

Introduction to Transformer (Part II)

Large Language Models: Introduction and Recent Advances

ELL881 · AIL821



Tanmoy Chakraborty
Associate Professor, IIT Delhi
<https://tanmoychak.com/>



Gemini 1.5 Pro surpasses GPT-4o and Claude 3.5 Sonnet in the Chatbot Arena* !

On August 2, 2024

[Update from LMSYS](#)

Category: Overall

Overall Questions: #models: 123 (100%) #votes: 1,575,901 (100%)

Google Gemini 1.5 Pro (0801) #1 in Arena!!

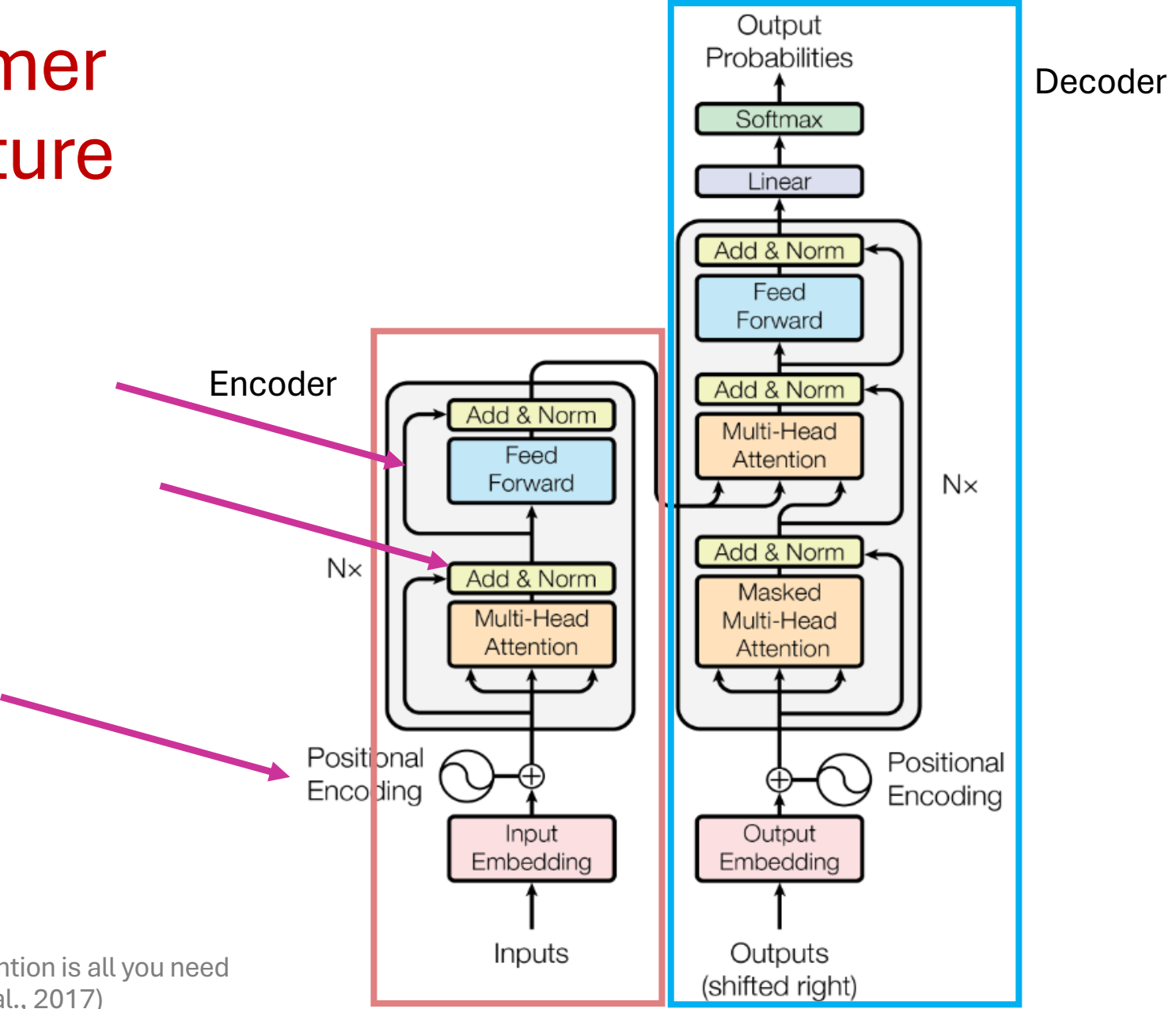
Rank* (UB)	Model	Arena Score	95% CI	Votes	Organization	License
1	Gemini-1.5-Pro-Exp-0801	1300	+6/-5	12672	Google	Proprietary
2	GPT-4o-2024-05-13	1286	+3/-2	69832	OpenAI	Proprietary
2	GPT-4o-mini-2024-07-18	1280	+6/-4	12047	OpenAI	Proprietary
4	Claude 3.5 Sonnet	1271	+3/-4	40174	Anthropic	Proprietary
4	Gemini-Advanced-0514	1266	+3/-4	50686	Google	Proprietary
4	Meta-Llama-3.1-405b-Instruct	1262	+6/-7	8454	Meta	Llama 3.1 Community
5	Gemini-1.5-Pro-API-0514	1261	+3/-3	62018	Google	Proprietary
6	Gemini-1.5-Pro-API-0409-Preview	1257	+4/-3	55673	Google	Proprietary
6	GPT-4-Turbo-2024-04-09	1257	+3/-2	79537	OpenAI	Proprietary
10	GPT-4-1106-preview	1251	+3/-2	90007	OpenAI	Proprietary
10	Claude 3 Opus	1248	+2/-3	151003	Anthropic	Proprietary
10	Athene-70b	1245	+6/-6	5511	NexusFlow	CC-BY-NC-4.0
10	Meta-Llama-3.1-70b-Instruct	1242	+10/-7	4450	Meta	Llama 3.1 Community

Prior to the release of **Gemini 1.5 Pro**, GPT-4o held the top score closely followed by Claude-3.5 Sonnet.

*Chatbot Arena is a benchmark platform for LLMs that features anonymous, randomized battles in a crowdsourced manner. The rankings are based on 1,000,000+ human votes for those battle based on “Elo rating system” (a rating system used in chess and similar competitive games)

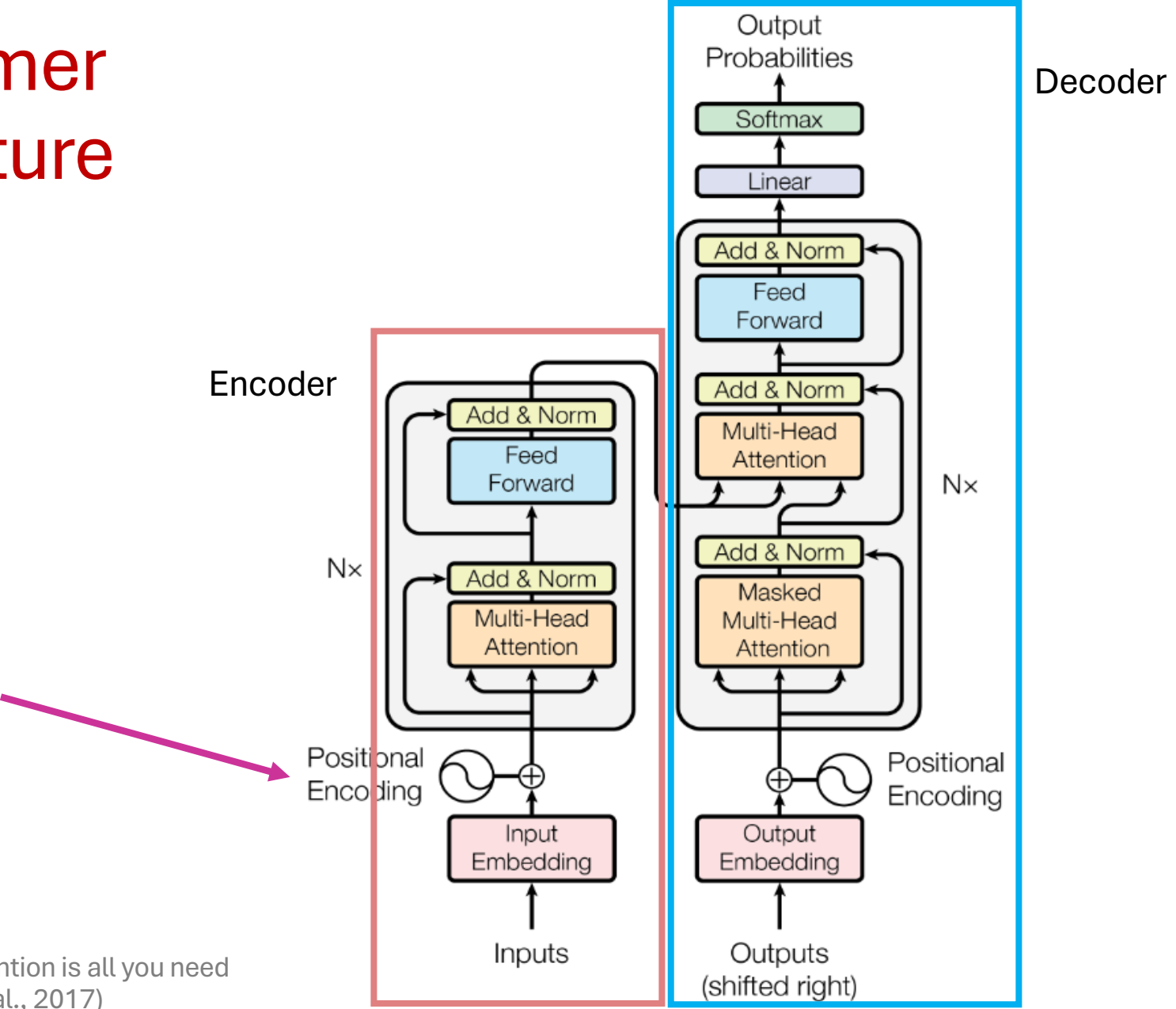
- ### Gemini 1.5 Pro Category Rankings:
- Overall: #1
 - Math: #1-3
 - Instruction-Following: #1-2
 - Coding: #3-5
 - Hard Prompts (English): #2-5

Transformer Architecture



Source of Image : Attention is all you need (Vaswani et al., 2017)

Transformer Architecture



Source of Image : Attention is all you need (Vaswani t al., 2017)

Positional Encoding

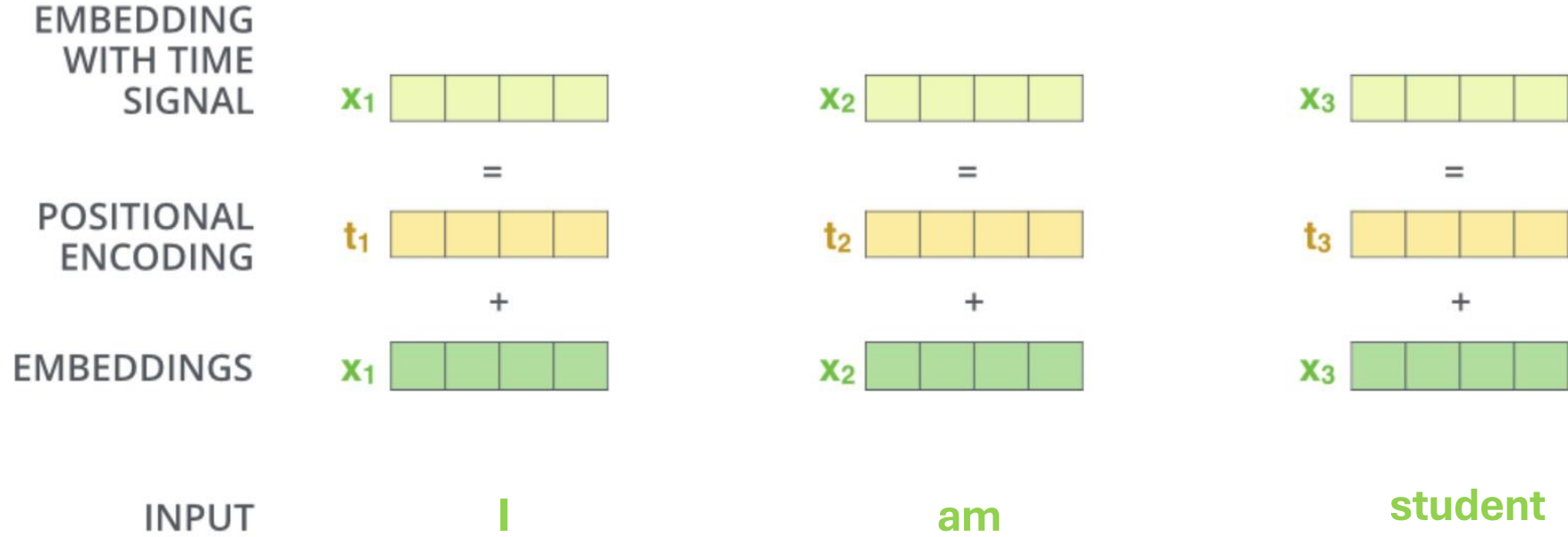
Position Information in Transformers: An Overview

Philipp Dufter, Martin Schmitt, Hinrich Schütze

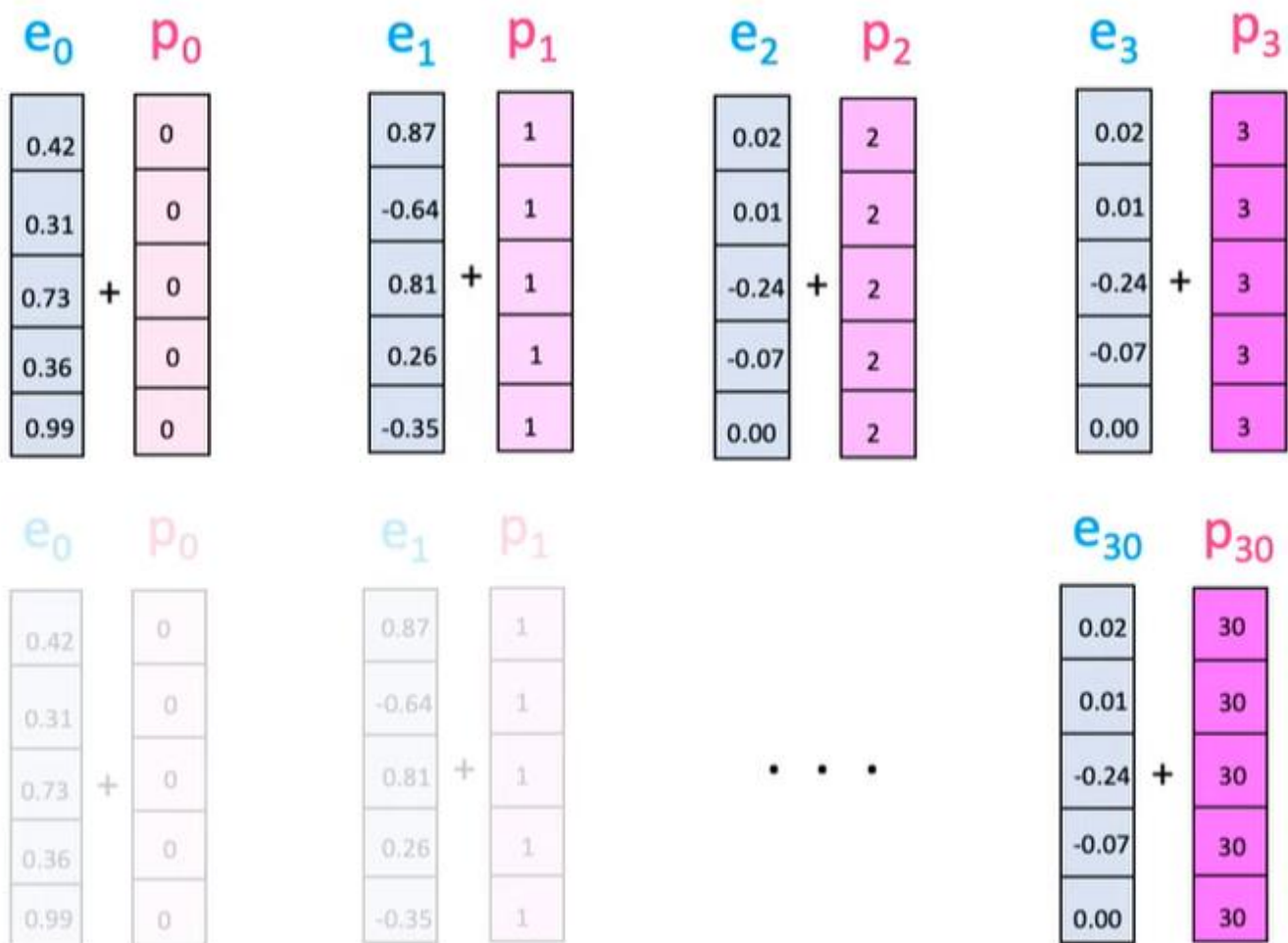
Abstract

Transformers are arguably the main workhorse in recent natural language processing research. By definition, a Transformer is invariant with respect to reordering of the input. However, language is inherently sequential and word order is essential to the semantics and syntax of an utterance. In this article, we provide an overview and theoretical comparison of existing methods to incorporate position information into Transformer models. The objectives of this survey are to (1) showcase that position information in Transformer is a vibrant and extensive research area; (2) enable the reader to compare existing methods by providing a unified notation and systematization of different approaches along important model dimensions; (3) indicate what characteristics of an application should be taken into account when selecting a position encoding; and (4) provide stimuli for future research.

