Multi-Head Low-Rank Attention

Songtao Liu^{1*} Hongwu Peng² Zhiwei Zhang¹ Zhengyu Chen³ Yue Guo⁴

¹The Pennsylvania State University ²University of Connecticut

³Carnegie Mellon University ⁴University of California, Los Angeles

Abstract

Long-context inference in large language models is bottlenecked by Key-Value (KV) cache loading during the decoding stage, where the sequential nature of generation requires repeatedly transferring the KV cache from off-chip to on-chip memory at each step. Recent architectures like Multi-Head Latent Attention (MLA) significantly reduce the KV cache size to $4.5d_h$ per token per layer while maintaining high model quality. However, when using tensor parallelism (TP) with sufficient devices for inference, MLA still decodes slower than Grouped-Query Attention (GQA) because its single latent vector cannot be sharded, forcing each device to load $4.5d_h$ versus $2d_h$ for GQA. In this work, we propose Multi-Head Low-Rank Attention (MLRA), a TP-friendly attention mechanism that slashes the per-device KV cache under TP to just $1.5d_h$. Extensive experiments show that MLRA achieves state-of-the-art perplexity and downstream task performance, while also delivering a $2.8\times$ decoding speedup over MLA.

1 Introduction

Inference-time scaling (OpenAI et al., 2024) is critical for large language models (LLMs) to produce high-quality responses. Both retrieval-augmented generation (RAG) (Lewis et al., 2020) and long chain-of-thought (CoT) reasoning (Wei et al., 2022) rely on maintaining long context before generating the final answer, substantially increasing the number of tokens that must be processed at each decoding step. Sequential token generation under Multi-Head Attention (MHA) (Vaswani et al., 2017) requires reloading the Key-Value (KV) cache from high-bandwidth memory every step, so data movement (Ivanov et al., 2021; Ootomo and Yokota, 2023; Gholami et al., 2024), not computation, dominates latency for long-context inference (Sadhukhan et al., 2025). The small amount of compute per step relative to this data movement leads to low arithmetic intensity (Williams et al., 2009) and poor GPU utilization (He and Zhai, 2024; Zadouri et al., 2025).

Recent work (Zadouri et al., 2025) proposes to analyze the performance and decoding efficiency of inference-aware attention mechanisms (Hu et al., 2024; Sun et al., 2025b; Zheng et al., 2025) along four axes: (1) model quality, (2) KV cache footprint per token, (3) arithmetic intensity, and (4) the degree of tensor parallelism (TP) (Pope et al., 2023) across attention heads. Multi-Query Attention (MQA) (Shazeer, 2019) shares single key and value heads across all query heads, reducing the KV cache to one head $(2d_h)$ yet maintaining the same floating point operations per second (FLOPs) as MHA. Despite the higher arithmetic intensity and smaller KV cache than MHA, it often leads to noticeable quality degradation. Grouped-Query Attention (GQA) (Ainslie et al., 2023) shares single key and value heads within a group of query heads, improving model quality over MQA while achieving higher arithmetic intensity than MHA. However, its KV cache of $2gd_h$ still scales with the group size (with g=8 in LLaMA-3-70B (Llama et al., 2024) and Grok-2), which can be memory intensive. As a result, reaching the fastest decoding $(2d_h)$ per device) typically requires g-way TP.

Multi-Head Latent Attention (MLA) (DeepSeek et al., 2024a) compresses the KV cache into a latent vector (4.5 d_h per token). By absorbing the up-projection matrices into the queries during decoding,

^{*}Correspondence to: Songtao Liu <skl5761@psu.edu>.

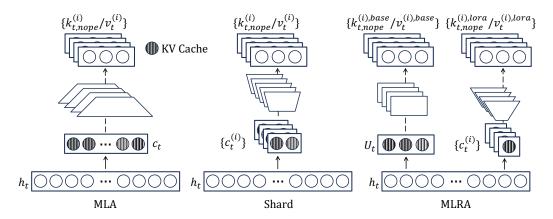


Figure 1: Illustration of MLA, Shard, and MLRA adopts a balanced design that achieves TP through the low-rank path while maintaining model quality via the base path.

its attention computation is similar to that of MQA with shared KV states, which increases arithmetic intensity by $2\times$ over MQA and delivers better model quality compared with GQA and MHA. However, MLA is unfriendly to TP because its single latent vector cannot be sharded. To mitigate this, Grouped Latent Attention (GLA) (Zadouri et al., 2025) compresses the KV cache into two latent vectors, each $2d_h$, plus an additional $0.5d_h$ for partial RoPE (Su et al., 2024). Nevertheless, its per-device $2.5d_h$ KV cache still leaves decoding slower than GQA when enough devices are available. In this work, we target a high TP degree with per-device KV cache below $2d_h$.

A simple and straightforward approach is to split MLA's single high-dimensional latent vector into h smaller latent vectors (illustrated in the middle of Figure 1), as in GLA, which splits the latent vector into two smaller vectors. However, this simple approach can cause a significant drop in model quality, as illustrated by the 354M-model training loss curves in Figure 2. To address this limitation, we introduce Multi-Head Low-Rank Attention (MLRA), a dual-path attention mechanism shown on the right of Figure 1, which can achieve a small per-device KV cache under TP and high model quality. In the low-rank path, the KV cache is compressed to h tiny latent vectors, each of which is upprojected to single key and value heads, and can be partitioned across devices to support TP. In the base path, we follow MLA by compressing the KV cache into a base latent vector whose dimension matches the perhead query (d_h) , then up-project it to produce h key

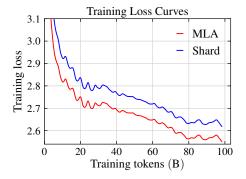


Figure 2: Training loss curves on FineWeb-Edu (100B tokens) for 354M models: MLA vs. Shard. Shard splits the latent vector into 16 smaller vectors.

and value heads. We then sum the attention outputs from the two paths. Therefore, our proposed MLRA can achieve a per-device KV cache of $1.5d_h$ with 4-way TP. Based on our 2.9B scale experiments, MLRA achieves the lowest perplexity (12.998 vs. 13.115 for MLA and 13.399 for GQA) and highest downstream accuracy (57.06% vs. 55.46% for MLA and 55.68% for GQA). Our kernel delivers a 1.05- $1.26 \times$ speedup over GQA in long-context decoding.

2 Preliminaries

2.1 Tensor Parallelism

Long-context decoding is bottlenecked by KV cache loading from high-bandwidth memory at each step. Given the sequential nature of generation, we can only accelerate inference by distributing the KV cache and model weights across devices via tensor parallelism. With high-bandwidth GPU interconnects such as NVLink, the synchronization overhead in TP can be substantially reduced.

2.2 Translation Equivariance

Equivariance is a fundamental property in geometric systems such as molecules, where vector features such as atomic forces or dipole moments must transform consistently with the coordinate system (Thomas et al., 2018; Weiler et al., 2018; Fuchs et al., 2020; Satorras et al., 2021). In the context of sequence models, a common transformation is sequence translation. Let $\mathbf{x}=(\mathbf{x}_1,\mathbf{x}_2,\ldots,\mathbf{x}_n)\in X$ be a sequence of tokens, and define a translation operator $T_s:X\to X$ that translates the entire sequence by s positions. Let $\phi:X\to Y$ be a function that maps a sequence to a matrix of attention scores $\phi(\mathbf{x})\in\mathbb{R}^{n\times n}$, where each element $\phi(\mathbf{x})_{i,j}=A(\mathbf{x}_i,\mathbf{x}_j)$ denotes the attention score between tokens \mathbf{x}_i and \mathbf{x}_j . We say that ϕ is translation equivariant if there exists a corresponding output-space transformation $S_s:Y\to Y$ such that:

$$\phi\left(T_s\left(\mathbf{x}\right)\right) = S_s\left(\phi\left(\mathbf{x}\right)\right), \quad \forall s. \tag{1}$$

This property ensures that the attention score between two tokens depends only on their relative positions, not their absolute positions. That is crucial for batch inference using left padding, where sequences of different lengths are offset to align ends. The first non-padding token of a sequence is no longer at position 0, yet attention scores remain invariant to this translation.

2.3 Rotary Position Embedding

Rotary Position Embedding (RoPE) (Su et al., 2024) is a positional encoding method designed to incorporate relative position information directly into the attention mechanism. We show in this section that RoPE is translation equivariant.

Theorem 1. Given two tokens with query \mathbf{q} and key \mathbf{k} at positions m and n, respectively, and applying RoPE to obtain RoPE (\mathbf{q}, m) and RoPE (\mathbf{k}, n), we aim to prove that RoPE is translation equivariant with respect to the inner product. In the definition of translation equivariance (Eq.(1)), this corresponds to setting the output transformation S_s as the identity map, i.e., $S_s = \mathbf{I}$. Specifically, we show that translating both positions by an offset s leaves the inner product unchanged:

$$\langle \text{RoPE}(\mathbf{q}, m + s), \text{RoPE}(\mathbf{k}, n + s) \rangle = \langle \text{RoPE}(\mathbf{q}, m), \text{RoPE}(\mathbf{k}, n) \rangle.$$

Proof can be found in Appendix A.

Remark 1. While RoPE is translation equivariant, applying a linear projection **W** after RoPE generally breaks this property. Specifically, suppose we define the inner product after RoPE and linear projection as:

$$\langle \text{RoPE}(\mathbf{q}, m) \mathbf{W}_{Q}, \text{ RoPE}(\mathbf{k}, n) \mathbf{W}_{K} \rangle = \mathbf{q} \mathbf{R}_{m} \mathbf{W}_{Q} \mathbf{W}_{K}^{\top} \mathbf{R}_{n}^{\top} \mathbf{k}^{\top}.$$
 (2)

The term $\mathbf{W}_{\mathrm{Q}}\mathbf{W}_{\mathrm{K}}^{\top}$ breaks translation equivariance by disrupting the expression's dependence on the relative position m-n. The property would only be preserved in the specific case where $\mathbf{W}_{\mathrm{Q}}\mathbf{W}_{\mathrm{K}}^{\top}=\mathbf{I}$, which would reduce the expression to its original form. However, since this constraint is difficult to enforce during training, translation equivariance is generally lost when applying a linear projection after RoPE.

3 Multi-Head Low-Rank Attention

3.1 Implementation

As noted in Remark 1, applying a linear projection after RoPE destroys translation equivariance. Therefore, we employ MLA's partial RoPE (DeepSeek et al., 2024a). We first apply a linear down-projection over the hidden states $\mathbf{H} \in \mathbb{R}^{n \times d}$ for an n-token sequence:

$$\left[\mathbf{U}_{\mathrm{KV}},\ \overline{\mathbf{C}}_{\mathrm{KV}},\ \widetilde{\mathbf{K}}_{\mathrm{RoPE}}\right] = \mathbf{H}\mathbf{A}_{\mathrm{KV}},\tag{3}$$

where $\mathbf{A}_{\mathrm{KV}} \in \mathbb{R}^{d \times \left(d_u + hr + d_h^R\right)}$ is a down-projection matrix. By default we set $d_u = d_h, r = \frac{3d_h}{h}$, and $d_h^R = 0.5d_h$ so that the total KV cache size matches that of MLA. The projected output is then partitioned as follows: a base latent head $\mathbf{U}_{\mathrm{KV}} \in \mathbb{R}^{n \times d_u}$ for the base path; $\overline{\mathbf{C}}_{\mathrm{KV}} \in \mathbb{R}^{n \times hr}$, reshaped to $\mathbf{C}_{\mathrm{KV}} \in \mathbb{R}^{n \times h \times r}$ to yield h tiny latent heads for the low-rank path; and $\widetilde{\mathbf{K}}_{\mathrm{ROPE}} \in \mathbb{R}^{n \times d_h^R}$, which is RoPE-encoded to obtain $\mathbf{K}_{\mathrm{ROPE}}$.

Base Path. We up-project the base latent head $\mathbf{U}_{\mathrm{KV}} \in \mathbb{R}^{n \times d_u}$ with the up-projection matrix $\mathbf{B}_{\mathrm{base}} \in \mathbb{R}^{d_u \times (2hd_h)}$ to produce concatenated keys and values:

$$\left[\overline{\mathbf{K}}_{\text{base}}, \overline{\mathbf{V}}_{\text{base}}\right] = \mathbf{U}_{\text{KV}} \, \mathbf{B}_{\text{base}} \in \mathbb{R}^{n \times (2hd_h)}.$$
 (4)

We then split and reshape this output to obtain h heads of keys $\widetilde{\mathbf{K}}_{\text{base}} \in \mathbb{R}^{n \times h \times d_h}$ and values $\widetilde{\mathbf{V}}_{\text{base}} \in \mathbb{R}^{n \times h \times d_h}$ for the base path.

Low-Rank Path. We up-project the h tiny latent heads using a head-wise matrix $\mathbf{B}_{lora} \in \mathbb{R}^{h \times r \times 2d_h}$ to form concatenated keys and values via Einstein summation over the latent dimension:

$$[\overline{\mathbf{K}}_{lora}, \overline{\mathbf{V}}_{lora}] = \operatorname{einsum}("\mathtt{nhr,hrd} \to \mathtt{nhd}", \mathbf{C}_{KV}, \mathbf{B}_{lora}),$$
 (5)

where the result is then split along the last dimension to obtain h heads of keys $\widetilde{\mathbf{K}}_{lora} \in \mathbb{R}^{n \times h \times d_h}$ and values $\widetilde{\mathbf{V}}_{lora} \in \mathbb{R}^{n \times h \times d_h}$ for the low-rank path.

