

加快語言模型 的生成速度

先備知識



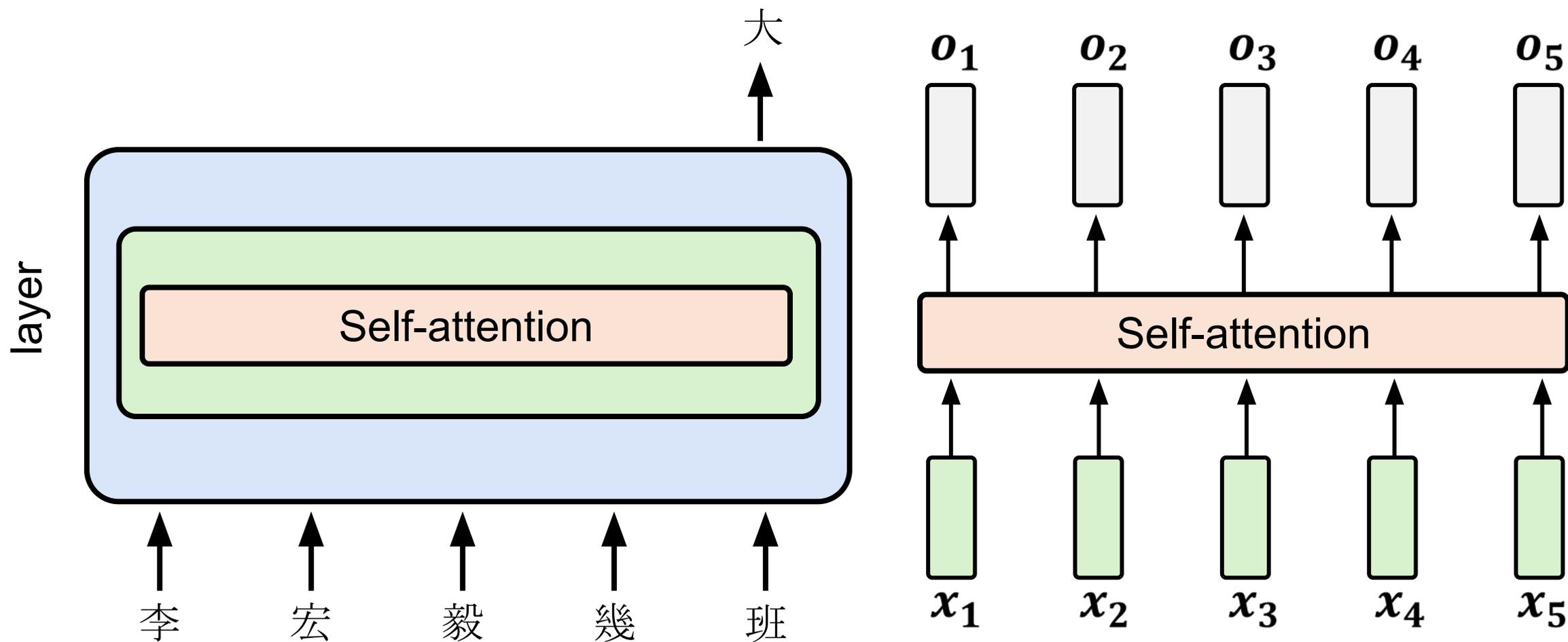
[https://youtu.be/8iFvM7WUUs8?
si=gymWr9Vurpb8ri2Z](https://youtu.be/8iFvM7WUUs8?si=gymWr9Vurpb8ri2Z)

一堂課看懂 語言模型內部運作

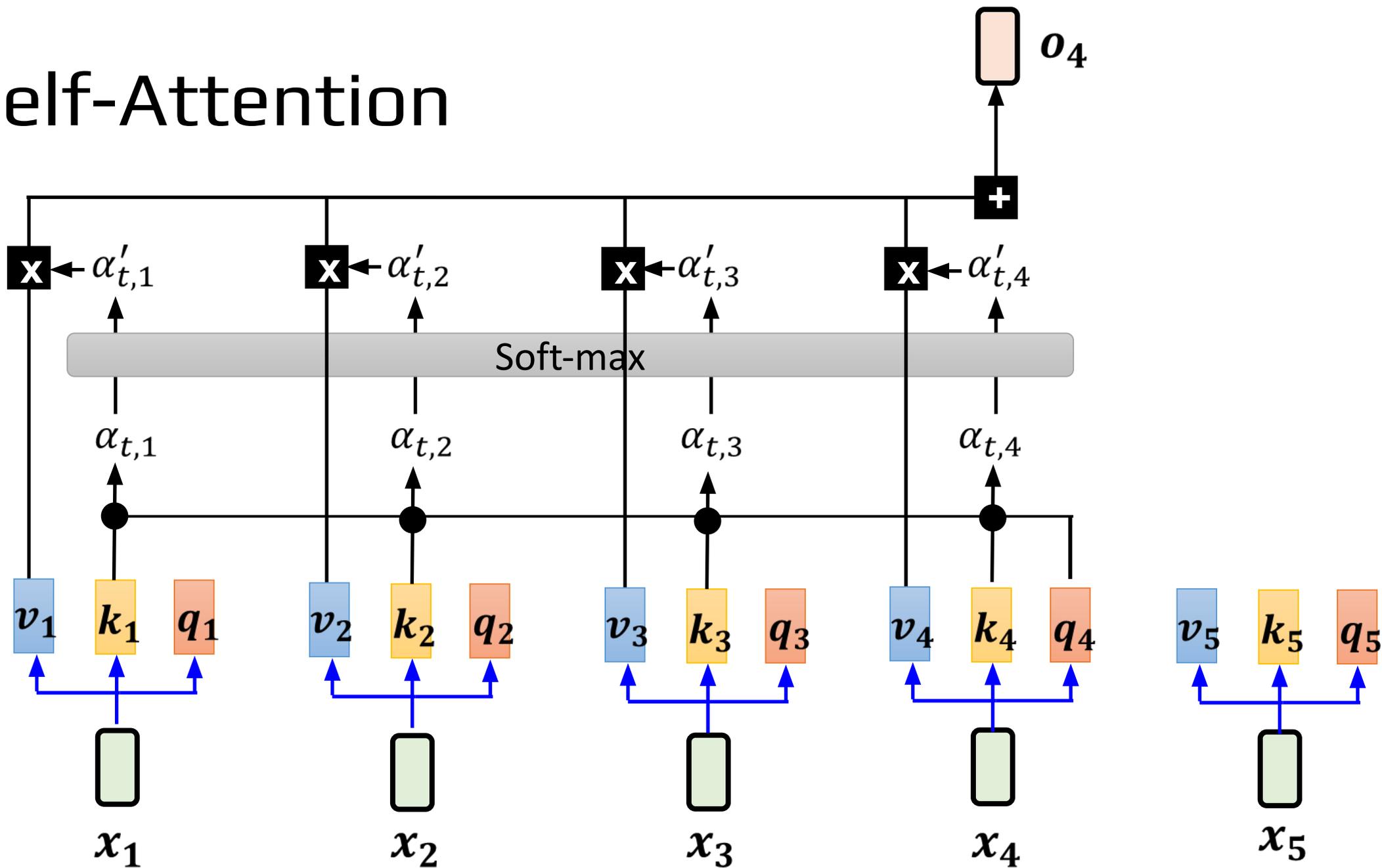
李宏毅

【生成式人工智慧與機器學習導論2025】第3講：解剖大型語言模型

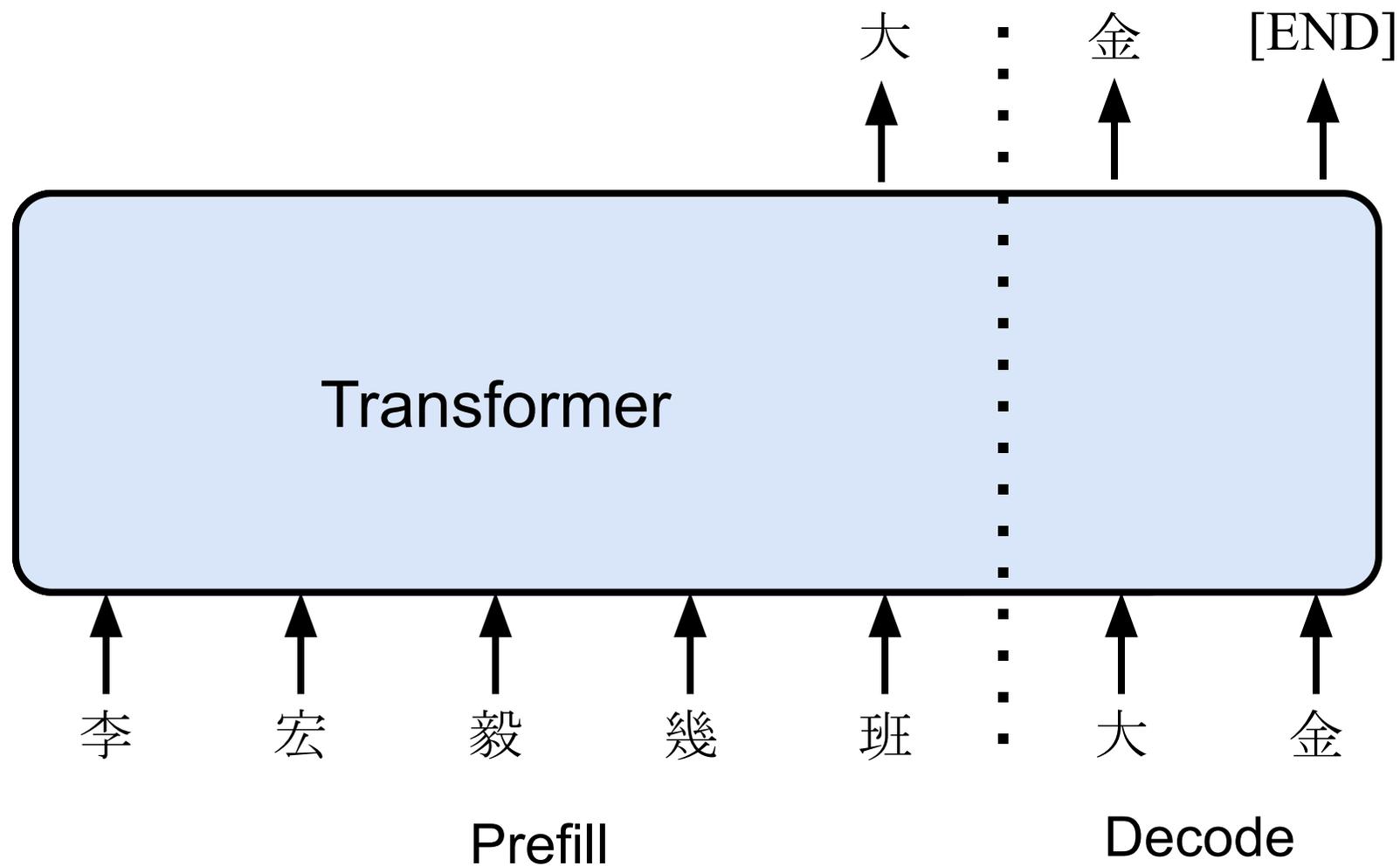
語言模型如何生成 (推論, Inference)



Self-Attention



語言模型如何生成 (推論, Inference)



加快推論速度的經典方法

Today

Flash Attention

KV Cache

Speculative Decoding

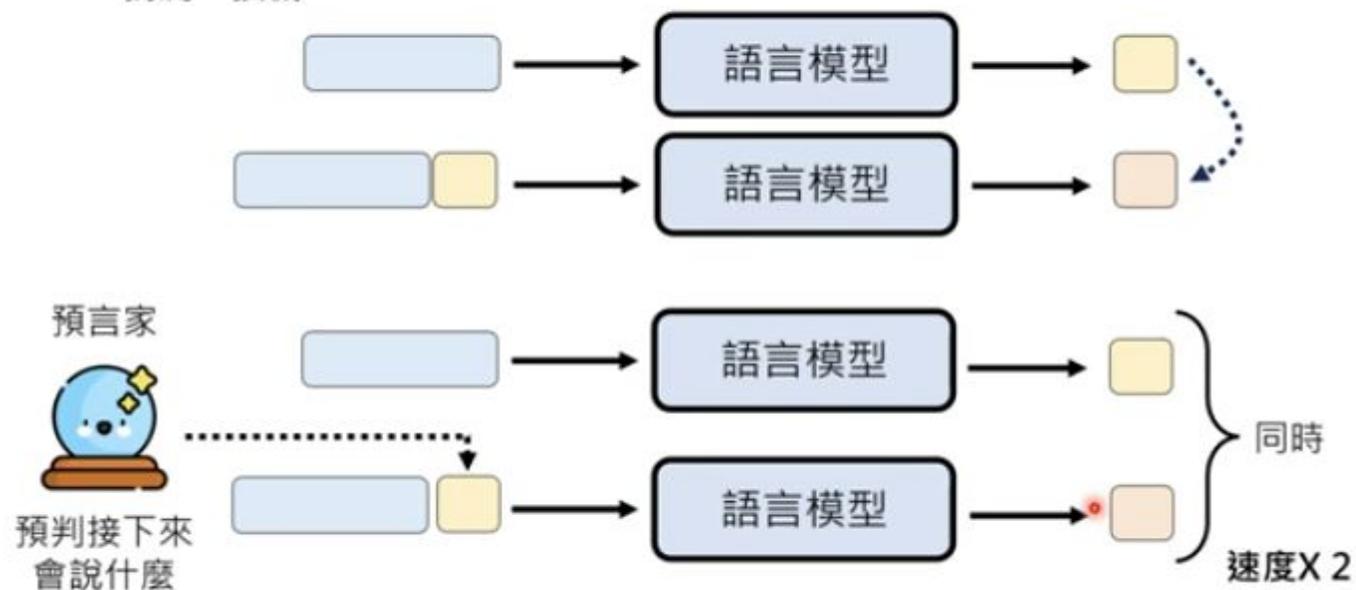
Speculative Decoding

<https://arxiv.org/abs/2211.17192>
<https://arxiv.org/abs/2302.01318>



Speculative Decoding

猜測、投機



【生成式AI導論 2024】第16講：可以加速所有語言模型生成速度的神奇外掛 — Speculative Decoding

<https://youtu.be/MAbGgsWKrg8?si=88lZbhk0GVJtN0Yi>

Flash Attention

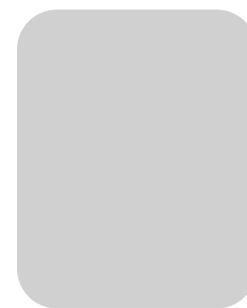
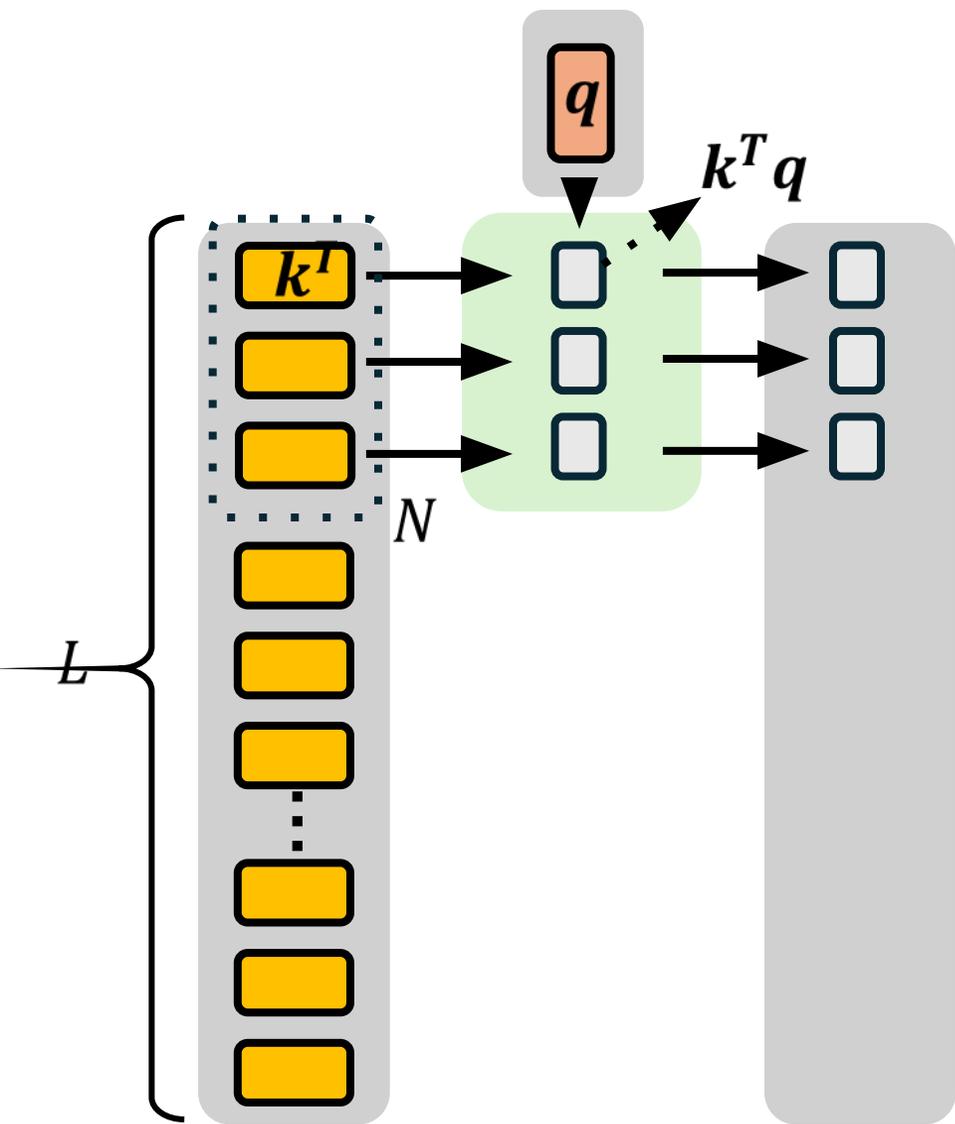
<https://arxiv.org/abs/2205.14135>



