

---

# A<sup>3</sup>: an Analytical Low-Rank Approximation Framework for Attention

---

Anonymous Authors<sup>1</sup>

## Abstract

Large language models have demonstrated remarkable performance; however, their massive parameter counts make deployment highly expensive. Low-rank approximation offers a promising compression solution, yet existing approaches have two main limitations: (1) They focus on minimizing the output error of individual linear layers, without considering the architectural characteristics of Transformers, and (2) they decompose a large weight matrix into two small low-rank matrices. Consequently, these methods often fall short compared to other compression techniques like pruning and quantization, and introduce runtime overhead such as the extra GEMM kernel launches and memory operations for decomposed small matrices. To address these limitations, we propose A<sup>3</sup>, a post-training low-rank approximation framework. A<sup>3</sup> splits a Transformer layer into three functional components, namely QK, OV, and MLP, and provides analytical solutions that reduces the hidden dimension size inside each component while minimizing the component’s functional loss. This approach directly reduces model sizes, KV cache sizes, and FLOPs without introducing any runtime overheads. Under the same reduction budget in computation and memory, our low-rank approximated LLaMA 3.1-70B achieves a perplexity of 4.69 on WikiText-2, outperforming the previous SoTA’s 7.87 by 3.18.

## 1. Introduction

Large language models (LLMs) have shown exceptional performance in various applications, including language understanding, code completion, and reasoning tasks (Vaswani et al., 2017; Brown et al., 2020; Chen et al., 2021a; Wei et al., 2022). However, these models usually contain bil-

ions of parameters, resulting in high computational costs and memory requirements. Linear layers and the attention mechanism contribute significantly to the model size and computational complexity, while the KV cache produced during generation further exacerbates the memory burden.

Low-rank approximation is a promising technique that breaks down a matrix into smaller sub-matrices, directly reducing computational complexity and memory usage without the need of additional specialized hardware support. Usually a trained linear layer  $\mathbf{W} \in \mathbb{R}^{m \times n}$  is approximated by  $\widetilde{\mathbf{W}}_r = \mathbf{A}_r \mathbf{B}_r$ , where  $\mathbf{A}_r \in \mathbb{R}^{m \times r}$  and  $\mathbf{B}_r \in \mathbb{R}^{r \times n}$  are two rank- $r$  matrices with  $r \ll m, n$ . At inference time, the original GEMM operation  $\mathbf{XW}$  is replaced by two smaller GEMM operations  $\mathbf{XA}_r$  and  $(\mathbf{XA}_r)\mathbf{B}_r$ . The challenge is to construct the optimal  $\mathbf{A}_r$  and  $\mathbf{B}_r$  that maintains the end-to-end model performance. Recent studies show that minimizing the layer output error instead of the weight error gives better model performance (Zhang et al., 2024a;b; Mozaffari & Dehnavi, 2024), thus various activation-aware methods have been proposed, such as SVD-LLM (Wang et al., 2024), ASVD (Yuan et al., 2023), FWSVD (Hsu et al., 2022). However, these methods usually target general linear layers, *which ineffectively save the FLOPs and memory proportional to  $\frac{m+n}{mn}r$*  (Note that  $\frac{m+n}{mn}r < r$  for any  $m, n > 2$ ). Moreover, these methods rarely consider the architectural characteristics of Transformer, and suffer from severe performance degradation when compared to pruning and quantization.

To address these limitations, we propose A<sup>3</sup>, a new analytical framework for post-training low-rank approximation. A<sup>3</sup> splits the Transformer architecture into three functional components: query-key (QK) component, output-value (OV) component, and multi-layer perceptron (MLP) component, and minimizes the functional loss of each component. This enables better end-to-end model performance than minimizing the error of individual linear layer outputs.

We highlight the following contributions of A<sup>3</sup>:

- We propose a three-part low-rank approximation setup for multi-head attention (MHA), which formulates the problem as three separate objectives: minimizing the functional loss of (1) QK’s attention score, (2) OV’s attention output, and (3) MLP’s layer output.

---

<sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

- We derive closed-form solutions for the three objectives, which reduces the hidden dimensions shared within each component: QK head dimension, OV head dimension, and MLP intermediate size. This naturally reduces model sizes, KV cache sizes, FLOPs, and avoids runtime overheads like extra GEMM operations. Moreover,  $A^3$  *trims both FLOPs/memory and information energy proportionally to the rank  $r$* , enabling a more effective trade-off between model performance and hardware efficiency.
- We have adapted  $A^3$  for use with diverse Transformer architectures, including group query attention (GQA) and rotary position embedding (RoPE), which allows the application of  $A^3$  across a broad spectrum of models. This overcomes the limitation of existing low-rank approximation methods, which can only be applied on the vanilla MHA architecture.
- We conduct extensive experiments on various LLMs, and show that  $A^3$  outperforms SoTA low-rank methods by a significant margin. For example, our compressed LLaMA 3.1-70B achieves a perplexity of 4.69 on WikiText-2, outperforming SoTA method’s 7.87 by 3.18. We also demonstrate further applications of  $A^3$ , such as its effectiveness in KV cache reduction, combination with quantization, and mixed-rank assignments for better performance.

## 2. Related Work

**Low-rank approximation for compressing Transformers** The linear layer takes a simple form but contributes most parameters to Transformer. For compressing large Transformer models like LLMs, low-rank approximation has been widely studied (Chen et al., 2021b; Saha et al., 2024). Usually a trained linear layer  $\mathbf{W} \in \mathbb{R}^{m \times n}$  is approximated by  $\mathbf{W} \approx \widetilde{\mathbf{W}}_r = \mathbf{A}_r \mathbf{B}_r$ , where  $\mathbf{A}_r \in \mathbb{R}^{m \times r}$  and  $\mathbf{B}_r \in \mathbb{R}^{r \times n}$  are two rank- $r$  matrices that effectively reduce the number of parameters and FLOPs with a small enough  $r$ . The problem is how to find the optimized  $\mathbf{A}_r$  and  $\mathbf{B}_r$  that maintains the model performance. If the objective is to minimize the Frobenius norm of the weight error,

$$\operatorname{argmin}_{\widetilde{\mathbf{W}}_r} \|\mathbf{W} - \widetilde{\mathbf{W}}_r\|_F^2 \quad \text{s.t.} \quad \operatorname{rank}(\widetilde{\mathbf{W}}_r) = r, \quad (1)$$

according to Eckart-Young theorem (Eckart & Young, 1936a), the optimal solution is to perform truncated singular value decomposition (SVD) on the weight matrix  $\mathbf{W}$ .

$$\mathbf{W} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T, \quad \widetilde{\mathbf{W}}_r = \operatorname{SVD}_r(\mathbf{W}) = \mathbf{U}_{:, :r} \mathbf{\Sigma}_{:k, :r} \mathbf{V}_{:r, :}^T, \quad (2)$$

where  $\mathbf{U} \in \mathbb{R}^{m \times k}$  and  $\mathbf{V} \in \mathbb{R}^{k \times n}$  are the left and right singular vectors and  $\mathbf{\Sigma} \in \mathbb{R}^{k \times k}$  is the diagonal matrix of singular values. Recent studies show that minimizing the

layer output error instead of the weight error gives better end-to-end model performance (Zhang et al., 2024a;b),

$$\operatorname{argmin}_{\widetilde{\mathbf{W}}_r} \|\mathbf{X} \mathbf{W} - \mathbf{X} \widetilde{\mathbf{W}}_r\|_2^2 \quad \text{s.t.} \quad \operatorname{rank}(\widetilde{\mathbf{W}}_r) = r, \quad (3)$$

where  $\mathbf{X} \in \mathbb{R}^{l \times d}$  denotes the activation of  $\mathbf{W}$ . The objective above minimizes the expected  $l_2$ -norm of layer output error. Recent studies find the optimal solution is:

$$\widetilde{\mathbf{W}}_r = (\mathbf{R}_{\mathbf{X}\mathbf{X}}^{\frac{1}{2}})^{-1} \operatorname{SVD}_r(\mathbf{R}_{\mathbf{X}\mathbf{X}}^{\frac{1}{2}} \mathbf{W}), \quad (4)$$

assuming  $\mathbf{R}_{\mathbf{X}\mathbf{X}}$  is positive definite, where  $\mathbf{R}_{\mathbf{X}\mathbf{X}} = \frac{1}{l} \mathbf{X}^T \mathbf{X}$  is the autocorrelation matrix with respect to  $\mathbf{X}$ , and  $\mathbf{R}_{\mathbf{X}\mathbf{X}}^{\frac{1}{2}}$  denotes the unique symmetric square root of  $\mathbf{R}_{\mathbf{X}\mathbf{X}}$ .

Interestingly, there are several works proposing or leveraging the solution in Equation (4) for various applications. DRONE (Chen et al., 2021b) and SVD-LLM (Wang et al., 2024) directly apply the solution to approximate all linear layers in Transformers. QERA (Zhang et al., 2024b) uses it to build high-precision low-rank terms to compensate for output quantization error, while CALDERA (Saha et al., 2024) further proposes iterative methods to quantize the low-rank terms, achieving performant sub-2.5-bit post-training quantization. Palu (Chang et al., 2024) reduces KV cache size by decomposing key and value weight matrices with the solution and caching smaller intermediate activations instead of original keys and values. ESPACE (Sakr & Khailany, 2024), a training based method, extends the base  $L_2$  objective by incorporating outlier, gradient information, and per-layer selection of the calibration strategy to maximize fine-tuning accuracy recovery. SLiM (Mozafari & Dehnavi, 2024) and Oats (Zhang & Pappan, 2024) incorporate sparsity with low-rank to achieve an overall compression gain. Beyond low-rank methods, quantization techniques such as Quarot (Ashkboos et al., 2024), AWQ (Lin et al., 2024), GPTQ (Frantar et al., 2022) and HQQ (Badri & Shaji, 2023) are often compatible with low-rank approaches for multiplicative gains such as Palu, SLiM and QERA. We show that  $A^3$  can also be effectively combined with quantization and creates a continuous spectrum of compression levels between discrete quantization points, yielding a way better Pareto frontier at extreme compression.

However, these works target general linear layers and minimize the linear layer output error without considering architectural characteristics. In this work, we step forward to the optimization for functional components. We propose analytical low-rank approximation methods of compressing the QK, OV, and MLP components that minimize the functional errors of attention scores, attention outputs, and MLP outputs, respectively, in a training-free manner.

### 3. The $\mathbf{A}^3$ Framework

In this section, for each component (QK, OV, MLP), we define the problem (optimization objectives), clarify the assumptions if any, and propose our analytical solutions. The proof for each lemma and theorem is provided in the Appendix. We also provide notation tables in Tables 4 and 5 in the appendix for the ease of reading.

#### 3.1. $\mathbf{A}^3$ -QK

In the QK component, each head computes its pre-softmax attention scores between queries and keys:

$$\mathbf{A}_i = \mathbf{Q}_i \mathbf{K}_i^T = \mathbf{X} \mathbf{W}_{q,i} \mathbf{W}_{k,i}^T \mathbf{X}^T = \mathbf{X} \mathbf{W}_{qk,i} \mathbf{X}^T, \quad (5)$$

where  $\mathbf{X}_q \in \mathbb{R}^{l_q \times d_m}$ ,  $\mathbf{X}_{kv} \in \mathbb{R}^{l_{kv} \times d_m}$  are the input of query layer and key/value layer respectively, and  $\mathbf{W}_{qk,i} := \mathbf{W}_{q,i} \mathbf{W}_{k,i}^T$  denotes the fused weight matrix of the  $i$ -th head. We seek for the low-rank approximation of  $\mathbf{W}_{qk,i}$  that minimizes the error of pre-softmax attention scores.

**Problem 1** (Minimization of the pre-softmax attention score error). Given a pretrained Transformer layer, for the  $i$ -th head of QK component  $\mathbf{A}_i = \mathbf{X}_q \mathbf{W}_{qk,i} \mathbf{X}_{kv}^T$  and its rank- $r$  approximated form  $\widetilde{\mathbf{A}}_i = \mathbf{X}_q \widetilde{\mathbf{W}}_{qk,i} \mathbf{X}_{kv}^T$ , approximating the head by minimizing the error between  $\mathbf{A}_i$  and  $\widetilde{\mathbf{A}}_i$  on a calibration set:

$$\begin{aligned} \operatorname{argmin}_{\widetilde{\mathbf{W}}_{qk,i}} & \left\| \mathbf{X}_q (\mathbf{W}_{qk,i} - \widetilde{\mathbf{W}}_{qk,i}) \mathbf{X}_{kv}^T \right\|_F^2 \\ \text{s.t.} & \quad \operatorname{rank}(\widetilde{\mathbf{W}}_{qk,i}) = r. \end{aligned} \quad (6)$$

where  $\mathbf{X}_q \in \mathbb{R}^{l_q \times d_m}$  and  $\mathbf{X}_{kv} \in \mathbb{R}^{l_{kv} \times d_m}$  denote the inputs for query and key/value projections, which can also be considered as the calibration set when  $l_q$  and  $l_{kv}$  are sufficiently large.

**Lemma 3.1** (Equivalent form of Problem 1). *The objective in Problem 1 is equivalent to:*

$$\operatorname{argmin}_{\widetilde{\mathbf{W}}_{qk,i}} \left\| \mathbf{R}_{\mathbf{X}_q \mathbf{X}_q}^{\frac{1}{2}} (\mathbf{W}_{qk,i} - \widetilde{\mathbf{W}}_{qk,i}) \mathbf{R}_{\mathbf{X}_{kv} \mathbf{X}_{kv}}^{\frac{1}{2}} \right\|_F^2, \quad (7)$$

where  $\mathbf{R}_{\mathbf{X}_q \mathbf{X}_q} := \frac{1}{l_q} \mathbf{X}_q^T \mathbf{X}_q$  and  $\mathbf{R}_{\mathbf{X}_{kv} \mathbf{X}_{kv}} := \frac{1}{l_{kv}} \mathbf{X}_{kv}^T \mathbf{X}_{kv}$  are the autocorrelation matrices of the query and key/value calibration activations respectively, and  $\mathbf{R}_{\mathbf{X}_q \mathbf{X}_q}^{\frac{1}{2}}$ ,  $\mathbf{R}_{\mathbf{X}_{kv} \mathbf{X}_{kv}}^{\frac{1}{2}}$  denote the corresponding unique symmetric matrix square roots.

The complete derivation of Lemma 3.1 is given in Section B.1.1.

