

# 各式各樣的 Attention

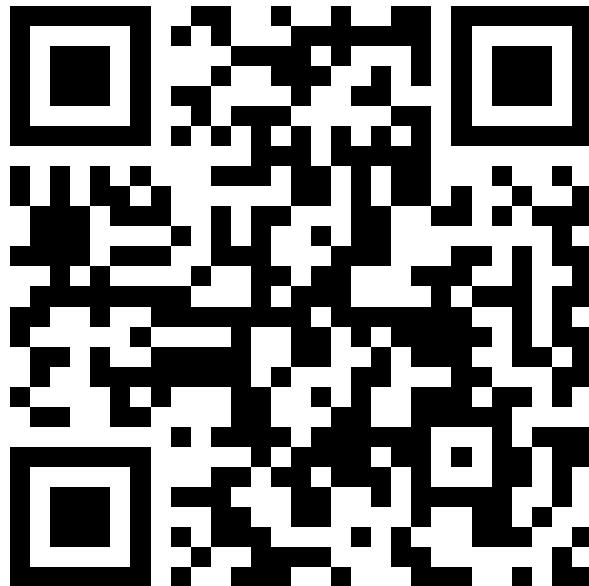
Hung-yi Lee 李宏毅

# Prerequisite



<https://youtu.be/hYdO9CscNes>

【機器學習2021】自注意力  
機制 (Self-attention) (上)



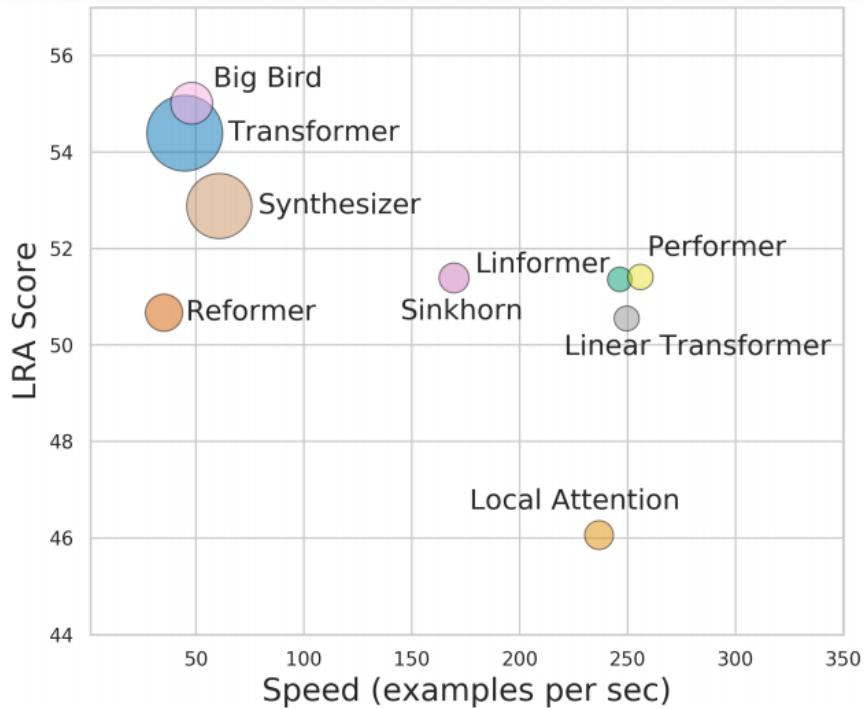
<https://youtu.be/gmsMY5kc-zw>

【機器學習2021】自注意力  
機制 (Self-attention) (下)

# To Learn More ...

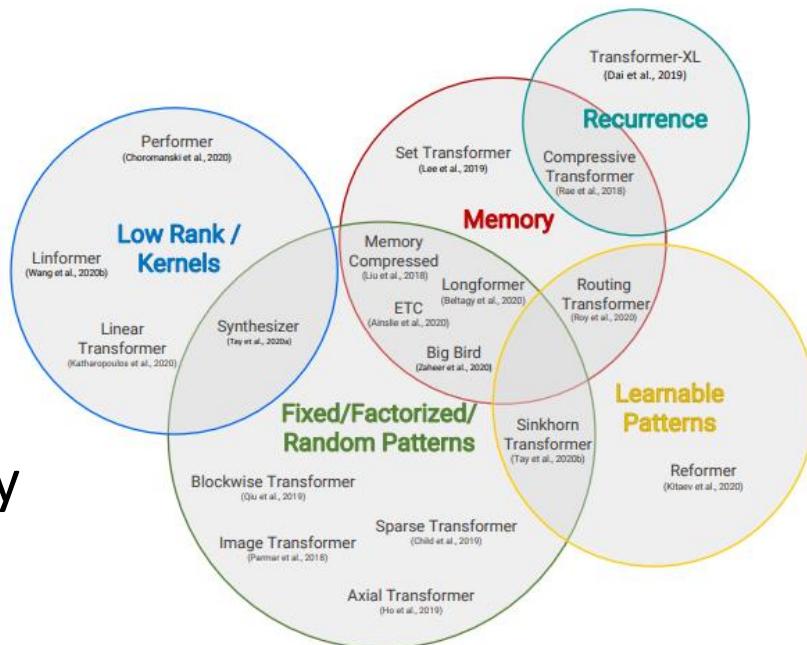
## Long Range Arena: A Benchmark for Efficient Transformers

<https://arxiv.org/abs/2011.04006>



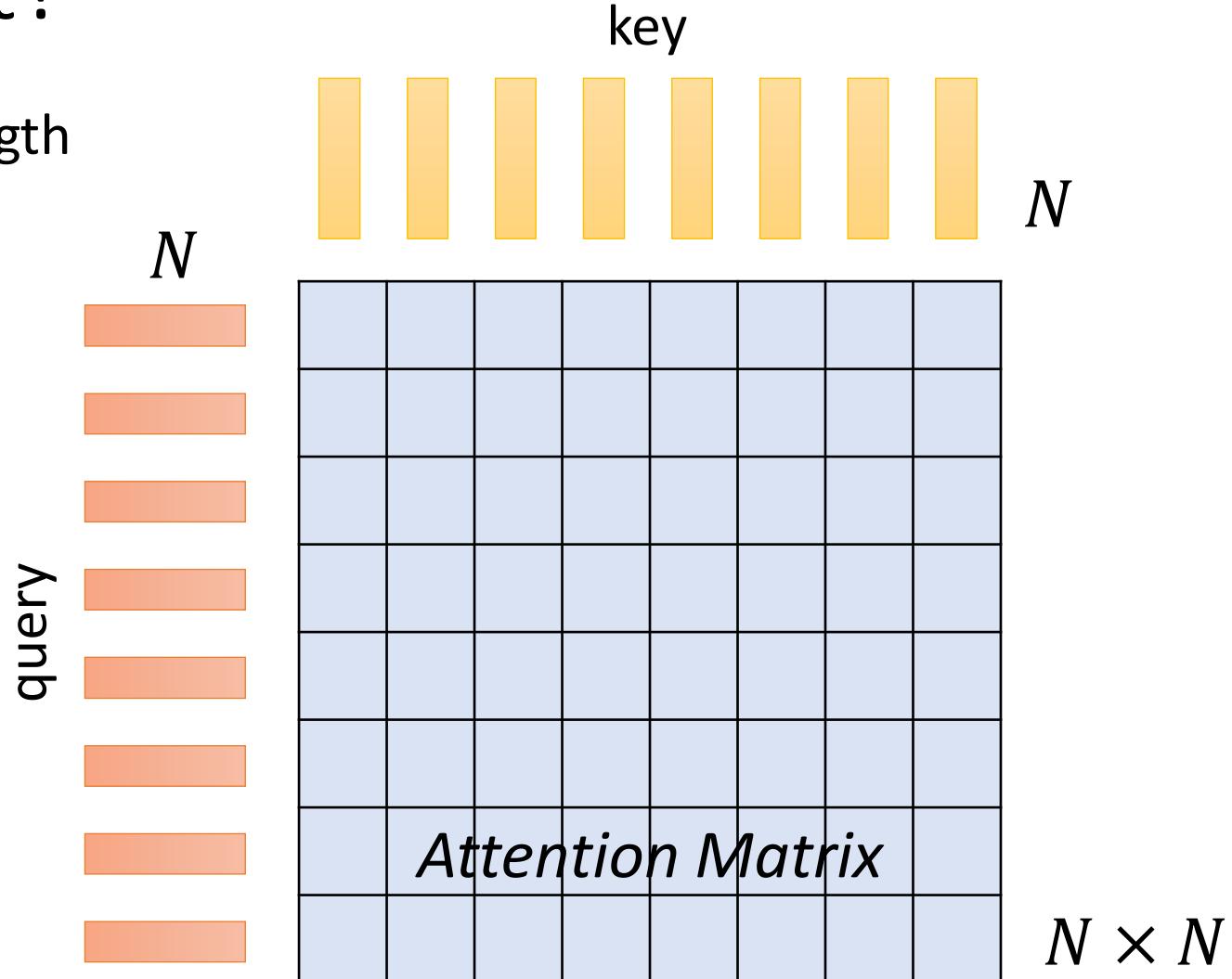
## Efficient Transformers: A Survey

<https://arxiv.org/abs/2009.06732>



# How to make self-attention efficient?

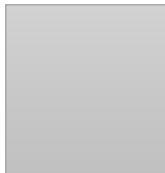
Sequence length  
 $= N$

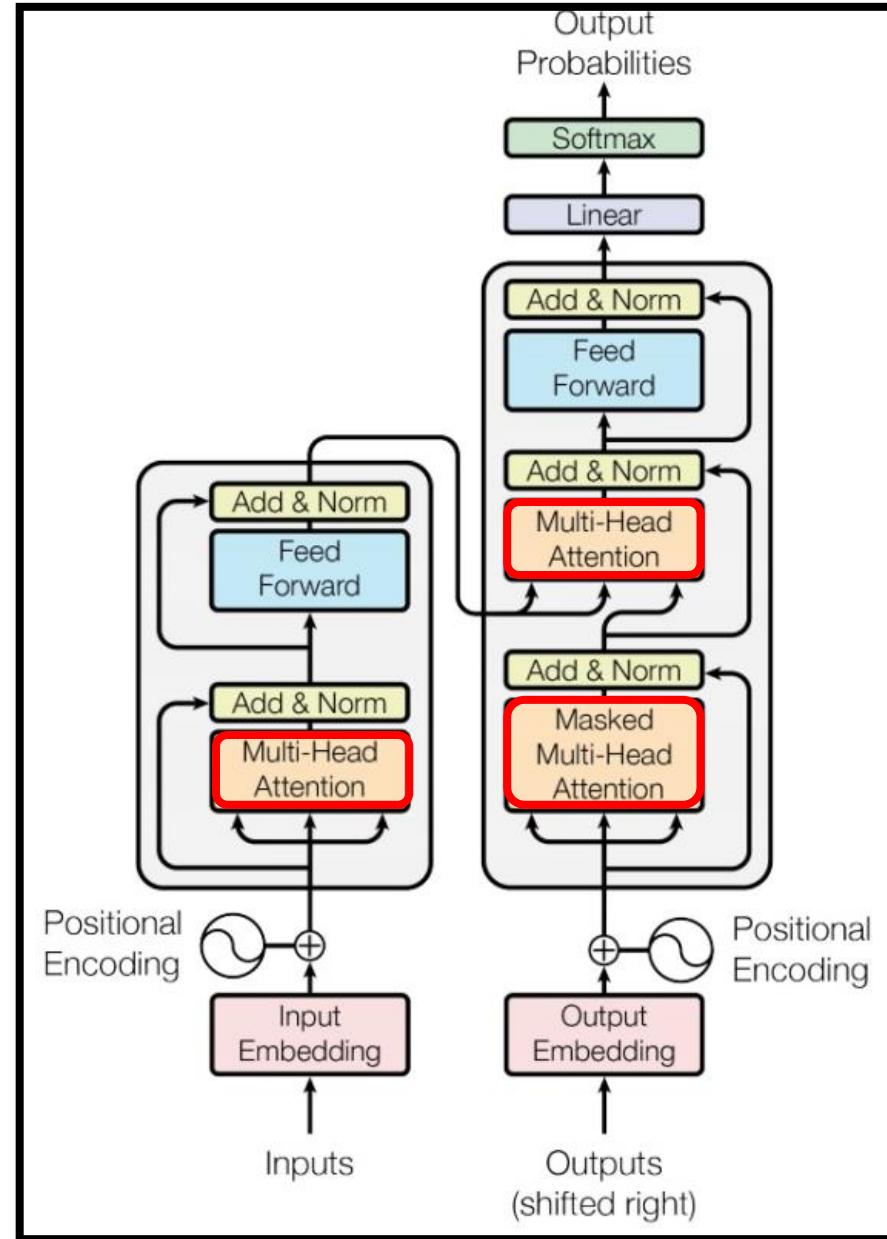


# Notice

- Self-attention is only a module in a larger network.
- Self-attention dominates computation when  $N$  is large.
- Usually developed for image processing

