

---

# 000 LORA MEETS RIEMANNION: MUON OPTIMIZER 001 FOR PARAMETRIZATION-INDEPENDENT LOW-RANK 002 ADAPTERS 003 004 005

006 **Anonymous authors**

007 Paper under double-blind review  
008  
009

## 010 ABSTRACT 011

012  
013 This work presents a novel, fully Riemannian framework for Low-Rank Adap-  
014 tation (LoRA) that geometrically treats low-rank adapters by optimizing them di-  
015 rectly on the fixed-rank manifold. This formulation eliminates the parametrization  
016 ambiguity present in standard Euclidean optimizers. Our framework integrates  
017 three key components to achieve this: (1) we derive Riemannion, a new Rie-  
018 mannian optimizer on the fixed-rank matrix manifold that generalizes the recently  
019 proposed Muon optimizer; (2) we develop a Riemannian gradient-informed LoRA  
020 initialization, and (3) we provide an efficient implementation without prominent  
021 overhead that uses automatic differentiation to compute arising geometric opera-  
022 tions while adhering to best practices in numerical linear algebra. Comprehensive  
023 experimental results on both LLM and diffusion model architectures demonstrate  
024 that our approach yields consistent and noticeable improvements in convergence  
025 speed and final task performance over both standard LoRA and its state-of-the-art  
026 modifications.

## 027 1 INTRODUCTION 028

029 Large language models (LLMs) have demonstrated remarkable capabilities across a wide range of  
030 natural language processing tasks Brown et al. (2020); Touvron et al. (2023a;b). However, the  
031 computational and storage costs associated with training and deploying such models at scale pose  
032 significant challenges. To reduce these costs, parameter-efficient fine-tuning techniques such as  
033 low-rank adaptation (LoRA) Hu et al. (2022) have emerged as a practical solution. LoRA enables  
034 efficient adaptation of pre-trained models by embedding learnable low-rank matrices into specific  
035 weight updates, allowing most of the original parameters to remain frozen. In particular, the main  
036 idea of LoRA is to fine-tune a pretrained model using a rank- $r$  correction matrix  $\Delta W$ :

$$037 W + \Delta W = W + AB^\top, \quad A \in \mathbb{R}^{m \times r}, \quad B \in \mathbb{R}^{n \times r},$$

038 where  $W$  remains constant during training and  $A, B$  are optimized via gradient-based optimization  
039 methods.  
040

041 Despite its efficiency, the dominant practice of optimizing the LoRA factors ( $A, B$ ) with Euclidean  
042 optimizers such as SGD (Robbins & Monro, 1951), Adam (Kingma & Ba, 2014), Adagrad (Duchi  
043 et al., 2011), RMSProp (Tieleman, 2012), etc. that misaligned with the geometry of the low-rank  
044 constraint. The same update  $\Delta W$  can be represented by infinitely many factorizations: for any  
045  $A \in \mathbb{R}^{m \times r}$ ,  $B \in \mathbb{R}^{n \times r}$  and any invertible matrix  $S \in \mathbb{R}^{r \times r}$ , we may write:

$$046 \Delta W = AB^\top = \tilde{A}\tilde{B}^\top, \quad \text{where } \tilde{A} = AS, \quad \tilde{B} = BS^{-\top}. \quad (1)$$

047  
048 Ideally, training should be *reparameterization (transformation) invariant*: the update to  $\Delta W$  must  
049 not depend on which factorization ( $A, B$ ) is used (Yen et al., 2024). Empirically, this lack of in-  
050 variance manifests as unbalanced learning where one factor dominates and the other stalls, fragile  
051 hyperparameter sensitivity, and path-dependent solutions. These issues have prompted geometry-  
052 aware formulations. Riemannian treatments of low-rank models operate on the *fixed-rank manifold*  
053 rather than the ambient factor space, projecting gradients to the tangent space and retracting back to  
the manifold. Such steps can be implemented efficiently when  $r \ll \min\{m, n\}$  and avoid forming

---

054 full-size matrices. Within the LoRA literature, existing Riemannian approaches either use stan-  
055 dard SGD-type optimizers Mo et al. (2025, LORO) or rely on the Adam (Zhang & Pilanci, 2024)  
056 optimizer for auxiliary matrices within the chosen parameterization, which deviate them from the  
057 Riemannian framework and introduce dependence on parameterization.

