

000  
001  
002  
003  
004  
005  
006  
007  
008  
009  
010  
011  
012  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030  
031  
032  
033  
034  
035  
036  
037  
038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053  

# UNBIASED GRADIENT LOW-RANK PROJECTION

Anonymous authors

Paper under double-blind review

## ABSTRACT

Memory-efficient optimization is critical for training increasingly large language models (LLMs). A popular strategy involves gradient low-rank projection, storing only the projected optimizer states, with GaLore being a representative example. However, a significant drawback of many such methods is their lack of convergence guarantees, as various low-rank projection approaches introduce inherent biases [at each step](#) relative to the original optimization algorithms [even after taking expectation over stochastic sampling](#), which contribute to performance gaps compared to full-parameter training. Aiming to tackle this problem, this paper investigates the layerwise sampling technique for debiasing low-rank projection mechanisms. In particular, an instantiation of the paradigm gives rise to a novel and [step-wise](#) unbiased low-rank optimization method built upon GaLore’s mechanism and the Muon algorithm, named GaLore Unbiased with Muon (GUM). We theoretically prove our method matches the convergence guarantees of the base Muon algorithm while preserving the memory efficiency of low-rank techniques. Empirical experiments on LLM fine-tuning and pretraining also demonstrate non-trivial improvements over GaLore and even better performance than full-parameter training. Further investigation shows that the improvement of this technique comes from a more uniform distribution of knowledge inside layers, leading to more efficient utilization of the model parameter space and better memorization.

## 1 INTRODUCTION

Large language models (LLMs) have demonstrated impressive performance across a diverse range of tasks, including conversation (Ouyang et al., 2022; Grattafiori et al., 2024b), mathematical reasoning (Guo et al., 2025), and agentic applications (Qin et al., 2025). The advancement of these powerful LLMs demands substantial GPU memory due to the large size of the underlying models. For example, training a 70B model with full parameters requires approximately 1.2 terabytes of GPU memory, which exceeds the capacity of even  $8 \times$ H100 GPUs.

To address this issue, memory-efficient training techniques such as GaLore (Zhao et al., 2024) have been introduced. GaLore projects gradients into a low-rank space, reducing the memory footprint of optimizer states during training. Specifically, it employs the top- $r$  components from Singular Value Decomposition (SVD) to define a compact low-rank space, into which the gradients are projected as  $R_t \leftarrow P_t^\top G_t$ . The optimization step is then performed in this low-rank space, enabling memory savings for the optimizer states. For example, the first and second moments in Adam (Kingma & Ba, 2014) are updated using  $\tilde{M}_t \leftarrow \beta_1 \tilde{M}_{t-1} + (1 - \beta_1)R_t$  and  $\tilde{V}_t \leftarrow \beta_2 \tilde{V}_{t-1} + (1 - \beta_2)R_t$ , where the low-rank projected gradient  $R_t$  replaces the original gradient  $G_t$ . After the optimization step, the parameter update is projected back to the original space.

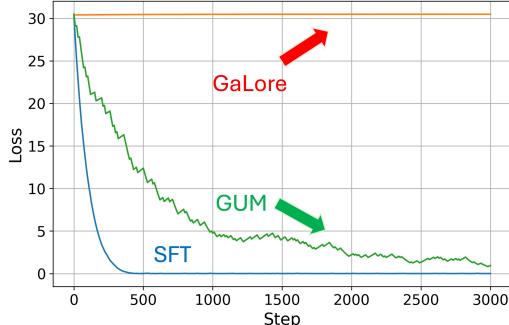


Figure 1: A counterexample of GaLore in linear regression with Muon optimizer (Jordan et al., 2024), where its debiased version GUM converges while GaLore fails to converge.

054 **Algorithm 1** Low-rank projection based gradient descent algorithms

---

```

055 1: Input: Initial weight  $W_0 \in \mathbb{R}^{m \times n}$  (suppose  $m \leq n$ ), number of iterations  $K$ , learning rate  $\eta$ ,
056   projection rank  $r$ .
057 2: for  $t = 0$  to  $K - 1$  do
058   3:    $G_t = G(W_t)$                                  $\triangleright$  Obtain the gradient at  $W_t$ 
059   4:    $P_t \leftarrow \text{get\_projector}()$             $\triangleright$  Obtain the projector  $P_t \in \mathbb{R}^{m \times r}$ 
060   5:    $\tilde{G}_t = P_t^\top G_t$                        $\triangleright$  Obtain projection of the gradient in the low-rank space
061   6:    $S_t \leftarrow \text{optimizer.update\_state}(\tilde{G}_t)$      $\triangleright$  Run the base algorithm with  $\tilde{G}_t$ 
062   7:    $W_{t+1} = W_t - \eta P_t S_t$             $\triangleright$  Project the update back and update the weights
063 end for
064
065

```

---

066 Nevertheless, most low-rank optimization methods introduce biased gradient estimations [at each step](#) 067 during training (Muhamed et al., 2024; Zhang et al., 2024a; He et al., 2024; Huang et al., 2025), which 068 can lead to suboptimal convergence behavior and measurable performance gaps compared to standard 069 full-parameter training. These biases arise because low-rank projections, while computationally and 070 memory efficient, do not fully preserve the direction and magnitude of the true gradient, especially 071 in high-dimensional parameter spaces. As a result, the optimization trajectory diverges from that 072 of full-precision training, potentially causing slower convergence, reduced final model quality, or 073 instability in certain regimes (Zhao et al., 2024; Ding et al., 2022; Zhang et al., 2024a; Huang et al., 074 2025). This limitation is particularly critical when pre-training large language models (LLMs), where 075 even small discrepancies in gradient estimation can propagate and amplify across many layers and 076 iterations.

077 To address this fundamental issue, we investigate the general debiasing technique using layerwise 078 sampling (Pan et al., 2024), which preserves the memory efficiency of training methods via randomly 079 freezing most of the layers. Specifically, the unique strength of layerwise sampling over the typical 080 low-rank projected algorithms of GaLore is analyzed both theoretically and empirically. The 081 introduction of the debiasing technique into GaLore gives rise to a new algorithm called **Galore** 082 **Unbiased with Muon (GUM)**, which demonstrates much better convergence guarantees and practical 083 performance in LLM training tasks. We summarize our major contributions as follows:

- 084 • We investigate the layerwise-sampling debiasing technique and propose a novel algorithm 085 called **Galore Unbiased with Muon (GUM)**, which unifies the strengths of GaLore and 086 Muon. GUM achieves the same theoretical convergence guarantees as Muon while retaining 087 the memory efficiency of GaLore, enabling scalable and effective training of large models.
- 088 • Empirical experiments in LLM training demonstrate that GUM consistently outperforms 089 GaLore in instruction-following, mathematical reasoning, and commonsense reasoning tasks 090 under the same memory budget. Surprisingly, in LLM pre-training experiments, GUM even 091 outperforms full-parameter trained AdamW by a non-trivial overall accuracy margin of 0.3%-1.1%, while obtaining on-par or better performance than AdamW in 6 out of 7 tasks.
- 092 • We analyze the underlying reasoning of GUM’s empirical improvements, discovering that 093 its high-rank update nature leads to a larger overall stable rank and more evenly distributed 094 singular values in model weights, which further induce a more long-tailed activation pattern 095 in trained models. This implies the performance gain is brought by more efficient utilization 096 of the model parameter space, in other words, better memorization.

099 **2 RELATED WORK**

101 **Parameter-Efficient Algorithms in Practice.** Parameter-Efficient Fine-Tuning (PEFT) methods 102 are widely adopted for training large-scale LLMs in practice. A typical approach is LoRA (Hu 103 et al., 2022), which freezes the original model and attaches a small trainable low-rank adapter, 104 thereby reducing memory consumption and improving training efficiency. However, LoRA has 105 been reported to exhibit a non-trivial performance gap compared to full-parameter training (Ding 106 et al., 2022; Lialin et al., 2023), due to its altered parameter space. These changes in parameter 107 space also introduce theoretical challenges in analyzing LoRA’s convergence properties with respect 108 to the original parameter space. To address the aforementioned deficiencies and extend LoRA to

108 larger-scale training settings, GaLore (Zhao et al., 2024) proposes a different approach, which projects  
 109 the gradients—rather than the parameters—into low-rank spaces. In doing so, the error between the  
 110 full gradients and the approximated gradients becomes numerically quantifiable, as they now operate  
 111 within the same parameter space.

112 Following GaLore, a number of low-rank projection-based algorithms have emerged, where the  
 113 key component follows a similar paradigm to Algorithm 1, but with different projection matrices  
 114  $P$ . GaLore utilizes the top- $r$  entries  $U[:, :r]$  from SVD, which is computationally expensive. To  
 115 address this issue, GRASS (Muhamed et al., 2024) derives a sparse projection matrix  $P$  based solely  
 116 on the row norms of the gradients. Specifically, each projection entry is sampled from a multimodal  
 117 distribution proportional to the row norms. GRASS has been reported to achieve performance  
 118 comparable to GaLore with lower computational cost, though no theoretical guarantees have been  
 119 provided regarding its convergence. LoGE (Zhang et al., 2024a) obtains the low-rank projection  $P$   
 120 by decomposing the original weight matrix  $W = BC$ , thereby implicitly allowing the backward  
 121 gradient to be low-rank. However, it is difficult to guarantee theoretical convergence due to the  
 122 empirical nature of the low-rank decomposition. GradNormLoRP (Huang et al., 2025) combines  
 123 ideas from LoRA and GaLore, resulting in a two-level projection  $P$  that further enhances memory  
 124 efficiency and reduces training cost. A variety of salience-aware sparse projections are also employed  
 125 in (Guo et al., 2020; Sung et al., 2021; Ansell et al., 2021; Das et al., 2023; Liu et al., 2024a), each  
 126 using different saliency metrics.

127 Despite the strong empirical performance across various practical settings, most of the aforementioned  
 128 methods lack guarantees regarding their theoretical convergence rates, which can be attributed to the  
 129 biasedness of the projected gradients. To bridge this gap, we investigate the debiasing technique of  
 130 layerwise-sampling that compensates for the errors introduced by low-rank projected updates, aiming  
 131 to improve their theoretical convergence guarantees while maintaining practical memory efficiency.

132 **Unbiased Optimization Methods.** The research on unbiased methods is an important part of  
 133 the optimization field, especially for distributed and memory-efficient optimization. This includes  
 134 methods of unbiased quantization (Alistarh et al., 2017; Suresh et al., 2017; Wang et al., 2022) and  
 135 unbiased sparsification (Wangni et al., 2018; Stich et al., 2018; Wang et al., 2018). The unbiased  
 136 property of these methods enables low communication/memory burden while maintaining guaranteed  
 137 convergence. For the recently popular low-rank projection-based methods, Fira (Chen et al., 2024)  
 138 provides an attempt to involve full-rank information by adding a scaled gradient projection to the  
 139 update, but without a rigorous theoretical justification of the approach. GoLore (He et al., 2024)  
 140 is probably the closest to building an unbiased algorithm. However, they employ a totally random  
 141 projection matrix for the algorithm to enable the convergence guarantee, which may fail to capture  
 142 the loss landscape properties and lead to slow convergence.

142 **Muon Optimizer.** Muon (Jordan et al., 2024) is a novel optimizer proposed recently, which  
 143 is gaining rapidly increasing attention because of its great potential in training large foundation  
 144 models (Liu et al., 2025a; Kimi, 2025), empirically outperforming AdamW on specific large-scale  
 145 tasks. On the theoretical side, An et al. (2025); Li & Hong (2025) proves its non-convex deterministic  
 146 and stochastic convergence, respectively, showing a strong theoretical guarantee for the optimizer.

### 148 3 ALGORITHM

#### 151 3.1 GALORE UNBIASED WITH MUON

153 As previously shown in Algorithm 1, the core of low-rank gradient methods is to only store the  
 154 low-rank projected optimizer states, i.e., related to  $\tilde{G}_t \in \mathbb{R}^{m \times r}$ , which is then projected back to the  
 155 weight space by multiplying  $P_t$  to update the weight  $W_t$ . The update conceptually shares similarities  
 156 with running the base optimizer using low-rank projected gradients  $P_t P_t^\top G_t$  instead of  $G_t$ .

157 This inspires the key idea of *debiasing*, that is, to compensate for biased errors introduced by the  
 158 low-rank projection  $P_t P_t^\top G_t$ . To implement this while retaining memory efficiency, we refer to the  
 159 main idea of LISA (Pan et al., 2024), which allows some of the blocks to be sampled uniformly with  
 160 probability  $q$  in each period. This compensated full-rank updates use  $G_t - P_t P_t^\top G_t$ , while other  
 161 blocks still do the original low-rank update. By carefully balancing the scaling constants for the  
 162 two different updates, the biased low-rank term can be canceled out in expectation, resulting in an

162

**Algorithm 2** GaLore Unbiased with Muon (GUM)

163

---

```

1: Input:  $\{W_{0,\ell} \in \mathbb{R}^{m_\ell \times n_\ell}\}$  with each  $\ell$  corresponding to the  $\ell$ -th block of parameters, number of
2: blocks  $N_L$ , sampling period  $K$ , rank  $r$  for each layer, full-rank update layer number  $\gamma$ 
3: for  $t = 0$  to  $T/K - 1$  do
4:   for  $\ell = 1$  to  $N_L$  do
5:     Initialize  $R_{t,0,\ell} = 0$                                  $\triangleright$  Restart the momentum to clear memory
6:      $G_{t,0,\ell} = G(W_{t,\ell})$                              $\triangleright$  Obtain the gradient of the  $\ell$ -th layer at  $W_t$ 
7:      $U_{t,\ell}, S_{t,\ell}, V_{t,\ell} = \text{SVD}(G_{t,0,\ell})$        $\triangleright$  Compute SVD of gradient obtained at  $W_{tK}$ 
8:      $P_{t,\ell} = U_{t,\ell}[:, :r]$                              $\triangleright$  Obtain GaLore projector  $P_{t,\ell} \in \mathbb{R}^{m_\ell \times r}$  (suppose  $m_\ell \leq n_\ell$ )
9:   end for
10:  Each block  $\ell$  is sampled to do full-rank updates with probability  $q_{t,\ell} \equiv q = \frac{\gamma}{N_L}$ 
11:  for  $k = 0$  to  $K - 1$  do
12:    Run (1) for all blocks sampled to compute low-rank update
13:    Run (2) for all blocks sampled to compute full-rank update
14:  end for
15: end for

```

---

178

unbiased estimation of gradients across iterations. Due to page limit, we present this general unbiased algorithm paradigm in Algorithm 3 in Appendix A.

181

For a practical instance of this paradigm, we consider applying GaLore as the low-rank projection method and Muon as the base algorithm, which gives birth to our proposed optimization algorithm, called **GaLore Unbiased with Muon (GUM)**, as presented in Algorithm 2.

184

In one training process, the algorithm contains separated periods just like the vanilla GaLore and LISA. During each period  $t$ , each block of parameters is sampled to do full-rank updates with probability  $q_{t,\ell}$ . In each iteration  $k$  in the period, we first compute the projection matrix  $P_{t,k,\ell}$  and sample the layers to do full-rank updates in this period.

185

186

187

188

If block  $\ell$  is sampled to do the low-rank update, we apply the following update adapted from Muon with  $G_{t,k,\ell} = G(W_{tK+k,\ell})$  as the gradient of block  $\ell$  at iteration  $k$  in period  $t$ :

191

192

$$R_{t,k,\ell} = \beta R_{t,k-1,\ell} + \frac{1}{1 - q_{t,\ell}} P_{t,\ell}^\top G_{t,k,\ell} \quad (1)$$

193

$$W_{Kt+k+1,\ell} = W_{Kt+k,\ell} + \eta_{t,k} P_{t,\ell} \text{NewtonSchulz}(R_{t,k,\ell})$$

194

Note that if we set  $q_{t,\ell} = 0$ , (1) is exactly GaLore with Muon as the base optimizer, which we will refer to as GaLore-Muon. In terms of memory consumption, we can see that the optimizer states requiring storage are the projection matrix  $P_{t,\ell} \in \mathbb{R}^{m_\ell \times r}$  and  $R_{t,k,\ell} \in \mathbb{R}^{r \times n_\ell}$ . Otherwise, the block is sampled to compute high-rank updates, and the compensated projection update is applied.

