

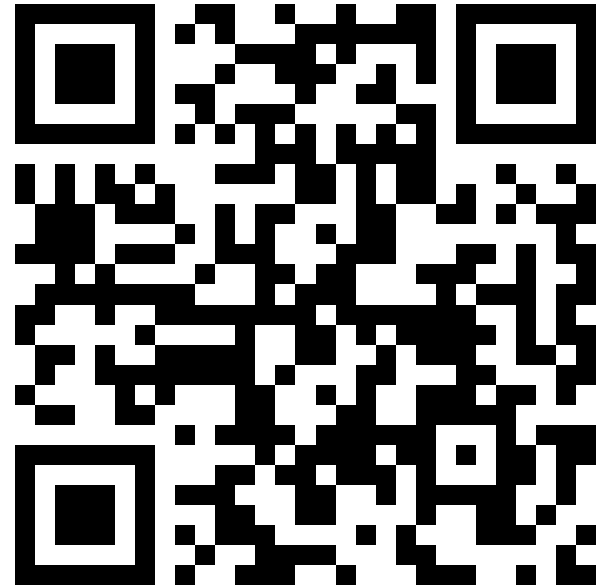
各式各樣的 Attention

Hung-yi Lee 李宏毅

Prerequisite



<https://youtu.be/hYdO9CscNes>

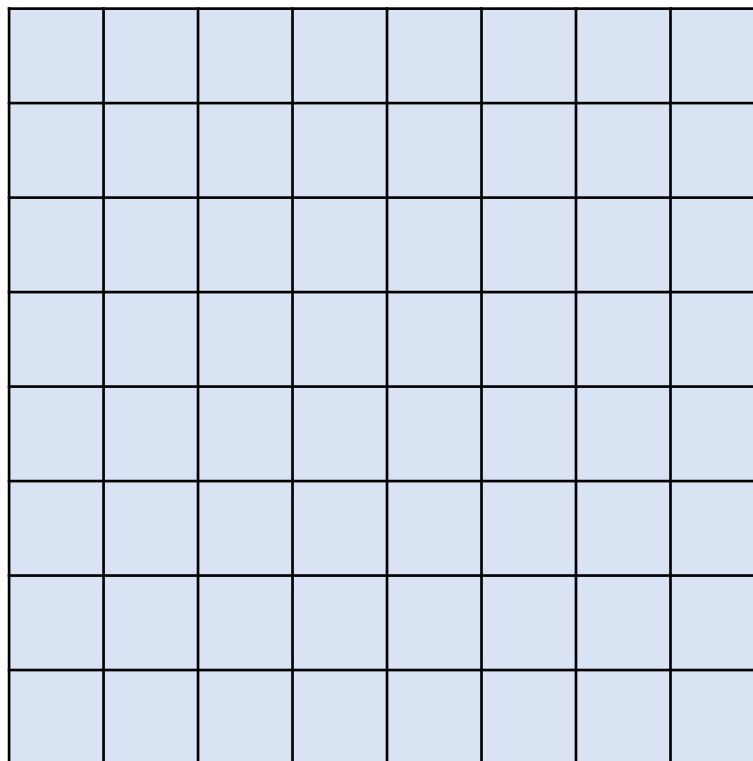
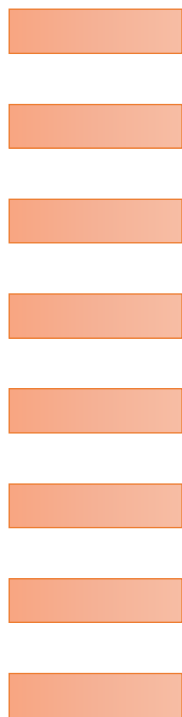


<https://youtu.be/gmsMY5kc-zw>

Review

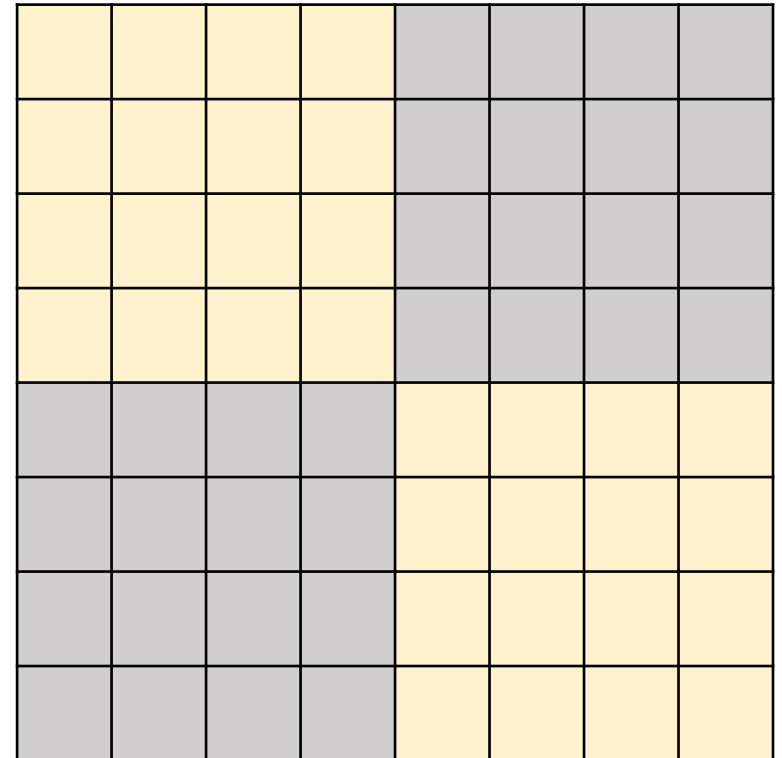
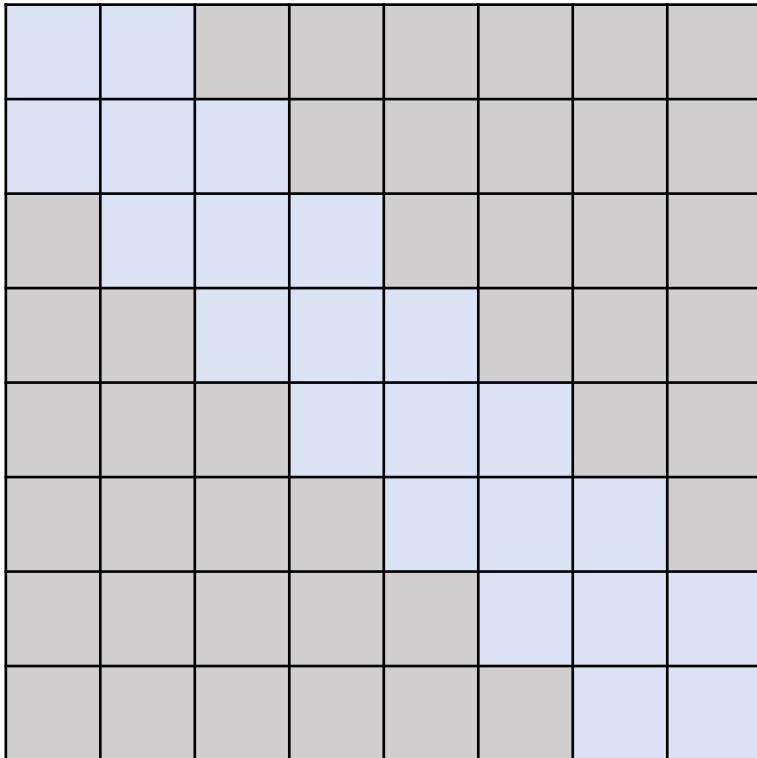
key

query



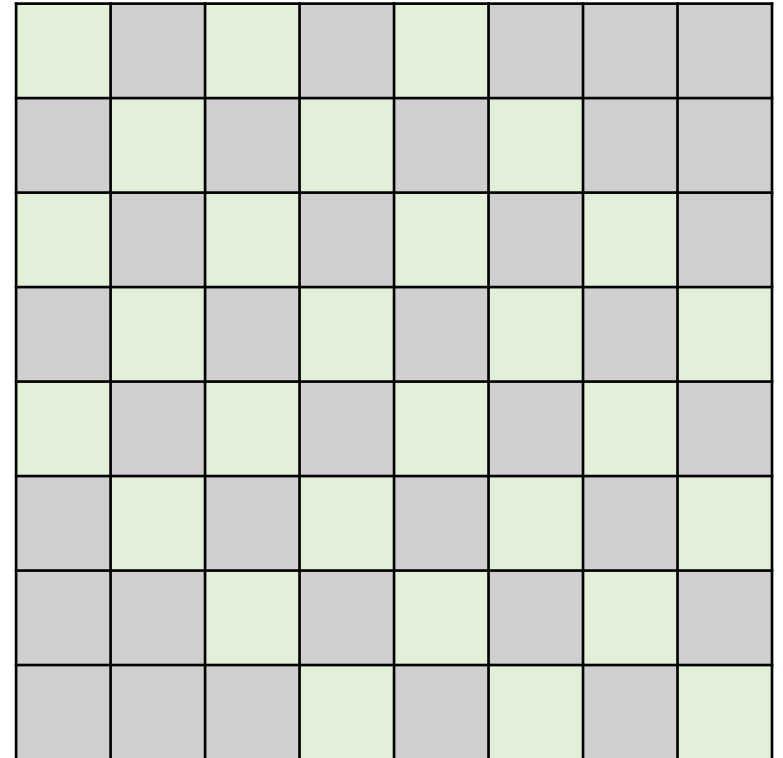
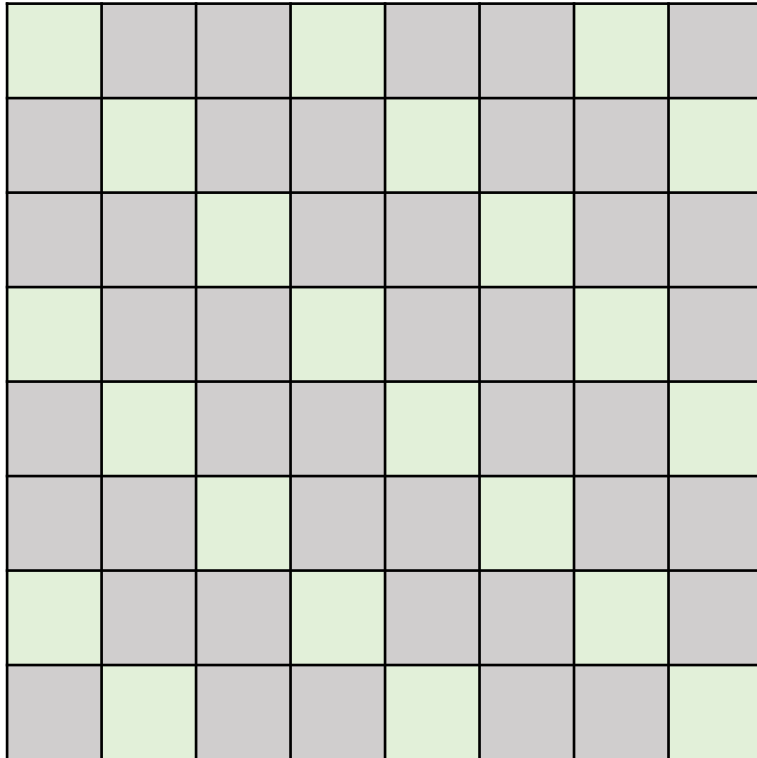
Attention Matrix

