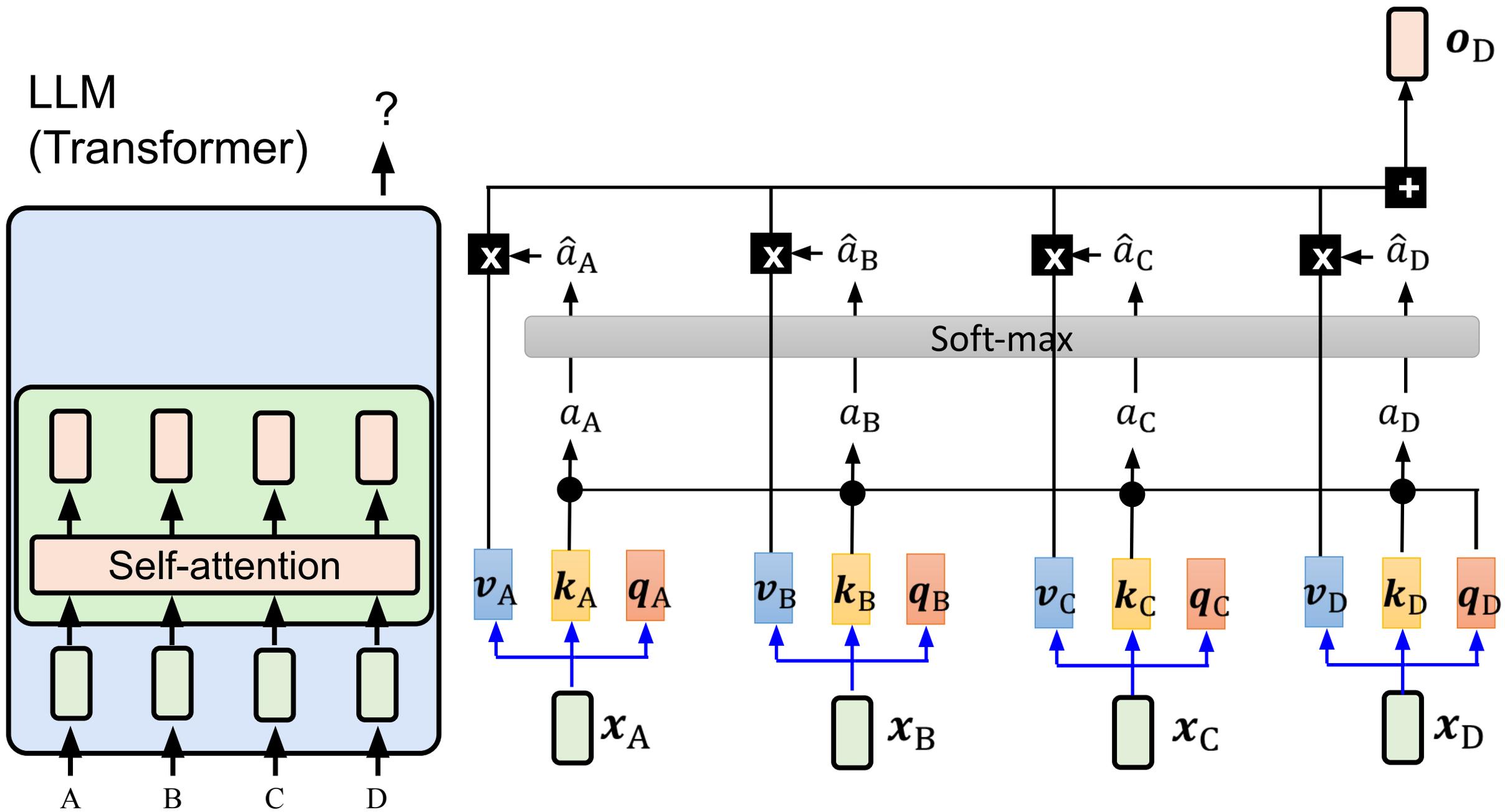


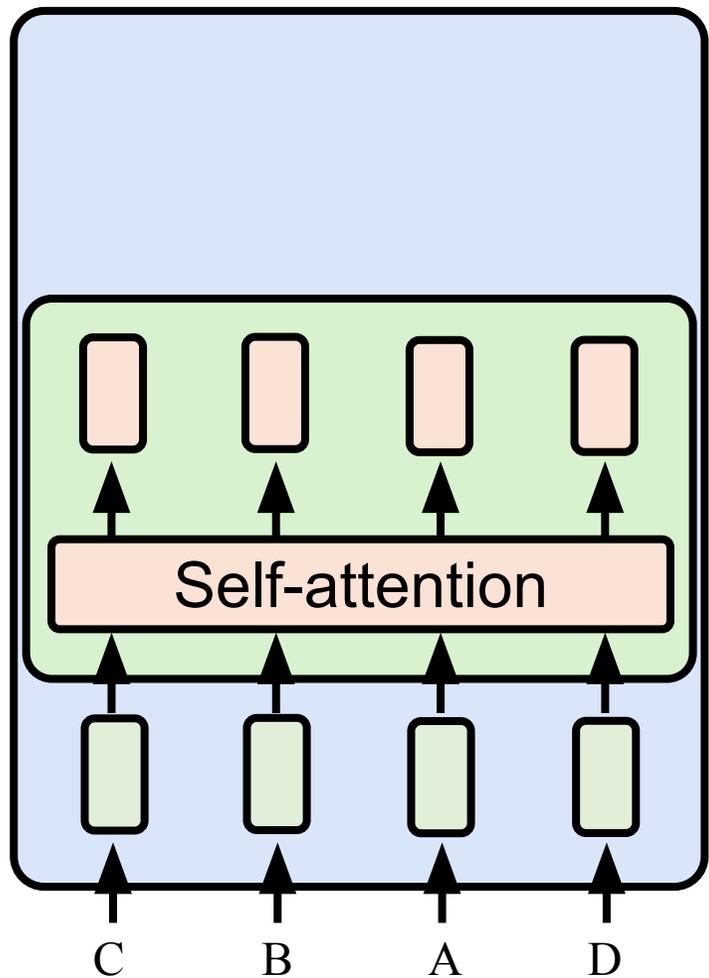
Transformer 怎麼知道 Token 的順序？

Positional Embedding

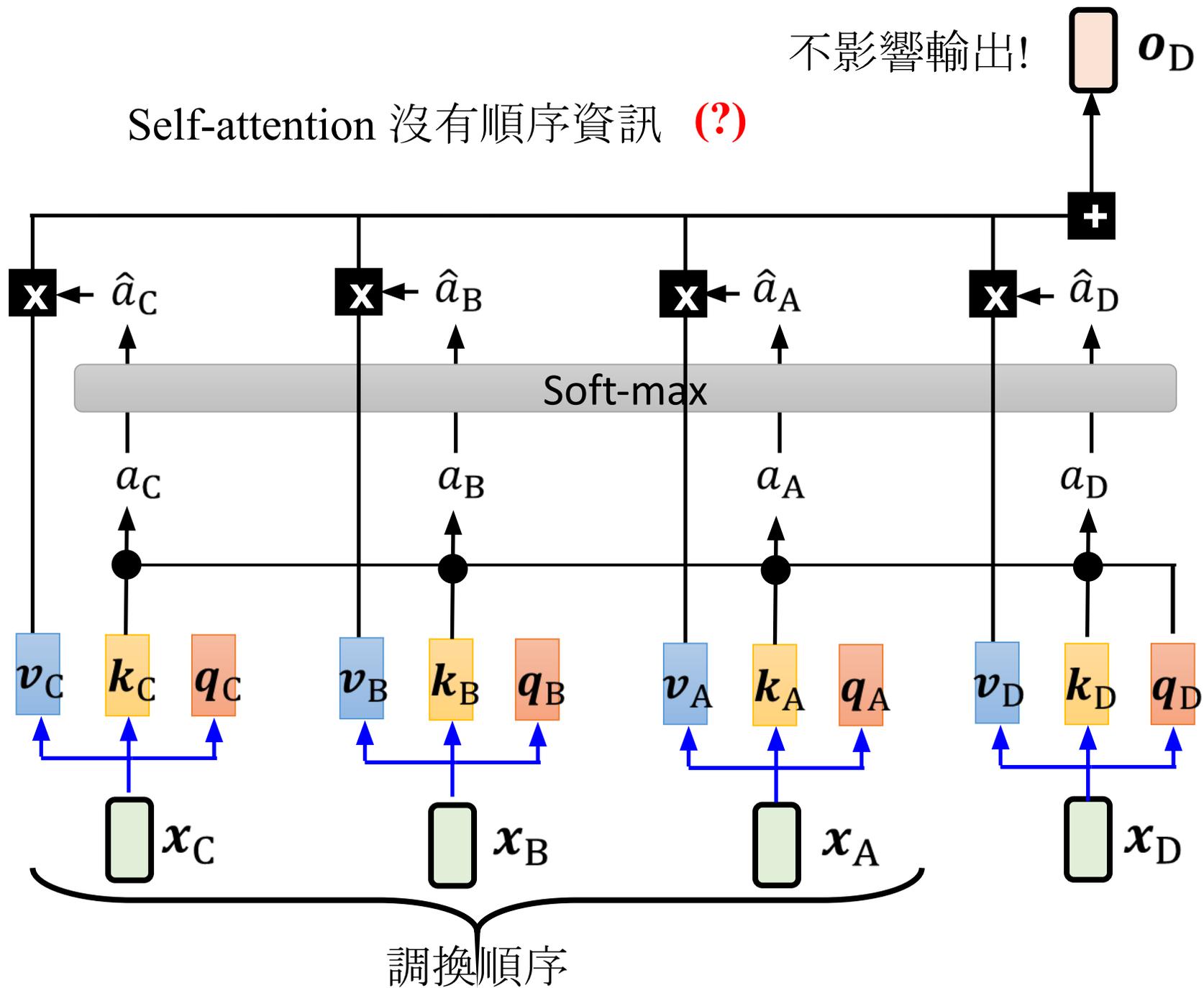


LLM
(Transformer)

?

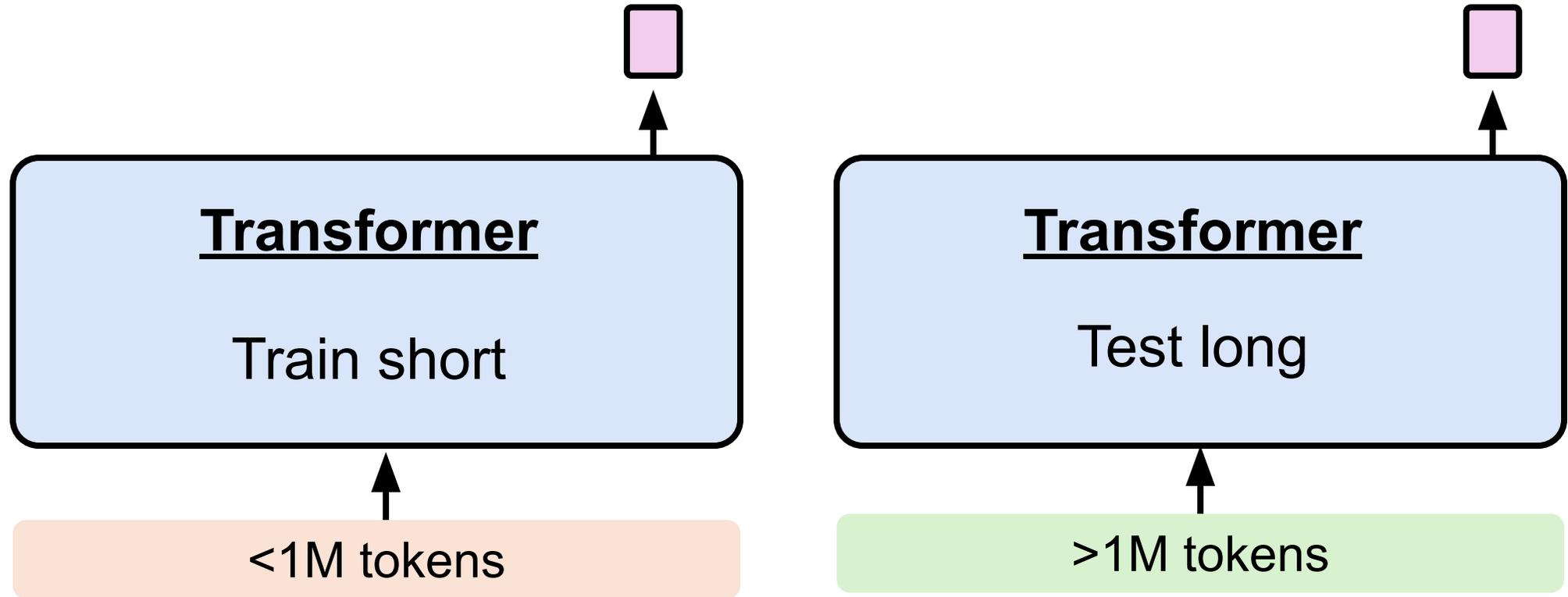


Self-attention 沒有順序資訊 (?)



Zero-shot Long Context

- Train short, test long



Outline

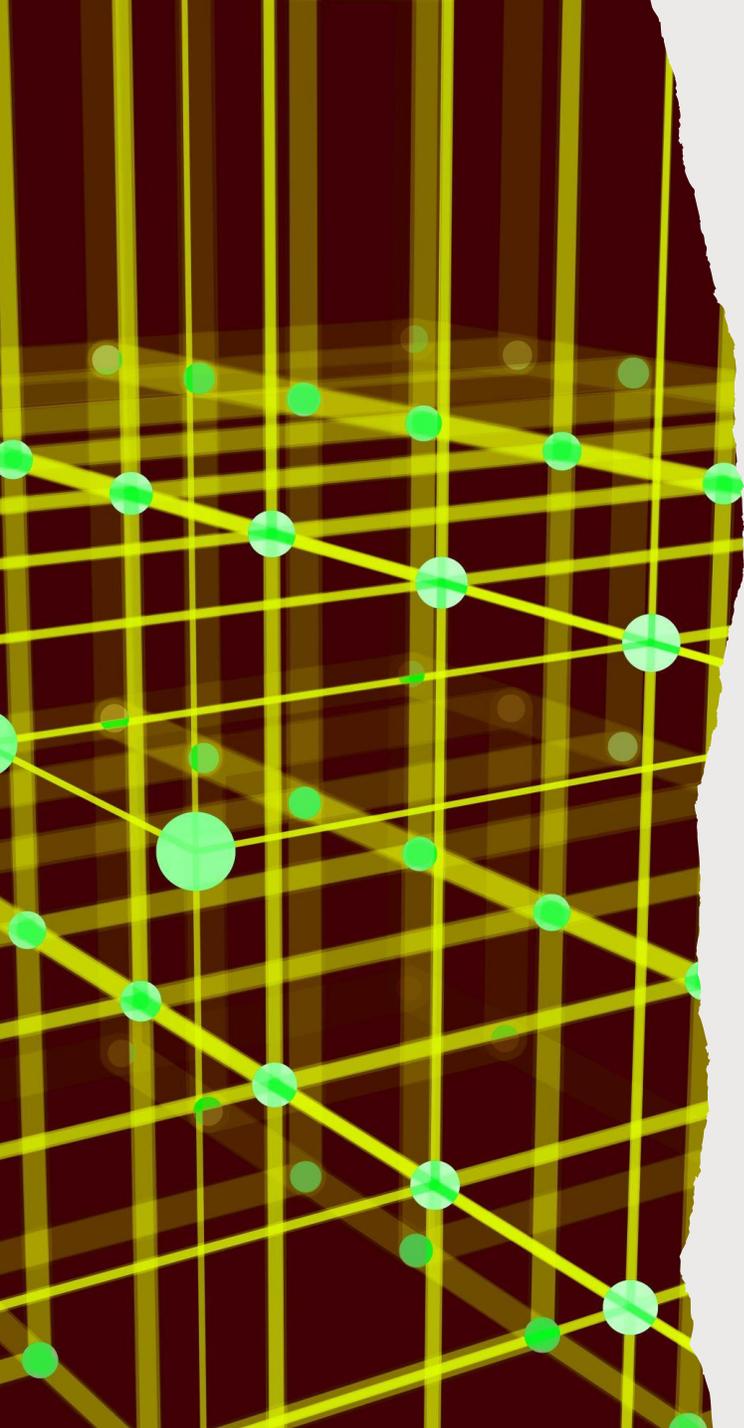
Absolute Positional Embedding

Relative Positional Embedding

RoPE

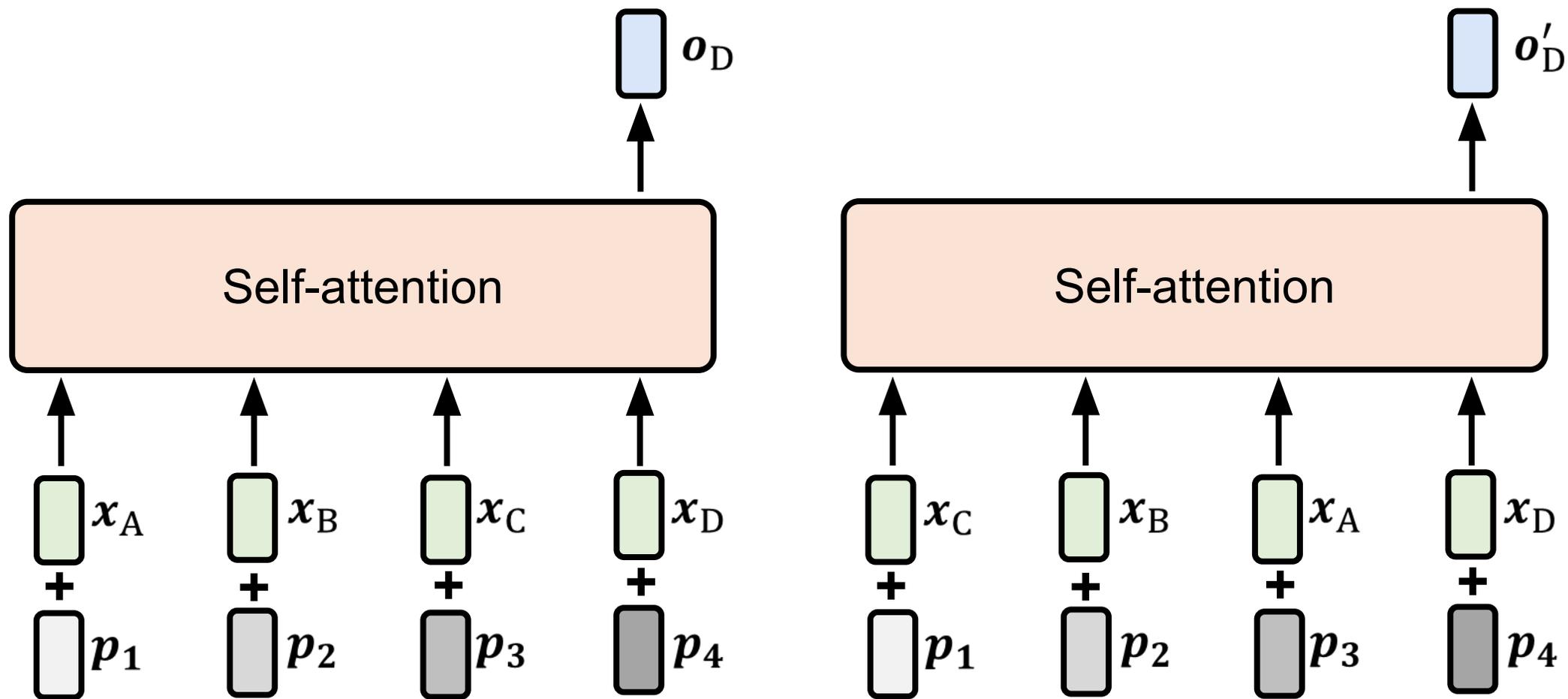
Train short, test long

No Positional Embedding?!

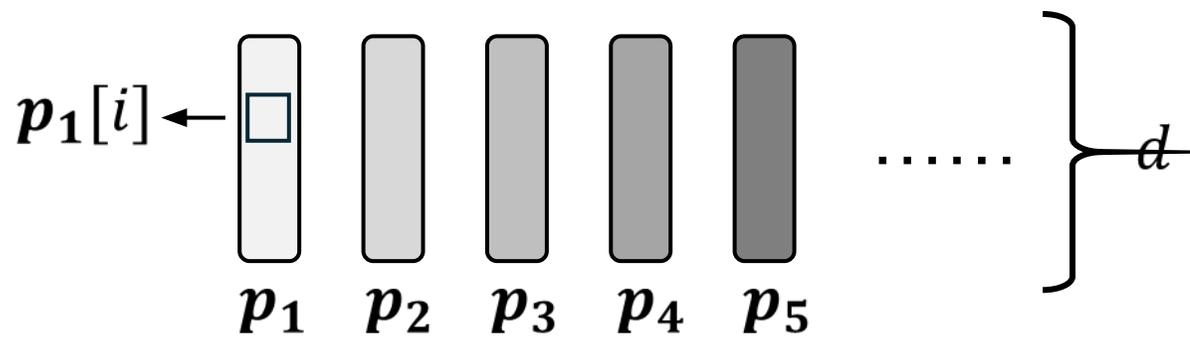


Absolute positional embedding

Absolute Positional Embedding

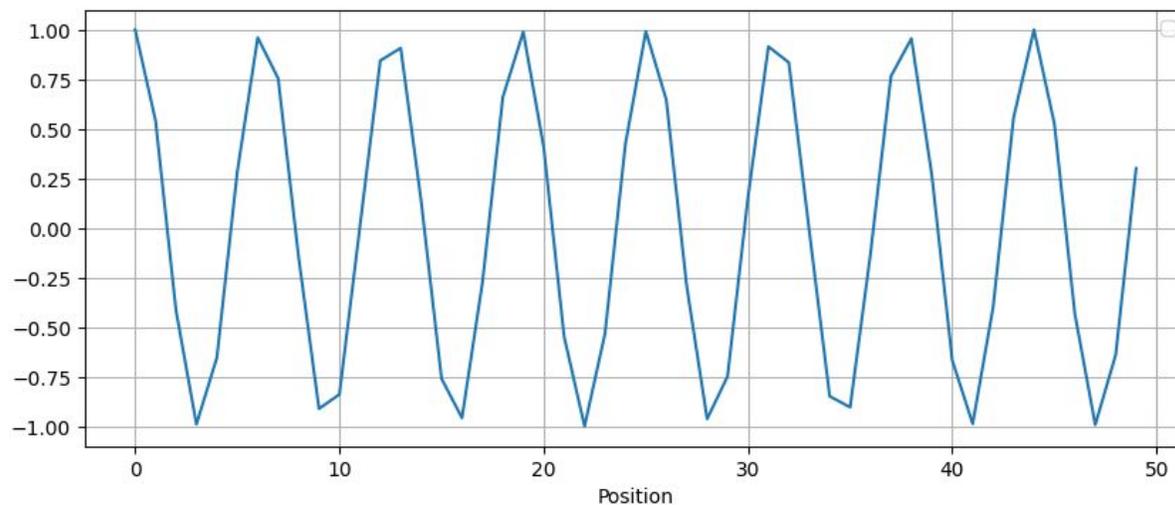
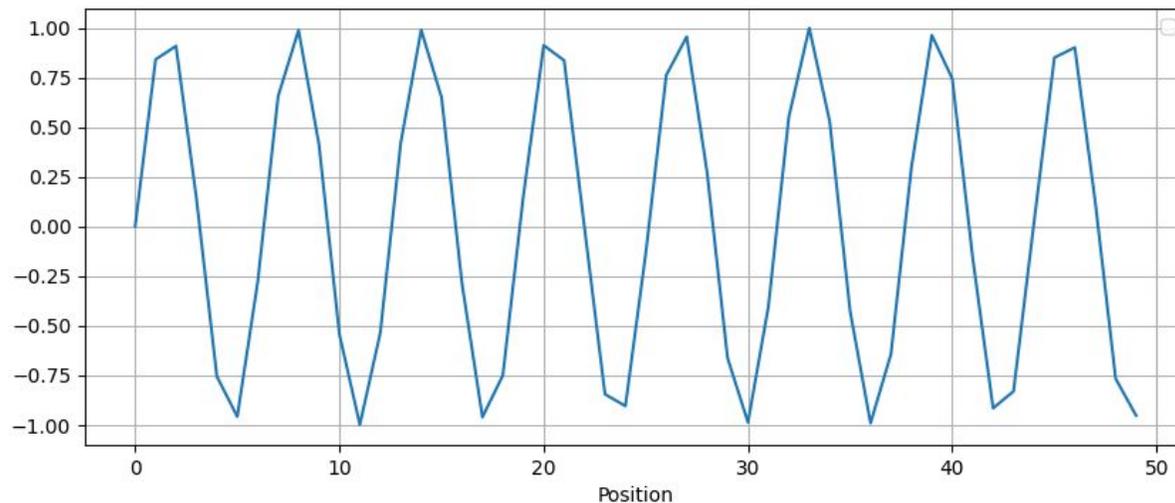