195

196

197

$$R_{t,k,\ell} = \beta R_{t,k-1,\ell} + \frac{1}{q_{t,\ell}} (G_{t,k,\ell} - P_{t,\ell} P_{t,\ell}^\top G_{t,k,\ell}) \quad (2)$$

198

$$W_{tK+k+1,\ell} = W_{tK+k,\ell} + \eta_{t,k} \text{NewtonSchulz}(R_{t,k,\ell})$$

199

In this case,  $P_{t,\ell} \in \mathbb{R}^{m_\ell \times r}$  and  $R_{t,k,\ell} \in \mathbb{R}^{m_\ell \times n_\ell}$  are required to be stored.

200

201

202

Summarizing both cases, the overall memory consumption comparison with the vanilla GaLore-Muon algorithm is obtained, as shown in Table 3. The memory consumption of GUM is higher than that of GaLore when using the same projection rank  $r$ , due to the use of probabilistic full-rank updates. However, as demonstrated in Section 5, by employing a smaller projection rank  $r'$  as a trade-off, the benefits of this additional memory consumption are sufficient to recover the performance loss and even achieve a smaller overall memory footprint.

203

204

205

We can show that this update is unbiased compared to the original Muon update.

206

**Lemma 1** (GUM is unbiased). *A single iteration of Algorithm 2 for  $W \in \mathbb{R}^{m \times n}$  is equivalent to*

207

208

209

210

211

212

213

214

215

$$\tilde{M}^+ = \beta \tilde{M} + \tilde{G}$$

$$W^+ = W - \eta \text{NewtonSchulz}(\tilde{M}^+)$$

with  $\mathbb{E}[\tilde{G}] = G \in \mathbb{R}^{m \times n}$ , where  $G$  denotes the gradient obtained at  $W$ , and the expectation is taken over stochastic sampling and layerwise sampling random variable  $\zeta$ , where

216     •  $\zeta = 0$  when layer  $\ell$  does a low-rank update, with probability  $1 - q$ ,  
 217     •  $\zeta = 1$  when layer  $\ell$  does a full-rank update, with probability  $q$ .

218

219  
 220 This unbiased technique is crucial for the convergence of the algorithm. As we will see in the  
 221 next subsection, GUM can recover similar convergence properties as the original Muon algorithm,  
 222 regardless of the employed projection matrix. This demonstrates substantial theoretical advantages  
 223 over the original biased GaLore-Muon algorithm.

224

## 225     4 CONVERGENCE ANALYSIS OF GUM

226

227 In this section, we present the convergence analysis of GUM. We consider the following assumptions  
 228 for the minimization problem  $\min_{W \in \mathbb{R}^{m \times n}} f(W)$  with  $m \leq n$ .

229 **Assumption 1** (Lower bounded). *There exists  $f^* > -\infty$  such that  $f(W) \geq f^*$  for all  $W \in \mathbb{R}^{m \times n}$ .*

230 **Assumption 2** (Smoothness).  *$f$  is  $L_{\text{op}}$ -smooth with respect to the spectral norm  $\|\cdot\|_{\text{op}}$ , i.e.,*

232     
$$\|\nabla f(W_1) - \nabla f(W_2)\|_* \leq L_{\text{op}} \|W_1 - W_2\|_{\text{op}},$$

233

234 for all  $W_1, W_2 \in \mathbb{R}^{m \times n}$ .  $\|\cdot\|_{\text{op}}$  and  $\|\cdot\|_*$  denotes the spectral norm and trace norm respectively.

235 **Assumption 3** (Gradient noise). *We assume the stochastic gradient  $G(W)$  obtained at  $W$  is unbiased  
 236 and there exists a matrix  $V \in \mathbb{R}^{m \times n}$  such that*

237     
$$\mathbb{E}[N(W)] = 0 \quad \text{and} \quad \mathbb{E}[N(W)N(W)^\top] \preceq VV^\top,$$

238

239 where  $N(W) \triangleq G(W) - \nabla f(W)$  and  $A \preceq B$  denotes that  $B - A$  is positive semidefinite.

240 Assumption 1 is standard in non-convex analysis. Based on the equivalence between norms, Assumption  
 241 2 implies nothing more than the standard smoothness assumption on Frobenius norm, but is  
 242 more suitable in analyzing GUM or Muon (Jordan et al., 2024). Assumption 3 can also imply the  
 243 standard bounded variance assumption by  $\mathbb{E}[\|N(W)\|_{\text{F}}^2] \leq \|V\|_{\text{F}}^2$ . The style of these assumptions  
 244 can be found in previous work on analyzing adaptive methods and Sign-based methods (Bernstein  
 245 et al., 2018; Crawshaw et al., 2022; Liu et al., 2024b; An et al., 2025), where the assumptions are  
 246 employed for more fine-grained analysis and analyzing the potential benefits of these optimizers.

247 **Assumption 4** (Exact Newton Schulz). *We consider the case where the Newton-Schulz iteration  
 248 computes the exact solution, i.e.,  $\text{NewtonSchulz}(X) = UV^\top$  with  $X = U\Sigma V^\top$  as the SVD of  $X$ .*

249 Assumption 4 is needed for analyzing Muon. As noted in Jordan et al. (2024); Liu et al. (2025a),  
 250 though the Newton-Schulz iteration adopted in Muon does not compute the exact  $UV^\top$  matrix, it  
 251 turns out that this error has little influence on the training curve. Then, based on the assumptions, we  
 252 can obtain the convergence guarantee for GUM.

253 **Theorem 1** (Non-convex Convergence). *Under Assumption 1-4, after running a total of  $T$  iterations  
 254 for Algorithm 2 with parameters set as (12), it holds that*

255     
$$\min_{0 \leq s \leq T-1} \mathbb{E}[\|\nabla f(W_s)\|] \leq \mathcal{O} \left( \frac{1}{\alpha} \sqrt{\frac{L_{\text{op}} \Delta}{T}} + \left( \frac{L_{\text{op}} \Delta \|V\|_*^2}{\alpha^5 T} \right)^{\frac{1}{4}} + \frac{\|V\|_*}{\sqrt{\alpha^3 T}} \right),$$

256

257 where  $\Delta \triangleq f(W_0) - f^*$  and  $\alpha \triangleq \min\{q, 1 - q\}$ .

258 The proof can be found in Appendix B. The convergence theorem for GUM leads to several important  
 259 observations. Firstly, when we set  $q$  to be an absolute constant, the convergence of GUM matches  
 260 exactly the convergence rate of Muon. In the deterministic case, it matches the convergence result  
 261 of Muon proven in An et al. (2025). When the noise  $V$  is the dominant term, it also matches  
 262 the  $\mathcal{O}(T^{-1/4})$  rate proven in Li & Hong (2025). Moreover, since we use more fine-grained and  
 263 appropriate assumptions to analyze GUM, Theorem 1 shows an even better dimensional dependence  
 264 than Li & Hong (2025). This consistency shows the power of the unbiased design, maintaining the  
 265 memory reduction of gradient low-rank methods without sacrificing the convergence guarantee.

270 **Remark 1.** In Theorem 1, the optimal choice is  $q = 0.5$  for fast convergence. Although any positive  
 271 constant  $q \in (0, 1)$  ensures the proved theoretical convergence, in practical settings, a large constant  
 272  $q$  may lead to huge memory consumption. This indicates a fundamental tradeoff between time and  
 273 space, which is controlled by the choice of  $q$ . If the memory requirement is constant, it is preferable  
 274 to choose the largest affordable  $q \leq 0.5$ . If the memory support is dynamic, e.g., the number of nodes  
 275 decreases at a later stage of training, it is better to schedule a diminishing  $q$  to adapt to this dynamic  
 276 memory.

277 As noted by He et al. (2024), GaLore using SGD with momentum (SGDM) as the base algorithm  
 278 converges in the deterministic non-convex setting, but can possibly diverge when the gradient noise  
 279 is large. We also empirically examine an extreme counterexample where GaLore-Muon doesn't  
 280 converge at all in Section 5. Clearly, GUM fixes this problem. GoLore (He et al., 2024) is also  
 281 designed to correct the convergence of GaLore. However, though GoLore shows a good convergence  
 282 guarantee when the base algorithm is SGDM, it employs a thoroughly random projection matrix  
 283 to do low-rank updates, failing to capture the potential gradient low-rank properties as the GaLore  
 284 projection matrix does. This can lead to a much slower convergence speed when applied to real  
 285 training tasks.

## 287 5 EXPERIMENTAL RESULTS

### 289 5.1 SYNTHETIC SETTINGS

291 To better illustrate how GaLore may fail due to the low-rank projection, we consider the following  
 292 synthetic noisy problem.

293 **Setup.** The settings of the experiment are generally the same as the synthetic experiment in He et al.  
 294 (2024). We consider the following noisy linear regression problem.

$$296 \min_{X \in \mathbb{R}^{n \times n}} f(x) \triangleq \frac{1}{2} \|AX\|_F^2 + \langle B, X \rangle, \quad \nabla f(X; \xi) = \nabla f(X) + \xi \sigma C,$$

298 where  $A = [I_{n-r} \ 0] \in \mathbb{R}^{(n-r) \times n}$ ,  $B = \begin{bmatrix} D & 0 \\ 0 & 0 \end{bmatrix} \in \mathbb{R}^{n \times n}$  with  $D \in \mathbb{R}^{(n-r) \times (n-r)}$  a Gaussian  
 299 random matrix,  $C = \begin{bmatrix} 0 & 0 \\ 0 & I_r \end{bmatrix} \in \mathbb{R}^{n \times n}$ ,  $\xi$  is a random variable with probability 0.5 to be 1 and  
 300 probability 0.5 to be 0, and  $\sigma$  is a constant controlling the noise level. It is straightforward to verify  
 301 that this is a smooth and convex optimization problem, with bounded gradient variance.

304 In our experiment, we specifically set  $n = 20$ ,  $r = 12$ ,  $\sigma = 100$  to construct a small-scale but noisy  
 305 problem. For the vanilla (biased) GaLore Muon algorithm, we set the projection rank to be 12 as  
 306 well. For GUM, we set  $r = 2$  and  $q_{t,\ell} = 0.5$ . We can see that in this case, the memory footprints of  
 307 the two algorithms are the same.

308 **Results.** The convergence result is shown in Figure 1. We adjust the minimum loss to 0 to better  
 309 visualize the difference. As we can see, GaLore  
 310 fails to converge at all, while GUM converges to  
 311 a comparable accuracy with the full-parameter  
 312 Muon baseline. The experiment shows a clear  
 313 benefit of the unbiased method, at least in noisy  
 314 settings.

316 Here is a more detailed analysis of why these  
 317 conditions lead to GaLore's failure. In this syn-  
 318 synthetic problem, the noise level is set to be large  
 319 and has rank  $r = 12$ , which is equal to the pro-  
 320 jection rank of GaLore. Since the noise is in a  
 321 dominant position, every time the  $r$  largest sin-  
 322 gular values of the stochastic gradient  $\nabla f(X; \xi)$  come from the noise, so do the corresponding  
 323 singular vectors and the GaLore projection matrix. This meaningless projection makes the training  
 process not even take a single effective step towards solving the problem. Therefore, this synthetic

315 Table 1: **Space complexity** comparison between  
 316 GaLore and GUM for a block  $W \in \mathbb{R}^{m \times m}$  with  
 317  $r' < r \leq m$  respectively. GUM uses a full-  
 318 rank update with probability  $q \in [0, 1]$ , where the  
 319 memory GUM has the same memory consumption  
 320 when  $q = 2(r - r')/(m - r')$ .

Method	Space Complexity
GaLore	$\mathcal{O}(2mr)$
GUM	$\mathcal{O}((2 - q)mr' + qm^2)$
SFT	$\mathcal{O}(m^2)$

324 experiment shows an extreme case in which GaLore can fail when the gradient noise is large. Also,  
 325 the experiment shows that GUM fixes the non-convergence problem with the same memory cost as  
 326 GaLore-Muon.

## 329 5.2 LLM FINE-TUNING SETTINGS

330 To verify the empirical effectiveness of the proposed algorithm in practice, we compare GUM with  
 331 GaLore in LLM fine-tuning settings.

332 **Setup.** The performance of the fine-tuned models is evaluated on two types of tasks: 1) IFEval  
 333 (Zhou et al., 2023), an instruction-following benchmark that assesses models’ adherence to explicit,  
 334 verifiable instructions, and 2) GSM8K (Cobbe et al., 2021a), a mathematical reasoning benchmark  
 335 that evaluates models’ problem-solving skills in grad-school level math questions.

336 For model choices, LLaMA3-8B (Grattafiori et al., 2024a), Qwen2.5-7B (Qwen et al., 2025), and  
 337 Gemma2-9B (Team et al., 2024) are adopted, which are commonly used in practical applications.

338 For training datasets, GPT-4-LLM is adopted on the instruction-following tasks of IFEval, which  
 339 consists of 54.6K high-quality GPT-4-generated instruction-response pairs across various instruction  
 340 categories. As for the mathematical reasoning task of GSM8K, a 2K-sized high-quality mixture<sup>1</sup> from  
 341 DART-Math (Tong et al., 2024), Ultra-Interact (Yuan et al., 2024), MathInstruct (Yue et al., 2023),  
 342 and Orca-Math (Mitra et al., 2024) is employed, which allows strong models such as Qwen-2.5-7B to  
 343 still obtain reasonable improvements after fine-tuning.

344 For hyperparameters, we adopt a rank of 512 for GaLore and 2 + 128 for GUM. The baselines include  
 345 Full-parameter Training with Muon (Jordan et al., 2024) (FT-Muon), Full-parameter Training with  
 346 AdamW (Loshchilov & Hutter, 2019) (FT-AdamW), Gradient Low-Rank Projection (GaLore) (Zhao  
 347 et al., 2024), [Fira](#) (Chen et al., 2024), [Golore](#) (He et al., 2024), [LDADAM](#) (Robert et al., 2025), and  
 348 [Apollo](#) (Zhu et al., 2025), where further details are available in Appendix C.

349 **Memory Efficiency.** We conducted peak GPU  
 350 memory experiments to evaluate GUM’s mem-  
 351 ory efficiency, demonstrating its comparable or  
 352 reduced memory footprint relative to GaLore.  
 353 Specifically, we focus on two key hyperparam-  
 354 eters: the rank and the number of selected lay-  
 355 ers for full-rank updates in GUM. To ensure a  
 356 fair comparison, all methods used a consistent  
 357 mini-batch size of 1, without employing addi-  
 358 tional GPU memory-saving techniques such as  
 359 offloading (Ren et al., 2021) or flash attention  
 360 (Dao et al., 2022; Dao, 2024).

361 As shown in Table 3, the GUM configuration  
 362 reaches comparable or better memory consump-  
 363 tion than GaLore. This improvement is not lim-  
 364 ited to a single case; consistent memory savings  
 365 are observed across multiple model architec-  
 366 tures.

367 **Results.** As shown in Table 2, GUM consistently outperforms GaLore in both tasks, highlighting its  
 368 robustness and general effectiveness.

369 A closer look at GSM8K results reveals that GUM achieves notable improvements and even outper-  
 370 forms full-parameter training methods, suggesting its strength in enhancing reasoning capabilities.  
 371 In Section 5.4, it will be revealed that this improvement is very likely to have originated from its  
 372 improvements in memorization, especially when the learned activations are required to be long-tailed.

373 **Table 3: Peak GPU memory usage** across differ-  
 374 ent model architectures and configurations, empha-  
 375 sizing the variations among them. As specified in  
 376 the table, the GUM configuration 2 + 128 involves  
 377 updating two layers with full-rank gradients, while  
 378 all other layers are updated with low-rank gradi-  
 379 ents of rank  $r = 128$ .

Model	GaLore	GUM Layers + Rank	
	512	4 + 128	2 + 128
LLaMA-3-8B	42G	41G	<b>40G</b>
Qwen-2.5-7B	41G	40G	<b>39G</b>
Gemma-2-9B	47G	46G	<b>44G</b>

375  
 376  
 377 <sup>1</sup>The dataset is from [https://huggingface.co/datasets/HanningZhang/scalebio\\_distill\\_qwen\\_math](https://huggingface.co/datasets/HanningZhang/scalebio_distill_qwen_math), generated using the same setting as Appendix A.2 of (Pan et al., 2025).

378  
 379 **Table 2: LLM Fine-tuning Results.** Trained models are evaluated on IFEval (instruction-following)  
 380 and GSM8K (mathematical reasoning). All experiments are conducted on a single H100 GPU.