Fixed Pattern

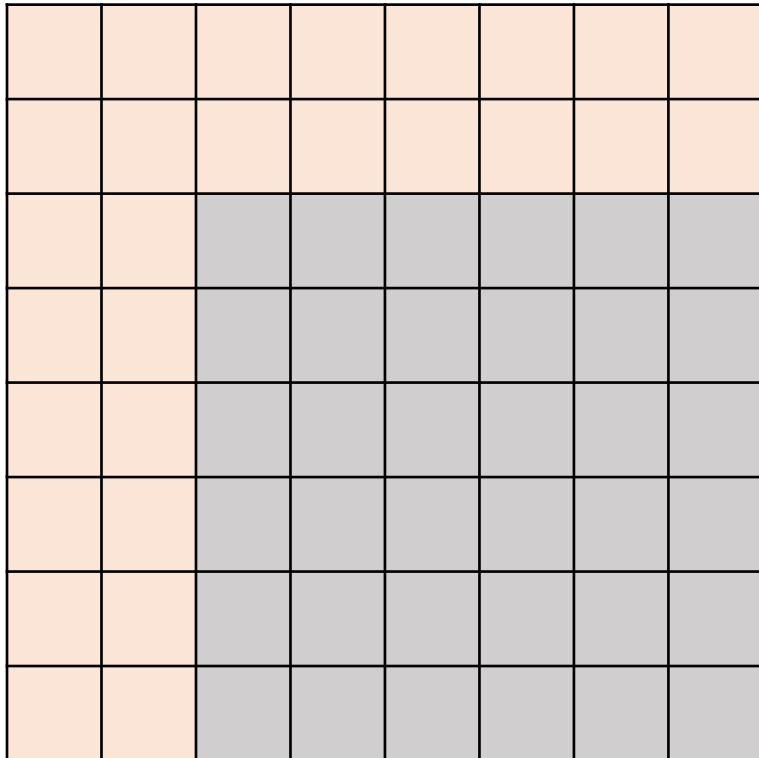


c.f. CNN

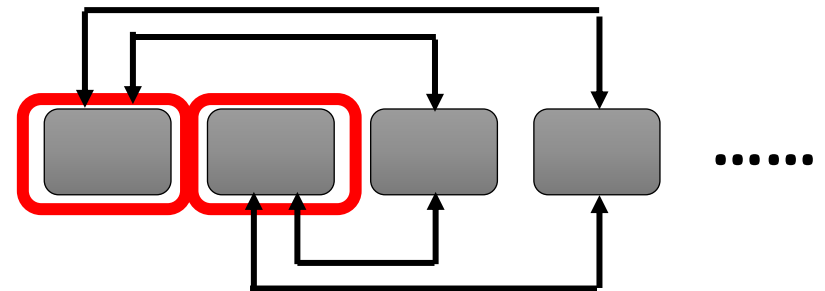
Fixed Pattern



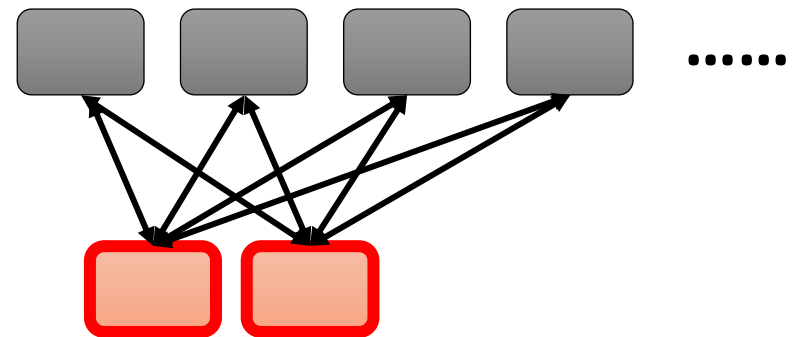
Fixed Pattern



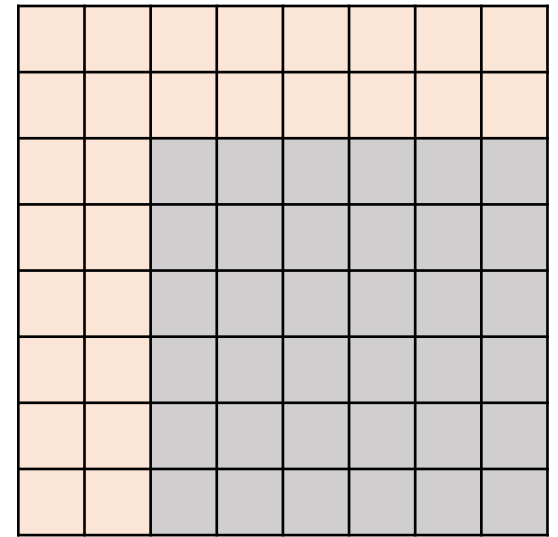
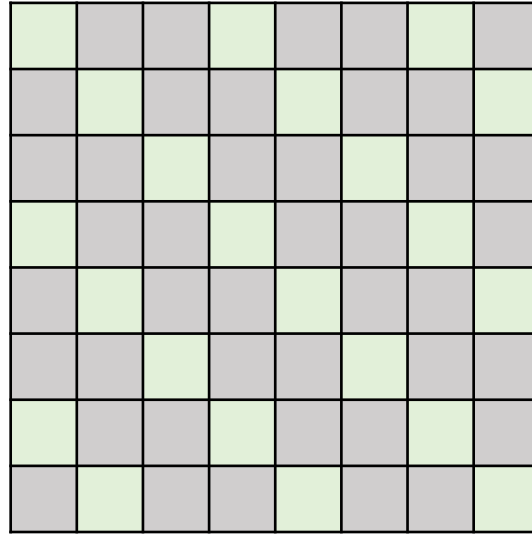
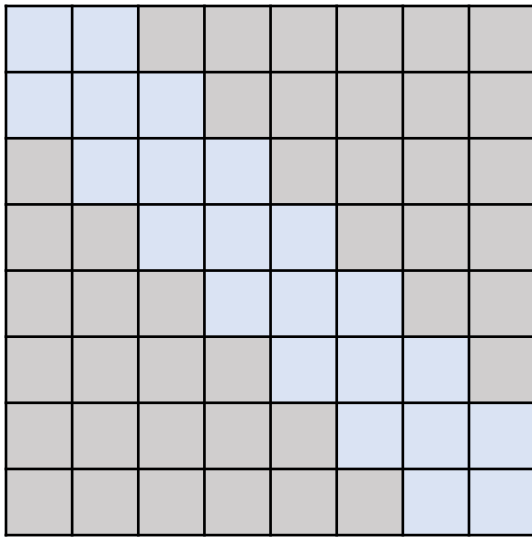
Internal transformer
construction (ITC)



Extended transformer
construction (TTC)



Fixed Pattern

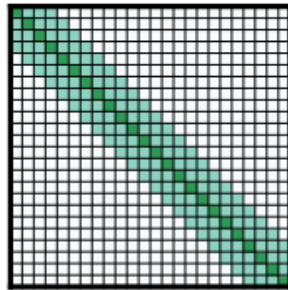


Different heads use different patterns.