Sinusoidal Positional Embedding

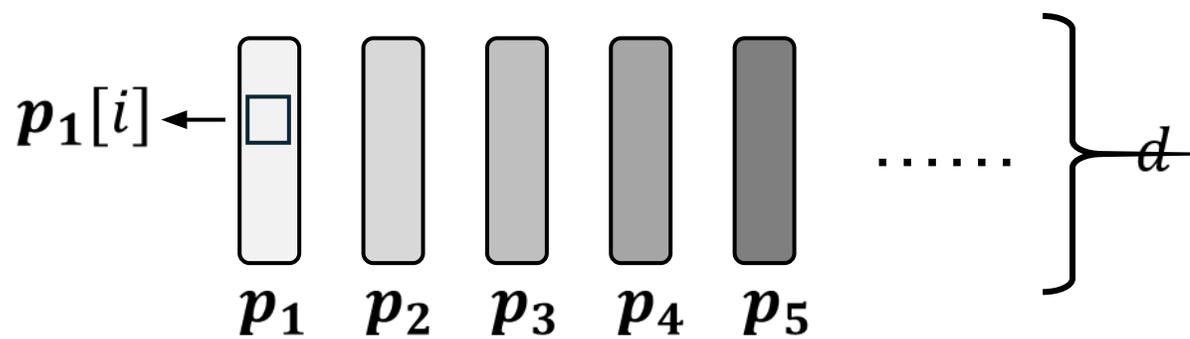


$$p_k[2i] = \sin\left(\frac{k}{10000^{2i/d}}\right)$$

$$p_k[2i + 1] = \cos\left(\frac{k}{10000^{2i/d}}\right)$$

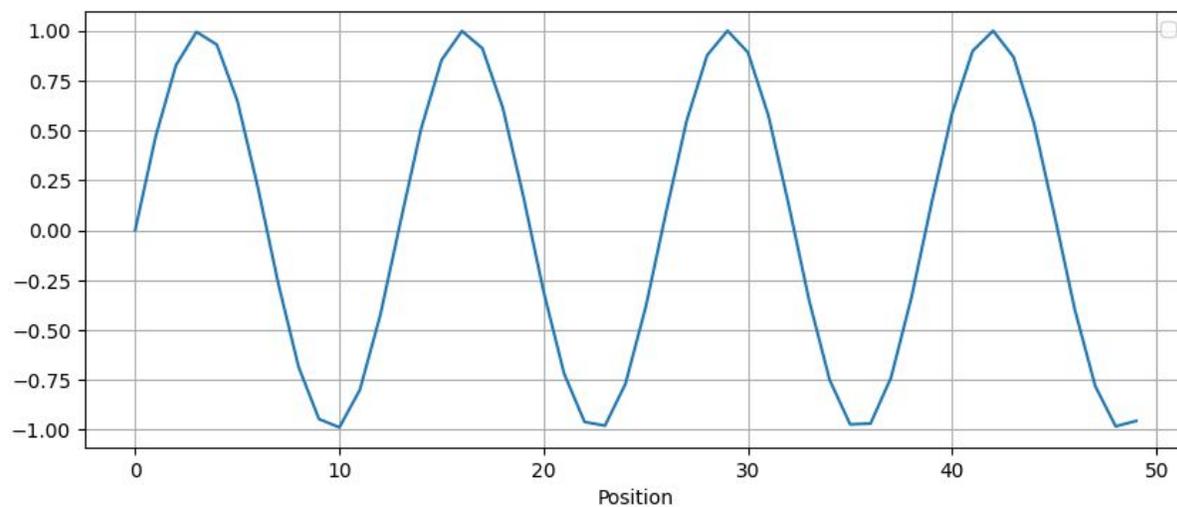
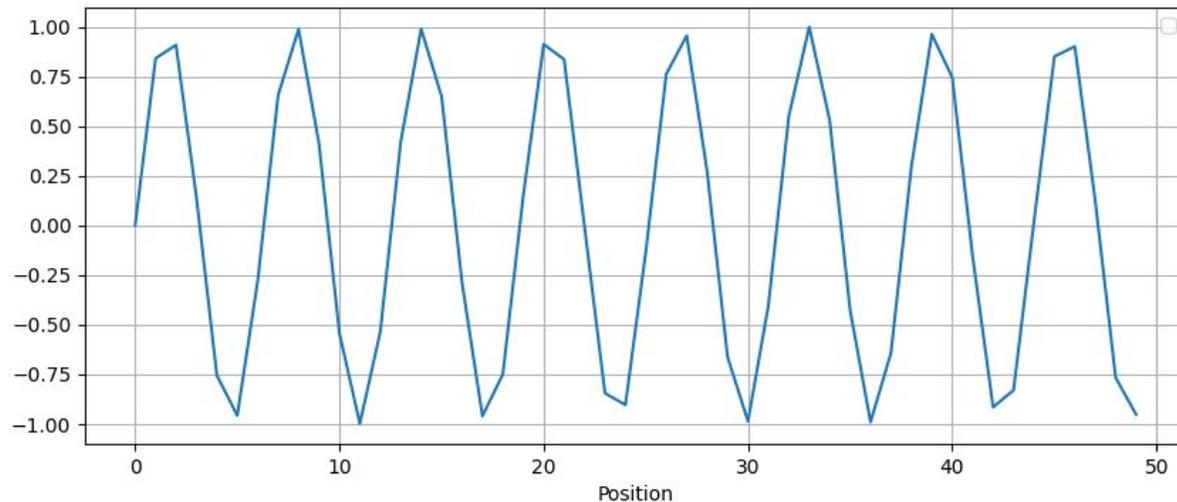


Sinusoidal Positional Embedding

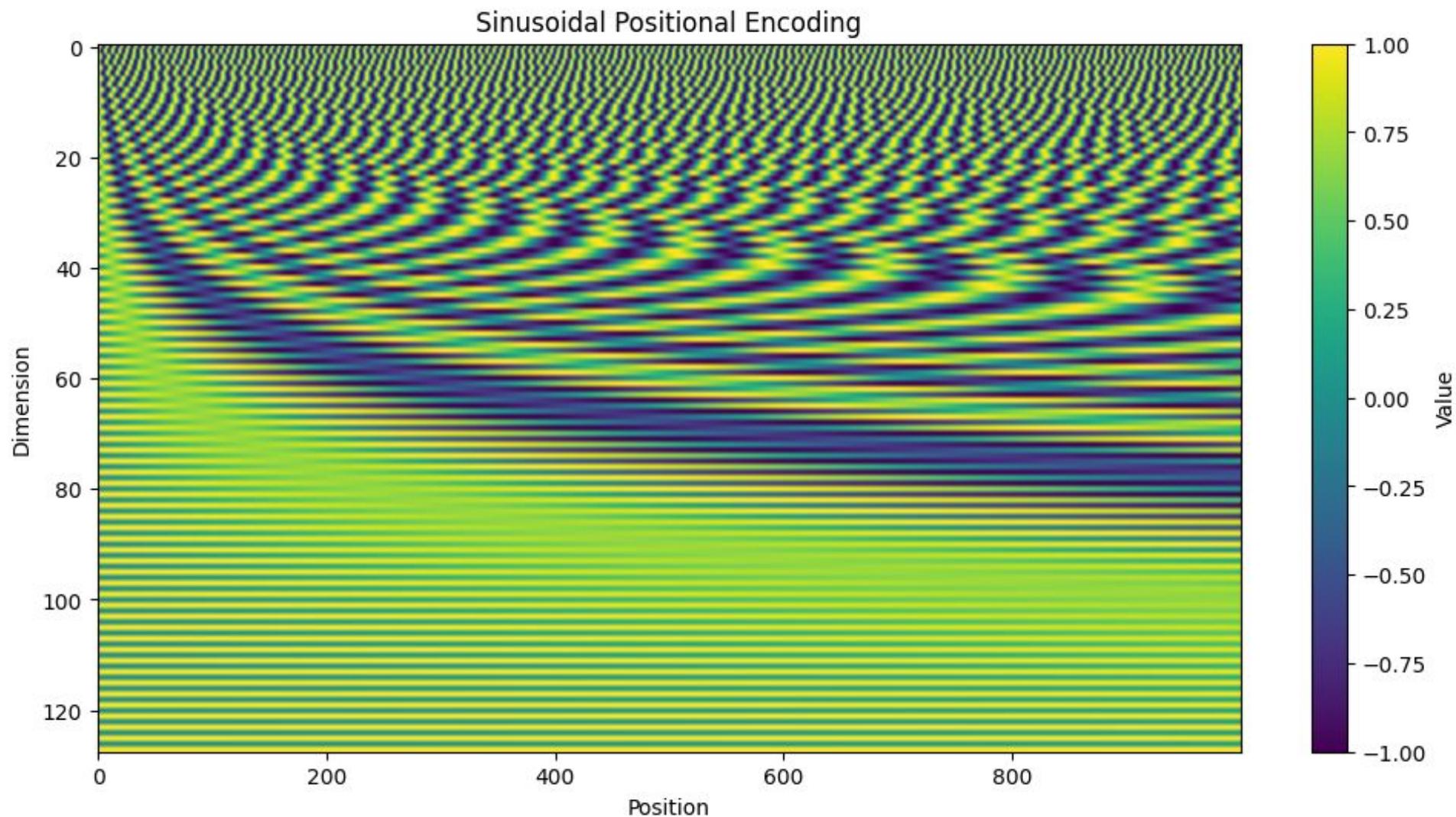


$$p_k[2i] = \sin\left(\frac{k}{10000^{2i/d}}\right)$$

$$p_k[2i + 1] = \cos\left(\frac{k}{10000^{2i/d}}\right)$$



Sinusoidal Positional Embedding



Sinusoidal Positional Embedding

$$\mathbf{p}_k[2i] = \sin\left(\frac{k}{10000^{2i/d}}\right) \quad \mathbf{p}_k[2i + 1] = \cos\left(\frac{k}{10000^{2i/d}}\right)$$

$$\frac{k}{10000^{2i/d}} = 2\pi$$

$$i = 0 \text{ (0,1 dim)}$$

6.3

$$k = 2\pi \cdot 10000^{2i/d}$$

$$i = 32 \text{ (64,65 dim)}$$

628.3

$$d = 128$$

$$i = 63 \text{ (126,127 dim)}$$

54410.1

3.5 Positional Encoding

Since our model contains no recurrence and no convolution, in order for the model to make use of the order of the sequence, we must inject some information about the relative or absolute position of the tokens in the sequence. To this end, we add "positional encodings" to the input embeddings at the bottoms of the encoder and decoder stacks. The positional encodings have the same dimension d_{model} as the embeddings, so that the two can be summed. There are many choices of positional encodings, learned and fixed [9].

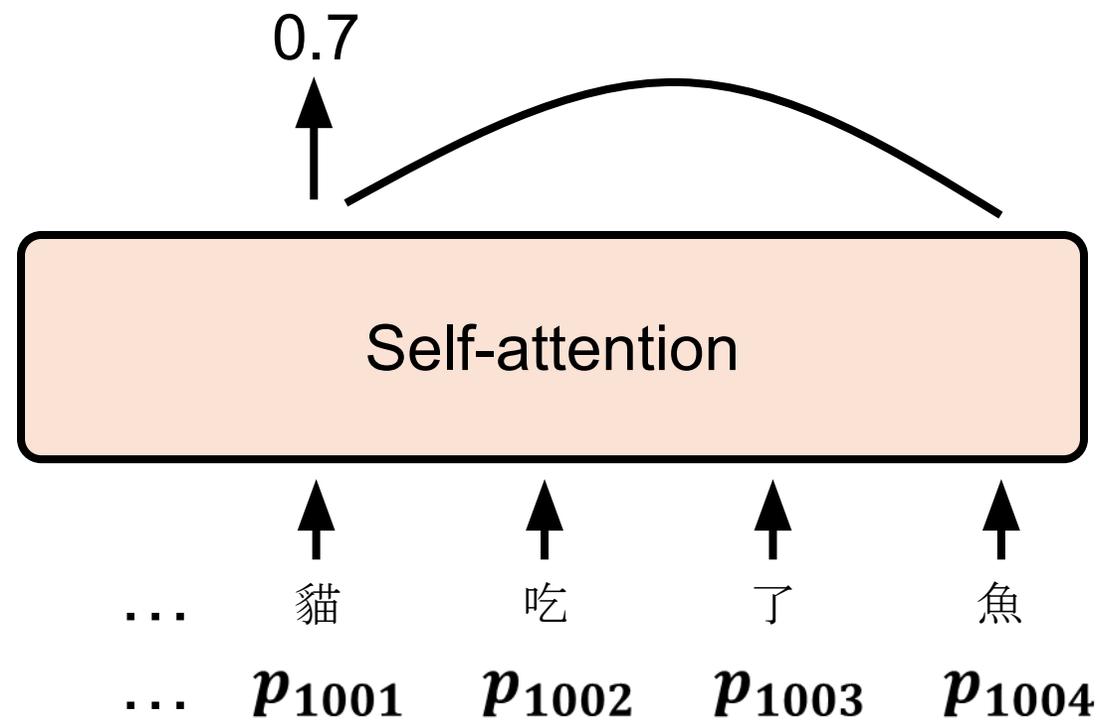
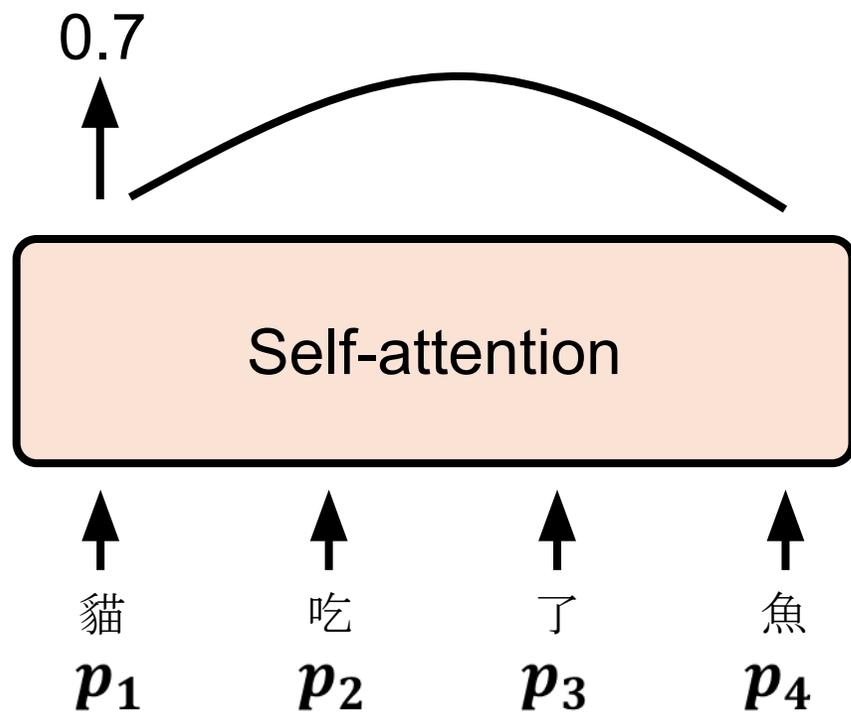
In this work, we use sine and cosine functions of different frequencies:

$$PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$
$$PE_{(pos,2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$

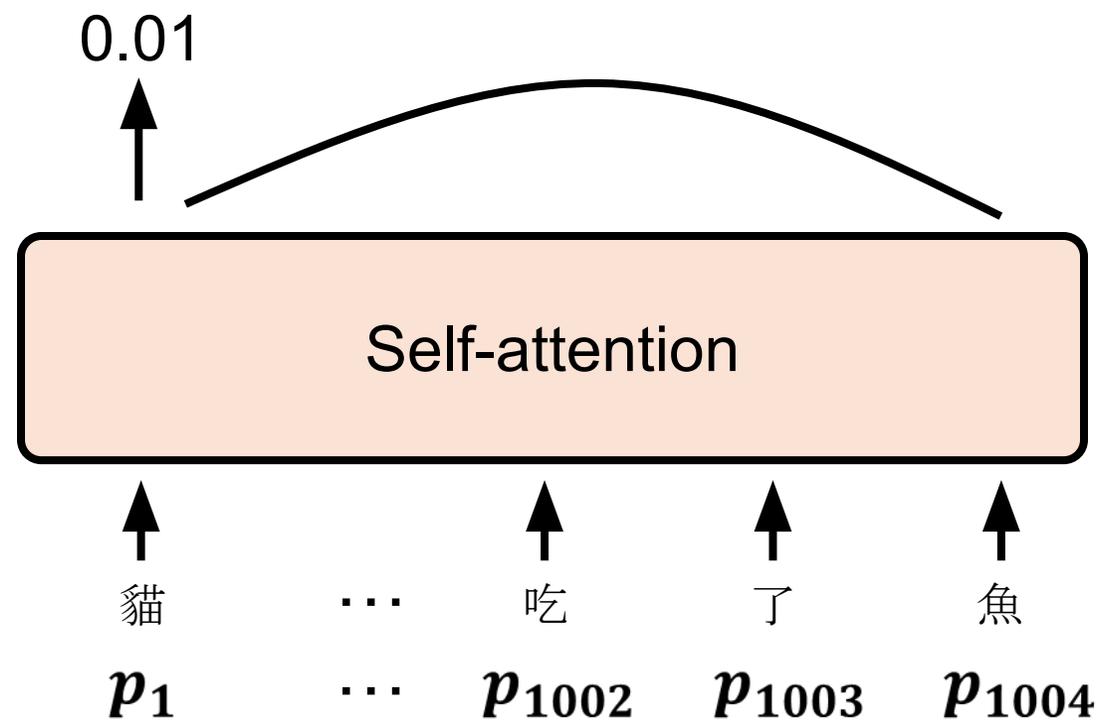
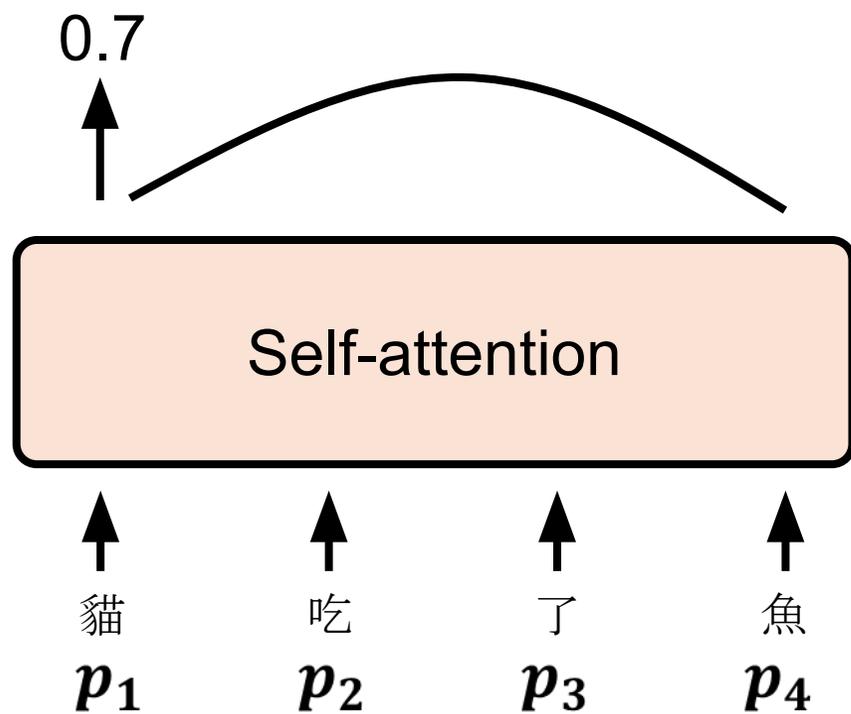
where pos is the position and i is the dimension. That is, each dimension of the positional encoding corresponds to a sinusoid. The wavelengths form a geometric progression from 2π to $10000 \cdot 2\pi$. We chose this function because we hypothesized it would allow the model to easily learn to attend by relative positions, since for any fixed offset k , PE_{pos+k} can be represented as a linear function of PE_{pos} .

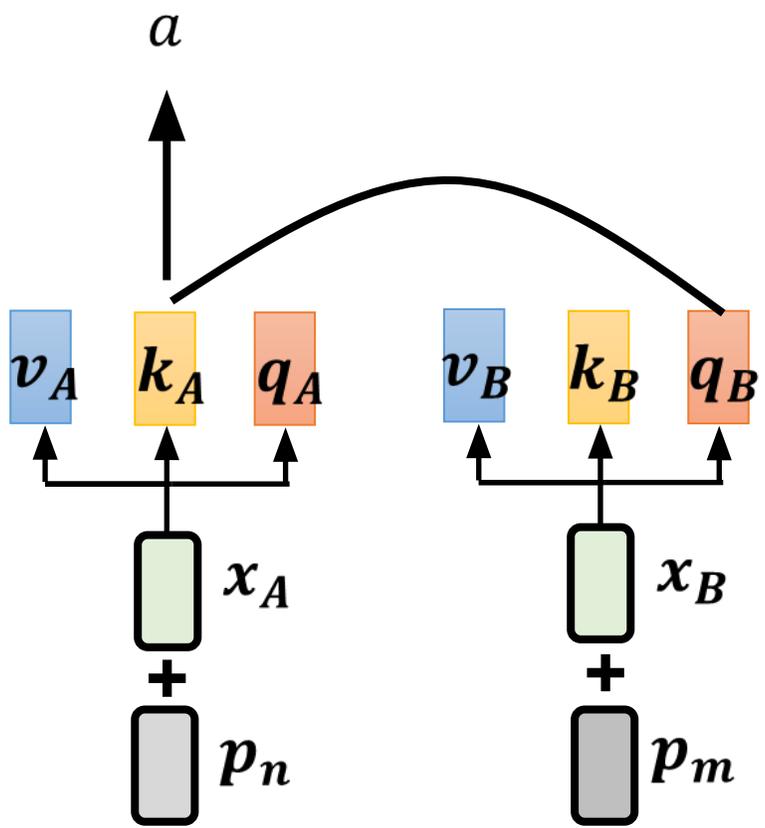
We also experimented with using learned positional embeddings [9] instead, and found that the two versions produced nearly identical results (see Table 3 row (E)). We chose the sinusoidal version because it may allow the model to extrapolate to sequence lengths longer than the ones encountered during training.

Relative positioning is crucial!



