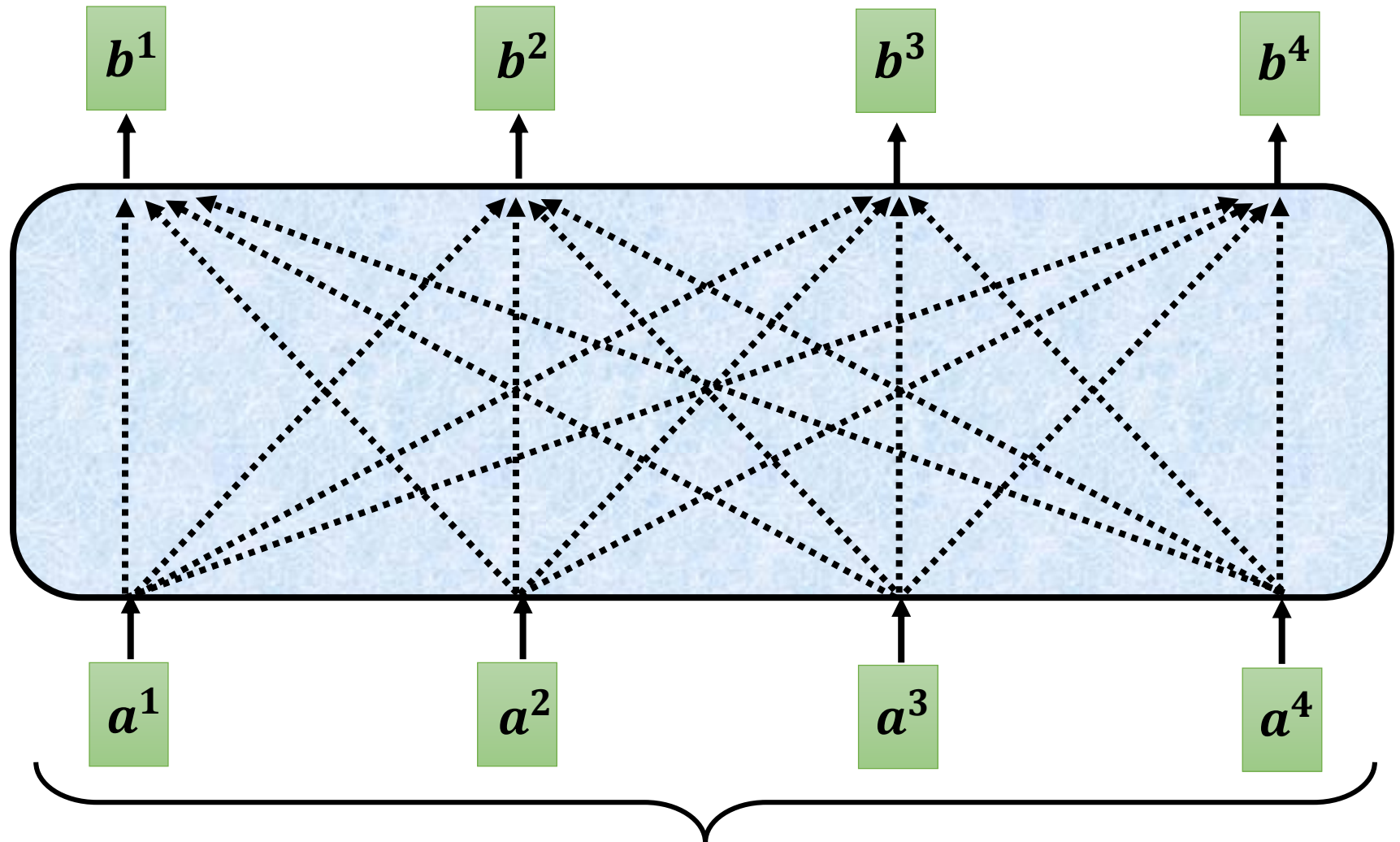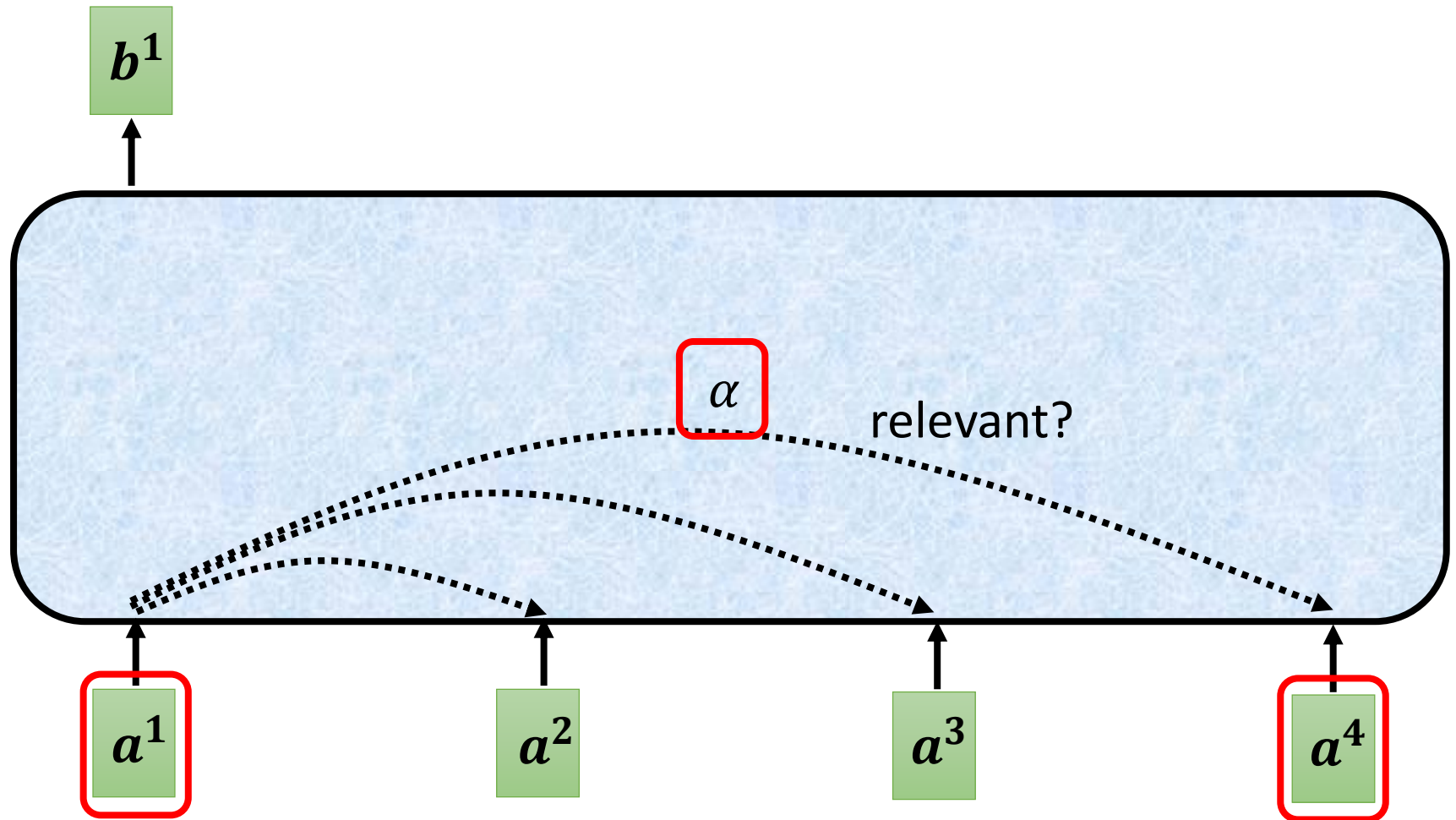# Self-attention



