

ACL2022 Findings Meta

Transformer Language Models without Positional Encodings Still Learn Positional Information

Adi Haviv^τ Ori Ram^τ Ofir Press^ω Peter Izsak^ℓ Omer Levy^{τμ}

^τTel Aviv University ^ωUniversity of Washington ^ℓIntel Labs ^μMeta AI

ACL2023 Honorable Mentions CMU&Princeton

Latent Positional Information is in the Self-Attention Variance of Transformer Language Models Without Positional Embeddings

Ta-Chung Chi[†]
Carnegie Mellon University

Ting-Han Fan
Princeton University

Li-Wei Chen
Carnegie Mellon University

Alexander I. Rudnický
Carnegie Mellon University

Peter J. Ramadge
Princeton University

Preprint Mila

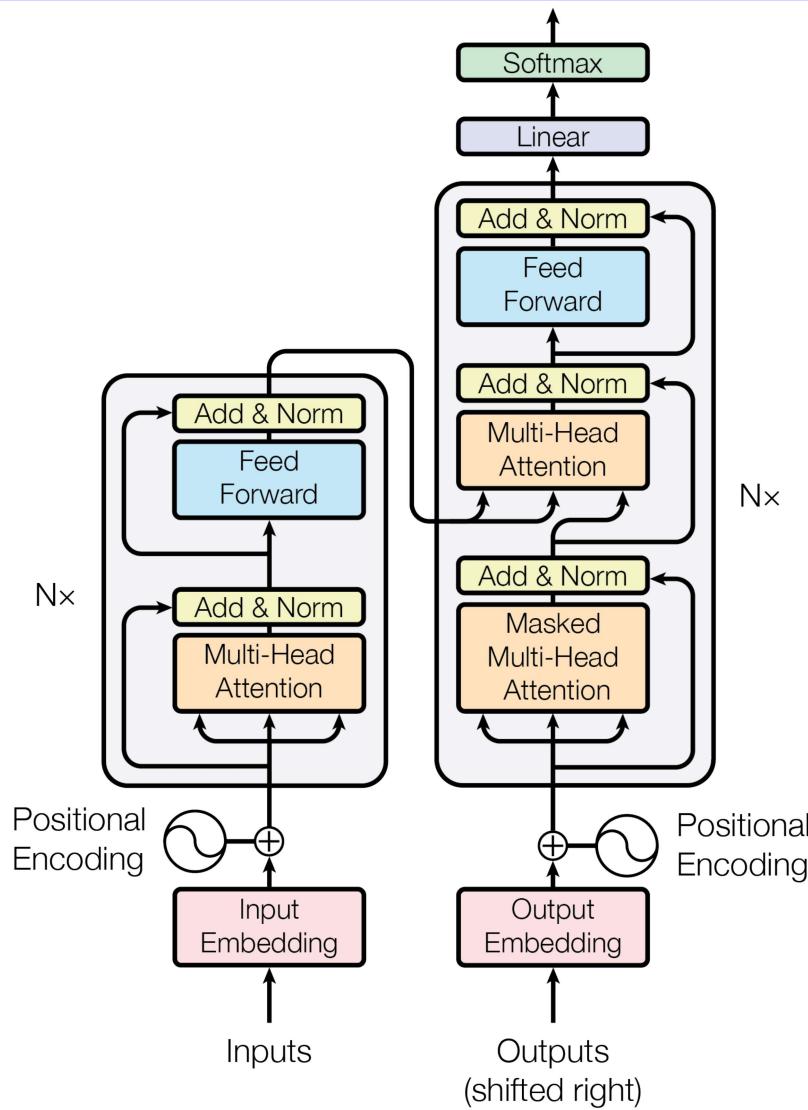
The Impact of Positional Encoding on Length Generalization in Transformers

Amirhossein Kazemnejad^{1,2}, Inkit Padhi³
Karthikeyan Natesan Ramamurthy³, Payel Das³, Siva Reddy^{1,2,4}

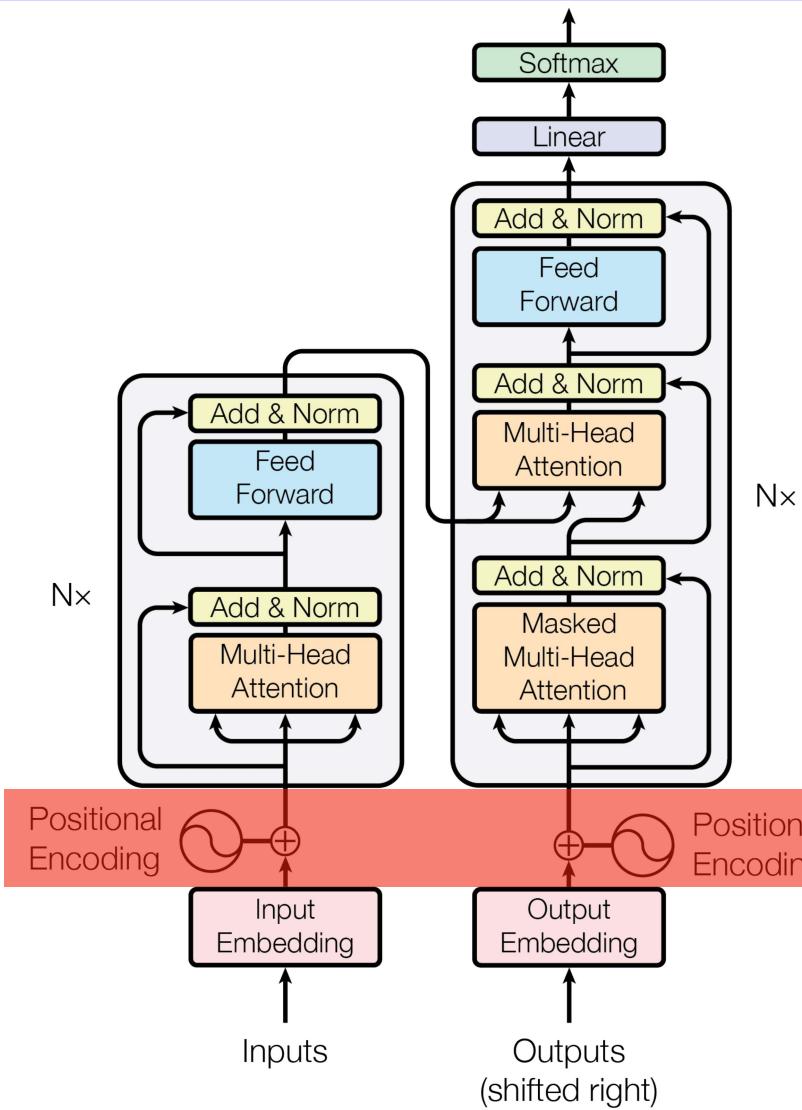
¹Mila - Québec AI Institute; ²McGill University;

³IBM Research; ⁴Facebook CIFAR AI Chair

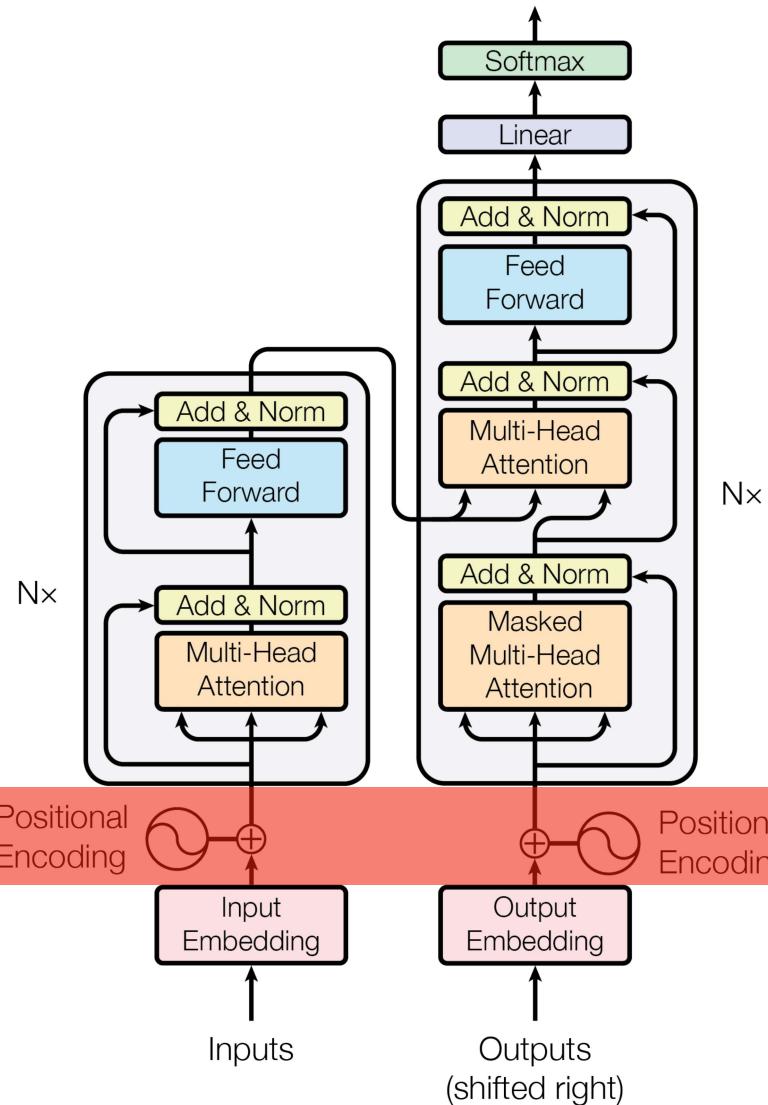
BG: PE in Transformer



BG: PE in Transformer



BG: PE in Transformer



Position Information in Transformers

纪 烹

Topic	Absolute	Reference Point Absolute & Relative	Relative
Sequential	(Devlin et al., 2019) (Kitayev et al., 2020) (Liu et al., 2020b) (Press et al., 2020) (Wang et al., 2020) (Dehghani et al., 2019)	(Shaw et al., 2018) (Ke et al., 2020) (Dufner et al., 2020) (He et al., 2020)	(Dai et al., 2019) (Raffel et al., 2020) (Wu et al., 2020) (Huang et al., 2020) (Shen et al., 2018) (Neishi & Yoshinaga, 2019)
Sinusoidal	(Vaswani et al., 2017) (Li et al., 2019)		(Yan et al., 2019)
Graphs		(Shiv & Quirk, 2019) (Wang et al., 2019)	(Zhu et al., 2019) (Cai & Lam, 2020) (Schmitt et al., 2021)
Decoder		(Takase & Okazaki, 2019) (Oka et al., 2020) (Bao et al., 2019)	
Crosslingual		(Artetxe et al., 2020) (Ding et al., 2020) (Liu et al., 2020a) (Liu et al., 2020c)	
Analysis	(Yang et al., 2019) (Wang & Chen, 2020)	(Rosendahl et al., 2019) (Wang et al., 2021)	

► [Position Information in Transformers: An Overview](#)

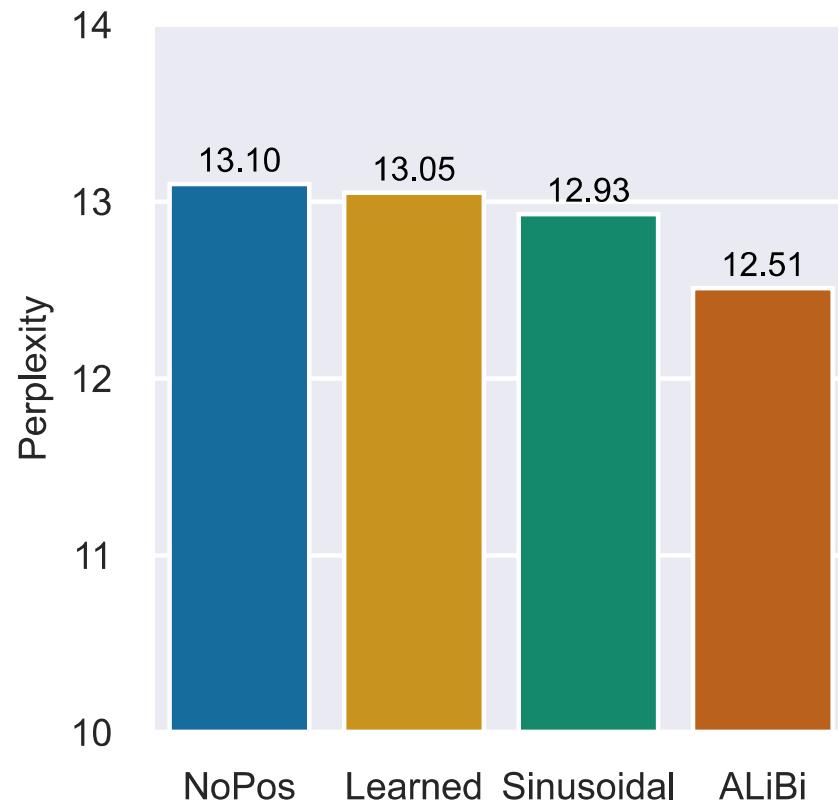
ACL2022 Findings

Transformer Language Models without Positional Encodings Still Learn Positional Information

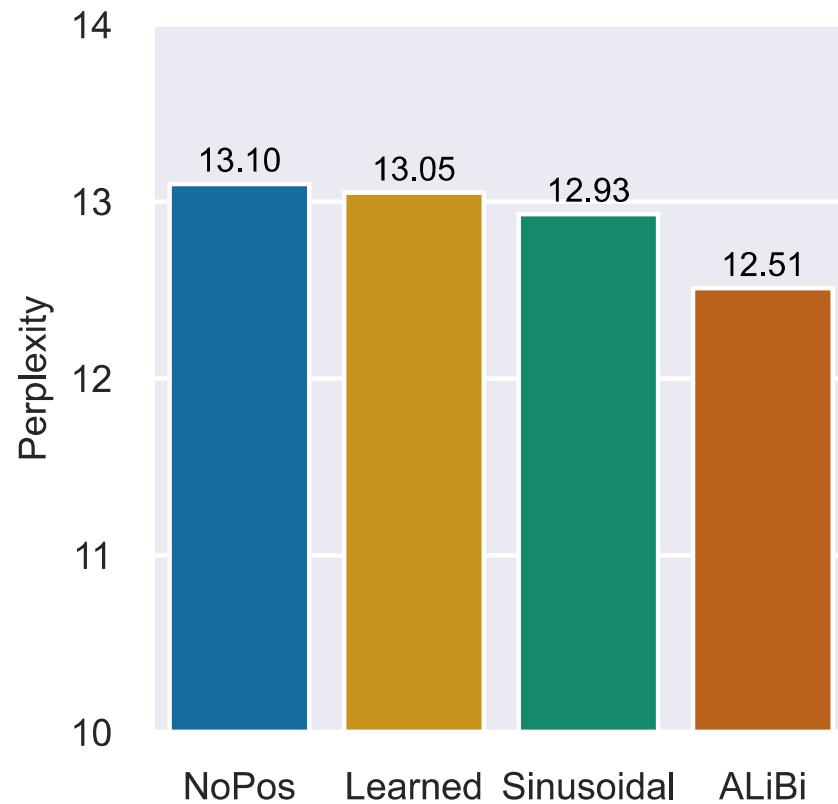
Adi Haviv^τ Ori Ram^τ Ofir Press^ω Peter Izsak^ι Omer Levy^{τμ}

^τTel Aviv University ^ωUniversity of Washington ^ιIntel Labs ^μMeta AI

1.3B GPT@ Pile



1.3B GPT@ Pile



😱 NoPE 🤝

1.3B GPT@wiki103 & Pile

	WikiText-103	The Pile
NoPos	20.97	13.10
Learned	20.42	13.05
Sinusoidal	20.16	12.93
ALiBi	19.71	12.51

1.3B GPT@wiki103 & Pile

	WikiText-103	The Pile
NoPos	20.97	13.10
Learned	20.42	13.05
Sinusoidal	20.16	12.93
ALiBi	19.71	12.51

2 datasets → NoPE 🤓

GPT@ Pile

Model Size	125M	350M	760M	1.3B
NoPos	22.15	16.87	14.29	13.10
Learned	22.04	16.84	14.21	13.05
Sinusoidal	21.49	16.58	14.04	12.93
ALiBi	19.94	15.66	13.53	12.51

Seq Length	256	512	1024	2048
NoPos	14.98	13.82	13.10	12.87
Learned	14.94	13.77	13.05	12.72
Sinusoidal	14.84	13.66	12.93	12.62
ALiBi	14.65	13.37	12.51	12.06

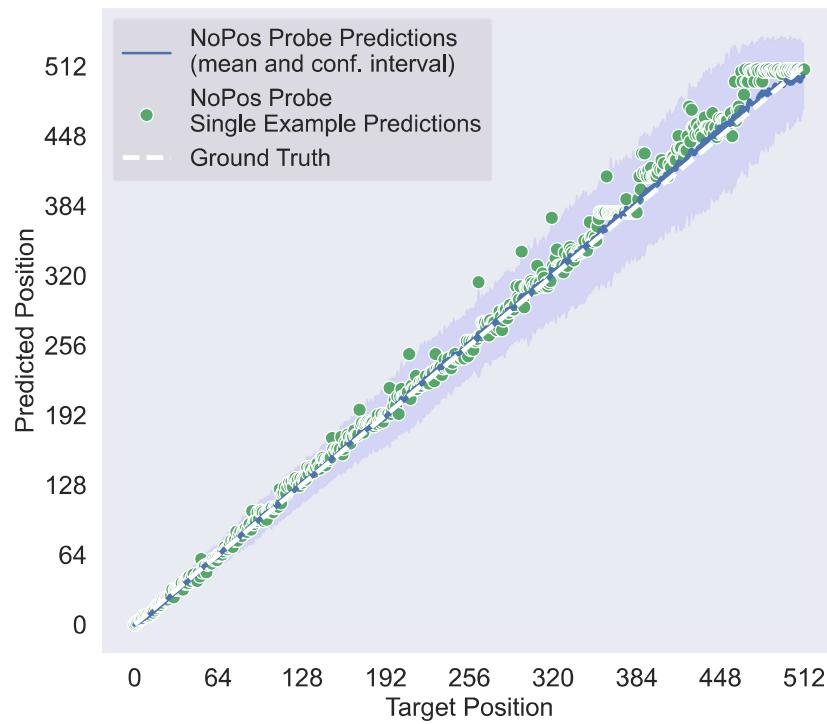
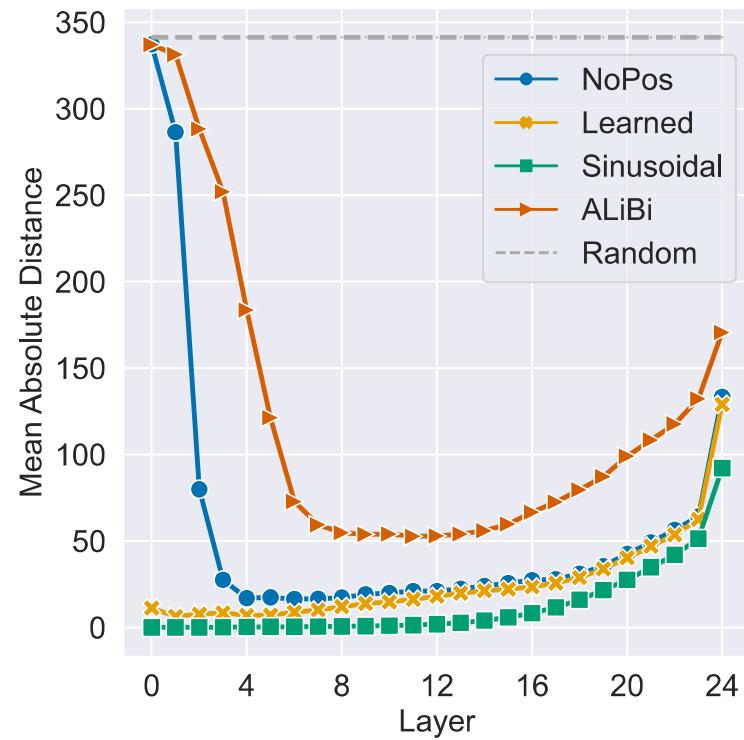
GPT@ Pile