Positional Encoding



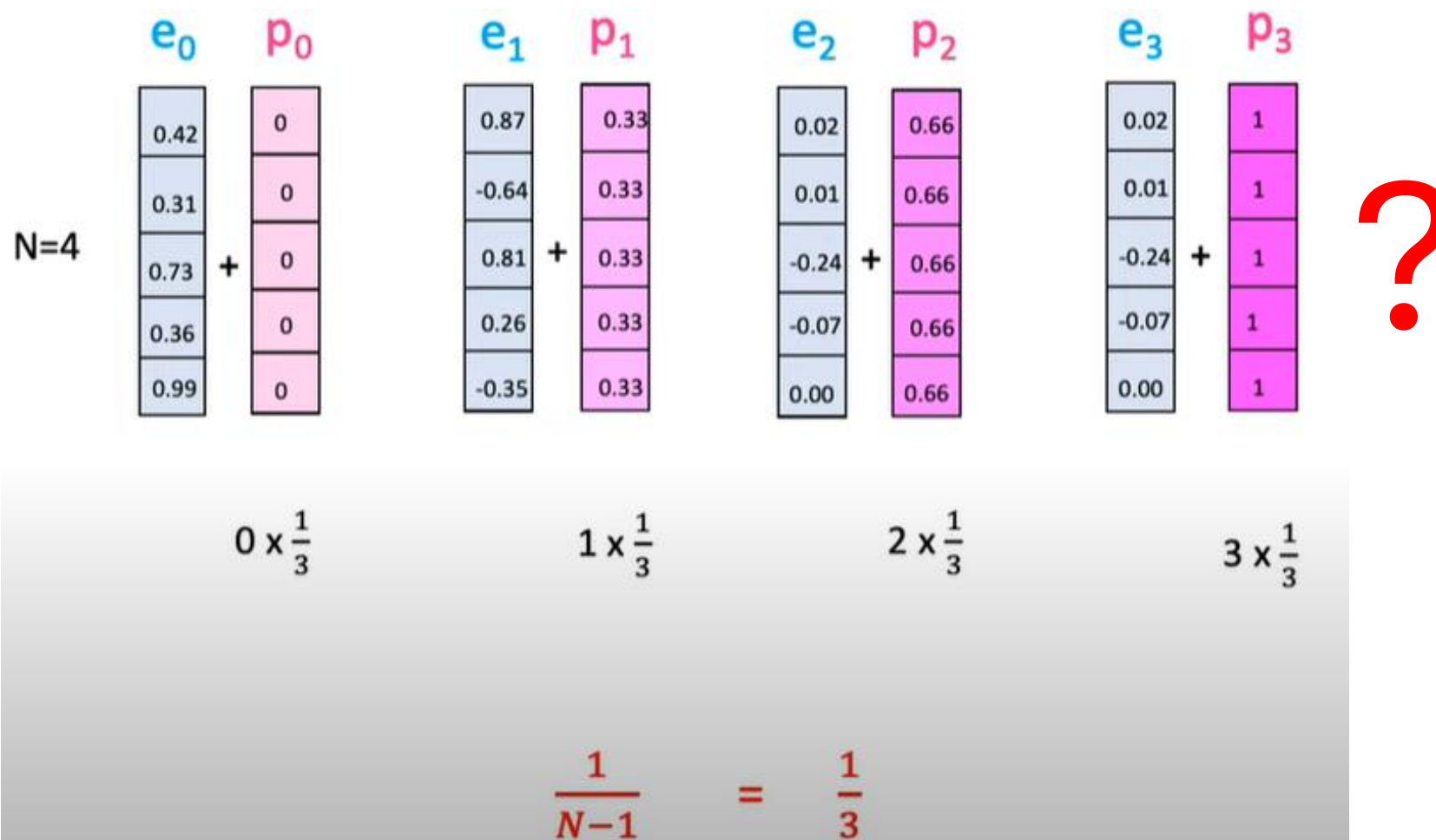
Option 1



Credits: <https://www.youtube.com/watch?v=dichlcUZfOw>



Option 2

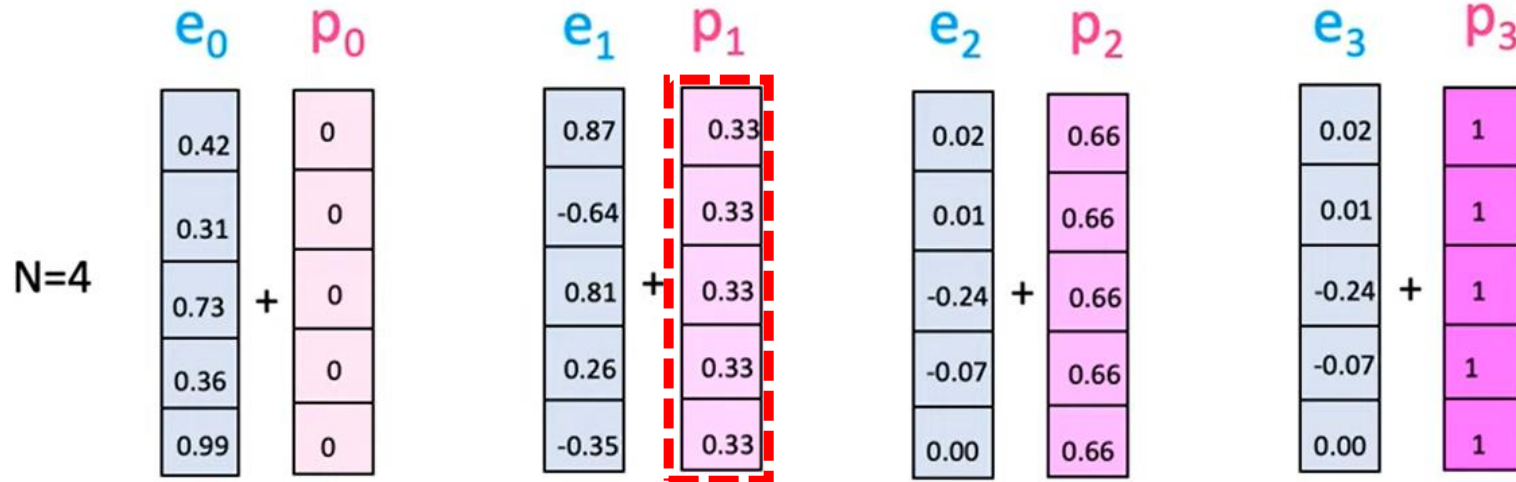


Credits: <https://www.youtube.com/watch?v=dichlcUZfOw>

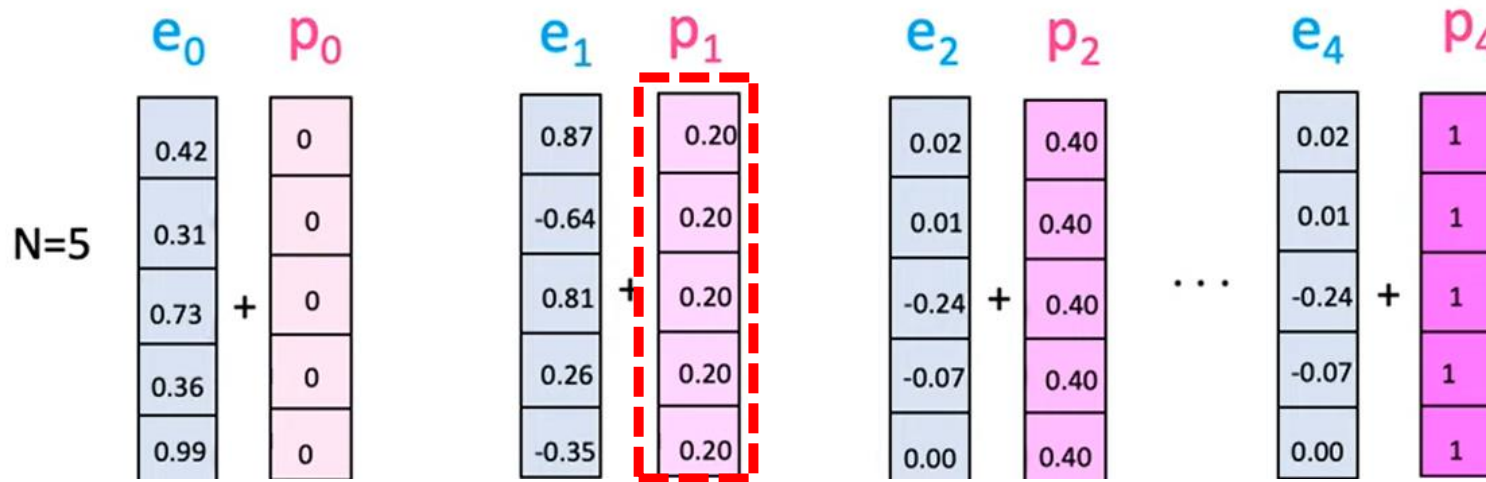


Option 2

Sentence 1



Sentence 2



The positional embedding vector at a given position should remain the same irrespective of the length of the sequence

Credits: <https://www.youtube.com/watch?v=dichlcUZfOw>



Creating Positional Encodings

- We could just concatenate a fixed value to each time step (e.g., 1, 2, 3, ... 1000) that corresponds to its position, but then what happens if we get a sequence with 5000 words at test time?
- We want something that can generalize to arbitrary sequence lengths. We also may want to make attending to *relative positions* (e.g., tokens in a local window to the current token) easier.
- Distance between two positions should be consistent with variable-length inputs



Intuitive Example

0 :	0	0	0	0	8 :	1	0	0	0
1 :	0	0	0	1	9 :	1	0	0	1
2 :	0	0	1	0	10 :	1	0	1	0
3 :	0	0	1	1	11 :	1	0	1	1
4 :	0	1	0	0	12 :	1	1	0	0
5 :	0	1	0	1	13 :	1	1	0	1
6 :	0	1	1	0	14 :	1	1	1	0
7 :	0	1	1	1	15 :	1	1	1	1

https://kazemnejad.com/blog/transformer_architecture_positional_encoding/



Transformer Positional Encoding

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

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

For $d_{model} = 512$,

Positional encoding is a 512-dimensional vector

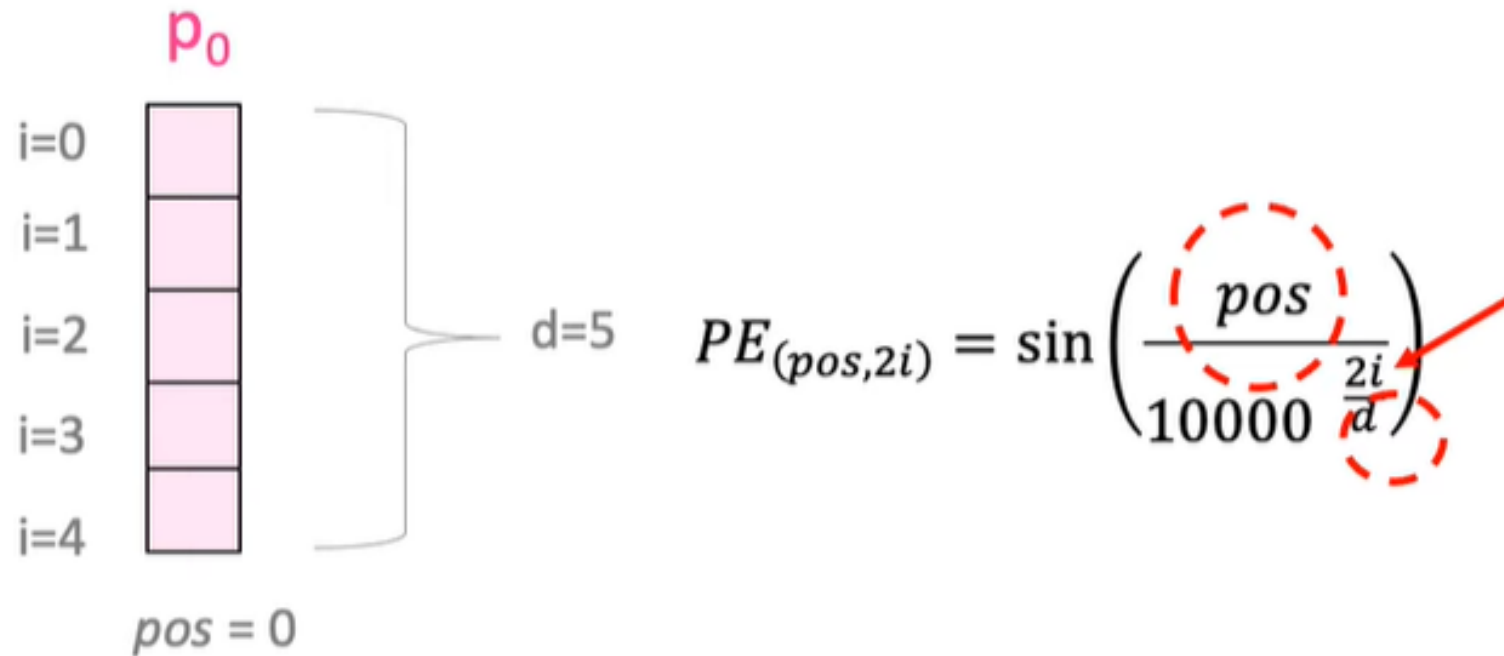
(Note: **Dimension of positional encoding is same as dimension of the word embeddings**)