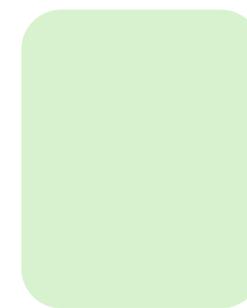
HBM



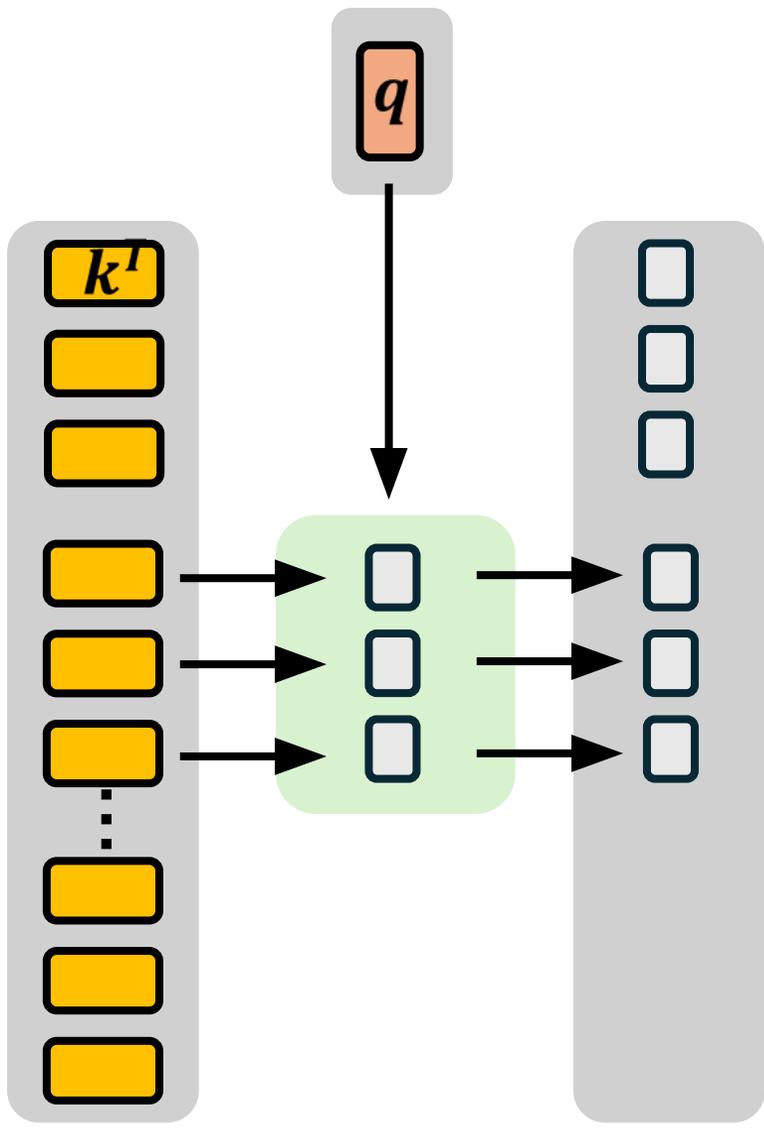
SRAM

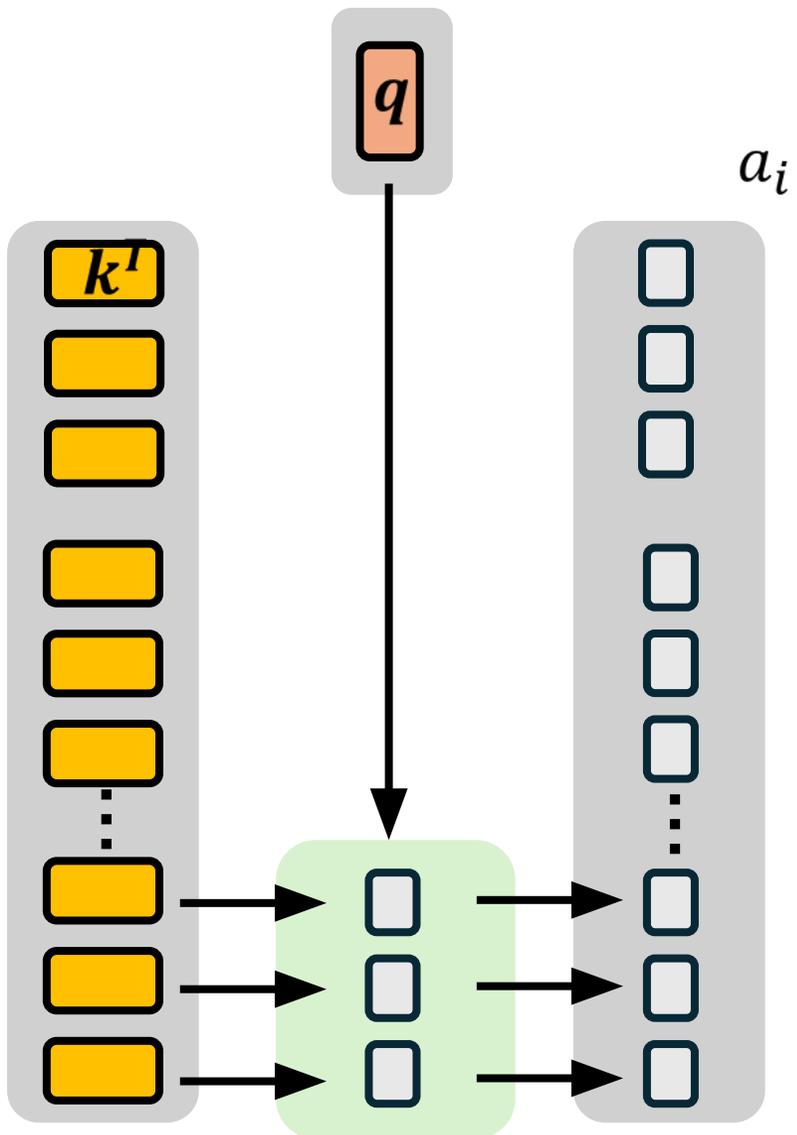


倉庫



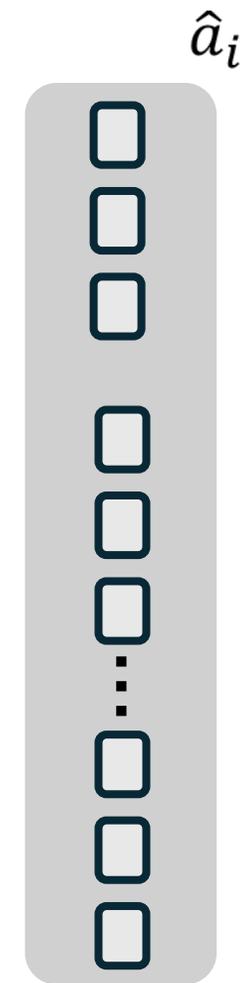
工作臺

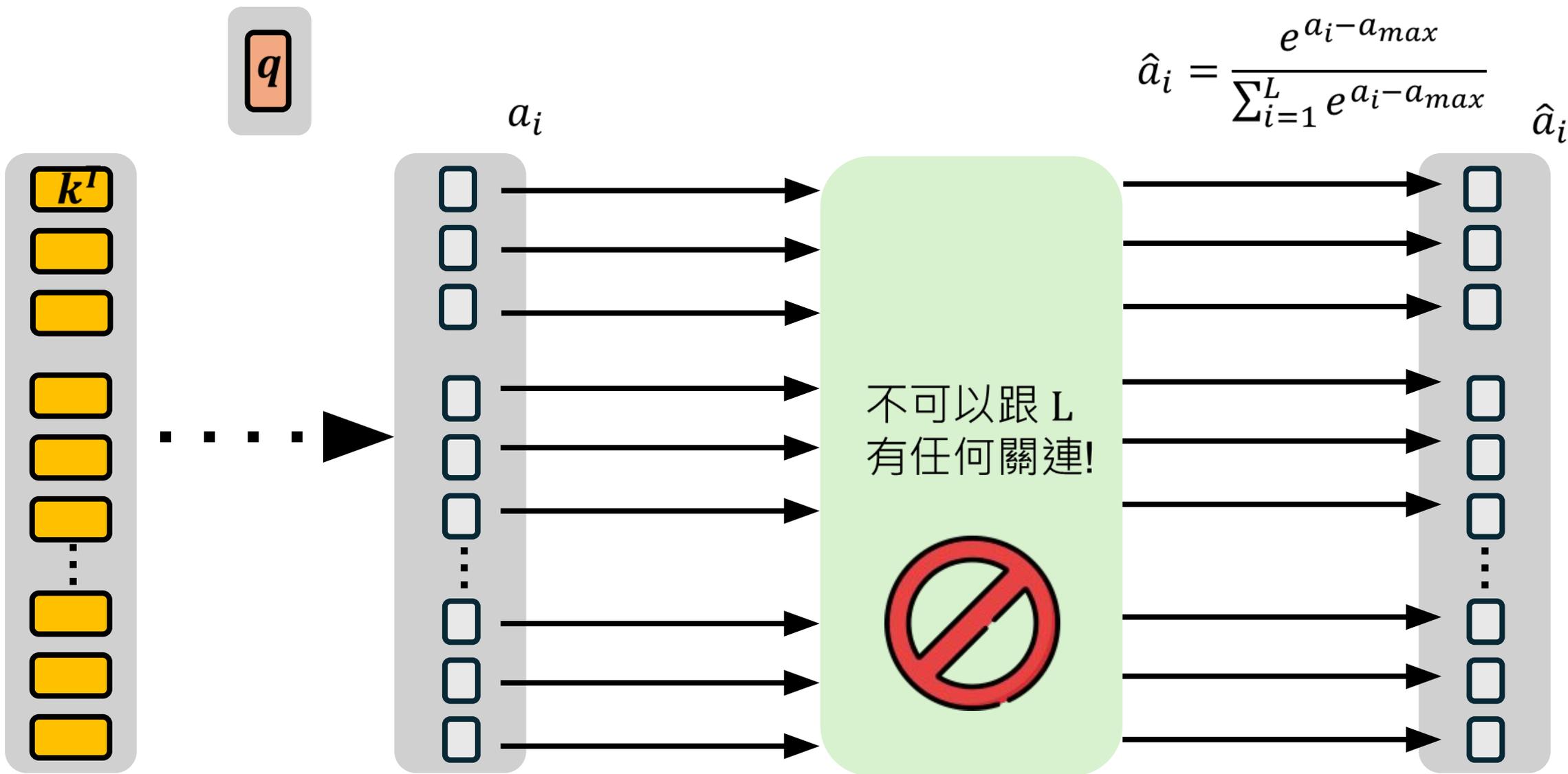


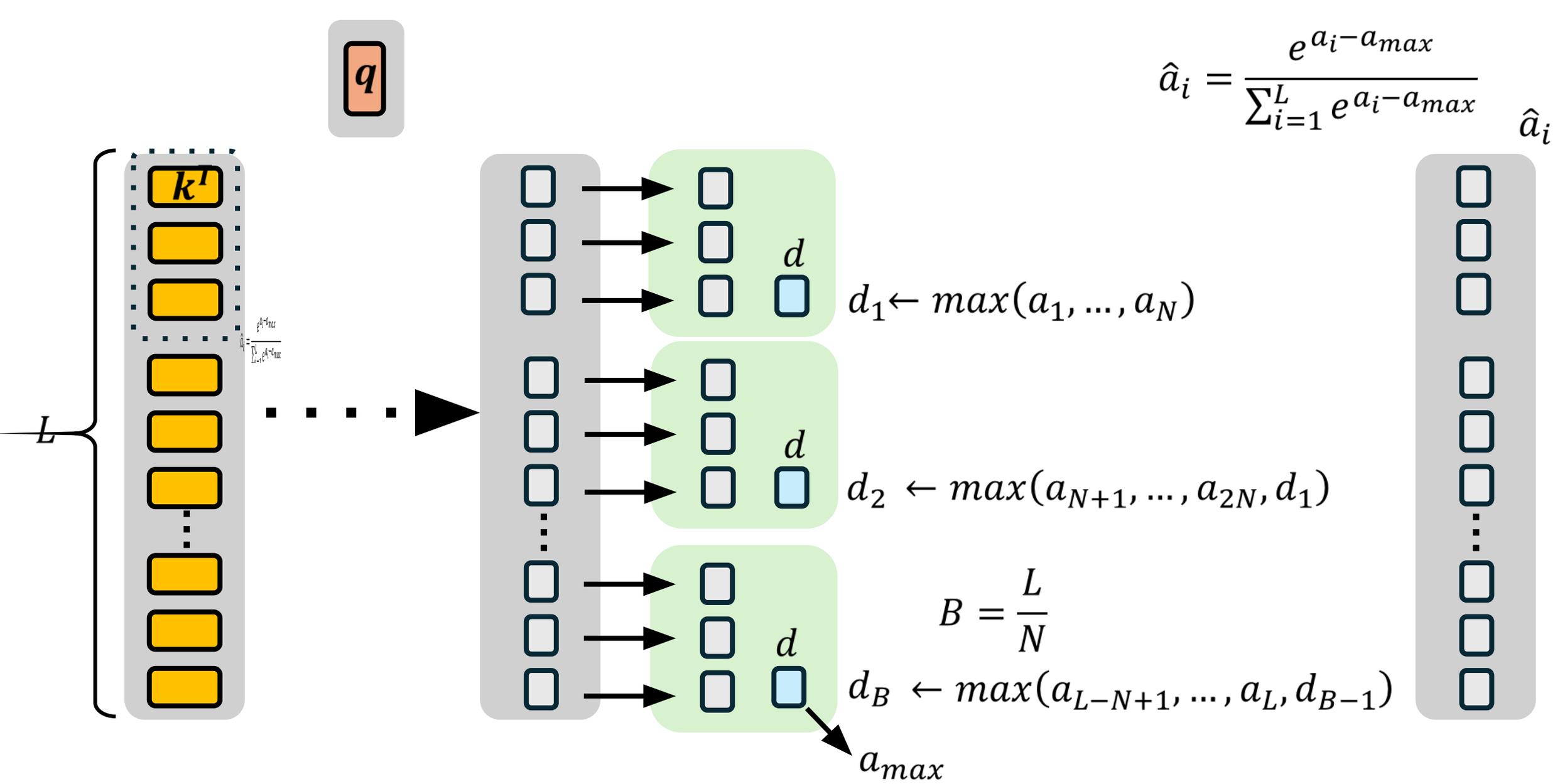


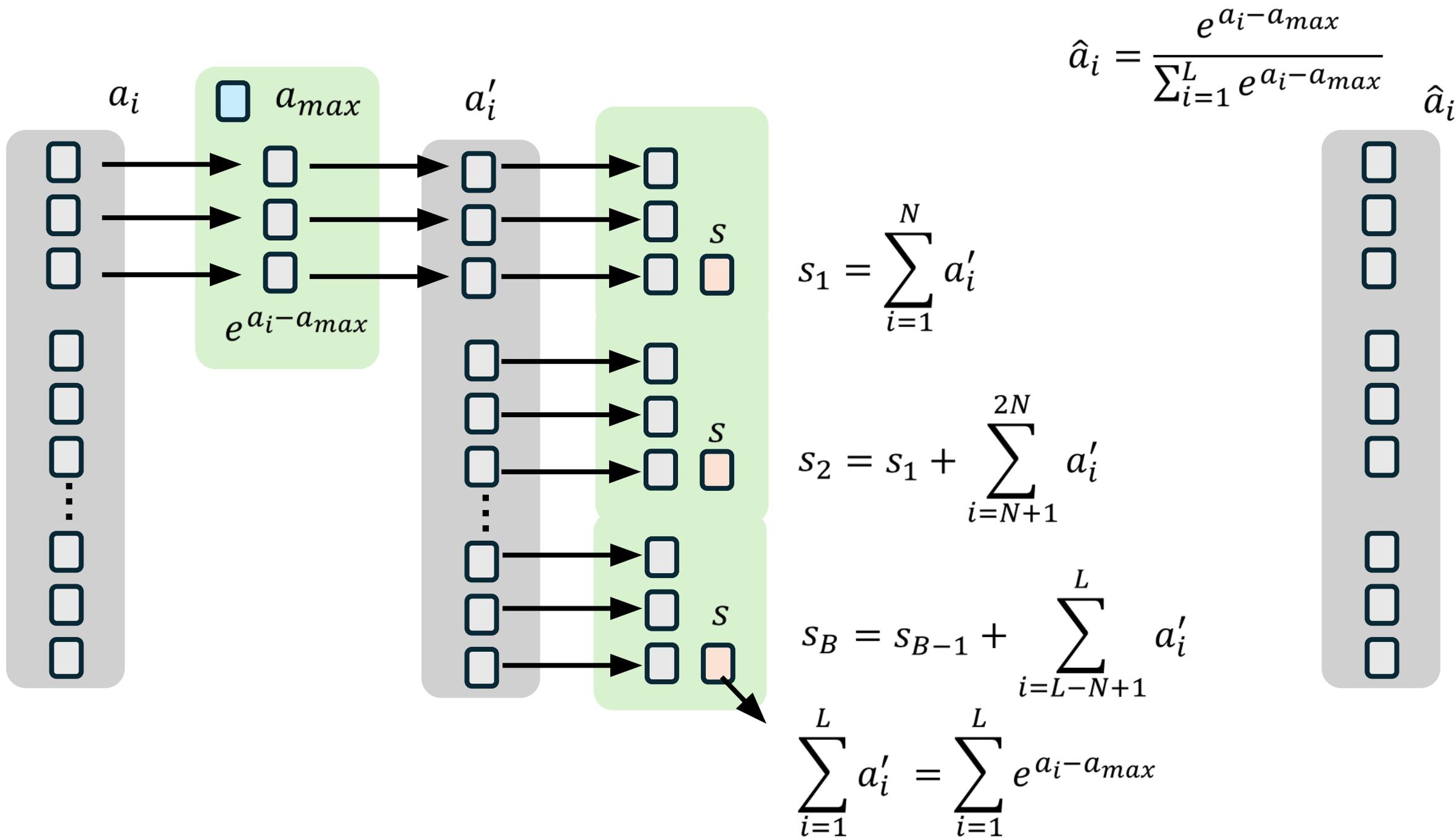
$$\hat{a}_i = \frac{e^{a_i}}{\sum_{i=1}^L e^{a_i}}$$

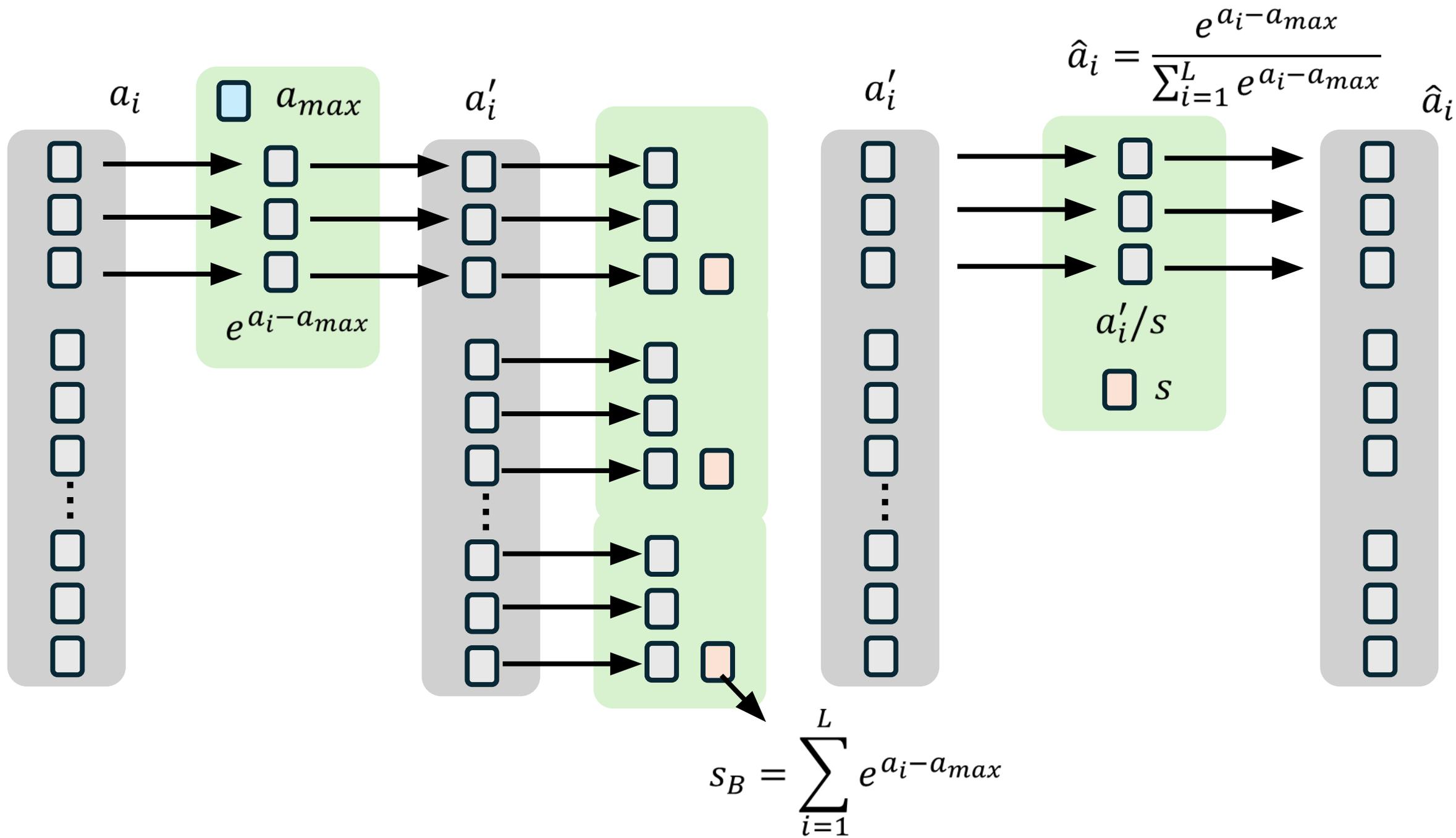
$$\hat{a}_i = \frac{e^{a_i - a_{max}}}{\sum_{i=1}^L e^{a_i - a_{max}}}$$