Can be either **input** or **a hidden layer**
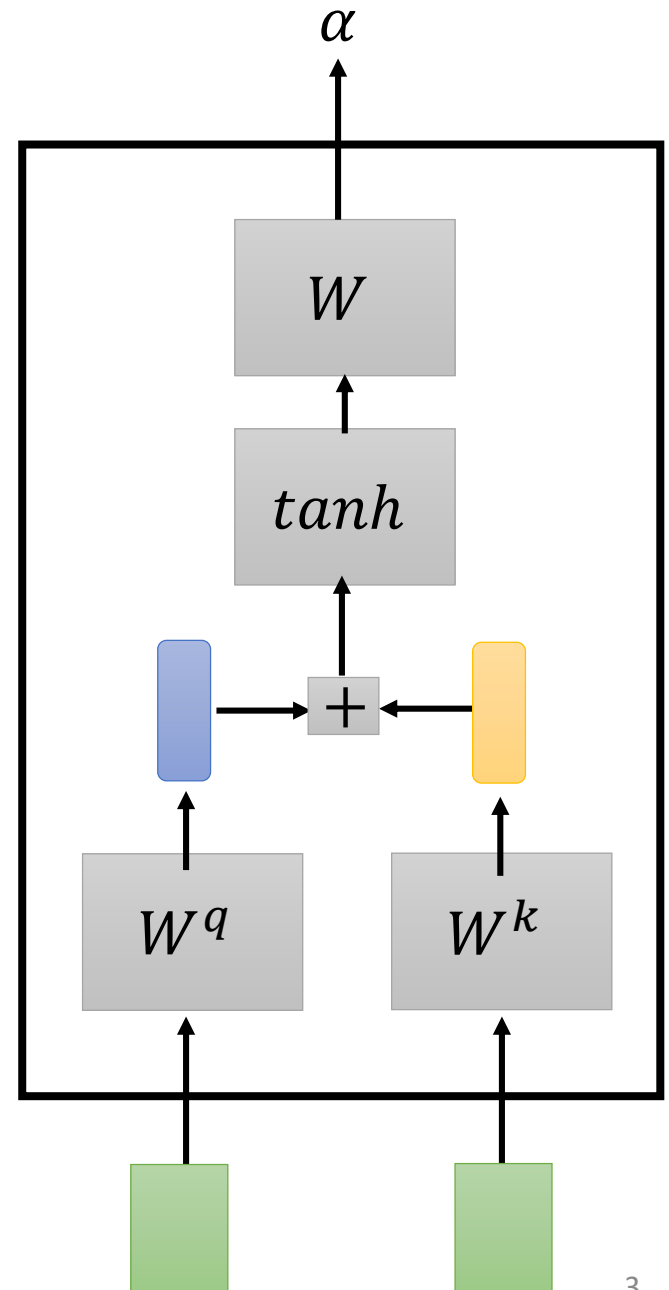