Relative positioning is crucial!





$$a = \mathbf{q}_B \cdot \mathbf{k}_A = (\mathbf{q}_B)^T \mathbf{k}_A$$

$$= \left(W_q (\mathbf{x}_B + \mathbf{p}_m) \right)^T W_k (\mathbf{x}_A + \mathbf{p}_n)$$

$$= (\mathbf{x}_B + \mathbf{p}_m)^T W_q^T W_k (\mathbf{x}_A + \mathbf{p}_n)$$

$$= \mathbf{x}_B^T W_q^T W_k \mathbf{x}_A \quad \text{只跟內容有關}$$

$$+ \mathbf{x}_B^T W_q^T W_k \mathbf{p}_n + \mathbf{p}_m^T W_q^T W_k \mathbf{x}_A$$

內容、位置
交互影響

$$+ \mathbf{p}_m^T W_q^T W_k \mathbf{p}_n \quad \text{只跟位置有關}$$

$m = 3, n = 1$ 和 $m = 1003, n = 1001$
沒有甚麼關聯

Back to Sinusoidal Positional Embedding

$$\mathbf{p}_{k+r} = M_r \mathbf{p}_k \quad \mathbf{p}_k[2i] = \sin\left(\frac{k}{10000^{2i/d}}\right) \quad \mathbf{p}_k[2i+1] = \cos\left(\frac{k}{10000^{2i/d}}\right)$$

$$\mathbf{p}_{k+r}[2i] = \sin\left(\frac{k+r}{z}\right) \quad \mathbf{p}_k[2i] = \sin\left(\frac{k}{z}\right)$$

$$\mathbf{p}_{k+r}[2i+1] = \cos\left(\frac{k+r}{z}\right) \quad \mathbf{p}_k[2i+1] = \cos\left(\frac{k}{z}\right)$$

$$\begin{aligned} \sin(a+b) &= \sin(a)\cos(b) + \cos(a)\sin(b) \\ \cos(a+b) &= \cos(a)\cos(b) - \sin(a)\sin(b) \end{aligned}$$

Back to Sinusoidal Positional Embedding

$$\mathbf{p}_{k+r} = M_r \mathbf{p}_k \quad \mathbf{p}_k[2i] = \sin\left(\frac{k}{10000^{2i/d} z}\right) \quad \mathbf{p}_k[2i+1] = \cos\left(\frac{k}{10000^{2i/d} z}\right)$$

$$\mathbf{p}_{k+r}[2i] = \sin\left(\frac{k+r}{z}\right) = \sin\left(\frac{k}{z}\right) \cos\left(\frac{r}{z}\right) + \cos\left(\frac{k}{z}\right) \sin\left(\frac{r}{z}\right)$$

$$\mathbf{p}_{k+r}[2i+1] = \cos\left(\frac{k+r}{z}\right) = \cos\left(\frac{k}{z}\right) \cos\left(\frac{r}{z}\right) - \sin\left(\frac{k}{z}\right) \sin\left(\frac{r}{z}\right)$$

$$\mathbf{p}_k[2i] = \sin\left(\frac{k}{z}\right)$$

$$\mathbf{p}_k[2i+1] = \cos\left(\frac{k}{z}\right)$$

$$\sin(a+b) = \sin(a)\cos(b) + \cos(a)\sin(b)$$

$$\cos(a+b) = \cos(a)\cos(b) - \sin(a)\sin(b)$$

$$\mathbf{p}_{k+r} = M_r \mathbf{p}_k \quad \mathbf{p}_k[2i] = \sin\left(\frac{k}{\underbrace{10000^{2i/d}}_z}\right) \quad \mathbf{p}_k[2i+1] = \cos\left(\frac{k}{\underbrace{10000^{2i/d}}_z}\right)$$

$$\begin{aligned} \mathbf{p}_{k+r}[2i] &= \sin\left(\frac{k+r}{z}\right) = \sin\left(\frac{k}{z}\right)\cos\left(\frac{r}{z}\right) + \cos\left(\frac{k}{z}\right)\sin\left(\frac{r}{z}\right) \\ &= \mathbf{p}_k[2i]\cos\left(\frac{r}{z}\right) + \mathbf{p}_k[2i+1]\sin\left(\frac{r}{z}\right) \end{aligned}$$

$$\mathbf{p}_{k+r}[2i+1] = \cos\left(\frac{k+r}{z}\right) = \cos\left(\frac{k}{z}\right)\cos\left(\frac{r}{z}\right) - \sin\left(\frac{k}{z}\right)\sin\left(\frac{r}{z}\right)$$

$$\mathbf{p}_k[2i] = \sin\left(\frac{k}{z}\right) = \mathbf{p}_k[2i+1]\cos\left(\frac{r}{z}\right) - \mathbf{p}_k[2i]\sin\left(\frac{r}{z}\right)$$

$$\mathbf{p}_k[2i+1] = \cos\left(\frac{k}{z}\right)$$

$$\mathbf{p}_{k+r} = M_r \mathbf{p}_k \quad \mathbf{p}_k[2i] = \sin\left(\frac{k}{\underbrace{10000^{2i/d}}_Z}\right) \quad \mathbf{p}_k[2i+1] = \cos\left(\frac{k}{\underbrace{10000^{2i/d}}_Z}\right)$$

$$\mathbf{p}_{k+r}[2i] = \mathbf{p}_k[2i] \cos\left(\frac{r}{Z}\right) + \mathbf{p}_k[2i+1] \sin\left(\frac{r}{Z}\right)$$

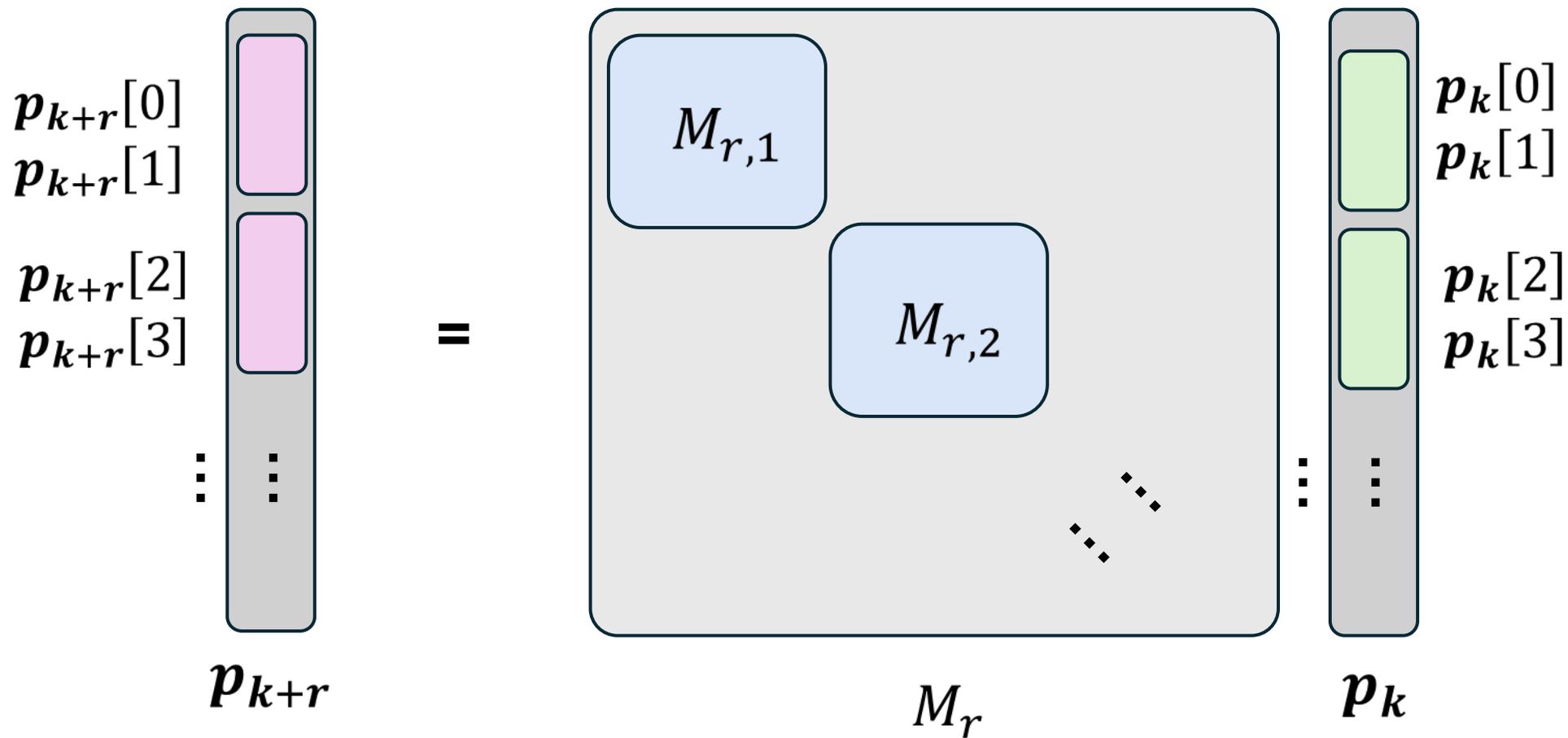
$$\mathbf{p}_{k+r}[2i+1] = \mathbf{p}_k[2i+1] \cos\left(\frac{r}{Z}\right) - \mathbf{p}_k[2i] \sin\left(\frac{r}{Z}\right)$$

$$\begin{bmatrix} \mathbf{p}_{k+r}[2i] \\ \mathbf{p}_{k+r}[2i+1] \end{bmatrix} = \begin{bmatrix} \cos\left(\frac{r}{Z}\right) & \sin\left(\frac{r}{Z}\right) \\ -\sin\left(\frac{r}{Z}\right) & \cos\left(\frac{r}{Z}\right) \end{bmatrix} \begin{bmatrix} \mathbf{p}_k[2i] \\ \mathbf{p}_k[2i+1] \end{bmatrix}$$

$$M_{r,i}$$

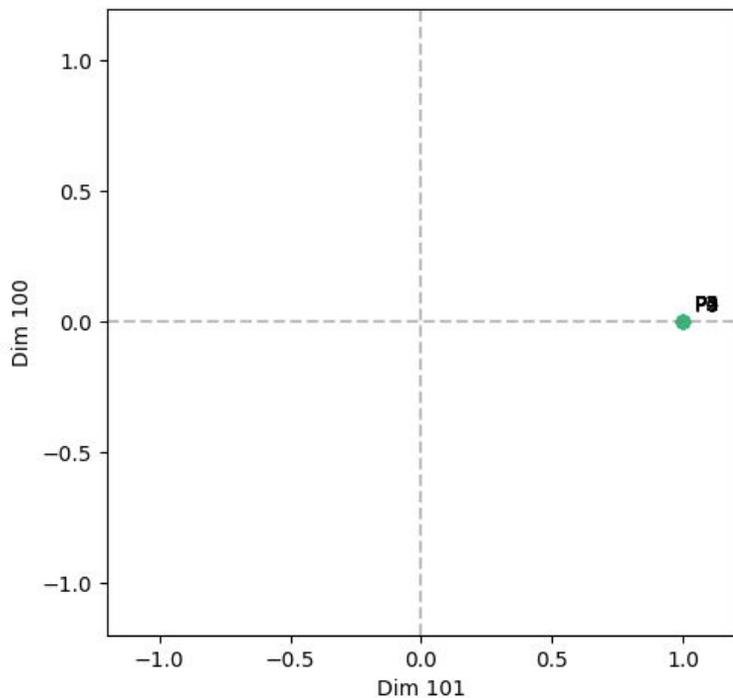
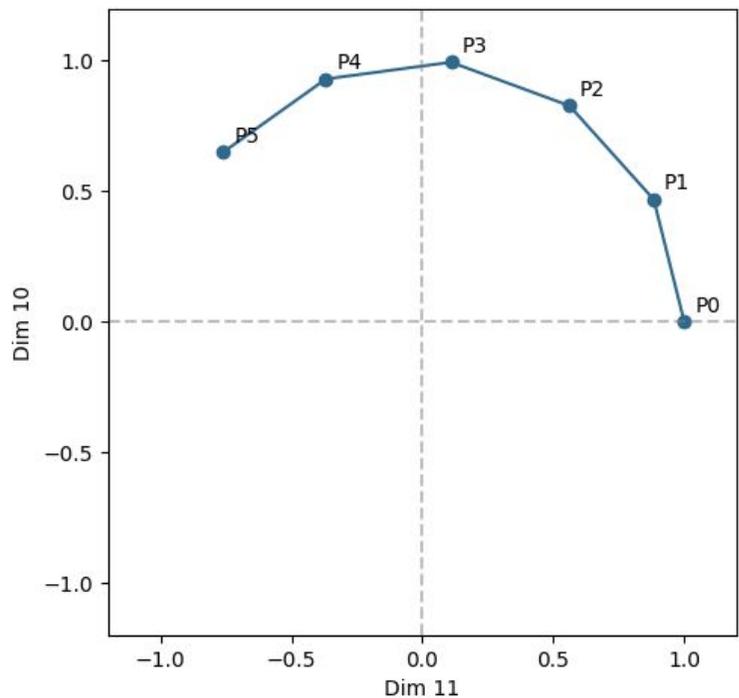
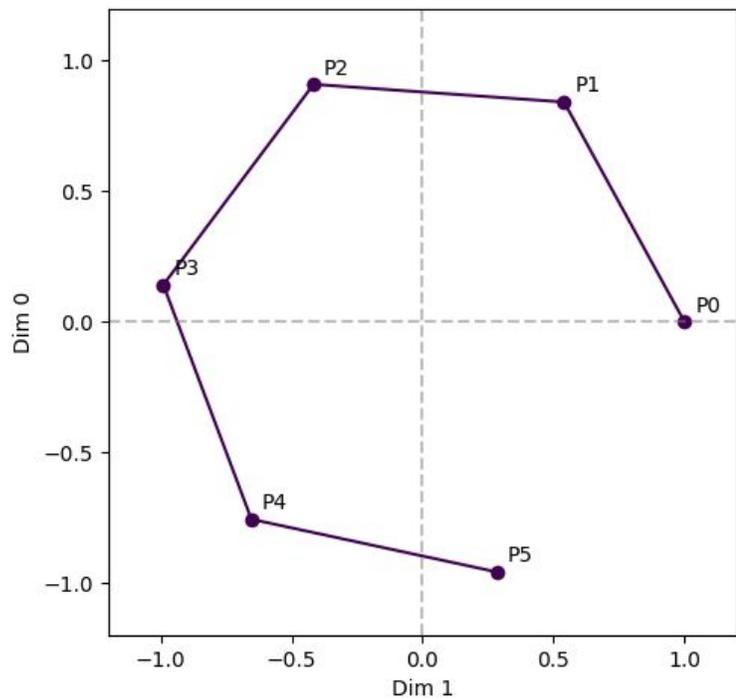
$$\mathbf{p}_{k+r} = M_r \mathbf{p}_k \quad \mathbf{p}_k[2i] = \sin\left(\frac{k}{10000^{2i/d}}\right) \quad \mathbf{p}_k[2i+1] = \cos\left(\frac{k}{10000^{2i/d}}\right)$$

$$\begin{bmatrix} \mathbf{p}_{k+r}[2i] \\ \mathbf{p}_{k+r}[2i+1] \end{bmatrix} = M_{r,i} \begin{bmatrix} \mathbf{p}_k[2i] \\ \mathbf{p}_k[2i+1] \end{bmatrix}$$

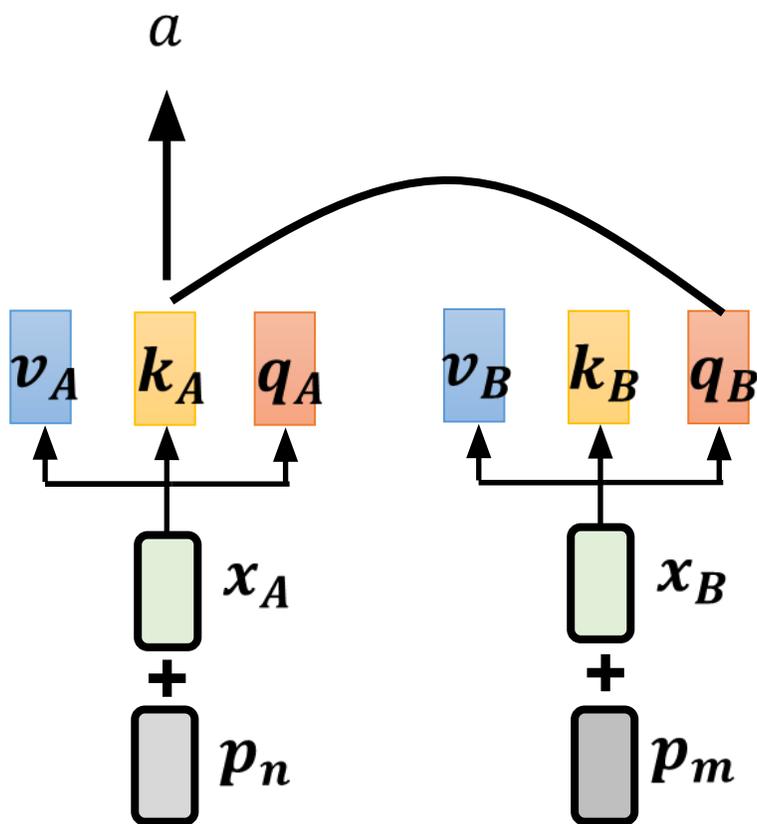


$$\begin{bmatrix} \mathbf{p}_{k+r}[2i] \\ \mathbf{p}_{k+r}[2i+1] \end{bmatrix} = \begin{bmatrix} \cos\left(\frac{r}{Z}\right) & \sin\left(\frac{r}{Z}\right) \\ -\sin\left(\frac{r}{Z}\right) & \cos\left(\frac{r}{Z}\right) \end{bmatrix} \begin{bmatrix} \mathbf{p}_k[2i] \\ \mathbf{p}_k[2i+1] \end{bmatrix}$$

$M_{r,i}$



$$\mathbf{p}_{k+r} = M_r \mathbf{p}_k$$



$$a = \mathbf{q}_B \cdot \mathbf{k}_A = (\mathbf{q}_B)^T \mathbf{k}_A$$

$$= \left(W_q (\mathbf{x}_B + \mathbf{p}_m) \right)^T W_k (\mathbf{x}_A + \mathbf{p}_n)$$

$$= (\mathbf{x}_B + \mathbf{p}_m)^T W_q^T W_k (\mathbf{x}_A + \mathbf{p}_n)$$

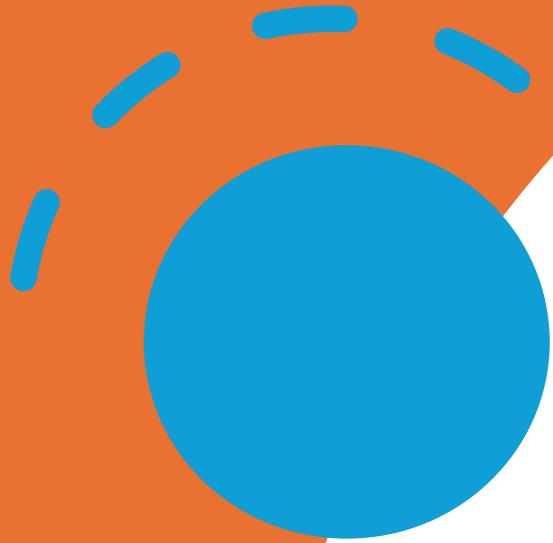
$$= \mathbf{x}_B^T W_q^T W_k \mathbf{x}_A \quad \text{只跟內容有關}$$

$$+ \mathbf{x}_B^T W_q^T W_k \mathbf{p}_n + \mathbf{p}_m^T W_q^T W_k \mathbf{x}_A$$

內容、位置
交互影響

$$+ \mathbf{p}_m^T W_q^T W_k \mathbf{p}_n \quad \text{只跟位置有關}$$

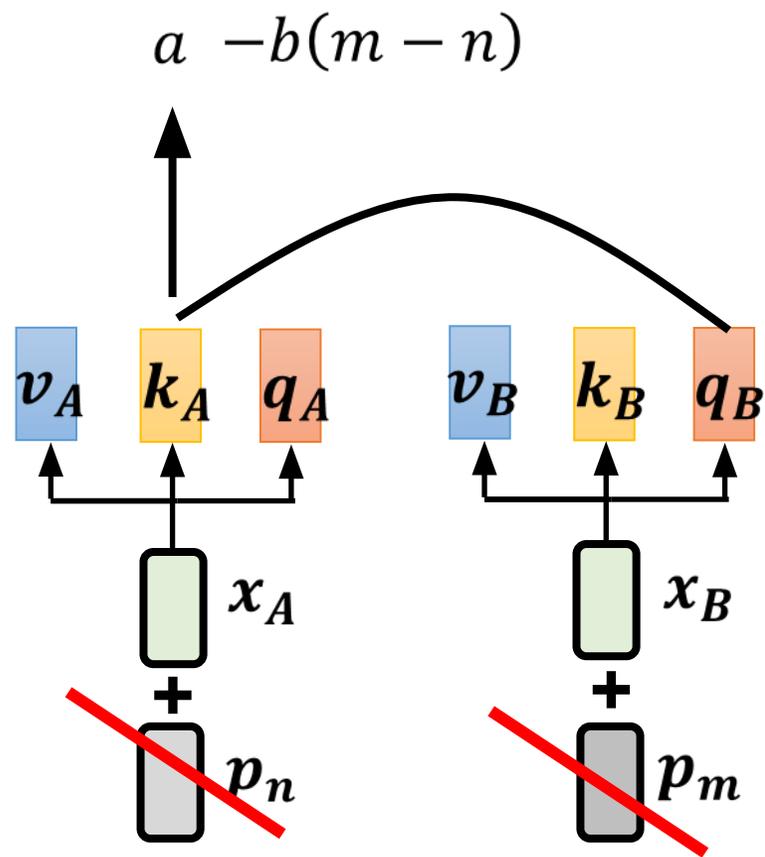
$$= (M_{m-n} \mathbf{p}_n)^T W_q^T W_k \mathbf{p}_n = (\mathbf{p}_n)^T M_{m-n} W_q^T W_k \mathbf{p}_n$$



Relative Positional Embedding

Attention with Linear Biases (ALiBi)

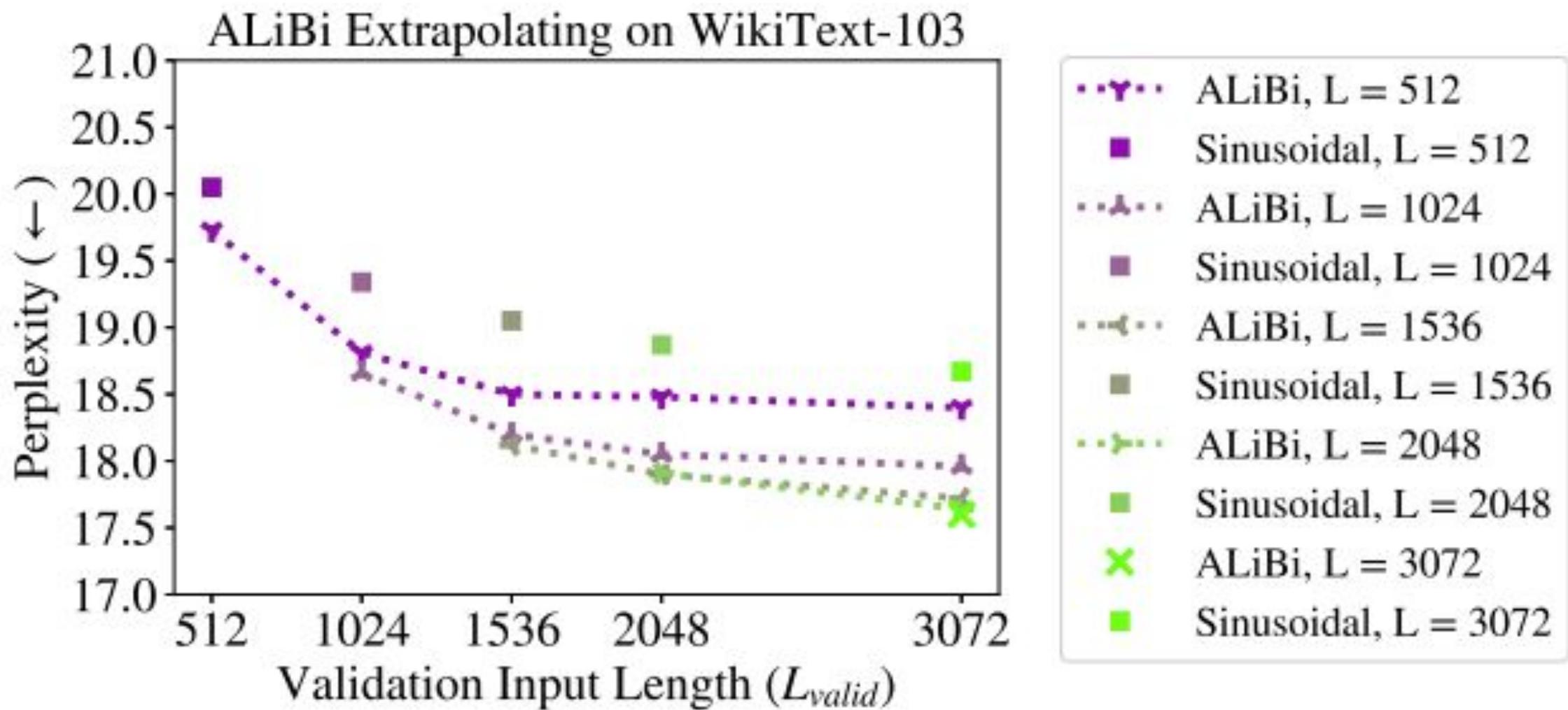
<https://arxiv.org/abs/2108.12409>



距離越遠, attention 越小,
就這樣!