Model Size	125M	350M	760M	1.3B
NoPos	22.15	16.87	14.29	13.10
Learned	22.04	16.84	14.21	13.05
Sinusoidal	21.49	16.58	14.04	12.93
ALiBi	19.94	15.66	13.53	12.51

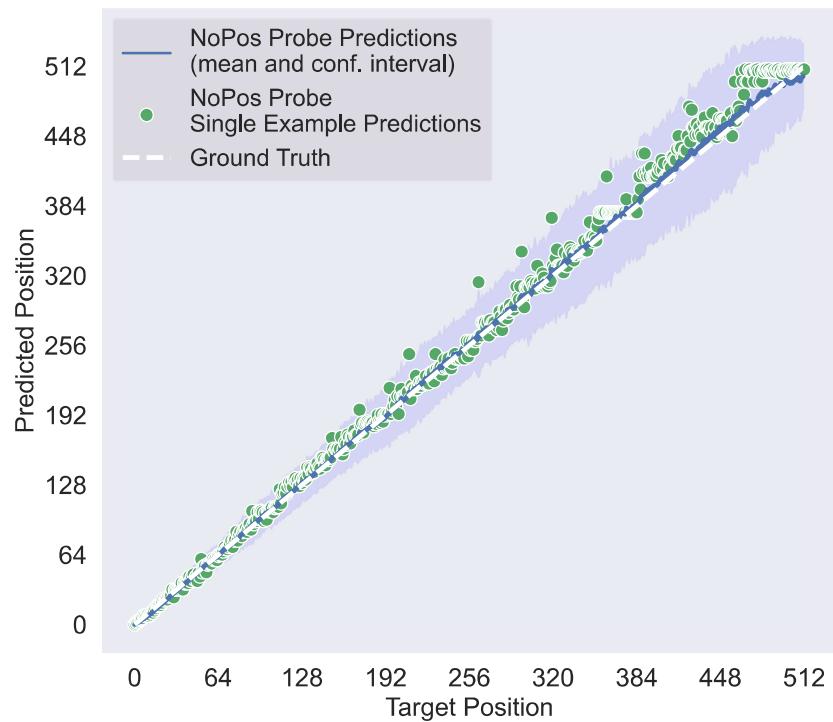
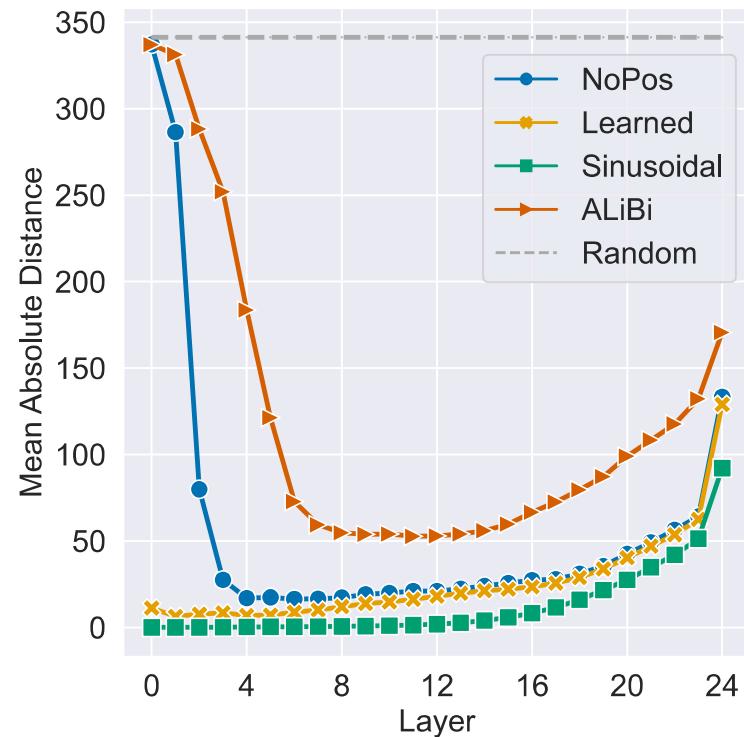
Seq Length	256	512	1024	2048
NoPos	14.98	13.82	13.10	12.87
Learned	14.94	13.77	13.05	12.72
Sinusoidal	14.84	13.66	12.93	12.62
ALiBi	14.65	13.37	12.51	12.06

Various size👉 &👉 length → NoPE 🤝

POS Probing



POS Probing



NoPos: POS info in hidden too

How about BERT-style

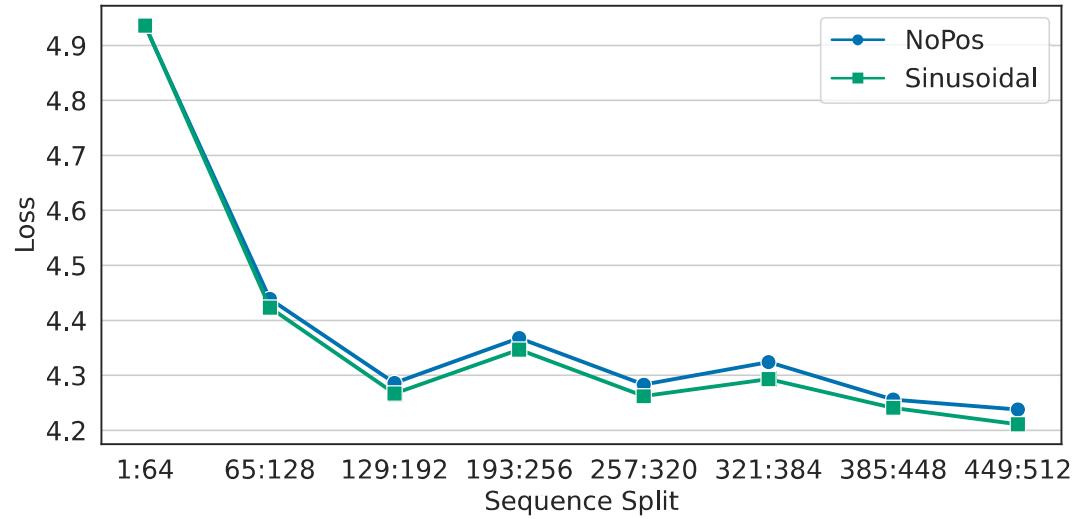
MLM Perplexity	
NoPos	147.18
Learned	4.06
Sinusoidal	4.07
ALiBi	4.00

How about BERT-style

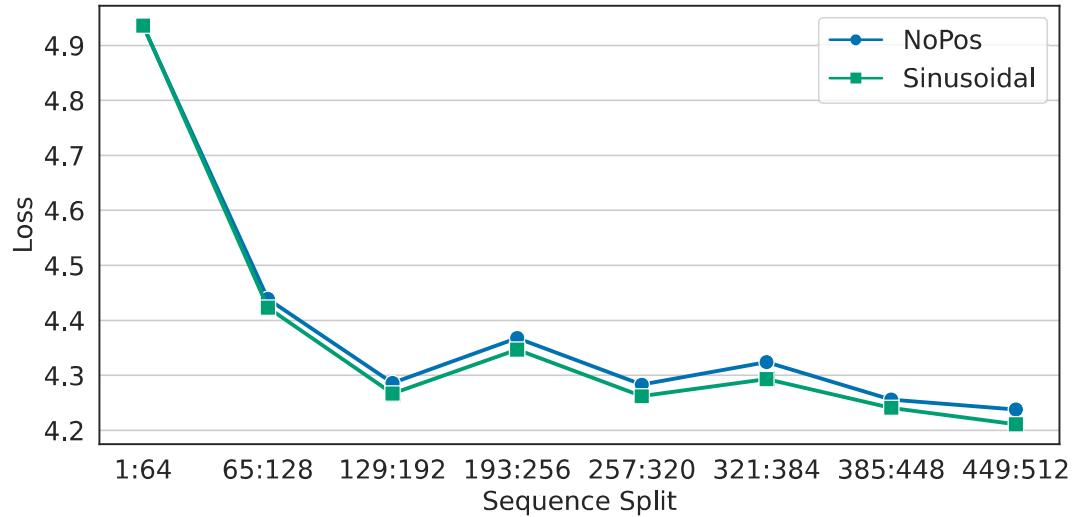
MLM Perplexity	
NoPos	147.18
Learned	4.06
Sinusoidal	4.07
ALiBi	4.00

GPT-style  BERT-style  T5-style ? 

Sequence split

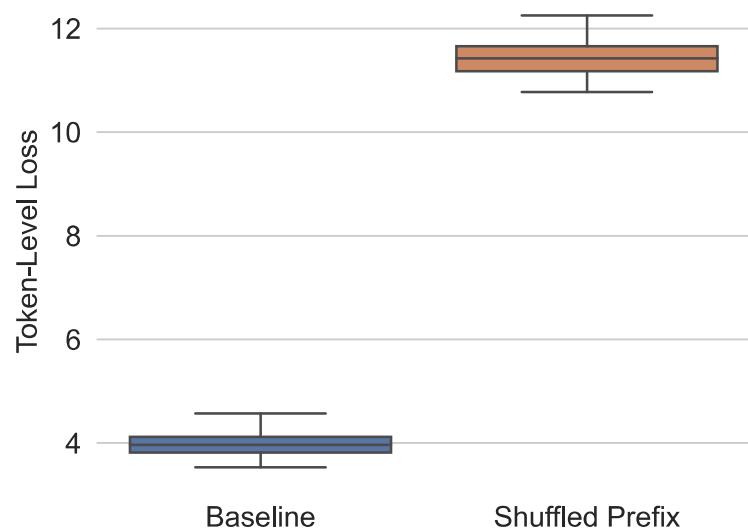


Sequence split



Slightly Worse

Shuffled Prefix



B Word Order Analysis

Is positional information necessary for language modeling, or does the order of the input tokens not matter? To answer this, we conduct the following experiment: instead of computing the loss on the complete sequence, we pick a specific token in the sequence. The next token prediction is conditioned on the previous tokens in the sequence, and so we shuffle the order of the tokens in the prefix and compute the loss only for that specific token. We repeat the experiment with the original, un-shuffled prefix sequence as the baseline and compare the results.

ACL2022 Findings

Transformer Language Models without Positional Encodings Still Learn Positional Information

Adi Haviv^τ Ori Ram^τ Ofir Press^ω Peter Izsak^ι Omer Levy^{τμ}

^τTel Aviv University ^ωUniversity of Washington ^ιIntel Labs ^μMeta AI

ACL2022 Findings

**Transformer Language Models without Positional Encodings
Still Learn Positional Information**

Adi Haviv^τ Ori Ram^τ Ofir Press^ω Peter Izsak^ι Omer Levy^{τμ}

^τTel Aviv University ^ωUniversity of Washington ^ιIntel Labs ^μMeta AI

😲 NoPE 🤝 But why? 🤔

ACL2023 Honorable Mentions

Latent Positional Information is in the Self-Attention Variance of Transformer Language Models Without Positional Embeddings

Ta-Chung Chi[†]

Carnegie Mellon University

Ting-Han Fan

Princeton University

Li-Wei Chen

Carnegie Mellon University

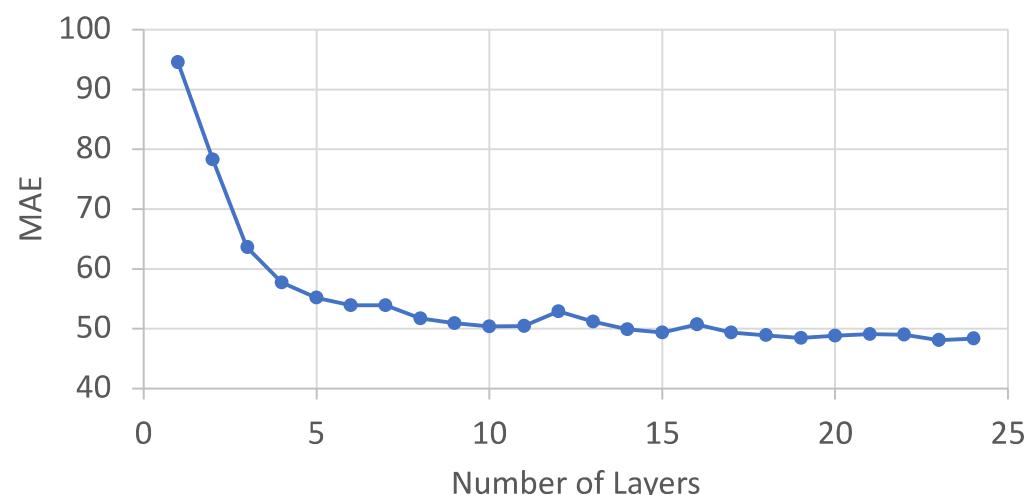
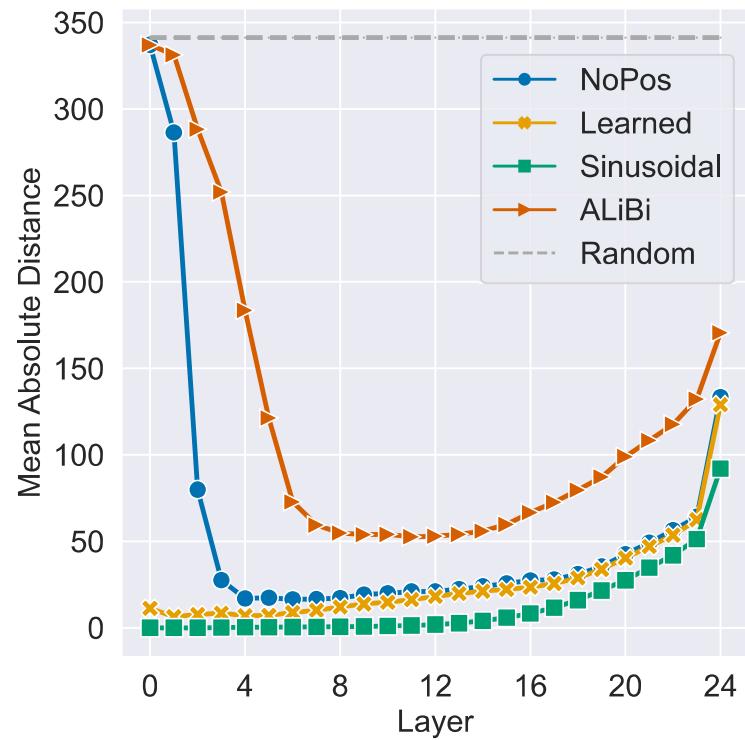
Alexander I. Rudnicky

Carnegie Mellon University

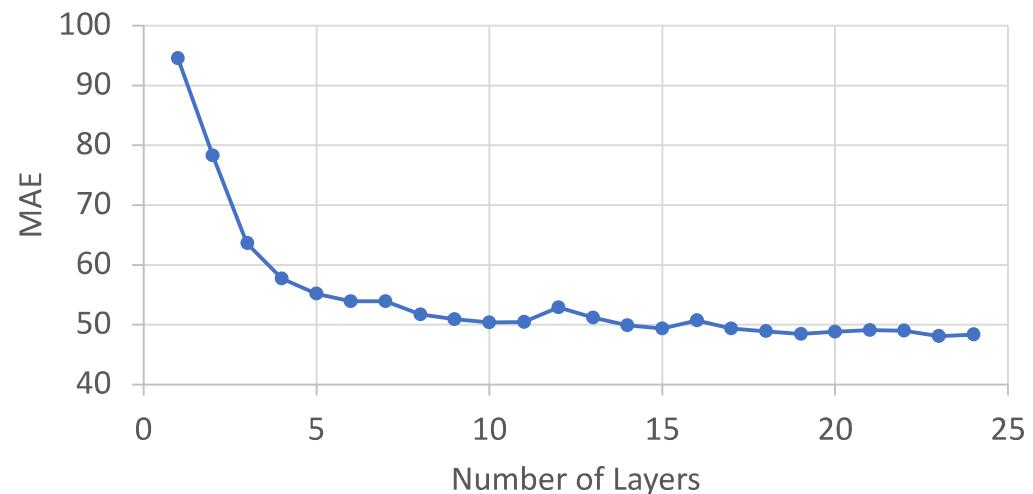
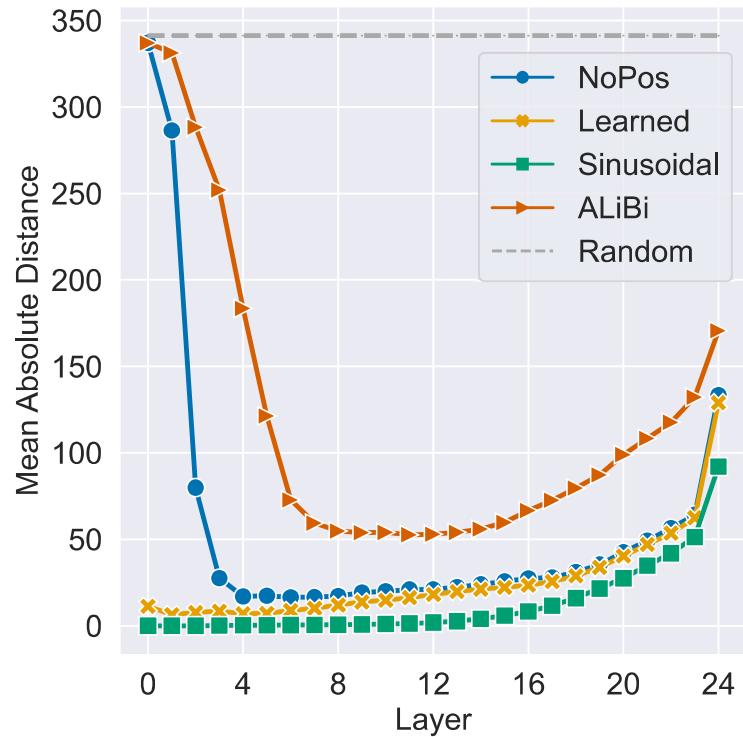
Peter J. Ramadge