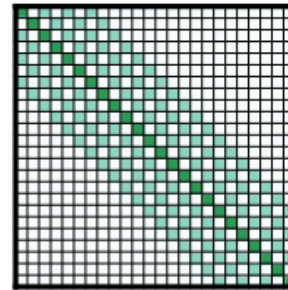
Fixed Pattern

- Longformer

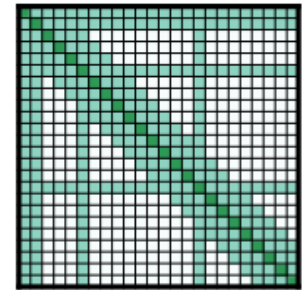
<https://arxiv.org/abs/2004.05150>



(b) Sliding window attention



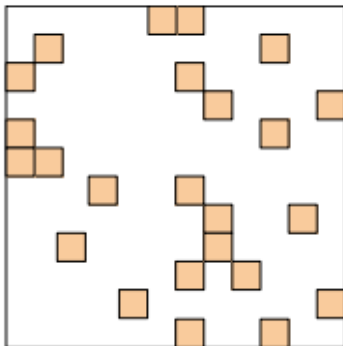
(c) Dilated sliding window



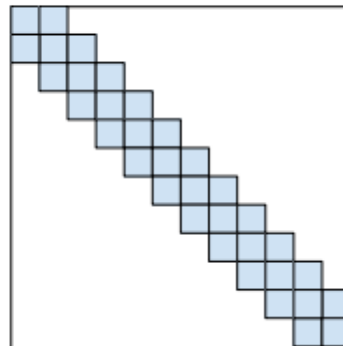
(d) Global+sliding window

- Big Bird

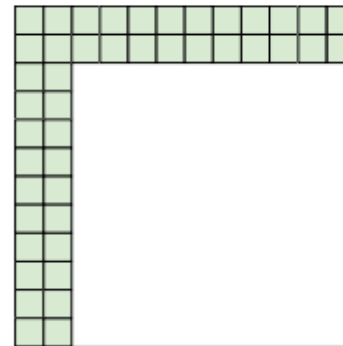
<https://arxiv.org/abs/2007.14062>



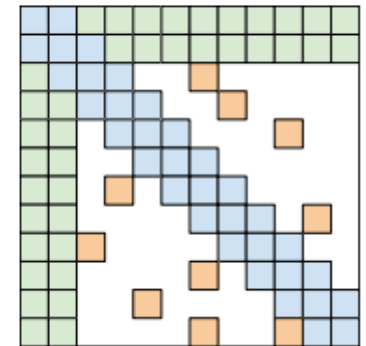
(a) Random attention



(b) Window attention



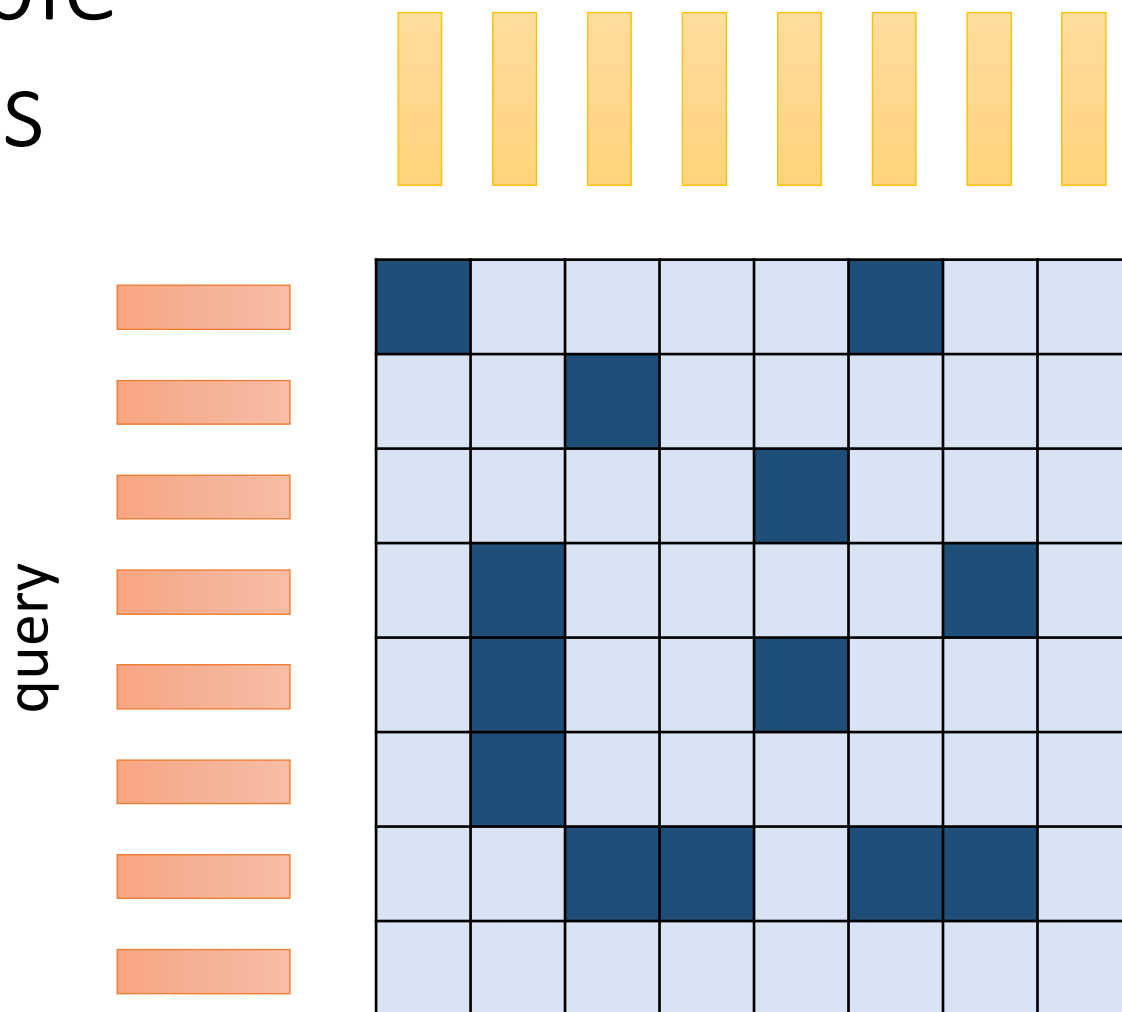
(c) Global Attention



(d) BIGBIRD

Learnable Patterns

key



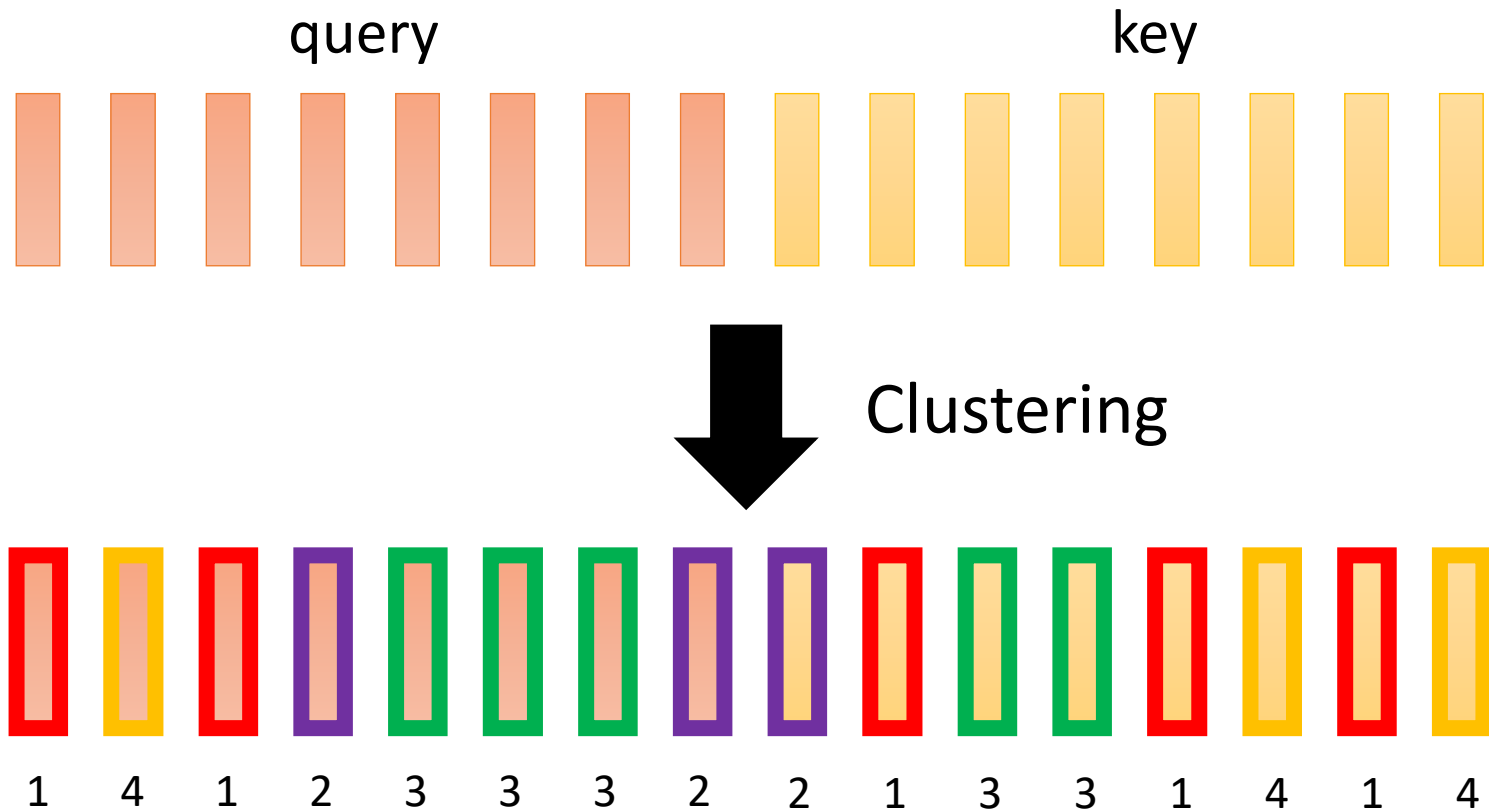
Learnable Patterns

Reformer

