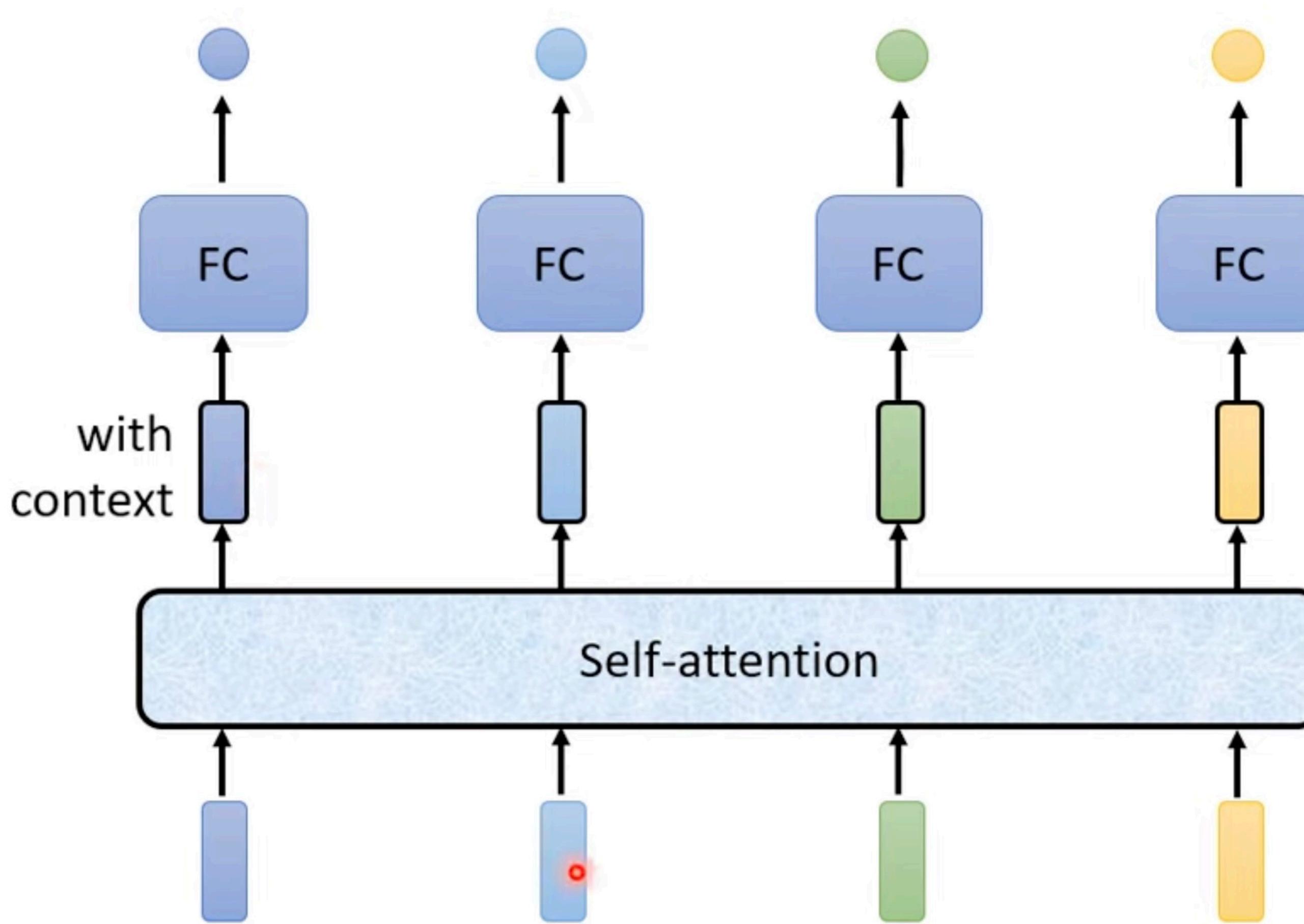
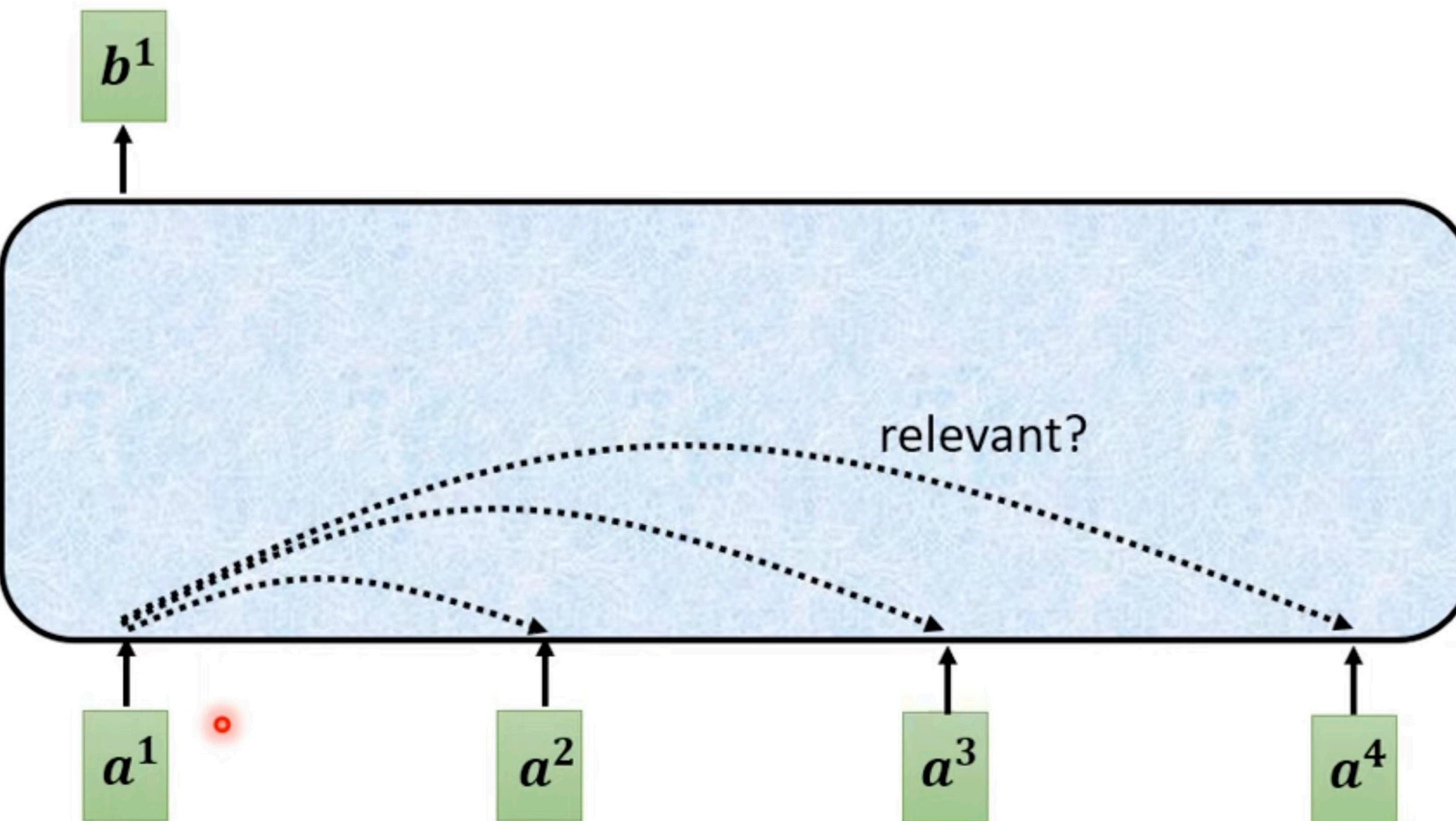