**Theorem 3.2** ( $\mathbf{A}^3$ -QK for MHA-NoPE). *The optimal solution to Problem 1 is*

$$\widetilde{\mathbf{W}}_{qk,i} = \left( \mathbf{R}_{\mathbf{X}_q \mathbf{X}_q}^{1/2} \right)^{-1} \operatorname{SVD}_r \left( \mathbf{R}_{\mathbf{X}_q \mathbf{X}_q}^{1/2} \mathbf{W}_{qk,i} \mathbf{R}_{\mathbf{X}_{kv} \mathbf{X}_{kv}}^{1/2} \right) \left( \mathbf{R}_{\mathbf{X}_{kv} \mathbf{X}_{kv}}^{1/2} \right)^{-1}. \quad (8)$$

where  $\operatorname{SVD}_r(\cdot)$  denotes the truncated SVD operator.

The proof of Theorem 3.2 is given in Appendix B.1.2. In practice, we apply Theorem 3.2 to all pairs of QK heads ( $i = 1, \dots, h_q$ ), and assign

$$\widetilde{\mathbf{W}}_{q,i} := \left( \mathbf{R}_{\mathbf{X}_q \mathbf{X}_q}^{\frac{1}{2}} \right)^{-1} \mathbf{U}_{:,k}, \quad \widetilde{\mathbf{W}}_{k,i}^T := \boldsymbol{\Sigma}_{:,k} \mathbf{V}_{:,k}^T \left( \mathbf{R}_{\mathbf{X}_{kv} \mathbf{X}_{kv}}^{\frac{1}{2}} \right)^{-1},$$

where  $\mathbf{U}_{:,k}$ ,  $\boldsymbol{\Sigma}_{:,k}$ , and  $\mathbf{V}_{:,k}^T$  are the truncated SVD components given by Theorem 3.2. This gives approximated query and key weights with a new smaller head dimension  $r < d_{qk}$ . Note that Theorem 3.2 is performed on each head separately, but the low-rank head weights can still be concatenated together and implemented as a single linear layer at inference time.

#### 3.2. $\mathbf{A}^3$ -OV

Expand the summation over all OV head outputs in ???. The matrix form of the attention layer output  $\mathbf{O} \in \mathbb{R}^{l_q \times d_m}$  can be expressed as

$$\mathbf{O} = \sum_{i=1}^{h_q} \mathbf{O}_i = \sum_{i=1}^{h_q} \mathbf{A}'_i \mathbf{X}_{kv} \mathbf{W}_{v,i} \mathbf{W}_{o,i} = \sum_{i=1}^{h_q} \mathbf{P}_i \mathbf{W}_{vo,i}, \quad (9)$$

where  $\mathbf{P}_i := \mathbf{A}'_i \mathbf{X}_{kv} \in \mathbb{R}^{l_q \times d_m}$  is the product between post-softmax attention score and the input matrix of key/value layer, and  $\mathbf{W}_{vo,i} := \mathbf{W}_{v,i} \mathbf{W}_{o,i} \in \mathbb{R}^{d_m \times d_m}$  denotes the fused weight matrix of the  $i$ -th head. Now each term  $\mathbf{O}_i$  takes the form of a linear layer  $\mathbf{O}_i = \mathbf{P}_i \mathbf{W}_{vo,i}$ . If  $\widetilde{\mathbf{O}}_i = \mathbf{P}_i \widetilde{\mathbf{W}}_{vo,i}$  denotes the approximated  $\mathbf{O}_i$ , the upper bound of the attention output error can be derived as follows:

$$\|\mathbf{O} - \widetilde{\mathbf{O}}\|_2^2 = \left\| \sum_{i=1}^{h_q} (\mathbf{O}_i - \widetilde{\mathbf{O}}_i) \right\|_2^2 \leq \sum_{i=1}^{h_q} \|\mathbf{O}_i - \widetilde{\mathbf{O}}_i\|_2^2. \quad (10)$$

Though  $\|\mathbf{O} - \widetilde{\mathbf{O}}\|_2^2$  can be directly minimized via matrix stacking and truncated SVD, its closed-form solution incurs a higher computational cost, as elaborated in Section B.2.3. Here, we relax the objective and treat minimizing each error term  $\|\mathbf{O}_i - \widetilde{\mathbf{O}}_i\|_2^2$  as an independent problem. Thus the optimal solution to  $\widetilde{\mathbf{W}}_{vo,i}$  is already given by Equation (4).

**Problem 2** (Minimization of per-head attention output error). Given a pretrained Transformer layer, for the  $i$ -th head of OV component  $\mathbf{O}_i = \mathbf{P}_i \mathbf{W}_{vo,i}$  and its approximated form  $\widetilde{\mathbf{O}}_i = \mathbf{P}_i \widetilde{\mathbf{W}}_{vo,i}$ , approximating the head by minimizing the head output error is to minimize the following loss:

$$\operatorname{argmin}_{\widetilde{\mathbf{W}}_{vo,i}} \left\| \mathbf{P}_i (\mathbf{W}_{vo,i} - \widetilde{\mathbf{W}}_{vo,i}) \right\|_F^2 \quad \text{s.t.} \quad \operatorname{rank}(\widetilde{\mathbf{W}}_{vo,i}) = r. \quad (11)$$

**Theorem 3.3** ( $\mathbf{A}^3$ -OV for MHA-NoPE). *The optimal solution to Problem 2 is*

$$\widetilde{\mathbf{W}}_{vo,i} = \left( \mathbf{R}_{\mathbf{P}_i \mathbf{P}_i}^{\frac{1}{2}} \right)^{-1} \operatorname{SVD}_r \left( \mathbf{R}_{\mathbf{P}_i \mathbf{P}_i}^{\frac{1}{2}} \mathbf{W}_{vo,i} \right), \quad (12)$$

where  $\mathbf{R}_{\mathbb{X}_{p_i} \mathbb{X}_{p_i}} = \frac{1}{l_q} \mathbf{P}_i^T \mathbf{P}_i$  is the autocorrelation matrix of  $\mathbf{P}_i$  and  $\mathbf{R}_{\mathbb{X}_{p_i} \mathbb{X}_{p_i}}^{\frac{1}{2}}$  denotes its unique symmetric matrix square root.

Similar to QK component, in practice we assign

$$\widetilde{\mathbf{W}}_{v,i} := \left( \mathbf{R}_{\mathbb{X}_{p_i} \mathbb{X}_{p_i}}^{\frac{1}{2}} \right)^{-1} \mathbf{U}_{:,i}, \quad \widetilde{\mathbf{W}}_{vo,i} := \Sigma_{:,k,:k} \mathbf{V}_{:,k}^T \left( \mathbf{R}_{\mathbb{X}_{p_i} \mathbb{X}_{p_i}}^{\frac{1}{2}} \right)^{-1},$$

to get the approximated value and output weights of head- $i$  with a smaller head dimension  $r < d_{vo}$ . Extension on GQA and RoPE is detailed in Appendix C.

### 3.3. A<sup>3</sup>-MLP

The non-linear activation function in MLP component prohibits us from directly applying SVD. Instead, we first derive an objective for minimizing the MLP output error, and uses CUR decomposition (Mahoney & Drineas, 2009) to find the low-rank form of MLP weights.

**Problem 3** (Minimization of MLP output error). Given a pretrained down projection layer  $\mathbf{Y}_{\text{mlp}} = \mathbf{X}_d \mathbf{W}_d$  in MLP and its approximated low-rank form  $\widetilde{\mathbf{Y}}_{\text{mlp}} = \widetilde{\mathbf{X}}_d \widetilde{\mathbf{W}}_d = \mathbf{X}_d \mathbf{U} \mathbf{W}_d$ , minimizing the MLP output error is to minimize the following loss:

$$\begin{aligned} \arg \min_{\mathbf{U} = \text{diag}(u_1, \dots, u_{d_{\text{inter}}})} \|\mathbf{X}_d \mathbf{U} \mathbf{W}_d - \mathbf{X}_d \mathbf{W}_d\|_F^2 \\ \text{s.t. } \text{rank}(\mathbf{U}) = r. \end{aligned} \quad (13)$$

where  $\mathbf{X}_d \in \mathbb{R}^{l_{\text{down}} \times d_{\text{inter}}}$  is the matrix of intermediate activation vectors, and  $\mathbf{U} \in \mathbb{R}^{d_{\text{inter}} \times d_{\text{inter}}}$  is a diagonal matrix determining which  $r$  columns of  $\mathbf{W}_d$  to keep.

**Lemma 3.4** (Equivalent form of Problem 3). *The objective in Problem 3 is equivalent to the following error on the calibration dataset:*

$$\begin{aligned} \arg \min_{\mathbf{U} = \text{diag}(u_1, u_2, \dots, u_{d_{\text{inter}}})} \left\| \mathbf{R}_{\mathbb{X}_d \mathbb{X}_d}^{\frac{1}{2}} \mathbf{U} \mathbf{W}_d - \mathbf{R}_{\mathbb{X}_d \mathbb{X}_d}^{\frac{1}{2}} \mathbf{W}_d \right\|_F^2 \\ \text{s.t. } \text{rank}(\mathbf{U}) = r. \end{aligned} \quad (14)$$

where  $\mathbf{R}_{\mathbb{X}_d \mathbb{X}_d} := \frac{1}{l_{\text{down}}} \mathbf{X}_d^T \mathbf{X}_d$  is the autocorrelation matrix of the calibration set  $\mathbf{X}_d$ .

The derivation of Lemma 3.4 is in Section B.3.1. This CUR approximation is a well-studied NP-hard problem, and various CUR methods have been proposed (Boutsidis & Woodruff, 2014; Drineas et al., 2006). We thus pick a simple but effective solution from (Drineas et al., 2006) and name this approach A<sup>3</sup>-MLP.

Following (Drineas et al., 2006), we build  $\mathbf{U}$  by sorting the F-norm of the outer product between the columns of  $\mathbf{R}_{\mathbb{X}_d \mathbb{X}_d}^{\frac{1}{2}}$  and the rows of  $\mathbf{W}_d$ :

$$\lambda_i = \|\mathbf{r}_i^T \mathbf{w}_i\|_F^2 = \|\mathbf{r}_i\|_2^2 \cdot \|\mathbf{w}_i\|_2^2, \quad (15)$$

where  $\mathbf{r}_i$  is the  $i$ -th column of  $\mathbf{R}_{\mathbb{X}_d \mathbb{X}_d}^{\frac{1}{2}}$  and  $\mathbf{w}_i$  is the  $i$ -th row of  $\mathbf{W}_d$ . Then  $\mathbf{U}$  is built by selecting the indexes that gives the top- $r$   $\lambda_i$ :

$$\begin{aligned} \mathbf{U} &= \text{diag}(u_1, u_2, \dots, u_{d_{\text{inter}}}), \\ u_i &= \frac{1}{r \lambda_i} \text{ if } i \in \text{top-}r(\lambda_i) \text{ else } 0. \end{aligned} \quad (16)$$

In practice, we compute  $\lambda_i$  for all  $i = 1, \dots, d_{\text{inter}}$ . Then we select the  $r$  rows of  $\mathbf{W}_d$  that have  $r$  largest non-zero  $\lambda_i$  to form  $\widetilde{\mathbf{W}}_d \in \mathbb{R}^{r \times d_m}$ . Accordingly,  $\widetilde{\mathbf{W}}_u, \widetilde{\mathbf{W}}_g \in \mathbb{R}^{d_m \times r}$  are formed by selecting the corresponding columns of  $\mathbf{W}_u$  and  $\mathbf{W}_g$  respectively.

Note that most related works, *e.g.*, the ones introduced in Section 2, target general linear layers and replace a weight matrix with two low-rank matrices  $\mathbf{X} \mathbf{W} \approx (\mathbf{X} \mathbf{A}_r) \mathbf{B}_r$ , which introduces one more GEMM operation per linear layer at inference time. In contrast, *all of our three solutions only reduce the hidden dimensions of the components ( $h_q$ ,  $d_{vo}$ , and  $d_{\text{inter}}$ ), resulting in the same number of GEMM operations with smaller problem sizes. This naturally enables reduced model sizes, saved the FLOPs of both linear layers and attention, compressed KV cache, without introducing any runtime overhead.*

## 4. Experiments

**Baselines** We compare A<sup>3</sup> against a range of baselines, including vanilla low-rank approximation using SVD and weight-magnitude-based column/row pruning, as well as SoTA approaches, including FWSVD (Hsu et al., 2022), ASVD (Yuan et al., 2023), SVD-LLM (Wang et al., 2024), SVD-LLM v2 (Wang et al., 2024), Palu (Chang et al., 2024), Wanda (Sun et al., 2023), and CLOVER (Meng et al., 2025). However, only several baselines support approximating all the three components (QK, OV, MLP), including SVD, FWSVD, ASVD, SVD-LLM, and SVD-LLM-v2. We conduct a comprehensive comparison against these methods. For other baselines, we align the components to approximate and present results in the ablation study. Unless otherwise specified, the compression ratio is defined in terms of parameter count to ensure consistency across all baselines.

**Models and benchmarks** Our evaluation covers vanilla Transformer and its variants, including MHA without RoPE, denoted as MHA-NoPE, (MPT (Team et al., 2023)), MHA-RoPE (LLaMA 1&2 (Touvron et al., 2023a;b)), and GQA-RoPE (LLaMA 3.1 (Grattafiori et al., 2024), Phi 3 (Abdin et al., 2024), Mistral 3 (Liu et al., 2026)). We evaluate on pretraining tasks (WikiText-2 (Merity et al., 2016), C4 (Raffel et al., 2020), and SlimPajama (Shen et al., 2023)) using SVD-LLM’s perplexity evaluation code snippet, and downstream tasks (ARC-Challenge, BoolQ, Winogrande, GSM8K (strict match), and MMLU) us-

ing lm-eval-harness (Gao et al., 2024). All experiments are post-training low-rank approximation *without fine-tuning*. We use 128 random 2048-token sequences from SlimPajama for all evaluations, except in Figure 1, where we calibrate on WikiText2 to match SVD-LLM’s setup (see Appendix D for details).

#### 4.1. Main Results

This section presents the main evaluation results where we compare  $A^3$  against all the baselines that can be applied to all of the three main components (QK, OV, MLP) in Transformer. We first simply evaluate on LLaMA-7B and eliminate less promising baselines, then conduct a comprehensive evaluation on more models and tasks. Lastly, we present profiling results to highlight  $A^3$ ’s improvement on hardware efficiency.

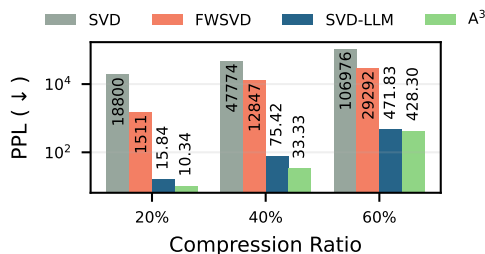


Figure 1. LLaMA-7b PPL on C4, compared to SVD, FWSVD and SVD-LLM.