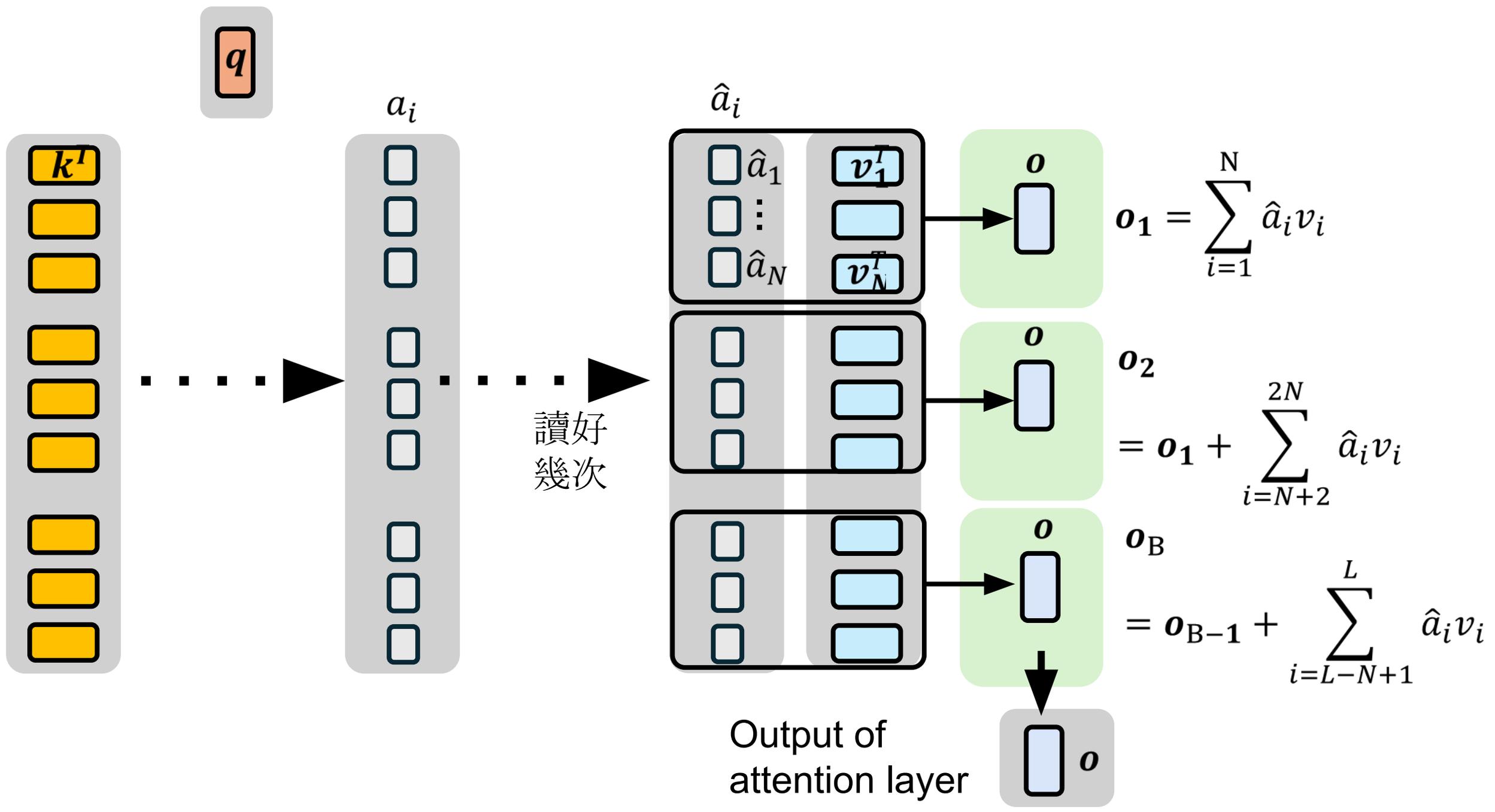






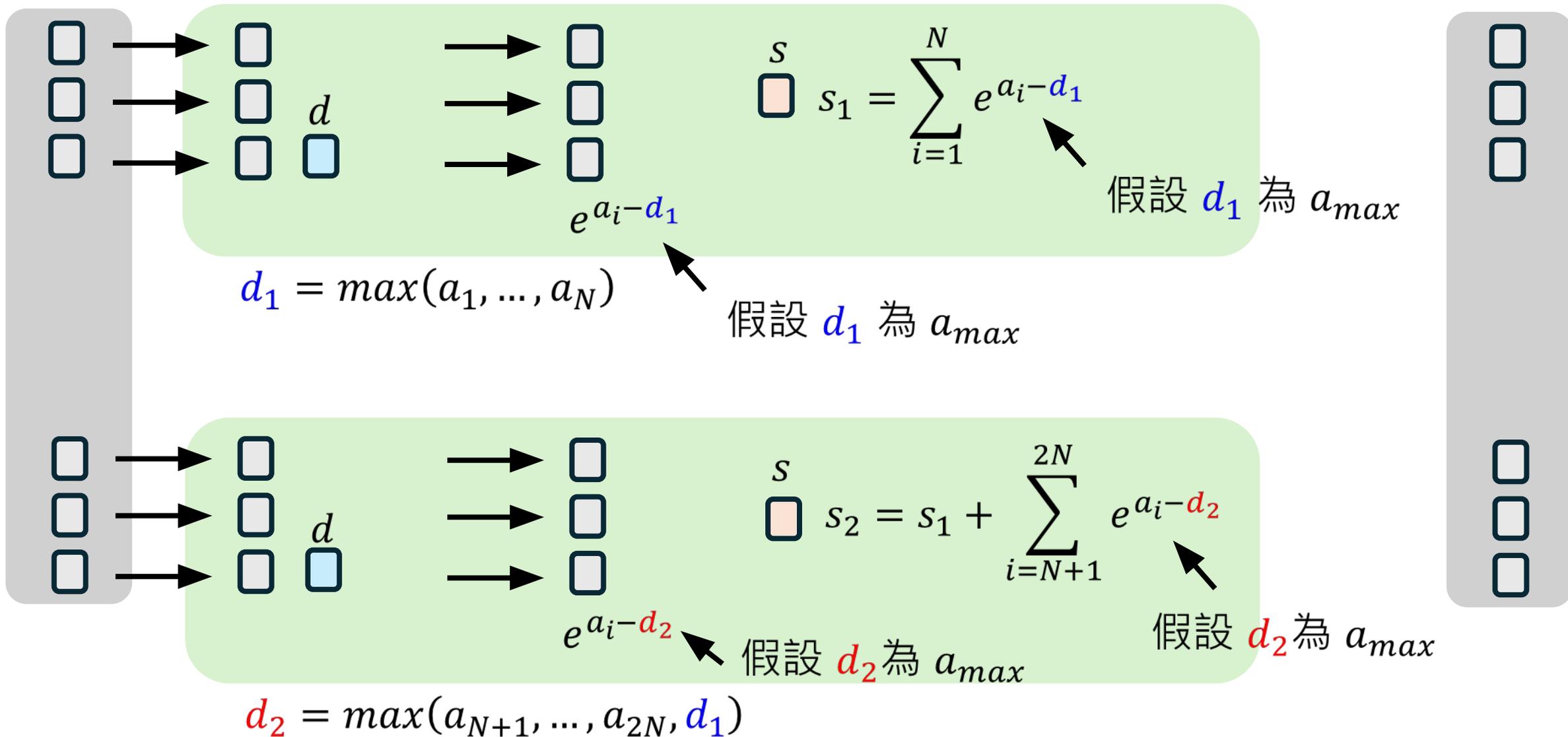






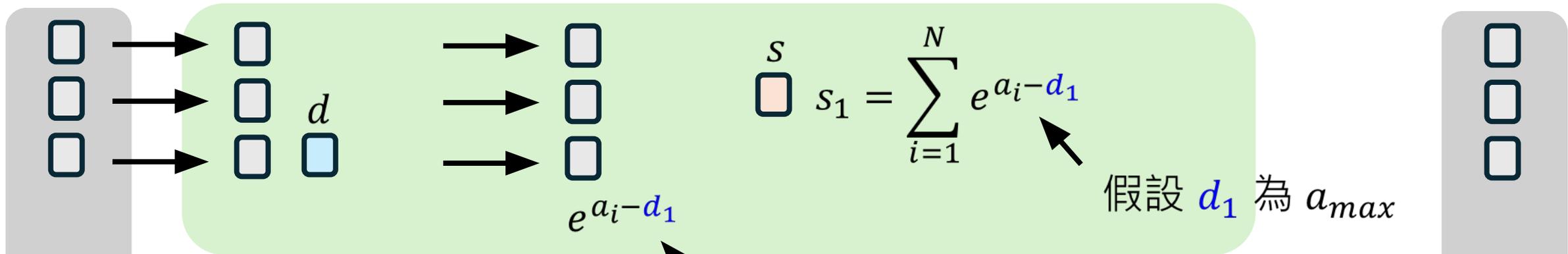
一次找出 a_{max} 和 $\sum_{i=1}^L e^{a_i - a_{max}}$

$$\hat{a}_i = \frac{e^{a_i - a_{max}}}{\sum_{i=1}^L e^{a_i - a_{max}}} \hat{a}_i$$



一次找出 a_{max} 和 $\sum_{i=1}^L e^{a_i - a_{max}}$

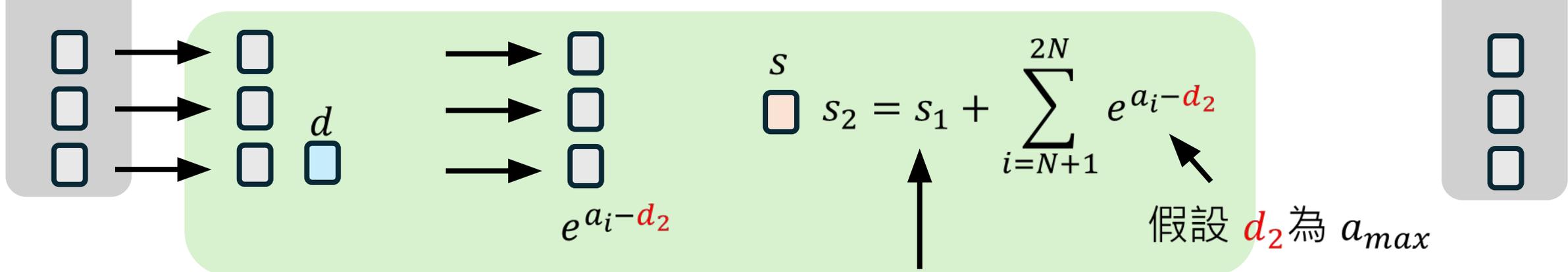
$$\hat{a}_i = \frac{e^{a_i - a_{max}}}{\sum_{i=1}^L e^{a_i - a_{max}}} \hat{a}_i$$



$$d_1 = \max(a_1, \dots, a_N)$$

假設 d_1 為 a_{max}

假設 d_1 為 a_{max}



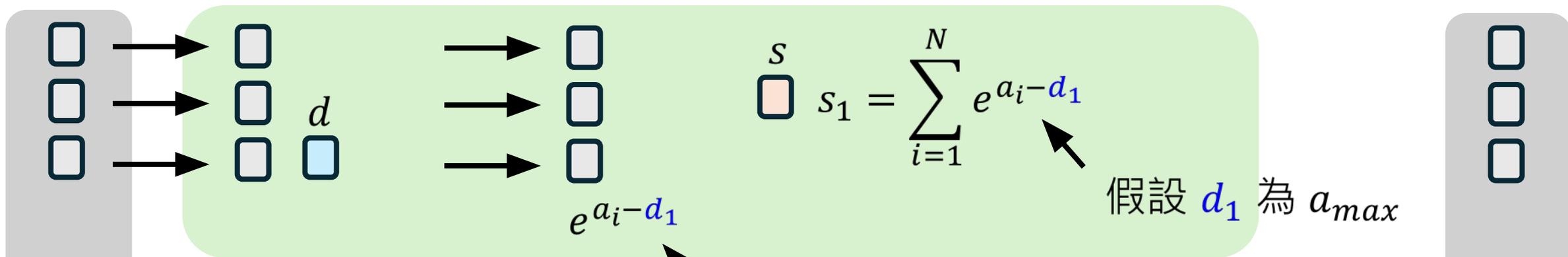
$$d_2 = \max(a_{N+1}, \dots, a_{2N}, d_1)$$

這是假設 l_1 為 a_{max} 計算的

假設 d_2 為 a_{max}

一次找出 a_{max} 和 $\sum_{i=1}^L e^{a_i - a_{max}}$

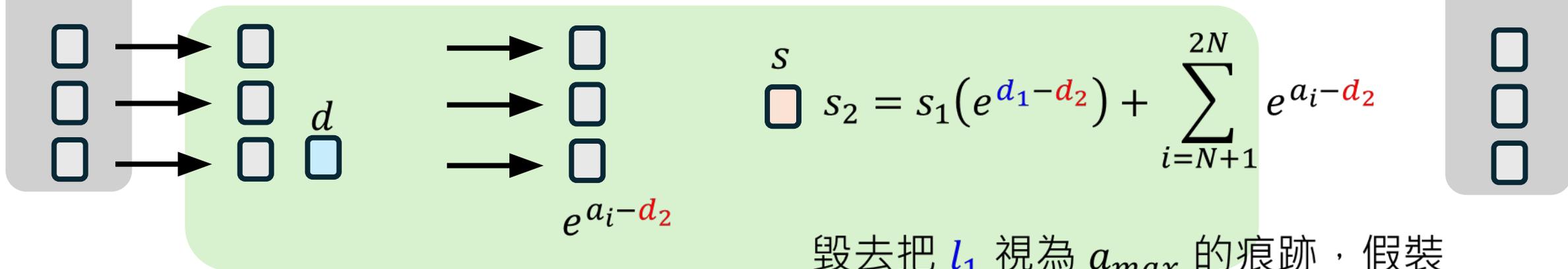
$$\hat{a}_i = \frac{e^{a_i - a_{max}}}{\sum_{i=1}^L e^{a_i - a_{max}}} \hat{a}_i$$



$$d_1 = \max(a_1, \dots, a_N)$$

假設 d_1 為 a_{max}

假設 d_1 為 a_{max}

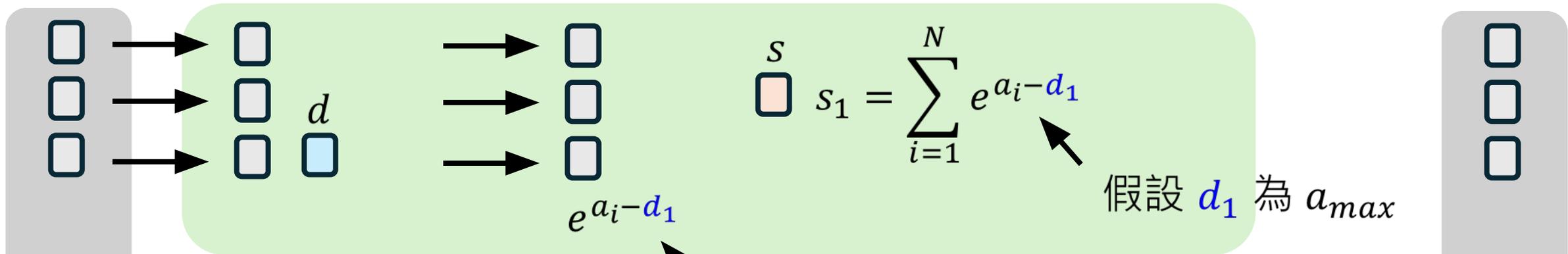


$$d_2 = \max(a_{N+1}, \dots, a_{2N}, d_1)$$

毀去把 l_1 視為 a_{max} 的痕跡，假裝已經是把 l_2 視為 a_{max}

一次找出 a_{max} 和 $\sum_{i=1}^L e^{a_i - a_{max}}$

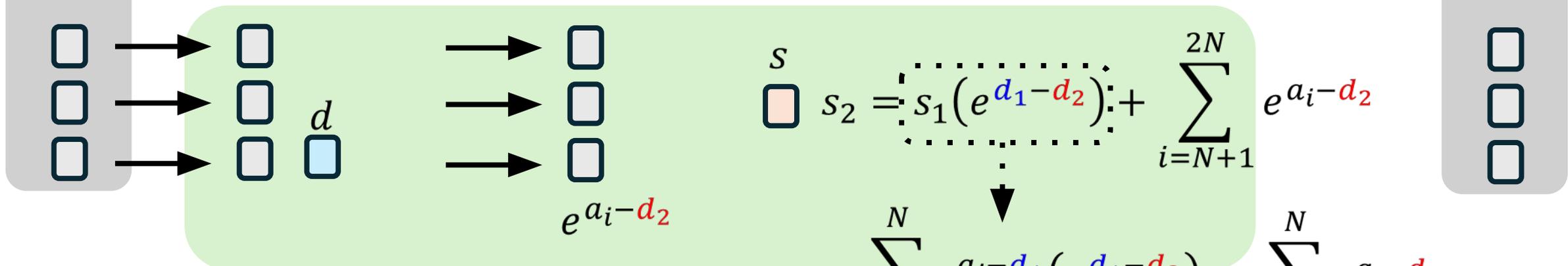
$$\hat{a}_i = \frac{e^{a_i - a_{max}}}{\sum_{i=1}^L e^{a_i - a_{max}}} \hat{a}_i$$



$$d_1 = \max(a_1, \dots, a_N)$$

假设 d_1 为 a_{max}

假设 d_1 为 a_{max}

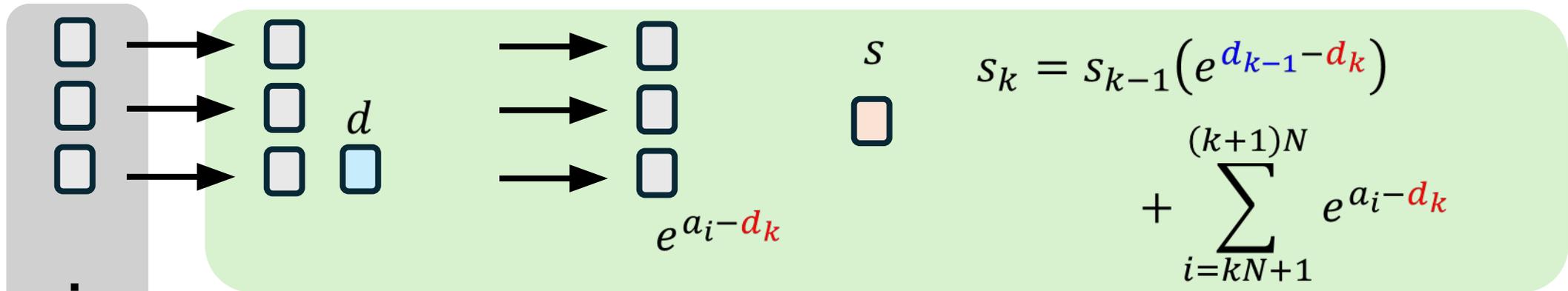


$$d_2 = \max(a_{N+1}, \dots, a_{2N}, d_1)$$

$$\sum_{i=1}^N e^{a_i - d_1} (e^{d_1 - d_2}) = \sum_{i=1}^N e^{a_i - d_2}$$

$$\hat{a}_i = \frac{e^{a_i - a_{max}}}{\sum_{i=1}^L e^{a_i - a_{max}}} \hat{a}_i$$

k-th chunk

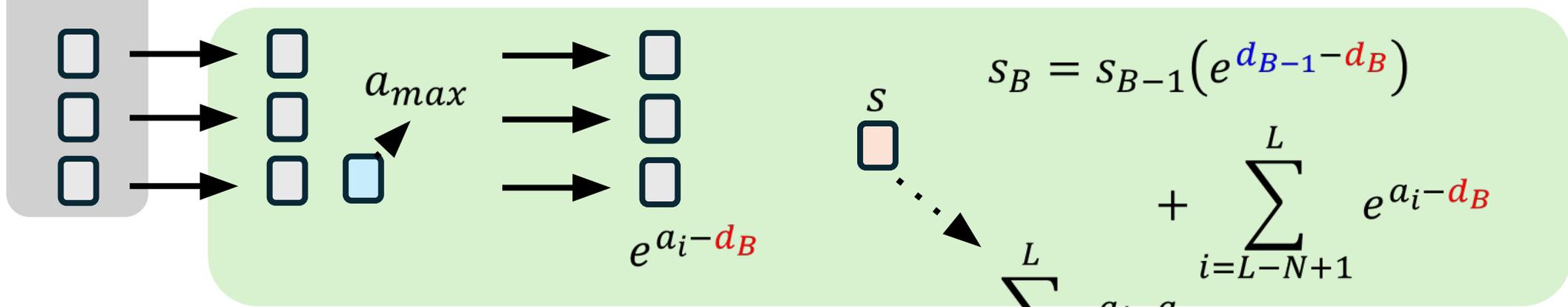


$$S_k = S_{k-1} (e^{d_{k-1} - d_k})$$

$$+ \sum_{i=kN+1}^{(k+1)N} e^{a_i - d_k}$$

$$d_k = \max(a_1, \dots, a_N, d_{k-1})$$

Last chunk (B-th chunk)

