

# RoPE

Rotary Position Embedding

# Papers

Attention is all you need:

- <https://arxiv.org/pdf/1706.03762>

RoFormer: Enhanced transformer with Rotary Position Embedding

- <https://arxiv.org/pdf/2104.09864>

Position Interpolation:

- <https://arxiv.org/pdf/2306.15595>

Transformers are RNNs: Fast Autoregressive Transformers with Linear Attention

- <https://arxiv.org/pdf/2006.16236>

# Videos

Efficient NLP: <https://www.youtube.com/watch?v=o29P0Kpobz0>

Code\_your\_own\_ai: <https://www.youtube.com/watch?v=GQP0tyITy54>

DeepLearning Hero : <https://www.youtube.com/watch?v=DvP8f7eWS7U>

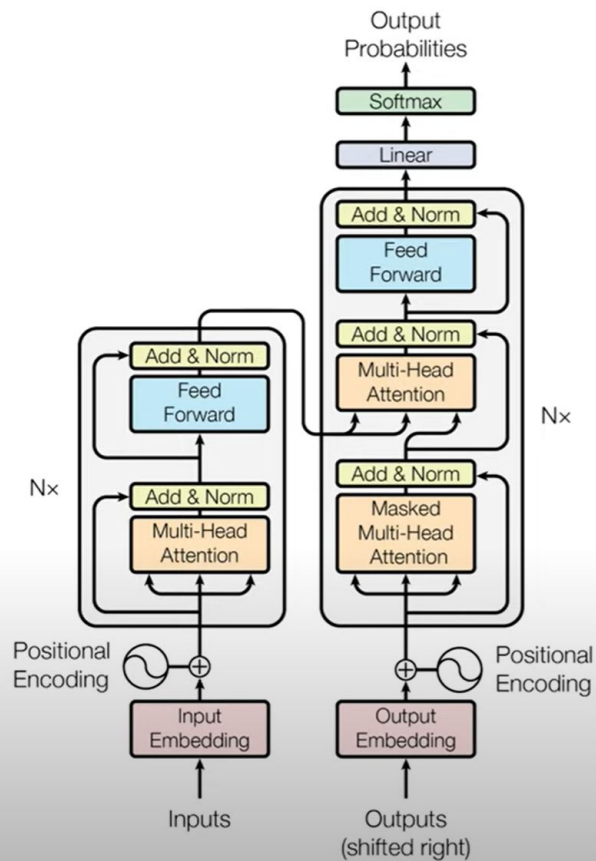
# Web resources

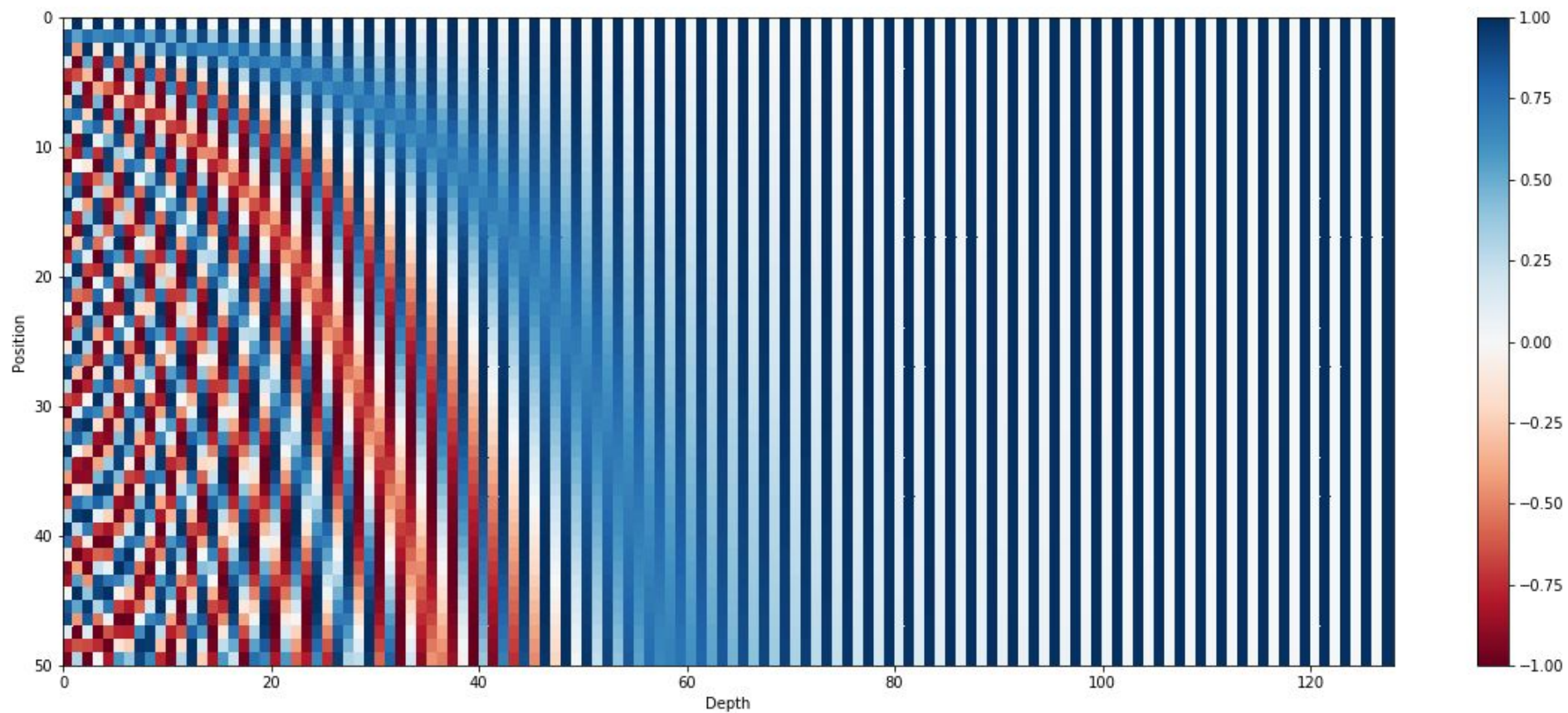