$$N = 256 * 256$$

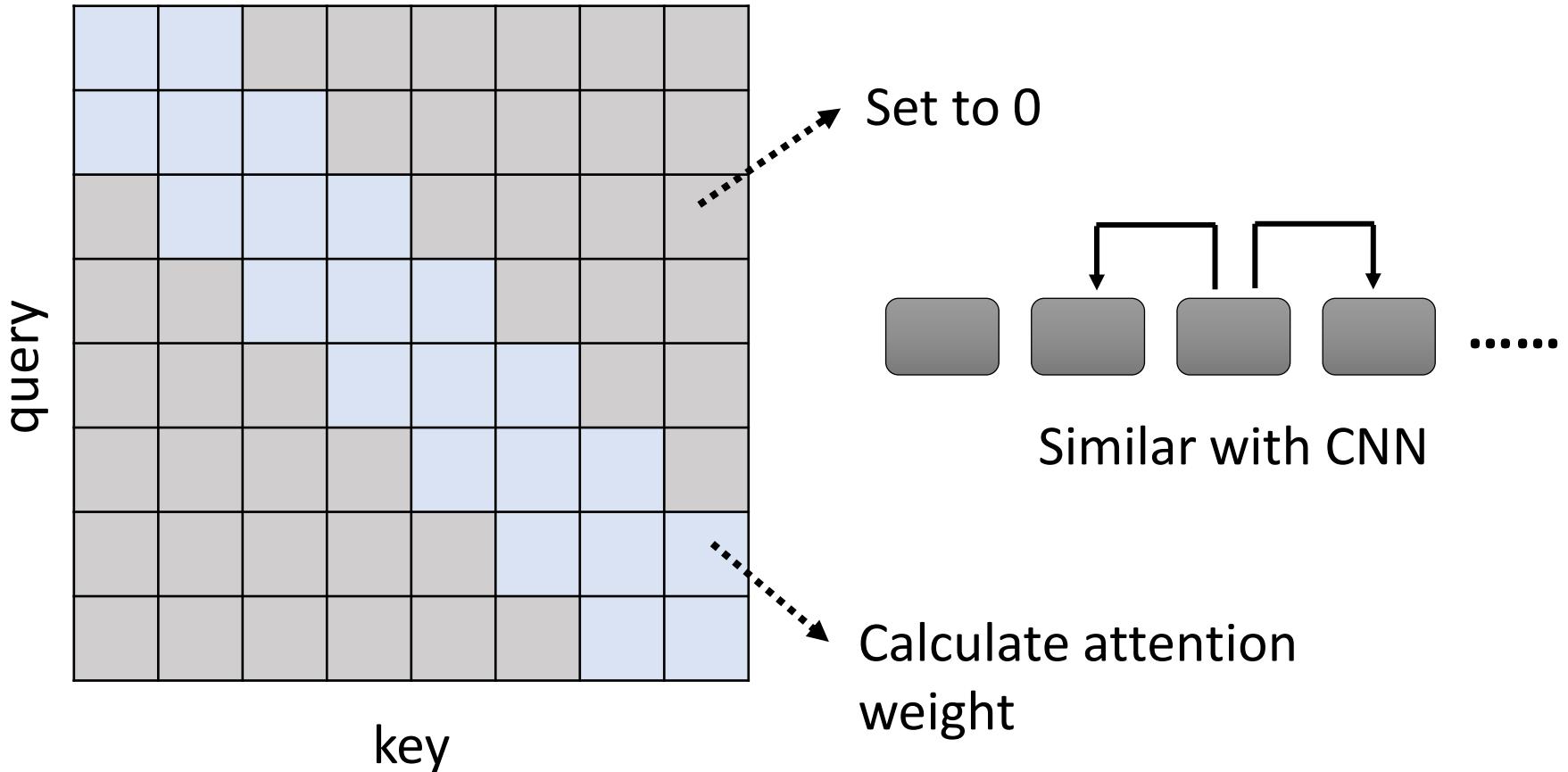
256        
256



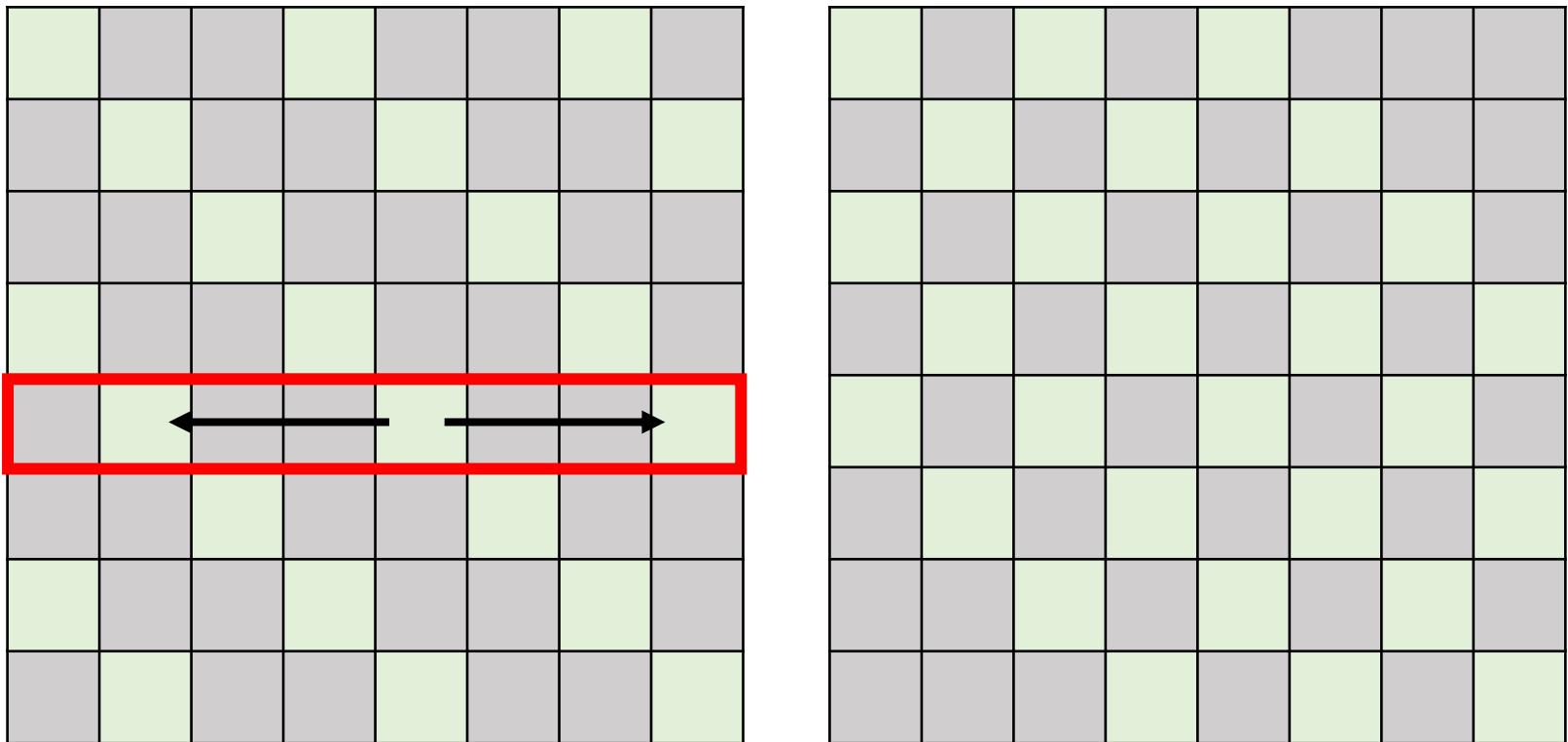
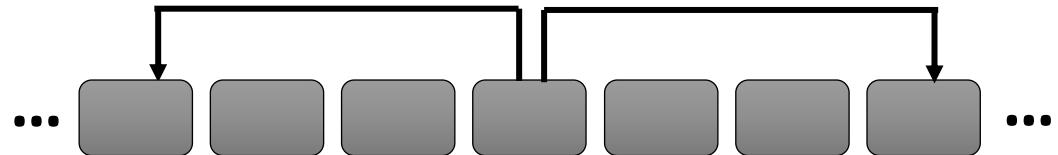
# Skip Some Calculations with Human Knowledge


Can we fill in some values with human knowledge?

# Local Attention / Truncated Attention



# Stride Attention

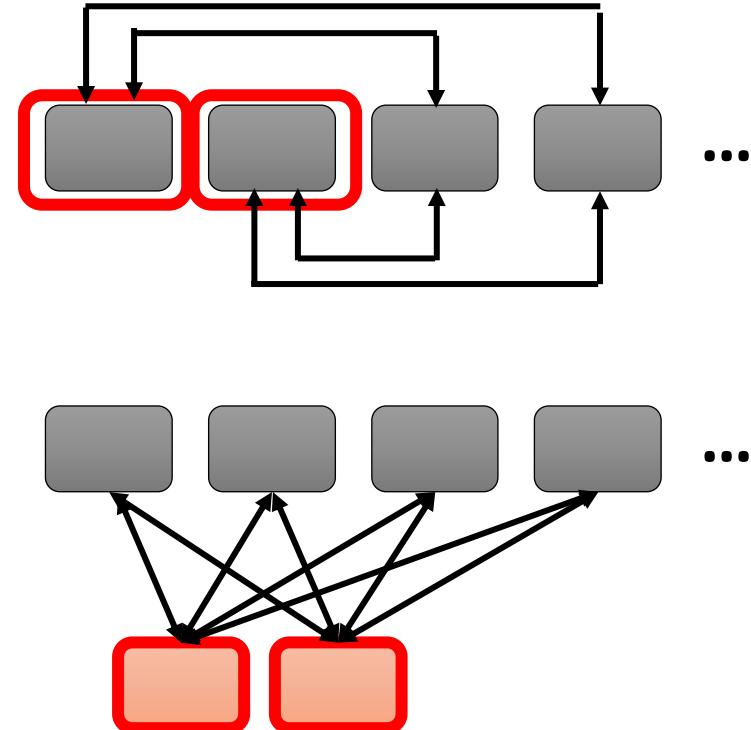
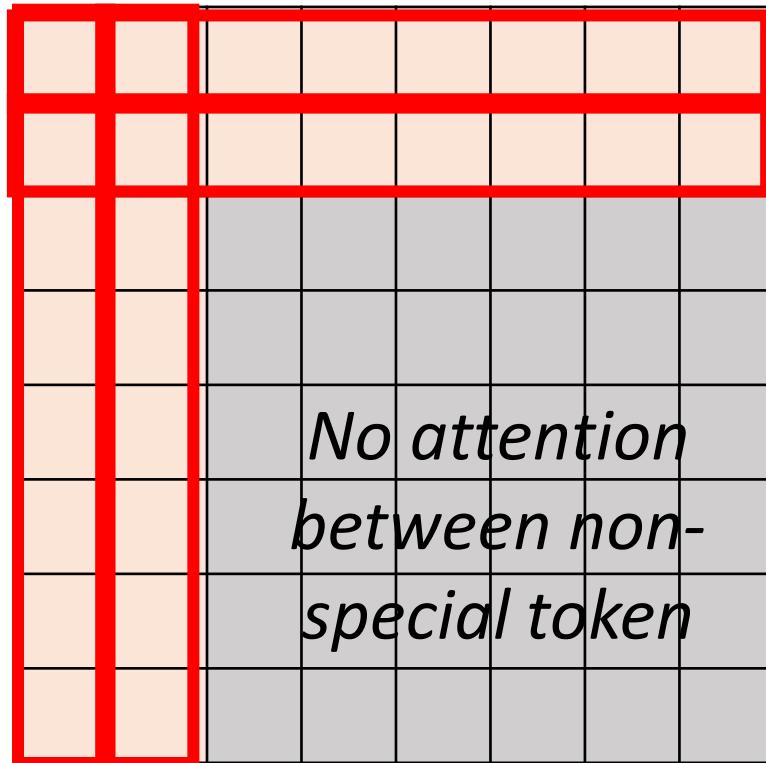


# Global Attention

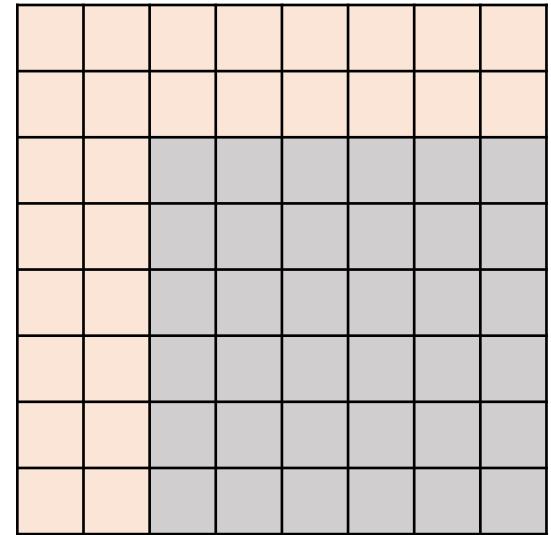
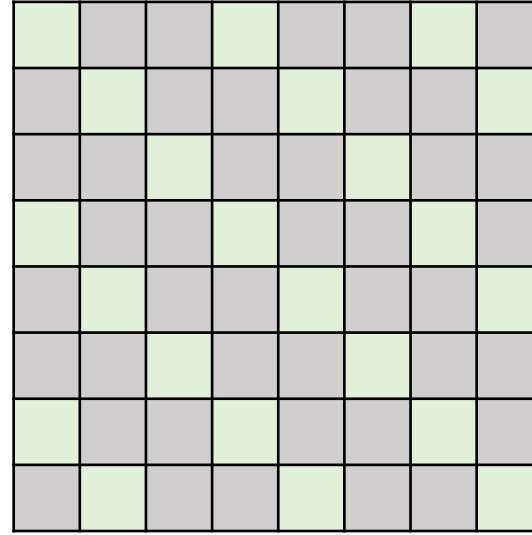
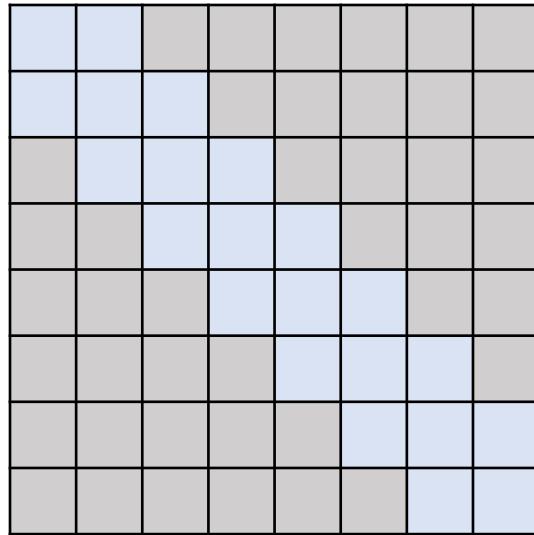
special token = “token中的里長伯”

Add special token into original sequence

- Attend to every token → collect global information
- Attended by every token → it knows global information



# Many Different Choices ...



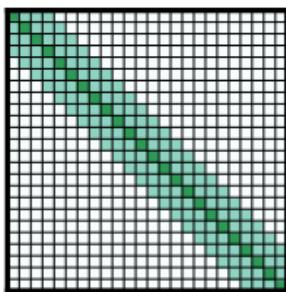
小孩子才做選擇 . . .