Princeton University

Probing



Probing



NoPos: POS info in hidden too

Proof



vector variance position

Proof

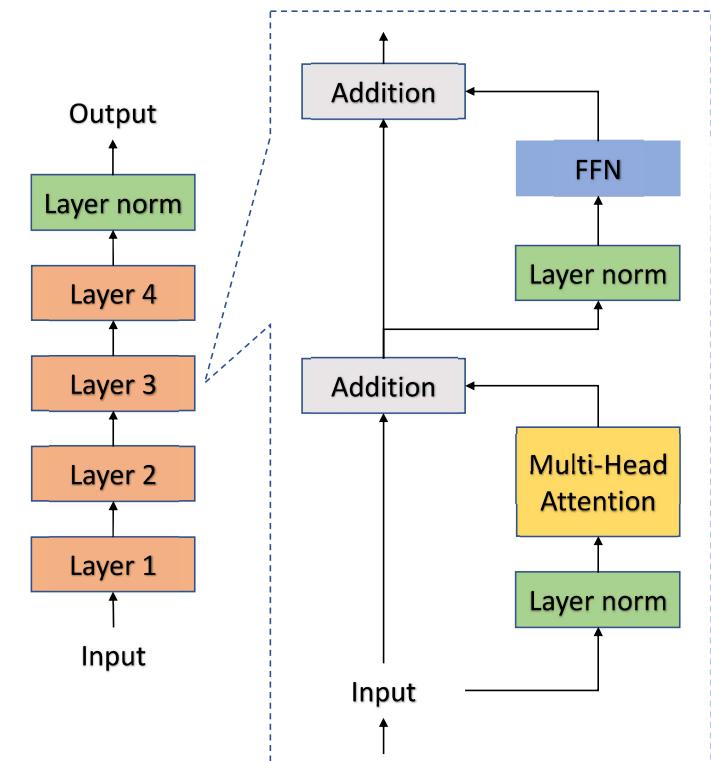


vector variance \leftrightarrow position

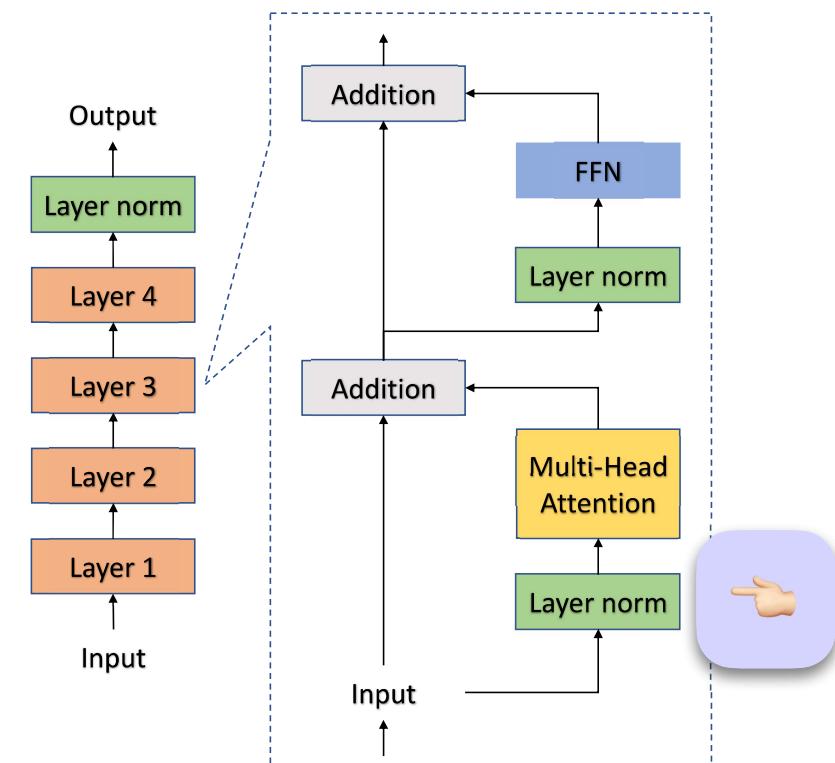
are: hidden dimension $d = 768$, number of attention heads $H = 12$, head dimension $d/H = 64$, sequence length $L = 512$, standard deviation for initialization $\sigma = 0.02$. All proofs of lemmas are deferred to Appendix A.

Given a sequence of randomly sampled input embeddings $\{\mathbf{x}_m\}_{m=1}^L$, where each element of $\mathbf{x}_m \in \mathbb{R}^d$ is sampled i.i.d from $N(0, \sigma^2)$, a TLM

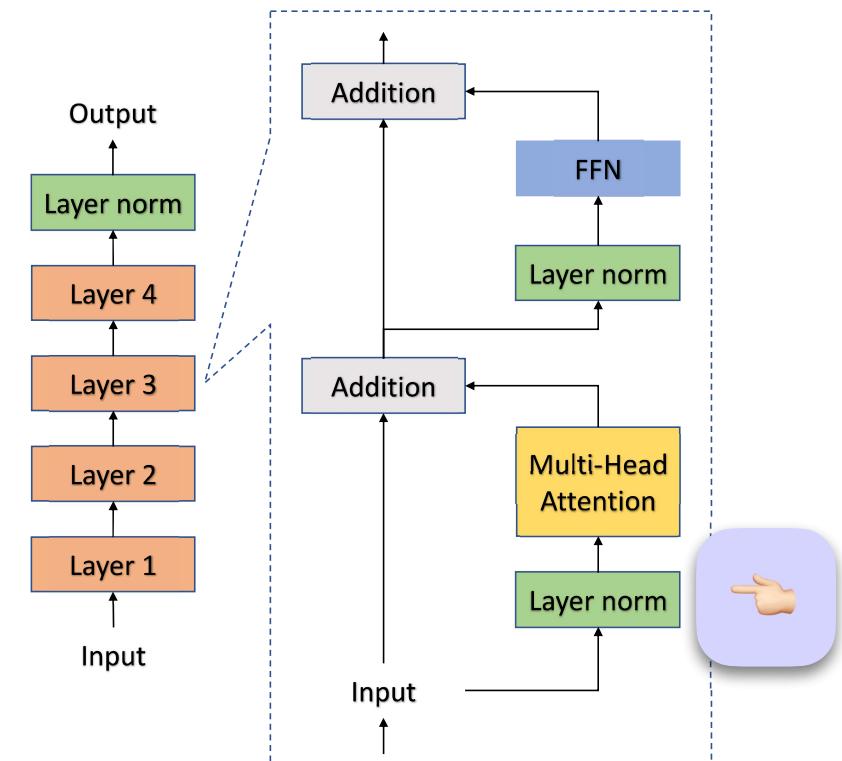
Layer Norm



Layer Norm



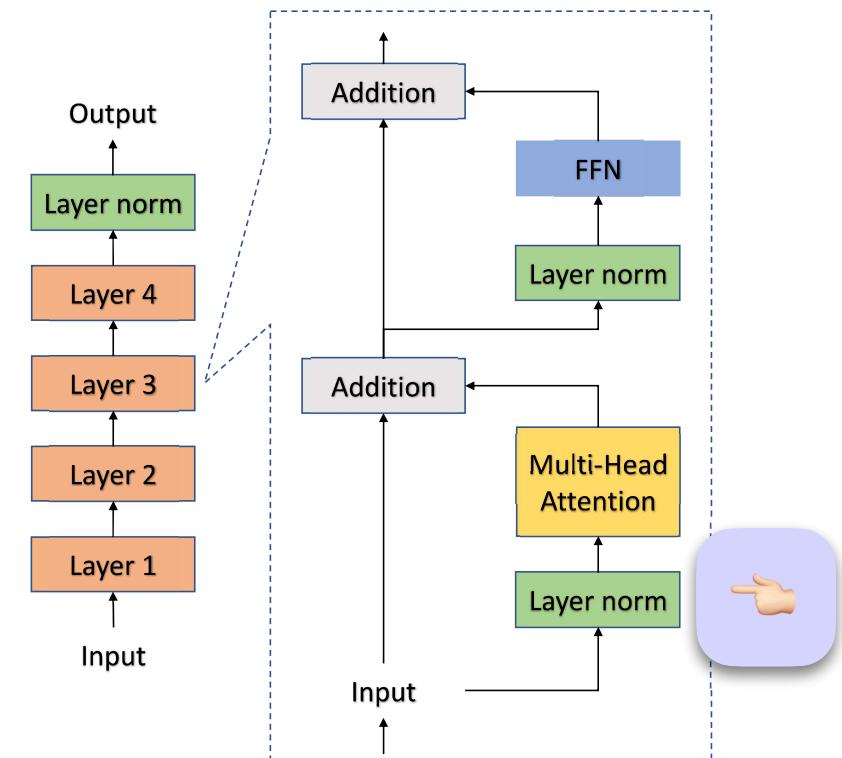
Layer Norm



$$\bar{x}_{m,:} = \frac{\sum_{i=1}^d x_{mi}}{d}, \quad S(\mathbf{x}_{m,:}) = \frac{\sum_{i=1}^d (\mathbf{x}_{mi} - \bar{x}_{m,:})^2}{d}$$

$$e_{mi} = \frac{\mathbf{x}_{mi} - \bar{x}_{m,:}}{\sqrt{S(\mathbf{x}_{m,:})}} * \gamma + \beta$$
$$\stackrel{(*)}{\approx} \frac{\mathbf{x}_{mi} - \mathbb{E}[\mathbf{x}_{mi}]}{\sqrt{\mathbb{V}[\mathbf{x}_{mi}]}} \sim N(0, 1),$$

Layer Norm



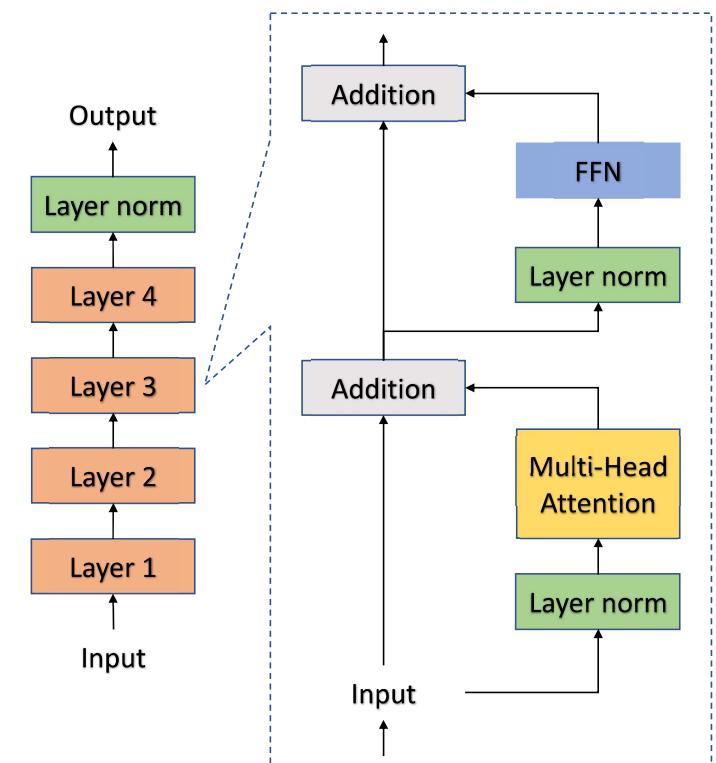
$$\bar{\mathbf{x}}_{m,:} = \frac{\sum_{i=1}^d \mathbf{x}_{mi}}{d}, \quad S(\mathbf{x}_{m,:}) = \frac{\sum_{i=1}^d (\mathbf{x}_{mi} - \bar{\mathbf{x}}_{m,:})^2}{d}$$

$$e_{mi} = \frac{\mathbf{x}_{mi} - \bar{\mathbf{x}}_{m,:}}{\sqrt{S(\mathbf{x}_{m,:})}} * \gamma + \beta$$

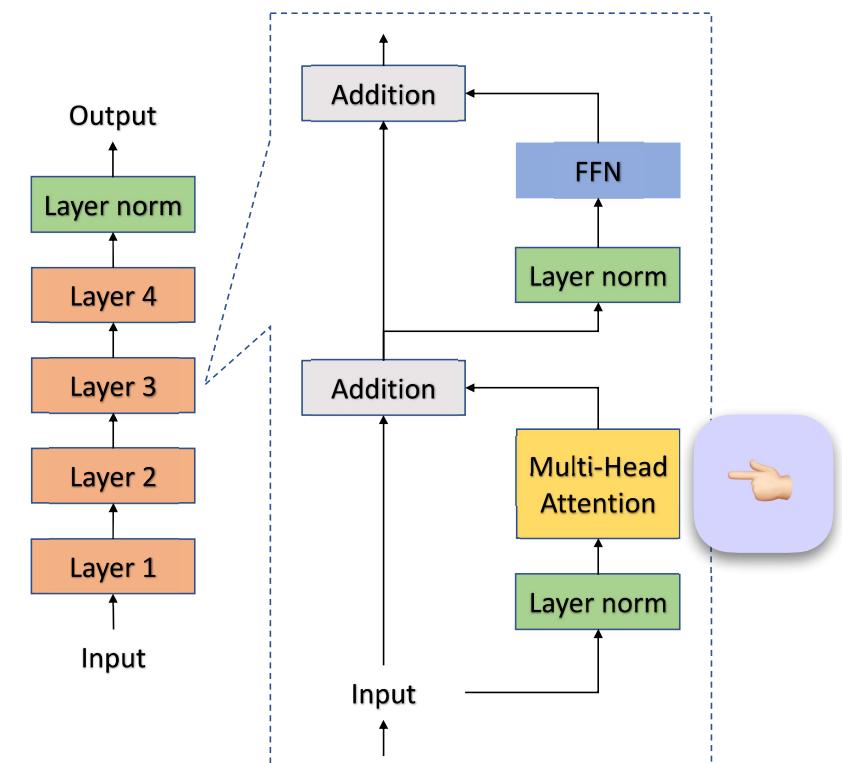
$$\stackrel{(*)}{\approx} \frac{\mathbf{x}_{mi} - \mathbb{E}[\mathbf{x}_{mi}]}{\sqrt{\mathbb{V}[\mathbf{x}_{mi}]}} \sim N(0, 1),$$

$\gamma = 1$ and $\beta = 0$

MHA

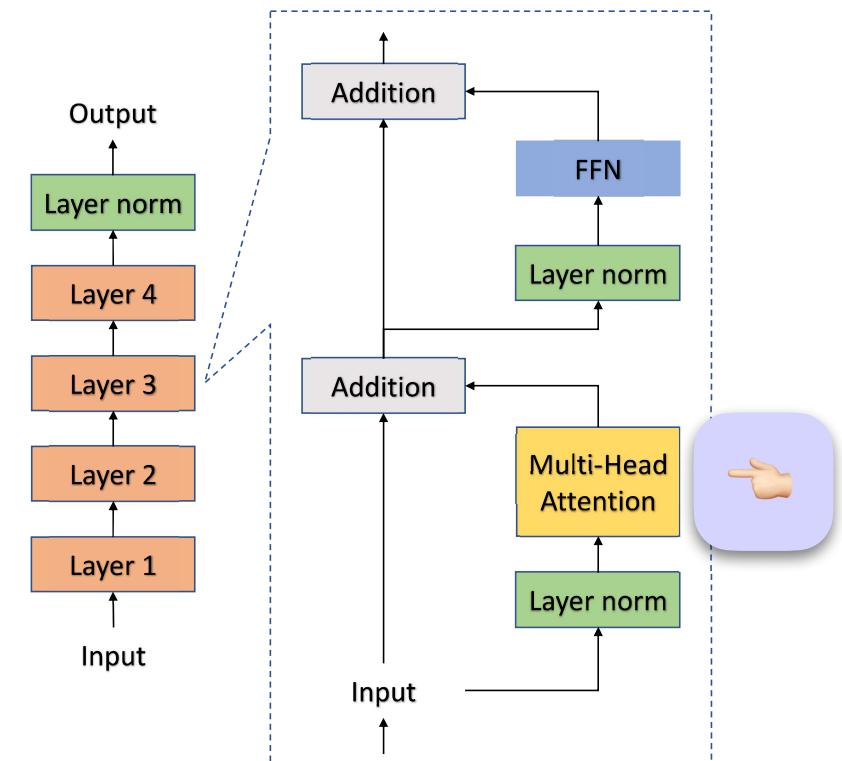


MHA

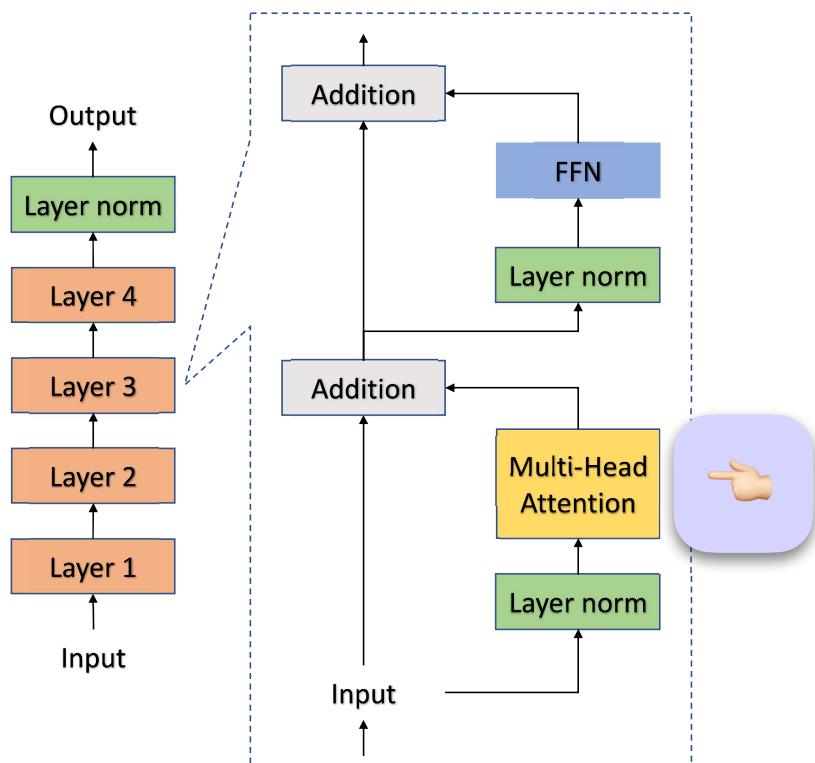


MHA

Lemma 1. q_m , k_m , and v_m have zero mean and $(d\sigma^2) \cdot I$ covariance matrix.



