CMPT 413/713: Natural Language Processing

# Modern LLM Architecture

Spring 2026

2026-02-26

# Lots of developments to the transformer block!



Transformers [2017]

https://magazine.sebastianraschka.com/p/the-big-llm-architecture-comparison

# PaLM

*PaLM: Scaling Language Modeling with Pathways, Chowdhery et al, Google, 2022*
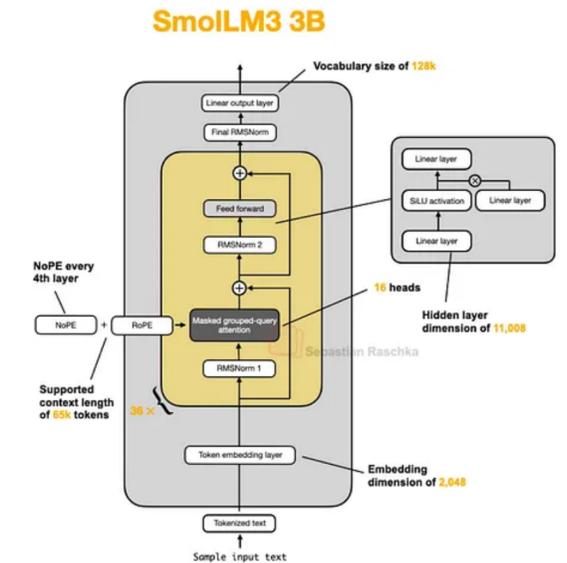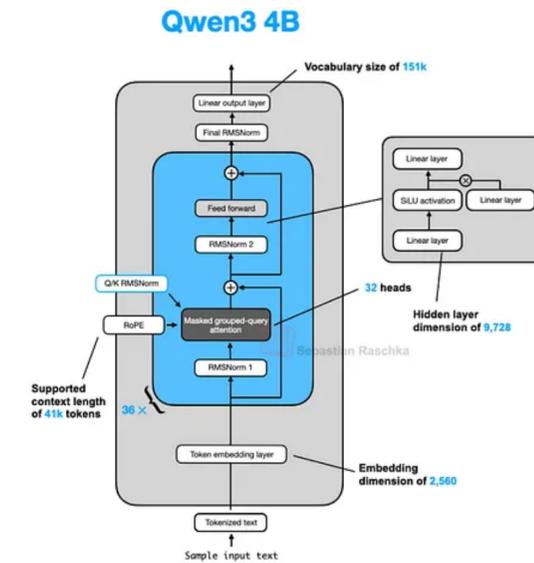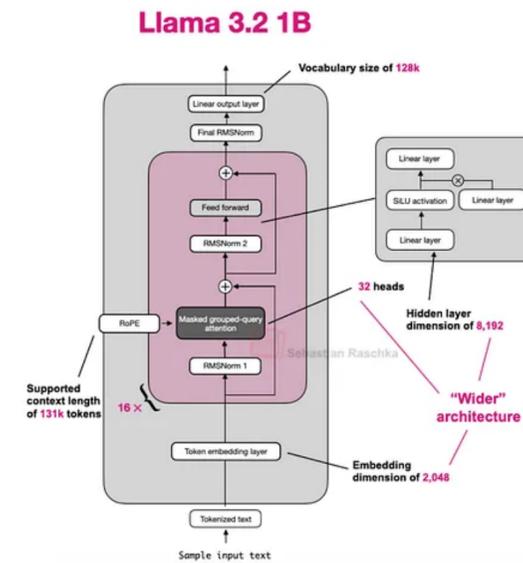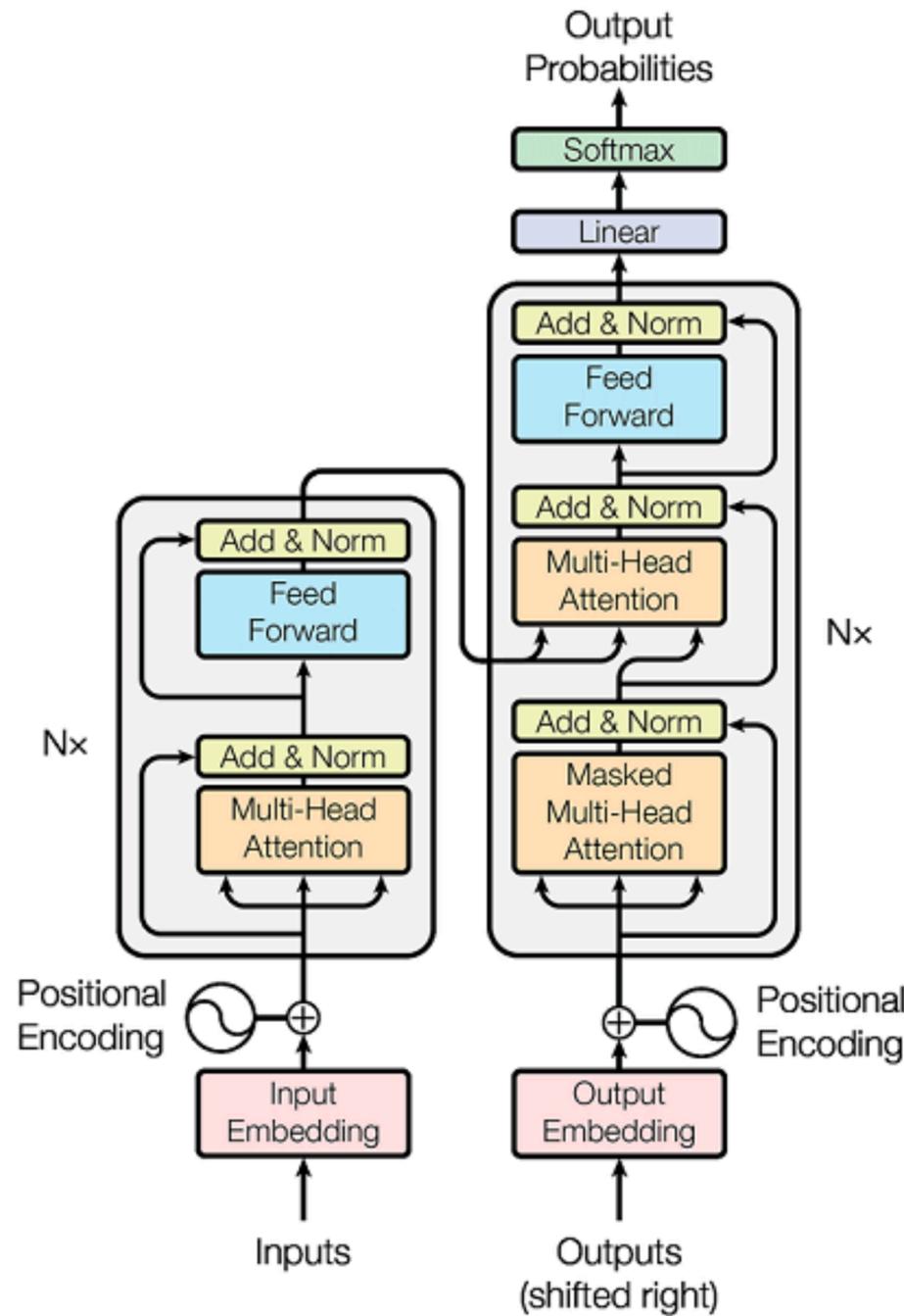
## Architecture

- SwiGLU activation: $\text{Swish}(xW) \otimes xV$
- Parallel layers
  - Serial: $y = x + \text{MLP}(\text{LayerNorm}(x + \text{Attention}(\text{LayerNorm}(x)))$
  - Parallel: $y = x + \text{MLP}(\text{LayerNorm}(x)) + \text{Attention}(\text{LayerNorm}(x))$
    - 15% faster training speed (degradation for small models 8B, but no degradation at 62B)
- Attention: Shared key-value across heads, query is still separately projected per head
- RoPE (rotary position) embeddings
- Shared input-output embeddings
- No biases: increased training stability
- Vocabulary: SentencePiece with 256k tokens

## Training data

- 780 billion tokens of natural language + source code from github

# LLaMa

| | Vaswani et al. | LLaMa |
|---|---|---|
| **Norm Position** | Post | Pre |
| **Norm Type** | LayerNorm | RMSNorm |
| **Non-linearity** | ReLU | SwiGLU |
| **Positional Encoding** | Sinosoidal | RoPE |
| **Attention** | Full Multi-Head Attention | Grouped Multi-Query Attention |



TRANSFORMERS ARCHITECTURE vs. LLAMA ARCHITECTURE

# Revisiting transformers

- Positional encoding

- Normalization

  - Layer normalization vs RMS

  - Post vs Pre-Layer norm

- Activation functions

- Attention variations

# Revisiting transformers

- **Positional encoding**

- Normalization

  - Layer normalization vs RMS

  - Post vs Pre-Layer norm

- Activation functions

- Attention variations

# Positional encoding

- Original transformer: fixed sinusoidal absolute embeddings

- Learned encoding

- Absolute vs relative

  - In most cases, it is the relative position between two words that matter (not their absolute position)

  - Relative encoding can be learned [Self-Attention with Relative Position Representations, Shaw et al. 2018]

- Rotary embeddings (RoPE)

# Learned encoding

- Advantage: Flexible, learned representations

- Disadvantage: bunch of extra parameters that need to be learned

- Disadvantage: impossible to extrapolate to longer sequences

# Learned encoding

## What do position embeddings learn?

- Visualize cosine similarity between position embeddings

- GPT-2 learned embeddings are quite good: can effectively predict absolute position using linear regression and relative ordering using logistic regression



What Do Position Embeddings Learn? [Wang and Chen 2020]

# Learned encoding

## What do position embeddings learn?

- Visualize cosine similarity between position embeddings

- GPT-2 learned embeddings are quite good: can effectively predict absolute position using linear regression and relative ordering using logistic regression

Absolute

| Type | PE | MAE |
|------|------|------|
| Learned | BERT | 34.14 |
| | RoBERTa | 6.06 |
| | GPT-2 | 1.03 |
| Pre-Defined | sinusoid | 0.0 |

Relative

| Type | PE | Error Rate |
|------|------|------|
| Learned | BERT | 19.72% |
| | RoBERTa | 7.23% |
| | GPT-2 | 1.56% |
| Pre-Defined | sinusoid | 5.08% |

What Do Position Embeddings Learn? [Wang and Chen 2020]

# Relative encoding

- Learnable relative embeddings

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

Self-Attention with Relative Position Representations
[Shaw et al. 2018]

- Modify attention scores to capture relative embedding

$$\boldsymbol{q}_m^\mathsf{T} \boldsymbol{k}_n = \boldsymbol{x}_m^\mathsf{T} \boldsymbol{W}_q^\mathsf{T} \boldsymbol{W}_k \boldsymbol{x}_n + \boldsymbol{x}_m^\mathsf{T} \boldsymbol{W}_q^\mathsf{T} \boldsymbol{W}_k \boldsymbol{p}_n + \boldsymbol{p}_m^\mathsf{T} \boldsymbol{W}_q^\mathsf{T} \boldsymbol{W}_k \boldsymbol{x}_n + \boldsymbol{p}_m^\mathsf{T} \boldsymbol{W}_q^\mathsf{T} \boldsymbol{W}_k \boldsymbol{p}_n$$

- Simplify to just learning a bias term

$$\boldsymbol{q}_m^\mathsf{T} \boldsymbol{k}_n = \boldsymbol{x}_m^\mathsf{T} \boldsymbol{W}_q^\mathsf{T} \boldsymbol{W}_k \boldsymbol{x}_n + b_{i,j}$$

Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer [Raffel et al. 2018]

# Attention with Linear Biases (ALiBi)

- No explicit position embedding

- Bias query-key attention scores with fixed penalty that is proportional to the distance

- Allows for better extrapolation to long sequences at test time



Train Short, Test Long: Attention with Linear Biases Enables Input Length Extrapolation [Press et al. 2021]

# Attention with Linear Biases (ALiBi)



Train Short, Test Long: Attention with Linear Biases Enables Input Length Extrapolation [Press et al. 2021]

# Rotary encoding

- Design absolute embeddings so the dot product result in function of relative position

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

- **Rotary Position Embedding (RoPE)**: Apply rotation to encode positional encoding (vs using addition).

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

- For dimension $d$, break down into $d/2$ 2D subspace.

  - For vector of length $d$, group into $d/2$ pairs of numbers - each pair is a vector in 2D.