**Different heads use different patterns.**

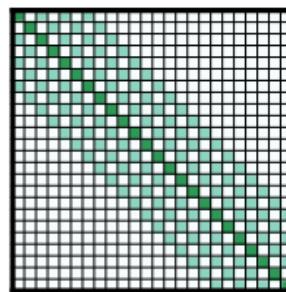
# Many Different Choices ...

- Longformer

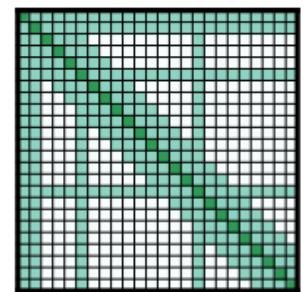
<https://arxiv.org/abs/2004.05150>



(b) Sliding window attention



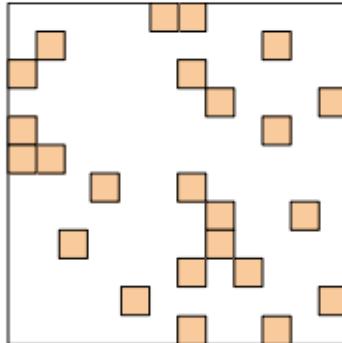
(c) Dilated sliding window



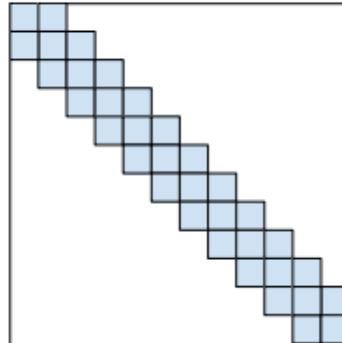
(d) Global+sliding window

- Big Bird

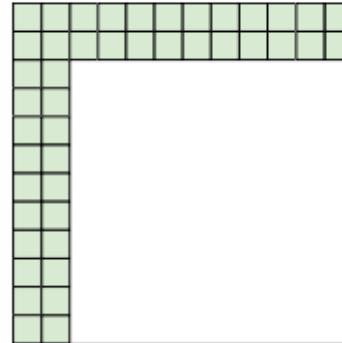
<https://arxiv.org/abs/2007.14062>



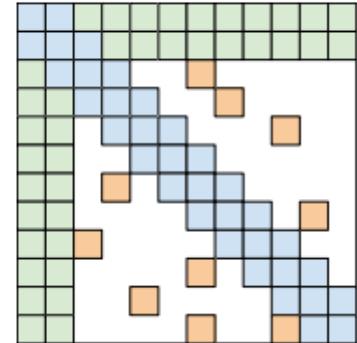
(a) Random attention



(b) Window attention

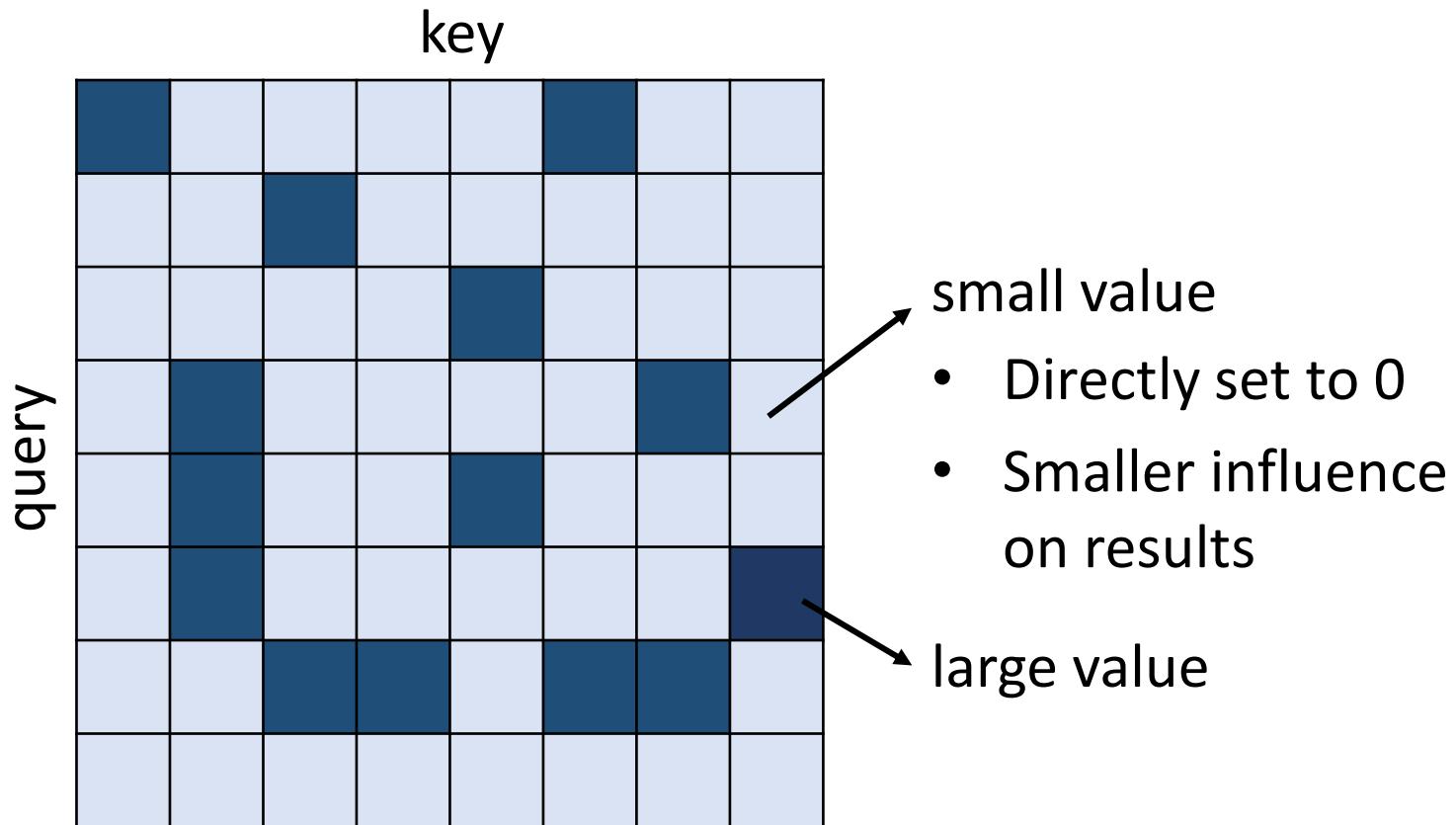


(c) Global Attention



(d) BIGBIRD

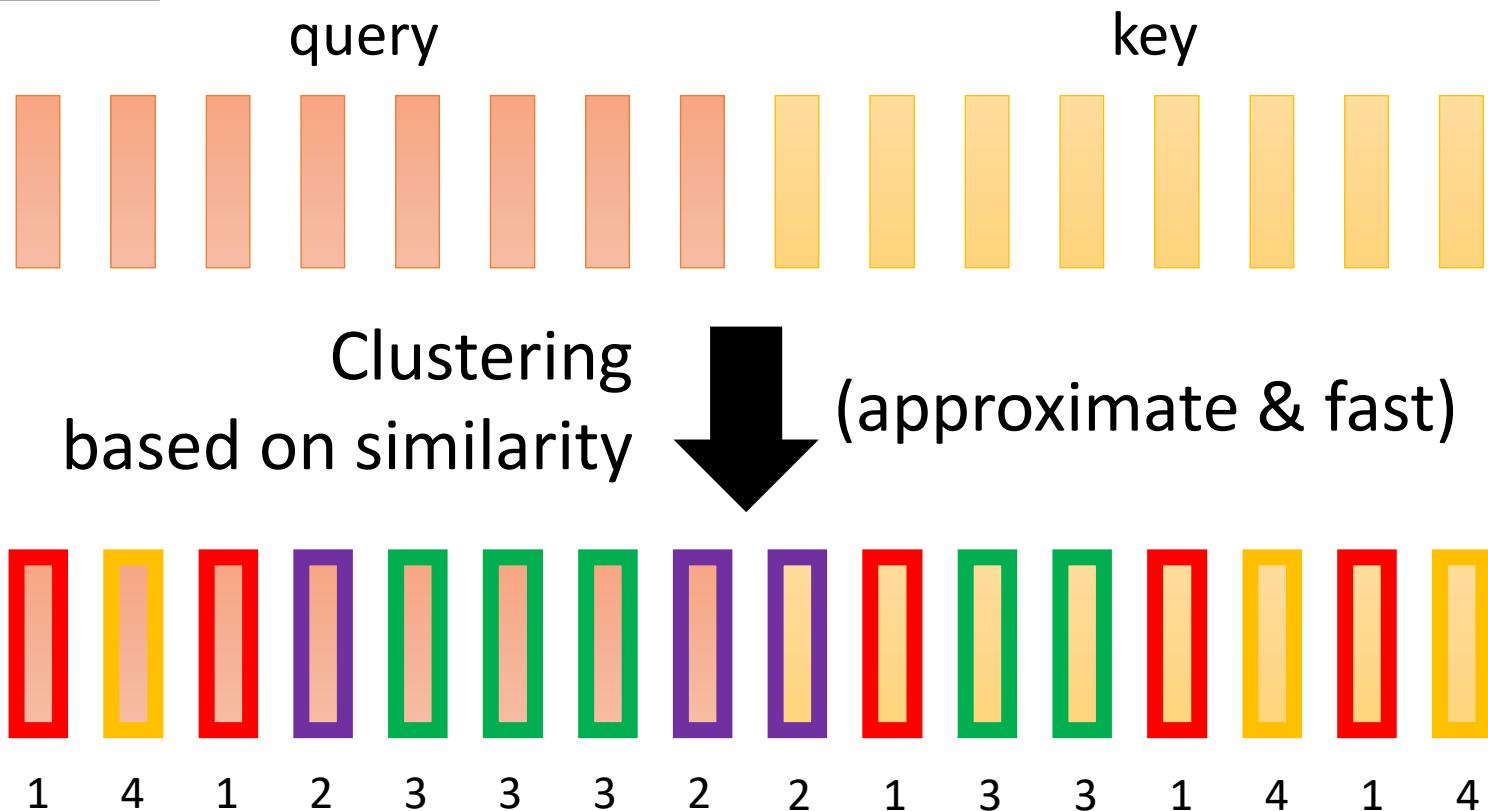
# Can we only focus on Critical Parts?



How to quickly estimate the portion with small attention weights?

# Clustering

## Step 1



Reformer

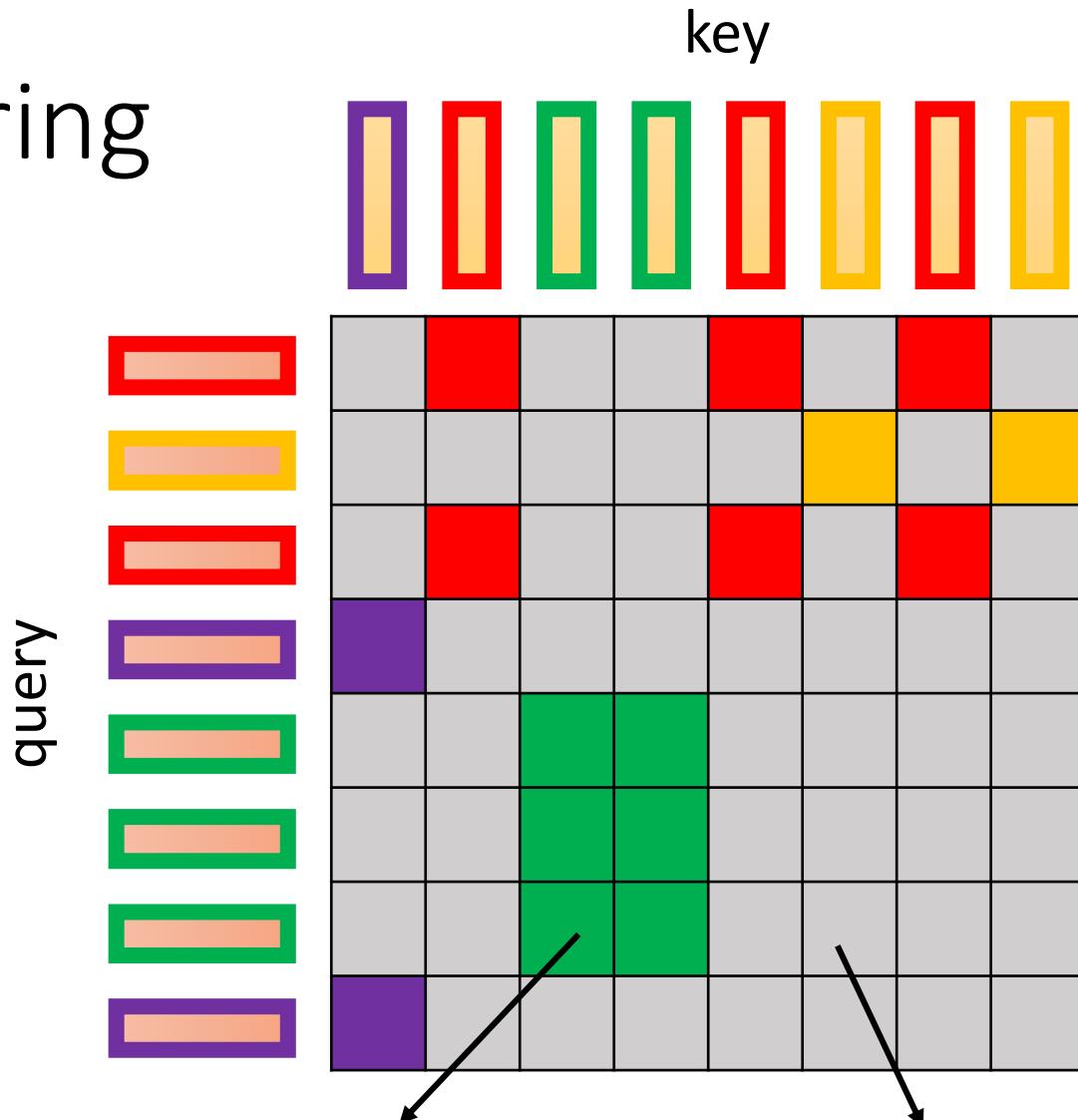
<https://openreview.net/forum?id=rkgNKkHtvB>