Self-attention



這邊就給你一個另外一個Vector

Self-attention



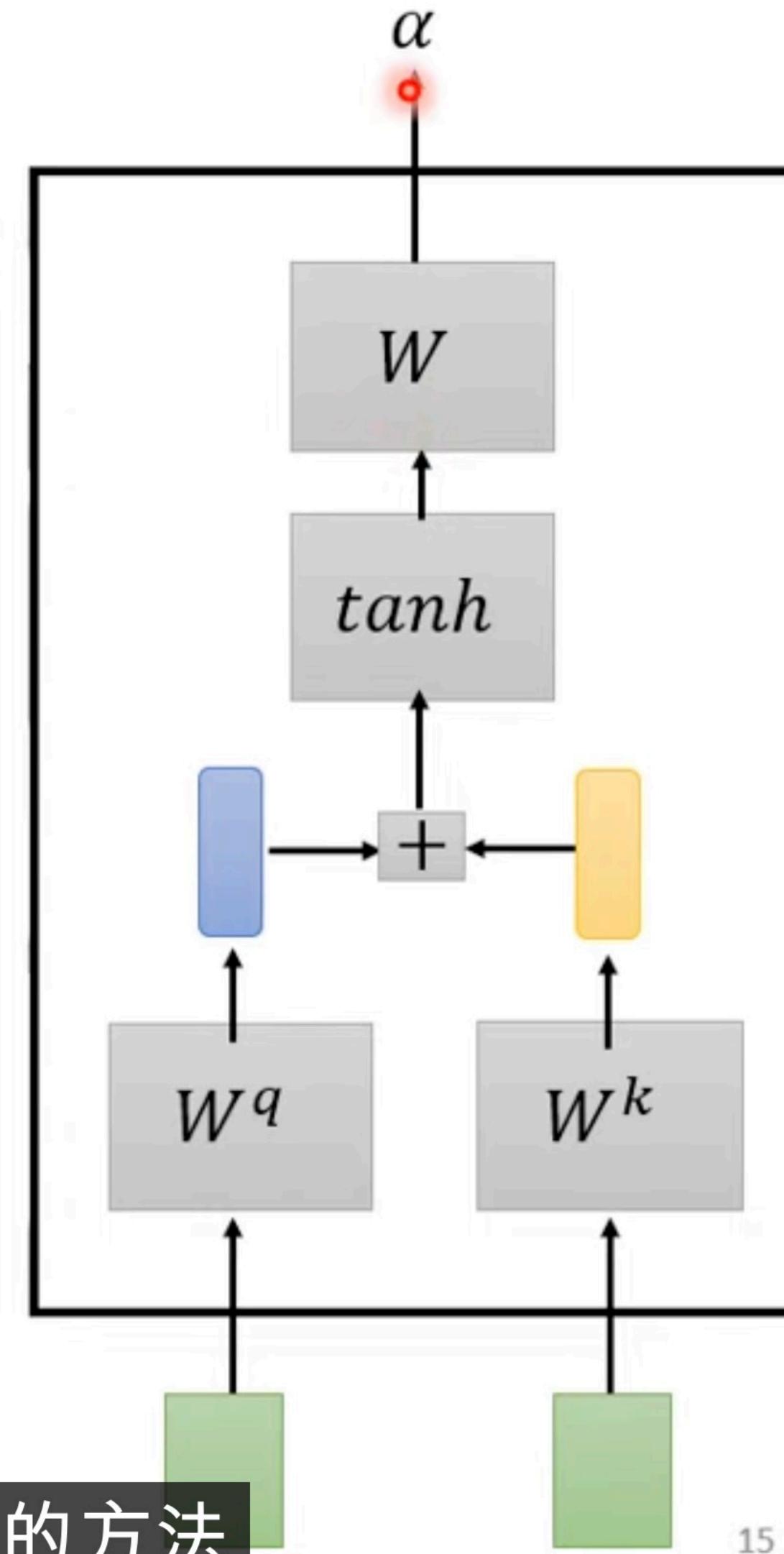
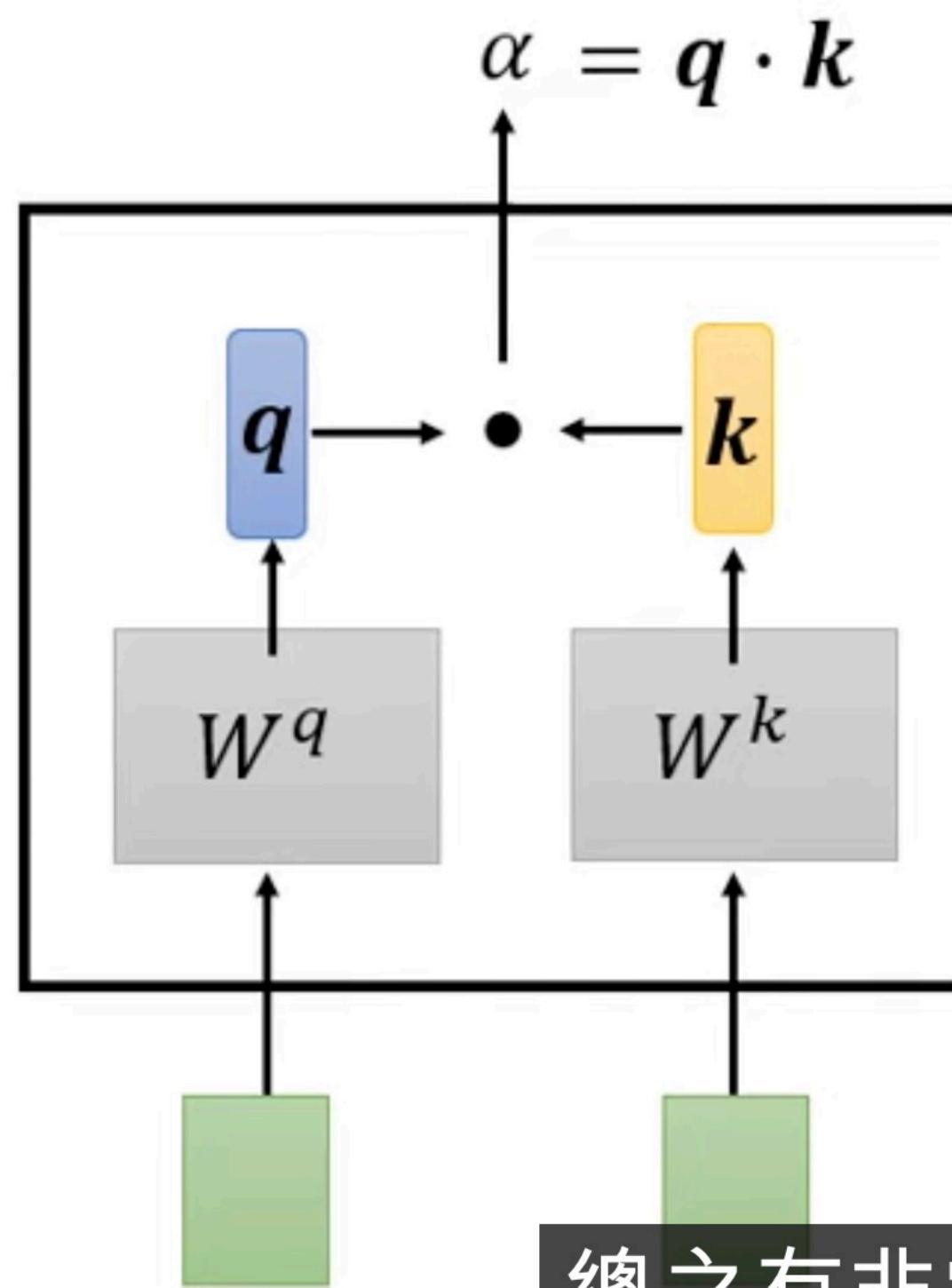
Find the relevant vectors in a sequence

