## 000 FIRA: CAN WE ACHIEVE FULL-RANK TRAINING OF 001 LLMs UNDER LOW-RANK CONSTRAINT? 002 003

Anonymous authors

004

010 011

013

014

015

017

018

019

021

023

025

026

027

028

029

030

031

032

034

Paper under double-blind review

# ABSTRACT

Low-rank training has emerged as a promising approach for reducing memory 012 usage in training Large Language Models (LLMs). Previous methods either rely on decomposing weight matrices (e.g., LoRA), or seek to decompose gradient matrices (e.g., GaLore) to ensure reduced memory consumption. However, both of them constrain the training in a low-rank subspace, thus inevitably leading to sub-optimal performance. This raises a question: whether it is possible to con-016 sistently preserve the low-rank constraint for memory efficiency, while achieving full-rank training (i.e., training with full-rank gradients of full-rank weights) to avoid inferior outcomes? In this paper, we propose a new plug-and-play training framework for LLMs called Fira, as the first attempt to achieve this goal. First, we observe an interesting phenomenon during LLM training: the scaling impact of adaptive optimizers (e.g., Adam) on the gradient norm remains similar from low-rank to full-rank training. Based on this observation, we propose a normbased scaling method, which utilizes the scaling impact of low-rank optimizers as substitutes for that of original full-rank optimizers to enable full-rank training. In this way, we can preserve the low-rank constraint in the optimizer while achieving full-rank training for better performance. Moreover, we find that there are sudden gradient rises during the optimization process, potentially causing loss spikes. To address this, we further put forward a norm-growth limiter to smooth the gradient via regulating the relative increase of gradient norms. Extensive experiments on the pre-training and fine-tuning of LLMs show that Fira outperforms both LoRA and GaLore, achieving performance that is comparable to or even better than fullrank training. For instance, our Fira can reduce the memory usage of optimizer states by 61.1%, while achieving improved performance for pre-training on the LLaMA 1B architecture. Notably, for pre-training on the LLaMA 7B architecture, our method uses an  $8 \times$  smaller rank than GaLore, yet outperforms it by a large margin.

037

## INTRODUCTION 1

In recent years, Large Language Models (LLMs) have achieved remarkable advancements in various 040 domains (Achiam et al., 2023; Sima et al., 2023; Feng et al., 2024). While the substantial increase in 041 model size contributes significantly to these advancements, it also introduces considerable memory 042 bottlenecks, especially for optimizer states (Zhao et al., 2024a). For instance, pre-training a LLaMA-043 7B model from scratch <sup>1</sup> requires at least 58 GB memory, allocated as follows: 14GB for loading 044 parameters, 14GB for weight gradients, 28GB for Adam (Kingma & Ba, 2014) optimizer states, and 045 2GB for activations (Zhao et al., 2024a). Notably, the optimizer states consume even more memory 046 than the parameters themselves. To address this, low-rank training has demonstrated its effectiveness 047 to reduce the memory usage of the optimizer states by conducting training in a low-rank subspace 048 (Zhao et al., 2024a; Hu et al., 2022).

049 The current low-rank training methods can be broadly divided into two categories: weight matrixbased and gradient matrix-based low-rank decomposition. For the weight matrix decomposition 051 methods, the most representative one is Low-Rank Adaptation (LoRA) (Hu et al., 2022), where 052 its basic idea is to use low-rank matrices as decomposed representations of the pre-trained weights

<sup>1</sup>Training the model with a single batch size and maximum sequence length of 2048 under BF16 precision.



Figure 1: This analyses three types of memory-efficient approaches at a macro level.

066 during training, as shown in Figure 1 (a). However, the optimization of LoRA is constrained in a 067 low-rank subspace of the weights. This will inevitably cause the reduction of representation capac-068 ity, leading to sub-optimal outcomes (Zhang et al., 2023b; Xia et al., 2024). Although the variant 069 ReLoRA (Lialin et al., 2024) attempts to extend the application of LoRA from fine-tuning to pretraining, by periodically updating high-rank weights with multiple low-rank updates. It still requires 071 full-rank weight training as a warm-up before low-rank training, thus rendering memory efficiency 072 unachievable (Zhao et al., 2024a). For the gradient matrix decomposition based methods, the typical one is the gradient low-rank projection (GaLore) proposed recently (Zhao et al., 2024a). In 073 contrast to LoRA, GaLore attempts to reduce the memory usage in optimizer states via decompos-074 ing the gradient matrix, as shown in Figure 1 (b). While GaLore supports the training of full-rank 075 weights, it leverages only low-rank gradients, restricting them to a low-rank subspace. Consequently, 076 any gradient information outside this subspace is lost, in contrast to training with full-rank gradi-077 ents. Note that since these methods constrain LLM training to a low-rank subspace, this inevitably 078 leads to sub-optimal results compared to full-rank training (i.e., training with full-rank gradients and 079 full-rank weights). This raises the question: Can we achieve full-rank training for LLMs while 080 consistently maintaining a low-rank constraint? 081

In light of this, we propose a new memory-efficient training framework for LLMs, called **Fira**, 082 which, to the best of our knowledge, is the first to achieve full-rank training while consistently 083 maintaining a low-rank constraint. To achieve this goal, a significant challenge is that the low-084 rank constraint makes it hard to preserve complete optimizer states (e.g., gradient momentum and 085 variance) of full-rank weights in the commonly-used adaptive optimizer (e.g., Adam). As a result, the adaptive optimizer fails to correct the full-rank raw gradient according to the optimizer states. 087 Without this correction, adaptive optimization algorithms would degrade into simple SGD, leading 088 to significantly reduced optimization performance (Kingma & Ba, 2014; Zhang et al., 2020). This point is further validated in Section 4.1 and Section 5.4. Fortunately, we observe an interesting phenomenon during LLM training: the scaling factor  $^2$  of the optimizer (e.g., Adam) is similar from 090 low-rank training to full-rank training. As illustrated in Figure 3, if we sort the weight matrices by 091 their average scaling factors, we can obtain a similar rank order. A detailed quantitative analysis of 092 this similarity is presented in Appendix A.5 and A.6.

Based on this observation, we put forward a norm-based scaling method that utilizes the scaling 094 factor of a weight matrix in low-rank training to replace the corresponding matrix's scaling factor in full-rank training. In this way, our scaling factor can also play a similar role in correcting the raw 096 gradient, as adaptive optimizers do. Therefore, we can enable full-rank training while preserving the low-rank constraint. Additionally, we observe sudden increases in the gradient during training, 098 which results in a spike in training loss (as depicted in Figure 4). This issue could lead to substantial 099 parameter updates, causing the loss function to reach a much higher value and undermining prior op-100 timization efforts (Goodfellow et al., 2016; Zhang et al., 2020). Despite the use of gradient clipping 101 techniques (Pascanu et al., 2013), this issue may not be adequately resolved, as shown in Figure 4. 102 To this end, we propose a norm-growth limiter, which aims to smooth the gradient by restricting 103 the magnitude of the gradient norm's increase. By employing our limiter, we adaptively convert 104 sudden rises in gradients into gradual increases, thereby facilitating a smooth update that mitigates 105 the problem of loss spikes.

106 107

The scaling factor  $\phi_t(R_t)$  is defined as  $\frac{||\psi(R_t)||}{||R_t||}$ , where  $||R_t||$  is the norm of the raw gradient,  $||\psi(R_t)||$  is the norm of the gradient corrected by the gradient correction function  $\psi$  of the optimizer (e.g., Adam).

Our main contributions can be summarized as follows:

- 1. We propose Fira, a plug-and-play memory-efficient training framework of LLMs, constituting the first attempt to enable full-rank training consistently under the low-rank constraint. We will release the source code and package of our Fira into a Python library for easy use.
- 2. We design two components in Fira: a norm-based scaling strategy that leverages the scaling effects of low-rank optimizers to facilitate full-rank training, and a norm-growth limiter to address the issue of loss spikes by limiting the growth of gradient norm.
- 3. Extensive experiments across various parameter counts (60M, 130M, 350M, 1B, 7B) validate the effectiveness of Fira in both pre-training and fine-tuning tasks. Our framework not only outperforms both LoRA and GaLore, but also achieves performance comparable to or better than full-rank training.
- 121 2 RELATED WORK

110

111

112

113 114

115 116

117

118

119 120

123 Low-rank Adaptation. Low-Rank Adaptation (LoRA) has been introduced by Hu et al. (2022) 124 as an efficient fine-tuning method for LLMs. The core idea of LoRA is to freeze the pre-trained 125 weights and introduce trainable low-rank matrices as decomposed representations of the pre-trained weights. In this way, the memory usage of training LLMs could be saved. Recently, a variety 126 of methods by extending LoRA have been proposed to further improve the performance (Zhang 127 et al., 2023c; Wen & Chaudhuri, 2023; Xia et al., 2024; Zhang et al., 2023b; Dettmers et al., 2024). 128 For instance, ReLoRA (Lialin et al., 2024) is proposed to extend the application of LoRA from fine-129 tuning to pre-training. However, it still requires full-rank warm-up training before low-rank training, 130 which prevents achieving memory efficiency. It is worth noting that while LoRA based methods 131 reduce memory usage by limiting training to a low-rank parameter subspace, they inevitably reduce 132 representation capacity (Xia et al., 2024). 133

Gradient Projection. Recent works (Zhang et al., 2023b; Xia et al., 2024; Valipour et al., 2022) 134 have indicated that LoRA may yield sub-optimal performance since its low-rank constraints in pa-135 rameters. Inspired by traditional projected gradient descent methods (Chen & Wainwright, 2015; 136 Chen et al., 2019), GaLore (Zhao et al., 2024a) has been proposed recently to mitigate this prob-137 lem. It enables full-parameter learning under low-rank constraints by projecting the gradient into 138 a low-rank subspace, reducing memory usage for optimizer states. However, while GaLore allows 139 memory-efficient full-parameter training, it confines the gradient to a low-rank subspace, discarding 140 the portion outside the subspace and resulting in significant information loss. Inspired by Galore, 141 several new memory-efficient methods have been developed. Flora (Hao et al., 2024) improves ef-142 ficiency by randomly generating projection matrices as an alternative to the SVD method. Besides, LISA (Pan et al., 2024) enhances memory efficiency by freezing certain layers during optimization. 143

144 System-Based Memory-Efficient Techniques. Many system-based techniques have been devel-145 oped to reduce memory usage in LLM training (Chen et al., 2016; Ren et al., 2021). However, most 146 of these methods achieve memory efficiency by compromising either time or precision. Gradient 147 checkpointing (Chen et al., 2016) is proposed to reduce memory usage by trading increased com-148 putational time for the re-computation of activations. Quantization (Dettmers et al., 2024) reduces memory consumption by using lower-bit data types, but at the cost of model precision. Memory 149 offloading (Zhang et al., 2023a; Ren et al., 2021) reduces GPU memory usage by using non-GPU 150 memory (e.g., CPU) as an extension. However, it introduces additional communication overhead, 151 such as CPU-GPU transfer time. It's important to note that our proposed method is complementary 152 to these approaches and can potentially be combined with them to further reduce memory usage. 153