[https://uvadlc-notebooks.readthedocs.io/en/latest/tutorial\\_notebooks/tutorial6/Transformers\\_and\\_MHAttention.html](https://uvadlc-notebooks.readthedocs.io/en/latest/tutorial_notebooks/tutorial6/Transformers_and_MHAttention.html)

<https://towardsdatascience.com/attention-is-all-you-need-e498378552f9>

<https://blog.eleuther.ai/rotary-embeddings/>

# Transformer





$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{model}})$$

$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{model}})$$

The Current Position      Dimension Index      Dimension = 512

Figure 4: Positional Encodings Formula ([source](#))

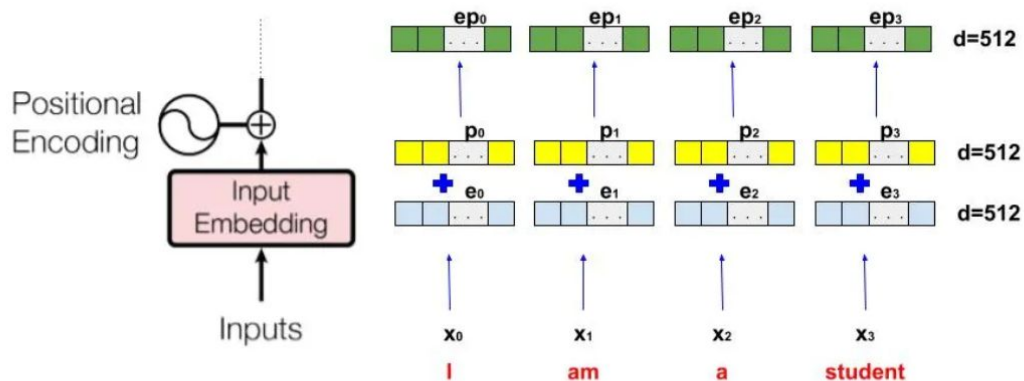


Figure 5: Adding Positional Encodings to the Embeddings to Generate Positional Embeddings (ep) (Image by Author)