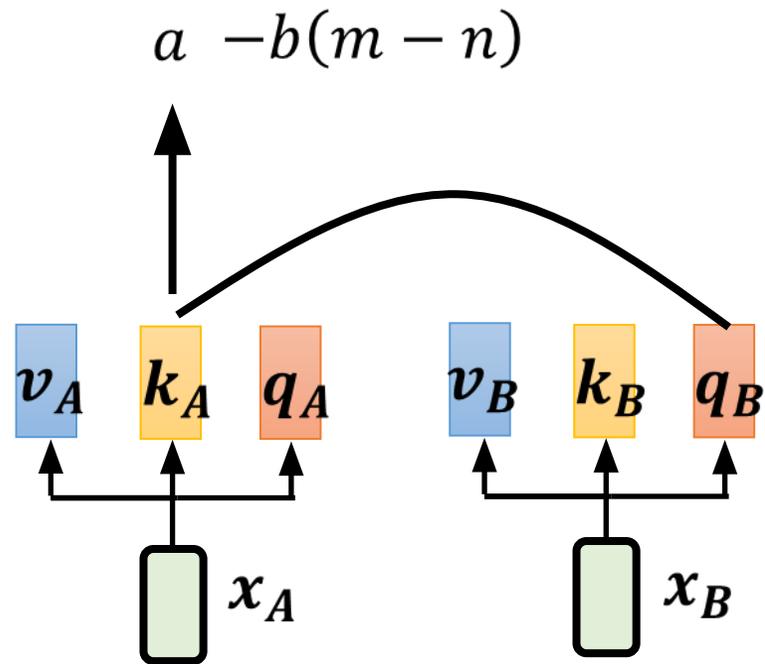
b 是手動設置

不同 attention head
設置不同值

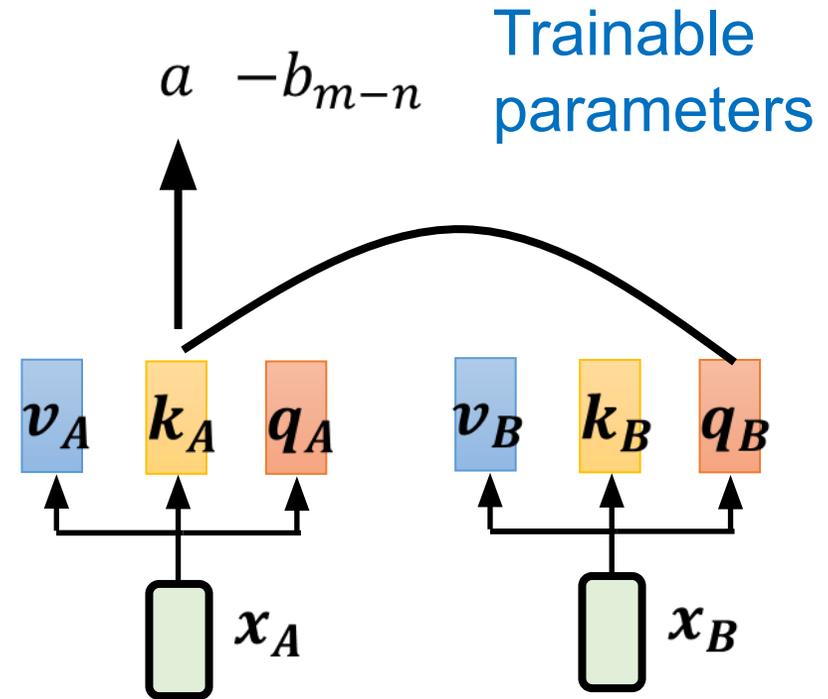


<https://arxiv.org/abs/2108.12409>

Text-to-Text Transfer Transformer (T5)



<https://arxiv.org/abs/2108.12409>



<https://arxiv.org/abs/1910.10683>

RoPE

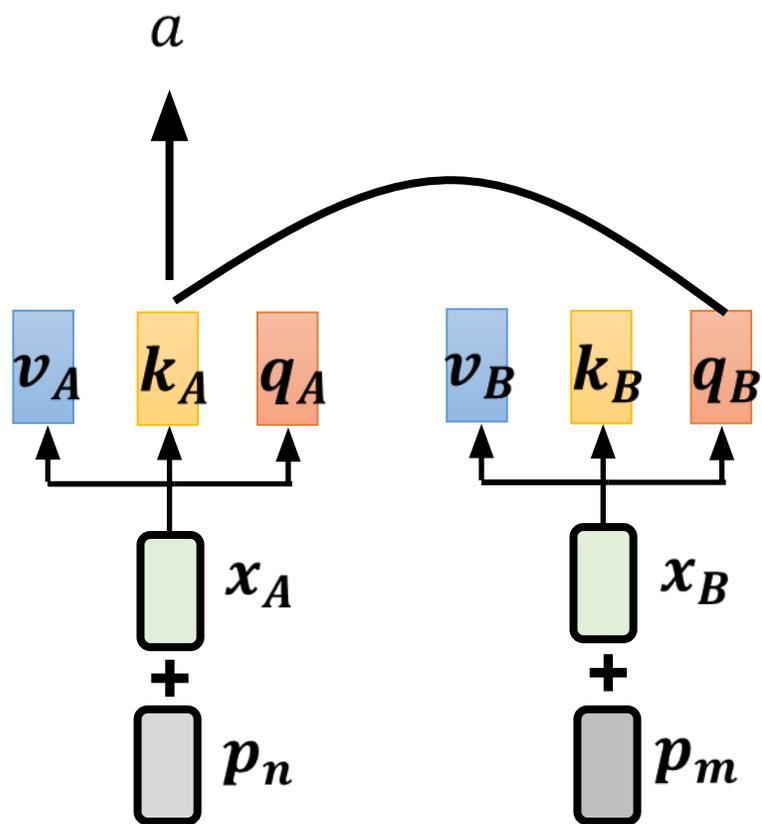
Rotary Position Embedding

Llama, Qwen, Gemma

<https://arxiv.org/abs/2104.09864>

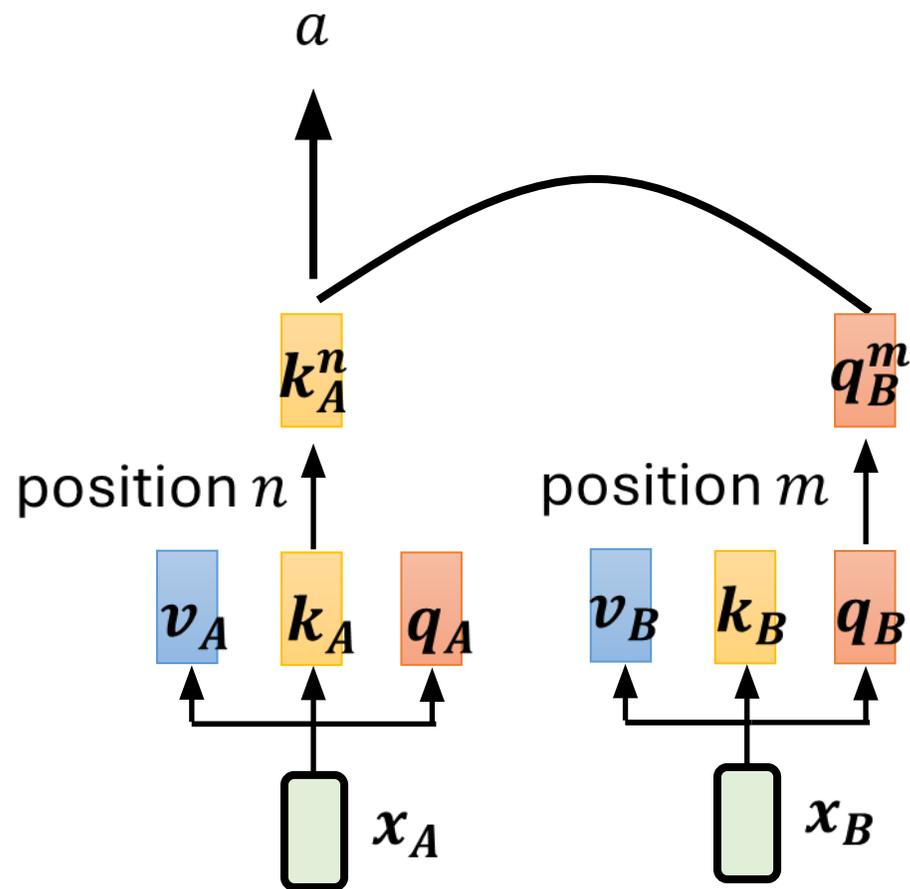
a

$$= \left(W_q(\mathbf{x}_A + \mathbf{p}_m) \right)^T W_k(\mathbf{x}_B + \mathbf{p}_n)$$

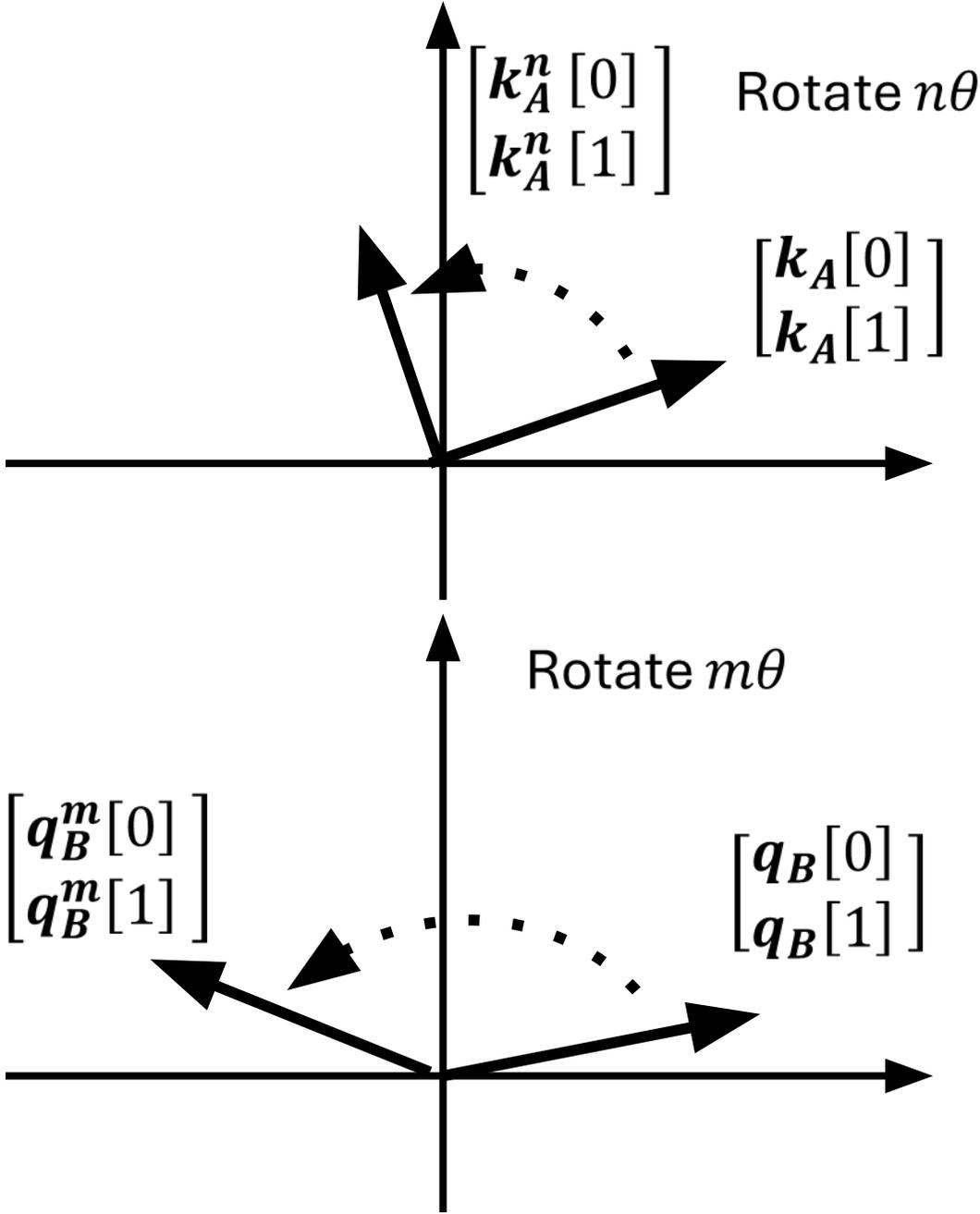
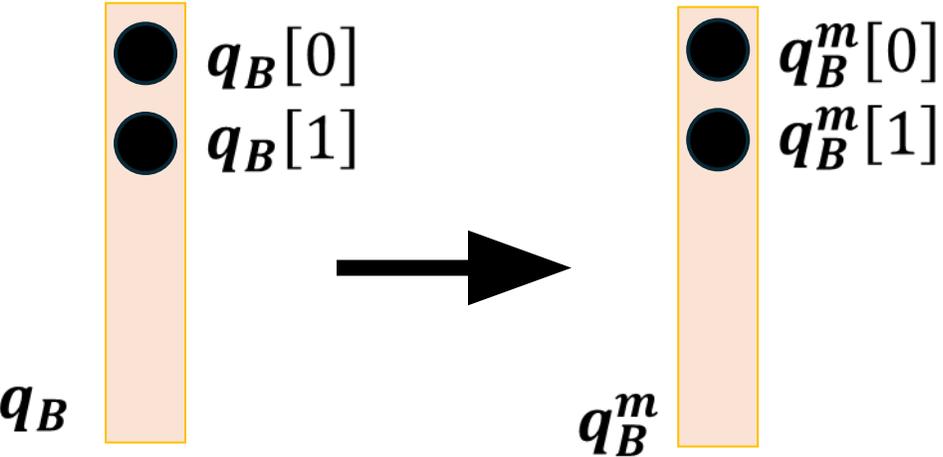
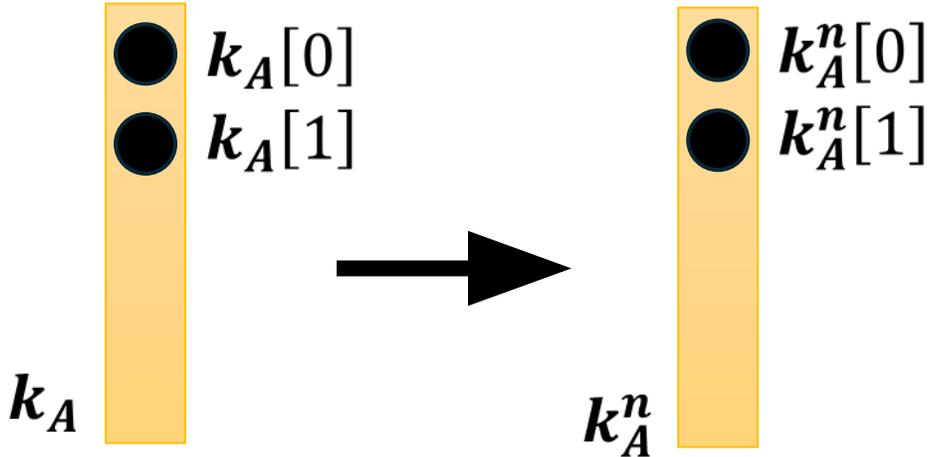


$$a = (\mathbf{k}_A^n)^T \mathbf{q}_B^m$$

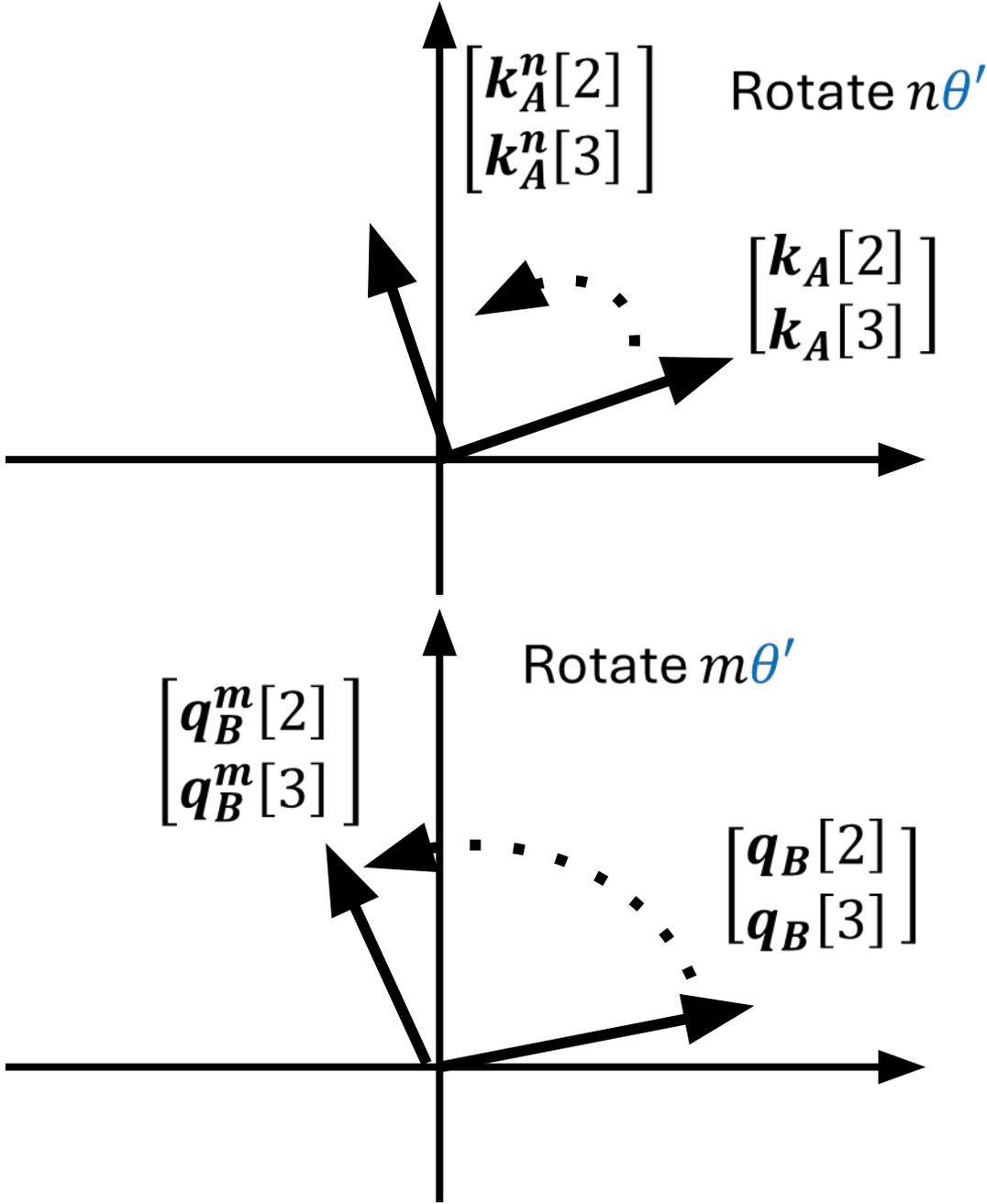
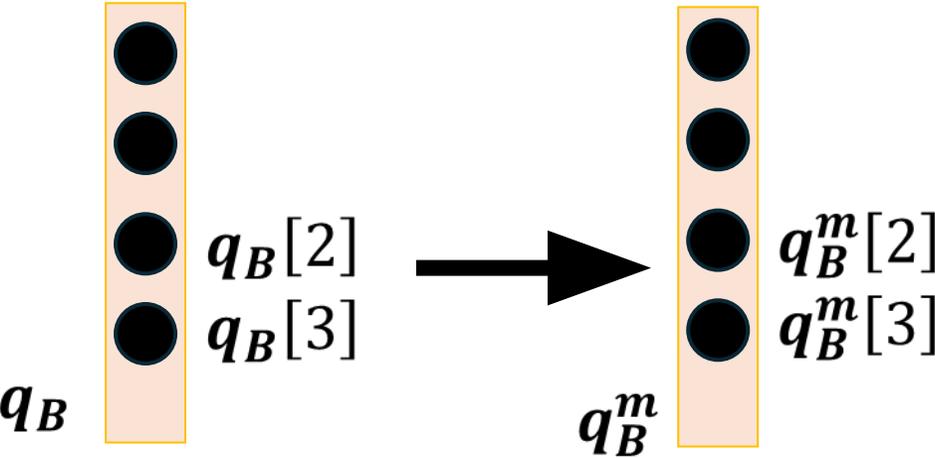
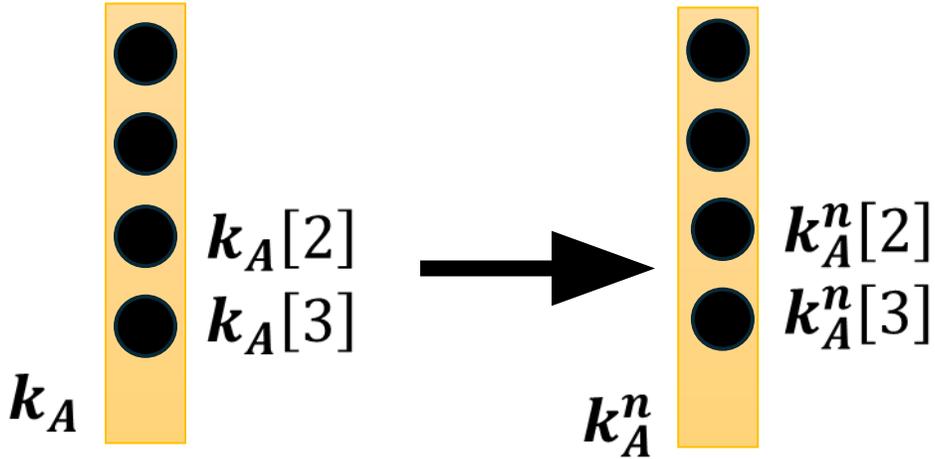
$$= (\mathbf{q}_A)^T R_{m-n} \mathbf{q}_B$$



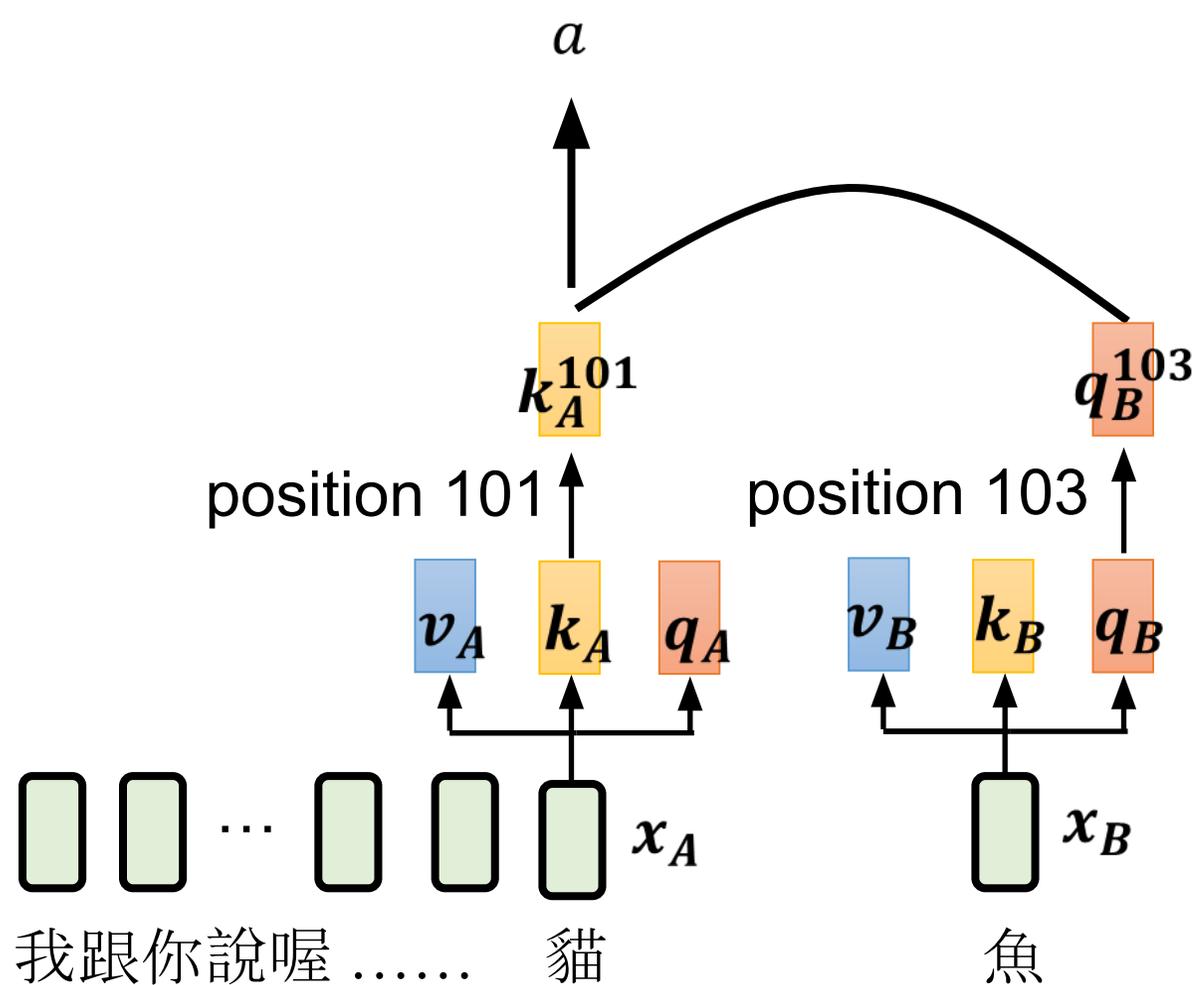
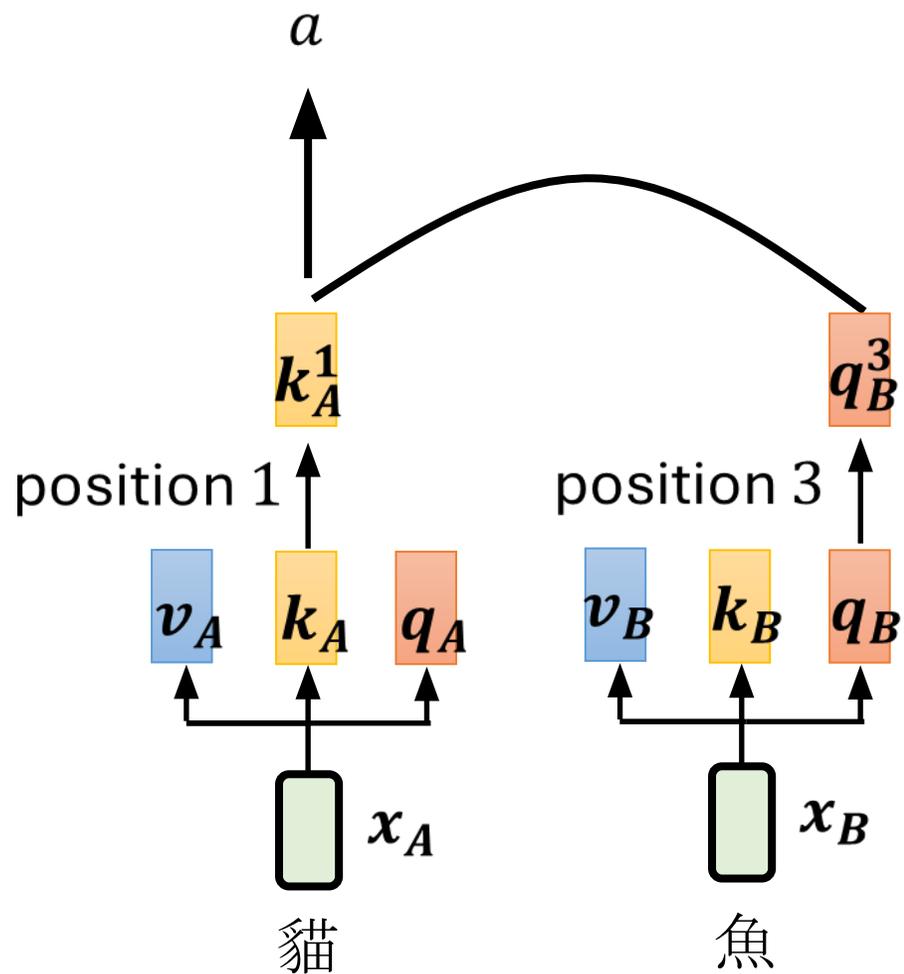
Rotation as Position