154 155

156

157

**3** PRELIMINARIES

# 3.1 REGULAR FULL-RANK TRAINING

At time step t, we denote the full-rank weight matrix as  $W_t \in \mathbb{R}^{m \times n}$ . The full-rank gradient can be represented as  $G_t = \nabla_W f_t(W_t) \in \mathbb{R}^{m \times n}$ , where f is the objective function. Then the regular full-rank training can be expressed as follows:

$$W_{t+1} = W_t - \eta \psi_t(G_t),\tag{1}$$

162 where  $\eta$  is the learning rate, and  $\psi_t$  is the gradient correction function of the optimizer (for vanilla 163 SGD,  $\psi_t(G_t) = G_t$ . Instead of vanilla SGD, adaptive optimizers (e.g., Adam (Kingma & Ba, 164 2014), AdamW (Loshchilov & Hutter, 2019)) are usually employed to correct the raw gradient for improving the training performance. However, this typically requires additional memory for storing 165 166 optimizer states used in gradient correction. For instance, Adam (Kingma & Ba, 2014) requires storing the optimizer states M and V, which consume 2mn of memory. The gradient correction 167 process is as follows: 168

$$M_t = \beta_1 M_{t-1} + (1 - \beta_1) G_t, \tag{2}$$

$$V_t = \beta_2 V_{t-1} + (1 - \beta_2) G_t^2, \tag{3}$$

178

189

190

195

199 200

201

208

209

215

169

 $\psi_t(G_t) = \frac{\sqrt{1-\beta_2^t}}{1-\beta_1^t} \cdot \frac{M_t}{\sqrt{V_t}+\epsilon},$ (4)

173 where all matrix operations are element-wise.  $\beta_1$  and  $\beta_2$  are Adam's hyper-parameters, and  $\epsilon$  is a 174 small constant (e.g.,  $1 \times 10^{-8}$ ) used for numerical stability. Since this regular full-rank training 175 typically consumes a large amount of memory for training LLMs, many representative low-rank 176 training methods, e.g., LoRA (Hu et al., 2022) and Galore (Zhao et al., 2024a), have been proposed 177 to reduce memory usage in recent years.

## 179 3.2 LOW-RANK ADAPTATION

The basic idea behind LoRA (Hu et al., 2022) is to use low-rank matrices as decomposed repre-181 sentations of the pre-trained weights during training, in order to reduce memory usage. Formally, 182 LoRA freezes the full-rank weight matrix  $W_0 \in \mathbb{R}^{m \times n}$  and incorporates two low-rank matrices  $A_t$ 183 and  $B_t$  for training as: 184

$$W_t = W_0 + B_t A_t, \tag{5}$$

185 where  $B_t \in \mathbb{R}^{m \times r}$ ,  $A_t \in \mathbb{R}^{r \times n}$ , and the rank  $r < \min(m, n)$ . While LoRA reduces memory usage 186 by limiting training to a low-rank subspace of the weight, it inevitably diminishes the representation 187 capacity of the weight matrix  $W_t$ . 188

# 3.3 GRADIENT LOW-RANK PROJECTION

191 In contrast to LoRA, GaLore (Zhao et al., 2024a) utilizes a projection matrix  $P_t \in \mathbb{R}^{m \times r}$  to project the full-rank gradient  $G_t \in \mathbb{R}^{m \times n}$  to a low-rank gradient  $R_t = P_t^{\top} G_t \in \mathbb{R}^{r \times n}$   $(m \le n)^3$ . By doing 192 193 so, the memory usage of optimizer states could be reduced. The parameter update in GaLore can be formulated as: 194

$$W_{t+1} = W_t - \eta P_t \psi_t(R_t), \tag{6}$$

where the projection matrix  $P_t$  can be obtained through singular value decomposition (SVD) of  $G_t$ 196 and can be updated every T steps: 197

$$G_t = U\Sigma V^\top \approx \sum_{i=1}^r \sigma_i u_i v_i^\top, \quad P_t = [u_1, u_2, \dots, u_r], \tag{7}$$

where  $u_i$  is the *i*-th column vector of the left singular matrix U. By selecting the first r columns of matrix U that correspond to the largest singular values, the projection matrix  $P_t$  effectively cap-202 tures the most significant directions in the gradient space, leading to faster convergence (Zhao et al., 203 2024a). The optimal switching frequency T is usually set to be between 50 to 1000, and the addi-204 tional computational overhead introduced by SVD is negligible (< 10%), as stated in (Zhao et al., 205 2024a). Since Galore restricts the gradient in the low-rank subspace, the gradient information out-206 side this subspace is lost, leading to inferior performance. 207

## 4 **PROPOSED METHOD**

210 To achieve full-rank training under low-rank constraints, our framework, named Fira, consists of two 211 important components: (i) a norm-based scaling method, enabling full-rank training by leveraging 212 the scaling effects of adaptive optimizers; (ii) a norm-growth limiter, which restricts the growth of the 213 gradient norm to prevent spikes in training loss. Next, we will elaborate on these two components. 214

<sup>&</sup>lt;sup>3</sup>For simplicity, we assume  $m \le n$ , following (Zhao et al., 2024a). If m > n,  $R_t = G_t Q_t \in \mathbb{R}^{m \times r}$ ,  $Q_t \in \mathbb{R}^{m \times r}$  $\mathbb{R}^{n \times r}$ .

# 216 4.1 NORM-BASED SCALING

The low-rank constraint makes it challenging to record complete optimizer states for correcting raw gradients in full-rank training. Fortunately, we find an interesting phenomenon in LLM training: the scaling factor at the matrix level remains similar from low-rank training to full-rank training. Based on this observation, we propose a norm-based scaling strategy that approximately corrects the raw gradient, similar to adaptive optimizers, thereby enabling full-rank training.

Challenge Analysis. Given the difficulty of incor-224 porating trainable low-rank weights into LoRA to 225 achieve full-rank weight training (Zhao et al., 2024a), we focus on investigating how to achieve full-rank 226 gradient training by extending the gradient projection 227 method, Galore, in this paper. In GaLore, the projec-228 tion matrix  $P_t \in \mathbb{R}^{m \times r}$  projects the full-rank gradient 229  $G_t \in \mathbb{R}^{m \times n}$  of the full-rank weight  $W_t \in \mathbb{R}^{m \times n}$ , to 230 the low-rank subspace gradient  $\overline{R_t} = P_t^{\top} G_t \in \mathbb{R}^{r \times n}$ . 231 The gradient outside this subspace can be represented 232 as:  $(I - P_t P_t^{\top})G_t = G_t - P_t R_t$ . In other words, the 233 full-rank gradient  $G_t$  can be divided into two terms: 234  $P_t R_t$  and  $(G_t - P_t R_t)$ .

235 In GaLore, the optimizer states only store the infor-236 mation of  $R_t$  instead of  $G_t$  to realize the low-rank 237 constraint. The term of  $(G_t - P_t R_t)$  is directly dis-238 carded in Galore due to the lack of corresponding op-239 timizer states of  $G_t$  for correction in optimizers. This 240 would lead to significant information loss especially 241 when  $r \ll d_{model}$ , where  $d_{model} = \min(m, n)$  is the full-rank dimension of models (This point can be ver-242

247

248

258

259



Figure 2: Training loss of different methods for pre-training LLaMA 60M on C4 dataset with  $r/d_{model} = 16/256$  and T = 200.

ified in our experiment section, as illustrated in Figure 6. In Figure 6, the validation perplexity of GaLore significantly increases at r = 4 compared to r = 128 when  $d_{model} = 256$ , indicating a substantial loss of information and decreased training performance). Intuitively, to capture the information of  $(G_t - P_t R_t)$ , we can directly add it based on Eq. (6) as follows:

$$W_{t+1} = W_t - \eta P_t \psi_t(R_t) - \eta (G_t - P_t R_t).$$
(8)

We denote the update strategy in Eq. (8) as GaLore-add. However, as illustrated in Figure 2, GaLoreadd exhibits almost no improvement compared to updates using Eq. (6) in GaLore. This phenomenon primarily arises because the term of  $(G_t - P_t R_t)$  doesn't have corresponding optimizer states for gradient correction. As a result, the optimization of  $(G_t - P_t R_t)$  uses vanilla SGD, yielding sub-optimal outputs. Besides, in GaLore-add,  $P_t\psi_t(R_t)$  employs the Adam optimizer for training while  $(G_t - P_t R_t)$  employs vanilla SGD. This gradient misalignment may also account for the lack of noticeable improvement.

Similarity of Scaling Factor. To tackle this challenge, we propose the concept of the *scaling factor*,
 which is defined as follows:

$$\phi_t(R_t) = \frac{||\psi_t(R_t)||}{||R_t||},$$
(9)

where the scaling factor  $\phi_t$  represents the magnitude of the correction applied by the adaptive optimizer to the gradient norm. Based on the scaling factor  $\phi_t$ , we observe an interesting phenomenon during LLM training: the scaling factors at the matrix level exhibit a high degree of similarity between low-rank and full-rank training. As shown in Figure 3, sorting weight matrices by their average scaling factors results in an almost similar order.

Based on this observation, we can use the scaling factors in low-rank training to replace those in full-rank training, even though the absolute magnitude of the scaling factors may vary between low-rank and full-rank training. This is because the absolute magnitude can be regarded as the learning rate, to which adaptive optimizers (e.g., Adam) are not sensitive (Zhao et al., 2024b; Liu et al., 2019).
Instead, the discrepancies (or relative orders) in the correction magnitudes of different parameters (or weight matrices) are of greater significance.



Figure 3: Scaling factor  $\phi_t(R_t)$  of different weight matrices for pre-training LLaMA 60M on the C4 dataset for 10K steps with varying ranks. Here x-axis is the number of recorded steps and y-axis is the ID of the weight matrices. Each row is labeled according to the order of its average scaling factor value. Case study and lager-scale quantitative analysis are provided in Appendix A.5 and Appendix A.6, respectively.

**Norm-based Scaling.** Inspired by this, we propose a norm-based scaling method that utilizes the scaling factor of a weight matrix in low-rank training as a substitute for the corresponding factor in full-rank training:

$$W_{t+1} = W_t - \eta P_t \psi_t(R_t) - \eta \phi_t(R_t) (G_t - P_t R_t).$$
(10)

By Eq. (10), we can approximately correct  $(G_t - P_t R_t)$  as adaptive optimizers do, so as to achieve full-rank training under low-rank constraints.

Furthermore, we can use a more fine-grained average of the scaling factor in Eq.(10), by considering each column in the weight matrix:

$$\phi_t(R_t)_i = \frac{||\psi(R_{t,:,i})||}{||R_{t,:,i}||}, \quad i = 1, 2, \dots, n,$$
(11)

where  $R_{t,i}$  is the *i*-th column of  $R_t$ , and  $\phi_t(R_t)_i$  is the *i*-th scaling factor.