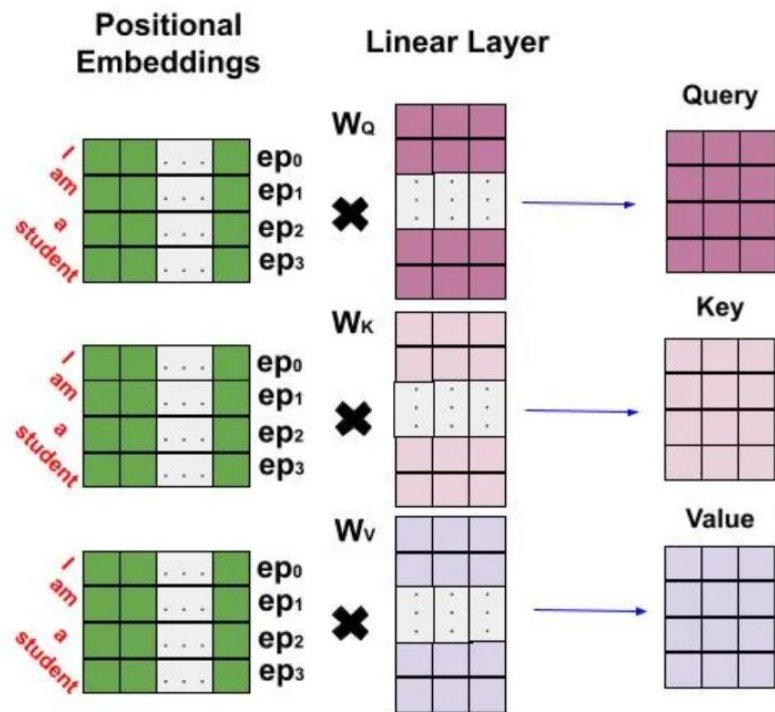
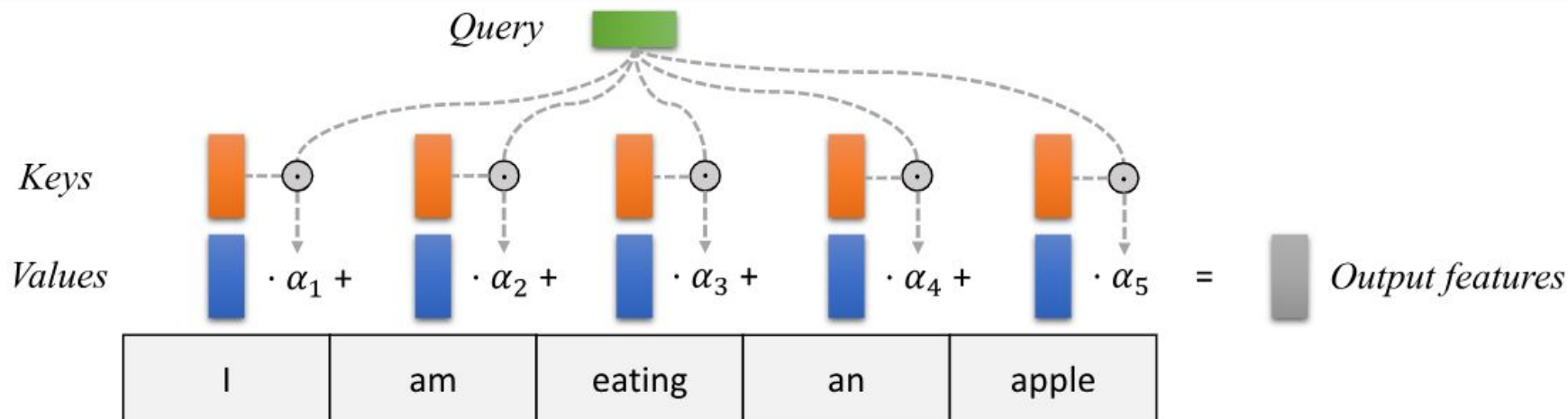