381 382 383 384 385 386 387 388 389 390 391 392 393 394 395 396 397 398 399 400 401 402 403 404	382 383 384 385 386 387 388 389 390 391 392 393 394 395 396 397 398 399 400 401 402 403 404	382 383 384 385 386 387 388 389 390 391 392 393 394 395 396 397 398 399 400 401 402 403 404	382 383 384 385 386 387 388 389 390 391 392 393 394 395 396 397 398 399 400 401 402 403 404	382 383 384 385 386 387 388 389 390 391 392 393 394 395 396 397 398 399 400 401 402 403 404		382 383 384 385 386 387 388 389 390 391 392 393 394 395 396 397 398 399 400 401 402 403 404	
				Prompt-level Strict-Accuracy	Prompt-level Loose-Accuracy		
382 383 384 385 386 387 388 389 390 391	382 383 384 385 386 387 388 389 390 391	382 383 384 385 386 387 388 389 390 391	382 383 384 385 386 387 388 389 390 391	FT-AdamW	23.66	25.14	57.39
				FT-Muon	23.11	26.06	57.65
				Apollo (Zhu et al., 2025)	19.04	21.63	56.03
				GaLore (Zhao et al., 2024)	21.07	22.74	57.38
				Fira (Chen et al., 2024)	21.81	23.73	56.41
				LDAdam (Robert et al., 2025)	22.74	24.40	57.92
				GoLore (He et al., 2024)	23.01	24.95	57.54
				GUM	22.37	24.03	58.45
				FT-AdamW	35.12	39.74	85.75
				FT-Muon	34.38	39.19	85.90
392 393 394 395 396 397 398 399 400 401	392 393 394 395 396 397 398 399 400 401	392 393 394 395 396 397 398 399 400 401	392 393 394 395 396 397 398 399 400 401	Apollo	31.61	36.41	85.67
				GaLore	33.09	37.71	86.28
				Fira	32.35	36.04	86.81
				LDAdam	28.10	30.31	83.78
				GoLore	30.87	35.67	86.66
				GUM	33.46	38.82	86.81
				FT-AdamW	OOM	OOM	OOM
				FT-Muon	28.47	32.16	76.92
				Apollo	25.14	28.10	75.28
				GaLore	30.31	33.64	77.18
402 403 404	402 403 404	402 403 404	402 403 404	Fira	29.21	33.64	75.44
				LDAdam	28.84	32.53	75.13
				GoLore	31.05	34.38	74.98
				GUM	33.27	36.60	77.48

### 5.3 LLM PRE-TRAINING SETTINGS

To provide stronger evidence for validating the effectiveness of GUM, a standard pre-training setting is introduced to compare different training methods’ performance.

**Setup.** To evaluate the improvements in commonsense reasoning, the following downstream tasks are employed: ARC (Clark et al., 2018), OpenBookQA (Mihaylov et al., 2018b), HellaSwag Zellers et al. (2019a), PIQA (Bisk et al., 2020), SIQA (Sap et al., 2019), and WinoGrande (Sakaguchi et al., 2021a), which are common choices for LLM pre-training (Hoffmann et al., 2022a; Groeneveld et al., 2024; Zhang et al., 2024b). For model choice, following the standard setting in Zhao et al. (2024), the experiments covered three model sizes—60M, 130M, and 350M parameters of LLaMA. For training datasets, we employ the widely-used C4 corpus (Raffel et al., 2023) under configurations guided by the Chinchilla scaling law (Hoffmann et al., 2022b): 1.5B tokens for 60M, 2B tokens for 130M, and 7B tokens for 350M. For baselines, in addition to Galore and full-parameter training methods, we include Fira (Chen et al., 2024) and Subtrack++ (He et al., 2024). Further details are available in Appendix C.3.

**Results.** The performance comparison presented in Table 4 clearly indicates that GUM achieves consistently better results than GaLore and, more surprisingly, even full-parameter training methods like AdamW and Muon. This improvement can largely be attributed to the unbiased low-rank update mechanism employed in GUM. The mechanism captures long-tailed gradient updates distributed across layers and thereby enhances model memorization.

### 5.4 UNDERSTANDING THE EFFECT OF LAYERWISE SAMPLING

In this section, we investigate the underlying reason why the proposed algorithm of GUM can yield empirical improvements over GaLore. In short, GUM’s high-rank gradient update leads to a more uniform singular value distribution in model parameters, which further results in more evenly

432  
 433 **Table 4: LLM Pre-training Results.** Trained models are evaluated on seven widely adopted  
 434 commonsense reasoning tasks. All experiments are conducted on H100 GPUs.

435 Model	436 Method	437 ARC-E	438 ARC-C	439 OBQA	440 HellaSwag	441 PIQA	442 SIQA	443 Winogrande	444 Avg.
445 LLaMA-60M	FT-AdamW	32.87	17.92	12.68	26.70	58.87	35.88	50.12	33.58
	FT-Muon	36.45	17.92	12.88	26.89	59.79	35.82	51.22	34.42
	GaLore	35.35	17.92	12.47	26.74	59.63	35.62	49.88	33.94
	Fira	35.02	18.94	12.27	26.75	58.71	36.24	50.28	34.03
	Subtrack++	37.96	16.92	13.08	27.07	60.45	37.05	51.67	34.89
	GUM	36.28	17.41	13.68	26.70	60.12	36.54	51.85	34.65
446 LLaMA-130M	FT-AdamW	37.08	18.86	13.48	27.04	59.14	36.18	51.07	34.69
	FT-Muon	38.34	18.00	13.08	27.67	62.68	37.00	49.33	35.16
	GaLore	36.49	18.00	13.28	27.08	60.34	35.36	50.20	34.39
	Fira	26.01	19.54	12.27	26.13	53.65	34.19	49.80	31.66
	Subtrack++	36.49	17.58	14.08	26.92	61.70	36.08	52.09	34.99
	GUM	38.01	18.34	14.69	27.32	61.26	36.44	52.49	35.51
447 LLaMA-350M	FT-AdamW	44.02	18.77	14.08	30.04	64.42	37.97	50.51	37.12
	FT-Muon	44.91	18.69	17.10	31.05	65.72	37.87	51.93	38.18
	GaLore	43.10	18.52	14.89	29.09	62.19	37.10	52.01	36.58
	Fira	42.38	18.77	15.49	29.27	63.00	37.97	51.85	36.96
	Subtrack++	40.45	18.43	14.49	28.50	63.06	37.72	50.25	36.13
	GUM	44.44	19.80	15.69	29.28	64.53	38.13	51.38	37.42

451 distributed activations for input samples. This implies the long-tailed knowledge is better preserved  
 452 in GUM-trained models, yielding better memorization.

453 **Setup.** We adopt the model of LLaMA-130M and benchmark of ARC-E (Clark et al., 2018), while  
 454 keeping other settings the same as in Section 5.3.

455 **Results.** As shown in Figure 2, the overall stable ranks  $\mathbb{E} [\|M\|_F^2 / \|M\|_2^2]$  of GaLore and GUM are  
 456 positively correlated with their performance in ARC-E, which provides direct evidence that higher  
 457 stable ranks are generally beneficial for improving commonsense reasoning.

458 On top of that, it is observed in Figure 3 that GUM not only improves the overall stable rank  
 459 of the trained model, but also shapes a set of more evenly distributed singular values in trained  
 460 models, which further leads to more long-tail distributed activation across all modules. This provides  
 461 indirect evidence and an intuitive explanation for the performance improvements in ARC-E: instead of  
 462 overusing a low-dimensional space or a limited number of modules, GUM-trained models demonstrate  
 463 a tendency to evenly distribute knowledge across all dimensions and modules, implying better  
 464 memorization. Additional evidence is available in Appendix D.2.

## 468 6 CONCLUSIONS

471 In this paper, we investigate the debiasing technique of layerwise sampling for memory-efficient  
 472 LLM training, whose combination with GaLore restores the theoretical convergence properties of  
 473 full-parameter training. Our proposed algorithm, GUM, demonstrates that it is possible to achieve  
 474 provable convergence in low-rank optimization without impairing its empirical performance and  
 475 memory efficiency. Further analysis shows that the empirical gains are brought by the inherent high-  
 476 rank updates, which lead to a higher overall stable rank and more uniformly distributed singular values  
 477 in model parameters, yielding more long-tailed activation patterns and implying better memorization.

## 478 480 ETHICS STATEMENT

481 After carefully reviewing the ethical regulations of the conference, to the best of our knowledge, this  
 482 work does not present any foreseeable ethical concerns. No negative societal or ethical impacts are  
 483 anticipated for the contribution of this work. The proposed algorithms are for general large language  
 484 model training, and do not involve anything about human subjects, potentially harmful insights,  
 485 potential conflicts of interest and sponsorship, discrimination/bias/fairness concerns, privacy and  
 security issues, legal compliance, or research integrity issues.

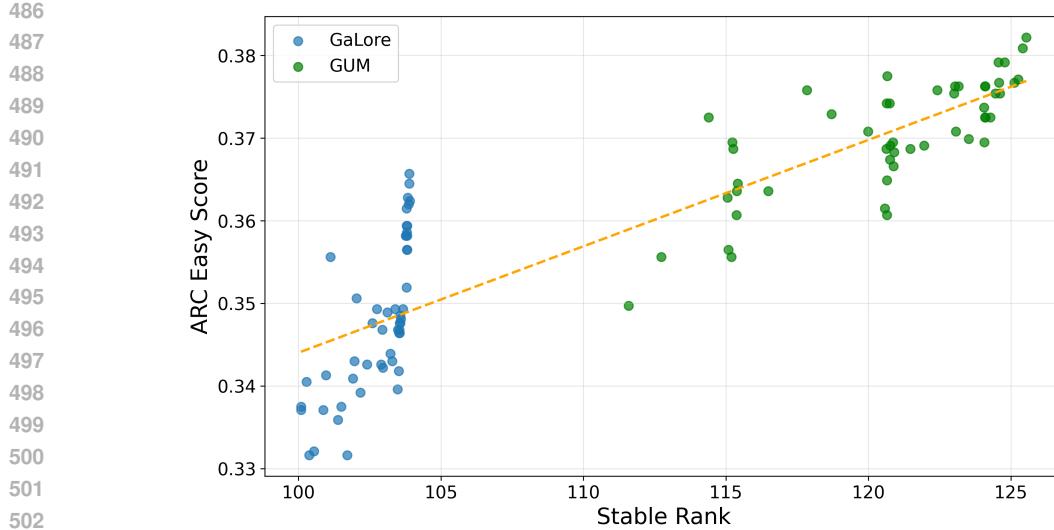


Figure 2: **Higher Stable Rank  $\rightarrow$  Better Performance.** A positive correlation is observed between the overall stable rank  $\mathbb{E} [\|M\|_F^2 / \|M\|_2^2]$  and ARC Easy score. Each dot represents a checkpoint during pre-training after 1,000 steps, saved every 20 steps.

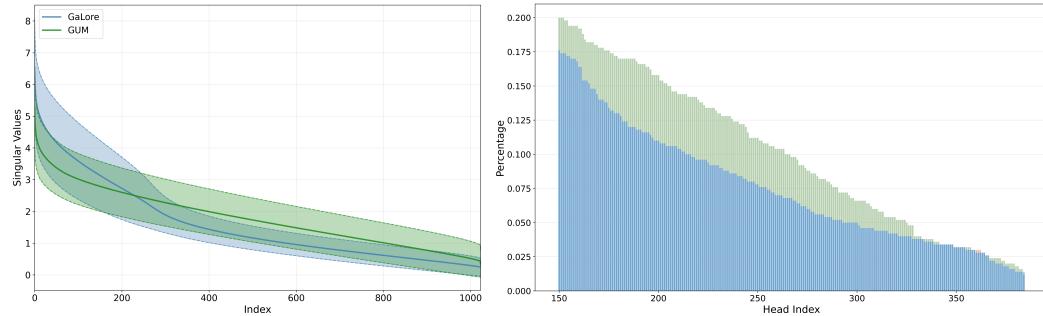


Figure 3: **Left: Updates  $\rightarrow$  Weights:** Singular value distribution across layers of GaLore and GUM, where GUM demonstrates a more even and long-tailed distribution of singular values. **Right: Weights  $\rightarrow$  Activations:** Tail distribution of modules that contain salient activations, where salient activations are defined as activations with top-k ( $k = 10,000$ ) attention scores over all modules. Randomly sampled 1K inputs from the C4 corpus are utilized as prompts. Blue parts correspond to GaLore’s tail distribution, while green parts stand for GUM’s further increase on top of GaLore.

## REPRODUCIBILITY STATEMENT

We have made efforts to ensure that our work is reproducible, with details provided in Section 5 and Appendix C.

## REFERENCES

Dan Alistarh, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. Qsgd: Communication-efficient sgd via gradient quantization and encoding. *Advances in neural information processing systems*, 30, 2017.

Kang An, Yuxing Liu, Rui Pan, Shiqian Ma, Donald Goldfarb, and Tong Zhang. Asgo: Adaptive structured gradient optimization. *arXiv preprint arXiv:2503.20762*, 2025.

Alan Ansell, Edoardo Maria Ponti, Anna Korhonen, and Ivan Vulić. Composable sparse fine-tuning for cross-lingual transfer. *arXiv preprint arXiv:2110.07560*, 2021.

540 Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization.  
 541 *arXiv preprint arXiv:1907.02893*, 2019.

542

543 Jeremy Bernstein, Yu-Xiang Wang, Kamyar Azizzadenesheli, and Animashree Anandkumar. signs gd:  
 544 Compressed optimisation for non-convex problems. In *International Conference on Machine  
 545 Learning*, pp. 560–569. PMLR, 2018.

546 Yonatan Bisk, Rowan Zellers, Omer Yakhini, and Yejin Choi. Piqa: Reasoning about physical  
 547 commonsense in natural language. *arXiv preprint*, 2020. URL <https://arxiv.org/abs/1911.11641>.

548

549 Xi Chen, Kaituo Feng, Changsheng Li, Xunhao Lai, Xiangyu Yue, Ye Yuan, and Guoren Wang.  
 550 Fira: Can we achieve full-rank training of llms under low-rank constraint?, 2024. URL <https://arxiv.org/abs/2410.01623>.

551

552 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick,  
 553 and Oyvind Tafjord. Think you have solved question answering? try ARC, the AI2 reasoning  
 554 challenge. In *Proceedings of the 2018 Conference of the North American Chapter of the Association  
 555 for Computational Linguistics: Human Language Technologies*, New Orleans, Louisiana, May  
 556 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-1074. URL <https://arxiv.org/abs/1803.05457>.

557

558

559 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,  
 560 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John  
 561 Schulman. Training verifiers to solve math word problems, 2021a.

562

563 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,  
 564 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John  
 565 Schulman. Training verifiers to solve math word problems, 2021b. URL <https://arxiv.org/abs/2110.14168>.

566

567 Michael Crawshaw, Mingrui Liu, Francesco Orabona, Wei Zhang, and Zhenxun Zhuang. Robustness  
 568 to unbounded smoothness of generalized signs gd. *Advances in neural information processing  
 569 systems*, 35:9955–9968, 2022.

570

571 Ashok Cutkosky and Harsh Mehta. Momentum improves normalized sgd. In *International conference  
 572 on machine learning*, pp. 2260–2268. PMLR, 2020.

573

574 Tri Dao. Flashattention-2: Faster attention with better parallelism and work partitioning. In  
 575 *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=mZn2Xyh9Ec>.

576

577 Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. FlashAttention: Fast and  
 578 memory-efficient exact attention with IO-awareness. In *Advances in Neural Information Processing  
 579 Systems*, 2022.

580

581 Sarkar Snigdha Sarathi Das, Ranran Haoran Zhang, Peng Shi, Wenpeng Yin, and Rui Zhang. Unified  
 582 low-resource sequence labeling by sample-aware dynamic sparse finetuning. *arXiv preprint  
 583 arXiv:2311.03748*, 2023.

584

585 Aaron Defazio, Xingyu Yang, Ahmed Khaled, Konstantin Mishchenko, Harsh Mehta, and Ashok  
 586 Cutkosky. The road less scheduled. *Advances in Neural Information Processing Systems*, 37:  
 587 9974–10007, 2024.

588

589 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of  
 590 deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and  
 591 Thamar Solorio (eds.), *Proceedings of the 2019 Conference of the North American Chapter of the  
 592 Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and  
 593 Short Papers)*, pp. 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational  
 594 Linguistics. doi: 10.18653/v1/N19-1423. URL <https://aclanthology.org/N19-1423>.

