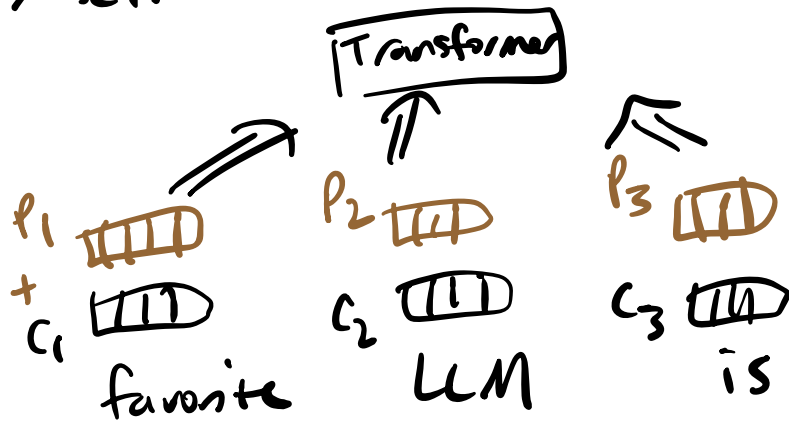


# Positional encodings

↳ self-attn is not inherently position-aware



↳ absolute position info

↳ additive embeddings

↳ new params assoc. w/ position

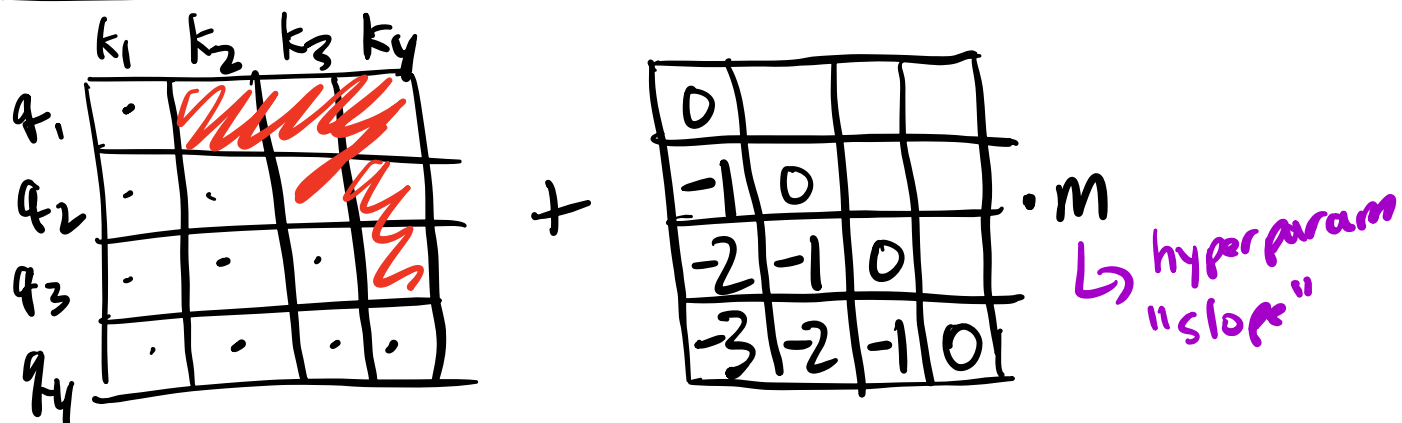
↳ lacks position extrapolation

↳ cannot generalize to sequences longer than max. training length

## relative position



ALiBi: no position embs added to input  $L$ ;



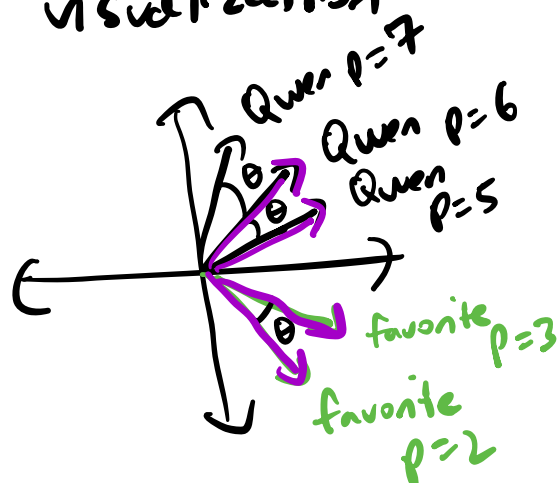
↳ each head can have a diff  $m$

↳ no extra params to train

↳ much better extrapolation to  
Seqs longer than max training length

RoPE: rotary position encoding

- 2d visualization



dot product  
between  $Qwen_5$  and  
 $favorite_2$   
is same as  
 $Qwen_6 \cdot favorite_3$

how do we rotate a vector?

↳ multiply by a rotation matrix

$$W_{R_{\theta,p}} = \begin{bmatrix} \cos(p\theta) & -\sin(p\theta) \\ \sin(p\theta) & \cos(p\theta) \end{bmatrix}$$

↑                      ↑  
rotation      position  
freq            in  
                  seq

how do we integrate rotation into self-atten?

RoPE:

$$q_{\text{Queen}} = W_{R_{\theta,p=5}} \cdot W_q \cdot c_{\text{Queen}}$$

$$k_{\text{favorite}} = W_{R_{\theta,p=2}} \cdot W_k \cdot c_{\text{favorite}}$$

$$q_{\text{Queen}} \cdot k_{\text{favorite}}$$

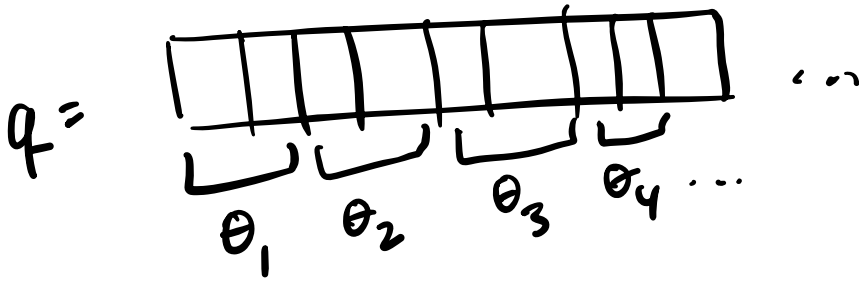
↳ due to rotation matrix properties

$$= (W_q c_{\text{Queen}})^T \cdot W_{R_{\theta,5-2}} \cdot (W_k c_{\text{favorite}})$$

↳  $\underset{=3}{\phantom{5-2}}$

↑ this depends only  
on the relative pos

↳ same dot product regardless of abs. pos,  
as long as  $p_{\text{Queen}} - p_{\text{favorite}}$  is the same



↳ each  $\theta_i$  is a constant  
that controls rotation freq