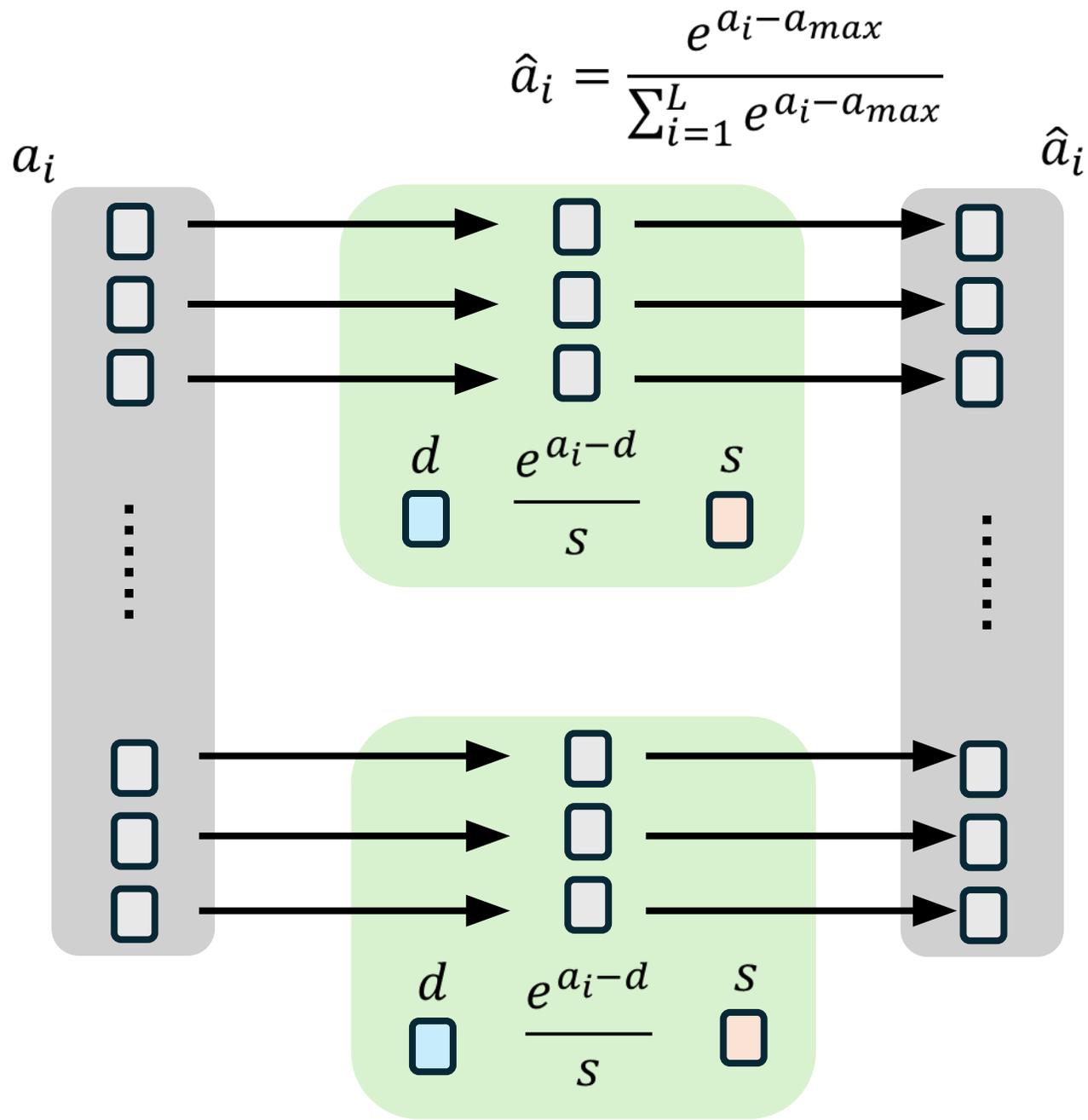
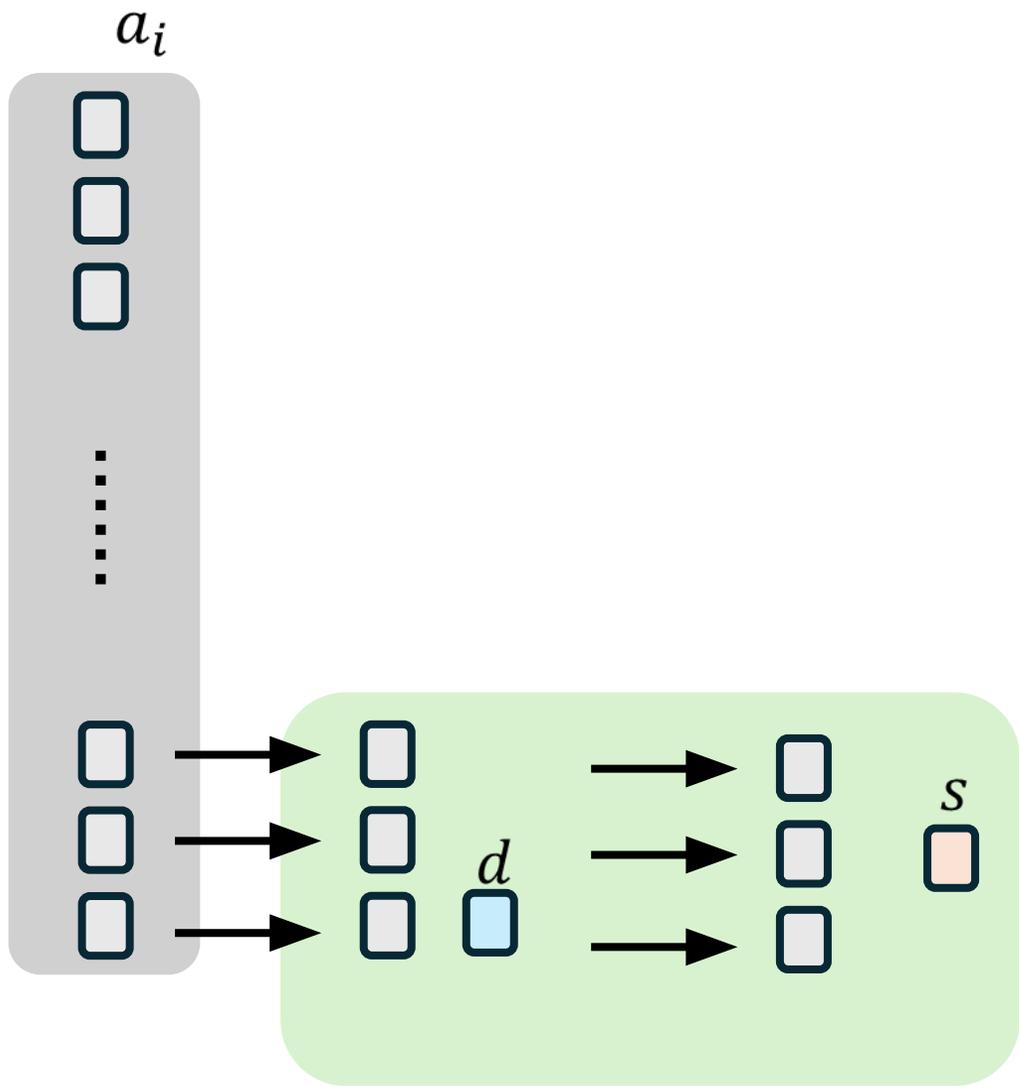


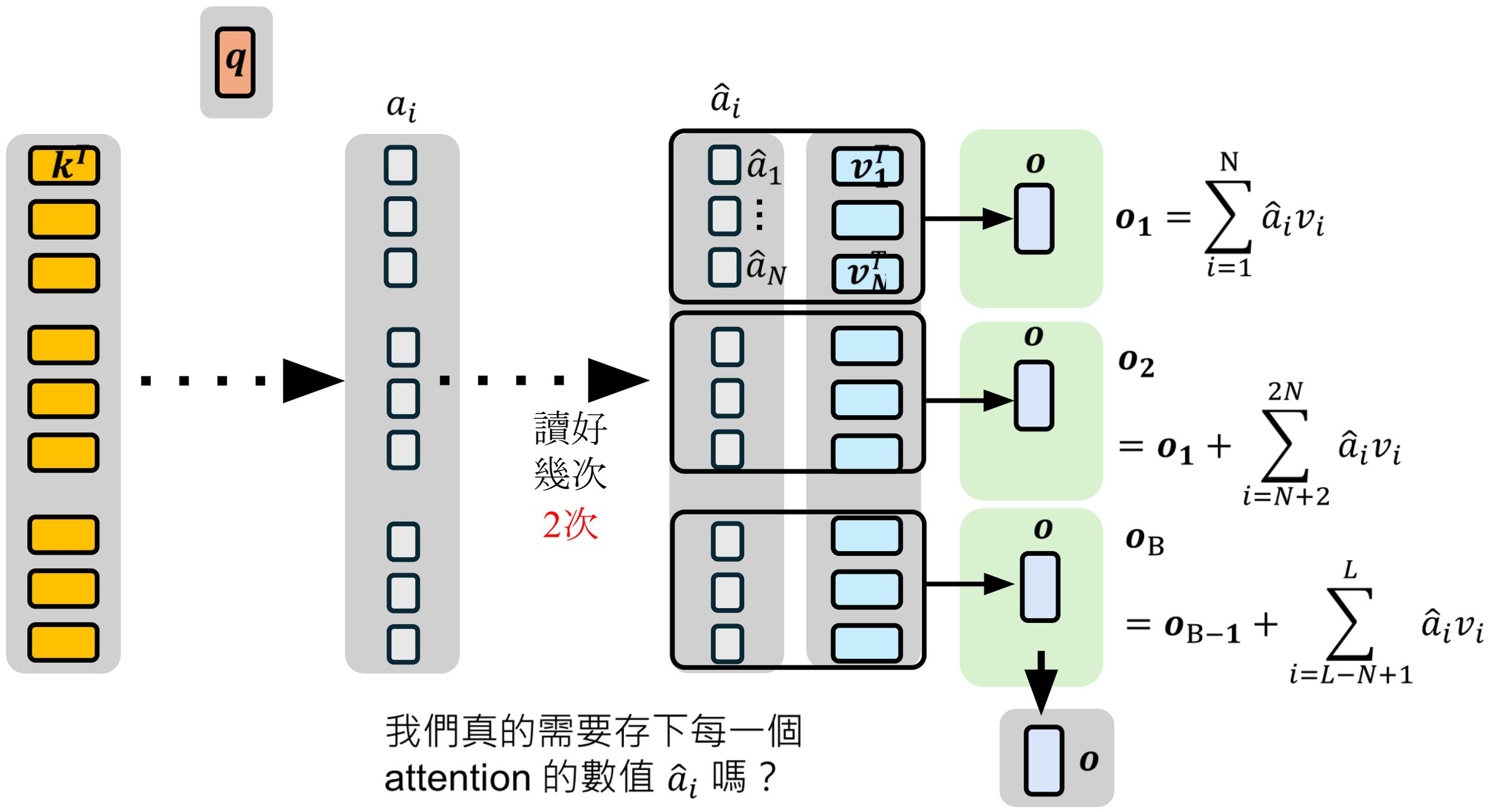
$$S_B = S_{B-1} (e^{d_{B-1} - d_B})$$

$$+ \sum_{i=L-N+1}^L e^{a_i - d_B}$$

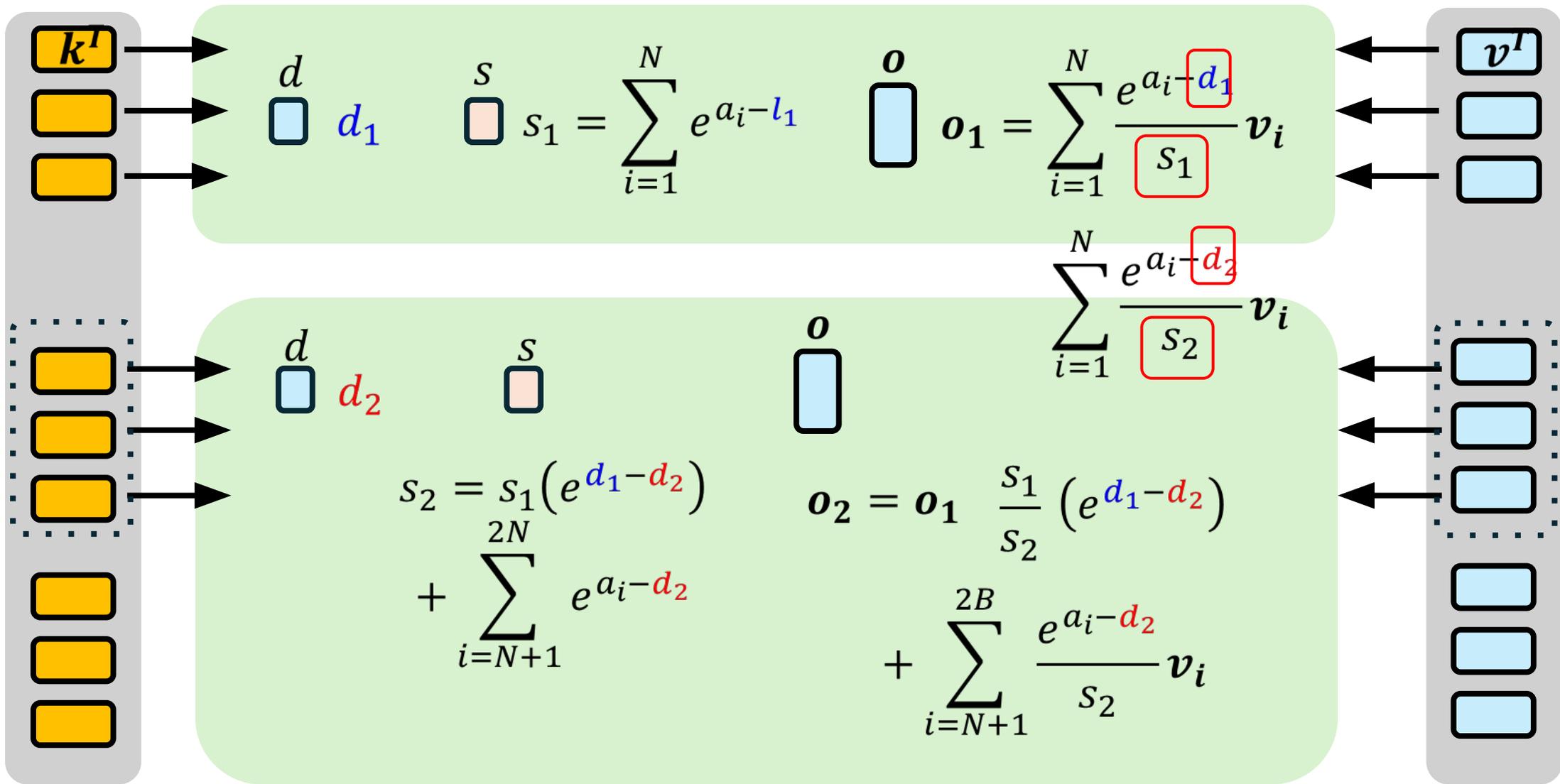
$$d_B = \max(a_{L-N+1}, \dots, a_L, d_{B-1})$$

