# Rotary Position Embeddings

*arxiv.org/abs/2104.09864*

Васильев Руслан

8 декабря 2021

# Positional Encoding

...

Encoder

"token **x** on position **k**" →

Input is sum of two embeddings: for token and position →

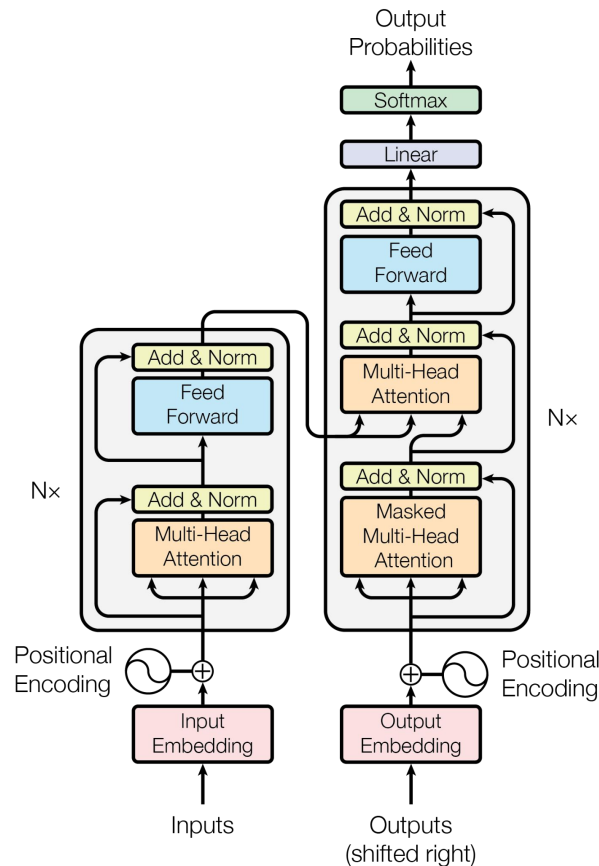tokens → Я "I"    видел "saw"    котю "cat"    <eos>

positions → 0    1    2    3

# Positional Encoding
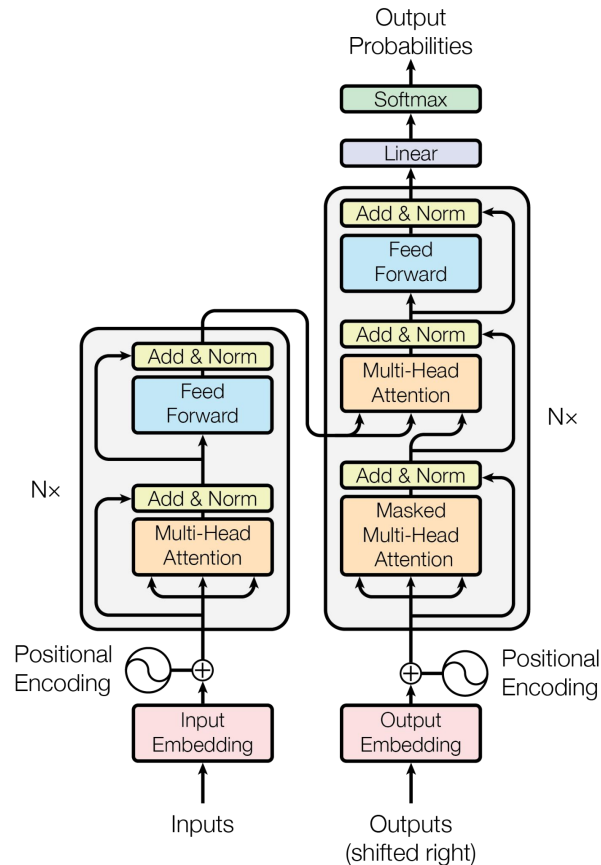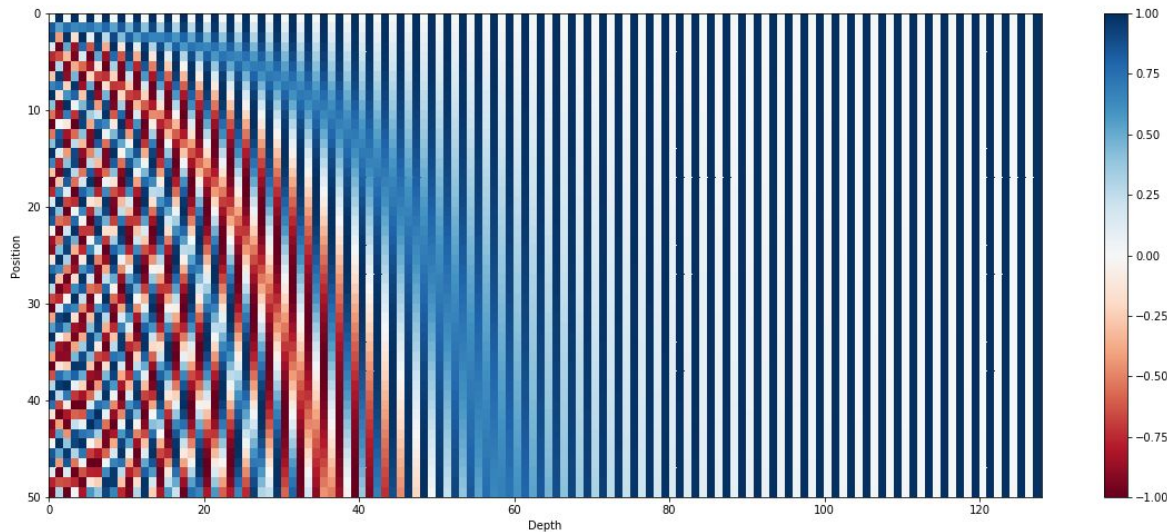
$$\vec{p_t}^{(i)} = f(t)^{(i)} := \begin{cases} \sin(\omega_k.t), & \text{if } i = 2k \\ \cos(\omega_k.t), & \text{if } i = 2k + 1 \end{cases}$$