Attention. To preserve translation equivariance, we concatenate the keys with partial RoPE \mathbf{K}_{RoPE} in both the base and low-rank paths.

$$\mathbf{K}_{\text{base}} = \operatorname{Concat}\left[\widetilde{\mathbf{K}}_{\text{base}}, \operatorname{repeat}\left(\mathbf{K}_{\text{RoPE}}, h\right)\right], \qquad \mathbf{K}_{\text{lora}} = \operatorname{Concat}\left[\alpha \widetilde{\mathbf{K}}_{\text{lora}}, \operatorname{repeat}\left(\mathbf{K}_{\text{RoPE}}, h\right)\right],$$
(6)

where repeat (\mathbf{K}_{ROPE} , h) replicates the partial RoPE across h heads to match the dimensionality required for concatenation. The scalar α is a tunable hyperparameter that controls the scale of the low-rank keys. Finally, we sum the base and low-rank attention outputs for each head i:

$$\mathbf{head}^{(i)} = \operatorname{Attention}\left(\mathbf{Q}^{(i)}, \mathbf{K}_{\text{base}}^{(i)}, \mathbf{V}_{\text{base}}^{(i)}\right) + \operatorname{Attention}\left(\mathbf{Q}^{(i)}, \mathbf{K}_{\text{lora}}^{(i)}, \alpha \mathbf{V}_{\text{lora}}^{(i)}\right), \tag{7}$$

where the computation of \mathbf{Q} can be found in Appendix B.

3.2 Decoding without KV Materialization

We refer to the DeepSeek official inference implementation (DeepSeek et al., 2024b) to illustrate how to "absorb" up-projection matrices into the queries and attention output to avoid explicit KV materialization in our MLRA decoding. Assume we have a query $\mathbf{q} = \operatorname{Concat}\left[\mathbf{q}_{\text{NoPE}}, \, \mathbf{q}_{\text{RoPE}}\right]$ and KV cache \mathbf{U}_{KV} , \mathbf{C}_{KV} , \mathbf{K}_{RoPE} . We split the up-projection matrices \mathbf{B}_{base} and \mathbf{B}_{lora} into per-head blocks, yielding $\mathbf{B}_{\text{base},\text{K}}^{(i)}, \mathbf{B}_{\text{base},\text{V}}^{(i)} \in \mathbb{R}^{d_u \times d_h}$ and $\mathbf{B}_{\text{lora},\text{K}}^{(i)}, \mathbf{B}_{\text{lora},\text{V}}^{(i)} \in \mathbb{R}^{r \times d_h}$. We likewise partition \mathbf{C}_{KV} along the head dimension as $\mathbf{C}_{\text{KV}} = \left[\mathbf{C}_{\text{KV}}^{(1)} \cdots \mathbf{C}_{\text{KV}}^{(h)}\right]$ with $\mathbf{C}_{\text{KV}}^{(i)} \in \mathbb{R}^{n \times r}$. For head i, letting $\tau = \frac{1}{\sqrt{d_h + d_h^R}}$, we compute the attention output for the base path:

$$Softmax \left(\tau \operatorname{Concat}\left[\mathbf{q}_{NoPE}^{(i)}, \mathbf{q}_{RoPE}^{(i)}\right] \left(\operatorname{Concat}\left[\mathbf{U}_{KV}\mathbf{B}_{base,K}^{(i)}, \mathbf{K}_{RoPE}\right]\right)^{\top}\right) \left(\mathbf{U}_{KV}\mathbf{B}_{base,V}^{(i)}\right)$$

$$= Softmax \left(\tau \mathbf{q}_{NoPE}^{(i)} \left(\mathbf{U}_{KV}\mathbf{B}_{base,K}^{(i)}\right)^{\top} + \tau \mathbf{q}_{RoPE}^{(i)}\mathbf{K}_{RoPE}^{\top}\right) \left(\mathbf{U}_{KV}\mathbf{B}_{base,V}^{(i)}\right)$$

$$= Softmax \left(\tau \mathbf{q}_{NoPE}^{(i)} \left(\mathbf{B}_{base,K}^{(i)}\right)^{\top} \mathbf{U}_{KV}^{\top} + \tau \mathbf{q}_{RoPE}^{(i)}\mathbf{K}_{RoPE}^{\top}\right) \left(\mathbf{U}_{KV}\mathbf{B}_{base,V}^{(i)}\right)$$

$$= Softmax \left(\tau \mathbf{q}_{NoPE}^{(i)} \left(\mathbf{B}_{base,K}^{(i)}\right)^{\top} \mathbf{U}_{KV}^{\top} + \tau \mathbf{q}_{RoPE}^{(i)}\mathbf{K}_{RoPE}^{\top}\right) \left(\mathbf{U}_{KV}\mathbf{B}_{base,V}^{(i)}\right)$$

$$= Softmax \left(\tau \mathbf{q}_{NoPE}^{(i)} \left(\mathbf{B}_{base,K}^{(i)}\right)^{\top} \mathbf{U}_{KV}^{\top} + \tau \mathbf{q}_{RoPE}^{(i)}\mathbf{K}_{RoPE}^{\top}\right) \mathbf{U}_{KV}\mathbf{B}_{base,V}^{(i)}.$$

$$(8)$$

Similarly, we compute the attention output for the low-rank path:

$$Softmax \left(\tau \operatorname{Concat}\left[\mathbf{q}_{NoPE}^{(i)}, \mathbf{q}_{RoPE}^{(i)}\right] \left(\operatorname{Concat}\left[\alpha \mathbf{C}_{KV}^{(i)} \mathbf{B}_{lora,K}^{(i)}, \mathbf{K}_{RoPE}\right]\right)^{\top}\right) \left(\alpha \mathbf{C}_{KV}^{(i)} \mathbf{B}_{lora,V}^{(i)}\right)$$

$$= Softmax \left(\tau \alpha \mathbf{q}_{NoPE}^{(i)} \left(\mathbf{C}_{KV}^{(i)} \mathbf{B}_{lora,K}^{(i)}\right)^{\top} + \tau \mathbf{q}_{RoPE}^{(i)} \mathbf{K}_{RoPE}^{\top}\right) \left(\alpha \mathbf{C}_{KV}^{(i)} \mathbf{B}_{lora,V}^{(i)}\right)$$

$$= Softmax \left(\tau \alpha \mathbf{q}_{NoPE}^{(i)} \left(\mathbf{B}_{lora,K}^{(i)}\right)^{\top} \left(\mathbf{C}_{KV}^{(i)}\right)^{\top} + \tau \mathbf{q}_{RoPE}^{(i)} \mathbf{K}_{RoPE}^{\top}\right) \left(\alpha \mathbf{C}_{KV}^{(i)} \mathbf{B}_{lora,V}^{(i)}\right)$$

$$= Softmax \left(\tau \alpha \mathbf{q}_{NoPE}^{(i)} \left(\mathbf{B}_{lora,K}^{(i)}\right)^{\top} \left(\mathbf{C}_{KV}^{(i)}\right)^{\top} + \tau \mathbf{q}_{RoPE}^{(i)} \mathbf{K}_{RoPE}^{\top}\right) \left(\alpha \mathbf{C}_{KV}^{(i)} \mathbf{B}_{lora,V}^{(i)}\right).$$

$$(9)$$

Below, we show how to avoid materializing h heads of keys and values during decoding by exploiting matrix multiplication associativity.

Step 1 (Keys/Logits). We exploit matrix-multiplication associativity by multiplying the query and the key up-projection matrix:

$$\tilde{\mathbf{q}}_{\text{base}}^{(i)} = \mathbf{q}_{\text{NoPE}}^{(i)} \left(\mathbf{B}_{\text{base},K}^{(i)} \right)^{\top} \in \mathbb{R}^{1 \times d_u}, \quad \boldsymbol{\lambda}_{\text{base}}^{(i)} = \operatorname{Softmax} \left(\tau \tilde{\mathbf{q}}_{\text{base}}^{(i)} \mathbf{U}_{\text{KV}}^{\top} + \tau \mathbf{q}_{\text{RoPE}}^{(i)} \mathbf{K}_{\text{RoPE}}^{\top} \right), \tag{10}$$

$$\tilde{\mathbf{q}}_{\text{lora}}^{(i)} = \alpha \mathbf{q}_{\text{NoPE}}^{(i)} \left(\mathbf{B}_{\text{lora},K}^{(i)} \right)^{\top} \in \mathbb{R}^{1 \times r}, \quad \boldsymbol{\lambda}_{\text{lora}}^{(i)} = \text{Softmax} \left(\tau \tilde{\mathbf{q}}_{\text{lora}}^{(i)} \left(\mathbf{C}_{\text{KV}}^{(i)} \right)^{\top} + \tau \mathbf{q}_{\text{RoPE}}^{(i)} \mathbf{K}_{\text{RoPE}}^{\top} \right). \quad (11)$$

This absorbs the key up-projection matrices into queries and avoids explicitly materializing keys $\mathbf{K}_{\text{base}}^{(i)} = \mathbf{U}_{\text{KV}} \mathbf{B}_{\text{base}, \text{K}}^{(i)} \in \mathbb{R}^{n \times 1 \times d_h}, \mathbf{K}_{\text{lora}}^{(i)} = \alpha \mathbf{C}_{\text{KV}}^{(i)} \mathbf{B}_{\text{lora}, \text{K}}^{(i)} \in \mathbb{R}^{n \times 1 \times d_h}.$ Note that \mathbf{U}_{KV} and \mathbf{K}_{RoPE} are shared across the h absorbed-query heads, so the base-path KV cache load is only $n \times (d_u + d_h^R)$.

Step 2 (Values/Output). After computing the attention scores, we use them to perform a weighted aggregation of compressed KV, then up-project the aggregated output:

$$\tilde{\mathbf{z}}_{\text{base}}^{(i)} = \boldsymbol{\lambda}_{\text{base}}^{(i)} \mathbf{U}_{\text{KV}} \in \mathbb{R}^{1 \times d_u}, \qquad \mathbf{z}_{\text{base}}^{(i)} = \tilde{\mathbf{z}}_{\text{base}}^{(i)} \mathbf{B}_{\text{base}, V}^{(i)} \in \mathbb{R}^{1 \times d_h}, \qquad (12)$$

$$\tilde{\mathbf{z}}_{\text{lora}}^{(i)} = \boldsymbol{\lambda}_{\text{lora}}^{(i)} \mathbf{C}_{\text{KV}}^{(i)} \in \mathbb{R}^{1 \times r}, \qquad \mathbf{z}_{\text{lora}}^{(i)} = \alpha \, \tilde{\mathbf{z}}_{\text{lora}}^{(i)} \mathbf{B}_{\text{lora}, V}^{(i)} \in \mathbb{R}^{1 \times d_h}. \qquad (13)$$

$$\tilde{\mathbf{z}}_{lora}^{(i)} = \boldsymbol{\lambda}_{lora}^{(i)} \mathbf{C}_{KV}^{(i)} \in \mathbb{R}^{1 \times r}, \qquad \mathbf{z}_{lora}^{(i)} = \alpha \, \tilde{\mathbf{z}}_{lora}^{(i)} \mathbf{B}_{lora,V}^{(i)} \in \mathbb{R}^{1 \times d_h}. \tag{13}$$

By deferring the up-projection for values with $\mathbf{B}_{\text{base,V}}^{(i)}$ and $\mathbf{B}_{\text{lora,V}}^{(i)}$, we never materialize $\mathbf{V}_{\text{base}}^{(i)} = \mathbf{U}_{\text{KV}}\mathbf{B}_{\text{base,V}}^{(i)} \in \mathbb{R}^{n \times 1 \times d_h}$ or $\mathbf{V}_{\text{lora}}^{(i)} = \mathbf{C}_{\text{KV}}^{(i)}\mathbf{B}_{\text{lora,V}}^{(i)} \in \mathbb{R}^{n \times 1 \times d_h}$. While the KV states are shared $(\mathbf{K} = \mathbf{V})$, FlashMLA (Jiashi Li, 2025) uses online softmax to load them once and avoid reloading.

3.3 Analysis

Translation Equivariance. We analyze the translation equivariance property of MLRA. Let $\mathbf{q}_m^{(i)} = \operatorname{Concat}\left(\mathbf{Q}_{\operatorname{NoPE}}\left[m,i,:\right],\ \mathbf{Q}_{\operatorname{RoPE}}\left[m,i,:\right]\right)$ and $\mathbf{k}_n^{(i)} = \operatorname{Concat}\left(\mathbf{K}_{\operatorname{NoPE}}\left[n,i,:\right],\ \mathbf{K}_{\operatorname{RoPE}}\left[n,:\right]\right)$ denote the query and key vectors for head i at positions m and n, respectively. $\mathbf{K}_{\operatorname{NoPE}}$ is either $\mathbf{K}_{\operatorname{base}}$ or \mathbf{K}_{lora} . The attention score for this head is given by the inner product:

$$\left\langle \mathbf{q}_{m}^{(i)}, \mathbf{k}_{n}^{(i)} \right\rangle = \left\langle \mathbf{Q}_{\text{NoPE}}\left[m, i, :\right], \mathbf{K}_{\text{NoPE}}\left[n, i, :\right] \right\rangle + \left\langle \mathbf{Q}_{\text{RoPE}}\left[m, i, :\right], \mathbf{K}_{\text{RoPE}}\left[n, :\right] \right\rangle. \tag{14}$$

Between the two terms, the second term with RoPE is position-dependent yet translation equivariant, due to RoPE's translation equivariance. The first term is position-independent and thus unchanged under joint translations of m and n. Therefore, the attention score $\langle \mathbf{q}_m^{(i)}, \mathbf{k}_n^{(i)} \rangle$ is invariant to translation in input positions. Although MLRA introduces positional inductive bias via partial RoPE and is translation equivariant, we refer to this property as semi-translation equivariance to distinguish it from the full RoPE with translation equivariance. Further details on other attention mechanisms and their translation equivariance analysis are provided in Appendix C.