Routing Transformer

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

# Clustering

Step 2

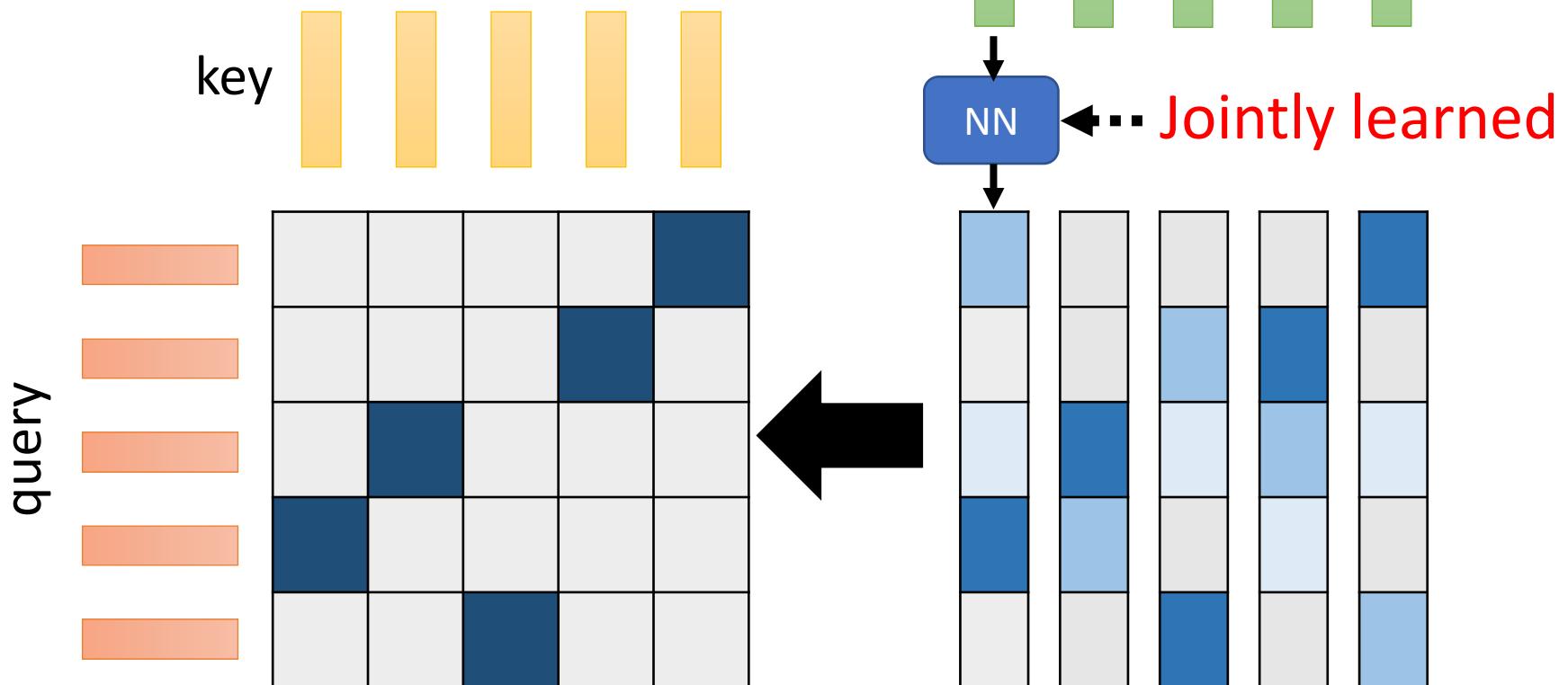


Belong to the same cluster, then  
calculate attention weight

Not the same cluster,  
set to 0

# Learnable Patterns

## Sinkhorn Sorting Network



A grid should be skipped or not is  
decided by another learned module

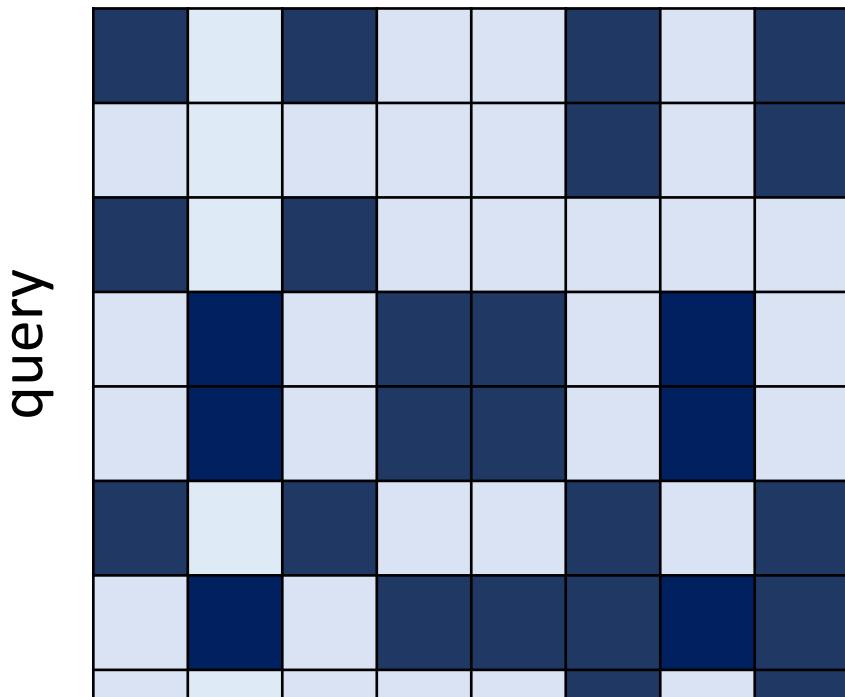
(simplified version)

# Do we need full attention matrix?

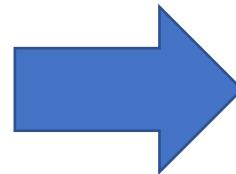
Linformer

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

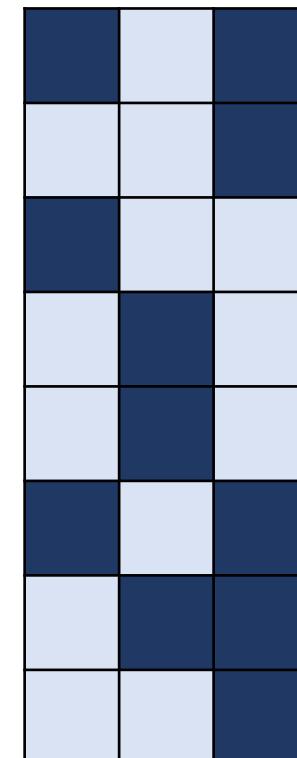
Many redundant columns

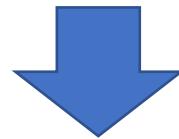
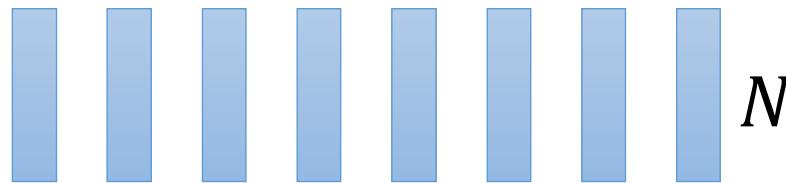


key

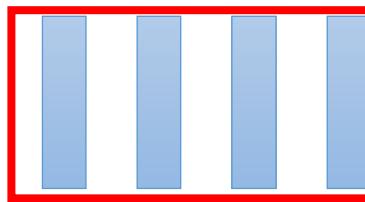


**Low Rank**





$N$   
value

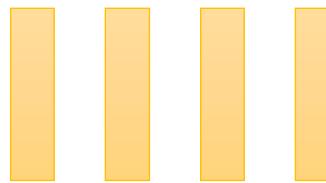


$K$



$K$

$N$   
key



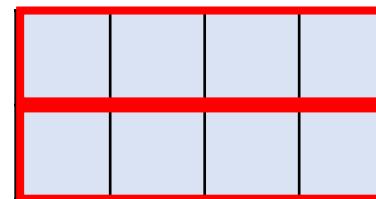
Representative  
keys

Can we reduce  
the number of  
queries?

