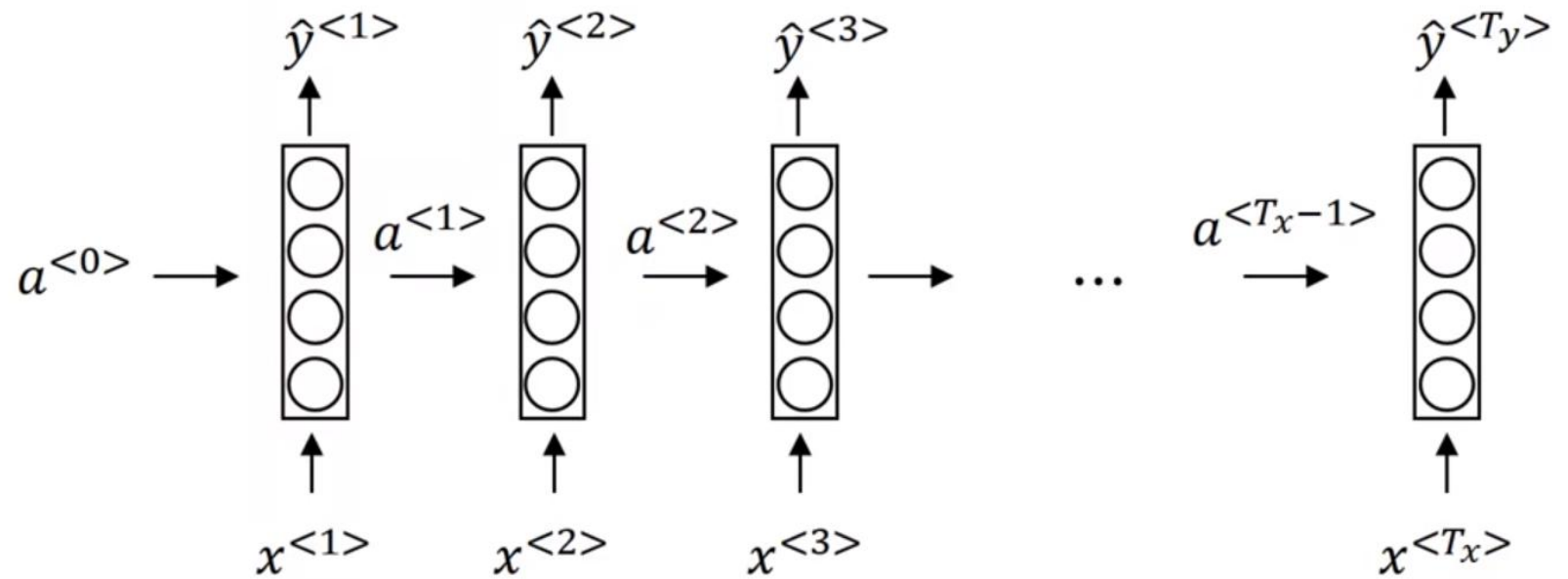# Attention

# RNN
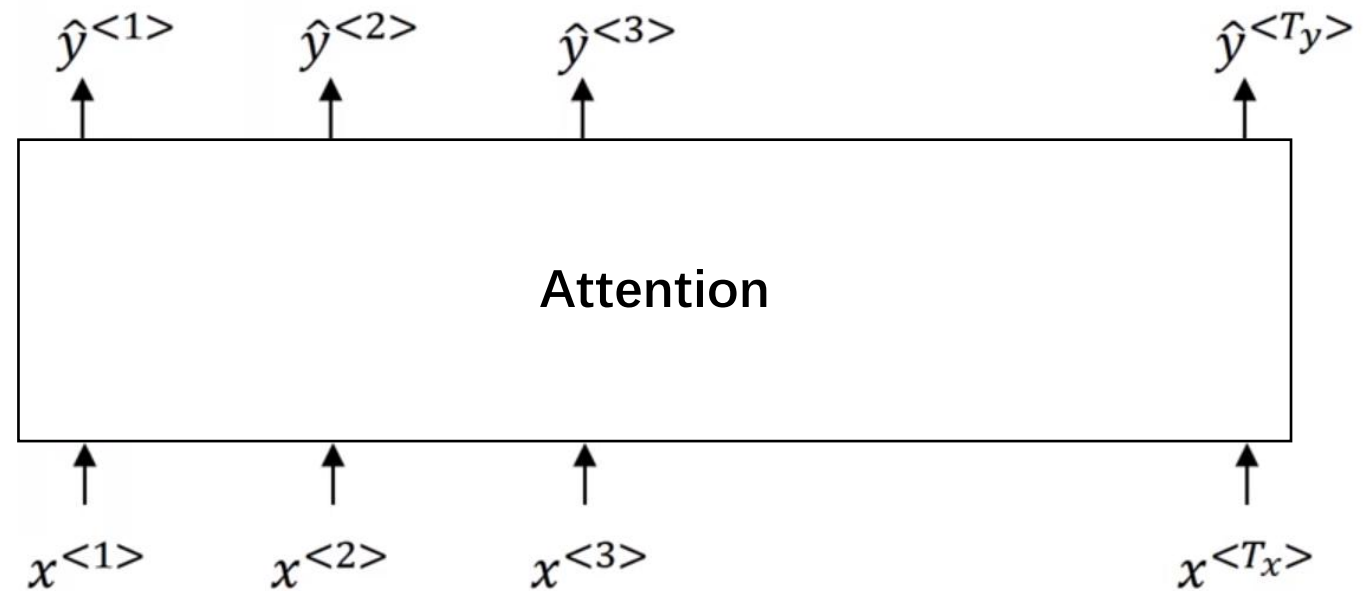
$$a^{<t>} = g(W_{aa}a^{<t-1>} + W_{ax}x^{<t>} + b_a)$$

$$\hat{y}^{<t>} = g(W_y a^{<t>} + b_y)$$
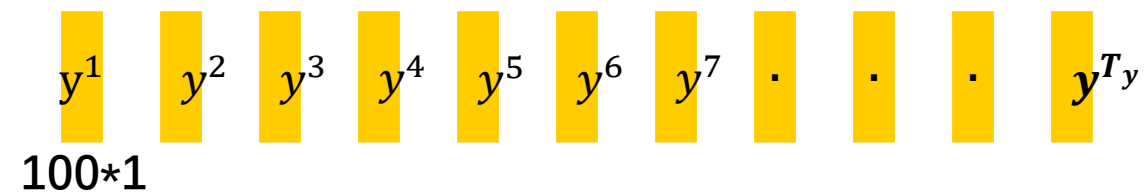
# Attention

$$Q = W_q X$$
$$K = W_k X$$
$$V = W_v X$$
$$A = K^T Q$$
$$A' = softmax(A)$$
$$Y = VA'$$

# Attention
## 向量矩阵化

$$x^1 \quad x^2 \quad x^3 \quad x^4 \quad x^5 \quad x^6 \quad x^7 \quad . \quad . \quad . \quad x^{T_x}$$

200*1

$$y^1 \quad y^2 \quad y^3 \quad y^4 \quad y^5 \quad y^6 \quad y^7 \quad . \quad . \quad . \quad y^{T_y}$$

100*1

# Attention
向量矩阵化

$X$

$$x^1 \quad x^2 \quad x^3 \quad x^4 \quad x^5 \quad x^6 \quad x^7 \quad . \quad . \quad . \quad x^{T_x}$$