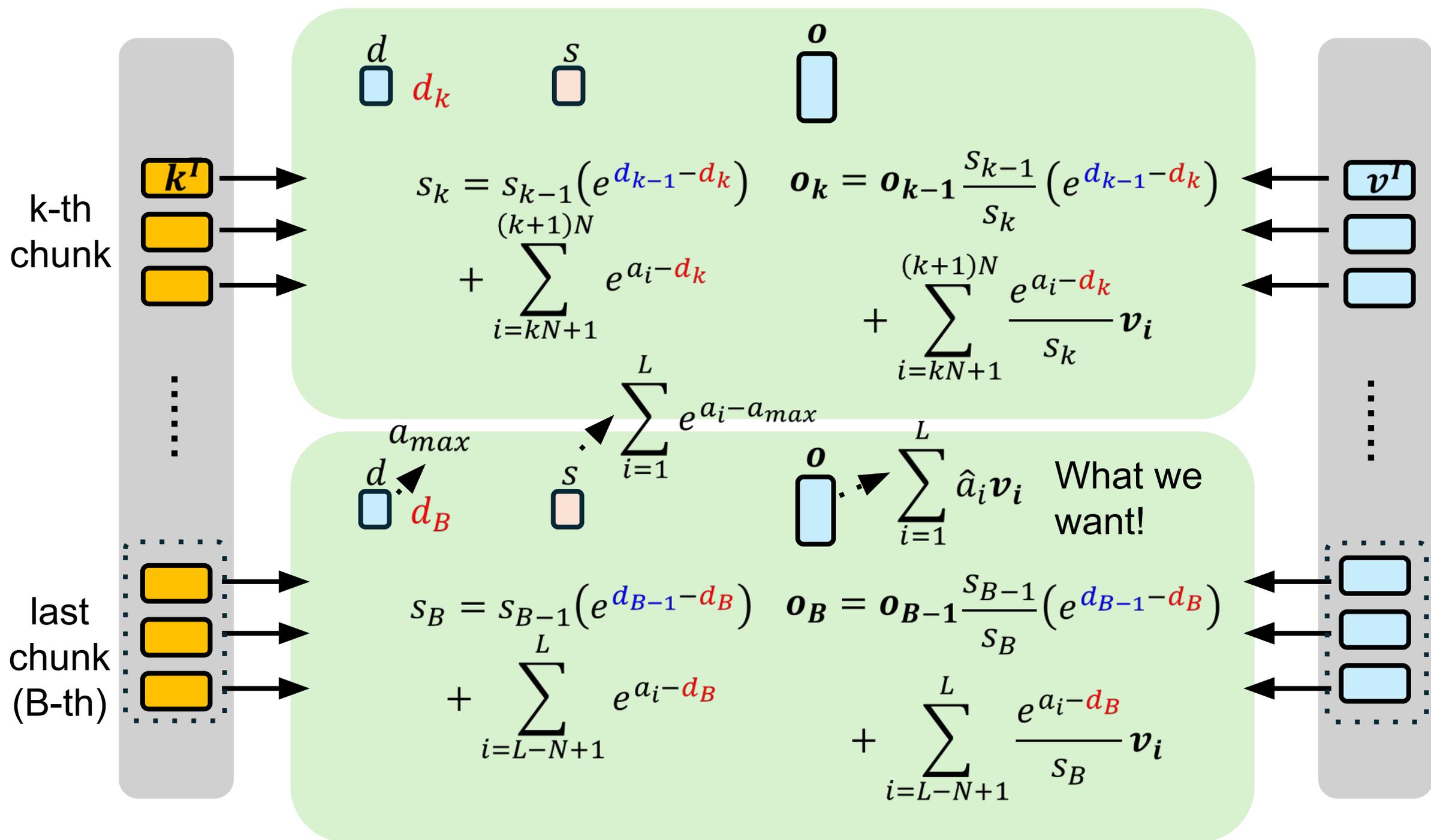
$$\sum_{i=1}^L e^{a_i - a_{max}}$$





我們真的需要存下每一個 attention 的數值 \hat{a}_i 嗎？



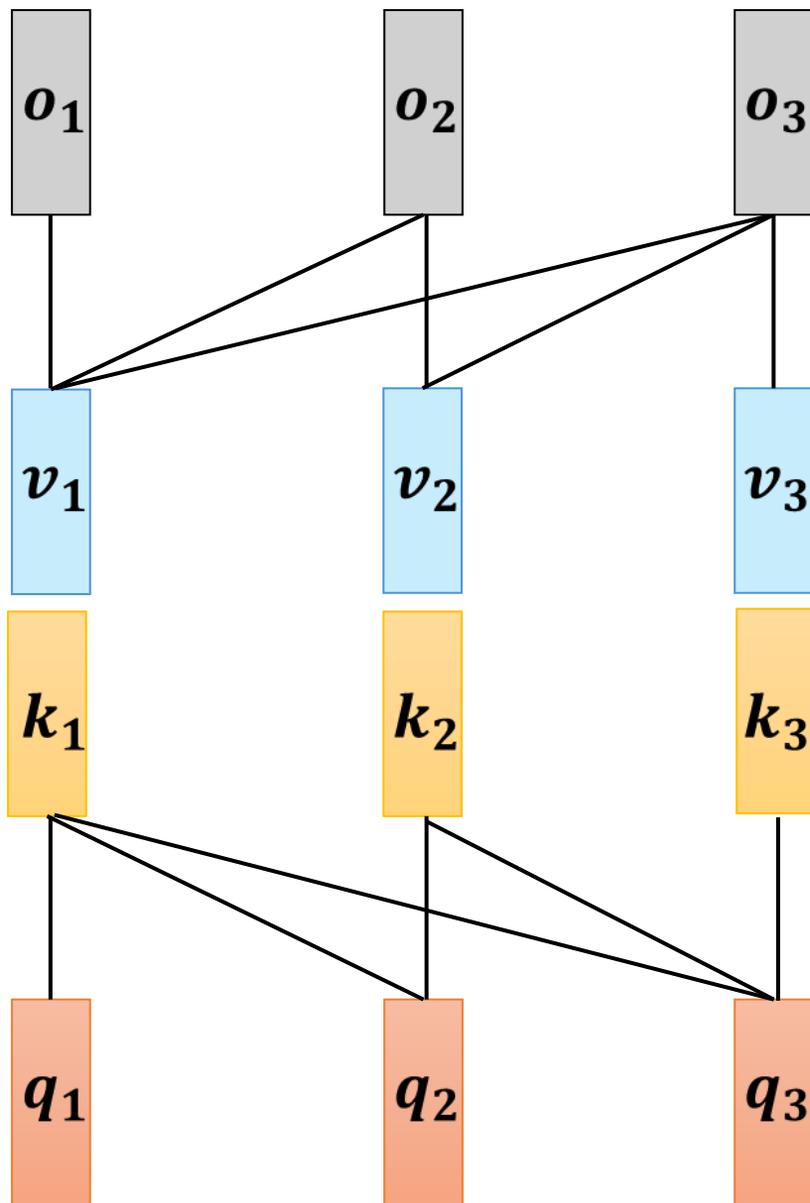


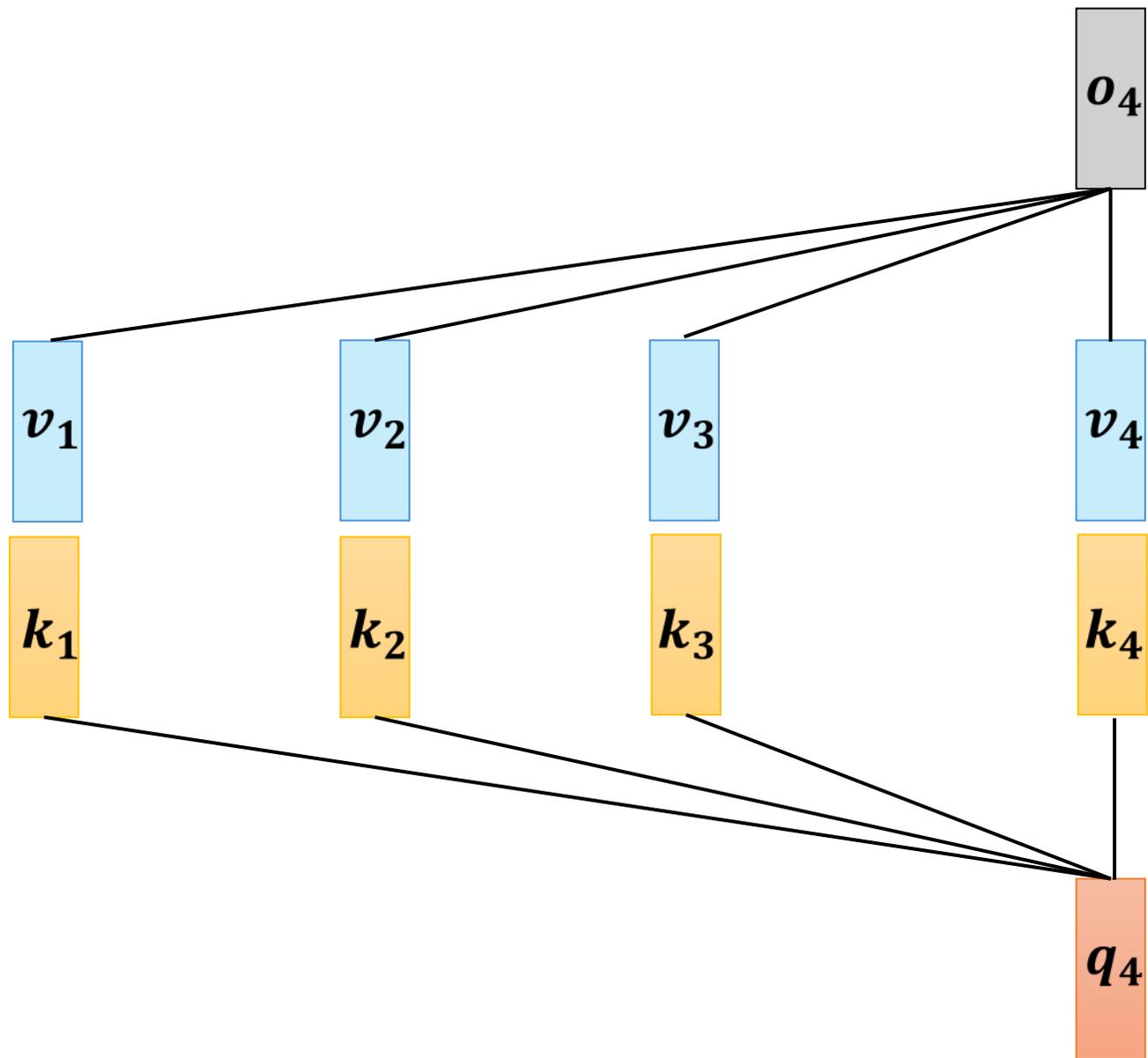
範例程式

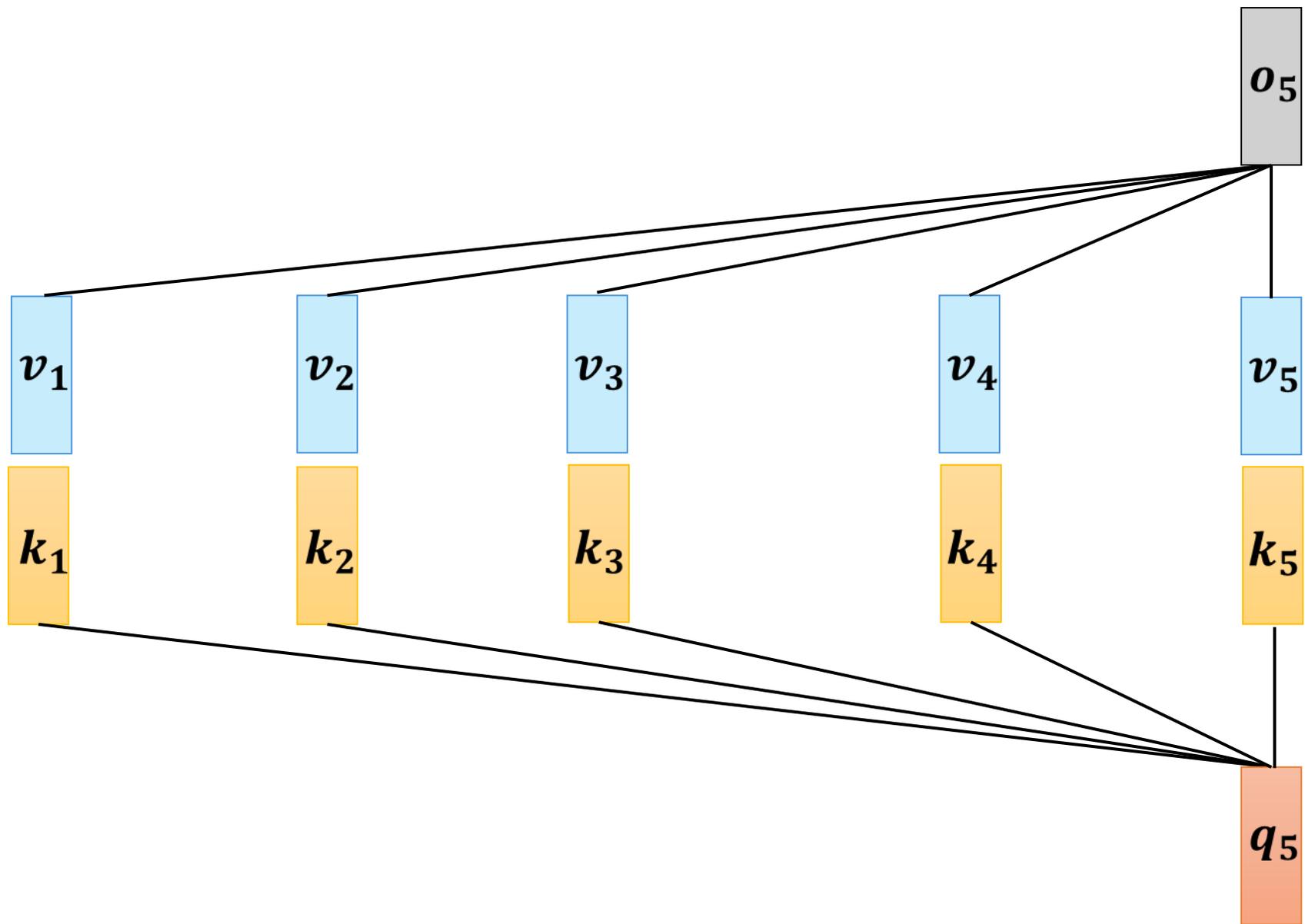
- https://colab.research.google.com/drive/1KoeKKIXSXI9b-pYg0kun3-uLQkP6p_hC?usp=sharing

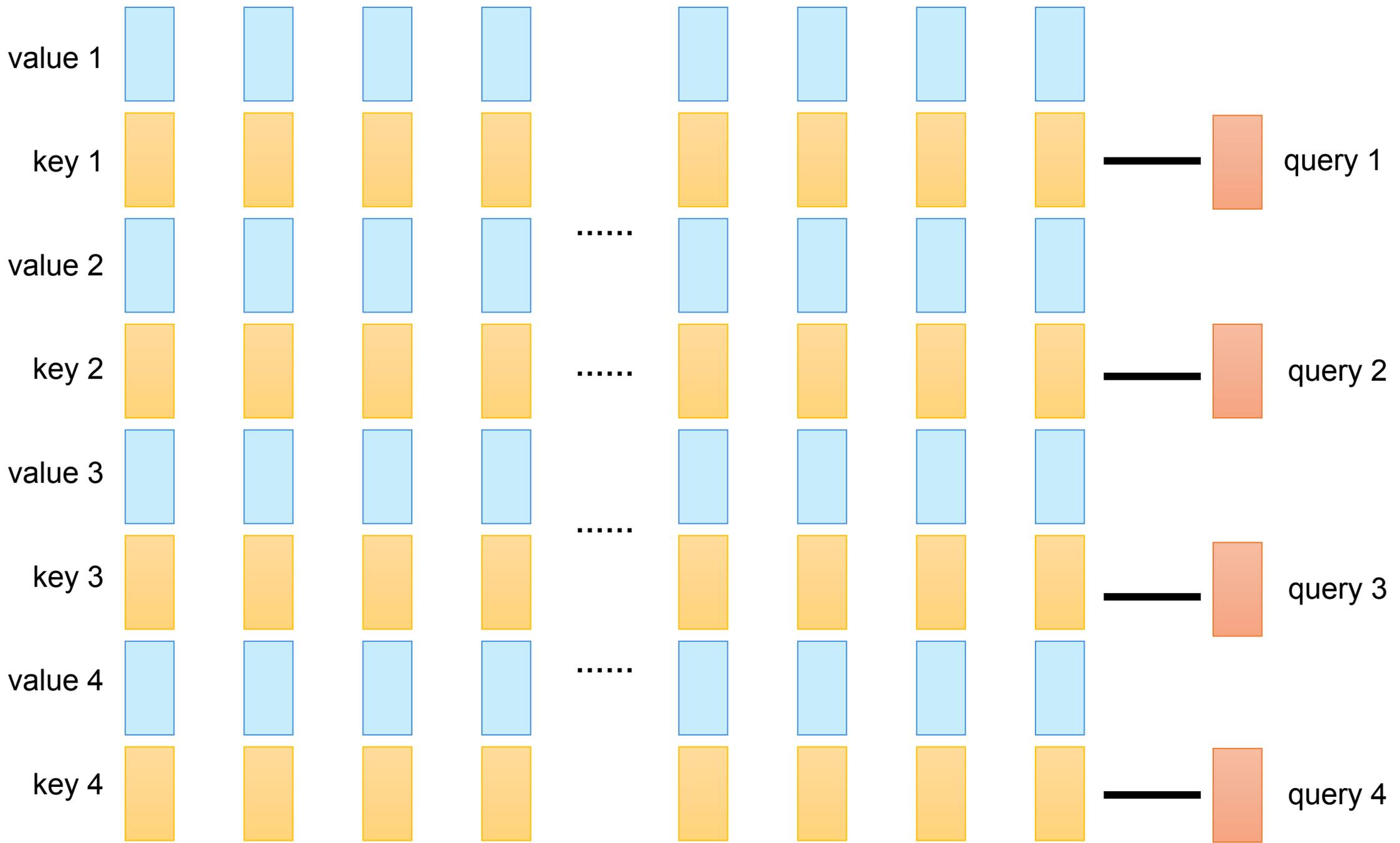


KV Cache











HBM



SRAM

Gemma 2

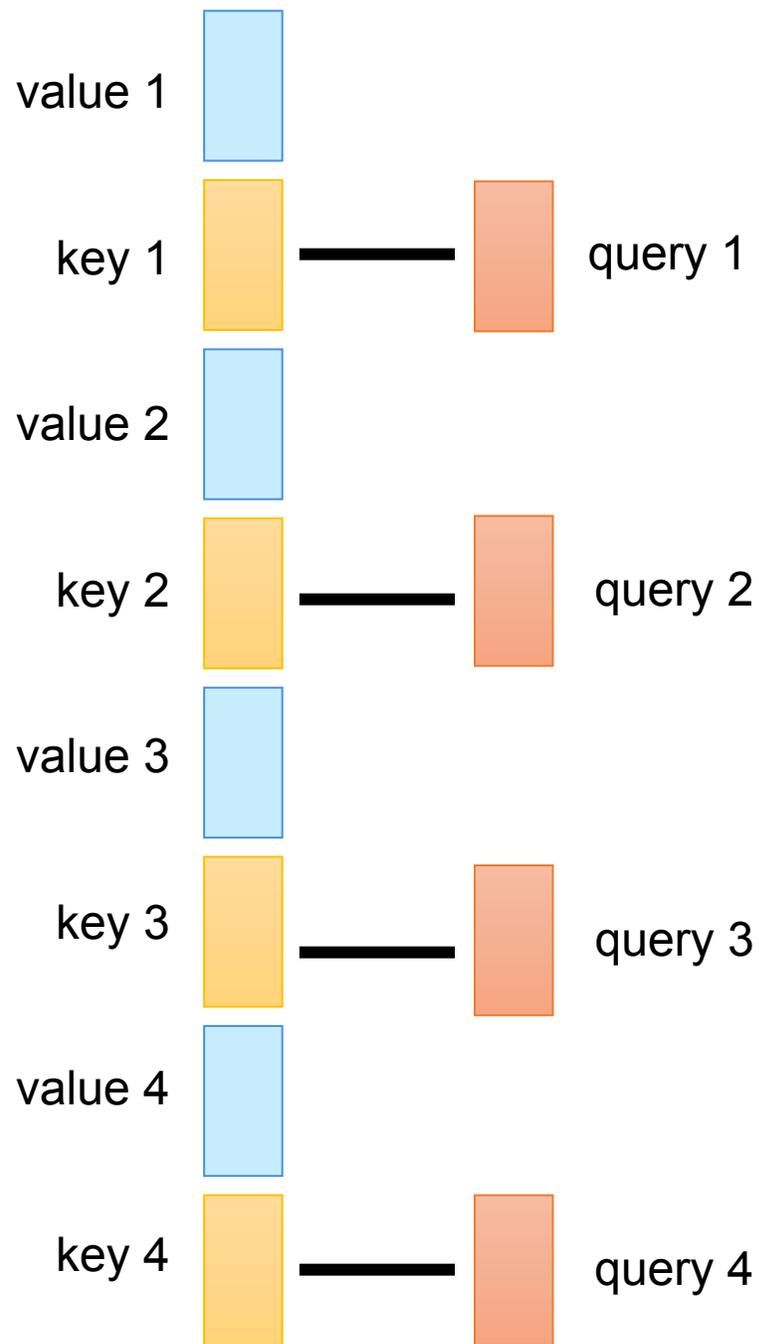
Parameters	2B	9B	27B
d_{model}	2304	3584	4608
Layers	26	42	46
Pre-norm	yes	yes	yes
Post-norm	yes	yes	yes
Non-linearity	GeGLU	GeGLU	GeGLU
Feedforward dim	18432	28672	73728
Head type	GQA	GQA	GQA
Num heads	8	16	32
Num KV heads	4	8	16
Head size	256	256	128
Global att. span	8192	8192	8192
Sliding window	4096	4096	4096
Vocab size	256128	256128	256128
Tied embedding	yes	yes	yes

Each token:

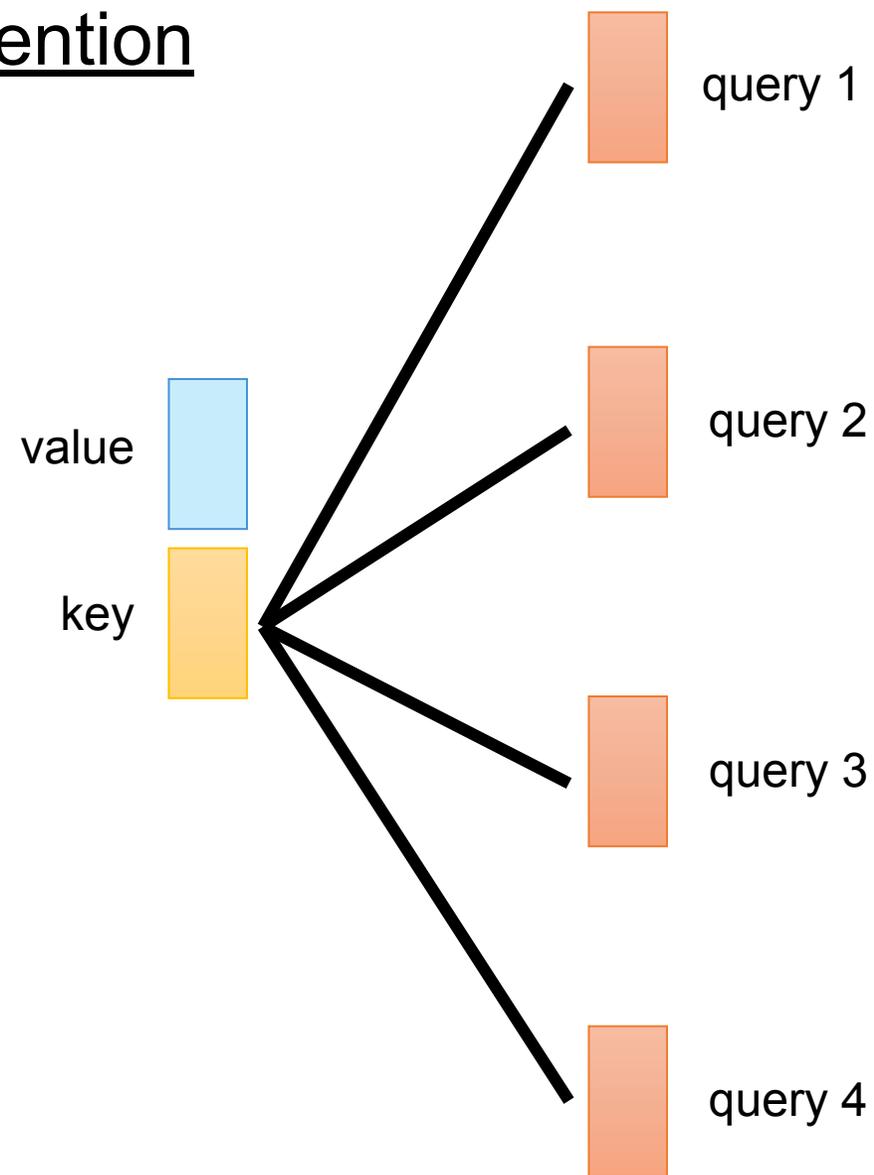
46 (layers) x 32 (heads) x 128
(dimension) x 2 (FP16) x 2 (value, key)
= 753664 bytes (736 KB, 0.72MB)

A100 (80GB):