Table 1: Comparison of parameters and KV cache among attention mechanisms. Some results are taken from Zhang et al. (2025) and Zadouri et al. (2025). For attention mechanism details, refer to Appendix C.

Method	KV Cache	# Parameters	# Query Heads	# KV Heads
MHA (Vaswani et al., 2017)	$2hd_h$	$4dhd_h$	h	h
MQA (Shazeer, 2019)	$2d_h$	$(2d + 2d_h) hd_h$	h	1
GQA (Ainslie et al., 2023)			h	g
MLA (DeepSeek et al., 2024a)	A (DeepSeek et al., 2024a) $ d_c + d_h^R \qquad \qquad d_c' \left(d + hd_h + hd_h^R \right) + dd_h^R \\ + d_c (d + 2hd_h) + dhd_h $		h	h
TPA (Zhang et al., 2025)	$2R_{KV}(h + d_h)$	$d(R_Q + 2R_{KV})(h + d_h) + dhd_h$	h	h
GLA (Zadouri et al., 2025)	(Zadouri et al., 2025) $ d_c + d_h^R \qquad \qquad d_c' \left(d + h d_h + h d_h^R \right) + d d_h^R \\ + d_c \left(d + h d_h \right) + d h d_h $		h	h
		$d(d_u + hr + d_h^R) + 2d_u h d_h + 2hr d_h + d(d'_u + hr') + d'_u h (d_h + d_h^R) + hr' (d_h + d_h^R) + dh d_h$	h	h
Method	Method KV Cache Per Token KV Cache Per Tok Per Device (2 GPU		KV Cache Per Token Per Device (4 GPUs)	KV Cache per token Per Device (8 GPUs)
MHA (Vaswani et al., 2017)	IHA (Vaswani et al., 2017) $64d_h$ $32d_h$		$16d_h$	$8d_h$
MQA (Shazeer, 2019)			$2d_h$	$2d_h$
GQA (Ainslie et al., 2023)			$4d_h$	$2d_h$
MLA (DeepSeek et al., 2024a)	$4.5d_h$	$4.5d_h$	$4.5d_h$	$4.5d_h$
TPA (Zhang et al., 2025)	FPA (Zhang et al., 2025) $4d_h + 256$ $4d_h + 128$		$4d_h + 64$	$4d_h + 32$
GLA (Zadouri et al., 2025)	GLA (Zadouri et al., 2025) $4.5d_h$ $2.5d_h$		$2.5d_h$	$2.5d_h$
2.07.70.1		A 2 4	2 2 4	

Table 2: Comparison of arithmetic intensity among attention mechanisms.

 $1.5d_{I}$

 $1.5d_{h}$

GLA (Zadouri et al., 2025) MLRA

 $4.5d_{\rm h}$

Method	MHA	MQA	GQA	MLA	TPA	GLA	MLRA
Arithmetic Intensity	$\frac{4nhd_h}{4nhd_h} \approx 1$	$\frac{\frac{4nhd_h}{4nd_h}}{\approx h}$	$\frac{\frac{4nhd_h}{4ngd_h}}{\approx \frac{h}{g}}$	$\frac{\frac{4nhd_c + 2nhd_h^R}{2n\left(d_c + d_h^R\right)}}{\approx 2h}$	$\frac{4nhR_{\text{KV}}d_h + 4nhd_h}{4nR_{\text{KV}}(h + d_h)} \approx \frac{3h}{4}$	$\frac{\frac{2nhd_c + 2nhd_h^R}{2n\left(d_c + d_h^R\right)}}{\approx h}$	$\frac{4nhr + 4nhd_h + 2nhd_h^R}{2n\left(hr + d_h + d_h^R\right)} \approx \frac{h}{2}$

KV Cache. We cache U_{KV} , C_{KV} , and K_{ROPE} during inference. Table 1 summarizes the total KV cache and parameters for the different attention mechanisms. While C_{KV} can be sharded across heads, \mathbf{U}_{KV} cannot. To support TP, we place \mathbf{U}_{KV} on one device and partition \mathbf{C}_{KV} across several other devices; consequently, $\mathbf{K}_{\mathrm{RoPE}}$ is replicated on every device. Therefore, KV cache under φ -way ($\varphi>1$) TP is $\frac{d_u+hr}{\varphi}+d_h^R$ per device for the low-rank path and $1.5d_h$ for the base path. We evaluate per-device KV cache under TP using LLaMA-3-70B (Llama et al., 2024) and Grok-2 as base models. Both models have 64 query heads and 8 key-value heads, so g = 8 for GQA. For GLA, we follow the paper's setup and consider two latent heads, each with a dimension of $2d_h$. For TPA, we follow the paper's setup and consider $R_{KV} = 2$.

Table 1 summarizes the per-device KV cache under TP as the number of devices varies. To support TP, the official MLA decoding implementation, FlashMLA (Jiashi Li, 2025), distributes upprojection matrices across devices by head. However, this approach replicates the KV cache on each device, so the per-device KV cache remains $4.5d_h$. TPA constructs its h key-value heads as linear combinations of $R_{\rm KV}$ shared heads. It supports TP only for the combination coefficients, while the shared heads are replicated across devices. So the per-device KV cache is $4d_h + \frac{256}{c}$. GLA splits the latent vector into two smaller latent vectors, and the per-device KV cache is $2.5d_h$ with 2-way TP. Note that for MLA with TP > 1 and for GLA with TP > 2, the KV cache is replicated, so the perdevice KV cache remains $4.5d_h$ and $2.5d_h$, respectively. For MLRA, when TP > 4, we can shard C_{KV} across additional devices and only replicate the partial RoPE. While GQA needs 8-way TP to reach $2d_h$ per device, MLRA achieves $1.5d_h$ with just 4-way TP, yielding the lowest per-device KV cache and faster decoding than other attention mechanisms.

Arithmetic Intensity. Arithmetic intensity (AI) is measured through the ratio of arithmetic operations to memory access (FLOPs per byte), which helps determine whether a workload is memorybound or compute-bound (Zadouri et al., 2025). We evaluate the arithmetic intensity of the attention mechanisms with $d_h = h = 128$ and $R_{KV} = 2$. Our analysis focuses on the long-context decoding, where the context length n dominates all other factors. We report AI of MLA, GLA, and MLRA without KV materialization. Implementation details of MLA decoding are provided in Appendix C.3. As shown in Table 2, our MLRA achieves relatively high AI. Under TP, the workload

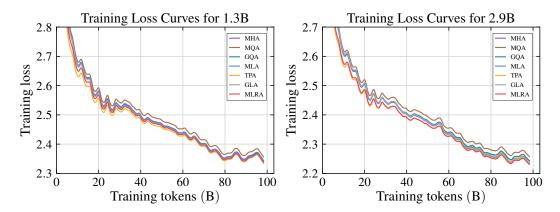


Figure 3: Training loss curves at the 1.3B and 2.9B scales for various attention mechanisms on the FineWeb-Edu 100B dataset.

is imbalanced across devices. In the base path, the AI is approximately 2h, matching that of MLA. In the low-rank path, if inference uses p devices, the AI is about $\frac{h}{p+6}$. This asymmetry motivates a heterogeneous deployment (Zhao et al., 2024b; Griggs et al., 2024; Jiang et al., 2024c, 2025). We can run the base path on a high-end GPU (e.g., H100) and distribute the low-rank path across multiple lower-cost GPUs (e.g., RTX 4090). Their reduced bandwidth and compute capabilities are sufficient given the low-rank path's low arithmetic intensity and the fact that its KV cache can be sharded, provided that the latency of the low-rank path does not exceed that of the base path.

4 Experiments

4.1 Experimental Setup

We train models at three scales: 354M, 1311M (1.3B), and 2873M (2.9B), on FineWeb-Edu-100B (Penedo et al., 2024). We report main results at the 1.3B and 2.9B scales, and use the 354M model for initialization ablations (Appendix E.2). Each model is pretrained from scratch on 98.3B tokens with an additional 0.1B tokens for validation. We adopt the LLaMA-3 architecture (Appendix D) and use six attention mechanisms as baselines: MHA (Vaswani et al., 2017), MQA (Shazeer, 2019), GQA (Ainslie et al., 2023), MLA (DeepSeek et al., 2024a), TPA (Zhang et al., 2025), and GLA (Zadouri et al., 2025). Architectural hyperparameters and training setup follow the GPT-3 configuration (Appendices E.1 and E.2). All models are trained on 8 NVIDIA H100 80G GPUs. Figure 3 shows the training loss curves for all attention mechanisms at the 1.3B and 2.9B scales.

4.2 Experimental Results

4.2.1 Model Quality

Validation Perplexity. We report validation perplexity for all attention mechanisms across six datasets (FineWeb-Edu (Penedo et al., 2024), Wikipedia, C4, Common Crawl, Pile (Gao et al., 2020), and Arxiv), using 100M tokens per dataset. As Table 3 shows, MLRA achieves the lowest average perplexity at the 2.9B (12.998) scale across six datasets. It achieves the best results on FineWeb-Edu, Wikipedia, and Pile, while placing second on C4, Common Crawl, and Arxiv. At the 1.3B scale, MLRA remains highly competitive (second-best average) and is best on Pile.

Downstream Evaluation. We evaluate zero-shot performance on standard benchmarks, including ARC-Easy (ARC-E) (Yadav et al., 2019), ARC-Challenge (ARC-C), BoolQ (Clark et al., 2019), HellaSwag (Zellers et al., 2019), PIQA (Bisk et al., 2020), and MMLU (Hendrycks et al., 2021), using the lm-evaluation-harness (Gao et al., 2024) package. We report normalized accuracy for ARC-Easy, ARC-Challenge, HellaSwag, and PIQA, and standard accuracy for the remaining tasks. As Table 4 shows, MLRA achieves the highest zero-shot average accuracy at both scales:

Table 3: Validation perplexity (lower is better) at the 1.3B and 2.9B scales on six datasets: FineWeb-Edu, Wikipedia, C4, CommonCrawl, Pile, and Arxiv. Best is **bold**; second best is underlined.

	1.3B												
Method	FineWeb-Edu	Wikipedia	C4	Common Crawl	Pile	Arxiv	Avg						
MHA	10.462	16.754	17.851	16.419	13.262	14.039	14.798						
MQA	10.592	17.039	18.113	16.580	12.701	14.161	14.864						
GQA	10.456	16.857	17.889	16.405	13.197	13.918	14.787						
MLA	10.449	15.881	17.758	16.143	12.240	13.657	14.355						
TPA	10.365	16.505	17.761	16.156	12.575	13.875	14.539						
GLA	<u>10.395</u>	16.548	17.725	16.057	12.463	13.595	14.464						
MLRA	10.402	16.170	17.787	16.186	12.234	13.705	14.414						
	2.9B												
Method	FineWeb-Edu	Wikipedia	C4	Common Crawl	Pile	Arxiv	Avg						
MHA	9.368	14.814	16.372	14.864	11.941	12.600	13.326						
MQA	9.499	15.232	16.664	15.173	11.886	12.972	13.571						
GQA	9.356	15.089	16.466	14.937	11.858	12.690	13.399						
MLA	9.280	14.901	16.278	14.673	11.190	12.370	13.115						
TPA	9.235	15.088	16.355	14.804	11.702	12.574	13.293						
GLA	9.298	<u>14.368</u>	16.381	14.792	11.242	12.417	13.083						
MLRA	9.224	14.214	16.287	14.702	11.147	12.412	12.998						

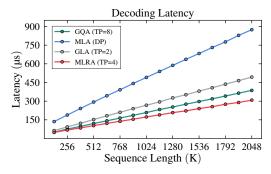
Table 4: Downstream evaluation at the 1.3B and 2.9B scales on six datasets: ARC-Easy, ARC-Challenge, BoolQ, HellaSwag, PIQA, and MMLU. Best is **bold**; second best is <u>underlined</u>.

•	-						
			1.	3B			
Method	ARC-E	ARC-C	BoolQ	HellaSwag	PIQA	MMLU	Avg
MHA	64.69	38.23	63.39	57.96	73.50	25.06	53.80
MQA	63.68	37.63	61.50	57.31	73.99	23.32	52.91
GQA	65.57	37.88	59.36	57.89	73.61	23.39	52.95
MLA	64.44	36.52	58.59	57.60	73.61	24.57	52.55
TPA	64.69	38.48	61.56	58.22	74.70	<u>25.37</u>	53.84
GLA	65.53	36.86	63.33	57.93	73.99	25.69	53.89
MLRA	66.08	37.29	63.09	57.76	74.27	25.26	53.96
			2.	9B			
Method	ARC-E	ARC-C	BoolQ	HellaSwag	PIQA	MMLU	Avg
MHA	68.14	41.13	63.39	61.97	75.30	25.79	55.95
MQA	69.49	40.27	63.12	61.11	75.14	23.69	55.47
GQA	69.74	40.36	62.05	61.83	75.52	24.58	55.68
MLA	70.58	41.04	59.24	62.20	75.03	24.65	55.46
TPA	70.12	40.70	64.50	62.27	76.33	27.04	56.83
GLA	69.65	40.96	63.55	62.44	<u>75.79</u>	25.35	56.29
MLRA	70.88	43.77	63.33	62.31	75.63	<u>26.43</u>	57.06

53.96% at 1.3B and 57.06% at 2.9B. At the 2.9B scale, it leads ARC-Easy and ARC-Challenge and places second on HellaSwag and MMLU. At the 1.3B scale, it leads ARC-Easy and places second on PIQA.

