# **Positional Embeddings and Relative Attention**



- Consider a phrase like:

The Space Needle is located in \_\_\_\_\_



- Consider the sentence:

The Space Needle is located in \_\_\_\_\_

- Now consider phrases:
  - The famous landmark known as the Space Needle is located in \_\_\_\_\_
  - I think the Space Needle is located in \_\_\_\_\_
  - The CN tower is located in Toronto and the Space Needle is located in \_\_\_\_\_







- Consider the sentence:

The Space Needle is located in \_\_\_\_\_



- Now consider phrases:

The famous landmark known as the Space Needle is located in \_

I think the Space Needle is located in \_\_\_\_\_

The CN tower is located in Toronto and the Space Needle is located in \_\_\_\_\_

If  $h_{-3}$  to  $h_{-1}$  = "is located in", we should pay attention to h.







- Paying attention to all the text is necessary
- However, in many cases, the important tokens are the same distance away.
- With relative attention, we can reinforce recurring relationships even with offsets!



- Paying attention to all the text is necessary
- However, in many cases, the important tokens are the same distance away.
- With relative attention, we can reinforce recurring relationships!
- of course, absolute attention *allows* the formation of relative attention.
- But positional encodings help extra...

for position encodings  $p_t$  at position t, and following position encoding  $p_{t+\Phi}$ 

 $p_{t} = \begin{bmatrix} \sin(\omega_{1}.t) \\ \cos(\omega_{1}.t) \\ \sin(\omega_{2}.t) \\ \cos(\omega_{2}.t) \\ \vdots \\ \sin(\omega_{d/2}.t) \\ \cos(\omega_{d/2}.t) \end{bmatrix}$ 

for position encodings  $p_t$  at position t, and following position encoding  $p_{t+\phi}$ 



for position encodings  $p_t$  at position *t*, and following position encoding  $p_{t+\Phi}$ 

	we can express encoding $p_{t+\Phi}$ as linearly transformed
$\left[ egin{array}{c} \sin(\omega_1,t) \ \cos(\omega_1,t) \end{array}  ight]$	p <sub>t</sub>
$\sin(\omega_2.t)\ \cos(\omega_2.t)$	$M p_t = p_{t+\Phi}$
$\vdots$ $\sin(\omega_{d/2}, t)$	$M. \left[ egin{smmatrix} \sin(\omega_k,t) \ \cos(\omega_k,t) \end{matrix}  ight] = \left[ egin{smmatrix} \sin(\omega_k,(t+\phi)) \ \cos(\omega_k,(t+\phi)) \end{matrix}  ight]$
$\left\lfloor \cos(\omega_{d/2},t) ight floor$	

 $p_t =$ 





#### Remember: matrices are linear transformations!



$$\begin{pmatrix} a & b \\ c & d \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} ax + by \\ cx + dy \end{pmatrix}$$

When we multiply by a 2x2 matrix, we are changing our basis from  $\{(1,0), (0,1)\}$  to  $\{(a, c), (b, d)\}$ 

$$M. \left[ egin{smmatrix} \sin(\omega_k.\,t) \ \cos(\omega_k.\,t) \end{matrix} 
ight] = \left[ egin{smmatrix} \sin(\omega_k.\,(t+\phi)) \ \cos(\omega_k.\,(t+\phi)) \end{matrix} 
ight]$$

$$M. \left[ egin{array}{c} \sin(\omega_k.\,t) \ \cos(\omega_k.\,t) \end{array} 
ight] = \left[ egin{array}{c} \sin(\omega_k.\,(t+\phi)) \ \cos(\omega_k.\,(t+\phi)) \end{array} 
ight]$$

$$egin{bmatrix} u_1 & v_1 \ u_2 & v_2 \end{bmatrix} \cdot egin{bmatrix} \sin(\omega_k,t) \ \cos(\omega_k,t) \end{bmatrix} = egin{bmatrix} \sin(\omega_k,(t+\phi)) \ \cos(\omega_k,(t+\phi)) \end{bmatrix}$$

$$M. egin{bmatrix} \sin(\omega_k,t)\ \cos(\omega_k,t) \end{bmatrix} = egin{bmatrix} \sin(\omega_k,(t+\phi))\ \cos(\omega_k,(t+\phi)) \end{bmatrix} \ \begin{bmatrix} u_1 & v_1\ u_2 & v_2 \end{bmatrix} . egin{bmatrix} \sin(\omega_k,t)\ \cos(\omega_k,t) \end{bmatrix} = egin{bmatrix} \sin(\omega_k,(t+\phi))\ \cos(\omega_k,(t+\phi)) \end{bmatrix}$$