Figure 7: Generating the Query, Key and Value (Image by Author)



$$\alpha_i = \frac{\exp(f_{\text{attn}}(\text{key}_i, \text{query}))}{\sum_j \exp(f_{\text{attn}}(\text{key}_j, \text{query}))}, \quad \text{out} = \sum_i \alpha_i \cdot \text{value}_i$$

he attention over a sequence of words as follows:



# Attention is all you need.

## 3.5 Positional Encoding

Since our model contains no recurrence and no convolution, in order for the model to make use of the order of the sequence, we must inject some information about the relative or absolute position of the tokens in the sequence. To this end, we add "positional encodings" to the input embeddings at the bottoms of the encoder and decoder stacks. The positional encodings have the same dimension  $d_{\text{model}}$  as the embeddings, so that the two can be summed. There are many choices of positional encodings, learned and fixed [9].

In this work, we use sine and cosine functions of different frequencies:

$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$
$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$

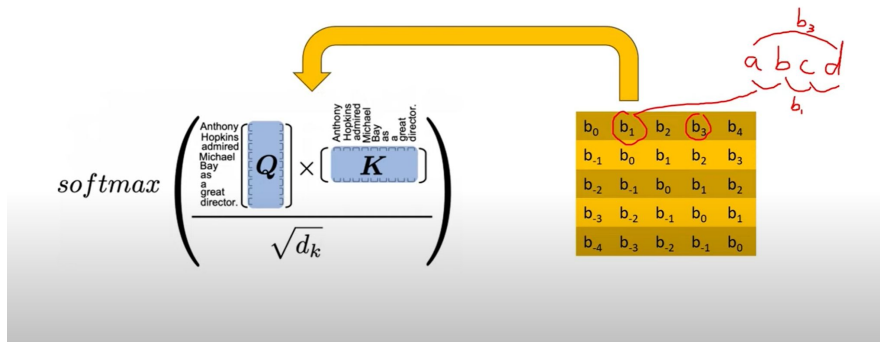
where  $pos$  is the position and  $i$  is the dimension. That is, each dimension of the positional encoding corresponds to a sinusoid. The wavelengths form a geometric progression from  $2\pi$  to  $10000 \cdot 2\pi$ . We chose this function because we hypothesized it would allow the model to easily learn to attend by relative positions, since for any fixed offset  $k$ ,  $PE_{pos+k}$  can be represented as a linear function of  $PE_{pos}$ .

We also experimented with using learned positional embeddings [9] instead, and found that the two versions produced nearly identical results (see Table 3 row (E)). We chose the sinusoidal version because it may allow the model to extrapolate to sequence lengths longer than the ones encountered during training.