$$\omega_k = \frac{1}{10000^{2k/d}}$$

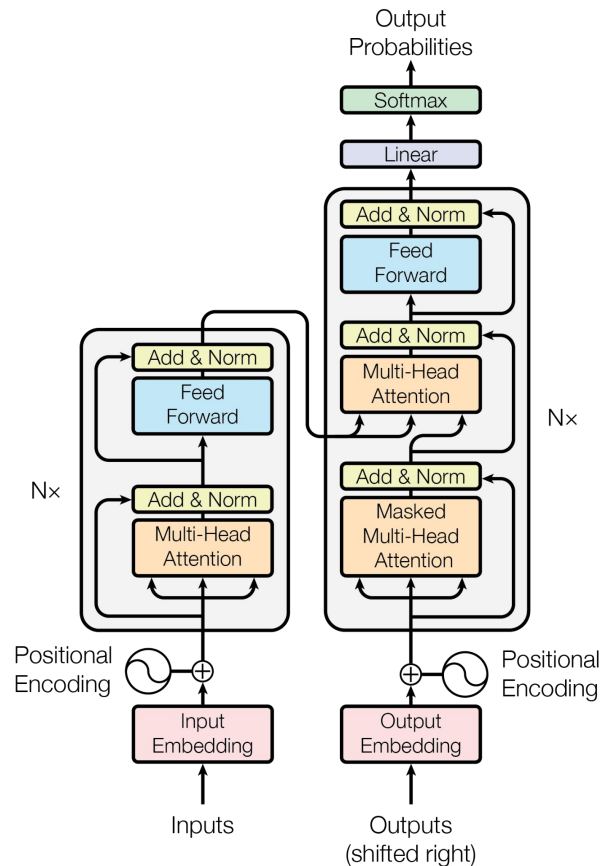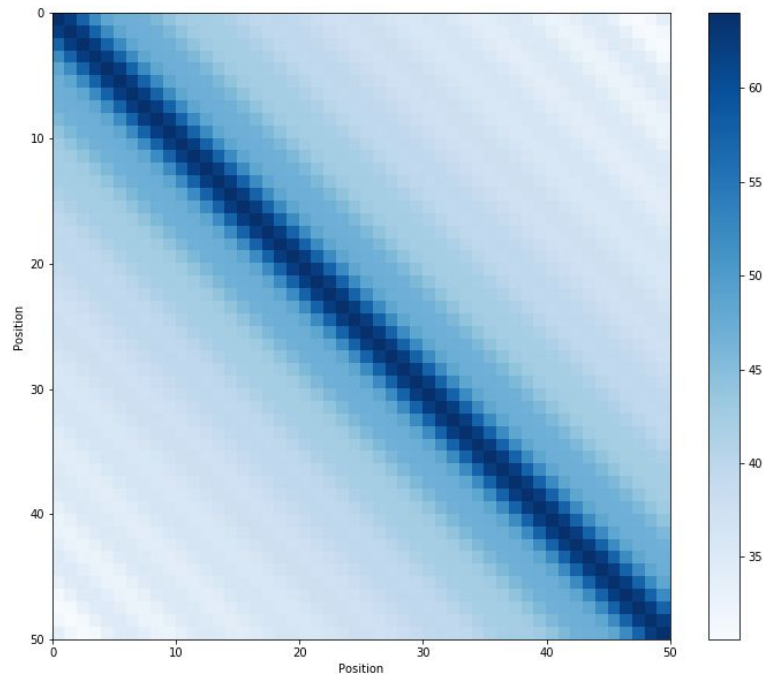