所需要用到的資訊

Self-attention

Additive

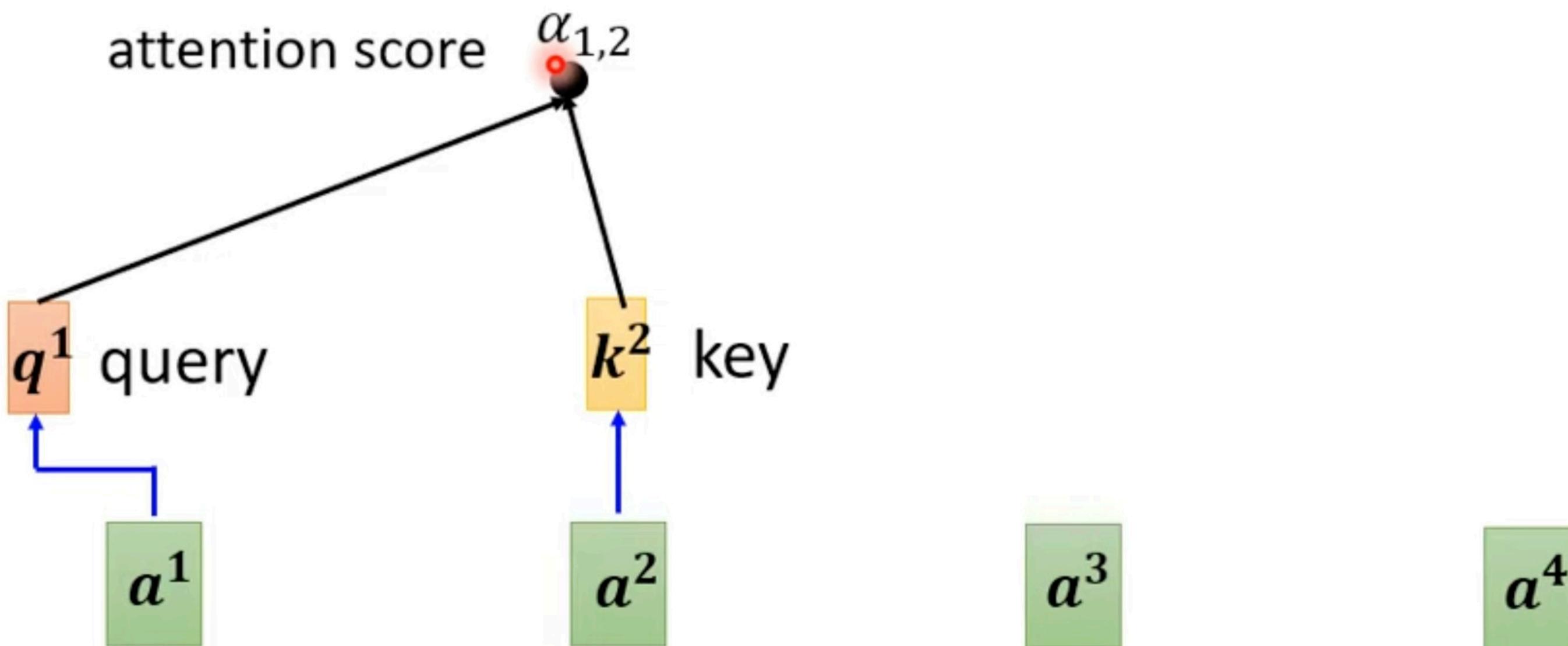
Dot-product



總之有非常多不同的方法

Self-attention

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

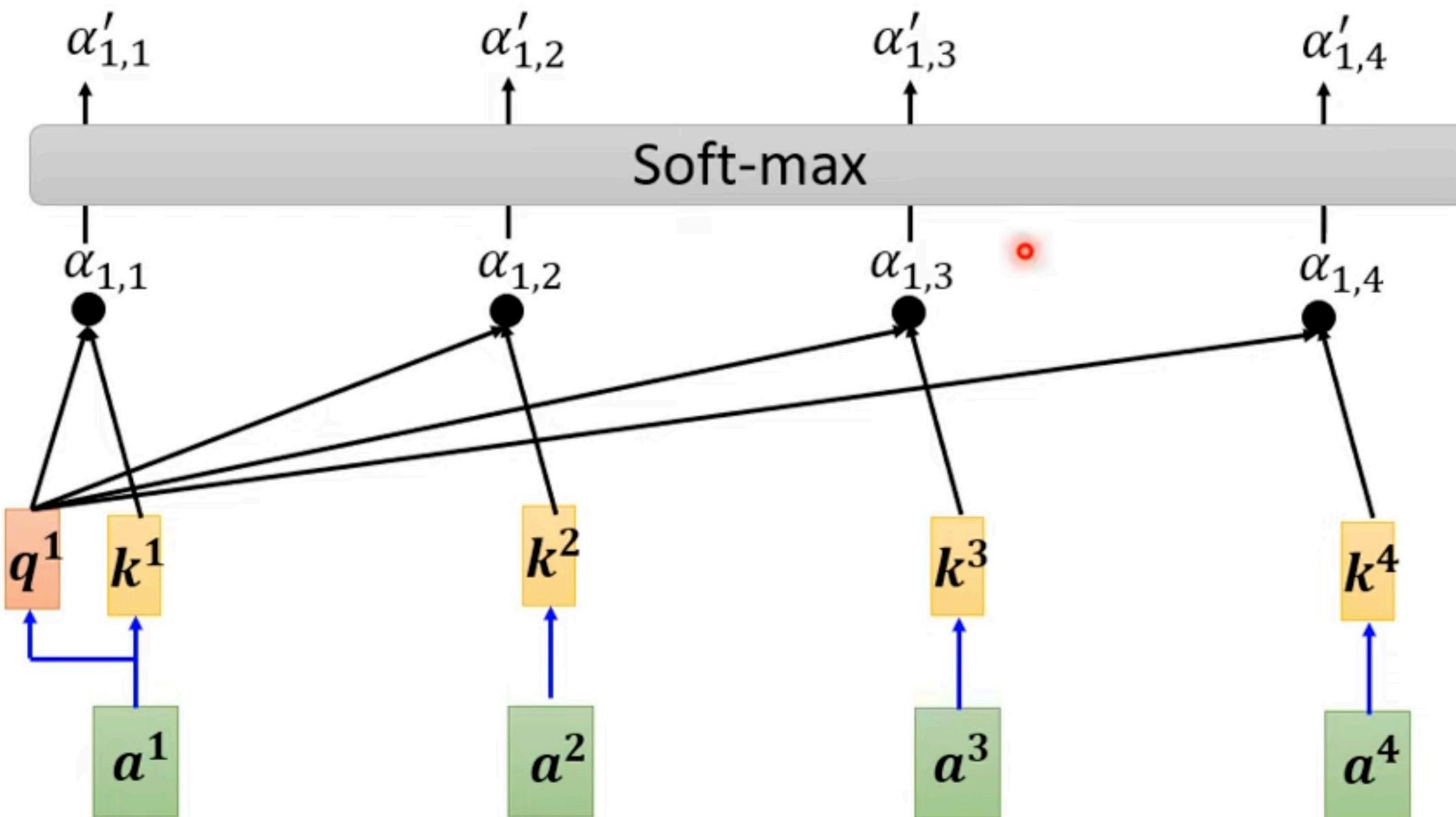


