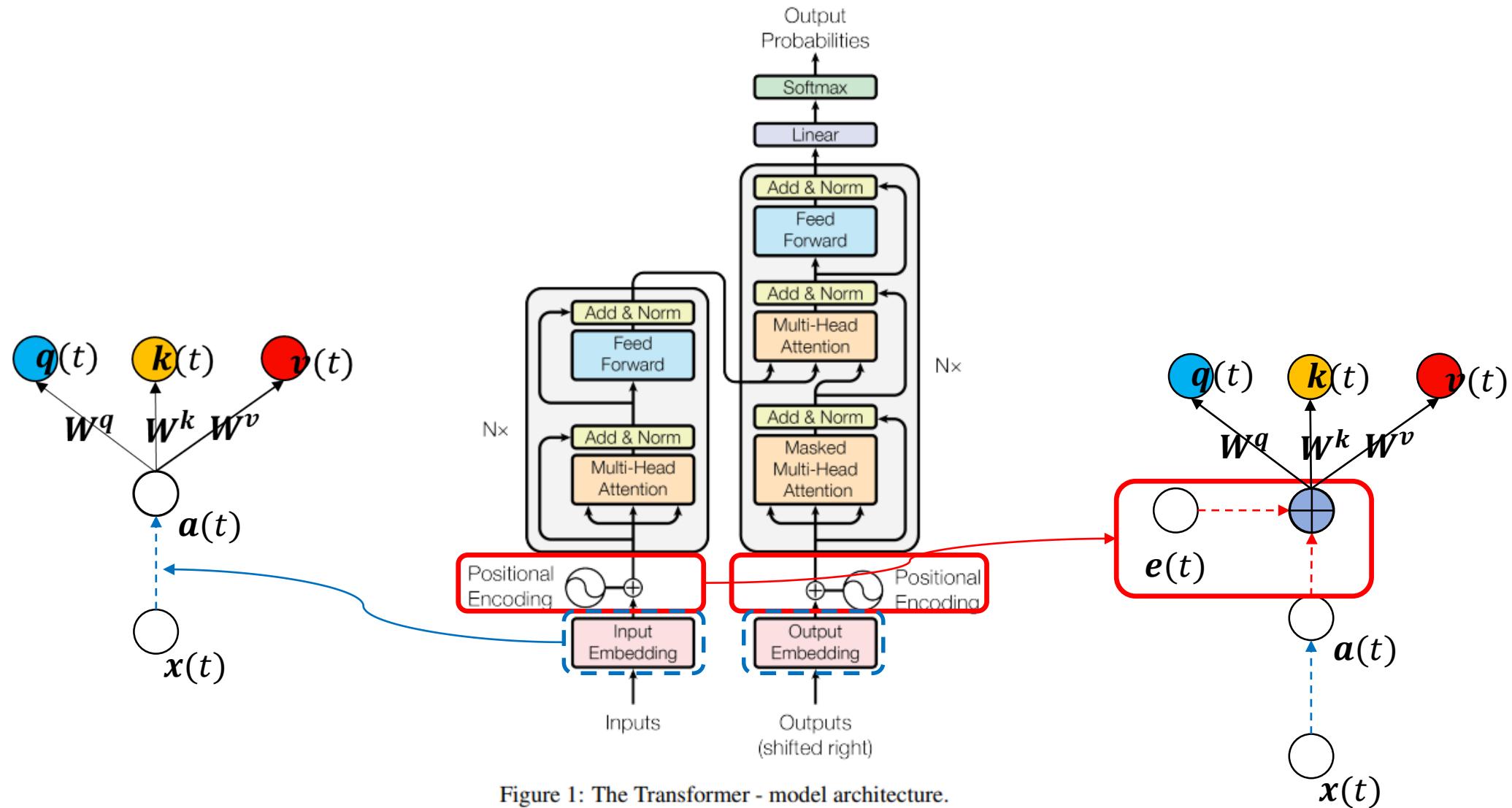
輸入是有訊續的 $x(0), x(1), \dots, x(T)$ ，但在前面介紹的Attention是沒有考慮順序性的。



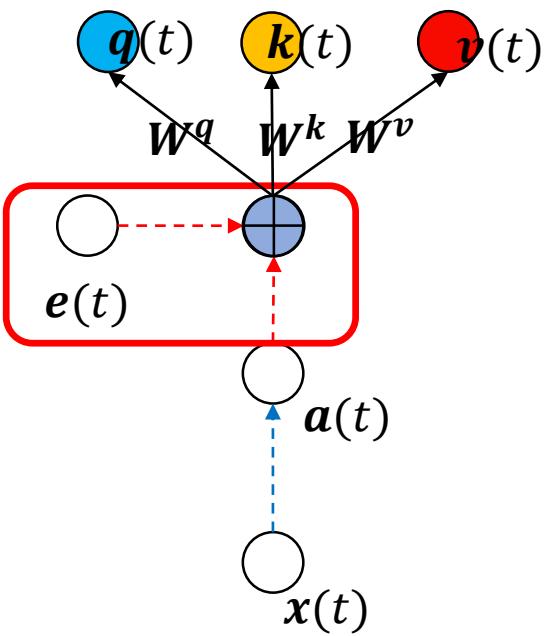
因為架構沒有順序(位置)資料，所以在Transformer中每一個位置給予一個positional vector( $e_i, \forall i = 0, 1, \dots, T$ )，整個過程稱為position encoding。



# Positional Encoding



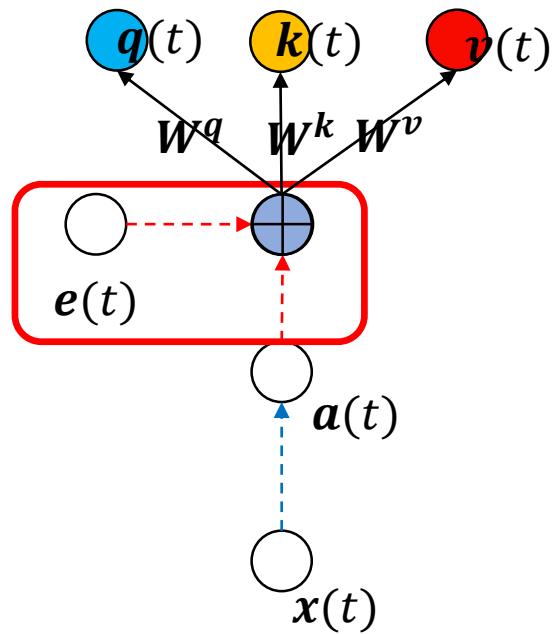
# Positional Encoding



$$\begin{aligned}
 & \mathbf{a}(1) + \mathbf{e}(1) = \begin{bmatrix} a_1 \\ a_2 \\ a_3 \\ a_4 \end{bmatrix}_{t=1} + \begin{bmatrix} e_1 \\ e_2 \\ e_3 \\ e_4 \end{bmatrix}_{t=1} \\
 & \mathbf{a}(2) + \mathbf{e}(2) = \begin{bmatrix} a_1 \\ a_2 \\ a_3 \\ a_4 \end{bmatrix}_{t=2} + \begin{bmatrix} e_1 \\ e_2 \\ e_3 \\ e_4 \end{bmatrix}_{t=2} \\
 & \mathbf{a}(3) + \mathbf{e}(3) = \begin{bmatrix} a_1 \\ a_2 \\ a_3 \\ a_4 \end{bmatrix}_{t=3} + \begin{bmatrix} e_1 \\ e_2 \\ e_3 \\ e_4 \end{bmatrix}_{t=3} \\
 & = \begin{bmatrix} \mathbf{e}(1) \\ e_1 \\ e_2 \\ e_3 \\ e_4 \end{bmatrix}_{t=1} = \begin{bmatrix} \mathbf{a}(1) \\ a_1 \\ a_2 \\ a_3 \\ a_4 \end{bmatrix}_{t=1} \quad = \begin{bmatrix} \mathbf{e}(2) \\ e_1 \\ e_2 \\ e_3 \\ e_4 \end{bmatrix}_{t=1} = \begin{bmatrix} \mathbf{a}(2) \\ a_1 \\ a_2 \\ a_3 \\ a_4 \end{bmatrix}_{t=2} \\
 & = \begin{bmatrix} \mathbf{e}(3) \\ e_1 \\ e_2 \\ e_3 \\ e_4 \end{bmatrix}_{t=1} = \mathbf{a}(T=3) \quad = \begin{bmatrix} a_1 \\ a_2 \\ a_3 \\ a_4 \end{bmatrix}_{t=3}
 \end{aligned}$$



# Positional Encoding



$$\begin{aligned}
 \mathbf{a}(1) + \mathbf{e}(1) &= \begin{bmatrix} 1.1 \\ 0.2 \\ 0.3 \\ 0.4 \end{bmatrix} \\
 \mathbf{e}(1) &= \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} & \mathbf{a}(1) &= \begin{bmatrix} 0.1 \\ 0.2 \\ 0.3 \\ 0.4 \end{bmatrix} \\
 \mathbf{a}(2) + \mathbf{e}(2) &= \begin{bmatrix} 0.2 \\ 1.3 \\ 0.2 \\ 0.3 \end{bmatrix} \\
 \mathbf{e}(2) &= \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} & \mathbf{a}(2) &= \begin{bmatrix} 0.2 \\ 0.3 \\ 0.2 \\ 0.3 \end{bmatrix} \\
 \mathbf{a}(3) + \mathbf{e}(3) &= \begin{bmatrix} 0.1 \\ 0.1 \\ 1.1 \\ 0.1 \end{bmatrix} \\
 \mathbf{e}(3) &= \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} & \mathbf{a}(3) &= \begin{bmatrix} 0.1 \\ 0.1 \\ 0.1 \\ 0.1 \end{bmatrix}
 \end{aligned}$$