i = a particular dimension of this vector

pos = position of the word in the sequence



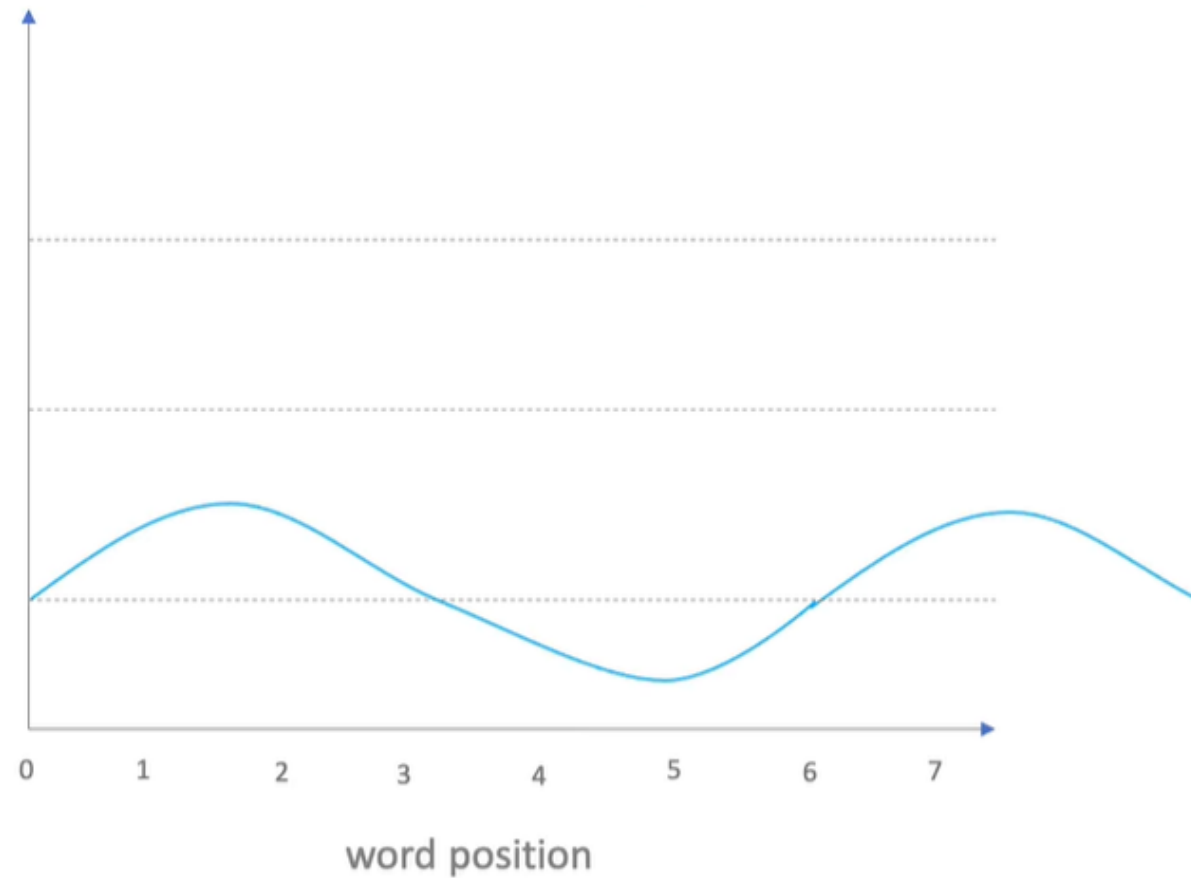
Transformer Positional Encoding



Credits: <https://www.youtube.com/watch?v=dichlcUZfOw>



Only Varying The Position

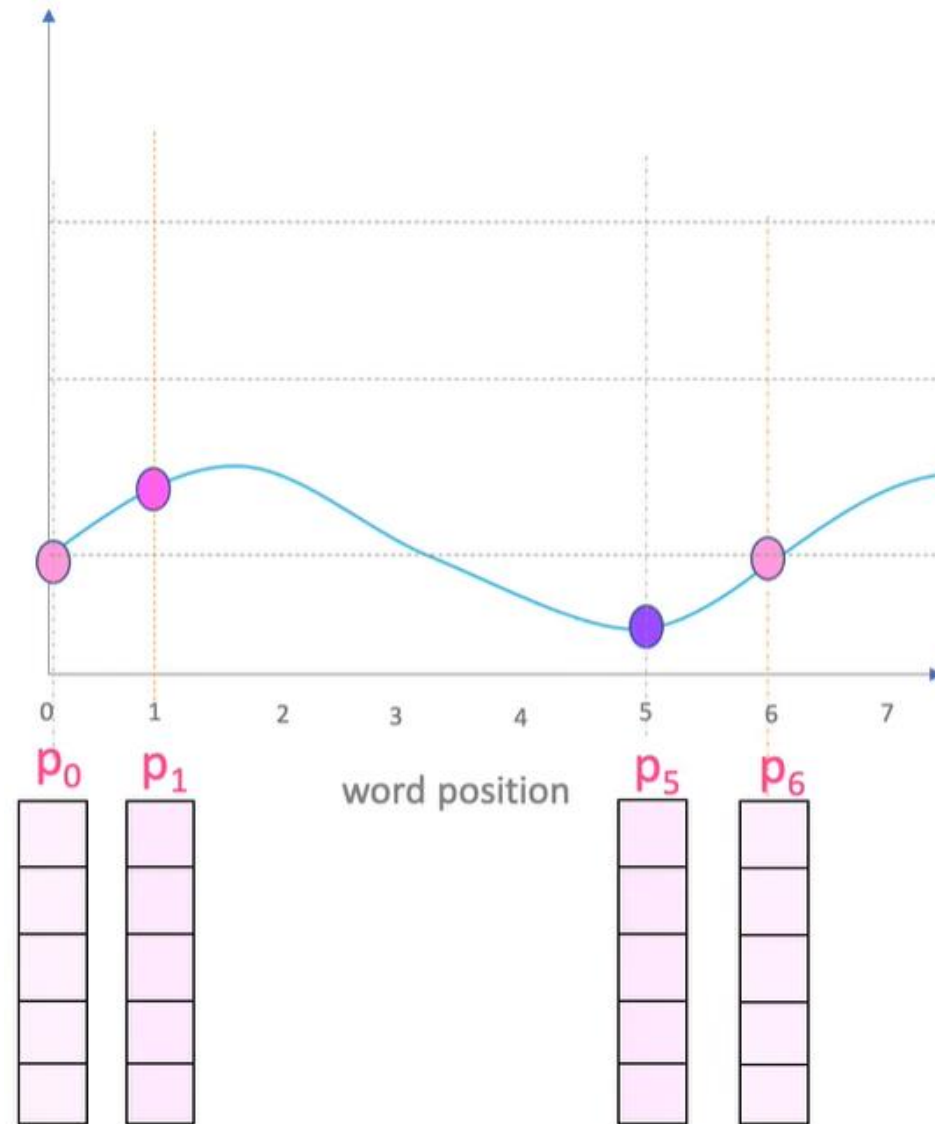


$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000} \frac{2i}{\pi}\right)$$

Credits: <https://www.youtube.com/watch?v=dichlcUZfOw>



Only Varying The Position



Pros

- It is independent of the length of the sequence

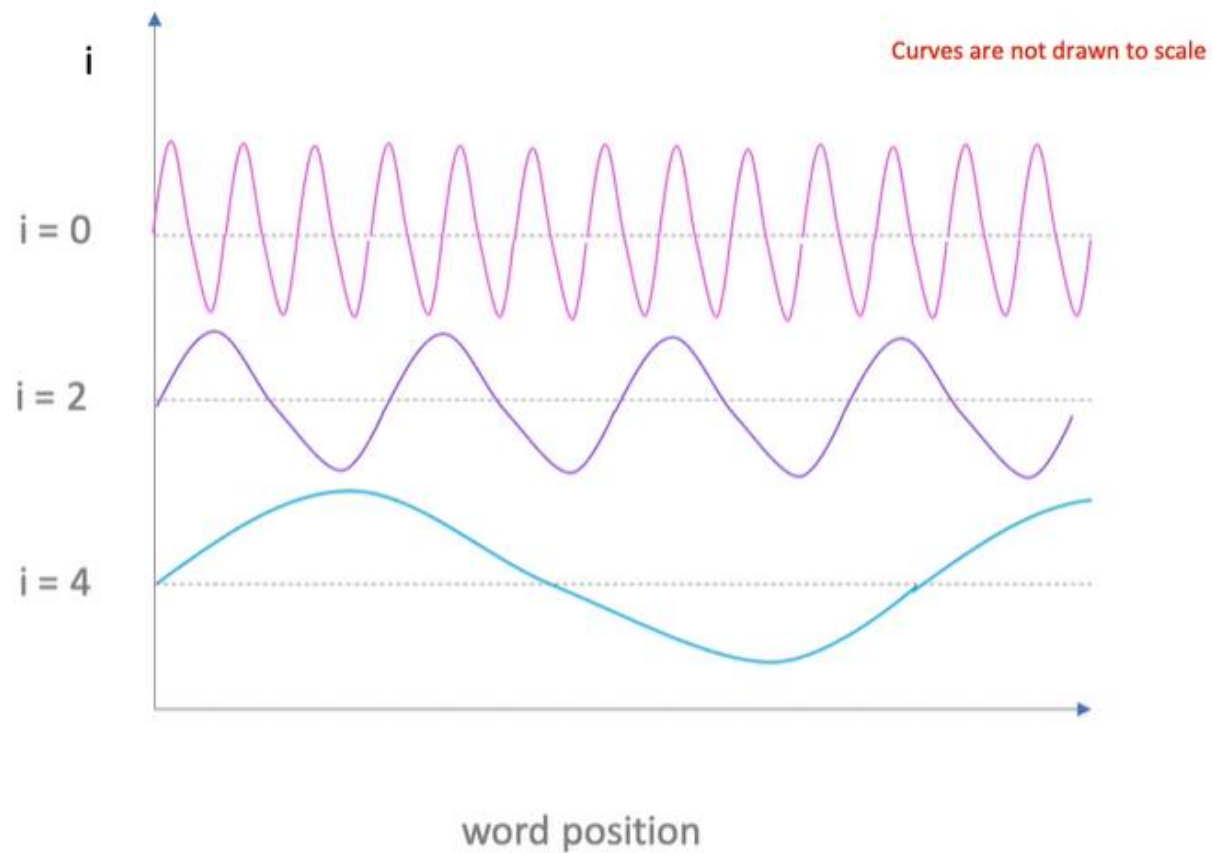
Cons

- p_0 and p_6 have same positional embeddings

Credits: <https://www.youtube.com/watch?v=dichlcUZfOw>



Varying Both Position and i

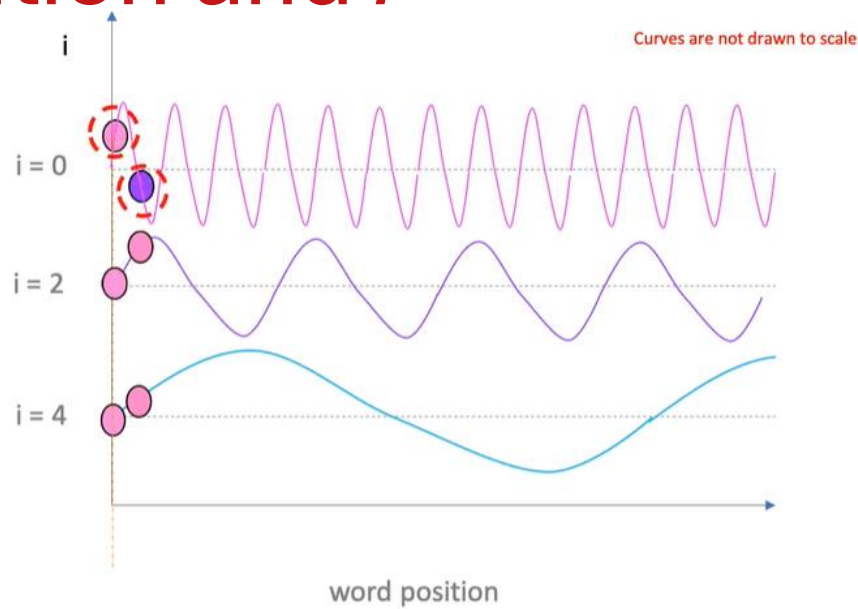


$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000} \frac{2i}{d}\right)$$

Credits: <https://www.youtube.com/watch?v=dichlcUZfOw>

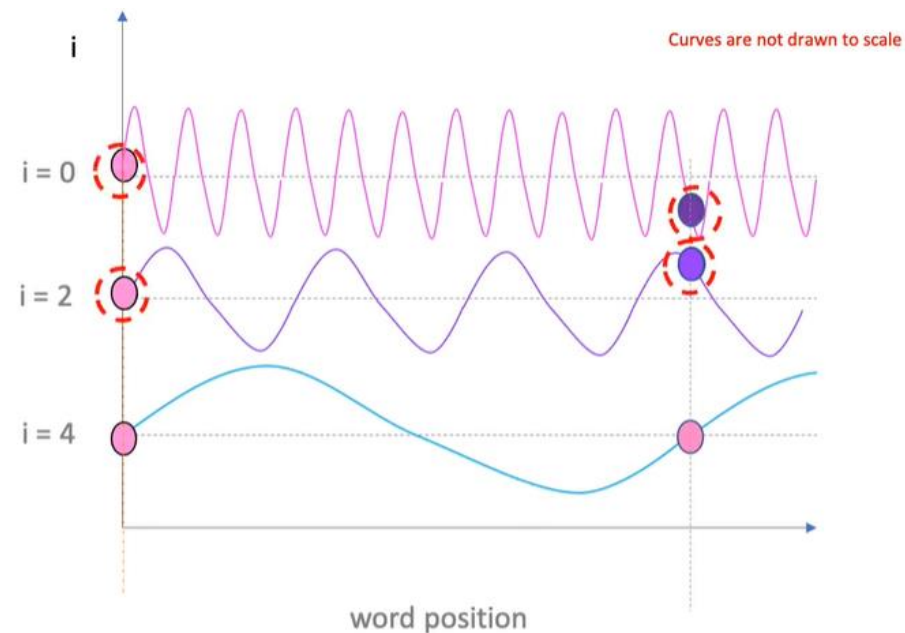


Varying Both Position and i



$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000} \frac{2i}{\alpha}\right)$$

If two points are close by on the curve, they will remain identical at low frequencies too. It is at the high frequency where their y-axis values differ and we may be able to take them apart.



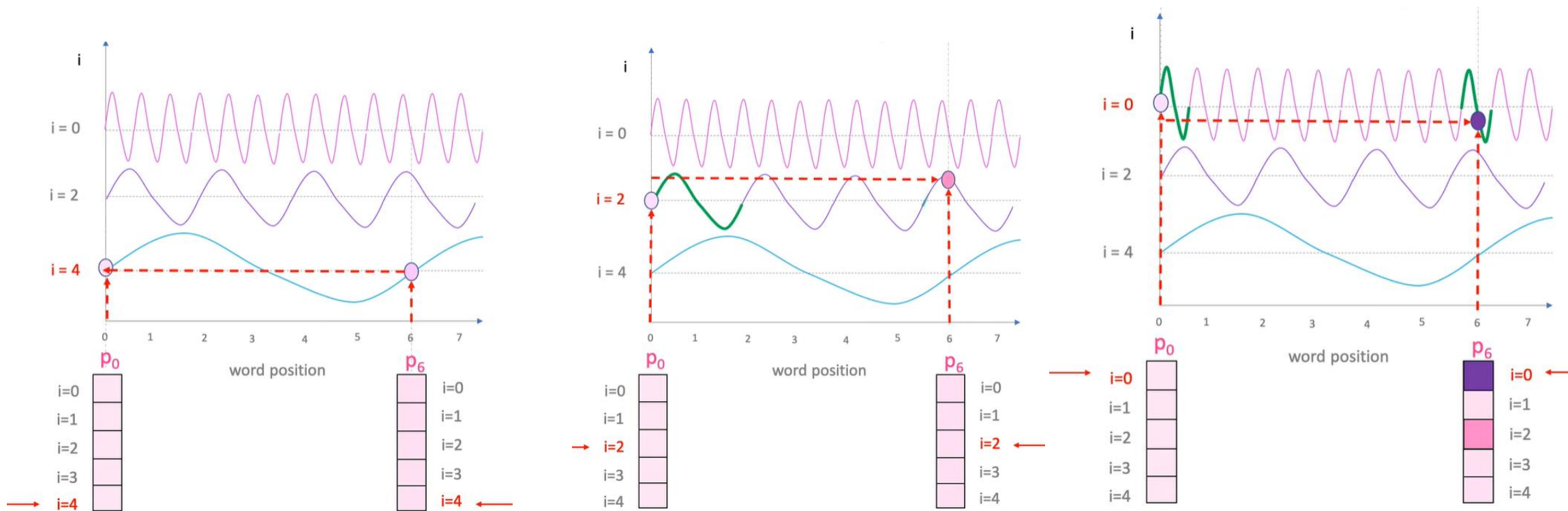
$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000} \frac{2i}{\alpha}\right)$$

For points which are far apart, we will see them falling apart on the y-axis quite early on

Credits: <https://www.youtube.com/watch?v=dichlcUZfOw>



Varying Both Position and i : Example



At lower frequency, the value is exactly the same

This will start to differ significantly as we move to high frequencies

Credits: <https://www.youtube.com/watch?v=dichlcUzfOw>



Example

For example, for word w at position $pos \in [0, L - 1]$ in the input sequence $\mathbf{w} = (w_0, \dots, w_{L-1})$, with 4-dimensional embedding e_w , and $d_{model} = 4$, the operation would be

$$\begin{aligned} e'_w &= e_w + \left[\sin \left(\frac{pos}{10000^0} \right), \cos \left(\frac{pos}{10000^0} \right), \sin \left(\frac{pos}{10000^{2/4}} \right), \cos \left(\frac{pos}{10000^{2/4}} \right) \right] \\ &= e_w + \left[\sin(pos), \cos(pos), \sin \left(\frac{pos}{100} \right), \cos \left(\frac{pos}{100} \right) \right] \end{aligned}$$

where the formula for positional encoding is as follows

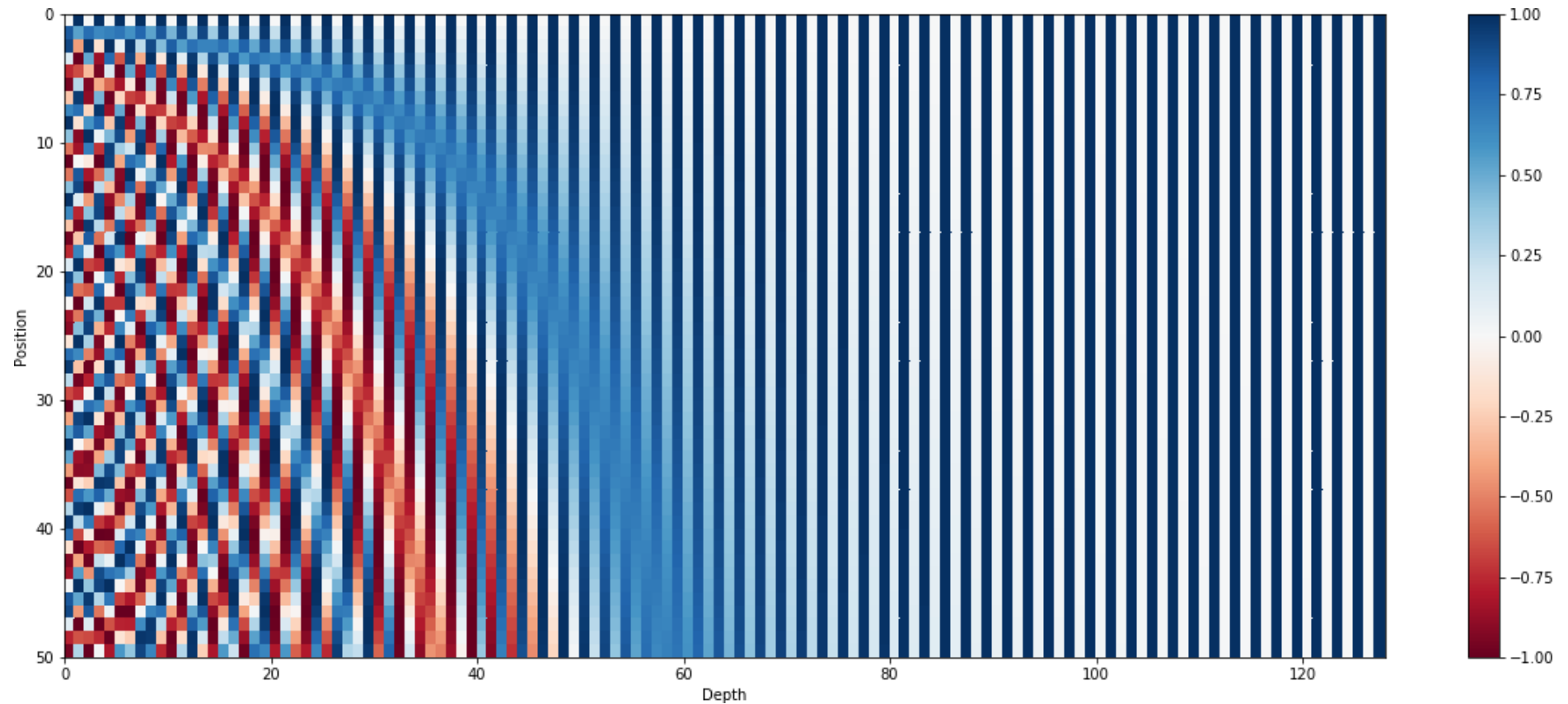
$$\begin{aligned} \text{PE}(pos, 2i) &= \sin \left(\frac{pos}{10000^{2i/d_{model}}} \right), \\ \text{PE}(pos, 2i + 1) &= \cos \left(\frac{pos}{10000^{2i/d_{model}}} \right). \end{aligned}$$

<https://datascience.stackexchange.com/questions/51065/what-is-the-positional-encoding-in-the-transformer-model>



What does this look like?