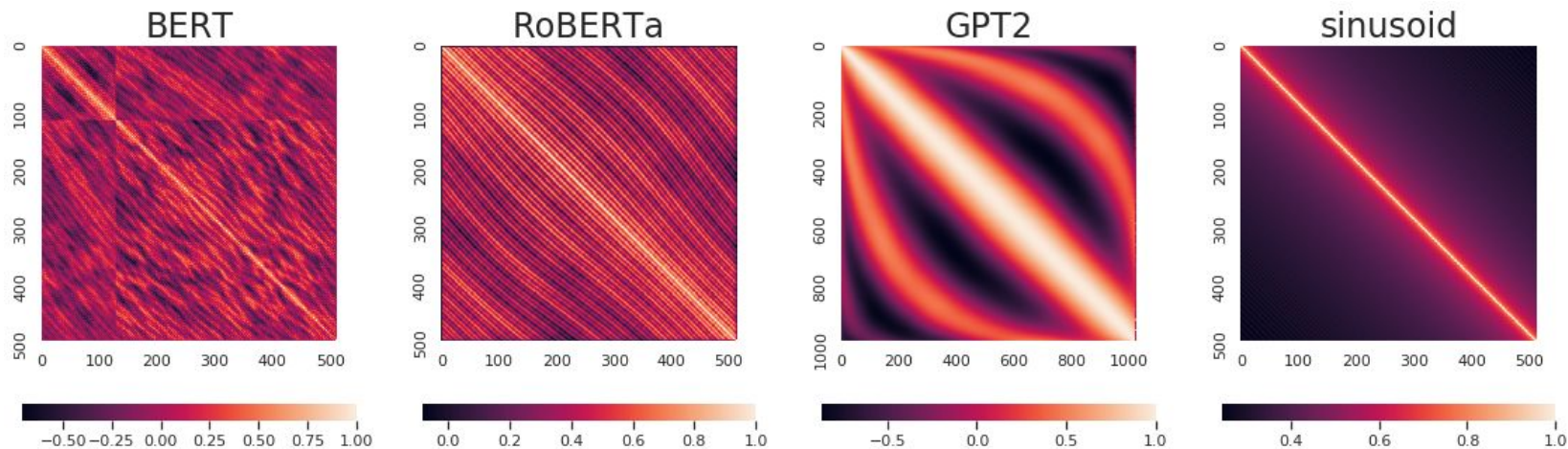$$\psi'(w_t) = \psi(w_t) + \vec{p_t}$$
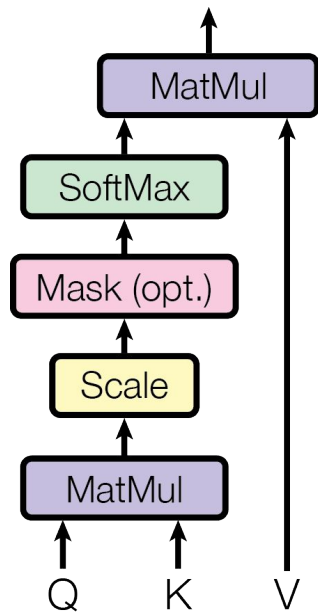