# Add & Norm

**Multi-Head Attention**

$$O = \text{MultiHeadAttent}(X)$$

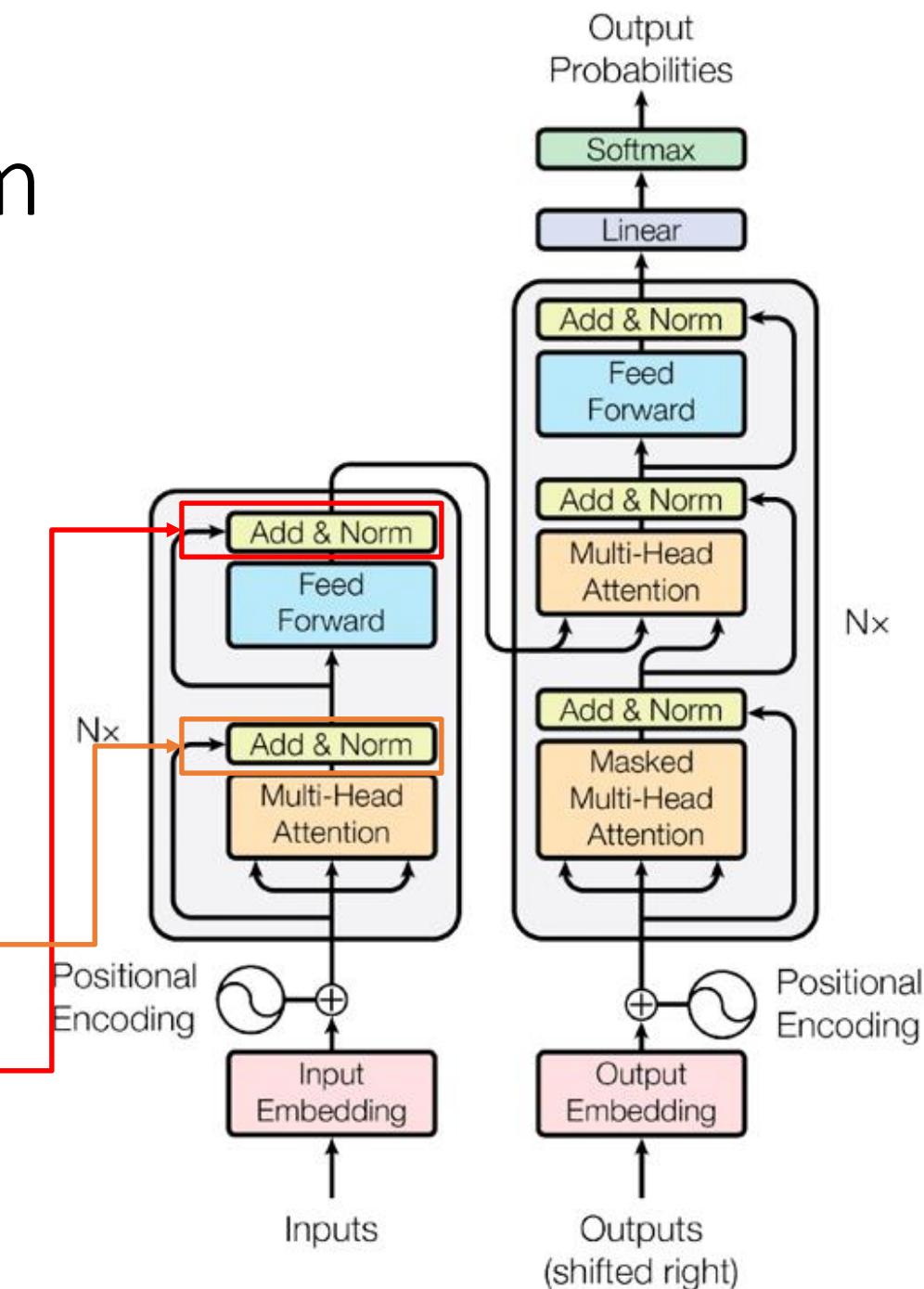
**Add**

$$X + \text{MultiHeadAttent}(X)$$

**Add & Norm**

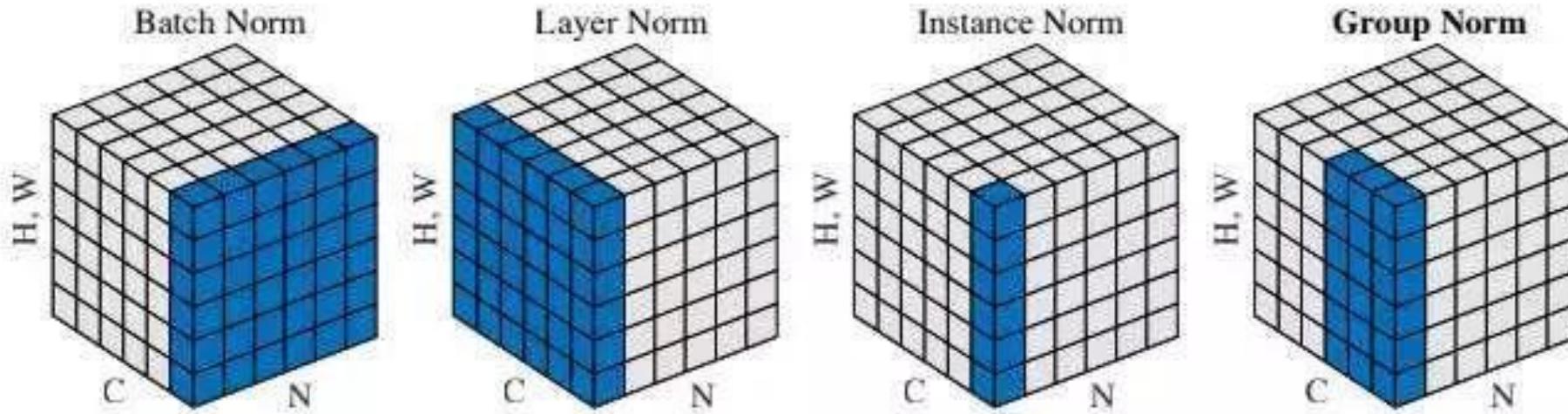
$$\text{LayerNorm}(X + \text{MultiHeadAttent}(X))$$