**Preliminary experiments** We first apply plain SVD, FWSVD, SVD-LLM, and  $A^3$  on LLaMA-7B with 20%, 40%, and 60% compression ratios and compare the perplexity (PPL ↓) results on WikiText-2 to find the most promising baselines. As shown in Figure 1, SVD-LLM and  $A^3$  achieve perplexities smaller than others by two to three orders of magnitude. We thus conduct further experiments on SVD-LLM and  $A^3$ .

**Pretraining tasks and downstream tasks** In Table 1, we include more models to compare  $A^3$  against SVD-LLM on WikiText-2, C4, and SlimPajama, covering MHA-RoPE and GQA-RoPE architectures.  $A^3$  outperforms SVD-LLM by a large margin most of the time. Remarkably,  $A^3$  achieves a perplexity of 4.69 on WikiText-2 with LLaMA 3.1-70B, which is 3.18 lower than SVD-LLM’s 7.87 (a perplexity reduction of 58.6%) at 10% compression ratio. We present the downstream task results in Table 2, where  $A^3$  consistently outperforms SVD-LLM in terms of average accuracy (↑) across all five tasks. We observe that the advantage of  $A^3$  is more pronounced when the compression ratio is small (10%). We attribute this to the adaptation of  $A^3$  for RoPE and GQA, which we will discuss in Section E. Results on Phi 3 and Minstral 3 are presented in Appendix I.

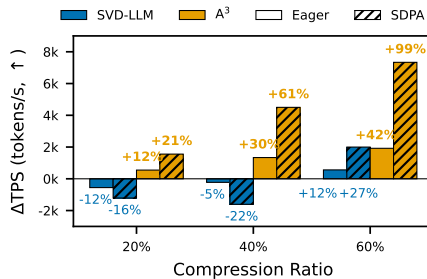


Figure 2. Performance comparisons in Tokens per Second (TPS) of  $A^3$  and SVD-LLM (LLaMA-2-13b, A100 40GB, batch size=2, sequence length=2048, attention backend=Eager/SDPA).

**Higher inference throughput** The low-rank approximation methods that target general linear layers only save the FLOPs of GEMM in linear layers, but induce runtime overhead like extra GEMM kernel launches and read/write for small matrices. In contrast,  $A^3$  saves the GEMM FLOPs in both linear layers and attention, without inducing these overheads. We profile prefilling throughput measured in TPS (tokens/sec) of LLaMA-2-13B on an A100 40GB, and visualize the speedup in Figure 2. SVD-LLM only has speedup for aggressive compression, while  $A^3$  always achieves a speedup, higher than SVD-LLM. More runtime analysis can be found in Appendix G.

**Stronger KV Cache Compression** Table 3 compares  $A^3$  with Clover and Palu across the full range of compression ratios on MPT-7B and MPT-30B. Although all three methods reduce parameters and KV-cache size by the **same** amount at a given compression ratio, their perplexity trends diverge significantly. Clover’s performance degrades sharply, even at only 20% compression, its perplexity rises above 40 on the MPT-7B model. Palu is closer to  $A^3$ , but it still underperforms by a large margin because its objective only reduces loss in the K and V projection outputs. In contrast,  $A^3$  consistently achieves the lowest perplexity across all datasets and model sizes. This performance gap widens at higher compression levels and with larger model sizes: on MPT-30B,  $A^3$  is the **only** method that stays below 50 perplexity at 80% compression.

**Ablation Study and other Baselines** We conduct ablation studies to evaluate  $A^3$ ’s impact on individual components (QK, OV, MLP). We also include baselines that can be applied to the target components to show  $A^3$ ’s advantage. For attention without RoPE, we compare  $A^3$  against CLOVER, Palu, and two simplified variants of  $A^3$ -QK. For RoPE, we compare against structure pruning methods adapted to this setting. Details are provided in Appendix E.

## 5. Discussion

In Appendix F, we showcase various applications of  $A^3$ , highlighting its compatibility with quantization, lora fine-

Table 1. A comparison of perplexity ( $\downarrow$ ) on WikiText2, C4, and SlimPajama.

Model	Method	10%			20%		
		WikiText-2	C4	SlimPajama	WikiText-2	C4	SlimPajama
LLaMA-2-7B (MHA-RoPE)	SVD-LLM	8.78 (+3.30)	11.73 (+4.14)	9.49 (+3.35)	11.58 (+6.1)	14.91 (+7.32)	11.93 (+5.79)
	A <sup>3</sup>	<b>5.96 (+0.48)</b>	<b>8.34 (+0.74)</b>	<b>6.68 (+0.54)</b>	<b>7.22 (+1.73)</b>	<b>9.91 (+2.31)</b>	<b>7.91 (+1.77)</b>
LLaMA-2-13B (MHA-RoPE)	SVD-LLM	7.09 (+2.19)	9.98 (+2.92)	7.95 (+2.26)	9.03 (+4.13)	12.35 (+5.29)	9.75 (+4.06)
	A <sup>3</sup>	<b>5.32 (+0.42)</b>	<b>7.65 (+0.59)</b>	<b>7.65 (+1.97)</b>	<b>6.24 (+1.34)</b>	<b>8.99 (+1.92)</b>	<b>7.15 (+1.47)</b>
LLaMA-3.1-8B (GQA-RoPE)	SVD-LLM	19.12 (+12.86)	19.37 (+9.33)	15.14 (+7.57)	42.28 (+36.02)	33.6 (+23.56)	27.44 (+19.86)
	A <sup>3</sup>	<b>7.93 (+1.67)</b>	<b>12.56 (+2.52)</b>	<b>9.52 (+1.94)</b>	<b>11.36 (+5.1)</b>	<b>17.87 (+10.29)</b>	<b>13.58 (+3.54)</b>
LLaMA-3.1-70B (GQA-RoPE)	SVD-LLM	7.87 (+5.07)	11.3 (+3.76)	8.43 (+2.94)	9.75 (+6.95)	<b>13.77 (+6.23)</b>	10.44 (+4.95)
	A <sup>3</sup>	<b>4.69 (+1.90)</b>	<b>8.83 (+1.31)</b>	<b>6.59 (+1.10)</b>	<b>8.32 (+5.52)</b>	13.94 (+6.40)	<b>10.02 (+4.53)</b>

Table 2. A comparison of downstream task accuracy ( $\uparrow$ ).

Model	CRatio	Method	ARC-c	BoolQ	Winogrande	GSM8k	MMLU	Avg.
LLaMA-2-7b (MHA-RoPE)	-	Original	0.4829	0.7777	0.7498	0.1387	0.4582	0.5158
	10%	SVD-LLM	0.3882	0.6749	0.6803	0.0129	0.3477	0.4166
	A <sup>3</sup>	<b>0.4761</b>	<b>0.7330</b>	<b>0.7435</b>	<b>0.1130</b>	<b>0.4398</b>	<b>0.4960</b>	
LLaMA-2-13b (MHA-RoPE)	20%	SVD-LLM	0.3139	0.6602	0.6464	0.0045	0.3119	0.3837
	A <sup>3</sup>	<b>0.4369</b>	<b>0.7174</b>	<b>0.7072</b>	<b>0.0751</b>	<b>0.3979</b>	<b>0.4621</b>	
	-	Original	0.5538	0.8086	0.7711	0.2343	0.5513	0.5774
LLaMA-3.1-8B (GQA-RoPE)	10%	SVD-LLM	0.4206	<b>0.8061</b>	0.7308	0.0902	0.4772	0.5000
	A <sup>3</sup>	<b>0.5213</b>	0.7865	<b>0.7743</b>	<b>0.1971</b>	<b>0.5324</b>	<b>0.5560</b>	
	20%	SVD-LLM	0.3472	<b>0.7877</b>	0.6898	0.0379	0.4318	0.4546
A <sup>3</sup>	<b>0.4727</b>	0.7654	<b>0.7364</b>	<b>0.1645</b>	<b>0.4804</b>	<b>0.5180</b>		
LLaMA-3.1-70B (GQA-RoPE)	-	Original	0.5401	0.8190	0.7822	0.4920	0.6535	0.6484
	10%	SVD-LLM	0.3575	0.7458	<b>0.7111</b>	0.0447	0.4708	0.4603
	A <sup>3</sup>	<b>0.4565</b>	<b>0.7884</b>	0.7072	<b>0.2388</b>	<b>0.5922</b>	<b>0.5500</b>	
LLaMA-3.1-8B (GQA-RoPE)	20%	SVD-LLM	0.2534	<b>0.6948</b>	<b>0.6440</b>	0.0113	0.3604	0.3880
	A <sup>3</sup>	<b>0.3345</b>	0.6823	0.6417	<b>0.0705</b>	<b>0.4649</b>	<b>0.4336</b>	
	-	Original	0.6536	0.8538	0.8445	0.8036	0.7864	0.7768
LLaMA-3.1-70B (GQA-RoPE)	10%	SVD-LLM	0.5742	0.8401	0.8051	0.5087	0.7181	0.6797
	A <sup>3</sup>	<b>0.6323</b>	<b>0.8532</b>	<b>0.8335</b>	<b>0.7453</b>	<b>0.7470</b>	<b>0.7508</b>	
	20%	SVD-LLM	<b>0.4957</b>	<b>0.8226</b>	<b>0.7727</b>	0.3040	<b>0.6620</b>	0.6025
A <sup>3</sup>	0.4667	0.8144	0.6875	<b>0.4951</b>	0.6145	<b>0.6071</b>		

Table 3. A comparison of perplexity ( $\downarrow$ ) on WikiText2, C4, and SlimPajama. CRatio indicates compression ratio on both KV-Cache and parameter count.

Model	CRatio	SlimPajama			C4			Wikikitext-2		
		Clover	Palu	A <sup>3</sup>	Clover	Palu	A <sup>3</sup>	Clover	Palu	A <sup>3</sup>
MPT-7B	20%	48.11	9.67	<b>8.88</b>	53.29	11.74	<b>10.77</b>	77.78	8.73	<b>8.05</b>
	40%	383	11.51	<b>9.90</b>	408	14.18	<b>12.20</b>	795	10.60	<b>9.19</b>
	60%	5397	25.73	<b>15.34</b>	4919	32.26	<b>18.71</b>	7895	25.09	<b>15.58</b>
	80%	15467	5270	<b>388</b>	11661	3210	<b>373</b>	14434	13714	<b>849</b>
MPT-30B	20%	11.52	7.91	<b>7.71</b>	14.53	9.87	<b>9.59</b>	13.07	7.04	<b>6.73</b>
	40%	18.00	8.99	<b>8.33</b>	22.43	11.30	<b>10.44</b>	23.47	8.40	<b>7.40</b>
	60%	54.97	15.59	<b>11.52</b>	70.65	18.91	<b>14.22</b>	95.45	18.88	<b>11.28</b>
	80%	779	211	<b>37.09</b>	732	253	<b>42.85</b>	1524	339	<b>46.72</b>

tuning and extensibility to mixed-rank allocation for additional performance gains. We find A<sup>3</sup> is orthogonal to weight-only quantization methods like HQQ (Badri & Shaji, 2023) and a simple mixed-rank A<sup>3</sup> outperforms ASVD and SVD-LLM v2. When combined with quantization, it yields a better Pareto frontier than using quantization alone at sub-4-bit setting. We also discuss the limitations of A<sup>3</sup> in

Appendix F, mainly caused by the sub-optimality of CUR decomposition, a compromise to RoPE. Additionally, we provide empirical diagnostics that link the reductions in individual local objectives for QK and OV to the overall end-to-end perplexity. Finally, we explore the impact of calibration set selection on overall performance, showing that a mixture of calibration datasets boosts accuracies on downstream tasks like Winogrande.

## 6. Conclusion

We propose A<sup>3</sup>, an analytical framework that decomposes the transformer into its core components, QK, OV, MLP, and compresses them by minimizing their respective errors. This method reduces model size, KV cache, and FLOPs without runtime overhead, while achieving SoTA performance.

## Impact Statement

This paper proposes an analytical framework that decomposes the transformer into its core components, QK, OV, MLP, and compresses them by minimizing their respective errors. This method reduces model size, KV cache, and FLOPs without runtime overhead. It lowers both training and inference costs, decreases computational demands, reduces power consumption, and minimizes carbon emissions.

## References

- Abdin, M., Aneja, J., Awadalla, H., Awadallah, A., Awan, A. A., et al. Phi-3 technical report: A highly capable language model locally on your phone, 2024. URL <https://arxiv.org/abs/2404.14219>.
- Ashkboos, S., Mohtashami, A., Croci, M. L., Li, B., Cameron, P., Jaggi, M., Alistarh, D., Hoefler, T., and Hensman, J. Quarot: Outlier-free 4-bit inference in rotated llms. *Advances in Neural Information Processing Systems*, 37:100213–100240, 2024.
- Badri, H. and Shaji, A. Half-quadratic quantization of large machine learning models. *Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob*, 2023.
- Barbero, F., Vitvitskyi, A., Perivolaropoulos, C., Pascanu, R., and Veličković, P. Round and round we go! what makes rotary positional encodings useful? *arXiv preprint arXiv:2410.06205*, 2024.
- Boutsidis, C. and Woodruff, D. P. Optimal cur matrix decompositions. In *Proceedings of the forty-sixth annual ACM symposium on Theory of computing*, pp. 353–362, 2014.
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33: 1877–1901, 2020.
- 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. Palu: Compressing kv-cache with low-rank projection. *arXiv preprint arXiv:2407.21118*, 2024.
- Chen, M., Tworek, J., Jun, H., Yuan, Q., Pinto, H. P. D. O., Kaplan, J., Edwards, H., Burda, Y., Joseph, N., Brockman, G., et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021a.
- Chen, P., Yu, H.-F., Dhillon, I., and Hsieh, C.-J. Drone: Data-aware low-rank compression for large nlp models. *Advances in neural information processing systems*, 34: 29321–29334, 2021b.
- Drineas, P., Kannan, R., and Mahoney, M. W. Fast monte carlo algorithms for matrices i: Approximating matrix multiplication. *SIAM Journal on Computing*, 36(1):132–157, 2006.
- Eckart, C. and Young, G. The approximation of one matrix by another of lower rank. *Psychometrika*, 1(3):211–218, 1936a.
- Eckart, C. and Young, G. The approximation of one matrix by another of lower rank. *Psychometrika*, 1(3):211–218, 1936b.
- Frantar, E., Ashkboos, S., Hoefler, T., and Alistarh, D. Gptq: Accurate post-training quantization for generative pre-trained transformers. *arXiv preprint arXiv:2210.17323*, 2022.
- Gao, L., Tow, J., Abbasi, B., Biderman, S., Black, S., DiPofi, A., Foster, C., Golding, L., Hsu, J., Le Noac’h, A., Li, H., McDonnell, 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. The language model evaluation harness, 07 2024. URL <https://zenodo.org/records/12608602>.
- Gelberg, Y., Eguchi, K., Akiba, T., and Cetin, E. Extending the context of pretrained llms by dropping their positional embeddings. *arXiv preprint arXiv:2512.12167*, 2025.
- Grattafiori, A., Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A., Vaughan, A., et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Hsu, Y.-C., Hua, T., Chang, S., Lou, Q., Shen, Y., and Jin, H. Language model compression with weighted low-rank factorization. *arXiv preprint arXiv:2207.00112*, 2022.
- Huang, Y., Zheng, K., Yu, Z., and Bouganis, C.-S. Itera-llm: Boosting sub-8-bit large language model inference via iterative tensor decomposition. In *2025 IEEE 33rd Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM)*, pp. 114–122. IEEE, 2025.
- Ji, T., Guo, B., Wu, Y., Guo, Q., Shen, L., Chen, Z., Qiu, X., Zhang, Q., and Gui, T. Towards economical inference: Enabling deepseek’s multi-head latent attention in any transformer-based llms. *arXiv preprint arXiv:2502.14837*, 2025.
- Lin, J., Tang, J., Tang, H., Yang, S., Chen, W.-M., Wang, W.-C., Xiao, G., Dang, X., Gan, C., and Han, S. Awq: Activation-aware weight quantization for on-device llm compression and acceleration. *Proceedings of machine learning and systems*, 6:87–100, 2024.