595

596 Shizhe Diao, Rui Pan, Hanze Dong, Ka Shun Shum, Jipeng Zhang, Wei Xiong, and Tong Zhang.  
 597 Lmflow: An extensible toolkit for finetuning and inference of large foundation models. *arXiv  
 598 preprint arXiv:2306.12420*, 2023.

594 Ning Ding, Yujia Qin, Guang Yang, Fuchao Wei, Zonghan Yang, Yusheng Su, Shengding Hu, Yulin  
 595 Chen, Chi-Min Chan, Weize Chen, et al. Delta tuning: A comprehensive study of parameter  
 596 efficient methods for pre-trained language models. *arXiv preprint arXiv:2203.06904*, 2022.

597

598 Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-aware minimization  
 599 for efficiently improving generalization. *arXiv preprint arXiv:2010.01412*, 2020.

600 Rong Ge, Sham M Kakade, Rahul Kidambi, and Praneeth Netrapalli. The step decay schedule: A  
 601 near optimal, geometrically decaying learning rate procedure for least squares. *Advances in neural*  
 602 *information processing systems*, 32, 2019a.

603

604 Rong Ge, Zhize Li, Weiyao Wang, and Xiang Wang. Stabilized svrg: Simple variance reduction for  
 605 nonconvex optimization. In *Conference on learning theory*, pp. 1394–1448. PMLR, 2019b.

606 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad  
 607 Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela  
 608 Fan, Anirudh Goyal, Anthony Hartshorn, and et al. The llama 3 herd of models, 2024a. URL  
 609 <https://arxiv.org/abs/2407.21783>.

610

611 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad  
 612 Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of  
 613 models. *arXiv preprint arXiv:2407.21783*, 2024b.

614 Dirk Groeneveld, Iz Beltagy, Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord,  
 615 Ananya Harsh Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, et al. Olmo: Accelerating  
 616 the science of language models. *arXiv preprint arXiv:2402.00838*, 2024.

617

618 Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv*  
 619 *preprint arXiv:2312.00752*, 2023.

620

621 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,  
 622 Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms  
 623 via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.

624

625 Demi Guo, Alexander M Rush, and Yoon Kim. Parameter-efficient transfer learning with diff pruning.  
 626 *arXiv preprint arXiv:2012.07463*, 2020.

627

628 Yifan Hao, Xingyuan Pan, Hanning Zhang, Chenlu Ye, Rui Pan, and Tong Zhang. Understanding  
 629 overadaptation in supervised fine-tuning: The role of ensemble methods. *arXiv preprint*  
 630 *arXiv:2506.01901*, 2025.

631

632 Yutong He, Pengrui Li, Yipeng Hu, Chuyan Chen, and Kun Yuan. Subspace optimization for large  
 633 language models with convergence guarantees. *arXiv preprint arXiv:2410.11289*, 2024.

634

635 Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza  
 636 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al.  
 637 Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022a.

638

639 Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza  
 640 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom  
 641 Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy,  
 642 Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre.  
 643 Training compute-optimal large language models, 2022b. URL <https://arxiv.org/abs/2203.15556>.

644

645 Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,  
 646 and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International*  
 647 *Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=nZeVKeFYf9>.

648

649 Zijian Hu, Jipeng Zhang, Rui Pan, Zhaozhuo Xu, Shanshan Han, Han Jin, Alay Diliphai Shah,  
 650 Dimitris Stripelis, Yuhang Yao, Salman Avestimehr, et al. Fox-1: Open small language model for  
 651 cloud and edge. *arXiv preprint arXiv:2411.05281*, 2024.

648 Jia-Hong Huang, Yixian Shen, Hongyi Zhu, Stevan Rudinac, and Evangelos Kanoulas. Gradient  
 649 weight-normalized low-rank projection for efficient llm training. In *Proceedings of the AAAI  
 650 Conference on Artificial Intelligence*, volume 39, pp. 24123–24131, 2025.

651 Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive variance  
 652 reduction. *Advances in neural information processing systems*, 26, 2013.

653 Keller Jordan, Yuchen Jin, Vlado Boza, Jiacheng You, Franz Cesista, Laker Newhouse, and Jeremy  
 654 Bernstein. Muon: An optimizer for hidden layers in neural networks, 2024. URL <https://kellerjordan.github.io/posts/muon/>.

655 Team Kimi. Kimi-vl technical report. *arXiv preprint arXiv:2504.07491*, 2025.

656 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint  
 657 arXiv:1412.6980*, 2014.

658 Jeffrey Li, Alex Fang, Georgios Smyrnis, Maor Ivgi, Matt Jordan, Samir Yitzhak Gadre, Hritik  
 659 Bansal, Etash Guha, Sedrick Scott Keh, Kushal Arora, et al. Datacomp-lm: In search of the  
 660 next generation of training sets for language models. *Advances in Neural Information Processing  
 661 Systems*, 37:14200–14282, 2024.

662 Jiaxiang Li and Mingyi Hong. A note on the convergence of muon and further. *arXiv preprint  
 663 arXiv:2502.02900*, 2025.

664 Vladislav Lialin, Namrata Shivagunde, Sherin Muckatira, and Anna Rumshisky. Relora: High-rank  
 665 training through low-rank updates, 2023.

666 Haotian Liu, Chunyuan Li, Qingsyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in  
 667 neural information processing systems*, 36:34892–34916, 2023.

668 Jingyuan Liu, Jianlin Su, Xingcheng Yao, Zhejun Jiang, Guokun Lai, Yulun Du, Yidao Qin, Weixin  
 669 Xu, Enzhe Lu, Junjie Yan, et al. Muon is scalable for llm training. *arXiv preprint arXiv:2502.16982*,  
 670 2025a.

671 Xinxin Liu, Aaron Thomas, Cheng Zhang, Jianyi Cheng, Yiren Zhao, and Xitong Gao. Refining  
 672 salience-aware sparse fine-tuning strategies for language models. *arXiv preprint arXiv:2412.13488*,  
 673 2024a.

674 Yuxing Liu, Rui Pan, and Tong Zhang. Adagrad under anisotropic smoothness. *arXiv preprint  
 675 arXiv:2406.15244*, 2024b.

676 Yuxing Liu, Yuze Ge, Rui Pan, An Kang, and Tong Zhang. Theoretical analysis on how learning rate  
 677 warmup accelerates convergence. *arXiv preprint arXiv:2509.07972*, 2025b.

678 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019. URL <https://arxiv.org/abs/1711.05101>.

679 Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct  
 680 electricity? a new dataset for open book question answering. In *EMNLP*, 2018a.

681 Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct  
 682 electricity? a new dataset for open book question answering. In *Proceedings of the 2018 Conference  
 683 on Empirical Methods in Natural Language Processing*, pp. 2381–2391, Brussels, Belgium,  
 684 October–November 2018b. Association for Computational Linguistics. doi: 10.18653/v1/D18-1260.  
 685 URL <https://aclanthology.org/D18-1260/>.

686 Arindam Mitra, Hamed Khanpour, Corby Rosset, and Ahmed Awadallah. Orca-math: Unlocking the  
 687 potential of slms in grade school math, 2024.

688 Siyuan Mu and Sen Lin. A comprehensive survey of mixture-of-experts: Algorithms, theory, and  
 689 applications. *arXiv preprint arXiv:2503.07137*, 2025.

690 Aashiq Muhamed, Oscar Li, David Woodruff, Mona Diab, and Virginia Smith. Grass: Compute effi-  
 691 cient low-memory llm training with structured sparse gradients. *arXiv preprint arXiv:2406.17660*,  
 692 2024.

693

702 Deanna Needell, Nathan Srebro, and Rachel Ward. Stochastic gradient descent, weighted sampling,  
 703 and the randomized kaczmarz algorithm. *Advances in neural information processing systems*, 27,  
 704 2014.

705 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong  
 706 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton,  
 707 Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan Leike, and  
 708 Ryan Lowe. Training language models to follow instructions with human feedback. In S. Koyejo,  
 709 S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information  
 710 Processing Systems*, volume 35, pp. 27730–27744. Curran Associates, Inc., 2022.

711 Rui Pan, Haishan Ye, and Tong Zhang. Eigencurve: Optimal learning rate schedule for sgd on  
 712 quadratic objectives with skewed hessian spectrums. *arXiv preprint arXiv:2110.14109*, 2021.

713 Rui Pan, Shizhe Diao, Jianlin Chen, and Tong Zhang. Extremebert: A toolkit for accelerating  
 714 pretraining of customized bert. *arXiv preprint arXiv:2211.17201*, 2022.

715 Rui Pan, Yuxing Liu, Xiaoyu Wang, and Tong Zhang. Accelerated convergence of stochastic heavy  
 716 ball method under anisotropic gradient noise. *arXiv preprint arXiv:2312.14567*, 2023.

717 Rui Pan, Xiang Liu, Shizhe Diao, Renjie Pi, Jipeng Zhang, Chi Han, and Tong Zhang. Lisa: layerwise  
 718 importance sampling for memory-efficient large language model fine-tuning. *Advances in Neural  
 719 Information Processing Systems*, 37:57018–57049, 2024.

720 Rui Pan, Dylan Zhang, Hanning Zhang, Xingyuan Pan, Minrui Xu, Jipeng Zhang, Renjie Pi, Xiaoyu  
 721 Wang, and Tong Zhang. ScaleBiO: Scalable bilevel optimization for LLM data reweighting. In  
 722 *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics*, pp.  
 723 31959–31982, 2025.

724 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of  
 725 the IEEE/CVF international conference on computer vision*, pp. 4195–4205, 2023.

726 Yujia Qin, Yining Ye, Junjie Fang, Haoming Wang, Shihao Liang, Shizuo Tian, Junda Zhang, Jiahao  
 727 Li, Yunxin Li, Shijue Huang, et al. Ui-tars: Pioneering automated gui interaction with native  
 728 agents. *arXiv preprint arXiv:2501.12326*, 2025.

729 Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan  
 730 Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang,  
 731 Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin  
 732 Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi  
 733 Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan,  
 734 Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report, 2025. URL  
 735 <https://arxiv.org/abs/2412.15115>.

736 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi  
 737 Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text  
 738 transformer, 2023. URL <https://arxiv.org/abs/1910.10683>.

739 Sahar Rajabi, Nayeema Nonta, and Sirisha Rambhatla. Subtrack++ : Gradient subspace tracking for  
 740 scalable LLM training. In *The Thirty-ninth Annual Conference on Neural Information Processing  
 741 Systems*, 2025. URL <https://openreview.net/forum?id=6geRIdlFWJ>.

742 Amrutha Varshini Ramesh, Vignesh Ganapathiraman, Issam H Laradji, and Mark Schmidt. Blockllm:  
 743 Memory-efficient adaptation of llms by selecting and optimizing the right coordinate blocks. *arXiv  
 744 preprint arXiv:2406.17296*, 2024.

745 Jie Ren, Samyam Rajbhandari, Reza Yazdani Aminabadi, Olatunji Ruwase, Shuangyan Yang, Minjia  
 746 Zhang, Dong Li, and Yuxiong He. Zero-offload: Democratizing billion-scale model training, 2021.  
 747 URL <https://arxiv.org/abs/2101.06840>.

748 Thomas Robert, Mher Safaryan, Ionut-Vlad Modoranu, and Dan Alistarh. LDAdam: Adaptive  
 749 optimization from low-dimensional gradient statistics. In *The Thirteenth International Conference  
 750 on Learning Representations*, 2025. URL <https://openreview.net/forum?id=Zkp1GuHerF>.

756 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
 757 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-  
 758 ence on computer vision and pattern recognition*, pp. 10684–10695, 2022.

759

760 Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An  
 761 adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106,  
 762 September 2021a. doi: 10.1145/3474381.

763 Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An  
 764 adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106,  
 765 2021b.

766

767 Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. Social IQa: Com-  
 768 monsense reasoning about social interactions. In *Proceedings of the 2019 Conference on  
 769 Empirical Methods in Natural Language Processing and the 9th International Joint Confer-  
 770 ence on Natural Language Processing (EMNLP-IJCNLP)*, pp. 4463–4473, Hong Kong, China,  
 771 November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1454. URL  
 772 <https://aclanthology.org/D19-1454/>.

773

774 Sebastian U Stich, Jean-Baptiste Cordonnier, and Martin Jaggi. Sparsified sgd with memory. *Advances  
 in neural information processing systems*, 31, 2018.

775

776 Tao Sun, Qingsong Wang, Dongsheng Li, and Bao Wang. Momentum ensures convergence of signsgd  
 777 under weaker assumptions. In *International Conference on Machine Learning*, pp. 33077–33099.  
 778 PMLR, 2023.

779

780 Yi-Lin Sung, Varun Nair, and Colin A Raffel. Training neural networks with fixed sparse masks.  
*Advances in Neural Information Processing Systems*, 34:24193–24205, 2021.

781

782 Ananda Theertha Suresh, X Yu Felix, Sanjiv Kumar, and H Brendan McMahan. Distributed mean  
 783 estimation with limited communication. In *International conference on machine learning*, pp.  
 784 3329–3337. PMLR, 2017.

785 Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya  
 786 Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, Johan  
 787 Ferret, Peter Liu, Pouya Tafti, Abe Friesen, Michelle Casbon, Sabela Ramos, Ravin Kumar,  
 788 Charline Le Lan, Sammy Jerome, Anton Tsitsulin, Nino Vieillard, Piotr Stanczyk, Sertan Girgin,  
 789 Nikola Momchev, Matt Hoffman, Shantanu Thakoor, Jean-Bastien Grill, Behnam Neyshabur,  
 790 Olivier Bachem, Alanna Walton, Aliaksei Severyn, Alicia Parrish, Aliya Ahmad, Allen Hutchison,  
 791 Alvin Abdagic, Amanda Carl, Amy Shen, Andy Brock, Andy Coenen, Anthony Laforge, Antonia  
 792 Paterson, Ben Bastian, Bilal Piot, Bo Wu, Brandon Royal, Charlie Chen, Chintu Kumar, Chris  
 793 Perry, Chris Welty, Christopher A. Choquette-Choo, Danila Sinopalnikov, David Weinberger,  
 794 Dimple Vijaykumar, Dominika Rogozińska, Dustin Herbison, Elisa Bandy, Emma Wang, Eric  
 795 Noland, Erica Moreira, Evan Senter, Evgenii Eltyshev, Francesco Visin, Gabriel Rasskin, Gary  
 796 Wei, Glenn Cameron, Gus Martins, Hadi Hashemi, Hanna Klimczak-Plucińska, Harleen Batra,  
 797 Harsh Dhand, Ivan Nardini, Jacinda Mein, Jack Zhou, James Svensson, Jeff Stanway, Jetha  
 798 Chan, Jin Peng Zhou, Joana Carrasqueira, Joana Iljazi, Jocelyn Becker, Joe Fernandez, Joost  
 799 van Amersfoort, Josh Gordon, Josh Lipschultz, Josh Newlan, Ju yeong Ji, Kareem Mohamed,  
 800 Kartikeya Badola, Kat Black, Katie Millican, Keelin McDonell, Kelvin Nguyen, Kiranbir Sodhia,  
 801 Kish Greene, Lars Lowe Sjoesund, Lauren Usui, Laurent Sifre, Lena Heuermann, Leticia Lago,  
 802 Lilly McNealus, Livio Baldini Soares, Logan Kilpatrick, Lucas Dixon, Luciano Martins, Machel  
 803 Reid, Manvinder Singh, Mark Iverson, Martin Görner, Mat Velloso, Mateo Wirth, Matt Davidow,  
 804 Matt Miller, Matthew Rahtz, Matthew Watson, Meg Risdal, Mehran Kazemi, Michael Moynihan,  
 805 Ming Zhang, Minsuk Kahng, Minwoo Park, Mofi Rahman, Mohit Khatwani, Natalie Dao, Nenshad  
 806 Bardoliwalla, Nesh Devanathan, Neta Dumai, Nilay Chauhan, Oscar Wahltinez, Pankil Botarda,  
 807 Parker Barnes, Paul Barham, Paul Michel, Pengchong Jin, Petko Georgiev, Phil Culliton, Pradeep  
 808 Kuppala, Ramona Comanescu, Ramona Merhej, Reena Jana, Reza Ardeshir Rokni, Rishabh  
 809 Agarwal, Ryan Mullins, Samaneh Saadat, Sara Mc Carthy, Sarah Cogan, Sarah Perrin, Sébastien  
 M. R. Arnold, Sebastian Krause, Shengyang Dai, Shruti Garg, Shruti Sheth, Sue Ronstrom, Susan  
 Chan, Timothy Jordan, Ting Yu, Tom Eccles, Tom Hennigan, Tomas Kociský, Tulsee Doshi,  
 Vihan Jain, Vikas Yadav, Vilobh Meshram, Vishal Dharmadhikari, Warren Barkley, Wei Wei,

810 Wenming Ye, Woohyun Han, Woosuk Kwon, Xiang Xu, Zhe Shen, Zhitao Gong, Zichuan Wei,  
 811 Victor Cotruta, Phoebe Kirk, Anand Rao, Minh Giang, Ludovic Peran, Tris Warkentin, Eli Collins,  
 812 Joelle Barral, Zoubin Ghahramani, Raia Hadsell, D. Sculley, Jeanine Banks, Anca Dragan, Slav  
 813 Petrov, Oriol Vinyals, Jeff Dean, Demis Hassabis, Koray Kavukcuoglu, Clement Farabet, Elena  
 814 Buchatskaya, Sebastian Borgeaud, Noah Fiedel, Armand Joulin, Kathleen Kenealy, Robert Dadashi,  
 815 and Alek Andreev. Gemma 2: Improving open language models at a practical size, 2024. URL  
 816 <https://arxiv.org/abs/2408.00118>.