$$q^1 = W^q a^1 \quad k^2 = W^k a^2$$

叫做Attention的Score

Self-attention

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



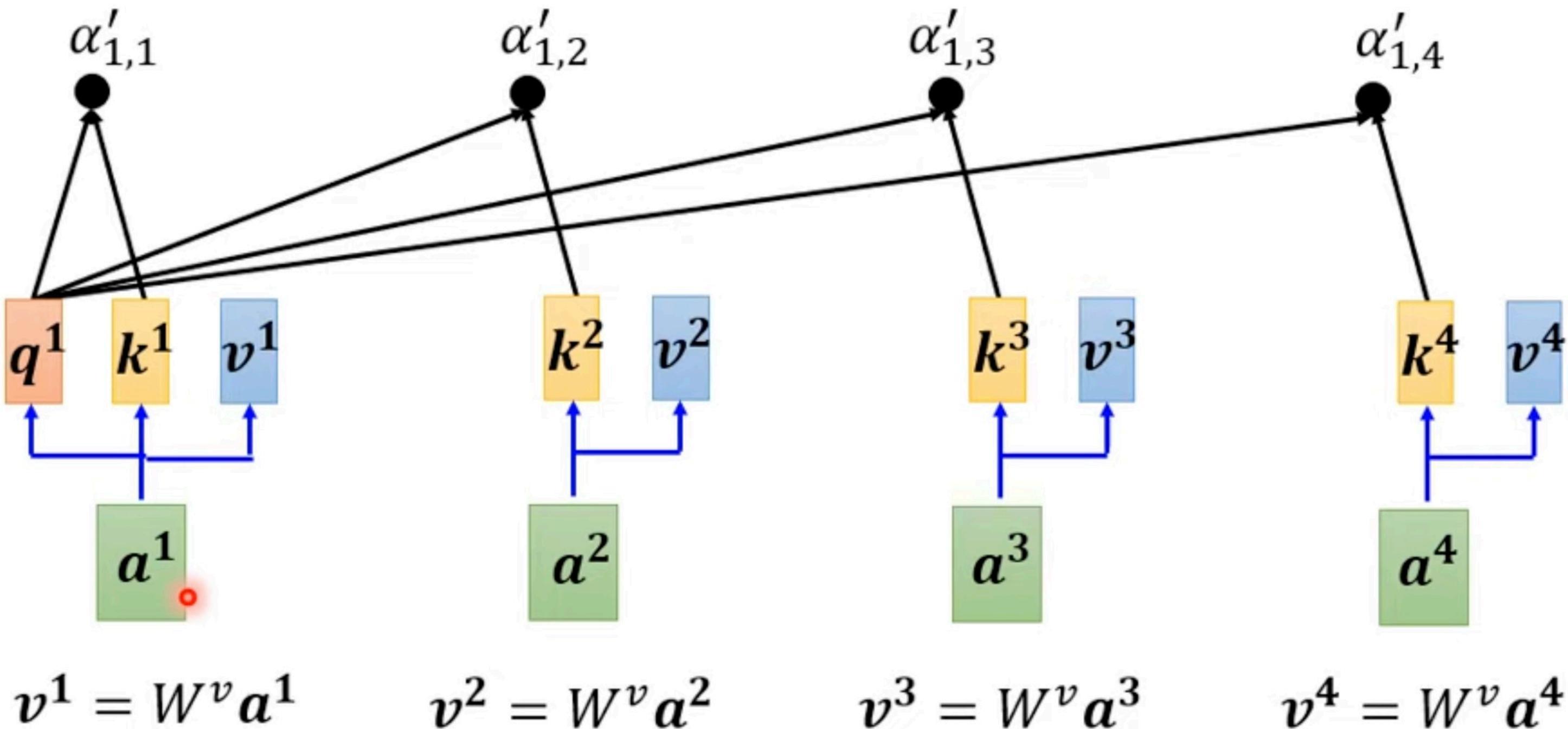
$$q^1 = W^q a^1 \quad k^2 = W^k a^2 \quad k^3 = W^k a^3 \quad k^4 = W^k a^4$$

