# ROFORMER-ENHANCED TRANSFORMER WITH ROTARY POSITION EMBEDDING

**3**

## 2.3 Relative position embedding

The authors of Shaw et al. [2018] applied different settings of Equation (1) as following: — Query bal. pre  — Query

$$f_q(\boldsymbol{x}_m) := \boldsymbol{W}_q(\boldsymbol{x}_m + \tilde{\boldsymbol{p}}_m)$$
$$f_k(\boldsymbol{x}_n, n) := \boldsymbol{W}_k(\boldsymbol{x}_n + \tilde{\boldsymbol{p}}_r^k)$$
$$f_v(\boldsymbol{x}_n, n) := \boldsymbol{W}_v(\boldsymbol{x}_n + \tilde{\boldsymbol{p}}_r^v)$$

(5)

where $\tilde{\boldsymbol{p}}_r^k, \tilde{\boldsymbol{p}}_r^v \in \mathbb{R}^d$ are trainable relative position embeddings. Note that $r = \text{clip}(m-n, r_{min}, r_{max})$ represents the relative distance between position $m$ and $n$. They clipped the relative distance with the hypothesis that precise relative position information is not useful beyond a certain distance. Keeping the form of Equation (3), the authors Dai et al. [2019] have proposed to decompose $\boldsymbol{q}_m^\top \boldsymbol{k}_n$ of Equation (2) as

$f_k(q,k)$ $\quad \boldsymbol{q}_m^\top \boldsymbol{k}_n = \boldsymbol{x}_m^\top \boldsymbol{W}_q^\top \boldsymbol{W}_k \boldsymbol{x}_n + \boldsymbol{x}_m^\top \boldsymbol{W}_q^\top \boldsymbol{W}_k \boldsymbol{p}_n + \boldsymbol{p}_m^\top \boldsymbol{W}_q^\top \boldsymbol{W}_k \boldsymbol{x}_n + \boldsymbol{p}_m^\top \boldsymbol{W}_q^\top \boldsymbol{W}_k \boldsymbol{p}_n$, (6)

the key idea is to replace the absolute position embedding $\boldsymbol{p}_n$ with its sinusoid-encoded relative counterpart $\tilde{\boldsymbol{p}}_{m-n}$, while the absolute position $\boldsymbol{p}_m$ in the third and fourth term with two trainable vectors $\mathbf{u}$ and $\mathbf{v}$ independent of the query positions. Further, $\boldsymbol{W}_k$ is distinguished for the content-based and location-based key vectors $\boldsymbol{x}_n$ and $\boldsymbol{p}_n$, denoted as $\boldsymbol{W}_k$ and $\widetilde{\boldsymbol{W}}_k$, resulting in:

$$\boldsymbol{q}_m^\top \boldsymbol{k}_n = \boldsymbol{x}_m^\top \boldsymbol{W}_q^\top \boldsymbol{W}_k \boldsymbol{x}_n + \boldsymbol{x}_m^\top \boldsymbol{W}_q^\top \widetilde{\boldsymbol{W}}_k \tilde{\boldsymbol{p}}_{m-n} + \mathbf{u}^\top \boldsymbol{W}_q^\top \boldsymbol{W}_k \boldsymbol{x}_n + \mathbf{v}^\top \boldsymbol{W}_q^\top \widetilde{\boldsymbol{W}}_k \tilde{\boldsymbol{p}}_{m-n}$$

(7)

---

**4**

The authors of He et al. [2020] argued that the relative positions of two tokens could only be fully modeled using the middle two terms of Equation (6). As a consequence, the absolute position embeddings $\boldsymbol{p}_m$ and $\boldsymbol{p}_n$ were simply replaced with the relative position embeddings $\tilde{\boldsymbol{p}}_{m-n}$:

$$\boldsymbol{q}_m^\top \boldsymbol{k}_n = \boldsymbol{x}_m^\top \boldsymbol{W}_q^\top \boldsymbol{W}_k \boldsymbol{x}_n + \boldsymbol{x}_m^\top \boldsymbol{W}_q^\top \boldsymbol{W}_k \tilde{\boldsymbol{p}}_{m-n} + \tilde{\boldsymbol{p}}_{m-n}^\top \boldsymbol{W}_q^\top \boldsymbol{W}_k \boldsymbol{x}_n$$

(10)

---

**4**

$$\langle f_q(\boldsymbol{x}_m, m), f_k(\boldsymbol{x}_n, n) \rangle = g(\boldsymbol{x}_m, \boldsymbol{x}_n, m-n).$$

(11)

The ultimate goal is to find an equivalent encoding mechanism to solve the functions $f_q(\boldsymbol{x}_m, m)$ and $f_k(\boldsymbol{x}_n, n)$ to conform the aforementioned relation.

---

**3**

### RoFormer

representation.

$$a_{m,n} = \frac{\exp(\frac{\boldsymbol{q}_m^\top \boldsymbol{k}_n}{\sqrt{d}})}{\sum_{j=1}^N \exp(\frac{\boldsymbol{q}_m^\top \boldsymbol{k}_j}{\sqrt{d}})}$$   $Q_1 \sqrt{d}$,

$$\boldsymbol{o}_m = \sum_{n=1}^N a_{m,n} \boldsymbol{v}_n$$   =)? no2 c-se/ inp?.