MHA

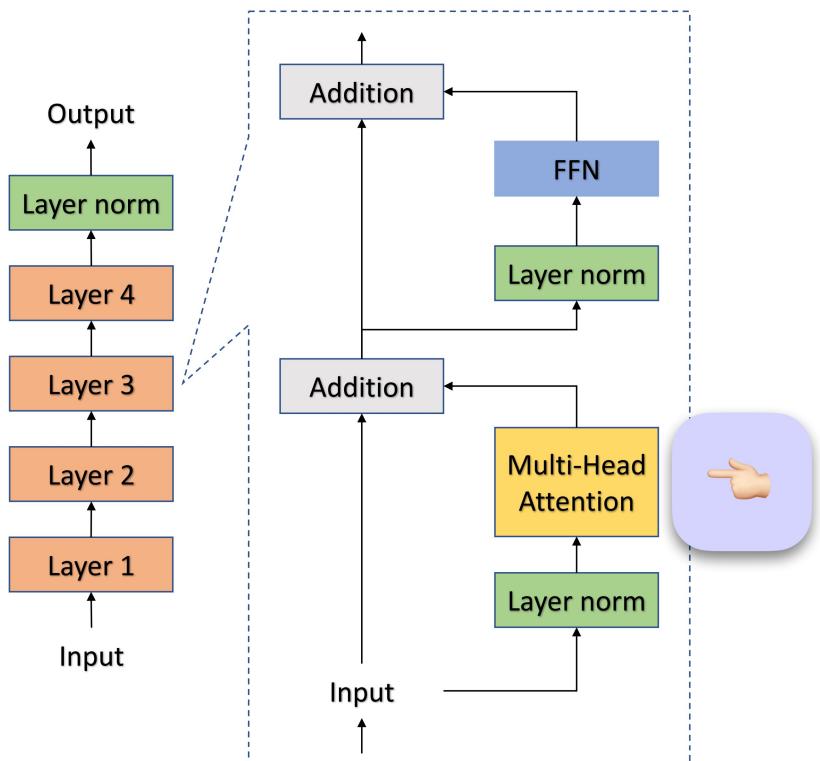


Lemma 1. q_m , k_m , and v_m have zero mean and $(d\sigma^2) \cdot I$ covariance matrix.

$$q_m = W_q e_m, \quad k_m = W_k e_m, \quad v_m = W_v e_m,$$

$W_q, W_k, W_v \in \mathbb{R}^{\frac{d}{H} \times d}$ sampled i.i.d from $N(0, \sigma^2)$

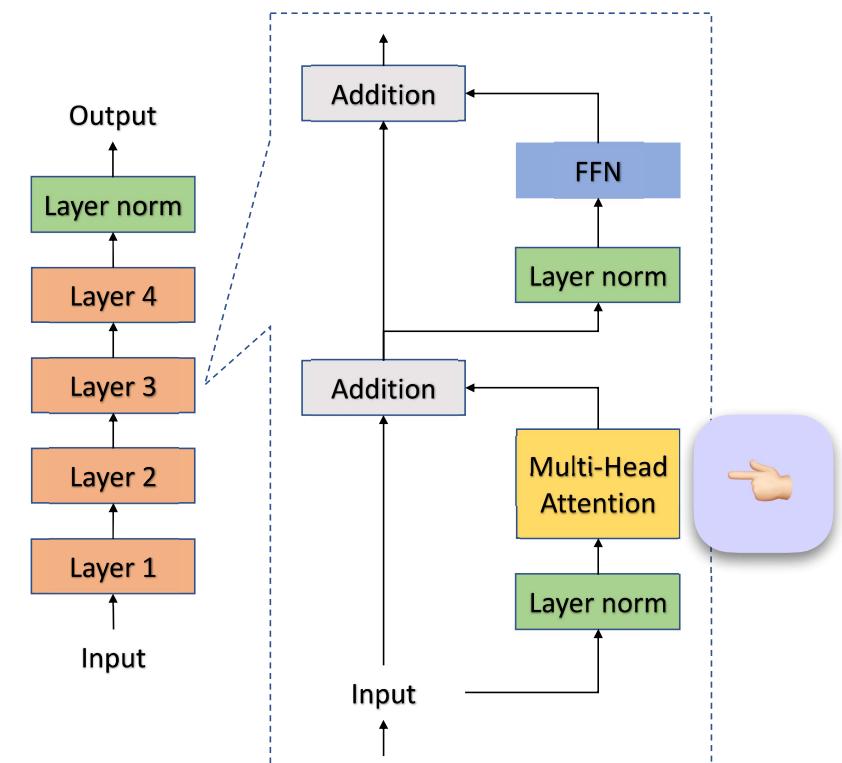
MHA: Proof



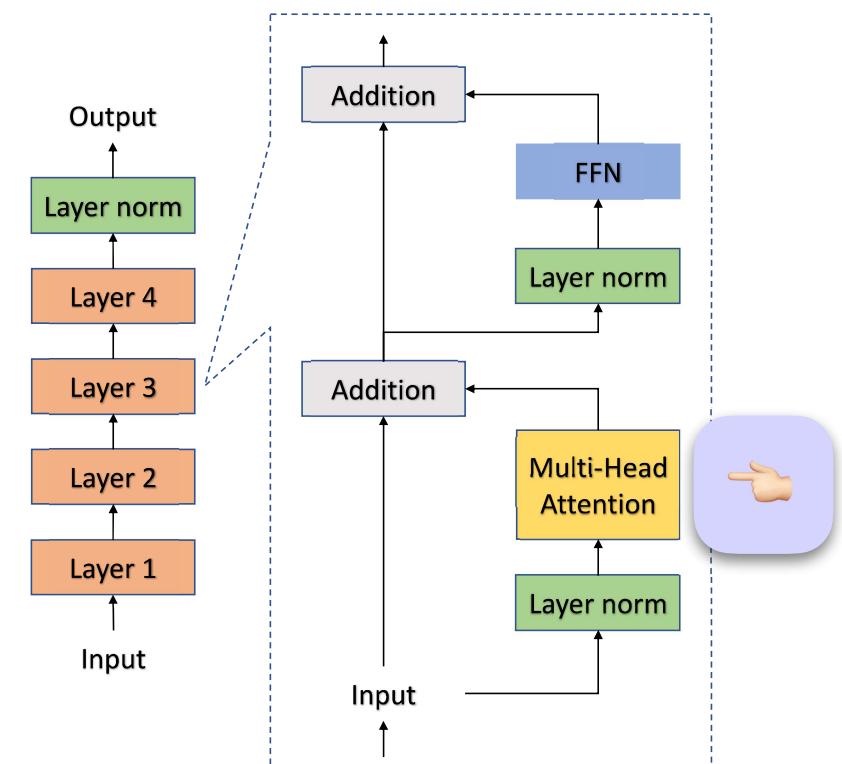
Lemma 1. q_m , k_m , and v_m have zero mean and $(d\sigma^2) \cdot I$ covariance matrix.

$$\begin{aligned}
 \text{cov}(\mathbf{v}_m, \mathbf{v}_n) &= \mathbb{E}[\mathbf{v}_m \mathbf{v}_n^\top] \\
 &= \mathbb{E}[\mathbf{W}_v \mathbf{e}_m \mathbf{e}_n^\top \mathbf{W}_v^\top] \\
 &= \mathbb{E} \left[\begin{bmatrix} \mathbf{r}_1^\top \mathbf{e}_m \\ \vdots \\ \mathbf{r}_{\frac{d}{H}}^\top \mathbf{e}_m \end{bmatrix} \begin{bmatrix} \mathbf{e}_n^\top \mathbf{r}_1 & \dots & \mathbf{e}_n^\top \mathbf{r}_{\frac{d}{H}} \end{bmatrix} \right] \\
 &= \left[\mathbb{E}[\mathbf{r}_i^\top \mathbf{e}_m \mathbf{e}_n^\top \mathbf{r}_j] \right]_{i,j=1}^{\frac{d}{H}} \\
 &= \left[\mathbb{E}[\text{Tr}(\mathbf{r}_j \mathbf{r}_i^\top \mathbf{e}_m \mathbf{e}_n^\top)] \right]_{i,j=1}^{\frac{d}{H}} \\
 &= \left[\text{Tr}(\mathbb{E}[\mathbf{r}_j \mathbf{r}_i^\top] \mathbb{E}[\mathbf{e}_m \mathbf{e}_n^\top]) \right]_{i,j=1}^{\frac{d}{H}} \\
 &\stackrel{(*)}{=} \left[\text{Tr}((\mathbb{1}_{i=j} \sigma^2) \cdot I_d \cdot \mathbb{1}_{m=n} \cdot I_d) \right]_{i,j=1}^{\frac{d}{H}} \\
 &= \left[\mathbb{1}_{i=j} \mathbb{1}_{m=n} d\sigma^2 \right]_{i,j=1}^{\frac{d}{H}} \\
 &= (\mathbb{1}_{m=n} d\sigma^2) \cdot I_{d/H}
 \end{aligned}$$

MHA

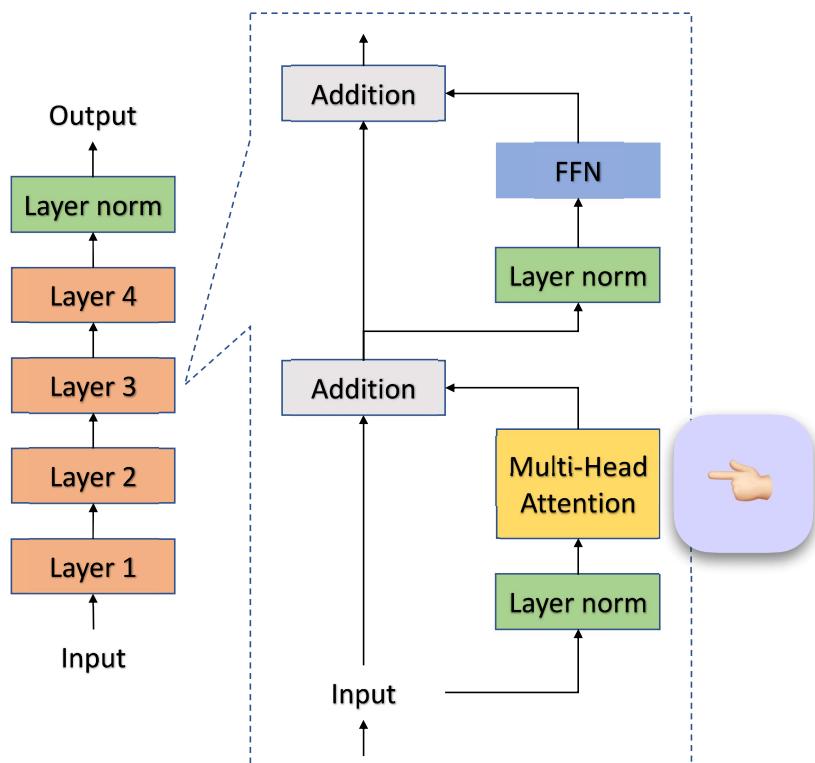


MHA



Lemma 2. l_{mn} has zero mean and $\frac{d^3\sigma^4}{H^2}$ variance.
 $l_{mn}/\sqrt{d/H}$ has $\frac{d^2\sigma^4}{H}$ variance.

MHA

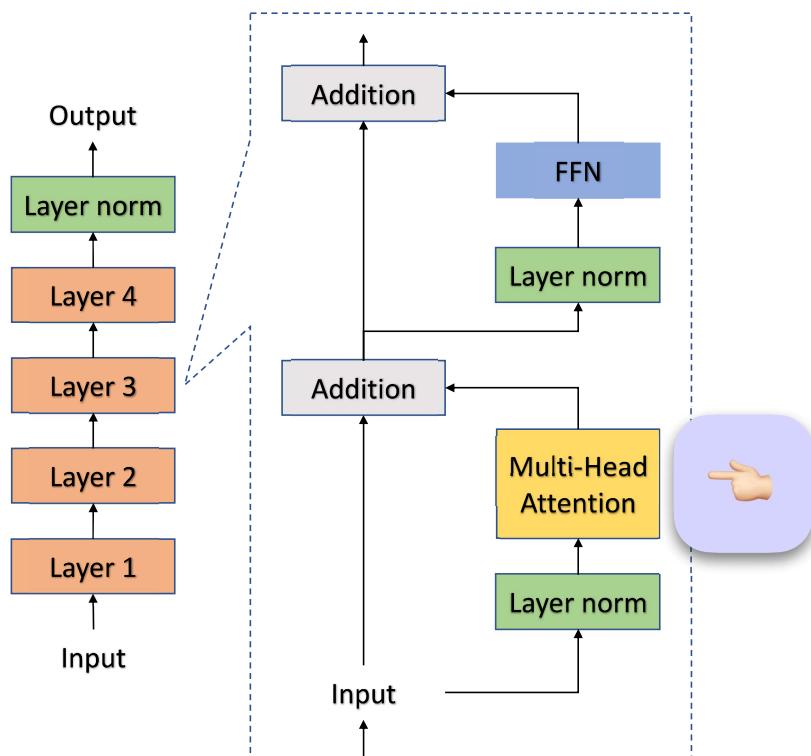


Lemma 2. l_{mn} has zero mean and $\frac{d^3\sigma^4}{H^2}$ variance.
 $l_{mn}/\sqrt{d/H}$ has $\frac{d^2\sigma^4}{H}$ variance.

$$l_{mn} = \begin{cases} \langle q_m, k_n \rangle, & \text{if } m \geq n \\ -\infty, & \text{otherwise} \end{cases}$$

$$a_{mn} = \frac{\exp(l_{mn}/\sqrt{d/H})}{\sum_{i=1}^L \exp(l_{mi}/\sqrt{d/H})}$$

MHA



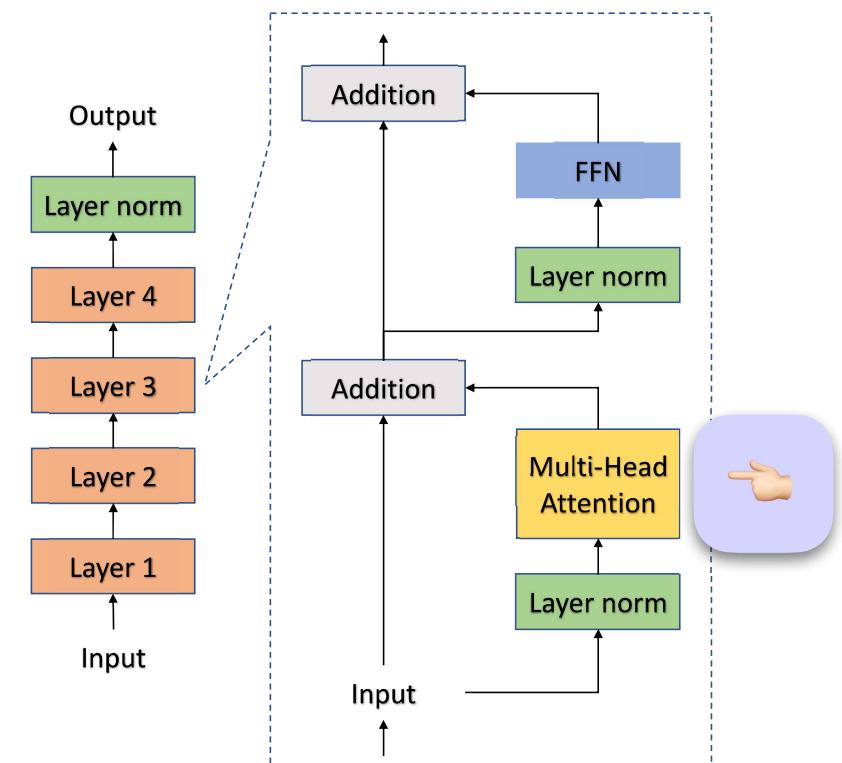
Lemma 2. l_{mn} has zero mean and $\frac{d^3\sigma^4}{H^2}$ variance.
 $l_{mn}/\sqrt{d/H}$ has $\frac{d^2\sigma^4}{H}$ variance.

$$l_{mn} = \begin{cases} \langle q_m, k_n \rangle, & \text{if } m \geq n \\ -\infty, & \text{otherwise} \end{cases}$$

$$a_{mn} = \frac{\exp(l_{mn}/\sqrt{d/H})}{\sum_{i=1}^L \exp(l_{mi}/\sqrt{d/H})}$$

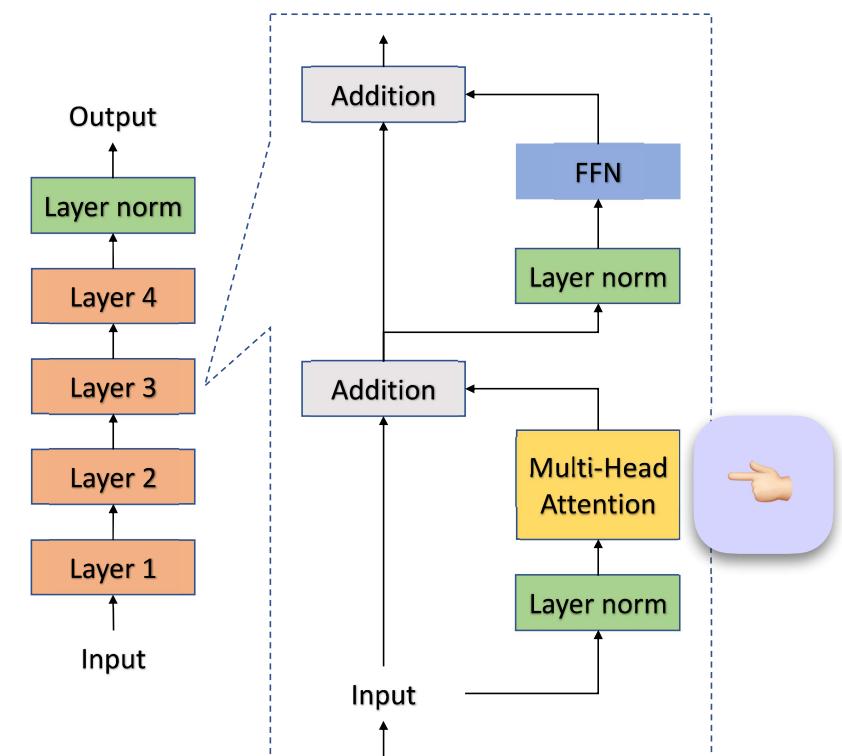
Property 1. When $\sigma^4 \ll \frac{H}{d^2}$, $l_{m,:}$ has small variance, making the attention weights $a_{m,:}$ almost evenly distributed among all positions.²

MHA



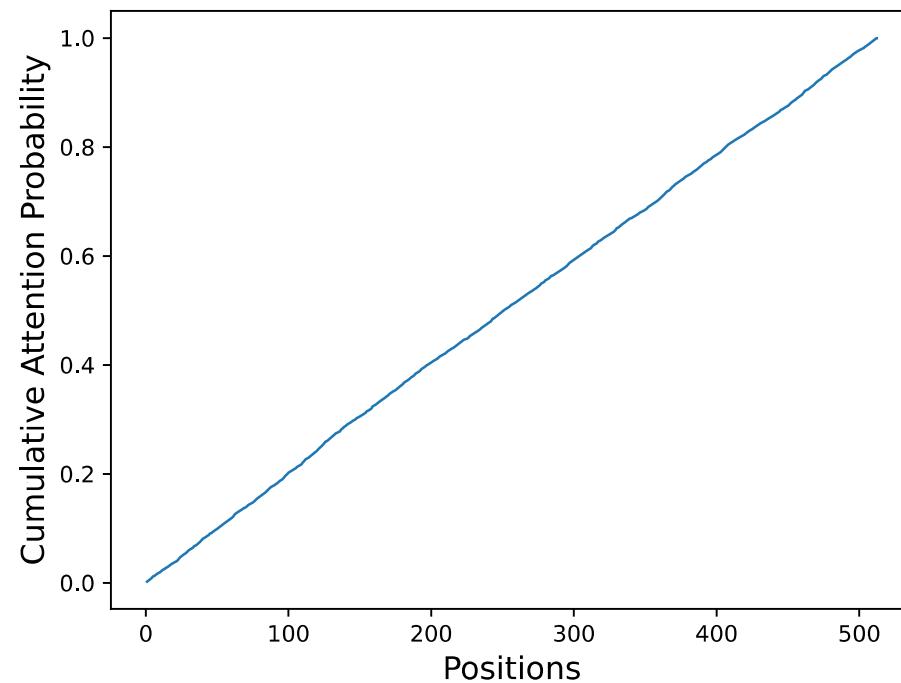
Property 1. When $\sigma^4 \ll \frac{H}{d^2}$, $l_{m,:}$ has small variance, making the attention weights $a_{m,:}$ almost evenly distributed among all positions.²