4.2.2 Efficiency

Decoding Speed. We benchmark decoding speed for GQA, MLA, GLA, and MLRA in long-context settings. All models use 64 heads with a head dimension of 128; for MLA, GLA, and MLRA, the partial RoPE dimension is 64. MLA is evaluated using DeepSeek's official implementation FlashMLA (Jiashi Li, 2025). GQA and GLA use FlashAttention kernels (Dao et al., 2022; Dao, 2024; Shah et al., 2024). We implement our MLRA kernel based on FlashAttention. We evaluate decoding speed across sequence lengths from 131,072 to 2,097,152 tokens (128K-2M). As shown in Figure 4, MLRA consistently outperforms all baselines at every length, yielding $1.05 \times -1.26 \times$ speedups over GQA. The gap grows with context length against GQA and GLA, while the speedup over MLA remains steady at about $2.8 \times$, indicating that MLRA with TP=4 substantially reduces long-context decoding latency.



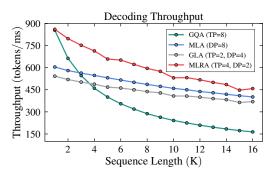


Figure 4: Decoding latency (lower is better) versus sequence length (batch=1) for GQA, MLA, GLA, and MLRA.

Figure 5: Decoding throughput versus sequence length (batch=128) for GQA, MLA, GLA, and MLRA.

Decoding Throughput. We evaluate decoding throughput for GQA, MLA, GLA, and MLRA on eight NVIDIA H100 80G GPUs, fixing the number of attention heads to 128 and the hidden size to 7168, following DeepSeekV3 (DeepSeek et al., 2024c). We set g = 16 for GQA. For MLA decoding deployment, there is a trade-off between data parallelism (DP) and tensor parallelism. With DP, we assign different requests to different devices, so attention parameters are replicated across devices and load can become imbalanced due to varying sequence lengths. With TP, the upprojection parameters are sharded by head, but the KV cache is replicated across devices. Following SGLang (Zheng et al., 2024), we otherwise strive to avoid KV cache duplication. Therefore, we use DP=8 for MLA, TP=2/DP=4 for GLA, TP=4/DP=2 for MLRA, and TP=8 for GQA. Throughput is reported for sequence lengths ranging from 1,024 to 16,384 tokens, and our end-to-end measurements include both the pre-attention stage that prepares inputs for the attention kernel and the attention computation itself. We accelerate pre-attention with torch.compile (Paszke et al., 2019) for MLA, GLA, and MLRA, and with custom Triton kernels for GQA. As shown in Figure 5, MLRA achieves the highest decoding throughput across both short and long sequence lengths. This suggests that MLRA with TP=4/DP=2 reduces parameter redundancy relative to MLA's DP=8, while introducing only modest partial RoPE duplication, thereby yielding higher throughput than MLA. For short sequences, GQA outperforms MLA and GLA because pre-attention dominates latency. However, MLRA remains competitive with GQA due to having even fewer query, key, and value parameters, as shown in Appendix E.1.

5 Conclusion

In this work, we propose Multi-Head Low-Rank Attention (MLRA), a TP-friendly dual-path attention mechanism. In the base path, we compress the KV cache into a single latent head and up-project it to produce multiple keys and values. In the low-rank path, we compress the KV cache into several tiny latent heads, each up-projected to generate its key and value. Our MLRA supports TP, reducing the per-device KV cache to $1.5d_h$, which is lower than that of GQA. At the 2.9 B scale, MLRA achieves state-of-the-art model quality across our benchmarks, including perplexity and downstream evaluations. We implement MLRA kernel based on FlashAttention. For long-context decoding (up to 2M tokens), MLRA achieves the lowest latency with 4-way TP and delivers the highest decoding throughput across sequence lengths from 1K to 16K tokens. Future work is to train MLRA at larger scales, such as 70B, to validate its effectiveness.

Acknowledgement

We thank Songlin Yang for helpful discussion.

References

- Ainslie, J., Lee-Thorp, J., de Jong, M., Zemlyanskiy, Y., Lebron, F., and Sanghai, S. (2023). GQA: Training generalized multi-query transformer models from multi-head checkpoints. In *Empirical Methods in Natural Language Processing*.
- Anagnostidis, S., Pavllo, D., Biggio, L., Noci, L., Lucchi, A., and Hofmann, T. (2023). Dynamic context pruning for efficient and interpretable autoregressive transformers. In Advances in Neural Information Processing Systems.
- Bisk, Y., Zellers, R., Bras, R. L., Gao, J., and Choi, Y. (2020). Piqa: Reasoning about physical commonsense in natural language. In *AAAI Conference on Artificial Intelligence*.
- Cai, R., Tian, Y., Wang, Z., and Chen, B. (2024). Lococo: Dropping in convolutions for long context compression. In *International Conference on Machine Learning*.
- Chang, C.-C., Lin, W.-C., Lin, C.-Y., Chen, C.-Y., Hu, Y.-F., Wang, P.-S., Huang, N.-C., Ceze, L., Abdelfattah, M. S., and Wu, K.-C. (2025). Palu: KV-cache compression with low-rank projection. In *International Conference on Learning Representations*.
- Chen, R., Wang, Z., Cao, B., Wu, T., Zheng, S., Li, X., Wei, X., Yan, S., Li, M., and Liang, Y. (2024a). Arkvale: Efficient generative LLM inference with recallable key-value eviction. In *Advances in Neural Information Processing Systems*.
- Chen, Y., Qian, S., Tang, H., Lai, X., Liu, Z., Han, S., and Jia, J. (2024b). LongloRA: Efficient fine-tuning of long-context large language models. In *International Conference on Learning Representations*.
- Clark, C., Lee, K., Chang, M.-W., Kwiatkowski, T., Collins, M., and Toutanova, K. (2019). BoolQ: Exploring the surprising difficulty of natural yes/no questions. In *North American Association for Computational Linguistics*.
- Dao, T. (2024). FlashAttention-2: Faster attention with better parallelism and work partitioning. In *International Conference on Learning Representations*.
- Dao, T., Fu, D. Y., Ermon, S., Rudra, A., and Re, C. (2022). Flashattention: Fast and memory-efficient exact attention with IO-awareness. In *Advances in Neural Information Processing Systems*.
- DeepSeek et al. (2024a). Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model. arXiv preprint arXiv:2405.04434.
- DeepSeek et al. (2024b). Deepseek-v3 technical report. https://github.com/deepseek-ai/DeepSeek-V3.
- DeepSeek et al. (2024c). Deepseek-v3 technical report. arXiv preprint arXiv:2412.19437.
- Dettmers, T., Pagnoni, A., Holtzman, A., and Zettlemoyer, L. (2023). Qlora: Efficient finetuning of quantized llms. In *Advances in Neural Information Processing Systems*.
- Dong, H., Yang, X., Zhang, Z., Wang, Z., Chi, Y., and Chen, B. (2024). Get more with less: Synthesizing recurrence with kv cache compression for efficient llm inference. In *International Conference on Machine Learning*.
- Fuchs, F., Worrall, D., Fischer, V., and Welling, M. (2020). Se (3)-transformers: 3d roto-translation equivariant attention networks. In *Advances in Neural Information Processing Systems*.
- Gao, L., Biderman, S., Black, S., Golding, L., Hoppe, T., Foster, C., Phang, J., He, H., Thite, A., Nabeshima, N., et al. (2020). The pile: An 800gb dataset of diverse text for language modeling. arXiv preprint arXiv:2101.00027.
- Gao, L., Tow, J., Abbasi, B., Biderman, S., Black, S., DiPofi, A., Foster, C., Golding, L., Hsu, J., Le Noac'h, A., Li, H., McDonell, K., Muennighoff, N., Ociepa, C., Phang, J., Reynolds, L., Schoelkopf, H., Skowron, A., Sutawika, L., Tang, E., Thite, A., Wang, B., Wang, K., and Zou, A. (2024). A framework for few-shot language model evaluation.

- Ge, S., Zhang, Y., Liu, L., Zhang, M., Han, J., and Gao, J. (2024). Model tells you what to discard: Adaptive KV cache compression for LLMs. In *International Conference on Learning Representations*.
- Gholami, A., Yao, Z., Kim, S., Hooper, C., Mahoney, M. W., and Keutzer, K. (2024). Ai and memory wall. *IEEE Micro*.
- Glorot, X. and Bengio, Y. (2010). Understanding the difficulty of training deep feedforward neural networks. In *International Conference on Artificial Intelligence and Statistics*.
- Griggs, T., Liu, X., Yu, J., Kim, D., Chiang, W.-L., Cheung, A., and Stoica, I. (2024). Mélange: Cost efficient large language model serving by exploiting gpu heterogeneity. arXiv preprint arXiv:2404.14527.
- He, J. and Zhai, J. (2024). Fastdecode: High-throughput gpu-efficient llm serving using heterogeneous pipelines. *arXiv preprint arXiv:2403.11421*.
- Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., and Steinhardt, J. (2021). Measuring massive multitask language understanding. In *International Conference on Learning Representations*.
- Hooper, C. R. C., Kim, S., Mohammadzadeh, H., Mahoney, M. W., Shao, S., Keutzer, K., and Gholami, A. (2024). KVQuant: Towards 10 million context length LLM inference with KV cache quantization. In *Advances in Neural Information Processing Systems*.
- Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., and Chen, W. (2022). LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*.
- Hu, J., Li, H., Zhang, Y., Wang, Z., Zhou, S., Zhang, X., Shum, H.-Y., and Jiang, D. (2024). Multimatrix factorization attention. *arXiv preprint arXiv:2412.19255*.
- Ivanov, A., Dryden, N., Ben-Nun, T., Li, S., and Hoefler, T. (2021). Data movement is all you need: A case study on optimizing transformers. In *Machine Learning and Systems*.
- Jiang, H., Li, Y., Zhang, C., Wu, Q., Luo, X., Ahn, S., Han, Z., Abdi, A. H., Li, D., Lin, C.-Y., Yang, Y., and Qiu, L. (2024a). MInference 1.0: Accelerating pre-filling for long-context LLMs via dynamic sparse attention. In *Advances in Neural Information Processing Systems*.
- Jiang, T., Huang, S., Luo, S., Zhang, Z., Huang, H., Wei, F., Deng, W., Sun, F., Zhang, Q., Wang, D., et al. (2024b). Mora: High-rank updating for parameter-efficient fine-tuning. arXiv preprint arXiv:2405.12130.
- Jiang, Y., Yan, R., Yao, X., Zhou, Y., Chen, B., and Yuan, B. (2024c). Hexgen: Generative inference of large language model over heterogeneous environment. In *International Conference on Machine Learning*.
- Jiang, Y., Yan, R., and Yuan, B. (2025). Hexgen-2: Disaggregated generative inference of LLMs in heterogeneous environment. In *International Conference on Learning Representations*.
- Jiashi Li, S. L. (2025). Flashmla: Efficient mla decoding kernels. https://github.com/deepseek-ai/FlashMLA.
- Kim, J.-H., Yeom, J., Yun, S., and Song, H. O. (2024). Compressed context memory for online language model interaction. In *International Conference on Learning Representations*.
- Kwon, W., Li, Z., Zhuang, S., Sheng, Y., Zheng, L., Yu, C. H., Gonzalez, J., Zhang, H., and Stoica, I. (2023). Efficient memory management for large language model serving with pagedattention. In *Symposium on Operating Systems Principles*.
- Lewis, P., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., Kuttler, H., Lewis, M., tau Yih, W., Rocktäschel, T., Riedel, S., and Kiela, D. (2020). Retrieval-augmented generation for knowledge-intensive nlp tasks. In *Advances in Neural Information Processing Systems*.