$$k^1 = W^k a^1$$

接下來這邊會作為一個Soft-Max

Self-attention

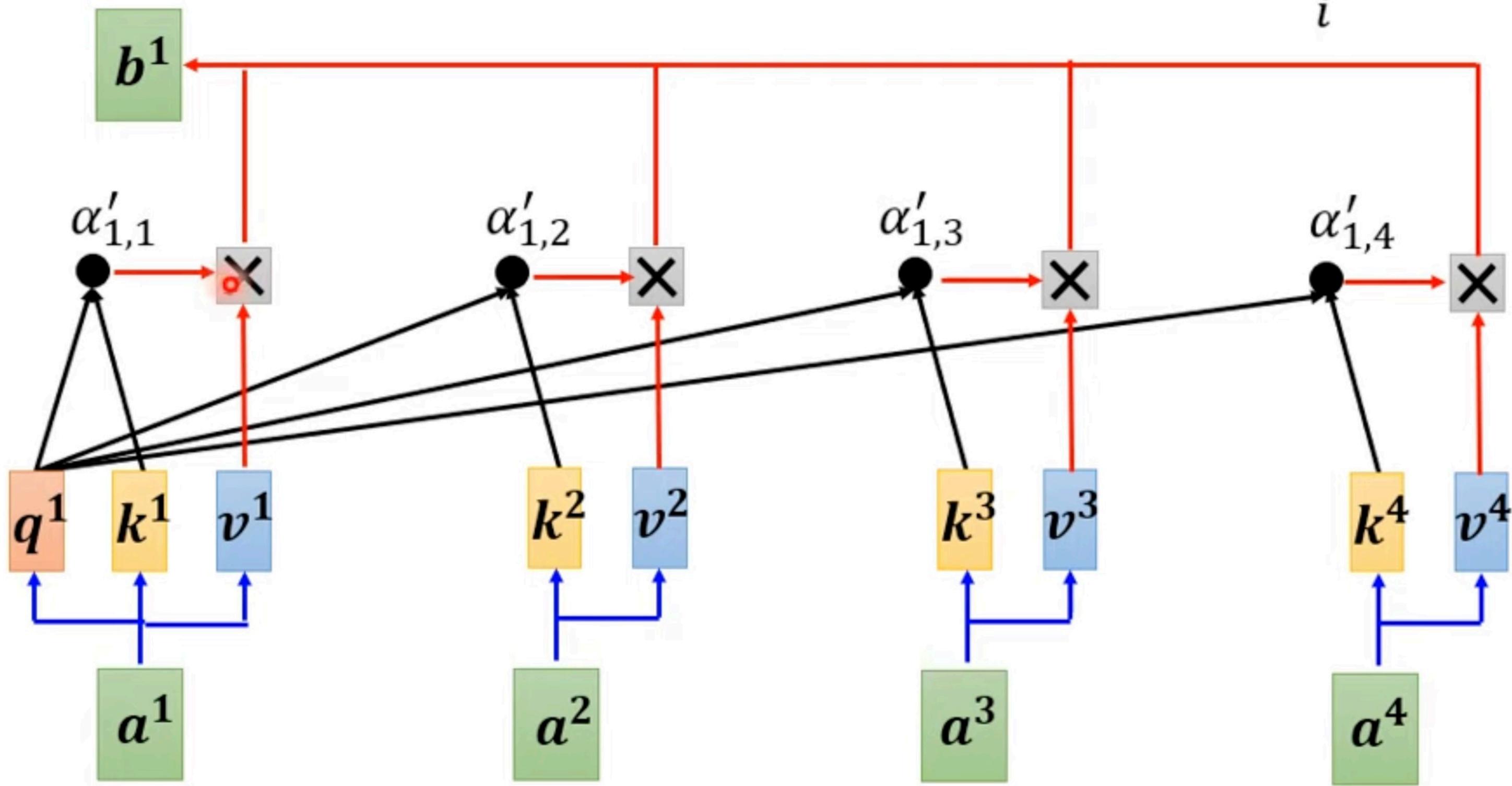
Extract information based
on attention scores