$$\text{LayerNorm}(X + \text{FeedForward}(X))$$

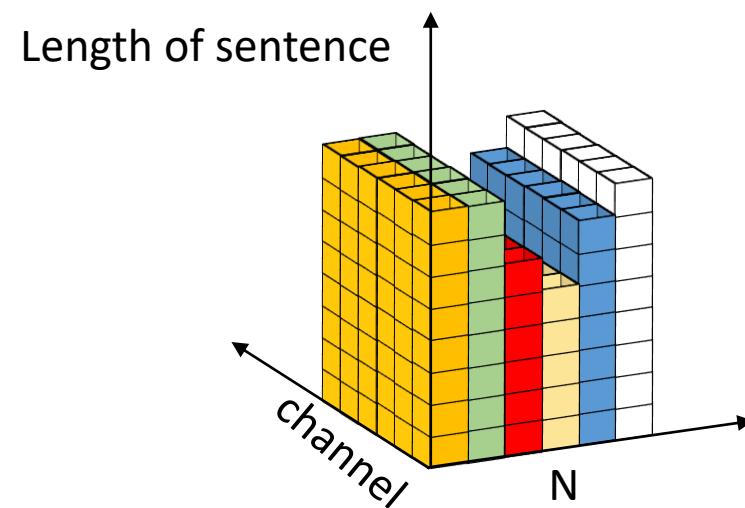


# Normalization

$$y = \frac{x - E[x]}{\sqrt{Var[x] + \epsilon}} * \gamma + \beta$$



BN: 針對每個channel的batch data所有分布來計算

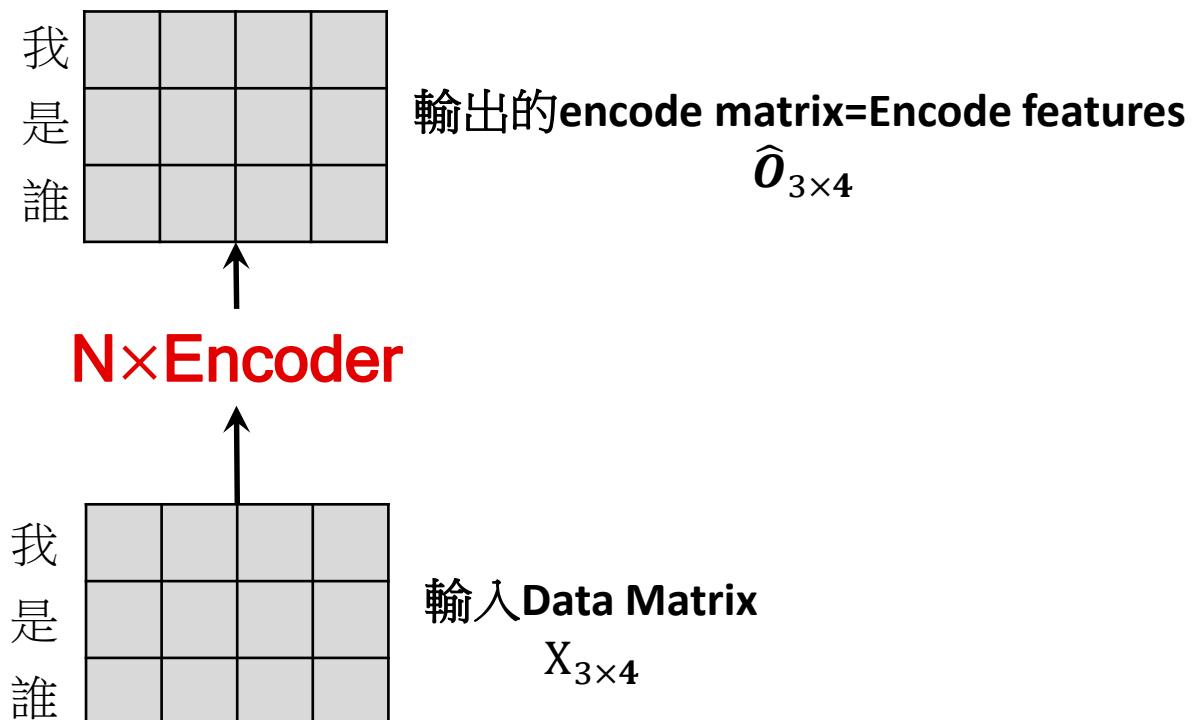


語言模型的句子長度不固定  
(每筆資料的Z軸實際大小不一)  
用Layer Norm會比較合適。

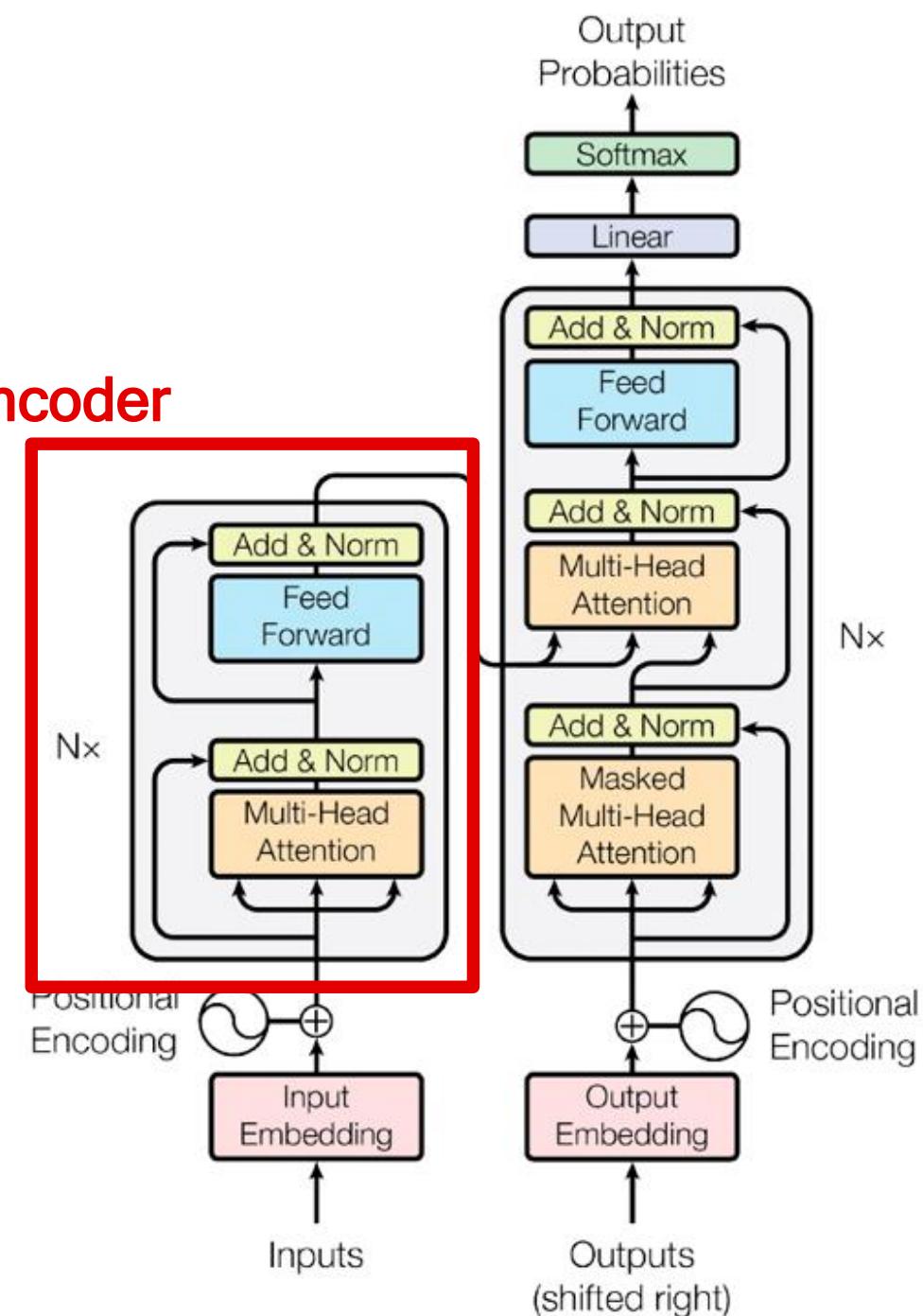
常用於RNN結構內



# Encoder Block



## Encoder





# Decoder Block

兩個Multi-Head Attention

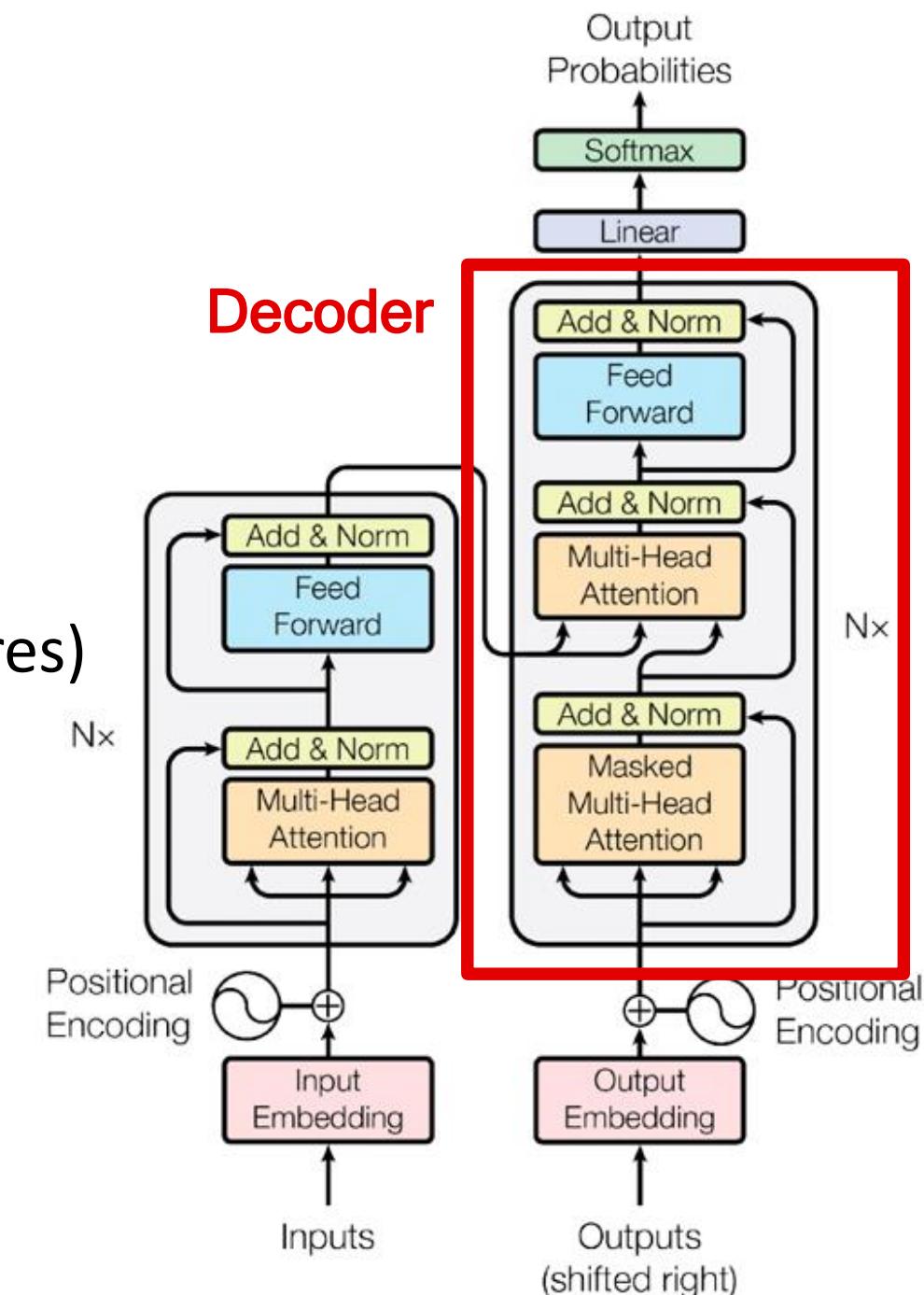
1. Masked Multi-Head Attention

2. Multi-Head Attention (包含Encode的features)

Q: Masked Multi-Head Attention 的輸出

K: Encode的features

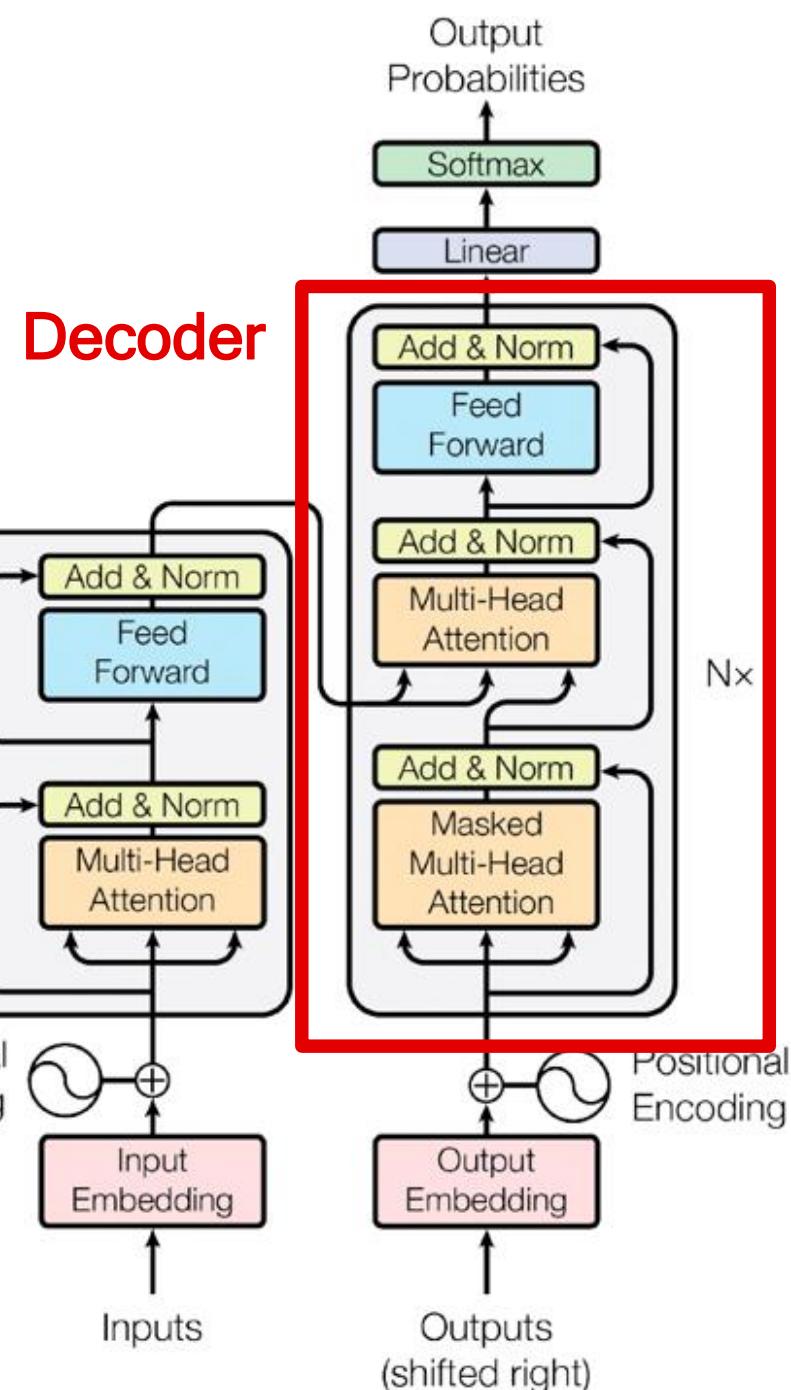
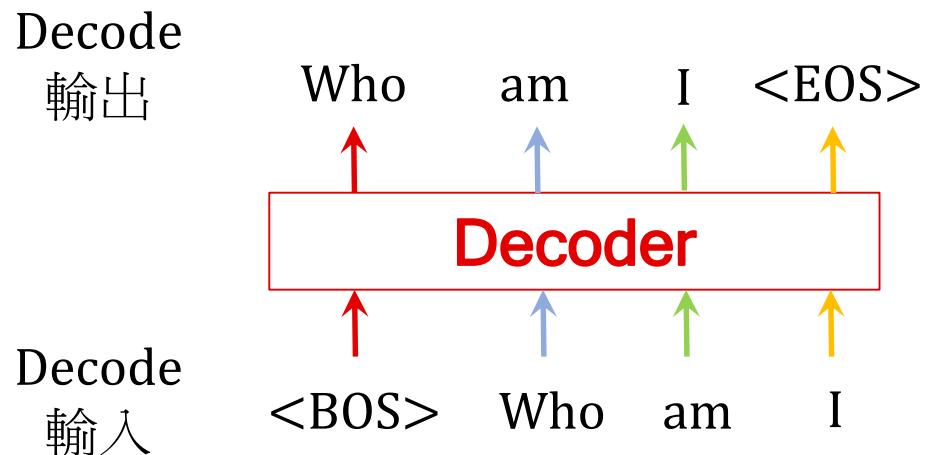
V: Encode的features



# 1. Masked Multi-Head Attention

Hugging Face is a startup based in New York City and Paris  
 p(word)

<https://huggingface.co/docs/transformers/perplexity>





Context1:

<BOS> I am Tommy <EOS>



1. <BOS> → I
2. <BOS> I → am
3. <BOS> I am → Tommy
4. <BOS> I am Tommy →<EOS>

還有很多組合

Context2:

<BOS> You are fine <EOS>



1. <BOS> → You
2. <BOS> You → are
3. <BOS> You are → fine
4. <BOS> You are fine →<EOS>

**N-gram**

$$p(x_1, x_2, \dots, x_L) = \prod_{i=1}^L p(x_i | x_1, \dots, x_{i-1}) \approx p(x_i | x_1, \dots, x_{i-1})$$

N=1

$$p(x_1, x_2, \dots, x_L) = \prod_{i=1}^L p(x_i | x_1, \dots, x_{i-1}) \approx p(x_i)$$

N=2

$$p(x_1, x_2, \dots, x_L) = \prod_{i=1}^L p(x_i | x_1, \dots, x_{i-1}) \approx p(x_i | x_{i-1})$$

N=3

$$p(x_1, x_2, \dots, x_L) = \prod_{i=1}^L p(x_i | x_1, \dots, x_{i-1}) \approx p(x_i | x_{i-1}, x_{i-2})$$





## 以2-gram為例

Context1:

<BOS> I am Tommy <EOS>



1. <BOS> → I
2. I → am
3. am → Tommy
4. Tommy → <EOS>

Context2:

<BOS> You are fine <EOS>



1. <BOS> → You
2. You → are
3. are → fine
4. fine → <EOS>

## 以3-gram為例

Context1:

<BOS> I am Tommy <EOS>



1. <BOS> I → am
2. I am → Tommy
3. am Tommy → <EOS>

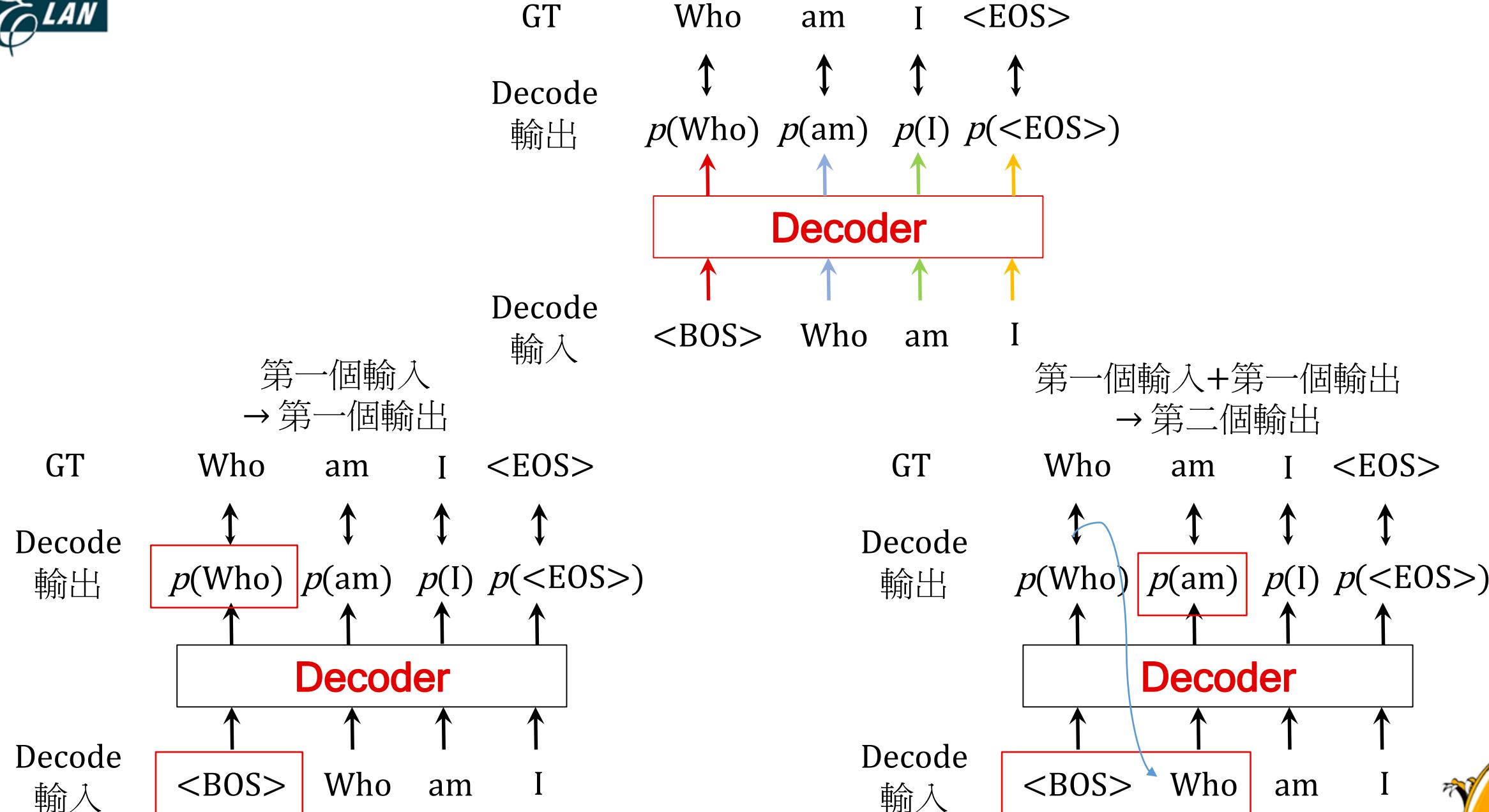
Context2:

<BOS> You are fine <EOS>



1. <BOS> You → are
2. You are → fine
3. are fine → <EOS>





# 1. Masked Multi-Head Attention

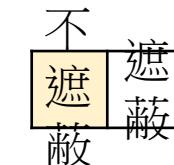
|         |  |  |  |  |
|---------|--|--|--|--|
| < BOS > |  |  |  |  |
| 我       |  |  |  |  |
| 是       |  |  |  |  |
| 誰       |  |  |  |  |
| < EOS > |  |  |  |  |

輸入矩陣  $X$ 

|         | $\wedge$ BOS $\wedge$ | $QK^T$ | $\wedge$ EOS $\wedge$ |
|---------|-----------------------|--------|-----------------------|
|         | $\vee$                | 我是誰    | $\vee$                |
| < BOS > |                       |        |                       |
| 我       |                       |        |                       |
| 是       |                       |        |                       |
| 誰       |                       |        |                       |
| < EOS > |                       |        |                       |

◎ Hadamard product

我是誰應該要寫 Who am I  
但不好呈現，所以我用中文取代



得到 Mask  $QK^T$  後和  $V$  內積後  
就可以繼續算輸出的  $O$

$$\odot \begin{array}{c} < \text{BOS} > \\ \text{我} \\ \text{是} \\ \text{誰} \\ < \text{EOS} > \end{array} \begin{array}{c} < \text{BOS} > \\ \text{我} \\ \text{是} \\ \text{誰} \\ < \text{EOS} > \end{array} \begin{array}{c} < \text{BOS} > \\ \text{我} \\ \text{是} \\ \text{誰} \\ < \text{EOS} > \end{array}$$

Mask 矩陣

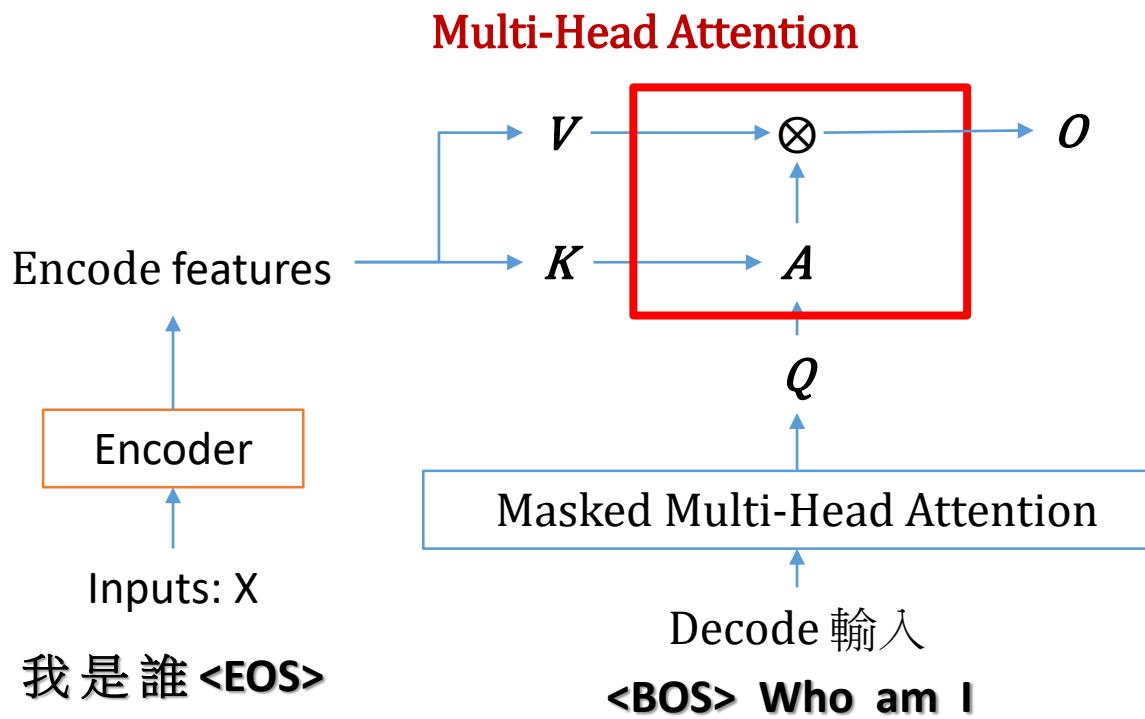
$$= \begin{array}{c} < \text{BOS} > \\ \text{我} \\ \text{是} \\ \text{誰} \\ < \text{EOS} > \end{array} \begin{array}{c} < \text{BOS} > \\ \text{我} \\ \text{是} \\ \text{誰} \\ < \text{EOS} > \end{array} \begin{array}{c} < \text{BOS} > \\ \text{我} \\ \text{是} \\ \text{誰} \\ < \text{EOS} > \end{array}$$

Mask  $QK^T$



## 2. Multi-Head Attention (包含Encode的features)

- Q: Masked Multi-Head Attention 的輸出
- K: Encode的features
- V: Encode的features



|         |  |  |  |
|---------|--|--|--|
| < BOS > |  |  |  |
| Who     |  |  |  |
| am      |  |  |  |
| I       |  |  |  |
| < EOS > |  |  |  |

Q

|                 |  |  |  |
|-----------------|--|--|--|
| 我               |  |  |  |
| 是               |  |  |  |
| 誰               |  |  |  |
| < EOS >         |  |  |  |
| Encode features |  |  |  |

$O_{T \times d}$  : 將V從K空間投影到Q空間

所以將Encode feature投影到Q空間

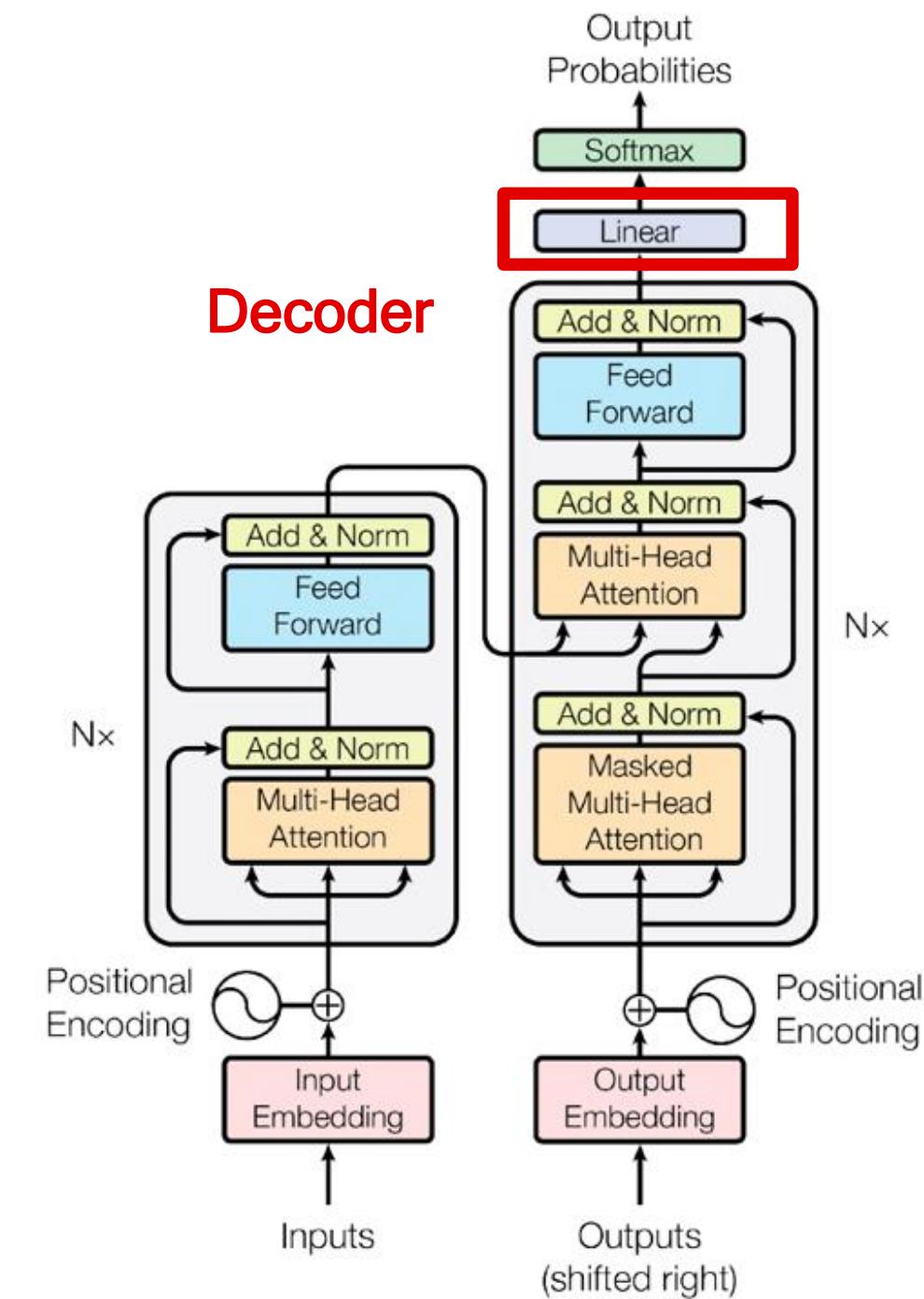
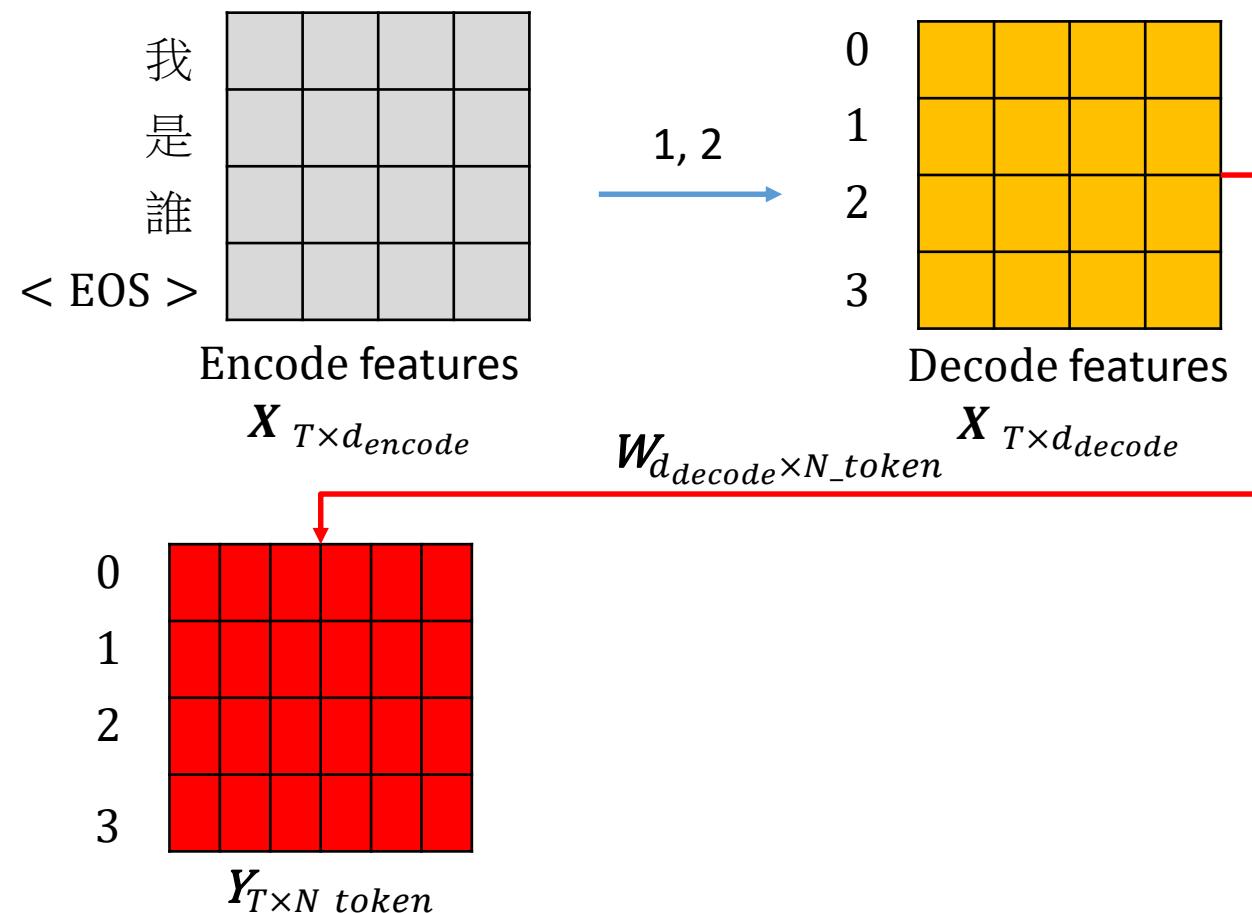
將中文的特徵投影到英文的特徵空間