arxiv.org/abs/1706.03762
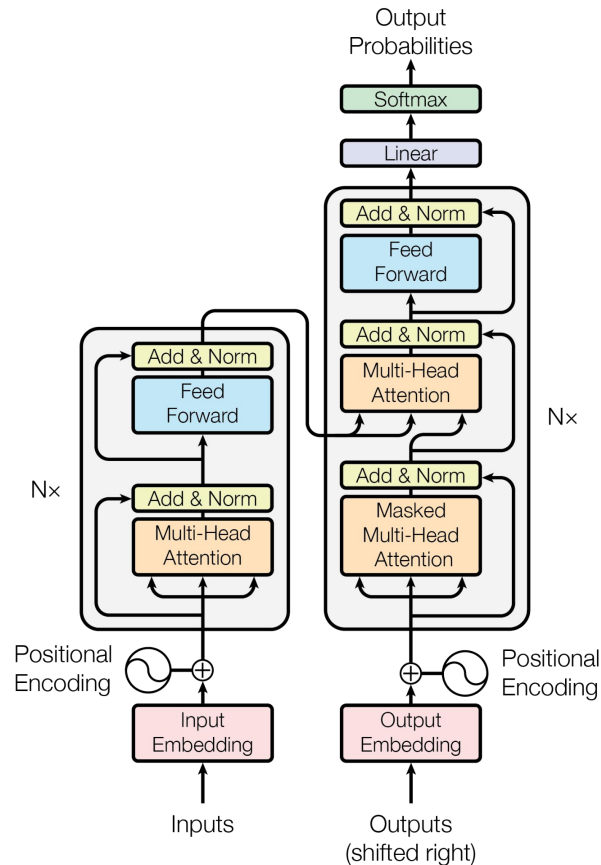
# Positional Encoding

# Positional Encoding

# Trainable Position Embedding

# Attention



$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

# Multihead Attention

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, ..., \text{head}_h)W^O$$

$$\text{where head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

# Attention

$$\boldsymbol{q}_m = f_q(\boldsymbol{x}_m, m)$$
$$\boldsymbol{k}_n = f_k(\boldsymbol{x}_n, n)$$
$$\boldsymbol{v}_n = f_v(\boldsymbol{x}_n, n)$$

$$\mathbf{o}_m = \sum_{n=1}^{N} a_{m,n} \boldsymbol{v}_n$$

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

# Relative Position Encoding

$$f_q(\boldsymbol{x}_m) := \boldsymbol{W}_q \boldsymbol{x}_m$$

$$f_k(\boldsymbol{x}_n, n) := \boldsymbol{W}_k(\boldsymbol{x}_n + \tilde{\boldsymbol{p}}_r^k)$$

$$f_v(\boldsymbol{x}_n, n) := \boldsymbol{W}_v(\boldsymbol{x}_n + \tilde{\boldsymbol{p}}_r^v)$$

$$r = \mathrm{clip}(m - n, r_{\min}, r_{\max})$$

| Model | Position Information | EN-DE BLEU | EN-FR BLEU |
|---|---|---|---|
| Transformer (base) | Absolute Position Representations | 26.5 | 38.2 |
| Transformer (base) | Relative Position Representations | **26.8** | **38.7** |
| Transformer (big) | Absolute Position Representations | 27.9 | 41.2 |
| Transformer (big) | Relative Position Representations | **29.2** | **41.5** |