Self-attention

Extract information based
on attention scores

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

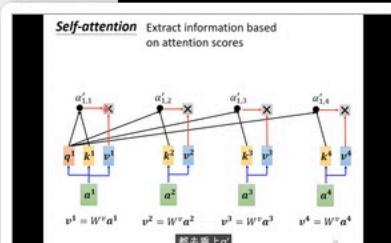


$$v^1 = W^v a^1$$

$$v^2 = W^v a^2$$

$$v^3 = W^v a^3$$

$$v^4 = W^v a^4$$



Self-attention

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

$$q^1 \boxed{q^2} \boxed{q^3} \boxed{q^4}$$

Q

$$W^q$$

$$a^1 \boxed{a^2} \boxed{a^3} \boxed{a^4}$$

I

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

$$k^1 \boxed{k^2} \boxed{k^3} \boxed{k^4}$$

K

$$W^k$$

$$a^1 \boxed{a^2} \boxed{a^3} \boxed{a^4}$$

I

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

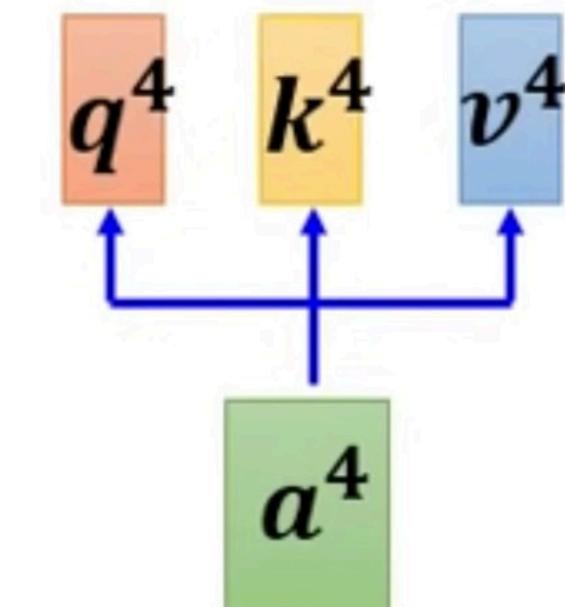
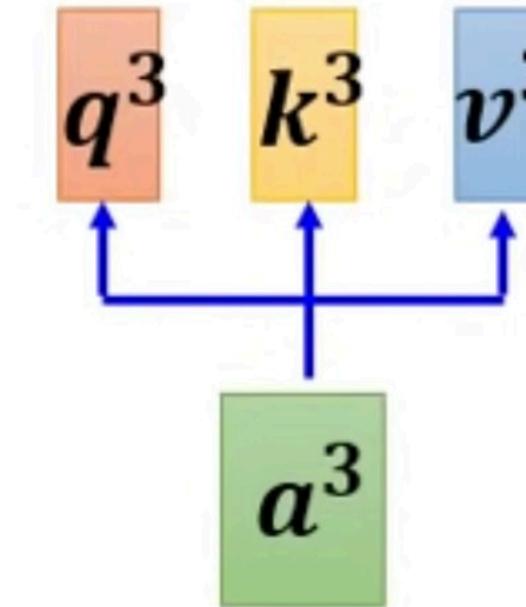
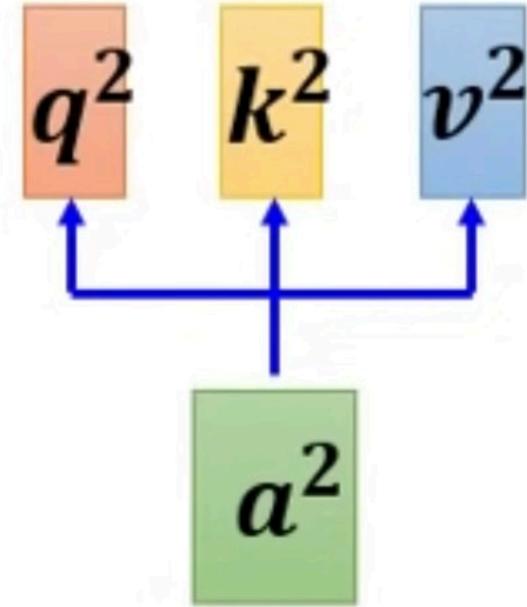
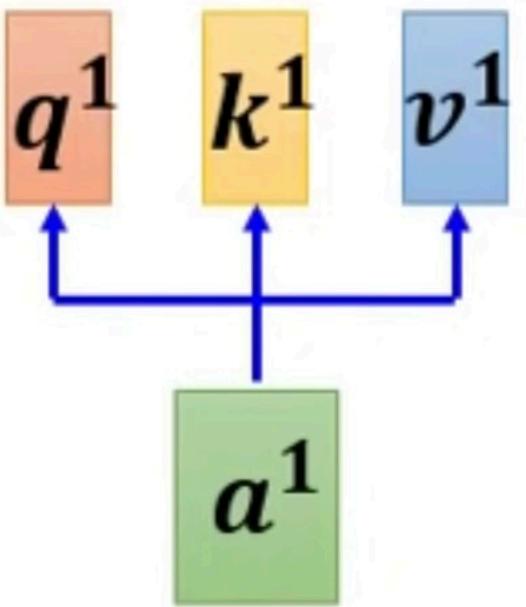
$$v^1 \boxed{v^2} \boxed{\color{red}{v^3}} \boxed{v^4}$$