(each row is the positional encoding of a 50-word sentence)



https://kazemnejad.com/blog/transformer_architecture_positional_encoding/



Despite the intuitive flaws, many models these days use ***learned positional embeddings*** (i.e., they cannot generalize to longer sequences, but this isn't a big deal for their use cases)

General Properties of Positional Embeddings

We define the function $\phi(\cdot, \cdot)$ to measure the proximity between positional embeddings.

- **Monotonicity:** The proximity of position embeddings positions decreases when positions are further apart.

$$\forall x, m, n \in \mathbb{N} : m > n \iff \phi(\vec{x}, \overrightarrow{x + m}) < \phi(\vec{x}, \overrightarrow{x + n})$$

- **Translation invariance:** The proximity of embedded positions are translation invariant.

$$\forall x_1, \dots, x_n, m \in \mathbb{N} : \phi(\vec{x}_1, \overrightarrow{x_1 + m}) = \phi(\vec{x}_2, \overrightarrow{x_2 + m}) = \dots = \phi(\vec{x}_n, \overrightarrow{x_n + m})$$

- **Symmetry:** The proximity of embedded positions is symmetric.

$$\forall x, y \in \mathbb{N} : \phi(\vec{x}, \vec{y}) = \phi(\vec{y}, \vec{x})$$

Rotary Positional Encoding (RoPE)



ROFORMER: ENHANCED TRANSFORMER WITH ROTARY POSITION EMBEDDING

Jianlin Su

Zhuiyi Technology Co., Ltd.
Shenzhen
bojonesu@wezhuiyi.com

Yu Lu

Zhuiyi Technology Co., Ltd.
Shenzhen
julianlu@wezhuiyi.com

Shengfeng Pan

Zhuiyi Technology Co., Ltd.
Shenzhen
nickpan@wezhuiyi.com

Ahmed Murtadha

Zhuiyi Technology Co., Ltd.
Shenzhen
mengjiayi@wezhuiyi.com

Bo Wen

Zhuiyi Technology Co., Ltd.
Shenzhen
brucewen@wezhuiyi.com

Yunfeng Liu

Zhuiyi Technology Co., Ltd.
Shenzhen
glenliu@wezhuiyi.com

Adopted by

- PaLM
- GPT-Neo and GPT-J
- LLaMa 1 and 2

November 9, 2023

Credit: Slides of RoPE is adopted from Manish Gupta



Absolute Positional Encoding

- Learned from data
- Sinusoidal function (like the original Transformer)



Absolute Position Embeddings

$$\mathbb{E}_N = \{\mathbf{x}_i\}_{i=1}^N$$

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

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

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

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

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

- w_i are tokens and x_i are embeddings

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

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

- Generate p_i using the sinusoidal function
 - Sinusoidal functions provide continuity between close positions.
 - This also allows for the model to scale to virtually unlimited input sequence length



Relative Positional Encoding

The dog chased the pig

Once upon a time, the pig
chased the dog



Relative Position Embeddings

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

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

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

- $\tilde{\mathbf{p}}_r^k, \tilde{\mathbf{p}}_r^v \in R^d$ are trainable relative position embeddings.
- $r = \text{clip}(m-n, r_{min}, r_{max})$ is relative distance between positions m and n .
 - Relative position info is not useful beyond a certain distance.



Relative Position Embeddings

Transformer-XL

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

- Decompose $q_m^T k_n$

$$q_m^T k_n = \mathbf{x}_m^T \mathbf{W}_q^T \mathbf{W}_k \mathbf{x}_n + \mathbf{x}_m^T \mathbf{W}_q^T \mathbf{W}_k \mathbf{p}_n + \mathbf{p}_m^T \mathbf{W}_q^T \mathbf{W}_k \mathbf{x}_n + \mathbf{p}_m^T \mathbf{W}_q^T \mathbf{W}_k \mathbf{p}_n$$

- Replace absolute position embedding p_n with relative \tilde{p}_{m-n}
- Replace absolute position embedding p_m with two trainable vectors u and v independent of query positions.
- W_k is distinguished for content-based and location-based key vectors x_n and p_n , denoted as W_k and \tilde{W}_k

$$q_m^T k_n = \mathbf{x}_m^T \mathbf{W}_q^T \mathbf{W}_k \mathbf{x}_n + \mathbf{x}_m^T \mathbf{W}_q^T \tilde{\mathbf{W}}_k \tilde{\mathbf{p}}_{m-n} + \mathbf{u}^T \mathbf{W}_q^T \mathbf{W}_k \mathbf{x}_n + \mathbf{v}^T \mathbf{W}_q^T \tilde{\mathbf{W}}_k \tilde{\mathbf{p}}_{m-n}$$



Relative Position Embeddings

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

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

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

- $\tilde{\mathbf{p}}_r^k, \tilde{\mathbf{p}}_r^v \in R^d$ are trainable relative position embeddings.
- $r = \text{clip}(m-n, r_{\min}, r_{\max})$ is relative distance between positions m and n .
 - Relative position info is not useful beyond a certain distance.

Transformer-XL

- Decompose $q_m^T k_n$

$$q_m^T k_n = \mathbf{x}_m^T \mathbf{W}_q^T \mathbf{W}_k \mathbf{x}_n + \mathbf{x}_m^T \mathbf{W}_q^T \mathbf{W}_k \mathbf{p}_n + \mathbf{p}_m^T \mathbf{W}_q^T \mathbf{W}_k \mathbf{x}_n + \mathbf{p}_m^T \mathbf{W}_q^T \mathbf{W}_k \mathbf{p}_n$$

- Replace abs pos embedding p_n with relative \tilde{p}_{m-n}
- Replace abs pos embedding p_m with two trainable vectors u and v independent of query positions.
- W_k is distinguished for content-based and location-based key vectors x_n and p_n , denoted as W_k and \tilde{W}_k

$$q_m^T k_n = \mathbf{x}_m^T \mathbf{W}_q^T \mathbf{W}_k \mathbf{x}_n + \mathbf{x}_m^T \mathbf{W}_q^T \tilde{\mathbf{W}}_k \tilde{\mathbf{p}}_{m-n} + \mathbf{u}^T \mathbf{W}_q^T \mathbf{W}_k \mathbf{x}_n + \mathbf{v}^T \mathbf{W}_q^T \tilde{\mathbf{W}}_k \tilde{\mathbf{p}}_{m-n}$$

- T5 uses a very simplified relative position embedding

$$q_m^T k_n = \mathbf{x}_m^T \mathbf{W}_q^T \mathbf{W}_k \mathbf{x}_n + b_{i,j}$$

- $b_{i,j}$ is a trainable bias.
- Another formulation

$$q_m^T k_n = \mathbf{x}_m^T \mathbf{W}_q^T \mathbf{W}_k \mathbf{x}_n + \mathbf{p}_m^T \mathbf{U}_q^T \mathbf{U}_k \mathbf{p}_n + b_{i,j}$$



Relative Position Embeddings

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

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

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

- $\tilde{\mathbf{p}}_r^k, \tilde{\mathbf{p}}_r^v \in R^d$ are trainable relative position embeddings.
- $r = \text{clip}(m-n, r_{\min}, r_{\max})$ is relative distance between positions m and n .
 - Relative position info is not useful beyond a certain distance.

Transformer-XL

- Decompose $q_m^T k_n$

$$q_m^T k_n = \mathbf{x}_m^T \mathbf{W}_q^T \mathbf{W}_k \mathbf{x}_n + \mathbf{x}_m^T \mathbf{W}_q^T \mathbf{W}_k \mathbf{p}_n + \mathbf{p}_m^T \mathbf{W}_q^T \mathbf{W}_k \mathbf{x}_n + \mathbf{p}_m^T \mathbf{W}_q^T \mathbf{W}_k \mathbf{p}_n$$

- Replace abs pos embedding p_n with relative \tilde{p}_{m-n}
- Replace abs pos embedding p_m with two trainable vectors u and v independent of query positions.
- W_k is distinguished for content-based and location-based key vectors x_n and p_n , denoted as W_k and \tilde{W}_k

- T5 uses a very simplified relative position embedding

$$q_m^T k_n = \mathbf{x}_m^T \mathbf{W}_q^T \mathbf{W}_k \mathbf{x}_n + b_{i,j}$$

- $b_{i,j}$ is a trainable bias.
- Another formulation

$$q_m^T k_n = \mathbf{x}_m^T \mathbf{W}_q^T \mathbf{W}_k \mathbf{x}_n + \mathbf{p}_m^T \mathbf{U}_q^T \mathbf{U}_k \mathbf{p}_n + b_{i,j}$$

- DeBERTa

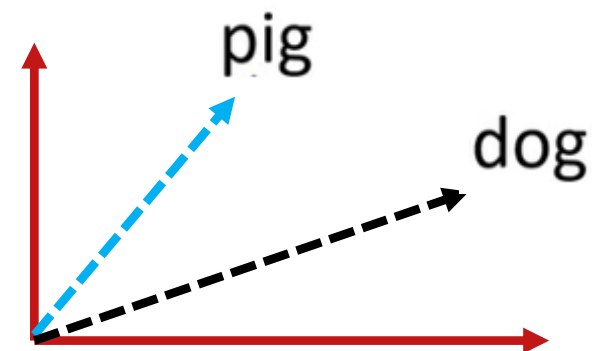
- Absolute position embeddings p_m and p_n are replaced with the relative position embeddings \tilde{p}_{m-n}

$$q_m^T k_n = \mathbf{x}_m^T \mathbf{W}_q^T \mathbf{W}_k \mathbf{x}_n + \mathbf{x}_m^T \mathbf{W}_q^T \mathbf{W}_k \tilde{\mathbf{p}}_{m-n} + \tilde{\mathbf{p}}_{m-n}^T \mathbf{W}_q^T \mathbf{W}_k \mathbf{x}_n$$



Combining both Relative and Absolute Positions

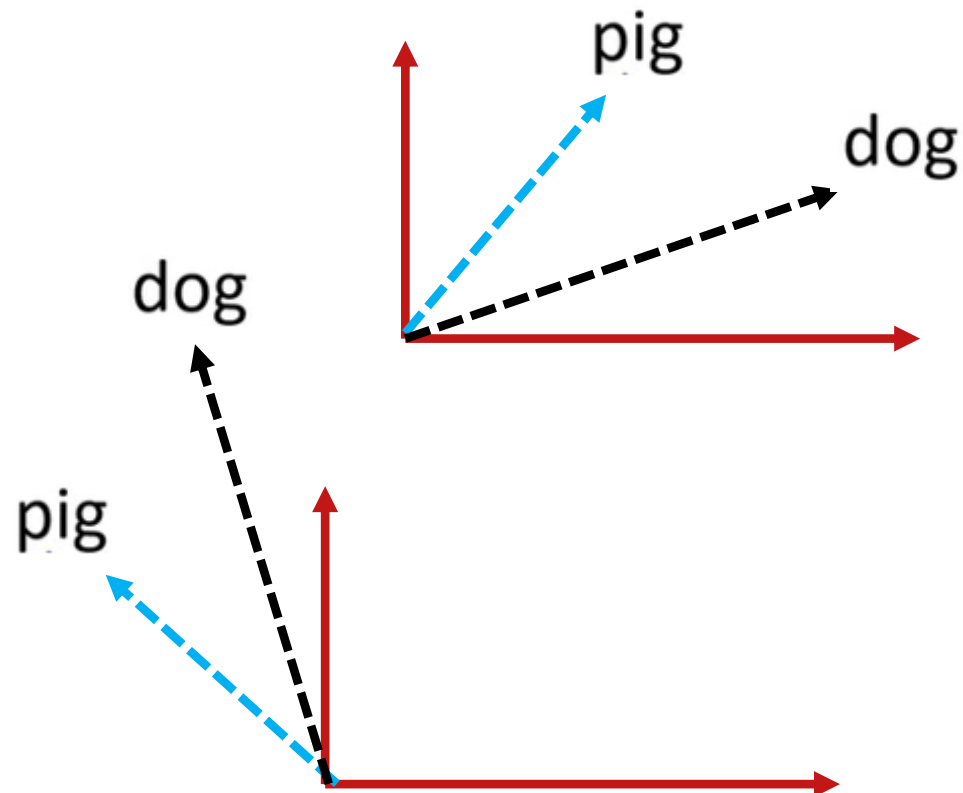
The dog chased the pig



Combining both Relative and Absolute Positions

The dog chased the pig

Once upon a time, the pig chased the dog



Rotary Position Embedding (RoPE)

- Encodes abs pos with a rotation matrix
- Incorporates explicit relative pos dependency in self-attention formulation.
- Require: $q_m^T k_n$ be a function (g) of only word embeddings $\mathbf{x}_m, \mathbf{x}_n$, and their relative position $m-n$

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

- What is a good choice for f_q and f_k ?
- For dimension $d = 2$, a solution is

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

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

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

- $\text{Re}[\cdot]$ is real part of a complex number
- $(\mathbf{W}_k \mathbf{x}_n)^*$ is conjugate complex number of $(\mathbf{W}_k \mathbf{x}_n)$.
- $\theta \in \mathbb{R}$ is a preset non-zero constant.

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

- Euler's formula: $e^{im\theta} = \cos(m\theta) + i \sin(m\theta)$
- The matrix form of $e^{im\theta}$ is $\begin{pmatrix} \cos m\theta & -\sin m\theta \\ \sin m\theta & \cos m\theta \end{pmatrix}$
- $\text{Re}[e^{im\theta}] = \cos(m\theta)$; $\text{Im}[e^{im\theta}] = \sin(m\theta)$
- This matrix captures the rotation by $m\theta$.
- Rotary Position Embedding
 - Rotate the affine-transformed word embedding vector by amount of angle multiples of its position index



Rotary Position Embedding (RoPE): General Form

- Generalizing from 2D to d-dimensions
 - Divide the d-dimension space into d/2 sub-spaces

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

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

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

- Applying RoPE to self-attention

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

$$\mathbf{R}_{\Theta, n-m}^d = (\mathbf{R}_{\Theta, m}^d)^\top \mathbf{R}_{\Theta, n}^d$$



Rotary Position Embedding (RoPE): General Form

- Generalizing from 2D to d-dimensions
 - Divide the d-dimension space into d/2 sub-spaces