- 385 Liu, A. H., Khandelwal, K., Subramanian, S., Jouault, V.,  
386 Rastogi, A., Sadé, A., Jeffares, A., Jiang, A., Cahill,  
387 A., Gavaudan, A., et al. Ministral 3. *arXiv preprint*  
388 *arXiv:2601.08584*, 2026.
- 389 Mahoney, M. W. and Drineas, P. Cur matrix decompositions  
390 for improved data analysis. *Proceedings of the National*  
391 *Academy of Sciences*, 106(3):697–702, 2009.
- 393 Meng, F., Tang, P., Jiang, F., and Zhang, M. Clover: Cross-  
394 layer orthogonal vectors pruning. In *Forty-second Inter-*  
395 *national Conference on Machine Learning*, 2025.
- 396 Merity, S., Xiong, C., Bradbury, J., and Socher, R.  
397 Pointer sentinel mixture models. *arXiv preprint*  
398 *arXiv:1609.07843*, 2016.
- 400 Mozaffari, M. and Dehnavi, M. M. Slim: One-shot quan-  
401 tized sparse plus low-rank approximation of llms. *URL*  
402 <https://arxiv.org/abs/2410.09615>, 2024.
- 403 Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S.,  
404 Matena, M., Zhou, Y., Li, W., and Liu, P. J. Exploring  
405 the limits of transfer learning with a unified text-to-text  
406 transformer. *Journal of machine learning research*, 21  
407 (140):1–67, 2020.
- 409 Saha, R., Sagan, N., Srivastava, V., Goldsmith, A., and  
410 Pilanci, M. Compressing large language models using  
411 low rank and low precision decomposition. *Advances*  
412 *in Neural Information Processing Systems*, 37:88981–  
413 89018, 2024.
- 414 Sakr, C. and Khailany, B. Espace: Dimensionality reduction  
415 of activations for model compression. *Advances in Neural*  
416 *Information Processing Systems*, 37:17489–17517, 2024.
- 418 Shen, Z., Tao, T., Ma, L., Neiswanger, W., Liu, Z., Wang,  
419 H., Tan, B., Hestness, J., Vassilieva, N., Soboleva, D.,  
420 et al. Slimpajama-dc: Understanding data combinations  
421 for llm training. *arXiv preprint arXiv:2309.10818*, 2023.
- 422 Sun, M., Liu, Z., Bair, A., and Kolter, J. Z. A simple and  
423 effective pruning approach for large language models.  
424 *arXiv preprint arXiv:2306.11695*, 2023.
- 426 Team, M. N. et al. Introducing mpt-7b: A new standard  
427 for open-source, commercially usable llms. *DataBricks*  
428 *(May, 2023) www.mosaicml.com/blog/mpt-7b*, 2023.
- 430 Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux,  
431 M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E.,  
432 Azhar, F., et al. Llama: Open and efficient foundation lan-  
433 guage models. *arXiv preprint arXiv:2302.13971*, 2023a.
- 434 Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi,  
435 A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P.,  
436 Bhosale, S., et al. Llama 2: Open foundation and fine-  
437 tuned chat models. *arXiv preprint arXiv:2307.09288*,  
438 2023b.
- 439 Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones,  
L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. At-  
tention is all you need. *Advances in neural information*  
*processing systems*, 30, 2017.
- Wang, X., Zheng, Y., Wan, Z., and Zhang, M. Svd-  
llm: Truncation-aware singular value decomposition  
for large language model compression. *arXiv preprint*  
*arXiv:2403.07378*, 2024.
- Wang, X., Alam, S., Wan, Z., Shen, H., and Zhang,  
M. Svd-llm v2: Optimizing singular value truncation  
for large language model compression. *arXiv preprint*  
*arXiv:2503.12340*, 2025.
- Wei, J., Wang, X., Schuurmans, D., Bosma, M., Xia, F., Chi,  
E., Le, Q. V., Zhou, D., et al. Chain-of-thought prompting  
elicits reasoning in large language models. *Advances in*  
*neural information processing systems*, 35:24824–24837,  
2022.
- Yuan, Z., Shang, Y., Song, Y., Wu, Q., Yan, Y., and Sun,  
G. Gsvd: Activation-aware singular value decomposition  
for compressing large language models. *arXiv preprint*  
*arXiv:2312.05821*, 2023.
- Zhang, C., Cheng, J., Constantinides, G. A., and Zhao, Y.  
Lqer: Low-rank quantization error reconstruction for llms.  
*arXiv preprint arXiv:2402.02446*, 2024a.
- Zhang, C., Wong, J. T., Xiao, C., Constantinides, G. A., and  
Zhao, Y. Qera: an analytical framework for quantization  
error reconstruction. *arXiv preprint arXiv:2410.06040*,  
2024b.
- Zhang, S. and Pappayan, V. Oats: Outlier-aware pruning  
through sparse and low rank decomposition. *arXiv*  
*preprint arXiv:2409.13652*, 2024.

## A. Notations

Table 4 includes the notations of matrices and vectors and Table 5 summarizes the notations of dimensions in this paper.

Table 4. Notation of matrices and vectors in this paper.

Notation	Description
$\mathbf{R}_{\mathbb{X}\mathbb{X}}$	The autocorrelation matrices of $\mathbf{X} \in \mathbb{R}^{l \times d}$ computed as $\frac{1}{l} \mathbf{X}^T \mathbf{X}$
$\mathbf{R}_{\mathbb{X}\mathbb{X}}^{\frac{1}{2}}$	The corresponding unique symmetric matrix square roots of $\mathbf{R}_{\mathbb{X}\mathbb{X}}$
$\mathbf{q}_{q,i}$	An input row vector to query projection of $i$ -th head
$\mathbf{k}_{k,i}$	An input row vector to key/value projection of $i$ -th head
$\mathbf{X}_q$	Input activation to the query layer
$\mathbf{X}_{kv}$	Input activation to the key/value layer
$\mathbf{W}_{q,i}$	Weight of query projection of $i$ -th head
$\mathbf{W}_{k,i}$	Weight of key projection of $i$ -th head
$\mathbf{W}_{qk,i}$	Fused weight of query/key projection of $i$ -th head
$\widetilde{\mathbf{W}}_{qk,i}$	Low-rank approximation of $\mathbf{W}_{qk,i}$
$\widetilde{\mathbf{W}}_{q,i}$	Approximated $\mathbf{W}_{q,i}$ , left low-rank matrix of $\widetilde{\mathbf{W}}_{qk,i}$
$\widetilde{\mathbf{W}}_{k,i}$	Approximated $\mathbf{W}_{k,i}$ , right low-rank matrix of $\widetilde{\mathbf{W}}_{qk,i}$
$\mathbf{Q}_i$	Query of $i$ -th head
$\mathbf{K}_i$	Key of $i$ -th head
$a_i$	A single attention score of $i$ -th head
$\mathbf{A}_i$	Pre-softmax attention score of $i$ -th head
$\mathbf{A}'_i$	Post-softmax attention score of $i$ -th head
$\mathbf{W}_{v,i}$	Weight of value projection of $i$ -th head
$\mathbf{W}_{o,i}$	Weight of output projection of $i$ -th head
$\mathbf{W}_{vo,i}$	Fused weight of value/output projection of $i$ -th head
$\widetilde{\mathbf{W}}_{vo,i}$	Low-rank approximation of $\mathbf{W}_{vo,i}$
$\mathbf{o}$	A row of attention output
$\mathbf{o}_i$	A row of attention output $i$ -th head
$\widetilde{\mathbf{o}}$	Low-rank approximation of $\mathbf{o}$
$\mathbf{O}$	Attention output matrix
$\mathbf{x}_d$	An input row vector of down projection layer in MLP
$\widetilde{\mathbf{x}}_d$	An input row vector of down projection layer in MLP after low-rank approximation
$\mathbf{y}_{\text{mlp}}$	Output vectors of down projection layer in MLP after low-rank approximation
$\mathbf{X}_{\text{mlp}}$	Input of MLP in transformer
$\mathbf{X}_d$	Input of down projection layer in MLP
$\mathbf{Y}_d$	Output of gate projection layer in MLP
$\mathbf{Y}_u$	Output of up projection layer in MLP
$\mathbf{Y}_{\text{mlp}}$	Output of down projection layer in MLP
$\mathbf{W}_u$	Weight of up projection layer in MLP
$\mathbf{W}_d$	Weight of down projection layer in MLP
$\mathbf{W}_g$	Weight of gate projection layer in MLP
$\widetilde{\mathbf{W}}_d$	Weight of up projection layer in low-rank approximated MLP
$\widetilde{\mathbf{W}}_d$	Weight of down projection layer in low-rank approximated MLP
$\widetilde{\mathbf{W}}_d$	Weight of gate projection layer in low-rank approximated MLP
$\mathbf{r}_i$	The $i$ -th column of $\mathbf{R}_{\mathbb{X}_d \mathbb{X}_d}^{\frac{1}{2}}$
$\mathbf{w}_i$	The $i$ -th row of $\mathbf{W}_d$

Table 5. Notation of dimensions in this paper.

Notation	Description
$l_q$	Query sequence length
$l_{kv}$	Key and value sequence length
$l_{\text{down}}$	Down layer sequence length
$d_m$	Model hidden size
$h_q$	Number of attention (query) heads
$h_{kv}$	Number of key and value heads
$g := \lfloor h_q/h_{kv} \rfloor$	Number of query heads per key/value head in GQA
$d_{vo}$	Head dimension shared by value and head output projection
$d_{qk}$	Head dimension shared by query and key
$d_{\text{inter}}$	Intermediate size of FFN

## B. Derivations for $\mathbf{A}^3$

### B.1. $\mathbf{A}^3$ -QK

#### B.1.1. EQUIVALENT OBJECTIVE

Here we provide the full derivation for Lemma 3.1 from Problem 1:

$$\begin{aligned} & \operatorname{argmin}_{\widetilde{\mathbf{W}}_{qk,i}} \|\mathbf{X}_q(\mathbf{W}_{qk,i} - \widetilde{\mathbf{W}}_{qk,i})\mathbf{X}_{kv}^T\|_F^2 \quad \text{s.t.} \quad \operatorname{rank}(\widetilde{\mathbf{W}}_{qk,i}) = r \\ & \Rightarrow \operatorname{argmin}_{\widetilde{\mathbf{W}}_{qk,i}} \|\mathbf{R}_{\mathbb{X}_q \mathbb{X}_q}^{\frac{1}{2}}(\mathbf{W}_{qk,i} - \widetilde{\mathbf{W}}_{qk,i})\mathbf{R}_{\mathbb{X}_{kv} \mathbb{X}_{kv}}^{\frac{1}{2}}\|_F^2. \end{aligned} \quad (17)$$

We begin with the right-hand side (RHS) of Equation (17). For clarity, we define some intermediate variables:

$$\begin{aligned} & \|\mathbf{X}_q(\mathbf{W}_{qk,i} - \widetilde{\mathbf{W}}_{qk,i})\mathbf{X}_{kv}^T\|_F^2 \\ & = \operatorname{Tr}(\mathbf{X}_{kv}(\mathbf{W}_{qk,i} - \widetilde{\mathbf{W}}_{qk,i})^T \mathbf{X}_q^T \mathbf{X}_q (\mathbf{W}_{qk,i} - \widetilde{\mathbf{W}}_{qk,i}) \mathbf{X}_{kv}^T) \\ & = \operatorname{Tr}(\mathbf{X}_{kv}^T \mathbf{X}_{kv} (\mathbf{W}_{qk,i} - \widetilde{\mathbf{W}}_{qk,i})^T \mathbf{X}_q^T \mathbf{X}_q (\mathbf{W}_{qk,i} - \widetilde{\mathbf{W}}_{qk,i})), \end{aligned} \quad (18)$$

If we assign  $\mathbf{R}_{\mathbb{X}_{kv} \mathbb{X}_{kv}} = \frac{1}{l_{kv}} \mathbf{X}_{kv}^T \mathbf{X}_{kv}$ , and  $\mathbf{R}_{\mathbb{X}_q \mathbb{X}_q} = \frac{1}{l_q} \mathbf{X}_q^T \mathbf{X}_q$ ,

$$\begin{aligned} LHS & = \operatorname{Tr}(l_{kv} l_q \mathbf{R}_{\mathbb{X}_{kv} \mathbb{X}_{kv}} ((\mathbf{W}_{qk,i} - \widetilde{\mathbf{W}}_{qk,i})^T \mathbf{R}_{\mathbb{X}_q \mathbb{X}_q} ((\mathbf{W}_{qk,i} - \widetilde{\mathbf{W}}_{qk,i}))) \\ & = l_{kv} l_q \|\mathbf{R}_{\mathbb{X}_{kv} \mathbb{X}_{kv}}^{\frac{1}{2}}(\mathbf{W}_{qk,i} - \widetilde{\mathbf{W}}_{qk,i})\mathbf{R}_{\mathbb{X}_q \mathbb{X}_q}^{\frac{1}{2}}\|_F^2 \\ & = \|\mathbf{R}_{\mathbb{X}_{kv} \mathbb{X}_{kv}}^{\frac{1}{2}}(\mathbf{W}_{qk,i} - \widetilde{\mathbf{W}}_{qk,i})\mathbf{R}_{\mathbb{X}_q \mathbb{X}_q}^{\frac{1}{2}}\|_F^2. \end{aligned} \quad (19)$$

the positive  $l_q l_{kv}$  can be dropped since they do not affect the minimizer.

