



# LLM Edge

after reviewing the fundamentals this note dives deep into some of the more specific knowledge gaps regarding LLMs

## THINK BEFORE YOU SPEAK:

### TRAINING LANGUAGE MODELS WITH PAUSE TOKENS

Sachin Goyal\*  
Machine Learning Department  
Carnegie Mellon University  
saching@andrew.cmu.edu

Ziwei Ji  
Google Research, NY  
ziwei.ji@google.com

Ankit Singh Rawat  
Google Research, NY  
ankitsrawat@google.com

Aditya Krishna Menon  
Google Research, NY  
adityakmenon@google.com

Sanjiv Kumar  
Google Research, NY  
sanjivk@google.com

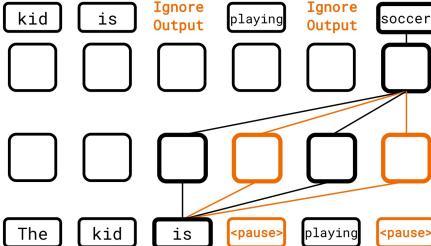
Vaishnavh Nagarajan†  
Google Research, NY  
vaishnavh@google.com

as for each  $(K+1)^{th}$  token the model uses  $K$  previous vectors, paper argues for some challenging tokens,  $K+M$  vectors might be needed.

even though this works, it's not known why :))

by appending dummy tokens to the input sequence, the generation of tokens is delayed.

during pre-training  
pauser inserted  
at random and  
the associated  
output is ignored  
for calculating  
the loss function.



(b) Pause-pretraining

the pause tokens don't give the model more "time" to think, but rather more computation to work with. pause tokens manipulate the model to do more computation than it normally does.

Note, that more tokens don't cause more activations or using more "notes". The computational flow in the paper means that more tokens provide more context or hidden vectors (token embeddings)

### Meet in the Middle: A New Pre-training Paradigm

Anh Nguyen\* Nikos Karampatziakis\* Weizhu Chen  
Microsoft Azure AI

regarding the unidirectional limitation of autoregressive models (referred on my notes about hallucination), this paper proposes a forward and a backward models that agree on a token in the middle. This is great for infill tasks where you don't generate text from scratch.

bidirectional models while better for infilling tasks, are less capable of in-context learning. regularizer to agree on some tokens

$$\sum_{x \in S} \sum_{i=1}^{|x|} -\log(\vec{p}(x_i | x_{<i})) - \log(\vec{p}(x_i | x_{>i})) + \beta D_{i,x}^{TV}(\vec{p} || \vec{p}).$$

during inference, both the forward and backward models generate sequences until they meet in the middle. In the best case they produce the same tokens. To avoid false positive (both agreeing on "the") n-gram is used and not just a single token.



# Infinite Context Window

how does Gemini manage a One million CW?

Leave No Context Behind:

Efficient Infinite Context Transformers with Infini-attention

Tsendsuren Munkhdalai, Manaal Faruqui and Siddharth Gopal  
Google  
tsendsuren@google.com

introduces memory of past context similar to how RNNs work.

## Why store K & V in the memory?

The K are important to finding the relevant V to a Q. Without a K, the model would not be able to retrieve.

there is nothing directly learnable about the memory matrix.

## Recap

the previous K and Vs are stored in a M matrix like [K|V]

now each Q<sub>s</sub> is compared to the Ks by dot product (or something) and through the resulting attention weights, related Vs are returned. This is actually a weighted sum of previous Vs.

This Could be key to Gemini's 1M Context Window!

so the CW isn't 1M in reality, but through a memory function, it effectively stores the previous 1M tokens.

an overview of how self-attention works



$K = XW_K$ ,  $V = XW_V$  and  $Q = XW_Q$ .  
 $V_V \in \mathbb{R}^{d_{model} \times d_{value}}$  and  $W_Q \in \mathbb{R}^{d_{model} \times d_k}$   
ion context is calculated as a weighted sum

$$A_{dot} = \text{softmax}\left(\frac{QK^T}{\sqrt{d_{model}}}\right)V.$$

$$A_{mem} = \frac{\sigma(Q)M_{s-1}}{\sigma(Q)z_{s-1}}.$$
 a piece of memory (A mem) is retrieved from the memory matrix M. σ is a non-lin activation

Once a memory is retrieved, the new memory is updated with new K and V. if the new info already exists, no update is required

$$M_s \leftarrow M_{s-1} + \sigma(K)^T(V - \underbrace{\frac{\sigma(K)M_{s-1}}{\sigma(K)z_{s-1}}}_{\text{a piece of memory}}).$$

the retrieved memory and the new A<sub>dot</sub> are aggregated through a learnable β weight.

$$A = \text{sigmoid}(\beta) \odot A_{mem} + (1 - \text{sigmoid}(\beta)) \odot A_{dot}.$$



# ROPE

## Rotary Positional Encoding

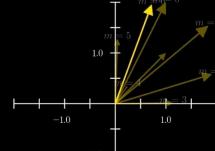
So what is wrong with the original sinusoidal PE?