## 4.2 NORM-GROWTH LIMITER

282

283

284

285

286 287 288

289

290

291 292

295

296

301 302

303

We find that there are suddenly sharp increases of the gradient during training, which could introduce loss spikes. As shown in Figure 4, Fira-w.o.-limiter (our method without using the proposed norm-growth limiter) experiences spikes in both gradient norm and training loss. In this section, we analyze the reasons for this issue and propose a norm-growth limiter which transforms abrupt gradient spikes into gradual, smooth increases.

309 Loss Spike Analysis. There are two main reasons for the spikes: (i) Switching the projection matrix 310  $P_t$  in gradient projection methods would cause instability during training. As illustrated in Fig-311 ure 2, both GaLore and GaLore-add exhibit significant training loss spikes at integer multiples of 312 T (i.e., the frequency of switching the projection matrix  $P_t$ ). This instability occurs because, when 313 switching projection matrices  $P_t$ , the optimizer retains states linked to the previous matrix, while the 314 current input gradient uses a new projection matrix, leading to significant misalignment. Furthermore, as shown in Figure 2, GaLore-add also exhibits training spikes, reinforcing our earlier claim 315 that directly incorporating  $(G_t - P_t R_t)$  may introduce instability and hinder training; (ii) Maintain-316 ing the original direction of the raw gradient  $(G_t - P_t R_t)$  may be insufficient for handling the sharp 317 loss landscapes in LLM training, unlike Adam (Zhang et al., 2020). Due to space constraints, further 318 analysis is provided in Appendix A.7. 319

Addressing Loss Spikes. To address this issue, a straightforward solution is to use gradient clipping
 techniques (Pascanu et al., 2013) to avoid loss spikes. However, clipping based on the absolute norm
 of gradient matrices fails to account for significant differences between them, leading to sub-optimal
 results. This point can be also verified in Figure 4 and Table 4. To this end, we propose a norm growth limiter method that constrains the ratio of the current gradient norm to the previous step's



Figure 4: Training loss and gradient norm of three variants of Fira for pre-training LLaMA 60M on C4 dataset.

norm to a fixed ratio  $\gamma$  when the gradient norm increases:

if 
$$\frac{||S_t||}{||S_{t-1}||} > \gamma$$
 then  $S_t \leftarrow \frac{S_t}{||S_t||} \cdot \gamma ||S_{t-1}||,$  (12)

where  $\gamma$  is a threshold ensuring that the rate of gradient growth does not exceed this value.  $S_t =$  $\phi_t(R_t)(G_t - P_tR_t)$  is the corrected gradient by applying our norm-based scaling. This approach limits the magnitude of gradient norm increases, converting sudden spikes into gradual rises and thus preventing loss spikes. Moreover, by constraining the relative increase of each gradient matrix's norm, our method is more flexible than the absolute norm clipping. As illustrated in Figure 2 and Figure 4, Fira with our proposed limiter improves the optimization performance without significant spikes.

# Algorithm 1 Fira with Adam

338

339

340 341

342 343

344 345

346

347

348

349

350

351

352

353 354 **Input:**  $\eta$  : step size,  $\{\beta_1, \beta_2\}$  : decay rates,  $W \in \mathbb{R}^{m \times n}$  with  $m \leq n$  : weight matrices, r : 355 rank, T : switching frequency,  $\alpha$  : hyper-parameter of Galore,  $\gamma$  : limiter threshold. 356 **Output:**  $W_t$ : resulting weight matrix 1:  $M_0, V_0 \in \mathbb{R}^{r \times n} \leftarrow 0, 0$   $t \leftarrow 0$ 357 ▷ Initialize moving 1st, 2nd moment and step 358 2: repeat  $G_t \in \mathbb{R}^{m \times n} \leftarrow \nabla_W f_t(W_t)$ 3: ▷ Calculate full-rank gradients of full-rank weights 359 4: if  $t \mod T = 0$  then 360  $U, \Sigma, V^{\top} \leftarrow \text{SVD}(G_t) \quad P_t \leftarrow U[:,:r] \triangleright \text{Initialize the projection matrix every } T \text{ steps}$ 5: 361 else  $P_t \leftarrow P_{t-1}$ 6:  $\triangleright$  Reuse the previous projection matrix 362 7: end if 363  $R_t, S_t \leftarrow P_t^{\top} G_t, (I - P_t P_t^{\top}) G_t \triangleright$  Divide gradients into two terms by gradient projection 8: 364  $M_t \leftarrow \beta_1 M_{t-1} + (1 - \beta_1) R_t$ 9:  $\triangleright \psi_t(R_t)$ : Apply Adam with low-rank gradients  $R_t$ 365  $V_t \leftarrow \beta_2 V_{t-1} + (1 - \beta_2) \hat{R}_t^2$ 10: 
$$\begin{split} & v_t \leftarrow \rho_2 v_{t-1} + (1 - \rho_2) \mathcal{K}_t^{\tau} \\ & N_t \leftarrow \frac{\sqrt{1 - \beta_1^t}}{1 - \beta_1^t} \cdot \frac{M_t}{\sqrt{V_t + \epsilon}} \\ & \overline{K} \leftarrow [\frac{||N_t[:, 1]||}{||R_t[:, 1]|| + \epsilon}, \frac{||N_t[:, 2]||}{||R_t[:, 2]|| + \epsilon}, \cdots, \frac{||N_t[:, n]||}{||R_t[:, n]|| + \epsilon}] \\ & S_t \leftarrow [k_1 S_t[:, 1], k_2 S_t[:, 2], \cdots, k_n S_t[:, n]] \end{split}$$
366 367 11: 368 369 12: ▷ Norm-Based Scaling 370 13: 371  $S_t \leftarrow S_t \cdot \gamma / \max\{\frac{||S_t||}{||S_{t-1}|| + \epsilon}, \gamma\}$ 372 14: ▷ Norm-Growth Limiter 373  $\tilde{G}_t \leftarrow \alpha \cdot (P_t N_t + S_t)$ 15: Project back and complete full-rank gradients 374  $W_t \leftarrow W_{t-1} - \eta \cdot \tilde{G}_t \quad t \leftarrow t+1$ ▷ Update the weight matrix 16: 375 17: **until** convergence criteria met 376 18: return  $W_T$ 377

378	Table 2: Comparison of different methods for pre-training LLaMA models of various sizes on the
379	C4 dataset. We report validation perplexity $(\downarrow)$ with a memory estimate of total parameters and
380	optimizer states. Results and memory estimates of all baselines are taken from Zhao et al. (2024a).
381	r refers to the rank and $d_{model}$ is the full-rank dimension of models.

	60M	130M	350M	1B
Full-Rank	34.06 (0.36G)	25.08 (0.76G)	18.80 (2.06G)	15.56 (7.80G)
<b>Fira</b> GaLore	<b>31.06</b> (0.24G) 34.88 (0.24G)	<b>22.73</b> (0.52G) 25.36 (0.52G)	<b>16.85</b> (1.22G) 18.95 (1.22G)	<b>14.31</b> (4.38G) 15.64 (4.38G)
LoRA ReLoRA	34.99 (0.36G) 37.04 (0.36G)	33.92 (0.80G) 29.37 (0.80G)	25.58 (1.76G) 29.08 (1.76G)	19.21 (6.17G) 18.33 (6.17G)
$r/d_{model}$ Training Tokens	128 / 256 1.1B	256 / 768 2.2B	256 / 1024 6.4B	512 / 2048 13.1B

# 4.3 OVERALL ALGORITHM

We present the overall algorithm of 396 Fira with Adam in Algorithm 1. Our 397 main components, the norm-based 398 scaling method and the norm-growth limiter, are straightforward to im-399 plement, requiring only 3 additional 400 lines of code. Moreover, Fira is a 401 plug-and-play framework which can 402 be easily integrated into the train-403 ing process without requiring signif-404 icant modifications. The plug-and-405

Table 1: Comparison between Fira, GaLore, and LoRA. Denote  $W_t \in \mathbb{R}^{m \times n}$   $(m \leq n)$ , rank r.

	Fira GaLo		LoRA
Weights Optimizer States	$ \begin{array}{c} mn \\ mr + 2nr + 1 \end{array} $	$mn \ mr + 2nr$	$\frac{mn+mr+nr}{2mr+2nr}$
Full-Rank Gradients Full-Rank Weights Pre-Training Fine-Tuning		×	√ × ✓

play Pytorch-like pseudo-code of Fira is provided in Appendix A.4.

407 It's worth noting that Fira only introduces one parameter  $||S_{t-1}||$  for each weight matrix in the 408 optimizer state, which is negligible, as shown in Table 1. Besides, in addition to the original hyper-409 parameters of optimizers and gradient projection methods, Fira only adds one hyper-parameter  $\gamma$ 410 in the norm-growth limiter. The hyper-parameter  $\gamma$  is set to 1.01 across all experiments, which 411 consistently yields satisfactory results.

412

414 415

416

417 418

419 420

394

# 5 EXPERIMENTS

In this section, we validate the effectiveness of Fira in pre-training and fine-tuning tasks of LLMs. In our experiments, we denote our method using the strategy of Eq. (10) as Fira-matrix, and denote our method additionally using the column-wise strategy of Eq. (11) as Fira.

# 5.1 MEMORY-EFFICIENT PRE-TRAINING

**Experimental Setup.** We follow the settings in Galore (Zhao et al., 2024a) to conduct the pre-421 training experiments. We compare Fira with GaLore (Zhao et al., 2024a), LoRA (Hu et al., 2022), 422 ReLoRA (Lialin et al., 2024), and full-rank training baselines. Adam optimizer is used for training 423 all baselines and our method on the C4 dataset in the BF16 format. The settings of these baselines 424 can be found in Zhao et al. (2024a). The dataset C4 is a colossal, cleaned version of Common 425 Crawl's web crawl corpus, which is widely used in LLM pre-training (Raffel et al., 2020). Fol-426 lowing Zhao et al. (2024a), we utilize LLaMA-based architectures equipped with RMSNorm and 427 SwiGLU activations (Zhang & Sennrich, 2019; Shazeer, 2020; Touvron et al., 2023). As in Zhao 428 et al. (2024a), our training protocol excludes data repetition and spans a sufficiently large dataset, encompassing a diverse array of model sizes (60M, 130M, 350M, 1B). To guarantee a fair compar-429 ison, we employ the same learning rate 0.01 as used in GaLore and maintain the same rank r for 430 each model size. The detailed settings of pre-training are provided in Appendix A.2. We use 8 A100 431 80G GPUs to conduct pre-training experiments.