V

$$W^v$$

$$a^1 \boxed{a^2} \boxed{a^3} \boxed{a^4}$$

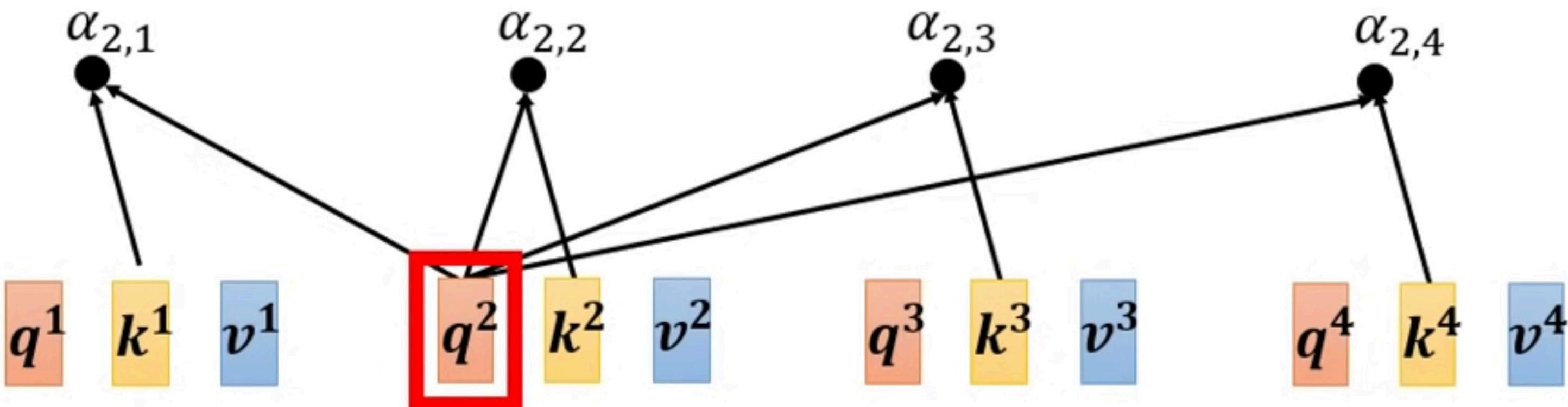
I



Self-attention

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

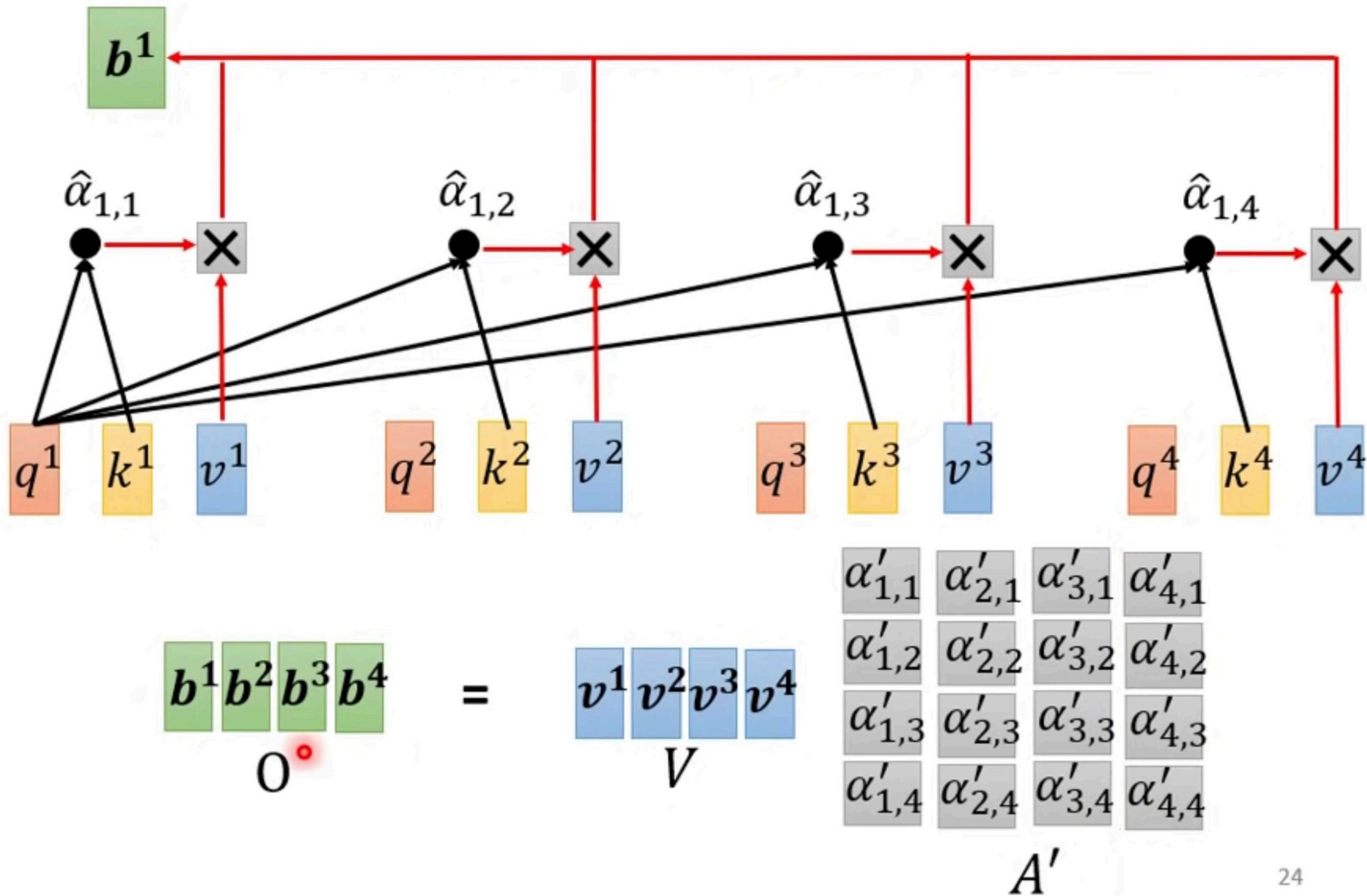
$$\begin{matrix} \alpha_{1,1} \\ \alpha_{1,2} \\ \alpha_{1,3} \\ \alpha_{1,4} \end{matrix} = \begin{matrix} k^1 \\ k^2 \\ k^3 \\ k^4 \end{matrix} \quad \begin{matrix} q^1 \\ q^1 \\ q^1 \\ q^1 \end{matrix}$$



$$\begin{array}{c} \begin{matrix} \alpha'_{1,1} & \alpha'_{2,1} & \alpha'_{3,1} & \alpha'_{4,1} \\ \alpha'_{1,2} & \alpha'_{2,2} & \alpha'_{3,2} & \alpha'_{4,2} \\ \alpha'_{1,3} & \alpha'_{2,3} & \alpha'_{3,3} & \alpha'_{4,3} \\ \alpha'_{1,4} & \alpha'_{2,4} & \alpha'_{3,4} & \alpha'_{4,4} \end{matrix} \leftarrow \text{softmax}^\circ \quad \begin{matrix} \alpha_{1,1} & \alpha_{2,1} & \alpha_{3,1} & \alpha_{4,1} \\ \alpha_{1,2} & \alpha_{2,2} & \alpha_{3,2} & \alpha_{4,2} \\ \alpha_{1,3} & \alpha_{2,3} & \alpha_{3,3} & \alpha_{4,3} \\ \alpha_{1,4} & \alpha_{2,4} & \alpha_{3,4} & \alpha_{4,4} \end{matrix} \\ A' \quad \text{softmax}^\circ \quad A \end{array} = \begin{matrix} k^1 \\ k^2 \\ k^3 \\ k^4 \end{matrix} \quad \begin{matrix} q^1 \\ q^2 \\ q^3 \\ q^4 \end{matrix} \quad Q$$

$$= \begin{matrix} k^1 \\ k^2 \\ k^3 \\ k^4 \end{matrix} \quad \begin{matrix} v^1 \\ v^2 \\ v^3 \\ v^4 \end{matrix} \quad K^T$$