#### B.1.2. ANALYTICAL SOLUTION

Here we provide the proof of Theorem 3.2:

*Proof.* We continue with Lemma 3.1.

$$\begin{aligned} & \operatorname{argmin}_{\widetilde{\mathbf{W}}_{qk,i}} \|\mathbf{R}_{\mathbb{X}_q \mathbb{X}_q}^{\frac{1}{2}}(\mathbf{W}_{qk,i} - \widetilde{\mathbf{W}}_{qk,i})\mathbf{R}_{\mathbb{X}_{kv} \mathbb{X}_{kv}}^{\frac{1}{2}}\|_F^2 \quad \text{s.t.} \quad \operatorname{rank}(\widetilde{\mathbf{W}}_{qk,i}) = r \\ & \Rightarrow \operatorname{argmin}_{\widetilde{\mathbf{W}}_{qk,i}} \|\mathbf{R}_{\mathbb{X}_q \mathbb{X}_q}^{\frac{1}{2}} \mathbf{W}_{qk,i} \mathbf{R}_{\mathbb{X}_{kv} \mathbb{X}_{kv}}^{\frac{1}{2}} - \mathbf{R}_{\mathbb{X}_q \mathbb{X}_q}^{\frac{1}{2}} \widetilde{\mathbf{W}}_{qk,i} \mathbf{R}_{\mathbb{X}_{kv} \mathbb{X}_{kv}}^{\frac{1}{2}}\|_F^2. \end{aligned} \quad (20)$$

Note that multiplication by the invertible matrix  $\mathbf{R}_{\mathbb{X}_q \mathbb{X}_q}^{\frac{1}{2}}$  and  $\mathbf{R}_{\mathbb{X}_{kv} \mathbb{X}_{kv}}^{\frac{1}{2}}$  does not change the rank of the matrix  $\mathbf{W}_{qk,i}$ . According

to the Eckart-Young-Mirsky theorem (Eckart & Young, 1936b), the optimal rank  $r$  approximation to  $(\mathbf{R}_{\mathbb{X}_q \mathbb{X}_q}^{\frac{1}{2}} \mathbf{W}_{\text{qk},i} \mathbf{R}_{\mathbb{X}_{\text{kv}} \mathbb{X}_{\text{kv}}}^{\frac{1}{2}})$  is the truncated SVD of  $(\mathbf{R}_{\mathbb{X}_q \mathbb{X}_q}^{\frac{1}{2}} \mathbf{W}_{\text{qk},i} \mathbf{R}_{\mathbb{X}_{\text{kv}} \mathbb{X}_{\text{kv}}}^{\frac{1}{2}})$ :

$$(\mathbf{R}_{\mathbb{X}_q \mathbb{X}_q}^{\frac{1}{2}} \mathbf{W}_{\text{qk},i} \mathbf{R}_{\mathbb{X}_{\text{kv}} \mathbb{X}_{\text{kv}}}^{\frac{1}{2}})_r = \mathbf{U}_{:, :r} \mathbf{\Sigma}_{:, :r} \mathbf{V}_{:, :r}^T, \quad (21)$$

where  $\mathbf{U} \mathbf{\Sigma} \mathbf{V}^T = \text{SVD}(\mathbf{R}_{\mathbb{X}_q \mathbb{X}_q}^{\frac{1}{2}} \mathbf{W}_{\text{qk},i} \mathbf{R}_{\mathbb{X}_{\text{kv}} \mathbb{X}_{\text{kv}}}^{\frac{1}{2}})$ . Thus the optimal rank- $k$  solution to  $\widetilde{\mathbf{W}}_{\text{qk},i}$  is:

$$\widetilde{\mathbf{W}}_{\text{qk},i} = \left(\mathbf{R}_{\mathbb{X}_q \mathbb{X}_q}^{\frac{1}{2}}\right)^{-1} \text{SVD}_r \left(\mathbf{R}_{\mathbb{X}_q \mathbb{X}_q}^{\frac{1}{2}} \mathbf{W}_{\text{qk},i} \mathbf{R}_{\mathbb{X}_{\text{kv}} \mathbb{X}_{\text{kv}}}^{\frac{1}{2}}\right) \left(\mathbf{R}_{\mathbb{X}_{\text{kv}} \mathbb{X}_{\text{kv}}}^{\frac{1}{2}}\right)^{-1}. \quad (22)$$

□

## B.2. A<sup>3</sup>-OV

### B.2.1. EQUIVALENT OBJECTIVE

Similarly, we provide the derivation of the equivalent objective in Problem 2:

$$\begin{aligned} \arg\min_{\widetilde{\mathbf{W}}_{\text{vo},i}} \|\mathbf{P}_i(\mathbf{W}_{\text{vo},i} - \widetilde{\mathbf{W}}_{\text{vo},i})\|_F^2 \quad \text{s.t.} \quad \text{rank}(\widetilde{\mathbf{W}}_{\text{vo},i}) = r \\ \Rightarrow \arg\min_{\widetilde{\mathbf{W}}_{\text{vo},i}} \|\mathbf{R}_{\mathbb{X}_{\mathbf{P}_i} \mathbb{X}_{\mathbf{P}_i}}^{\frac{1}{2}}(\mathbf{W}_{\text{vo},i} - \widetilde{\mathbf{W}}_{\text{vo},i})\|_F^2. \end{aligned} \quad (23)$$

*Proof.* We begin with the right-hand side (RHS) of Equation (23).

$$\begin{aligned} \|\mathbf{O}_i - \widetilde{\mathbf{O}}_i\|_F^2 &= \|\mathbf{P}_i(\mathbf{W}_{\text{vo},i} - \widetilde{\mathbf{W}}_{\text{vo},i})\|_F^2 \\ &= \text{Tr}((\mathbf{W}_{\text{vo},i} - \widetilde{\mathbf{W}}_{\text{vo},i})^T \mathbf{P}_i^T \mathbf{P}_i (\mathbf{W}_{\text{vo},i} - \widetilde{\mathbf{W}}_{\text{vo},i})), \end{aligned} \quad (24)$$

If we assign  $\mathbf{R}_{\mathbb{X}_{\mathbf{P}_i} \mathbb{X}_{\mathbf{P}_i}} = \frac{1}{l_q} \mathbf{P}_i^T \mathbf{P}_i$ ,

$$\begin{aligned} \text{LHS} &= \text{Tr}(l_q(\mathbf{W}_{\text{vo},i} - \widetilde{\mathbf{W}}_{\text{vo},i})^T \mathbf{R}_{\mathbb{X}_{\mathbf{P}_i} \mathbb{X}_{\mathbf{P}_i}} (\mathbf{W}_{\text{vo},i} - \widetilde{\mathbf{W}}_{\text{vo},i})) \\ &= l_q \|\mathbf{R}_{\mathbb{X}_{\mathbf{P}_i} \mathbb{X}_{\mathbf{P}_i}}^{1/2}(\mathbf{W}_{\text{vo},i} - \widetilde{\mathbf{W}}_{\text{vo},i})\|_F^2 \\ &= \|\mathbf{R}_{\mathbb{X}_{\mathbf{P}_i} \mathbb{X}_{\mathbf{P}_i}}^{1/2}(\mathbf{W}_{\text{vo},i} - \widetilde{\mathbf{W}}_{\text{vo},i})\|_F^2. \end{aligned} \quad (25)$$

the positive  $l_q$  can be dropped since they do not affect the minimizer. □

### B.2.2. ANALYTICAL SOLUTION

Here we provide the proof of Theorem 3.3.

*Proof.* We continue with Equation 25:

$$\begin{aligned} \arg\min_{\widetilde{\mathbf{W}}_{\text{vo},i}} \|\mathbf{P}_i(\mathbf{W}_{\text{vo},i} - \widetilde{\mathbf{W}}_{\text{vo},i})\|_F^2 \quad \text{s.t.} \quad \text{rank}(\widetilde{\mathbf{W}}_{\text{vo},i}) = r \\ \Rightarrow \arg\min_{\widetilde{\mathbf{W}}_{\text{vo},i}} \|\mathbf{R}_{\mathbb{X}_{\mathbf{P}_i} \mathbb{X}_{\mathbf{P}_i}}^{\frac{1}{2}}(\mathbf{W}_{\text{vo},i} - \widetilde{\mathbf{W}}_{\text{vo},i})\|_F^2. \end{aligned} \quad (26)$$

Note that multiplication by the invertible matrix  $\mathbf{R}_{\mathbb{X}_{\mathbf{P}_i} \mathbb{X}_{\mathbf{P}_i}}$  does not change the rank of the matrix  $\mathbf{W}_{\text{vo},i}$ . According to the Eckart-Young-Mirsky theorem (Eckart & Young, 1936b), the optimal rank  $r$  approximation to  $(\mathbf{R}_{\mathbb{X}_{\mathbf{P}_i} \mathbb{X}_{\mathbf{P}_i}}^{\frac{1}{2}} \mathbf{W}_{\text{vo},i})$  is the truncated SVD of  $(\mathbf{R}_{\mathbb{X}_{\mathbf{P}_i} \mathbb{X}_{\mathbf{P}_i}}^{\frac{1}{2}} \mathbf{W}_{\text{vo},i})$ :

$$(\mathbf{R}_{\mathbb{X}_{\mathbf{P}_i} \mathbb{X}_{\mathbf{P}_i}}^{\frac{1}{2}} \mathbf{W}_{\text{vo},i})_r = \mathbf{U}_{:, :r} \mathbf{\Sigma}_{:, :r} \mathbf{V}_{:, :r}^T, \quad (27)$$

where  $U\Sigma V^T = \text{SVD}\left(\mathbf{R}_{\mathbb{X}_{P_i}\mathbb{X}_{P_i}}^{\frac{1}{2}} \mathbf{W}_{\text{vo},i}\right)$ . Thus the optimal rank- $k$  solution to  $\widetilde{\mathbf{W}}_{\text{vo},i}$  is:

$$\widetilde{\mathbf{W}}_{\text{vo},i} = \left(\mathbf{R}_{\mathbb{X}_q\mathbb{X}_q}^{\frac{1}{2}}\right)^{-1} \text{SVD}_r\left(\mathbf{R}_{\mathbb{X}_{P_i}\mathbb{X}_{P_i}}^{\frac{1}{2}} \mathbf{W}_{\text{vo},i}\right). \quad (28)$$

□

### B.2.3. ALTERNATIVE SOLUTION TO $A^3$ -OV

Here we elaborate on the alternative solution to Problem 2 by directly minimizing  $\|\mathbf{O} - \widetilde{\mathbf{O}}\|_F^2$  through matrix stacking. With matrix stacking, we can write the overall attention output as two matrix multiplications:

$$\mathbf{O} = \sum_{i=1}^{h_q} \mathbf{O}_i = \sum_{i=1}^{h_q} \mathbf{P}_i \mathbf{W}_{\text{vo},i} = [\mathbf{P}_1 \quad \mathbf{P}_2 \quad \dots \quad \mathbf{P}_{h_q}] \begin{bmatrix} \mathbf{W}_{\text{vo},1} \\ \mathbf{W}_{\text{vo},2} \\ \vdots \\ \mathbf{W}_{\text{vo},h_q} \end{bmatrix} = \mathbf{P}_{\text{cat}} \mathbf{W}_{\text{vo,cat}}, \quad (29)$$

where  $\mathbf{P}_{\text{cat}} \in \mathbb{R}^{l_q \times d_m d_{h_{kv}}}$ ,  $\mathbf{W}_{\text{vo,cat}} \in \mathbb{R}^{d_m d_{h_{kv}} \times d_m}$  denote the concatenated attention score weighted values and the fused value/output projection.

The overall term  $\mathbf{O}$  takes the form of a linear layer  $\mathbf{O} = \mathbf{P}_{\text{cat}} \mathbf{W}_{\text{vo,cat}}$ . If  $\widetilde{\mathbf{O}} = \mathbf{P}_{\text{cat}} \widetilde{\mathbf{W}}_{\text{vo,cat}}$  denotes the approximated  $\mathbf{o}$ , the optimal solution to  $\widetilde{\mathbf{W}}_{\text{vo,cat}}$  is already given by Equation (4).

**Problem 4** (Minimization of overall attention output error). Given a pretrained Transformer layer, the attention output of OV component  $\mathbf{O} = \mathbf{P}_{\text{cat}} \mathbf{W}_{\text{vo,cat}}$  and its approximated form  $\widetilde{\mathbf{O}} = \mathbf{P}_{\text{cat}} \widetilde{\mathbf{W}}_{\text{vo,cat}}$ , approximating the fused value/output projection by minimizing the output error is to minimize the following error:

$$\|\mathbf{P}_{\text{cat}}(\mathbf{W}_{\text{vo,cat}} - \widetilde{\mathbf{W}}_{\text{vo,cat}})\|_F^2 \quad \text{s.t.} \quad \text{rank}(\widetilde{\mathbf{W}}_{\text{vo,cat}}) = r. \quad (30)$$

**Theorem B.1** ( $A^3$ -OV-overall for MHA-NoPE). *The optimal solution to Problem 4 is*

$$\widetilde{\mathbf{W}}_{\text{vo,cat}} = \left(\mathbf{R}_{\mathbb{X}_{P_{\text{cat}}}\mathbb{X}_{P_{\text{cat}}}}^{\frac{1}{2}}\right)^{-1} \text{SVD}_r\left(\mathbf{R}_{\mathbb{X}_{P_{\text{cat}}}\mathbb{X}_{P_{\text{cat}}}}^{\frac{1}{2}} \mathbf{W}_{\text{vo,cat}}\right), \quad (31)$$

where  $\mathbf{R}_{\mathbb{X}_{P_{\text{cat}}}\mathbb{X}_{P_{\text{cat}}}}$  is the autocorrelation matrix respect to  $\mathbf{P}_{\text{cat}}$  and  $\mathbf{R}_{\mathbb{X}_{P_{\text{cat}}}\mathbb{X}_{P_{\text{cat}}}}^{\frac{1}{2}}$  denotes its unique symmetric matrix square root.