  - Encode positional encoding $m$ on each pair by rotating by an angle $m\theta$

RoFormer: Enhanced Transformer with Rotary Position Embedding [Su et al. 2021]

# Rotary encoding

- Design absolute embeddings so the dot product result in function of relative position

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

- **Rotary Position Embedding (RoPE)**: Apply rotation to encode positional encoding (vs using addition).
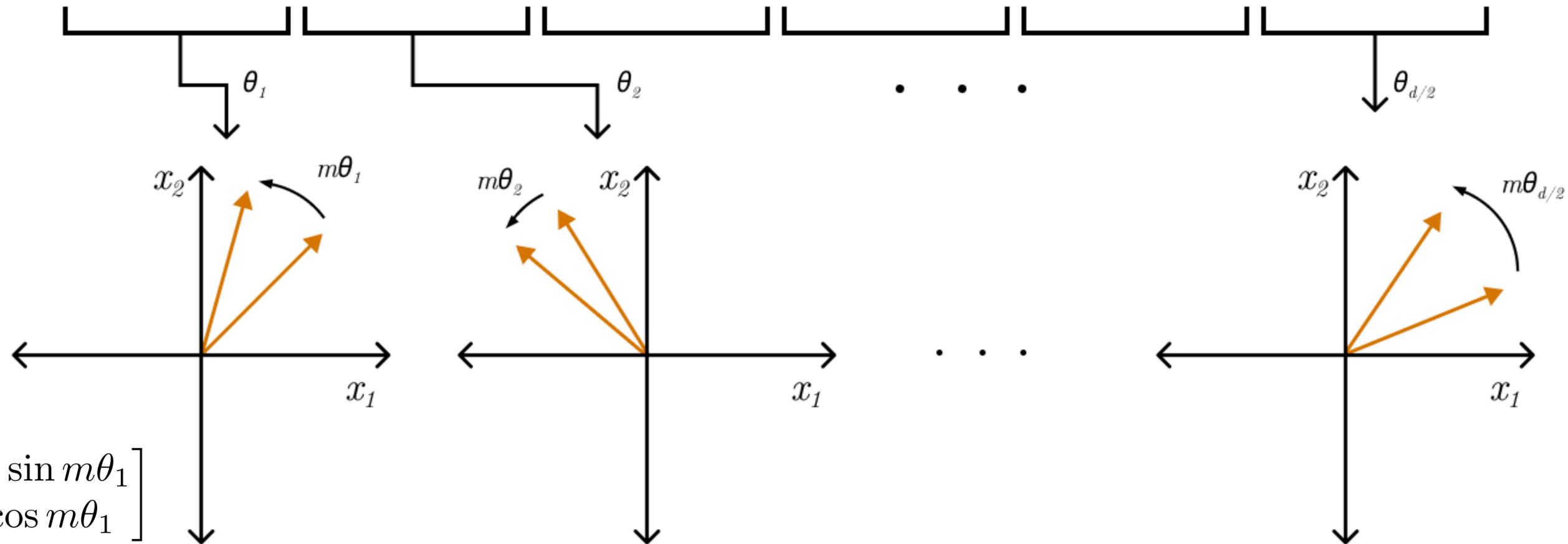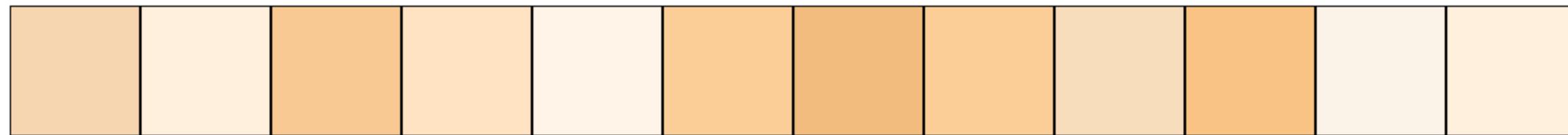
$$f_{\{q,k\}}(\mathbf{x}_m, m) = \mathbf{R}_{\Theta,m}^d \mathbf{W}_{\{q,k\}} \mathbf{x}_m$$

$$R_{\Theta,m}^d = \begin{bmatrix} \cos m\theta_1 & -\sin m\theta_1 & 0 & 0 & \cdots & 0 & 0 \\ \sin m\theta_1 & \cos m\theta_1 & 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \cos m\theta_2 & -\sin m\theta_2 & \cdots & 0 & 0 \\ 0 & 0 & \sin m\theta_2 & \cos m\theta_2 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & \cos m\theta_{\frac{d}{2}} & -\sin m\theta_{\frac{d}{2}} \\ 0 & 0 & 0 & 0 & \cdots & \sin m\theta_{\frac{d}{2}} & \cos m\theta_{\frac{d}{2}} \end{bmatrix}$$

RoFormer: Enhanced Transformer with Rotary Position Embedding [Su et al. 2021]

# Rotary encoding



$$\begin{bmatrix} \cos m\theta_1 & -\sin m\theta_1 \\ \sin m\theta_1 & \cos m\theta_1 \end{bmatrix}$$

# Rotary encoding

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

## Why it works?

- For dimension $d$, break down into $d/2$ 2D subspace.

  - For vector of length $d$, group into $d/2$ pairs of numbers - each pair is a vector in 2D.

  - Encode positional encoding $m$ on each pair by rotating by an angle $m\theta$

$$R_{\Theta,m}^{d} = \begin{bmatrix} \cos m\theta_1 & -\sin m\theta_1 & 0 & 0 & \cdots & 0 & 0 \\ \sin m\theta_1 & \cos m\theta_1 & 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \cos m\theta_2 & -\sin m\theta_2 & \cdots & 0 & 0 \\ 0 & 0 & \sin m\theta_2 & \cos m\theta_2 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & \cos m\theta_{\frac{d}{2}} & -\sin m\theta_{\frac{d}{2}} \\ 0 & 0 & 0 & 0 & \cdots & \sin m\theta_{\frac{d}{2}} & \cos m\theta_{\frac{d}{2}} \end{bmatrix}$$

RoFormer: Enhanced Transformer with Rotary Position Embedding [Su et al. 2021]
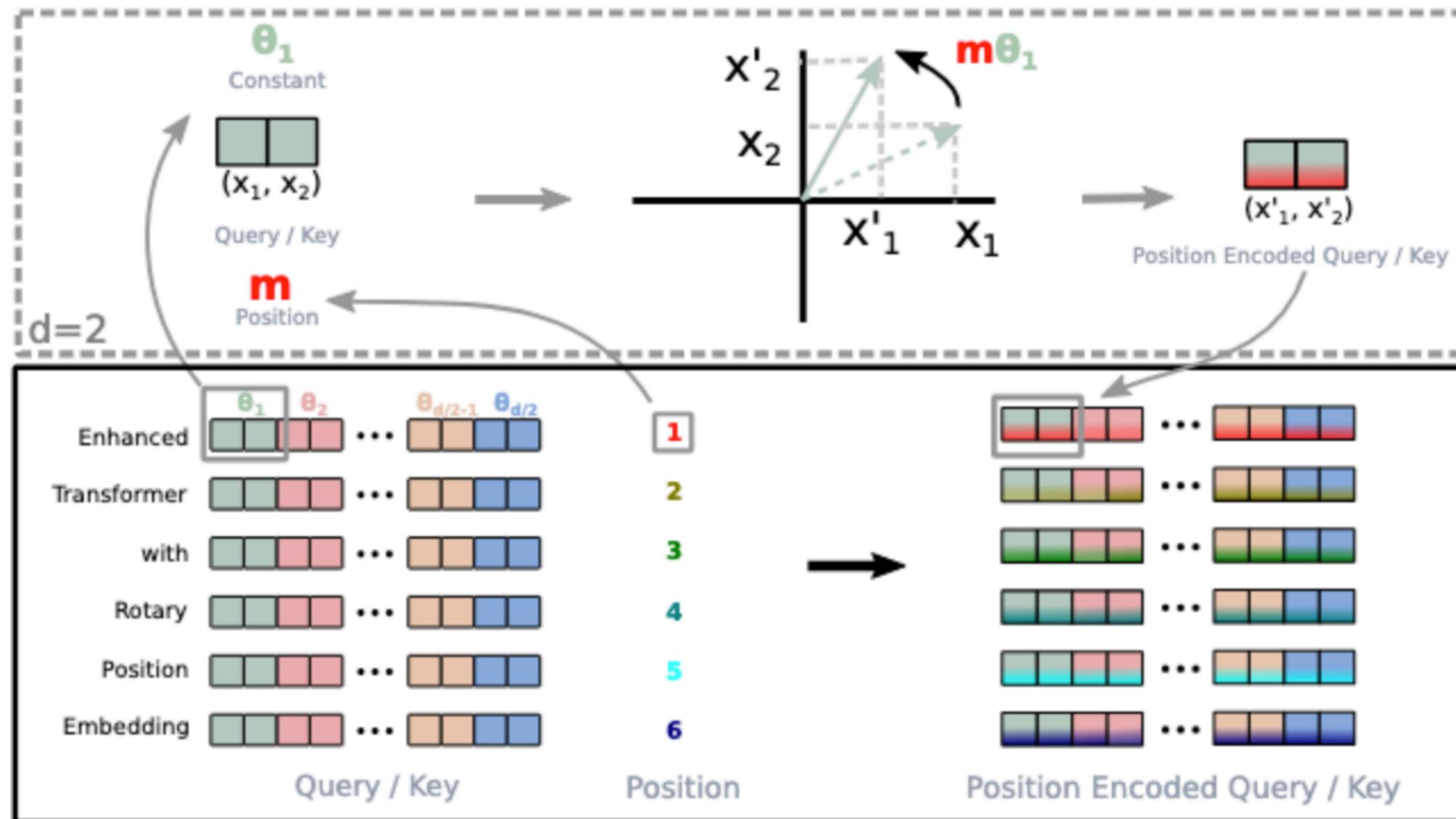
# Rotary encoding
## More efficient form

- With just element wise multiply and addition

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

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

RoFormer: Enhanced Transformer with Rotary Position Embedding [Su et al. 2021]

# Rotary encoding



RoFormer: Enhanced Transformer with Rotary Position Embedding [Su et al. 2021]

# Rotary encoding
## Benefits

- Absolute encoding that captures relative information

- Computationally efficient

- Integrates well into KV cache

RoFormer: Enhanced Transformer with Rotary Position Embedding [Su et al. 2021]

# Rotary encoding

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

## Dot product

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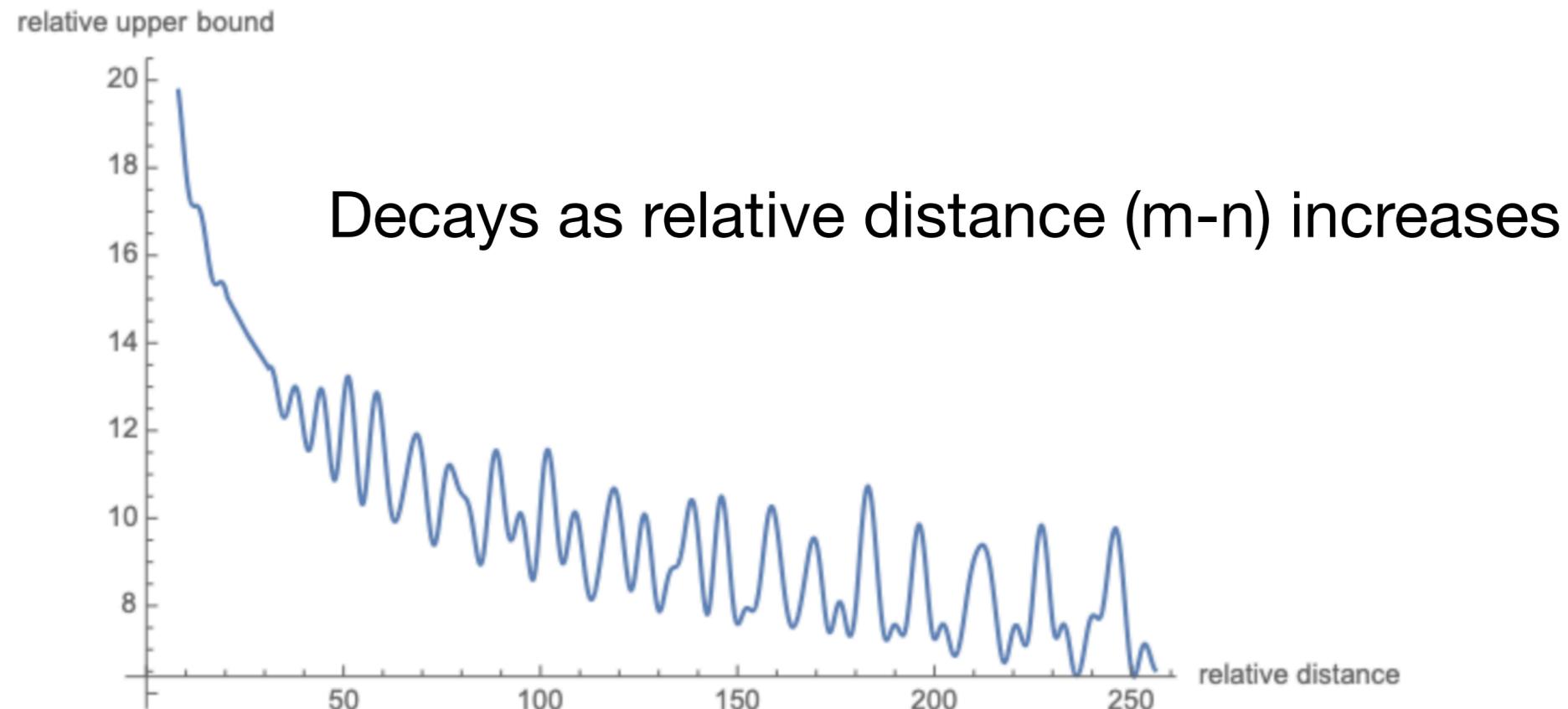


Decays as relative distance (m-n) increases

Figure 2: Long-term decay of RoPE.

RoFormer: Enhanced Transformer with Rotary Position Embedding [Su et al. 2021]

# Current LLMs use RoPE

- Advantage of RoPE may not hold for large models

| Model size | Pos. Emb. | Performance | | |
|---|---|---|---|---|
| | | zs-main ↑ | zs-small ↑ | ppl-pile ↓ |
| *Likely* threshold (1-$\sigma$) | | ±1.1 | ±0.4 | ±0.002 |
| 1B | ALiBi | 49.2 | 43.1 | 0.895 |
| | URPE | 49.6 | 43.1 | 0.885 |
| | RoPE | **50.0** | **44.2** | **0.883** |
| 3B | ALiBi | **54.4** | 49.8 | 0.807 |
| | RoPE | **54.4** | **50.5** | **0.799** |

The Falcon Series of Open Language Models
[Almazrouei et al. 2023]

- But can depend on value of theta

  - Recent models use $\theta = 500K$ (vs initial $\theta = 10K$) for long context

    LLaMa, OLMo

# Comparing RoPE variants

Adjusted base frequency (ABF)

- RoPE with $\theta = 10K$ vs $\theta = 500K$

| PE | Books | CC | Wikipedia |
|---|---|---|---|
| RoPE | 6.548 | 6.816 | 3.802 |
| RoPE PI | 6.341 | 6.786 | 3.775 |
| RoPE ABF | **6.323** | **6.780** | **3.771** |
| xPos ABF | 6.331 | **6.780** | **3.771** |

Effective Long-Context Scaling of Foundation Models
[Xiong et al. 2023]

- Position interpolation (PI)

$$\mathbf{f}'(\mathbf{x}, m) = \mathbf{f}\left(\mathbf{x}, \frac{mL}{L'}\right)$$