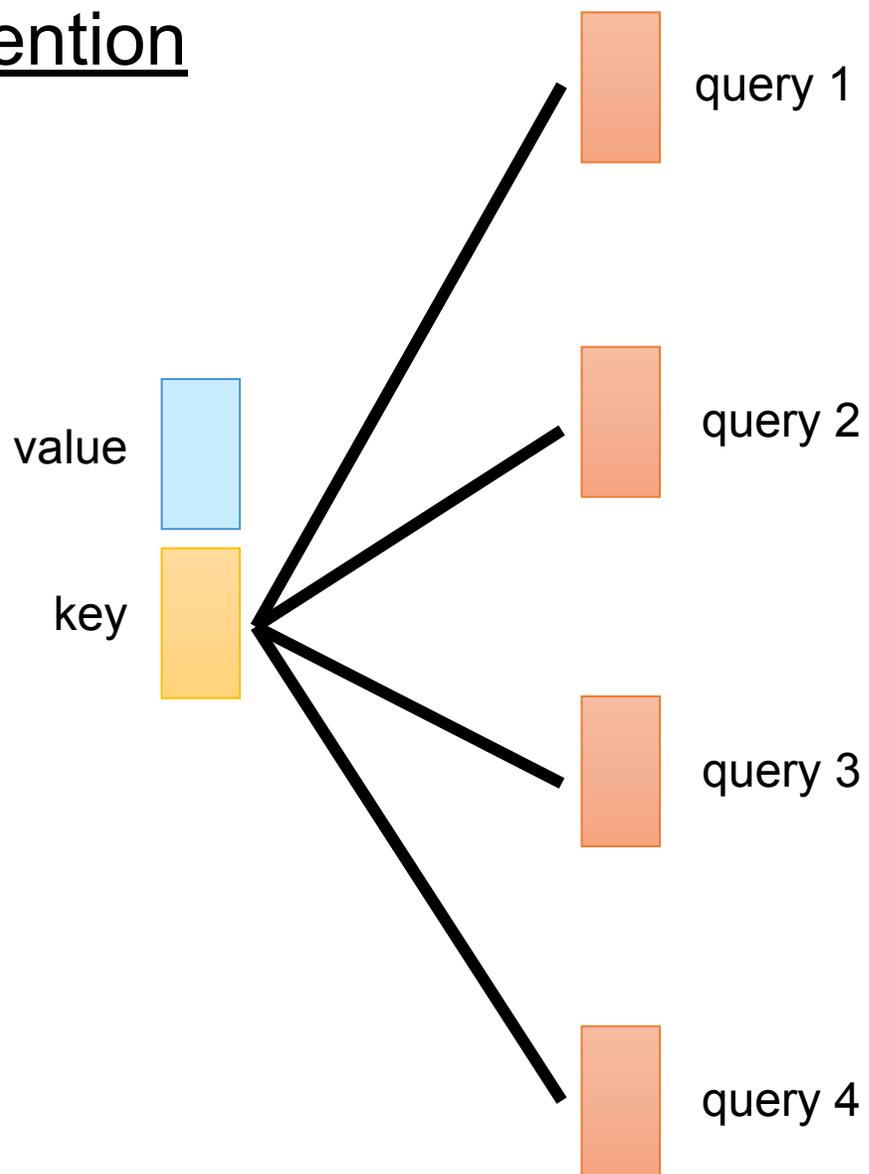
About 114k tokens



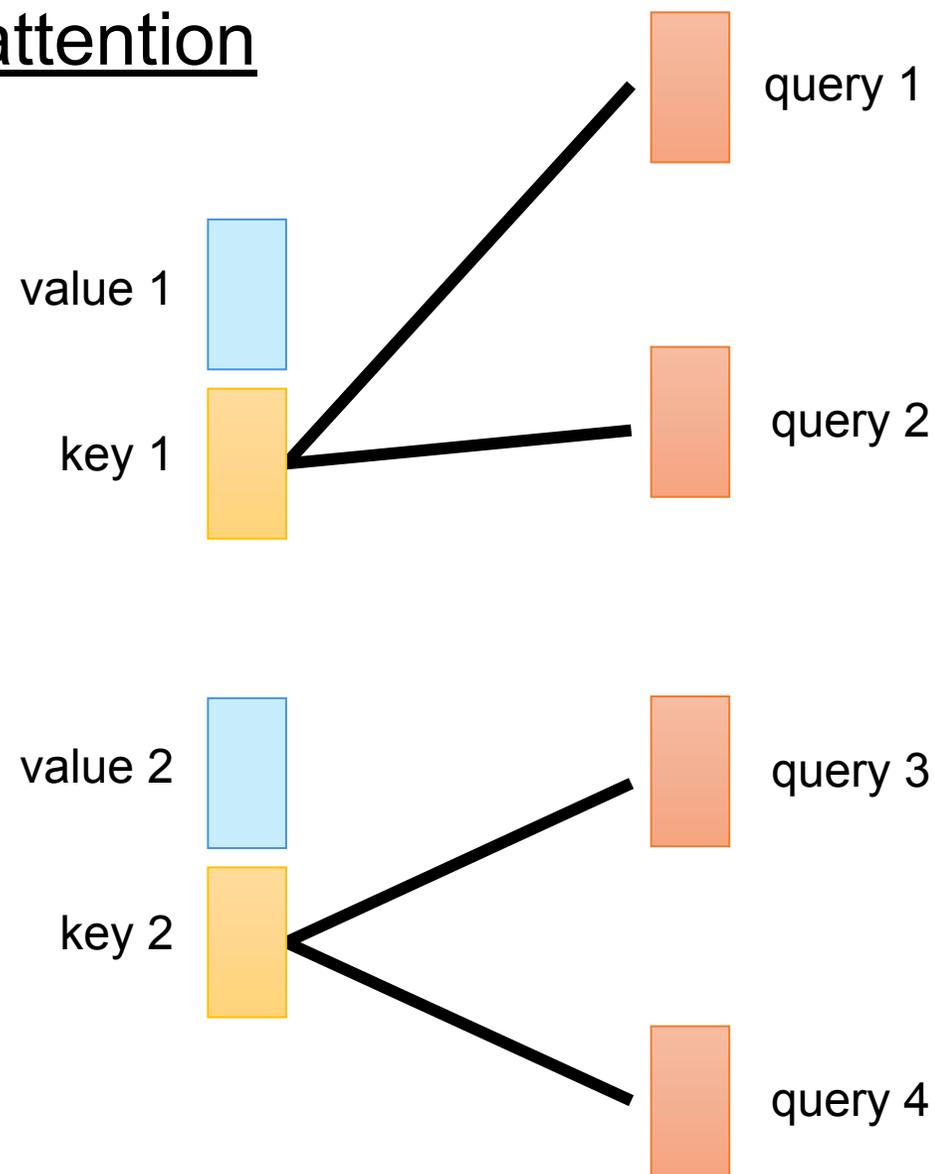
Multi-query attention



Multi-query attention

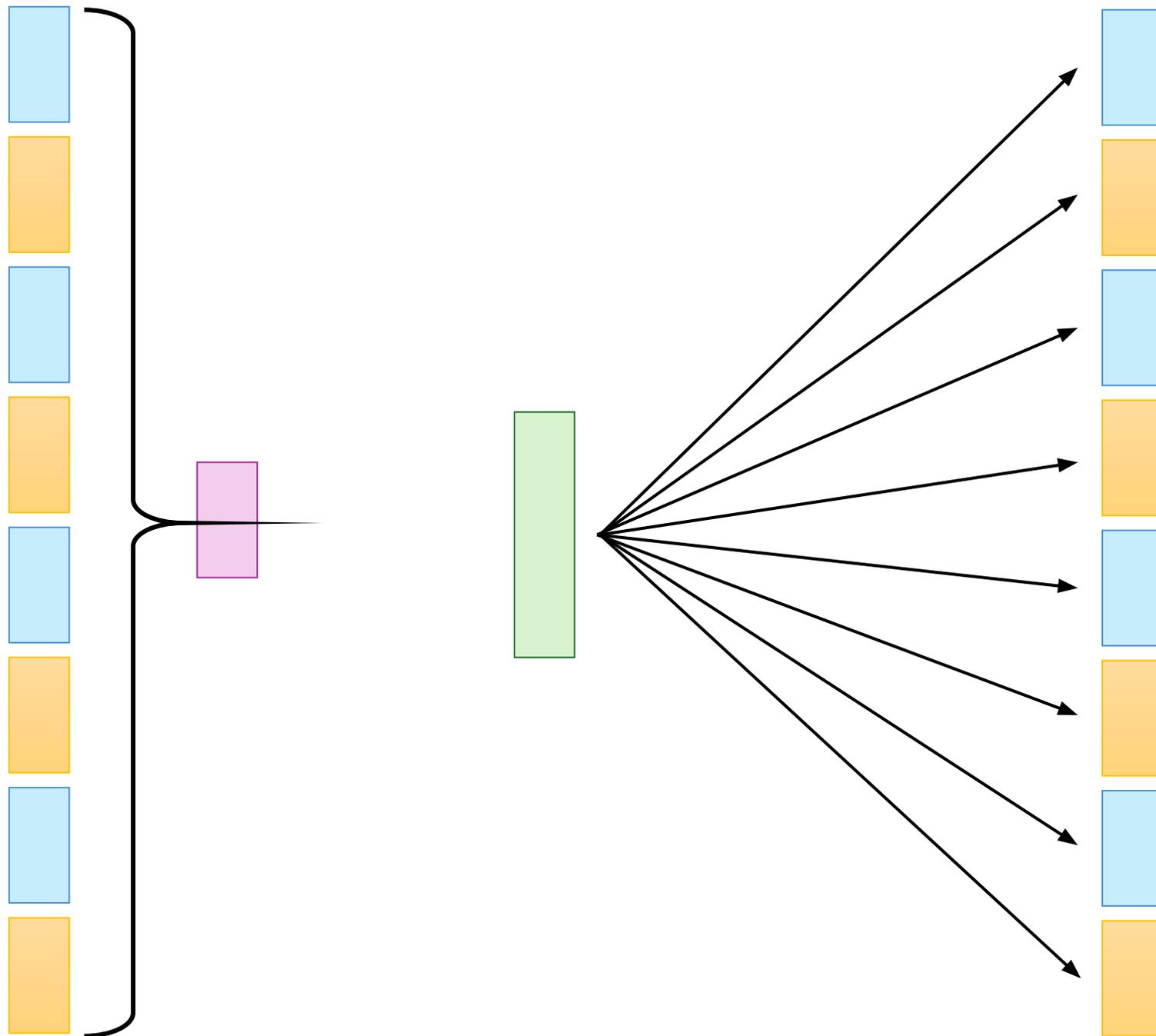


Group-query attention



Llama, Gemma, etc.