# Decoder Block

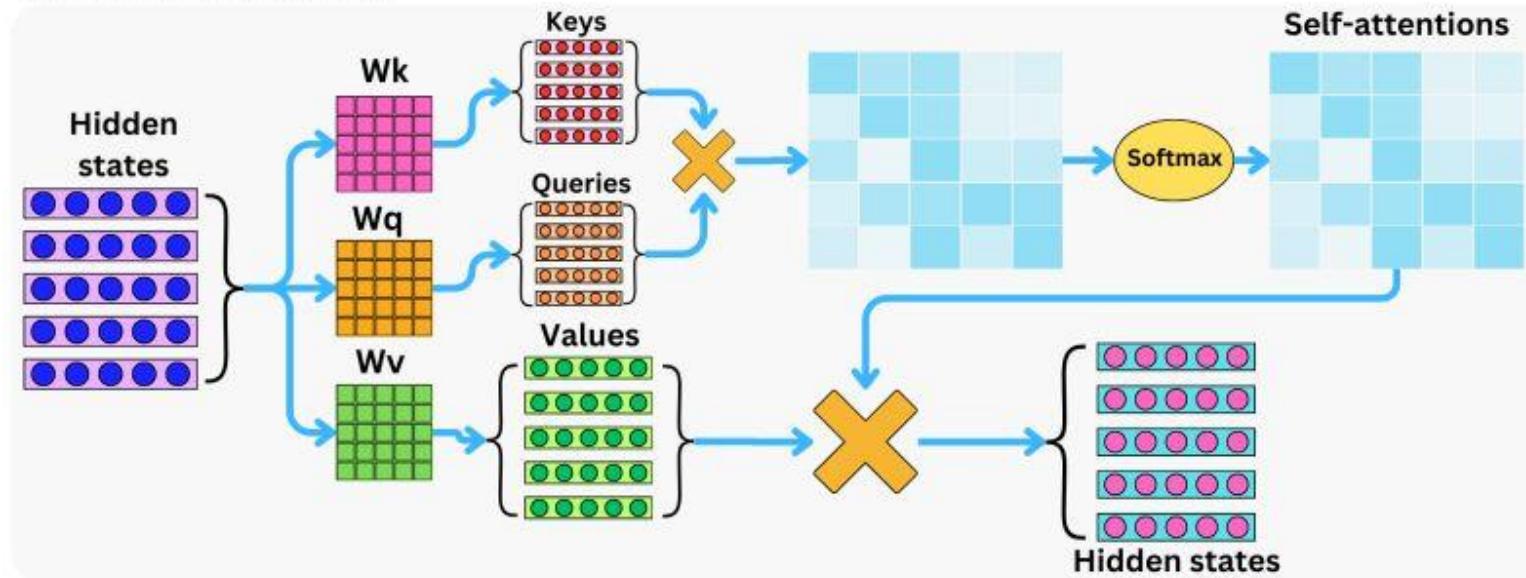
兩個Multi-Head Attention

1. Masked Multi-Head Attention
2. Multi-Head Attention (包含Encode的features)

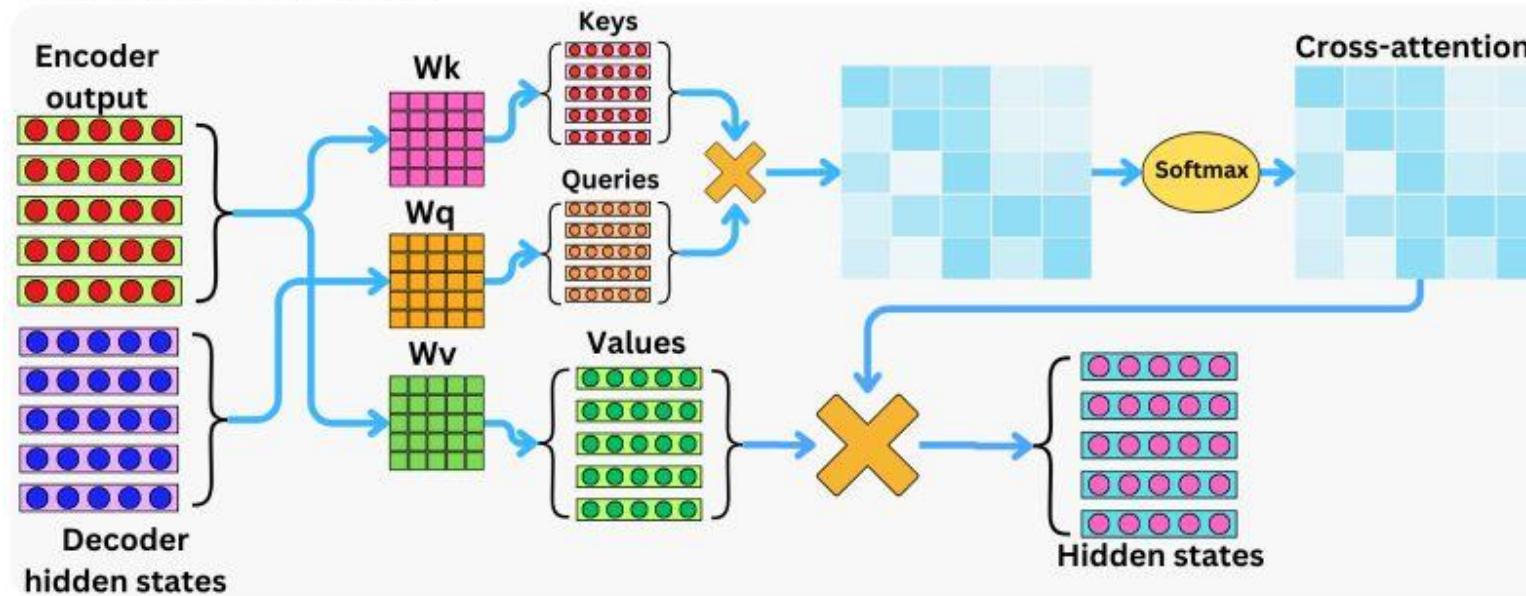


# Self-Attention VS Cross-Attention

## The Self-Attentions



## The Cross-Attentions



# Cross-attention

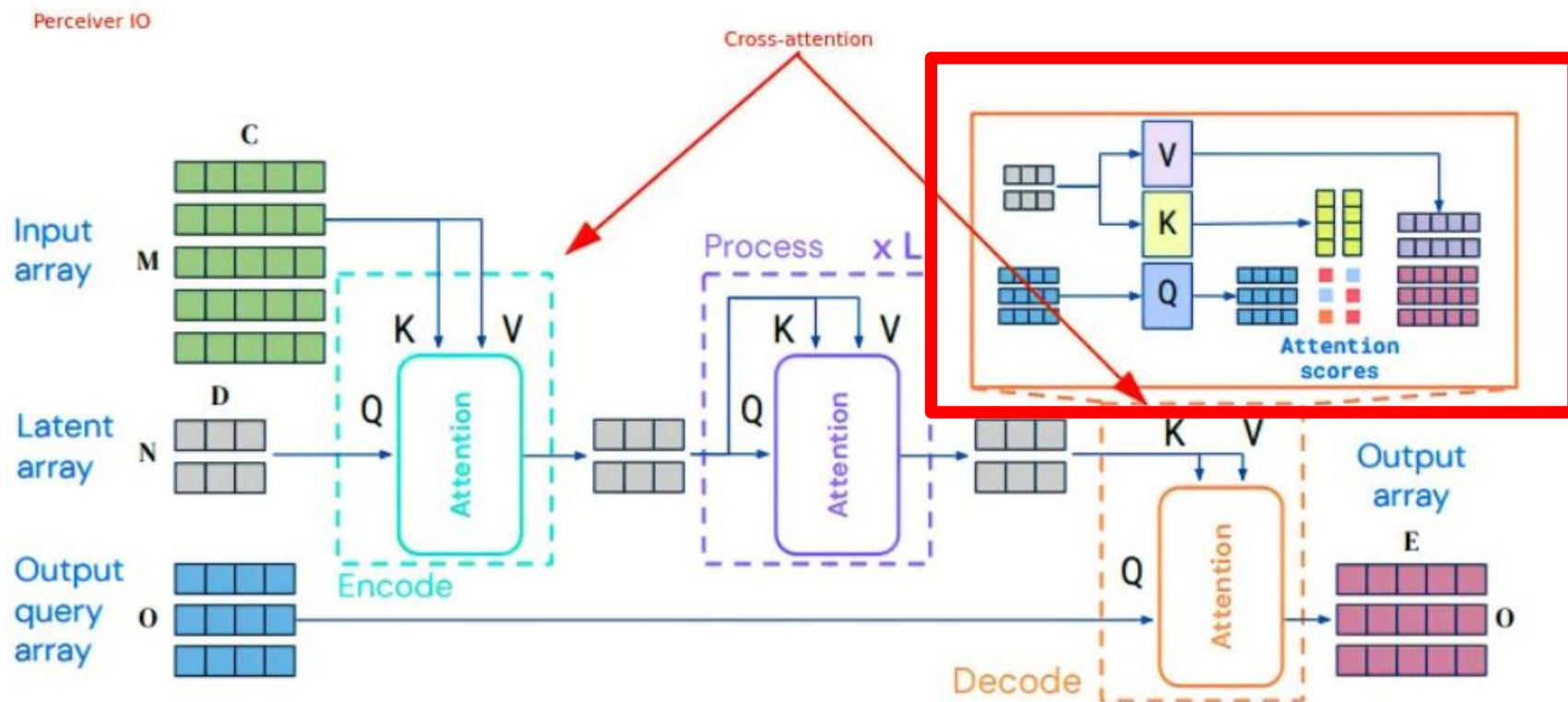
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