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

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

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

- Applying RoPE to self-attention

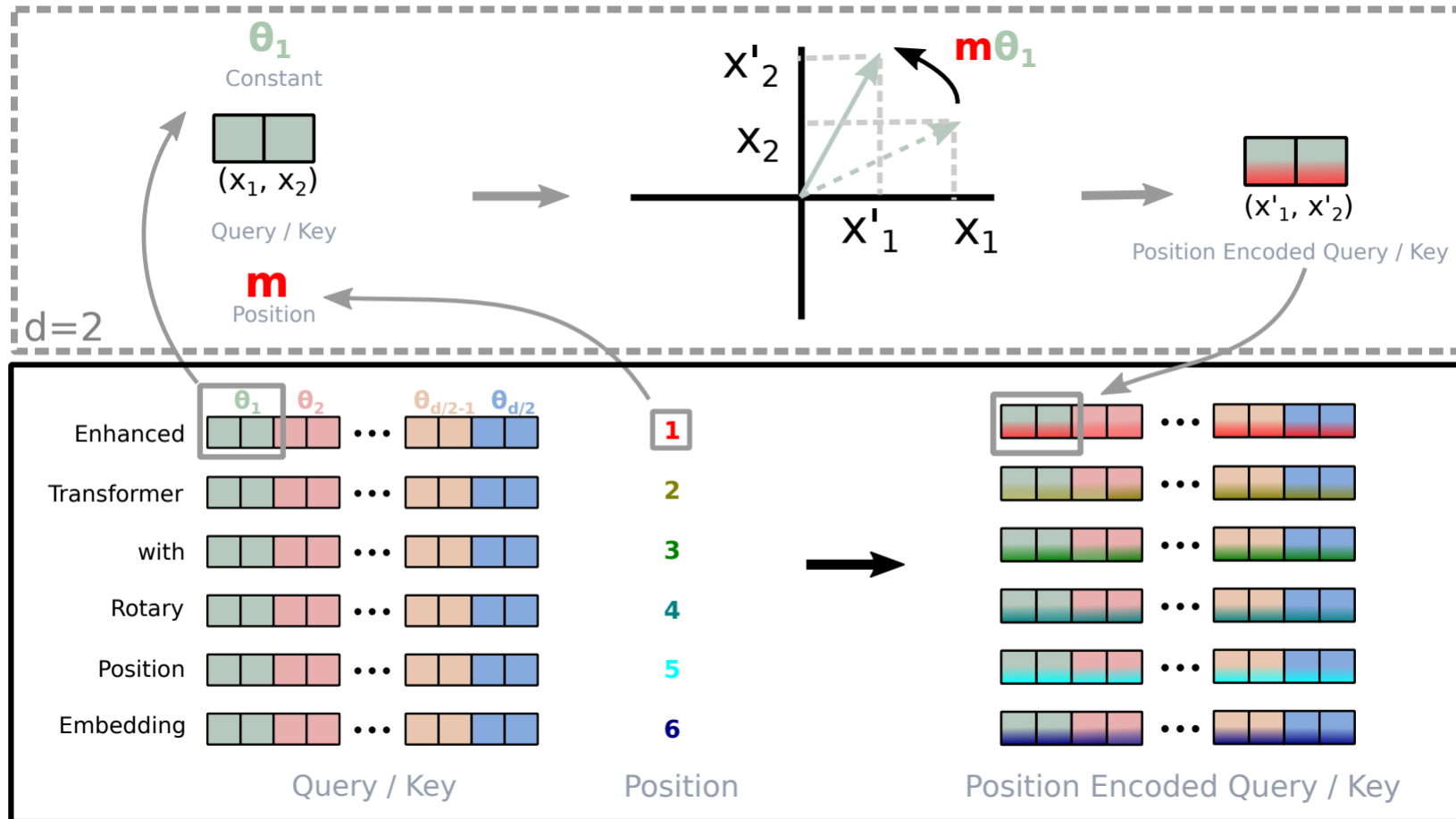
$$\mathbf{q}_m^\top \mathbf{k}_n = (\mathbf{R}_{\Theta, m}^d \mathbf{W}_q \mathbf{x}_m)^\top (\mathbf{R}_{\Theta, n}^d \mathbf{W}_k \mathbf{x}_n) = \mathbf{x}_m^\top \mathbf{W}_q \mathbf{R}_{\Theta, n-m}^d \mathbf{W}_k \mathbf{x}_n$$

$$\mathbf{R}_{\Theta, n-m}^d = (\mathbf{R}_{\Theta, m}^d)^\top \mathbf{R}_{\Theta, n}^d$$

- In contrast to earlier position embedding methods, RoPE is multiplicative.
- RoPE naturally incorporates relative pos info through rotation matrix product instead of altering terms in the expanded formulation of additive position encoding when applied with self-attention.
- RoPE
 - Represents token embeddings as complex numbers
 - Represents their positions as pure rotations
 - If we shift both the query and key by the same amount, changing absolute position but not relative position, this will lead both representations to be additionally rotated in the same manner, thus the angle between them will remain unchanged and thus the dot product will also remain unchanged.



RoPE Implementation



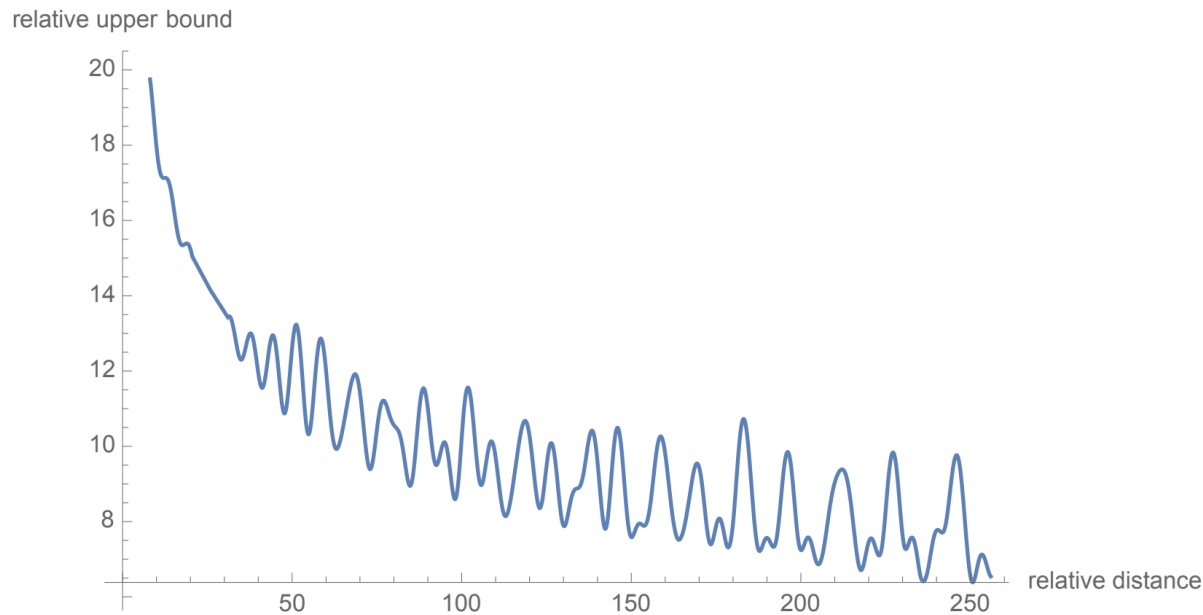
$$\Theta = \{\theta_i = 10000^{-2(i-1)/d}, i \in [1, 2, \dots, d/2]\}$$

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



Properties of RoPE

- Long-term decay
 - Inner-product decays when the relative position increase.
 - A pair of tokens with a long relative distance should have less connection.



- Computational efficient realization of rotary matrix multiplication

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

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



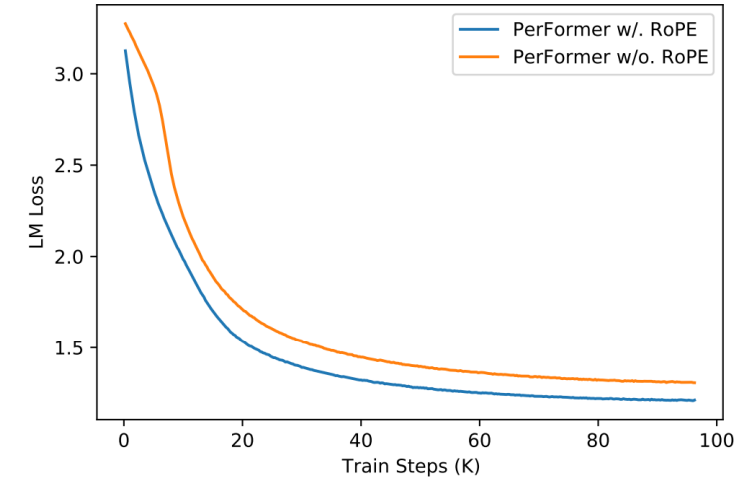
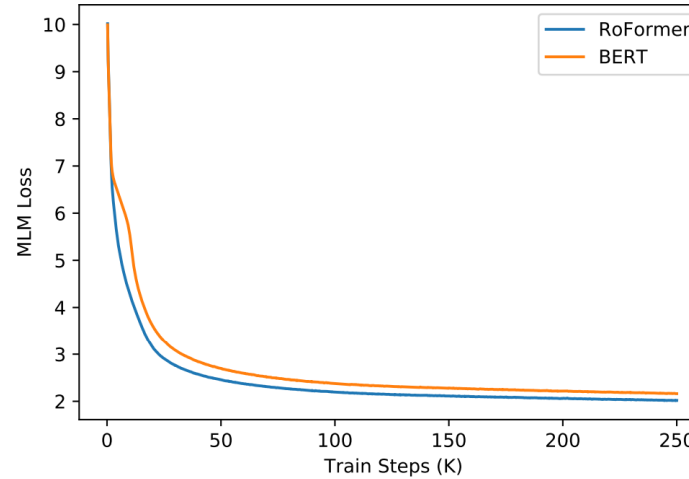
$$\mathbf{R}_{\Theta, m}^d \mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ \vdots \\ x_{d-1} \\ x_d \end{pmatrix} \otimes \begin{pmatrix} \cos m\theta_1 \\ \cos m\theta_1 \\ \cos m\theta_2 \\ \cos m\theta_2 \\ \vdots \\ \cos m\theta_{d/2} \\ \cos m\theta_{d/2} \end{pmatrix} + \begin{pmatrix}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



RoPE Performance

Model	BLEU
Transformer-base Vaswani et al. [2017]	27.3
RoFormer	27.5

WMT 2014 English-to-German translation task



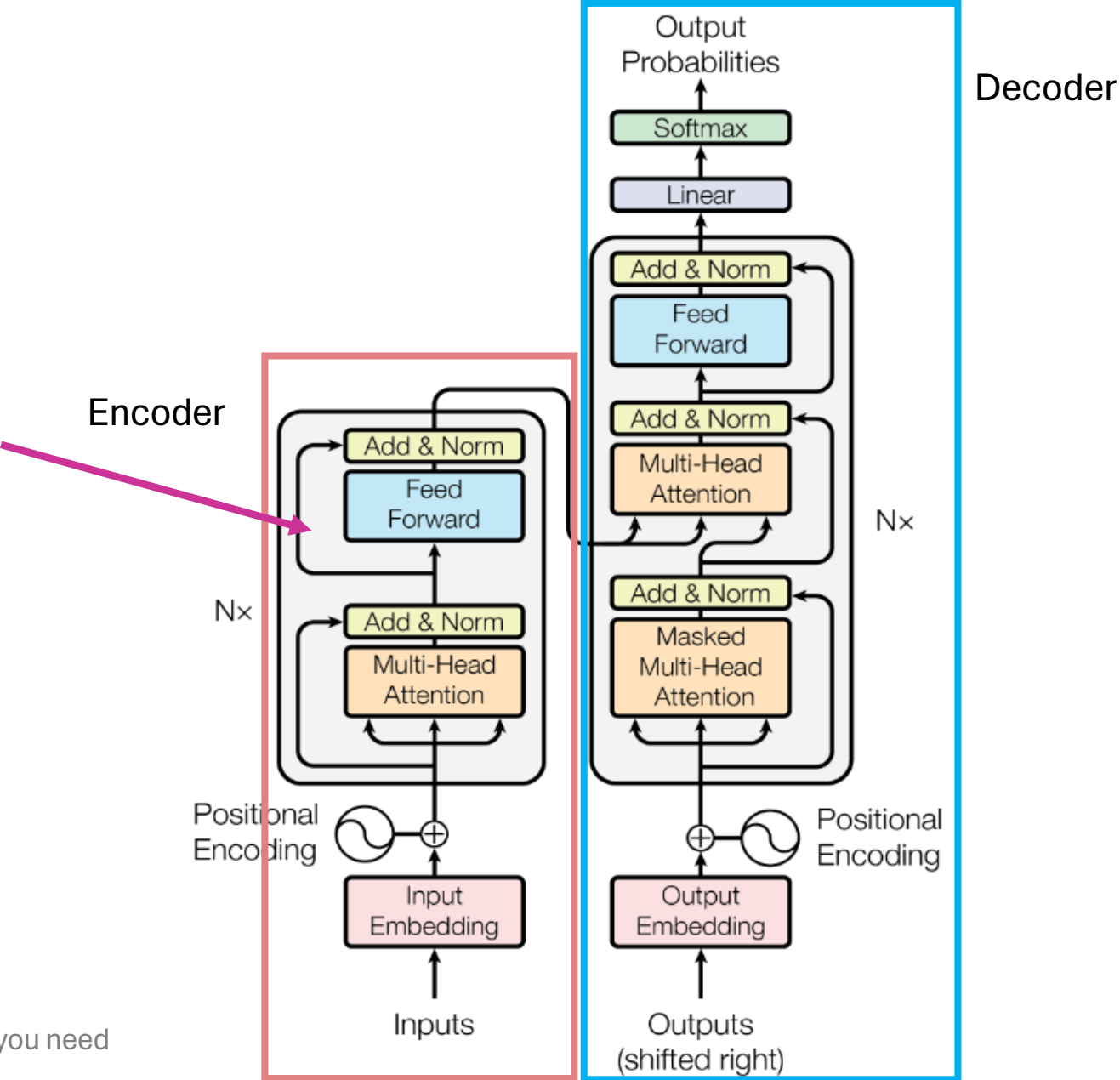
Language modeling pre-training. Left: training loss. Right: training loss for PerFormer with and without RoPE.

Table 2: Comparing RoFormer and BERT by fine tuning on downstream GLEU tasks.

Model	MRPC	SST-2	QNLI	STS-B	QQP	MNLI(m/mm)
BERT Devlin et al. [2019]	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



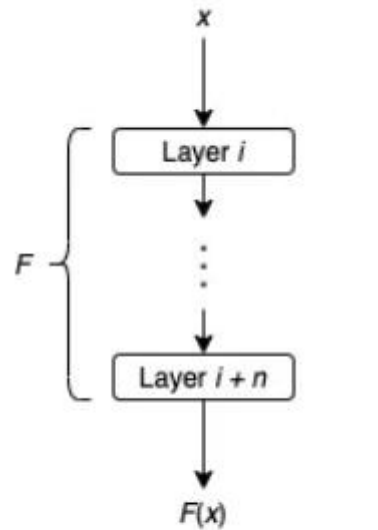
Transformer Architecture



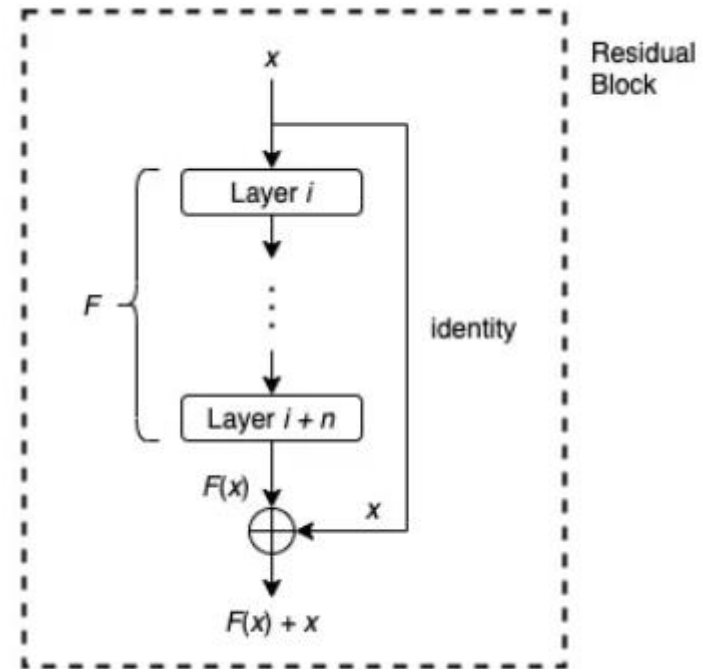
Source of Image : Attention is all you need (Vaswani et al., 2017)



Residual Connection



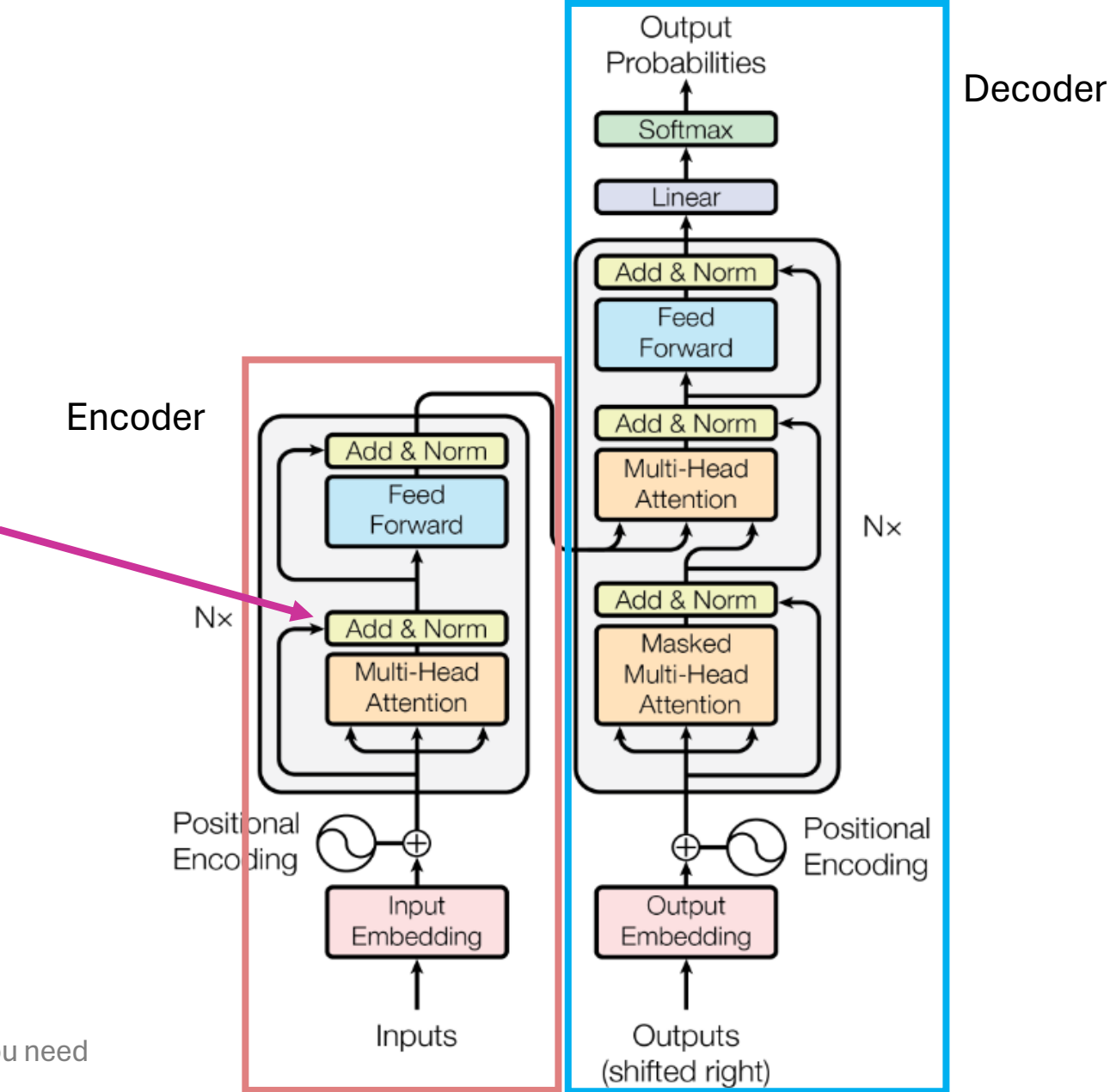
Traditional Feedforward
without Residual Connection



With Residual Connection



Transformer Architecture



Source of Image : Attention is all you need (Vaswani et al., 2017)



Batch and Layer Normalization



Normalization

- Example: student loans with the age of the student and the tuition as two input features
 - two values are on totally different scales.
 - the age of a student will have a median value in the range 18 to 25 years
 - the tuition could take on values in the range \$20K - \$50K for a given academic year.