MHA

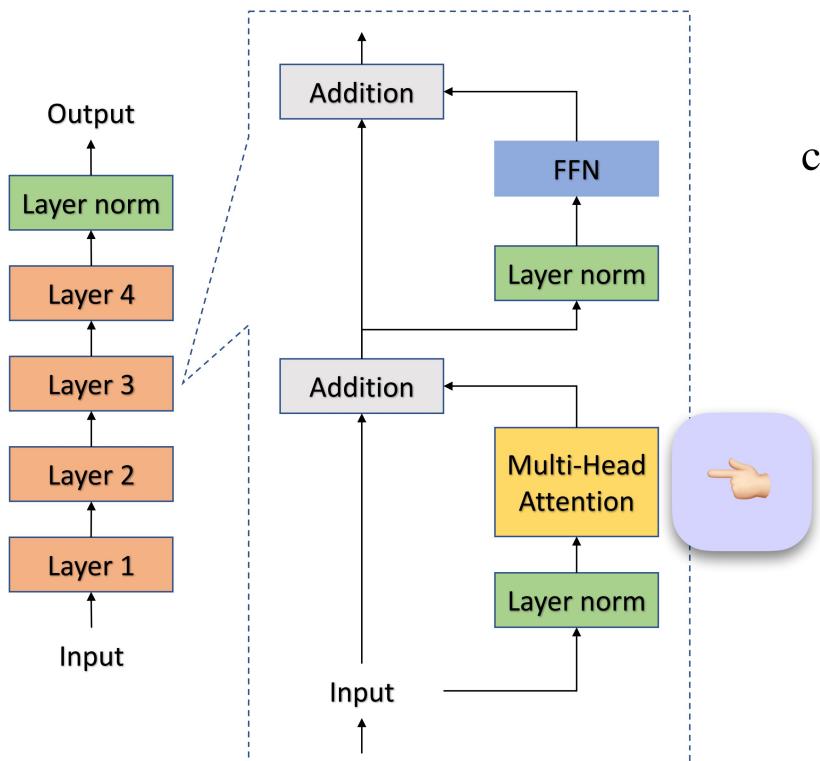


Property 1. When $\sigma^4 \ll \frac{H}{d^2}$, $l_{m,:}$ has small variance, making the attention weights $a_{m,:}$ almost evenly distributed among all positions.²

$$\frac{768^2 \cdot 0.02^4}{12} \approx 0.0079$$



MHA: Proof



$\text{cov}(l_{mn}, l_{mp})$

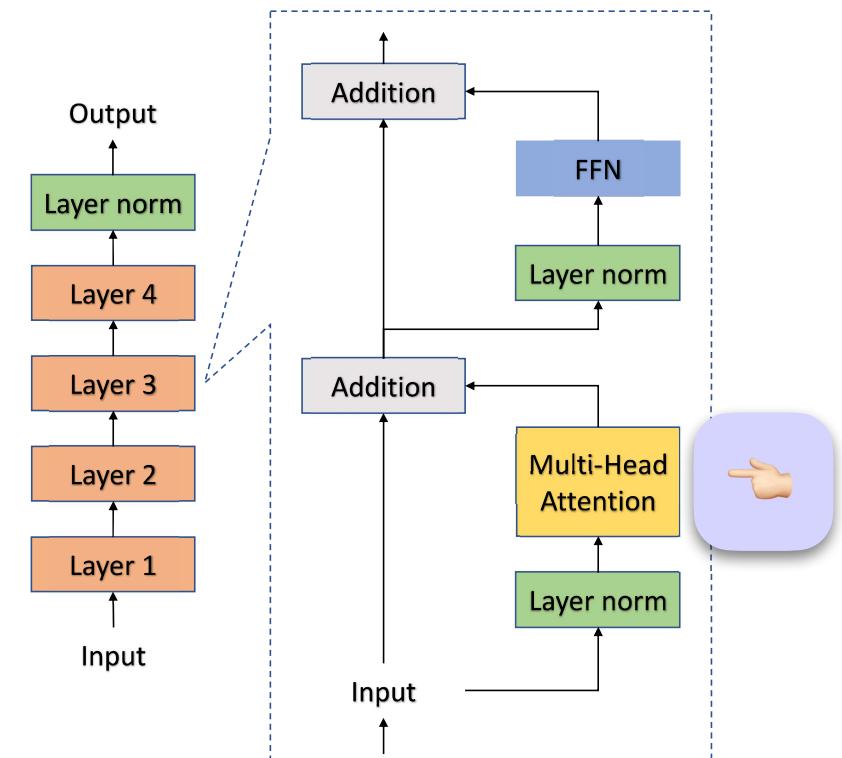
$$\begin{aligned}
 &= \mathbb{E}[(e_m^\top W_q^\top W_k e_n)(e_m^\top W_q^\top W_k e_p)^\top] \\
 &= \mathbb{E}[\text{Tr}(e_m^\top W_q^\top W_k e_n e_p^\top W_k^\top W_q e_m)] \\
 &= \mathbb{E}[\text{Tr}(e_m e_m^\top W_q^\top W_k e_n e_p^\top W_k^\top W_q)] \\
 &= \text{Tr}(\mathbb{E}[e_m e_m^\top] \mathbb{E}[W_q^\top W_k e_n e_p^\top W_k^\top W_q]) \\
 &= \mathbb{E}[\text{Tr}(e_n e_p^\top W_k^\top W_q W_q^\top W_k)] \\
 &= \text{Tr}(\mathbb{E}[e_n e_p^\top] \mathbb{E}[W_k^\top W_q W_q^\top W_k]) \\
 &= (\mathbb{1}_{n=p}) \text{Tr}(\mathbb{E}[W_q W_q^\top] \mathbb{E}[W_k W_k^\top]) \\
 &\stackrel{(*)}{=} (\mathbb{1}_{n=p}) \text{Tr}\left(\left(\frac{d}{H}\sigma^2 \cdot I\right)\left(\frac{d}{H}\sigma^2 \cdot I\right)\right) \\
 &= (\mathbb{1}_{n=p}) \frac{d^3\sigma^4}{H^2}
 \end{aligned}$$

$$\mathbb{E}[W_q W_q^\top] = \mathbb{E} \left[\begin{bmatrix} r_1^\top \\ \vdots \\ r_{\frac{d}{H}}^\top \end{bmatrix} \begin{bmatrix} r_1 & \dots & r_{\frac{d}{H}} \end{bmatrix} \right]$$

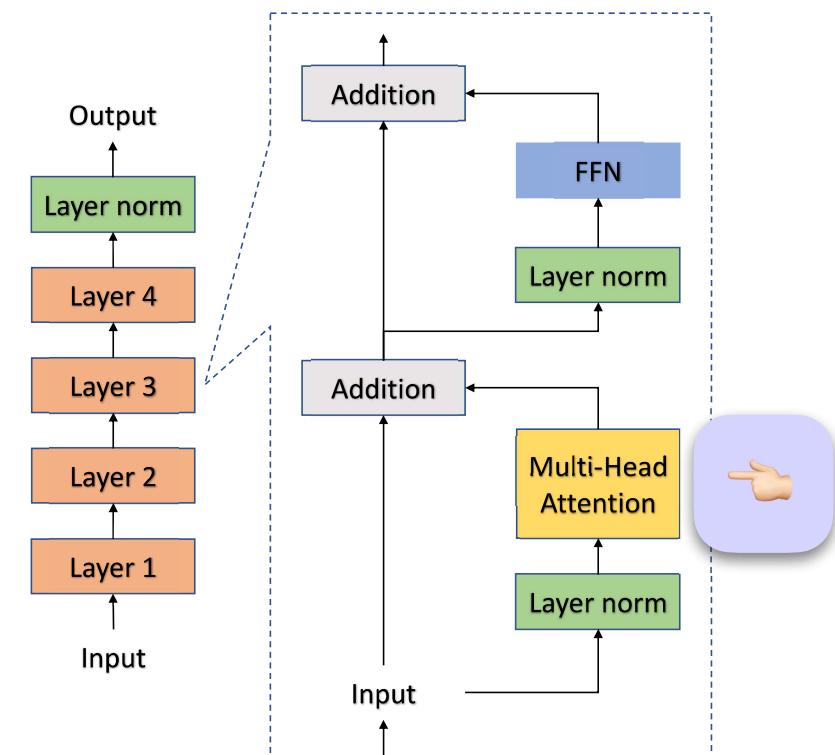
$$= \left[\mathbb{E}[r_i^\top r_j] \right]_{i,j=1}^{\frac{d}{H}} = \frac{d}{H}\sigma^2 \cdot I$$

MHA

Lemma 3. \mathbf{o}_m has zero mean and $\frac{d^2\sigma^4}{m}\mathbf{I}$ covariance matrix.

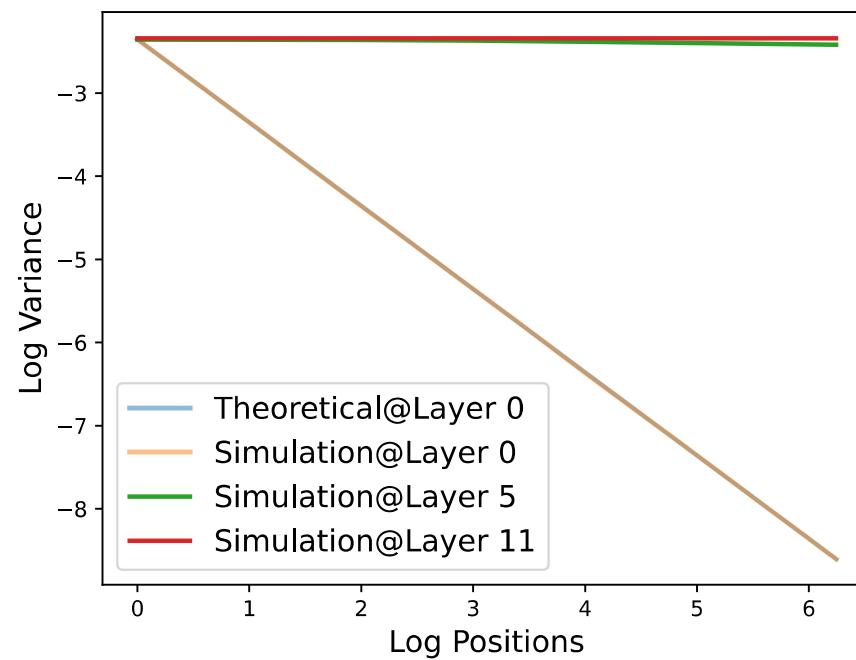


MHA

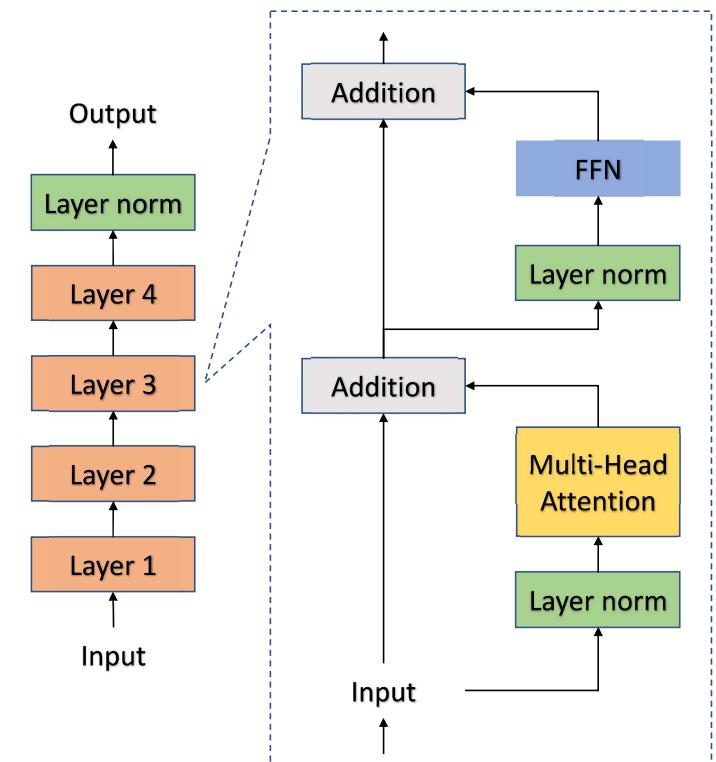


Lemma 3. \mathbf{o}_m has zero mean and $\frac{d^2\sigma^4}{m}I$ covariance matrix.

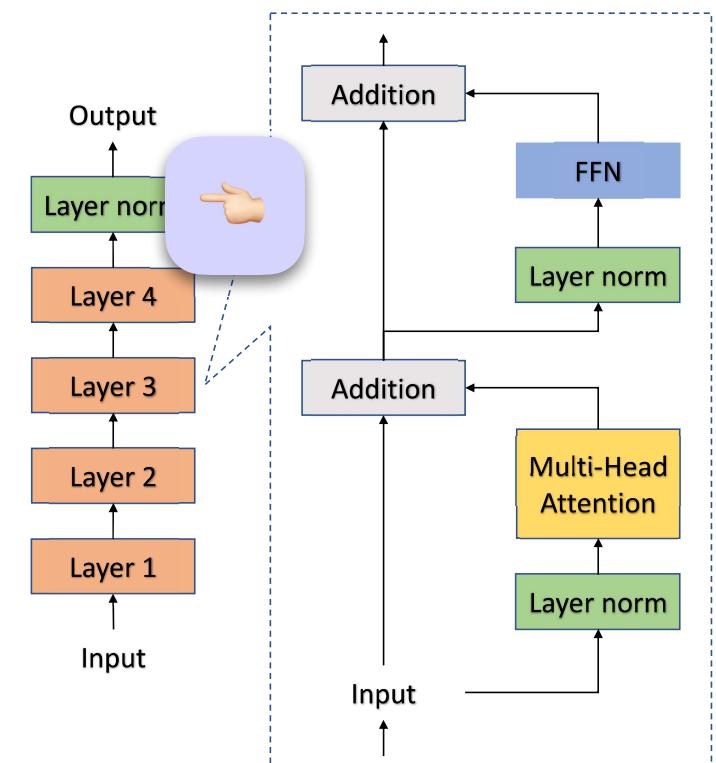
$$\mathbf{o}_m = \mathbf{W}_o \left(\bigoplus_{h=1}^H \sum_{n=1}^L a_{mn}^{(h)} \mathbf{v}_n^{(h)} \right),$$



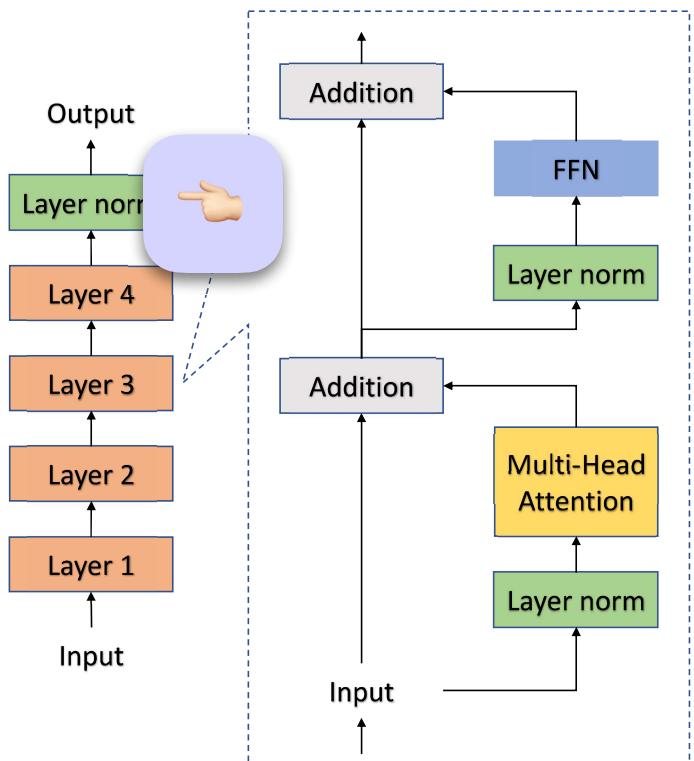
Final Layer Norm



Final Layer Norm



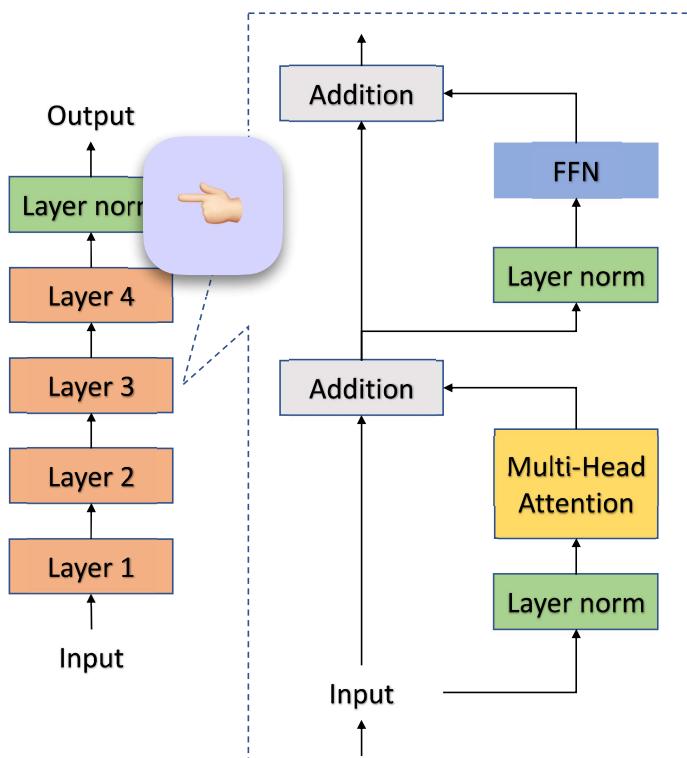
Final Layer Norm



Lemma 4. *The variance of the j -th dimension of y_m is:*

$$\frac{m\sigma^2 + \sum_i (\mathbf{W}_{o,j} : \mathbf{W}_{v,:i})^2}{m\sigma^2 + d^2\sigma^4},$$

Final Layer Norm



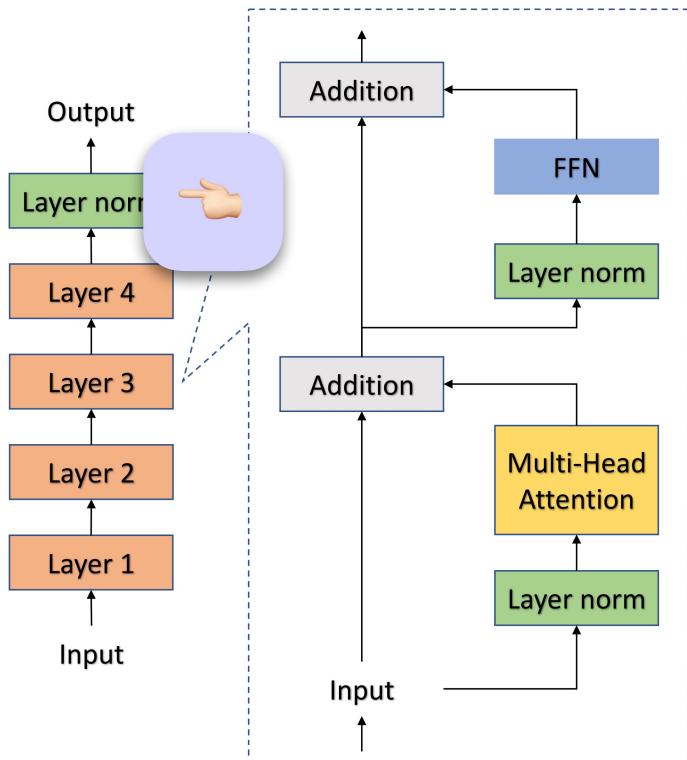
Lemma 4. *The variance of the j -th dimension of \mathbf{y}_m is:*

$$\frac{m\sigma^2 + \sum_i (\mathbf{W}_{o,j} : \mathbf{W}_{v,:i})^2}{m\sigma^2 + d^2\sigma^4},$$

$$\mathbf{y}'_{mi} \approx \frac{\mathbf{y}_{mi} - \mathbb{E}[\mathbf{y}_{mi}]}{\sqrt{\mathbb{V}[\mathbf{y}_{mi}]}} \approx \frac{\mathbf{x}_{mi} + \mathbf{W}_o \mathbf{W}_v \frac{\sum_n \mathbf{e}_{ni}}{m}}{\sqrt{\sigma^2 + \frac{d^2\sigma^4}{m}}},$$

$$\begin{aligned} \mathbb{E}[\mathbf{y}_{mi}] &= 0, \quad \mathbb{V}[\mathbf{y}_{mi}] = \mathbb{V}[\mathbf{x}_{mi}] + \mathbb{V}[\mathbf{o}_{mi}] \\ &= \sigma^2 + \frac{d^2\sigma^4}{m} \end{aligned}$$