200*1

200*800

$Y$

$$y^1 \quad y^2 \quad y^3 \quad y^4 \quad y^5 \quad y^6 \quad y^7 \quad . \quad . \quad . \quad y^{T_y}$$

100*1

100*800

# Attention

$$Q = W_q X$$
$$K = W_k X$$
$$V = W_v X$$
$$A = K^T Q$$
$$A' = softmax(A)$$
$$Y = V A'$$

# Attention

$x^1$

$x^2$

$x^3$

$x^4$

# Attention

$$q^i = W_q x^i$$
$$k^i = W_k x^i$$
$$v^i = W_v x^i$$

# Attention

$$\alpha_{i,j} = q^i \cdot k^j$$

# Attention

$$\alpha'_{1,i} = \frac{e^{\alpha_{1,i}}}{\sum_j e^{\alpha_{1,j}}} \qquad \alpha'_{1,2} = \frac{e^{\alpha_{1,2}}}{e^{\alpha_{1,1}} + e^{\alpha_{1,2}} + e^{\alpha_{1,3}} + e^{\alpha_{1,4}}}$$
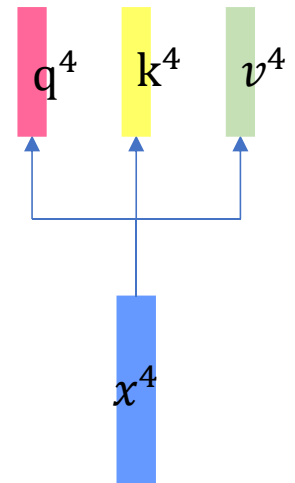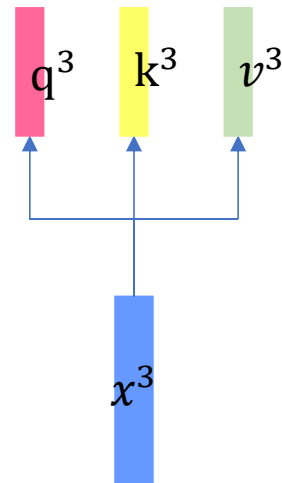
# Attention