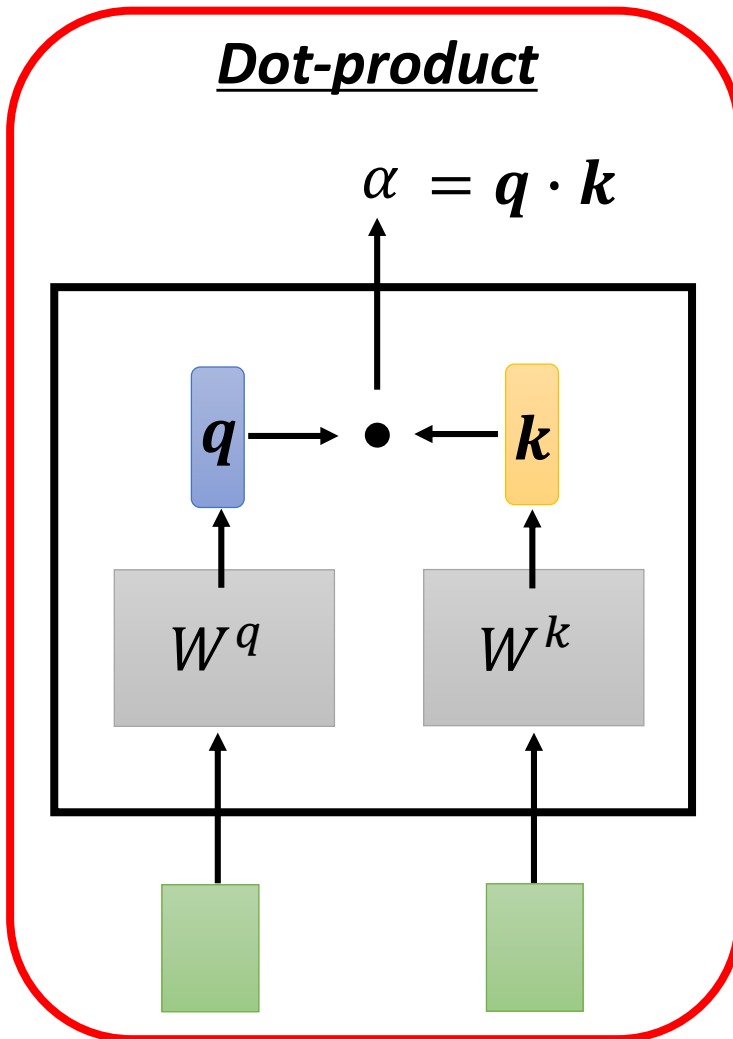
# _Self-attention_
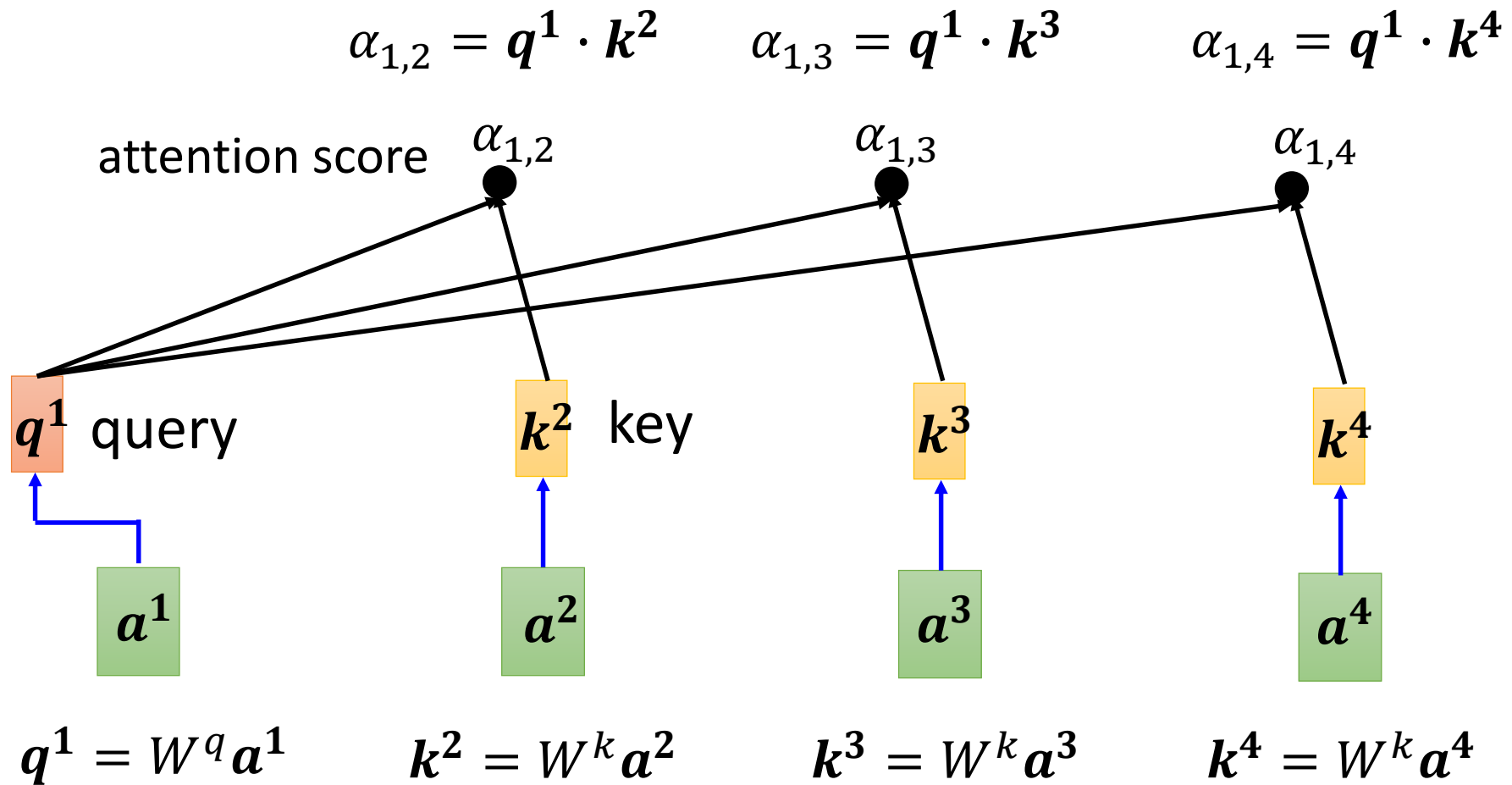