Similar to QK component, in practice we assign

$$\mathbf{L}_{vo} := \left(\mathbf{R}_{\mathbb{X}_{P_{\text{cat}}}\mathbb{X}_{P_{\text{cat}}}}^{\frac{1}{2}}\right)^{-1} \mathbf{U}_{:,k}, \quad \mathbf{R}_{vo} := \Sigma_{:k,:k} \mathbf{V}_{:k,:}^T \left(\mathbf{R}_{\mathbb{X}_{P_{\text{cat}}}\mathbb{X}_{P_{\text{cat}}}}^{\frac{1}{2}}\right)^{-1},$$

to get the approximated fused value/output weights with two low-rank matrices. In terms of attention head,  $\mathbf{L}_{vo}$  can be viewed as concatenating all value heads with head dimension of  $r$  together, and  $\mathbf{R}_{vo}$  can be viewed as a shared output head across the values:

$$\mathbf{L}_{vo} := \begin{bmatrix} \mathbf{L}_{v,1} \\ \mathbf{L}_{v,2} \\ \vdots \\ \mathbf{L}_{v,h_q} \end{bmatrix}, \quad \mathbf{R}_{vo} = \mathbf{R}_o,$$

where  $\mathbf{L}_{v,i} \in \mathbb{R}^{d_m \times r}$  and  $\mathbf{R}_o \in \mathbb{R}^{r \times d_m}$ . Attention output can then be computed as follows:

$$\mathbf{O} = \sum_{i=1}^{h_q} \mathbf{P}_i \mathbf{X}_{kv} \mathbf{L}_{v,i} \mathbf{R}_o. \quad (32)$$

Although this solution can theoretically save more parameters than minimizing the per-head attention output loss under the same model performance, it requires significantly more KV-cache storage, even exceeding that of the uncompressed model. This is because the shared rank  $r$  across all heads typically needs to be larger than  $d_{vo}$  of one head to maintain competitive performance. As a result, the KV-cache size increases by a factor of  $\frac{r}{d_{vo}}$  compared to uncompressed models.

### B.3. A<sup>3</sup>-MLP

#### B.3.1. EQUIVALENT OBJECTIVE

Here we provide the derivation of Lemma 3.4 in Problem 3:

$$\begin{aligned} & \operatorname{argmin}_{\mathbf{U}=\operatorname{diag}(u_1, u_2, \dots, u_{d_{\text{inter}}})} \|\mathbf{X}_d \mathbf{U} \mathbf{W}_d - \mathbf{X}_d \mathbf{W}_d\|_F^2 \quad \text{s.t.} \quad \operatorname{rank}(\mathbf{U}) = r, \\ & \Rightarrow \operatorname{argmin}_{\mathbf{U}=\operatorname{diag}(u_1, u_2, \dots, u_{d_{\text{inter}}})} \|\mathbf{R}_{\mathbb{X}_d \mathbb{X}_d}^{\frac{1}{2}} \mathbf{U} \mathbf{W}_d - \mathbf{R}_{\mathbb{X}_d \mathbb{X}_d}^{\frac{1}{2}} \mathbf{W}_d\|_F^2. \end{aligned} \quad (33)$$

*Proof.* We begin with the right-hand side (RHS) of Equation (33). For clarity, we define some intermediate variables:

$$\mathbf{V} := (\mathbf{U} \mathbf{W}_{\text{down}} - \mathbf{W}_{\text{down}}) = [\mathbf{v}_1^T, \mathbf{v}_2^T, \dots, \mathbf{v}_{d_{\text{inter}}}^T]^T, \quad \mathbf{x}_{\text{down}} = [x_1 \quad x_2 \quad \dots \quad x_{d_{\text{inter}}}] . \quad (34)$$

We continue by substituting Equation (34) to RHS of Equation (33):

$$\begin{aligned} & \|\mathbf{X}_d \mathbf{U} \mathbf{W}_d - \mathbf{X}_d \mathbf{W}_d\|_F^2 \\ &= \|\mathbf{X}_d \mathbf{V}\|_F^2 \\ &= \operatorname{Tr}((\mathbf{X}_d \mathbf{V})^T \mathbf{X}_d \mathbf{V}) \\ &= \operatorname{Tr}(\mathbf{V}^T \mathbf{X}_d^T \mathbf{X}_d \mathbf{V}) \\ &= \operatorname{Tr}(\mathbf{X}_d^T \mathbf{X}_d \mathbf{V} \mathbf{V}^T), \end{aligned} \quad (35)$$

If we assign  $\mathbf{R}_{\mathbb{X}_d \mathbb{X}_d} = \frac{1}{l_{\text{down}}} \mathbf{X}_d^T \mathbf{X}_d$ ,

$$\begin{aligned} LHS &= \operatorname{Tr}(\mathbf{R}_{\mathbb{X}_d \mathbb{X}_d} \mathbf{V} \mathbf{V}^T) \\ &= \operatorname{Tr}(l_{\text{down}} \mathbf{R}_{\mathbb{X}_d \mathbb{X}_d}^{\frac{1}{2}} \mathbf{V} \mathbf{V}^T (\mathbf{R}_{\mathbb{X}_d \mathbb{X}_d}^{\frac{1}{2}})^T) \\ &= l_{\text{down}} \|\mathbf{R}_{\mathbb{X}_d \mathbb{X}_d}^{\frac{1}{2}} \mathbf{V}\|_F^2 \\ &= l_{\text{down}} \|\mathbf{R}_{\mathbb{X}_d \mathbb{X}_d}^{\frac{1}{2}} \mathbf{U} \mathbf{W}_d - \mathbf{R}_{\mathbb{X}_d \mathbb{X}_d}^{\frac{1}{2}} \mathbf{W}_d\|_F^2 \\ &= \|\mathbf{R}_{\mathbb{X}_d \mathbb{X}_d}^{\frac{1}{2}} \mathbf{U} \mathbf{W}_d - \mathbf{R}_{\mathbb{X}_d \mathbb{X}_d}^{\frac{1}{2}} \mathbf{W}_d\|_F^2. \end{aligned} \quad (36)$$

the positive  $l_{\text{down}}$  can be dropped since they do not affect the minimizer.  $\square$

### B.4. Extending A<sup>3</sup> for GQA

Similar to GQA's QK Component, we can apply joint SVD to the OV component. The concatenation is done along the second dimension:

$$\operatorname{SVD} \left( \left[ \mathbf{R}_{\mathbb{X}_{kv} \mathbb{X}_{kv}}^{\frac{1}{2}} \widetilde{\mathbf{W}}_{\text{vo},1} \quad \dots \quad \mathbf{R}_{\mathbb{X}_{kv} \mathbb{X}_{kv}}^{\frac{1}{2}} \widetilde{\mathbf{W}}_{\text{vo},g} \right] \right) = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T. \quad (37)$$

Then we make the assignment below to build the shared value head weights  $\widetilde{\mathbf{W}}_{\text{v, shared}}$  and the  $i$ -th approximated output head weights  $\widetilde{\mathbf{W}}_{\text{o},i}$  for this group.

$$\widetilde{\mathbf{W}}_{\text{o},i} := \mathbf{V}_{:,id_m:(i+1)d_m}^T \left( \mathbf{R}_{\mathbb{X}_{kv} \mathbb{X}_{kv}}^{\frac{1}{2}} \right)^{-1}, \quad \widetilde{\mathbf{W}}_{\text{v, shared}} := \mathbf{U}_{:, :r} \mathbf{\Sigma}_{:, :r}. \quad (38)$$

Note that in Equation (37) we use  $\mathbf{R}_{\mathbb{X}_{kv} \mathbb{X}_{kv}}$  instead of  $\mathbf{R}_{\mathbb{X}_{p_i} \mathbb{X}_{p_i}}$ , because  $\mathbf{R}_{\mathbb{X}_{p_i} \mathbb{X}_{p_i}}$  is head-specific, which prevents us from building a shared  $\widetilde{\mathbf{W}}_{\text{v, shared}}$  for all heads in the same GQA group. We name these two methods A<sup>3</sup>-QK-CR and A<sup>3</sup>-OV-CR for GQA's QK and OV components respectively.

## B.5. Extending $\mathbf{A}^3$ for RoPE

Instead of sorting by the product of  $l_2$ -norm,  $\lambda_i = \|\mathbf{L}_{:,i}\|_2^2 \cdot \|\mathbf{R}_{i,:}\|_2^2$ , for  $i = 0, 1, \dots, d_{\text{qk}} - 1$ , we sort by the sum of  $\lambda_i$  of adjacent pairs:

$$\lambda_{2i} = \|\mathbf{L}_{:,2i}\|_2^2 \cdot \|\mathbf{R}_{2i,:}\|_2^2 + \|\mathbf{L}_{:,2i+1}\|_2^2 \cdot \|\mathbf{R}_{2i+1,:}\|_2^2 \quad \text{for } i = 0, 1, \dots, \frac{d_{\text{qk}}}{2} - 1. \quad (39)$$

This will drop pairs of less important columns in  $\mathbf{L}$  and rows in  $\mathbf{R}$ , as well as the corresponding pairs of RoPE frequencies. This requires a frequency index array for each QK pair, and indexing the RoPE constants at runtime. Usually the head dimension  $d_{\text{qk}}$  is 64 or 128, so the RoPE frequency indices can be saved in a compact INT8 array. However, to achieve high throughput/low latency, a custom kernel is needed to fuse indexing and rotation together, which is out of the scope of this paper.

This RoPE extension can be combined with the GQA extension, which means the sorting in Equation (39) is done on the concatenated  $\mathbf{L}$  and  $\mathbf{R}$  matrices in Equation (40).

## C. Adapting $\mathbf{A}^3$ for GQA and RoPE

The  $\mathbf{A}^3$ -QK and  $\mathbf{A}^3$ -OV methods described above are designed for vanilla multi-head attention (MHA), as are most related works discussed in Section 2. However, modern large language models typically employ Transformer variants such as GQA and RoPE. In this subsection, we extend  $\mathbf{A}^3$  to support GQA and RoPE, thereby broadening its applicability to contemporary model architectures.

**Joint SVD for GQA ( $\mathbf{A}^3$ -QK and  $\mathbf{A}^3$ -OV for GQA-NoPE)** In GQA, a key head is shared with multiples query heads in the same QK group. This grouping prevents us from applying Theorem 3.2 to each QK head independently. Inspired by (Ji et al., 2025), we first concatenate the scaled error matrices in Equation (8) within the same QK group and apply joint SVD:

$$\text{SVD} \left( \begin{bmatrix} \mathbf{R}_{\mathbb{X}_q \mathbb{X}_q}^{\frac{1}{2}} \mathbf{W}_{\text{qk},1} \mathbf{R}_{\mathbb{X}_{\text{kv}} \mathbb{X}_{\text{kv}}}^{\frac{1}{2}} \\ \vdots \\ \mathbf{R}_{\mathbb{X}_q \mathbb{X}_q}^{\frac{1}{2}} \mathbf{W}_{\text{qk},g} \mathbf{R}_{\mathbb{X}_{\text{kv}} \mathbb{X}_{\text{kv}}}^{\frac{1}{2}} \end{bmatrix} \right) = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T, \quad (40)$$

where  $g := \lfloor h_q/h_{\text{kv}} \rfloor$  is the number of query heads in this group. Then for this group, we assign

$$\begin{aligned} \widetilde{\mathbf{W}}_{\text{qk},i} &:= \left( \mathbf{R}_{\mathbb{X}_q \mathbb{X}_q}^{\frac{1}{2}} \right)^{-1} \mathbf{U}_{id_m:(i+1)d_m, :r}, \\ \widetilde{\mathbf{W}}_{\text{k, shared}} &:= \mathbf{\Sigma}_{:r, :r} \mathbf{V}_{:r, :d_m}^T \left( \mathbf{R}_{\mathbb{X}_{\text{kv}} \mathbb{X}_{\text{kv}}}^{\frac{1}{2}} \right)^{-1}. \end{aligned} \quad (41)$$

to build the  $i$ -th head's approximated query weights  $\widetilde{\mathbf{W}}_{\text{q},i}$  and the shared head key weights  $\widetilde{\mathbf{W}}_{\text{k, shared}}$ . The subscript with colons, e.g.,  $\mathbf{\Sigma}_{:r, :r}$ , denotes array slicing. Similarly, we can apply joint SVD to the OV component by concatenating the scaled error matrices along the column dimension. A detailed description can be found in Section B.4.

**CUR Approximation for MHA with RoPE ( $\mathbf{A}^3$ -QK for RoPE)** Recent work have shown that not all RoPE frequencies are helpful for the model performance (Barbero et al., 2024; Ji et al., 2025). DroPE (Gelberg et al., 2025) further show that dropping RoPE after training can extend the context length. However, most models still rely on RoPE, motivating our adaptation of  $\mathbf{A}^3$ -QK to RoPE-based attention. RoPE inserts a position-dependent operation before the dot product between queries and keys:

$$\text{RoPE}(\mathbf{q}_{q,i}, m, \mathbf{k}_{k,i}, n) = \mathbf{q}_{q,i} \mathbf{\Phi}_m \mathbf{\Phi}_n^T \mathbf{k}_{k,i}^T, \quad (42)$$

where  $m$  and  $n$  are the position indexes of  $\mathbf{q}_{q,i}$  and  $\mathbf{k}_{k,i}$ ,  $\mathbf{\Phi}_m, \mathbf{\Phi}_n \in \mathbb{R}^{d_{\text{qk}} \times d_{\text{qk}}}$  are matrices rotating adjacent pairs of query elements and key elements. To deal with these pairwise rotations, we use CUR approximation to solve the problem in Lemma 3.1. Similar to  $\mathbf{A}^3$ -MLP, we seek for a rank- $r$  CUR approximation of  $(\mathbf{R}_{\mathbb{X}_q \mathbb{X}_q}^{\frac{1}{2}} \mathbf{W}_{q,i})(\mathbf{W}_{k,i}^T \mathbf{R}_{\mathbb{X}_{\text{kv}} \mathbb{X}_{\text{kv}}}^{\frac{1}{2}})$  that extracts the most important head dimensions as well as RoPE frequencies. Assign  $\mathbf{L} := \mathbf{R}_{\mathbb{X}_q \mathbb{X}_q}^{\frac{1}{2}} \mathbf{W}_{q,i}$  and  $\mathbf{R} = \mathbf{W}_{k,i}^T \mathbf{R}_{\mathbb{X}_{\text{kv}} \mathbb{X}_{\text{kv}}}^{\frac{1}{2}}$ , the

770 objective is

$$\begin{aligned}
 & \arg \min_{U=\text{diag}(u_1, u_2, \dots, u_{d_{\text{qk}}})} \|LUR - LR\|_F^2 \\
 & \text{s.t. } \text{rank}(U) = r.
 \end{aligned}
 \tag{43}$$

775 Instead of sorting by the product of  $l_2$ -norm, *i.e.*,  $\lambda_i = \|\mathbf{L}_{:,i}\|_2 \cdot \|\mathbf{R}_{i,:}\|_2$ , for  $i = 0, 1, \dots, d_{\text{qk}} - 1$ , we sort by the sum of  
 776  $\lambda_i$  of adjacent pairs (Check Section B.5). This will drop pairs of less important columns in  $\mathbf{L}$  and rows in  $\mathbf{R}$ , as well as the  
 777 corresponding pairs of RoPE frequencies.