817 Yuxuan Tong, Xiwen Zhang, Rui Wang, Ruidong Wu, and Junxian He. Dart-math: Difficulty-aware  
 818 rejection tuning for mathematical problem-solving. 2024. URL <https://arxiv.org/abs/2407.13690>.

819 Bokun Wang, Mher Safaryan, and Peter Richtárik. Theoretically better and numerically faster  
 820 distributed optimization with smoothness-aware quantization techniques. *Advances in Neural  
 821 Information Processing Systems*, 35:9841–9852, 2022.

822 Hongyi Wang, Scott Sievert, Shengchao Liu, Zachary Charles, Dimitris Papailiopoulos, and Stephen  
 823 Wright. Atomo: Communication-efficient learning via atomic sparsification. *Advances in neural  
 824 information processing systems*, 31, 2018.

825 Zihan Wang, Rui Pan, Jiarui Yao, Robert Csordas, Linjie Li, Lu Yin, Jiajun Wu, Tong Zhang, Manling  
 826 Li, and Shiwei Liu. Chain-of-experts: Unlocking the communication power of mixture-of-experts  
 827 models. *arXiv preprint arXiv:2506.18945*, 2025.

828 Jianqiao Wangni, Jialei Wang, Ji Liu, and Tong Zhang. Gradient sparsification for communication-  
 829 efficient distributed optimization. *Advances in Neural Information Processing Systems*, 31, 2018.

830 Mengzhou Xia, Tianyu Gao, Zhiyuan Zeng, and Danqi Chen. Sheared llama: Accelerating language  
 831 model pre-training via structured pruning. *arXiv preprint arXiv:2310.06694*, 2023.

832 Lifan Yuan, Ganqu Cui, Hanbin Wang, Ning Ding, Xingyao Wang, Jia Deng, Boji Shan, Huimin Chen,  
 833 Ruobing Xie, Yankai Lin, Zhenghao Liu, Bowen Zhou, Hao Peng, Zhiyuan Liu, and Maosong Sun.  
 834 Advancing llm reasoning generalists with preference trees, 2024.

835 Xiang Yue, Xingwei Qu, Ge Zhang, Yao Fu, Wenhao Huang, Huan Sun, Yu Su, and Wenhui Chen.  
 836 Mammoth: Building math generalist models through hybrid instruction tuning. *arXiv preprint  
 837 arXiv:2309.05653*, 2023.

838 Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine  
 839 really finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association for  
 840 Computational Linguistics*, pp. 4791–4800, Florence, Italy, July 2019a. Association for Com-  
 841 putational Linguistics. doi: 10.18653/v1/P19-1472. URL <https://aclanthology.org/P19-1472/>.

842 Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine  
 843 really finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association for  
 844 Computational Linguistics*, 2019b.

845 Luoming Zhang, Zhenyu Lou, Yangwei Ying, Cheng Yang, and Hong Zhou. Efficient fine-tuning of  
 846 large language models via a low-rank gradient estimator. *Applied Sciences*, 15(1):82, 2024a.

847 Peiyuan Zhang, Guangtao Zeng, Tianduo Wang, and Wei Lu. Tinyllama: An open-source small  
 848 language model, 2024b.

849 Yuchen Zhang and Lin Xiao. Stochastic primal-dual coordinate method for regularized empirical risk  
 850 minimization. *Journal of Machine Learning Research*, 18(84):1–42, 2017.

851 Jiawei Zhao, Zhenyu Zhang, Beidi Chen, Zhangyang Wang, Anima Anandkumar, and Yuandong  
 852 Tian. Galore: Memory-efficient llm training by gradient low-rank projection. *arXiv preprint  
 853 arXiv:2403.03507*, 2024.

854 Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny  
 855 Zhou, and Le Hou. Instruction-following evaluation for large language models, 2023. URL  
 856 <https://arxiv.org/abs/2311.07911>.

864 Hanqing Zhu, Zhenyu Zhang, Wenyan Cong, Xi Liu, Sem Park, Vikas Chandra, Bo Long, David Z.  
865 Pan, Zhangyang Wang, and Jinwon Lee. APOLLO: SGD-like memory, adamw-level performance.  
866 In *Eighth Conference on Machine Learning and Systems*, 2025. URL <https://openreview.net/forum?id=mJrPkdcZDj>.  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917

---

**Algorithm 3** An unbiased version of Algorithm 1

---

1: **Input:**  $\{W_{0,\ell} \in \mathbb{R}^{m_\ell \times n_\ell}\}$  with each  $\ell$  corresponding to the  $\ell$ -th block of parameters, number of blocks  $N_L$ , sampling period  $K$ , projection rank  $r$

2: **for**  $t = 0$  to  $T/K - 1$  **do**

3:     **delete\_optimizer\_states()**                   ▷ Delete the optimizer states to clear memory

4:     Each block  $\ell$  is sampled to do full-rank updates with probability  $q_{t,\ell}$

5:     **for**  $k = 0$  to  $K - 1$  **do**

6:         **for**  $\ell = 1$  to  $N_\ell$  **do**

7:              $G_{t,k,\ell} = G(W_{tK+k-1,\ell})$                    ▷ Obtain the gradient of the  $\ell$ -th layer at  $W_t$

8:              $P_{t,k,\ell} \leftarrow \text{get\_projector}()$                    ▷ Obtain the projector  $P_{t,k,\ell} \in \mathbb{R}^{m_\ell \times r}$

9:              $\tilde{G}_{t,k,\ell} = \begin{cases} \frac{1}{q_{t,\ell}}(I_{m_\ell} - P_{t,k,\ell}P_{t,k,\ell}^\top)G_{t,k,\ell}, & \text{if block } \ell \text{ is sampled to be full-rank} \\ \frac{1}{1-q_{t,\ell}}P_{t,k,\ell}^\top G_{t,k,\ell}, & \text{else} \end{cases}$

10:              $S_{t,k,\ell} \leftarrow \text{optimizer.update\_state}(\tilde{G}_{t,k,\ell})$                    ▷ Run the base algorithm with  $\tilde{G}_{t,k,\ell}$

11:              $W_{tK+k+1,\ell} = \begin{cases} W_{tK+k,\ell} - \eta S_{t,k,\ell}, & \text{if block } \ell \text{ is sampled to be full-rank} \\ W_{tK+k,\ell} - \eta P_{t,k,\ell}S_{t,k,\ell}, & \text{else} \end{cases}$

12:         **end for**

13:     **end for**

14: **end for**

A GENERAL UNBIASED LOW-RANK GRADIENT METHOD PARADIGM

Here, we present our unbiased algorithm paradigm in Algorithm 3. The key idea of the algorithm is to compensate for biased errors introduced by the low-rank projection  $P_t P_t^\top G_t$ . To implement this while retaining memory efficiency, we refer to the main idea of LISA (Pan et al., 2024), which allows some of the blocks to be sampled uniformly with probability  $q$  in each period. This compensated full-rank updates use  $G_t - P_t P_t^\top G_t$ , while other blocks still do the original low-rank update. By carefully balancing the scaling constants for the two different updates, the biased low-rank term can be canceled out in expectation, resulting in an unbiased estimation of gradients across iterations. This unbiased version of the algorithm is presented in Algorithm 3.

In one training process, the algorithm contains separated periods just like the vanilla GaLore algorithm (Zhao et al., 2024) and LISA (Pan et al., 2024). During each period  $t$ , each block of parameters is sampled to do full-rank updates with probability  $q_{t,\ell}$ . In each iteration  $k$  in the period, we first compute the projection matrix  $P_{t,k,\ell}$ . Note that a lot of strategies for selecting projection matrices and sampling importance can be applied here (Guo et al., 2020; Sung et al., 2021; Ansell et al., 2021; Das et al., 2023; Muhamed et al., 2024; Ramesh et al., 2024; Liu et al., 2024a). Then, the blocks not sampled to do full-rank updates run basically the same low-rank update with Algorithm 1, while the full-rank blocks directly run the base optimizer with the compensated gradient  $\tilde{G}_{t,k,\ell} = (I_m - P_{t,k,\ell}P_{t,k,\ell}^\top)G_{t,k,\ell}$ .

We note that the proposed debiasing technique Algorithm 3 works generally when the following properties are satisfied:

- **Property I.** The columns of the low-rank projection matrix  $P_t \in \mathbb{R}^{m \times r}$  with  $r \leq m$  are orthonormal, i.e.,  $P_t^\top P_t = I_{r \times r}$ .
- **Property II.** The projection and optimization updates are commutable, which means that  $S_t = P_t \text{optimizer}.\text{update\_state}(\tilde{G}_t) = \text{optimizer}.\text{update\_state}(P_t \tilde{G}_t)$ . Optimizers satisfying this property typically treat the update parameters as matrices instead of vectors, and only conduct matrix operations in the update. Two standard examples include SGD and Muon (Jordan et al., 2024).

If the two properties are satisfied, we can show that Algorithm 3 is unbiased compared to the base optimizer, since it is equivalent to running the base optimizer with an unbiased estimation of the gradient at each iteration.

972     **Lemma 2** (Unbiased update of Algorithm 3). *When Property I and II are satisfied, a single iteration  
973     of Algorithm 3 for  $W \in \mathbb{R}^{m \times n}$  is equivalent to*

974                      $S \leftarrow \text{optimizer.update\_state}(\hat{G})$   
975                      $W^+ = W - \eta S$   
976

977     with  $\mathbb{E}[\hat{G}] = G \in \mathbb{R}^{m \times n}$ , where  $G$  denotes the gradient obtained at  $W$ .

980     **B PROOFS OF SECTION 3 AND 4**

982     **B.1 PROOF OF LEMMA 2 AND 1**

984     *Proof of Lemma 2.* A single step of Algorithm 3 writes:

985                      $\tilde{G} = \begin{cases} \frac{1}{q}(I - PP^\top)G, & \text{with probability } q \\ \frac{1}{1-q}P^\top G, & \text{with probability } 1 - q \end{cases}$   
986                      $S = \text{optimizer.update\_state}(\tilde{G})$   
987                      $W^+ = \begin{cases} W - \eta S, & \text{with probability } q \\ W - \eta PS, & \text{with probability } 1 - q \end{cases}$   
988

992     where  $G$  is the gradient at  $W$  and  $P$  is the projection matrix obtained at  $W$ . Based on the commutative  
993     property, we know that

994                      $W^+ = W - \eta PS = W - \eta \text{optimizer.update\_state}(P\tilde{G}),$   
995

996     which means that the update step is equivalent to

997                      $\hat{G} = \begin{cases} \frac{1}{q}(I - PP^\top)G, & \text{with probability } q \\ \frac{1}{1-q}PP^\top G, & \text{with probability } 1 - q \end{cases}$   
998                      $S = \text{optimizer.update\_state}(\hat{G})$   
999                      $W^+ = W - \eta S$   
1000

1003     Since we have  $\hat{G}$  is an unbiased estimation of  $G$ :

1004                      $\mathbb{E}[\hat{G}] = q \cdot \frac{1}{q}(I - PP^\top)G + (1 - q) \cdot \frac{1}{1-q}PP^\top G = G,$   
1005

1006     we finish the proof that Algorithm 3 is unbiased compared to the base optimizer.  $\square$

1008     *Proof of Lemma 1.* Based on Lemma 2, we only need to prove that GUM satisfies the two properties  
1009     to finish the proof of Lemma 1.

1011     **Property I.** Denote the projection matrix at one specific iteration  $P$ . Since  $P$  is obtained from the  
1012     SVD, we have  $P \in \mathbb{R}^{m \times r}$  and  $P^\top P = I_r$ .

1013     **Property II.** The base algorithm of GUM is Muon (Jordan et al., 2024). To prove the commutative  
1014     property, we only need to prove that the Newton-Schulz iteration is commutable with  $P$ . In each  
1015     iteration of the Newton Schulz iteration `NewtonSchulz(X)`, we compute

1017                      $X^+ = aX + bXX^\top X + cXX^\top XX^\top X,$

1018     where  $a, b, c \in \mathbb{R}$  are absolute constants. Then consider `NewtonSchulz(PX)`, we get

1020                      $X^+ = aPX + bPX(PX)^\top (PX) + cPX(PX)^\top (PX)(PX)^\top (PX)$   
1021                      $= P(aX + bXX^\top X + cXX^\top XX^\top X),$

1023     where the second equality is because of Property I. Therefore, we obtain that

1024                      $\text{NewtonSchulz}(PX) = P \cdot \text{NewtonSchulz}(X),$

1025     which finishes the proof of Property II and thus the unbiased property of GUM.  $\square$

1026 B.2 PROOF OF THEOREM 1  
1027

1028 We first state the notations in the following proof writing. For simplicity, we assume that the total  
1029 iteration number  $T = K\tau$ . For  $k = 0, \dots, K-1$  in a specific period  $t = 0, \dots, \tau-1$ , Algorithm 2  
1030 is mathematically equivalent to the following formulation:

$$1031 \tilde{G}_{t,k} = \begin{cases} \frac{1}{1-q_t} P_t P_t^\top G_{t,k}, & \text{if } \xi_t = 0 \\ \frac{1}{q_t} (I - P_t P_t^\top) G_{t,k}, & \text{else} \end{cases}$$

$$1034 \tilde{M}_{t,k} = \beta \tilde{M}_{t,k-1} + (1 - \beta) \tilde{G}_{t,k}$$

$$1035 W_{tK+k+1} = W_{tK+k} - \eta \text{NewtonSchulz}(\tilde{M}_{t,k})$$

1037 where  $G_{t,k}$  is the stochastic gradient obtained at  $W_{tK+k}$  and  $\xi_t \sim \text{Bernoulli}(q_t)$  is the indicator  
1038 random variable such that  $\xi_t = 1$  means using full-rank update in period  $t$ . We assume that the  
1039 full-rank probability  $q_t \equiv q$  and step size  $\eta_t \equiv \eta$  are constants. The equivalence of Algorithm 2  
1040 and this formulation is shown by Lemma 1. At the beginning of each period, we initialize  $P_t$  from  
1041  $G_{t,0}$  and set  $\tilde{M}_{t,-1} = 0$ . Also, we denote  $\nabla f_{t,k} \triangleq \nabla f(W_{tK+k})$  and  $\text{msign}(X) \triangleq UV^\top$  for  
1042  $X = U\Sigma V^\top$  as the SVD of  $X$ . Under Assumption 4, we have  $\text{NewtonSchulz}(X) = \text{msign}(X)$ .  
1043 Note that here in the theoretical proof, we consider the damping, i.e., the  $1 - \beta$  term in the update of  
1044  $\tilde{M}_{t,k}$ . Since we initialize  $\tilde{M}_{t,k} = 0$  in each period, this damping will not affect the algorithm because  
1045 the Newton-Schulz iteration is irrelevant to the input scale.