Multi-head
Latent
Attention

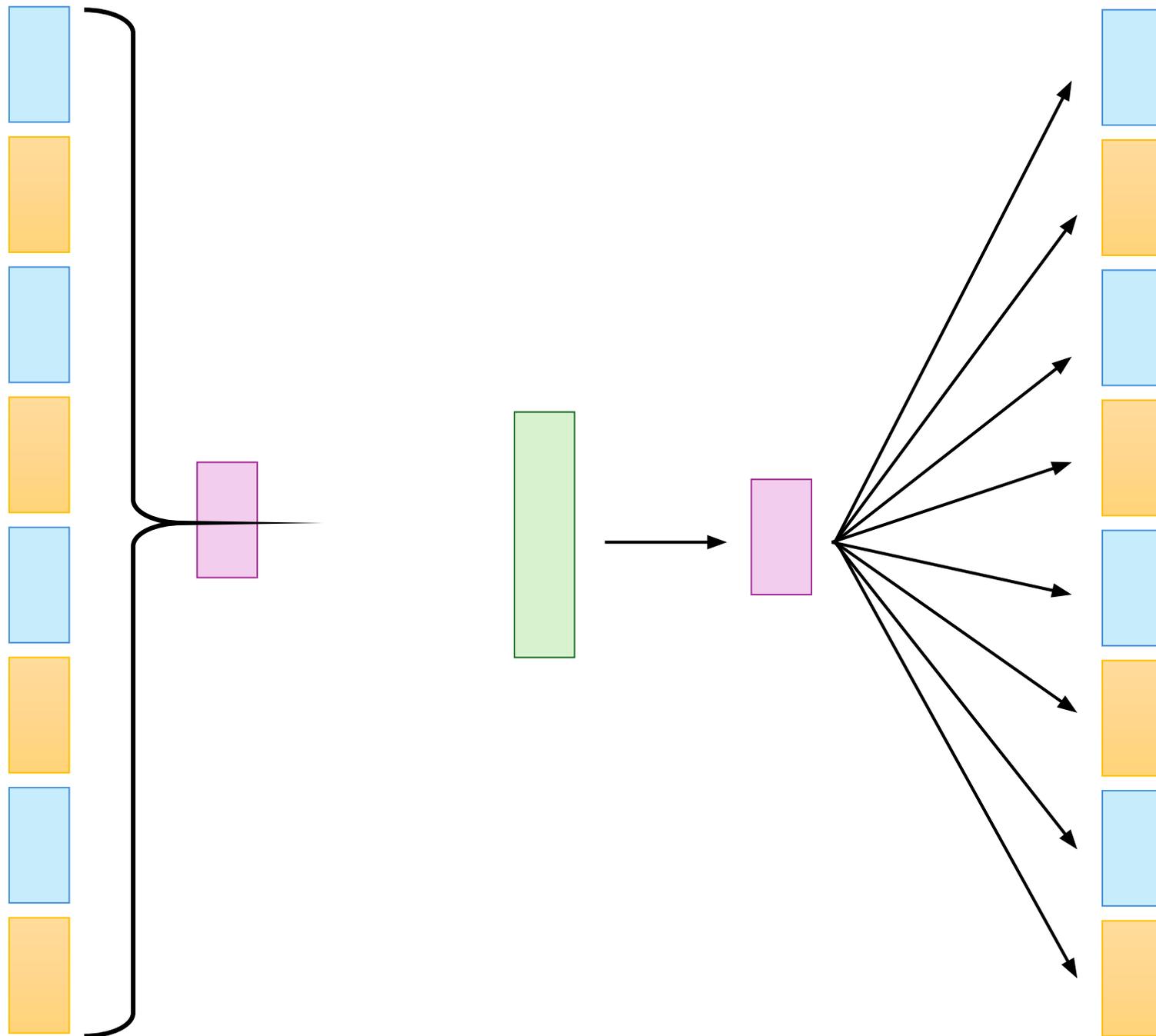


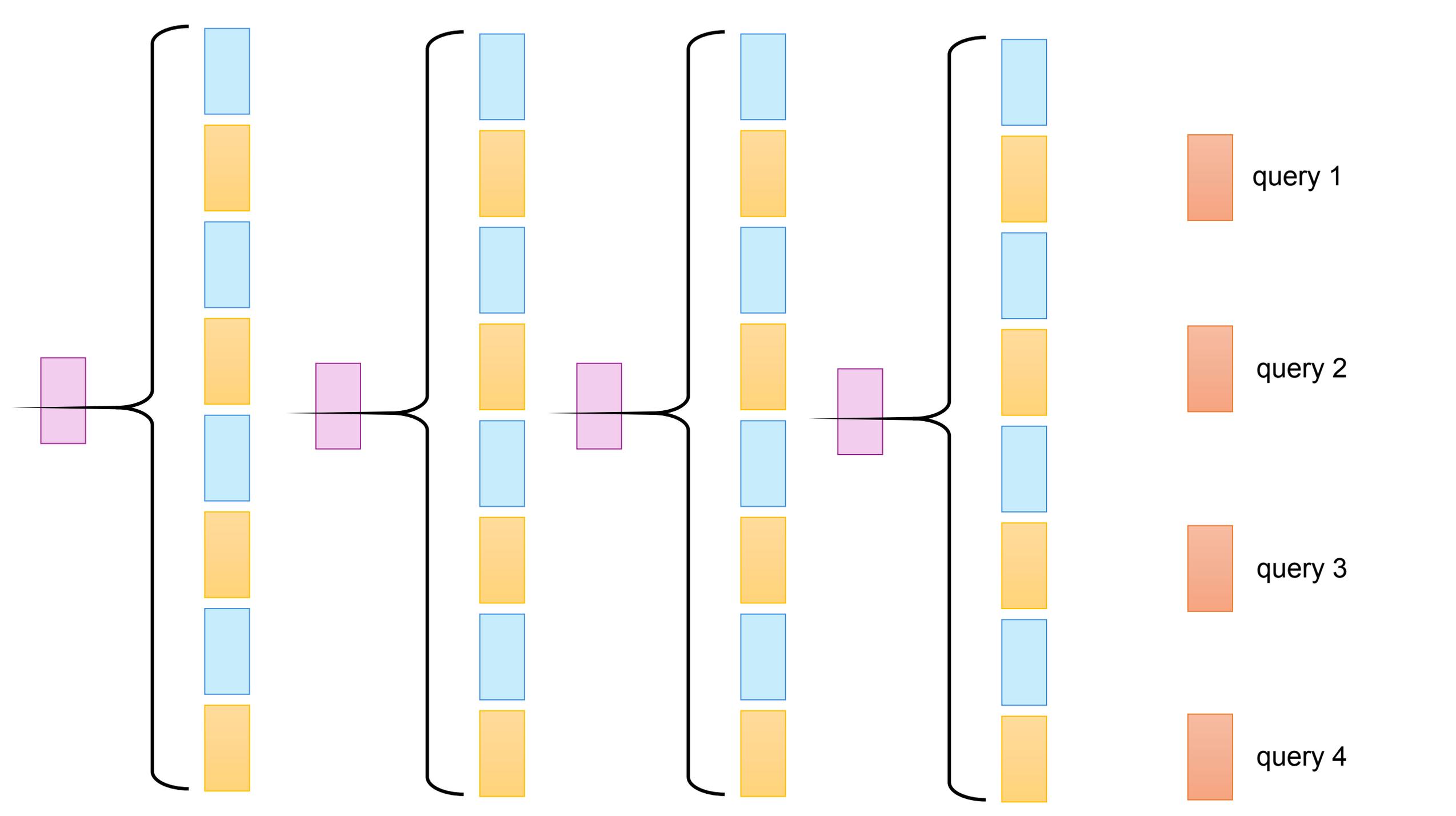
<https://arxiv.org/abs/2405.04434>

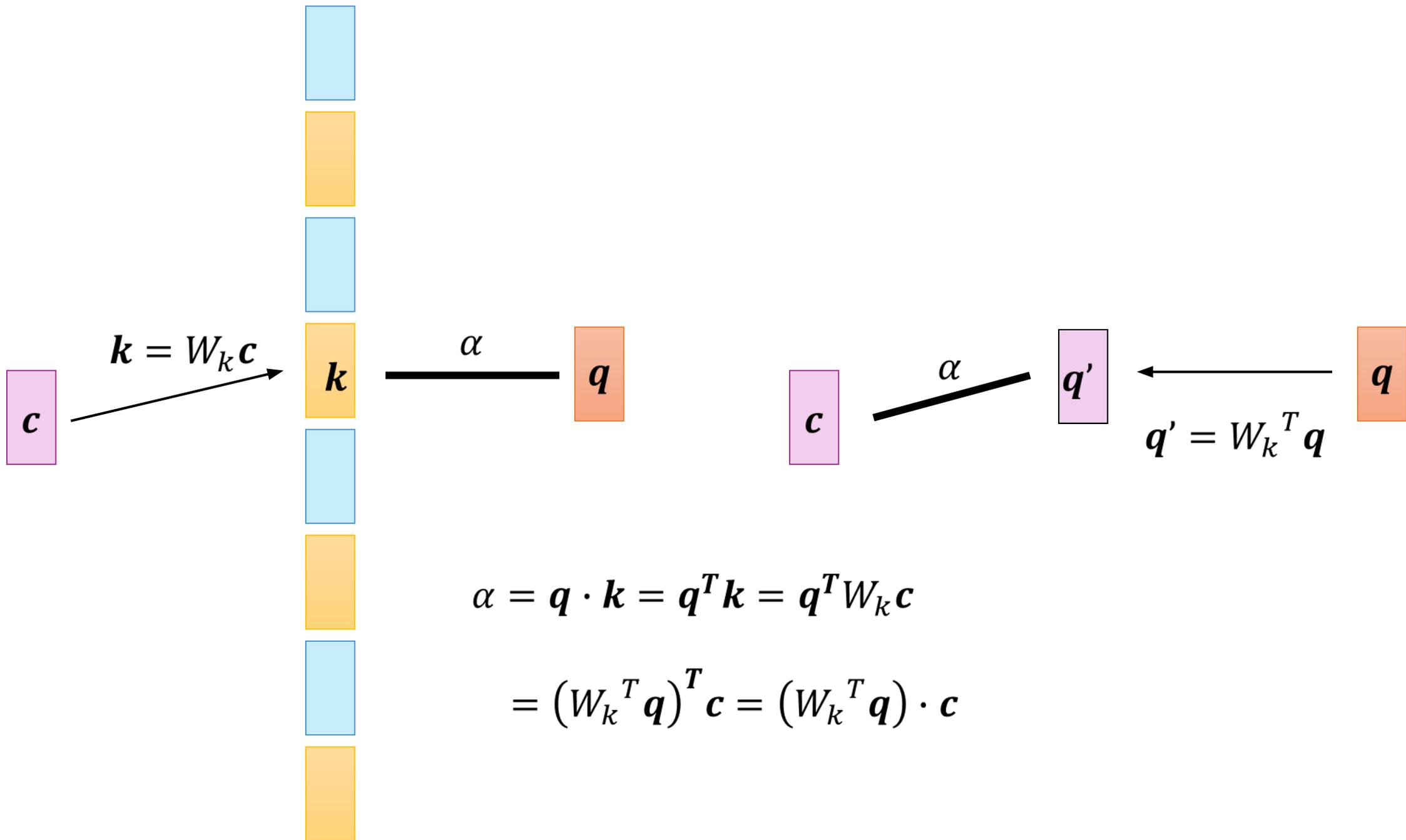
Multi-head
Latent
Attention

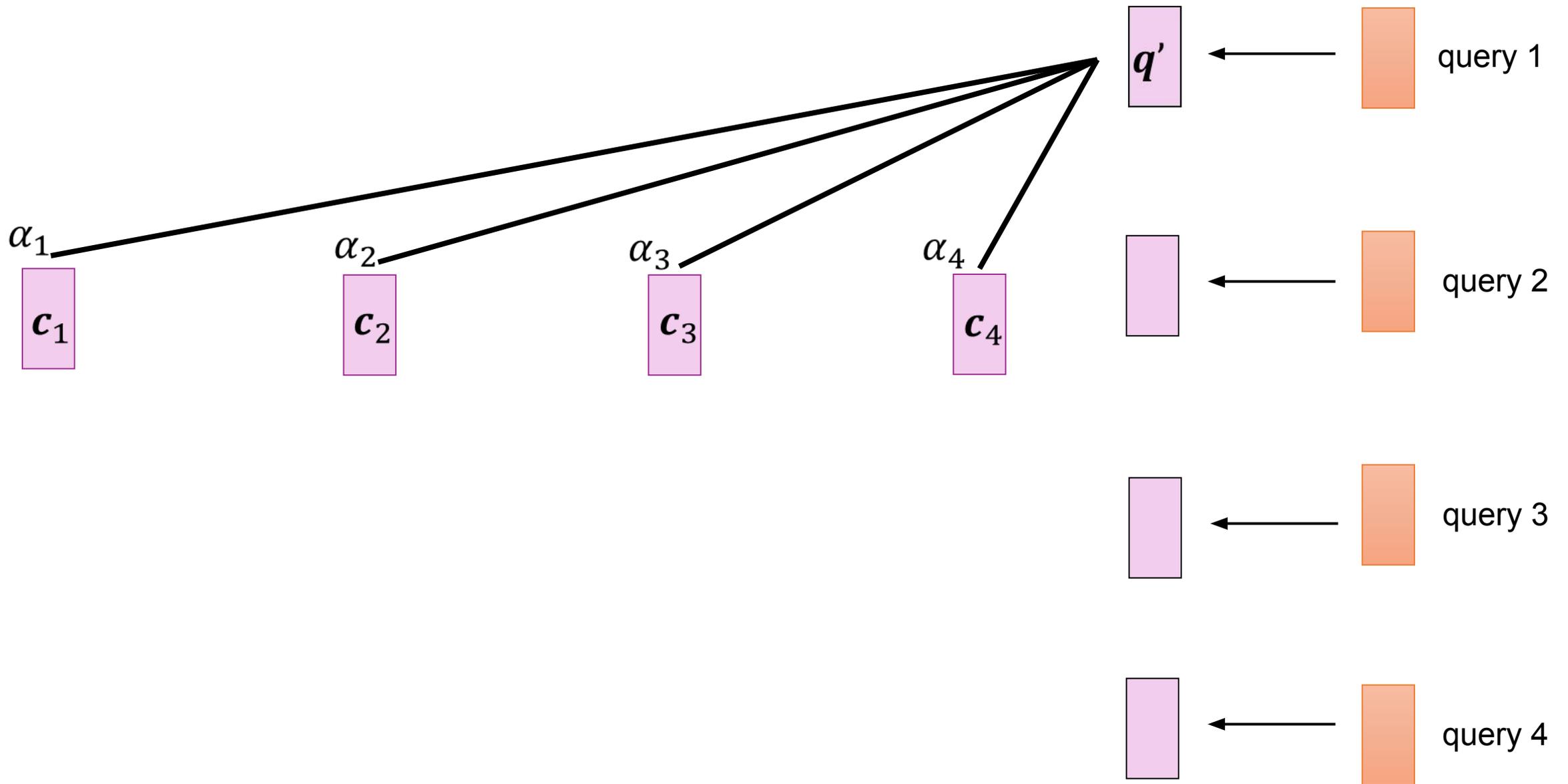
DeepSeek

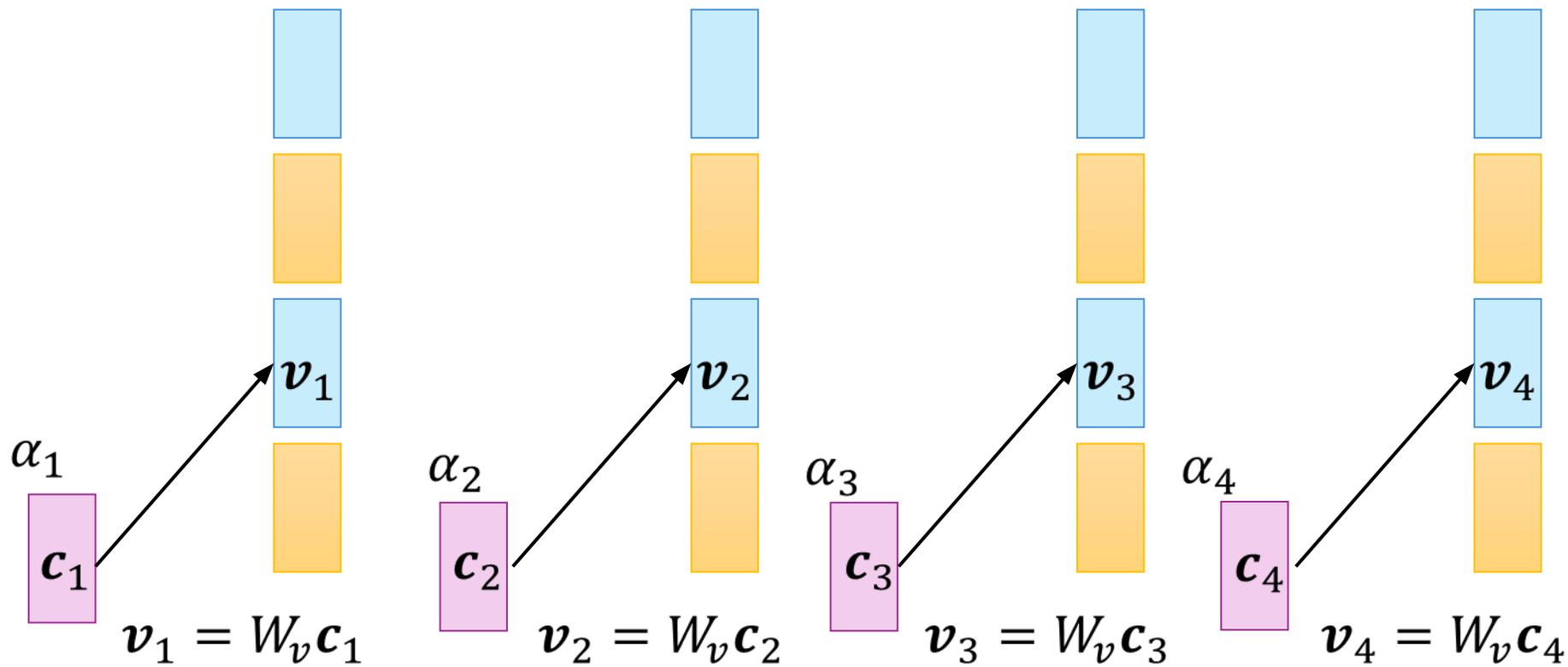
<https://arxiv.org/abs/2405.04434>





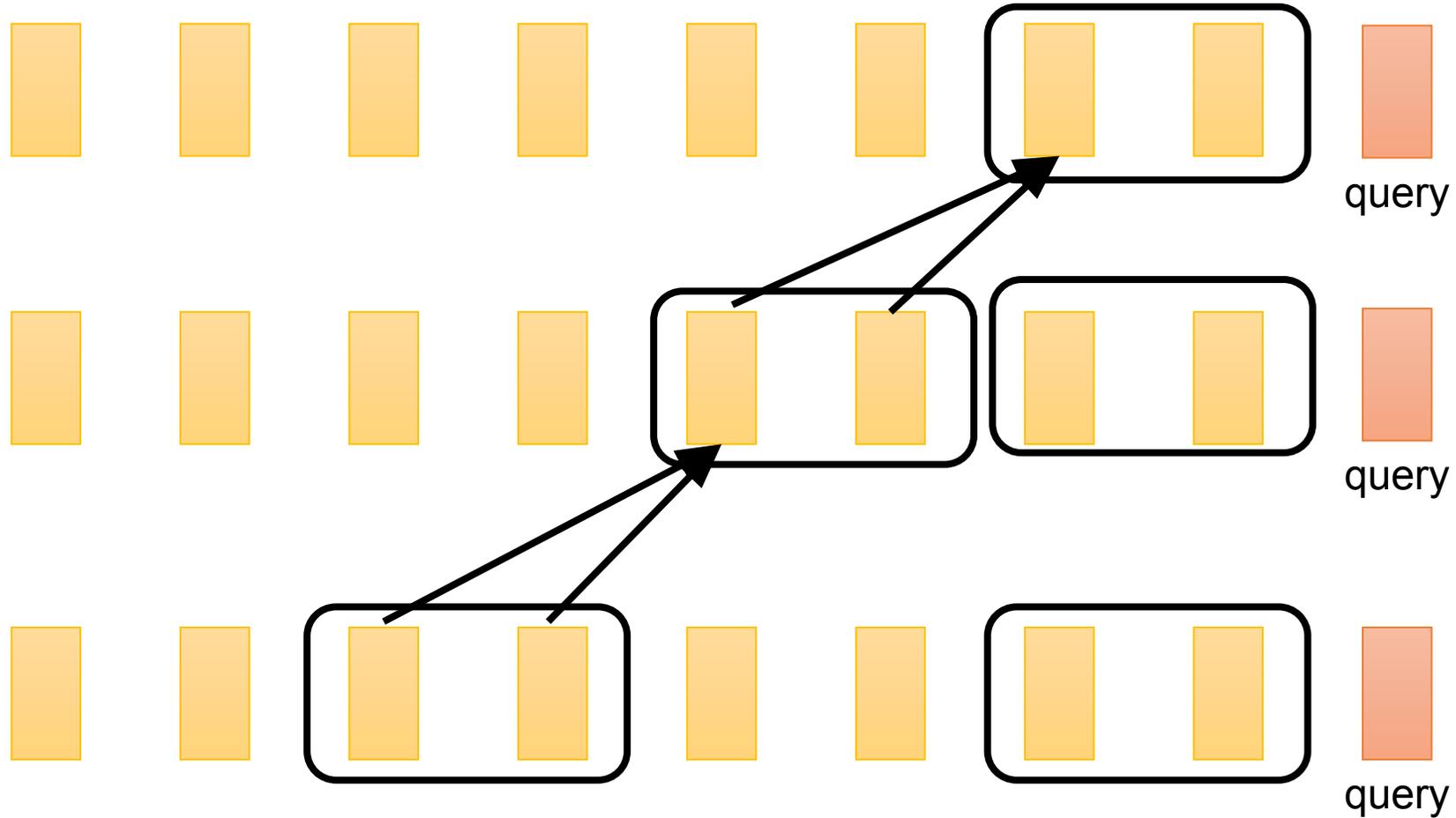






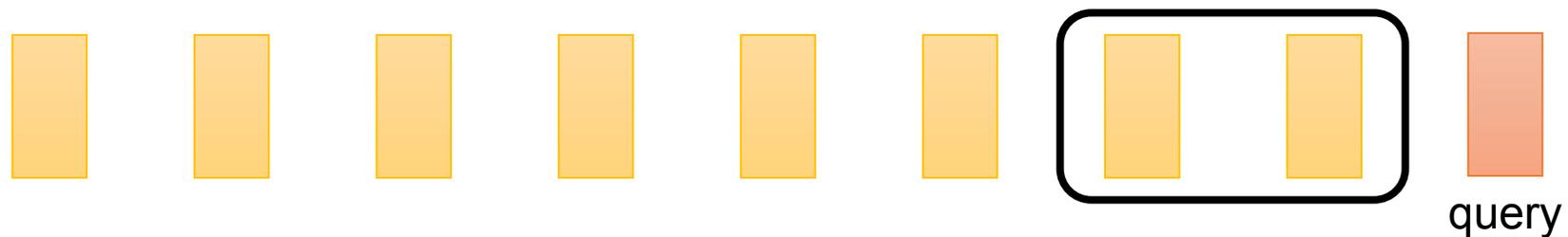
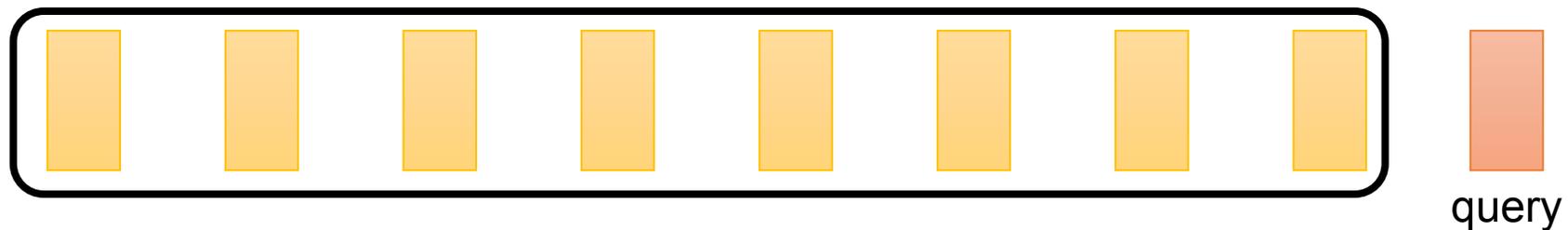
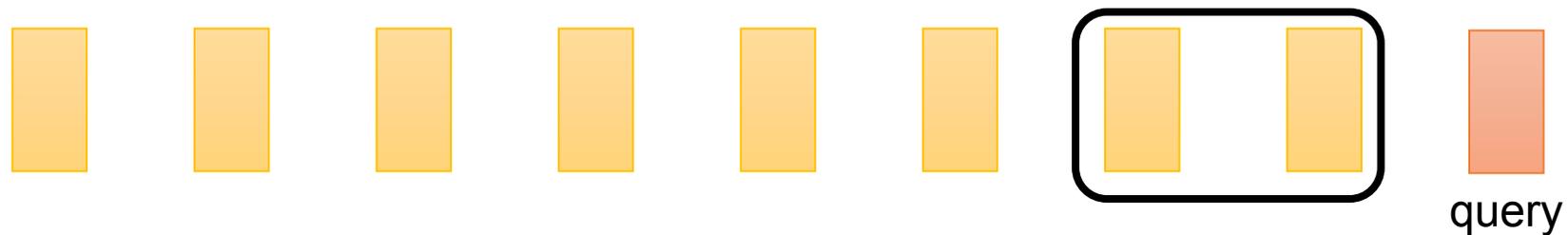
$$\begin{aligned}\mathbf{o} &= \alpha_1 \mathbf{v}_1 + \alpha_2 \mathbf{v}_2 + \alpha_3 \mathbf{v}_3 + \alpha_4 \mathbf{v}_4 \\ &= \alpha_1 W_v \mathbf{c}_1 + \alpha_2 W_v \mathbf{c}_2 + \alpha_3 W_v \mathbf{c}_3 + \alpha_4 W_v \mathbf{c}_4 \\ &= W_v (\alpha_1 \mathbf{c}_1 + \alpha_2 \mathbf{c}_2 + \alpha_3 \mathbf{c}_3 + \alpha_4 \mathbf{c}_4)\end{aligned}$$

Sliding Window Attention



Mistral 7B

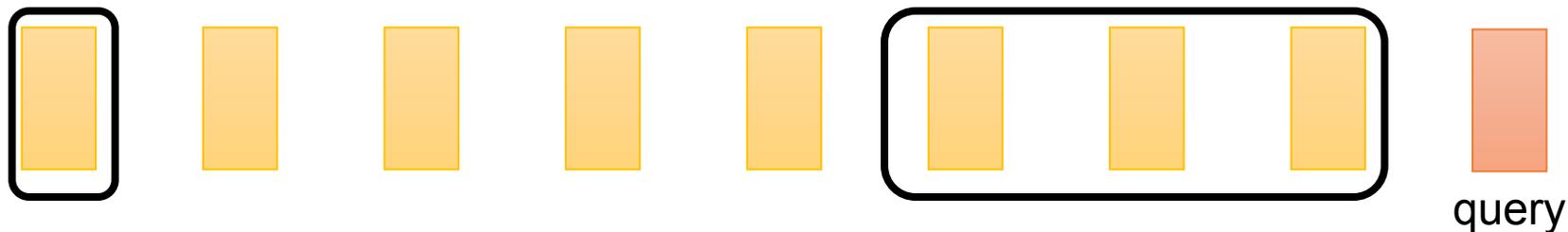
Sliding Window Attention



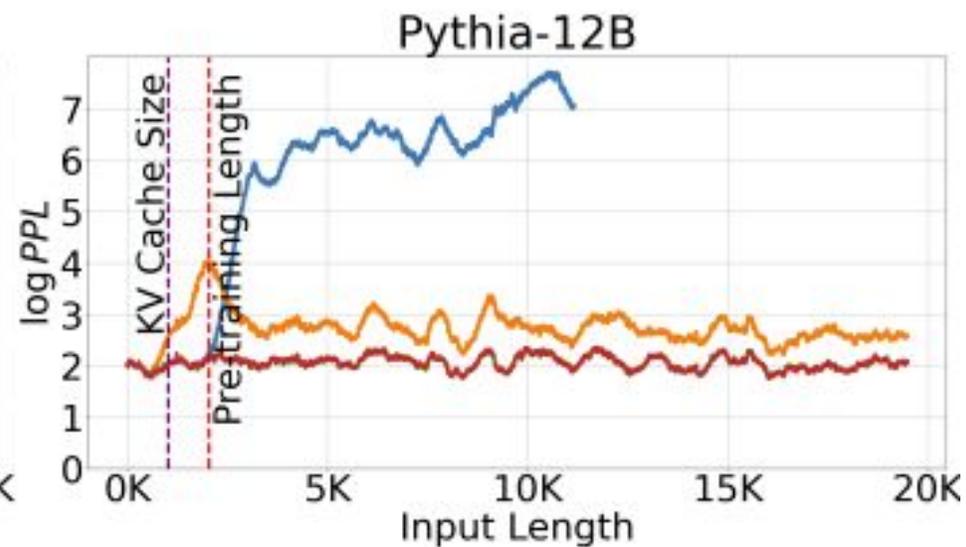
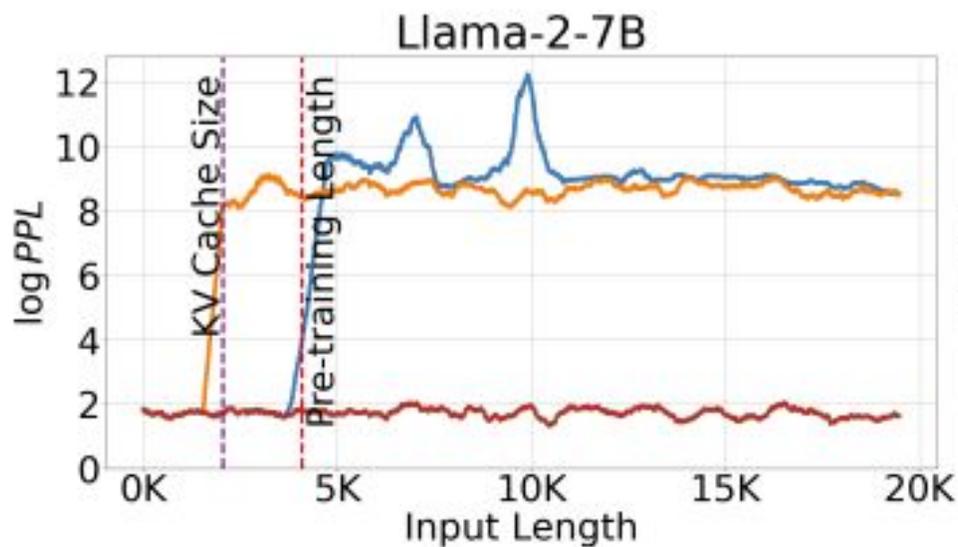
GPT-OSS

Streaming LLM

<https://arxiv.org/abs/2309.17453>



- Dense Attention
- Window Attention
- StreamingLLM

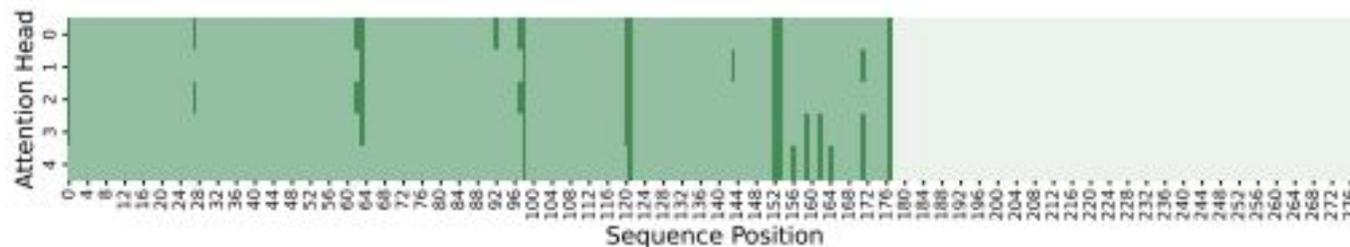


Pruning KV Cache

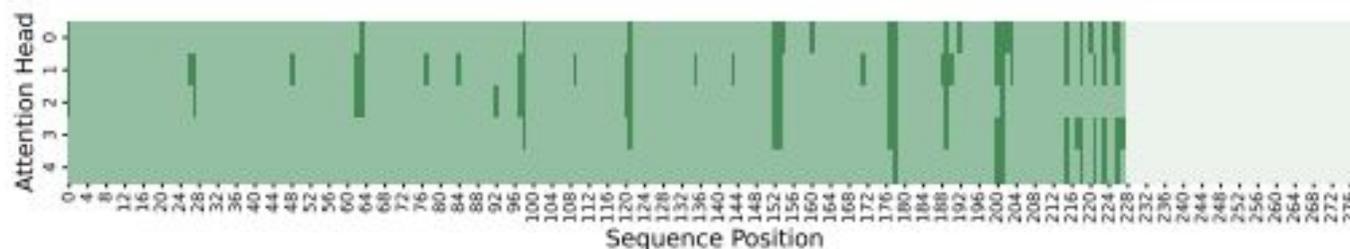
Scissorhands: <https://arxiv.org/abs/2305.17118>

H2O: <https://arxiv.org/abs/2306.14048>

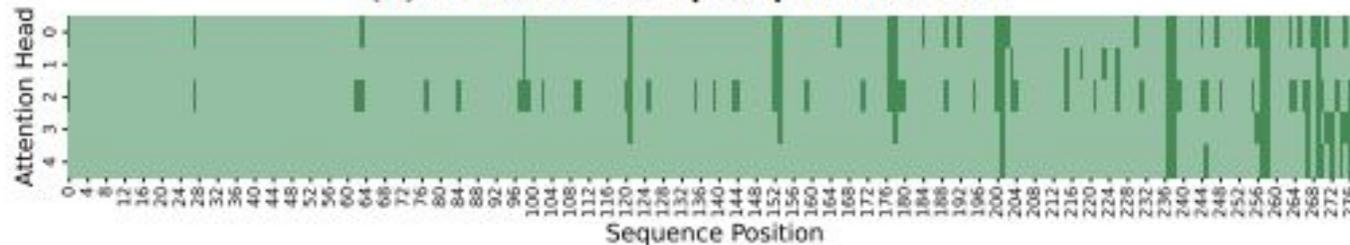
- 每次只有小部分的 Tokens 有 Attention
- 少數 token 會反覆吸走大量 attention



(a) Attention map at position 178

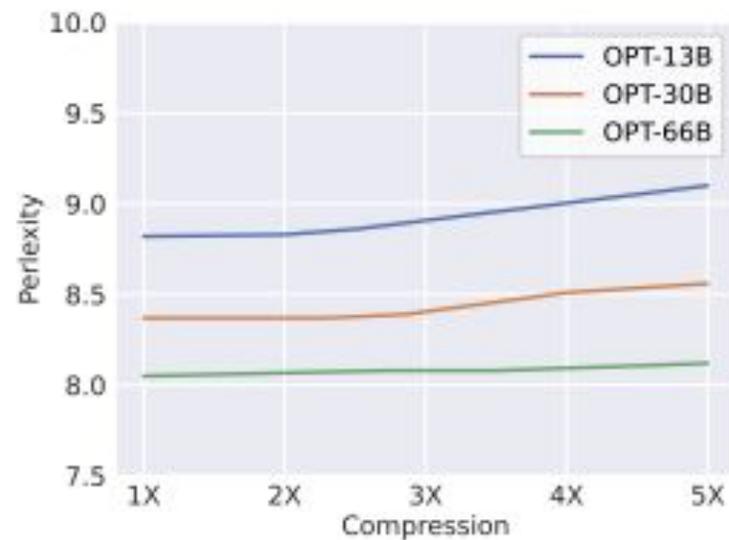


(b) Attention map at position 228

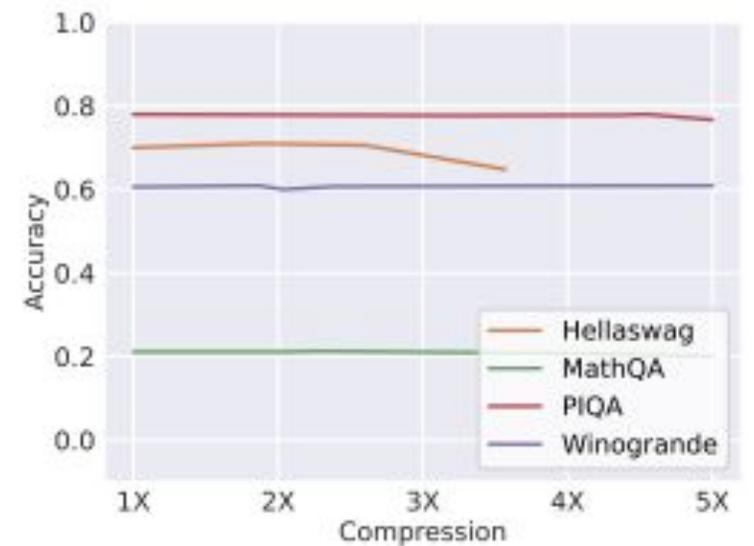


(c) Attention map at position 278

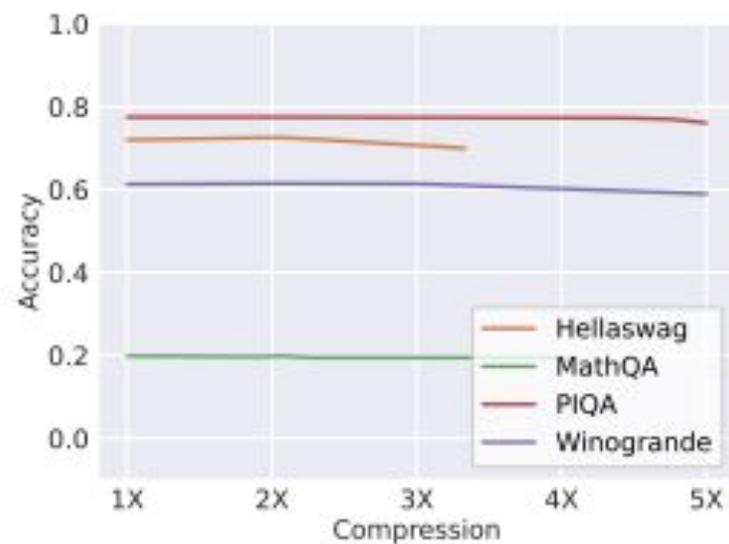
Scissorhands:
<https://arxiv.org/abs/2305.17118>