The really cool part: 
$$sin(x+y) = sin(x)cos(y) + cos(x)sin(y)$$

$$M. egin{bmatrix} \sin(\omega_k.\,t)\ \cos(\omega_k.\,t) \end{bmatrix} = egin{bmatrix} \sin(\omega_k.\,(t+\phi))\ \cos(\omega_k.\,(t+\phi)) \end{bmatrix} \ \begin{bmatrix} u_1 & v_1\ u_2 & v_2 \end{bmatrix} . egin{bmatrix} \sin(\omega_k.\,t)\ \cos(\omega_k.\,t) \end{bmatrix} = egin{bmatrix} \sin(\omega_k.\,(t+\phi))\ \cos(\omega_k.\,(t+\phi)) \end{bmatrix}$$

The really cool part: sin(x+y) = sin(x)cos(y) + cos(x)sin(y)

The right hand side produces  $u_1 sin(x) + v_1 cos(x)$ 

$$M. egin{bmatrix} \sin(\omega_k.\,t)\ \cos(\omega_k.\,t) \end{bmatrix} = egin{bmatrix} \sin(\omega_k.\,(t+\phi))\ \cos(\omega_k.\,(t+\phi)) \end{bmatrix} \ \begin{bmatrix} u_1 & v_1\ u_2 & v_2 \end{bmatrix} . egin{bmatrix} \sin(\omega_k.\,t)\ \cos(\omega_k.\,t) \end{bmatrix} = egin{bmatrix} \sin(\omega_k.\,(t+\phi))\ \cos(\omega_k.\,(t+\phi)) \end{bmatrix}$$

The really cool part: sin(x+y) = sin(x)cos(y) + cos(x)sin(y)

The right hand side produces  $u_1 sin(x) + v_1 cos(x)$ 

$$M. egin{bmatrix} \sin(\omega_k.\,t)\ \cos(\omega_k.\,t) \end{bmatrix} = egin{bmatrix} \sin(\omega_k.\,(t+\phi))\ \cos(\omega_k.\,(t+\phi)) \end{bmatrix} \ u_1 \quad v_1\ u_2 \quad v_2 \end{bmatrix}. egin{bmatrix} \sin(\omega_k.\,t)\ \cos(\omega_k.\,t) \end{bmatrix} = egin{bmatrix} \sin(\omega_k.\,(t+\phi))\ \cos(\omega_k.\,(t+\phi)) \end{bmatrix}$$

The really cool part: sin(x+y) = sin(x)cos(y) + cos(x)sin(y)The right hand side produces  $u_1 sin(x) + v_1 cos(x)$ 

So we're going to end up with only terms involving y...

$$M. egin{bmatrix} \sin(\omega_k,t)\ \cos(\omega_k,t) \end{bmatrix} = egin{bmatrix} \sin(\omega_k,(t+\phi))\ \cos(\omega_k,(t+\phi)) \end{bmatrix} \ u_1 \quad v_1\ u_2 \quad v_2 \end{bmatrix} . egin{bmatrix} \sin(\omega_k,t)\ \cos(\omega_k,t) \end{bmatrix} = egin{bmatrix} \sin(\omega_k,(t+\phi))\ \cos(\omega_k,(t+\phi)) \end{bmatrix}$$

The really cool part: sin(x+y) = sin(x)cos(y) + cos(x)sin(y)

The right hand side produces  $u_1 sin(x) + v_1 cos(x)$ 

So we're going to end up with only terms involving y...