432 **Result Analysis.** As shown in Table 2, Fira consistently outperforms low-rank training baselines 433 by a large margin under the same rank constraint, and even surpasses full-rank training. Following 434 Zhao et al. (2024a), we estimate the memory reduction of the optimizer states via the same memory 435 estimation method introduced in Zhao et al. (2024a). From Table 2, our Fira saves 61.1% memory 436 usage of the optimizer states when pre-training the LLaMA 1B architecture compared to full-rank training, while Fira achieving better results. Compared to full-rank training, Fira's superior perfor-437 mance may be attributed to the following reason: the gradient direction in the norm-based scaling 438 method is determined by the current state, rather than by historical gradients in Adam. Therefore, 439 Fira introduces a higher degree of randomness in training, which can enhance the model's ability to 440 escape the local optima, leading to better training performance (Zhou et al., 2020). 441

## 5.2 SCALING UP TO LLAMA 7B PRE-TRAINING.

444 To validate the scalability of our method, we scale up 445 by pre-training the LLaMA 7B model with the full-446 rank dimension  $d_{\text{model}} = 4096$ . We compare Fira 447 with the GaLore baseline, which generally achieves 448 the best performance among low-rank training base-449 lines, as shown in Table 2. As illustrated in Figure 5, our method demonstrates a significant improvement 450 over GaLore for pre-training LLaMA 7B, while using 451 an  $8 \times$  smaller rank. This highlights Fira's effective-452 ness, suggesting it could be a viable solution for large-453 scale LLM pre-training. 454



## 454 455 456

## 5.3 Memory-Efficient Fine-Tuning

Figure 5: Pre-training LLaMA 7B with different methods on the C4 dataset.

457 Experimental Setup. Following Hu et al. (2023), we perform the fine-tuning task to compare Fira 458 with LoRA, GaLore, Flora, ReLoRA, Full-rank training, and other baseline methods, including 459 Prefix-tuning (Prefix) (Li & Liang, 2021), Series Adapter (Series) (Houlsby et al., 2019), and Paral-460 lel Adapter (Parallel) (He et al., 2021), on the LLaMA-7B model for commonsense reasoning tasks. 461 This task consists of eight sub-tasks, each with its own designated training and testing sets. Follow-462 ing the approach of Hu et al. (2023), we combine the training datasets from all eight sub-tasks into a unified training set, while evaluating each sub-task individually using its respective testing dataset. 463 In the fine-tuning task, the rank r is set to 32 and the learning rate is set to 1e-4. The detailed 464 settings of fine-tuning are provided in Appendix A.3. We adopt RTX 4090 GPUs for fine-tuning 465 experiments. 466

467 Result Analysis. As shown in Table 3, our Fira achieves the highest performance on 4 out of 8
468 datasets, demonstrating better or comparable performance compared to the baseline methods. No469 tably, GaLore struggles to adapt to the HellaSwag and WinoGrande datasets, resulting in a significant
470 decline in scores. In contrast, our Fira adapts to these tasks well and achieves the highest scores on
471 WinoGrande. In terms of memory efficiency, our method uses comparable or even less memory than
472 the low-rank training methods LoRA and GaLore. These results illustrate the effectiveness of our
473 method for the fine-tuning of LLMs.

Table 3: Accuracy (<sup>↑</sup>) of various fine-tuning methods on eight commonsense reasoning datasets with
LLaMA 7B. Results for all baseline methods, except GaLore, are taken from Hu et al. (2023).

Method	Memory	BoolQ	PIQA	SIQA	HellaSwag	WinoGrande	ARC-e	ARC-c	OBQA	Avg
Fira	14.26G	69.4	82.6	78.0	76.8	81.2	82.2	64.4	80.8	76.9
Prefix	14.05G	64.3	76.8	73.9	42.1	72.1	72.9	54.0	60.6	64.6
Series	14.42G	63.0	79.2	76.3	67.9	75.7	74.5	57.1	72.4	70.8
Parallel	15.49G	67.9	76.4	78.8	69.8	78.9	73.7	57.3	75.2	72.2
LoRA	14.35G	68.9	80.7	77.4	78.1	78.8	77.8	61.3	74.8	74.7
ReLoRA	A 14.35G	68.9	81.2	77.8	46.0	79.4	80.2	64.2	79.6	72.2
Flora	14.26G	50.1	77.5	74.2	53.8	45.5	79	64.6	74.8	64.9
GaLore	14.26G	69.5	82.0	75.1	32.2	18.0	80.7	65.8	78.0	62.7
Full-ran	k 42.00G	64.2	68.1	68.0	42.3	66.5	55.6	43.9	60.0	58.6

# 486 5.4 ABLATION STUDY

488 In this section, we conduct an ablation study to assess the effectiveness of each component in our 489 method. We adopt the same settings in Section 5.1 for pre-training the LLaMA 60M model. We 490 design four variants of our method for the ablation study: (1) Fira-w.o.-scaling: our Fira without 491 using the scaling factor to correct the gradient (i.e., setting  $\phi_t(R_t)$  to a fixed value of 1). (2) Fira-492 matrix: our Fira using the scaling factor at the matrix level instead of at the column level. (3) Firaw.o.-limiter: our Fira without using norm-growth limiter to avoid training loss spikes. (4) Fira-493 494 gradient-clipping: our Fira using gradient clipping to avoid loss spikes instead of our proposed norm-growth limiter. 495

496 Table 4 presents the results. It can be found that Fira 497 outperforms Fira-w.o.-scaling, thereby demonstrating 498 the effectiveness of our proposed norm-based scaling 499 method for gradient correction. This also suggests that directly incorporating the raw gradient outside the 500 subspace without correction will lead to sub-optimal 501 results. Besides, Fira yields better performance than 502 Fira-matrix, illustrating that a more fine-grained consideration of the scaling factor is beneficial. Further-504 more, Fira demonstrates improved performance over 505

Table 4: Ablation study on the C4 dataset.

Method	Perplexity (↓)
Fira-w.oscaling	37.06
Fira-matrix	31.52
Fira-w.olimiter	32.22
Fira-gradient-clipping	31.22
Fira	31.06

Fira-w.o.-limiter and Fira-gradient-clipping, indicating the effectiveness of our proposed normgrowth limiter in addressing the issue of training loss spikes.

508 509

510

526 527 528

529

6

CONCLUSION

# 5.5 PERFORMANCE UNDER VARYING RANKS

511 In this section, we illustrate the advantages of our Fira 512 over Galore under a lower rank. We adjust various 513 rank configurations within the set  $\{4, 16, 64, 128\}$  and  $d_{model} = 256$ , and then assess the performance of pre-514 training the LLaMA 60M model on the C4 dataset as 515 outlined in Section 5.1. The validation perplexity of 516 Fira and GaLore after 10K steps across different ranks 517 is depicted in Figure 6. From Figure 6, we can ob-518 serve that Fira consistently surpasses GaLore across 519 all rank configurations. Notably, even when the ranks 520 are set very low (4 and 16), Fira still achieves perfor-521 mance comparable to full-rank training. In contrast, 522 the performance of GaLore significantly declines in 523 these cases. These results highlight the superiority of 524 our proposed Fira at lower ranks and its effectiveness 525 in reducing memory usage.

70 Fira GaLore 60 -40 -416 64 128rank r

Figure 6: Validation perplexity of Fira and GaLore for varying ranks when pretraining LLaMA 60M on the C4 dataset with  $d_{model} = 256$ .

530 In this paper, we present a plug-and-play memory-efficient training framework for LLMs, called 531 Fira, as the first attempt to facilitate full-rank training consistently under low-rank constraints. First, 532 we find a notable phenomenon in LLM training: the scaling effect of adaptive optimizers on the gra-533 dient norm remains similar between low-rank and full-rank training. Building on this observation, 534 we propose a norm-based scaling method that applies the scaling effect of low-rank optimizers in place of full-rank optimizers to facilitate full-rank training. This allows us to maintain the low-rank 536 constraint within the optimizer while still benefiting from the advantages of full-rank training for 537 improved performance. Additionally, we observe there are sudden spikes in gradient values during optimization, which could lead to spikes in loss. To mitigate this, we propose a norm-growth limiter 538 that smooths gradients by regulating the relative increase in gradient norms. Extensive experiments in both pre-training and fine-tuning of LLMs demonstrate the effectiveness of our proposed Fira.

# 540 REFERENCES

547

567

568

569

- Hervé Abdi. The kendall rank correlation coefficient. *Encyclopedia of measurement and statistics*, 2:508–510, 2007.
- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical
  report. *arXiv preprint arXiv:2303.08774*, 2023.
- David Balduzzi, Marcus Frean, Lennox Leary, JP Lewis, Kurt Wan-Duo Ma, and Brian McWilliams. The shattered gradients problem: If resnets are the answer, then what is the question? In *International conference on machine learning*, pp. 342–350. PMLR, 2017.
- Han Chen, Garvesh Raskutti, and Ming Yuan. Non-convex projected gradient descent for general ized low-rank tensor regression. *Journal of Machine Learning Research*, 20(5):1–37, 2019.
- Tianqi Chen, Bing Xu, Chiyuan Zhang, and Carlos Guestrin. Training deep nets with sublinear memory cost. arXiv preprint arXiv:1604.06174, 2016.
- Yudong Chen and Martin J Wainwright. Fast low-rank estimation by projected gradient descent:
   General statistical and algorithmic guarantees. *arXiv preprint arXiv:1509.03025*, 2015.
- Tim Dettmers, Mike Lewis, Sam Shleifer, and Luke Zettlemoyer. 8-bit optimizers via block-wise
   quantization. In International Conference on Learning Representations, 2022. URL https:
   //openreview.net/forum?id=shpkpVXzo3h.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning of quantized llms. *Advances in Neural Information Processing Systems*, 36, 2024.
- Kaituo Feng, Changsheng Li, Xiaolu Zhang, JUN ZHOU, Ye Yuan, and Guoren Wang. Keypoint based progressive chain-of-thought distillation for llms. In *International Conference on Machine Learning*, 2024.
  - Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. *Deep learning*, volume 1. MIT Press, 2016.
- Yongchang Hao, Yanshuai Cao, and Lili Mou. Flora: Low-rank adapters are secretly gradient
   compressors. In *Forty-first International Conference on Machine Learning*, 2024.
- Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. Towards a unified view of parameter-efficient transfer learning. *arXiv preprint arXiv:2110.04366*, 2021.
- 575 Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for nlp.
   577 In *International conference on machine learning*, pp. 2790–2799. PMLR, 2019.
- Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022.
- Zhiqiang Hu, Lei Wang, Yihuai Lan, Wanyu Xu, Ee-Peng Lim, Lidong Bing, Xing Xu, Soujanya
  Poria, and Roy Ka-Wei Lee. Llm-adapters: An adapter family for parameter-efficient fine-tuning
  of large language models. In *The 2023 Conference on Empirical Methods in Natural Language Processing*, 2023.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Teven Le Scao, Thomas Wang, Daniel Hesslow, Stas Bekman, M Saiful Bari, Stella Biderman, Hady Elsahar, Niklas Muennighoff, Jason Phang, Ofir Press, et al. What language model to train if you have one million gpu hours? In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pp. 765–782, 2022.
- 593 Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. *arXiv* preprint arXiv:2101.00190, 2021.