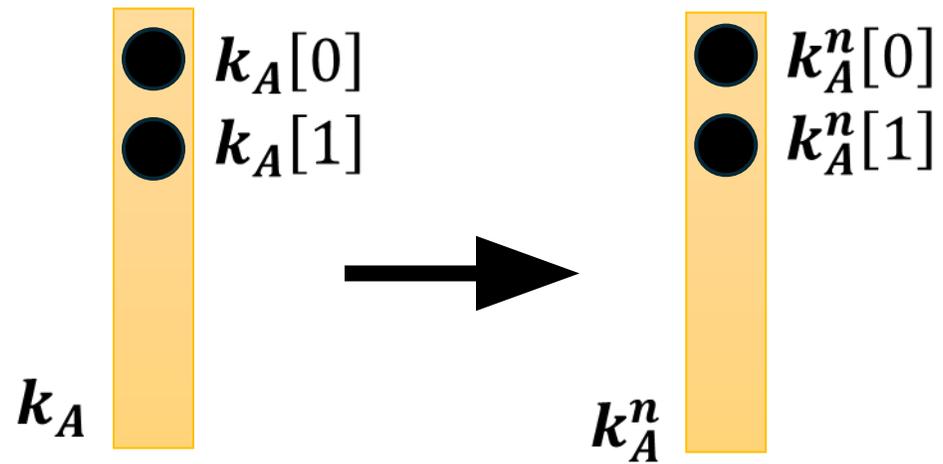


Rotation as Position



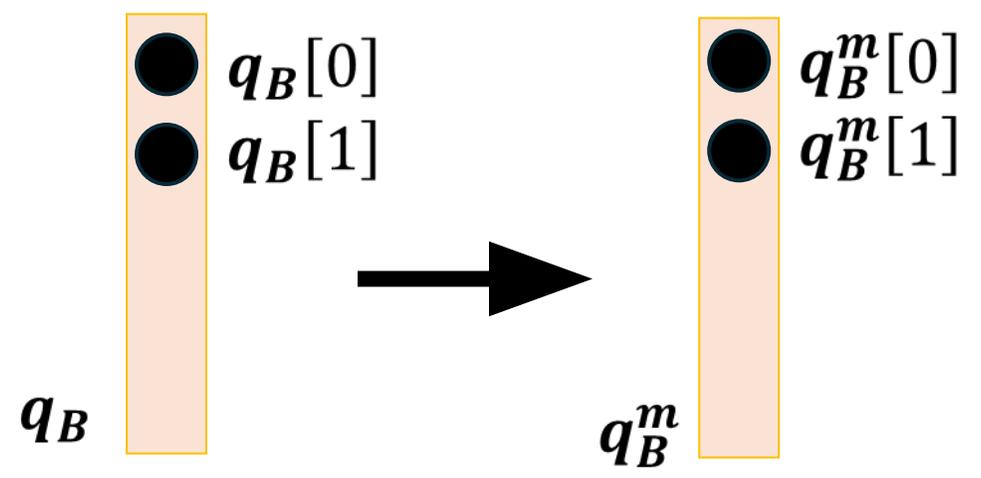
$$k_A^n \cdot q_B^m = k_A^{n+r} \cdot q_B^{m+r}$$





Rotate $n\theta$

$$\begin{bmatrix} k_A^n[0] \\ k_A^n[1] \end{bmatrix} = \begin{bmatrix} \cos(n\theta) & -\sin(n\theta) \\ \sin(n\theta) & \cos(n\theta) \end{bmatrix} \begin{bmatrix} k_A[0] \\ k_A[1] \end{bmatrix}$$



Rotate $m\theta$

$$\begin{bmatrix} q_B^m[0] \\ q_B^m[1] \end{bmatrix} = \begin{bmatrix} \cos(m\theta) & -\sin(m\theta) \\ \sin(m\theta) & \cos(m\theta) \end{bmatrix} \begin{bmatrix} q_B[0] \\ q_B[1] \end{bmatrix}$$

Rotate $n\theta$

$$\begin{bmatrix} \mathbf{k}_A^n[0] \\ \mathbf{k}_A^n[1] \end{bmatrix} = \begin{bmatrix} \cos(n\theta) & -\sin(n\theta) \\ \sin(n\theta) & \cos(n\theta) \end{bmatrix} \begin{bmatrix} \mathbf{k}_A[0] \\ \mathbf{k}_A[1] \end{bmatrix}$$

Rotate $m\theta$

$$\begin{bmatrix} \mathbf{q}_B^m[0] \\ \mathbf{q}_B^m[1] \end{bmatrix} = \begin{bmatrix} \cos(m\theta) & -\sin(m\theta) \\ \sin(m\theta) & \cos(m\theta) \end{bmatrix} \begin{bmatrix} \mathbf{q}_B[0] \\ \mathbf{q}_B[1] \end{bmatrix}$$

$$\begin{bmatrix} \mathbf{k}_A^n[0] \\ \mathbf{k}_A^n[1] \end{bmatrix} \cdot \begin{bmatrix} \mathbf{q}_B^m[0] \\ \mathbf{q}_B^m[1] \end{bmatrix} = \begin{bmatrix} \mathbf{k}_A^n[0] \\ \mathbf{k}_A^n[1] \end{bmatrix}^T \begin{bmatrix} \mathbf{q}_B^m[0] \\ \mathbf{q}_B^m[1] \end{bmatrix}$$

$$= \left(\begin{bmatrix} \cos(n\theta) & -\sin(n\theta) \\ \sin(n\theta) & \cos(n\theta) \end{bmatrix} \begin{bmatrix} \mathbf{k}_A[0] \\ \mathbf{k}_A[1] \end{bmatrix} \right)^T \begin{bmatrix} \cos(m\theta) & -\sin(m\theta) \\ \sin(m\theta) & \cos(m\theta) \end{bmatrix} \begin{bmatrix} \mathbf{q}_B[0] \\ \mathbf{q}_B[1] \end{bmatrix}$$

$$= \begin{bmatrix} \mathbf{k}_A[0] \\ \mathbf{k}_A[1] \end{bmatrix}^T \begin{bmatrix} \cos(n\theta) & -\sin(n\theta) \\ \sin(n\theta) & \cos(n\theta) \end{bmatrix}^T \begin{bmatrix} \cos(m\theta) & -\sin(m\theta) \\ \sin(m\theta) & \cos(m\theta) \end{bmatrix} \begin{bmatrix} \mathbf{q}_B[0] \\ \mathbf{q}_B[1] \end{bmatrix}$$

Rotate $-n\theta$

Rotate $m\theta$

$$= \begin{bmatrix} \mathbf{k}_A[0] \\ \mathbf{k}_A[1] \end{bmatrix}^T \begin{bmatrix} \cos((m-n)\theta) & -\sin((m-n)\theta) \\ \sin((m-n)\theta) & \cos((m-n)\theta) \end{bmatrix} \begin{bmatrix} \mathbf{q}_B[0] \\ \mathbf{q}_B[1] \end{bmatrix}$$

Rotate $(m-n)\theta$

Rotate $n\theta$

$$\begin{bmatrix} \mathbf{k}_A^n[0] \\ \mathbf{k}_A^n[1] \end{bmatrix} = \begin{bmatrix} \cos(n\theta) & -\sin(n\theta) \\ \sin(n\theta) & \cos(n\theta) \end{bmatrix} \begin{bmatrix} \mathbf{k}_A[0] \\ \mathbf{k}_A[1] \end{bmatrix}$$

Rotate $m\theta$

$$\begin{bmatrix} \mathbf{q}_B^m[0] \\ \mathbf{q}_B^m[1] \end{bmatrix} = \begin{bmatrix} \cos(m\theta) & -\sin(m\theta) \\ \sin(m\theta) & \cos(m\theta) \end{bmatrix} \begin{bmatrix} \mathbf{q}_B[0] \\ \mathbf{q}_B[1] \end{bmatrix}$$

$$\begin{bmatrix} \mathbf{k}_A^n[0] \\ \mathbf{k}_A^n[1] \end{bmatrix} \cdot \begin{bmatrix} \mathbf{q}_B^m[0] \\ \mathbf{q}_B^m[1] \end{bmatrix} = \begin{bmatrix} \mathbf{k}_A[0] \\ \mathbf{k}_A[1] \end{bmatrix}^T \begin{bmatrix} \cos((m-n)\theta) & -\sin((m-n)\theta) \\ \sin((m-n)\theta) & \cos((m-n)\theta) \end{bmatrix} \begin{bmatrix} \mathbf{q}_B[0] \\ \mathbf{q}_B[1] \end{bmatrix}$$

Rotate $(m-n)\theta$

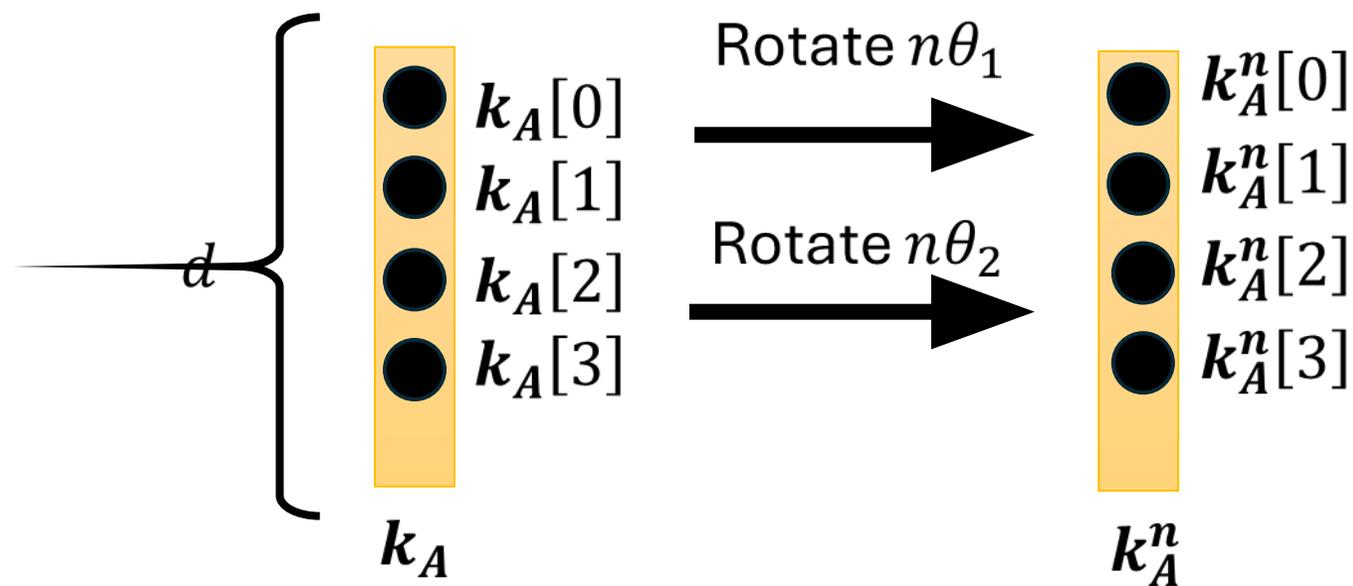
$$\begin{bmatrix} \mathbf{k}_A^{n+r}[0] \\ \mathbf{k}_A^{n+r}[1] \end{bmatrix} \cdot \begin{bmatrix} \mathbf{q}_B^{m+r}[0] \\ \mathbf{q}_B^{m+r}[1] \end{bmatrix}$$

$$= \begin{bmatrix} \mathbf{k}_A[0] \\ \mathbf{k}_A[1] \end{bmatrix}^T \begin{bmatrix} \cos(((m+r)-(n+r))\theta) & -\sin(((m+r)-(n+r))\theta) \\ \sin(((m+r)-(n+r))\theta) & \cos(((m+r)-(n+r))\theta) \end{bmatrix} \begin{bmatrix} \mathbf{q}_B[0] \\ \mathbf{q}_B[1] \end{bmatrix}$$

$$\text{Rotate } ((m+r)-(n+r))\theta = \text{Rotate } (m-n)\theta$$

旋轉角度怎麼設定

<https://arxiv.org/abs/2104.09864>



$$\theta_i = \frac{1}{10000^{2i/d}}$$

$$i = 0, 1, \dots, \frac{d}{2} - 1$$

RoPE 沒有越遠越小

<https://arxiv.org/pdf/2410.06205>

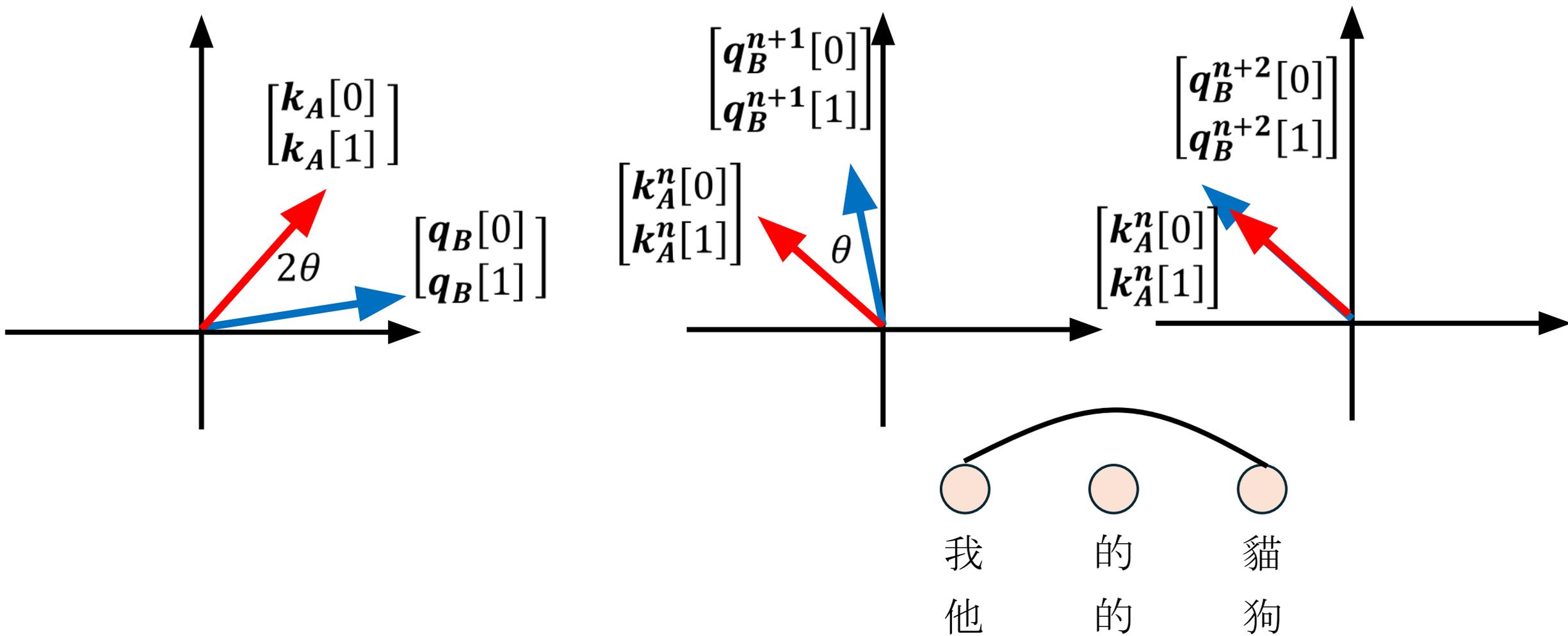
$$\begin{bmatrix} \mathbf{k}_A^n[0] \\ \mathbf{k}_A^n[1] \end{bmatrix} \cdot \begin{bmatrix} \mathbf{q}_B^m[0] \\ \mathbf{q}_B^m[1] \end{bmatrix} = \begin{bmatrix} \mathbf{k}_A[0] \\ \mathbf{k}_A[1] \end{bmatrix}^T \begin{bmatrix} \cos((m-n)\theta) & -\sin((m-n)\theta) \\ \sin((m-n)\theta) & \cos((m-n)\theta) \end{bmatrix} \begin{bmatrix} \mathbf{q}_B[0] \\ \mathbf{q}_B[1] \end{bmatrix}$$

Rotate $(m-n)\theta$



<https://colab.research.google.com/drive/1rWDtAkScrb2K3tcprSTzwwuRQo5bGiyKJ?usp=sharing>

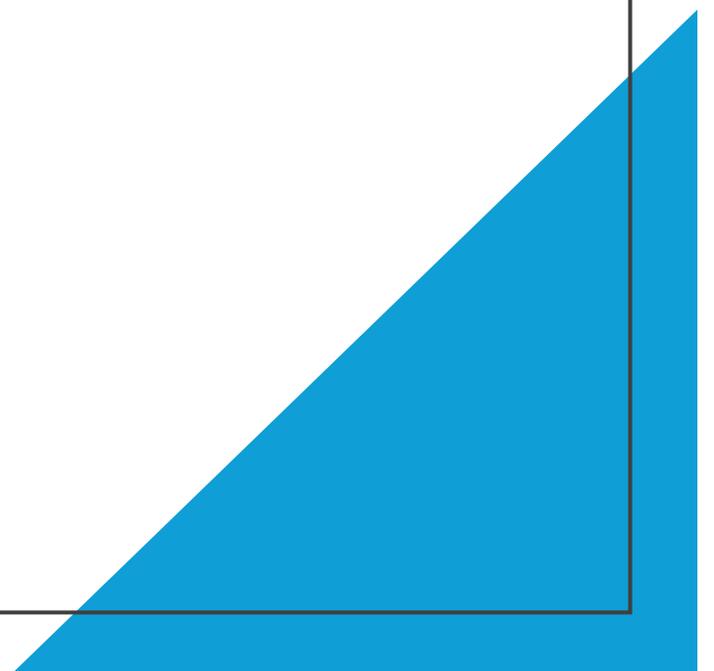
RoPE 沒有越遠越小，不一定是一件壞事



Train Short, Test Long

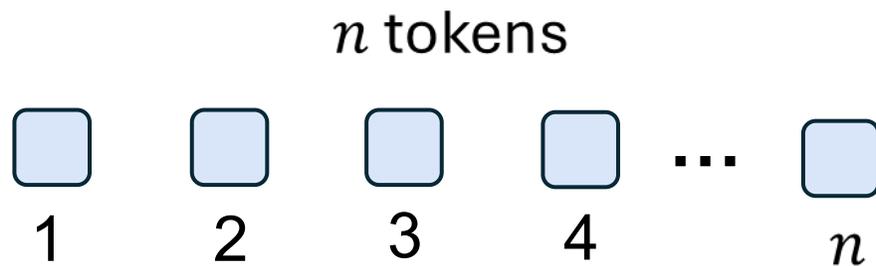
Reference:

<https://amaarora.github.io/posts/2025-09-21-rope-context-extension.html>

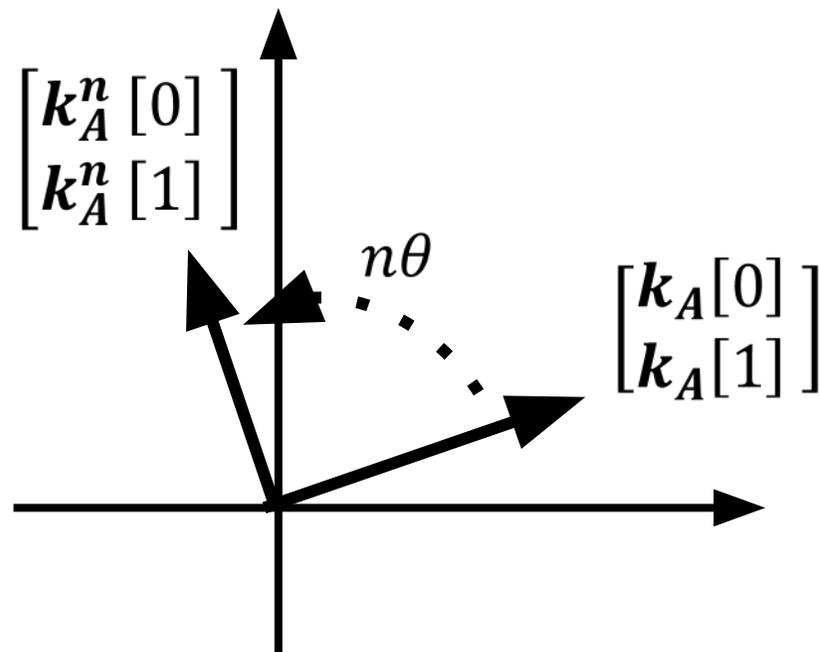
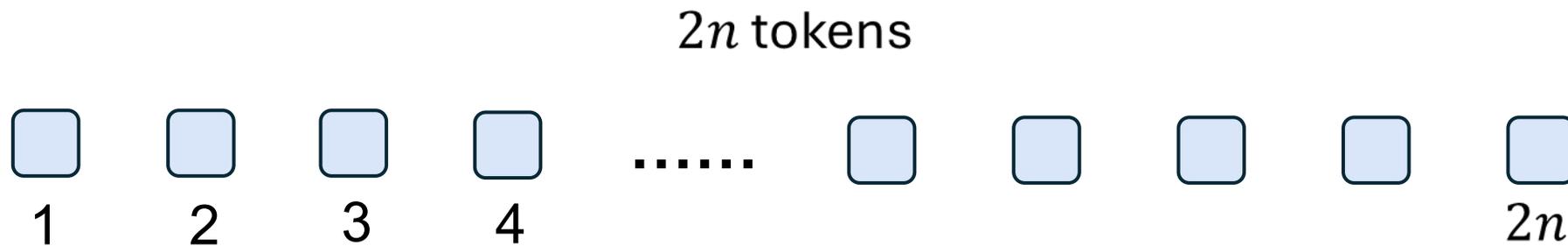