1046 To help simplify the convergence proof, we also denote the residual of the projector as  $R_t \in$   
1047  $\mathbb{R}^{m \times (m-r)}$ , i.e., we take  $U_t = [P_t \ R_t] \in \mathbb{R}^{m \times m}$ , which satisfies that  $P_t^\top R_t = 0$ ,  $R_t^\top P_t = 0$ . Note  
1048 that since we consider only the case  $m \leq n$  here, we have  $U_t U_t^\top = P_t P_t^\top + R_t R_t^\top = I$ . We further  
1049 define

$$1050 Q_t \triangleq \begin{cases} P_t, & \text{if } \xi_t = 0 \\ R_t, & \text{else} \end{cases} \quad (3)$$

1052 and the following auxiliary sequence

$$1054 M_{t,k} = \beta M_{t,k-1} + (1 - \beta) G_{t,k} \quad (4)$$

1055 with  $M_{t,-1} = 0$ , which is the exponential moving average of the real gradient. With these definitions,  
1056 we have

$$1058 \text{msign}(\tilde{M}_{t,k}) = \text{msign}(Q_t Q_t^\top M_{t,k}) = Q_t \text{msign}(Q_t^\top M_{t,k}), \quad (5)$$

1059 where the equation is based on the fact that  $Q_t^\top Q_t = I$ .

1061 We first make use of the smoothness assumption to obtain a one-step analysis.

1062 **Lemma 3** (One-step descent). *Under Assumption 2 and 4 and setting  $\eta_t \equiv \eta$ , for  $t = 0, \dots, \tau-1$  and  
1063  $k = 0, \dots, K-1$ , it holds that*

$$1065 f(W_{tK+k+1}) \leq f(W_{tK+k}) - \eta \|Q_t^\top \nabla f_{t,k}\|_* + \frac{1}{2} \eta^2 L_{\text{op}} + 2\eta \|M_{t,k} - \nabla f_{t,k}\|_*, \quad (6)$$

1067 where  $Q_t$  is defined as (3).

1068 *Proof.* Based on Assumption 2, we have the descent property

$$1071 f(W_{tK+k+1}) \leq f(W_{tK+k}) + \langle \nabla f_{t,K}, W_{tK+k+1} - W_{tK+k} \rangle + \frac{L_{\text{op}}}{2} \|W_{tK+k+1} - W_{tK+k}\|_{\text{op}}^2$$

$$1073 = f(W_{tK+k}) - \eta \langle \nabla f_{t,K}, \text{msign}(\tilde{M}_{t,k}) \rangle + \frac{L_{\text{op}}\eta^2}{2} \|\text{msign}(\tilde{M}_{t,k})\|_{\text{op}}^2$$

$$1075 = f(W_{tK+k}) - \eta \langle M_{t,K}, \text{msign}(\tilde{M}_{t,k}) \rangle + \frac{L_{\text{op}}\eta^2}{2}$$

$$1077 + \eta \langle M_{t,K} - \nabla f_{t,k}, \text{msign}(\tilde{M}_{t,k}) \rangle$$

$$1079 \leq f(W_{tK+k}) - \eta \langle M_{t,K}, \text{msign}(\tilde{M}_{t,k}) \rangle + \frac{L_{\text{op}}\eta^2}{2} + \eta \|M_{t,K} - \nabla f_{t,k}\|_*,$$

1080 where the last inequality is based on the fact that  $\|\cdot\|_*$  and  $\|\cdot\|_{\text{op}}$  are dual norms and  
 1081  $\|\text{msign}(\tilde{M}_{t,k})\|_{\text{op}} = 1$ . Then we further deal with the second term on the right hand side:  
 1082

$$\begin{aligned}
 1084 \quad & -\langle M_{t,K}, \text{msign}(\tilde{M}_{t,k}) \rangle \stackrel{(5)}{=} -\langle M_{t,K}, Q_t \text{msign}(Q_t^\top M_{t,k}) \rangle \\
 1085 \quad & = -\langle Q_t^\top M_{t,K}, \text{msign}(Q_t^\top M_{t,k}) \rangle = -\|Q_t^\top M_{t,k}\|_* \\
 1086 \quad & \leq -\|Q_t^\top \nabla f_{t,k}\|_* + \|Q_t^\top (M_{t,k} - \nabla f_{t,k})\|_* \\
 1087 \quad & \leq -\|Q_t^\top \nabla f_{t,k}\|_* + \|M_{t,k} - \nabla f_{t,k}\|_*, \\
 1088 \quad & \leq -\|Q_t^\top \nabla f_{t,k}\|_* + \|M_{t,k} - \nabla f_{t,k}\|_*, \\
 1089 \end{aligned}$$

1090 where the last inequality is based on that  $Q_t Q_t^\top \preceq I$ . Then combining the inequalities, we can finish  
 1091 the proof.  $\square$   
 1092

1093 Based on Lemma 3, we could find that a key to proving the convergence is the  $\|M_{t,k} - \nabla f_{t,k}\|_*$   
 1094 term. Let us define the following auxiliary sequences:  
 1095

$$1096 \quad \epsilon_{t,k} \triangleq M_{t,k} - \nabla f_{t,k}, \quad S_{t,k} \triangleq \nabla f_{t,k-1} - \nabla f_{t,k}, \quad N_{t,k} \triangleq G_{t,k} - \nabla f_{t,k} \quad (7)$$

1097 and additionally set  $\nabla f_{t,-1} \triangleq \nabla f_{t,0}$  for all  $t = 0, \dots, \tau - 1$ . Then we consider decomposing the  
 1098 desired  $\epsilon_t$  based on the properties of moving average sequences.  
 1099

1100 **Lemma 4** (Decompose  $\epsilon_{t,k}$ ). *For  $t = 0, \dots, \tau - 1$  and  $k = 0, \dots, K - 1$ , it holds that*

$$1102 \quad \epsilon_{t,k} = \sum_{i=1}^k \beta^{k-i+1} S_{t,i} + (1 - \beta) \sum_{i=0}^k \beta^{k-i} N_{t,i} - \beta^k \nabla f_{t,0}. \quad (8)$$

1105 *Proof.* From the definition of  $M_{t,k}$  in (4), we know that

$$1107 \quad M_{t,k} = \beta M_{t,k-1} + (1 - \beta) G_{t,k},$$

1108 which implies that

$$\begin{aligned}
 1110 \quad \epsilon_t &= \beta(M_{t,k-1} - \nabla f_{t,k-1}) + \beta(\nabla f_{t,k-1} - \nabla f_{t,k}) + (1 - \beta)(G_{t,k} - \nabla f_{t,k}) \\
 1111 &= \beta \epsilon_{t,k-1} + \beta S_{t,k} + (1 - \beta) N_{t,k}.
 \end{aligned}$$

1112 Then by applying the equality recursively and noting that  $\epsilon_{t,0} = (1 - \beta)G_{t,0} - \nabla f_{t,0} = (1 - \beta)N_{t,0} - \beta \nabla f_{t,0}$ , we conclude the proof.  $\square$   
 1113  
 1114

1115 Then we produce the next lemma to state the variance contraction properties of momentum for Muon,  
 1116 which has been explored for Normalized SGD (Cutkosky & Mehta, 2020) and SignSGD (Sun et al.,  
 1117 2023), and also for Muon (Li & Hong, 2025), but with different assumptions.  
 1118

1119 **Lemma 5** (Variance Contraction). *Under Assumption 3, for  $t = 0, \dots, \tau - 1$  and  $k = 0, \dots, K - 1$ ,  
 1120 it holds that*

$$1121 \quad \mathbb{E} \left[ \left\| (1 - \beta) \sum_{i=0}^k \beta^{k-i} N_{t,i} \right\|_* \right] \leq \|V\|_* \sqrt{(1 - \beta^{2k})(1 - \beta)}. \quad (9)$$

1124 *Proof.* Based on Lemma 8 in An et al. (2025), for an arbitrary symmetric positive definite matrix  
 1125  $H \in \mathbb{R}^{m \times m}$ , it holds that

$$\begin{aligned}
 1127 \quad & \mathbb{E} \left[ \left\| \sum_{i=0}^k \beta^{k-i} N_{t,i} \right\|_* \right] \leq \mathbb{E} \left[ \sqrt{\|H\|_* \text{tr} \left( \left( \sum_{i=0}^k \beta^{k-i} N_{t,i} \right)^\top H^{-1} \left( \sum_{i=0}^k \beta^{k-i} N_{t,i} \right) \right)} \right] \\
 1128 & = \mathbb{E} \left[ \sqrt{\|H\|_* \text{tr} \left( \left( \sum_{i=0}^k \beta^{k-i} N_{t,i} \right) \left( \sum_{i=0}^k \beta^{k-i} N_{t,i} \right)^\top H^{-1} \right)} \right]
 \end{aligned}$$

$$\begin{aligned}
&\leq \sqrt{\|H\|_* \mathbb{E} \left[ \text{tr} \left( \left( \sum_{i=0}^k \beta^{k-i} N_{t,i} \right) \left( \sum_{i=0}^k \beta^{k-i} N_{t,i} \right)^\top H^{-1} \right) \right]} \\
&= \sqrt{\|H\|_* \mathbb{E} \left[ \text{tr} \left( \left( \sum_{i=0}^k \beta^{2(k-i)} N_{t,i} N_{t,i}^\top \right) H^{-1} \right) \right]},
\end{aligned}$$

where the last inequality is based on the fact that  $\mathbb{E}[\sqrt{X}] \leq \sqrt{\mathbb{E}[X]}$  and the last equality is based on the assumption that  $N_{t,i}$  and  $N_{t,j}$  are independent for  $i \neq j$ , which implies  $\mathbb{E}[\text{tr}(N_{t,i} N_{t,j}^\top H)] = 0$ . Then taking  $H = (VV^\top)^{1/2}$  leads to

$$\begin{aligned}
\sqrt{\|H\|_* \mathbb{E} \left[ \text{tr} \left( \left( \sum_{i=0}^k \beta^{2(k-i)} N_{t,i} N_{t,i}^\top \right) H \right) \right]} &= \sqrt{\|V\|_* \mathbb{E} \left[ \text{tr} \left( \sum_{i=0}^k \beta^{2(k-i)} N_{t,i} N_{t,i}^\top (VV^\top)^{-\frac{1}{2}} \right) \right]} \\
&\leq \sqrt{\|V\|_* \sum_{i=0}^k \beta^{2(k-i)} \text{tr} \left( VV^\top (VV^\top)^{-\frac{1}{2}} \right)} \\
&\leq \|V\|_* \sqrt{\frac{1 - \beta^{2k}}{1 - \beta^2}},
\end{aligned}$$

where the first inequality is based on Assumption 3 and the second inequality is by algebra. Then, combining the inequalities and multiplying  $1 - \beta$  gives the result.  $\square$

**Lemma 6** (Bound  $\mathbb{E}[\|\epsilon_{t,k}\|_*]$ ). *Under Assumption 3, for  $t = 0, \dots, \tau - 1$  and  $k = 0, \dots, K - 1$ , it holds that*

$$\mathbb{E}[\|\epsilon_{t,k}\|_*] \leq \frac{1 - \beta^k}{1 - \beta} L_{\text{op}} \eta + \sqrt{(1 - \beta^{2k})(1 - \beta)} \|V\|_* + \beta^k \mathbb{E}[\|\nabla f_{t,0}\|_*]. \quad (10)$$

*Proof.* Based on Lemma 4, it holds that

$$\begin{aligned}
\mathbb{E}[\|\epsilon_{t,k}\|_*] &= \mathbb{E} \left[ \left\| \sum_{i=1}^k \beta^{k-i+1} S_{t,i} + (1 - \beta) \sum_{i=0}^k \beta^{k-i} N_{t,i} - \beta^k \nabla f_{t,0} \right\|_* \right] \\
&\leq \sum_{i=1}^k \beta^{k-i+1} \|S_{t,i}\|_* + \left\| (1 - \beta) \sum_{i=0}^k \beta^{k-i} N_{t,i} \right\|_* + \|\beta^k \nabla f_{t,0}\|_*
\end{aligned}$$

where the inequality is based on the triangular inequality. For the first term in the RHS, it holds that

$$\|S_{t,i}\|_* = \|\nabla f_{t,i-1} - \nabla f_{t,i}\|_* \leq L_{\text{op}} \|W_{tK+i-1} - W_{tK+i}\| = L_{\text{op}} \eta.$$

Thus we have

$$\begin{aligned}
\mathbb{E}[\|\epsilon_{t,k}\|_*] &\leq \sum_{i=1}^k \beta^{k-i+1} L_{\text{op}} \eta + \mathbb{E} \left[ \left\| (1 - \beta) \sum_{i=0}^k \beta^{k-i} N_{t,i} \right\|_* \right] + \mathbb{E}[\|\beta^k \nabla f_{t,0}\|_*] \\
&\stackrel{(9)}{\leq} \sum_{i=1}^k \beta^{k-i+1} L_{\text{op}} \eta + \sqrt{(1 - \beta^{2k})(1 - \beta)} \|V\|_* + \beta^k \mathbb{E}[\|\nabla f_{t,0}\|_*] \\
&\leq \frac{1 - \beta^k}{1 - \beta} L_{\text{op}} \eta + \sqrt{(1 - \beta^{2k})(1 - \beta)} \|V\|_* + \beta^k \mathbb{E}[\|\nabla f_{t,0}\|_*],
\end{aligned}$$

which concludes the proof.  $\square$

We need to further determine the expected projected gradient for  $\nabla f_{t,0}$ .

**Lemma 7** (Expected projected gradient). *For  $t = 0, \dots, \tau - 1$ , it holds that*

$$\mathbb{E}[\|Q_t^\top \nabla f_{t,0}\|_*] \geq \min\{q, 1 - q\} \mathbb{E}[\|\nabla f_{t,0}\|_*]. \quad (11)$$

1188 *Proof.* Based on the algorithm, we know that  $\xi_t$  and  $W_{tK}$  are independent, which means that  
1189

$$\mathbb{E} [\|Q_t^\top \nabla f_{t,0}\|_*] = (1-q)\mathbb{E} [\|P_t^\top \nabla f_{t,0}\|_*] + q\mathbb{E} [\|R_t^\top \nabla f_{t,0}\|_*].$$

1190 Because we have  $U_t = [P_t \ R_t]$  that satisfies  $U_t^\top U_t = U_t U_t^\top = I$ , it holds for any  $X \in \mathbb{R}^{m \times n}$  that  
1191

$$\begin{aligned} \|P_t^\top X\|_* + \|R_t^\top X\|_* &= \text{tr} \left( (X^\top P_t P_t^\top X)^{\frac{1}{2}} \right) + \text{tr} \left( (X^\top R_t R_t^\top X)^{\frac{1}{2}} \right) \\ &\geq \text{tr} \left( (X^\top (P_t P_t^\top + R_t R_t^\top) X)^{\frac{1}{2}} \right) \\ &= \text{tr} \left( (X^\top X)^{\frac{1}{2}} \right) = \|X\|_*. \end{aligned}$$

1192 Therefore, we have  
1193

$$\begin{aligned} \mathbb{E} [\|Q_t^\top \nabla f_{t,0}\|_*] &= (1-q)\mathbb{E} [\|P_t^\top \nabla f_{t,0}\|_*] + q\mathbb{E} [\|R_t^\top \nabla f_{t,0}\|_*] \\ &\geq \min \{q, 1-q\} (\mathbb{E} [\|P_t^\top \nabla f_{t,0}\|_*] + \mathbb{E} [\|R_t^\top \nabla f_{t,0}\|_*]) \\ &\geq \min \{q, 1-q\} \mathbb{E} [\|\nabla f_{t,0}\|_*], \end{aligned}$$

1194 which completes the proof.  $\square$   
1195

1196 With the lemmas in hand, we are able to prove Theorem 1.  
1197

1198 *Proof of Theorem 1.* Based on Lemma 3, for  $t = 0, \dots, \tau - 1$  and  $k = 0, \dots, K - 1$ , it holds that  
1199

$$\begin{aligned} f(W_{tK+k+1}) &\stackrel{(6)}{\leq} f(W_{tK+k}) - \eta \|Q_t^\top \nabla f_{t,k}\|_* + \frac{1}{2}\eta^2 L_{\text{op}} + 2\eta \|M_{t,k} - \nabla f_{t,k}\|_* \\ &= f(W_{tK+k}) - \eta \|Q_t^\top \nabla f_{t,k}\|_* + \frac{1}{2}\eta^2 L_{\text{op}} + 2\eta \|\epsilon_{t,k}\|_*, \end{aligned}$$