<https://openreview.net/forum?id=rkgNKkHtvB>

Routing Transformer

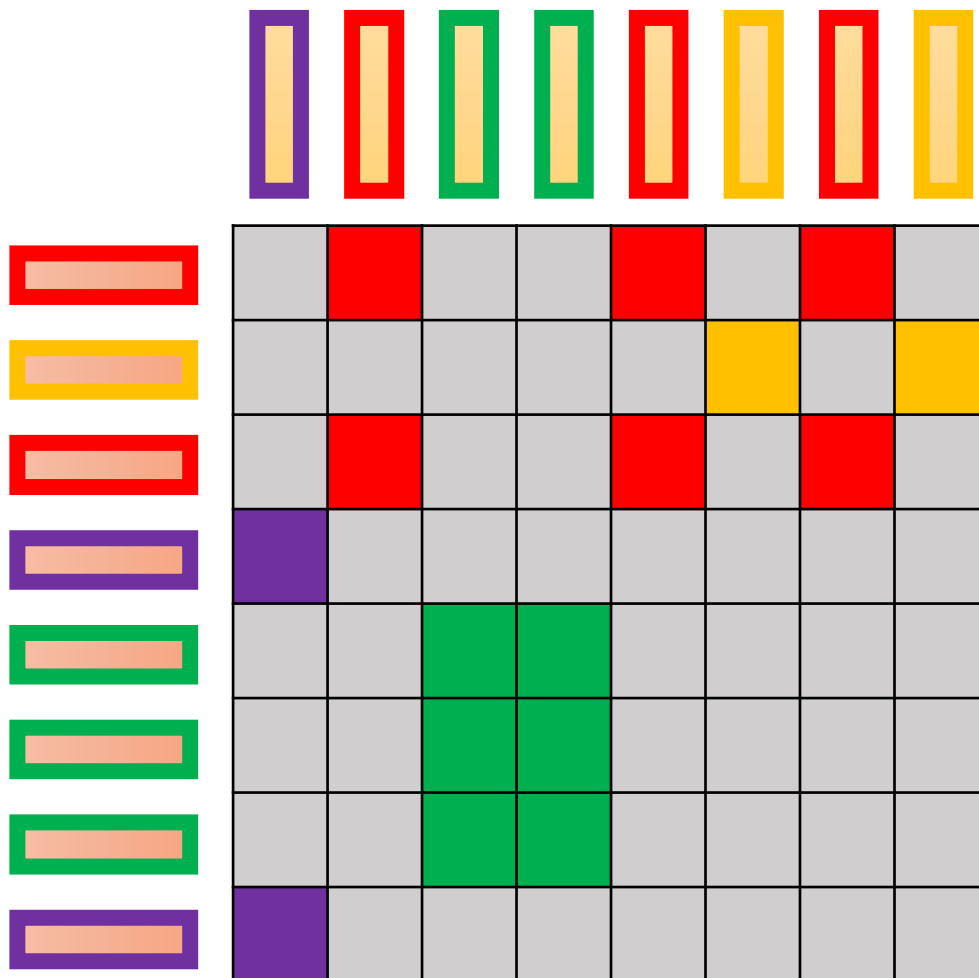
<https://arxiv.org/abs/2003.05997>



Learnable Patterns

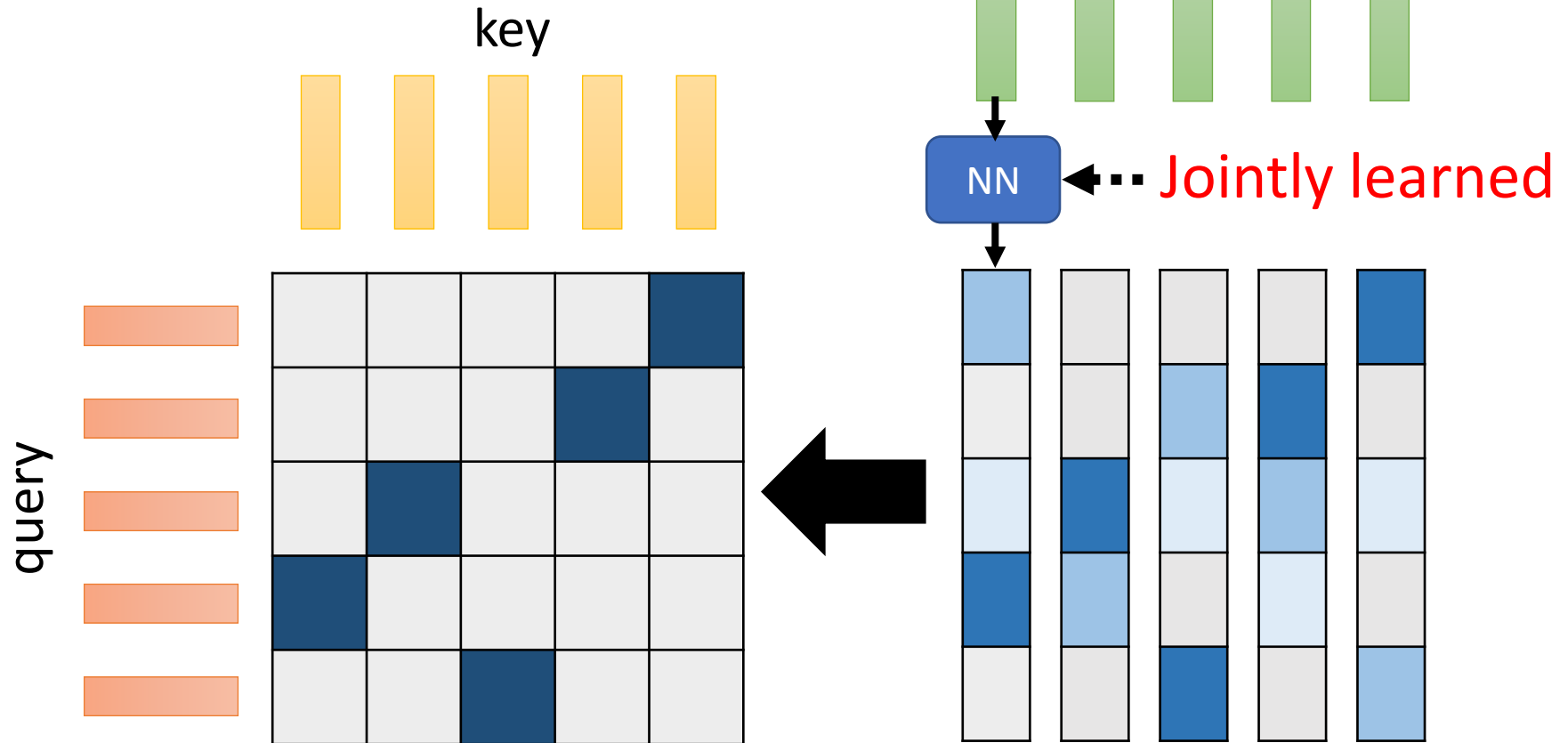
key

query



Learnable Patterns

Sinkhorn Sorting Network

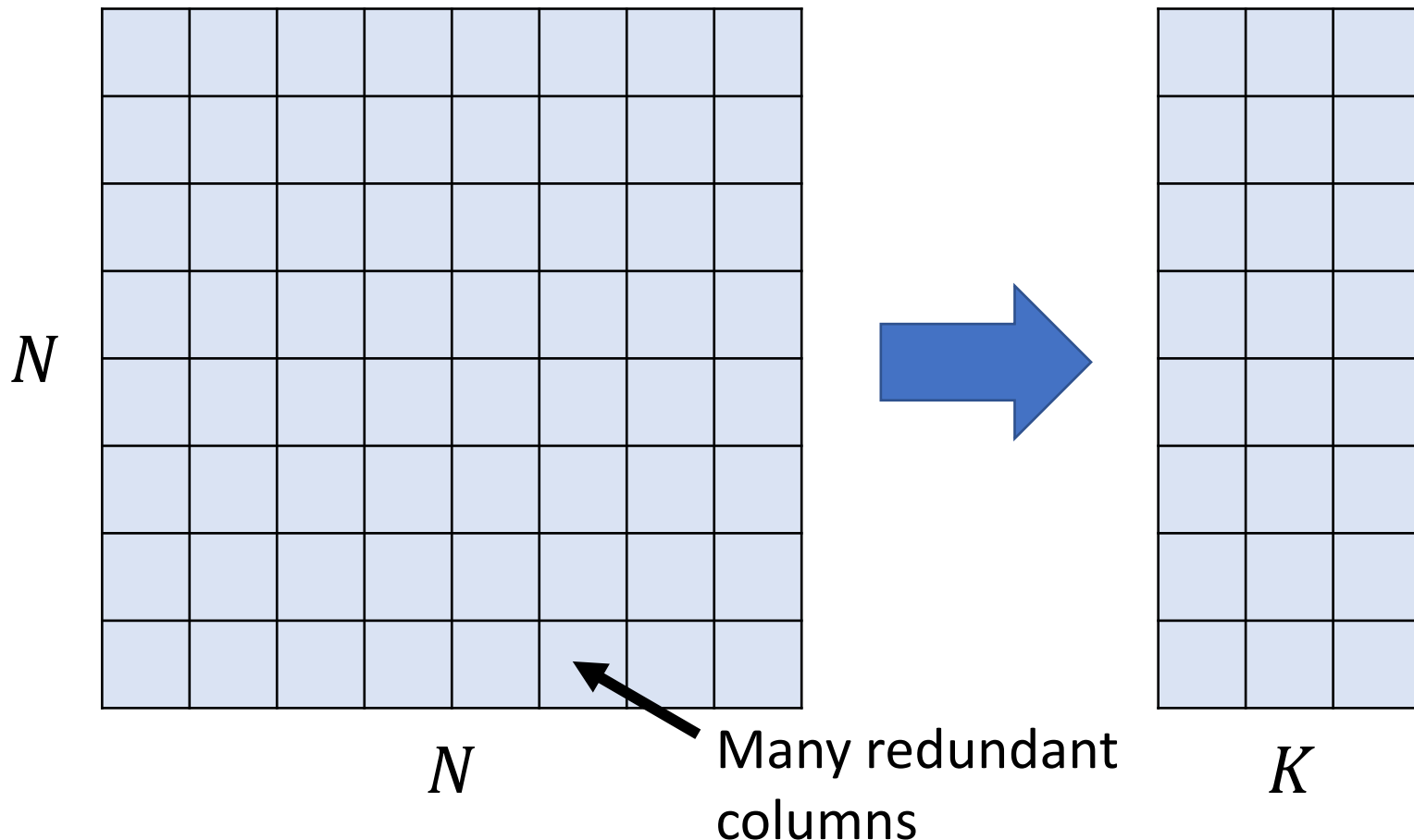


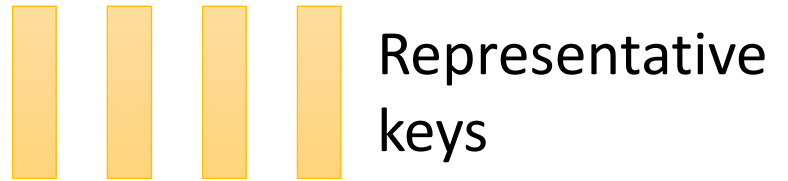
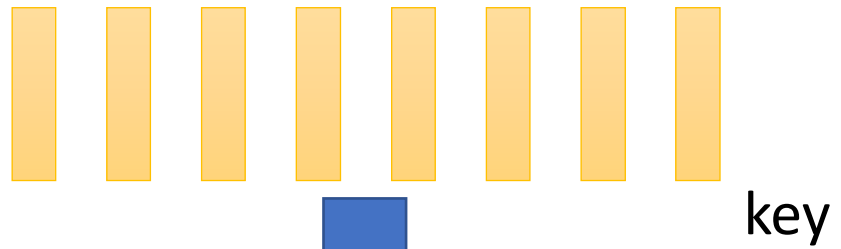
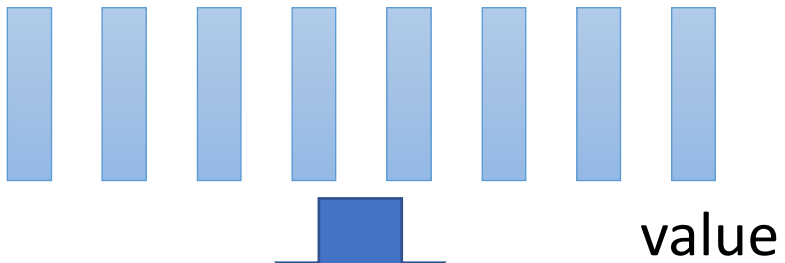
(simplified version)

Do we need full attention matrix?

Linformer

<https://arxiv.org/abs/2006.04768>

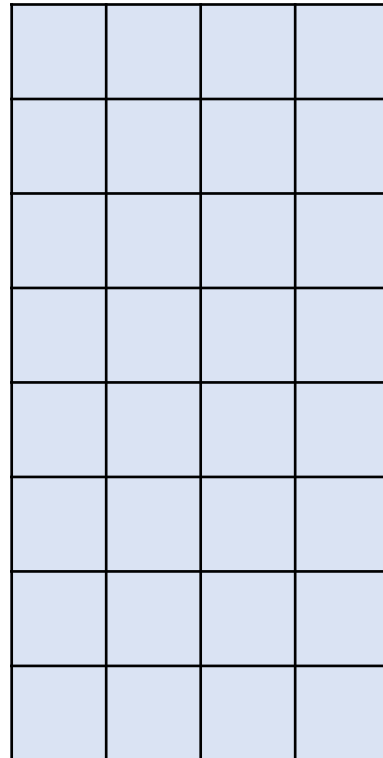




query



Reduce
Keys



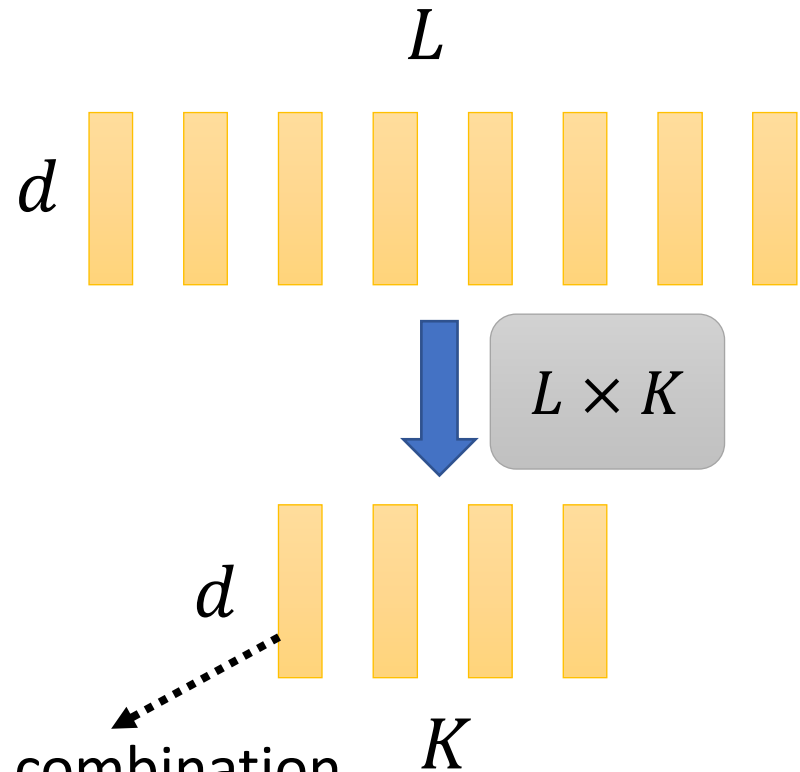
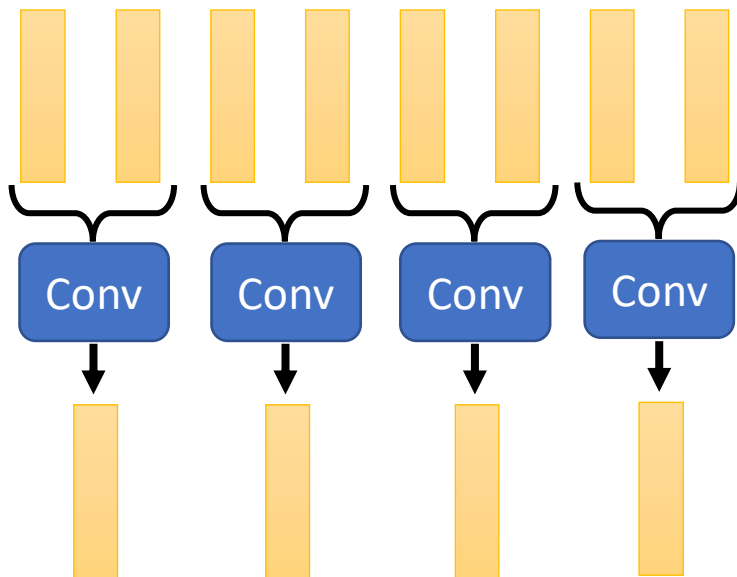
Reduce Number of Keys