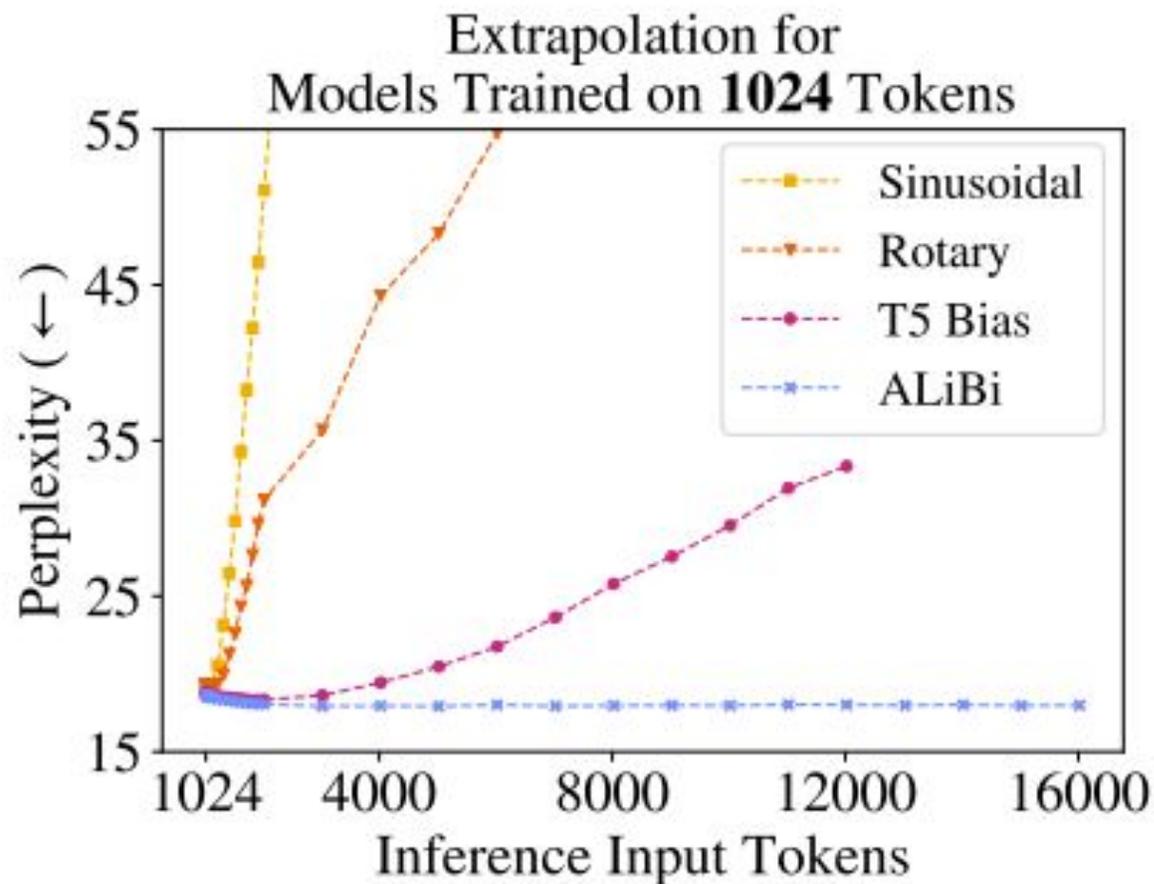
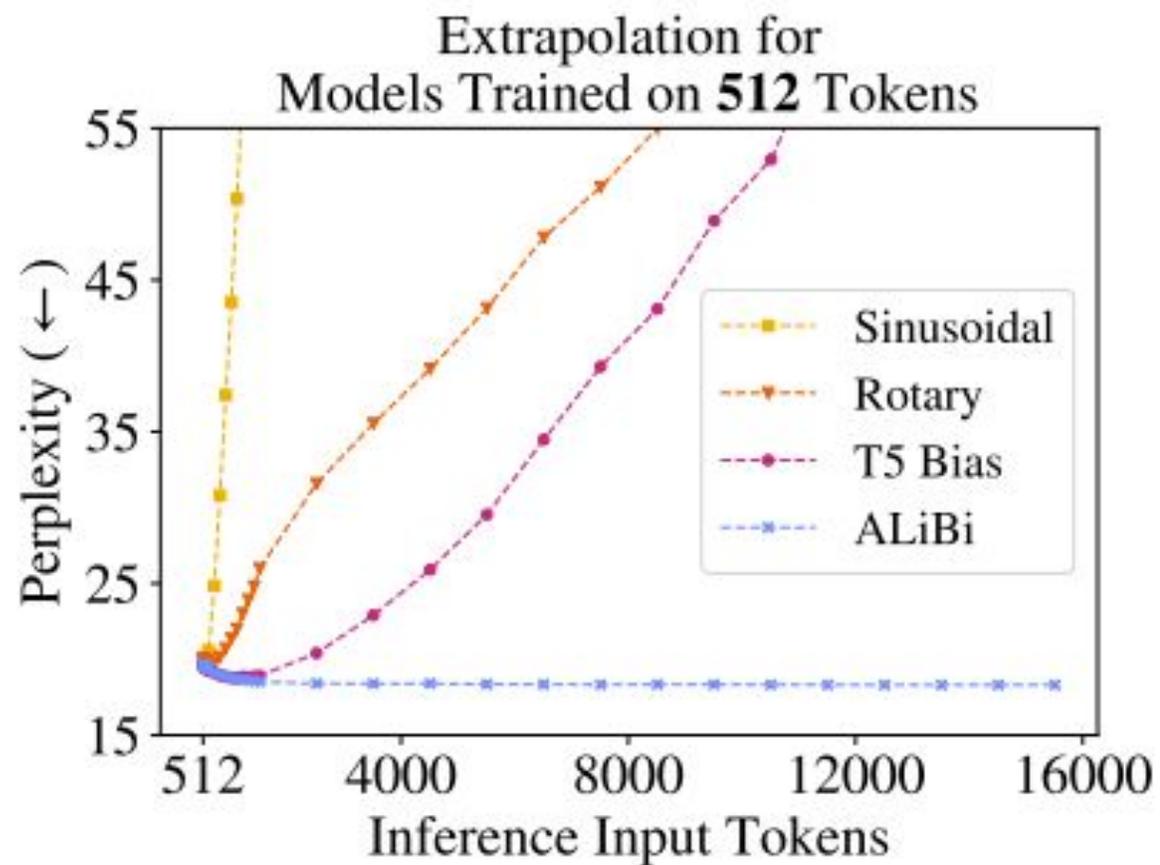
Train Short, Test Long

Train

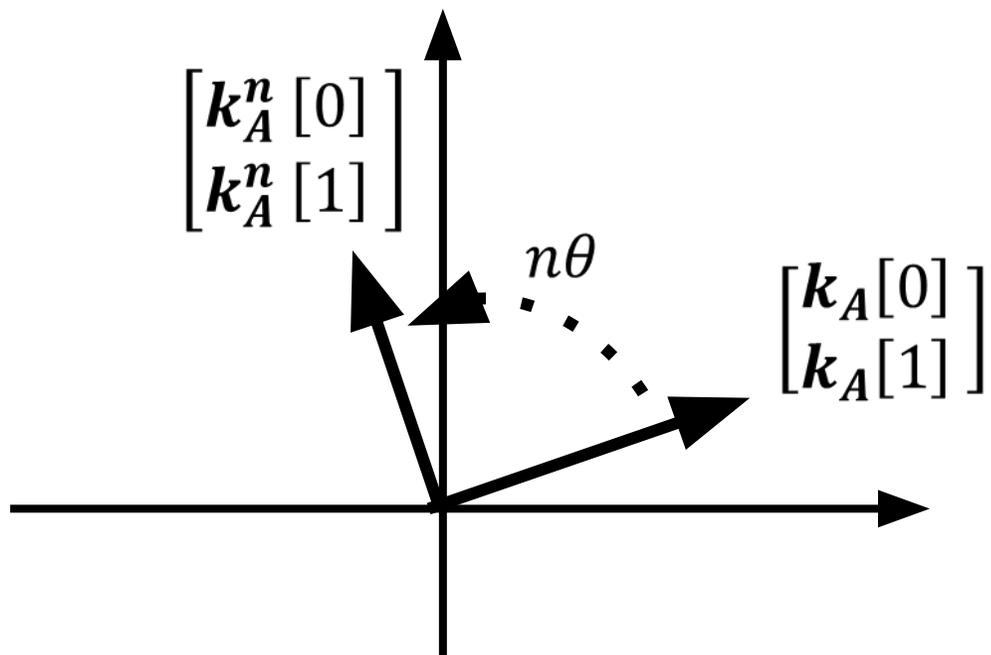


Test

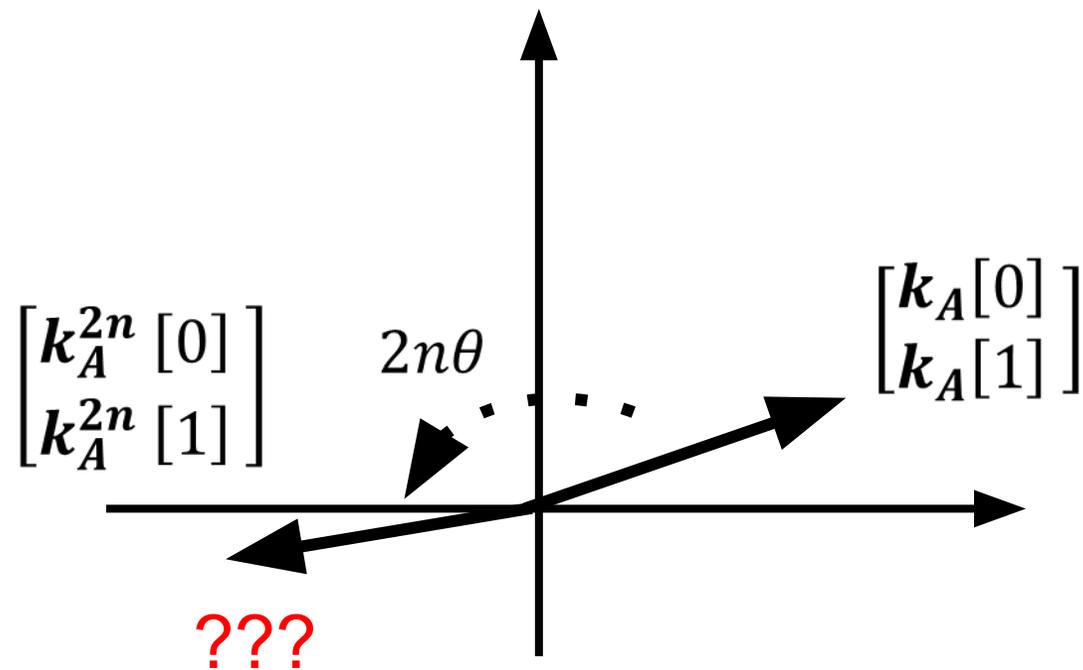




Why Fail?



Training: rotate at most $n\theta$

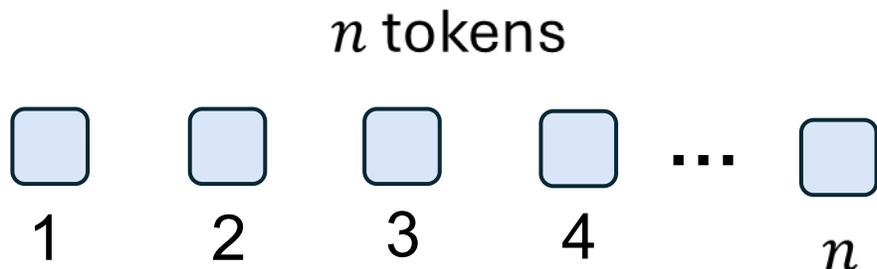


Testing: rotate $2n\theta$

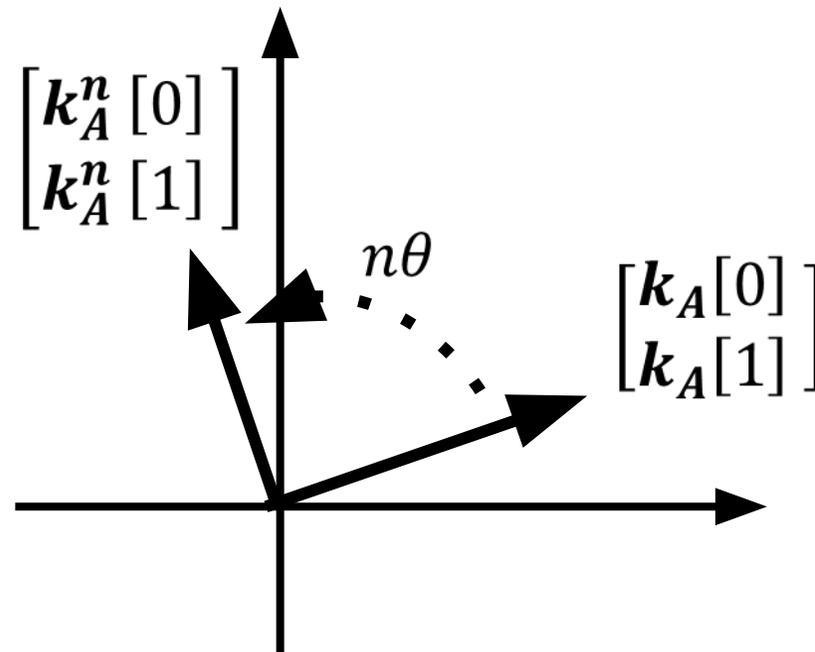
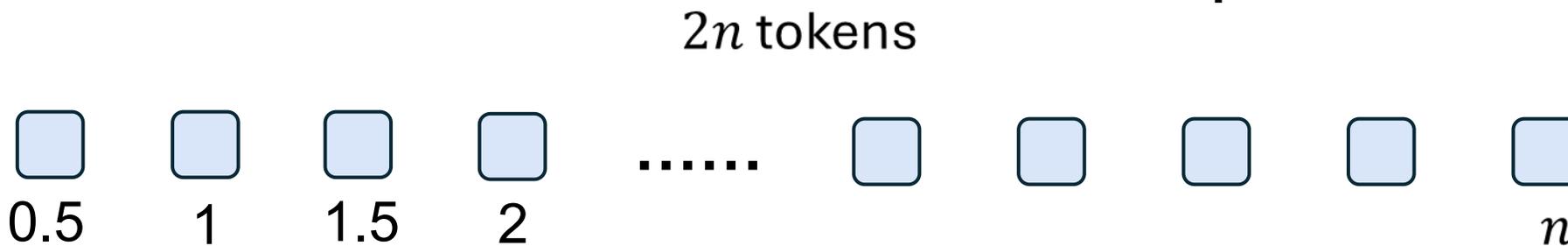
Position Interpolation

仍然需要微調模型參數

Train



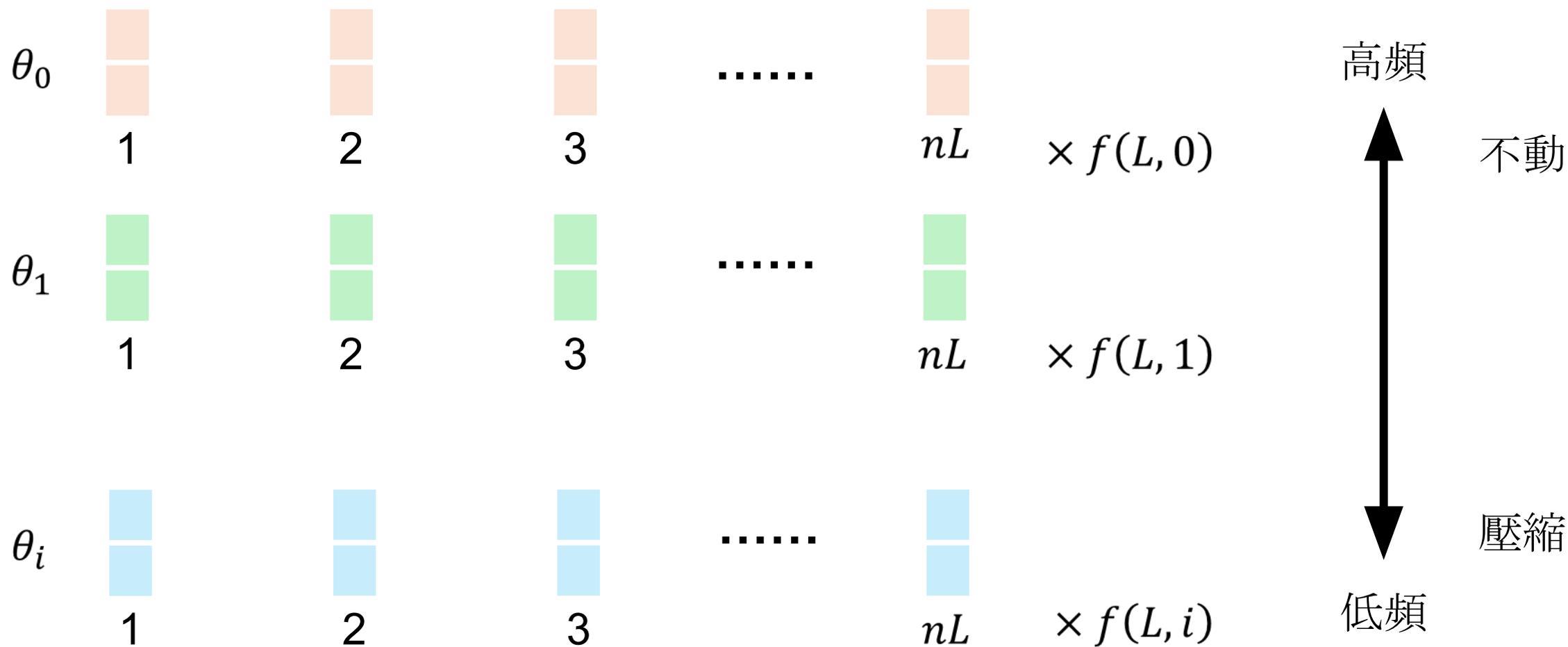
Test



伸長 L 倍，position index 乘上 $1/L$

Frequency-Based Approach

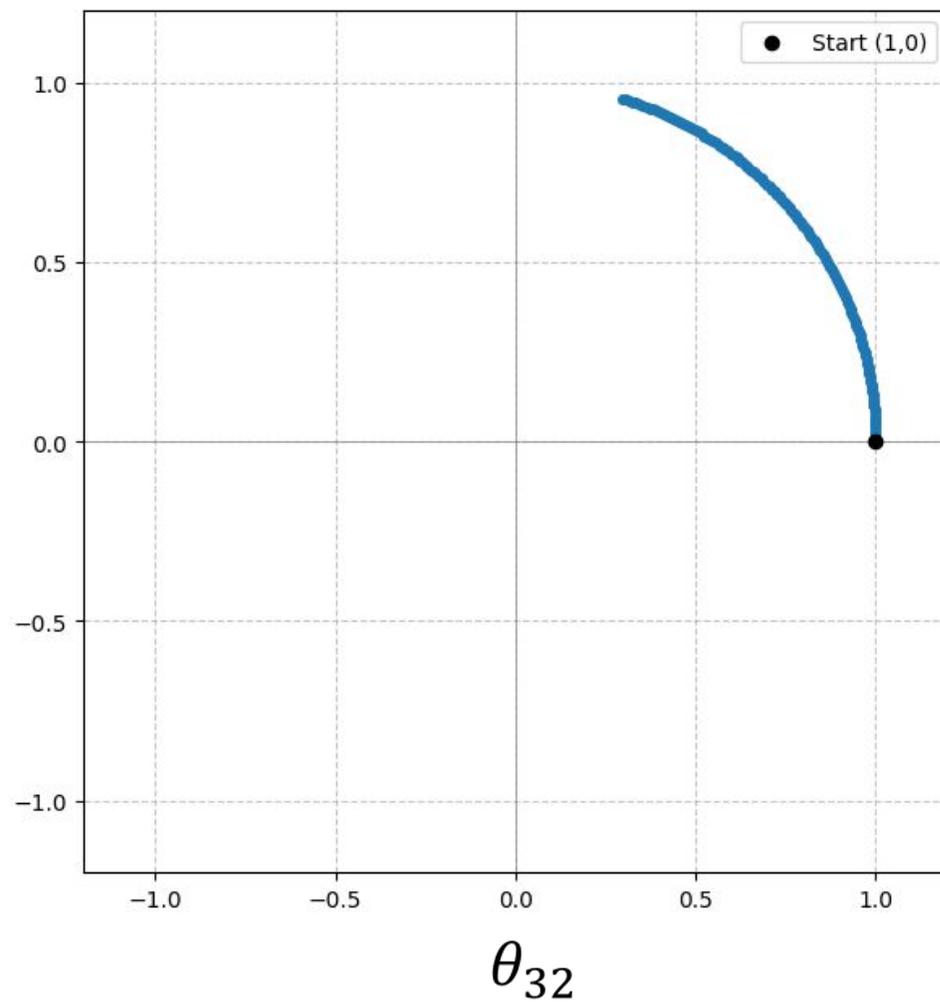
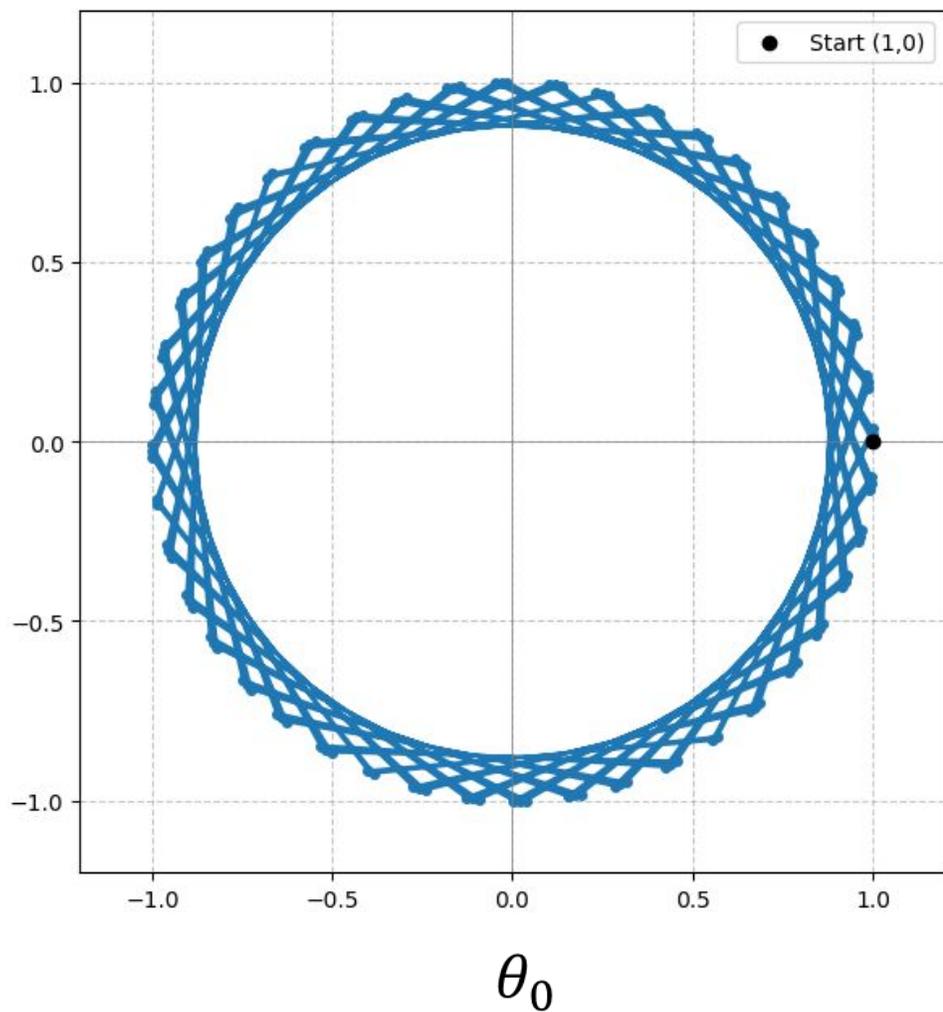
$$\theta_i = \frac{1}{10000^{2i/d}}$$
$$i = 0, 1, \dots, \frac{d}{2} - 1$$