Where x is  $w_k$  t, our only absolute position!

$$M. \left[ egin{array}{c} \sin(\omega_k,t) \ \cos(\omega_k,t) \end{array} 
ight] = \left[ egin{array}{c} \sin(\omega_k,(t+\phi)) \ \cos(\omega_k,(t+\phi)) \end{array} 
ight]$$

$$egin{bmatrix} u_1 & v_1 \ u_2 & v_2 \end{bmatrix} \cdot egin{bmatrix} \sin(\omega_k.\,t) \ \cos(\omega_k.\,t) \end{bmatrix} = egin{bmatrix} \sin(\omega_k.\,(t+\phi)) \ \cos(\omega_k.\,(t+\phi)) \end{bmatrix}$$

$$egin{bmatrix} u_1 & v_1 \ u_2 & v_2 \end{bmatrix} \cdot egin{bmatrix} \sin(\omega_k.\,t) \ \cos(\omega_k.\,t) \end{bmatrix} = egin{bmatrix} \sin(\omega_k.\,\phi) + \cos(\omega_k.\,t)\sin(\omega_k.\,\phi) \ \cos(\omega_k.\,t)\cos(\omega_k.\,\phi) - \sin(\omega_k.\,t)\sin(\omega_k.\,\phi) \end{bmatrix}$$

$$M. \begin{bmatrix} \sin(\omega_k, t) \\ \cos(\omega_k, t) \end{bmatrix} = \begin{bmatrix} \sin(\omega_k, (t + \phi)) \\ \cos(\omega_k, (t + \phi)) \end{bmatrix}$$
$$\begin{bmatrix} u_1 & v_1 \\ u_2 & v_2 \end{bmatrix} \cdot \begin{bmatrix} \sin(\omega_k, t) \\ \cos(\omega_k, t) \end{bmatrix} = \begin{bmatrix} \sin(\omega_k, (t + \phi)) \\ \cos(\omega_k, (t + \phi)) \end{bmatrix}$$
$$\begin{bmatrix} u_1 & v_1 \\ u_2 & v_2 \end{bmatrix} \cdot \begin{bmatrix} \sin(\omega_k, t) \\ \cos(\omega_k, t) \end{bmatrix} = \begin{bmatrix} \sin(\omega_k, t) \cos(\omega_k, \phi) + \cos(\omega_k, t) \sin(\omega_k, \phi) \\ \cos(\omega_k, t) \cos(\omega_k, \phi) - \sin(\omega_k, t) \sin(\omega_k, \phi) \end{bmatrix}$$

$$M. \begin{bmatrix} \sin(\omega_k, t) \\ \cos(\omega_k, t) \end{bmatrix} = \begin{bmatrix} \sin(\omega_k, (t + \phi)) \\ \cos(\omega_k, (t + \phi)) \end{bmatrix}$$
$$\begin{bmatrix} u_1 & v_1 \\ u_2 & v_2 \end{bmatrix} \cdot \begin{bmatrix} \sin(\omega_k, t) \\ \cos(\omega_k, t) \end{bmatrix} = \begin{bmatrix} \sin(\omega_k, (t + \phi)) \\ \cos(\omega_k, (t + \phi)) \end{bmatrix}$$
$$\begin{bmatrix} \overbrace{u_1} & v_1 \\ u_2 & v_2 \end{bmatrix} \cdot \begin{bmatrix} \sin(\omega_k, t) \\ \cos(\omega_k, t) \end{bmatrix} = \begin{bmatrix} \sin(\omega_k, t) \cos(\omega_k, \phi) + \cos(\omega_k, t) \sin(\omega_k, \phi) \\ \cos(\omega_k, t) \cos(\omega_k, \phi) - \sin(\omega_k, t) \sin(\omega_k, \phi) \end{bmatrix}$$

$$egin{aligned} &u_1\sin(\omega_k.\,t)+v_1\cos(\omega_k.\,t)\ &u_2\sin(\omega_k.\,t)+v_2\cos(\omega_k.\,t)\ \end{aligned}$$

$$M. \begin{bmatrix} \sin(\omega_k, t) \\ \cos(\omega_k, t) \end{bmatrix} = \begin{bmatrix} \sin(\omega_k, (t + \phi)) \\ \cos(\omega_k, (t + \phi)) \end{bmatrix}$$
 $\begin{bmatrix} u_1 & v_1 \\ u_2 & v_2 \end{bmatrix} \cdot \begin{bmatrix} \sin(\omega_k, t) \\ \cos(\omega_k, t) \end{bmatrix} = \begin{bmatrix} \sin(\omega_k, (t + \phi)) \\ \cos(\omega_k, (t + \phi)) \end{bmatrix}$ 
 $\begin{bmatrix} \overline{u_1} & v_1 \\ u_2 & v_2 \end{bmatrix} \cdot \begin{bmatrix} \sin(\omega_k, t) \\ \cos(\omega_k, t) \end{bmatrix} = \begin{bmatrix} \sin(\omega_k, t) \cos(\omega_k, \phi) + \cos(\omega_k, t) \sin(\omega_k, \phi) \\ \cos(\omega_k, t) \cos(\omega_k, \phi) - \sin(\omega_k, t) \sin(\omega_k, \phi) \end{bmatrix}$ 