Normalization works by mapping all values of a feature to be in the range [0,1] using the transformation

$$x_{norm} = \frac{x - x_{min}}{x_{max} - x_{min}}$$

Suppose a particular input feature x has values in the range $[x_{min}, x_{max}]$. When x is equal to x_{min} , x_{norm} is equal to 0 and when x is equal to x_{max} , x_{norm} is equal to 1. So for all values of x between x_{min} and x_{max} , x_{norm} maps to a value between 0 and 1.



Standardization

- Example: student loans with the age of the student and the tuition as two input features
 - two values are on totally different scales.
 - the age of a student will have a median value in the range 18 to 25 years
 - the tuition could take on values in the range \$20K - \$50K for a given academic year.

Standardization transforms the input values such that they follow a distribution with zero mean and unit variance.

$$x_{std} = \frac{x - \mu}{\sigma}$$

In practice, this process of *standardization* is also referred to as *normalization*



Batch Normalization

- For a network with hidden layers, the output of layer $k-1$ serves as the input to layer k
- Split the dataset into multiple batches and run the mini-batch gradient descent.
- The mini-batch gradient descent algorithm optimizes the parameters of the neural network by batch-wise processing of the dataset, one batch at a time.
- The input distribution at a particular layer keeps changing across batches.
- **Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift refers to this change in distribution of the input to a particular layer across batches as internal covariate shift.**
- For instance, if the distribution of data at the input of layer K keeps changing across batches, the network will take longer to train.



Batch Normalization

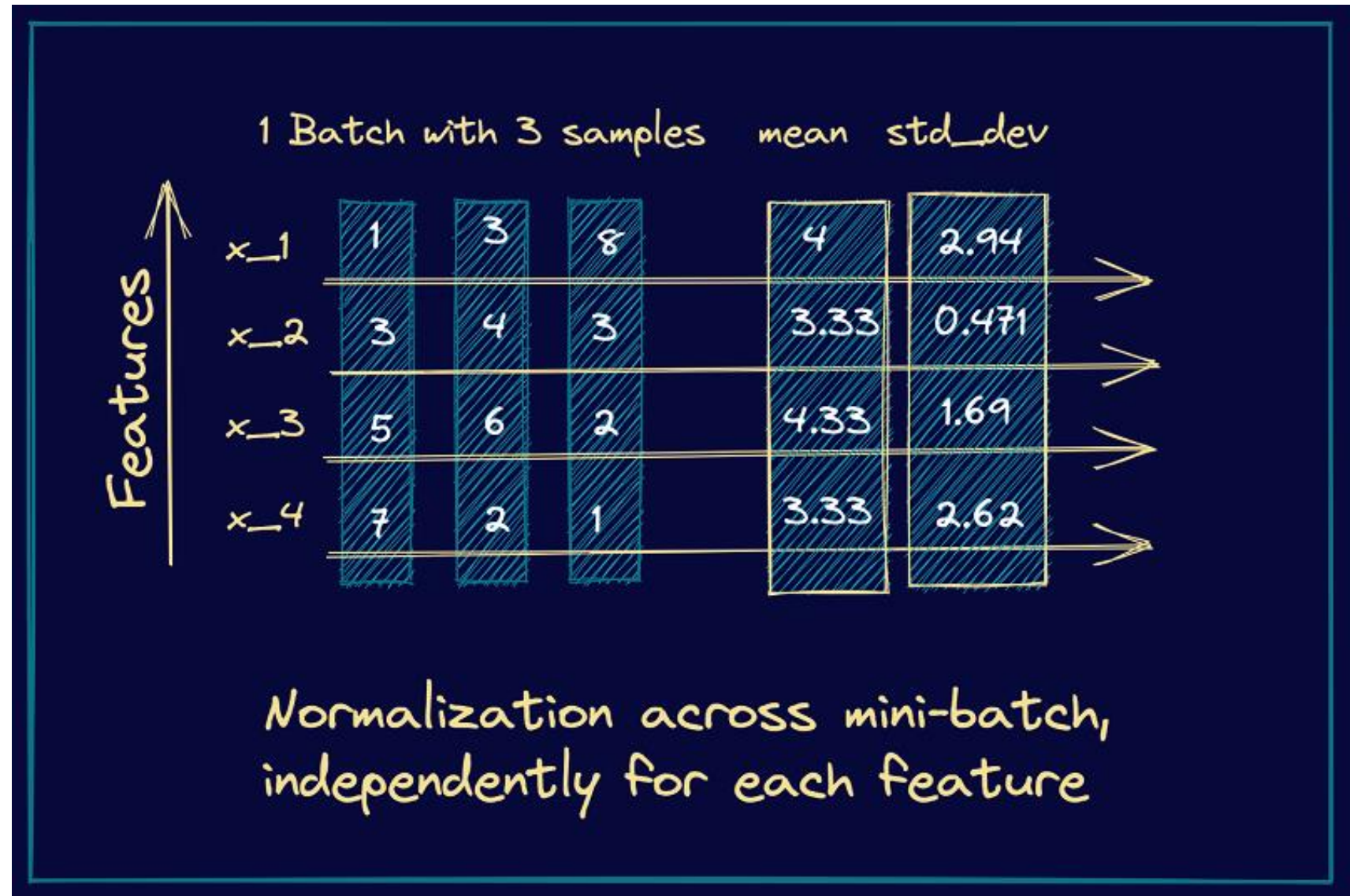


Image Source: <https://www.pinecone.io/learn/batch-layer-normalization/>



Batch Normalization

- Forcing all the pre-activations to be zero and unit standard deviation across all batches can be too restrictive.
- It may be the case that the fluctuant distributions are necessary for the network to learn certain classes better.

$$\mu_b = \frac{1}{B} \sum_{i=1}^B x_i \quad (1)$$

$$\sigma_b^2 = \frac{1}{B} \sum_{i=1}^B (x_i - \mu_b)^2 \quad (2)$$

$$\hat{x}_i = \frac{x_i - \mu_b}{\sqrt{\sigma_b^2}} \quad (3)$$

$$\text{or } \hat{x}_i = \frac{x_i - \mu_b}{\sqrt{\sigma_b^2 + \epsilon}} \quad (3)$$

Adding ϵ helps when σ_b^2 is small

$$y_i = \mathcal{BN}(x_i) = \gamma \cdot x_i + \beta \quad (4)$$

Trainable parameters



Batch Normalization: Limitations

- In batch normalization, we use **batch statistics**: mean and standard deviation for current mini-batch.
 - When **batch size is small, the sample mean and sample standard deviation are not representative enough** of the actual distribution and the network cannot learn anything meaningful.
- As batch normalization depends on batch statistics for normalization, it is **less suited for sequence models**.
 - This is because, in sequence models, we may have sequences of potentially different lengths and smaller batch sizes corresponding to longer sequences.



Layer Normalization

- All neurons in a particular layer effectively have the same distribution across all features for a given input

If each input has d features, it's a d -dimensional vector. If there are B elements in a batch, the normalization is done along the length of the d -dimensional vector and not across the batch of size B .

- Normalizing *across all features* but for each of the inputs to a specific layer removes the dependence on batches.
- This makes layer normalization well suited for sequence models such as Transformers and RNNs.



Layer Normalization

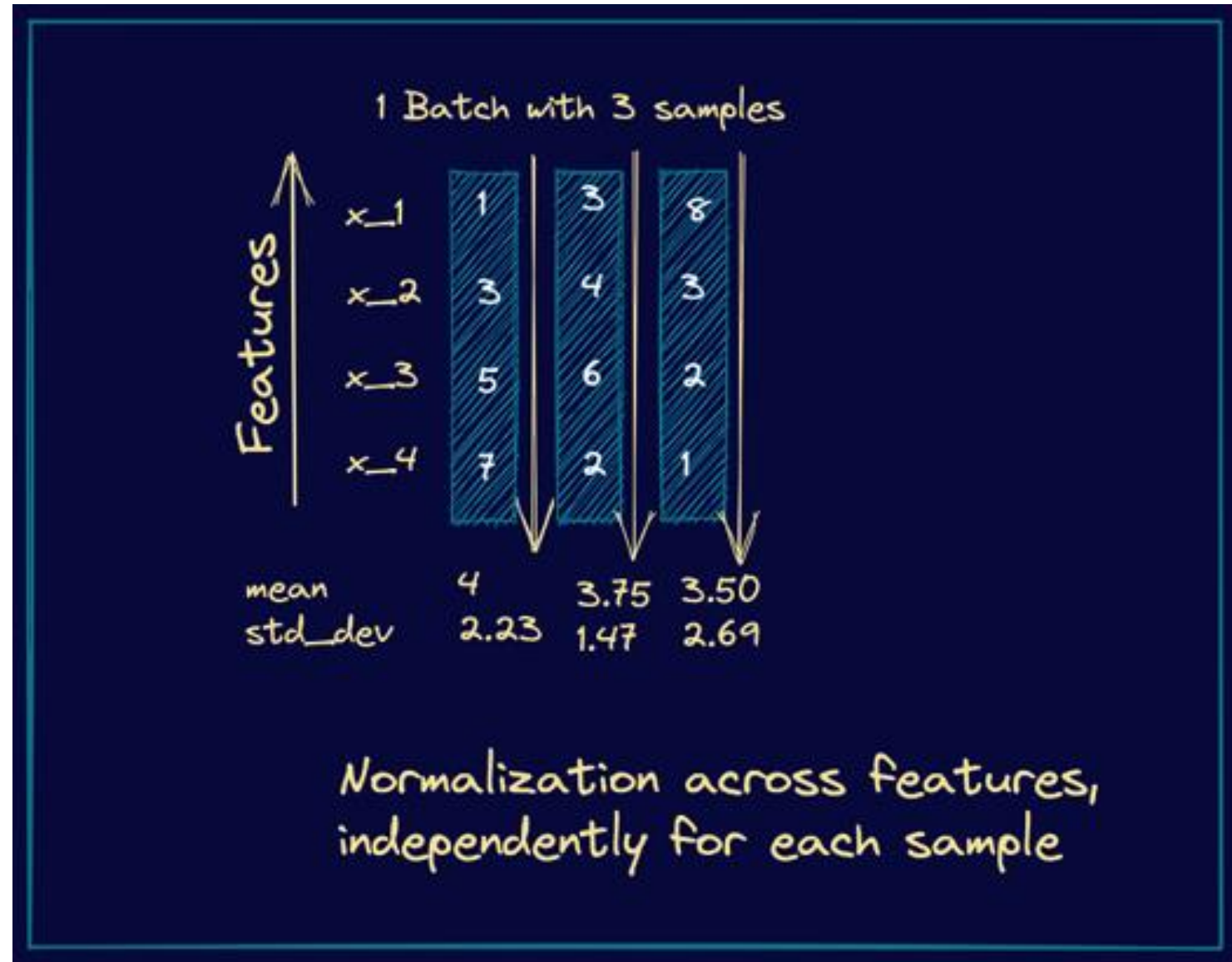


Image Source: <https://www.pinecone.io/learn/batch-layer-normalization/>



Layer Normalization

$$\mu_l = \frac{1}{d} \sum_{i=1}^d x_i \quad (1)$$

$$\sigma_l^2 = \frac{1}{d} \sum_{i=1}^d (x_i - \mu_l)^2 \quad (2)$$

$$\hat{x}_i = \frac{x_i - \mu_l}{\sqrt{\sigma_l^2}} \quad (3)$$

$$\text{or } \hat{x}_i = \frac{x_i - \mu_l}{\sqrt{\sigma_l^2 + \epsilon}} \quad (3)$$

Adding ϵ helps when σ_l^2 is small

$$y_i = \mathcal{LN}(x_i) = \gamma \cdot x_i + \beta \quad (4)$$

Trainable parameters

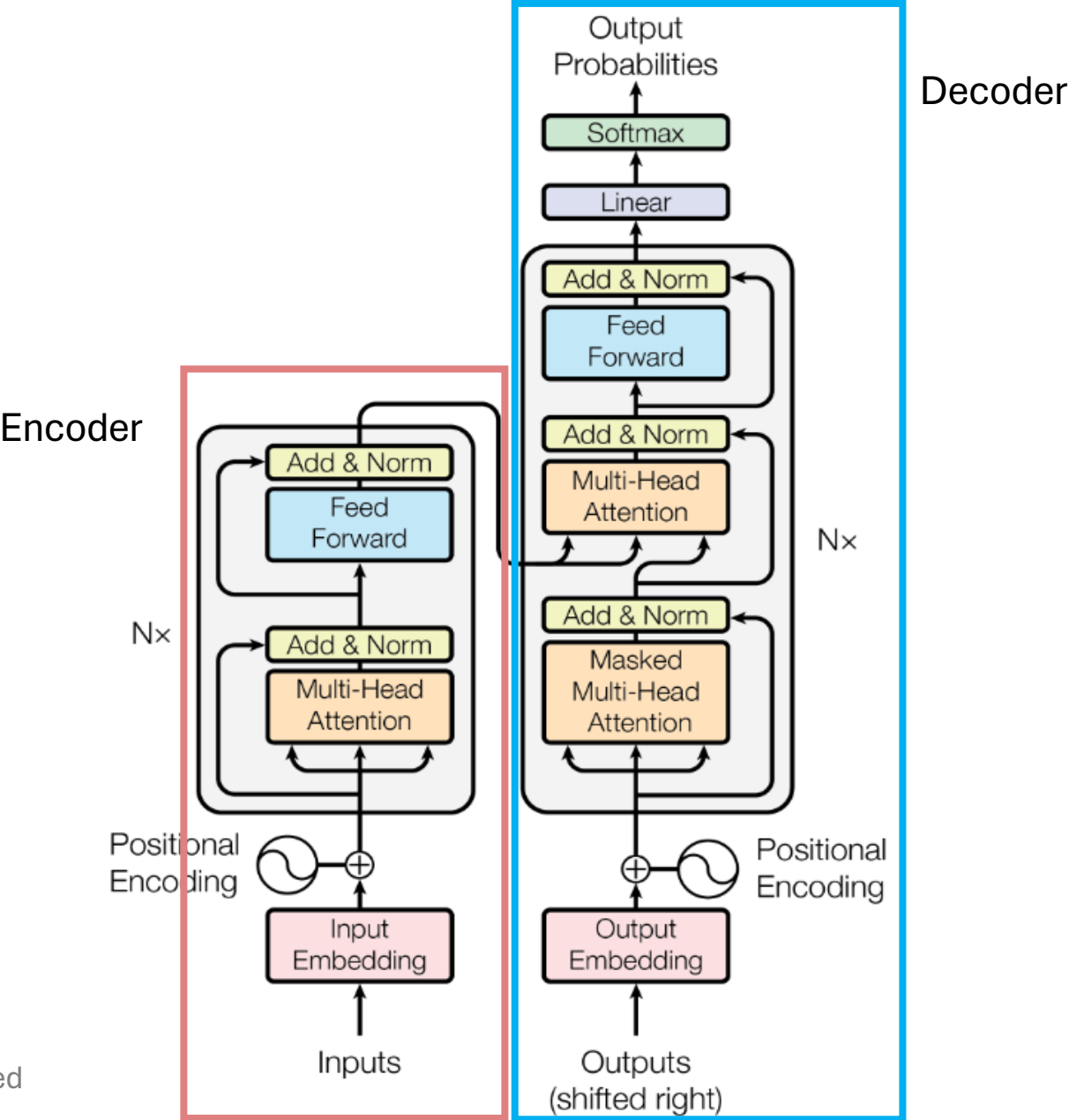


Batch Normalization vs Layer Normalization

- Batch normalization normalizes each feature independently across the mini-batch. Layer normalization normalizes each of the inputs in the batch independently across all features.
- As batch normalization is dependent on batch size, it's not effective for small batch sizes. Layer normalization is independent of the batch size, so it can be applied to batches with smaller sizes as well.
- Batch normalization requires different processing at training and inference times. As layer normalization is done along the length of input to a specific layer, the same set of operations can be used at both training and inference times.