arxiv.org/abs/1803.02155

# Formulation

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

# Solution (2D)

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

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

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

# General Solution

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

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

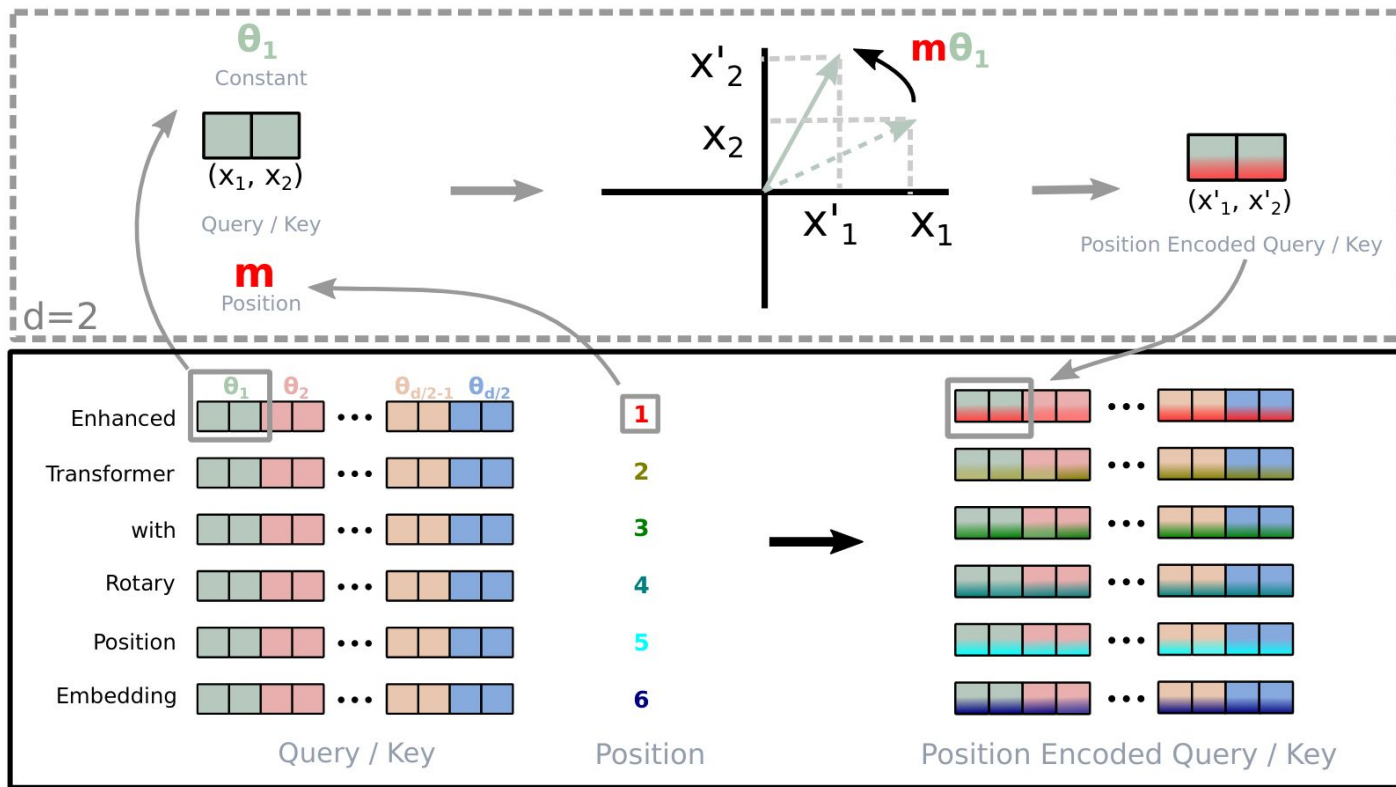$$\Theta = \{\theta_i = 10000^{-2(i-1)/d}, i \in [1, 2, ..., d/2]\}$$

# General Solution

$$f_{\{q,k\}}(\boldsymbol{x}_m, m) = \boldsymbol{R}^d_{\Theta,m} \boldsymbol{W}_{\{q,k\}} \boldsymbol{x}_m \qquad \boldsymbol{R}^d_{\Theta,n-m} = (\boldsymbol{R}^d_{\Theta,m})^{\mathsf{T}} \boldsymbol{R}^d_{\Theta,n}$$