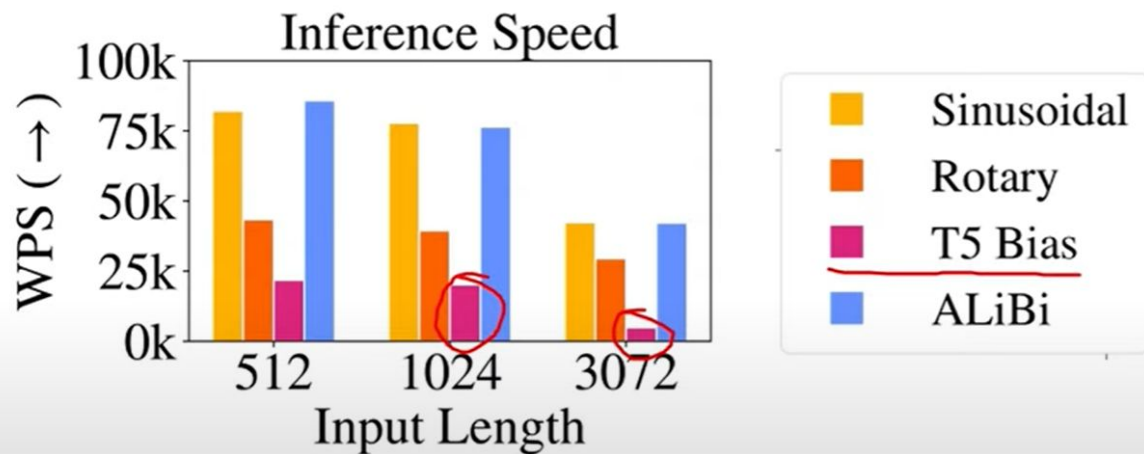
# Relative position embedding

## Efficient NLP

Relative Positional Embeddings: T5 Model



# Relative embeddings are slow!



**Why?**

- Extra step in self-attention layer

# Position embedding approaches

## 2.2 Absolute position embedding

A typical choice of Equation (1) is

$$f_{t:t \in \{a.k.v\}}(\mathbf{x}_i, i) := \mathbf{W}_{t:t \in \{a.k.v\}}(\mathbf{x}_i + \mathbf{p}_i),$$

## 2.3 Relative position embedding

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

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

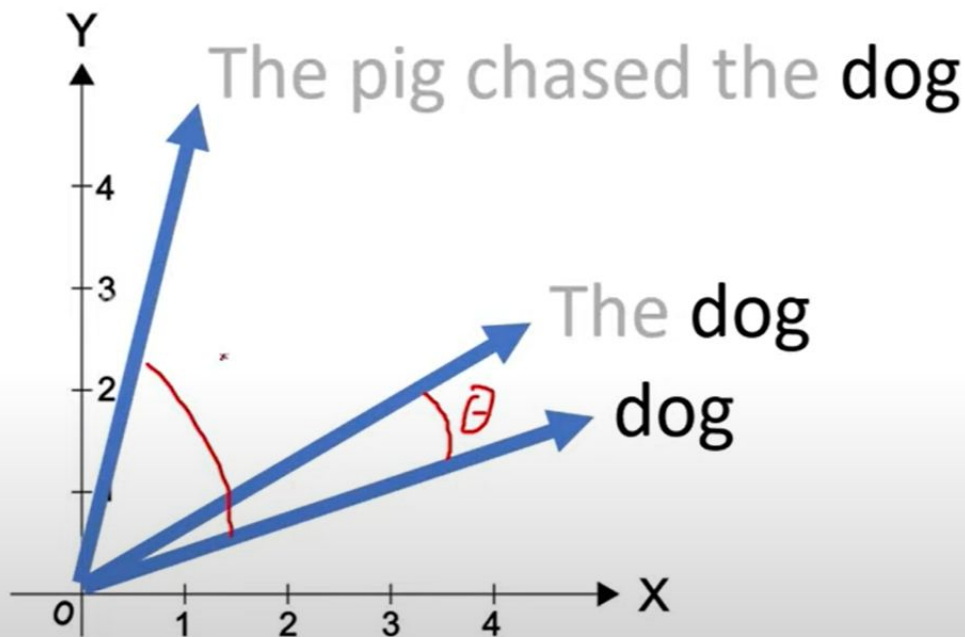
# Rotary Position Embedding

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

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

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

# Rotary Positional Embeddings



$$\mathbf{q}_m = f_q(\mathbf{x}_m, m)$$

$$\mathbf{k}_n = f_k(\mathbf{x}_n, n)$$

$$\mathbf{v}_n = f_v(\mathbf{x}_n, n),$$

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

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