Find the relevant vectors in a sequence
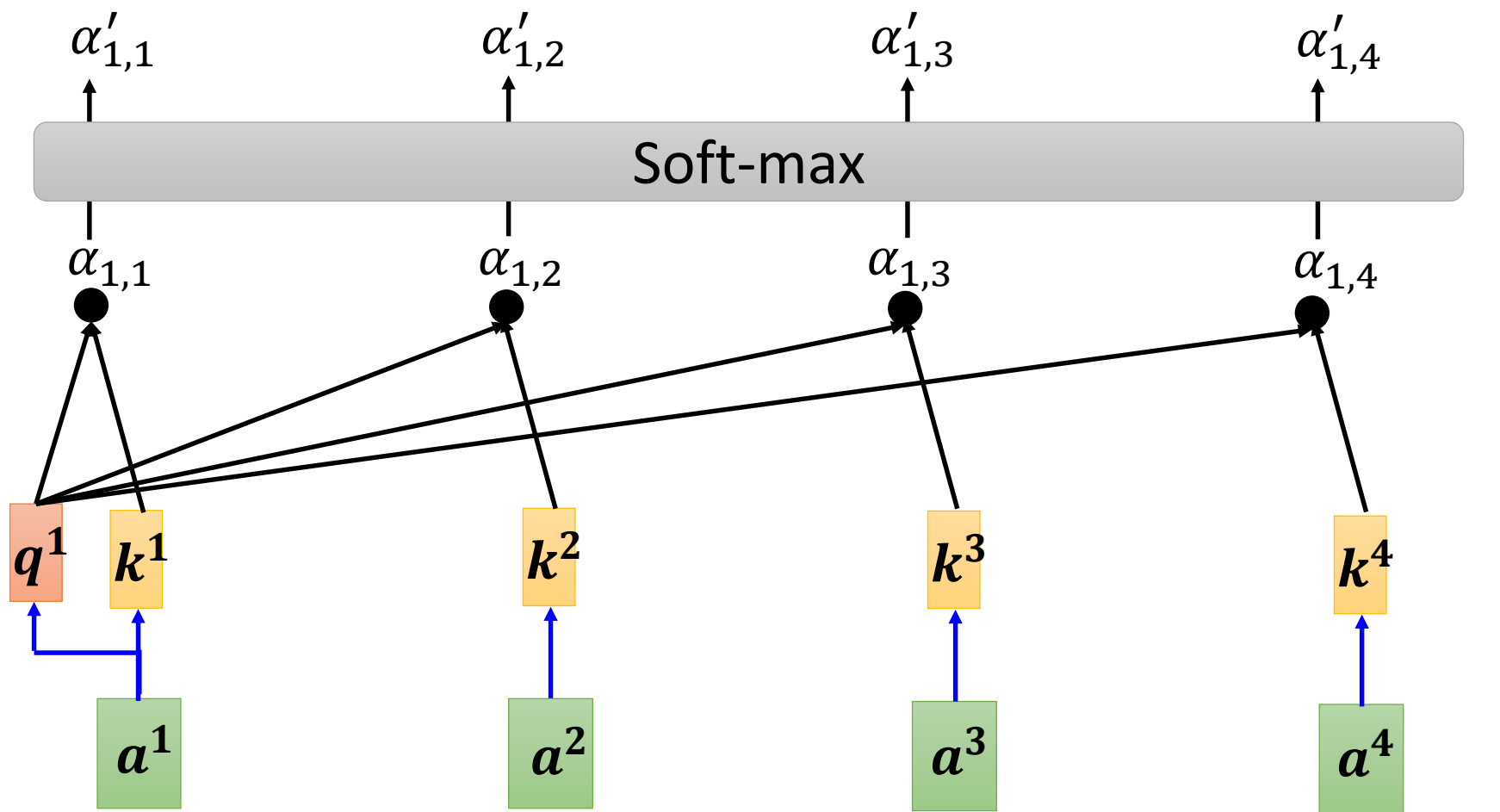
# Self-attention