$$\boldsymbol{q}^{\mathsf{T}}_m \boldsymbol{k}_n = (\boldsymbol{R}^d_{\Theta,m} \boldsymbol{W}_q \boldsymbol{x}_m)^{\mathsf{T}} (\boldsymbol{R}^d_{\Theta,n} \boldsymbol{W}_k \boldsymbol{x}_n) = \boldsymbol{x}^{\mathsf{T}} \boldsymbol{W}_q R^d_{\Theta,n-m} \boldsymbol{W}_k \boldsymbol{x}_n$$
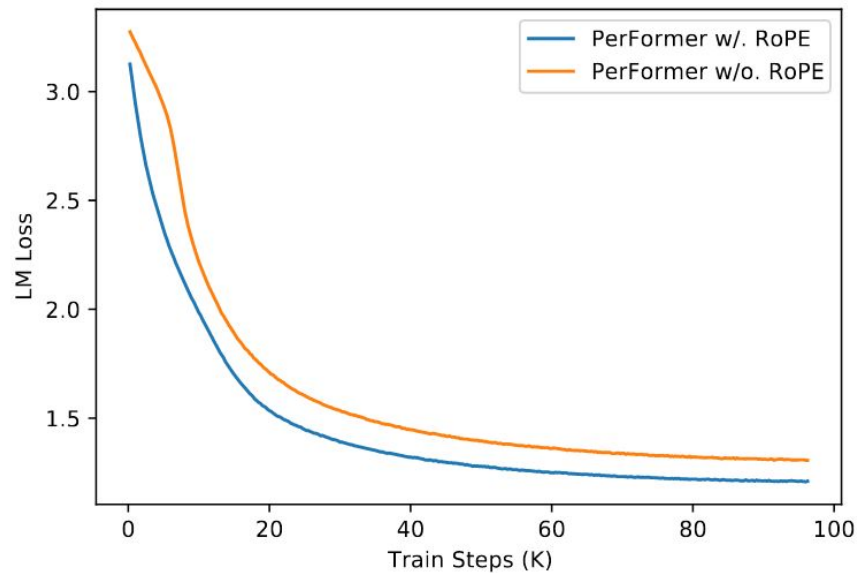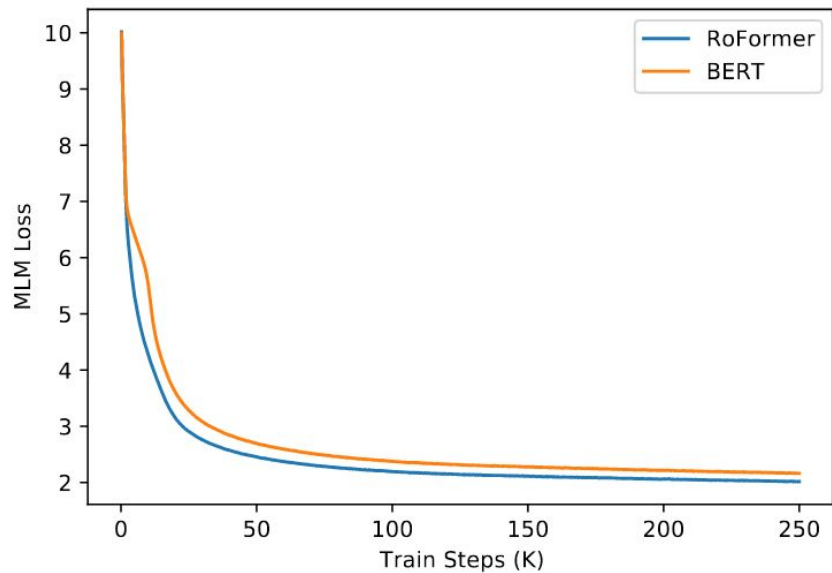
# Implementation

# Implementation

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

# Experiments: Pretraining
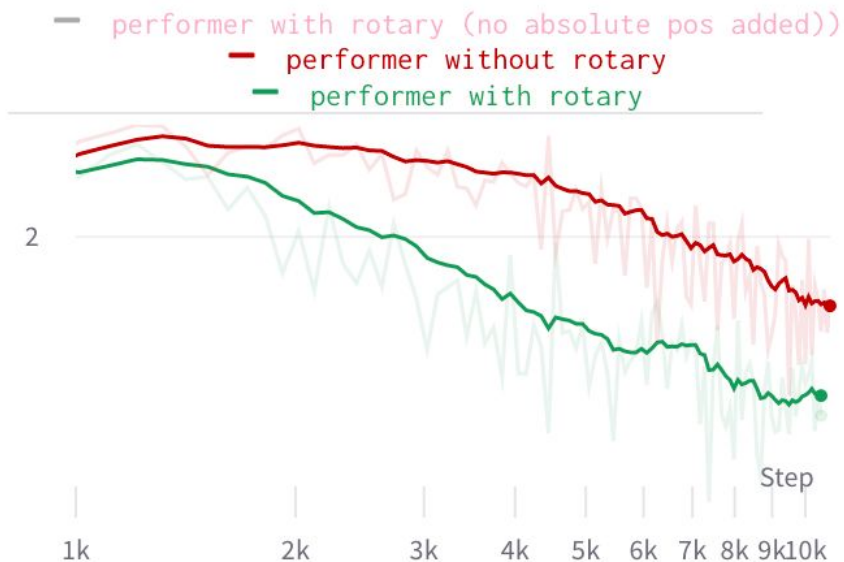
# Experiments: Downstream Tasks

| Model | MRPC | SST-2 | QNLI | STS-B | QQP | MNLI(m/mm) |
|-------|------|-------|------|-------|------|------------|
| BERT[8] | 88.9 | 93.5 | 90.5 | 85.8 | 71.2 | 84.6/83.4 |
| RoFormer | **89.5** | 90.7 | 88.0 | **87.0** | **86.4** | 80.2/79.8 |

| Model | validation | test |
|-------|------------|------|
| BERT-512 | 64.13% | 67.77% |
| WoBERT-512 | 64.07% | 68.10% |
| **RoFormer-512** | 64.13% | 68.29% |
| **RoFormer-1024** | **66.07**% | **69.79**% |

| Model | BLEU |
|-------|------|
| Transformer-base[37] | 27.3 |
| RoFormer | **27.5** |