778 Our adaptation for RoPE can be combined with the joint SVD for GQA, allowing  $A^3$  to be applied to various models. The  
 779 evaluation in Section 4 includes standard MHA, MHA with RoPE, and GQA with RoPE.  
 780

## 781 D. Detailed Experiment Setup

783 **Calibration** We concatenate the texts in SlimPajama and randomly sample 128 sequences of 2048 tokens for calibration.  
 784 We only calibrate on WikiText-2 for Table 8. SlimPajama is a pretraining dataset of high-quality corpus, better capturing the  
 785 statistics of auto-correlation than WikiText2. We calibrate the auto-correlation matrix using BF16 models, but accumulate  
 786 the outer product in FP64.  
 787

788 **Approximation** Since the autocorrelation matrix is symmetric and positive semi-definite, we used SVD to calculate its in-  
 789 verse and matrix square root, which improves the numerical stability. For GQA models, we also use `torch.svd_lowrank`  
 790 instead of `torch.linalg.svd` for faster solving. In cases where the autocorrelation matrices are ill-conditioned, we  
 791 follow the approach of GPTQ (Frantar et al., 2022) and add a small damping term to the zero eigenvalues. In all of our  
 792 experiments across different calibration datasets, the autocorrelation matrices were always invertible.  
 793

794 **Downstream evaluation** We use 0-shot prompting for BoolQ and OpenBookQA, 5-shot for Winogrande, GSM8K, and  
 795 MMLU, 25-shot for ARC-c. Other evaluation parameters are kept as the default provided by `lm-eval-harness`.  
 796

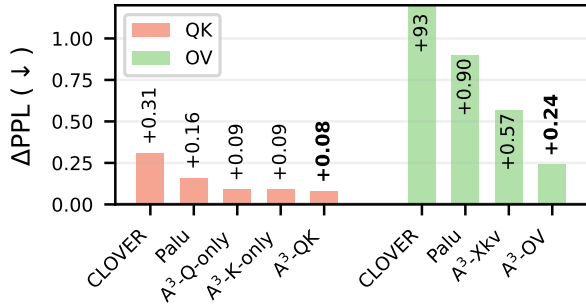
798 **Server specs** We run all the experiments on two GPU boxes, one with two NVIDIA H100s, and the other with 8 NVIDIA  
 799 A100s. In total, we spent around 1200 GPU hours on running  $A^3$ , and 800 hours on baselines. Specifically, for MHA models,  
 800 most of the GPU hours were spent on calibration and VO solving. For GQA models, `torch.svd_lowrank` speeds up  
 801 the VO solving, with most of the GPU hours on calibration and FFN solving. For the ASVD baseline, most of the GPU  
 802 hours were on the mixed-rank search, while other baselines took most of the time on calibration and approximation. Since  
 803 all the calibration, approximation, and evaluation were run on GPUs, our experiments were not bottlenecked by CPUs.  
 804

## 805 E. Ablation Study and other Baselines

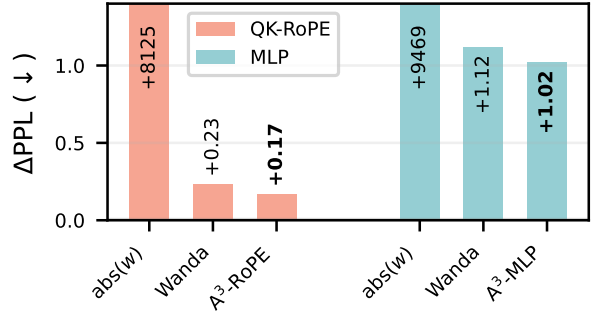
807 We conduct ablation studies to evaluate  $A^3$ 's impact on individual components (QK, OV, MLP). We also include baselines  
 808 that can be applied to the target components to show  $A^3$ 's advantage.  
 809

810 **Attention without RoPE** Theorem 3.2 ( $A^3$ -QK) and Theorem 3.3 ( $A^3$ -OV) provide optimal solutions for MHA-NoPE's  
 811 QK and OV components without the need of adaptation. We evaluate the increased perplexity ( $\Delta\text{PPL}\downarrow$ ) of  $A^3$ -QK and  $A^3$ -OV  
 812 on MPT-7B (MHA-NoPE) in Figure 3a, comparing against CLOVER (Meng et al., 2025) and Palu (Chang et al., 2024),  
 813 with compression ratio=20%. CLOVER is equivalent to  $A^3$ -QK but assumes  $\mathbf{R}_{\mathbf{x}_q, \mathbf{x}_q}$  and  $\mathbf{R}_{\mathbf{x}_{kv}, \mathbf{x}_{kv}}$  are identity matrices (no  
 814 activation information). The bars are grouped by the component being approximated (QK, OV). We add two bars representing  
 815 simplified version  $A^3$ -QK in the QK group,  $A^3$ -Q-only and  $A^3$ -K-only.  $A^3$ -Q-only (K-only) replaces the autocorrelation matrix  
 816 of key (query) in Equation (8) with an identity matrix. We also add a simplified version  $A^3$ -OV in the OV group,  $A^3$ - $\mathbf{x}_{kv}$ ,  
 817 which replaces the autocorrelation matrix of  $\mathbf{p}_i$  in Equation (12) with the autocorrelation matrix of  $\mathbf{x}_{kv}$ . The auto-correlation  
 818 matrices of these simplified versions of  $A^3$  are cheaper to calibrate. We observe that  $A^3$ -QK-SVD and  $A^3$ -OV-SVD and their  
 819 simplified versions outperform CLOVER and Palu by a clear margin.  
 820

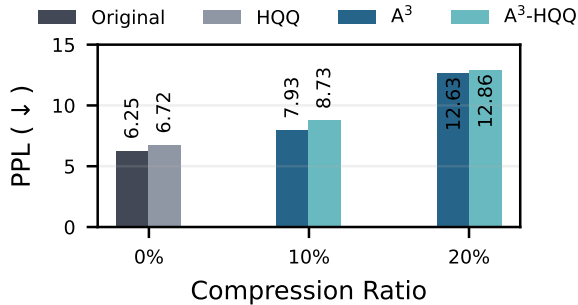
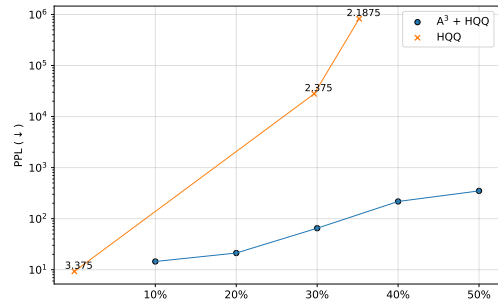
821 **Attention with RoPE** In Section C we propose using CUR approximation to solve Problem 1 for attention with RoPE,  
 822 which follows a similar approach as  $A^3$ -MLP in Section 3.3. Here we compare against structured pruning baselines that  
 823 can be adapted for this problem, including `abs(w)` and Wanda (Sun et al., 2023). `abs(w)` represents the classic pruning  
 824 method that drops weights with smaller magnitudes, while Wanda sorts by the product between the weight magnitude and



(a) QK and OV (MHA-NoPE).



(b) QK-RoPE and MLP (MHA-RoPE).

Figure 3. Ablation study of  $A^3$  components. (a) QK and OV on MPT-7B. (b) QK-RoPE and MLP on LLaMA-2-7B.(a) Perplexity (↓) on WikiText-2 using HQQ (4 bits) for both original and  $A^3$ -applied LLaMA-3.1-8B.(b) Perplexity (↓) on WikiText-2 for HQQ +  $A^3$  relative to HQQ alone in sub-4-bit on LLaMA-3.1-8B.Figure 4. Comparison of perplexity results on WikiText-2 across different HQQ and  $A^3$  settings for LLaMA-3.1-8B.

the average  $l_2$ -norm of the activation row <sup>1</sup>. Figure 3b illustrates  $\Delta PPL$  of LLaMA-2-7B on WikiText2, indicating the advantage of  $A^3$  over  $\text{abs}(w)$  and Wanda.

## F. Discussion

**Quantization compatibility** Here we show that  $A^3$  can be combined together with quantization. Figure 4a presents the perplexity of quantized LLaMA-3.1-8B with HQQ 4-bit quantization, before and after applying  $A^3$ . The small amount of extra model performance degradation caused by  $A^3$  indicates its orthogonality to quantization.

In extreme compression regimes, combining low-rank methods with quantization enables a continuous range of compression levels beyond the discrete choices offered by quantization alone, leading to a substantially improved Pareto frontier. As shown in Figure 4b, sub-3-bit quantization by itself destabilizes the model, whereas  $A^3$  combined with quantization achieves a markedly better accuracy-compression trade-off than quantization alone. Table 10 reports the  $\Delta PPL$  of MPT-30B on WikiText2 under extreme KV-cache compression. While 4-bit HQQ quantization introduces only minor degradation, more aggressive 2-bit quantization leads to a sharp increase in perplexity. In contrast, integrating  $A^3$  with 4-bit HQQ enables higher compression ratios with limited performance loss, and fine-tuning consistently mitigates the remaining degradation, achieving  $\Delta PPL = 0.99$  at an overall  $6.67\times$  compression ratio. Fine-tuning set up is described in the [Fine-tuning paragraph](#).

<sup>1</sup>In the case of Equation (43),  $\text{abs}(w)$  drops columns (rows) by the column (row) sum of magnitudes of  $W_{q,i}$  ( $W_{k,i}^T$ ), while Wanda assumes non-diagonal elements in  $R_{x_q, x_q}$  and  $R_{x_{kv}, x_{kv}}$  are all zeros.

Table 6. A comparison of perplexity (↓) on WikiText-2 and C4 of LLaMA-7B under 20% mixed-rank compression rate.

Method	WikiText-2	C4
Original	5.68	7.65
$A^3$	7.21	10.01
SVD-LLM v2	10.53	13.00
ASVD	10.45	13.1
$A^3$ -mix	<b>7.11</b>	<b>9.86</b>

This section presents a simple low-rank quantization setup. For more advanced low-rank plus quantization methods, we refer the reader to QERA (Zhang et al., 2024b), CALDERA (Saha et al., 2024) and ITERA (Huang et al., 2025).

**Mixed-rank allocation** SVD-LLM v2 (Wang et al., 2025) and ASVD (Yuan et al., 2023) have demonstrated that different layers exhibit varying sensitivity to rank reduction. Here we show that  $A^3$  can also benefit from mixed-rank allocation, achieving performance gain with minimal effort. Specifically, we conduct a simple search over rank allocations for each decoder layer in LLaMA-7b. As shown in Table 6,  $A^3$ -mixed outperforms both the uniform  $A^3$  and other mixed-rank approaches.

**Limitation of CUR decomposition** One limitation of  $A^3$  is its reliance on CUR decomposition for RoPE-based attention and MLP, which does not guarantee an optimal solution like SVD. When targeting a small compression ratio (e.g., CRatio=10%), CUR decomposition provides a good trade-off between performance and compression. As the compression ratio increases, its performance degrades much faster than SVD and is eventually surpassed by SVD-based approaches. However, we argue that even for SVD-based approaches, the model performance under a compression ratio larger than 10% is already very poor. For example, the C4 perplexity of LLaMA-3.1-70B with CRatio=20% is 13.77 for SVD-LLM, which is even higher than the original LLaMA-3.1-8B. A retraining is needed in this case, but is out of the scope of this paper.

**Choice of calibration datasets** We compared the model performance when calibrating on SlimPajama and WikiText-2. We find calibrating on SlimPajama gives closer perplexities on Wikitext-2 and C4, regardless of the compression level (Appendix Table 7). However, calibrating on WikiText-2 has a widening perplexity gap between WikiText-2 and C4 as the compression ratio increased, especially for SVD-LLM, which potentially indicates overfitting. We hypothesize that this may contribute to cases where SVD-LLM appears to perform better on particular downstream tasks, like Winogrande, in Table 2. To validate this, we used a more diverse calibration set with samples from SlimPajama and PTB in Table 8. The results show that with this mixture of calibration sets,  $A^3$  achieves a higher accuracy than SVD-LLM on Winogrande.

Table 7. Performance of LLaMA-7b compressed by SVD-LLM and  $A^3$  under different compression ratio using calibration data sampled from Slimpajama (our setting) and WikiText-2 datasets. The performance are reported by the average and difference in perplexity ( $\downarrow$ ) of Wikitext-2 and C4 datasets.

Method	10%		20%		40%	
	Avg	$ \Delta $	Avg	$ \Delta $	Avg	$ \Delta $
SVD-LLM (SlimPajama)	9.50	2.54	9.50	2.54	30.84	1.00
$A^3$ (SlimPajama)	7.24	2.24	8.52	2.79	25.22	9.83
SVD-LLM (WikiText-2)	9.62	5.07	11.89	7.90	44.58	61.69
$A^3$ (WikiText-2)	7.24	2.25	8.50	3.68	12.82	18.95

Table 8. Performance evaluation of LLaMA-3.1-8b (20% compression) using SVD-LLM and  $A^3$  with two calibration datasets: SlimPajama and a 50/50 SlimPajama-PTB mixture. Metrics include perplexity ( $\downarrow$ ) of Wikitext-2, C4 and slimPajama, and accuracy ( $\uparrow$ ) of BoolQ, Winogrande, ARC-c (with their average).

Method	WikiText-2	C4	SlimPajama	BoolQ	Winogrande	ARC-c	Avg.
SVD-LLM (SlimPajama)	42.28	33.6	27.44	0.6948	0.644	0.2534	0.5307
$A^3$ (SlimPajama)	11.36	17.87	13.58	0.6823	0.6417	0.3345	0.5528
SVD-LLM (SlimPajama+PTB)	36.76	36.62	31.79	0.7349	0.6717	0.2671	0.5579
$A^3$ (SlimPajama+PTB)	11.47	18.57	14.30	0.7220	0.6938	0.3524	0.5894

**Fine-tuning performance** Here we provide the  $A^3$  performance on Wikitext-2 (PPL  $\downarrow$ ) with simple LoRA fine-tuning setting. Following the SVD-LLM setup, we applied LoRA (rank 8) on  $A^3$  with 50K Alpaca-cleaned samples over 2 epochs. Tables 9 and 10 show  $A^3$  benefits notably from even basic fine-tuning due to its strong initialization.

**Relation of local objective reduction to end-to-end perplexity** Here we provide a diagnostic analysis to better understand the gap between individual local objective reduction and their combined effect on end-to-end perplexity. Table 11 compares

Table 9. Comparison of  $A^3$  and  $A^3$  + Fine-Tuning across compression ratios for Llama-2-7b.

Compression Ratio	$A^3$	$A^3$ + Fine-Tuning
20%	7.22 (+1.73)	6.94 (+1.45)
40%	32.04 (+24.31)	10.53 (+5.04)

Table 10. MPT-30B  $\Delta$ PPL relative to baseline (no compression) on WikiText2.