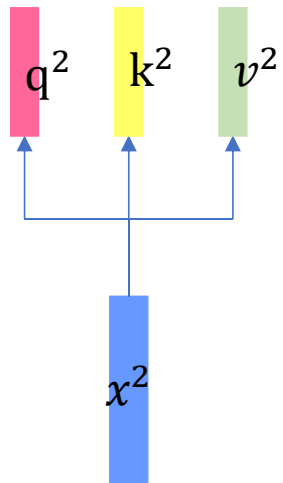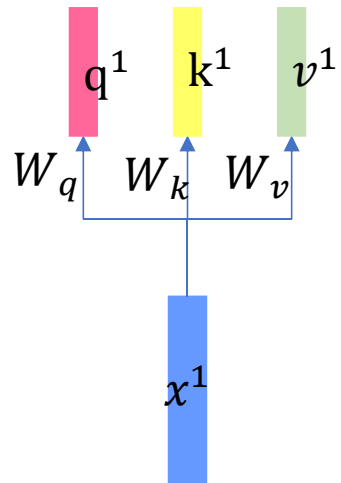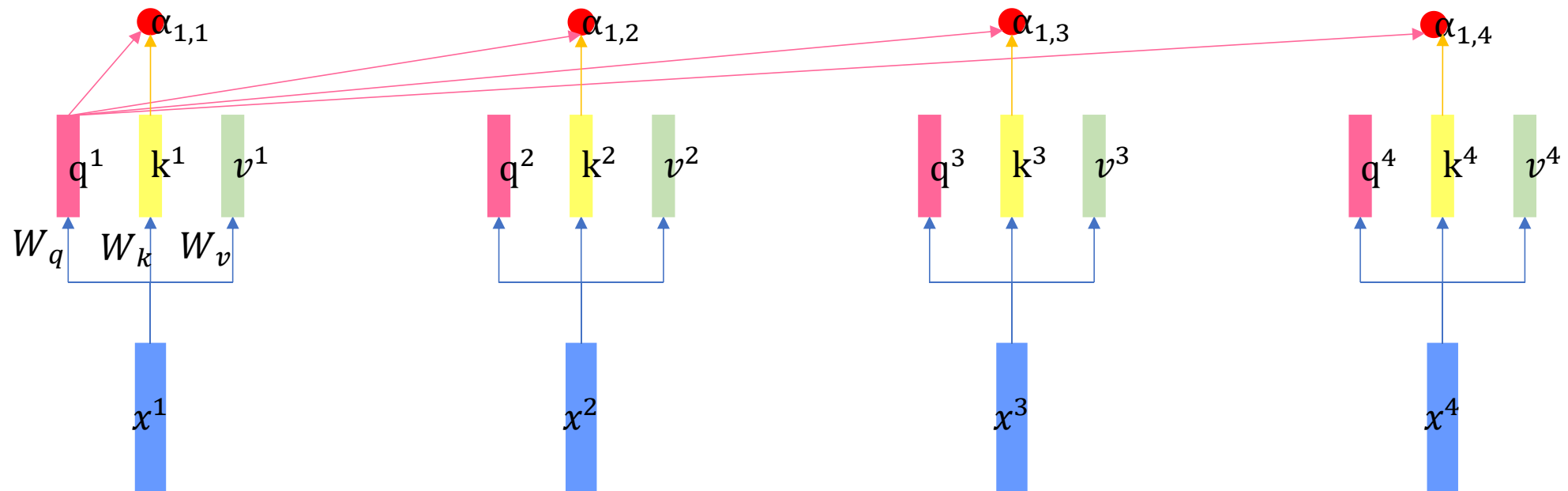$$y^1 = \sum_i \alpha'_{1,i} * v^i$$

# Attention

$x^1$

$x^2$

$x^3$

$x^4$

# Attention

$$q^i = W_q x^i$$
$$k^i = W_k x^i$$
$$v^i = W_v x^i$$

$$\begin{array}{cc} q^1 & q^2 & q^3 & q^4 \end{array} \quad = \quad W_q \quad \begin{array}{cccc} x^1 & x^2 & x^3 & x^4 \end{array}$$

**2*1**                     **2*3**     **3*1**

$q^1$   $k^1$   $v^1$          $q^2$   $k^2$   $v^2$          $q^3$   $k^3$   $v^3$          $q^4$   $k^4$   $v^4$

$W_q$   $W_k$   $W_v$

$x^1$                  $x^2$                  $x^3$                  $x^4$

# Attention

$$q^i = W_q x^i$$
$$k^i = W_k x^i$$
$$v^i = W_v x^i$$

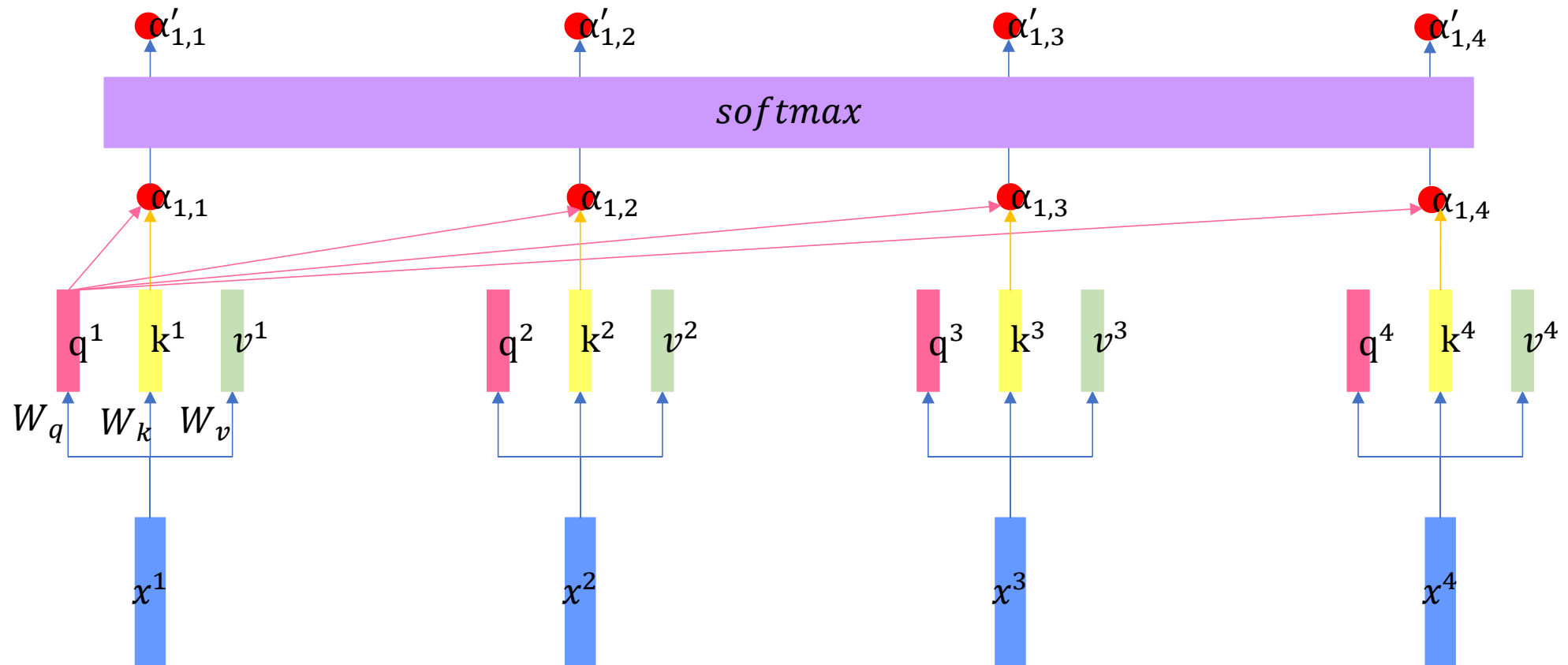$Q$ $^4$ $=$ $W_q$ $X$ $^4$ $\qquad Q = W_q X$