**Dot-product**

$$\alpha = \boldsymbol{q} \cdot \boldsymbol{k}$$

**Additive**

$$\alpha$$

3

# _Self-attention_

$$\alpha_{1,2} = \boldsymbol{q^1} \cdot \boldsymbol{k^2} \qquad \alpha_{1,3} = \boldsymbol{q^1} \cdot \boldsymbol{k^3} \qquad \alpha_{1,4} = \boldsymbol{q^1} \cdot \boldsymbol{k^4}$$

attention score $\quad\alpha_{1,2}\qquad\qquad\alpha_{1,3}\qquad\qquad\alpha_{1,4}$

$\boldsymbol{q^1}$ query $\qquad \boldsymbol{k^2}$ key $\qquad\qquad \boldsymbol{k^3}\qquad\qquad\qquad \boldsymbol{k^4}$

$\boldsymbol{a^1} \qquad\qquad \boldsymbol{a^2} \qquad\qquad\qquad \boldsymbol{a^3} \qquad\qquad\qquad \boldsymbol{a^4}$

$$\boldsymbol{q^1} = W^q \boldsymbol{a^1} \qquad \boldsymbol{k^2} = W^k \boldsymbol{a^2} \qquad \boldsymbol{k^3} = W^k \boldsymbol{a^3} \qquad \boldsymbol{k^4} = W^k \boldsymbol{a^4}$$

## Self-attention

$$\alpha'_{1,i} = exp(\alpha_{1,i})/\sum_j exp(\alpha_{1,j})$$

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

Soft-max

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

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

$a^1$ $a^2$ $a^3$ $a^4$

$q^1 = W^q a^1$ $\quad$ $k^2 = W^k a^2$ $\quad$ $k^3 = W^k a^3$ $\quad$ $k^4 = W^k a^4$
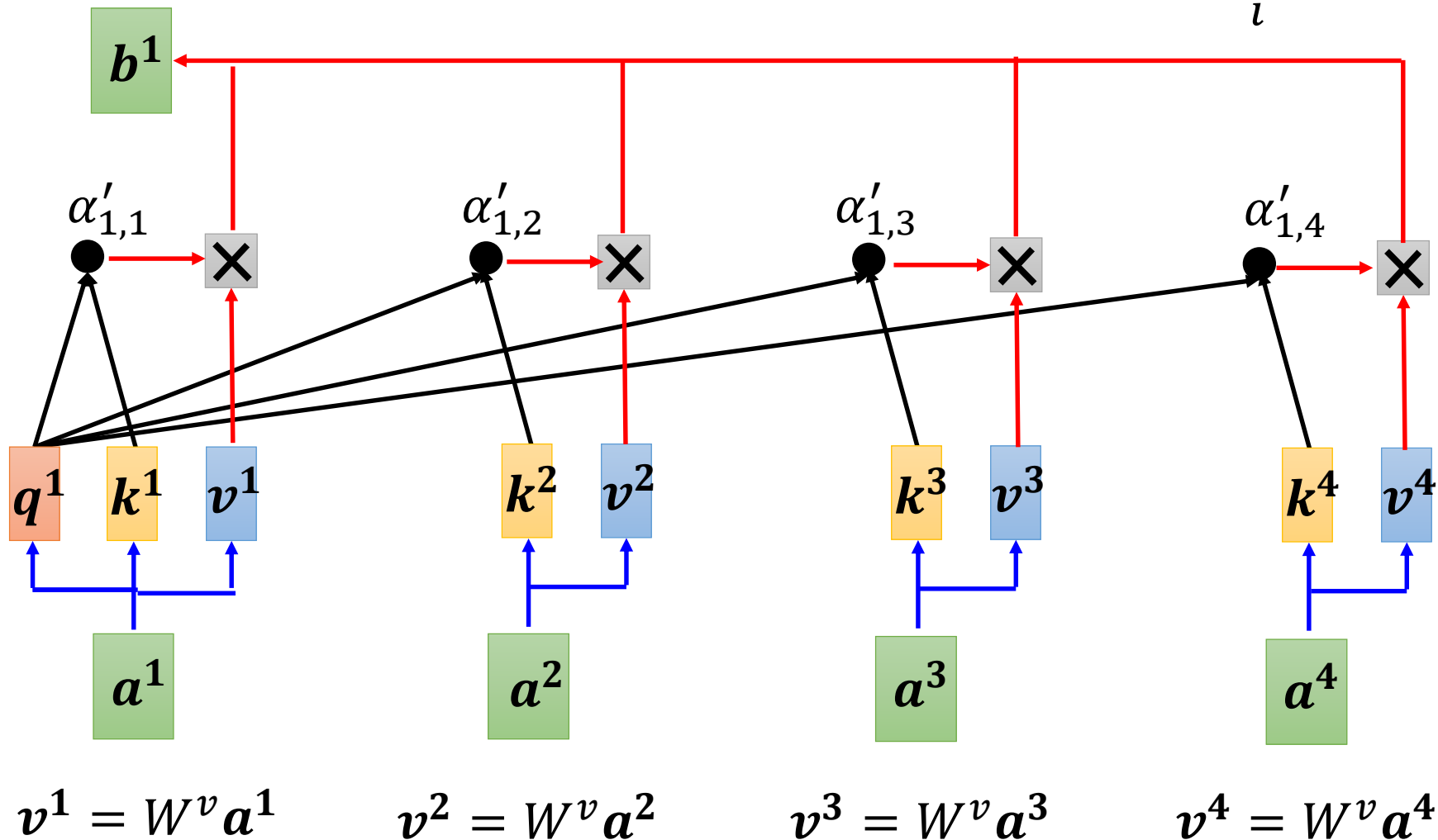
$k^1 = W^k a^1$