- Li, X., Ren, W., Qin, W., Wang, L., Zhao, T., and Hong, R. (2025). Analyzing and reducing catastrophic forgetting in parameter efficient tuning. In *International Conference on Acoustics*, *Speech and Signal Processing*.
- Li, Y., Huang, Y., Yang, B., Venkitesh, B., Locatelli, A., Ye, H., Cai, T., Lewis, P., and Chen, D. (2024). SnapKV: LLM knows what you are looking for before generation. In *Advances in Neural Information Processing Systems*.
- Lialin, V., Muckatira, S., Shivagunde, N., and Rumshisky, A. (2024). ReloRA: High-rank training through low-rank updates. In *International Conference on Learning Representations*.
- Liang, Y.-S. and Li, W.-J. (2024). Inflora: Interference-free low-rank adaptation for continual learning. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- Lin, C.-H., Gao, S., Smith, J. S., Patel, A., Tuli, S., Shen, Y., Jin, H., and Hsu, Y.-C. (2025). ModeGPT: Modular decomposition for large language model compression. In *International Conference on Learning Representations*.
- Liu, Y., Li, H., Cheng, Y., Ray, S., Huang, Y., Zhang, Q., Du, K., Yao, J., Lu, S., Ananthanarayanan, G., Maire, M., Hoffmann, H., Holtzman, A., and Jiang, J. (2024a). Cachegen: Kv cache compression and streaming for fast large language model serving. In *Conference of the ACM Special Interest Group on Data Communication*.
- Liu, Z., Desai, A., Liao, F., Wang, W., Xie, V., Xu, Z., Kyrillidis, A., and Shrivastava, A. (2023). Scissorhands: Exploiting the persistence of importance hypothesis for LLM KV cache compression at test time. In *Advances in Neural Information Processing Systems*.
- Liu, Z., Yuan, J., Jin, H., Zhong, S., Xu, Z., Braverman, V., Chen, B., and Hu, X. (2024b). KIVI: A tuning-free asymmetric 2bit quantization for KV cache. In *International Conference on Machine Learning*.
- Llama et al. (2024). The llama 3 herd of models. arXiv preprint arXiv:2407.21783.
- Loshchilov, I. and Hutter, F. (2016). Sgdr: Stochastic gradient descent with warm restarts. *arXiv* preprint arXiv:1608.03983.
- Loshchilov, I. and Hutter, F. (2017). Decoupled weight decay regularization. *arXiv preprint* arXiv:1711.05101.
- Malladi, S., Wettig, A., Yu, D., Chen, D., and Arora, S. (2023). A kernel-based view of language model fine-tuning. In *International Conference on Machine Learning*.
- Meng, F., Tang, P., Tang, X., Yao, Z., Sun, X., and Zhang, M. (2025). Transmla: Multi-head latent attention is all you need. *arXiv preprint arXiv:2502.07864*.
- Nawrot, P., Łańcucki, A., Chochowski, M., Tarjan, D., and Ponti, E. (2024). Dynamic memory compression: Retrofitting LLMs for accelerated inference. In *International Conference on Machine Learning*.
- Ootomo, H. and Yokota, R. (2023). Reducing shared memory footprint to leverage high throughput on tensor cores and its flexible api extension library. In *International Conference on High Performance Computing in Asia-Pacific Region*.
- OpenAI et al. (2024). Openai o1 system card. arXiv preprint arXiv:2412.16720.
- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al. (2019). Pytorch: An imperative style, high-performance deep learning library. In Advances in Neural Information Processing Systems, volume 32.
- Penedo, G., Kydlíček, H., allal, L. B., Lozhkov, A., Mitchell, M., Raffel, C., Werra, L. V., and Wolf, T. (2024). The fineweb datasets: Decanting the web for the finest text data at scale. In *Advances in Neural Information Processing Systems Datasets and Benchmarks Track*.
- Pope, R., Douglas, S., Chowdhery, A., Devlin, J., Bradbury, J., Heek, J., Xiao, K., Agrawal, S., and Dean, J. (2023). Efficiently scaling transformer inference. In *Machine Learning and Systems*.

- Sadhukhan, R., Chen, J., Chen, Z., Tiwari, V., Lai, R., Shi, J., Yen, I. E.-H., May, A., Chen, T., and Chen, B. (2025). Magicdec: Breaking the latency-throughput tradeoff for long context generation with speculative decoding. In *International Conference on Learning Representations*.
- Satorras, V. G., Hoogeboom, E., and Welling, M. (2021). E (n) equivariant graph neural networks. In *International Conference on Machine Learning*.
- Shah, J., Bikshandi, G., Zhang, Y., Thakkar, V., Ramani, P., and Dao, T. (2024). Flashattention-3: Fast and accurate attention with asynchrony and low-precision. In *Advances in Neural Information Processing Systems*.
- Shazeer, N. (2019). Fast transformer decoding: One write-head is all you need. *arXiv preprint arXiv:1911.02150*.
- Shi, Y., Wei, J., Wu, Y., Ran, R., Sun, C., He, S., and Yang, Y. (2024). Loldu: Low-rank adaptation via lower-diag-upper decomposition for parameter-efficient fine-tuning. *arXiv* preprint arXiv:2410.13618.
- Su, J., Ahmed, M., Lu, Y., Pan, S., Bo, W., and Liu, Y. (2024). Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*.
- Sun, H., Chang, L.-W., Bao, W., Zheng, S., Zheng, N., Liu, X., Dong, H., Chi, Y., and Chen, B. (2025a). ShadowKV: KV cache in shadows for high-throughput long-context LLM inference. In *International Conference on Machine Learning*.
- Sun, L., Deng, C., Jiang, J., Wu, X., Zhang, H., Chen, L., Ni, L., and Wang, J. (2025b). Gta: Grouped-head latent attention. *arXiv preprint arXiv:2506.17286*.
- Sun, Y., Dong, L., Zhu, Y., Huang, S., Wang, W., Ma, S., Zhang, Q., Wang, J., and Wei, F. (2024). You only cache once: Decoder-decoder architectures for language models. In *Advances in Neural Information Processing Systems*.
- Tang, J., Zhao, Y., Zhu, K., Xiao, G., Kasikci, B., and Han, S. (2024). QUEST: Query-aware sparsity for efficient long-context LLM inference. In *International Conference on Machine Learning*.
- Tang, X., Meng, F., Tang, P., Wang, Y., Yin, D., Sun, X., and Zhang, M. (2025). Tpla: Tensor parallel latent attention for efficient disaggregated prefill & decode inference. arXiv preprint arXiv:2508.15881.
- Thomas, N., Smidt, T., Kearnes, S., Yang, L., Li, L., Kohlhoff, K., and Riley, P. (2018). Tensor field networks: Rotation-and translation-equivariant neural networks for 3d point clouds. *arXiv* preprint arXiv:1802.08219.
- Vaswani, A., Shazeer, N. M., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. (2017). Attention is all you need. In *Advances in Neural Information Processing Systems*.
- Wang, X., Zheng, Y., Wan, Z., and Zhang, M. (2025). SVD-LLM: Truncation-aware singular value decomposition for large language model compression. In *International Conference on Learning Representations*.
- Wei, J., Wang, X., Schuurmans, D., Bosma, M., Xia, F., Chi, E., Le, Q. V., and Zhou, D. (2022). Chain-of-thought prompting elicits reasoning in large language models. In *Advances in Neural Information Processing Systems*.
- Weiler, M., Geiger, M., Welling, M., Boomsma, W., and Cohen, T. S. (2018). 3d steerable cnns: Learning rotationally equivariant features in volumetric data. In *Advances in Neural Information Processing Systems*.
- Williams, S., Waterman, A., and Patterson, D. (2009). Roofline: an insightful visual performance model for multicore architectures. *Communications of the ACM*.
- Xiao, G., Tang, J., Zuo, J., junxian guo, Yang, S., Tang, H., Fu, Y., and Han, S. (2025). Duoattention: Efficient long-context LLM inference with retrieval and streaming heads. In *International Conference on Learning Representations*.

- Xiao, G., Tian, Y., Chen, B., Han, S., and Lewis, M. (2024). Efficient streaming language models with attention sinks. In *International Conference on Learning Representations*.
- Yadav, V., Bethard, S., and Surdeanu, M. (2019). Quick and (not so) dirty: Unsupervised selection of justification sentences for multi-hop question answering. In *Empirical Methods in Natural Language Processing*.
- Yue, Y., Yuan, Z., Duanmu, H., Zhou, S., Wu, J., and Nie, L. (2024). Wkvquant: Quantizing weight and key/value cache for large language models gains more. *arXiv preprint arXiv:2402.12065*.
- Zadouri, T., Strauss, H., and Dao, T. (2025). Hardware-efficient attention for fast decoding. In *Conference on Language Modeling*.
- Zellers, R., Holtzman, A., Bisk, Y., Farhadi, A., and Choi, Y. (2019). Hellaswag: Can a machine really finish your sentence? In *Association for Computational Linguistics*.
- Zeng, Y. and Lee, K. (2024). The expressive power of low-rank adaptation. In *International Conference on Learning Representations*.
- Zhang, H. (2024). Sinklora: Enhanced efficiency and chat capabilities for long-context large language models. *arXiv preprint arXiv:2406.05678*.
- Zhang, Q., Chen, M., Bukharin, A., He, P., Cheng, Y., Chen, W., and Zhao, T. (2023a). Adaptive budget allocation for parameter-efficient fine-tuning. In *International Conference on Learning Representations*.
- Zhang, Y., Du, Y., Luo, G., Zhong, Y., Zhang, Z., Liu, S., and Ji, R. (2024). Cam: Cache merging for memory-efficient LLMs inference. In *International Conference on Machine Learning*.
- Zhang, Y., Liu, Y., Yuan, H., Qin, Z., Yuan, Y., Gu, Q., and Yao, A. C.-C. (2025). Tensor product attention is all you need. *arXiv preprint arXiv:2501.06425*.
- Zhang, Z., Sheng, Y., Zhou, T., Chen, T., Zheng, L., Cai, R., Song, Z., Tian, Y., Re, C., Barrett, C., Wang, Z., and Chen, B. (2023b). H2o: Heavy-hitter oracle for efficient generative inference of large language models. In *Advances in Neural Information Processing Systems*.
- Zhao, H., Ni, B., Fan, J., Wang, Y., Chen, Y., Meng, G., and Zhang, Z. (2024a). Continual forgetting for pre-trained vision models. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- Zhao, J., Wan, B., Wu, C., Peng, Y., and Lin, H. (2024b). Poster: Llm-pq: Serving llm on heterogeneous clusters with phase-aware partition and adaptive quantization. In *Symposium on Principles and Practice of Parallel Programming*.
- Zheng, C., Sun, J., Gao, Y., Wang, Y., Wang, P., Xiong, J., Ren, L., Cheng, H., Kulkarni, J., Shen, Y., et al. (2025). Sas: Simulated attention score. *arXiv preprint arXiv:2507.07694*.
- Zheng, L., Yin, L., Xie, Z., Sun, C., Huang, J., Yu, C. H., Cao, S., Kozyrakis, C., Stoica, I., Gonzalez, J. E., Barrett, C., and Sheng, Y. (2024). SGLang: Efficient execution of structured language model programs. In *Advances in Neural Information Processing Systemss*.
- Zhu, J., Greenewald, K., Nadjahi, K., Borde, H. S. D. O., Gabrielsson, R. B., Choshen, L., Ghassemi, M., Yurochkin, M., and Solomon, J. (2024). Asymmetry in low-rank adapters of foundation models. In *International Conference on Machine Learning*.

Appendix

A	Theorem	16
В	3 Query Computation in MLRA	16
C	C Attention Mechanism	17
	C.1 Multi-Head Attention (MHA)	. 17
	C.2 Multi-Query Attention (MQA) and Grouped-Query Attention (GQA)	. 17
	C.3 Multi-Head Latent Attention (MLA)	. 18
	C.4 Tensor Product Attention (TPA)	. 19
D	Llama-3 Architecture	20
E	Experimental Setup	20
	E.1 Model Hyperparameters	. 20
	E.2 Training Setup	. 22
F	Related Work	22

A Theorem

Theorem 1. Given two tokens with query \mathbf{q} and key \mathbf{k} at positions m and n, respectively, and applying RoPE to obtain RoPE (\mathbf{q}, m) and RoPE (\mathbf{k}, n) , we aim to prove that RoPE is translation equivariant with respect to the inner product. In the definition of translation equivariance (Eq.(1)), this corresponds to setting the output transformation S_s as the identity map, i.e., $S_s = \mathbf{I}$. Specifically, we show that translating both positions by an offset s leaves the inner product unchanged:

$$\langle \text{RoPE}(\mathbf{q}, m + s), \text{RoPE}(\mathbf{k}, n + s) \rangle = \langle \text{RoPE}(\mathbf{q}, m), \text{RoPE}(\mathbf{k}, n) \rangle.$$

Proof. We compute the inner product under RoPE as:

$$\left(\mathbf{q}\mathbf{R}_{m}\right)\left(\mathbf{k}\mathbf{R}_{n}\right)^{\top} = \mathbf{q}\mathbf{R}_{m}\mathbf{R}_{n}^{\top}\mathbf{k}^{\top} = \operatorname{Re}\left[\sum_{i=1}^{d/2}\mathbf{q}_{\left[2i-1:2i+1\right]}\mathbf{k}_{\left[2i-1:2i+1\right]}^{*}e^{i(m-n)\theta_{i}}\right],$$

where \mathbf{R}_m is the rotary matrix at position m. Now consider translating both tokens by an offset s. The difference in positions becomes:

$$(m+s) - (n+s) = m - n,$$

which leaves the term $e^{i(m-n)\theta_i}$ unchanged. Therefore, the inner product remains invariant under a translation of both tokens, and we have:

$$\phi(T_t(\mathbf{x})) = \phi(\mathbf{x}), \text{ with } S_s = \mathbf{I},$$

which proves that RoPE is translation equivariant with respect to the inner product.