Final Layer Norm: Proof



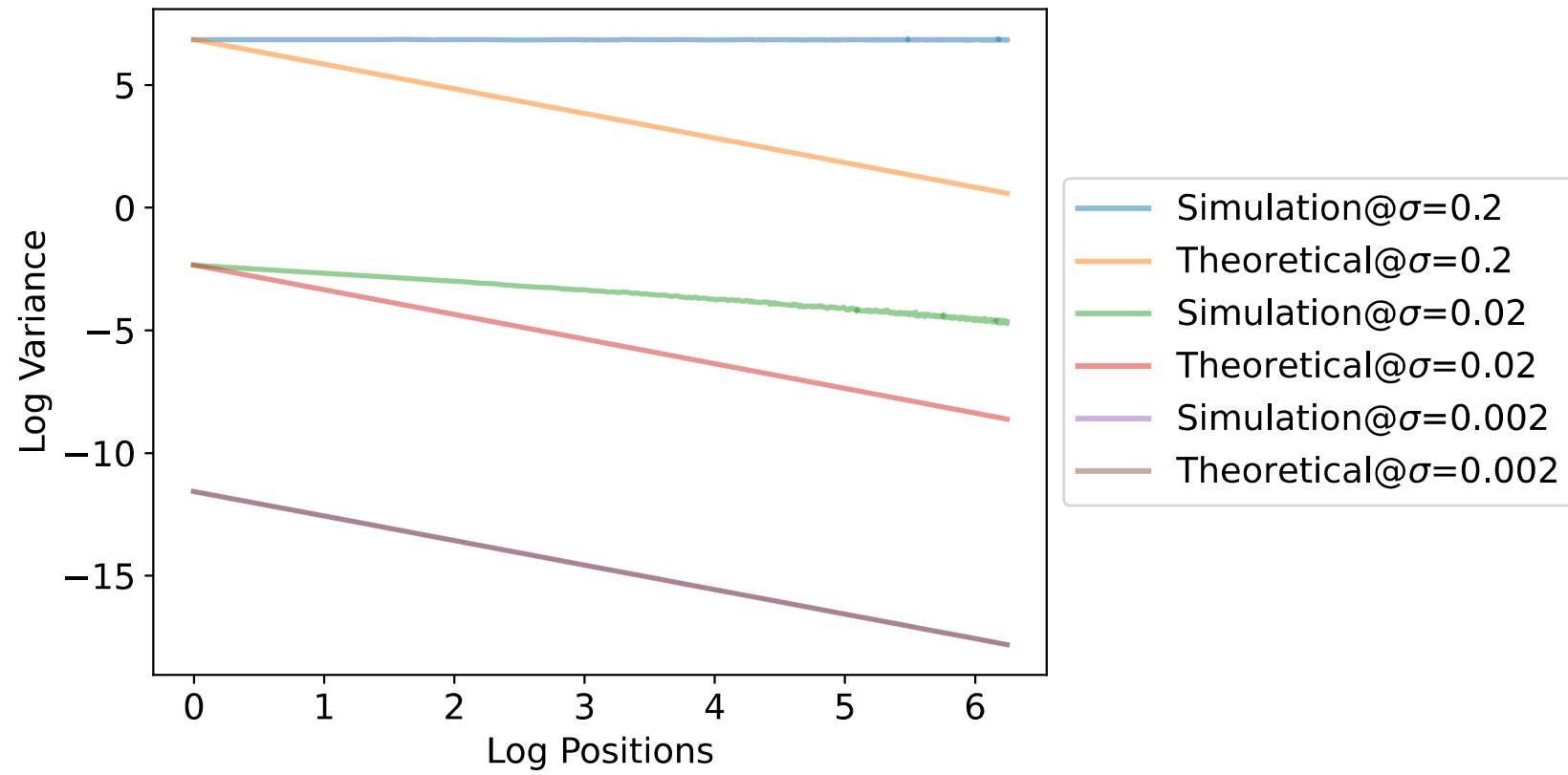
Lemma 4. *The variance of the j -th dimension of y_m is:*

$$\frac{m\sigma^2 + \sum_i (\mathbf{W}_{o,j} : \mathbf{W}_{v,:i})^2}{m\sigma^2 + d^2\sigma^4},$$

$$\begin{aligned} & \text{cov}(\mathbf{o}_m, \mathbf{o}_m) \\ & \approx \mathbb{E} \left[\mathbf{W}_o \frac{\sum_{i=1}^m \mathbf{v}_i}{m} \frac{\sum_{j=1}^m \mathbf{v}_j^\top}{m} \mathbf{W}_o^\top \right] \\ & = \frac{1}{m^2} \sum_{i,j=1}^m \mathbb{E}[\mathbf{W}_o \mathbf{v}_i \mathbf{v}_j^\top \mathbf{W}_o^\top] \\ & = \frac{1}{m^2} \sum_{i,j=1}^m \mathbb{E} \left[\begin{bmatrix} \mathbf{r}_1^\top \mathbf{v}_i \\ \vdots \\ \mathbf{r}_d^\top \mathbf{v}_i \end{bmatrix} \begin{bmatrix} \mathbf{v}_j^\top \mathbf{r}_1 & \dots & \mathbf{v}_j^\top \mathbf{r}_d \end{bmatrix} \right] \\ & = \frac{1}{m^2} \sum_{i,j=1}^m \left[\mathbb{E}[\mathbf{r}_k^\top \mathbf{v}_i \mathbf{v}_j^\top \mathbf{r}_l] \right]_{k,l=1}^d \\ & = \frac{1}{m^2} \sum_{i,j=1}^m \left[\mathbb{E}[\text{Tr}(\mathbf{r}_l \mathbf{r}_k^\top \mathbf{v}_i \mathbf{v}_j^\top)] \right]_{k,l=1}^d \\ & = \frac{1}{m^2} \sum_{i,j=1}^m \left[\text{Tr}(\mathbb{E}[\mathbf{r}_l \mathbf{r}_k^\top] \mathbb{E}[\mathbf{v}_i \mathbf{v}_j^\top]) \right]_{k,l=1}^d \\ & \stackrel{(*)}{=} \frac{1}{m^2} \sum_{i,j=1}^m [\text{Tr}((\mathbb{1}_{k=l}\sigma^2) \cdot I \cdot (\mathbb{1}_{i=j}d\sigma^2) \cdot I)]_{k,l=1}^d \\ & = \frac{d^2\sigma^4}{m} I \end{aligned}$$

Relaxing

? What if Property 1 does not hold?



Relaxing

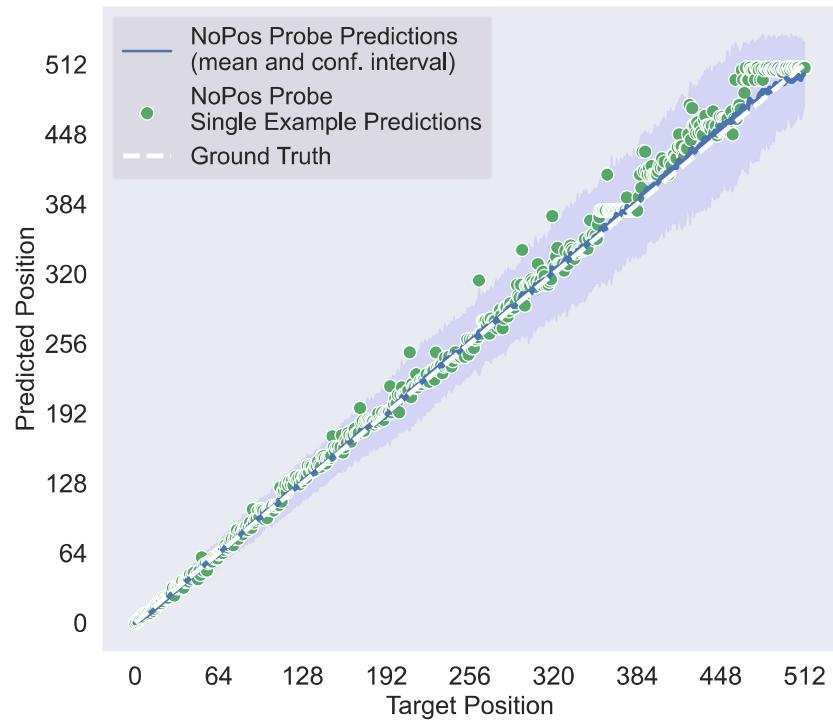
? Other Initialization Schemes

$$(\mathbf{W}_q, \mathbf{W}_k, \mathbf{W}_v, \mathbf{W}_o)$$

i.i.d. samples from a distribution

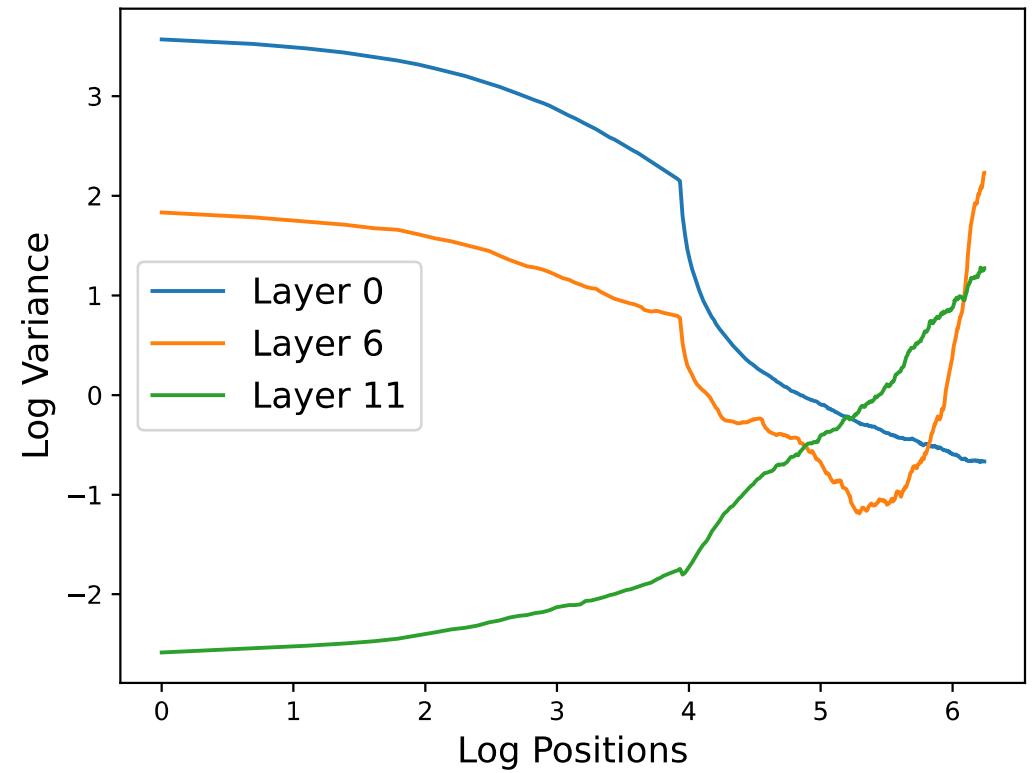
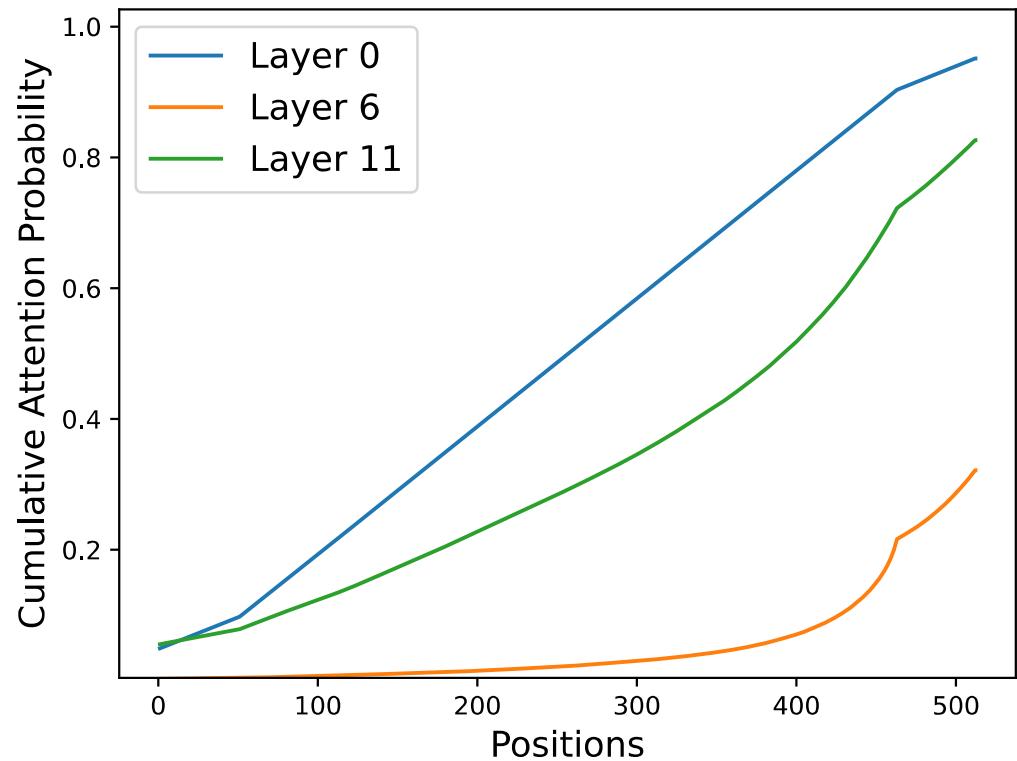
with zero mean and finite variance.

Discussions



Lemma 3. \mathbf{o}_m has zero mean and $\frac{d^2\sigma^4}{m}\mathbf{I}$ covariance matrix.

Discussions: Pretrained



Discussions

? Why does BERT fail to converge NoPE

Lemma 3. o_m has zero mean and $\frac{d^2\sigma^4}{L}I$ covariance matrix.

ACL2023 Honorable Mentions

Latent Positional Information is in the Self-Attention Variance of Transformer Language Models Without Positional Embeddings

Ta-Chung Chi[†]

Carnegie Mellon University

Ting-Han Fan

Princeton University

Li-Wei Chen

Carnegie Mellon University

Alexander I. Rudnicky

Carnegie Mellon University

Peter J. Ramadge

Princeton University

ACL2023 Honorable Mentions

Latent Positional Information is in the Self-Attention Variance of Transformer Language Models Without Positional Embeddings

Ta-Chung Chi[†]

Carnegie Mellon University

Ting-Han Fan

Princeton University

Li-Wei Chen

Carnegie Mellon University

Alexander I. Rudnicky

Carnegie Mellon University

Peter J. Ramadge

Princeton University



NoPos 🤯 why → variance

ACL2023 Honorable Mentions

Latent Positional Information is in the Self-Attention Variance of Transformer Language Models Without Positional Embeddings

Ta-Chung Chi[†]

Carnegie Mellon University

Ting-Han Fan

Princeton University

Li-Wei Chen

Carnegie Mellon University

Alexander I. Rudnicky

Carnegie Mellon University

Peter J. Ramadge

Princeton University



NoPos why → variance

initialization pretrained; Downstream tasks

ACL2023 Honorable Mentions

Latent Positional Information is in the Self-Attention Variance of Transformer Language Models Without Positional Embeddings

Ta-Chung Chi[†]

Carnegie Mellon University

Ting-Han Fan

Princeton University

Li-Wei Chen

Carnegie Mellon University

Alexander I. Rudnicky

Carnegie Mellon University

Peter J. Ramadge

Princeton University

😱 NoPos ✅ why → variance

initialization 🙏 pretrained 😊; Downstream tasks 🤔

<S> $\{x_m\}_{m=1}^L$

苏剑林 发表于 August 11th, 2023

在我的实验中，NoPE没有任何长度外推性，甚至连RoPE都不如，所以我无法评价

回复评论

Preprint **Mila**

The Impact of Positional Encoding on Length Generalization in Transformers

Amirhossein Kazemnejad^{1,2}, Inkit Padhi³
Karthikeyan Natesan Ramamurthy³, Payel Das³, Siva Reddy^{1,2,4}

¹Mila - Québec AI Institute; ²McGill University;

³IBM Research; ⁴Facebook CIFAR AI Chair

Preprint **Mila**

The Impact of Positional Encoding on Length Generalization in Transformers

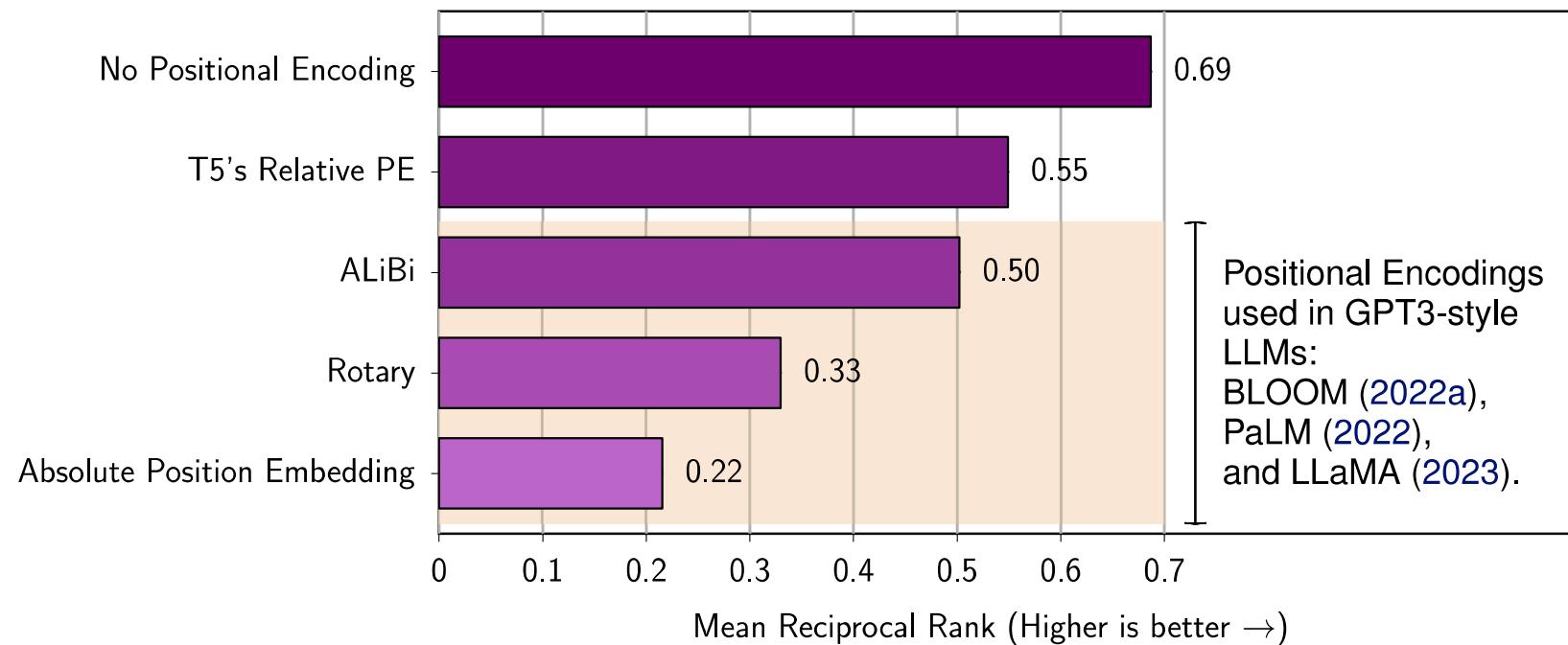
Amirhossein Kazemnejad^{1,2}, Inkit Padhi³
Karthikeyan Natesan Ramamurthy³, Payel Das³, Siva Reddy^{1,2,4}