Self-attention



Self-attention

$$\begin{array}{ccc} Q & = & W^q \\ K & = & W^k \\ V & = & W^v \end{array}$$

I I I

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

Attention Matrix

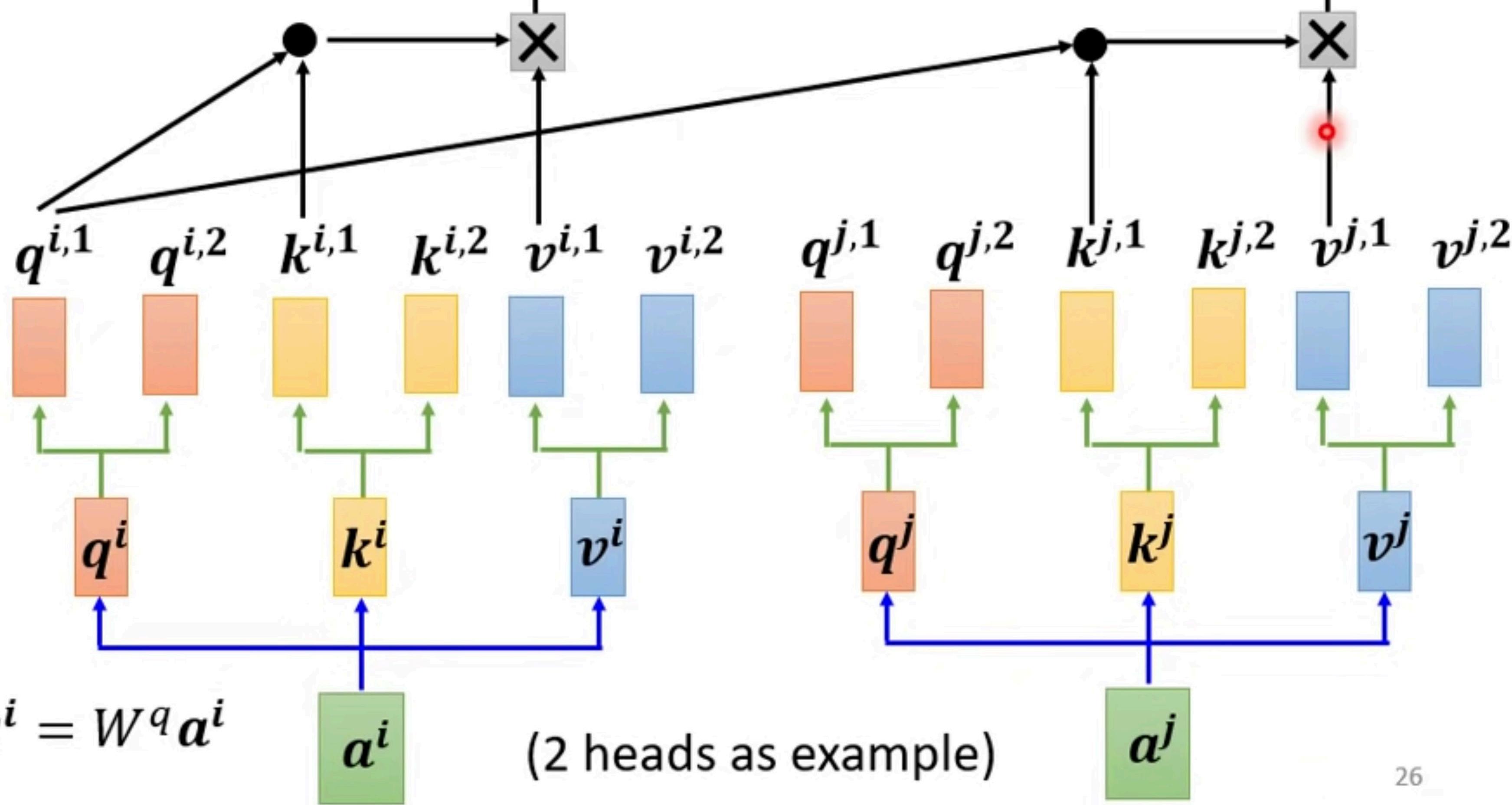
$$O = V A'$$

Multi-head Self-attention

Different types of relevance

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

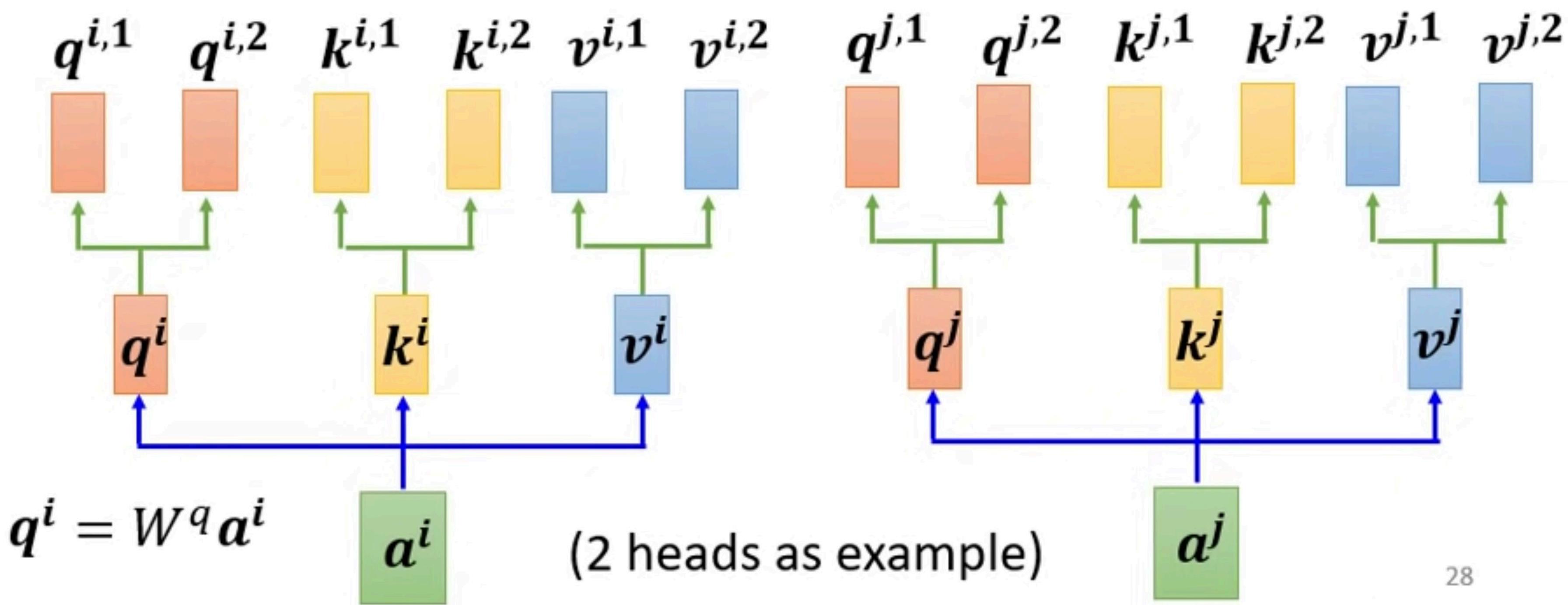
$$q^{i,2} = W^{q,2} q^i$$



Multi-head Self-attention

Different types of relevance

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



Positional Encoding

Each column represents a positional vector e^i

- No position information in self-attention.
- Each position has a unique positional vector e^i
- **hand-crafted**
- **learned from data**