query



change output  
sequence length



output

output

output

output

output

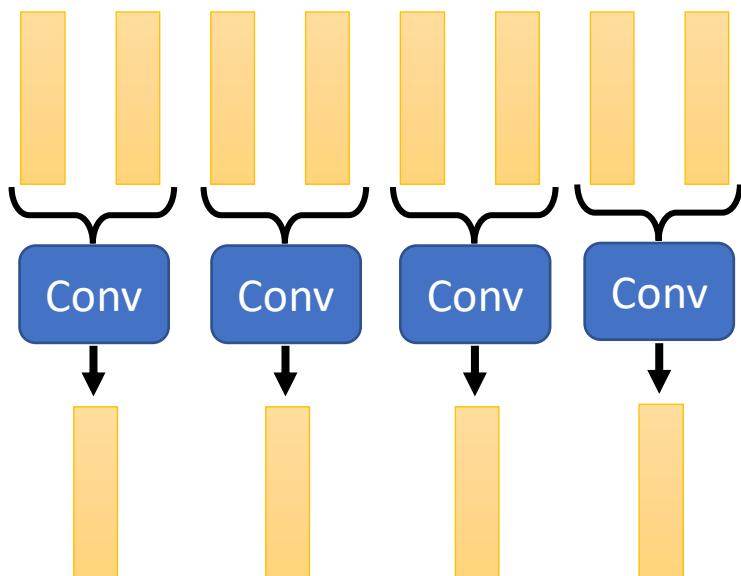
output

output

# Reduce Number of Keys

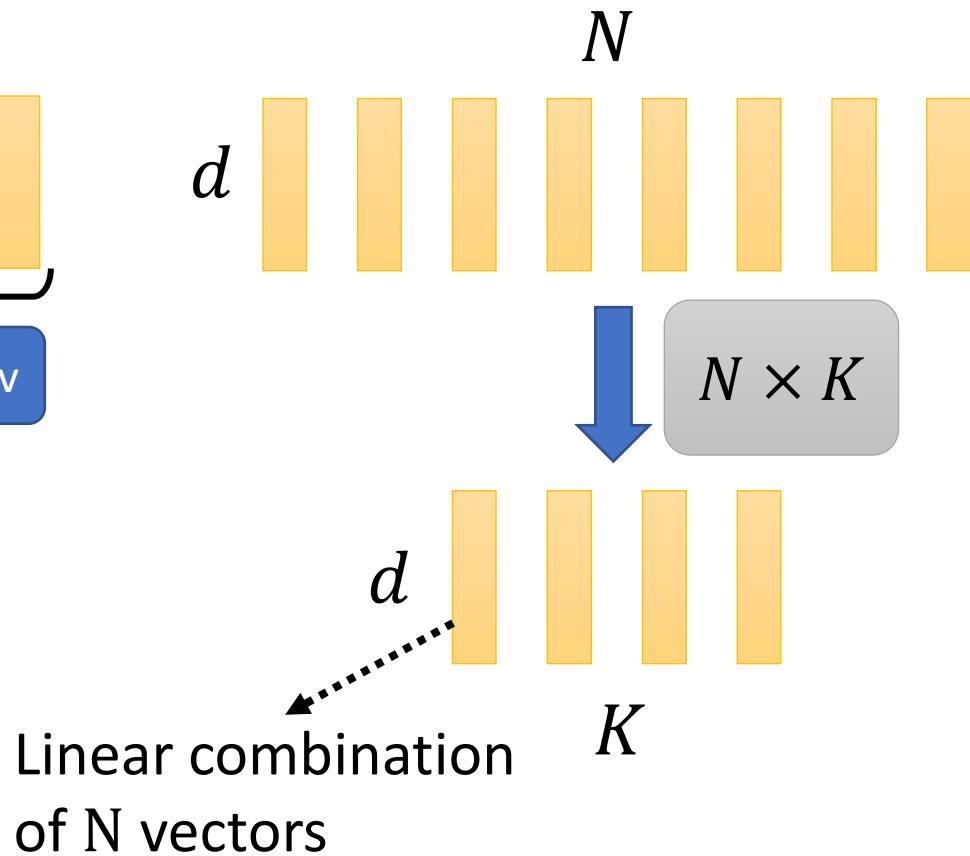
## Compressed Attention

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



## Linformer

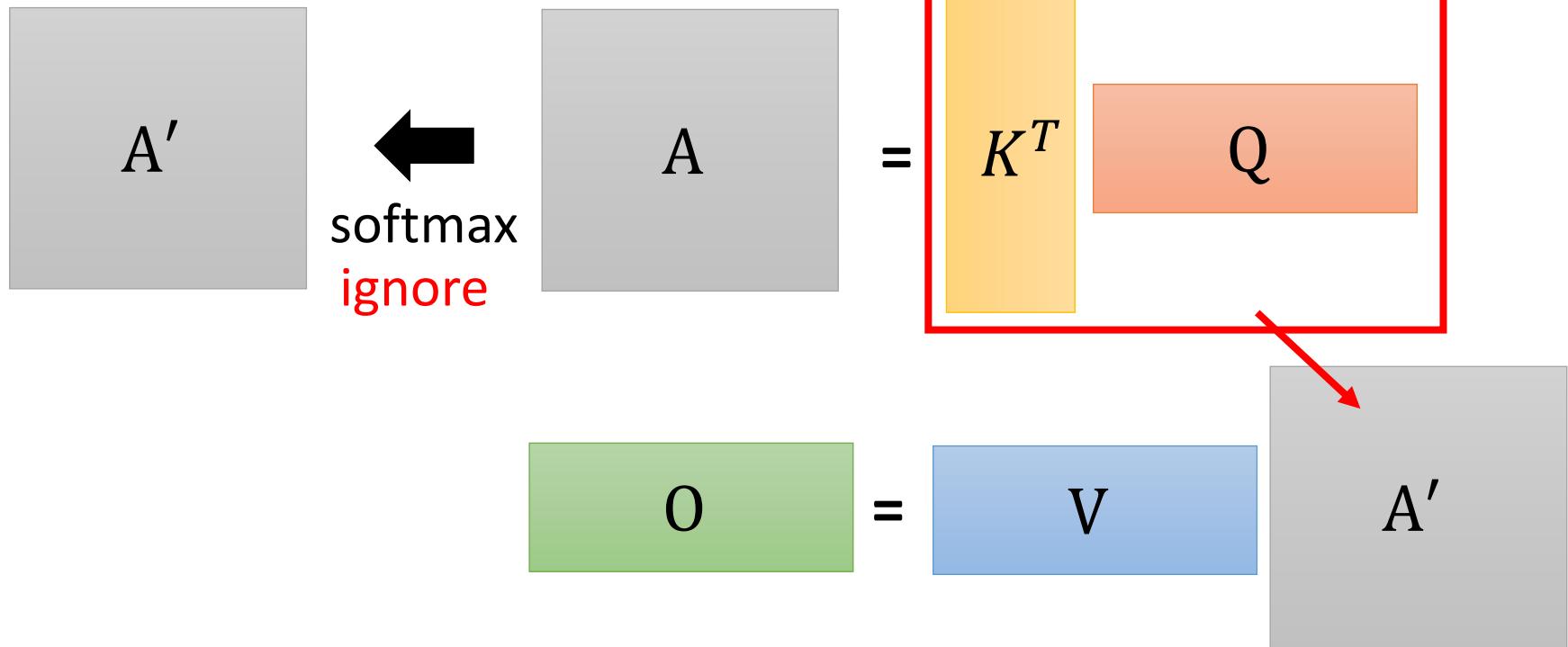
<https://arxiv.org/abs/2006.04768>



# **Attention Mechanism is three-matrix Multiplication**

Review

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



# ***Attention Mechanism is three-matrix Multiplication***

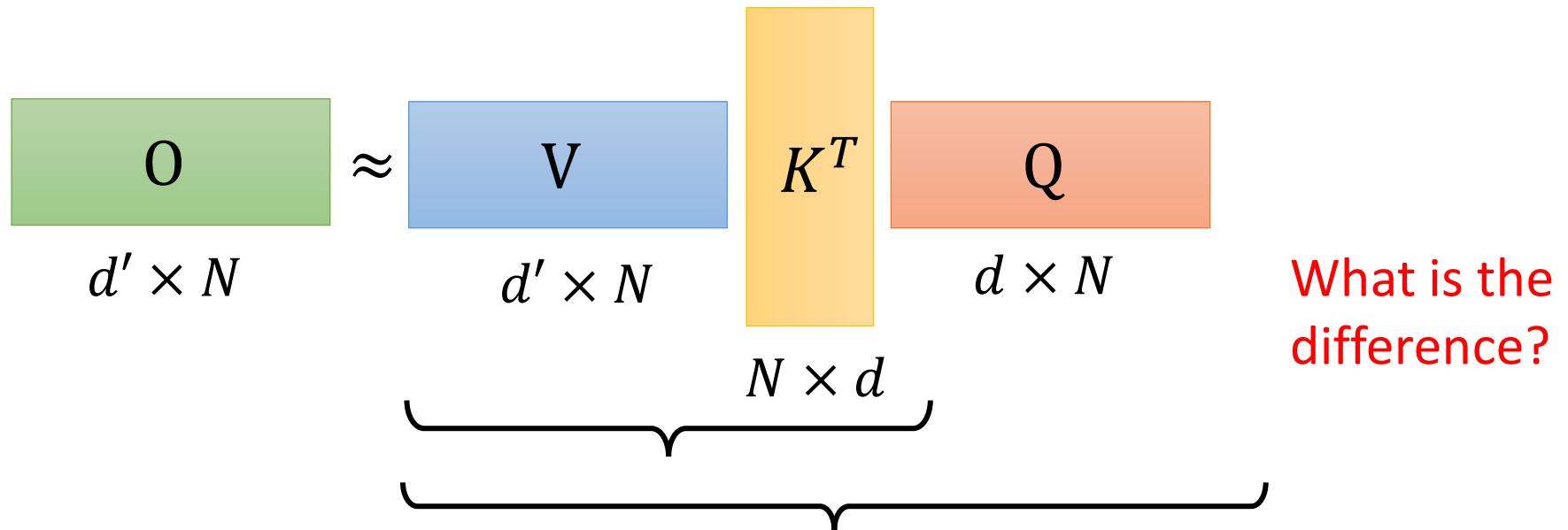
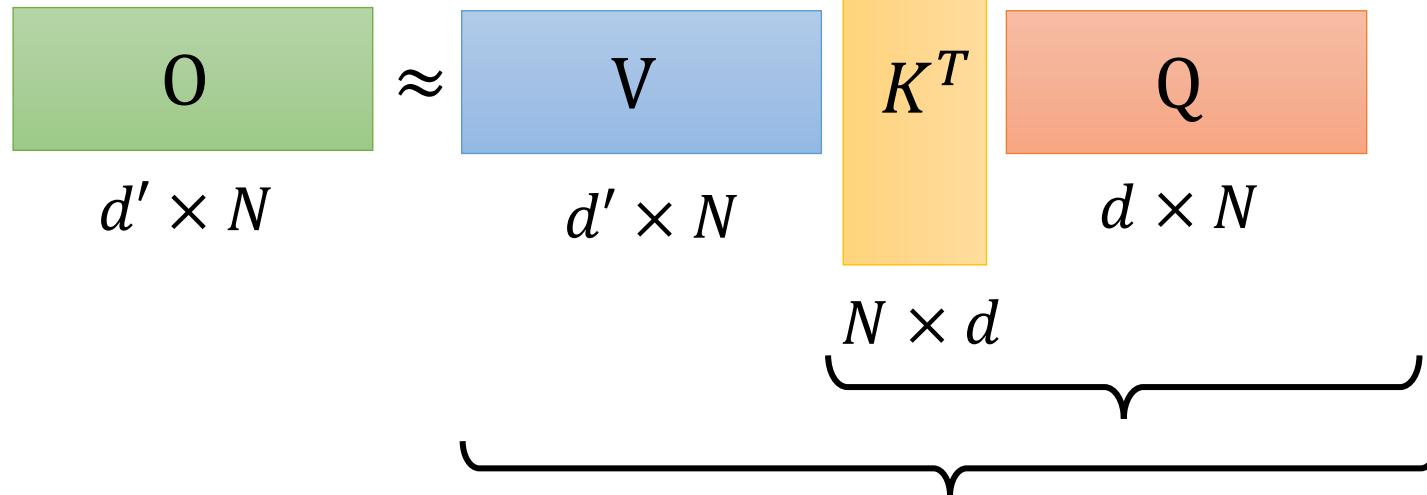
Review

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

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

$N \times d$

The diagram illustrates the dimensions of the matrices involved in the attention mechanism. The matrix  $\text{O}$  is  $d' \times N$ . The matrix  $\text{V}$  is also  $d' \times N$ . The matrix  $K^T$  is  $N \times d$ . The matrix  $\text{Q}$  is  $d \times N$ . Brackets indicate that the width of  $\text{O}$  and  $\text{V}$  is  $d'$ , and the height of  $K^T$  and  $\text{Q}$  is  $N$ .