Compressed Attention

<https://arxiv.org/abs/1801.10198>

Linformer

<https://arxiv.org/abs/2006.04768>

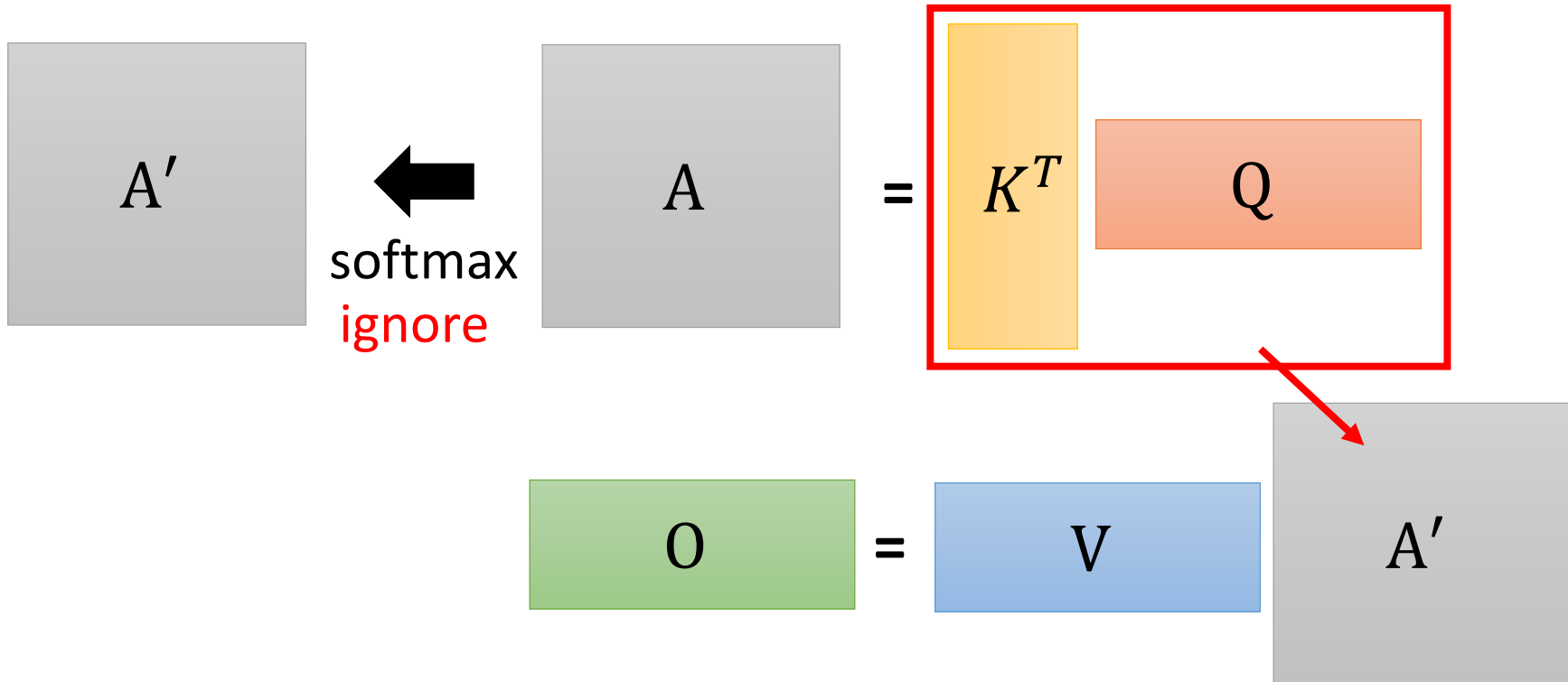


Linear combination
of L vectors

Attention Mechanism is three-matrix Multiplication

Review

$$\begin{array}{lcl} d \times N & \text{Q} & = W^q \text{I} \\ d \times N & \text{K} & = W^k \text{I} \\ d \times N & \text{V} & = W^v \text{I} \end{array}$$



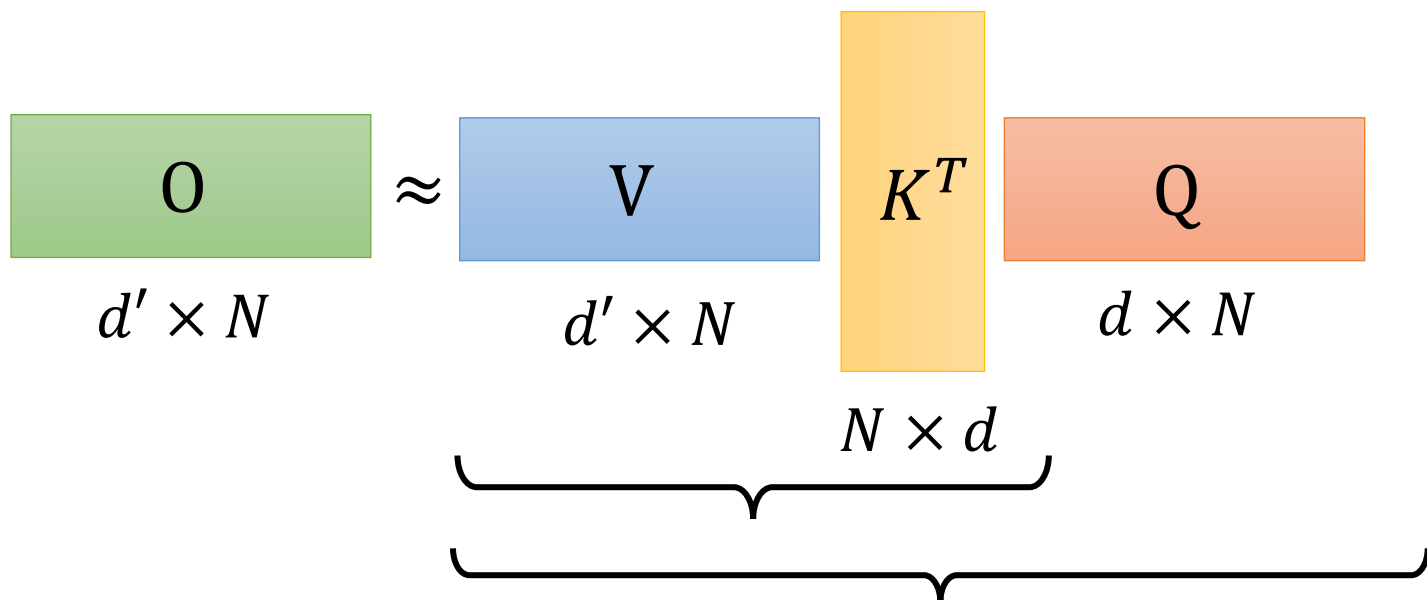
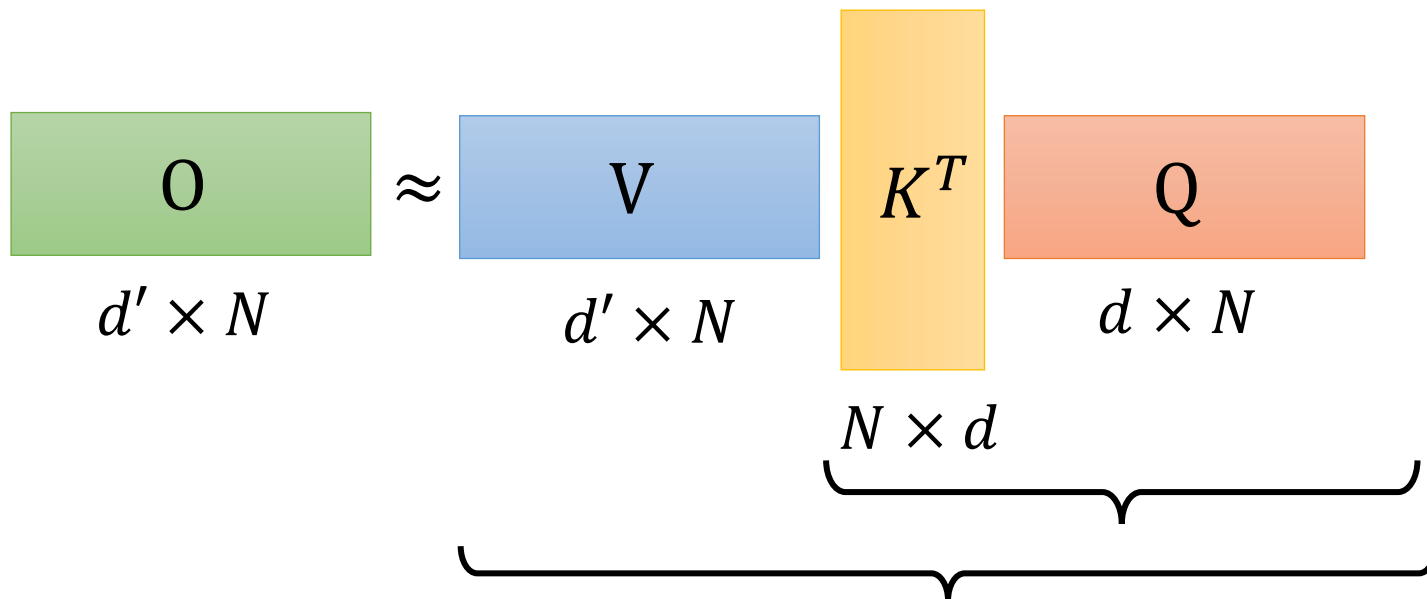
Attention Mechanism is three-matrix Multiplication

Review

$$\begin{array}{lcl} d \times N & \boxed{\text{Q}} & = \boxed{W^q} \boxed{\text{I}} \\ d \times N & \boxed{\text{K}} & = \boxed{W^k} \boxed{\text{I}} \\ d' \times N & \boxed{\text{V}} & = \boxed{W^v} \boxed{\text{I}} \end{array}$$

$$\begin{array}{ccccc} \boxed{\text{O}} & \approx & \boxed{\text{V}} & \boxed{K^T} & \boxed{\text{Q}} \\ d' \times N & & d' \times N & N \times d & d \times N \end{array}$$

⏟
⏟



Review Linear Algebra

Practical Issue

$k=1$ $m=1000$

$n=1$ $p=1000$

- Let A and B be $k \times m$ matrices, C be an $m \times n$ matrix, and P and Q be $n \times p$ matrices
 - $A(CP) = (AC)P$



$k \times m$ $m \times n$ $n \times p$



$k \times m$ $m \times p$

$m \times n \times p$

10^6

$k \times m \times p$

10^6



$k \times m$ $m \times n$ $n \times p$



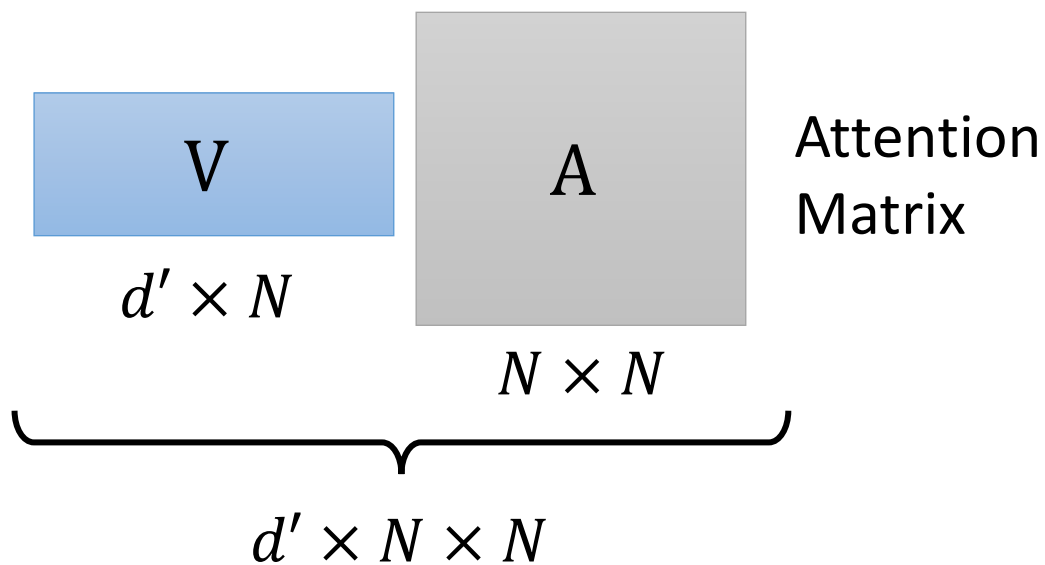
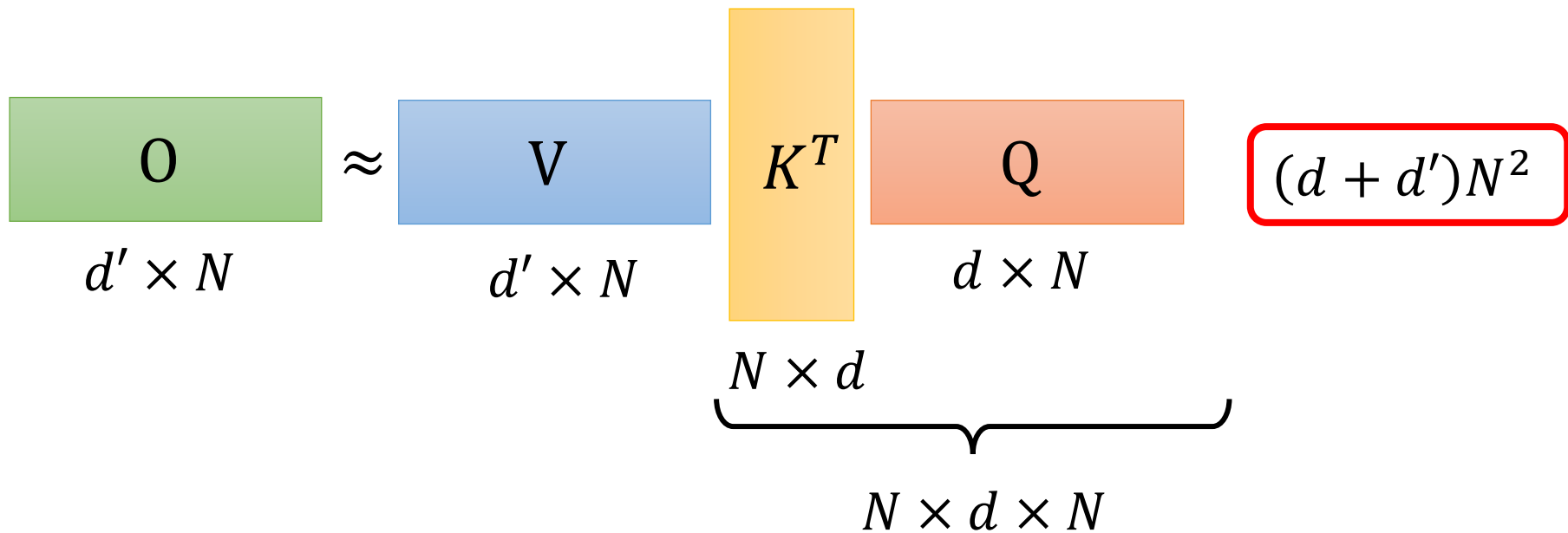
$k \times n$ $n \times p$