B Query Computation in MLRA

The hidden states $\mathbf{H} \in \mathbb{R}^{n \times d}$ are first passed through a linear down-projection $\mathbf{A}_Q \in \mathbb{R}^{d \times \left(d'_u + hr'\right)}$ yielding two parts

$$\left[\mathbf{U}_{\mathsf{Q}},\ \overline{\mathbf{C}}_{\mathsf{Q}}\right] = \mathbf{H}\mathbf{A}_{\mathsf{Q}} \in \mathbb{R}^{n \times \left(d'_u + hr'\right)}$$

By default, we set $d'_u = 2d_h$ and $r' = \frac{6d_h}{h}$. The base latent head $\mathbf{U}_{\mathbf{Q}} \in \mathbb{R}^{n \times d'_u}$ is up-projected by a weight matrix $\mathbf{B}^{\mathbf{Q}}_{\mathrm{base}} \in \mathbb{R}^{d'_u \times h \left(d_h + d_h^R\right)}$ producing

$$\overline{\mathbf{Q}}_{\mathrm{base}} = \mathbf{U}_{\mathrm{Q}} \, \mathbf{B}_{\mathrm{base}}^{\mathrm{Q}} \in \mathbb{R}^{n \times h \left(d_h + d_h^R\right)},$$

which reshapes and splits into

$$\mathbf{Q}_{\mathrm{base,NoPE}} \in \mathbb{R}^{n \times h \times d_h}, \qquad \mathbf{Q}_{\mathrm{base,RoPE}} \in \mathbb{R}^{n \times h \times d_h^R}.$$

 $\overline{\mathbf{C}}_{0}$ is reshaped to expose the head dimension to obtain h latent heads,

$$\mathbf{C}_{\mathbf{Q}} = \text{reshape}\left(\overline{\mathbf{C}}_{\mathbf{Q}}, [n, h, r']\right) \in \mathbb{R}^{n \times h \times r'},$$

and up-projected head-wise by a weight matrix $\mathbf{B}_{lora}^Q \in \mathbb{R}^{h \times r' \times \left(d_h + d_h^R\right)}$. This produces an intermediate query tensor

$$\mathbf{Q}_{\mathrm{lora}} = \mathrm{einsum}\left(exttt{"nhr,hrd}
ightarrow exttt{nhd"}, \mathbf{C}_{\mathrm{Q}}, \ \mathbf{B}_{\mathrm{lora}}^{\mathrm{Q}}
ight) \in \mathbb{R}^{n imes h imes \left(d_h + d_h^R
ight)},$$

which is split into

$$\mathbf{Q}_{\text{lora,NoPE}} \in \mathbb{R}^{n \times h \times d_h}, \qquad \mathbf{Q}_{\text{lora,RoPE}} \in \mathbb{R}^{n \times h \times d_h^R}.$$

The low-rank and base paths are fused:

$$\mathbf{Q}_{\text{NoPE}} = \gamma \mathbf{Q}_{\text{lora,NoPE}} + \mathbf{Q}_{\text{base,NoPE}}, \qquad \widetilde{\mathbf{Q}}_{\text{RoPE}} = \gamma \mathbf{Q}_{\text{lora,RoPE}} + \mathbf{Q}_{\text{base,RoPE}}$$

RoPE is applied to $\widetilde{\mathbf{Q}}_{RoPE}$ to obtain \mathbf{Q}_{RoPE} . This result is then concatenated with the other query components to form the final query:

$$\mathbf{Q} = \operatorname{Concat}\left[\mathbf{Q}_{\operatorname{NoPE}}, \; \mathbf{Q}_{\operatorname{RoPE}}\right] \in \mathbb{R}^{n \times h \times \left(d_h + d_h^R\right)}.$$

C Attention Mechanism

C.1 Multi-Head Attention (MHA)

Consider a sequence of n tokens with hidden states $\mathbf{H} \in \mathbb{R}^{n \times d}$. We first project these hidden states into query, key, and value representations using projection matrices $\mathbf{W}_0, \mathbf{W}_K, \mathbf{W}_V \in \mathbb{R}^{d \times (hd_h)}$:

$$\overline{\mathbf{Q}} = \text{RoPE}(\mathbf{H}\mathbf{W}_{\mathbf{Q}}), \quad \overline{\mathbf{K}} = \text{RoPE}(\mathbf{H}\mathbf{W}_{\mathbf{K}}), \quad \overline{\mathbf{V}} = \mathbf{H}\mathbf{W}_{\mathbf{V}},$$

where $\overline{\mathbf{Q}}, \overline{\mathbf{K}}, \overline{\mathbf{V}} \in \mathbb{R}^{n \times (hd_h)}$, h is the number of attention heads, and d_h is the dimensionality of each head. Next, we reshape these matrices to separate the head dimension:

$$\mathbf{Q} = \operatorname{reshape}(\overline{\mathbf{Q}}, [n, h, d_h]), \quad \mathbf{K} = \operatorname{reshape}(\overline{\mathbf{K}}, [n, h, d_h]), \quad \mathbf{V} = \operatorname{reshape}(\overline{\mathbf{V}}, [n, h, d_h]),$$

such that $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{n \times h \times d_h}$. In MHA, the keys \mathbf{K} and values \mathbf{V} from previous tokens are cached to accelerate inference.

Remark 2. Let $\mathbf{Q}, \mathbf{K} \in \mathbb{R}^{n \times h \times d_h}$ denote the RoPE-encoded queries and keys after projection and reshaping, and define per-head query and key vectors as:

$$\mathbf{q}_{m}^{(i)} := \mathbf{Q}[m, i, :], \quad \mathbf{k}_{n}^{(i)} := \mathbf{K}[n, i, :],$$

where m and n are token positions and i is the attention head index. Then, for any translation t, it follows from Theorem 1 that:

$$\left\langle \mathbf{q}_{m+t}^{(i)},\ \mathbf{k}_{n+t}^{(i)}\right\rangle =\left\langle \mathbf{q}_{m}^{(i)},\ \mathbf{k}_{n}^{(i)}\right\rangle .$$

C.2 Multi-Query Attention (MQA) and Grouped-Query Attention (GQA)

Both MQA and GQA reduce the number of key and value heads compared to MHA, while maintaining the full number of query heads. MQA takes this to the extreme by using a single key-value head for all query heads, whereas GQA partitions the query heads into groups that each share a key-value head. Given a sequence of n tokens with hidden states $\mathbf{H} \in \mathbb{R}^{n \times d}$, the queries are computed using the same projection as in MHA:

$$\overline{\mathbf{Q}} = \text{RoPE}\left(\mathbf{HW}_{\mathbf{Q}}\right), \quad \overline{\mathbf{Q}} \in \mathbb{R}^{n \times (hd_h)}.$$

To reduce KV cache footprint, we use projection matrices \mathbf{W}_{K} , $\mathbf{W}_{V} \in \mathbb{R}^{d \times d_{h}g}$, where g < h (e.g., g = 1 for MQA), to compute:

$$\overline{\mathbf{K}}_{C} = \text{RoPE}\left(\mathbf{H}\mathbf{W}_{K}\right), \quad \overline{\mathbf{V}}_{C} = \mathbf{H}\mathbf{W}_{V}, \quad \overline{\mathbf{K}}_{C}, \overline{\mathbf{V}}_{C} \in \mathbb{R}^{n \times d_{h}g}.$$

These are then reshaped into per-head form:

$$\mathbf{K}_{\mathrm{C}} = \mathrm{reshape}(\overline{\mathbf{K}}_{\mathrm{C}}, \; [n,g,d_h]), \quad \mathbf{V}_{\mathrm{C}} = \mathrm{reshape}(\overline{\mathbf{V}}_{\mathrm{C}}, \; [n,g,d_h]).$$

We cache $\overline{\mathbf{K}}_{\mathbf{C}}$ and $\overline{\mathbf{V}}_{\mathbf{C}}$ during inference. During decoding, we first repeat them along the head axis to match the h query heads (repeat factor r = h/g):

$$\mathbf{K} = \text{RepeatInterleave}(\mathbf{K}_{\mathbf{C}}, r, \text{dim} = 1) \in \mathbb{R}^{n \times h \times d_h}, \quad \mathbf{V} = \text{RepeatInterleave}(\mathbf{V}_{\mathbf{C}}, r, \text{dim} = 1).$$

Remark 3. Let $\mathbf{Q} \in \mathbb{R}^{n \times h \times d_h}$ and $\mathbf{K}_C \in \mathbb{R}^{n \times g \times d_h}$ denote the RoPE-encoded queries and keys, respectively. For attention head i, we define the query and key vectors at positions m and n as:

$$\mathbf{q}_m^{(i)} := \mathbf{Q}[m,i,:], \quad \mathbf{k}_n^{(i)} := \mathbf{K}_{\mathcal{C}}[n,j,:] \text{ for some } j \in \{1,\ldots,g\}.$$

Since each $\mathbf{k}_n^{(j)}$ is obtained via RoPE, and RoPE preserves inner products under joint position translations, we have, for any offset t:

$$\left\langle \mathbf{q}_{m+t}^{(i)},\ \mathbf{k}_{n+t}^{(i)}\right\rangle =\left\langle \mathbf{q}_{m}^{(i)},\ \mathbf{k}_{n}^{(i)}\right\rangle ,$$

where $\mathbf{k}_n^{(i)} = \mathbf{K}_C[n, j, :]$ for some $j \in \{1, \dots, g\}$.

C.3 Multi-Head Latent Attention (MLA)

Given a sequence of n tokens with hidden states $\mathbf{H} \in \mathbb{R}^{n \times d}$, MLA first computes queries as:

$$\begin{split} \mathbf{C}_{Q} &= \mathbf{H} \mathbf{A}_{Q}, \quad \overline{\mathbf{Q}}_{NoPE} = \mathbf{C}_{Q} \mathbf{B}_{Q}, \quad \overline{\mathbf{Q}}_{RoPE} = RoPE \left(\mathbf{C}_{Q} \mathbf{B}_{QR} \right), \\ \mathbf{A}_{Q} &\in \mathbb{R}^{d \times d_{c}'}, \quad \mathbf{B}_{Q} \in \mathbb{R}^{d_{c}' \times (hd_{h})}, \quad \mathbf{B}_{QR} \in \mathbb{R}^{d_{c}' \times d_{h}^{R}h}, \end{split}$$

where $d_c' \ll hd_h$ is a low-rank latent dimension.

We then reshape the query to separate heads:

$$\mathbf{Q}_{\text{NoPE}} = \text{reshape}(\overline{\mathbf{Q}}_{\text{C}}, [n, h, d_h]), \quad \mathbf{Q}_{\text{RoPE}} = \text{reshape}(\overline{\mathbf{Q}}_{\text{RoPE}}, [n, h, d_h^R]),$$

where $\mathbf{Q}_{\text{NoPE}} \in \mathbb{R}^{n \times h \times d_h}$ and $\mathbf{Q}_{\text{RoPE}} \in \mathbb{R}^{n \times h \times d_h^R}$. These are concatenated along the last dimension to form the final query:

$$\mathbf{Q} = \operatorname{Concat}\left[\mathbf{Q}_{\operatorname{NoPE}}, \ \mathbf{Q}_{\operatorname{RoPE}}\right] \in \mathbb{R}^{n \times h \times (d_h + d_h^R)}$$
.

To reduce the KV cache, MLA obtains shared compressed KV states via a down-projection:

$$\begin{split} \mathbf{C}_{\mathsf{KV}} &= \mathbf{H} \mathbf{A}_{\mathsf{KV}}, \quad \mathbf{A}_{\mathsf{KV}} \in \mathbb{R}^{d \times d_c}, \quad \mathbf{C}_{\mathsf{KV}} \in \mathbb{R}^{n \times d_c}, \\ \mathbf{K}_{\mathsf{RoPE}} &= \mathsf{RoPE}\left(\mathbf{H} \mathbf{A}_{\mathsf{KR}}\right), \quad \mathbf{A}_{\mathsf{KR}} \in \mathbb{R}^{d \times d_h^R}, \quad \mathbf{K}_{\mathsf{RoPE}} \in \mathbb{R}^{n \times d_h^R}. \end{split}$$

where $d_c \ll h d_h$ is the latent dimension. Both \mathbf{C}_{KV} and $\mathbf{K}_{\mathrm{RoPE}}$ are stored in the KV cache during inference.

At runtime, MLA reconstructs the full keys and values using learned up-projection matrices:

$$\overline{\mathbf{K}}_{\mathrm{NoPE}} = \mathbf{C}_{\mathrm{KV}} \mathbf{B}_{\mathrm{K}}, \quad \overline{\mathbf{V}} = \mathbf{C}_{\mathrm{KV}} \mathbf{B}_{\mathrm{V}}, \quad \mathbf{B}_{\mathrm{K}}, \mathbf{B}_{\mathrm{V}} \in \mathbb{R}^{d_c \times (hd_h)}$$

These are reshaped into per-head form:

$$\mathbf{K}_{\text{NoPE}} = \text{reshape}(\overline{\mathbf{K}}_{\text{NoPE}}, [n, h, d_h]), \quad \mathbf{V} = \text{reshape}(\overline{\mathbf{V}}, [n, h, d_h]),$$

where $\mathbf{K}_{\text{NoPE}} \in \mathbb{R}^{n \times h \times d_h}$ and $\mathbf{V} \in \mathbb{R}^{n \times h \times d_h}$.

To obtain per-head position-aware keys, MLA repeats the position encoding across heads:

$$\mathbf{K} = \operatorname{Concat} \left[\mathbf{K}_{\text{NoPE}}, \operatorname{repeat} \left(\mathbf{K}_{\text{RoPE}}, h \right) \right].$$

The translation equivariance analysis is the same as in MLRA.