將V從K空間投影到Q空間

K空間: 文字 Q空間: 影像

K空間: Sensor 1 Q空間: Sensor 2

Self-attention: 輸入只有同一個序列資料  
 Cross-attention: 輸入來至不同的序列資料



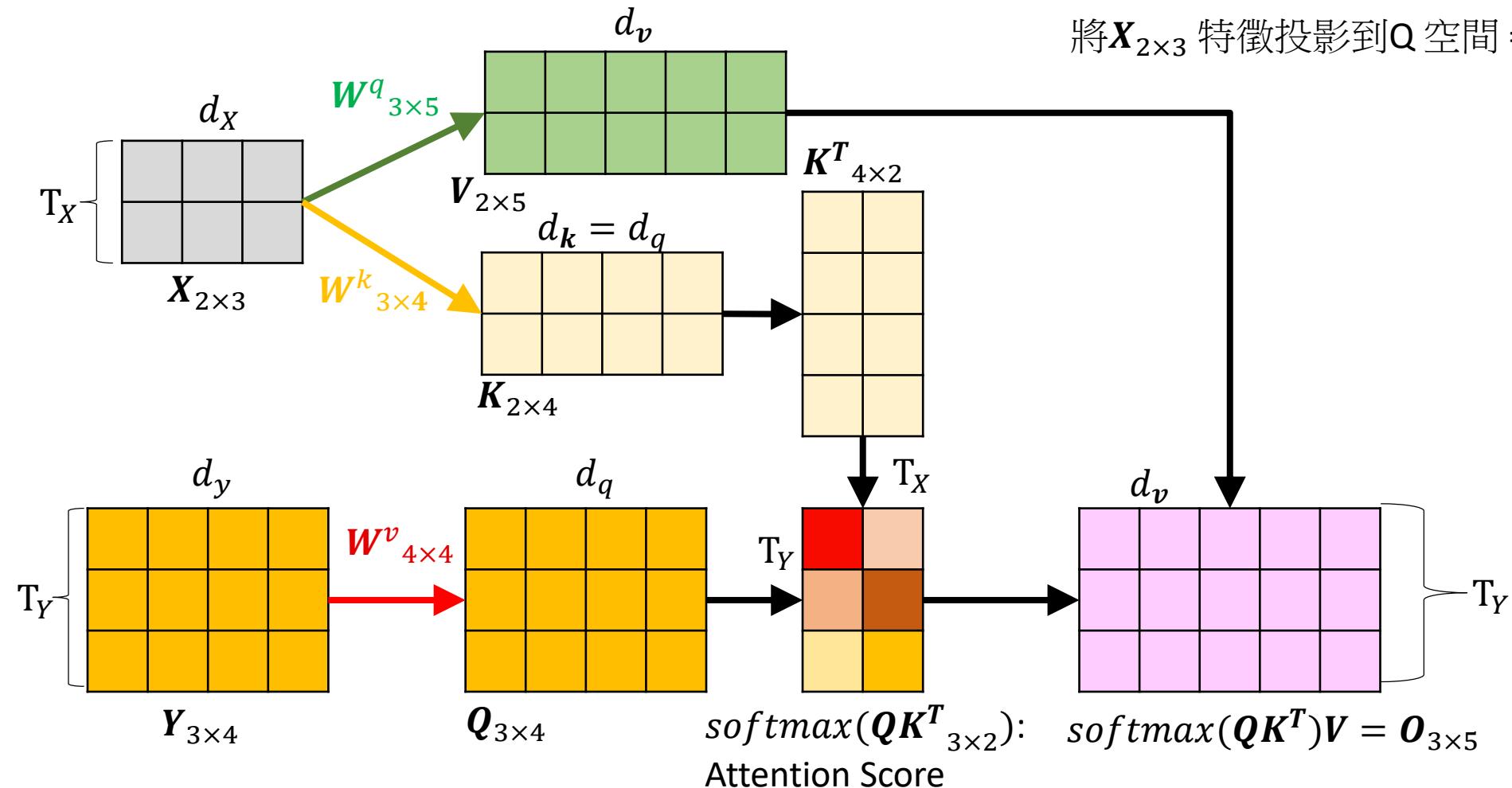
# Cross-attention

兩個不同的輸入源:  $X_{2 \times 3}$   $Y_{3 \times 4}$   
 輸出的序列  
 長度會跟  $Y$  一致  
 維度會跟  $V$  一致(可在投影時控制)

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

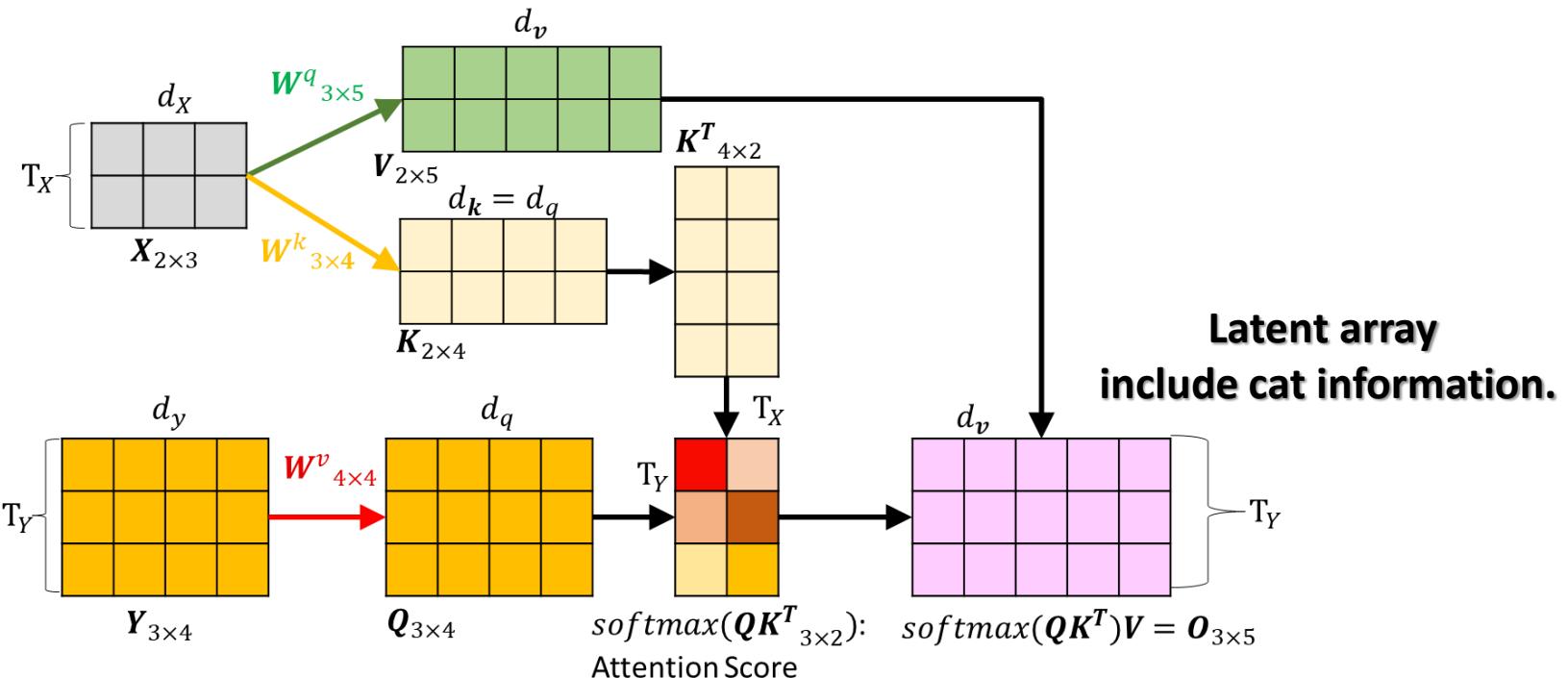
將  $V$  從  $K$  空間投影到  $Q$  空間

將  $X_{2 \times 3}$  特徵投影到  $Q$  空間 =  $O_{3 \times 5}$

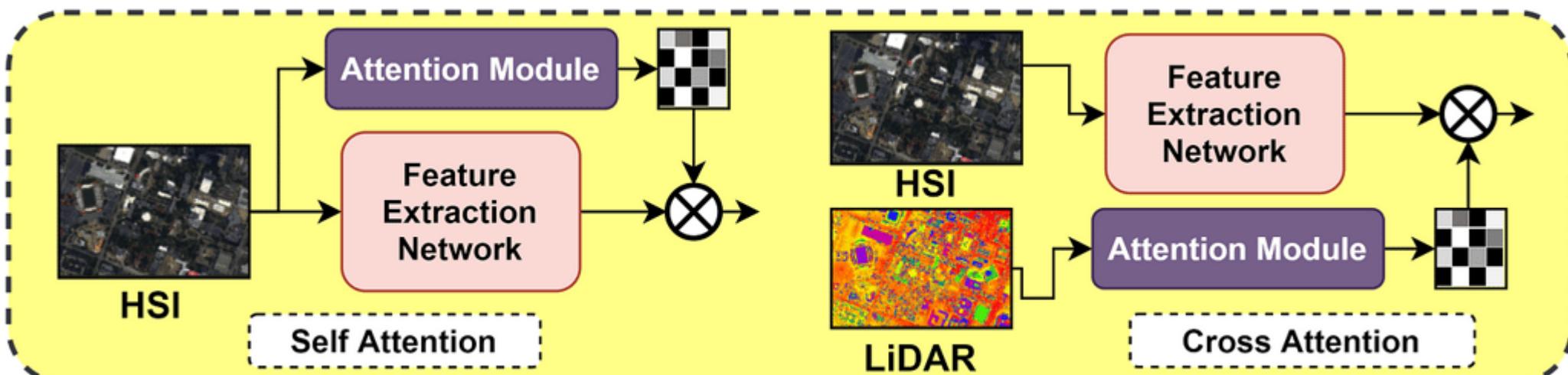


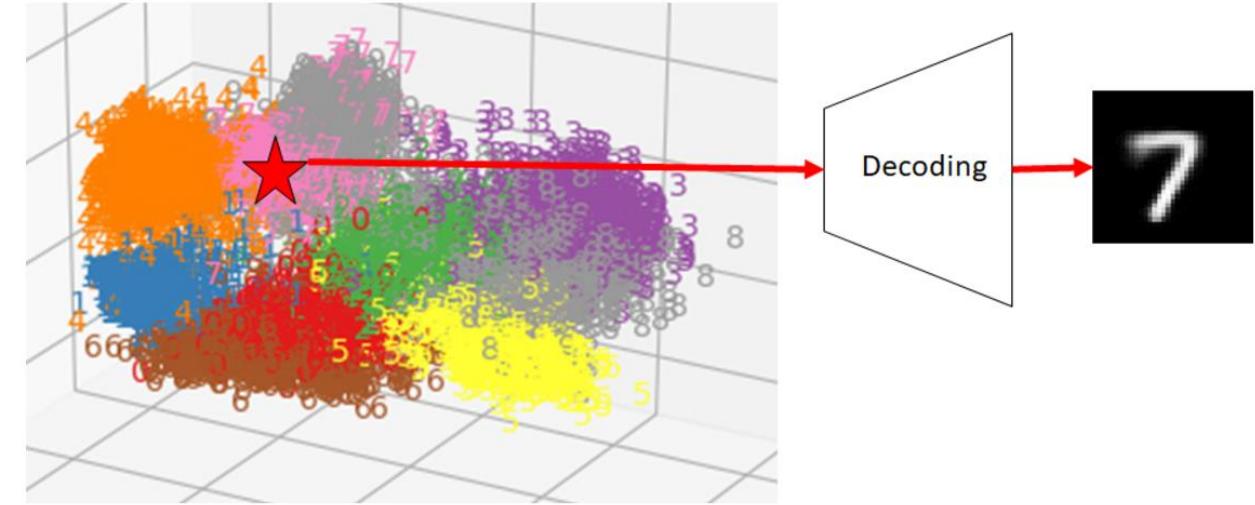
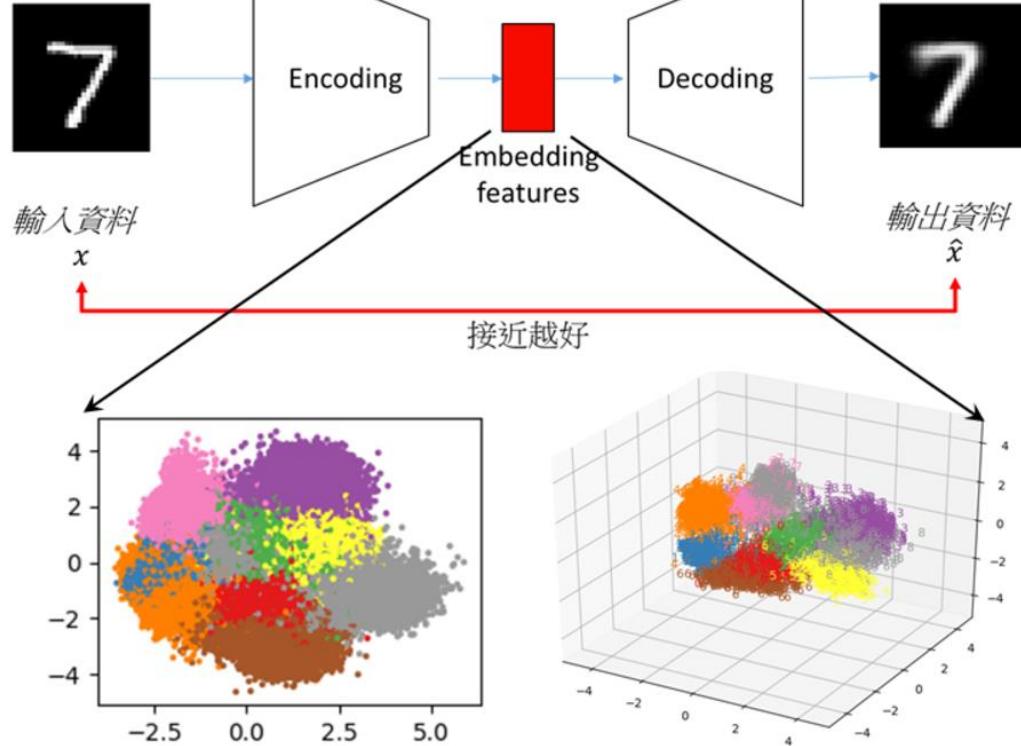
## Diffusion Model