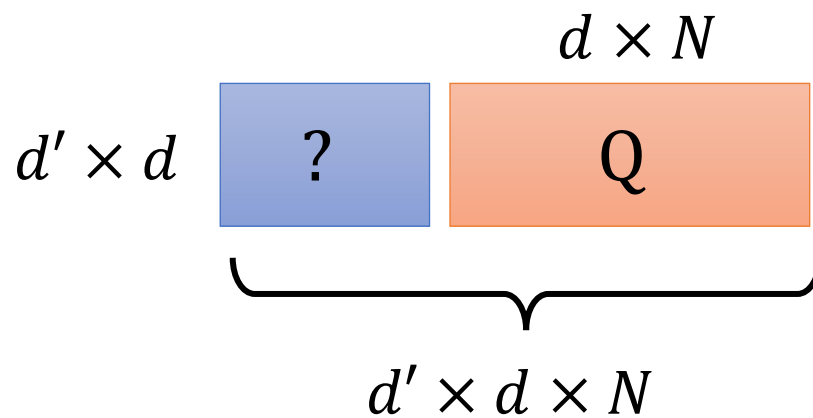
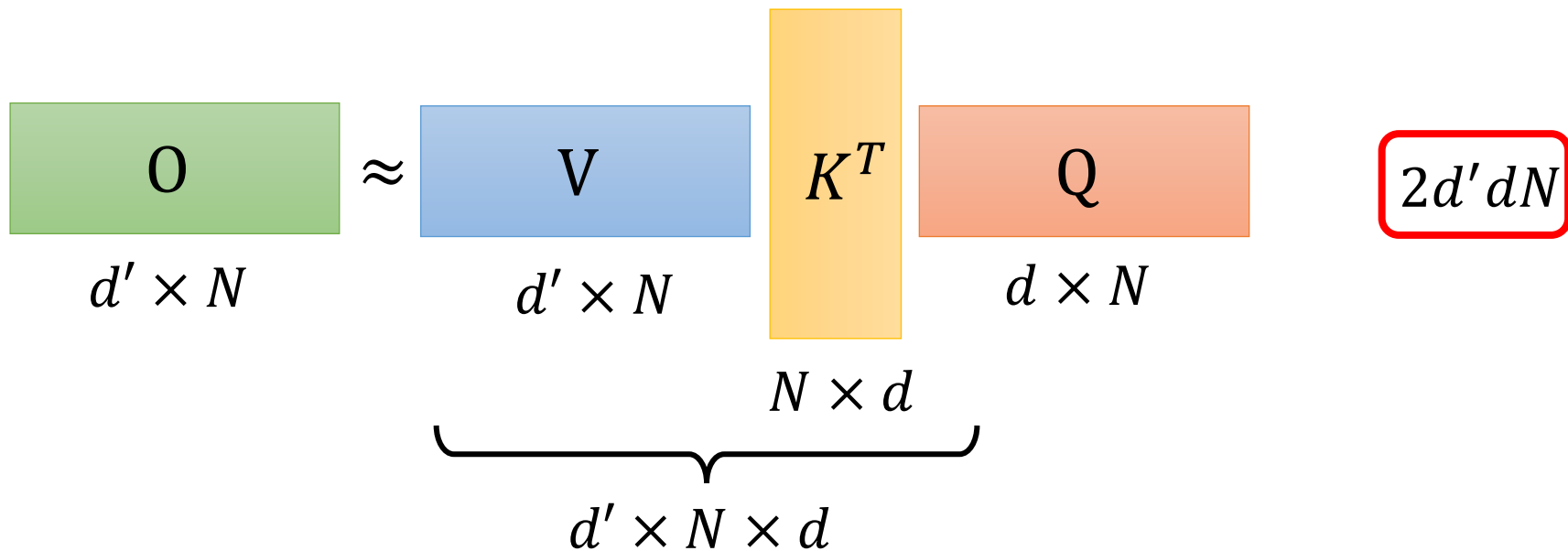
$k \times m \times n$

10^3

$k \times n \times p$

10^3

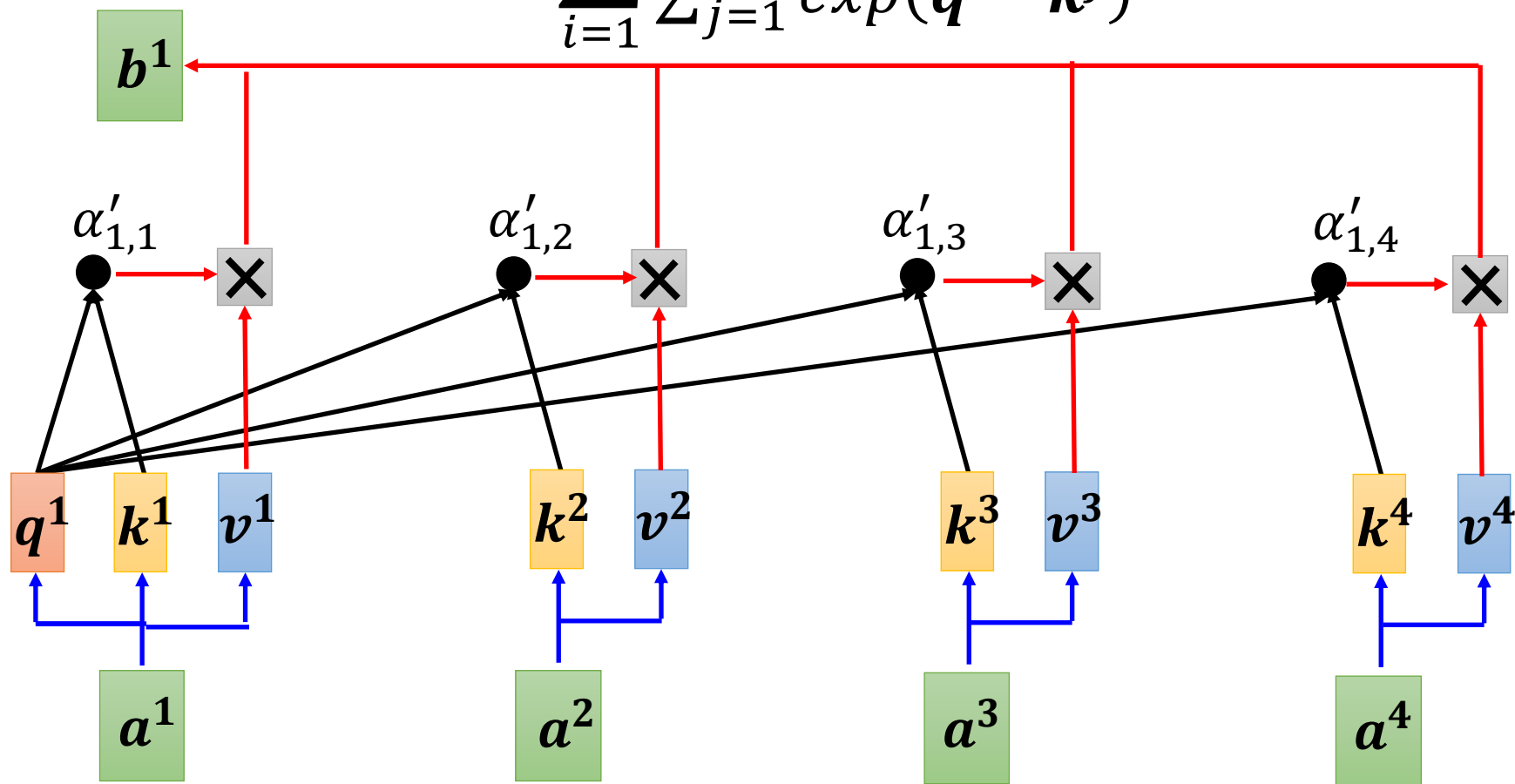




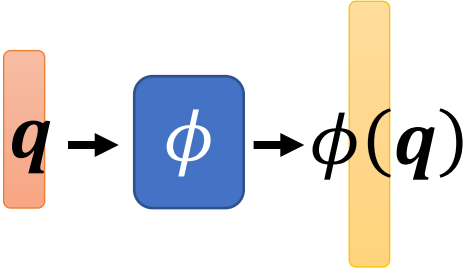
$$\begin{array}{ccccccc}
 \boxed{O} & \approx & \boxed{V} & \boxed{K^T} & \boxed{Q} & & \boxed{(d + d')N^2} \\
 d \times N & & d \times N & N \times d & d \times N & & \\
 & & & \underbrace{\hspace{10em}} & & & \\
 & & \underbrace{\hspace{15em}} & & & &
 \end{array}$$

$$\begin{array}{ccccccc}
 \boxed{O} & \approx & \boxed{V} & \boxed{K^T} & \boxed{Q} & & \boxed{2d'dN} \\
 d \times N & & d \times N & N \times d & d \times N & & \\
 & & & \underbrace{\hspace{10em}} & & & \\
 & & \underbrace{\hspace{15em}} & & & &
 \end{array}$$

$$\begin{aligned}
 \mathbf{b}^1 &= \sum_{i=1}^N \alpha'_{1,i} \mathbf{v}^i = \sum_{i=1}^N \frac{\exp(\alpha_{1,i})}{\sum_{j=1}^N \exp(\alpha_{1,j})} \mathbf{v}^i \\
 &= \sum_{i=1}^N \frac{\exp(\mathbf{q}^1 \cdot \mathbf{k}^i)}{\sum_{j=1}^N \exp(\mathbf{q}^1 \cdot \mathbf{k}^j)} \mathbf{v}^i
 \end{aligned}$$