(2)

---

**3**

$$\begin{cases} \boldsymbol{p}_{i,2t} & = \sin(k/10000^{2t/d}) \\ \boldsymbol{p}_{i,2t+1} & = \cos(k/10000^{2t/d}) \end{cases}$$

---

**4**

$$f_q(\boldsymbol{x}_m, m) = (\boldsymbol{W}_q \boldsymbol{x}_m) e^{im\theta}$$   $e^{\theta I}.$   $e^{i\cdot\theta}:$ $\cos\theta + i\sin\theta.$
$$f_k(\boldsymbol{x}_n, n) = (\boldsymbol{W}_k \boldsymbol{x}_n) e^{in\theta}$$   (12)
$$g(\boldsymbol{x}_m, \boldsymbol{x}_n, m-n) = \text{Re}[(\boldsymbol{W}_q \boldsymbol{x}_m)(\boldsymbol{W}_k \boldsymbol{x}_n)^* e^{i(m-n)\theta}]$$

where $\text{Re}[\cdot]$ is the real part of a complex number and $(\boldsymbol{W}_k \boldsymbol{x}_n)^*$ represents the conjugate complex number of $(\boldsymbol{W}_k \boldsymbol{x}_n)$. $\theta \in \mathbb{R}$ is a preset non-zero constant. We can further write $f_{\{q,k\}}$ in a multiplication matrix:

$$f_{\{q,k\}}(\boldsymbol{x}_m, m) = \begin{pmatrix} \cos m\theta & -\sin m\theta \\ \sin m\theta & \cos m\theta \end{pmatrix} \begin{pmatrix} W_{(q,k)}^{(11)} & W_{(q,k)}^{(12)} \\ W_{(q,k)}^{(21)} & W_{(q,k)}^{(22)} \end{pmatrix} \begin{pmatrix} x_m^{(1)} \\ x_m^{(2)} \end{pmatrix}$$   $g \Rightarrow$ matrix. h2D (xnk

(13)

---

**5**

$$f_{\{q,k\}}(\boldsymbol{x}_m, m) = \boldsymbol{R}_{\Theta,m}^d \boldsymbol{W}_{\{q,k\}} \boldsymbol{x}_m$$   (14)

where   $\begin{pmatrix} \cos\theta_1 & \sin\theta_1 \\ -\sin\theta_1 & \cos\theta_1 \end{pmatrix}$

$$\boldsymbol{R}_{\Theta,m}^d = \begin{pmatrix} \cos m\theta_1 & -\sin m\theta_1 & 0 & 0 & \cdots & 0 & 0 \\ \sin m\theta_1 & \cos m\theta_1 & 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \cos m\theta_2 & -\sin m\theta_2 & \cdots & 0 & 0 \\ 0 & 0 & \sin m\theta_2 & \cos m\theta_2 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & \cos m\theta_{d/2} & -\sin m\theta_{d/2} \\ 0 & 0 & 0 & 0 & \cdots & \sin m\theta_{d/2} & \cos m\theta_{d/2} \end{pmatrix}$$   (15)

$\cos(m\theta)_1 - \sin(m\theta)_2$,
$\sin(m\theta)_1, \cos(m\theta)_2$

is the rotary matrix with pre-defined parameters $\Theta = \{\theta_i = 10000^{-2(i-1)/d}, i \in [1, 2, ..., d/2]\}$. A graphic illustration of RoPE is shown in Figure (1). Applying our RoPE to self-attention in Equation (2), we obtain:

$$\boldsymbol{q}_m^\top \boldsymbol{k}_n = (\boldsymbol{R}_{\Theta,m}^d \boldsymbol{W}_q \boldsymbol{x}_m)^\top (\boldsymbol{R}_{\Theta,n}^d \boldsymbol{W}_k \boldsymbol{x}_n) = \boldsymbol{x}^\top \boldsymbol{W}_q \boldsymbol{R}_{\Theta,n-m}^d \boldsymbol{W}_k \boldsymbol{x}_n$$

(16)

---

**5**

tion.



Figure 1: Implementation of Rotary Position Embedding(RoPE).

---

**7**

$$\boldsymbol{R}_{\Theta,m}^d \boldsymbol{x} = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ \vdots \\ x_{d-1} \\ x_d \end{pmatrix} \otimes \begin{pmatrix} \cos m\theta_1 \\ \cos m\theta_1 \\ \cos m\theta_2 \\ \cos m\theta_2 \\ \vdots \\ \cos m\theta_{d/2} \\ \cos m\theta_{d/2} \end{pmatrix} + \begin{pmatrix} -x_2 \\ x_1 \\ -x_4 \\ x_3 \\ \vdots \\ -x_{d-1} \\ x_d \end{pmatrix} \otimes \begin{pmatrix} \sin m\theta_1 \\ \sin m\theta_1 \\ \sin m\theta_2 \\ \sin m\theta_2 \\ \vdots \\ \sin m\theta_{d/2} \\ \sin m\theta_{d/2} \end{pmatrix}$$