Position interpolation:
Rescale from new length $L'$
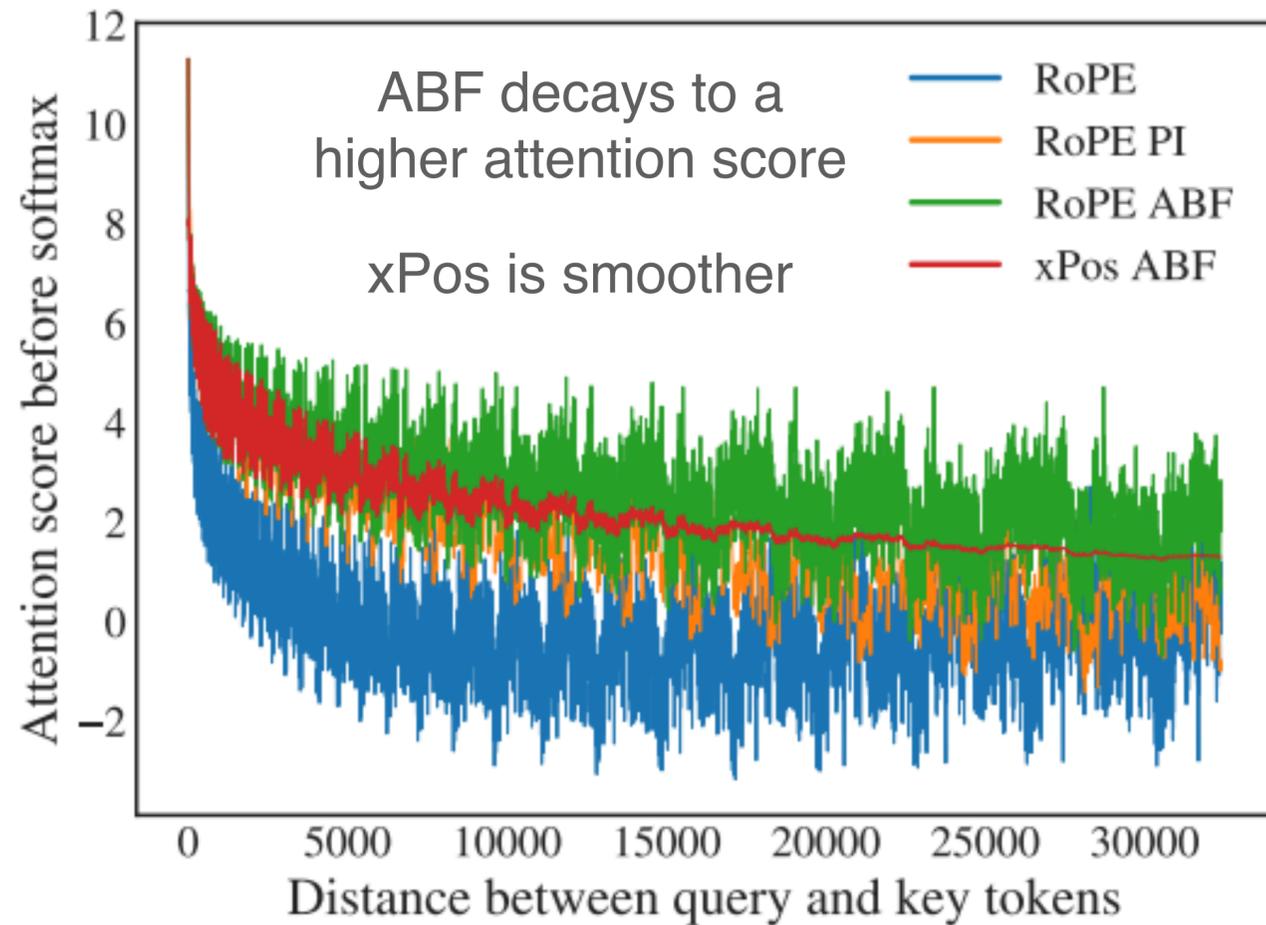to original length $L$

- xPos - Variant of RoPE

A length-extrapolatable transformer [Sun et al. 2022]

Extending Context Window of Large Language
Models via Positional Interpolation [Chen et al. 2023]

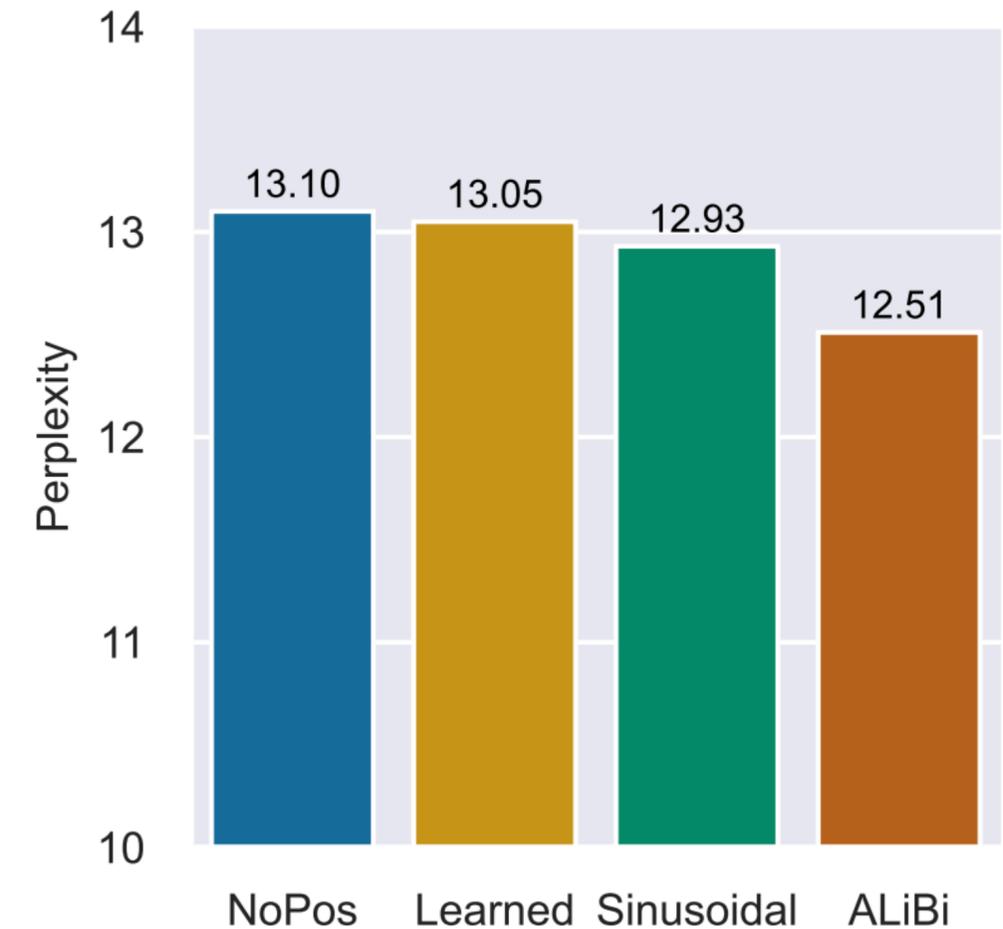# Current LLMs use RoPE

Adjusted base frequency (ABF)

- Comparing RoPE with $\theta = 10K$ vs $\theta = 500K$



First sentence retrieval

Effective Long-Context Scaling of Foundation Models [Xiong et al. 2023]

# How important is positional encoding?
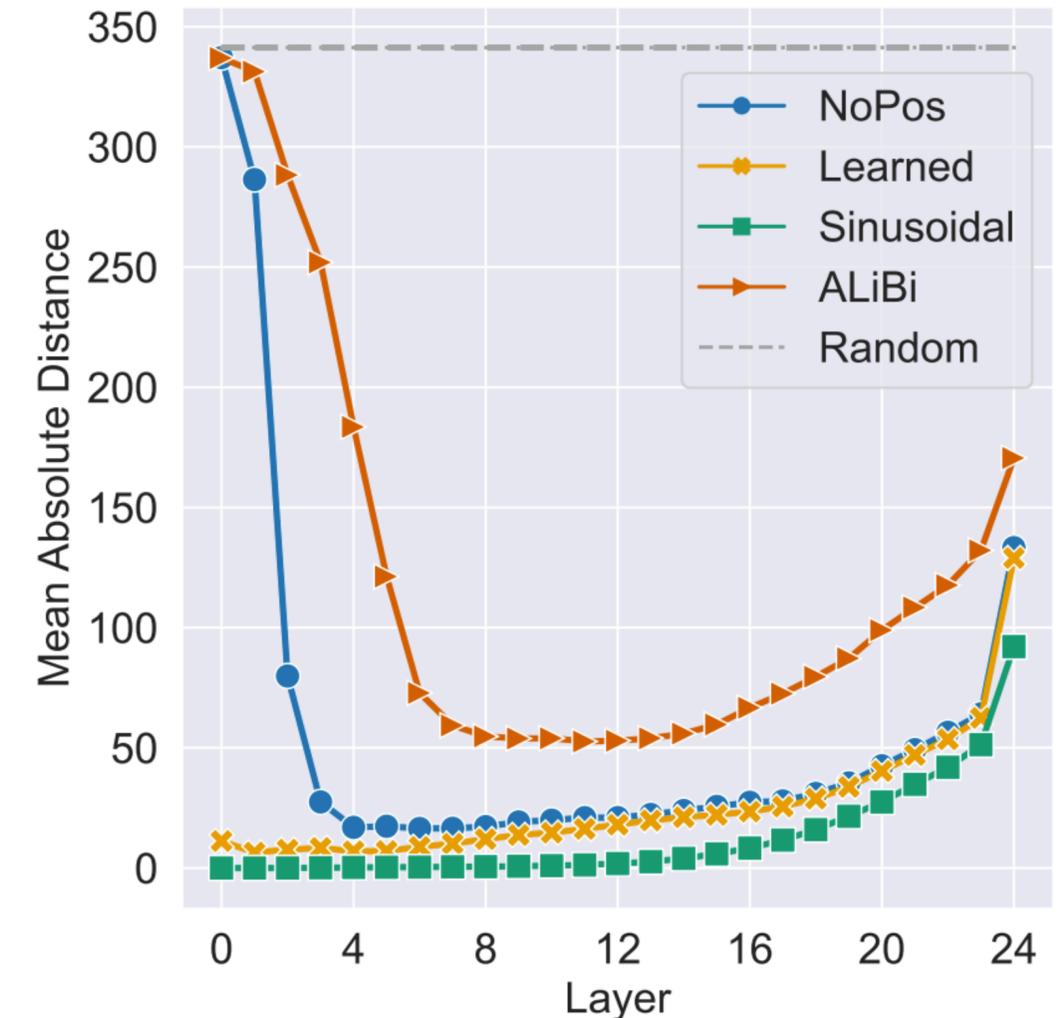
## Can we remove them?

- Train 1.3B parameter autoregressive LM on the Pile

- LM perplexity comparable to learned embedding

- Hypothesis: causal attention allows for encoding of absolute position of tokens



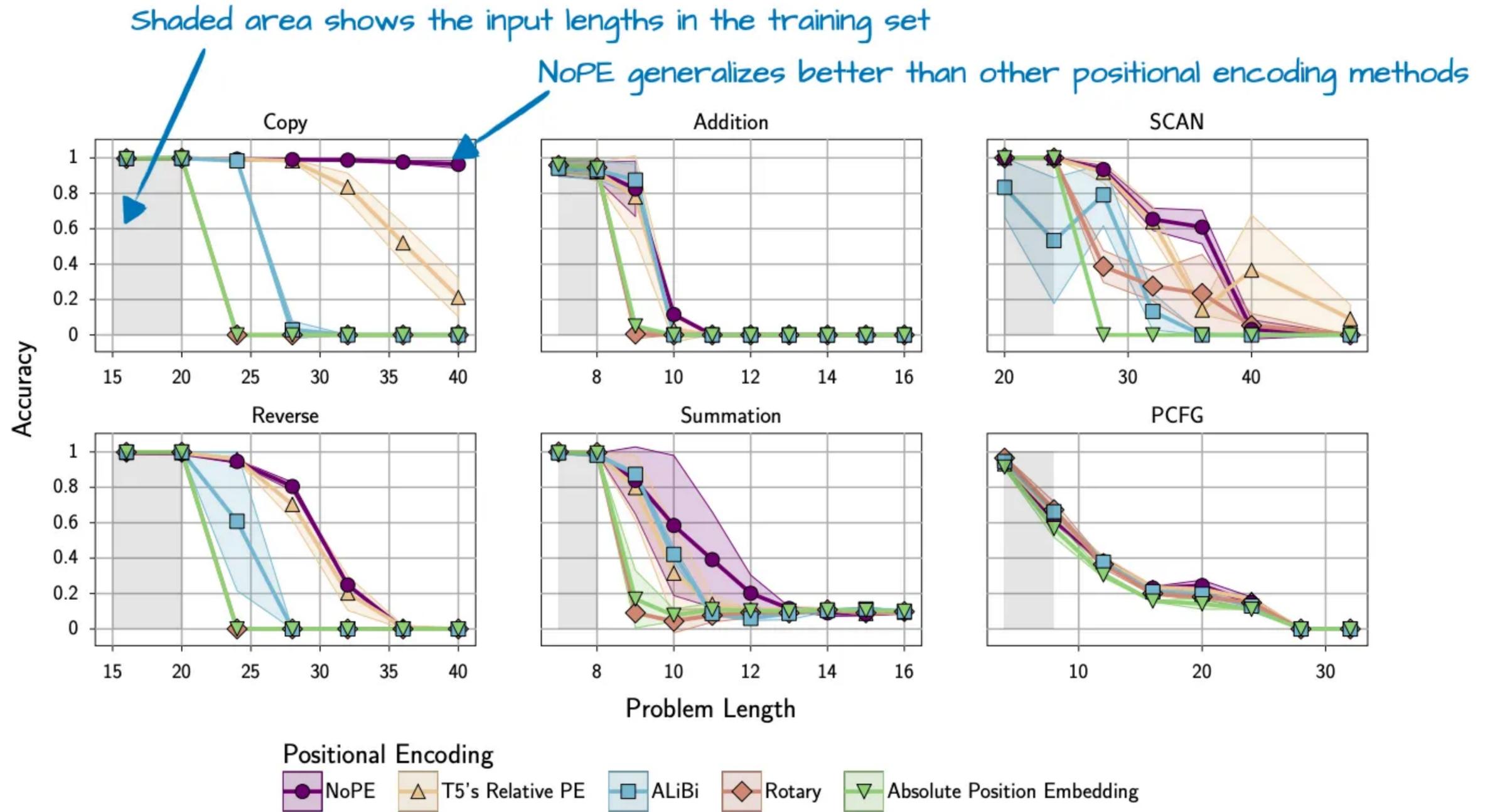Transformer Language Models without Positional Encodings Still Learn Positional Information [Haviv et al. 2022]

# How important is positional encoding?

## Can we remove them?

- Model can learn to encode absolute position

- Train 2-layer FF ReLU network (hooked up to transformer layer output representations) to predict absolute position (0 to 1023) of each token (as multi class classification)

- LM weights are frozen

- Later layers can predict position



Transformer Language Models without Positional Encodings Still Learn Positional Information [Haviv et al. 2022]

# NoPE (no positional embedding)

- Interleave NoPE layers with layers with positional embeddings



Shaded area shows the input lengths in the training set

NoPE generalizes better than other positional encoding methods

Positional Encoding: NoPE, T5's Relative PE, ALiBi, Rotary, Absolute Position Embedding

*The Impact of Positional Encoding on Length Generalization in Transformers* [Kazemnejad et al. 2023]
(Annotated version from https://magazine.sebastianraschka.com/p/the-big-llm-architecture-comparison)

# Revisiting transformers

- Positional encoding

- **Normalization**

  - Layer normalization vs RMS

  - Post vs Pre-Layer norm

- Activation functions

- Attention variations

# Add & Norm
## Residual Connections and Layer Norm

- Combine residual connection and layer norm into a single "Add & Norm" component

- Three choices:

  - Pre-norm (input): $\mathbf{z}^{\ell+1} = f(\text{LN}(\mathbf{z}^\ell)) + \mathbf{z}^\ell$  https://arxiv.org/abs/2002.04745

  - Res-Post-norm (output): $\mathbf{z}^{\ell+1} = \text{LN}(f(\mathbf{z}^\ell)) + \mathbf{z}^\ell$  https://arxiv.org/pdf/2111.09883

    - (Does not combine well with dropout)

  - Post-norm: $\mathbf{z}^{\ell+1} = \text{LN}(f(\mathbf{z}^\ell) + \mathbf{z}^\ell)$  Original transformer

- Pre-norm leads to faster training.  Res-Post-norm has less spikes.