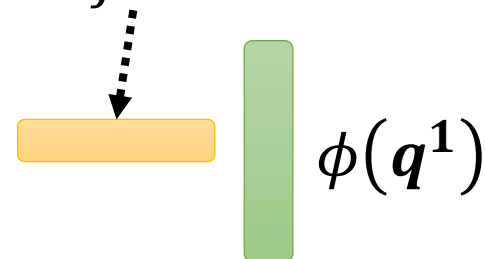
$$\mathbf{b}^1 = \sum_{i=1}^N \alpha'_{1,i} \mathbf{v}^i = \sum_{i=1}^N \frac{\exp(\alpha_{1,i})}{\sum_{j=1}^N \exp(\alpha_{1,j})} \mathbf{v}^i$$



$$\mathbf{q} \rightarrow \phi \rightarrow \phi(\mathbf{q}) = \sum_{i=1}^N \frac{\exp(\mathbf{q}^1 \cdot \mathbf{k}^i)}{\sum_{j=1}^N \exp(\mathbf{q}^1 \cdot \mathbf{k}^j)} \mathbf{v}^i$$

$$\exp(\mathbf{q} \cdot \mathbf{k}) \approx \phi(\mathbf{q}) \cdot \phi(\mathbf{k}) = \sum_{i=1}^N \frac{\phi(\mathbf{q}^1) \cdot \phi(\mathbf{k}^i)}{\sum_{j=1}^N \phi(\mathbf{q}^1) \cdot \phi(\mathbf{k}^j)} \mathbf{v}^i$$

$$= \frac{\sum_{i=1}^N [\phi(\mathbf{q}^1) \cdot \phi(\mathbf{k}^i)] \mathbf{v}^i}{\sum_{j=1}^N \phi(\mathbf{q}^1) \cdot \phi(\mathbf{k}^j)} = \frac{\sum_{i=1}^N [\phi(\mathbf{q}^1) \cdot \phi(\mathbf{k}^i)] \mathbf{v}^i}{\phi(\mathbf{q}^1) \cdot \sum_{j=1}^N \phi(\mathbf{k}^j)}$$



$$\mathbf{b}^1 = \sum_{i=1}^N \alpha'_{1,i} \mathbf{v}^i = \frac{\sum_{i=1}^N [\phi(\mathbf{q}^1) \cdot \phi(\mathbf{k}^i)] \mathbf{v}^i}{\phi(\mathbf{q}^1) \cdot \sum_{j=1}^N \phi(\mathbf{k}^j)}$$

$$\sum_{i=1}^N [\phi(\mathbf{q}^1) \cdot \phi(\mathbf{k}^i)] \mathbf{v}^i \quad \phi(\mathbf{q}^1) = \begin{bmatrix} q_1^1 \\ q_2^1 \\ \vdots \end{bmatrix} \quad \phi(\mathbf{k}^1) = \begin{bmatrix} k_1^1 \\ k_2^1 \\ \vdots \end{bmatrix}$$

$$= [\phi(\mathbf{q}^1) \cdot \phi(\mathbf{k}^1)] \mathbf{v}^1 + [\phi(\mathbf{q}^1) \cdot \phi(\mathbf{k}^2)] \mathbf{v}^2 + \dots$$

$$= (q_1^1 k_1^1 + q_2^1 k_2^1 + \dots) \mathbf{v}^1 + (q_1^1 k_1^2 + q_2^1 k_2^2 + \dots) \mathbf{v}^2 + \dots$$

$$= q_1^1 k_1^1 \mathbf{v}^1 + q_2^1 k_2^1 \mathbf{v}^1 + \dots + q_1^1 k_1^2 \mathbf{v}^2 + q_2^1 k_2^2 \mathbf{v}^2 + \dots + \dots$$

$$= q_1^1 (k_1^1 \mathbf{v}^1 + k_1^2 \mathbf{v}^2 + \dots) + q_2^1 (k_2^1 \mathbf{v}^1 + k_2^2 \mathbf{v}^2 + \dots)$$

$$\mathbf{b}^1 = \sum_{i=1}^N \alpha'_{1,i} \mathbf{v}^i = \frac{\sum_{i=1}^N [\phi(\mathbf{q}^1) \cdot \phi(\mathbf{k}^i)] \mathbf{v}^i}{\phi(\mathbf{q}^1) \cdot \sum_{j=1}^N \phi(\mathbf{k}^j)}$$

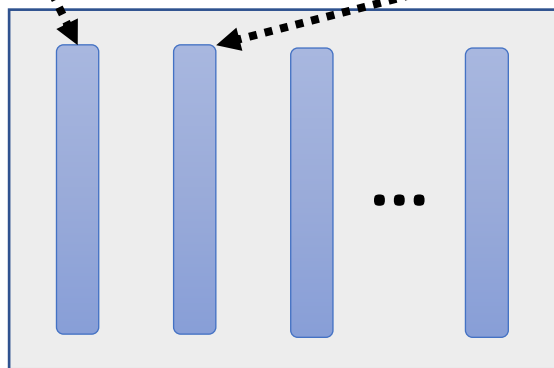
$$\sum_{i=1}^N [\phi(\mathbf{q}^1) \cdot \phi(\mathbf{k}^i)] \mathbf{v}^i$$

$$\phi(\mathbf{q}^1) = \begin{bmatrix} q_1^1 \\ q_2^1 \\ \vdots \end{bmatrix}$$

$$\phi(\mathbf{k}^1) = \begin{bmatrix} k_1^1 \\ k_2^1 \\ \vdots \end{bmatrix}$$

$$= q_1^1 (k_1^1 \mathbf{v}^1 + k_1^2 \mathbf{v}^2 + \dots) + q_2^1 (k_2^1 \mathbf{v}^1 + k_2^2 \mathbf{v}^2 + \dots)$$

$$\sum_{i=1}^N k_1^i \mathbf{v}^i$$

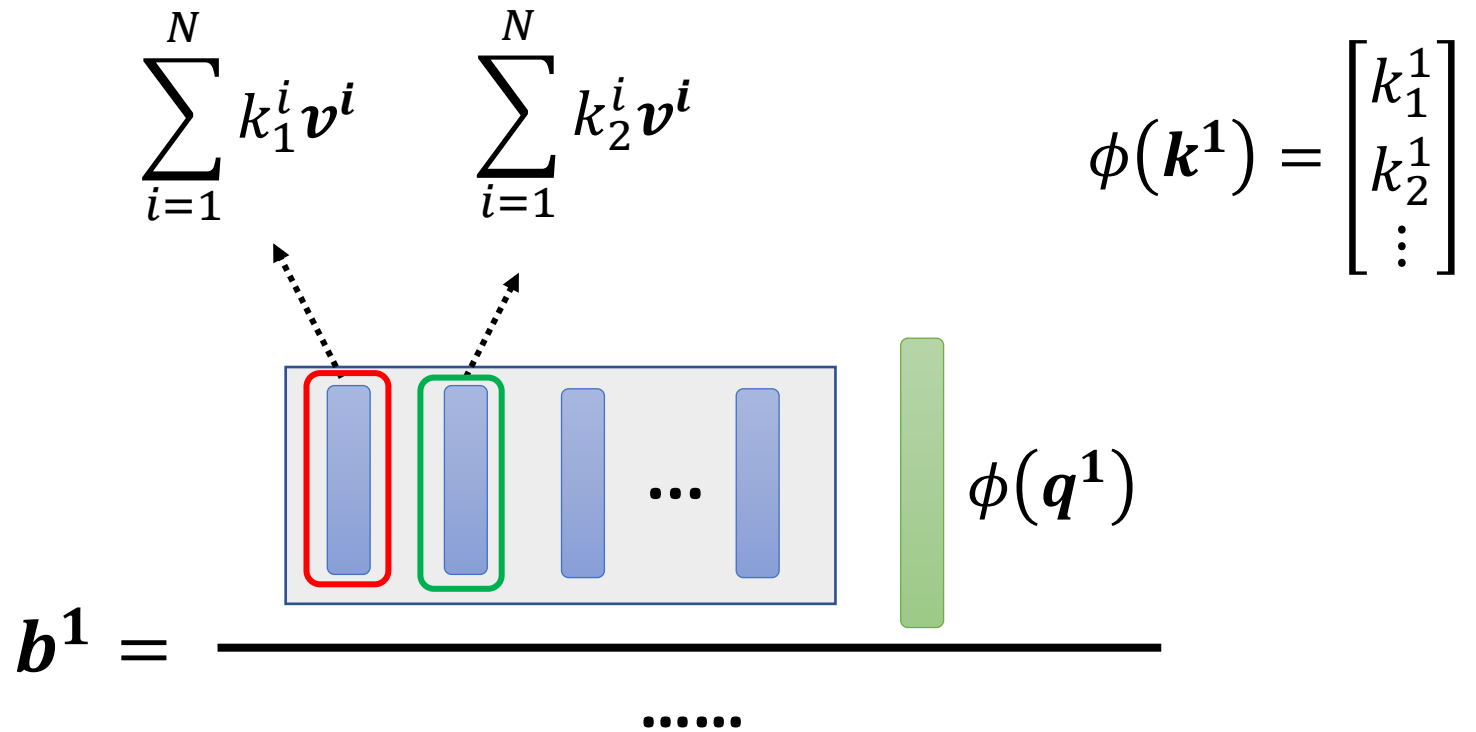
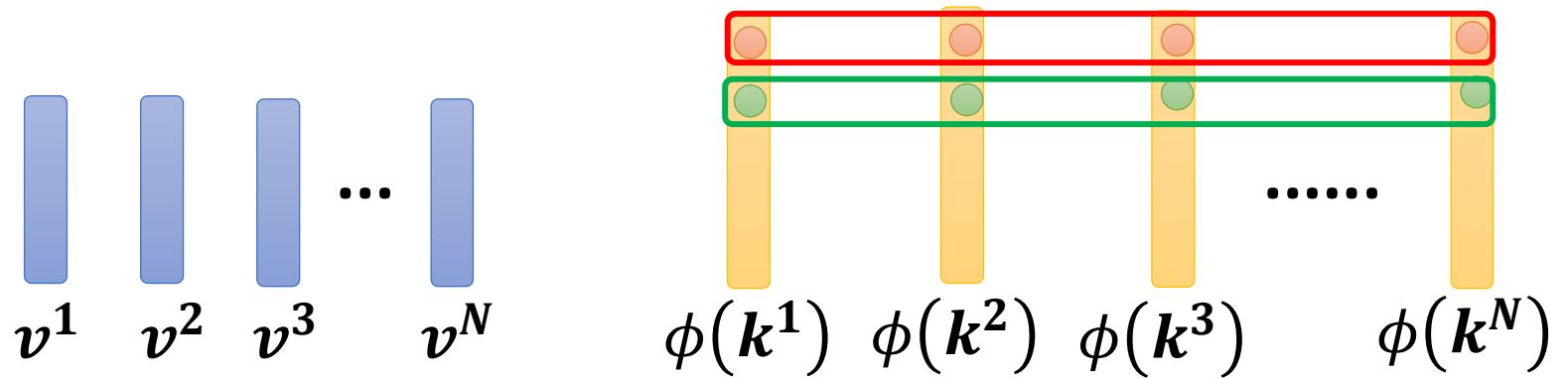