¹Mila - Québec AI Institute; ²McGill University;

³IBM Research; ⁴Facebook CIFAR AI Chair

😱 NoPE 🤓 But why? 🤔

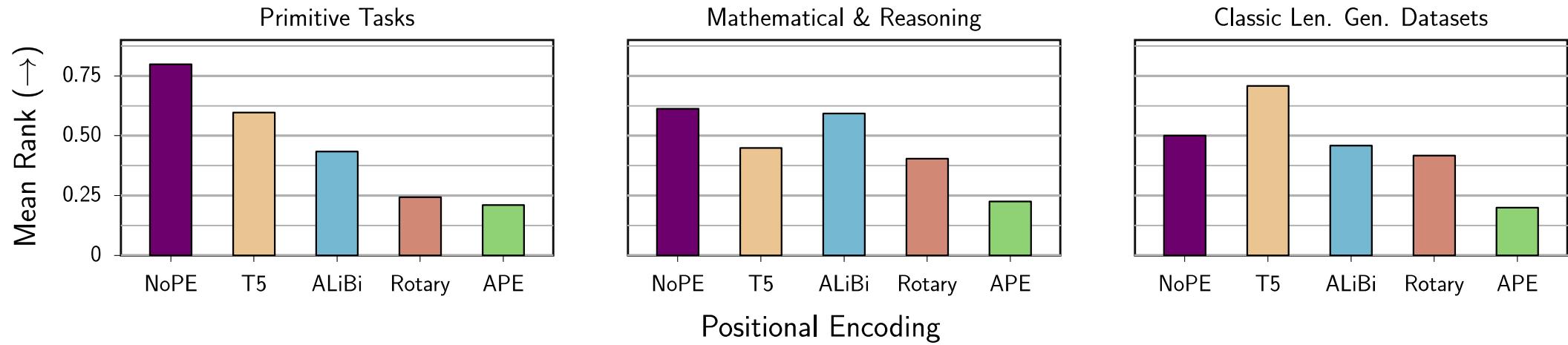
NoPE vs PE



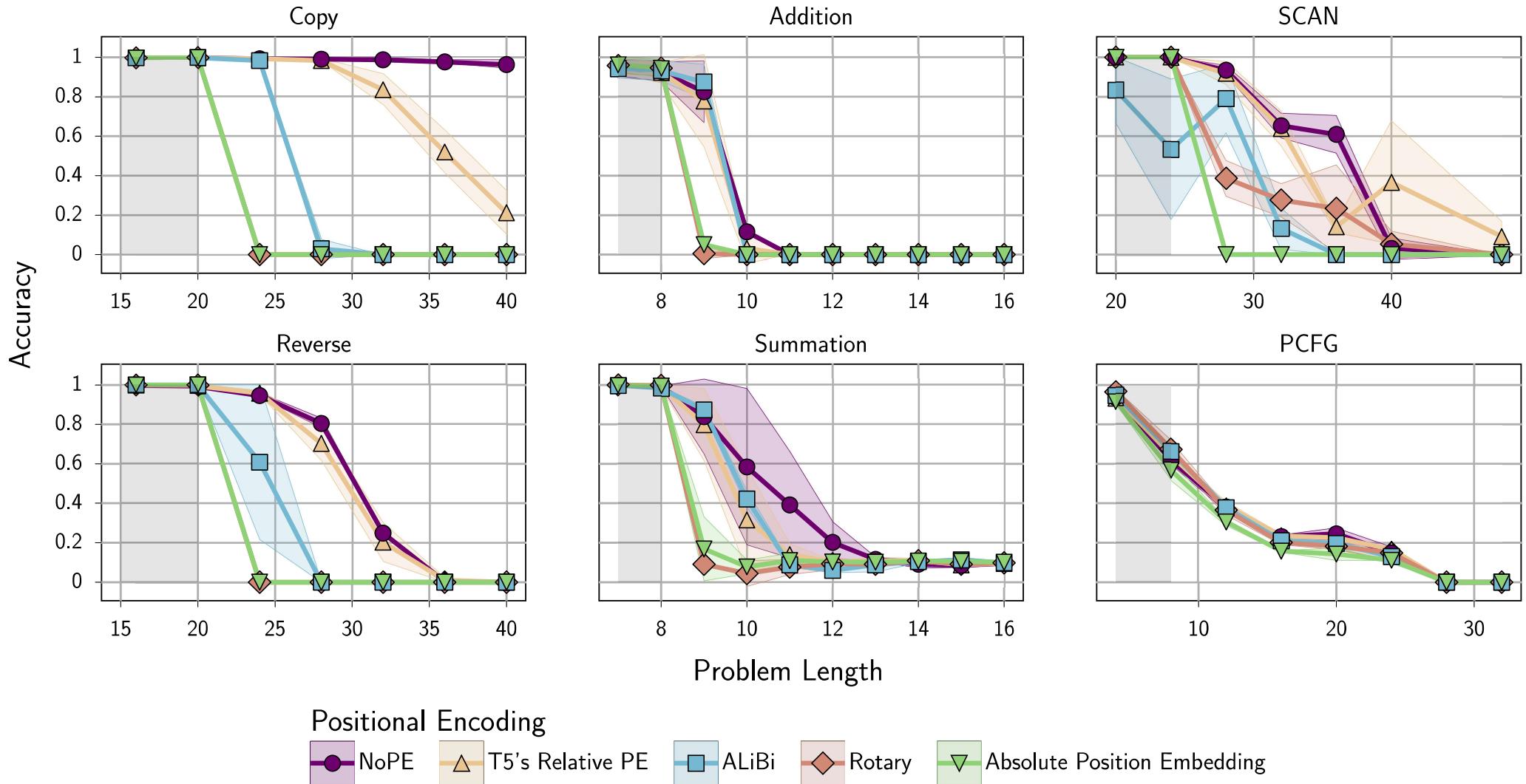
Length Generalization

Task	Input Example	Output Example
Primitive Tasks		
Copy	Copy the following words: <w1> <w2> <w3> <w4> <w5>	<w1> <w2> <w3> <w4> <w5>
Reverse	Reverse the following words: <w1> <w2> <w3> <w4> <w5>	<w5> <w4> <w3> <w2> <w1>
Mathematical and Algorithmic Tasks		
Addition	Compute: 5 3 7 2 6 + 1 9 1 7 ?	The answer is 5 5 6 4 3.
Polynomial Eval.	Evaluate x = 3 in (3 x ** 0 + 1 x ** 1 + 1 x ** 2) % 10 ?	The answer is 5.
Sorting	Sort the following numbers: 3 1 4 1 5 ?	The answer is 1 1 3 4 5.
Summation	Compute: (1 + 2 + 3 + 4 + 7) % 10 ?	The answer is 7.
Parity	Is the number of 1's even in [1 0 0 1 1] ?	The answer is No.
LEGO	If a = -1; b = -a; c = +b; d = +c. Then what is c?	The answer is +1.
Classical Length Generalization Datasets		
SCAN	jump twice and run left	JUMP JUMP TURN_LEFT RUN
PCFG	shift prepend K10 R1 K12 , E12 F16	F16 K10 R1 K12 E12

Length Generalization



Length Generalization



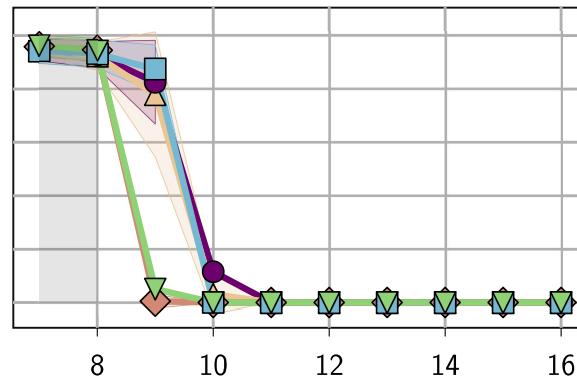
Length Generalization

Length Generalization

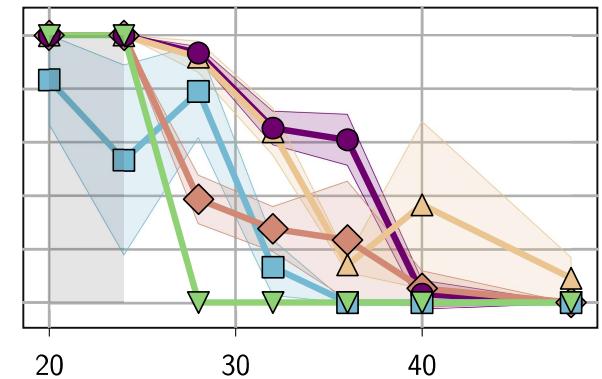


PE

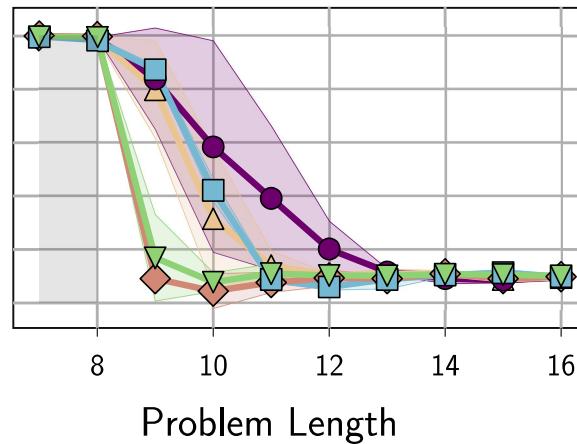
Addition



SCAN

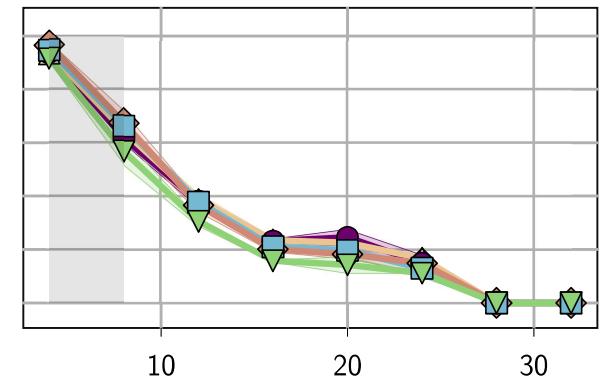


Summation



Problem Length

PCFG



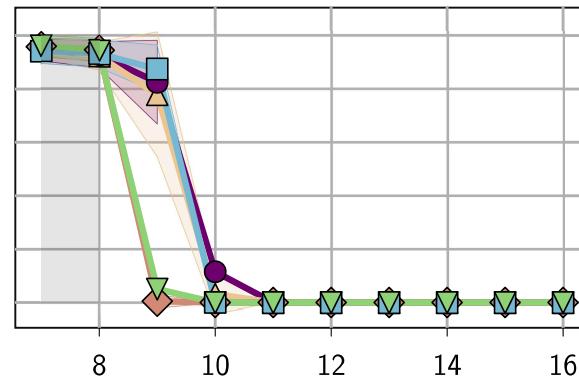
Length Generalization

Length Generalization

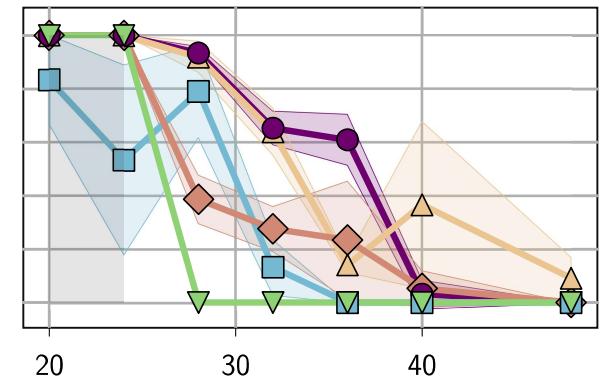


PE

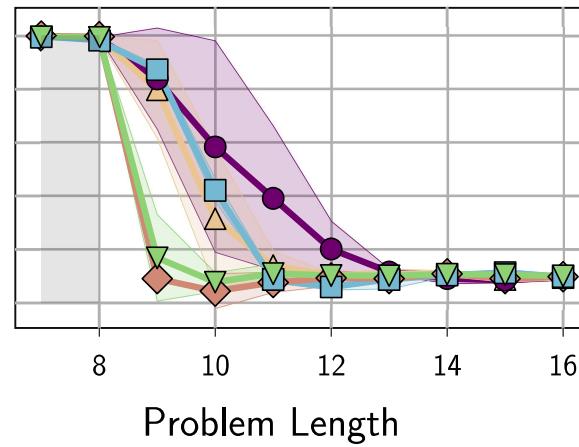
Addition



SCAN

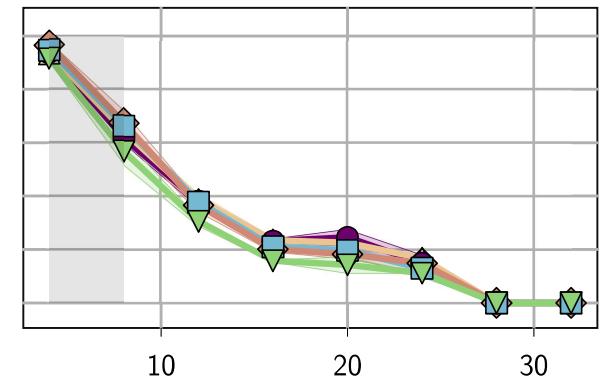


Summation

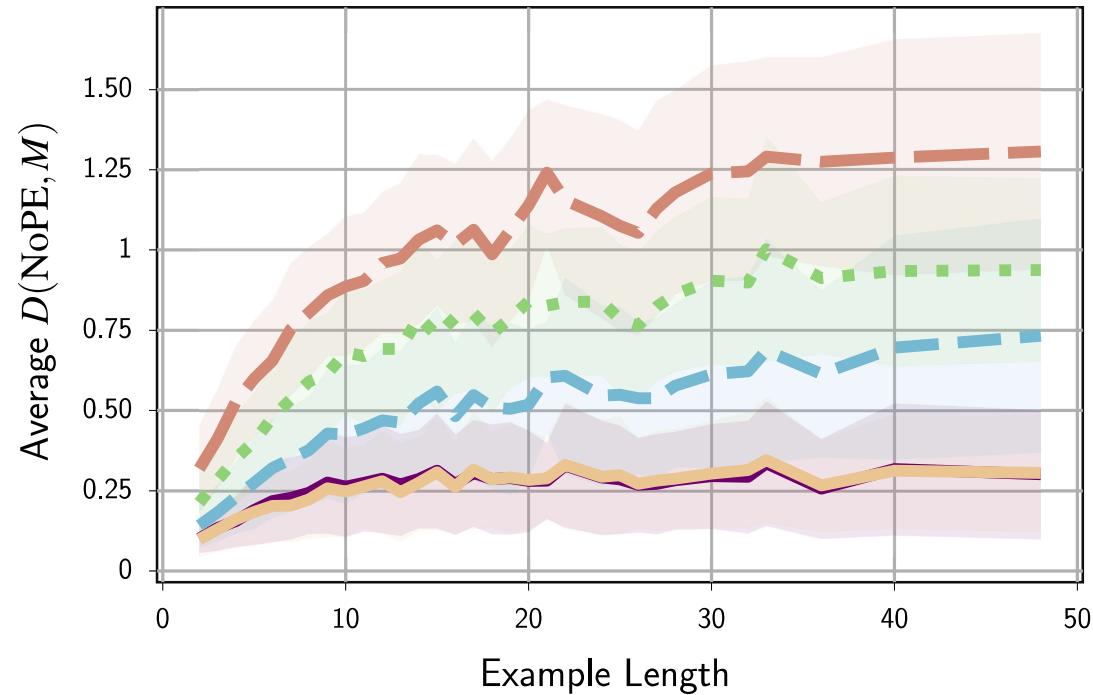
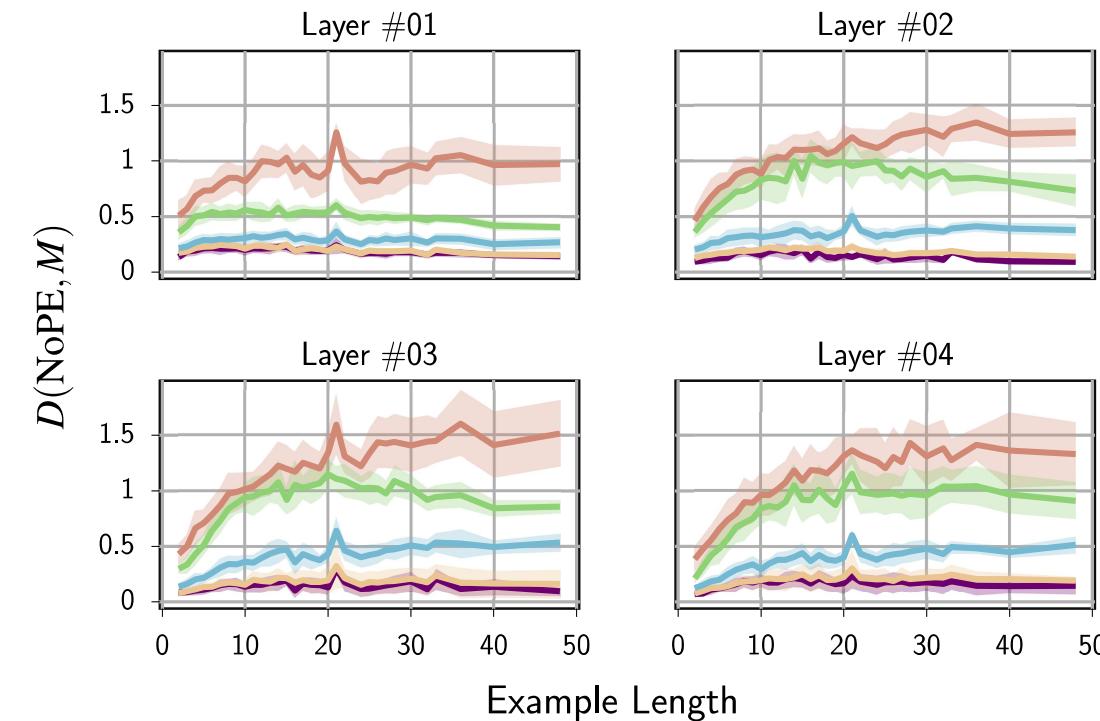


Problem Length

PCFG



NoPE 🤝 T5-rel



Positional Encoding (M)

NoPE'	T5's Relative PE	ALiBi	Rotary	Absolute Position Embedding
-------	------------------	-------	--------	-----------------------------

$$D_{\text{AT}}(\mathbf{P}, \mathbf{Q}) = \frac{1}{T} \sum_{t=1}^T D_{\text{JSD}}(\mathbf{P}_t || \mathbf{Q}_t)$$

$$D^{(l)}(\mathbf{A}, \mathbf{B}) = \min_{(\mathbf{P}, \mathbf{Q}) \in A_l \times B_l} D_{\text{AT}}(\mathbf{P}, \mathbf{Q})$$

