

# RoFormer: Enhanced Transformer with Rotary Position Embedding

Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo  
Wen, Yunfeng Liu

발표자: 나 리 나

## 기존의 Positional Embeddings

- Absolute Positional Embeddings
  - 입력 시퀀스의 각 토큰에 고유한 위치 식별자를 제공

# 기존의 Positional Embeddings

- Relative Positional Embeddings

- 토큰 간의 상대적 거리에 대한 정보를 인코딩



$$A_{i,j}^{abs} = \underbrace{E_{x_i}^T W_q^T W_k E_{x_j}}_{\text{Content-to-Content (a)}} + \underbrace{E_{x_i}^T W_q^T W_k U_j}_{\text{Position-to-Content (b)}} + \underbrace{U_i^T W_q^T W_k E_{x_j}}_{\text{Content-to-Position (c)}} + \underbrace{U_i^T W_q^T W_k U_j}_{\text{Position-to-Position (d)}}$$

$$A_{i,j}^{abs} = \left( W_q (E_{x_i} + U_i) \right)^T \left( W_k (E_{x_j} + U_j) \right)$$

## Rotary Positional Embeddings(RoPE)

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

- RoPE는 각 위치에 대한 회전 행렬을 사용하여 절대 위치를 인코딩하고, 상대적 위치 의존성을 명시적으로 셀프-어텐션 공식에 통합
- 토큰의 위치에 따라 인베딩이 회전함으로써 위치 정보를 효과적으로 모델링

## Rotary Positional Embeddings(RoPE): A 2D Case

$$f_q(\mathbf{x}_m, m) = (\mathbf{W}_q \mathbf{x}_m) e^{im\theta}$$

$$f_k(\mathbf{x}_n, n) = (\mathbf{W}_k \mathbf{x}_n) e^{in\theta}$$

$$g(\mathbf{x}_m, \mathbf{x}_n, m - n) = \text{Re}[(\mathbf{W}_q \mathbf{x}_m)(\mathbf{W}_k \mathbf{x}_n)^* e^{i(m-n)\theta}]$$

$$f_{\{q,k\}}(\mathbf{x}_m, m) = \begin{pmatrix} \cos m\theta & -\sin m\theta \\ \sin m\theta & \cos m\theta \end{pmatrix} \begin{pmatrix} W_{\{q,k\}}^{(11)} & W_{\{q,k\}}^{(12)} \\ W_{\{q,k\}}^{(21)} & W_{\{q,k\}}^{(22)} \end{pmatrix} \begin{pmatrix} x_m^{(1)} \\ x_m^{(2)} \end{pmatrix}$$

- 벡터가 2차원 평면에 있을 때의 처리 방식을 설명

## Rotary Positional Embeddings(RoPE): General form

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

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

- General form은 2차원을 넘어서 더 높은 차원에서의 위치 정보를 처리할 수 있도록 확장

## Rotary Positional Embeddings(RoPE)

- 이전 기법과의 비교
  - Additive vs Multiplicative
  - sinusoid individually vs mix pairs of coordinates

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

# Rotary Positional Embeddings(RoPE)

- 장점

- Long-term decay

- 상대 위치가 멀어질수록 내적 값이 감소하는 효과

- RoPE with linear attention

- 시퀀스 길이 차원을 감소시키지 않고도 확장성을 유지

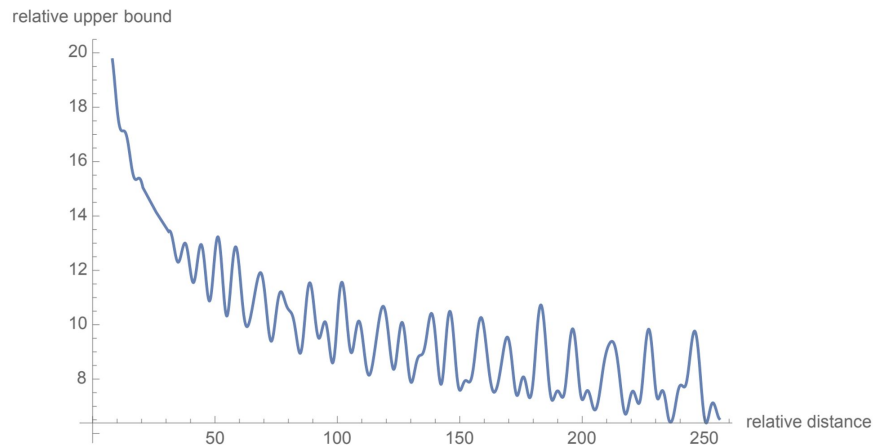


Figure 2: Long-term decay of RoPE.