5

# *Self-attention*

Extract information based on attention scores

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



$$v^1 = W^v a^1 \qquad v^2 = W^v a^2 \qquad v^3 = W^v a^3 \qquad v^4 = W^v a^4$$

# Self-attention

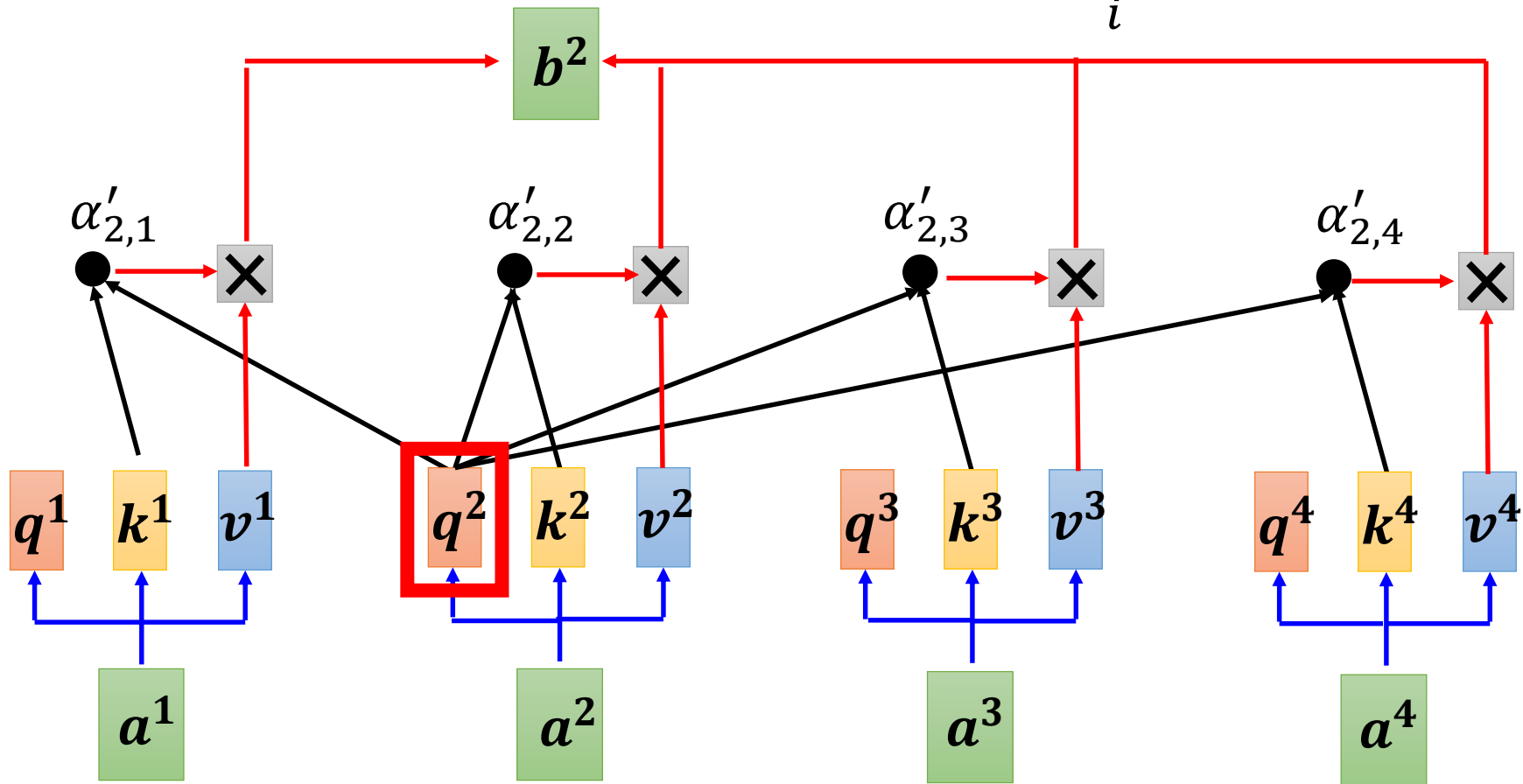parallel

$b^1$  $b^2$  $b^3$  $b^4$

$a^1$  $a^2$  $a^3$  $a^4$

Can be either **input** or **a hidden layer**

# Self-attention

$$b^2 = \sum_i \alpha'_{2,i} v^i$$

# Self-attention

$$q^i = W^q a^i$$

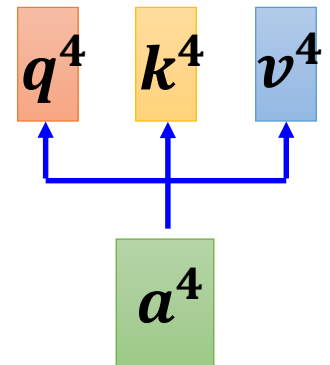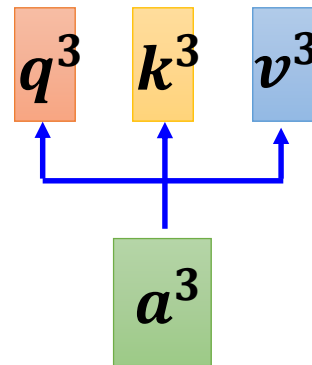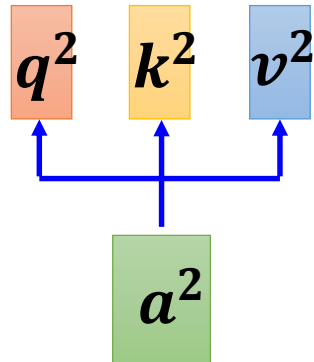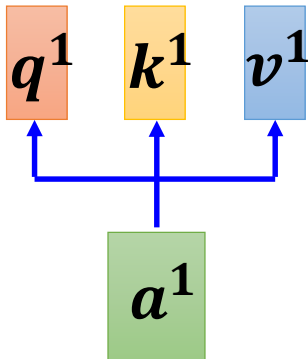$$\boxed{q^1\,q^2\,q^3\,q^4} = \boxed{W^q}\;\boxed{a^1\,a^2\,a^3\,a^4}$$