- $u_1 \sin(\omega_k, t) + v_1 \cos(\omega_k, t) = \cos(\omega_k, \phi) \sin(\omega_k, t) + \sin(\omega_k, \phi) \cos(\omega_k, t)$ (1)
- $u_2\sin(\omega_k,t)+v_2\cos(\omega_k,t)=-\sin(\omega_k,\phi)\sin(\omega_k,t)+\cos(\omega_k,\phi)\cos(\omega_k,t)$  (2)

$$\mathbf{h} \qquad M. \begin{bmatrix} \sin(\omega_k, t) \\ \cos(\omega_k, t) \end{bmatrix} = \begin{bmatrix} \sin(\omega_k, (t+\phi)) \\ \cos(\omega_k, (t+\phi)) \end{bmatrix}$$

$$\begin{bmatrix} u_1 & v_1 \\ u_2 & v_2 \end{bmatrix} \cdot \begin{bmatrix} \sin(\omega_k, t) \\ \cos(\omega_k, t) \end{bmatrix} = \begin{bmatrix} \sin(\omega_k, (t+\phi)) \\ \cos(\omega_k, (t+\phi)) \end{bmatrix}$$

$$\begin{bmatrix} u_1 & v_1 \\ u_2 & v_2 \end{bmatrix} \cdot \begin{bmatrix} \sin(\omega_k, t) \\ \cos(\omega_k, t) \end{bmatrix} = \begin{bmatrix} \sin(\omega_k, t) \cos(\omega_k, \phi) + \cos(\omega_k, t) \sin(\omega_k, \phi) \\ \cos(\omega_k, t) \cos(\omega_k, \phi) - \sin(\omega_k, t) \sin(\omega_k, \phi) \end{bmatrix}$$

$$\begin{bmatrix} u_1 \sin(\omega_k, t) + v_1 \cos(\omega_k, t) = \cos(\omega_k, \phi) \sin(\omega_k, t) + \sin(\omega_k, \phi) \cos(\omega_k, t) \sin(\omega_k, \phi) \end{bmatrix}$$

$$\begin{bmatrix} u_1 \sin(\omega_k, t) + v_1 \cos(\omega_k, t) = \cos(\omega_k, \phi) \sin(\omega_k, t) + \sin(\omega_k, \phi) \cos(\omega_k, t) \sin(\omega_k, \phi) \end{bmatrix}$$

\_\_\_\_

$$M. \begin{bmatrix} \sin(\omega_k, t) \\ \cos(\omega_k, t) \end{bmatrix} = \begin{bmatrix} \sin(\omega_k, (t + \phi)) \\ \cos(\omega_k, (t + \phi)) \end{bmatrix}$$
$$\begin{bmatrix} u_1 & v_1 \\ u_2 & v_2 \end{bmatrix} \cdot \begin{bmatrix} \sin(\omega_k, t) \\ \cos(\omega_k, t) \end{bmatrix} = \begin{bmatrix} \sin(\omega_k, (t + \phi)) \\ \cos(\omega_k, (t + \phi)) \end{bmatrix}$$
$$\begin{bmatrix} u_1 & v_1 \\ u_2 & v_2 \end{bmatrix} \cdot \begin{bmatrix} \sin(\omega_k, t) \\ \cos(\omega_k, t) \end{bmatrix} = \begin{bmatrix} \sin(\omega_k, t) \cos(\omega_k, \phi) + \cos(\omega_k, t) \sin(\omega_k, \phi) \\ \cos(\omega_k, t) \cos(\omega_k, \phi) - \sin(\omega_k, t) \sin(\omega_k, \phi) \end{bmatrix}$$
$$\begin{bmatrix} u_1 \sin(\omega_k, t) + v_1 \cos(\omega_k, t) = \cos(\omega_k, \phi) \sin(\omega_k, t) + \sin(\omega_k, \phi) \cos(\omega_k, t) \quad (1) \\ u_2 \sin(\omega_k, t) + v_2 \cos(\omega_k, t) = -\sin(\omega_k, \phi) \sin(\omega_k, t) + \cos(\omega_k, \phi) \cos(\omega_k, t) \quad (2) \end{bmatrix}$$