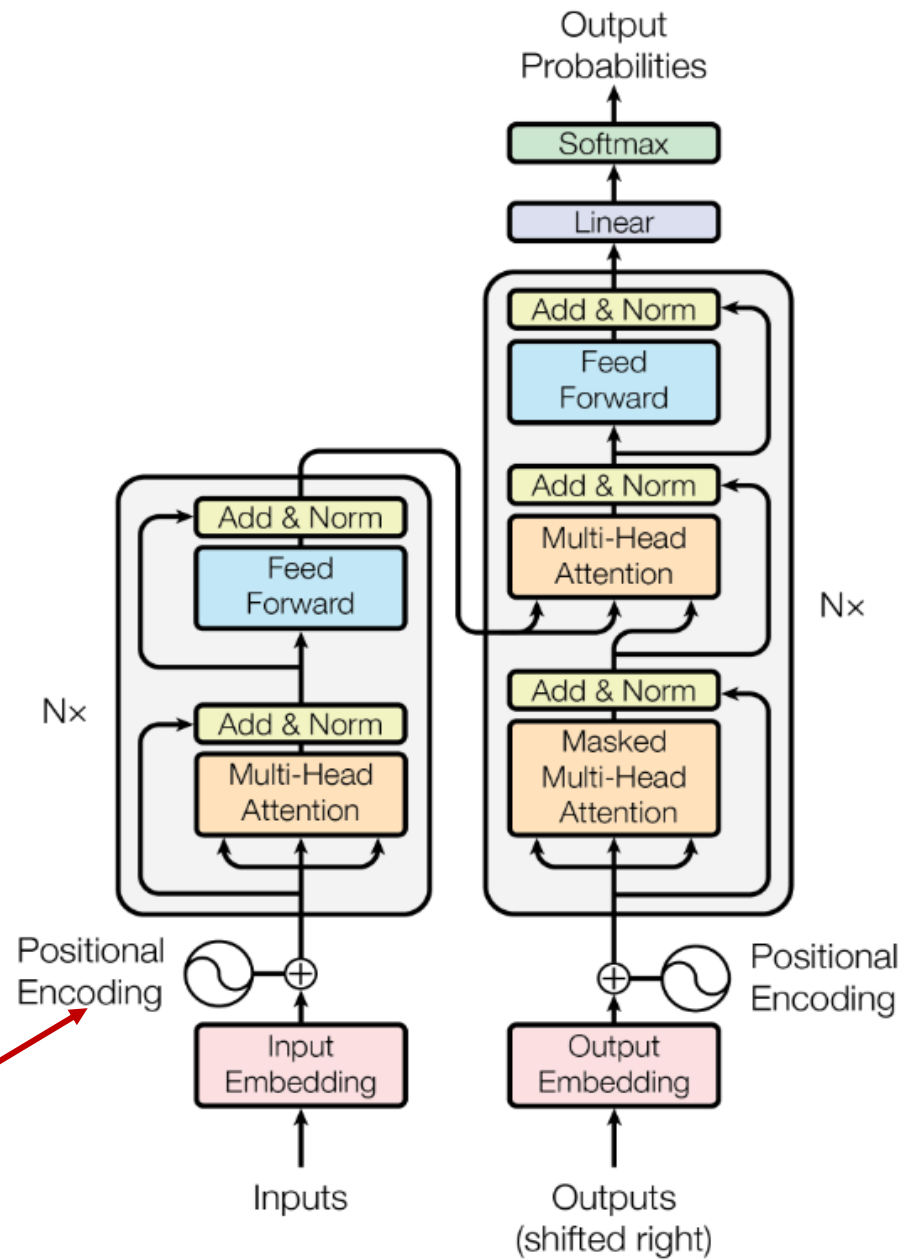
Transformer Architecture



Source of Image : Attention is all you need (Vaswani et al., 2017)



Transformer Architecture

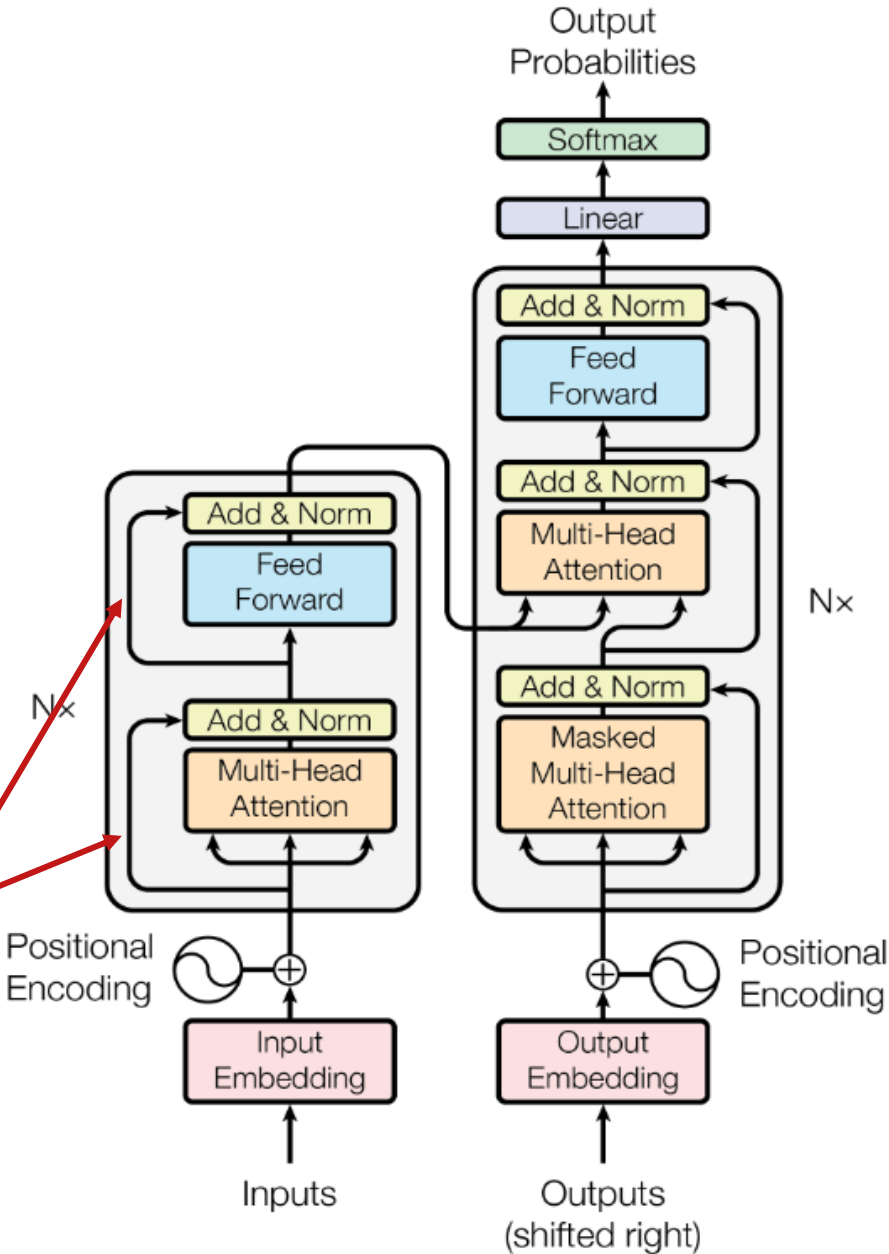


Position embeddings are *added* to each word embedding. Otherwise, since we have no recurrence, our model is unaware of the position of a word in the sequence!



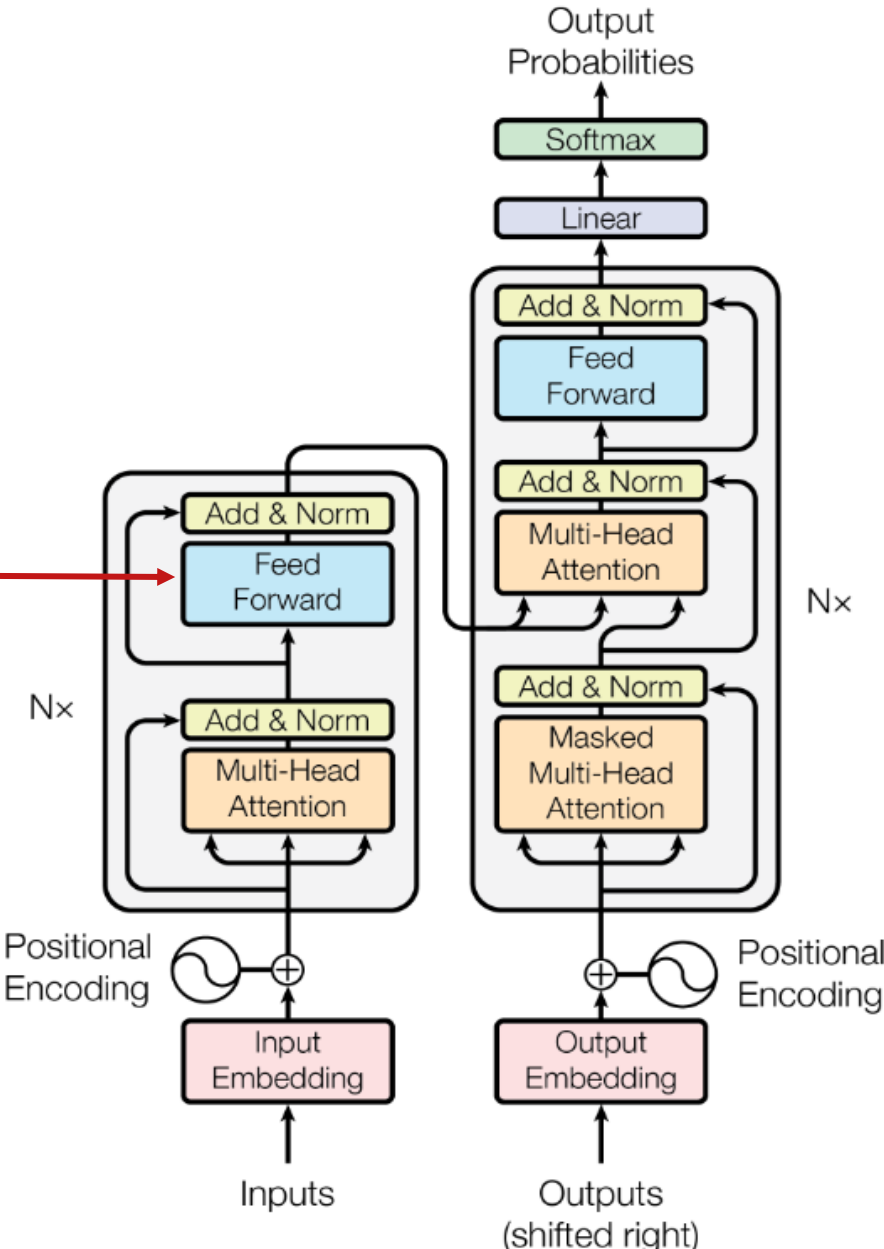
Transformer Architecture

Residual connections, which mean that we add the input to a particular block to its output, help improve gradient flow



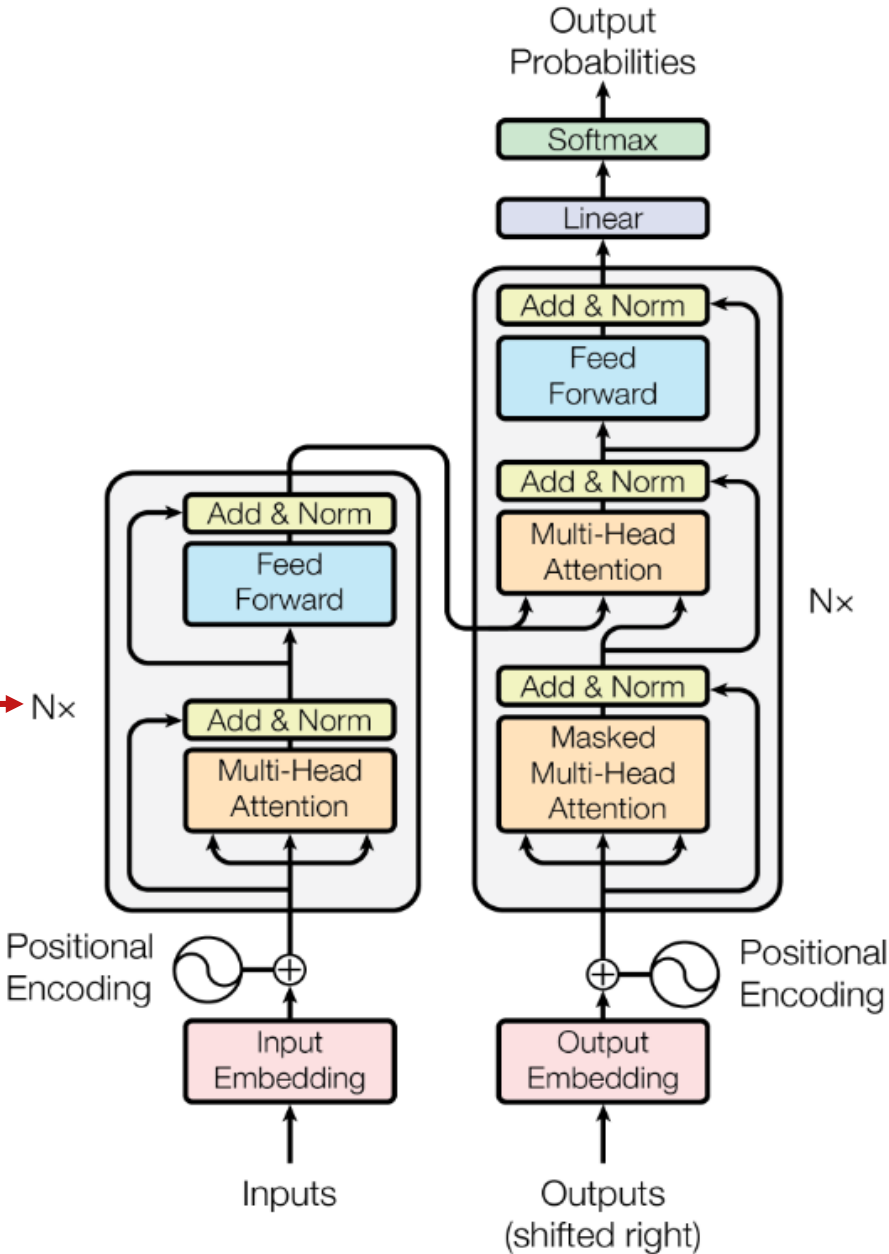
Transformer Architecture

A feed-forward layer on top of the attention-weighted averaged value vectors allows us to add more parameters / nonlinearity

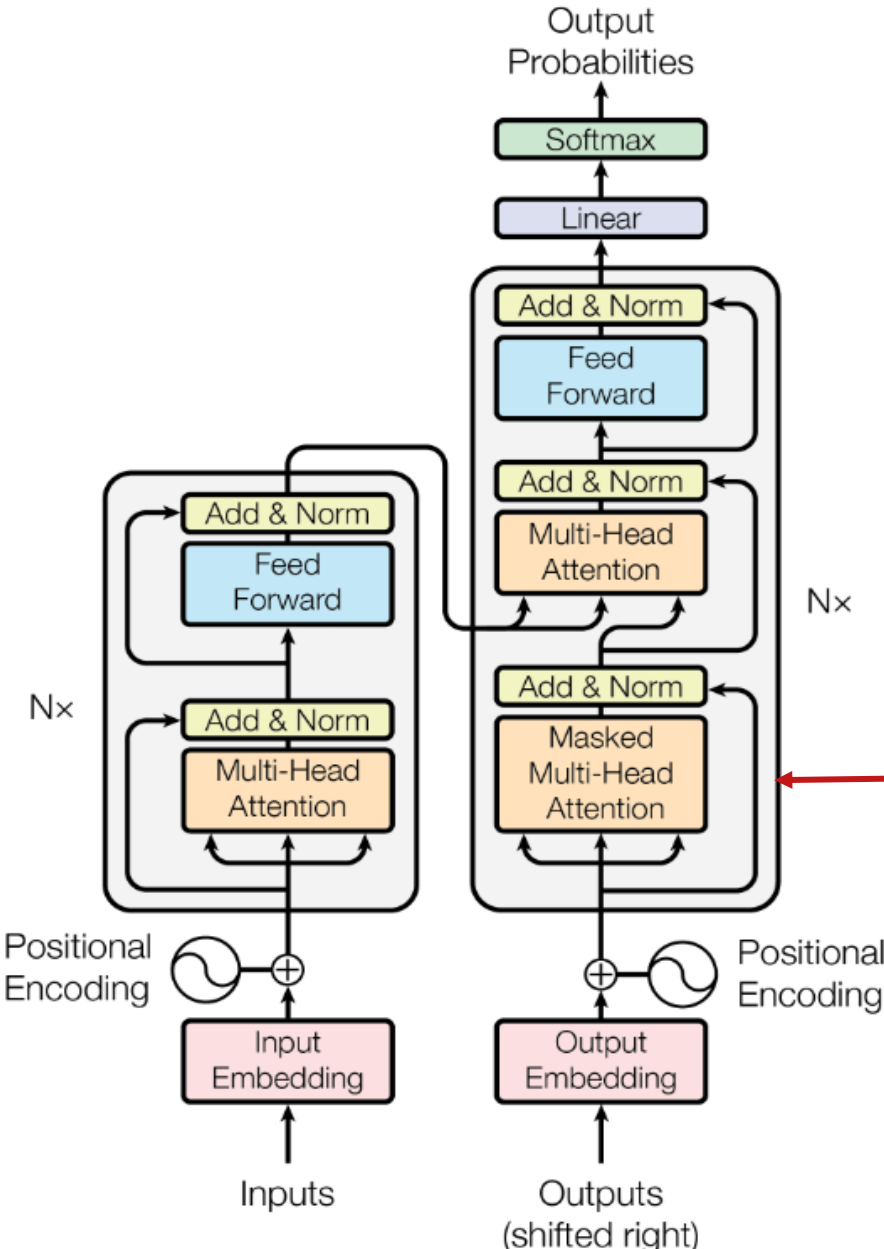


Transformer Architecture

We stack as many of these **Transformer blocks** on top of each other as we can (bigger models are generally better given enough data!)