# 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
  - $A(CP) = (AC)P$



$k \times m$



$m \times n$



$n \times p$



$k \times m$



$m \times p$

$m \times n \times p$

$10^6$

$k \times m \times p$

$10^6$



$k \times m$



$m \times n$



$n \times p$



$k \times n$



$n \times p$

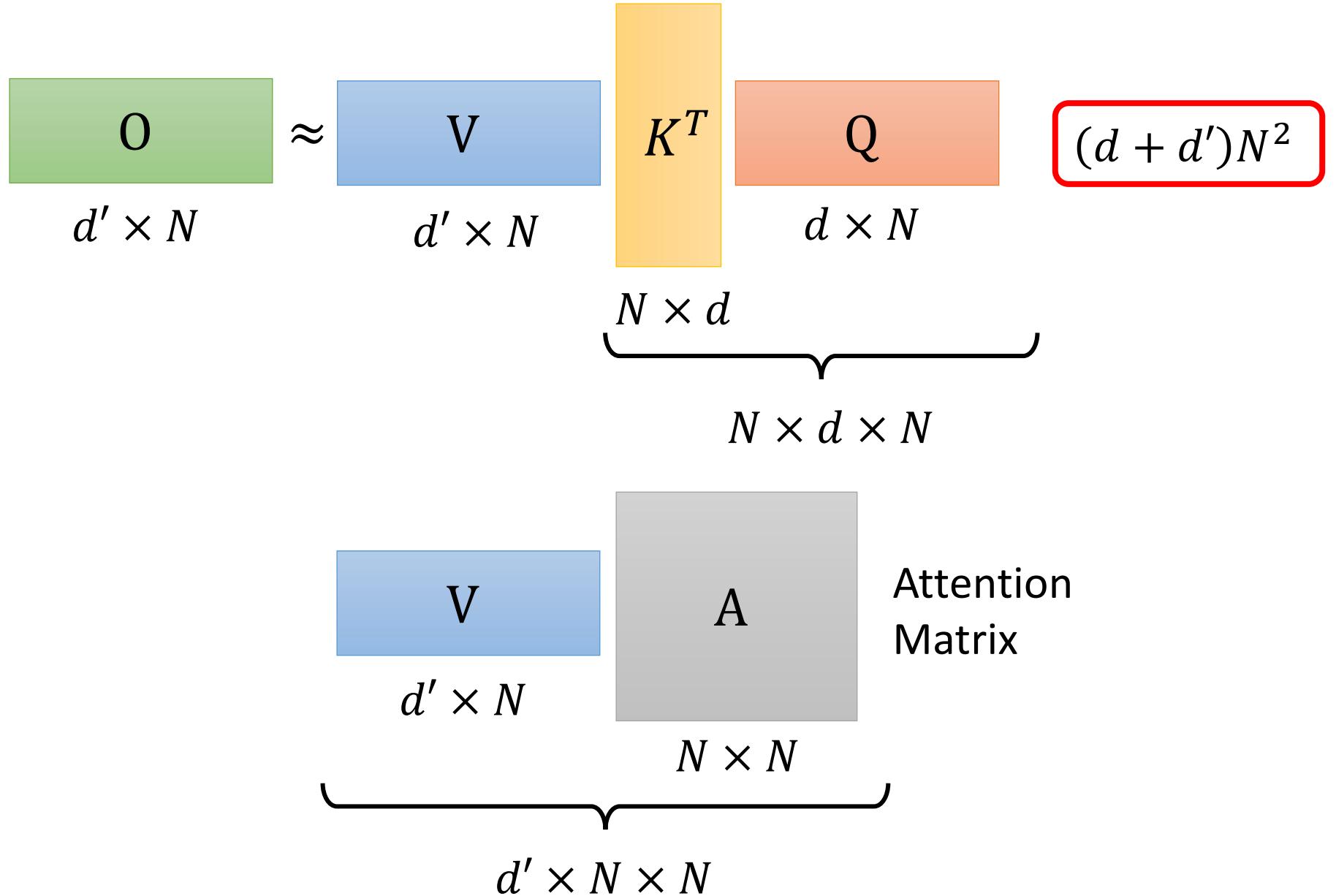
$k \times m \times n$

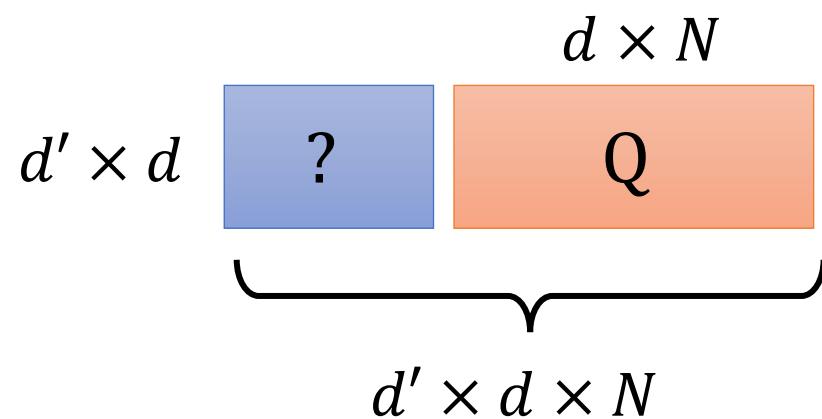
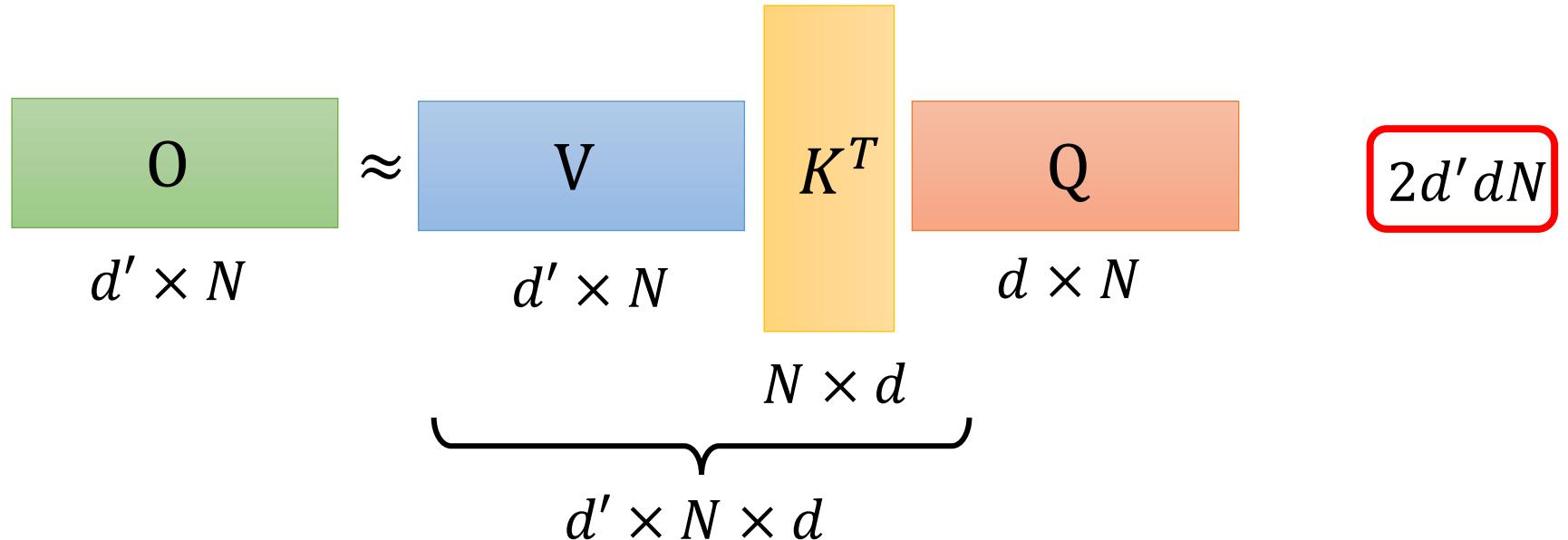
$10^3$

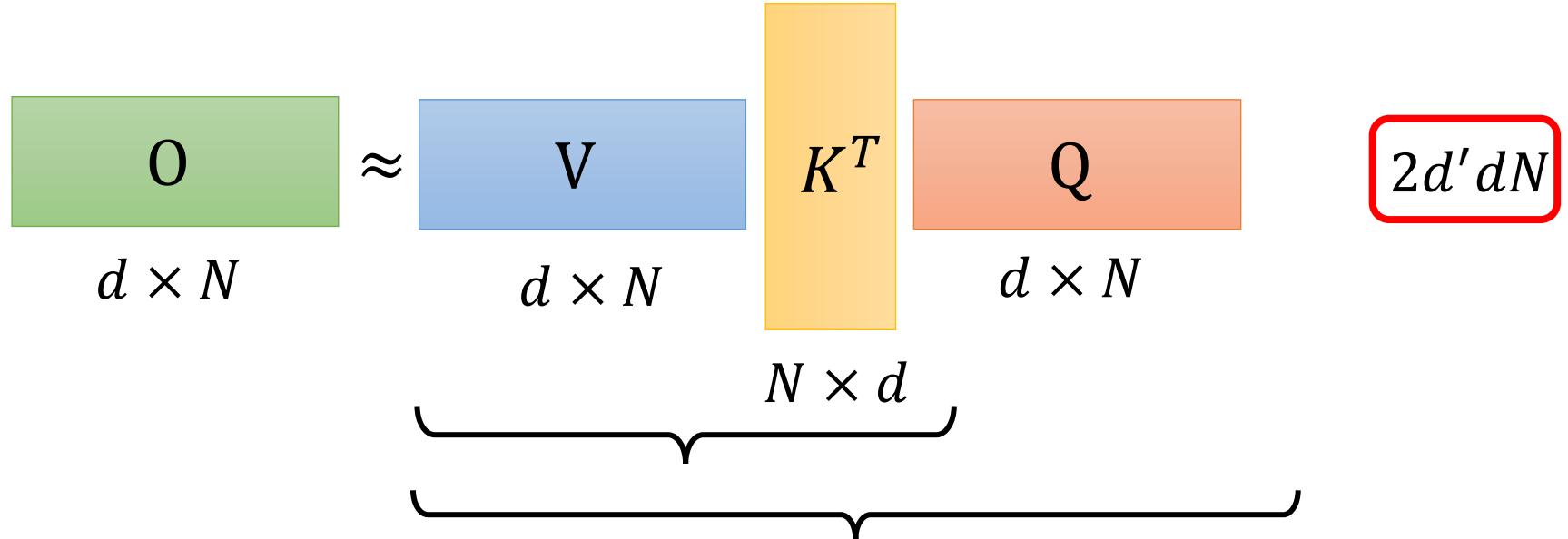
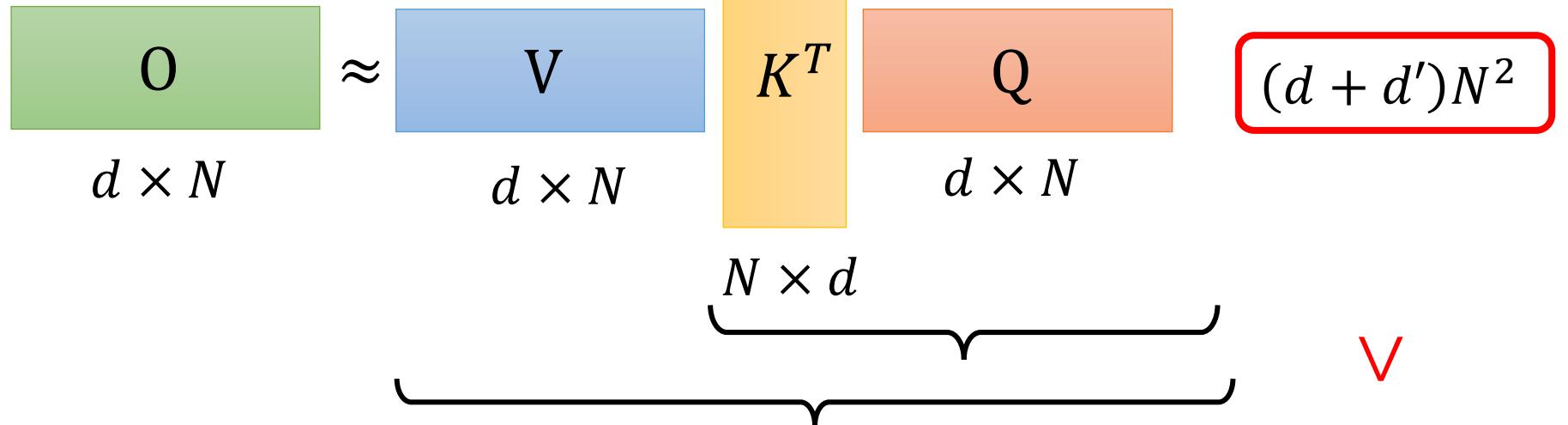
$k \times n \times p$

$10^3$





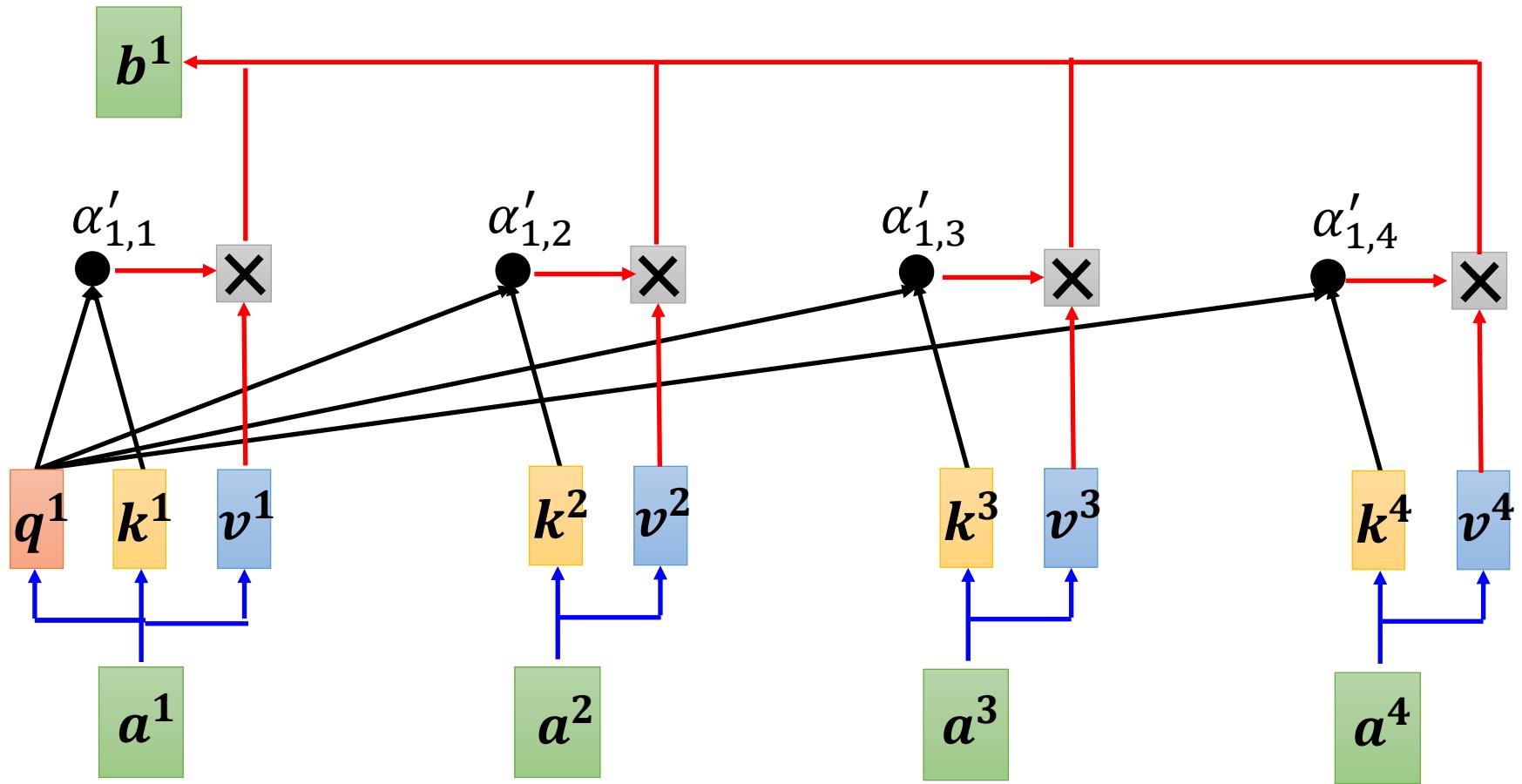




Let's put softmax back ...