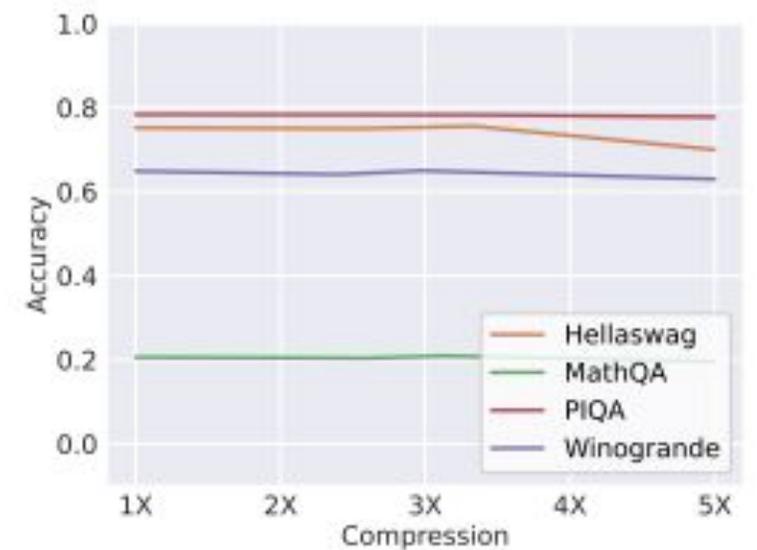
(a) Language Modeling



(b) OPT-6B Five shot

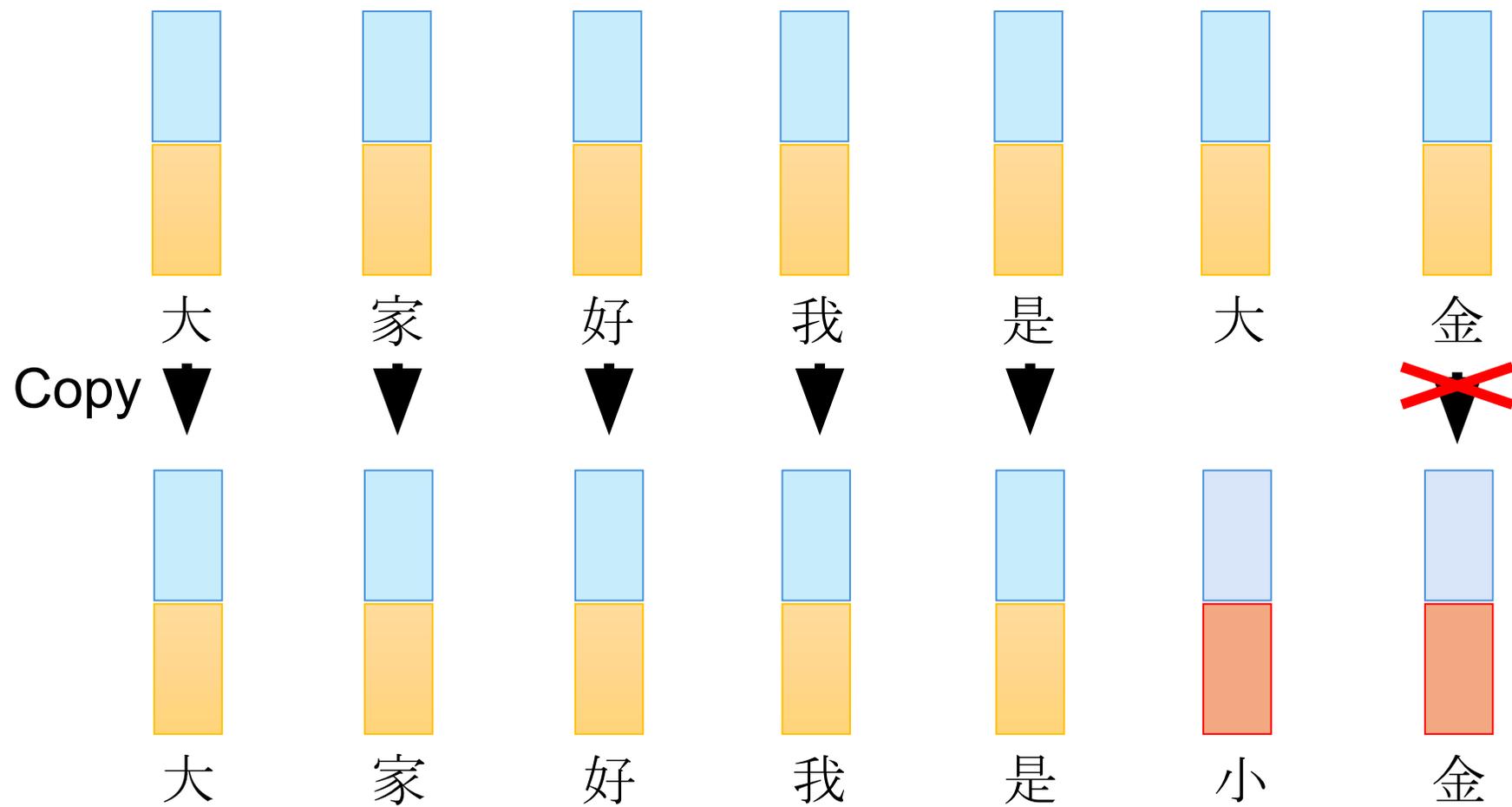


(c) OPT-13B Five shot



(d) OPT-30B Five shot

跨對話的 Cache



Text tokens

Prices per 1M tokens.

Standard

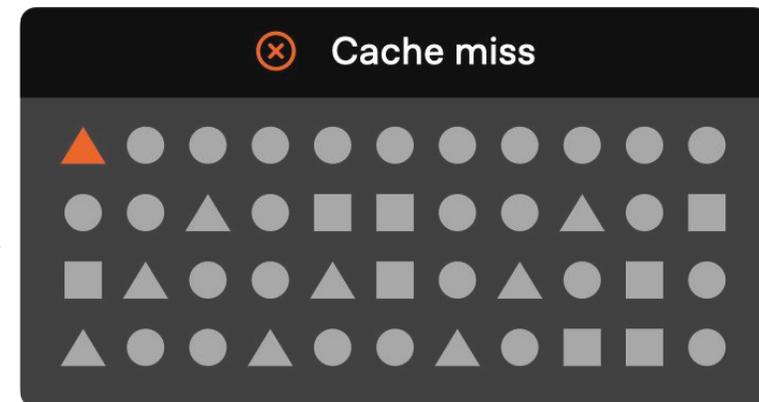
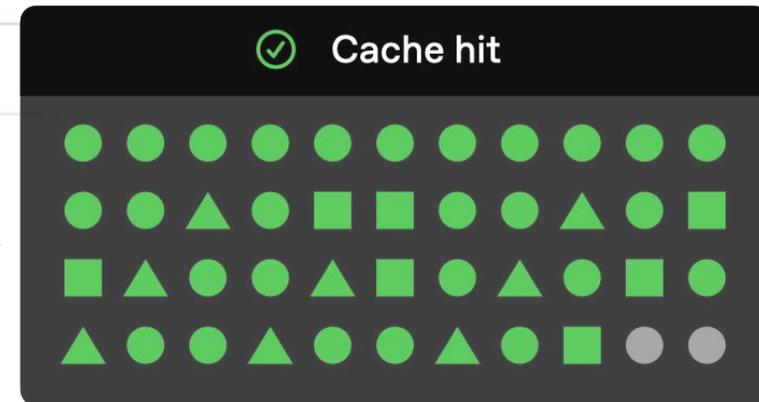
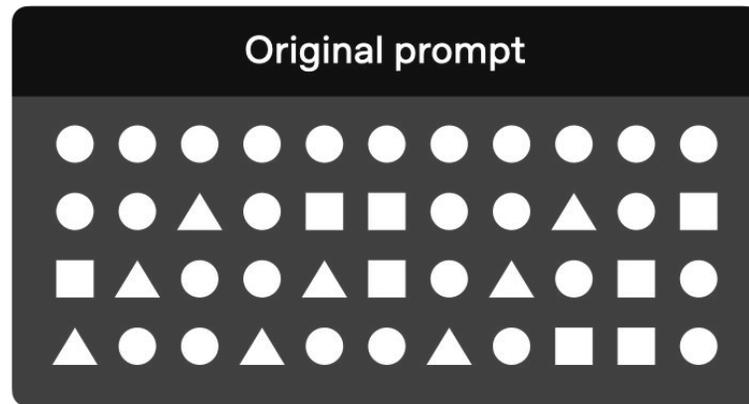
Batch

Flex

Priority

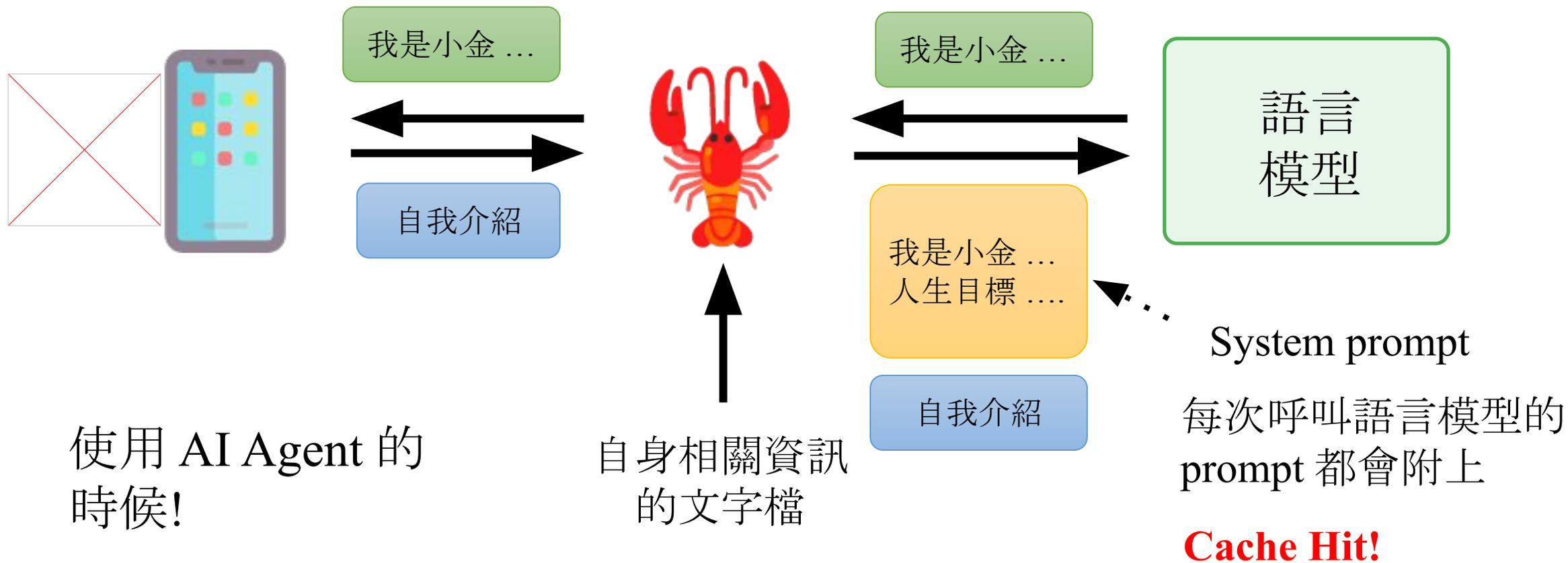
Our latest models

Model	Short context ⓘ			Long context ⓘ		
	Input	Cached input	Output	Input	Cached input	Output
gpt-5.4	\$2.50	\$0.25	\$15.00	\$5.00	\$0.50	\$22.50
gpt-5.4-pro	\$30.00	-	\$180.00	\$60.00	-	\$270.00

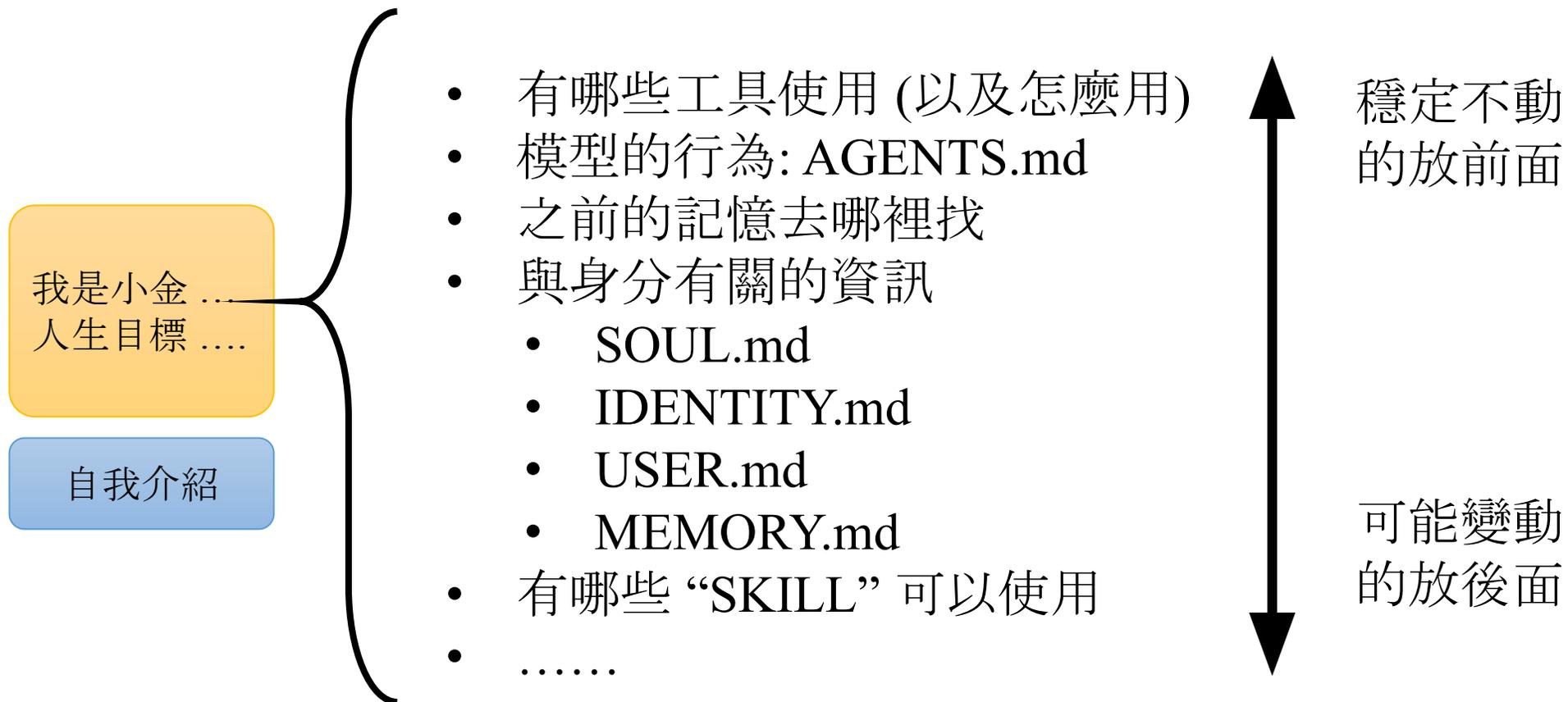


<https://developers.openai.com/api/docs/guides/prompt-caching/>
<https://developers.openai.com/api/docs/pricing>

甚麼時候跨對話的 Cache 可以發揮作用



甚麼時候跨對話的 Cache 可以發揮作用



幫我訂從台北到波士頓的班機

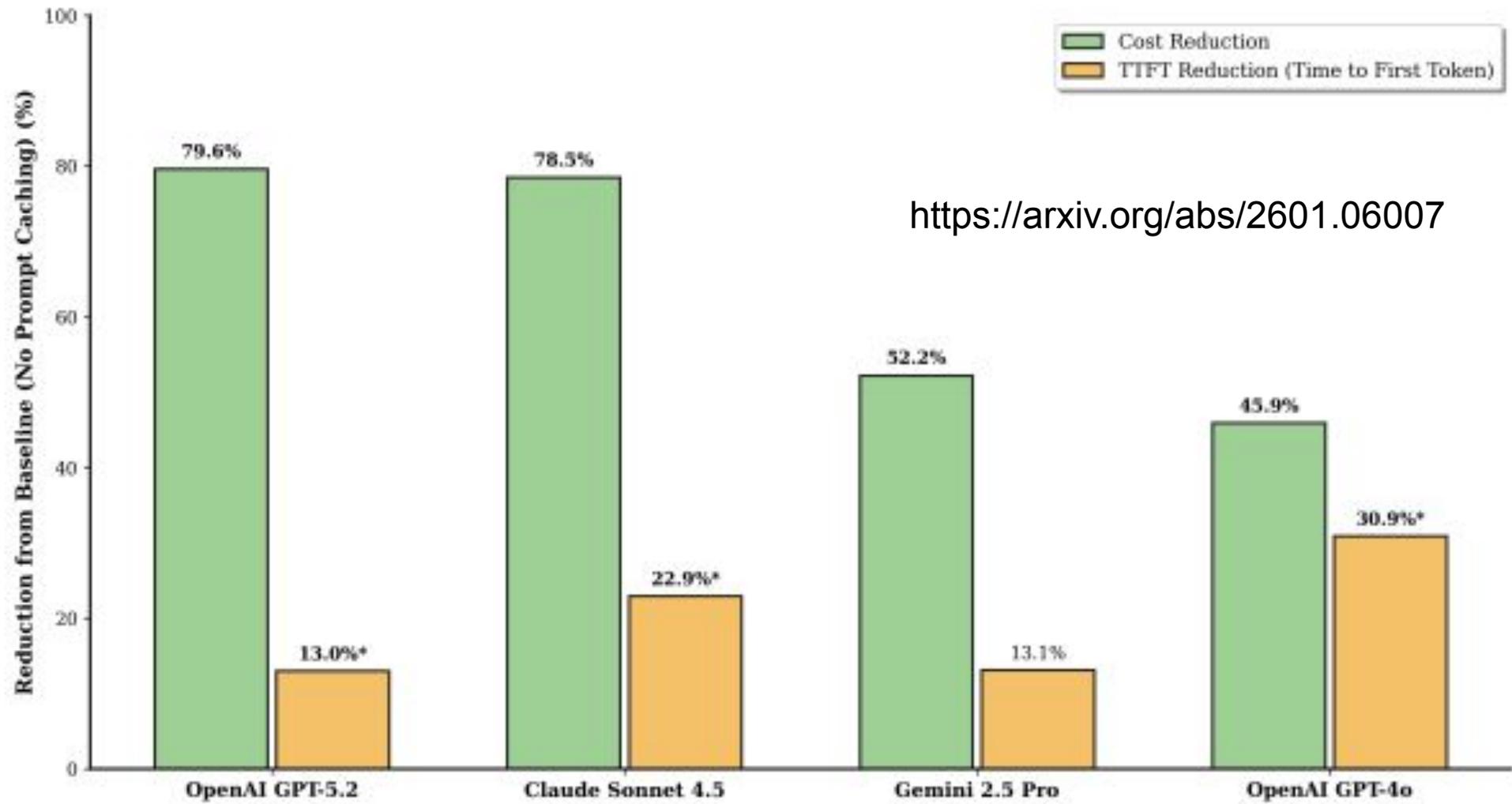
Cache hit little

幫我訂從舊金山到紐約的班機

幫我訂從 x 到 y 的班機 x =台北, y =波士頓

Cache hit more

幫我訂從 x 到 y 的班機 x =舊金山, y =紐約



<https://arxiv.org/abs/2601.06007>

結論

	方法	是否改變原有的 Attention	是否需要訓練模型	其他代價
Flash Attention	少搬資料	X	X	一點額外運算+一點點燒腦
KV Cache	儲存已經算出來的 key 和 value	X	X	占用記憶體
Multi-query attention	多個 query 共享 key 和 value	O	O	可能明顯傷害模型能力
Group-query attention		O	O	
Multi-head Latent Attention	壓縮 key 和 value	O	O	
Sliding Window Attention	改變 Attention 範圍	O	?	
Streaming LLM		O	?	
Pruning KV Cache	丟棄 key 和 value	O	X	可能明顯傷害模型能力
Speculative Decoding	用小模型來預言生成結果	X (理論上)	X	小模型還是需要耗費額外算力