$$egin{array}{ll} u_1 = & \cos(\omega_k,\phi) & v_1 = \sin(\omega_k,\phi) \ u_2 = -\sin(\omega_k,\phi) & v_2 = \cos(\omega_k,\phi) \end{array}$$

$$M. \begin{bmatrix} \sin(\omega_k, t) \\ \cos(\omega_k, t) \end{bmatrix} = \begin{bmatrix} \sin(\omega_k, (t+\phi)) \\ \cos(\omega_k, (t+\phi)) \end{bmatrix}$$
$$\begin{bmatrix} u_1 & v_1 \\ u_2 & v_2 \end{bmatrix} \cdot \begin{bmatrix} \sin(\omega_k, t) \\ \cos(\omega_k, t) \end{bmatrix} = \begin{bmatrix} \sin(\omega_k, (t+\phi)) \\ \cos(\omega_k, (t+\phi)) \end{bmatrix}$$
$$\begin{bmatrix} u_1 & v_1 \\ u_2 & v_2 \end{bmatrix} \cdot \begin{bmatrix} \sin(\omega_k, t) \\ \cos(\omega_k, t) \end{bmatrix} = \begin{bmatrix} \sin(\omega_k, t) \cos(\omega_k, \phi) + \cos(\omega_k, t) \sin(\omega_k, \phi) \\ \cos(\omega_k, t) \cos(\omega_k, \phi) - \sin(\omega_k, t) \sin(\omega_k, \phi) \end{bmatrix}$$
$$\begin{bmatrix} u_1 \sin(\omega_k, t) + v_1 \cos(\omega_k, t) = \cos(\omega_k, \phi) \sin(\omega_k, t) + \sin(\omega_k, \phi) \cos(\omega_k, t) \quad (1) \\ u_2 \sin(\omega_k, t) + v_2 \cos(\omega_k, t) = -\sin(\omega_k, \phi) \sin(\omega_k, t) + \cos(\omega_k, \phi) \cos(\omega_k, t) \quad (2) \end{bmatrix}$$

$$M_{\phi,k} = egin{bmatrix} \cos(\omega_k,\phi) & \sin(\omega_k,\phi) \ -\sin(\omega_k,\phi) & \cos(\omega_k,\phi) \end{bmatrix}$$

$$M_{\phi,k} = egin{bmatrix} \cos(\omega_k,\phi) & \sin(\omega_k,\phi) \ -\sin(\omega_k,\phi) & \cos(\omega_k,\phi) \end{bmatrix}$$

$$M_{\phi,k} = egin{bmatrix} \cos(\omega_k,\phi) & \sin(\omega_k,\phi) \ -\sin(\omega_k,\phi) & \cos(\omega_k,\phi) \end{bmatrix}$$

And one more interesting thing...

$$M_{\phi,k} = egin{bmatrix} \cos(\omega_k,\phi) & \sin(\omega_k,\phi) \ -\sin(\omega_k,\phi) & \cos(\omega_k,\phi) \end{bmatrix}$$

And one more interesting thing...



 $\boldsymbol{p}_{t\!+\!\Phi} \text{ is } \boldsymbol{p}_t \text{ rotated } \text{by an amount based on } \boldsymbol{\Phi}$ 

$$M_{\phi,k} = egin{bmatrix} \cos(\omega_k,\phi) & \sin(\omega_k,\phi) \ -\sin(\omega_k,\phi) & \cos(\omega_k,\phi) \end{bmatrix}$$

And one more interesting thing...



 $\boldsymbol{p}_{t\!+\!\Phi} \text{ is } \boldsymbol{p}_t \text{ rotated } \text{by an amount based on } \boldsymbol{\Phi}$ 

## Summary / final notes

- Positional embeddings are added to each token embedding

- relative attention: if  $p_{t-1}$ ,  $p_{t-2}$ ,  $p_{t-3}$  were "is located in", pay lots of attention to  $p_{t-4}$ 

- By using the positional embeddings of words, it becomes very easy for a transformer to represent relative attention...
- Because  $p_{t-1}$ ,  $p_{t-2}$ ,  $p_{t-3}$  are simple linear transforms of  $p_t$ !