# Experiments: Performer

# Experiments: LM



**validation lm loss value**

— rpe    — rotary    — learned

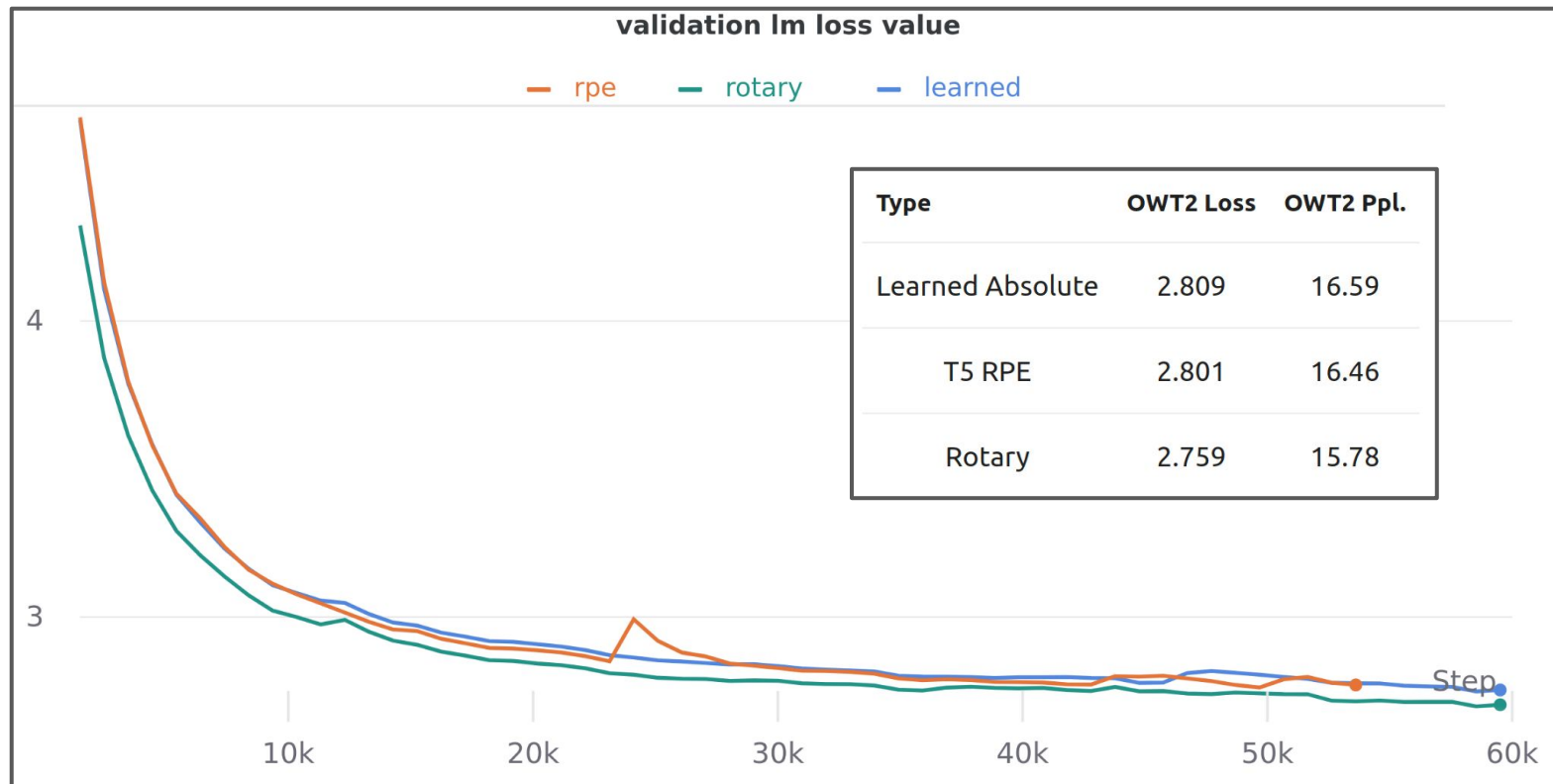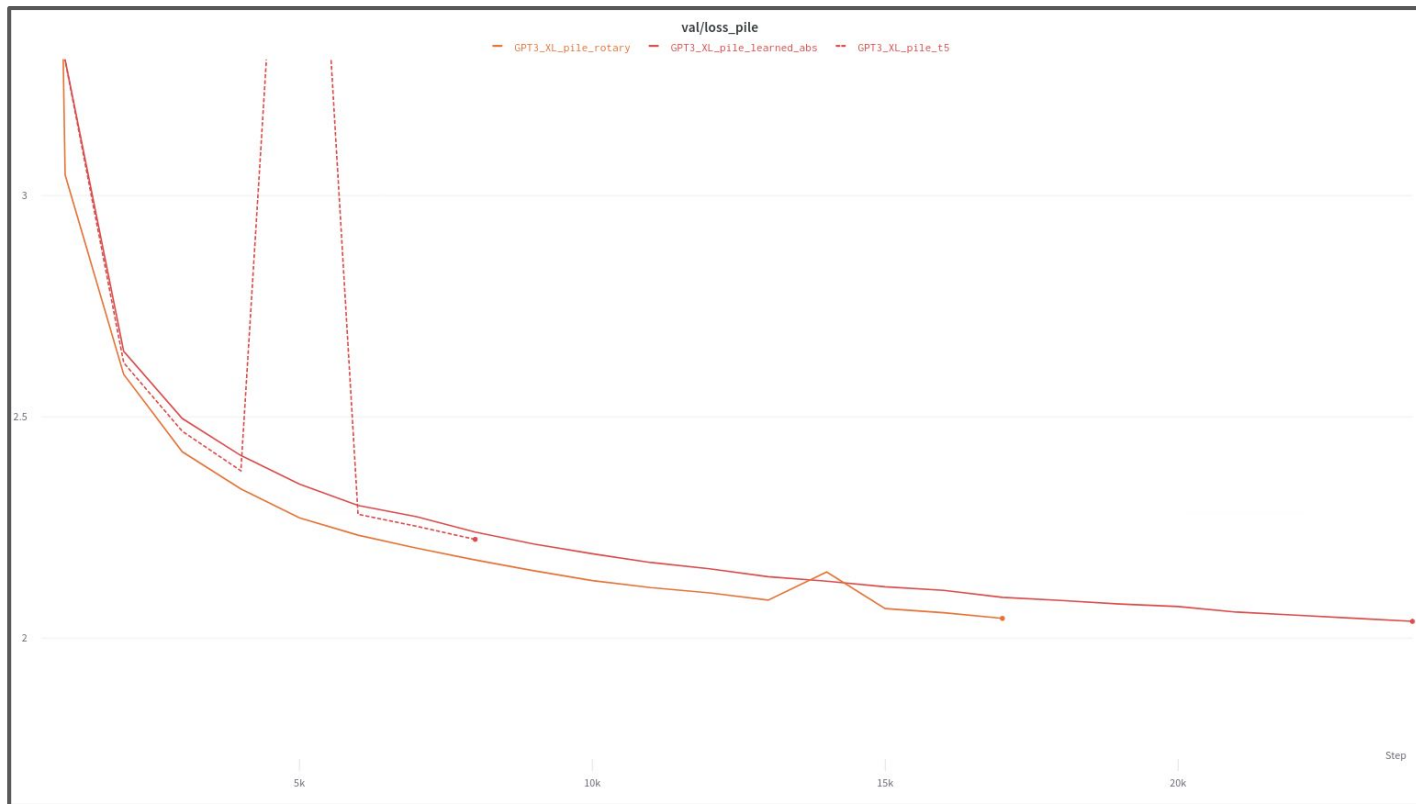| Type | OWT2 Loss | OWT2 Ppl. |
|---|---|---|
| Learned Absolute | 2.809 | 16.59 |
| T5 RPE | 2.801 | 16.46 |
| Rotary | 2.759 | 15.78 |

# Experiments: LM

# Conclusion: Rotary Embeddings

- Relative position in self-attention encoded through rotation matrix
- No training
- Faster convergence
- Greater stability