Decoding without KV materialization. Given a query $\mathbf{q} = \operatorname{Concat}\left[\mathbf{q}_{\operatorname{NoPE}}, \mathbf{q}_{\operatorname{RoPE}}\right] \in \mathbb{R}^{h \times \left(d_h + d_h^R\right)}$ and KV cache: $\mathbf{C}_{\operatorname{KV}} \in \mathbb{R}^{n \times d_c}, \mathbf{K}_{\operatorname{RoPE}} \in \mathbb{R}^{n \times d_h}, \operatorname{MLA}$ avoids KV materialization and instead operates in the compressed space. Let $\mathbf{B}_{\operatorname{K}}, \mathbf{B}_{\operatorname{V}} \in \mathbb{R}^{d_c \times (hd_h)}$ be partitioned per head as $\mathbf{B}_{\operatorname{K}} = \left[\mathbf{B}_{\operatorname{K}}^{(1)} \cdots \mathbf{B}_{\operatorname{K}}^{(h)}\right], \mathbf{B}_{\operatorname{V}} = \left[\mathbf{B}_{\operatorname{V}}^{(1)} \cdots \mathbf{B}_{\operatorname{V}}^{(h)}\right]$ with $\mathbf{B}_{\operatorname{K}}^{(i)}, \mathbf{B}_{\operatorname{V}}^{(i)} \in \mathbb{R}^{d_c \times d_h}$. For head i we compute with scale $\tau = \frac{1}{\sqrt{d_h + d_h^R}}$:

$$\begin{split} &\operatorname{Softmax}\left(\tau\operatorname{Concat}\left[\mathbf{q}_{\operatorname{NoPE}}^{(i)},\mathbf{q}_{\operatorname{RoPE}}^{(i)}\right]\left(\operatorname{Concat}\left[\mathbf{C}_{\operatorname{KV}}\mathbf{B}_{\operatorname{K}}^{(i)},\mathbf{K}_{\operatorname{RoPE}}\right]\right)^{\top}\right)\left(\mathbf{C}_{\operatorname{KV}}\mathbf{B}_{\operatorname{V}}^{(i)}\right),\\ &=\operatorname{Softmax}\left(\tau\mathbf{q}_{\operatorname{NoPE}}^{(i)}\left(\mathbf{C}_{\operatorname{KV}}\mathbf{B}_{\operatorname{K}}^{(i)}\right)^{\top}+\tau\mathbf{q}_{\operatorname{RoPE}}^{(i)}\mathbf{K}_{\operatorname{RoPE}}^{\top}\right)\left(\mathbf{C}_{\operatorname{KV}}\mathbf{B}_{\operatorname{V}}^{(i)}\right),\\ &=\operatorname{Softmax}\left(\tau\mathbf{q}_{\operatorname{NoPE}}^{(i)}\left(\mathbf{B}_{\operatorname{K}}^{(i)}\right)^{\top}\mathbf{C}_{\operatorname{KV}}^{\top}+\tau\mathbf{q}_{\operatorname{RoPE}}^{(i)}\mathbf{K}_{\operatorname{RoPE}}^{\top}\right)\left(\mathbf{C}_{\operatorname{KV}}\mathbf{B}_{\operatorname{V}}^{(i)}\right),\\ &=\operatorname{Softmax}\left(\tau\mathbf{q}_{\operatorname{NoPE}}^{(i)}\left(\mathbf{B}_{\operatorname{K}}^{(i)}\right)^{\top}\mathbf{C}_{\operatorname{KV}}^{\top}+\tau\mathbf{q}_{\operatorname{RoPE}}^{(i)}\mathbf{K}_{\operatorname{RoPE}}^{\top}\right)\mathbf{C}_{\operatorname{KV}}\mathbf{B}_{\operatorname{V}}^{(i)}. \end{split}$$

We compute attention without materializing per-head K, V by exploiting matrix multiplication associativity.

Step 1 (Keys/Logits): We absorb the key's up-projection matrix into the query, resulting in a score computation that mimics MOA with shared KV states. For head *i*.

$$\tilde{\mathbf{q}}^{(i)} = \mathbf{q}_{\text{NoPE}}^{(i)} \big(\mathbf{B}_{\text{K}}^{(i)} \big)^{\!\top}, \qquad \boldsymbol{\lambda}^{(i)} = \operatorname{Softmax} \Big(\tau \tilde{\mathbf{q}}^{(i)} \mathbf{C}_{\text{KV}}^{\!\top} + \tau \mathbf{q}_{\text{RoPE}}^{(i)} \mathbf{K}_{\text{RoPE}}^{\!\top} \Big).$$

This avoids forming $\mathbf{K}_{NoPE}^{(i)} = \mathbf{C}_{KV} \mathbf{B}_{K}^{(i)}$, and the RoPE keys \mathbf{K}_{RoPE} are shared across heads (no head-wise repeat).

Step 2 (Values/Output): Aggregate with shared KV, then up-project once:

$$\tilde{\mathbf{z}}^{(i)} = \boldsymbol{\lambda}^{(i)} \, \mathbf{C}_{\mathrm{KV}}, \qquad \mathbf{z}^{(i)} = \tilde{\mathbf{z}}^{(i)} \, \mathbf{B}_{\mathrm{V}}^{(i)}.$$

By deferring the right multiplication by $\mathbf{B}_{\mathrm{V}}^{(i)}$, we never materialize $\mathbf{V}^{(i)} = \mathbf{C}_{\mathrm{KV}} \mathbf{B}_{\mathrm{V}}^{(i)}$. Consequently, decoding maintains only $\mathbf{C}_{\mathrm{KV}} \in \mathbb{R}^{n \times d_c}$ and $\mathbf{K}_{\mathrm{RoPE}} \in \mathbb{R}^{n \times d_h^R}$ caches instead of per-head full \mathbf{K}, \mathbf{V} .

C.4 Tensor Product Attention (TPA)

TPA achieves KV cache compression through low-rank factorization. It decomposes the keys and values into two components: a set of r shared basis vectors (or "attention heads") and token-specific coefficient matrices. The keys and values are reconstructed by linearly combining the shared basis vectors according to these coefficients. This allows storing only the coefficients and basis vectors in the cache.

Given a sequence of n tokens with hidden states $\mathbf{H} \in \mathbb{R}^{n \times d}$, TPA first computes the query, key, and value as follows:

$$\begin{split} & \overline{\mathbf{A}}_{Q} = \mathbf{H} \mathbf{W}_{AQ}, \quad \mathbf{W}_{AQ} \in \mathbb{R}^{d \times R_{Q}h}, \quad \overline{\mathbf{A}}_{Q} \in \mathbb{R}^{n \times R_{Q}h}, \\ & \overline{\mathbf{B}}_{Q} = \mathbf{H} \mathbf{W}_{BQ}, \quad \mathbf{W}_{BQ} \in \mathbb{R}^{d \times R_{Q}d_{h}}, \quad \overline{\mathbf{B}}_{Q} \in \mathbb{R}^{n \times R_{Q}d_{h}}, \\ & \overline{\mathbf{A}}_{K} = \mathbf{H} \mathbf{W}_{AK}, \quad \mathbf{W}_{AK} \in \mathbb{R}^{d \times R_{KV}h}, \quad \overline{\mathbf{A}}_{K} \in \mathbb{R}^{n \times R_{KV}h}, \\ & \overline{\mathbf{B}}_{K} = \mathbf{H} \mathbf{W}_{BK}, \quad \mathbf{W}_{BK} \in \mathbb{R}^{d \times R_{KV}d_{h}}, \quad \overline{\mathbf{B}}_{K} \in \mathbb{R}^{n \times R_{KV}d_{h}}, \\ & \overline{\mathbf{A}}_{V} = \mathbf{H} \mathbf{W}_{AV}, \quad \mathbf{W}_{AV} \in \mathbb{R}^{d \times R_{KV}d_{h}}, \quad \overline{\mathbf{A}}_{V} \in \mathbb{R}^{n \times R_{KV}d_{h}}, \\ & \overline{\mathbf{B}}_{V} = \mathbf{H} \mathbf{W}_{BV}, \quad \mathbf{W}_{BV} \in \mathbb{R}^{d \times R_{KV}d_{h}}, \quad \overline{\mathbf{B}}_{V} \in \mathbb{R}^{n \times R_{KV}d_{h}}, \end{split}$$

We reshape the projections into 3D tensors:

$$\begin{split} \mathbf{A}_{\mathrm{Q}} &= \mathrm{reshape}\left(\overline{\mathbf{A}}_{\mathrm{Q}},\; [n,\; R_{\mathrm{Q}},\; h]\right), \quad \mathbf{B}_{\mathrm{Q}} = \mathrm{reshape}\left(\overline{\mathbf{B}}_{\mathrm{Q}},\; [n,\; R_{\mathrm{Q}},\; d_h]\right), \\ \mathbf{A}_{\mathrm{K}} &= \mathrm{reshape}\left(\overline{\mathbf{A}}_{\mathrm{K}},\; [n,\; R_{\mathrm{KV}},\; h]\right), \quad \mathbf{B}_{\mathrm{K}} = \mathrm{reshape}\left(\overline{\mathbf{B}}_{\mathrm{K}},\; [n,\; R_{\mathrm{KV}},\; d_h]\right), \\ \mathbf{A}_{\mathrm{V}} &= \mathrm{reshape}\left(\overline{\mathbf{A}}_{\mathrm{V}},\; [n,\; R_{\mathrm{KV}},\; h]\right), \quad \mathbf{B}_{\mathrm{V}} = \mathrm{reshape}\left(\overline{\mathbf{B}}_{\mathrm{V}},\; [n,\; R_{\mathrm{KV}},\; d_h]\right), \end{split}$$

so that $\mathbf{A}_Q \in \mathbb{R}^{n \times R_Q \times h}$, $\mathbf{B}_Q \in \mathbb{R}^{n \times R_Q \times d_h}$, $\mathbf{A}_K \in \mathbb{R}^{n \times R_{KV} \times h}$, $\mathbf{B}_K \in \mathbb{R}^{n \times R_{KV} \times d_h}$, $\mathbf{A}_V \in \mathbb{R}^{n \times R_{KV} \times d_h}$, and $\mathbf{B}_V \in \mathbb{R}^{n \times R_{KV} \times d_h}$. The final query, key, and value are computed as:

$$\mathbf{Q}_{[i,:,:]} = \frac{1}{R_{\mathbf{Q}}} \mathbf{A}_{\mathbf{Q}} [i,:,:]^{\top} \operatorname{RoPE} \left(\mathbf{B}_{\mathbf{Q}} [i,:,:] \right), \quad \mathbf{Q}_{[i,:,:]} \in \mathbb{R}^{h \times d_{h}},$$

$$\mathbf{K}_{[i,:,:]} = \frac{1}{R_{\mathbf{KV}}} \mathbf{A}_{\mathbf{K}} [i,:,:]^{\top} \operatorname{RoPE} \left(\mathbf{B}_{\mathbf{K}} [i,:,:] \right), \quad \mathbf{K}_{[i,:,:]} \in \mathbb{R}^{h \times d_{h}},$$

$$\mathbf{V}_{[i,:,:]} = \frac{1}{R_{\mathbf{KV}}} \mathbf{A}_{\mathbf{V}} [i,:,:]^{\top} \mathbf{B}_{\mathbf{V}} [i,:,:], \quad \mathbf{V}_{[i,:,:]} \in \mathbb{R}^{h \times d_{h}},$$

where A_K , B_K , B_V , and A_V are cached during inference.

Remark 4. We analyze the translation equivariance of TPA by focusing on a single attention head. Let $\mathbf{q}_m \in \mathbb{R}^{d_h}$ and $\mathbf{k}_n \in \mathbb{R}^{d_h}$ denote the query and key vectors for the *i*-th head at positions m and n, respectively. In TPA, they are computed as:

$$\mathbf{q}_{m}^{(i)} = \frac{1}{R_{\mathcal{Q}}} \sum_{r_{o}=1}^{R_{\mathcal{Q}}} \alpha_{\mathcal{Q}}^{(m,r_{q},i)} \cdot \tilde{\mathbf{b}}_{\mathcal{Q}}^{(m,r_{q})}, \quad \mathbf{k}_{n}^{(i)} = \frac{1}{R_{\mathit{KV}}} \sum_{r_{kv}=1}^{R_{\mathit{KV}}} \alpha_{\mathit{K}}^{(n,r_{kv},i)} \cdot \tilde{\mathbf{b}}_{\mathit{K}}^{(n,r_{kv})},$$

where we define the shorthand notations:

$$\begin{split} &\alpha_{Q}^{(m,r_{q},i)} := \mathbf{A}_{Q}\left[m,r_{q},i\right], \quad \alpha_{K}^{(n,r_{kv},i)} := \mathbf{A}_{K}\left[n,r_{kv},i\right], \\ &\tilde{\mathbf{b}}_{Q}^{(m,r_{q})} := \operatorname{RoPE}\left(\mathbf{B}_{Q}\left[m,r_{q},:\right]\right), \quad \tilde{\mathbf{b}}_{K}^{(n,r_{kv})} := \operatorname{RoPE}\left(\mathbf{B}_{K}\left[n,r_{kv},:\right]\right). \end{split}$$

Then, the inner product between query and key becomes:

$$\begin{split} \left\langle \mathbf{q}_{m}^{(i)}, \ \mathbf{k}_{n}^{(i)} \right\rangle &= \left\langle \frac{1}{R_{\mathcal{Q}}} \sum_{r_{q}=1}^{R_{\mathcal{Q}}} \alpha_{\mathcal{Q}}^{(m,r_{q},i)} \cdot \tilde{\mathbf{b}}_{\mathcal{Q}}^{(m,r_{q})}, \ \frac{1}{R_{\mathit{KV}}} \sum_{r_{kv}=1}^{R_{\mathit{KV}}} \alpha_{\mathit{K}}^{(n,r_{kv},i)} \cdot \tilde{\mathbf{b}}_{\mathit{K}}^{(n,r_{kv})} \right\rangle \\ &= \frac{1}{R_{\mathcal{Q}} R_{\mathit{KV}}} \sum_{r_{q}=1}^{R_{\mathcal{Q}}} \sum_{r_{kv}=1}^{R_{\mathit{KV}}} \alpha_{\mathcal{Q}}^{(m,r_{q},i)} \cdot \alpha_{\mathit{K}}^{(n,r_{kv},i)} \cdot \left\langle \tilde{\mathbf{b}}_{\mathcal{Q}}^{(m,r_{q})}, \ \tilde{\mathbf{b}}_{\mathit{K}}^{(n,r_{kv})} \right\rangle. \end{split}$$

Since α_Q and α_K are position-independent and RoPE satisfies translation equivariance, we have, for any r_q and r_{kv} :

$$\left\langle \tilde{\mathbf{b}}_{Q}^{(m+t,r_{q})},\; \tilde{\mathbf{b}}_{K}^{(n+t,r_{kv})} \right\rangle = \left\langle \tilde{\mathbf{b}}_{Q}^{(m,r_{q})},\; \tilde{\mathbf{b}}_{K}^{(n,r_{kv})} \right\rangle.$$

Therefore, the query-key inner product is invariant under joint position translations:

$$\left\langle \mathbf{q}_{m+t}^{(i)},\ \mathbf{k}_{n+t}^{(i)}\right\rangle =\left\langle \mathbf{q}_{m}^{(i)},\ \mathbf{k}_{n}^{(i)}\right\rangle .$$

This confirms that TPA preserves translation equivariance.