**2\*4** **2\*3** **3\*4**

$$\begin{pmatrix} 1 \\ 2 \end{pmatrix} \begin{pmatrix} 3 \\ 1 \end{pmatrix} \begin{pmatrix} 3 \\ 0 \end{pmatrix} \begin{pmatrix} 2 \\ 6 \end{pmatrix} \quad = \quad \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \\ 2 \end{pmatrix} \begin{pmatrix} 3 \\ 2 \\ 1 \end{pmatrix} \begin{pmatrix} 3 \\ 0 \\ 0 \end{pmatrix} \begin{pmatrix} 2 \\ 9 \\ 6 \end{pmatrix}$$

q$^1$   k$^1$   $v^1$      q$^2$   k$^2$   $v^2$      q$^3$   k$^3$   $v^3$      q$^4$   k$^4$   $v^4$

$W_q$   $W_k$   $W_v$

$x^1$          $x^2$          $x^3$          $x^4$

# Attention

$$q^i = W_q x^i$$
$$k^i = W_k x^i$$
$$v^i = W_v x^i$$

$Q$ ₄ = $W_q$ $X$ ₄

$K$ ₄ = $W_k$ $X$ ₄

$V$ ₄ = $W_v$ $X$ ₄

$$Q = W_q X$$
$$K = W_k X$$
$$V = W_v X$$

q¹  k¹  $v^1$ 

q²  k²  $v^2$ 

$W_q$  $W_k$  $W_v$ 

q³  k³  $v^3$ 

q⁴  k⁴  $v^4$ 

$x^1$

$x^2$

$x^3$

$x^4$

# Attention

$$\alpha_{2,1} = k^1 \ q^2$$

$$\alpha_{2,2} = k^2 \ q^2$$

$$\alpha_{2,3} = k^3 \ q^2$$

$$\alpha_{2,4} = k^4 \ q^2$$

**1*1**     **1*2**  **2*1**

# Attention

$$\alpha_{2,1}$$
$$\alpha_{2,2}$$
$$\alpha_{2,3}$$
$$\alpha_{2,4}$$

$= \begin{matrix} k^1 \\ k^2 \\ k^3 \\ k^4 \end{matrix} \quad q^2$

**4*1**  **4*2**  **2*1**



$\alpha_{2,1}$  $\alpha_{2,2}$  $\alpha_{2,3}$  $\alpha_{2,4}$

$q^1$ k$^1$ $v^1$  $q^2$ k$^2$ $v^2$  $q^3$ k$^3$ $v^3$  $q^4$ k$^4$ $v^4$

$W_q$  $W_k$  $W_v$

$x^1$  $x^2$  $x^3$  $x^4$

# Attention

$\alpha_{2,1}$

$\alpha_{2,2}$

$\alpha_{2,3}$

$\alpha_{2,4}$

**4\*1**

$=$

$k^1$
$k^2$
$k^3$
$k^4$

**4\*2**

$K^T$

$q^2$

**2\*1**



$\alpha_{2,1}$   $\alpha_{2,2}$   $\alpha_{2,3}$   $\alpha_{2,4}$

$q^1$  $k^1$  $v^1$      $q^2$  $k^2$  $v^2$      $q^3$  $k^3$  $v^3$      $q^4$  $k^4$  $v^4$

$W_q$   $W_k$   $W_v$

$x^1$        $x^2$        $x^3$        $x^4$

# Attention



$\alpha_{2,1}$ $\alpha_{3,1}$
$\alpha_{2,2}$ $\alpha_{3,2}$
$\alpha_{2,3}$ $\alpha_{3,3}$
$\alpha_{2,4}$ $\alpha_{3,4}$

**4*2**

=

$k^1$
$k^2$
$k^3$
$k^4$

**4*2**

$K^T$

$q^2$ $q^3$

**2*2**

$\alpha_{2,1}$    $\alpha_{2,2}$    $\alpha_{2,3}$    $\alpha_{2,4}$

$q^1$ $k^1$ $v^1$    $q^2$ $k^2$ $v^2$    $q^3$ $k^3$ $v^3$    $q^4$ $k^4$ $v^4$

$W_q$  $W_k$  $W_v$

$x^1$    $x^2$    $x^3$    $x^4$

# Attention



$$\begin{matrix} \alpha_{2,1} & \alpha_{3,1} & \alpha_{4,1} \\ \alpha_{2,2} & \alpha_{3,2} & \alpha_{4,2} \\ \alpha_{2,3} & \alpha_{3,3} & \alpha_{4,3} \\ \alpha_{2,4} & \alpha_{3,4} & \alpha_{4,4} \end{matrix} \quad = \quad \begin{matrix} k^1 \\ k^2 \\ k^3 \\ k^4 \end{matrix} \qquad q^2 \ q^3 \ q^4$$