613

620

638

594	Vladislav Lialin, Sherin Muckatira, Namrata Shivagunde, and Anna Rumshisky. Relora: High-
595	rank training through low-rank updates. In The Twelfth International Conference on Learning
596	Representations, 2024.
597	

- Liyuan Liu, Haoming Jiang, Pengcheng He, Weizhu Chen, Xiaodong Liu, Jianfeng Gao, and Jiawei
   Han. On the variance of the adaptive learning rate and beyond. In *International Conference on Learning Representations*, 2019.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019.
- Robert Nau. Forecasting with moving averages. *Fuqua School of Business, Duke University*, pp. 1–3, 2014.
- Rui Pan, Xiang Liu, Shizhe Diao, Renjie Pi, Jipeng Zhang, Chi Han, and Tong Zhang. LISA: Layerwise importance sampling for memory-efficient large language model fine-tuning. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=L8ifDX5XNq.
- Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. On the difficulty of training recurrent neural networks. In *International conference on machine learning*, pp. 1310–1318. Pmlr, 2013.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi
   Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text
   transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
- Jie Ren, Samyam Rajbhandari, Reza Yazdani Aminabadi, Olatunji Ruwase, Shuangyan Yang, Minjia Zhang, Dong Li, and Yuxiong He. {Zero-offload}: Democratizing {billion-scale} model
  training. In 2021 USENIX Annual Technical Conference (USENIX ATC 21), pp. 551–564, 2021.
- Philip Sedgwick. Spearman's rank correlation coefficient. *Bmj*, 349, 2014.
- Noam Shazeer. Glu variants improve transformer. *arXiv preprint arXiv:2002.05202*, 2020.
- Chonghao Sima, Katrin Renz, Kashyap Chitta, Li Chen, Hanxue Zhang, Chengen Xie, Ping Luo,
   Andreas Geiger, and Hongyang Li. Drivelm: Driving with graph visual question answering. *arXiv preprint arXiv:2312.14150*, 2023.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- Mojtaba Valipour, Mehdi Rezagholizadeh, Ivan Kobyzev, and Ali Ghodsi. Dylora: Parameter efficient tuning of pre-trained models using dynamic search-free low-rank adaptation. *arXiv preprint arXiv:2210.07558*, 2022.
- Yeming Wen and Swarat Chaudhuri. Batched low-rank adaptation of foundation models. *arXiv* preprint arXiv:2312.05677, 2023.
- <sup>637</sup> Kirk M Wolter and Kirk M Wolter. *Introduction to variance estimation*, volume 53. Springer, 2007.
- 639 Wenhan Xia, Chengwei Qin, and Elad Hazan. Chain of lora: Efficient fine-tuning of language 640 models via residual learning. *arXiv preprint arXiv:2401.04151*, 2024.
- Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu, Wendi Zheng, Xiao Xia, Weng Lam Tam, Zixuan Ma, Yufei Xue, Jidong Zhai, Wenguang Chen, Zhiyuan Liu, Peng Zhang, Yuxiao Dong, and Jie Tang. GLM-130b: An open bilingual pre-trained model. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=-Aw0rrrPUF.
- 647 Biao Zhang and Rico Sennrich. Root mean square layer normalization. Advances in Neural Information Processing Systems, 32, 2019.

648 649 650	Haoyang Zhang, Yirui Zhou, Yuqi Xue, Yiqi Liu, and Jian Huang. G10: Enabling an efficient unified gpu memory and storage architecture with smart tensor migrations. In <i>Proceedings of the 56th Annual IEEE/ACM International Symposium on Microarchitecture</i> , pp. 395–410, 2023a.
651 652 653 654	Jingzhao Zhang, Tianxing He, Suvrit Sra, and Ali Jadbabaie. Why gradient clipping accelerates training: A theoretical justification for adaptivity. In <i>International Conference on Learning Representations</i> , 2020.
655 656 657	Longteng Zhang, Lin Zhang, Shaohuai Shi, Xiaowen Chu, and Bo Li. Lora-fa: Memory-efficient low-rank adaptation for large language models fine-tuning. <i>arXiv preprint arXiv:2308.03303</i> , 2023b.
658 659 660 661	Qingru Zhang, Minshuo Chen, Alexander Bukharin, Nikos Karampatziakis, Pengcheng He, Yu Cheng, Weizhu Chen, and Tuo Zhao. Adalora: Adaptive budget allocation for parameter- efficient fine-tuning. <i>arXiv preprint arXiv:2303.10512</i> , 2023c.
662 663 664	Jiawei Zhao, Zhenyu Zhang, Beidi Chen, Zhangyang Wang, Anima Anandkumar, and Yuandong Tian. Galore: Memory-efficient llm training by gradient low-rank projection. In <i>International Conference on Machine Learning</i> , 2024a.
665 666	Rosie Zhao, Depen Morwani, David Brandfonbrener, Nikhil Vyas, and Sham Kakade. Deconstruct- ing what makes a good optimizer for language models. <i>arXiv preprint arXiv:2407.07972</i> , 2024b.
668 669 670 671	Pan Zhou, Jiashi Feng, Chao Ma, Caiming Xiong, Steven Chu Hong Hoi, et al. Towards theoreti- cally understanding why sgd generalizes better than adam in deep learning. <i>Advances in Neural</i> <i>Information Processing Systems</i> , 33:21285–21296, 2020.
672	
673	
674	
675	
676	
677	
678	
620	
691	
682	
683	
684	
685	
686	
687	
688	
689	
690	
691	
692	
693	
694	
695	
696	
697	
698	
699	
700	

### APPENDIX А

### THEORETICAL ANALYSIS A.1

Error Upper Bound for the Approximation of Scaling Factors. In Section 4.1, we use the scaling factors of low-rank gradients to approximate that of full-rank gradients. To quantify the effectiveness of this approximation, we will derive its error upper bound theoretically, and verify our analysis experimentally. The error of the approximation  $\kappa(r)$  can be written as:

$$\kappa(r) = \left|\phi_t^2(G_t) - \phi_t^2(R_t)\right|,$$
(13)

where rank  $r \leq n$ , and  $R_t = P_t^{\top} G_t \in \mathbb{R}^{r \times n}$ . To simplify the proof, we consider that r components  $(g_1, \ldots, g_r)$  of low-rank gradients are directly sampled from n components  $(g_1, \ldots, g_n)$  of full-rank gradients. Under these conditions, the error can be rewritten as: 

$$\kappa(r) = \left| \frac{\sum_{i=1}^{n} \psi_i^2(g_i)}{\sum_{i=1}^{n} g_i^2} - \frac{\sum_{i=1}^{r} \psi_i^2(g_i)}{\sum_{i=1}^{r} g_i^2} \right|.$$
(14)

Assumption 1. (Bounded Scaling Factors): we assume that the adaptive optimizer scales each gradient component  $q_i$  by the scaling factor that lies within known bounds. Specifically, there exist constants  $c_{\min}$  and  $c_{\max}$  such that for all *i*: 

$$c_{\min} \le \left| \frac{\psi_i(g_i)}{g_i} \right| \le c_{\max}.$$
 (15)

This implies:

$$c_{\min}^2 \le \frac{\psi_i^2(g_i)}{g_i^2} \le c_{\max}^2.$$
 (16)

**Theorem 1.** (Error Upper Bound for Approximation) Under the assumption that  $\frac{\psi_i^2(g_i)}{g_i^2}$  are bounded between constants  $c_{\min}^2$  and  $c_{\max}^2$  for all components *i*, the approximation error  $\kappa(r)$  satisfies:

$$\kappa(r) \le (c_{\max}^2 - c_{\min}^2) \cdot (1 - \frac{\sum_{i=1}^r g_i^2}{\sum_{i=1}^n g_i^2}).$$
(17)

**Proof.** We first define the following quantities: Total Gradient Norm  $s_n = \sum_{i=1}^n g_i^2$ , Partial Gradient Norm (first r components)  $s_r = \sum_{i=1}^r g_i^2$ , Remaining Gradient Norm  $s_{n-r} = s_n - s_r = \sum_{i=r+1}^n g_i^2$ , Total Corrected Gradient Norm  $S_n = \sum_{i=1}^n \psi_i^2(g_i)$ , Partial Adjusted Gradient Norm (first r components)  $S_r = \sum_{i=1}^r \psi_i^2(g_i)$ , Remaining Adjusted Gradient Norm  $S_{n-r} = S_n - S_r =$  $\sum_{i=r+1}^{n} \psi_i^2(g_i).$ 

Then, our estimation error  $\kappa(r)$  can be rewritten using these definitions: 

$$\kappa(r) = \left| \frac{S_n}{s_n} - \frac{S_r}{s_r} \right|. \tag{18}$$

First, we rewrite the estimation error  $\kappa(r)$  in terms of  $S_r$ ,  $S_{n-r}$ ,  $s_r$ , and  $s_{n-r}$ :

 $\kappa$ 

$$(r) = \left| \frac{S_r + S_{n-r}}{s_r + s_{n-r}} - \frac{S_r}{s_r} \right|.$$
(19)

Then, compute the difference in the numerator:

$$\kappa(r) = \left| \frac{(S_r + S_{n-r})s_r - S_r(s_r + s_{n-r})}{(s_r + s_{n-r})s_r} \right|.$$
(20)

Simplify the numerator, thus, the estimation error becomes: 

C e  $\kappa(r)$ 

$$r) = \frac{|S_{n-r}s_r - S_rs_{n-r}|}{s_n s_r}.$$
(21)

After that, we factor out  $s_{n-r}$ :

$$\kappa(r) = \frac{s_{n-r}}{s_n} \cdot \left| \frac{S_{n-r}}{s_{n-r}} - \frac{S_r}{s_r} \right|.$$
(22)

From our bounded assumption, we have:

$$c_{\min}^2 \le \frac{S_r}{s_r} \le c_{\max}^2$$
 and  $c_{\min}^2 \le \frac{S_{n-r}}{s_{n-r}} \le c_{\max}^2$ . (23)

Therefore, the maximum possible difference between  $\frac{S_{n-r}}{s_{n-r}}$  and  $\frac{S_r}{s_r}$  is:

$$\left|\frac{S_{n-r}}{s_{n-r}} - \frac{S_r}{s_r}\right| \le c_{\max}^2 - c_{\min}^2.$$
(24)

Finally, since  $\frac{s_{n-r}}{s_n} = (1 - \frac{\sum_{i=1}^r g_i^2}{\sum_{i=1}^n g_i^2})$ , the approximation error  $\kappa(r)$  is bounded above by:

$$\kappa(r) \le (c_{\max}^2 - c_{\min}^2) \cdot (1 - \frac{\sum_{i=1}^r g_i^2}{\sum_{i=1}^n g_i^2}).$$
(25)

From this theory, we can find that the error upper bound on the approximation of scaling factors is mainly determined by two aspects, and we can verify them experimentally:

- Variability of Scaling Factor  $(c_{\max}^2 c_{\min}^2)$ : This term represents the maximum variation in the scaling factors of different gradient components. For further validation, we designed *Fira-only-scaling*, a variant of Fira. It directly applies the low-rank scaling factors to the full-rank gradients by changing the Eq. (10) from  $W_{t+1} = W_t \eta P_t \psi_t(R_t) \eta \phi_t(R_t)(G_t P_tR_t)$  to  $W_{t+1} = W_t \eta \phi_t(R_t)G_t$ . In this way, we are able to exclude the influence of the original Adam term  $P_t\psi_t(R_t)$  and better analyze the effectiveness of our approximation. As shown in Table 5, Fira-only-scaling (column-level) gains better performance than Fira-only-scaling (matrix-level) for its more fine-grained consideration of the scaling factor, which also means a smaller maximum variation  $(c_{\max}^2 c_{\min}^2)$ .
- Effectiveness of Gradient Sampling  $(1 \sum_{i=1}^{r} g_i^2)$ : This term represents the proportion of the gradients norm contributed by the sampled low-rank r components from full-rank n components. As shown in Table 6, we conducted ablation experiments *Fira-only-scaling-w.o.-svd*, i.e., Fira-only-scaling without SVD in low-rank gradient sampling. As we can see, SVD is capable of sampling more prominent low-rank gradients, which leads to a reduction in the upper bound of error and enhanced performance. Similarly, as shown in Table 7, employing a higher rank enables the sampling of a greater proportion of the gradients norm, resulting in reducing error upper bound and improved performance.