# Layer normalization

## Pre- vs Post- layer normalization

- Moving layer normalization to before the MHA and FFN block helps improve training stability

- Warm-up of learning rate is critical for training when using Post-LN
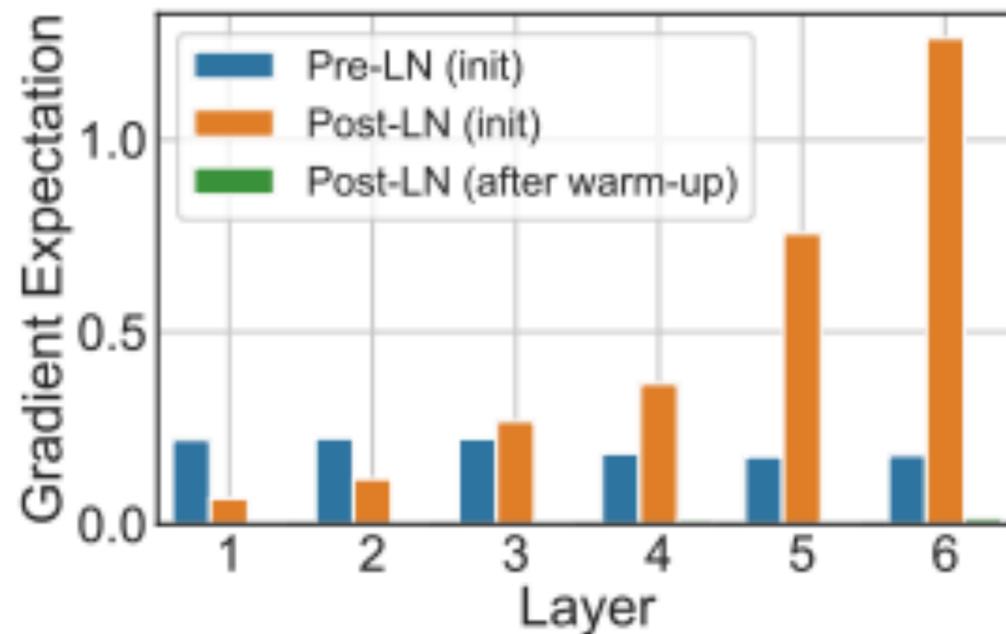
Original transformer
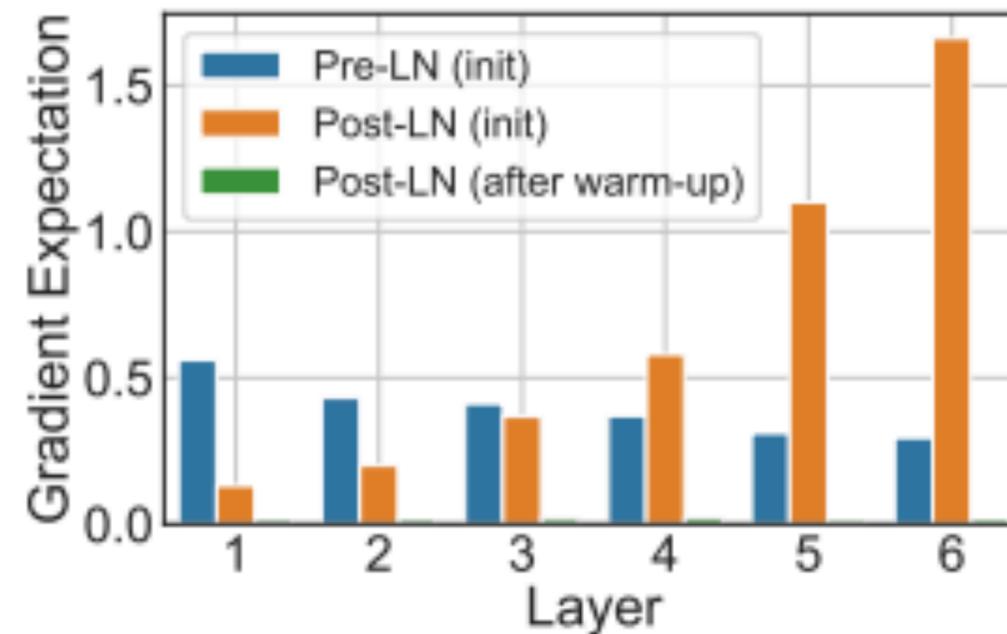


Post                          Pre

On Layer Normalization in the Transformer Architecture [Xiong et al. 2020]

# Layer normalization

## Pre- vs Post- layer normalization

- Warmup very important for Post-LN

- Gradient for Post-LN varies across layers during warmup, and is very small after warmup



(a) $W^1$ in the FFN sub-layers   (b) $W^2$ in the FFN sub-layers

On Layer Normalization in the Transformer Architecture [Xiong et al. 2020]

# Replace LayerNorm with RMSNorm

**LayerNorm**

- changes input features to have mean 0 and variance 1 per layer.

- Adds two more learnable parameters (vectors)

  - Gain **g** and bias **b**

$$\mu^l = \frac{1}{H} \sum_{i=1}^{H} x_i^l \qquad \sigma^l = \sqrt{\frac{1}{H} \sum_{i=1}^{H} (x_i^l - \mu^l)^2}$$

$$h_i = \frac{g_i}{\sigma_i}(x_i - \mu_i) + b_i$$

**RMSNorm**

- Simplifies LayerNorm by removing the mean and bias terms

$$\text{RMS}(\mathbf{x}) = \sqrt{\frac{1}{n} \sum_{i=1}^{n} x_i^2}$$

$$\text{RMSNorm}(\mathbf{x}) = \frac{\mathbf{x}}{\text{RMS}(\mathbf{x})} \cdot \mathbf{g}$$

Root Mean Square Layer Normalization
[Zhang and Sennrich, 2019]

# Revisiting transformers

- Positional encoding

- Normalization

  - Layer normalization vs RMS

  - Post vs Pre-Layer norm

- **Activation functions**

- Attention variations

# Activation function


ReLU Function

- Vaswani et al.: ReLU

$$\mathrm{ReLU}(\mathbf{x}) = \max(0, \mathbf{x})$$

## **Gated linear unit (GLU)**

- Output of a linear transformation of x is modulated by a gate

- LLaMa: Swish/SiLU (Hendricks and Gimpel 2016)

$$\mathrm{GLU} = \sigma(xW) \otimes xV$$

$$\mathrm{SwiGLU} = \mathrm{Swish}(xW) \otimes xV$$

$$\mathrm{Swish}(\mathbf{x}; \beta) = \mathbf{x} \odot \sigma(\beta\mathbf{x})$$


Swish Activation Function

## **Squared ReLU**

Outperforms ReLU, competitive with SwiGLU

- Sparse computations
- Used in Primer

# Revisiting transformers

- Positional encoding

- Normalization

  - Layer normalization vs RMS

  - Post vs Pre-Layer norm

- Activation functions

- **Attention variations**

# QK norm

- Normalize query and key to unit vectors

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

$\longrightarrow$

- Variants using LayerNorm

- Training stability

$\gamma$ is learnable parameter

$$A(Q, K, V) = \text{softmax}\left(\gamma \cdot Q'K'^T\right) V$$

$$Q_i' = \frac{Q_i}{\|Q_i\|}, K_j' = \frac{K_j}{\|K_j\|}$$

Query-Key Normalization for Transformers [Henry et al. 2020]

# Parallel attention and QK norm

- Parallel attention

$$y' = \text{LayerNorm}(x),$$
$$y = x + \text{MLP}(y') + \text{Attention}(y').$$

- QK norm

$$\text{softmax}\left[\frac{1}{\sqrt{d}}\text{LN}(XW^Q)(\text{LN}(XW^K))^T\right]$$



An encoder layer with parallel Attention-MLP blocks

Scaling Vision Transformers to 22 Billion Parameters [Dehghani et al. 2023]

# QK norm

- Prevents exploding attention logits (which leads to almost one-hot attention weights with near-zero entropy)



Figure 1: Effect of query/key normalization on an 8B parameter model.

Scaling Vision Transformers to 22 Billion Parameters [Dehghani et al. 2023]

# Sliding window attention

- Limited context / local attention

- Interleave with global attention layers

**Regular causal self-attention mask**

|     | The | cat | sat | on | the |
|-----|-----|-----|-----|-----|-----|
| The | 1 | 0 | 0 | 0 | 0 |
| cat | 1 | 1 | 0 | 0 | 0 |
| sat | 1 | 1 | 1 | 0 | 0 |
| on  | 1 | 1 | 1 | 1 | 0 |
| the | 1 | 1 | 1 | 1 | 1 |

Longformer: The long-document transformer.

Using a causal attention mask, the current token can only attend previous tokens (+ itself)

**Sliding window attention**

|     | The | cat | sat | on | the |
|-----|-----|-----|-----|-----|-----|
| The | 1 | 0 | 0 | 0 | 0 |
| cat | 1 | 1 | 0 | 0 | 0 |
| sat | 1 | 1 | 1 | 0 | 0 |
| on  | 0 | 1 | 1 | 1 | 0 |
| the | 0 | 0 | 1 | 1 | 1 |

Not attended to save computation

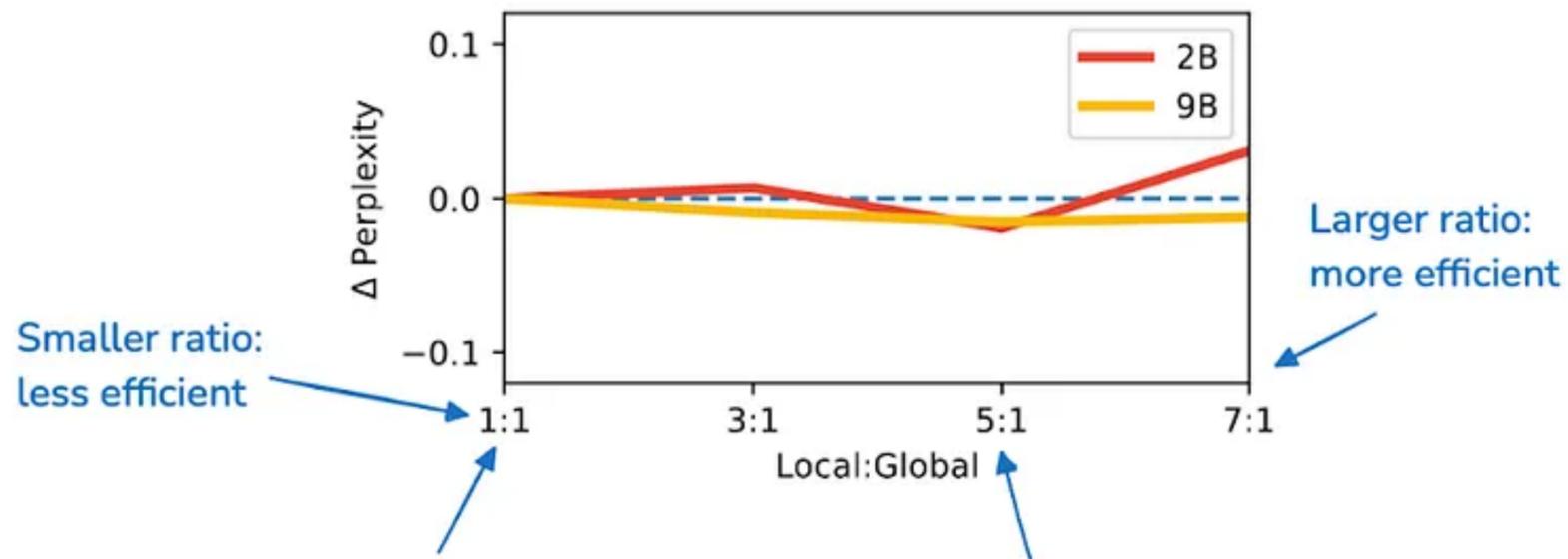Using a causal attention mask, the current token can only attend previous tokens **within a certain limit**

https://magazine.sebastianraschka.com/p/the-big-llm-architecture-comparison

*Longformer: The Long-Document Transformer* [Beltagy et al. 2020]

# Sliding window attention



Gemma 3 [Gemma Team. 2025]

(Annotated version from https://magazine.sebastianraschka.com/p/the-big-llm-architecture-comparison)

# Sliding window attention



Gemma 3 [Gemma Team. 2025]
(Annotated version from https://magazine.sebastianraschka.com/p/the-big-llm-architecture-comparison)