D Llama-3 Architecture

Given hidden states $\mathbf{H} \in \mathbb{R}^{n \times d}$ for a sequence of n tokens, we first compute the attention output

$$\mathbf{O}_{\text{attn}} = \text{Attention}(\mathbf{H}) \in \mathbb{R}^{n \times (hd_h)}$$

then project back to the model dimension and add a residual:

$$\mathbf{H} \leftarrow \mathbf{H} + \mathbf{O}_{attn} \mathbf{W}_{O,attn}, \qquad \mathbf{W}_{O,attn} \in \mathbb{R}^{(hd_h) \times d}$$

Next, an MLP block (gated form) is applied:

$$\mathbf{O}_{\text{mln}} = \sigma \left(\mathbf{H} \mathbf{W}_1 \right) \odot \left(\mathbf{H} \mathbf{W}_2 \right), \qquad \mathbf{W}_1, \mathbf{W}_2 \in \mathbb{R}^{d \times d_1}.$$

followed by the output projection and residual:

$$\mathbf{H} \leftarrow \mathbf{H} + \mathbf{O}_{mlp} \mathbf{W}_{O,mlp}, \qquad \mathbf{W}_{O,mlp} \in \mathbb{R}^{d_i \times d},$$

where $\sigma(\cdot)$ is an elementwise nonlinearity function such as SiLU and \odot denotes elementwise multiplication.

E Experimental Setup

E.1 Model Hyperparameters

We adopt a Llama-3 architecture and a GPT-2 tokenizer with a 50k vocabulary. Following Zadouri et al. (2025), we initialize our MHA baseline using the GPT-3 configuration at each target model size; this baseline serves as the anchor for parameter count. For all other attention variants, we widen the MLP layers until each model's total parameters match the MHA baseline. We report architecture hyperparameters for MHA, MQA, GQA, MLA, TPA, MLA, and MLRA in Tables 5, 6, 7, 8, 9, 10, 11. For GQA we set $g=\frac{h}{4}$. For MLA, we follow the paper's setup with $d'_c=12d_h$, $d'_c=4d_h$, and $d^R_h=0.5d_h$. For GLA, we follow the paper's setup with $d'_c=8d_h$, $d'_c=4d_h$, and $d^R_h=0.5d_h$. For TPA, we follow the paper's setup with $R_Q=6$ and $R_{KV}=2$. For MLRA, we set $d'_u=2d_h$, $d'_u=d_h$, $r'=\frac{3d_h}{h}$, and $r=\frac{6d_h}{h}$.

Table 5: Model configuration for the three model sizes for MHA in our experiments.

Model Size	# Parameters	# Layers	h	d	d_h	d_i
354M	353.94M	24	16	1024	64	2736
1.3B	1311.48M	24	16	2048	128	5464
2.9B	2872.59M	24	24	3072	128	8192

Table 6: Model configuration for the three model sizes for MQA in our experiments.

Model Size	# Parameters	# Layers	h	\boldsymbol{g}	d	d_h	d_i
354M	353.94M	24	16	1	1024	64	3376
1.3B	1311.48M	24	16	1	2048	128	6744
2.9B	2872.00M	24	24	1	3072	128	10152

Table 7: Model configuration for the three model sizes for GQA in our experiments.

Model Size	# Parameters	# Layers	h	\boldsymbol{g}	d	d_h	d_i
354M	353.94M	24	16	4	1024	64	3248
1.3B	1311.48M	24	16	4	2048	128	6488
2.9B	2872.59M	24	24	6	3072	128	9728

Table 8: Model configuration for the three model sizes for MLA in our experiments.

Model Size	# Parameters	# Layers	h	d_c'	d_c	d	d_h	d_h^R	d_i
354M	353.58M	24	16	768	256	1024	64	32	2848
1.3B	1311.13M	24	16	1536	512	2048	128	64	5696
2.9B	2872.05M	24	24	1536	512	3072	128	64	9448

Table 9: Model configuration for the three model sizes for TPA in our experiments.

Model Size	# Parameters	# Layers	h	$R_{ m Q}$	$R_{ m KV}$	d	d_h	d_i
354M	354.14M	24	16	6	2	1024	64	3496
1.3B	1311.48M	24	16	6	2	2048	128	7032
2.9B	2873.18M	24	24	6	2	3072	128	10760

Table 10: Model configuration for the three model sizes for GLA in our experiments.

Model Size	# Parameters	# Layers	h	d_c'	d_c	d	d_h	d_h^R	d_i
354M	353.96M	24	16	512	256	1024	64	32	3152
1.3B	1311.51M	24	16	1024	512	2048	128	64	6296
2.9B	2872.63M	24	24	1024	512	3072	128	64	10048

Table 11: Model configuration for the two model sizes for MLRA in our experiments.

Model Size	# Parameters	# Layers	h	d_u'	d_u	r'	r	d	d_h	d_h^R	d_i	α	$\overline{\gamma}$
1.3B	1311.88M	24	16	256	128	48	24	2048	128	64	6728	2	2
2.9B	2872.01M	24	24	256	128	32	16	3072	128	64	10488	2	1

E.2 Training Setup

Training Configuration. For pretraining, we adopt GPT-3 settings. We use AdamW (Loshchilov and Hutter, 2017) with $(\beta_1,\beta_2)=(0.9,0.95),\epsilon=10^{-8}$, weight decay 0.1 , and gradient clipping at 1.0 . The learning rate is linearly warmed up for the first 2,000 steps, then annealed with cosine decay (Loshchilov and Hutter, 2016) to 10% of the peak. Peak learning rates are $3\times10^{-4},2\times10^{-4},$ and 1.6×10^{-4} for the 354M , 1.3B, and 2.9B models, respectively. We use a context length of 2,048 tokens and a global batch of 240 sequences, yielding 491,520 tokens per step (\approx 0.5M). We train for 200,000 steps, totaling 98.3B tokens. By contrast, GPT-3 uses 1.0M tokens/step for 1.3B and 2.9B experiments.

8 8			ı		
Model Size	Micro-batch Size	Batch Size	Peak Learning Rate	Total Steps	
354M	120	240	3×10^{-4}	200,000	
1.3B	40	240	2×10^{-4}	200,000	
2.9B	8	240	1.6×10^{-4}	200,000	

Table 12: Training configuration for the three model sizes in our experiments.

Initialization. For MHA, MQA, GQA, MLA, and GLA, all learnable parameters are initialized with a normal distribution $\mathcal{N}(0, \sigma = 0.02)$.

For TPA, we follow the paper's setup. The matrices \mathbf{W}_{AQ} , \mathbf{W}_{BQ} , \mathbf{W}_{AK} , \mathbf{W}_{BK} , \mathbf{W}_{AV} , \mathbf{W}_{BV} use Xavier uniform initialization (Glorot and Bengio, 2010): each entry is sampled from \mathcal{U} (-bound, bound) with bound = $\sqrt{\frac{6}{d_{\mathrm{in}}+d_{\mathrm{out}}}}$, where d_{in} and d_{out} are the input and output dimensions of the respective matrix. The output projections $\mathbf{W}_{O,\,\mathrm{attn}}$ and $\mathbf{W}_{O,\,\mathrm{mlp}}$ are zero-initialized, while all remaining parameters use $\mathcal{N}(0,\sigma=0.02)$.

In MLRA, we also zero-initialize $W_{O, attn}$ and $W_{O, mlp}$; all other parameters are initialized with $\mathcal{N}(0, \sigma = 0.02)$.

We conduct an ablation study on the baseline mechanisms comparing zero vs. $\mathcal{N}(0,\sigma=0.02)$ initialization for $\mathbf{W}_{O,attn}$ and $\mathbf{W}_{O,mlp}$. Medium-size models are trained for 100B tokens using the configuration in Appendix E.2. We evaluate perplexity on 6 datasets: FineWeb-Edu (Penedo et al., 2024), Wikipedia, C4, Common Crawl, Pile (Gao et al., 2020), and ArXiv, each evaluated using 100M tokens. Results are reported in Table 13. We find that $\mathcal{N}(0,0.02)$ performs better for MHA, MQA, GQA, MLA, and GLA, whereas zero initialization is better for TPA.

		•	` ` ` '		,	-,	~, _F
Method	Initialization	FineWeb-Edu	Wikipedia	C4	Common Crawl	Pile	ArXiv Avg
MHA	$\begin{array}{ c c } \mathcal{N}(0, \sigma = 0.02) \\ 0 \end{array}$	13.159	20.369	21.732	20.230	15.974	17.836 18.217
MHA		13.604	20.579	22.432	20.838	15.812	18.255 18.586
MQA	$\begin{array}{c c} \mathcal{N}(0, \sigma = 0.02) \\ 0 \end{array}$	13.375	20.261	22.157	20.625	15.380	17.947 18.291
MQA		13.563	21.468	22.471	21.129	16.157	18.632 18.903
GQA	$\begin{array}{ c c } \hline \mathcal{N}(0, \sigma = 0.02) \\ 0 \\ \hline \end{array}$	13.158	19.862	21.810	20.239	16.253	17.969 18.215
GQA		13.352	20.360	22.074	20.604	15.730	17.861 18.330
MLA	$\left \begin{array}{c} \mathcal{N}(0, \sigma = 0.02) \\ 0 \end{array} \right $	13.095	19.693	21.604	19.860	14.914	17.176 17.724
MLA		13.184	20.599	21.768	20.136	15.609	17.228 18.087
TPA	$\begin{array}{ c c }\hline 0\\ \mathcal{N}(0, \sigma = 0.02)\end{array}$	13.140	19.856	21.770	20.187	15.138	17.554 17.941
TPA		13.486	20.511	22.336	20.788	15.279	17.849 18.375
GLA	$\begin{array}{ c c } \mathcal{N}(0, \sigma = 0.02) \\ 0 \end{array}$	13.057	18.932	21.575	19.926	14.477	17.140 17.518
GLA		13.133	19.651	21.684	19.975	14.680	17.217 17.723

Table 13: Ablation study on initialization ($\mathcal{N}(0, \sigma = 0.02)$ vs. zero) for $\mathbf{W}_{O, \text{attn}}$ and $\mathbf{W}_{O, \text{mlp}}$.

F Related Work

KV Cache Compression. Recent works (Liu et al., 2023; Anagnostidis et al., 2023; Zhang et al., 2023b; Ge et al., 2024; Xiao et al., 2024; Kim et al., 2024; Zhang et al., 2024; Nawrot et al., 2024;

Tang et al., 2024; Liu et al., 2024b; Dong et al., 2024; Yue et al., 2024; Cai et al., 2024; Liu et al., 2024a; Hooper et al., 2024; Sun et al., 2024; Chen et al., 2024a; Jiang et al., 2024a; Li et al., 2024; Xiao et al., 2025; Sun et al., 2025a; Meng et al., 2025; Tang et al., 2025) don't introduce new attention mechanisms; instead, they compress the KV cache for the pretrained models. Some of these works (Liu et al., 2024b; Hooper et al., 2024; Yue et al., 2024) use quantization to store the KV cache in low-bit formats. Some other approaches (Zhang et al., 2023b; Xiao et al., 2024; Li et al., 2024; Xiao et al., 2025) retain important tokens and discard others to compress the KV cache.

Low-Rank Approximation. Low-rank approximation approaches (Hu et al., 2022; Malladi et al., 2023; Zhang et al., 2023a; Dettmers et al., 2023; Lialin et al., 2024; Zhu et al., 2024; Zeng and Lee, 2024; Chen et al., 2024b; Zhang, 2024; Liang and Li, 2024; Zhao et al., 2024a; Shi et al., 2024; Jiang et al., 2024b; Lin et al., 2025; Wang et al., 2025; Chang et al., 2025; Li et al., 2025) are widely used to compress representations to a low-dimensional latent, then up-project to recover full representations. These methods greatly reduce trainable parameters (Hu et al., 2022; Dettmers et al., 2023) during fine-tuning and decrease the number of parameters (Lin et al., 2025; Wang et al., 2025) for pretrained models.

System for Attention. FlashAttention (Dao et al., 2022; Dao, 2024; Shah et al., 2024) uses tiling and online softmax to minimize reads and writes between high-bandwidth memory and on-chip SRAM, shifting attention from a memory bottleneck to a compute bottleneck. FlashMLA (Jiashi Li, 2025) avoids explicit KV materialization during decoding by absorbing the key and value up-projection matrices into the query and attention output. The following attention computation is similar to MQA with shared KV states. Inspired by classical virtual memory and paging in operating systems, PagedAttention (Kwon et al., 2023) and vLLM use block-level memory management and preemptive request scheduling to reduce fragmentation and redundant duplication.