Table 5: Ablation on the level of scaling factors for the variant Fira-only-scaling.

Level	Perplexity $(\downarrow)$
Column	31.68
Matrix	32.05

Table 6: Ablation on SVD for the variant Fira-only-scaling.

806 807	Method	Perplexity (↓)
808	Fira-only-scaling	31.68
809	Fira-only-scaling-w.osvd	32.22

Table 7:	Ablation	on rank	for the	variant	Fira-only	-scaling.

Rank	4	16	64	128
Perplexity $(\downarrow)$	35.91	32.90	31.93	31.68

837

838

839 840

841 842

858 859

860

861 862

863

810

811

812 813

817 Variance of Scaling Factors. The variance of adaptive learning rates is significantly elevated dur-818 ing the early stage of training, often necessitating a warm-up to mitigate this variance and stabilize 819 training (Liu et al., 2019). As illustrated in Figure 7, the scaling factor in Fira exhibits a similar 820 pattern, characterized by substantial variance during the early stage of training, which also neces-821 sitates a warm-up. However, the addition of an extra warm-up hyper-parameter for Fira would be 822 inefficient. Therefore, it is crucial to investigate whether the original warm-up would have mitigated the variance in Fira efficiently. In the subsequent theoretical analysis, we show that, during the early 823 training phase, the variance of the scaling factor of Fira is less than or equal to that of the adaptive 824 learning rate. This finding suggests that the existing warm-up strategy is sufficient to mitigate the 825 variance of Fira, thereby eliminating the need for an additional warm-up hyper-parameter. 826

Consider independent random vectors  $\{g^{(i)}\}_{i=1}^{n}$ , where each  $g^{(i)} = (g_{1}^{(i)}, g_{2}^{(i)}, \dots, g_{t}^{(i)})$ . Here, the superscript *i* indicates the index of the weight matrix to which the vector belongs, while the subscript *j* (where *j* ranges from 1 to *t*) denotes training iterations with each parameter. Following (Liu et al., 2019), we assume the adaptive learning rate of Adam  $\psi(.) = \sqrt{\frac{1-\beta_{2}^{t}}{(1-\beta_{2})\sum_{i=1}^{t}\beta_{2}^{t-i}g_{i}^{2}}},$ and  $g_{j}^{(i)} \sim \mathcal{N}(0, \sigma^{2})$  for all *i* and *j* in the early stage. Additionally, approximate the distribution of the exponential moving average as the distribution of the simple average,  $p(\psi(.)) =$  $p(\sqrt{\frac{1-\beta_{2}^{t}}{(1-\beta_{2})\sum_{i=1}^{t}\beta_{2}^{t-i}g_{i}^{2}}) \approx p(\sqrt{\frac{t}{\sum_{i=1}^{t}g_{i}^{2}}})$  (Nau, 2014), and then  $\psi^{2}(.) \sim$  Scale-inv- $\mathcal{X}^{2}(\rho, \frac{1}{\sigma^{2}})$ .

**Theorem 2.** (Variance of Scaling Factors) In the early stages of training, if  $\psi^2(\cdot) \sim$ Scale-inv- $\mathcal{X}^2(\rho, \frac{1}{\sigma^2})$ , and  $g_j^{(i)} \sim \mathcal{N}(0, \sigma^2)^4$  for all i, j, then for all  $\rho > 4$ , the scaling factor  $\phi^2 = \frac{\sum_{i=1}^n \psi_i^2(g_t^{(i)})^2}{\sum_{i=1}^n (g_t^{(i)})^2}$  satisfies  $\operatorname{Var}[\phi^2] \leq \operatorname{Var}[\psi^2]$ . If we approximate  $\sqrt{\psi^2}$  and  $\sqrt{\phi^2}$  to the first order, we have  $\operatorname{Var}[\phi] \leq \operatorname{Var}[\psi]$ .

*Proof.* We express  $\phi^2$  as a weighted sum:

$$\phi^2 = \sum_{i=1}^n w_i \psi_i^2,$$
(26)

where the weights are defined as:

$$w_i = \frac{(g_t^{(i)})^2}{\sum_{j=1}^n (g_t^{(j)})^2}.$$
(27)

Each  $w_i$  is a non-negative random variable satisfying  $\sum_{i=1}^{n} w_i = 1$ .

In the context of adaptive optimization algorithms like Adam, the squared gradients  $\psi_i^2$  accumulate information from past iterations to adapt the learning rate for each parameter. With  $\beta_2 = 0.999$ , the moving average of the squared gradients places significant weight on historical data, making  $\psi_i^2$ dependent mainly on past gradients, yielding:

$$\psi_i^2 \approx \psi_i^2(g_1^{(i)}, \dots, g_{t-1}^{(i)}).$$
 (28)

Since  $\psi_i^2$  primarily depend on past gradients  $g_1^{(i)}, \ldots, g_{t-1}^{(i)}$ , and  $w_i$  depend solely on the current gradients  $g_t^{(i)}$ , we can consider  $\psi_i^2$  and  $w_i$  to be independent random variables.

<sup>&</sup>lt;sup>4</sup>The assumption of a mean-zero normal distribution is valid at the outset of training, as the weights are sampled from normal distributions with a mean of zero (Balduzzi et al., 2017)



Figure 7: Scaling factor  $\phi_t(R_t)$  during the early stage of training (1K iterations of total 10K iterations).



Figure 8: The simulation of variance of the scaling factor Var[ $\phi$ ] across different rank settings. The adaptive learning rate  $\psi$  is equivalent to  $\phi$  when the rank equals 1.

Consequently, we can express the variance of  $\phi^2$  as:

$$\operatorname{Var}[\phi^2] = \operatorname{Var}\left[\sum_{i=1}^n w_i \psi_i^2\right].$$
(29)

Using the law of total variance, we have:

$$\operatorname{Var}[\phi^{2}] = \operatorname{E}\left[\operatorname{Var}\left[\sum_{i=1}^{n} w_{i}\psi_{i}^{2} \mid w_{1}, \dots, w_{n}\right]\right] + \operatorname{Var}\left(\operatorname{E}\left[\sum_{i=1}^{n} w_{i}\psi_{i}^{2} \mid w_{1}, \dots, w_{n}\right]\right).$$
(30)

Since  $\psi_i^2$  are independent of the  $w_i$ , we find:

$$\mathbf{E}\left[\sum_{i=1}^{n} w_i \psi_i^2 \mid w_1, \dots, w_n\right] = \mathbf{E}[\psi_i^2] \sum_{i=1}^{n} w_i = \mathbf{E}[\psi_i^2],$$
(31)

$$\operatorname{Var}\left[\sum_{i=1}^{n} w_{i}\psi_{i}^{2} \mid w_{1}, \dots, w_{n}\right] = \sum_{i=1}^{n} w_{i}^{2}\operatorname{Var}[\psi_{i}^{2}].$$
(32)

Thus, the variance simplifies to:

$$\operatorname{Var}[\phi^2] = \operatorname{Var}[\psi^2] \operatorname{E}\left[\sum_{i=1}^n w_i^2\right].$$
(33)

The second term,  $\operatorname{Var}\left(\operatorname{E}\left[\sum_{i=1}^{n} w_{i}\psi_{i}^{2} \mid w_{i}\right]\right)$ , is zero since  $\operatorname{E}[\phi^{2} \mid w_{i}] = \operatorname{E}[\psi_{i}^{2}]$  is constant.

Let  $X_i = (g_t^{(i)})^2$ , where each  $X_i \sim \sigma^2 \chi_1^2$ . Then, we can express the weights as:

$$w_i = \frac{X_i}{\sum_{j=1}^n X_j}.$$
(34)

Since  $X_i/\sigma^2 \sim \chi_1^2$ , and each  $w_i$  is the ratio of  $X_i$  to the sum of all  $X_j$ , the vector  $(w_1, \ldots, w_n)$ follows a Dirichlet distribution with parameters  $\alpha_i = \frac{\nu_i}{2} = \frac{1}{2}$ , where  $\nu_i = 1$  is the degrees of freedom of  $\chi_1^2$ .

911 For a Dirichlet distribution, the expected value of  $w_i^2$  is given by:

$$E[w_i^2] = \frac{\alpha_i(\alpha_i + 1)}{(\sum_{k=1}^n \alpha_k) (\sum_{k=1}^n \alpha_k + 1)}.$$
(35)

915 Substituting  $\alpha_i = \frac{1}{2}$  and  $\sum_{k=1}^n \alpha_k = \frac{n}{2}$  yields: 

$$\mathbf{E}[w_i^2] = \frac{\frac{1}{2} \cdot \frac{3}{2}}{\frac{n}{2} \cdot \left(\frac{n}{2} + 1\right)} = \frac{3}{4} \cdot \frac{4}{n(n+2)} = \frac{3}{n(n+2)}.$$
(36)

Thus, summing over all i gives: 919

$$E\left[\sum_{i=1}^{n} w_i^2\right] = n \cdot E[w_i^2] = \frac{3}{n+2}.$$
(37)

Finally, substituting this result back into the variance expression:

$$\operatorname{Var}[\phi^2] = \operatorname{Var}[\psi^2] \cdot \frac{3}{n+2}.$$
(38)

Since  $n \ge 1$ , it follows that:

$$\frac{3}{n+2} \le 1,\tag{39}$$

which implies:

$$\operatorname{Var}[\phi^2] \le \operatorname{Var}[\psi^2]. \tag{40}$$

Given  $\rho > 4$  and  $\psi^2(\cdot) \sim$  Scale-inv- $\mathcal{X}^2(\rho, \frac{1}{\sigma^2})$ , the variance of  $\psi^2(\cdot)$  exists (Liu et al., 2019).