$$\underset{Q}{\quad}\qquad\qquad\underset{I}{\quad}$$

$$k^i = W^k a^i$$

$$\boxed{k^1\,k^2\,k^3\,k^4} = \boxed{W^k}\;\boxed{a^1\,a^2\,a^3\,a^4}$$

$$\underset{K}{\quad}\qquad\qquad\underset{I}{\quad}$$

$$v^i = W^v a^i$$

$$\boxed{v^1\,v^2\,v^3\,v^4} = \boxed{W^v}\;\boxed{a^1\,a^2\,a^3\,a^4}$$

$$\underset{V}{\quad}\qquad\qquad\underset{I}{\quad}$$

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

# *Self-attention*

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

$$\alpha_{1,3} = \boxed{k^3} \boxed{q^1} \quad \alpha_{1,4} = \boxed{k^4} \boxed{q^1}$$

$$\begin{bmatrix} \alpha_{1,1} \\ \alpha_{1,2} \\ \alpha_{1,3} \\ \alpha_{1,4} \end{bmatrix} = \begin{bmatrix} k^1 \\ k^2 \\ k^3 \\ k^4 \end{bmatrix} \boxed{q^1}$$

# *Self-attention*

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

$$\alpha_{1,3} = \boxed{k^3}\ \boxed{q^1} \quad \alpha_{1,4} = \boxed{k^4}\ \boxed{q^1}$$
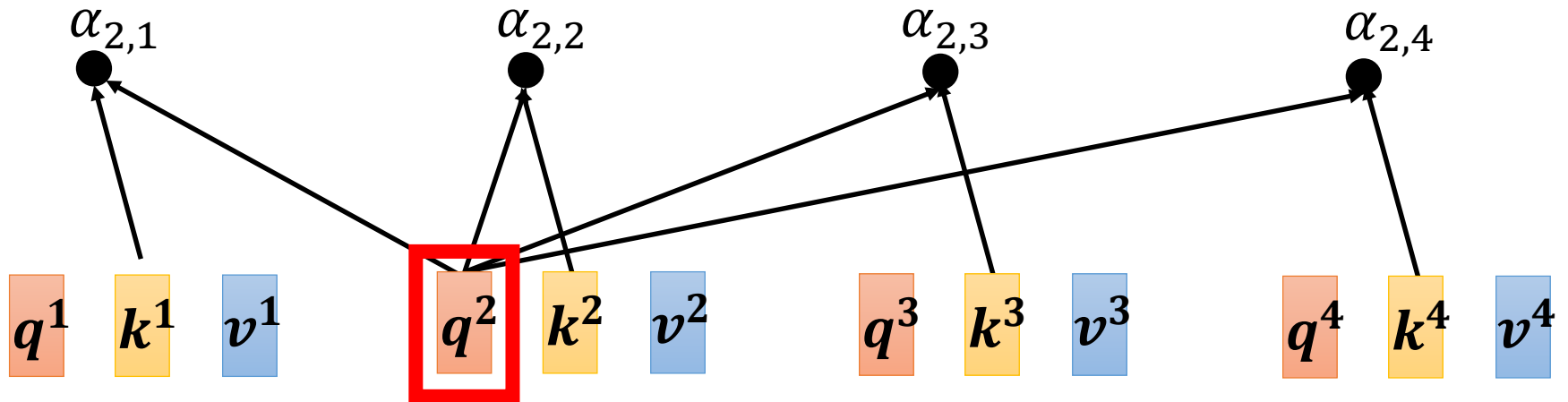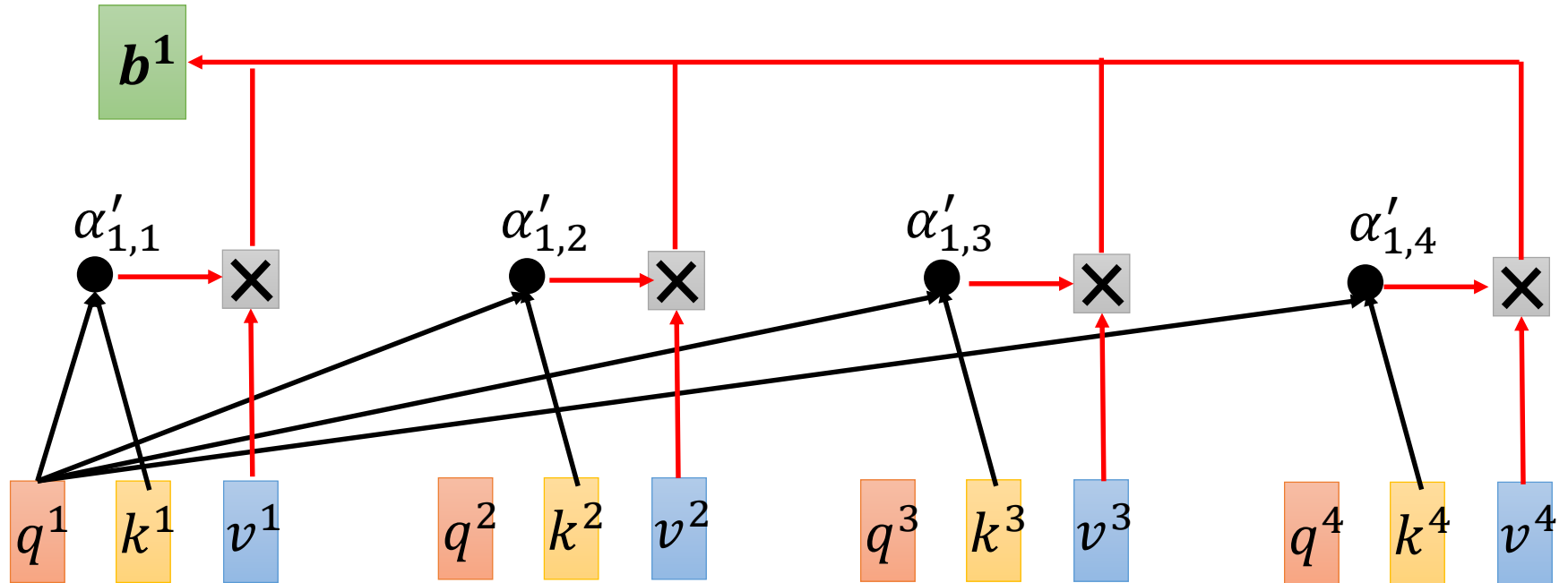
$$\begin{bmatrix} \alpha_{1,1} \\ \alpha_{1,2} \\ \alpha_{1,3} \\ \alpha_{1,4} \end{bmatrix} = \begin{bmatrix} k^1 \\ k^2 \\ k^3 \\ k^4 \end{bmatrix} \boxed{q^1}$$