幫我生成一隻貓



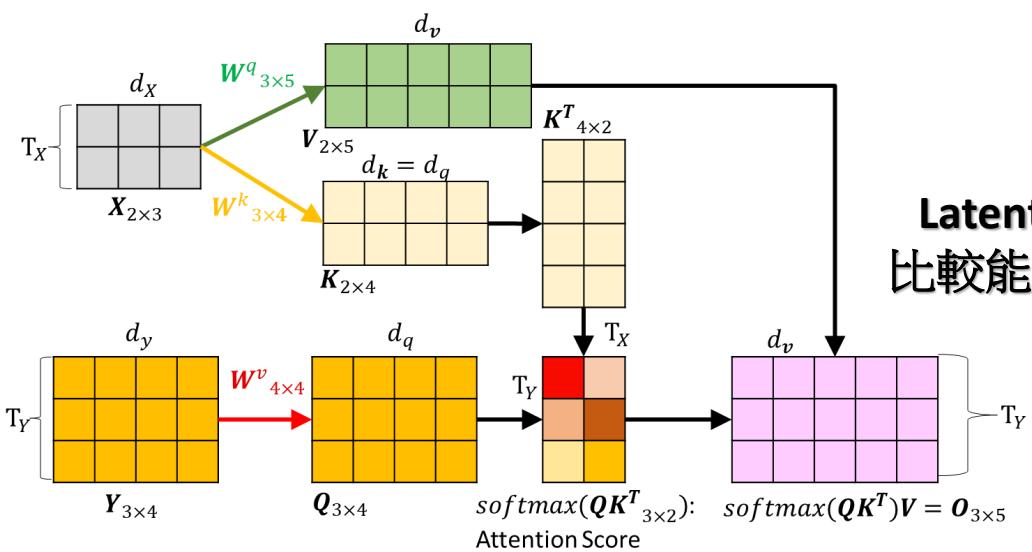
<https://www.genspark.ai/spark/cross-attention-in-transformer-architecture/58187f91-5956-3099-a0ec-2c084e963efb>





**Generate latent**

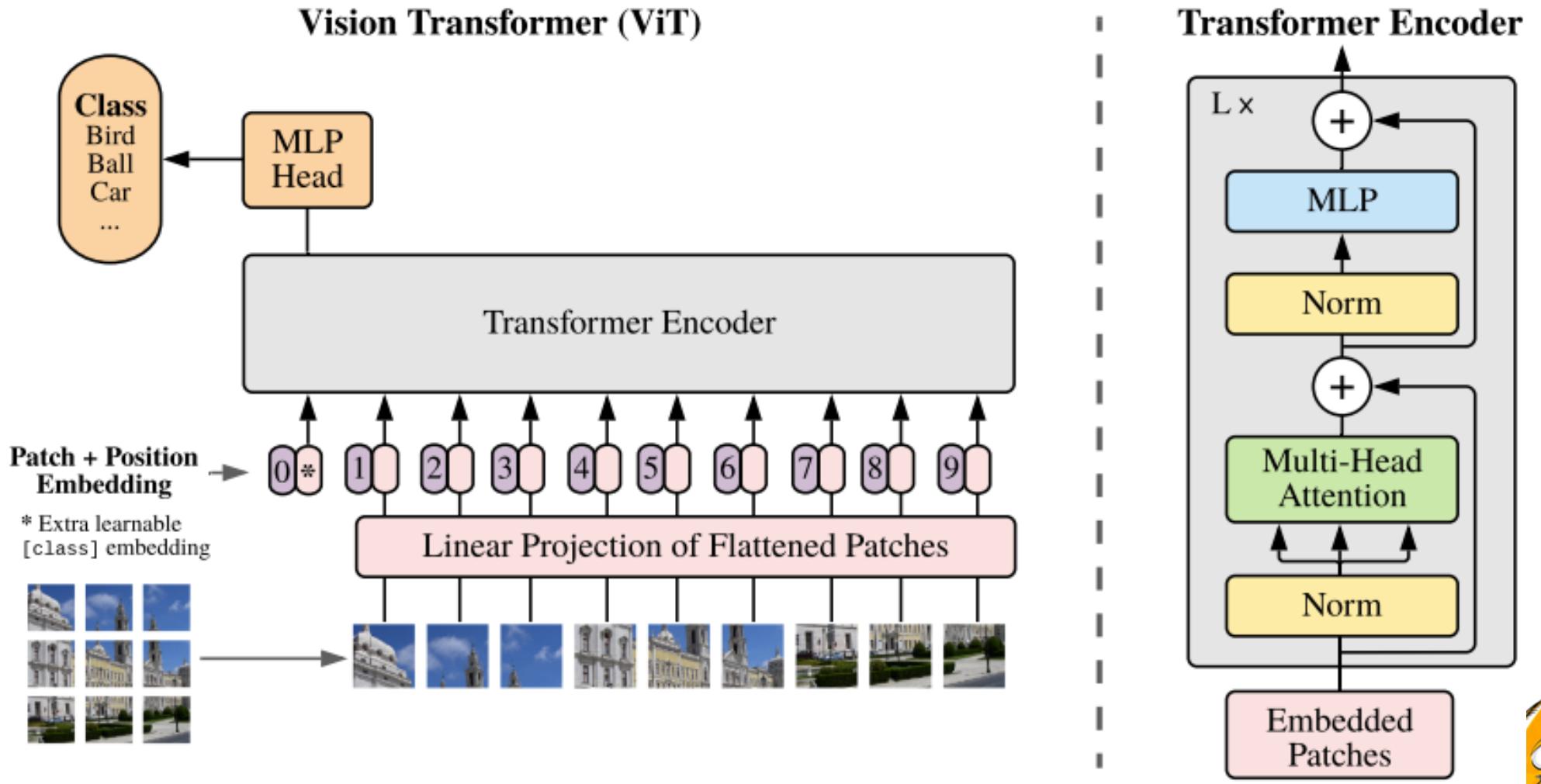
幫我數字7



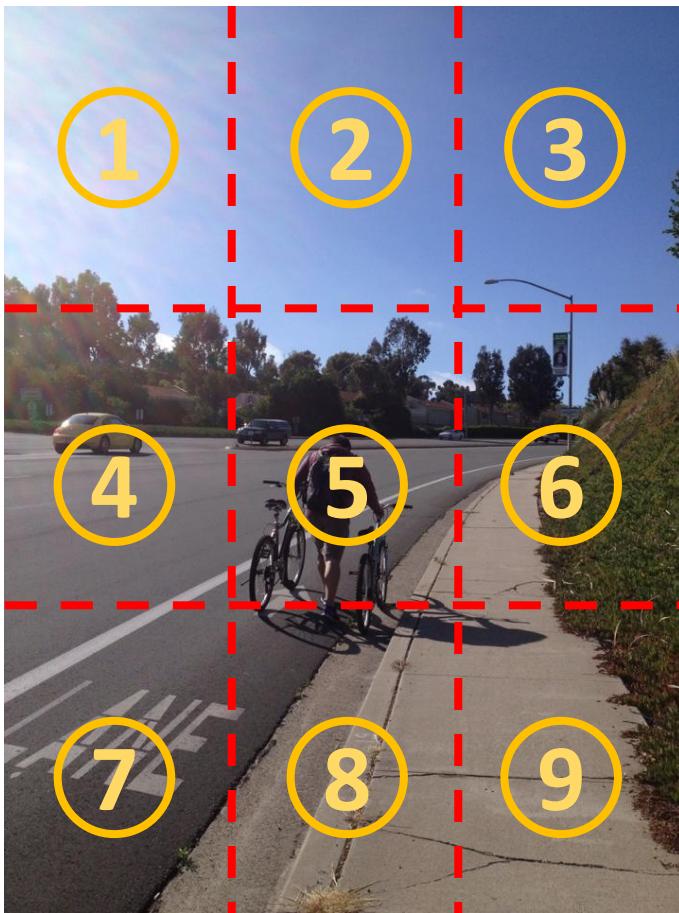
**Latent array加權後  
比較能取到七的位置**



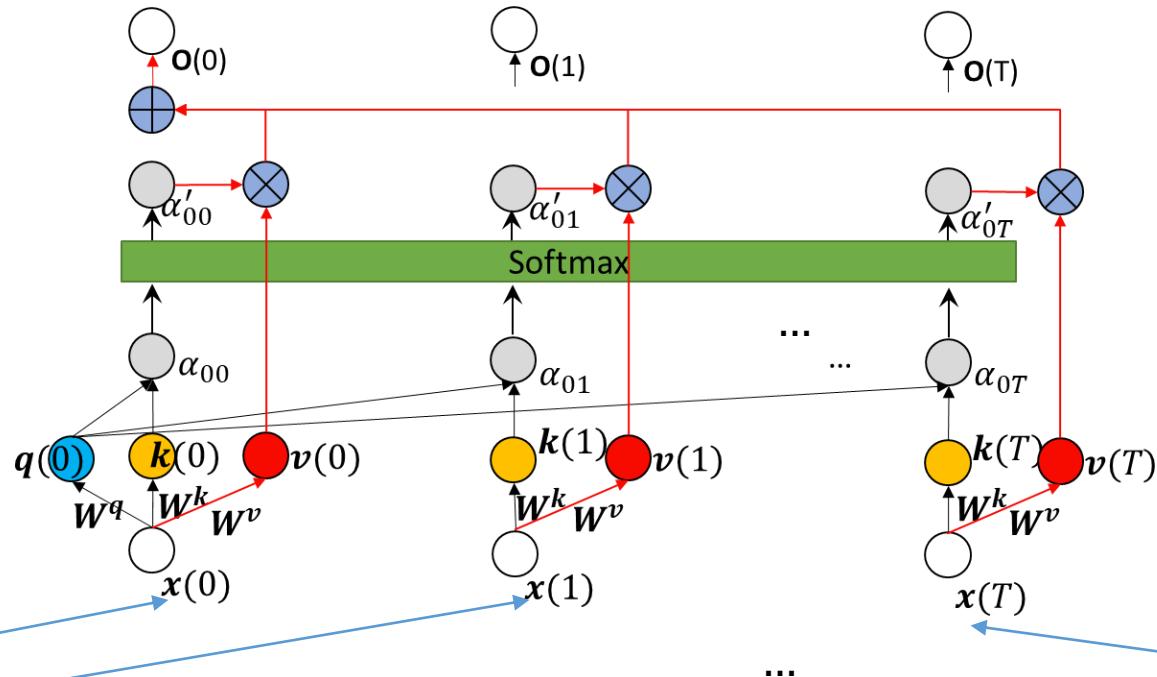
# Vision Transformer(ViT)



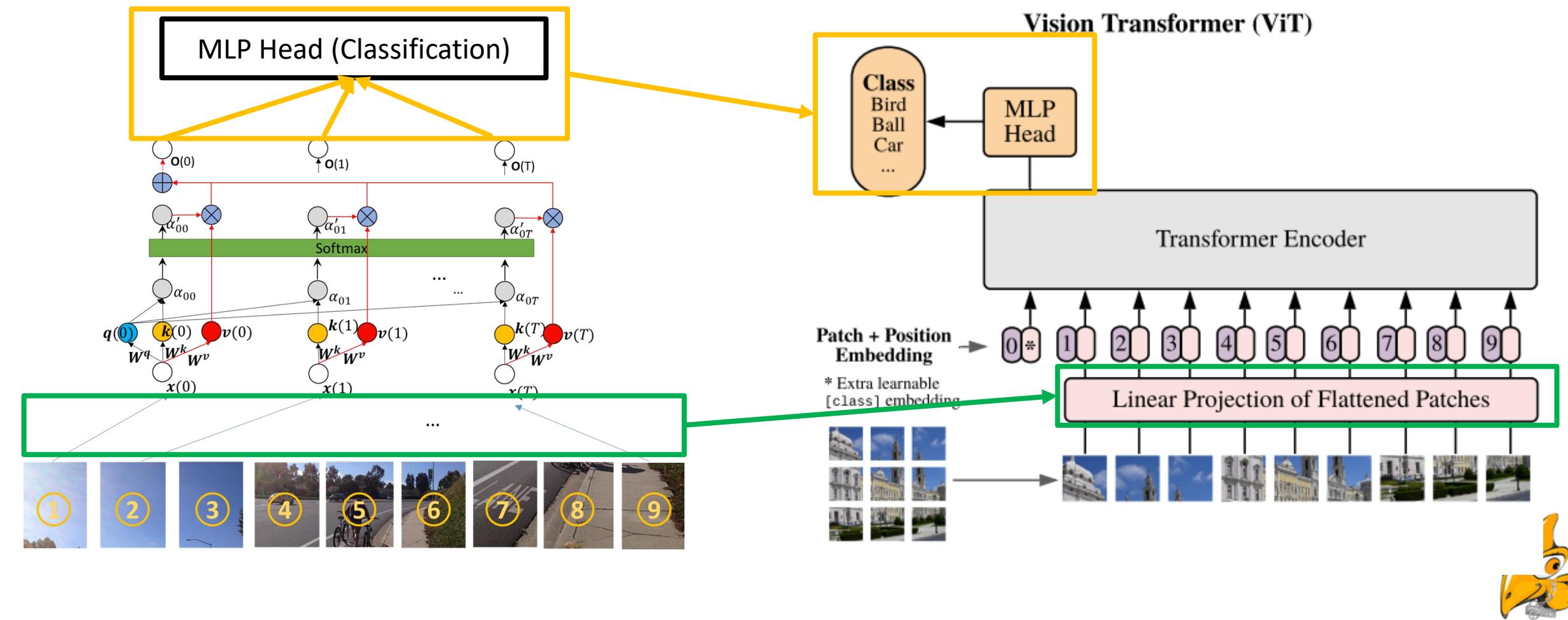
# Vision Transformer(ViT)