Frequency-Based Approach

$$\theta_i = \frac{1}{10000^{2i/d}}$$

$$i = 0, 1, \dots, \frac{d}{2} - 1$$

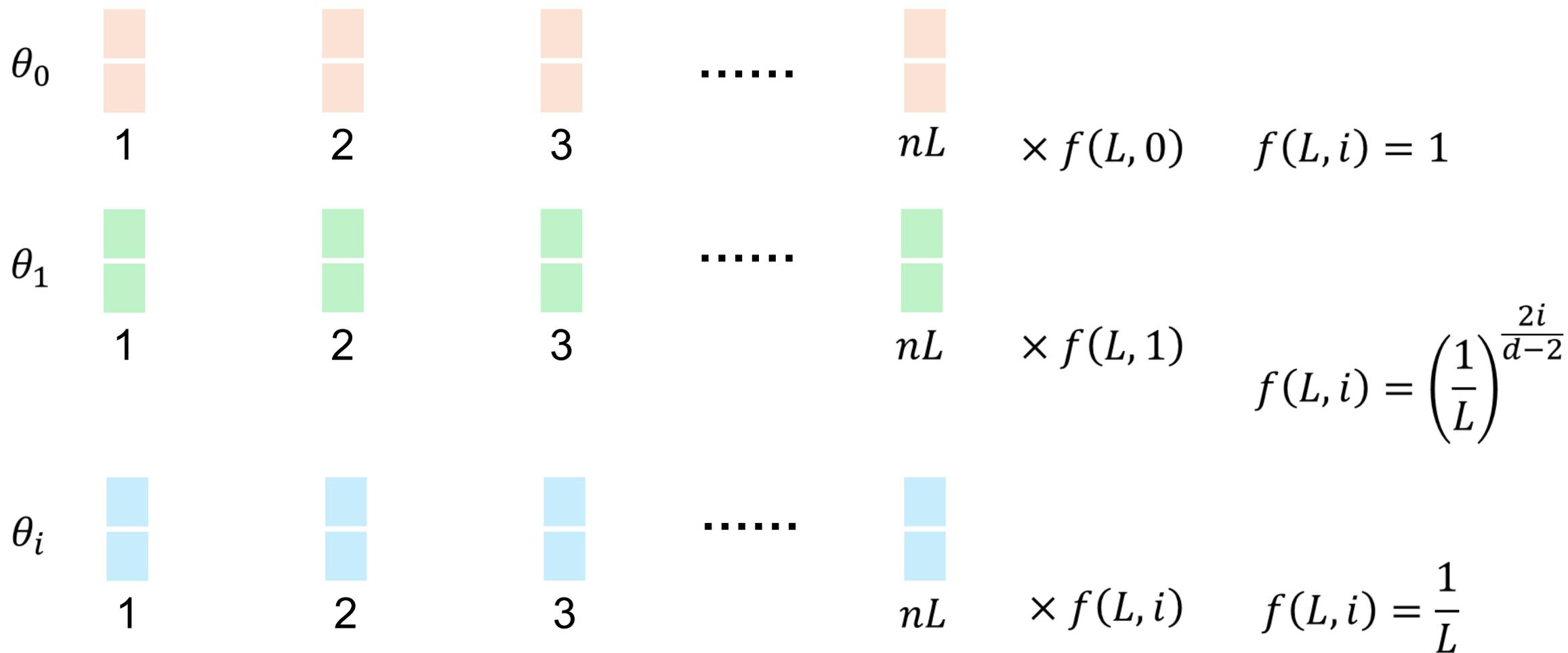


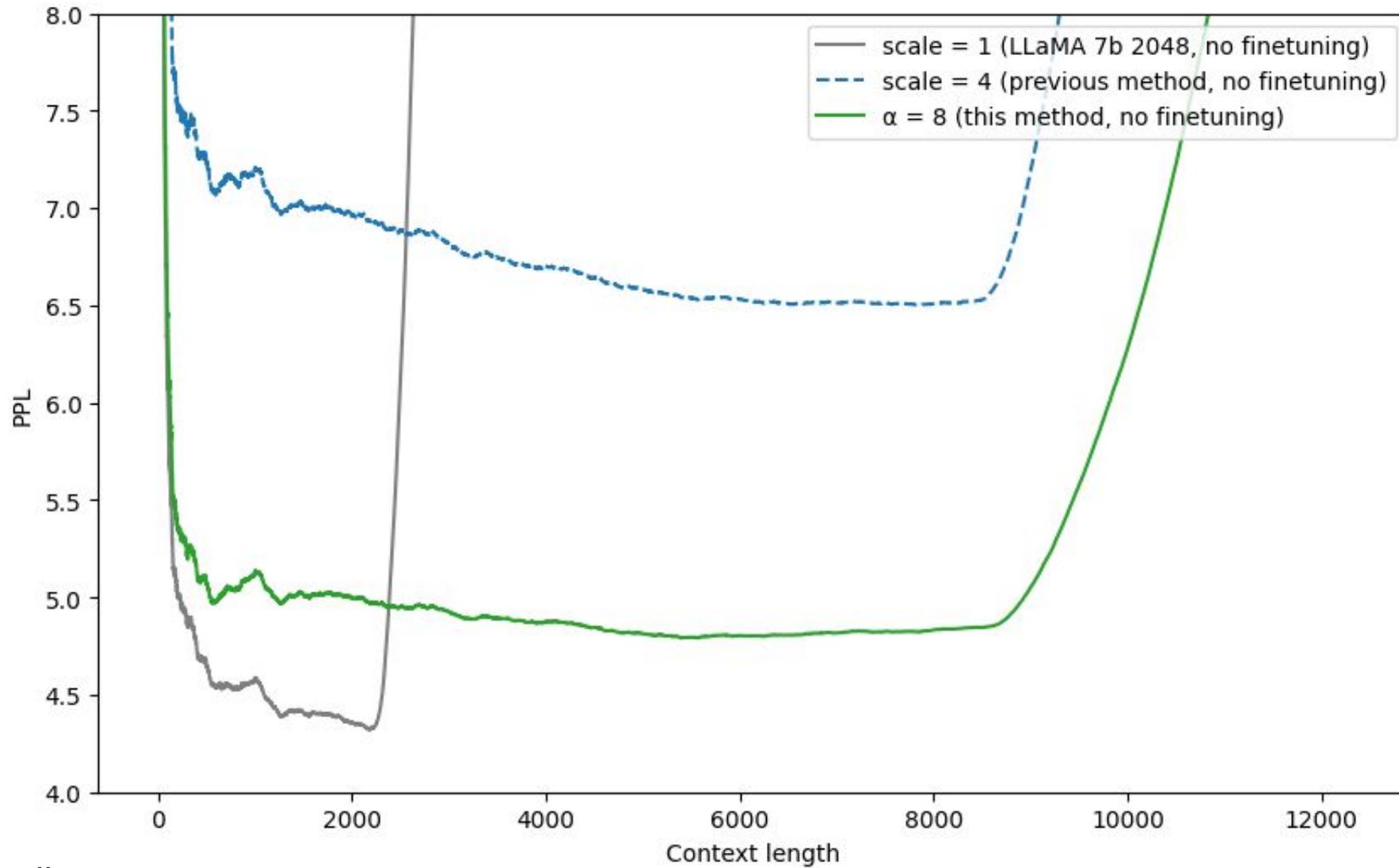
Frequency-Based Approach

NTK-Aware Scaling

$$\theta_i = \frac{1}{10000^{2i/d}}$$

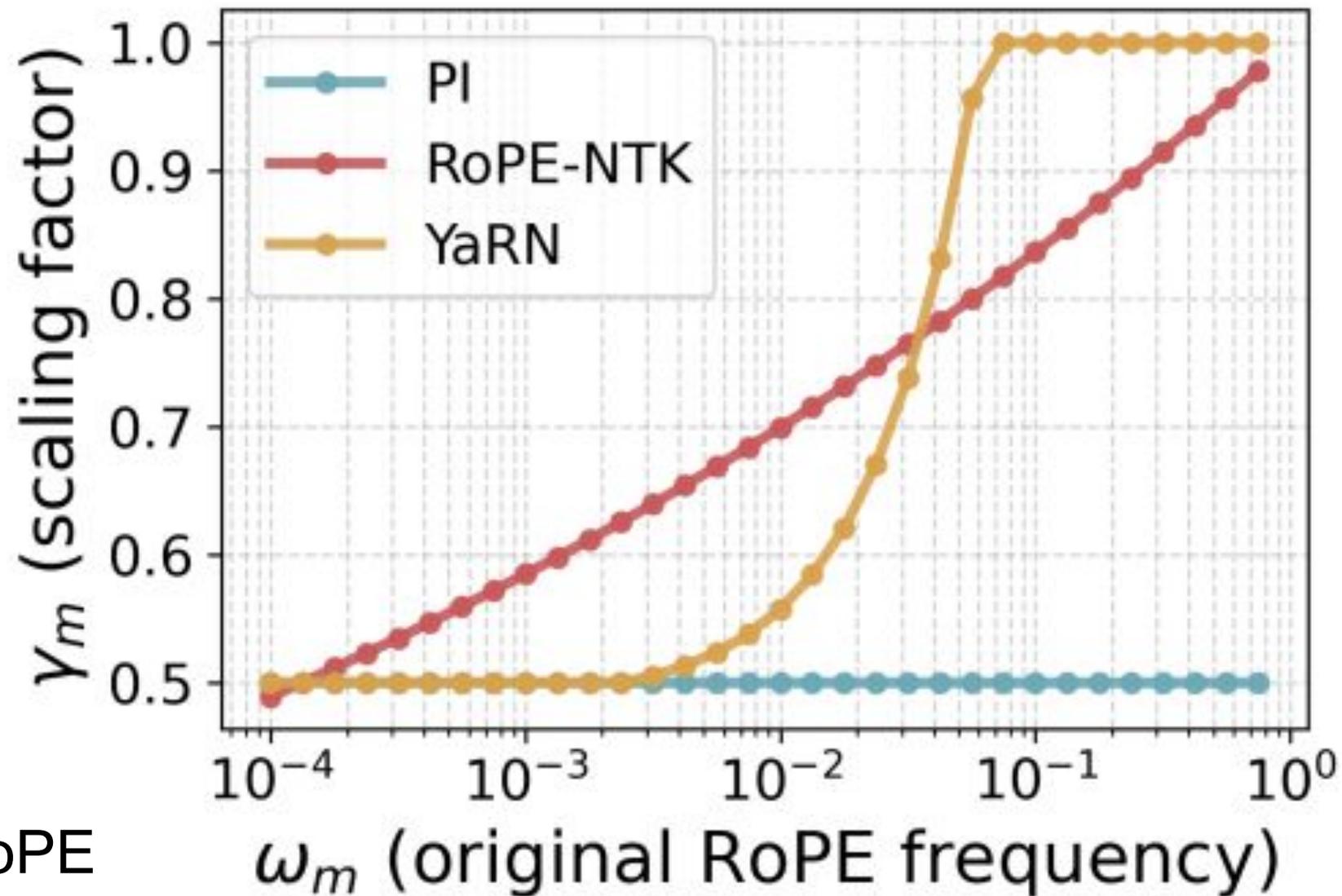
$$i = 0, 1, \dots, \frac{d}{2} - 1$$





NTK-Aware Scaling

https://www.reddit.com/r/LocalLLaMA/comments/14lz7j5/ntkaware_scaled_rope_allows_llama_models_to_have/



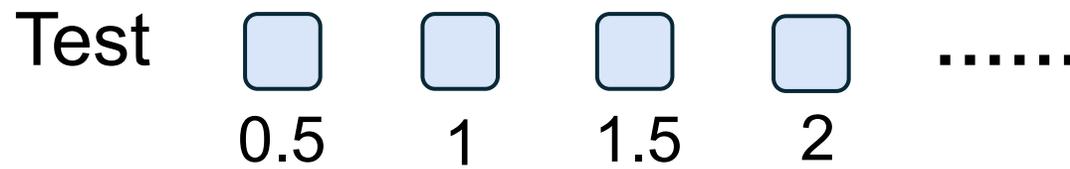
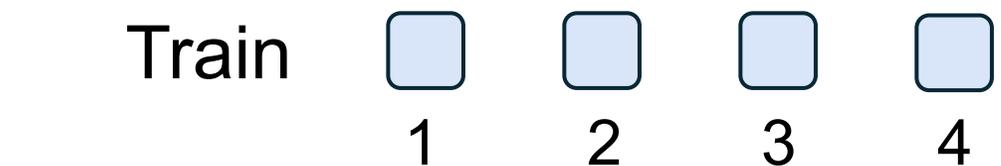
YaRN (Yet another RoPE extension method)

<https://arxiv.org/abs/2309.00071>

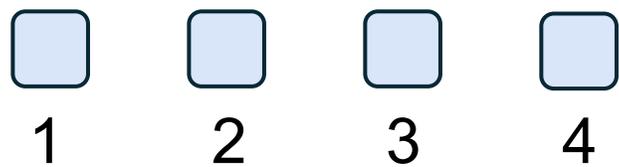
Source of figure: <https://arxiv.org/pdf/2512.12167>

Dynmaic Scaling

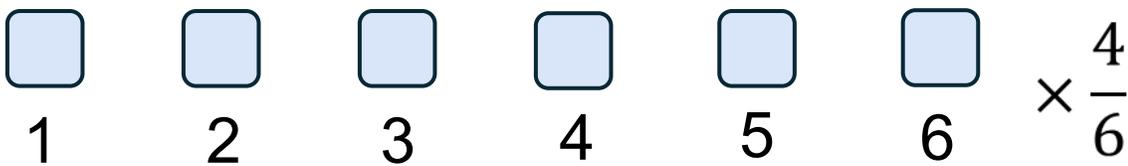
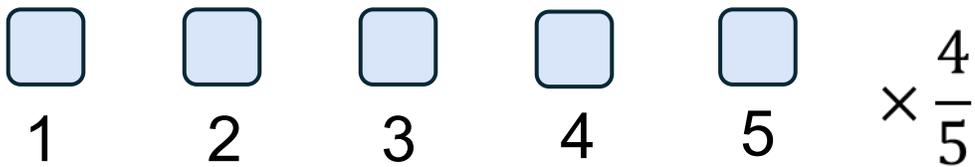
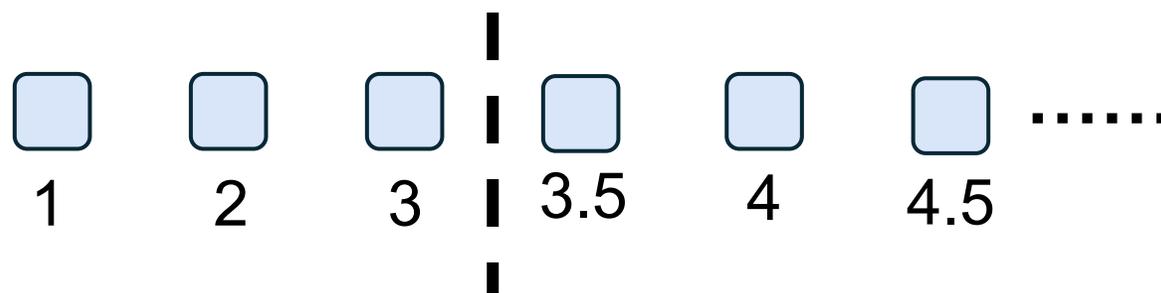
https://www.reddit.com/r/LocalLLaMA/comments/14mrgpr/dynamically_scaled_rope_further_increases/

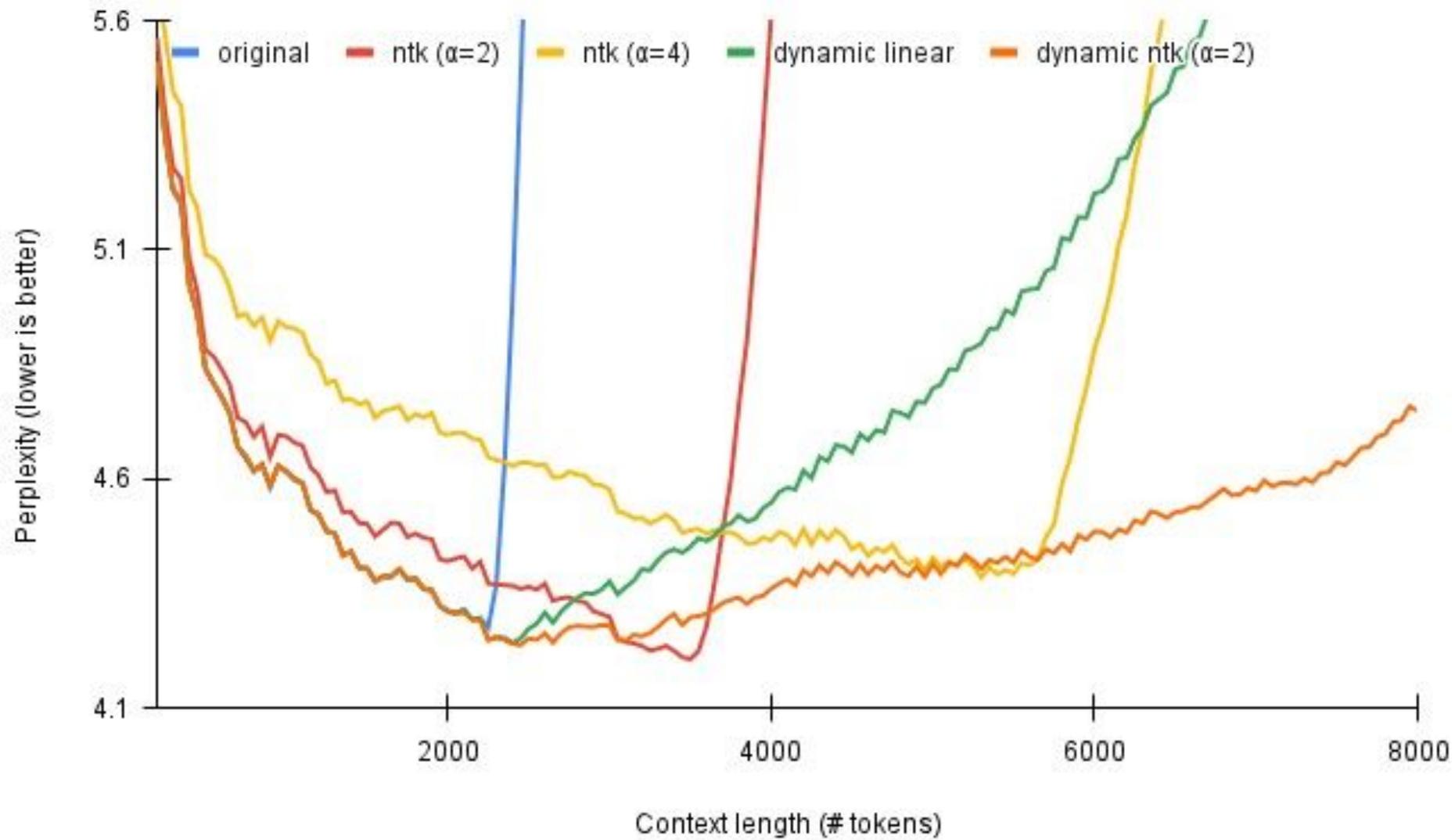


Test



Test





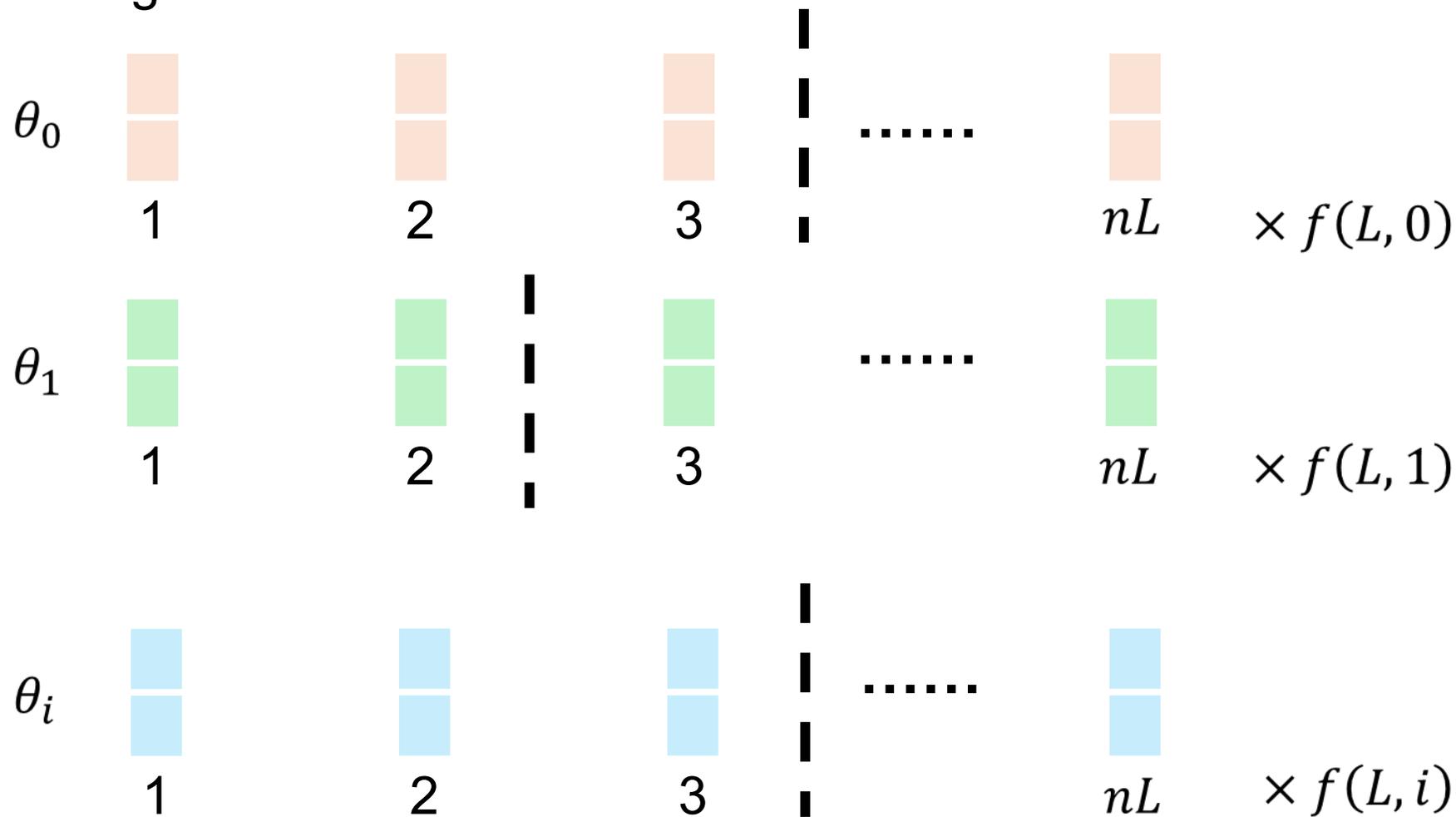
https://www.reddit.com/r/LocalLLaMA/comments/14mrgpr/dynamically_scaled_rope_further_increases/