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

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

$$\begin{bmatrix} \alpha'_{1,1} & \alpha'_{2,1} & \alpha'_{3,1} & \alpha'_{4,1} \\ \alpha'_{1,2} & \alpha'_{2,2} & \alpha'_{3,2} & \alpha'_{4,2} \\ \alpha'_{1,3} & \alpha'_{2,3} & \alpha'_{3,3} & \alpha'_{4,3} \\ \alpha'_{1,4} & \alpha'_{2,4} & \alpha'_{3,4} & \alpha'_{4,4} \end{bmatrix} \longleftarrow \begin{bmatrix} \alpha_{1,1} & \alpha_{2,1} & \alpha_{3,1} & \alpha_{4,1} \\ \alpha_{1,2} & \alpha_{2,2} & \alpha_{3,2} & \alpha_{4,2} \\ \alpha_{1,3} & \alpha_{2,3} & \alpha_{3,3} & \alpha_{4,3} \\ \alpha_{1,4} & \alpha_{2,4} & \alpha_{3,4} & \alpha_{4,4} \end{bmatrix} = \begin{bmatrix} k^1 \\ k^2 \\ k^3 \\ k^4 \end{bmatrix} \begin{bmatrix} q^1 & q^2 & q^3 & q^4 \end{bmatrix}$$

$A'$  softmax  $A$  $K^T$  $Q$

# _Self-attention_



$b^1 b^2 b^3 b^4 = v^1 v^2 v^3 v^4$
$O \qquad\qquad V$

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

$A'$

# *Self-attention*

$$Q = W^q I$$

$$K = W^k I$$

$$V = W^v I$$

Parameters to be learned

$$A' \leftarrow A = K^T Q$$

Attention Matrix

$$O = V A'$$

13

# Multi-head Self-attention　Different types of relevance

$$q^{i,1} = W^{q,1} q^i$$

$$q^{i,2} = W^{q,2} q^i$$



$b^{i,1}$

$q^{i,1} \quad q^{i,2} \quad k^{i,1} \quad k^{i,2} \quad v^{i,1} \quad v^{i,2} \quad q^{j,1} \quad q^{j,2} \quad k^{j,1} \quad k^{j,2} \quad v^{j,1} \quad v^{j,2}$

$q^i \quad k^i \quad v^i \quad q^j \quad k^j \quad v^j$

$$q^i = W^q a^i$$

$a^i$　(2 heads as example)　$a^j$

# *Multi-head Self-attention*   Different types of relevance

$q^{i,1} = W^{q,1} q^i$

$q^{i,2} = W^{q,2} q^i$

$b^{i,1}$

$b^{i,2}$

$q^{i,1}$ $q^{i,2}$ $k^{i,1}$ $k^{i,2}$ $v^{i,1}$ $v^{i,2}$ $q^{j,1}$ $q^{j,2}$ $k^{j,1}$ $k^{j,2}$ $v^{j,1}$ $v^{j,2}$

$q^i$ $k^i$ $v^i$ $q^j$ $k^j$ $v^j$

$q^i = W^q a^i$   $a^i$   (2 heads as example)   $a^j$

# _Multi-head Self-attention_ Different types of relevance

$$b^i = W^O \begin{bmatrix} b^{i,1} \\ b^{i,2} \end{bmatrix}$$

$q^{i,1}$ $q^{i,2}$ $k^{i,1}$ $k^{i,2}$ $v^{i,1}$ $v^{i,2}$    $q^{j,1}$ $q^{j,2}$ $k^{j,1}$ $k^{j,2}$ $v^{j,1}$ $v^{j,2}$

$q^i$ $k^i$ $v^i$    $q^j$ $k^j$ $v^j$

$q^i = W^q a^i$

$a^i$    (2 heads as example)    $a^j$