**4*3**    **4*2**    **2*3**

$K^T$

$\alpha_{2,1}$    $\alpha_{2,2}$    $\alpha_{2,3}$    $\alpha_{2,4}$

$q^1$ $k^1$ $v^1$    $q^2$ $k^2$ $v^2$    $q^3$ $k^3$ $v^3$    $q^4$ $k^4$ $v^4$

$W_q$ $W_k$ $W_v$

$x^1$    $x^2$    $x^3$    $x^4$

# Attention

$$\begin{matrix} \alpha_{1,1} & \alpha_{2,1} & \alpha_{3,1} & \alpha_{4,1} \\ \alpha_{1,2} & \alpha_{2,2} & \alpha_{3,2} & \alpha_{4,2} \\ \alpha_{1,3} & \alpha_{2,3} & \alpha_{3,3} & \alpha_{4,3} \\ \alpha_{1,4} & \alpha_{2,4} & \alpha_{3,4} & \alpha_{4,4} \end{matrix}$$

**4*4**

$=$

$k^1 \quad k^2 \quad k^3 \quad k^4$

**4*2**

$K^T$

$q^1 \quad q^2 \quad q^3 \quad q^4$

**2*4**

$\alpha_{2,1}$  $\alpha_{2,2}$  $\alpha_{2,3}$  $\alpha_{2,4}$

$q^1$  $k^1$  $v^1$

$q^2$  $k^2$  $v^2$

$q^3$  $k^3$  $v^3$

$q^4$  $k^4$  $v^4$

$W_q \quad W_k \quad W_v$

$x^1$

$x^2$

$x^3$

$x^4$

# Attention

# Attention

$A$ $=$ K$^T$ $Q$                    $A = K^T Q$

**4*1**          **4*2**     **2*1**

# Attention

# Attention

# Attention

$$y^i = \sum_j \alpha'_{i,j} * v^j$$

# Attention

$$y^i = \sum_j \alpha'_{i,j} * v^j$$

# Attention

$$y^i = \sum_j \alpha'_{i,j} * v^j$$

$y^2$

$\alpha'_{2,1}$ * $\alpha'_{2,2}$ * $\alpha'_{2,3}$ * $\alpha'_{2,4}$ *

$softmax$

$\alpha_{2,1}$ $\alpha_{2,2}$ $\alpha_{2,3}$ $\alpha_{2,4}$

$$y^2 = \alpha'_{2,1}v^1 + \alpha'_{2,2}v^2 + \alpha'_{2,3}v^3 + \alpha'_{2,4}v^4$$

$$\begin{pmatrix} 1 \\ 2 \end{pmatrix} = 0.1 \begin{pmatrix} 1 \\ 0 \end{pmatrix} + 0.2 \begin{pmatrix} 0 \\ 2 \end{pmatrix} + 0.3 \begin{pmatrix} 3 \\ 0 \end{pmatrix} + 0.4 \begin{pmatrix} 0 \\ 4 \end{pmatrix}$$

$$\begin{pmatrix} 1 \\ 2 \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \end{pmatrix} \begin{pmatrix} 0 \\ 2 \end{pmatrix} \begin{pmatrix} 3 \\ 0 \end{pmatrix} \begin{pmatrix} 0 \\ 4 \end{pmatrix} \begin{pmatrix} 0.1 \\ 0.2 \\ 0.3 \\ 0.4 \end{pmatrix}$$

# Attention

$$y^i = \sum_j \alpha'_{i,j} * v^j$$



$$y^2 = \alpha'_{2,1}v^1 + \alpha'_{2,2}v^2 + \alpha'_{2,3}v^3 + \alpha'_{2,4}v^4$$

$$\begin{pmatrix} 1 \\ 2 \end{pmatrix} = 0.1 \begin{pmatrix} 1 \\ 0 \end{pmatrix} + 0.2 \begin{pmatrix} 0 \\ 2 \end{pmatrix} + 0.3 \begin{pmatrix} 3 \\ 0 \end{pmatrix} + 0.4 \begin{pmatrix} 0 \\ 4 \end{pmatrix}$$

$$\begin{pmatrix} 1 \\ 2 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 3 & 0 \\ 0 & 2 & 0 & 4 \end{pmatrix} \begin{pmatrix} 0.1 \\ 0.2 \\ 0.3 \\ 0.4 \end{pmatrix} \qquad y^2 = \begin{pmatrix} v^1 & v^2 & v^3 & v^4 \end{pmatrix} \begin{pmatrix} \alpha'_{2,1} \\ \alpha'_{2,2} \\ \alpha'_{2,3} \\ \alpha'_{2,4} \end{pmatrix}$$

# Attention

$$y^i = \sum_j \alpha'_{i,j} * v^j$$

$$\begin{pmatrix} 1 \\ 2 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 3 & 0 \\ 0 & 2 & 0 & 4 \end{pmatrix} \begin{pmatrix} 0.1 \\ 0.2 \\ 0.3 \\ 0.4 \end{pmatrix}$$

$$y^2 = \begin{pmatrix} v^1 & v^2 & v^3 & v^4 \end{pmatrix} \begin{pmatrix} \alpha'_{2,1} \\ \alpha'_{2,2} \\ \alpha'_{2,3} \\ \alpha'_{2,4} \end{pmatrix}$$

$$y^2 \quad = \quad \boxed{v^1} \; \boxed{v^2} \; \boxed{v^3} \; \boxed{v^4} \quad \begin{matrix} \alpha'_{2,1} \\ \alpha'_{2,2} \\ \alpha'_{2,3} \\ \alpha'_{2,4} \end{matrix}$$