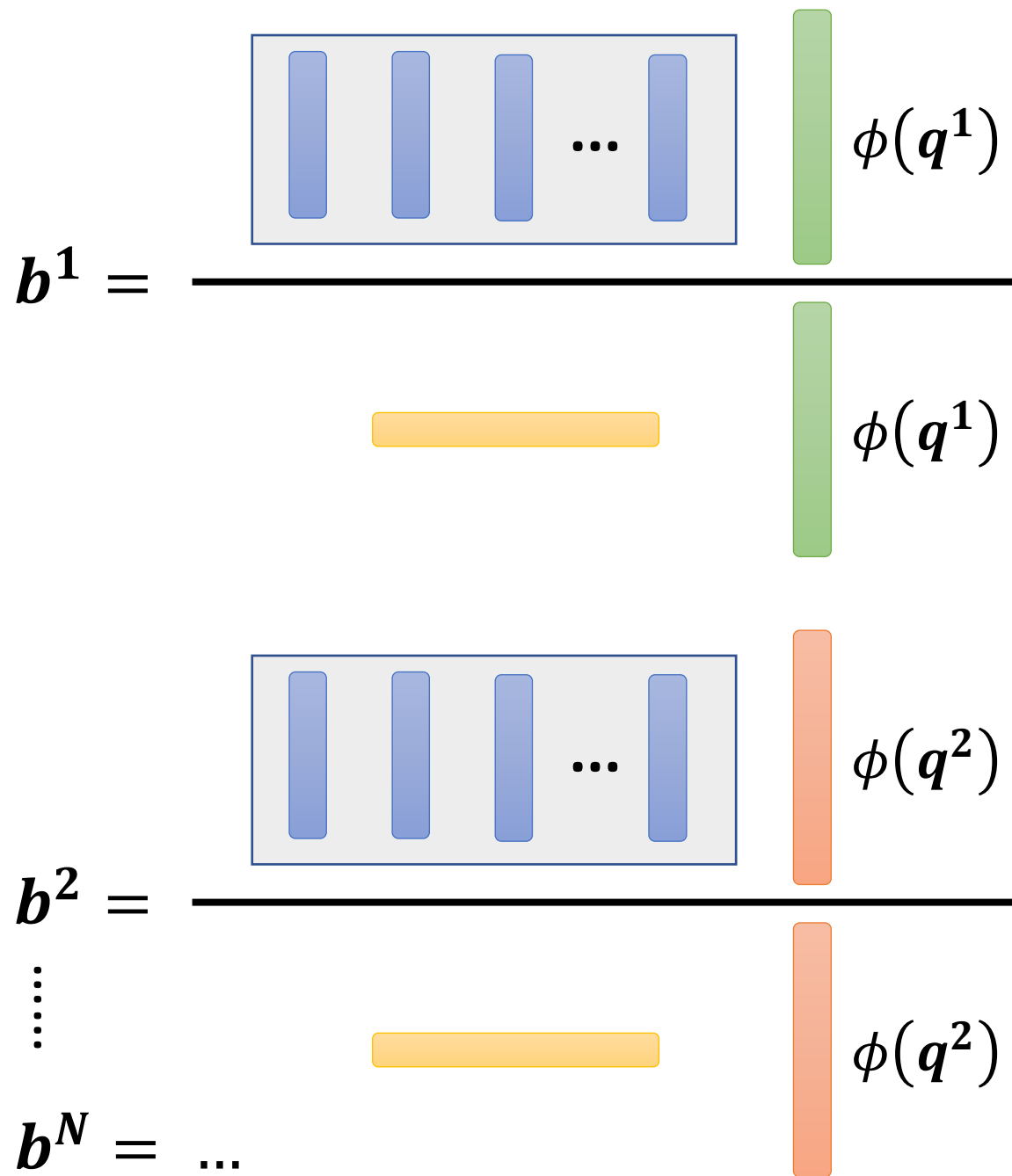
$$\phi(\mathbf{q}^1)$$



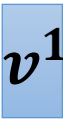



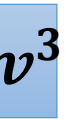

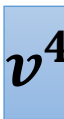
$$\sum_{i=1}^N k_2^i \mathbf{v}^i$$



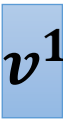



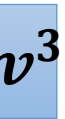

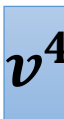
$$b^1 = \frac{\sum_{i=1}^N k_1^i \mathbf{v}^i \quad \sum_{i=1}^N k_2^i \mathbf{v}^i}{\sum_{j=1}^N \phi(\mathbf{k}^j)}$$

The diagram illustrates the calculation of b^1 . A horizontal line separates the numerator from the denominator. The numerator consists of two sums of weighted vectors, each with a dotted arrow pointing to a blue bar in a sequence of bars. The denominator is a sum of $\phi(\mathbf{k}^j)$, with an orange bar above it. To the right, two green bars represent $\phi(\mathbf{q}^1)$ above and below the line.



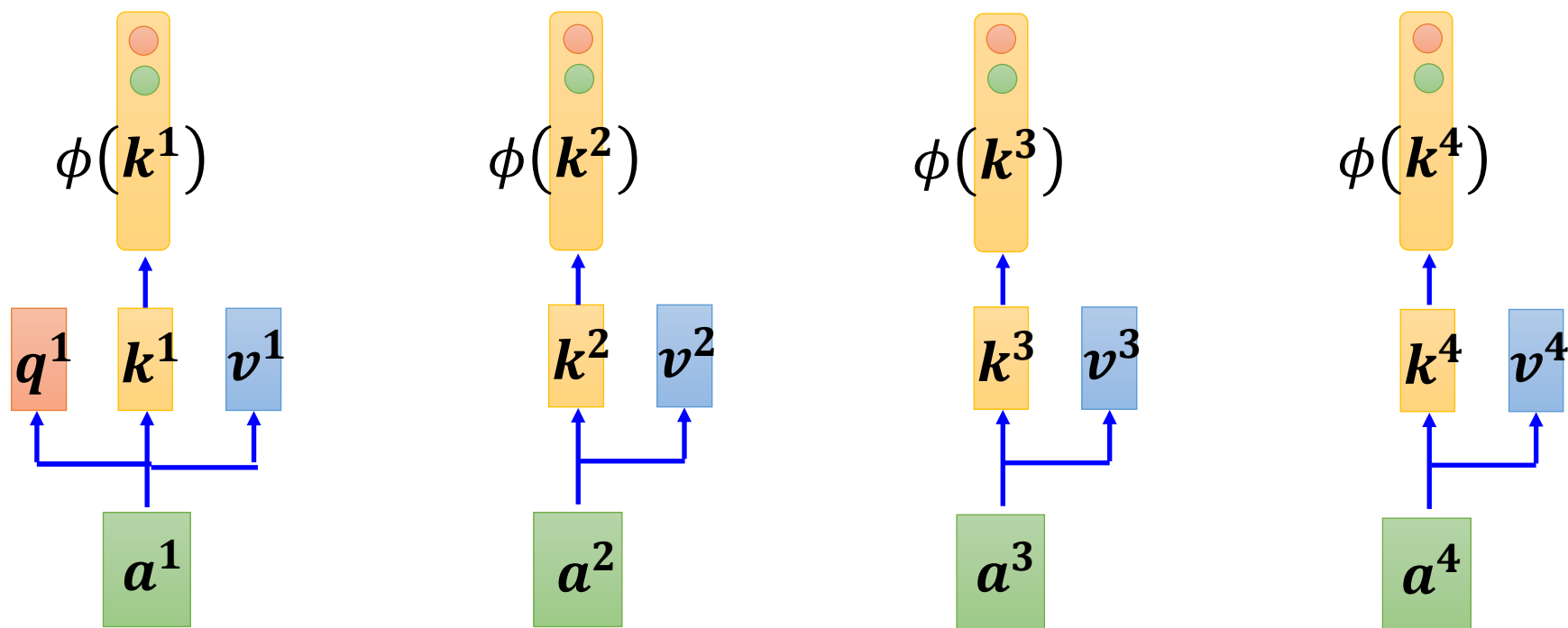


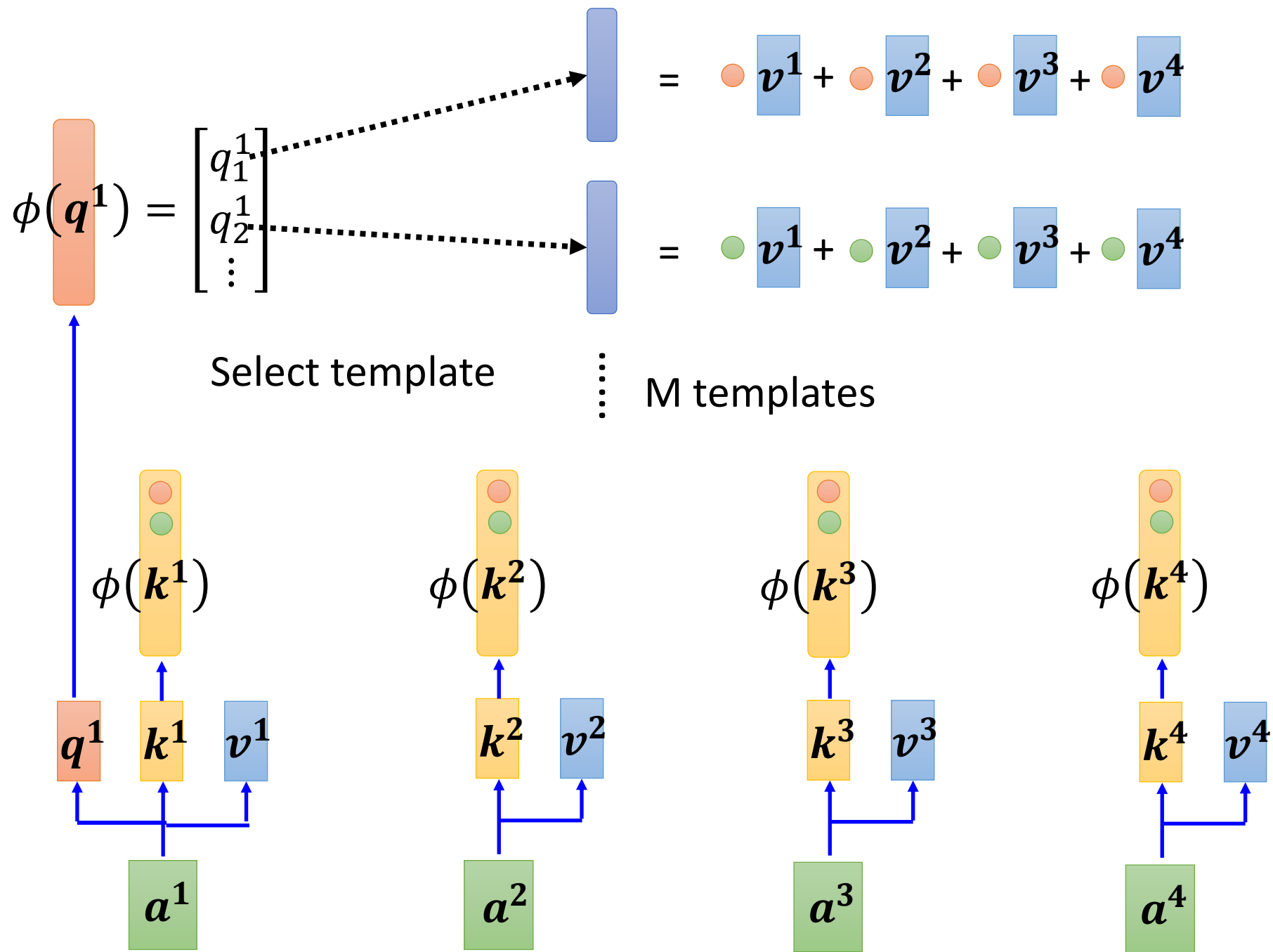
template  =   v^1 +   v^2 +   v^3 +   v^4

 =   v^1 +   v^2 +   v^3 +   v^4

M templates

M dimensions





Realization

- Efficient attention

<https://arxiv.org/pdf/1812.01243.pdf>

- Linear Transformer

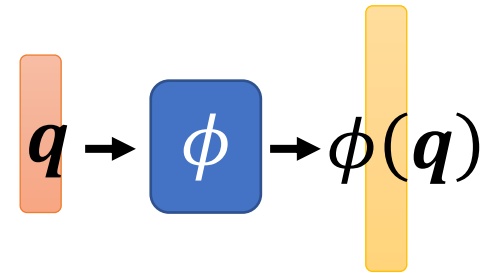
<https://linear-transformers.com/>

- Random Feature Attention

<https://arxiv.org/pdf/2103.02143.pdf>

- Performer

<https://arxiv.org/pdf/2009.14794.pdf>



$$\begin{aligned} \exp(\mathbf{q} \cdot \mathbf{k}) \\ \approx \phi(\mathbf{q}) \cdot \phi(\mathbf{k}) \end{aligned}$$