Frequency-Based + Dynamic

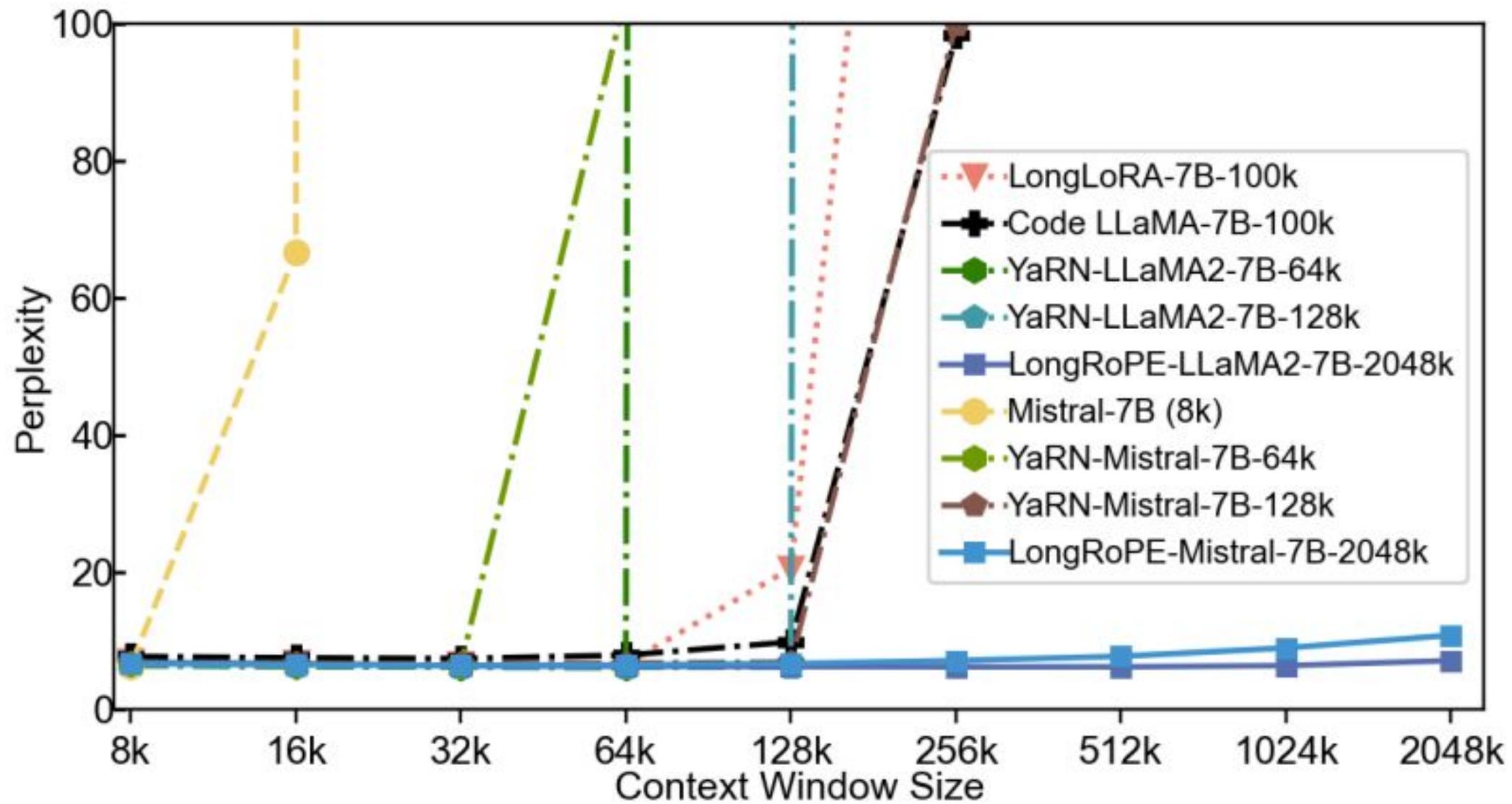
LongRoPE

$$\theta_i = \frac{1}{10000^{2i/d}}$$

$$i = 0, 1, \dots, \frac{d}{2} - 1$$

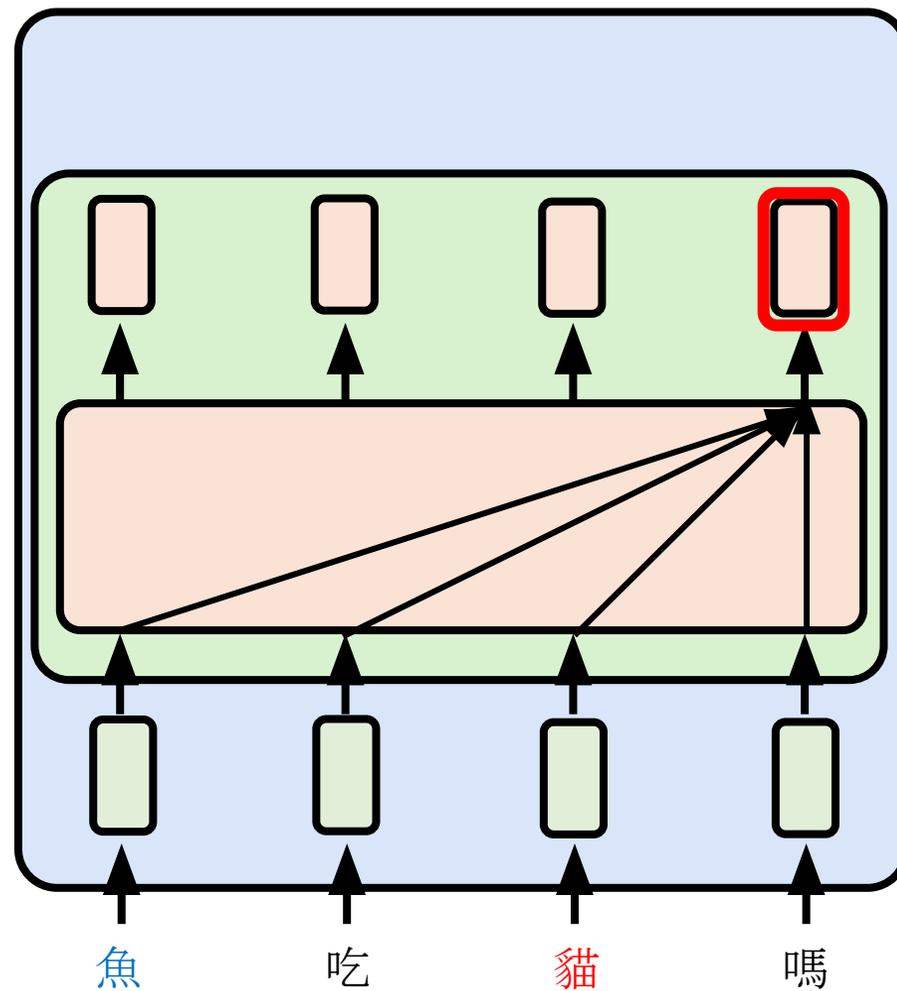
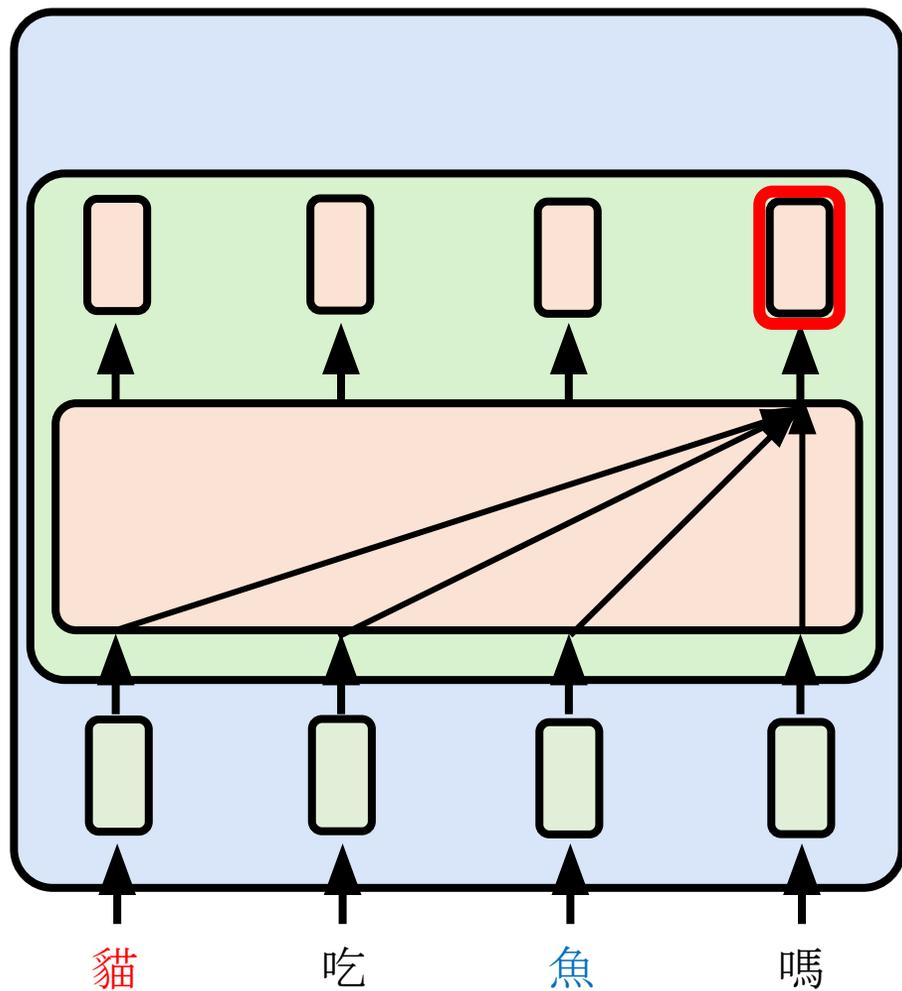


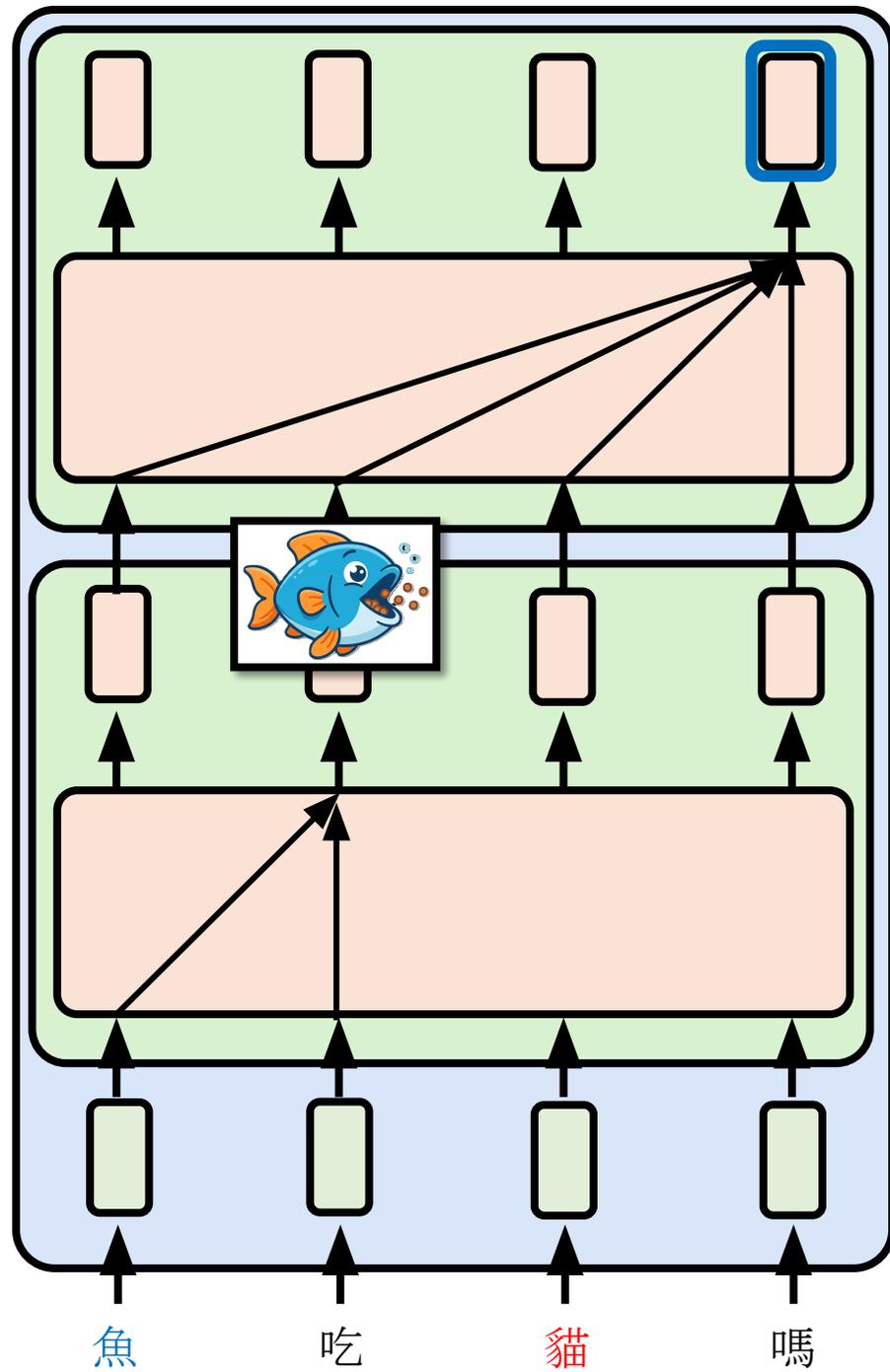
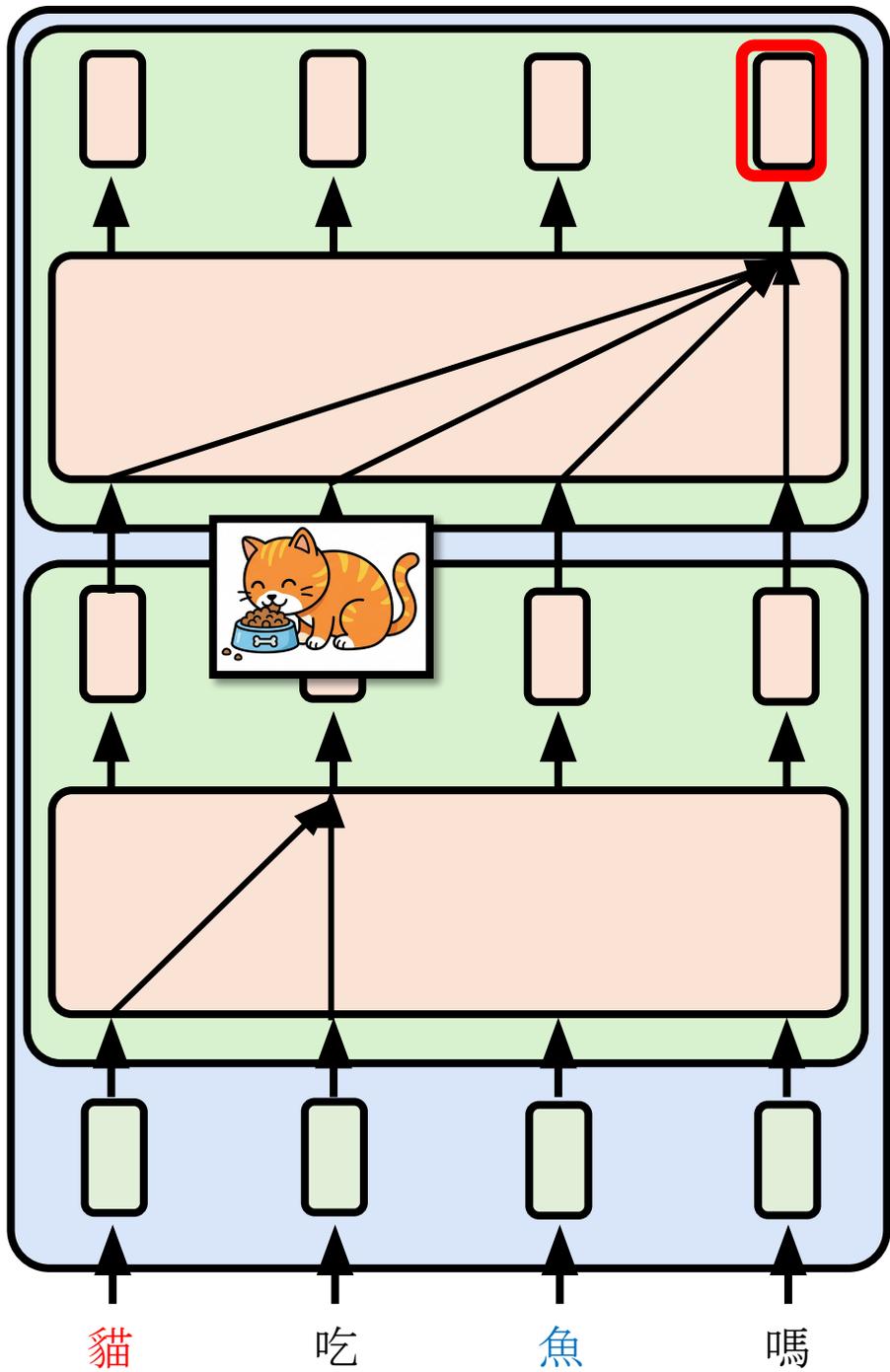
**Use
evolutionary
search**

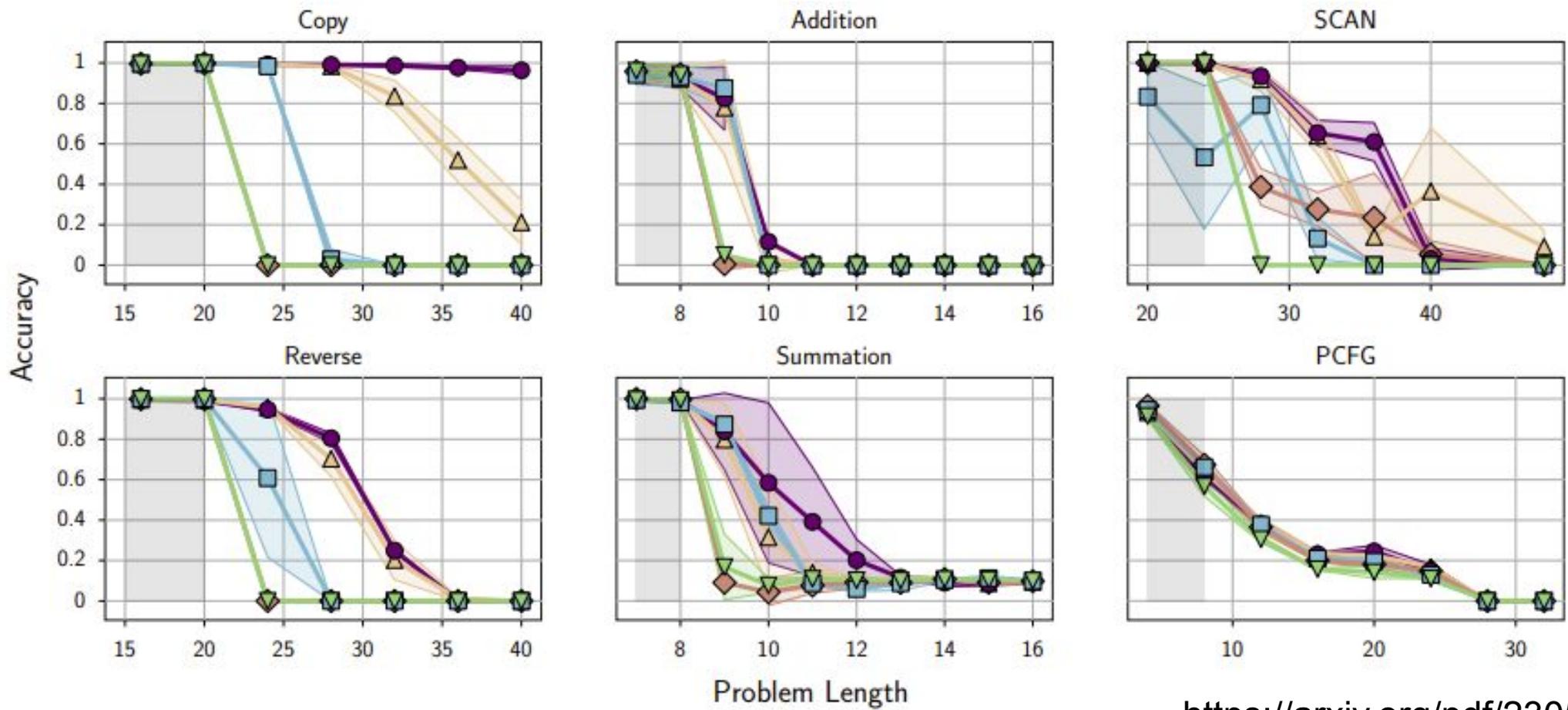


**No Positional
Embedding?!**

Self-attention 真的沒有位置資訊？

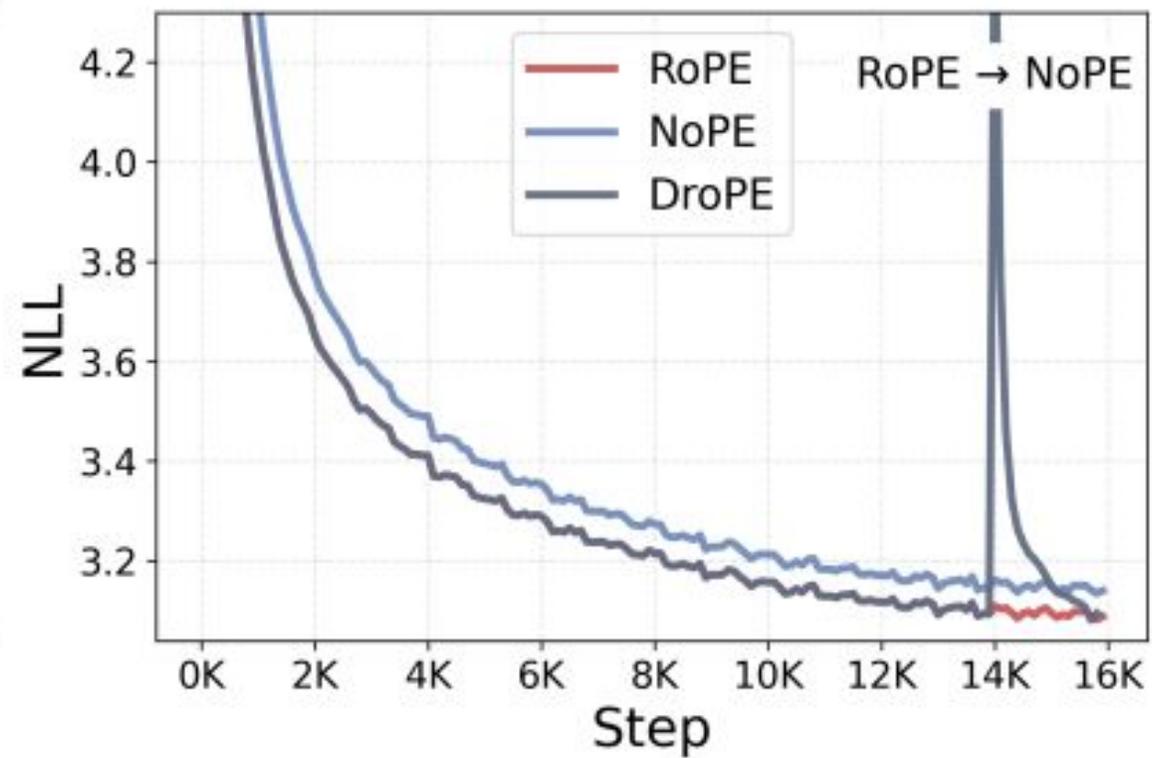
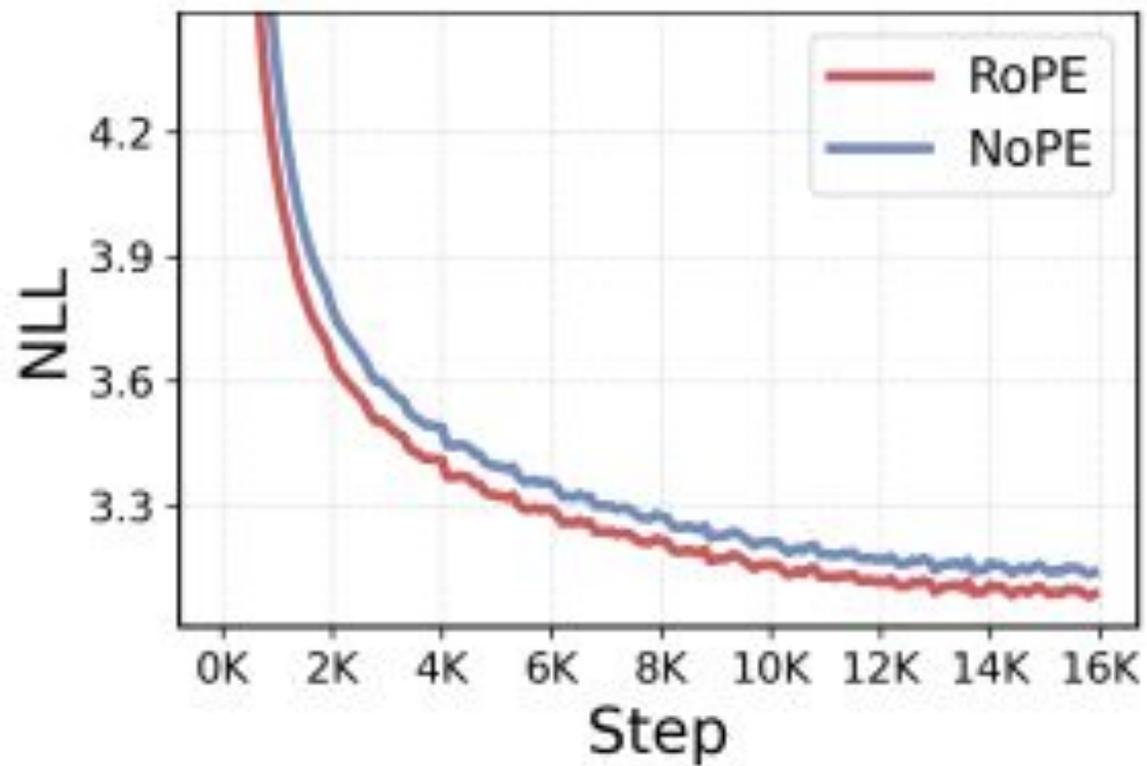


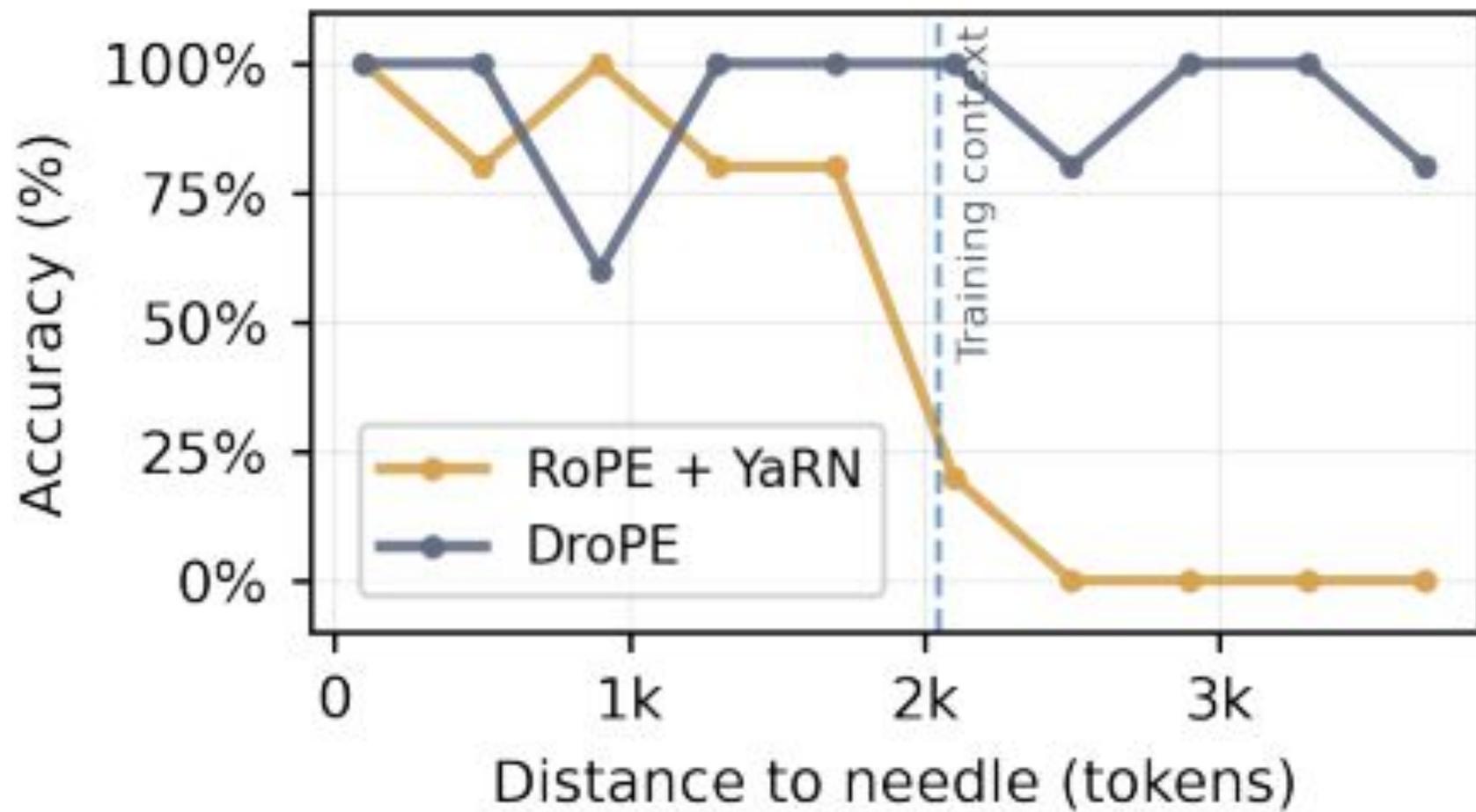




Positional Encoding
 ● NoPE ▲ T5's Relative PE ■ ALiBi ◆ Rotary ▼ Absolute Position Embedding

<https://arxiv.org/pdf/2305.19466>





Concluding Remarks

Absolute Positional Embedding

Relative Positional Embedding

RoPE

Train short, test long

No Positional Embedding?!