# Attention

$$y^i = \sum_j \alpha'_{i,j} * v^j$$

$$\begin{pmatrix} 1 \\ 2 \end{pmatrix} \begin{pmatrix} 0.9 \\ 0 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 3 & 0 \\ 0 & 2 & 0 & 4 \end{pmatrix} \begin{pmatrix} 0.1 \\ 0.2 \\ 0.3 \\ 0.4 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ 0.3 \\ 0 \end{pmatrix}$$

$$y^2 y^3 = \begin{pmatrix} v^1 & v^2 & v^3 & v^4 \end{pmatrix} \begin{pmatrix} \alpha'_{2,1} \\ \alpha'_{2,2} \\ \alpha'_{2,3} \\ \alpha'_{2,4} \end{pmatrix} \begin{pmatrix} \alpha'_{3,1} \\ \alpha'_{3,2} \\ \alpha'_{3,3} \\ \alpha'_{3,4} \end{pmatrix}$$

$$y^2 \quad y^3 \qquad = \qquad v^1 \quad v^2 \quad v^3 \quad v^4 \qquad \begin{matrix} \alpha'_{2,1} & \alpha'_{3,1} \\ \alpha'_{2,2} & \alpha'_{3,2} \\ \alpha'_{2,3} & \alpha'_{3,3} \\ \alpha'_{2,4} & \alpha'_{3,4} \end{matrix}$$

# Attention

$$y^i = \sum_j \alpha'_{i,j} * v^j$$

$$\begin{pmatrix} 1 \\ 2 \end{pmatrix} \begin{pmatrix} 0.9 \\ 0 \end{pmatrix} \begin{pmatrix} \textcolor{red}{0} \\ \textcolor{red}{1.6} \end{pmatrix} = \begin{pmatrix} 1 & 0 & 3 & 0 \\ 0 & 2 & 0 & 4 \end{pmatrix} \begin{pmatrix} 0.1 \\ 0.2 \\ 0.3 \\ 0.4 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ 0.3 \\ 0 \end{pmatrix} \begin{pmatrix} \textcolor{red}{0} \\ \textcolor{red}{0} \\ \textcolor{red}{0} \\ \textcolor{red}{0.4} \end{pmatrix}$$

$$y^2 y^3 \textcolor{red}{y^4} = \begin{pmatrix} v^1 & v^2 & v^3 & v^4 \end{pmatrix} \begin{pmatrix} \alpha'_{2,1} \\ \alpha'_{2,2} \\ \alpha'_{2,3} \\ \alpha'_{2,4} \end{pmatrix} \begin{pmatrix} \alpha'_{3,1} \\ \alpha'_{3,2} \\ \alpha'_{3,3} \\ \alpha'_{3,4} \end{pmatrix} \begin{pmatrix} \textcolor{red}{\alpha'_{4,1}} \\ \textcolor{red}{\alpha'_{4,2}} \\ \textcolor{red}{\alpha'_{4,3}} \\ \textcolor{red}{\alpha'_{4,4}} \end{pmatrix}$$

$y^2$  $y^3$  $y^4$   $=$   $v^1$  $v^2$  $v^3$  $v^4$

$$\begin{array}{ccc} \alpha'_{2,1} & \alpha'_{3,1} & \alpha'_{4,1} \\ \alpha'_{2,2} & \alpha'_{3,2} & \alpha'_{4,2} \\ \alpha'_{2,3} & \alpha'_{3,3} & \alpha'_{4,3} \\ \alpha'_{2,4} & \alpha'_{3,4} & \alpha'_{4,4} \end{array}$$

# Attention

$$y^i = \sum_j$$

$$\begin{pmatrix} \color{red}{0.1} \\ \color{red}{0} \end{pmatrix} \begin{pmatrix} 1 \\ 2 \end{pmatrix} \begin{pmatrix} 0.9 \\ 0 \end{pmatrix} \begin{pmatrix} 0 \\ 1.6 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 3 & 0 \\ 0 & 2 & 0 & 4 \end{pmatrix} \begin{pmatrix} \color{red}{0.1} \\ \color{red}{0} \\ \color{red}{0} \\ \color{red}{0} \end{pmatrix} \begin{pmatrix} 0.1 \\ 0.2 \\ 0.3 \\ 0.4 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ 0.3 \\ 0 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0.4 \end{pmatrix}$$

$$\color{red}{y^1}\, y^2 y^3 y^4 = \begin{pmatrix} v^1 & v^2 & v^3 & v^4 \end{pmatrix} \begin{pmatrix} \color{red}{\alpha'_{1,1}} \\ \color{red}{\alpha'_{1,2}} \\ \color{red}{\alpha'_{1,3}} \\ \color{red}{\alpha'_{1,4}} \end{pmatrix} \begin{pmatrix} \alpha'_{2,1} \\ \alpha'_{2,2} \\ \alpha'_{2,3} \\ \alpha'_{2,4} \end{pmatrix} \begin{pmatrix} \alpha'_{3,1} \\ \alpha'_{3,2} \\ \alpha'_{3,3} \\ \alpha'_{3,4} \end{pmatrix} \begin{pmatrix} \alpha'_{4,1} \\ \alpha'_{4,2} \\ \alpha'_{4,3} \\ \alpha'_{4,4} \end{pmatrix}$$

$$y^1 \quad y^2 \quad y^3 \quad y^4 \quad = \quad v^1 \quad v^2 \quad v^3 \quad v^4 \quad \begin{matrix} \alpha'_{1,1} & \alpha'_{2,1} & \alpha'_{3,1} & \alpha'_{4,1} \\ \alpha'_{1,2} & \alpha'_{2,2} & \alpha'_{3,2} & \alpha'_{4,2} \\ \alpha'_{1,3} & \alpha'_{2,3} & \alpha'_{3,3} & \alpha'_{4,3} \\ \alpha'_{1,4} & \alpha'_{2,4} & \alpha'_{3,4} & \alpha'_{4,4} \end{matrix}$$