NoPE still has abs-PE

input sequence $\mathbf{x} = [\langle \text{bos} \rangle, x_1, \dots, x_T]$

Theorem 1 (Absolute Encoding). *Let \mathbf{x} be an input sequence of length $T + 1$ to the model. Then, the first layer of f_θ can recover absolute positions $[1, \dots, T + 1]$ in the hidden state $\mathbf{H}^{(1)}$. That is, there exist \mathbf{W}_Q , \mathbf{W}_K , \mathbf{W}_V , \mathbf{W}_O , \mathbf{W}_1 , and \mathbf{W}_2 such that the self-attention and feedforward operations in the first layer compute absolute positions and write it to the next hidden state.*

NoPE still has abs-PE: Proof

$$\mathbf{W}_E = \begin{bmatrix} 1 & 1 & 1 & \dots & 1 \\ 1 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ e_{4,1} & e_{4,2} & e_{4,3} & \dots & e_{4,V} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ e_{d,1} & e_{d,2} & e_{d,2} & \dots & e_{d,V} \end{bmatrix}_{d \times V}$$

$$\mathbf{H}^{(0)} = \mathbf{W}_E \mathbf{X} = \begin{bmatrix} 1 & 1 & 1 & \dots & 1 \\ 1 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ e_{4,1} & e_{4,2} & e_{4,3} & \dots & e_{4,V} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ e_{d,1} & e_{d,2} & e_{d,2} & \dots & e_{d,V} \end{bmatrix}_{d \times (T+1)}$$

NoPE still has abs-PE: Proof

$$\mathbf{W}_K = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 1 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & \dots & 0 \end{bmatrix}_{h \times d}$$

$$\mathbf{k}_i = \mathbf{W}_K \mathbf{h}_i^{(0)}$$
$$\mathbf{k}_1 = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} \quad \mathbf{k}_2 = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} \quad \dots \quad \mathbf{k}_t = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}$$

NoPE still has abs-PE: Proof

$$\mathbf{W}_K = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 1 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & \dots & 0 \end{bmatrix}_{h \times d}$$

$$\mathbf{k}_i = \mathbf{W}_K \mathbf{h}_i^{(0)}$$
$$\mathbf{k}_1 = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} \quad \mathbf{k}_2 = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} \quad \dots \quad \mathbf{k}_t = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}$$

$$\mathbf{q}_t = \mathbf{W}_Q \mathbf{h}_t^{(0)}:$$

$$\mathbf{q}_t = [q_1, q_2, q_3, \dots, q_h]^\top$$

NoPE still has abs-PE: Proof

$$\mathbf{W}_K = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 1 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & \dots & 0 \end{bmatrix}_{h \times d} \quad \mathbf{k}_i = \mathbf{W}_K \mathbf{h}_i^{(0)}$$

$$\mathbf{k}_1 = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} \quad \mathbf{k}_2 = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} \quad \dots \quad \mathbf{k}_t = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}$$

$$\mathbf{q}_t = \mathbf{W}_Q \mathbf{h}_t^{(0)} :$$

$$\mathbf{q}_t = [q_1, q_2, q_3, \dots, q_h]^\top$$

$$\begin{aligned} \boldsymbol{\alpha} &= [\langle \mathbf{q}_t, \mathbf{k}_1 \rangle, \langle \mathbf{q}_t, \mathbf{k}_2 \rangle, \dots, \langle \mathbf{q}_t, \mathbf{k}_t \rangle]^\top \\ &= [\alpha^*, \alpha^*, \dots, \alpha^*]^\top \end{aligned}$$

NoPE still has abs-PE: Proof

$$\mathbf{W}_K = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 1 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & \dots & 0 \end{bmatrix}_{h \times d} \quad \mathbf{k}_i = \mathbf{W}_K \mathbf{h}_i^{(0)}$$

$$\mathbf{k}_1 = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} \quad \mathbf{k}_2 = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} \quad \dots \quad \mathbf{k}_t = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}$$

$$\mathbf{q}_t = \mathbf{W}_Q \mathbf{h}_t^{(0)} :$$

$$\mathbf{q}_t = [q_1, q_2, q_3, \dots, q_h]^\top$$

$$\begin{aligned} \boldsymbol{\alpha} &= [\langle \mathbf{q}_t, \mathbf{k}_1 \rangle, \langle \mathbf{q}_t, \mathbf{k}_2 \rangle, \dots, \langle \mathbf{q}_t, \mathbf{k}_t \rangle]^\top \\ &= [\alpha^*, \alpha^*, \dots, \alpha^*]^\top \end{aligned}$$

$$\hat{\boldsymbol{\alpha}} = \text{softmax}(\boldsymbol{\alpha}) = \left[\frac{1}{t}, \frac{1}{t}, \dots, \frac{1}{t} \right]^\top$$

NoPE still has abs-PE: Proof

$$\mathbf{W}_V = \begin{bmatrix} 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 0 \end{bmatrix}_{h \times d} \quad \mathbf{v}_i = \mathbf{W}_V \mathbf{h}_i^{(0)}$$
$$\mathbf{v}_1 = \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \quad \mathbf{v}_2 = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \quad \dots \quad \mathbf{v}_t = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

NoPE still has abs-PE: Proof

$$\mathbf{W}_V = \begin{bmatrix} 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 0 \end{bmatrix}_{h \times d} \quad \mathbf{v}_i = \mathbf{W}_V \mathbf{h}_i^{(0)}$$

$$\mathbf{v}_1 = \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \quad \mathbf{v}_2 = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \quad \dots \quad \mathbf{v}_t = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

$$\mathbf{W}_O = \begin{bmatrix} 0 & 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & 0 & \dots & 0 \\ 1 & 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & \dots & 0 \end{bmatrix}_{d \times h}$$

NoPE still has abs-PE: Proof

$$\mathbf{W}_V = \begin{bmatrix} 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 0 \end{bmatrix}_{h \times d} \quad \mathbf{v}_i = \mathbf{W}_V \mathbf{h}_i^{(0)}$$

$$\mathbf{v}_1 = \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \quad \mathbf{v}_2 = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \quad \dots \quad \mathbf{v}_t = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

$$\mathbf{W}_O = \begin{bmatrix} 0 & 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & 0 & \dots & 0 \\ 1 & 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & \dots & 0 \end{bmatrix}_{d \times h}$$

$$\mathbf{o}_t = \mathbf{W}_O \left(\sum_{i \leq t} \hat{\alpha}_i \mathbf{v}_i \right) = \mathbf{W}_O \left(\frac{1}{t} \sum_{i \leq t} \mathbf{v}_i \right) = \mathbf{W}_O \begin{pmatrix} 1/t \\ 0 \\ \vdots \\ 0 \end{pmatrix}_h = \begin{pmatrix} 0 \\ 0 \\ 1/t \\ 0 \\ \vdots \\ 0 \end{pmatrix}_d$$

NoPE still has rel-PE

Theorem 2 (Relative Encoding). Suppose that the hidden state $\mathbf{H}^{(1)}$ contains absolute positional information, as stated in [Theorem 1](#), and assume that it is not overwritten by any subsequent layers. Then, the self-attention in all subsequent layers can implement a relative positional encoding: there exists a parameterization of f_θ such that, for $l \geq 2$, the attention dot product between query \mathbf{q}_t and key \mathbf{k}_i at positions t and i ($t \geq i$) can be expressed as:

$$\langle \mathbf{q}_t, \mathbf{k}_i \rangle = f_{\text{cnt}}(\mathbf{q}_t, \mathbf{k}_i) + f_{\text{rel}}(t - i) \quad (1)$$

where f_{cnt} is a function of their content, and f_{rel} is a function of their relative distance.

NoPE still has rel-PE: Proof

$$\mathbf{H}^{(l)} = \begin{bmatrix} 1 & 1 & 1 & 1 & \dots & 1 \\ 1 & 0 & 0 & 0 & \dots & 0 \\ 1 & 2 & 3 & 4 & \dots & T+1 \\ h_{4,1} & h_{4,2} & h_{4,3} & h_{4,4} & \dots & h_{4,T+1} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ h_{d,1} & h_{d,2} & h_{d,3} & h_{d,4} & \dots & h_{d,T+1} \end{bmatrix}_{d \times (T+1)}$$

$$\mathbf{W}_Q = \begin{bmatrix} 1 & 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & -1 & 0 & \dots & 0 \\ w_{3,1} & w_{3,2} & w_{3,3} & w_{3,4} & \dots & w_{3,d} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ w_{h,1} & w_{h,2} & w_{h,3} & w_{h,4} & \dots & w_{h,d} \end{bmatrix}_{h \times d}$$

$\mathbf{q}_t = \mathbf{W}_Q \mathbf{h}_t^{(l)}$:
 $\mathbf{q}_t = [1, -t, q_3, \dots, q_h]^\top$

NoPE still has rel-PE: Proof

$$\mathbf{W}_V = \begin{bmatrix} 0 & 0 & 1 & 0 & \dots & 0 \\ 1 & 0 & 0 & 0 & \dots & 0 \\ w'_{3,1} & w'_{3,2} & w'_{3,3} & w'_{3,4} & \dots & w'_{3,d} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ w'_{h,1} & w'_{h,2} & w'_{h,3} & w'_{h,4} & \dots & w'_{h,d} \end{bmatrix}_{h \times d}$$

$$\mathbf{k}_1 = \begin{pmatrix} 1 \\ 1 \\ k_{3,1} \\ \vdots \\ k_{h,1} \end{pmatrix} \quad \mathbf{k}_2 = \begin{pmatrix} 2 \\ 1 \\ k_{3,2} \\ \vdots \\ k_{h,2} \end{pmatrix} \quad \mathbf{k}_3 = \begin{pmatrix} 3 \\ 1 \\ k_{3,3} \\ \vdots \\ k_{h,3} \end{pmatrix} \quad \dots \quad \mathbf{k}_t = \begin{pmatrix} t \\ 1 \\ k_{3,t} \\ \vdots \\ k_{h,t} \end{pmatrix}$$

$$\mathbf{k}_i = [i, 1, k_{3,i}, \dots, k_{h,i}]^\top$$

NoPE still has rel-PE: Proof

$$\mathbf{q}_t = [1, -t, q_3, \dots, q_h]^\top \quad \mathbf{k}_i = [i, 1, k_{3,i}, \dots, k_{h,i}]^\top$$

$$\begin{aligned}\langle \mathbf{q}_t, \mathbf{k}_i \rangle &= 1 \cdot i + (-t) \cdot 1 + q_3 \cdot k_{3,i} + \dots + q_h \cdot k_{h,i} \\ &= i - t + \sum_{j=3}^h q_j \cdot k_{j,i} \\ &= \left(\sum_{j=3}^h q_j \cdot k_{j,i} \right) - (t - i) \\ &= f_{\text{cnt}}(\mathbf{q}_t, \mathbf{k}_i) + f_{\text{rel}}(t - i)\end{aligned}$$

CoT

Input (x)
 Compute $5\ 3\ 7\ 2\ 6 + 1\ 9\ 1\ 7 =$

Output ($s; y$)
 <scratch>

\mathcal{I} For digits 6 and 7,

\mathcal{C} We have $(6 + 7 + \text{carry}) \% 10 = 13 \% 10 = 13 \% 10$

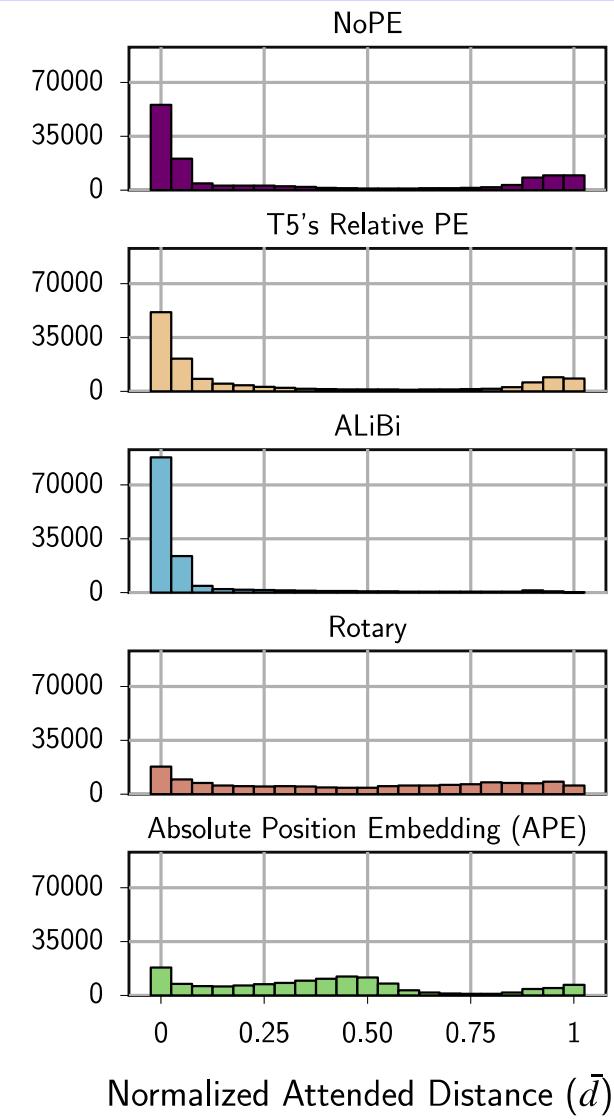
\mathcal{O} Which is equal to 3 .

\mathcal{V} We update carry to $13 // 10 = 1$.

\mathcal{R} So, the remaining input is
 $5\ 3\ 7\ 2 + 1\ 9\ 1$

...

</scratch>
 The answer is 5 5 6 4 3.



ACL2022 Findings Meta

Transformer Language Models without Positional Encodings Still Learn Positional Information

Adi Haviv^τ Ori Ram^τ Ofir Press^ω Peter Izsak^λ Omer Levy^{τμ}

^τTel Aviv University ^ωUniversity of Washington ^λIntel Labs ^μMeta AI

ACL2023 Honorable Mentions CMU&Princeton

Latent Positional Information is in the Self-Attention Variance of Transformer Language Models Without Positional Embeddings

Ta-Chung Chi[†]

Carnegie Mellon University

Ting-Han Fan

Princeton University

Li-Wei Chen

Carnegie Mellon University

Alexander I. Rudnicky

Carnegie Mellon University

Peter J. Ramadge

Princeton University

Preprint Mila

The Impact of Positional Encoding on Length Generalization in Transformers

Amirhossein Kazemnejad^{1,2}, Inkit Padhi³
Karthikeyan Natesan Ramamurthy³, Payel Das³, Siva Reddy^{1,2,4}

¹Mila - Québec AI Institute; ²McGill University;

³IBM Research; ⁴Facebook CIFAR AI Chair

ACL2022 Findings Meta

Transformer Language Models without Positional Encodings Still Learn Positional Information

Adi Haviv^τ Ori Ram^τ Ofir Press^ω Peter Izsak^λ Omer Levy^{τμ}

^τTel Aviv University ^ωUniversity of Washington ^λIntel Labs ^μMeta AI

NoPE 

ACL2023 Honorable Mentions CMU&Princeton

Latent Positional Information is in the Self-Attention Variance of Transformer Language Models Without Positional Embeddings

Ta-Chung Chi [†] Carnegie Mellon University	Ting-Han Fan Princeton University	Li-Wei Chen Carnegie Mellon University
Alexander I. Rudnicky Carnegie Mellon University	Peter J. Ramadge Princeton University	

NoPE 

Preprint Mila

The Impact of Positional Encoding on Length Generalization in Transformers

Amirhossein Kazemnejad^{1,2}, Inkit Padhi³
Karthikeyan Natesan Ramamurthy³, Payel Das³, Siva Reddy^{1,2,4}
¹Mila - Québec AI Institute; ²McGill University;
³IBM Research; ⁴Facebook CIFAR AI Chair

NoPE 

ACL2022 Findings Meta

Transformer Language Models without Positional Encodings Still Learn Positional Information

Adi Haviv^τ Ori Ram^τ Ofir Press^ω Peter Izsak^β Omer Levy^{τμ}

^τTel Aviv University ^ωUniversity of Washington ^βIntel Labs ^μMeta AI

NoPE 

ACL2023 Honorable Mentions CMU&Princeton

Latent Positional Information is in the Self-Attention Variance of Transformer Language Models Without Positional Embeddings

Ta-Chung Chi[†]
Carnegie Mellon University

Ting-Han Fan
Princeton University

Li-Wei Chen
Carnegie Mellon University

Alexander I. Rudnicky
Carnegie Mellon University

Peter J. Ramadge
Princeton University

NoPE 

why? 

Preprint Mila

The Impact of Positional Encoding on Length Generalization in Transformers

Amirhossein Kazemnejad^{1,2}, Inkit Padhi³
Karthikeyan Natesan Ramamurthy³, Payel Das³, Siva Reddy^{1,2,4}

¹Mila - Québec AI Institute; ²McGill University;

³IBM Research; ⁴Facebook CIFAR AI Chair

NoPE 

why? 

ACL2022 Findings Meta

Transformer Language Models without Positional Encodings Still Learn Positional Information

Adi Haviv^τ Ori Ram^τ Ofir Press^ω Peter Izsak^β Omer Levy^{τμ}

^τTel Aviv University ^ωUniversity of Washington ^βIntel Labs ^μMeta AI

NoPE 

ACL2023 Honorable Mentions CMU&Princeton

Latent Positional Information is in the Self-Attention Variance of Transformer Language Models Without Positional Embeddings

Ta-Chung Chi[†]
Carnegie Mellon University

Ting-Han Fan
Princeton University

Li-Wei Chen
Carnegie Mellon University

Alexander I. Rudnicky
Carnegie Mellon University

Peter J. Ramadge
Princeton University

NoPE 

why? 

Preprint Mila

The Impact of Positional Encoding on Length Generalization in Transformers

Amirhossein Kazemnejad^{1,2}, Inkit Padhi³
Karthikeyan Natesan Ramamurthy³, Payel Das³, Siva Reddy^{1,2,4}

¹Mila - Québec AI Institute; ²McGill University;

³IBM Research; ⁴Facebook CIFAR AI Chair

NoPE 

why? 

 Task

ACL2022 Findings Meta

Transformer Language Models without Positional Encodings Still Learn Positional Information

Adi Haviv^τ Ori Ram^τ Ofir Press^ω Peter Izsak^β Omer Levy^{τμ}

^τTel Aviv University ^ωUniversity of Washington ^βIntel Labs ^μMeta AI

NoPE 

ACL2023 Honorable Mentions CMU&Princeton

Latent Positional Information is in the Self-Attention Variance of Transformer Language Models Without Positional Embeddings

Ta-Chung Chi[†]
Carnegie Mellon University

Ting-Han Fan
Princeton University

Li-Wei Chen
Carnegie Mellon University

Alexander I. Rudnicky
Carnegie Mellon University

Peter J. Ramadge
Princeton University

NoPE 

why? 

Preprint Mila

The Impact of Positional Encoding on Length Generalization in Transformers

Amirhossein Kazemnejad^{1,2}, Inkit Padhi³
Karthikeyan Natesan Ramamurthy³, Payel Das³, Siva Reddy^{1,2,4}

¹Mila - Québec AI Institute; ²McGill University;

³IBM Research; ⁴Facebook CIFAR AI Chair

NoPE 

why? 

Task 