# Grouped multi-query attention

- Reduce number of heads used for keys and values

- Shared values and keys across heads



Fast Transformer Decoding: One Write-Head is All You Need [Shazeer 2019]

GQA: Training Generalized Multi-Query Transformer Models from Multi-Head Checkpoints [Ainslie et al. 2023]

# Multi-head latent attention (MLA)

$$c_n^{KV} = W^{DKV} x_n$$

- Use low rank matrix $C^{KV}$

$$K = W^{UK} C_{1:n}^{KV}$$

$$V = W^{UV} C_{1:n}^{KV}$$

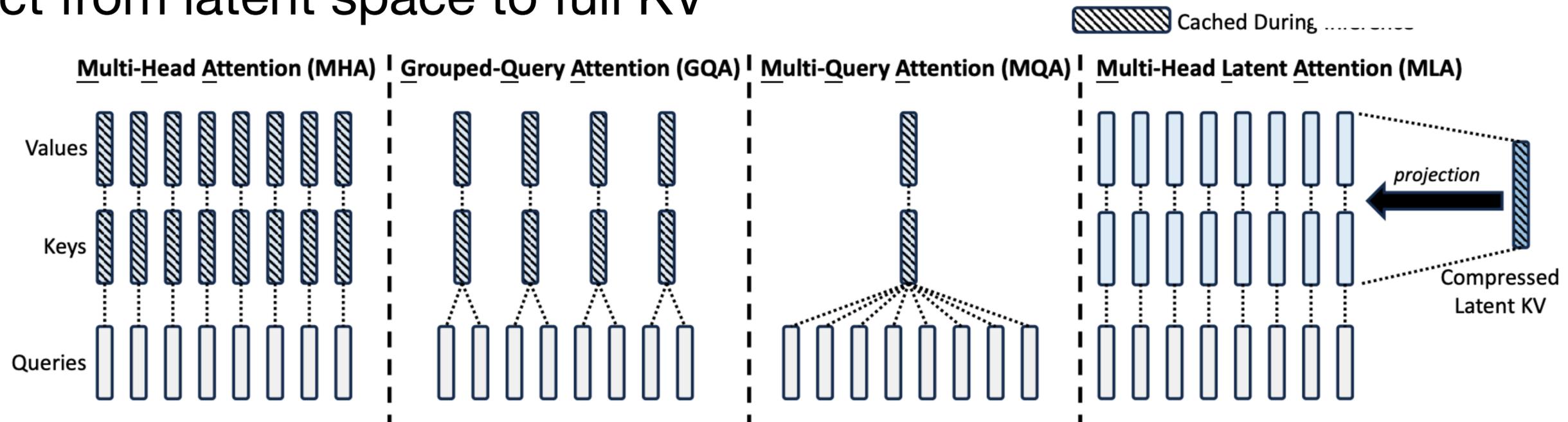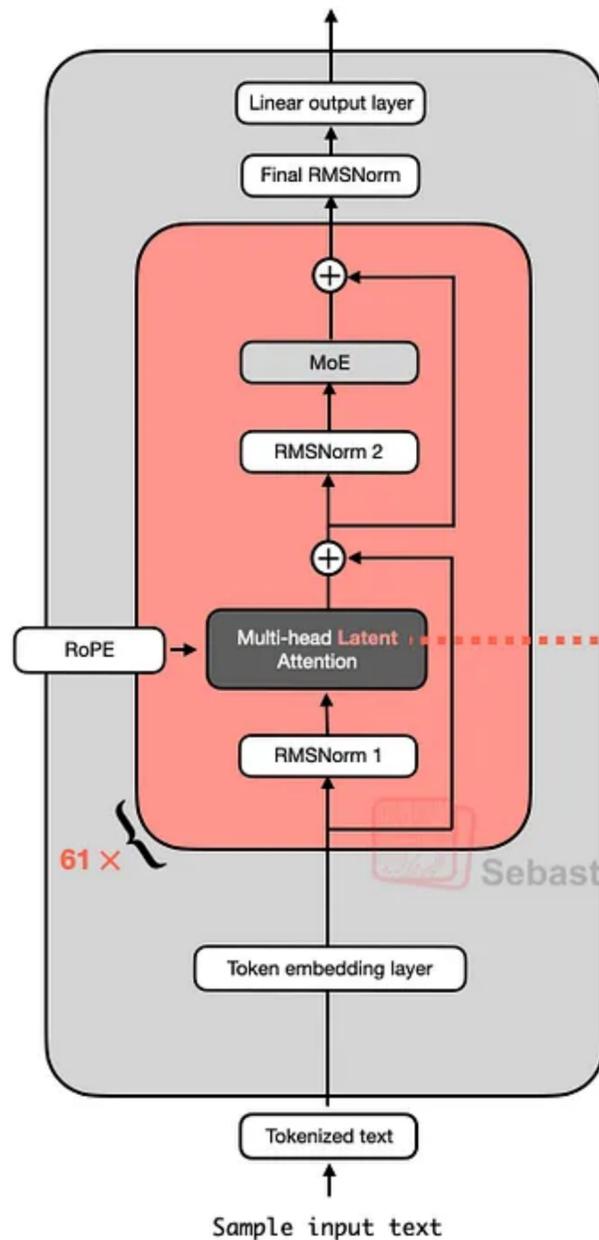- Project from latent space to full KV



Figure 3 | Simplified illustration of Multi-Head Attention (MHA), Grouped-Query Attention (GQA), Multi-Query Attention (MQA), and Multi-head Latent Attention (MLA). Through jointly compressing the keys and values into a latent vector, MLA significantly reduces the KV cache during inference.
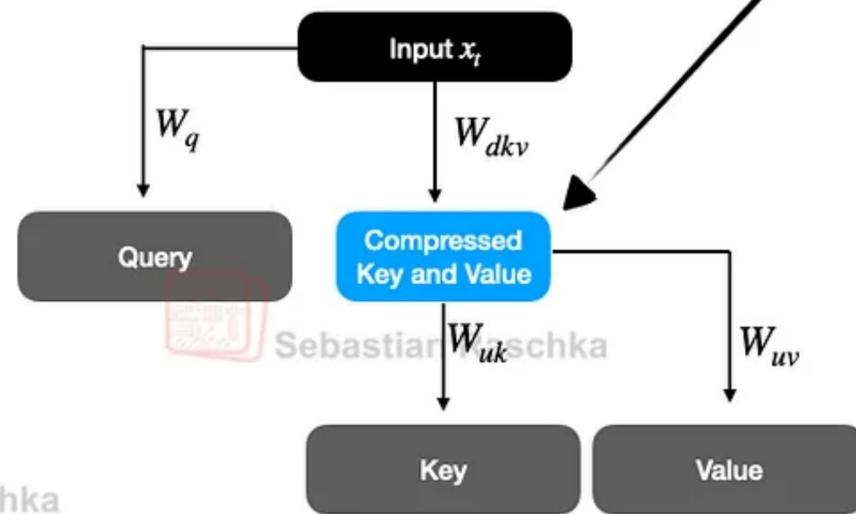
DeepSeek-V2: A Strong, Economical, and Efficient Mixture-of-Experts Language Model [DeepSeek AI 2024]

# Multi-head latent attention (MLA)



DeepSeek V3/R1

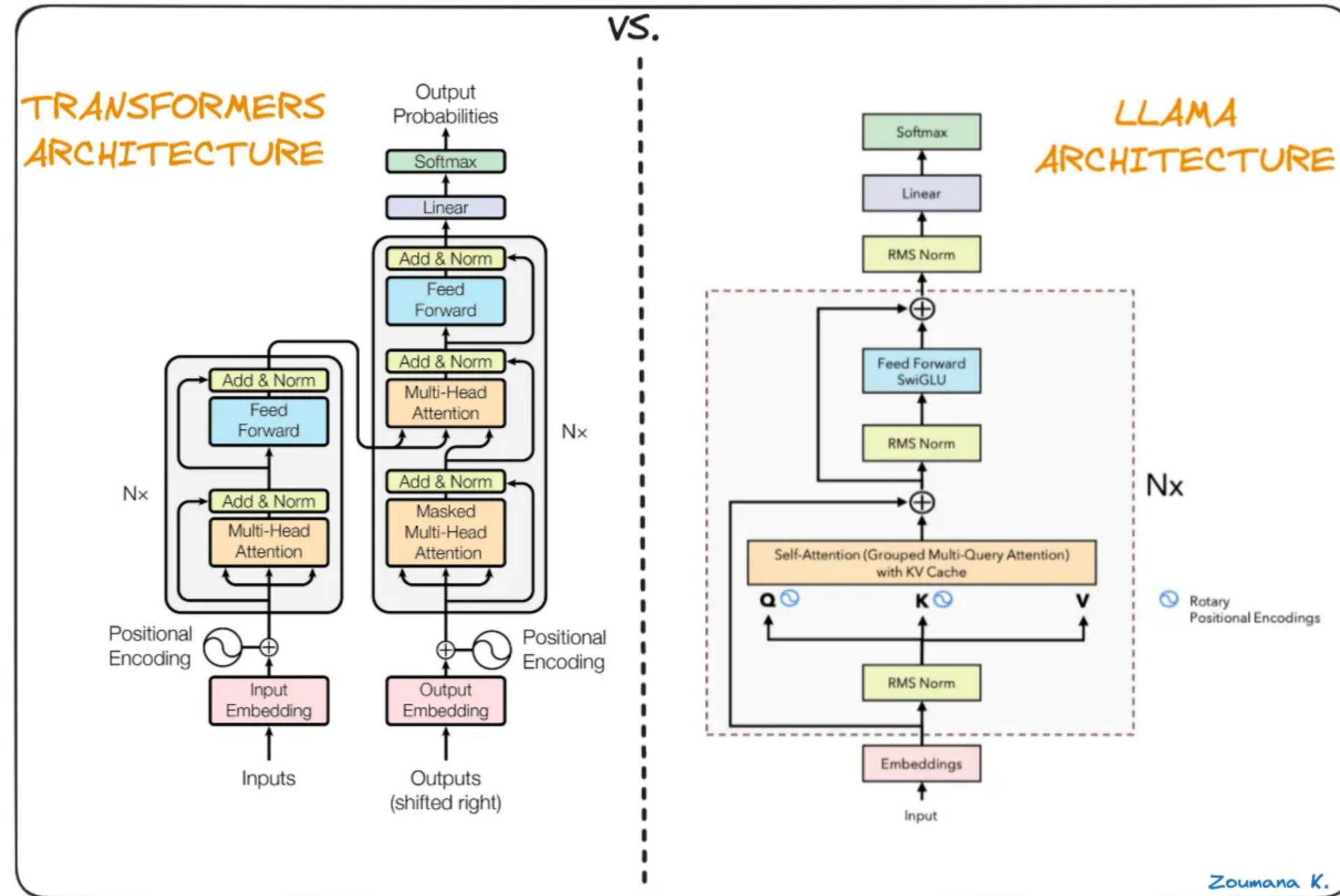Key idea: Reduces memory usage in KV cache

# Gated value embedding

- Add gated value embedding (at alternating layers)

- Additional expressivity (more parameters) with minimal compute overhead

```python
ve = value_embeds[layer_idx](token_ids)  # (B, T, kv_dim)
gate = 2 * sigmoid(ve_gate(x[:, :, :32]))  # range (0, 2)
v = v + gate * ve
```
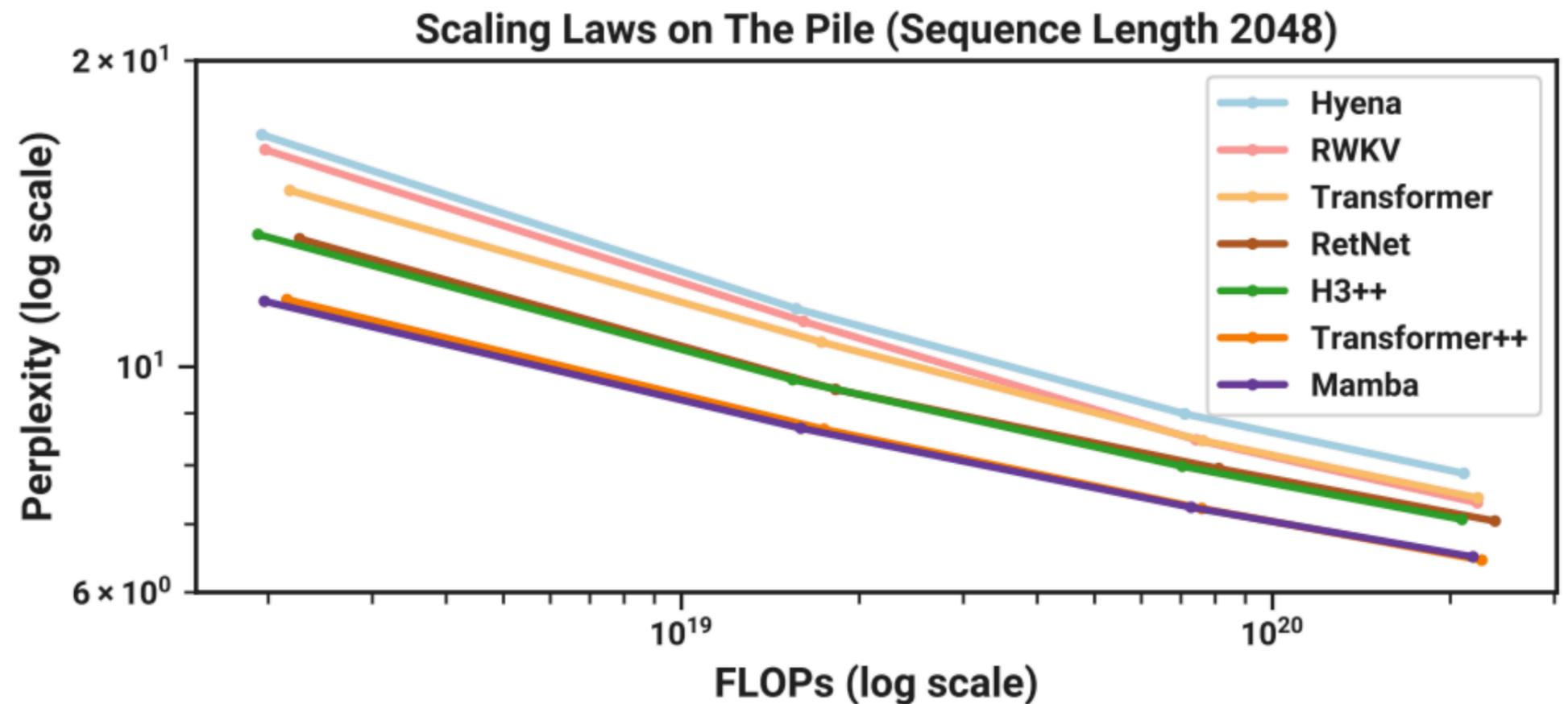
# LLaMa

| | Vaswani et al. | LLaMa |
|---|---|---|
| Norm Position | Post | Pre |
| Norm Type | LayerNorm | RMSNorm |
| Non-linearity | ReLU | SwiGLU |
| Positional Encoding | Sinosoidal | RoPE |
| Attention | Full Multi-Head Attention | Grouped Multi-Query Attention |



TRANSFORMERS ARCHITECTURE vs. LLAMA ARCHITECTURE

# Impact of these changes

- Transformer: Vaswani et al.