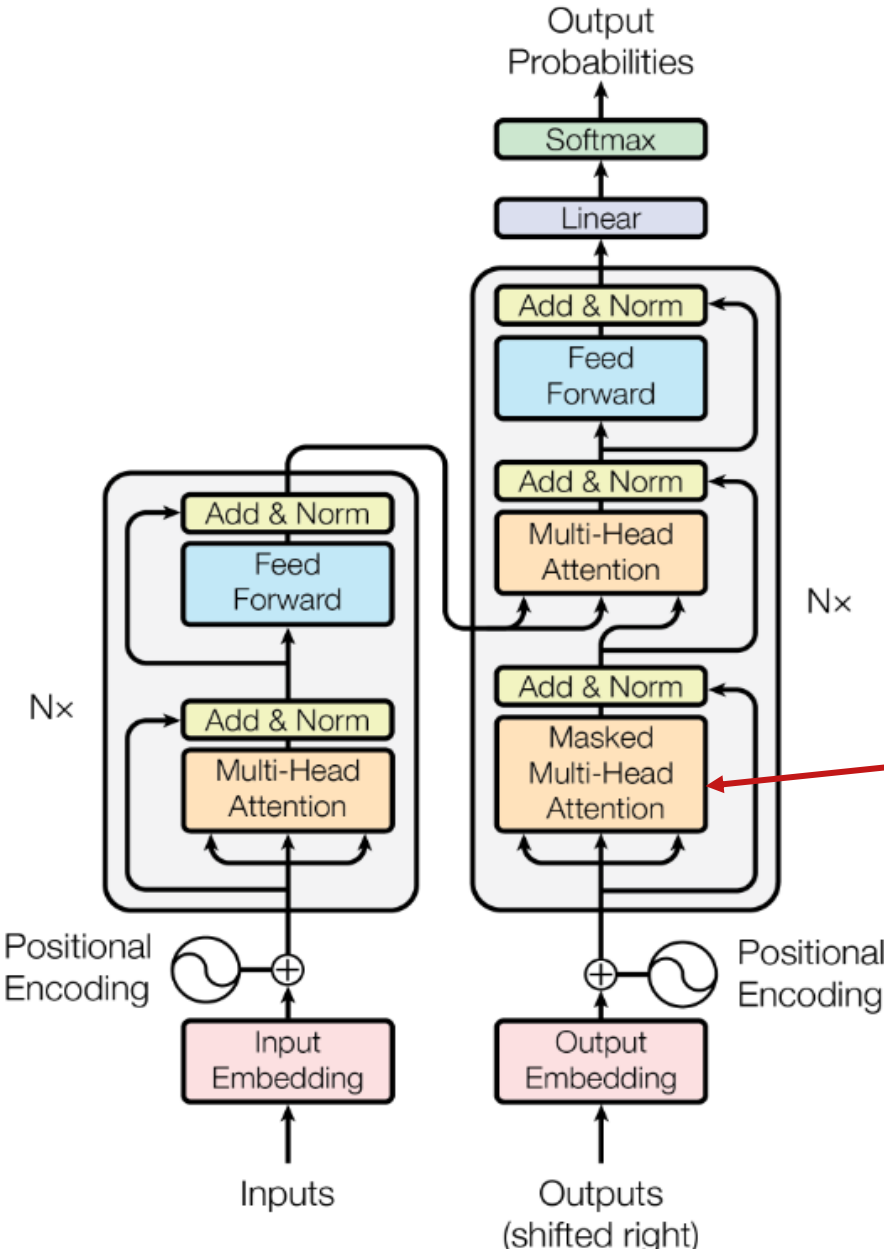
Transformer Architecture



Moving onto the decoder, which takes in English sequences that have been shifted to the right (e.g., *<START> schools opened their*)



Transformer Architecture

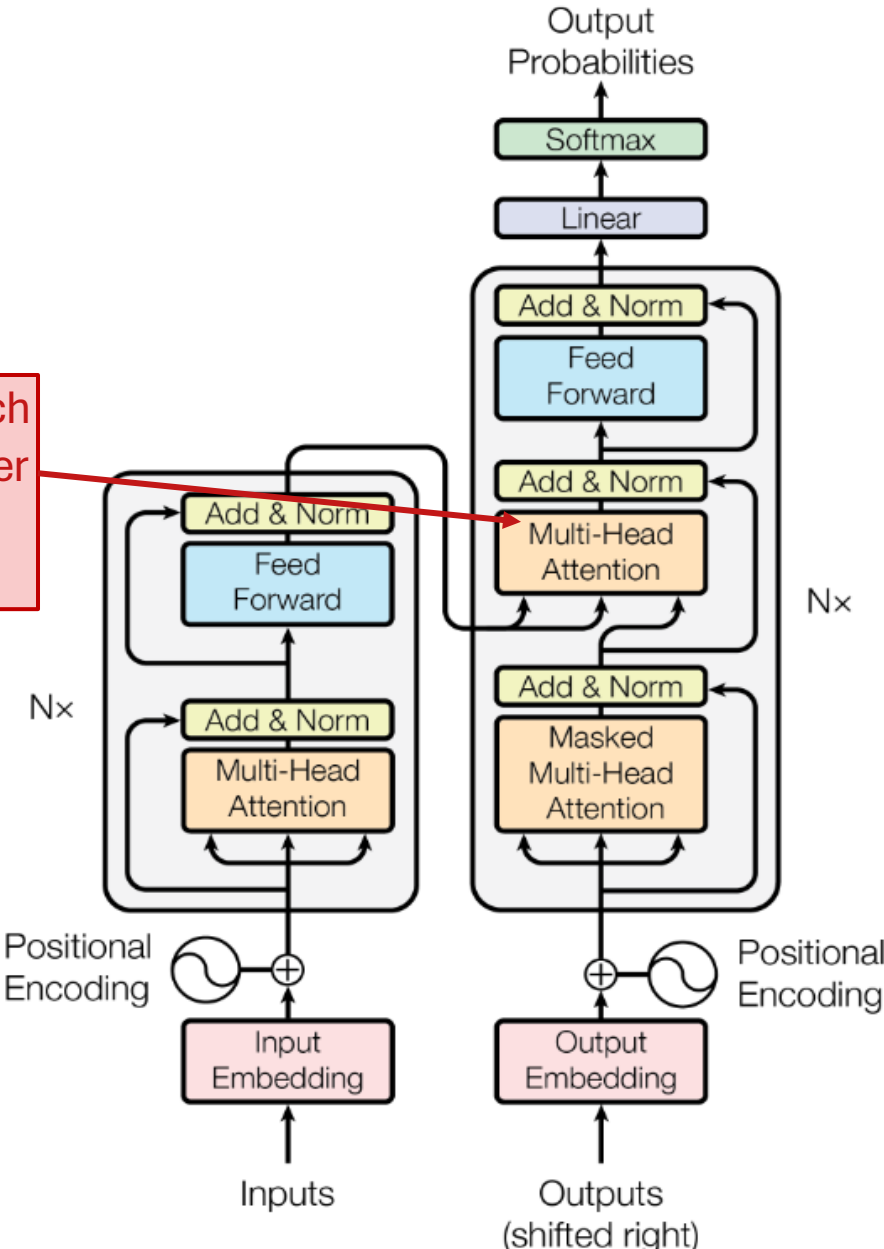


We first have an instance of **masked self attention**. Since the decoder is responsible for predicting the English words, we need to apply masking as we saw before.



Transformer Architecture

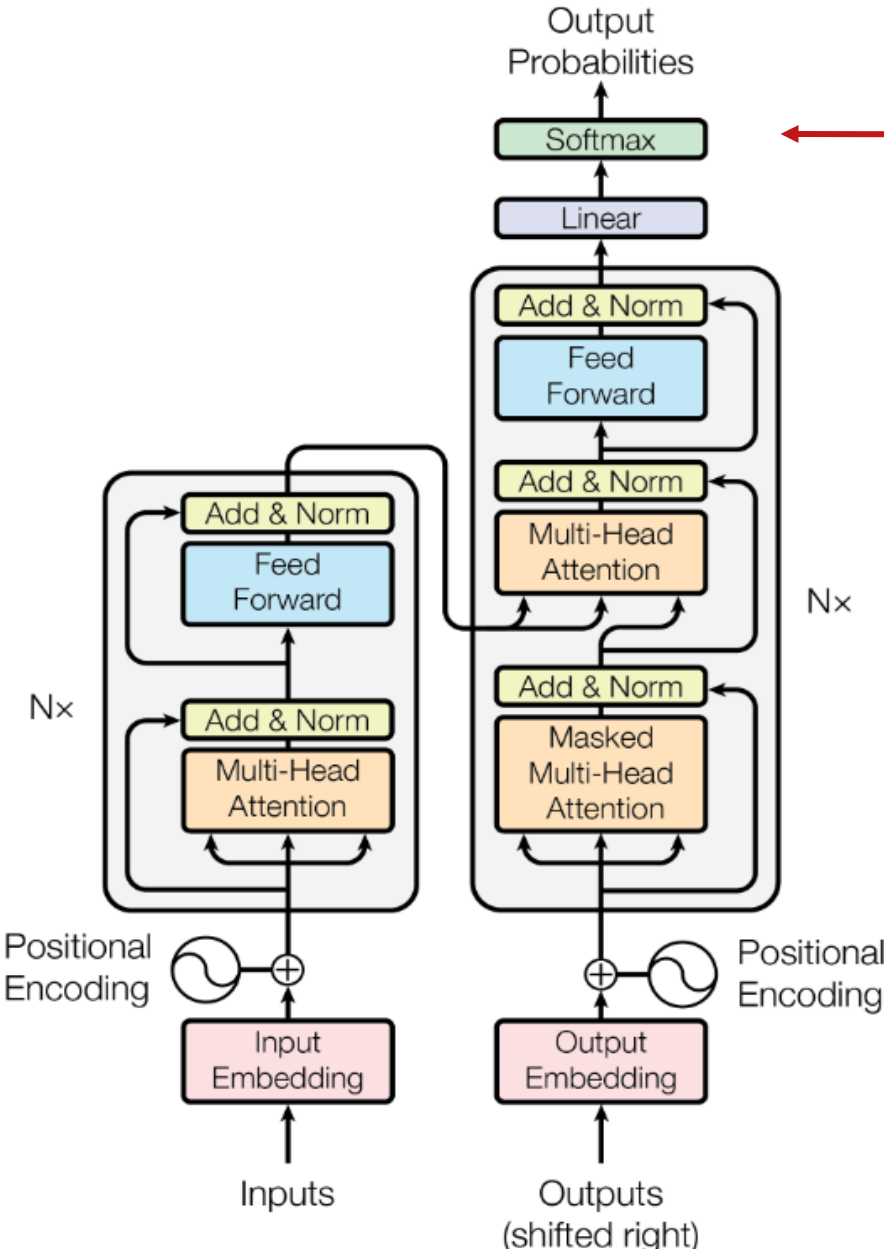
Now, we have **cross attention**, which connects the decoder to the encoder by enabling it to attend over the encoder's final hidden states.



Source of Image : Attention is all you need (Vaswani et al., 2017)



Transformer Architecture



After stacking a bunch of these decoder blocks, we finally have our familiar **softmax** layer to predict the next English word.

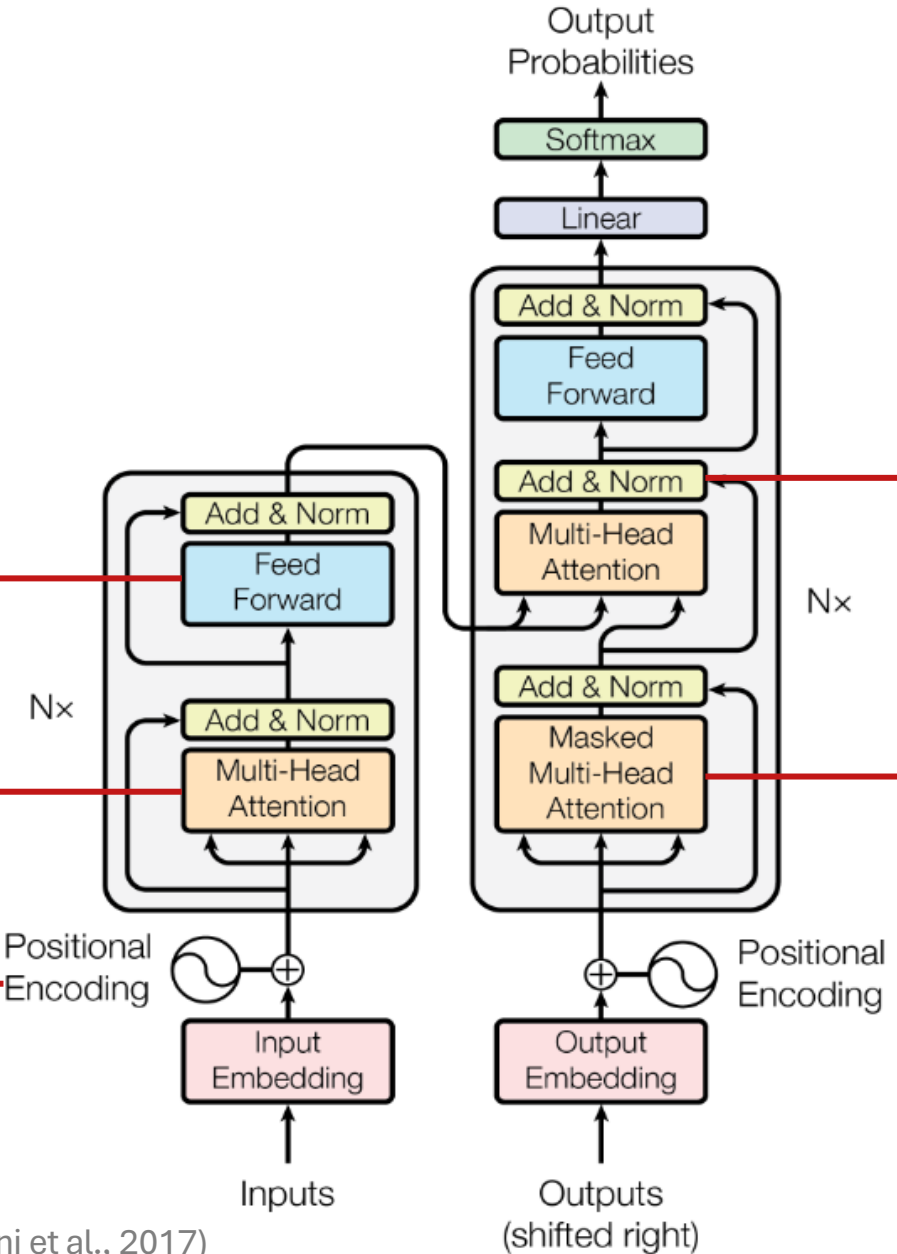


Transformer Architecture

Adding non-linearities

Allows querying multiple positions at each layer

Adds positional information



Reduces covariance shift and makes the system stable

Prevents attention lookups into the future while decoding

Source of Image : Attention is all you need (Vaswani et al., 2017)