1200 where  $Q_t$  is defined in (3) and  $\epsilon_{t,k}$  is defined in (7). Then, after rearrangement and summation over  $k$   
1201 and taking expectation, we have  
1202

$$\begin{aligned} \sum_{k=0}^{K-1} \eta \mathbb{E} [\|Q_t^\top \nabla f_{t,k}\|_*] &\leq \mathbb{E} [f(W_{tK}) - f(W_{(t+1)K})] + \frac{1}{2}\eta^2 K L_{\text{op}} + 2\eta \sum_{k=0}^{K-1} \mathbb{E} [\|\epsilon_{t,k}\|_*] \\ &\stackrel{(10)}{\leq} \mathbb{E} [f(W_{tK}) - f(W_{(t+1)K})] + \frac{1}{2}\eta^2 K L_{\text{op}} \\ &\quad + 2\eta \sum_{k=0}^{K-1} \left( \frac{1-\beta^k}{1-\beta} L_{\text{op}} \eta + \left( \sqrt{(1-\beta^{2k})(1-\beta)} + \beta^k \right) \|V\|_* \right) \\ &\leq \mathbb{E} [f(W_{tK}) - f(W_{(t+1)K})] + \eta^2 K L_{\text{op}} \left( \frac{1}{2} + \frac{2(1-\beta^K)}{1-\beta} \right) \\ &\quad + 2\eta \left( \sqrt{(1-\beta^{2K})(1-\beta)} \|V\|_* + \beta^K \mathbb{E} [\|\nabla f_{t,0}\|_*] \right). \end{aligned}$$

1203 Since  $W_{tK+k}$  is dependent on  $Q_t$ , it would be difficult to bound  $\mathbb{E} [\|Q_t^\top \nabla f_{t,k}\|_*]$  for  $k \geq 1$ . We  
1204 therefore consider  
1205

$$\begin{aligned} \sum_{k=0}^{K-1} \eta \mathbb{E} [\|Q_t^\top \nabla f_{t,k}\|_*] &\geq \sum_{k=0}^{K-1} \eta \mathbb{E} [\|Q_t^\top \nabla f_{t,0}\|_*] - \sum_{k=0}^{K-1} \eta \mathbb{E} [\|Q_t^\top (\nabla f_{t,k} - \nabla f_{t,0})\|_*] \\ &\geq \sum_{k=0}^{K-1} \eta \mathbb{E} [\|Q_t^\top \nabla f_{t,0}\|_*] - \sum_{k=1}^{K-1} \eta \mathbb{E} [\|\nabla f_{t,k} - \nabla f_{t,0}\|_*] \\ &\geq \sum_{k=0}^{K-1} \eta \mathbb{E} [\|Q_t^\top \nabla f_{t,0}\|_*] - \sum_{k=1}^{K-1} \eta \sum_{l=1}^k \mathbb{E} [\|\nabla f_{t,l} - \nabla f_{t,l-1}\|_*] \\ &\geq \sum_{k=0}^{K-1} \eta \mathbb{E} [\|Q_t^\top \nabla f_{t,0}\|_*] - \sum_{k=1}^{K-1} \eta L_{\text{op}} \sum_{l=1}^k \mathbb{E} [\|W_{tK+l} - W_{tK+l-1}\|_{\text{op}}] \end{aligned}$$

$$\geq K\eta\mathbb{E} [\|Q_t^\top \nabla f_{t,0}\|_*] - \frac{K^2}{2}\eta^2 L_{\text{op}},$$

where the first and third inequalities are based on the triangular inequality and the second inequality is based on that  $Q_t Q_t^\top \preceq I$ . The second last inequality uses Assumption 2. Then we combine the above inequalities and further sum up over  $t$  and use Assumption 1 to obtain that

$$\begin{aligned} \sum_{t=0}^{\tau-1} K\mathbb{E} [\|Q_t^\top \nabla f_{t,0}\|_*] &\leq \frac{f(W_0) - f^*}{\eta} + \eta K\tau L_{\text{op}} \left( \frac{K+1}{2} + \frac{2(1-\beta^K)}{1-\beta} \right) \\ &\quad + 2\tau K \sqrt{(1-\beta^{2K})(1-\beta)} \|V\|_* + \sum_{t=0}^{\tau-1} \frac{2(1-\beta^K)}{1-\beta} \mathbb{E} [\|\nabla f_{t,0}\|_*]. \end{aligned}$$

Combining Lemma 7, we have

$$K\mathbb{E} [\|Q_t^\top \nabla f_{t,0}\|_*] - \frac{2(1-\beta^K)}{1-\beta} \mathbb{E} [\|\nabla f_{t,0}\|_*] \geq \frac{K\alpha}{2} \mathbb{E} [\|\nabla f_{t,0}\|_*]$$

where  $\alpha \triangleq \min\{q, 1-q\}$  and we take  $\alpha > \frac{2}{K}$  and  $1-\beta \geq \frac{2}{K\alpha}$ . Thus, we can obtain that

$$\begin{aligned} \frac{\alpha}{2\tau} \sum_{t=0}^{\tau-1} \mathbb{E} [\|\nabla f_{t,0}\|_*] &\leq \frac{f(W_0) - f^*}{\eta T} + \eta L_{\text{op}} \left( \frac{K+1}{2} + \frac{2}{1-\beta} \right) + 2\sqrt{1-\beta} \|V\|_* \\ &\leq \frac{f(W_0) - f^*}{\eta T} + \eta L_{\text{op}} \left( \frac{K+1}{2} + K\alpha \right) + 2\sqrt{1-\beta} \|V\|_* \end{aligned}$$

By choosing the hyperparameter as

$$\eta = \sqrt{\frac{TL_{\text{op}} \left( \frac{K+1}{2} + K\alpha \right)}{f(W_0) - f^*}}, \quad \beta = 1 - \frac{2}{K\alpha}, \quad K = \max \left\{ 1, \min \left\{ \frac{\sigma\sqrt{T}}{\sqrt{\alpha L(f(W_0) - f^*)}}, T \right\} \right\}, \quad (12)$$

we can obtain that

$$\min_{t=0, \dots, \frac{T}{K}-1} \mathbb{E} [\|\nabla f(W_{tK})\|] \leq \mathcal{O} \left( \frac{1}{\alpha} \sqrt{\frac{L_{\text{op}} \Delta}{T}} + \left( \frac{L_{\text{op}} \Delta \|V\|_*^2}{\alpha^5 T} \right)^{\frac{1}{4}} + \frac{\|V\|_*}{\sqrt{\alpha^3 T}} \right),$$

with  $\Delta \triangleq f(W_0) - f^*$ , which finishes the proof.  $\square$

## C TRAINING SETUP AND HYPERPARAMETERS

### C.1 FINE-TUNING SETUP

In our experiments, we slightly modify the full-rank update rule (2) for GUM by multiplying  $(1-q_{t,\ell})$  on  $-P_{t,\ell} P_{t,\ell}^\top G_{t,k,\ell}$ . This modification still preserves the unbiased property while being able to recover the original full-parameter Muon algorithm by setting  $q_{t,\ell} = 1$ .

We utilize LMFlow (Diao et al., 2023)<sup>2</sup> to perform full-parameter fine-tuning, GaLore tuning, and GUM tuning. We set the number of training epochs for all fine-tuning scenarios to 1. All experiments were conducted on a single NVIDIA H100 GPU with 80 GB of memory.

We explored a range of learning rates from  $8 \times 10^{-6}$  to  $1 \times 10^{-4}$ , applying this range to Full Parameter Training, GaLore, and GUM. For GaLore, we fixed the rank  $r = 512$  and applied it uniformly across all layers. In the case of GUM, the number of layers ( $\gamma$ ) selected for full-rank updates was set to 2 for all models. The sampling interval  $K$ , which defines the number of update steps between each layer selection, was varied between 10 and 300, depending on factors such as dataset size, batch size, and total training steps. The models covered in this paper can be found in Table 5.

<sup>2</sup><https://github.com/OptimalScale/LMFlow>

1296

1297

1298

1299

1300

1301

1302

1303

1304

1305

## C.2 FINE-TUNING HYPERPARAMETERS

1306

We began our study by conducting a grid search over two key hyperparameters: (i) the learning rate and (ii) the number of sampling layers used for full-rank updates. Given the strong empirical performance of the GaLore method, we fixed the rank to  $r = 512$ . The learning rate was explored within the range  $\{8 \times 10^{-6}, 2 \times 10^{-5}, 4 \times 10^{-5}, 6 \times 10^{-5}, 8 \times 10^{-5}, 1 \times 10^{-4}\}$ , applied consistently across full parameter training, GaLore, and GUM. For GaLore, we followed the official Transformers implementation<sup>3</sup>, using the default settings and aligning the learning rate with the full parameter training. With respect to the number of sampling layers, and in accordance with Table 3, we selected values that did not exceed the GPU memory cost of GaLore. As a result,  $\gamma = 2$  was used in most GUM configurations. The sampling period  $K$  was uniformly set to 200 for all models. A detailed summary of the optimal hyperparameter values identified for each setting is provided in Table 6.

1316

1317

1318

Table 6: Optimal settings for each method were determined through hyperparameter search: FT (Full-parameter Training)-AdamW, FT-Muon, GaLore, and GUM.

1319

1320

Model	FT-AdamW		FT-Muon		GaLore		GUM		
	lr	lr	lr	Rank	lr	$\gamma$	$K$		
LLaMA-3-8B	$3 \times 10^{-5}$	$7 \times 10^{-5}$	$9 \times 10^{-5}$	512	$1 \times 10^{-4}$	2	200		
Qwen-2.5-7B	$1 \times 10^{-5}$	$5 \times 10^{-5}$	$7 \times 10^{-5}$	512	$7 \times 10^{-5}$	2	200		
Gemma-2-9B	—	$4 \times 10^{-5}$	$4 \times 10^{-5}$	512	$6 \times 10^{-5}$	2	200		

1321

1322

1323

1324

1325

1326

1327

1328

1329

## C.3 PRE-TRAINING HYPERPARAMETERS

1330

1331

In our experiments, we utilize C-optim<sup>4</sup> for the pre-training. Following standard protocol, we fixed the LLaMA context length to 1024 tokens. Similar to the fine-tuning setup, we made a grid search on learning rate and the number of sampling layers. The sampling period  $K$  was set to 100 for 130M and 350M models, 50 for the 60M model. A detailed summary of the optimal hyperparameter values identified for each setting is provided in Table 7.

1336

1337

1338

1339

Table 7: Optimal settings for each method were determined through hyperparameter search: AdamW, Muon, Fira, GaLore, and GUM.

1340

1341

1342

1343

1344

1345

1346

1347

1348

1349

<sup>3</sup><https://github.com/jiaweizhao/GaLore>

<sup>4</sup><https://github.com/kyleliang919/C-Optim>

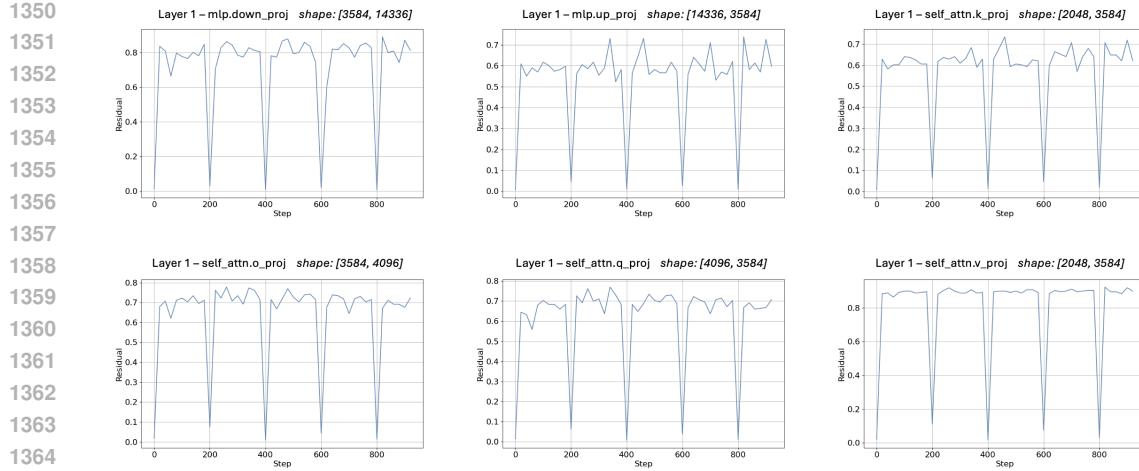


Figure 4: Residual ( $\chi_t = \|G_t^u - G_t^p\|_F / \|G_t^u\|_F$ ) between GaLore’s projected and original gradients across different blocks during **Gemma-2-9B** fine-tuning. High residuals persist throughout training (except for the iterations with projector updates), revealing systematic bias in GaLore updates.

## D ADDITIONAL EXPERIMENTAL RESULTS

### D.1 BIAS IN GALORE

To further illustrate how significant the bias in low-rank projection methods is, we analyze the residuals between low-rank projected gradients and the original full-rank gradients across multiple layers during the fine-tuning of the **Gemma-2-9B** model on the **GPT-4-LLM** dataset. The residual is computed as follows:

$$\chi_t = \frac{\|G_t^u - G_t^p\|_F}{\|G_t^u\|_F}, \quad (13)$$

where  $G_t^u$  represents the original gradient at iteration  $t$  without projection, and  $G_t^p$  denotes the low-rank projected gradients in GaLore-Muon. We can see that  $\chi_t$  presents the relative error between the original gradients and the projected gradients at iteration  $t$ , showing how much the projection operation makes the gradient estimation biased from the original one. We measure this relative error for each block of parameters along the trajectory of the GaLore-Muon algorithm every 20 iterations. The projector update frequency is set to 200, and the projection rank is 512. We use a batch size of 16 and a learning rate of  $7 \times 10^{-5}$ . For demonstration purposes, we specifically select the self-attention and MLP weights at layer 10.

As depicted in Figure 4, the relative error shows a periodic curve. It is relatively small (around 0 – 20%) in the iteration  $t$  such that  $t$  is a multiple of the update frequency 200, where the projector is updated based on the gradient. Since the GaLore projector is chosen as the singular vectors of the largest singular values of the current gradient, it is a good low-rank projector for the current gradient, which results in this small error.<sup>5</sup> However, we can see that the relative error rapidly increases after this and achieves even higher than 60 – 80% in less than 20 iterations. This implies that although the low-rank projection of GaLore doesn’t hurt much in the first iteration, it makes little sense for the following gradients, since the projection produces a really high relative error. Such a high relative error demonstrates a remarkably significant bias between the low-rank projected gradients and the original gradients, and between GaLore and the original gradient algorithm, highlighting the need to derive an unbiased low-rank projection algorithm.

<sup>5</sup>Note that while the projector is good for the stochastic gradient used in the algorithm, it can still be a large obstacle to the convergence, as shown in Figure 1.

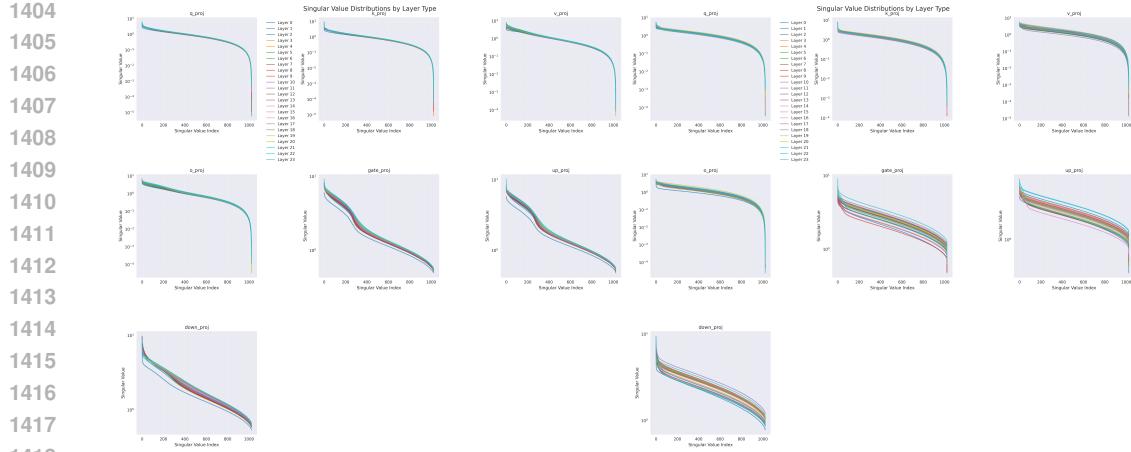


Figure 5: **Detailed Singular Value Distribution.** **Left:** GaLore. **Right:** GUM. It can be observed that GaLore has a sudden magnitude drop in the tail distribution of singular values in `gate_proj` and `up_proj` modules. GUM generally demonstrates smoother and more long-tailed singular value distributions. Furthermore, GUM has a differentiated spectrum across different layers, while this phenomenon is much weaker in GaLore.