It's Chaotic! The model has to memorize all the PEs

Rather than figure out a pattern.

	$i = 0$	$i = 1$
$m = 0$	0.0	1.0
$m = 1$	0.88	0.61
$m = 2$	0.94	-0.46
$m = 3$	0.17	-0.99
$m = 4$	-0.81	-0.73
$m = 5$	-0.97	0.3
$m = 6$	-0.31	0.97
$m = 7$	0.73	0.82
$m = 8$	0.99	-0.13

$$q, k_m = W_{\{q,k\}} \cdot (x_m + PE(m))$$

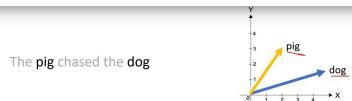


Sinusoidal is an absolute PE, meaning that position 1 vs. 87 have different PE relations than 2 vs. 88.

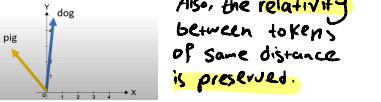
With **relative Positional embeddings** the PE between two tokens of relatively same position is the same.

the **RoPE** method uses rotations to encode the position of each token.

### Rotary Positional Embeddings



Once upon a time, the pig chased the dog



the formula breaks the dims to chunks of two dimensions,

$$\left( \begin{array}{ccccccccc} \cos \theta_1 & -\sin \theta_1 & 0 & 0 & \cdots & 0 & 0 \\ \sin \theta_1 & \cos \theta_1 & 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \cos \theta_2 & -\sin \theta_2 & \cdots & 0 & 0 \\ 0 & 0 & \sin \theta_2 & \cos \theta_2 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ & & & & & \cos \theta_{d/2} & -\sin \theta_{d/2} \\ & & & & & \sin \theta_{d/2} & \cos \theta_{d/2} \end{array} \right)$$

Sparse matrix  
and Computationally  
efficient

### RoFORMER: ENHANCED TRANSFORMER WITH ROTARY POSITION EMBEDDING

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

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

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

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

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

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

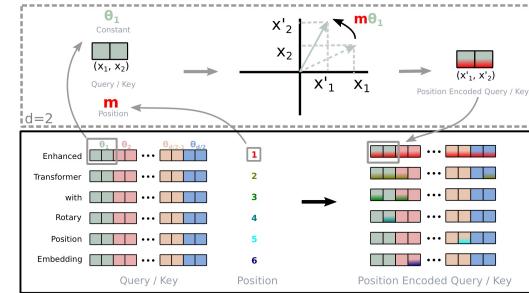
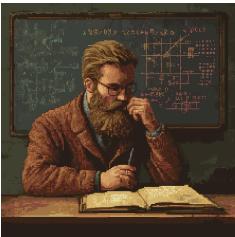
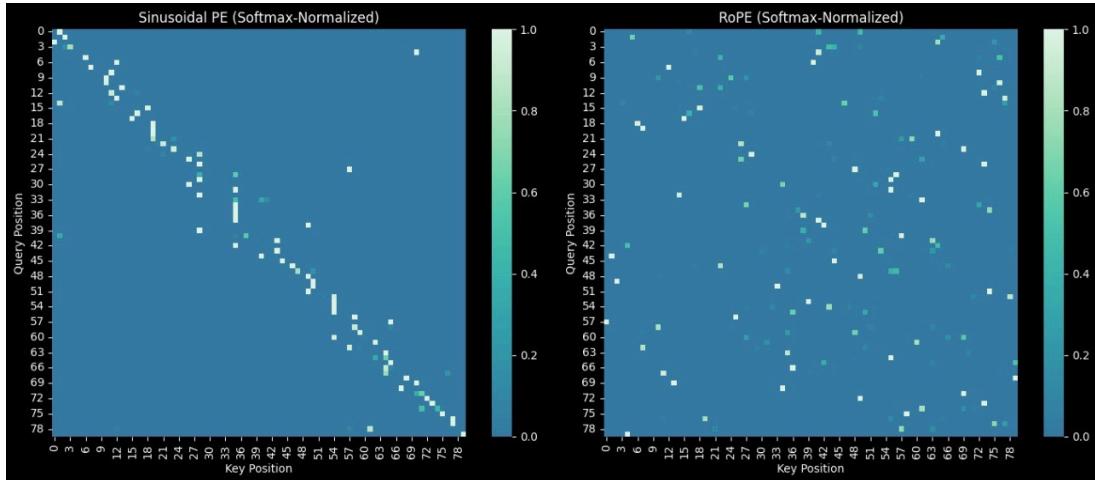


Figure 1: Implementation of Rotary Position Embedding(RoPE).