Method	Attention FLOPs/Parameter/KV Cache Compression Ratio	Without Fine-Tuning	Fine-Tuning
Dense	1.00x	0	0
Pure 4-bit HQQ Quantization	4.00x	+0.11	/
Pure 2-bit HQQ Quantization	8.00x	+12.78	+2.80
4-bit HQQ + $A^3$ @ 20%	5.00x	+0.99	+0.59
4-bit HQQ + $A^3$ @ 40%	6.67x	+1.15	+0.99
4-bit HQQ + $A^3$ @ 60%	10.00x	+18.15	+2.74

how locally compressing QK and OV impacts the global perplexity for two model architectures: MPT-7B and LLaMA-3.1-8B. For small compression ratios, the increase in perplexity is approximately equal to the sum of the individual contributions from QK and OV, particularly for standard  $A^3$ -QK and  $A^3$ -OV configurations without RoPE in MPT-7b. At larger compression ratios, the sum of contributions from QK and OV remains roughly on the same order of magnitude as the observed end-to-end perplexity.

Table 11. Perplexity changes under local and joint compression. Shows the effect of compressing QK and OV individually, their summed contribution (QK + OV), and **Both**, the end-to-end perplexity when both are compressed jointly, across different compression ratios for LLaMA-3.1-8B (left) and MPT-7B (right).

LLaMA-3.1-8B	5%	10%	15%	20%	40%	MPT-7B	5%	10%	15%	20%	40%
QK	0.07	0.16	0.32	0.56	13.28	QK	-0.004	0.005	0.040	0.092	0.75
OV	0.27	0.39	0.58	0.78	2.79	OV	0.048	0.097	0.166	0.248	0.98
QK + OV	0.34	0.55	0.90	1.34	16.07	QK + OV	0.045	0.103	0.206	0.340	1.73
Both	0.35	0.59	1.00	1.58	25.07	Both	0.044	0.102	0.197	0.313	1.52

## G. Runtime Analysis

### G.1. $A^3$ compression translates to gains in runtime performance

Here we provide more details for the runtime performance  $A^3$  including the rank, theoretical FLOPs for one decoder layer, throughput and peak allocated GPU memory across different compression ratio on 1 H100 GPU.

To compute theoretical FLOPs for a single decoder block during prefill, we sum attention, FFN, normalization, and residual costs:

$$FLOPs_{total} = FLOPs_{attn} + FLOPs_{mlp} + FLOPs_{norm} + FLOPs_{resid} ,$$

Attention ( $Q, K, V, O$  projections +  $QK^T$  + softmax +  $AV$ ):

$$FLOPs_{attn} = 8BLHD + 4BL^2AD + BL^2A ,$$

Feedforward network (3 projections, each 2 FLOPs per multiply-add):

$$FLOPs_{mlp} = 6BLIH ,$$

Normalization and residuals (each counted twice):

$$FLOPs_{norm} = 2BLH \quad FLOPs_{resid} = 2BLH ,$$

Final total:

$$\text{FLOPs}_{\text{total}} = 8BLHD + 4BL^2AD + BL^2A + 6BLIH + 4BLH.$$

Where:  $B$  = batch size,  $L$  = sequence length,  $H$  = hidden size,  $D$  = head dimension,  $A$  = number of attention heads,  $I$  = MLP intermediate size.

Since  $A^3$  compresses the model by proportionally reducing both  $D$  and  $I$ , the overall theoretical FLOPs reduction is approximately equal to the compression ratio, as many operations scale with  $D$  and  $I$ . However, it is not exactly equal due to additional operations such as normalization, residual connections, and softmax. Table 12 highlights that  $A^3$  compression ratio can directly translate to gains in runtime performance that closely align with the theoretical expectation without requiring specialized kernels.

Table 12. Analysis of runtime performance in LLaMA-2-13b across varying compression ratios and methods.

Compression	Ranks (qk/vo/mlp)	Method	Theoretical FLOPs	Throughput (token/s)	Peak Memory (MB)	Theoretical FLOPs	Speedup	Peak Memory
Original	128/128/13824	Eager	$2.77 \times 10^{12}$	7285	35004	1.00x	1.00x	1.00x
				12319	32917	1.00x	1.69x	0.94x
20%	128/128/13824	Eager	$2.16 \times 10^{12}$	8077	28114	0.78x	1.11x	0.80x
				15096	26037	0.78x	2.07x	0.74x
40%	128/128/13824	Eager	$1.56 \times 10^{12}$	9350	21336	0.56x	1.28x	0.61x
				20237	19270	0.56x	2.78x	0.55x
60%	128/128/13824	Eager	$1.08 \times 10^{12}$	10405	16139	0.39x	1.43x	0.46x
				25554	14078	0.39x	3.51x	0.40x

## G.2. $A^3$ Throughput Evidence at Scale

To evaluate the robustness of  $A^3$ 's throughput gains, Figure 5 shows throughput gains across a variety of settings, varying GPU type, batch size, model size, compression ratio, sequence length, and attention kernel. The experiments include:

- **GPU & Model Size:** Single A6000 (Llama-3.2-1B/Llama-3.2-3B/Llama-3.1-8B), Single H100 (Llama-3.2-3B/Llama-2-13B/Qwen3-32B)
- **Batch Size:** 1, 2, 4 for A6000; 1, 4, 8 for H100
- **Compression Ratio:** 20%, 40%
- **Sequence Length:** 1024, 2048
- **Attention Kernel:** eager, SDPA

Aligned with the analysis in Appendix G.1,  $A^3$  consistently achieves speedups close to the theoretical gains, as it reduces the effective problem size without introducing overhead or additional kernel launches in SVD-LLM. This is also reflected in the decoding throughput in Table 13.

Larger models and SDPA attention benefit more from  $A^3$ , since the reduced head dimension not only compresses the linear layers but also decreases attention FLOPs.

## H. Offline Compression Cost Analysis

Table 14 shows the compression time for both MHA (Llama-2-7b) and GQA (Llama-3-8b). Note that the compression is done offline before deployment and inference.  $A^3$ -OV compression for  $A^3$  on MHA models is more time-consuming because Equation 12 must be applied separately to each attention head. To ensure numerical stability, we use the weight truncation method from SVD-LLM-v2 (Wang et al., 2025), which involves two SVD operations per head: one to compute  $R^{1/2}$ , and another for the primary decomposition.

Figure 6 shows how offline compression time grows with model size. We evaluate 20% compression on a single H100 GPU for Llama-3.2-1B, Llama-3.2-3B, Llama-3.1-8B, Qwen3-32B, and Llama-3-70B. Overall, compression time

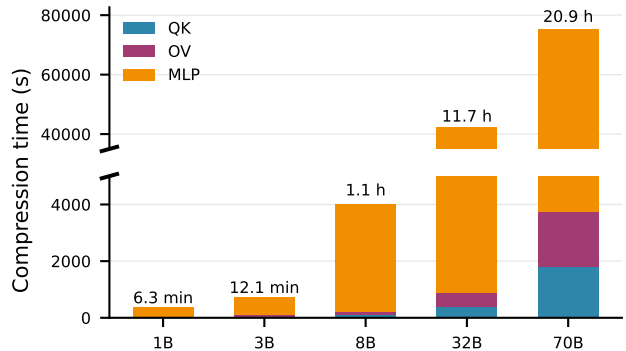


Figure 6.  $A^3$  offline compression time breakdown across model sizes

Table 13. Decoding throughput (TPS) for LLaMA-2-13B under different compression ratios on NVIDIA B200 using SDPA. Baselines correspond to different prefill/decode configurations. A<sup>3</sup> consistently improves throughput, while SVD-LLM often reduces it.

Baseline	Compression Ratio	SVD-LLM ( $\Delta$ TPS)	SVD-LLM (Gain %)	A <sup>3</sup> ( $\Delta$ TPS)	A <sup>3</sup> (Gain %)
548 TPS (1024/1024)	0.2	-39.8	-7.3%	+34.1	+6.2%
	0.4	-41.9	-7.6%	+39.8	+7.1%
	0.6	-48.8	-9.5%	+86.0	+15.7%
397.3 TPS (512/1536)	0.2	-0.5	-0.1%	+118.4	+29.8%
	0.4	-2.3	-0.6%	+118.5	+29.8%
	0.6	-5.0	-1.3%	+115.7	+29.1%
520.9 TPS (512/1536)	0.2	-0.5	-0.1%	+84.1	+16.1%
	0.4	-7.9	-1.5%	+81.9	+15.7%
	0.6	-8.9	-1.7%	+114.6	+21.2%

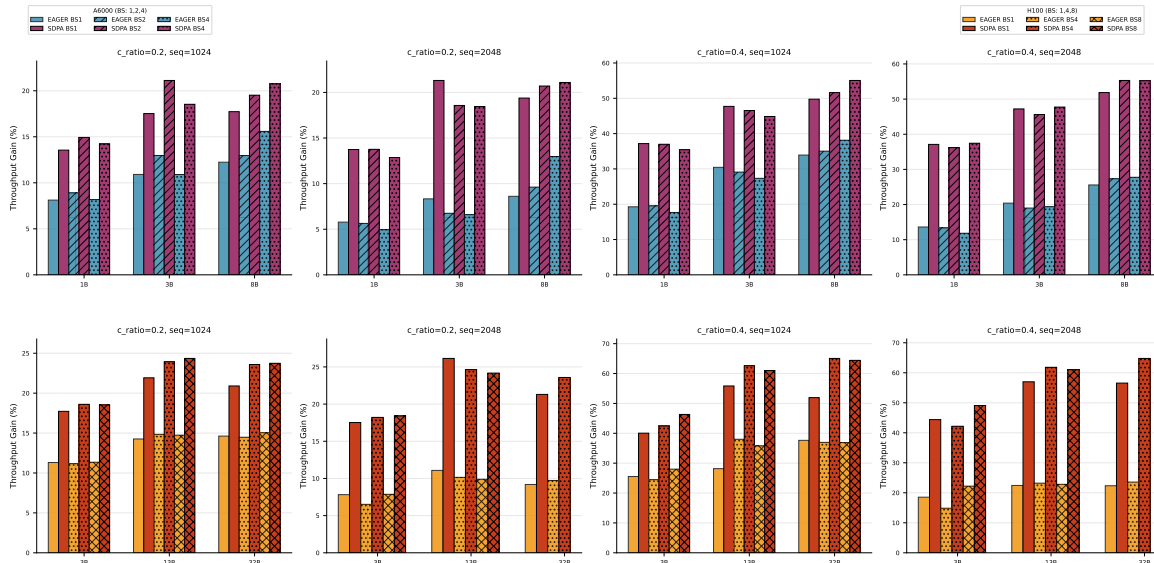


Figure 5. Throughput gains of A<sup>3</sup> relative to the uncompressed model across a range of settings. **C\_ratio** denotes the compression ratio, and **Seq** denotes the sequence length. The top row shows results on a single A6000, while the bottom row shows results on a single H100.

increases rapidly with model scale because many operations in A<sup>3</sup> have computational complexity greater than  $O(n^2)$ , where  $n$  is the model’s hidden dimension

It is worth noting that the compression of each layer is fully independent, so in practice A<sup>3</sup> can be parallelized across  $n$  GPUs to achieve roughly  $n$  times speedup.

For GQA’s QK projections with RoPE, A<sup>3</sup>-QK becomes increasingly expensive as hidden size grows. By the time the model reaches 70B parameters, A<sup>3</sup>-QK accounts for more than 90% of total compression time due to the cost of joint decomposition.

## I. Performance on Phi and Mistral Model

To evaluate the performance of A<sup>3</sup> across other model families, we provide additional results for the Phi and Mistral models following Table 1 and Table 2 metrics. Thanks to its optimization-based design, A<sup>3</sup> continues to perform very well for this model family.

Table 14. Time (in minutes) for different model components.

Model	A <sup>3</sup> -QK (min)	A <sup>3</sup> -OV (min)	A <sup>3</sup> -MLP (min)
LLaMA-7B MHA	00:49	46:55	12:55
LLaMA-8B GQA	00:56	01:30	23:00

Table 15. A comparison of phi-3-medium-4k-instruct (14B) perplexity ( $\downarrow$ ) on WikiText2, C4, and SlimPajama.

Compression	Method	Wikitext-2	SlimPajama	C4
10%	SVD-LLM	6.81 (+2.50)	8.40 (+1.69)	10.47 (+1.66)
	A <sup>3</sup>	<b>5.44 (+1.14)</b>	<b>7.28 (+0.58)</b>	<b>9.48 (+0.67)</b>
20%	SVD-LLM	8.14 (+3.83)	9.67 (+2.96)	11.90 (+3.10)
	A <sup>3</sup>	<b>6.40 (+2.10)</b>	<b>8.16 (+1.46)</b>	<b>10.59 (+1.79)</b>

Table 16. A comparison of downstream task accuracy ( $\uparrow$ ) of phi-3-medium-4k-instruct (14B).

Compression	Method	ARC Challenge	BoolQ	OpenbookQA	GSM8K (Strict)	MMLU	Avg
-	Original	0.6672	0.8850	0.4120	0.8279	0.7797	0.7144
10%	SVD-LLM	0.5751	0.8703	0.3720	0.6179	0.7134	0.6297
	A <sup>3</sup>	0.6118	0.8841	0.3800	0.7589	0.7340	<b>0.6738</b>
20%	SVD-LLM	0.5034	0.8618	0.3260	0.4913	0.6773	0.5720
	A <sup>3</sup>	0.5273	0.8645	0.3480	0.6073	0.6715	<b>0.6037</b>

Table 17. Perplexity ( $\downarrow$ ) on WikiText-2, C4, and SlimPajama across compression ratios for Ministral-3 models. Values in parentheses denote absolute degradation over the uncompressed baseline. A<sup>3</sup> consistently outperforms SVD-LLM.

Model	Method	WikiText-2 (10%)	C4 (10%)	SlimPajama (10%)	WikiText-2 (20%)	C4 (20%)	SlimPajama (20%)
Ministral-3-3B	SVD-LLM	19.97 (+12.27)	24.11 (+11.52)	18.27 (+9.03)	34.76 (+27.06)	33.83 (+21.25)	26.49 (+17.25)
	A <sup>3</sup>	11.82 (+4.12)	16.17 (+3.59)	12.15 (+2.91)	16.94 (+9.24)	25.58 (+13.00)	18.53 (+9.29)
Ministral-3-8B	SVD-LLM	15.42 (+8.82)	19.67 (+8.52)	14.72 (+6.58)	26.71 (+20.11)	26.72 (+15.58)	20.28 (+12.14)
	A <sup>3</sup>	7.94 (+1.33)	13.11 (+1.97)	9.66 (+1.52)	10.42 (+3.82)	16.99 (+5.85)	12.66 (+4.53)