# Attention

$$y^i = \sum_j$$

$$\begin{pmatrix} 0.1 & 1 & 0.9 & 0 \\ 0 & 2 & 0 & 1.6 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 3 & 0 \\ 0 & 2 & 0 & 4 \end{pmatrix} \begin{pmatrix} 0.1 & 0.1 & 0 & 0 \\ 0 & 0.2 & 0 & 0 \\ 0 & 0.3 & 0.3 & 0 \\ 0 & 0.4 & 0 & 0.4 \end{pmatrix}$$

$$y^1 y^2 y^3 y^4 = (v^1 \quad v^2 \quad v^3 \quad v^4) \begin{pmatrix} \alpha'_{1,1} & \alpha'_{2,1} & \alpha'_{3,1} & \alpha'_{4,1} \\ \alpha'_{1,2} & \alpha'_{2,2} & \alpha'_{3,2} & \alpha'_{4,2} \\ \alpha'_{1,3} & \alpha'_{2,3} & \alpha'_{3,3} & \alpha'_{4,3} \\ \alpha'_{1,4} & \alpha'_{2,4} & \alpha'_{3,4} & \alpha'_{4,4} \end{pmatrix}$$

$$\boldsymbol{Y = VA'}$$

$Y$ $^4$ $=$ $V$ $^{,4}$ $A'$

# Attention

$$Q = \begin{bmatrix} & & Q & & \end{bmatrix} = \begin{bmatrix} W_q \end{bmatrix} \begin{bmatrix} & & X & & \end{bmatrix}$$

$$K = \begin{bmatrix} & & K & & \end{bmatrix} = \begin{bmatrix} W_k \end{bmatrix} \begin{bmatrix} & & X & & \end{bmatrix}$$

$$V = \begin{bmatrix} & & V & & \end{bmatrix} = \begin{bmatrix} W_v \end{bmatrix} \begin{bmatrix} & & X & & \end{bmatrix}$$

$$\boldsymbol{Q = W_q X}$$
$$\boldsymbol{K = W_k X}$$
$$\boldsymbol{V = W_v X}$$

100*800                    100*200          200*800

# Attention



$A$     =     $K^T$     $Q$

800*800        800*100       100*800

$$A = K^T Q$$
$$A' = softmax(A)$$

# Attention

$$A' = softmax(\qquad A \qquad)$$

800*800                      800*800

$$A' = softmax(A)$$

# Self-Attention

$$Y = VA'$$

with matrices labeled:
- $Y$ — 100*800
- $V$ — 100*800
- $A'$ — 800*800

$$Y = VA'$$

天涯若比邻，Attention is all you need

**位置编码**

天涯若比邻，Attention is all you need

$y^2$

$\alpha'_{2,1}$ * $\alpha'_{2,2}$ * $\alpha'_{2,3}$ * $\alpha'_{2,4}$ *

$softmax$

$\alpha_{2,1}$ $\alpha_{2,2}$ $\alpha_{2,3}$ $\alpha_{2,4}$

$q^1$ $k^1$ $v^1$ $q^4$ $k^4$ $v^4$

$$PE(pos, 2i) = \sin(pos/10000^{2i/d_{model}})$$

$$PE(pos, 2i + 1) = \cos(pos/10000^{2i/d_{model}})$$

$W_q$ $W_k$ $W_v$

$e^1$ + $x^1$ $e^2$ + $x^2$ $e^3$ + $x^3$ $e^4$ + $x^4$

$y^2$

天涯若比邻，**Attention is all you need**

$\alpha'_{2,1}$ * $\alpha'_{2,2}$ * $\alpha'_{2,3}$ * $\alpha'_{2,4}$ *

$softmax$

$\alpha_{2,1}$ $\alpha_{2,2}$ $\alpha_{2,3}$ $\alpha_{2,4}$

$q^1$ $k^1$ $v^1$ $q^2$ $q^4$ $k^4$ $v^4$

$$PE(pos, 2i) = \sin(pos/10000^{2i/d_{model}})$$

$W_q$ $W_k$ $W_v$

$$PE(pos, 2i + 1) = \cos(pos/10000^{2i/d_{model}})$$

$pos = 2$
$d_{model} = 4$

$e^1$ + $x^1$ $e^2$ + $x^2$ $e^3$ + $x^3$ $e^4$ + $x^4$

天涯若比邻，Attention is all you need

$$PE(2,2i) = \sin(2/10000^{2i/4})$$

$$PE(2,2i+1) = \cos(2/10000^{2i/4})$$

$$\begin{pmatrix} ? \\ ? \\ ? \\ ? \end{pmatrix}$$

$softmax$

$y^2$

$\alpha'_{2,1}$  $*$   $\alpha'_{2,2}$  $*$   $\alpha'_{2,3}$  $*$   $\alpha'_{2,4}$  $*$