- Transformer++: basically LLaMa



- Stronger architecture is ≈10x more efficient!
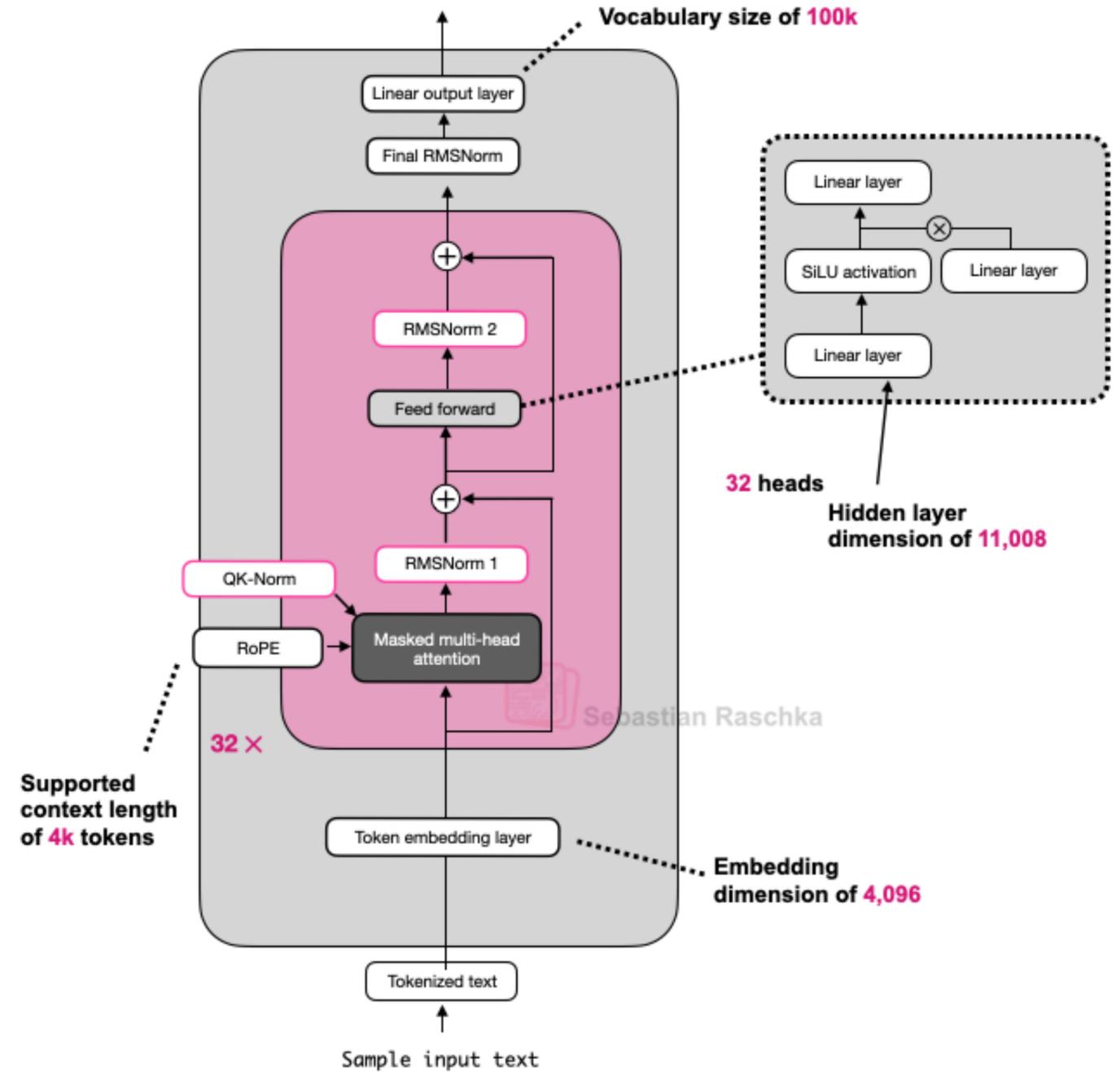
# OLMo 2 architecture

| | OLMo 1 (0224) | OLMo-0424 | OLMo 2 |
|---|---|---|---|
| Biases | None | None | None |
| Activation | SwiGLU | SwiGLU | SwiGLU |
| RoPE $\theta$ | $1 \cdot 10^4$ | $1 \cdot 10^4$ | $5 \cdot 10^5$ |
| QKV Normalization | None | Clip to 8 | QK-Norm |
| Layer Norm | non-parametric | non-parametric | RMSNorm |
| Layer Norm Applied to | Inputs | Inputs | Outputs |
| Z-Loss Weight | 0 | 0 | $10^{-5}$ |
| Weight Decay on Embeddings | Yes | Yes | No |

2 OLMo 2 Furious [AI2 OLMo team, 2024]

Llama 3 8B

Vocabulary size of 128k

Linear output layer

Final RMSNorm

Feed forward

Linear layer

SiLU activation    Linear layer

Linear layer

RMSNorm 2

Masked grouped-query attention

RoPE

32 heads

Hidden layer dimension of 14,336

RMSNorm 1

Supported context length of 128k tokens

32 ×

Token embedding layer

Embedding dimension of 4,096

Tokenized text

Sample input text

OLMo 2 7B

Vocabulary size of 100k

Linear output layer

Final RMSNorm

Linear layer

SiLU activation    Linear layer

Linear layer

RMSNorm 2

Feed forward

RMSNorm 1

QK-Norm

Masked multi-head attention

RoPE

32 heads

Hidden layer dimension of 11,008

Supported context length of 4k tokens

32 ×

Token embedding layer

Embedding dimension of 4,096

Tokenized text

Sample input text

Sebastian Raschka

https://magazine.sebastianraschka.com/p/the-big-llm-architecture-comparison

# Development of Open LLMs

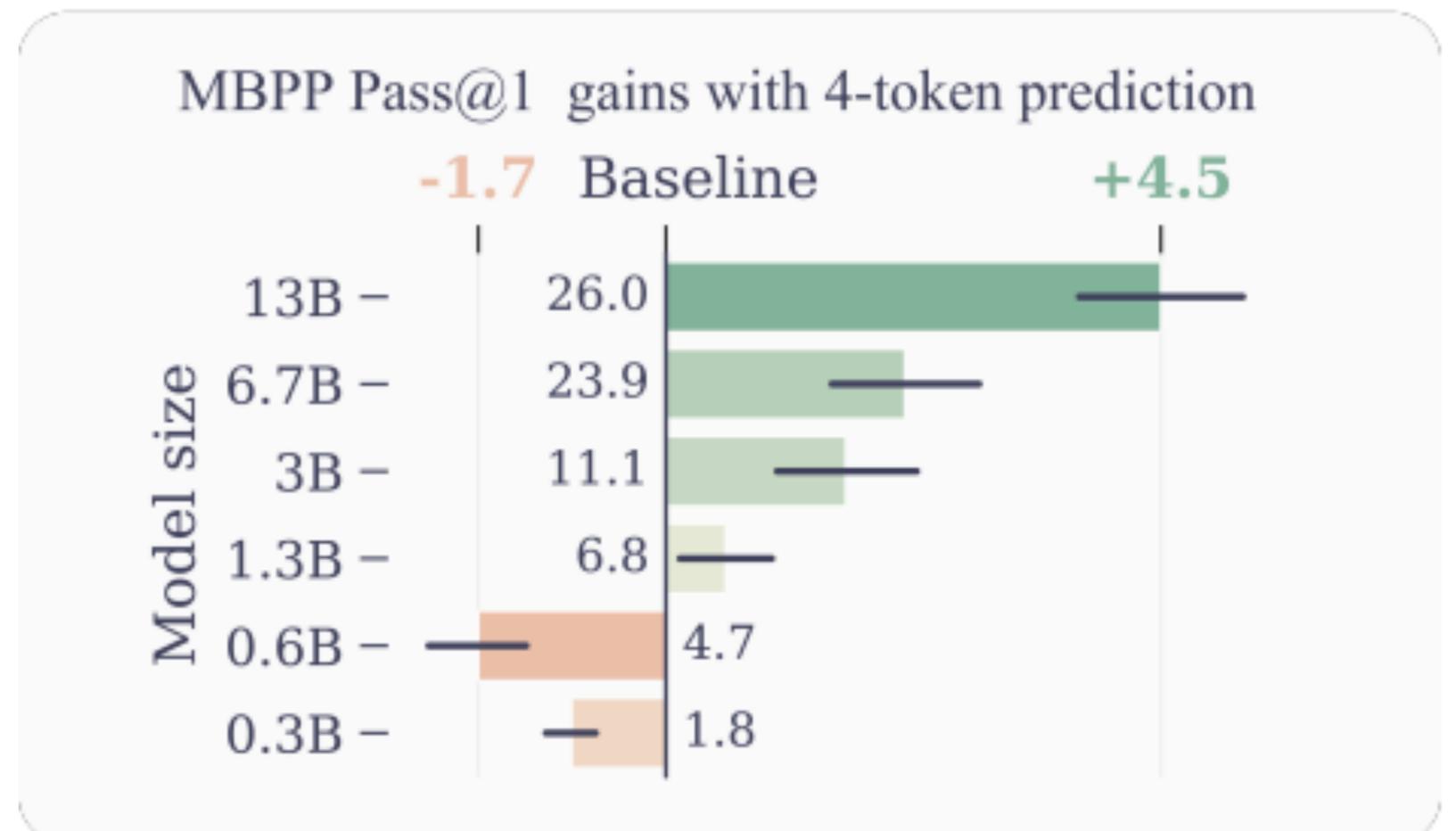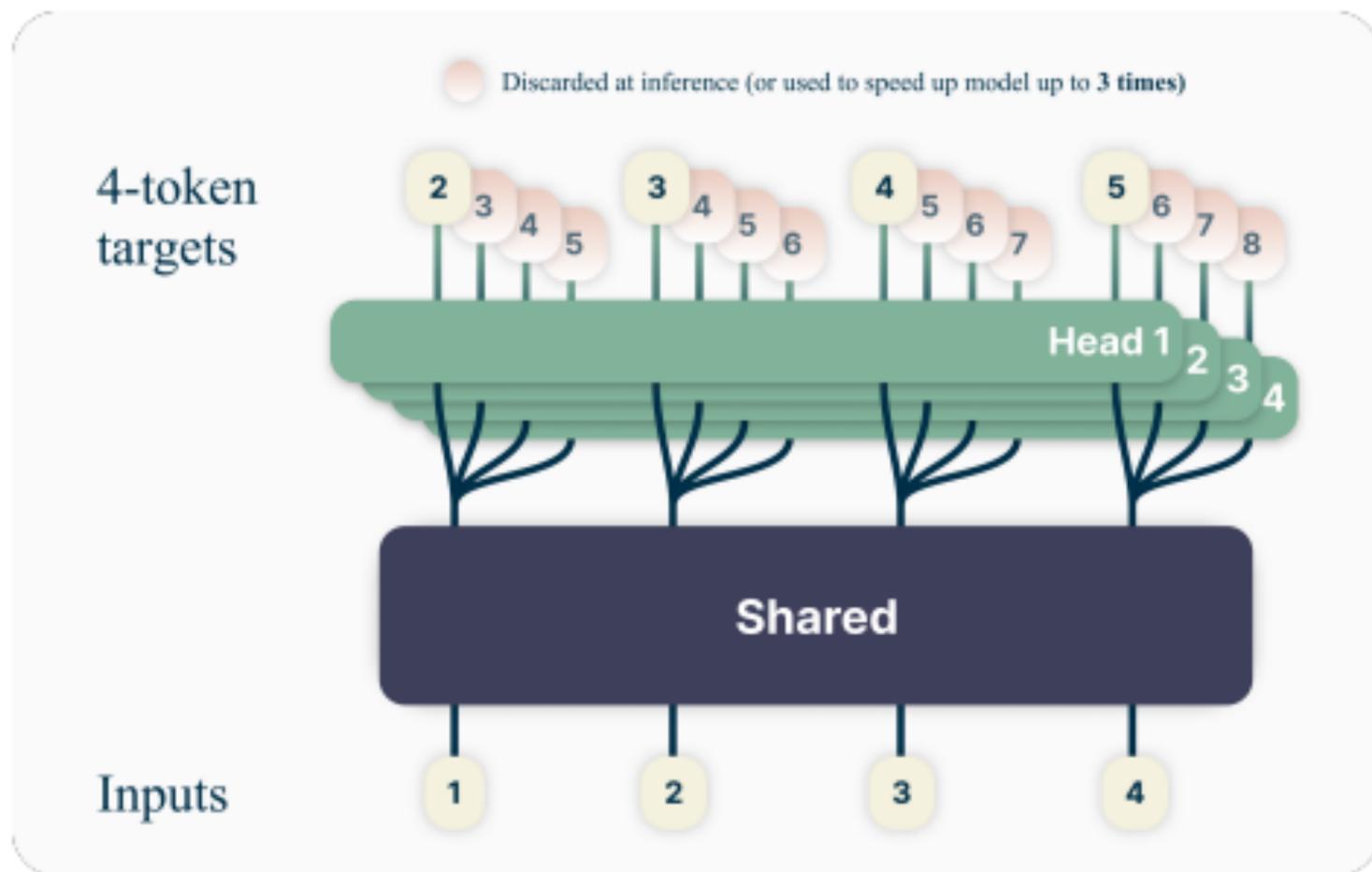2 OLMo 2 Furious [AI2 OLMo team, 2024]