Since  $\psi_i^2$  and  $w_i$  are independent and  $\sum_{i=1}^n E[w_i] = 1$ :

$$\mathbf{E}[\phi^2] = \mathbf{E}\left[\sum_{i=1}^n w_i \psi_i^2\right] = \sum_{i=1}^n \mathbf{E}[w_i] \, \mathbf{E}[\psi_i^2] = \mathbf{E}[\psi_i^2] \sum_{i=1}^n \mathbf{E}[w_i] = \mathbf{E}[\psi_i^2],\tag{41}$$

Thus, we have shown that:

$$\operatorname{Var}[\phi^2] \le \operatorname{Var}[\psi^2], \quad \text{and} \quad \operatorname{E}[\phi^2] = \operatorname{E}[\psi_i^2]. \tag{42}$$

Follow Liu et al. (2019), we approximate  $\sqrt{\psi^2}$  and  $\sqrt{\phi^2}$  to the first order (Wolter & Wolter, 2007)

$$\operatorname{Var}[\psi] \approx \frac{\operatorname{Var}[\psi^2]}{4\operatorname{E}[\psi^2]}, \quad \text{and} \quad \operatorname{Var}[\phi] \approx \frac{\operatorname{Var}[\phi^2]}{4\operatorname{E}[\phi^2]}.$$
 (43)

which implies:

$$\operatorname{Var}[\phi] \le \operatorname{Var}[\psi]. \tag{44}$$

To further examine our theorem, we conduct simulations to calculate the variance of the scaling factor  $\phi$  at ranks within the set  $\{1, 5, 10, 50, 100\}$ . The adaptive learning rate  $\psi$  is equivalent to that of  $\phi$  when the rank equals 1. As shown in Figure 8, the variance decreases as the rank increases, supporting our above theorem  $Var[\phi] \leq Var[\psi]$ . Furthermore, we observe a surprisingly large variance during the early stage, which corroborated our initial experiments. Consequently, we conclude that our method is efficient without requiring an additional warm-up.

# A.2 DETAILED PRE-TRAINING SETTING

This section provides an overview of the LLaMA architectures and the hyper-parameters employed during pre-training. To ensure a fair comparison, we adopt the same settings as Zhao et al. (2024a).
Table 8 presents the hyper-parameters of the LLaMA architectures across various sizes. For all architectures, we utilize a maximum sequence length of 256 and a batch size of 131K tokens. Furthermore, we implement a learning rate warm-up during the initial 10% of training steps and employ cosine annealing for the learning rate schedule, which decreases to 10% of the initial learning rate.

For all methods except Fira and GaLore, we tune the optimal learning rate from the set {0.01, 0.005, 0.001, 0.0005, 0.0001} across model sizes ranging from 60M to 1B, selecting their best validation perplexity to report. In contrast, both Fira and GaLore employ the same learning rate 0.01 and a subspace change frequency *T* of 200 without tuning. Additionally, the scale factor  $\alpha$  is considered a fractional learning rate (Zhao et al., 2024a). Furthermore, a relatively large learning rate may result in spikes of the training loss (Zhao et al., 2024a). To address this issue, for models with a size of less than 1B, we set  $\alpha$  to 0.25, while for models exceeding 1B, we adjust  $\alpha$  to 0.0625.

973	Table 8: Hyper-parameters of LLaMA architectures for pre-training.							
974		Demonstra	TEddam	Tutomo di sta	IIda	T	Ctores	Data Amanut (Talama)
975		Params	Hidden	Intermediate	Heads	Layers	Steps	Data Amount (Tokens)
970		60M	512	1376	8	8	10K	1.3 B
977		130M	768	2048	12	12	20K	2.6 B
970		350M	1024	2736	16	24	60K	7.8 B
979			2048	5461	24	32	100K	13.1 B
900		/ B	4090	11008	32	32	150K	19.7 В
080								
983		Tab	le 0. Hype	er-narameter co	nfiguratio	ns of fine-	tuning I	LaMA-7B for Fira
984		140	ne y. myp	r parameter con	iniguiatio	iis or inic	tuning D	
985				Hyper-par	ameters	Sett	ting	-
986				Rank r		3	$\frac{v}{r}$	_
987				$\alpha$		6	2 4	
988				Dropout		0.0	05	
989				Base optir	Base optimizer		mW	
990				LR	LR LR Scheduler		-4	
991				LR Sched			ear	
992				Batch size	Batch size		6	
993				warm-up S	warm-up Steps		00	
994				Epochs			3	
995				Where		Q,K,V,U	p,Down	_
996								
997		Deeler						
998	A.3	DETAIL	ED FINE-1	UNING SETTIN	NG			
999 1000 1001 1002 1003	We f signe hype	ine-tune the d for LLI r-paramete	he pre-trai M fine-tun er configui	ned LLaMA-7E ing, which incl ations.	B model f lude eight	for comme sub-tasks	onsense r s (Hu et	easoning tasks benchmark de- al., 2023). Table 9 shows the
1004	A.4	PLUG-A	ND-PLAY	FRAMEWORK I	FOR FIRA			
1005								
1006	Algo	rithm 2 P	lug-and-p	ay framework f	for Fira, P	ytorch-lik	e.	
1007	1: 1	or weight	in model.	parameters() <b>do</b>				
1008	2:	grad =	weight.gra	ad	1)			
1009	3:	sub_gra	ad, outer_g	rad = <b>project</b> (g	grad)	⊾ A dor	tivo anti	▷ Gradient projection.
1010	4: 5.	SUD_ad	apt = adag Jira – Fira	(sub grad sub	adapt ou	ter grad)	nive opu	Apply Fire to outer grad
1011	5. 6.	weight	undate –	nroiect hack(e	uh grad).	+ outer Fi	ra	⊳ full-rank training
1012	7:	weight	data += v	veight undate	<u>5</u> 1)			
1013	8: 6	end for						
1014								

972

# A.5 CASE QUANTITATIVE ANALYSIS OF SCALING FACTOR SIMILARITIES

1018 In this section, we analyze the similarities among the three ranking (i.e., order) sequences  $R_1$ ,  $R_2$ , 1019 and  $R_3$  depicted in Figure 3, which highlight the discrepancies in scaling factors across different weight matrices. Each ranking comprises ten distinct items. We will employ Kendall's Tau cor-1020 relation coefficient (Abdi, 2007) and Spearman's rank correlation coefficient (Sedgwick, 2014) to 1021 evaluate their concordance and divergence. 1022

1023 The three ranking sequences are defined as follows: 1024

1025

$$R_1 = (7, 6, 1, 2, 4, 8, 5, 10, 9, 3) \tag{45}$$

1000		
1026		
1027	$R_2 = (7, 8, 2, 1, 5, 4, 6, 10, 9, 3) \tag{4}$	<del>1</del> 6)
1028	$-\circ_2  (\cdot, \circ, -, -, \circ, -, \circ, -, \circ, \circ, \circ) $	,
1029		
1030		47
1031	$R_3 = (6, 8, 2, 1, 5, 4, 7, 10, 9, 3) \tag{4}$	+/)
1032		
1033	A.5.1 KENDALL'S TAU CORRELATION COEFFICIENT	
1034		т.
1035	Kendall's Tau, a non-parametric statistic, measures the ordinal association between two rankings	. It
1036	quantifies the degree to which the presence of one ranking implies a similar ranking in another. I	he
1037	formula for Kendall's Tau is given by:	
1038		
1039	C - D	
10/0	$\tau = \frac{1}{n(n-1)} \tag{4}$	18)
1040	2	
1041	,	
1042	where:	
1043	C is the number of concordent point, which are point of characteristic where the realized	far
1044 1045	both items agree,	101
1046	• D is the number of discordant pairs, where the ranks disagree,	
1047	• $n$ is the total number of observations	
1048		
1049	Kendall's Tau ranges from $-1$ to $+1$ with 1 indicating perfect agreement between the rankin	σς
1050	0 indicating no correlation and -1 indicating perfect disagreement. This measure is particula	rlv
1051	robust against outliers, making it useful for assessing the strength of relationships in ordinal data.	Â
1052	p-value, calculated to evaluate the statistical significance of the observed correlation, tests the n	ull
1053	hypothesis that there is no association between the two variables. A low p-value (typically less the	an
1054	0.05) suggests rejecting the null hypothesis, indicating that the observed correlation is statistica	llv
1055	significant.	2
1056		
1057	$\Delta$ 5.2 Sdeadman Rank Coddel ation Coefficient	
1058	A.J.2 DI LARMAN KANK CORRELATION COEFFICIENT	
1059	The Spearman rank correlation coefficient is another non-parametric measure that assesses	the
1060	strength and direction of association between two ranked variables. It is calculated as follows:	
1061	5	
1001		
1002	$o = 1 - \frac{6\sum d_i^2}{2} \tag{2}$	19)
1063	$\nu = 1$ $n(n^2 - 1)$ (-	. / )
1064		

1065 where:

1067 1068

1069

1077

1079

•  $d_i$  represents the differences between the ranks of each observation,

• *n* is the number of observations.

Spearman's coefficient also ranges from -1 to +1, with similar interpretations as Kendall's Tau. A coefficient of 1 indicates perfect positive correlation, -1 indicates perfect negative correlation, and 0 indicates no correlation. Spearman's method excels when analyzing datasets that fail to meet the normality assumptions requisite for parametric tests. To assess the significance of the Spearman correlation, a p-value is calculated alongside the coefficient. This p-value tests the null hypothesis of no correlation between the rankings of the variables. A small p-value (often less than 0.05) indicates a statistically significant correlation, providing strong evidence against the null hypothesis.