Let's see the practical difference of ROPE vs Sinusoidal. The code is in my repo of "ml-retreat".



now we can see clearly that sinusoidal puts more emphasis on nearby tokens while RoPE captures long-range dependencies. Why? The Sinusoidal, being an absolute PE method assigns a PE Score to tokens based on their positions. It puts more emphasis on closer tokens because it understands local context where the absolute positional difference is large enough. ROPE has the main idea to capture how far tokens are from each other rather than their specific absolute positions.

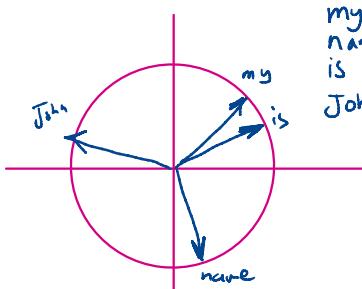
There's a big difference in how they're applied

Sens: Token Emb + Pos Emb

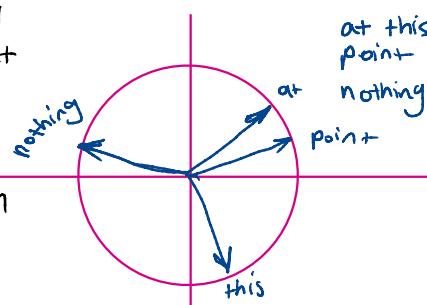
RoPE: Rotated Q, Rotated K

- ① RoPE doesn't directly modify token embs
- ② because of that, the V is not modified

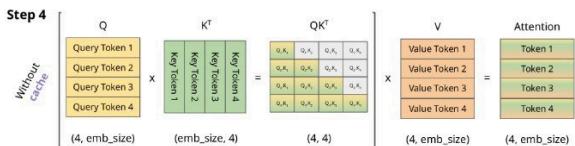
3 RoPE is based on a **Rotation Transformation** on the Q, K. So regardless of position, the positional Encoding **does not decay**. The transformation also does not change size of the token embeddings. Unlike Sinusoidal, RoPE **doesn't mess with the embeddings**.



the model learns that each degree relates to a position



# KV Cache



Comparison of scaled dot-product attention with and without KV caching. emb\_size means embedding size.  
Image created by the author.

during the process of calculating attention scores, many operations are repeated. KV cache stores some intermediate values and eliminate the need to perform every calculations again.

## Optimizations

$K = W_k \cdot X$  is optimized as the  $K$  of previous tokens is stored and concatenated to new  $K$ .

$$K_{\text{new}} = W_k \cdot X_{\text{new}}$$

$$K = \text{concat}(K_{\text{cache}}, K_{\text{new}})$$

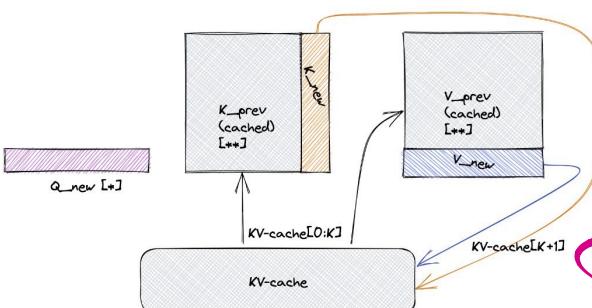
$V = W_v \cdot X$  is optimized the same way

$$V_{\text{new}} = W_v \cdot X_{\text{new}}$$

$V$  is  $V_{\text{cache}}$  and  $V_{\text{new}}$  concatenated

Scaled dot Product is performed for the new token only.

$$\text{Attention Score} = Q_{\text{new}} \times K_{\text{cache + new}}^T$$



# Mixture of Experts (MoE)

We have multiple smaller models each expert in a set of tasks.

Also the Gating function choosing the expert for each input

\* it's not an entirely novel idea

## Mixtral of Experts

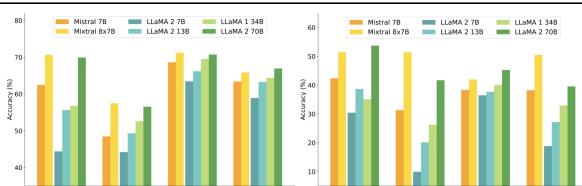
Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lampe, Lélio Renard Lavau, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang, Timothée Lacroix, William El Sayed

# MistralAI

the output is a weighted sum of all expert outputs. If  $G(n)$  is sparse, only one or few experts are needed for computation.

$G(n)$  takes input, applies a simple classification via a FF layer,  $G(x) := \text{Softmax}(\text{TopK}(x \cdot W_g))$ , Select the TopK logits and applies Softmax to normalize.

normal feed forward layer



it can outperform vanilla models of bigger size.

the assignment to experts is per token not input meaning that a sequence could be handed over multiple experts.

Also, the authors didn't notice any patterns of expert assignments on a topic.

This Wraps LLM Edge