# Other aspects of variation

- Training objective

- Scaling up with Mixture-of-Experts

- Optimizer

- Efficient training

# Multi-token prediction

- Predict next tokens with shared trunk



- Gains due to multi-token prediction

# Multi-token prediction
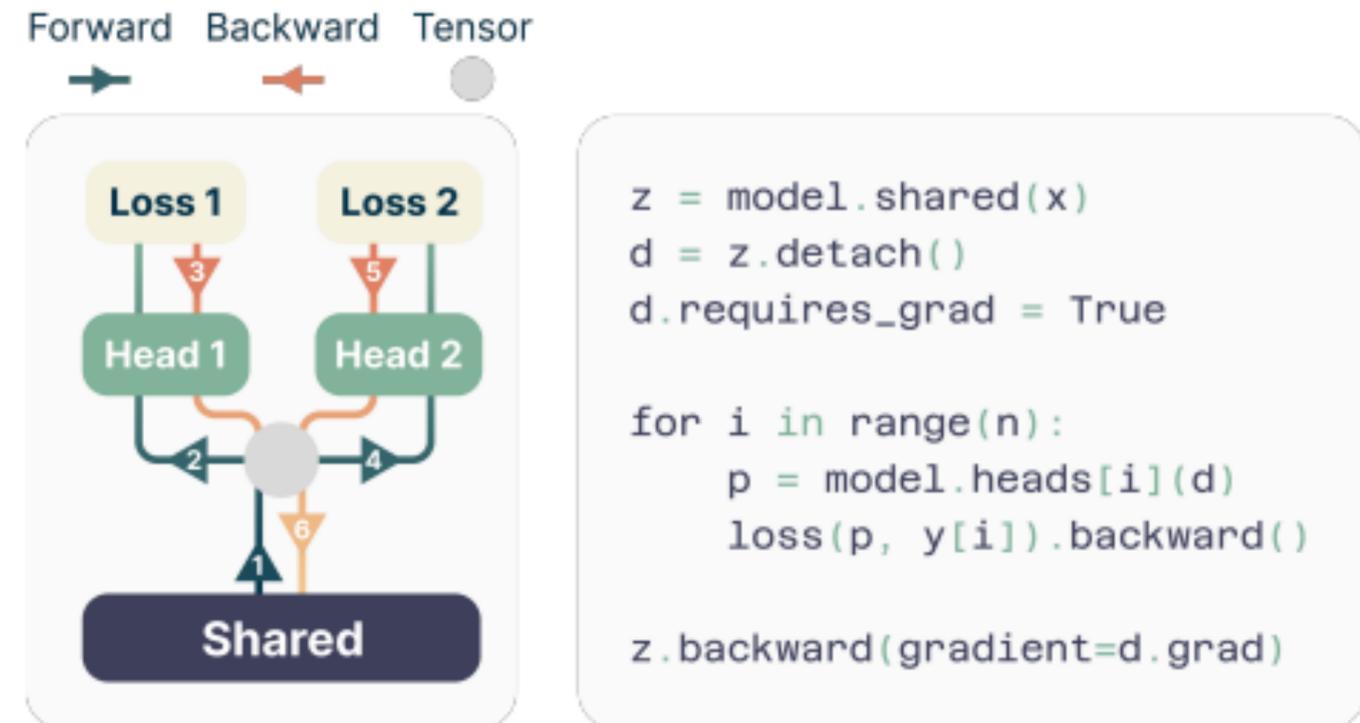
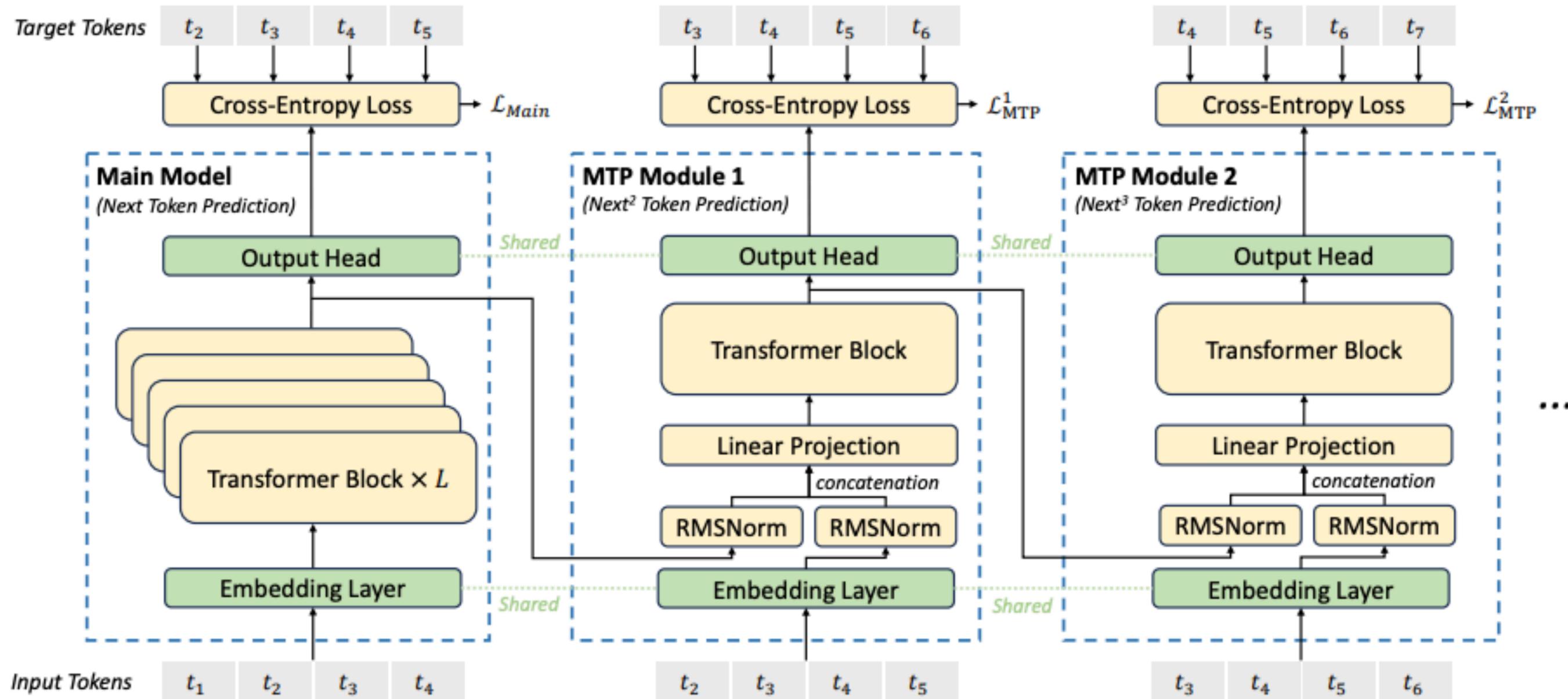- Predict in sequence to avoid large memory demands



Figure 2: **Order of the forward/backward in an $n$-token prediction model with $n = 2$ heads.** By performing the forward/backward on the heads in sequential order, we avoid materializing all unembedding layer gradients in memory simultaneously and reduce peak GPU memory usage.

# Multi-token prediction
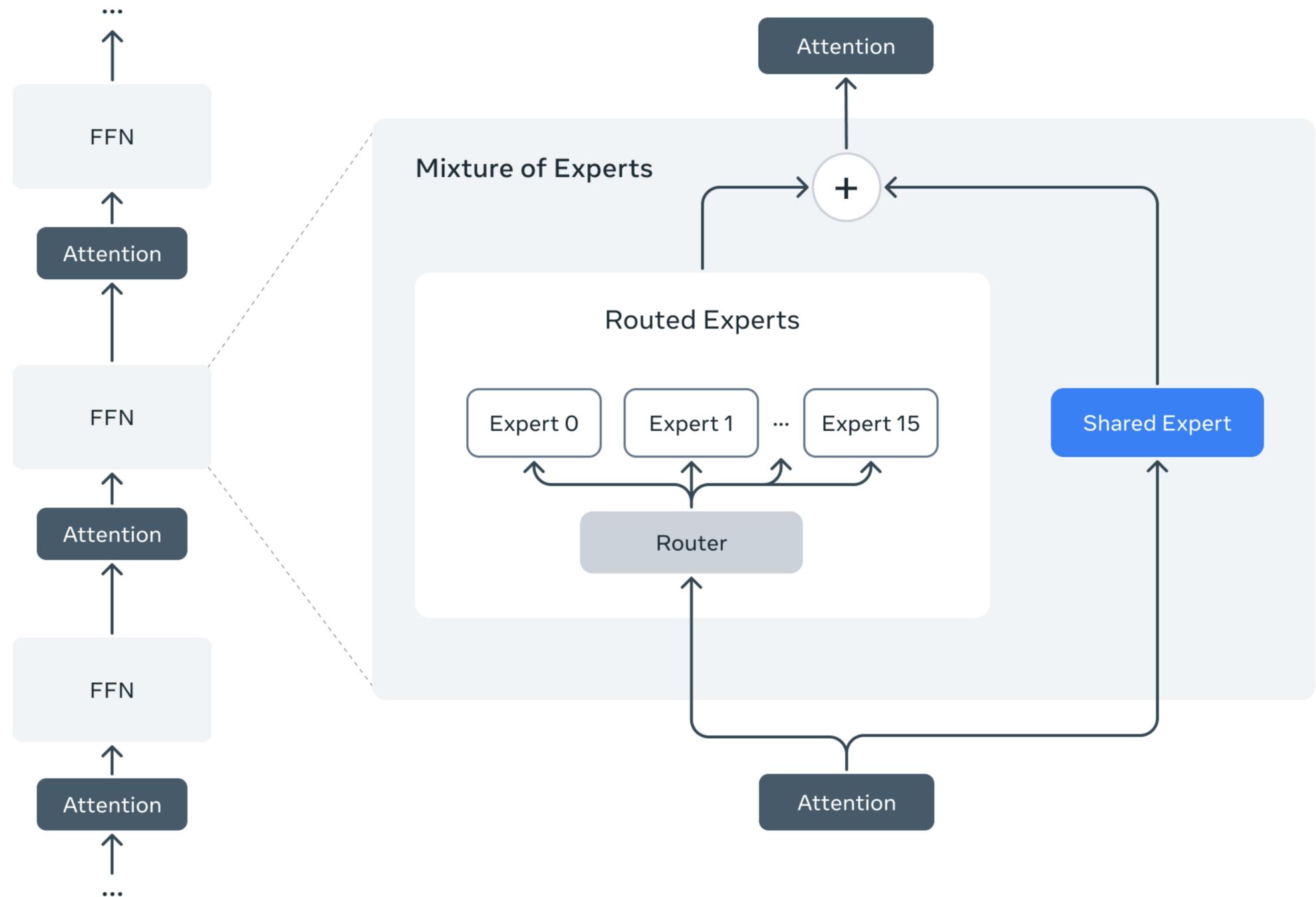
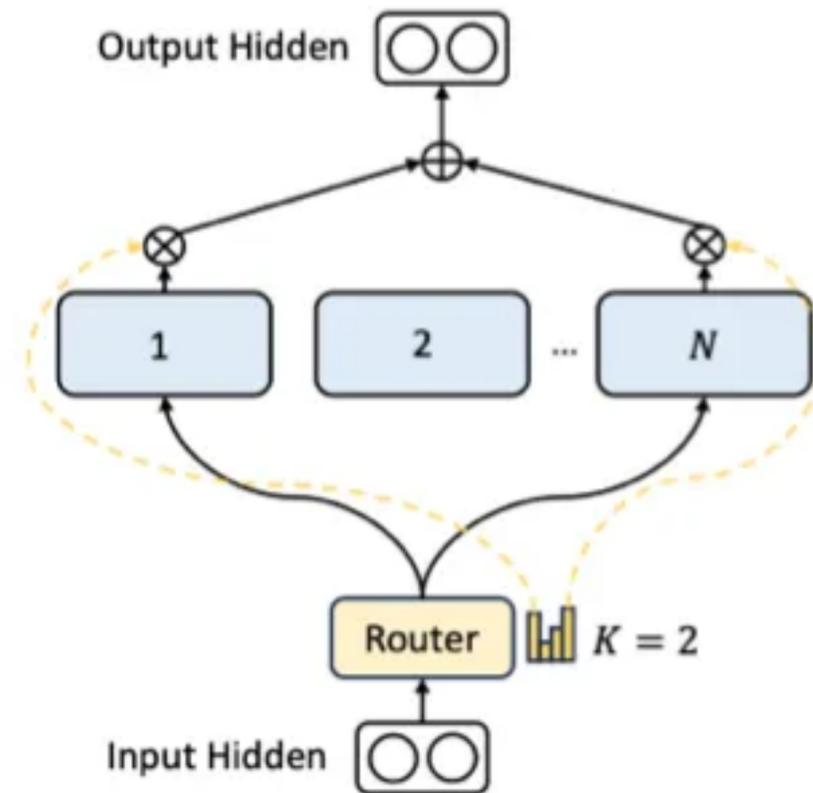- D sequential modules for predicting D tokens