1078 A.5.3 RESULTS

The results of the Spearman and Kendall correlation coefficients are summarized in Table 10:

1081	Table 10: S	pearman a	nd Kendall co	rrelation co	efficients with	correspondin	ig p-values.
1082		Sample	Spear	man	Kend	lall	
1083		Sample	Coefficient	P-value	Coefficient	P-value	
1085		$R_1R_2$	0.8545	0.0016	0.7333	0.0022	
1086		$R_1R_3$	0.8303	0.0029	0.6889	0.0047	
1087		$R_2R_3$	0.9879	9.31e-08	0.9556	5.51e-06	
1088							
1089	A.5.4 ANALYSI	S AND CO	NCLUSIONS				
1090							
1092 1093	below:	alysis indi	cates strong p	ositive relat	tionships amoi	ig the sample	es, as summarized
1094 1095	• For $R_1 R_2$ 0.7333 (p	$_{2}$ : The Spectrum (2) = 0.0022	earman coeffic both indicati	ient is 0.854	45 (p = 0.0016 nd statistically	) and the Ker significant co	ndall coefficient is orrelations.
1096	• For B. B.	• The Spe	arman coeffic	$i_{\text{ent}} \circ f \cap S3$	n = 0.0020	) and the Ker	dall coefficient of
1097	0.6889 (p	0 = 0.0047	further confi	m significat	nt positive asso	ociations.	idan coefficient of
1098	• For $R_2 R_3$	: The Spe	arman coeffic	ient of 0.98	79 (n = 9.31e)	08) and the K	Kendall coefficient
1100	of 0.9556	p = 5.51	e-06) suggest	an almost p	perfect correlat	tion, accompa	anied by high sta-
1101	tistical sig	gnificance.					
1102	In concept on volu	na laga tha	n 0 05 indiaat	a atmom a are	danaa againat	tha null huna	thesis suggesting
1103	that the observed c	correlation	s are statistica	llv significa	nt and unlikely	v to have occi	urred by chance
1104		aa avhihit	strong comple	tions The	aanaistanay a	maga hath ag	under of enance.
1105	underscores the re	eliability o	of these finding	gs. suggest	ing robust rela	ationships the	at warrant further
1107	exploration.	j -				F	
1108							
1109							
1110	A.6 LARGE-SC	ALE QUAN	NTITATIVE AN	NALYSIS OF	SCALING FA	CTOR SIMIL	ARITIES
1112	In the previous sec	ction we d	conduct a case	auantitativ	e analysis of t	he similarity	of scaling factors
1113	between low-rank	and full-ra	ank training a	cross 10 we	ight matrices of	of the LLaM	A 60M model. To
1114	further substantiate	e our findi	ngs, we expan	nd the scope	e of our experi	ment to inclu	ide all matrices of
1115	LLaMA models ra	anging from	m 60M to 1B	. Additional	lly, we add va	lue-based me	trics of similarity
1116	original order-base	ed metrics	. We use the	same low-ra	ink setup as in	Table 2. Th	en. we train these
1118	models and assess	s the simil	arity of scalir	ng factors a	veraged over	10,000 steps.	Additionally, to
1119	evaluate the effect	iveness of	a column-lev	el fine-grai	ned strategy f	or scaling fac	ctors, we perform
1120	a column-level qua	antitative s	similarity anal	lysis. Due t	o the computation of the computa	tional challes	nges posed by the
1121	cally, in the LLaM	A 1B mod	lel, over 10,00	0 columns a	are sampled.		anarysis. Speem-
1122	Both Speerman K	andall an	d Deerson com	alation cost	fficients range	from 1 to 1	1 A coefficient of
1123	1 signifies a perfec	endan, and t positive	correlation. a	id -1 signifi	es a perfect ne	gative correla	ation. The p-value
1124	helps us determine	e whether	the observed	correlation	is statistically	significant c	or if it could have
1126	occurred by rando	m chance.	For instance	, a p-value	less than 0.05	means there	is less than a 5%
1127	Generally a p yel	e observed	1  correlation  h	hat a signif	chance if the	re was actual	lly no correlation.
1128	score also ranges	from -1 to	1. Vectors w	ith scores of	close to 1 are	very similar.	As for MSE. the
1129	smaller the value t	he higher	similarity. As	shown in Ta	ble 11 and 12,	we can obse	rve the significant
1130	similarity of scalin	g factors b	etween low-ra	ank and full-	-rank LLM trai	ining (all coe	fficient and cosine
1131	similarity close to	1, while p-	-value and MS	E close to 0	). Thus, it is li	kely that the o	observed behavior
1132	provides a robust e	experiment	tal basis for o	unresung ac	norm-based so	ange of scena caling in Fira	and helps explain
1133	its effectiveness.	permen		- proposed	iioiiii ouood ba		and helps explain

Table 10: Spearman and Kendall correlation coefficients with corresponding p-values.

Table 11: Spearman, Kendall, and Pearson correlation coefficients (p-values) at both the Matrix and Column levels for pre-training LLaMA models ranging from 60M to 1B parameters, averaged over 10,000 steps. 

Size	Matrix Level			Column Level			
5120	Spearman	Kendall	Pearson	Spearman	Kendall	Pearson	
60M	0.9972 (2e-62)	0.9662 (7e-26)	0.9891 (1e-46)	0.9372 (0.0)	0.7942 (0.0)	0.8723 (0.0)	
130M	0.9925 (2e-76)	0.9409 (9e-37)	0.9813 (2e-60)	0.8698 (0.0)	0.6830 (0.0)	0.7805 (0.0)	
350M	0.9770 (3e-113)	0.8848 (5e-65)	0.9766 (1e-112)	0.9091 (0.0)	0.7400 (0.0)	0.8272 (0.0)	
1B	0.9469 (1e-83)	0.8249 (1e-56)	0.9457 (6e-83)	0.8331 (0.0)	0.6513 (0.0)	0.8112 (0.0)	

Table 12: Cosine Similarity and MSE at both the Matrix and Column levels for pre-training LLaMA models ranging from 60M to 1B parameters, averaged over 10,000 steps.

Size	Matrix Leve	1	Column Level		
5120	Cosine Similarity	MSE	Cosine Similarity	MSE	
60M	0.9922	3e-04	0.9273	3e-05	
130M	0.9901	2e-04	0.9046	2e-05	
350M	0.9893	1e-04	0.9174	1e-05	
1B	0.9795	2e-04	0.9229	1e-05	

### A.7 ADDITIONAL ANALYSIS OF SPIKES

Maintaining the direction of the raw gradient without correction might be unable to effectively deal with the steep loss landscapes of LLM training like Adam (Zhang et al., 2020). The steep loss landscapes are likely to cause abrupt increases in raw gradients. When the raw gradients increase abruptly, the gradients' norm after norm-based scaling may also increase abruptly, as illustrated in Figure 4. This arises from the fact that the norm-based scaling method only adjusts the average gradient norm of the gradient at the matrix level, failing to make fine-grained adjustments to each parameter, unlike the optimizer Adam. As a result, a significant parameter update may occur, under-mining previous optimization efforts, i.e. training loss spikes (Goodfellow et al., 2016; Zhang et al., 2020).

There are a lot of methods also proposed to solve loss spiking and stabilize training, e.g., embed-ding normalization (Le Scao et al., 2022), gradient shrink (Zeng et al., 2023), tensor-wise scaling (Dettmers et al., 2022). However, it is crucial to clarify that our Fira does not conflict these stabiliza-tion methods. For instance, when training the LLaMA model with Fira, it inherently incorporates stabilization methods like RMSNorm. Our norm-growth limiter is mainly aimed at addressing the gradient stability capability that our norm-based scaling method lacks compared to Adam. As shown in Figure 9, when we directly use Adam to pre-train the llama model, there will be no loss spike. However, since Fira maintains the original direction of the raw gradient  $(G_t - P_t R_t)$ , similar to SGD, it may lack the capability to navigate the sharp loss landscapes in LLM training, thus leading to an additional loss spike. 

In addition, for more comprehensive comparisons of our norm-growth limiter, we design two ad-ditional gradient stabilization variants to solve the loss spike: Gradient Shrink  $(|S_t| = |S_t| \cdot \alpha +$  $|S_{t-1}| \cdot (1-\alpha)$ , and Tensor-Wise Scaling  $(|S_t| = |S_t| \cdot \alpha)$ , where  $S_t = \phi_t(R_t)(G_t - P_tR_t)$  is the corrected gradient by applying our norm-based scaling. As shown in Figure 10 and Table 13, Fira outperforms other gradient stabilization methods. For further analysis, Gradient Shrink fails to solve the loss spike, while Tensor-Wise Scaling solves the loss spike but led to sub-optimal results.



Table 13: Validation perplexity ( $\downarrow$ ) of Fira across different gradient stabilization methods.

memory-efficient training methods for both pre-training and fine-tuning. As illustrated in Tables 14 and 15, Fira achieves superior memory efficiency compared to full-rank training without significantly reducing throughput. Although Fira's throughput is slightly lower than that of other memoryefficient methods, it delivers exceptional performance. During pre-training, methods like LoRA necessitate maintaining higher-rank adapters compared to full-rank training. In practice, maintain-1229 ing these higher-rank adapters outweighs the benefits of fewer trainable parameters, thus leading to 1230 more memory and less throughput. Furthermore, since full fine-tuning of LLaMA 7B's memory re-1231 quirements exceeds the A100's 80GB capacity, we utilize DeepSeed's Zero2 technology to mitigate 1232 its memory usage.

1233 1234

1188

Table 14: Real memory usage and normalized throughput when pre-training LLaMA 1B on the C4 1236 dataset.

Method	Fira	Galore	Flora	LoRA	ReLoRA	Full-rank
Memory (GB)	54.6	54.6	54.5	59.0	59.0	58.5
Normalized Throughput (%)	94.2	95.9	95.9	67.4	67.4	100

1243	Table 15: Real memory usage and normalized throughput when fine-tuning LLaMA 7B on com-
1244	monsense reasoning datasets.
1245	

Method	Fira	Galore	Flora	LoRA	ReLoRA	Full-rank
Memory (GB) Normalized Throughput (%)	23.4 156.1	23.4 201.1	23.3 210.3	23.7 232.8	23.7 232.8	>80 100







Figure 13: Cosine similarity trends over training iterations ( $r/d_{model} = 128/256$ ).



Figure 12: Euclidean distance trends over training iterations  $(r/d_{model} = 16/256)$ .



Figure 14: Cosine similarity trends over training iterations  $(r/d_{model} = 16/256)$ .

# A.10 ADDITIONAL EXPERIMENTS ON THE SIMILARITY TRENDS.

We conduct additional experiments on similarity trends of scaling factors from the same full-rank
 Adam dynamic using two rank settings: 16/256 and 128/256.

1291 As shown in Figure 11, 12, 13, and 14, the similarity exhibits fluctuations during the initial training 1293 phase but achieves a relatively steady pattern with high similarity in the later iterations. Under lower 1294 rank setting  $r/d_{model} = 16/256$ , there is negligible reduction in similarity. Besides,  $\phi_t(G_t)$  and  $\phi_t(G_t - P_t R_t)$  demonstrate significantly higher similarity owing to their closely aligned dimensions.

This observation further validates that  $\phi_t(R_t) \approx \phi_t(G_t) \approx \phi_t(G_t - P_t R_t)$ .

#### ADDITIONAL COMPARISONS OF PERPLEXITY TRENDS. A.11

In this section, we compare the perplexity trends of Fira, Fira-only-scaling, SGD, and Adam. As illustrated in 15 (a), the performance of vanilla SGD is significantly inferior, highlighting its inad-equacy for directly training LLMs. As depicted in 15 (b) and (c), while Adam demonstrates faster convergence during the initial stages, both Fira and Fira-only-scaling achieve superior performance in the later stages. This maybe because Fira applies an adaptive strategy only at the matrix-level while maintaining the original gradient direction within a weight matrix. In this way, Fira may introduce a higher degree of randomness in training and a better ability to escape the local optima.



Figure 15: Comparisons of perplexity ( $\downarrow$ ) trends for pre-training LLaMA 60M on C4 dataset.