# Vision Transformer(ViT)



# Vision Transformer(ViT)

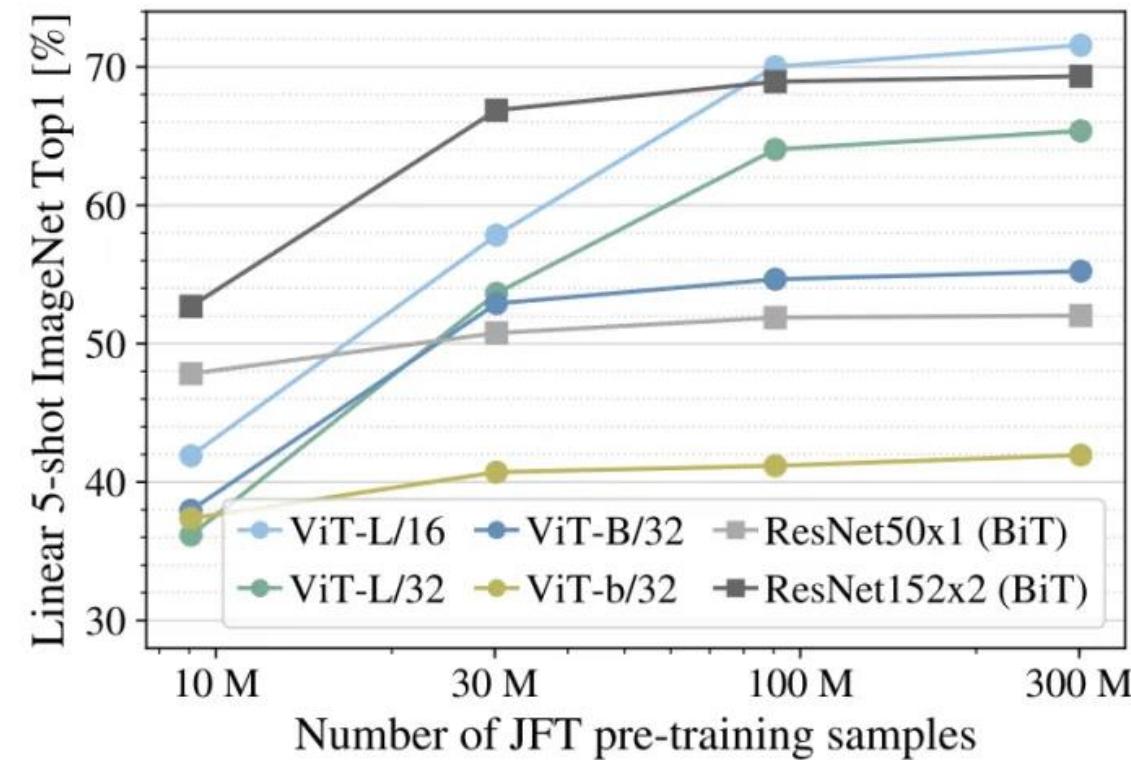


| Model     | Layers | Hidden size $D$ | MLP size | Heads | Params |
|-----------|--------|-----------------|----------|-------|--------|
| ViT-Base  | 12     | 768             | 3072     | 12    | 86M    |
| ViT-Large | 24     | 1024            | 4096     | 16    | 307M   |
| ViT-Huge  | 32     | 1280            | 5120     | 16    | 632M   |

Table 1: Details of Vision Transformer model variants.

ViT-L/16表示的是large model使用的是16\*16的patch

|                    | Ours-JFT<br>(ViT-H/14) | Ours-JFT<br>(ViT-L/16) | Ours-I21K<br>(ViT-L/16) | BiT-L<br>(ResNet152x4) | Noisy Student<br>(EfficientNet-L2) |
|--------------------|------------------------|------------------------|-------------------------|------------------------|------------------------------------|
| ImageNet           | <b>88.55</b> ± 0.04    | 87.76 ± 0.03           | 85.30 ± 0.02            | 87.54 ± 0.02           | 88.4 / 88.5*                       |
| ImageNet ReaL      | <b>90.72</b> ± 0.05    | 90.54 ± 0.03           | 88.62 ± 0.05            | 90.54                  | 90.55                              |
| CIFAR-10           | <b>99.50</b> ± 0.06    | 99.42 ± 0.03           | 99.15 ± 0.03            | 99.37 ± 0.06           | —                                  |
| CIFAR-100          | <b>94.55</b> ± 0.04    | 93.90 ± 0.05           | 93.25 ± 0.05            | 93.51 ± 0.08           | —                                  |
| Oxford-IIIT Pets   | <b>97.56</b> ± 0.03    | 97.32 ± 0.11           | 94.67 ± 0.15            | 96.62 ± 0.23           | —                                  |
| Oxford Flowers-102 | 99.68 ± 0.02           | <b>99.74</b> ± 0.00    | 99.61 ± 0.02            | 99.63 ± 0.03           | —                                  |
| VTAB (19 tasks)    | <b>77.63</b> ± 0.23    | 76.28 ± 0.46           | 72.72 ± 0.21            | 76.29 ± 1.70           | —                                  |
| TPUv3-core-days    | 2.5k                   | 0.68k                  | 0.23k                   | 9.9k                   | 12.3k                              |

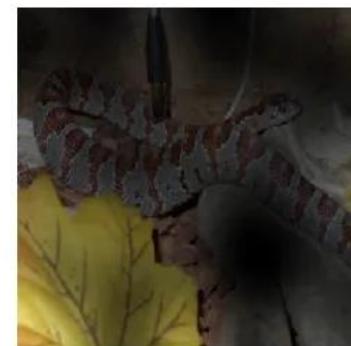
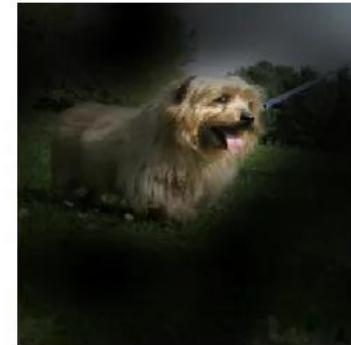


Big Transform: BiT





# Input      Attention



# CNN and Self-Attention

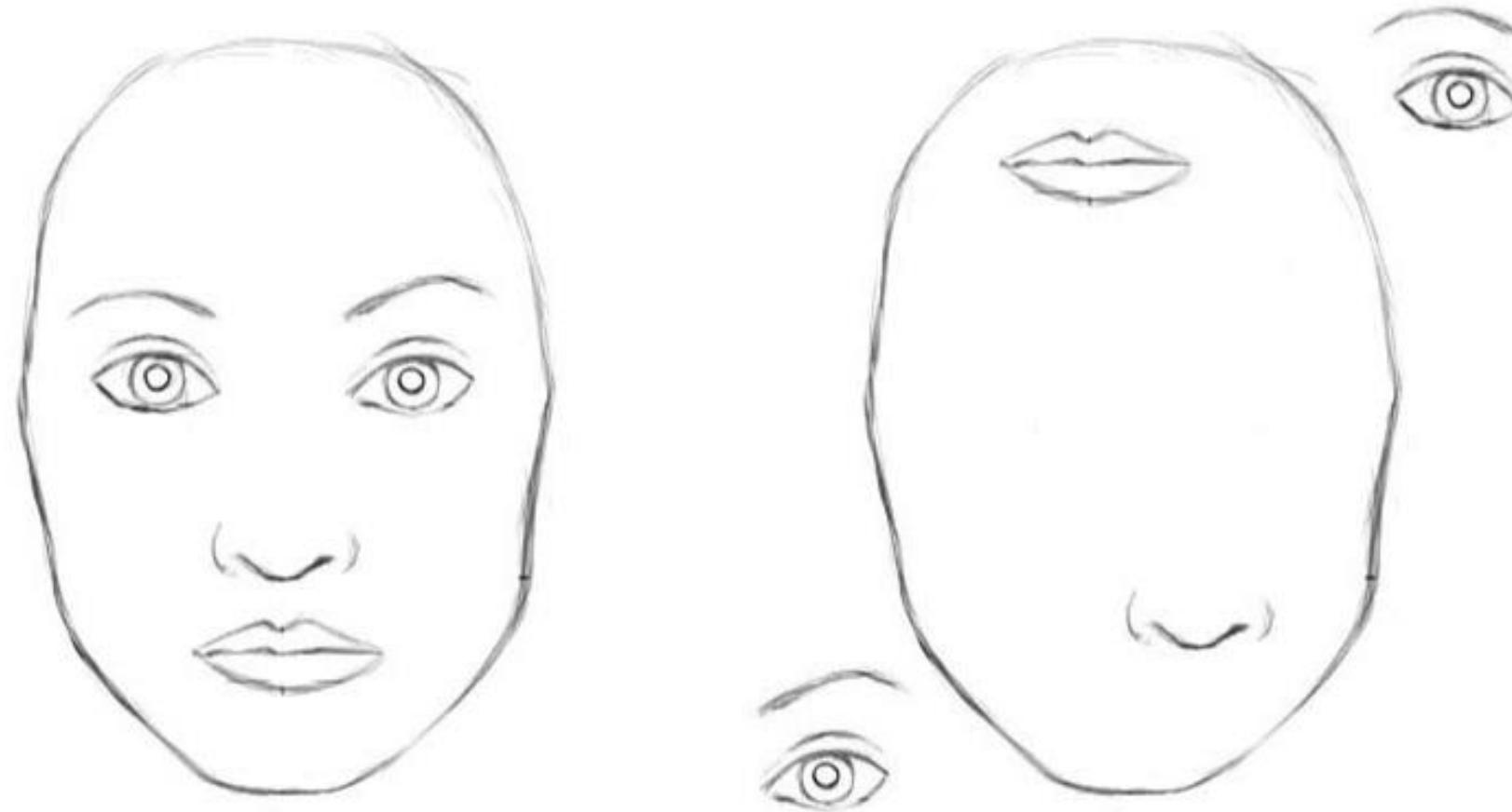


Fig 1: Both [images](#) are similar for a CNN as it does not consider the relative positioning of facial features

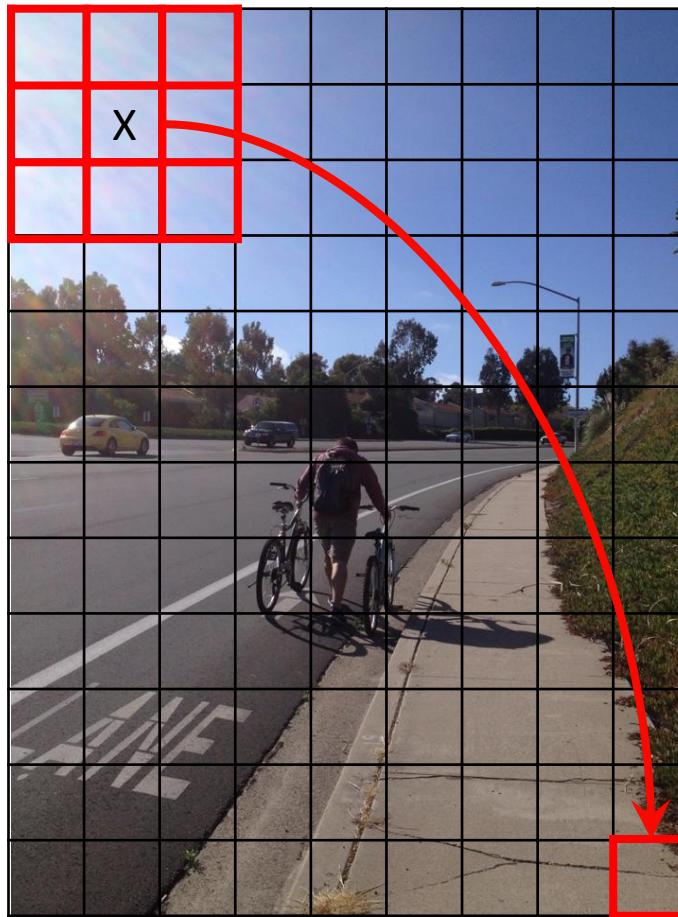
<https://towardsdatascience.com/is-this-the-end-for-convolutional-neural-networks-6f944dccc2e9>



# CNN and Self-Attention

CNN

Receptive field



Self - Attention