# Mixture of Experts

- Allows for large capacity with low computational cost per token

- Router that selects which expert (individual neural networks) to use



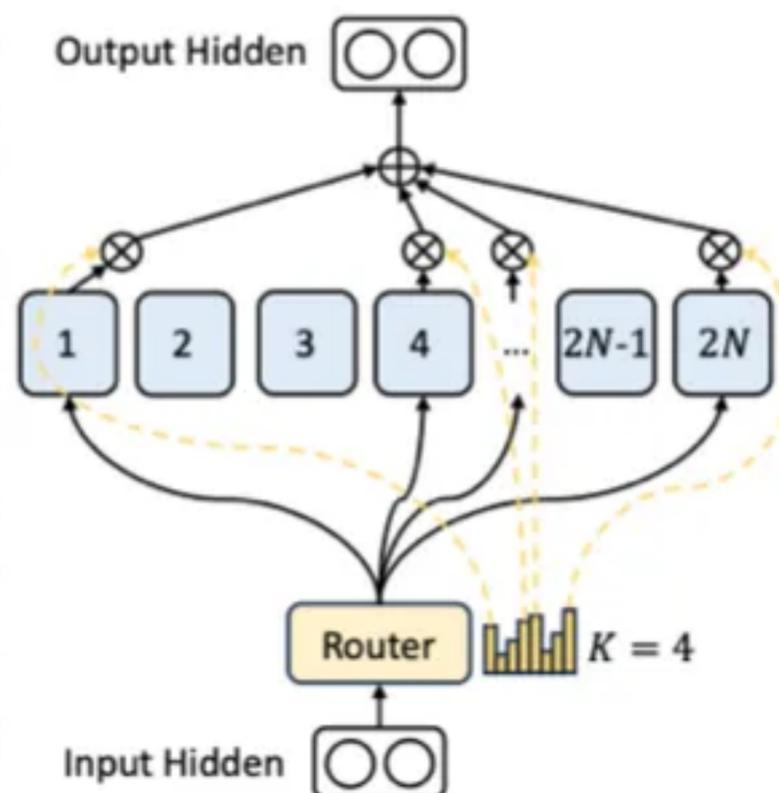https://ai.meta.com/blog/llama-4-multimodal-intelligence/
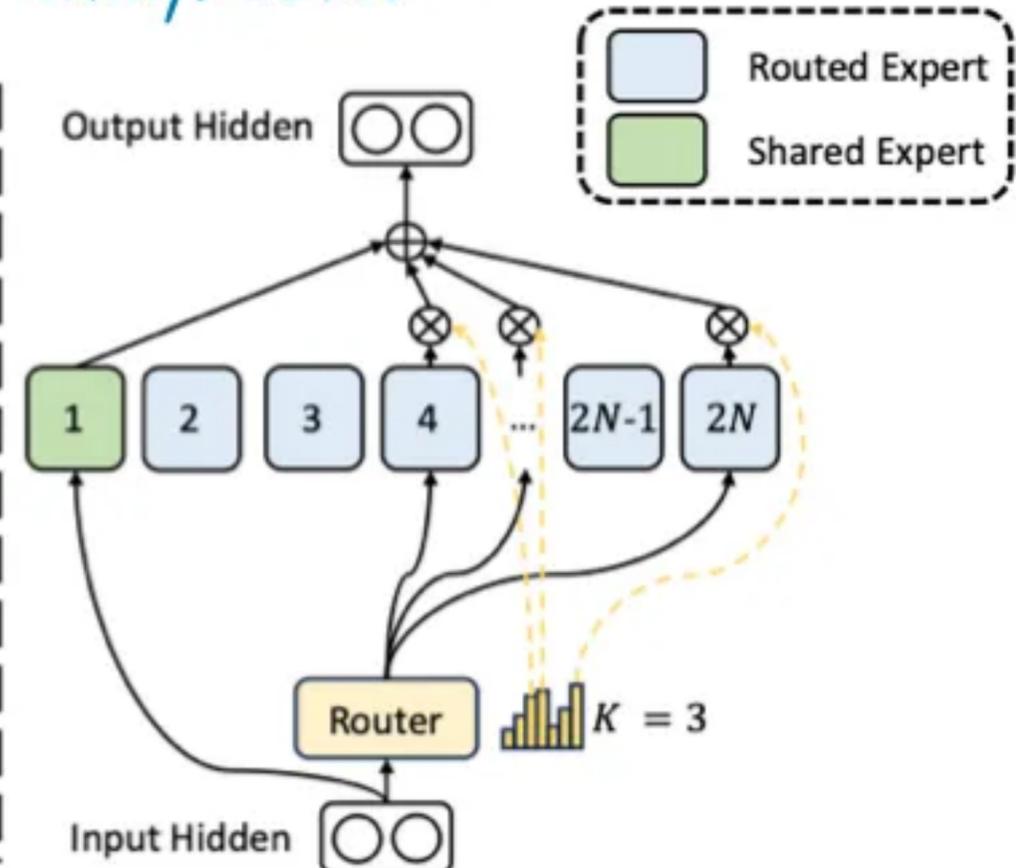
# Mixture of experts over time



Early MoE: Has bigger and fewer experts, and activates only a few experts (here: 2)

Fine-grained MoE uses more but smaller experts, and activates more experts (here: 4)

MoE with shared expert: also uses many small experts, but adds a shared expert that is always active
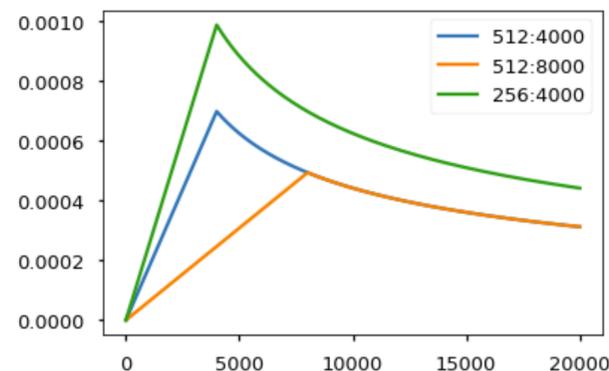
Routed Expert

Shared Expert

Output Hidden

Input Hidden

Router $K = 2$

Router $K = 4$

Router $K = 3$

(a) Conventional Top-2 Routing ➡ (b) + Fine-grained Expert Segmentation ➡ (c) + Shared Expert Isolation (DeepSeekMoE)

# Optimizer

- SGD: Update in the direction of reducing loss

- Adam: Add momentum and normalize by the stddev of the outputs

- AdamW: properly applies weight decay for regularization to Adam

- Adam with learning rate schedule:

$$\text{learning\_rate} = d_{\text{model}}^{-0.5} \cdot \min(\text{step\_num}^{-0.5}, \text{step\_num}^{-0.5} \cdot \text{warmup\_steps}^{-1.5})$$
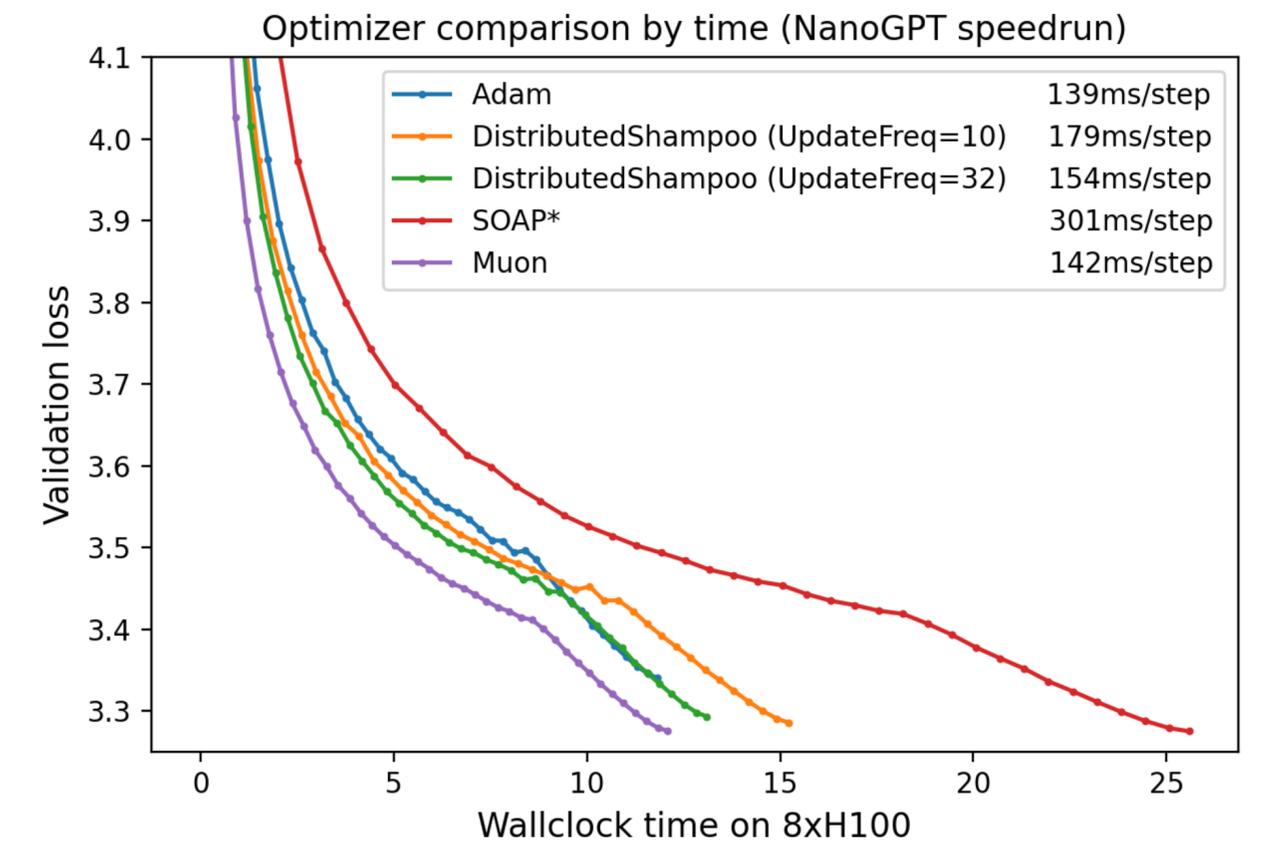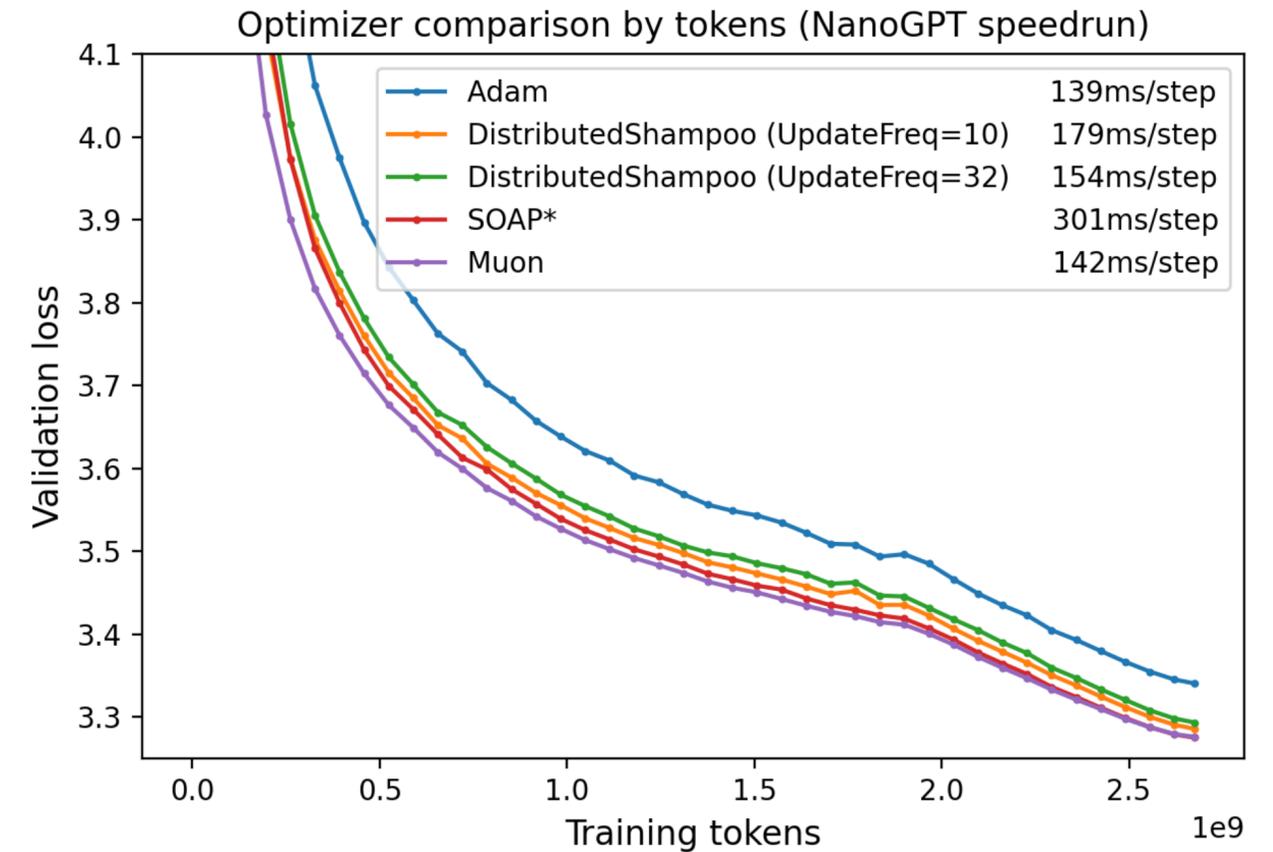
# Muon Optimizer

- Orthogonalize the update

---
**Algorithm 2** Muon

---
**Require:** Learning rate $\eta$, momentum $\mu$
1: Initialize $B_0 \leftarrow 0$
2: **for** $t = 1, \ldots$ **do**
3:      Compute gradient $G_t \leftarrow \nabla_\theta \mathcal{L}_t(\theta_{t-1})$
4:      $B_t \leftarrow \mu B_{t-1} + G_t$
5:      $O_t \leftarrow \text{NewtonSchulz5}(B_t)$
6:      Update parameters $\theta_t \leftarrow \theta_{t-1} - \eta O_t$
7: **end for**
8: **return** $\theta_t$

---

- Use for weight matrices, use AdamW for embeddings

https://kellerjordan.github.io/posts/muon/



Optimizer comparison by tokens (NanoGPT speedrun)

| | |
|---|---|
| Adam | 139ms/step |
| DistributedShampoo (UpdateFreq=10) | 179ms/step |
| DistributedShampoo (UpdateFreq=32) | 154ms/step |
| SOAP* | 301ms/step |
| Muon | 142ms/step |

Optimizer comparison by time (NanoGPT speedrun)

| | |
|---|---|
| Adam | 139ms/step |
| DistributedShampoo (UpdateFreq=10) | 179ms/step |
| DistributedShampoo (UpdateFreq=32) | 154ms/step |
| SOAP* | 301ms/step |
| Muon | 142ms/step |

*SOAP is under active development. Future versions will significantly improve the wallclock overhead.

# Resources

- LLM architectures

  - https://magazine.sebastianraschka.com/p/the-big-llm-architecture-comparison

  - https://magazine.sebastianraschka.com/p/a-dream-of-spring-for-open-weight

- Andrej Karpathy - NanoGPT journey

  - https://github.com/karpathy/nanochat/discussions/481