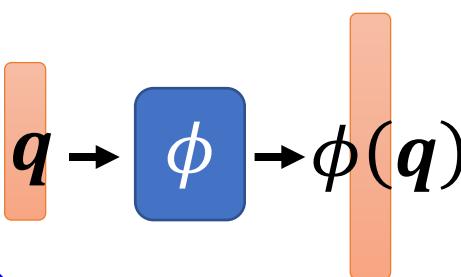
**Warning of math**

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



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

$$\exp(q \cdot k) \approx \phi(q) \cdot \phi(k)$$



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

$$\phi(q^1) \cdot \sum_{j=1}^N \phi(k^j)$$

↗

$\phi(q^1)$

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

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

$$\begin{aligned}
&= [\phi(\mathbf{q}^{\textcolor{green}{1}}) \cdot \phi(\mathbf{k}^{\textcolor{red}{1}})] \mathbf{v}^{\textcolor{red}{1}} + [\phi(\mathbf{q}^{\textcolor{green}{1}}) \cdot \phi(\mathbf{k}^{\textcolor{red}{2}})] \mathbf{v}^{\textcolor{red}{2}} + \dots \\
&= (q_1^{\textcolor{green}{1}} k_1^{\textcolor{red}{1}} + q_2^{\textcolor{green}{1}} k_2^{\textcolor{red}{1}} + \dots) \mathbf{v}^{\textcolor{red}{1}} + (q_1^{\textcolor{green}{1}} k_1^{\textcolor{red}{2}} + q_2^{\textcolor{green}{1}} k_2^{\textcolor{red}{2}} + \dots) \mathbf{v}^{\textcolor{red}{2}} + \dots \\
&= \underline{q_1^{\textcolor{green}{1}} k_1^{\textcolor{red}{1}} \mathbf{v}^{\textcolor{red}{1}}} + \underline{q_2^{\textcolor{green}{1}} k_2^{\textcolor{red}{1}} \mathbf{v}^{\textcolor{red}{1}}} + \dots + \underline{q_1^{\textcolor{green}{1}} k_1^{\textcolor{red}{2}} \mathbf{v}^{\textcolor{red}{2}}} + \underline{q_2^{\textcolor{green}{1}} k_2^{\textcolor{red}{2}} \mathbf{v}^{\textcolor{red}{2}}} + \dots + \dots \\
&= q_1^{\textcolor{green}{1}} (k_1^{\textcolor{red}{1}} \mathbf{v}^{\textcolor{red}{1}} + k_1^{\textcolor{red}{2}} \mathbf{v}^{\textcolor{red}{2}} + \dots) + q_2^{\textcolor{green}{1}} (k_2^{\textcolor{red}{1}} \mathbf{v}^{\textcolor{red}{1}} + k_2^{\textcolor{red}{2}} \mathbf{v}^{\textcolor{red}{2}} + \dots)
\end{aligned}$$

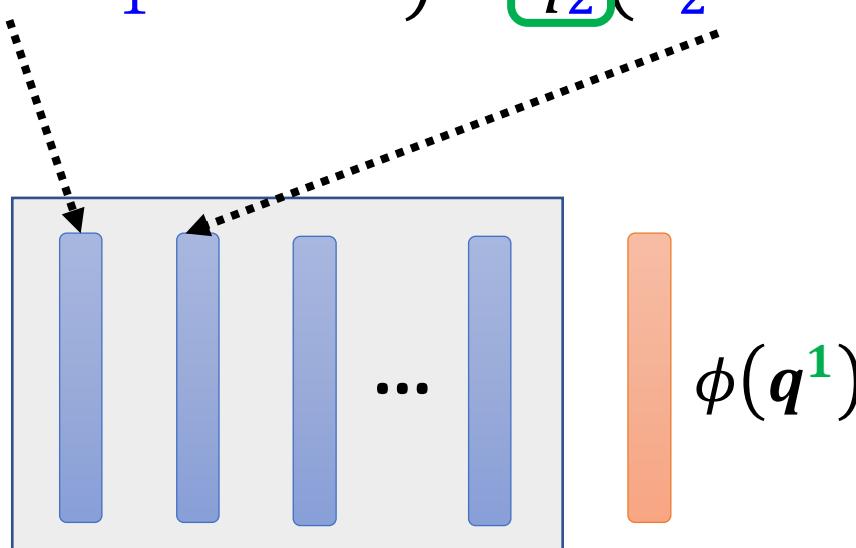
$$\mathbf{b}^{\textcolor{green}{1}} = \sum_{\textcolor{red}{i}=1}^N \alpha'_{\textcolor{green}{1},\textcolor{red}{i}} v^{\textcolor{red}{i}} = \frac{\sum_{\textcolor{red}{i}=1}^N [\phi(\mathbf{q}^{\textcolor{green}{1}}) \cdot \phi(\mathbf{k}^{\textcolor{red}{i}})] v^{\textcolor{red}{i}}}{\phi(\mathbf{q}^{\textcolor{green}{1}})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

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

$$\phi(\mathbf{q}^{\textcolor{green}{1}}) = \begin{bmatrix} q_1^{\textcolor{green}{1}} \\ q_2^{\textcolor{green}{1}} \\ \vdots \end{bmatrix} \quad \phi(\mathbf{k}^{\textcolor{red}{1}}) = \begin{bmatrix} k_1^{\textcolor{red}{1}} \\ k_2^{\textcolor{red}{1}} \\ \vdots \end{bmatrix}$$

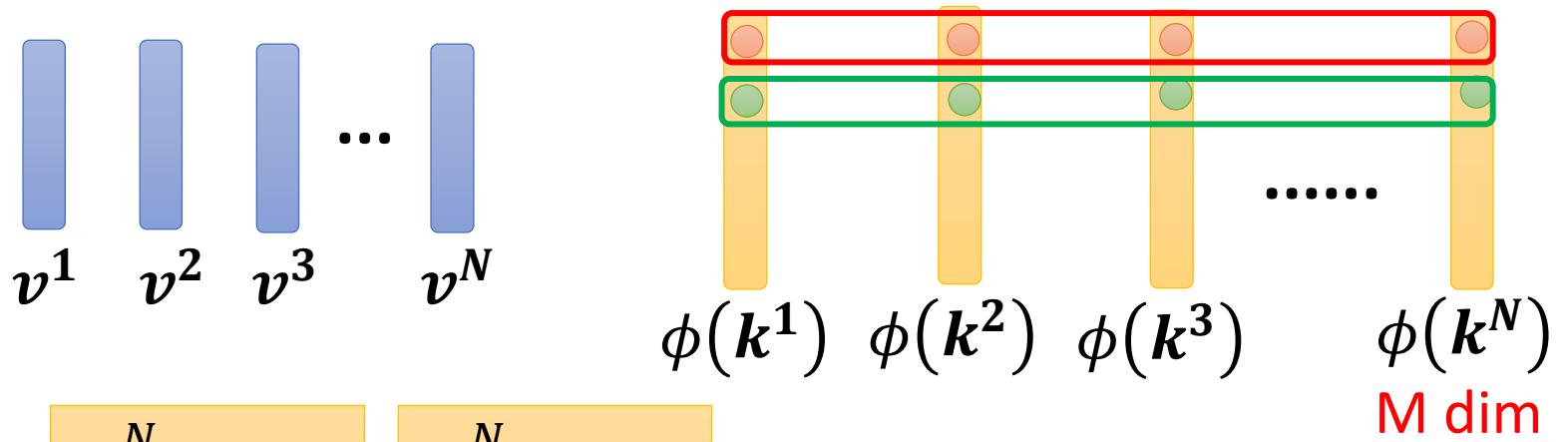
$$= \boxed{q_1^{\textcolor{green}{1}}} (k_1^{\textcolor{red}{1}} v^{\textcolor{red}{1}} + k_1^{\textcolor{red}{2}} v^{\textcolor{red}{2}} + \dots) + \boxed{q_2^{\textcolor{green}{1}}} (k_2^{\textcolor{red}{1}} v^{\textcolor{red}{1}} + k_2^{\textcolor{red}{2}} v^{\textcolor{red}{2}} + \dots)$$

$$\sum_{\textcolor{red}{j}=1}^N k_1^{\textcolor{red}{j}} v^{\textcolor{red}{j}}$$



$$\sum_{\textcolor{red}{j}=1}^N k_2^{\textcolor{red}{j}} v^{\textcolor{red}{j}}$$

$\mathbf{M}$  vectors

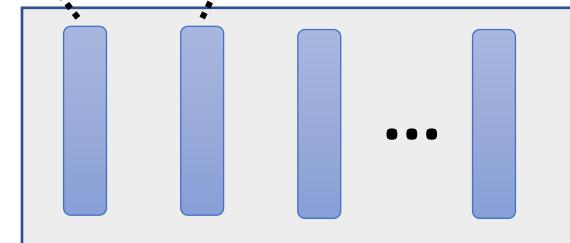


$$\sum_{j=1}^N k_1^j v^j$$

$$\sum_{j=1}^N k_2^j v^j$$

**M vectors**

$$b_1^{(1)} =$$



$$\sum_{j=1}^N \phi(k^j)$$

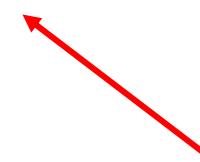
$$\phi(q_1^{(1)})$$

**M dim**

$$\phi(q_1^{(1)})$$

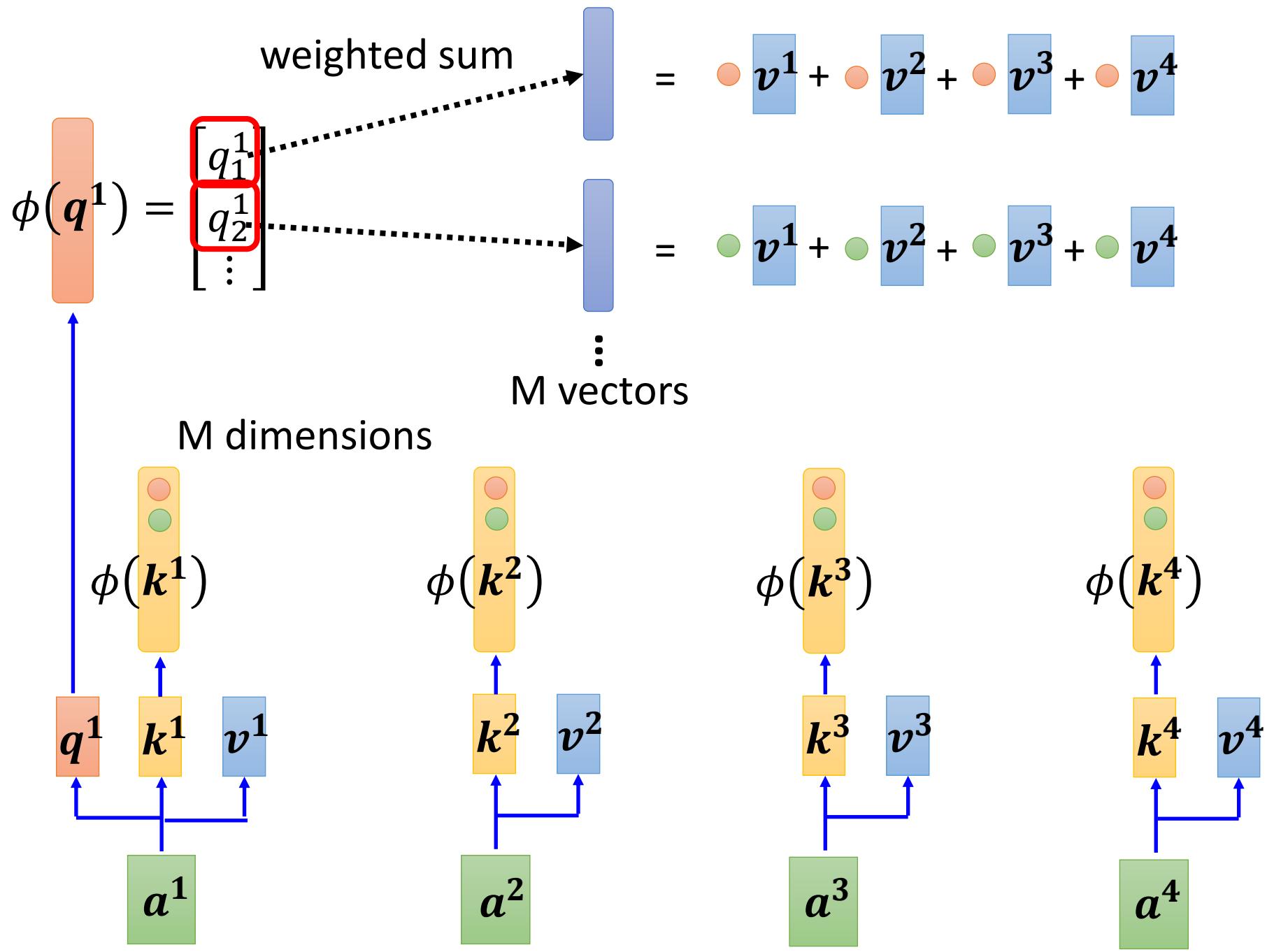
$$b^1 = \frac{\phi(q^1)}{\phi(q^1)}$$
$$b^2 = \frac{\phi(q^2)}{\phi(q^2)}$$
$$\vdots$$
$$b^N = \dots$$