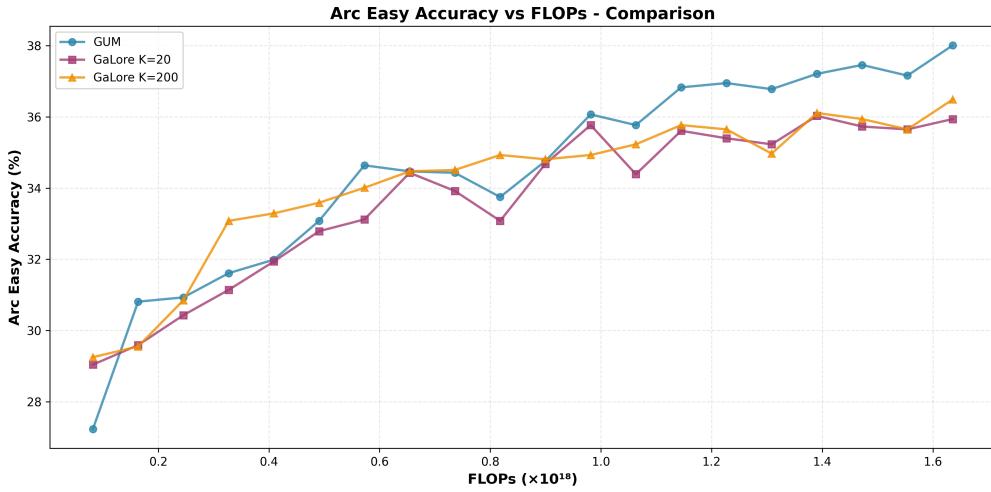


Figure 6: **Computational Cost Comparison.** The quality-vs.-time curve of GaLore with  $K = 20/200$  projector refreshing period and GUM.

## D.2 SINGULAR VALUE DISTRIBUTION OF MODEL WEIGHTS

As shown in Figure 5, GUM demonstrates a smoother and more long-tailed singular value distribution than GaLore, especially in modules of `gate_proj` and `up_proj`. The spectrums are also more differentiated and have a non-trivial diversity across layers in GUM.

## D.3 COMPUTATIONAL COST COMPARISON

To compare the computational cost between GaLore and GUM, additional experiments on LLaMA-130M are conducted, following the same setting in Section 5.3. Results of different projector refreshing periods  $K = 20/200$  for GaLore are also included to understand the effect of projector staleness. As shown in Figure 6, GUM is more computationally efficient than GaLore, and the projector refreshing period has little effect on GaLore’s computational efficiency.

1458 Table 8: **Ablation studies** for different choice of #sampled full-rank layers, sampling period  $K$ , and  
 1459 rank  $r$ . The default setup is #sampled full-rank layers = 4,  $K$ =200, and  $r$ =128.

#Sampled full-rank layers	Prompt-level Strict Accuracy	Prompt-level Loose Accuracy	$K$	Prompt-level Strict Accuracy	Prompt-level Loose Accuracy	$r$	Prompt-level Strict Accuracy	Prompt-level Loose Accuracy
1	30.87	34.20	20	29.21	33.83	32	30.98	<b>36.89</b>
2	31.61	35.12	100	30.68	35.21	64	31.24	36.60
4	33.27	36.60	200	<b>33.27</b>	<b>36.60</b>	128	<b>33.27</b>	36.60
6	<b>33.39</b>	<b>36.75</b>	500	29.39	33.27			
10	32.36	36.16						
20	28.36	30.76						

1467  
 1468 Table 9: **LLM Fine-tuning with More Baselines.** Trained models are evaluated on IFEval  
 1469 (instruction-following) and GSM8K (mathematical reasoning) with Qwen-2.5-7B model. All experi-  
 1470 ments are conducted on a single H100 GPU.

Method	IFEval		GSM8K
	Prompt-level Strict-Accuracy	Prompt-level Loose-Accuracy	Accuracy
LDAdamW (Robert et al., 2025)	28.10	30.31	83.78
Apollo (Zhu et al., 2025)	31.61	36.41	85.67
Subtrack++ (Rajabi et al., 2025)	29.76	34.01	86.66
<b>GUM</b>	<b>33.46</b>	<b>38.82</b>	<b>86.81</b>

#### D.4 ABLATION STUDIES

To better understand the tradeoff between sampling probability, sampling period, and ranks, additional ablation studies are conducted following the setting of Section 5.2. All experiments are conducted on IFEval benchmark with Gemma-2-9B.

As shown in the table above, the best choice of sampled layers is 6, where the performance starts to degrade when more full-rank layers are introduced. This is consistent with the observation in Pan et al. (2024), where this sampling-style training is conjectured to introduce an implicit regularization effect for supervised fine-tuning tasks.

For the best sampling period  $K$ , our choice of  $K = 200$  is already optimal, where a smaller  $K$  may compromise the momentum and decelerate the training process, while a larger  $K$  results in insufficient sampling of all layers.

For the best rank  $r$ , increasing the rank from 32 to 128 leads to overall performance improvements, especially in prompt-level strict accuracy. This means the higher-rank update captures more details for following the given instruction.

#### D.5 COMPARISON WITH MORE BASELINES

To further highlight GUM’s performance, we have included comparisons with LDAdamW (Robert et al., 2025), Apollo (Zhu et al., 2025), and SubTrack++ (Rajabi et al., 2025) on IFEval and GSM8K using the Qwen-2.5-7B model.

#### D.6 LARGER-SCALE PRE-TRAINING

To verify GUM’s effectiveness in pre-training on larger-sized model tasks, we conducted an additional pre-training experiment on a 7B-sized LLaMA model. Due to computational resource constraints, we follow the setup in SubTrack++ (Rajabi et al., 2025), and report results using the same number of tokens for pre-training.

As shown in the table, GUM outperforms GaLore, AdamW, and Fira in 7B-sized models as well, with better or no-worse performance on 4 out of 7 tasks, and a higher overall accuracy.

1512  
 1513 **Table 10: LLM Pre-training with 7B-sized LLaMA.** Trained models are evaluated on seven widely  
 1514 adopted commonsense reasoning tasks. All experiments are conducted on H100 GPUs.

Method	ARC-E	ARC-C	OBQA	HellaSwag	PIQA	SIQA	Winogrande	Avg.
GaLore	26.39	19.80	12.27	25.85	52.77	34.60	49.88	31.65
AdamW	26.68	18.94	13.48	26.11	51.69	34.54	50.70	31.73
Fira	26.46	20.90	13.01	25.65	52.34	34.75	50.54	31.95
<b>GUM</b>	<b>26.64</b>	<b>20.90</b>	<b>12.27</b>	<b>25.83</b>	<b>53.75</b>	<b>35.26</b>	<b>50.91</b>	<b>32.22</b>

## D.7 LONGER TRAINING

To verify GUM’s effectiveness in longer training scenarios, we extended the number of epochs and conducted additional experiments on the IFEval benchmark with Gemma-2-9B.

1525  
 1526 **Table 11: LLM Fine-tuning with Longer Training.** Trained models are evaluated on IFEval  
 1527 (instruction-following) with Gemma-2-9B model. All experiments are conducted on a single H100  
 1528 GPU.

#Epoch	IFEval	
	Prompt-level	Prompt-level
	Strict-Accuracy	Loose-Accuracy
1	33.27	36.60
3	29.57	31.42
5	26.25	28.28

As shown in Table 11, the performance degrades with an increasing number of epochs, indicating overfitting. So the number of epochs is sufficient for this supervised fine-tuning setting.

To further investigate GUM’s performance in effectively longer training settings, we conducted additional experiments on LLaMA-60M in Table 4, increasing the data amount to 5B tokens (originally 1.5B).

As shown in Table 12, GUM still outperforms Fira and GaLore by a non-trivial margin, demonstrating the effectiveness of GUM under longer training settings.

## E FURTHER DISCUSSION

### E.1 RELATIONSHIP WITH MUON OPTIMIZER

It is worth noticing that main focus of GUM is not the Muon optimizer (Jordan et al., 2024), but a technique for debiasing existing low-rank training methods like GaLore (Zhao et al., 2024), which is empirically orthogonal to the underlying optimizers such as AdamW and Muon.

Regarding Muon’s properties, there are several points worth mentioning:

- There is a fundamental tradeoff between Muon and AdamW across different model sizes.
  - Generally, Muon favors *deep and thin* networks, while AdamW has memory advantages in *large-scale wide* networks. On one hand, Muon may incur higher memory cost for extremely large hidden layers, since Muon requires matrix–matrix multiplication in the Newton–Schulz5 update, whereas AdamW only requires matrix–vector multiplication operations. On the other hand, Muon has only one momentum term, while AdamW has an additional second moment, which incurs extra memory consumption.
  - For commonly used  $\sim$ 7B-sized models like Gemma-2-9B, AdamW empirically requires more GPU memory, as shown in Table 2, where AdamW triggers an out-of-memory error while Muon does not.

1566  
 1567 **Table 12: LLM Pre-training with LLaMA-60M with long training (1.5B → 5B tokens).** Trained  
 1568 models are evaluated on seven widely adopted commonsense reasoning tasks. All experiments are  
 1569 conducted on H100 GPUs.

Method	ARC-E	ARC-C	OBQA	HellaSwag	PIQA	SIQA	Winogrande	Avg.
GaLore	37.50	18.60	11.07	27.07	60.66	37.56	49.57	34.58
Fira	37.16	17.66	14.29	27.38	60.72	37.62	50.83	35.09
<b>GUM</b>	<b>37.79</b>	<b>17.75</b>	<b>14.29</b>	<b>27.30</b>	<b>61.26</b>	<b>37.81</b>	<b>51.67</b>	<b>35.41</b>

1570  
 1571  
 1572  
 1573  
 1574  
 1575  
 1576  
 1577 • Muon has been successfully applied to Mixture-of-Experts training (Kimi, 2025) and  
 1578 outperforms AdamW, as shown in (Liu et al., 2025a).  
 1579 • **Why do we choose Muon as the base optimizer in GUM?** Muon performs well in large  
 1580 LLMs, as demonstrated by Kimi K2 (Kimi, 2025), which has 1T total parameters and 32B  
 1581 activated parameters. On the empirical side, Muon is demonstrated to perform better than  
 1582 AdamW in pre-training tasks (Liu et al., 2025a; Kimi, 2025). In addition, Muon incurs less  
 1583 memory consumption for common ~7B-sized models since it has no second moments. On  
 1584 the theoretical side, Muon satisfies properties I and II in Lemma 2, allowing the unbiasedness  
 1585 to be proven.

## F BROADER IMPACTS

1586  
 1587 Memory-efficient training techniques are critical for scalable LLM development and for democratizing  
 1588 customized LLMs for broader societal use. Improving theoretical guarantees provides insights for the  
 1589 invention of new methods with enhanced performance, leading to reduced computational resource  
 1590 consumption and lower carbon dioxide emissions.

## G LIMITATIONS

1591 The technique of sampled high-rank updates inherently introduces high variance into the per-iteration  
 1592 updates when the sampling probability is low, which leads to instability in the training procedure  
 1593 and requires more careful tuning of the hyperparameters. To alleviate this issue, standard theoretical  
 1594 tools for variance reduction can be employed (Johnson & Zhang, 2013; Needell et al., 2014; Ge  
 1595 et al., 2019b), which we leave for future work here. The analysis can also be combined with other  
 1596 acceleration (Zhang & Xiao, 2017; Ge et al., 2019a; Pan et al., 2021; 2023; Defazio et al., 2024; Liu  
 1597 et al., 2025b) and generalization techniques (Arjovsky et al., 2019; Foret et al., 2020; Hao et al., 2025),  
 1598 whose properties are worth investigating as open problems. The algorithm’s empirical performance  
 1599 and computational cost in other types of models (Devlin et al., 2019; Rombach et al., 2022; Pan  
 1600 et al., 2022; Liu et al., 2023; Gu & Dao, 2023; Hu et al., 2024; Wang et al., 2025; Mu & Lin, 2025)  
 1601 and applications (Xia et al., 2023; Peebles & Xie, 2023; Pan et al., 2025) also remain as interesting  
 1602 questions.

## H THE USE OF LARGE LANGUAGE MODELS

1603 ChatGPT and GPT-5 were adopted to polish the writing of the paper, where all revised sentences  
 1604 were double-checked by the authors. OpenAI Deep Research was utilized for finding dataset licenses.

## I LICENSES

1605 For mathematical reasoning tasks in LLM fine-tuning, the training dataset comes from 4 different  
 1606 sources: DART-Math (Tong et al., 2024), UltraInteract (Yuan et al., 2024), MathInstruct (Yue et al.,  
 1607 2023), and Orca-Math (Mitra et al., 2024), with their licenses listed in Table 13. Other datasets and  
 1608 benchmarks are also available in the same table.

1620	Training Datasets	#Samples	Kind	License
1621	teknium/GPT4-LLM-Cleaned <sup>6</sup>	55K	Instruction	CC BY-NC 4.0
1622	DART-Math <sup>7</sup> (Tong et al., 2024)	591K	Math	MIT
1623	openbmb/UltraInteract_sft <sup>8</sup> (Yuan et al., 2024)	289K	Reasoning	MIT
1624	TIGER-Lab/MathInstruct <sup>9</sup> (Yue et al., 2023)	262K	Reasoning	MIT
1625	microsoft/orca-math-word-problems-200k <sup>10</sup> (Mitra et al., 2024)	200K	Math	MIT
1626	C4 corpus <sup>11</sup> (Raffel et al., 2023)	>1B	Commonsense	ODC-BY
1627	IFEval (Zhou et al., 2023)	0.5K	Instruction	Apache-2.0
1628	GSM8K (Cobbe et al., 2021b)	7.5K	Math	MIT
1629	ARC-E (Clark et al., 2018)	5.2K	Instruction	CC-BY-SA-4.0
1630	ARC-C (Clark et al., 2018)	2.6K	Instruction	CC-BY-SA-4.0
1631	HellaSwag (Zellers et al., 2019b)	10K	Commonsense reasoning	MIT
1632	PIQA <sup>12</sup> (Bisk et al., 2020)	3K	Commonsense reasoning	Academic Free License v. 3.0
1633	SIQA (Sap et al., 2019)	2.2K	Commonsense reasoning	CC-BY-4.0 (Li et al., 2024)
1634	Winogrande <sup>13</sup> (Sakaguchi et al., 2021b)	1.8K	Commonsense reasoning	CC-BY
1635	OBQA (Mihaylov et al., 2018a)	5.9K	Commonsense reasoning	(permissive open license) <sup>14</sup>

Table 13: Licenses of training datasets and benchmarks. Here, the number of samples for benchmarks only counts the test set.

For code repositories, LMFlow (Diao et al., 2023) is released under Apache-2.0 license.

<sup>6</sup><https://huggingface.co/datasets/teknium/GPT4-LLM-Cleaned>

<sup>7</sup><https://huggingface.co/datasets/hkust-nlp/dart-math-uniform>

<sup>8</sup>[https://huggingface.co/datasets/openbmb/UltraInteract\\_sft](https://huggingface.co/datasets/openbmb/UltraInteract_sft)

<sup>9</sup><https://huggingface.co/datasets/TIGER-Lab/MathInstruct>

<sup>10</sup><https://huggingface.co/datasets/microsoft/orca-math-word-problems-200k>

<sup>11</sup><https://huggingface.co/datasets/allenai/c4>

<sup>12</sup><https://github.com/ybisk/ybisk.github.io/tree/master/piqa>

<sup>13</sup><https://github.com/allenai/winogrande>

<sup>14</sup>The OpenBookQA dataset is released under a permissive open license, making it freely available for academic research. In practice, sources indicate that the dataset is in the public domain or under a very permissive license. For example, a Kaggle distribution of OpenBookQA explicitly labels it CC0 1.0 Universal (Public Domain) (<https://www.kaggle.com/datasets/thedevastator/openbookqa-a-new-dataset-for-advanced-question-a>). Similarly, a curated dataset list reports OpenBookQA's license as Apache 2.0 (<https://github.com/lmmlzn/Awesome-LLMs-Datasets>). Both of these licenses allow unrestricted use, redistribution, and modification of the data, including for academic purposes.