$\alpha_{2,1}$   $\alpha_{2,2}$   $\alpha_{2,3}$   $\alpha_{2,4}$

$q^1$  $k^1$  $v^1$   $q^2$  $k^2$  $v^2$   $q^3$  $k^3$  $v^3$   $q^4$  $k^4$  $v^4$

$W_q$  $W_k$  $W_v$

$e^1$ $+$ $x^1$        $e^2$ $+$ $x^2$   $pos = 2$   $e^3$ $+$ $x^3$        $e^4$ $+$ $x^4$
$d_{model} = 4$

**位置编码**

天涯若比邻，**Attention is all you need**

$y^2$

$\alpha'_{2,1}$ * $\alpha'_{2,2}$ * $\alpha'_{2,3}$ * $\alpha'_{2,4}$ *

$softmax$

$\alpha_{2,1}$ $\alpha_{2,2}$ $\alpha_{2,3}$ $\alpha_{2,4}$

$$PE(2,2i) = \sin(2/10000^{2i/4})$$

$$\left( cos(2/10000^{2*1/4} \right.$$

$$PE(2,2i+1) = \cos(2/10000^{2i/4})$$

$$\left. \right)$$

$q^1$ $k^1$ $v^1$ $q^2$ $k^2$ $v^2$ $q^3$ $k^3$ $v^3$ $q^4$ $k^4$ $v^4$

$W_q$ $W_k$ $W_v$

$e^1$ + $x^1$ $e^2$ + $x^2$ $e^3$ + $x^3$ $e^4$ + $x^4$

$$pos = 2$$
$$d_{model} = 4$$

**位置编码**

$y^2$

天涯若比邻，Attention is all you need

$\alpha'_{2,1}$ *     $\alpha'_{2,2}$ *     $\alpha'_{2,3}$ *     $\alpha'_{2,4}$ *

$$softmax$$

$\alpha_{2,1}$     $\alpha_{2,2}$     $\alpha_{2,3}$     $\alpha_{2,4}$

$q^1 \quad k^1 \quad v^1$     $q^2 \quad k^2 \quad v^2$     $q^3 \quad k^3 \quad v^3$

$$PE(2,2i) = \sin(2/10000^{2i/4})$$

$$PE(2,2i+1) = \cos(2/10000^{2i/4})$$

$$\begin{pmatrix} cos(2/10000^{2*1/4}) \\ sin(2/10000^{2*2/4}) \end{pmatrix}$$

$W_q \quad W_k \quad W_v$

$e^1$ + $x^1$     $e^2$ + $x^2$

$$pos = 2$$
$$d_{model} = 4$$

$e^3$ + $x^3$     $e^4$ + $x^4$

**位置编码**

天涯若比邻，Attention is all you need

$y^2$

$\alpha'_{2,1}$ $*$    $\alpha'_{2,2}$ $*$    $\alpha'_{2,3}$ $*$    $\alpha'_{2,4}$ $*$

$softmax$

$\alpha_{2,1}$    $\alpha_{2,2}$    $\alpha_{2,3}$    $\alpha_{2,4}$

$q^1$   $k^1$   $v^1$    $q^2$   $PE(2,2i) = \sin(2/10000^{2i/4})$   $k^3$   $v^3$

$\begin{pmatrix} cos(2/10000^{2*1/4}) \\ sin(2/10000^{2*2/4}) \\ cos(2/10000^{2*3/4}) \end{pmatrix}$

$W_q$   $W_k$ , $W_v$    $PE(2,2i+1) = \cos(2/10000^{2i/4})$

$e^1$ $+$ $x^1$     $e^2$ $+$ $x^2$   $pos = 2$   $e^3$ $+$ $x^3$     $e^4$ $+$ $x^4$

$d_{model} = 4$

**位置编码**

$y^2$

天涯若比邻，**Attention is all you need**

$\alpha'_{2,1}$ *     $\alpha'_{2,2}$ *     $\alpha'_{2,3}$ *     $\alpha'_{2,4}$ *

$softmax$

$\alpha_{2,1}$     $\alpha_{2,2}$     $\alpha_{2,3}$     $\alpha_{2,4}$

$q^1$   $k^1$   $v^1$     $q^2$   $k^2$   $v^2$     $q^3$   $k^3$   $v^3$

$$PE(2,2i) = \sin(2/10000^{2i/4})$$

$$PE(2,2i+1) = \cos(2/10000^{2i/4})$$

$$\begin{pmatrix} cos(2/10000^{2*1/4}) \\ \sin(2/10000^{2*2/4}) \\ cos(2/10000^{2*3/4}) \\ \sin(2/10000^{2*4/4}) \end{pmatrix}$$

$W_q$   $W_k$   $W_v$

$e^1$ + $x^1$     $e^2$ + $x^2$    $\begin{matrix} pos = 2 \\ d_{model} = 4 \end{matrix}$   $e^3$ + $x^3$     $e^4$ + $x^4$

**位置编码**



天涯若比邻，Attention is all you need

$$\begin{pmatrix} cos(2/10000^{2*1/4} \\ sin(2/10000^{2*2/4}) \\ cos(2/10000^{2*3/4}) \\ sin(2/10000^{2*4/4}) \end{pmatrix} = \begin{pmatrix} 0.99998333341666 \\ -0.02463979854800 \\ 0.99999166666944 \\ 0.99970691178397 \end{pmatrix}$$

$pos = 2$
$d_{model} = 4$

# Attention



$$Q = W_q X$$
(1) $$K = W_k X$$
$$V = W_v X$$

(2) $$A = K^T Q$$

(3) $$A' = softmax(A)$$

(4) $$Y = VA'$$