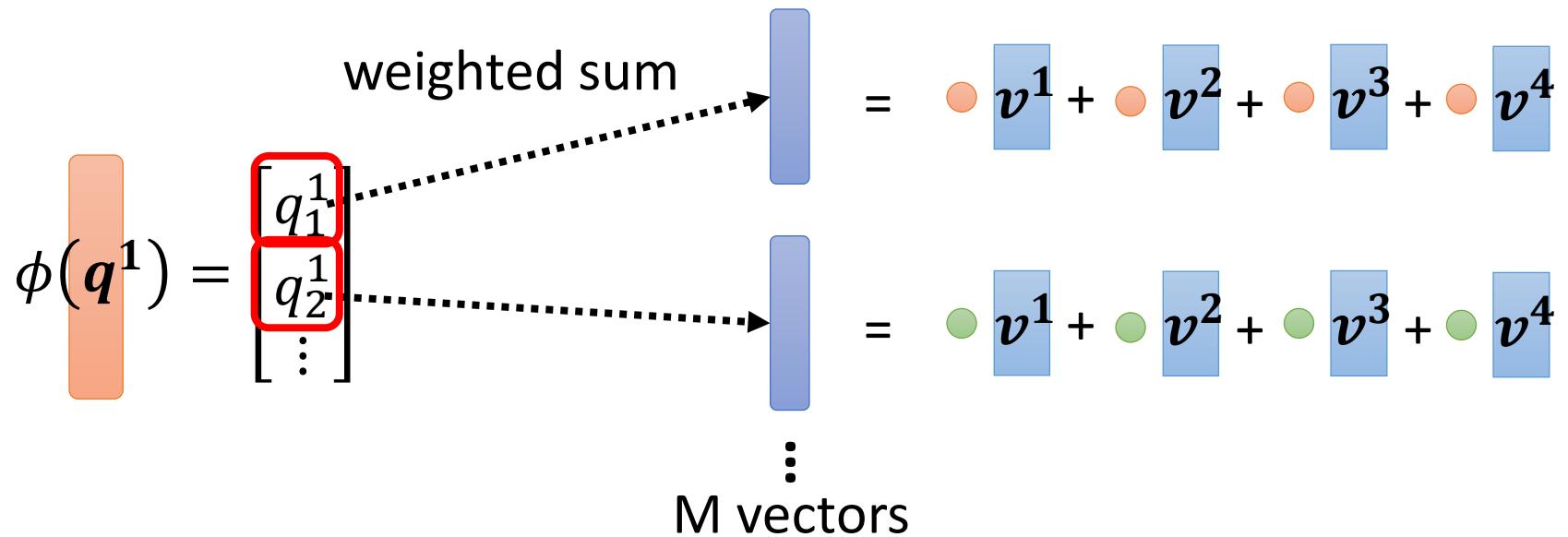
Don't compute again



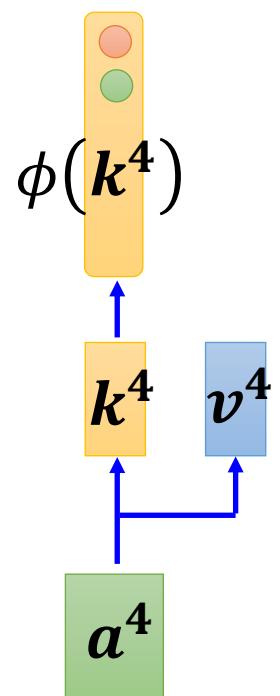
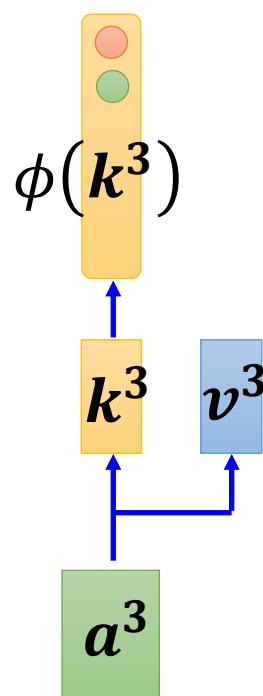
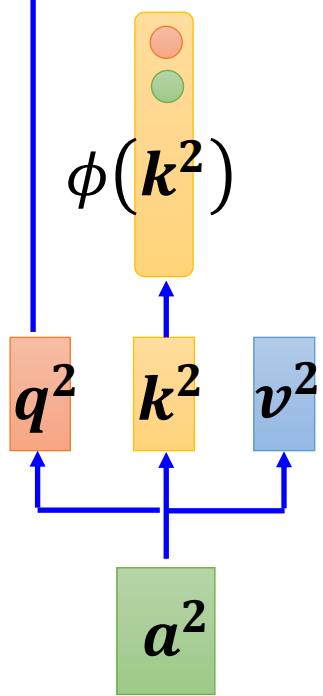
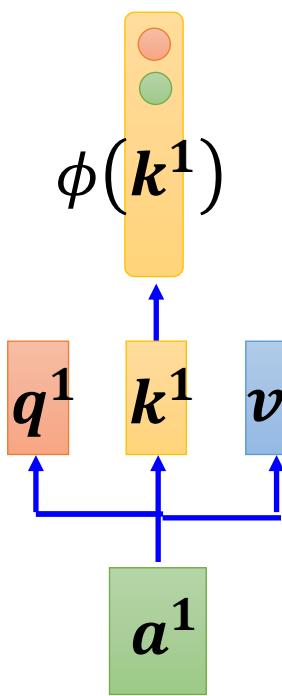
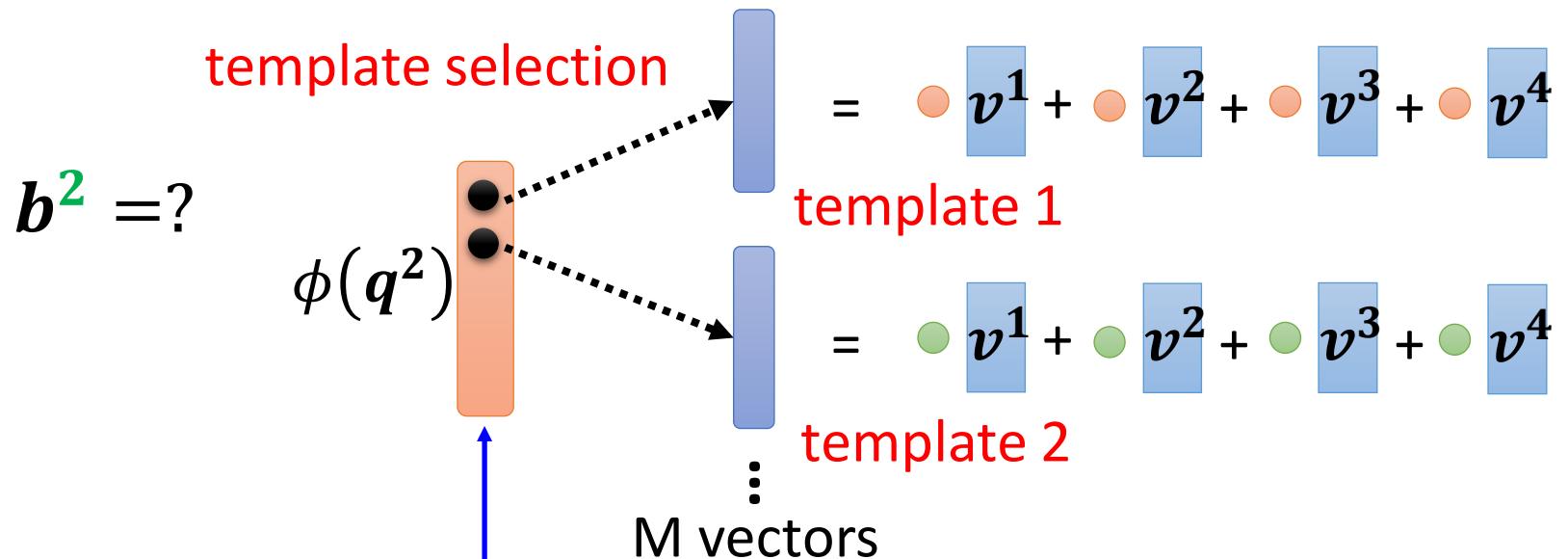
Let's put softmax back ...

End of warning

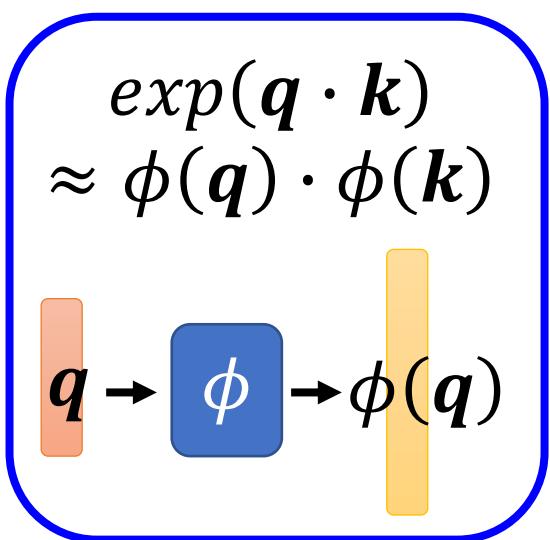




$$b^1 = \frac{\sum_{j=1}^N \phi(k^j)}{\phi(q^1)}$$



# 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

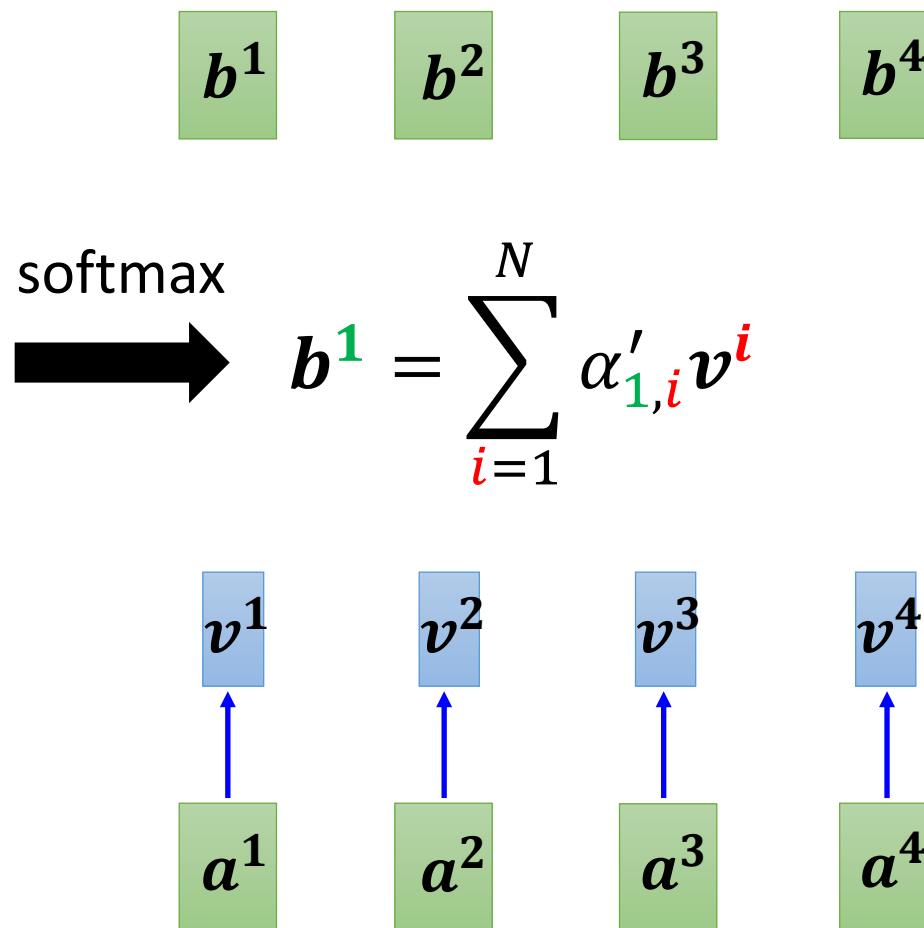
<https://arxiv.org/pdf/2009.14794.pdf>

# Do we need q and k to compute attention? Synthesizer!

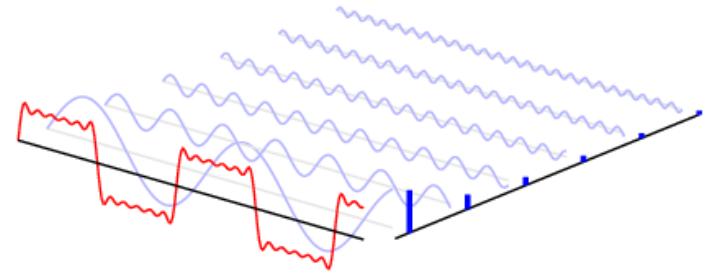
$\alpha_{1,1}$	$\alpha_{1,2}$	$\alpha_{1,3}$	$\alpha_{1,4}$
$\alpha_{1,2}$	$\alpha_{2,2}$	$\alpha_{2,3}$	$\alpha_{2,4}$
$\alpha_{1,3}$	$\alpha_{2,3}$	$\alpha_{3,3}$	$\alpha_{3,4}$
$\alpha_{1,4}$	$\alpha_{2,4}$	$\alpha_{3,4}$	$\alpha_{4,4}$

~~From q and k?~~

They are network  
parameters!



# Attention-free?



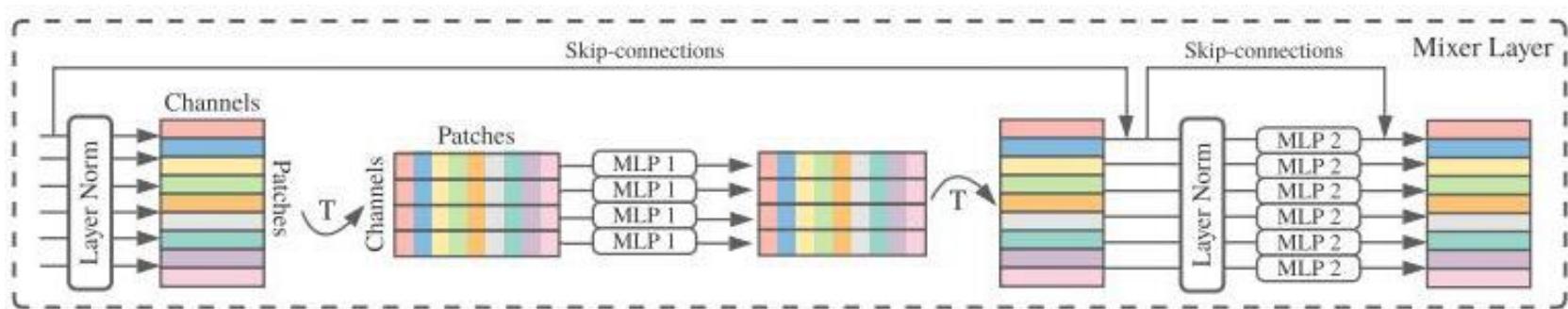
- Fnet: Mixing tokens with fourier transforms

<https://arxiv.org/abs/2105.03824>

- Pay Attention to MLPs <https://arxiv.org/abs/2105.08050>

- MLP-Mixer: An all-MLP Architecture for Vision

<https://arxiv.org/abs/2105.01601>



# Summary

- Human knowledge
  - Local Attention, Big Bird
- Clustering
  - Reformer
- Learnable Pattern
  - Sinkhorn
- Representative key
  - Linformer
- $k, q$  first  $\rightarrow v, k$  first
  - Linear Transformer, Performer
- New framework
  - Synthesizer