058 In this work, we introduce a *fully Riemannian* framework for training LoRA that optimizes the  
059 adapter  $X = \Delta W$  directly on the fixed-rank matrix manifold

$$060 \mathcal{M}_r = \{X \in \mathbb{R}^{m \times n} : \text{rank}(X) = r\},$$

062 eliminating factorization ambiguity by construction. Central to our approach is *Riemannion*, a new  
063 Riemannian optimizer on  $\mathcal{M}_r$  that *generalizes the recently proposed Muon optimizer* (Jordan et al.,  
064 2024) to the fixed-rank setting. In contrast to prior Riemannian LoRA variants that port Adam-  
065 like mechanics to the manifold with ad hoc choices, our design inherits Muon’s geometry-aligned  
066 normalization, yielding transformation invariance of the learned update. We further propose a Rie-  
067 mannian gradient-informed initialization that places the initial adapter at a good location on  $\mathcal{M}_r$ ,  
068 and we provide a practical, low-overhead implementation that assembles projections, retractions,  
069 and vector transports via automatic differentiation, also following best practices from numerical  
070 linear algebra. Extensive experiments on LLM and diffusion architectures show consistent gains  
071 in convergence speed and final task performance over standard LoRA and recent state-of-the-art  
072 modifications.

073 Our contributions are as follows:

- 074 • **Riemannion: Muon on the fixed-rank manifold.** We derive *Riemannion*, the first opti-  
075 mizer that *generalizes Muon* to the manifold  $\mathcal{M}_r$  of fixed rank matrices.
- 076 • **Riemannian gradient-informed initialization.** We propose an initialization strategy  
077 which yields best alignment between the initial Riemannian gradient and the Euclidean  
078 gradient. We also propose an efficient way for this strategy by using a randomized SVD  
079 algorithm with implicit matrix multiplication (Section 5). Finally, we show the connection  
080 of this initialization to LoRA-GA.
- 081 • **Efficient implementation with automatic differentiation.** We pay special attention to  
082 numerical implementation to make the method robust without any prominent overhead  
083 compared to vanilla LoRA at small ranks.
- 084 • **Comprehensive empirical validation.** We showcase the performance of our framework  
085 for fine-tuning LLMs and in subject-driven generation using diffusion models. Among  
086 positive effects that we observe are: boost in target metrics, improved convergence, and  
087 reduction of variance.

## 089 2 RELATED WORK

091 The problem of an optimal initial guess selection for low-rank LLM adaptation has been addressed  
092 in a sequence of works: the authors Meng et al. (2024, PiSSA) have suggested a heuristic that  
093 involves using a low-rank truncated SVD of pretrained parameters as an initial point for LoRA and  
094 its orthogonal complement as frozen layer’s parameters, so that the tuning process starts without  
095 changing the starting value of the loss function. A similar approach was implemented by Wang  
096 et al. (2025, MiLoRA) with the main difference of optimizing the smallest singular components of  
097 unadapted parameter matrix. A context-aware initialization was considered in (Yang et al., 2024,  
098 CorDA) and (Parkina & Rakhuba, 2025, COALA) proposes a numerically robust inversion-free  
099 framework for low-rank weighted approximations for this setting. Another idea is to initialize LoRA  
100 with a subset of left and right singular vectors of a doubled-rank truncated SVD of the loss function  
101 gradient at the starting parameters, proposed by Wang et al. (2024, LoRA-GA). We show direct  
102 connection of this method to our Riemannian initialization strategy and propose how to additionally  
103 significantly accelerate the computation of SVD using our approach. Attempting to overcome the  
104 asymmetry in the initialization of vanilla LoRA fine-tuning process, (Hayou et al., 2024, LoRA+)  
105 introduced a scale-free step size selection for LoRA factors.

106 Riemannian optimization is widely used for algorithms on matrix manifolds and allows for exploit-  
107 ing task geometry or imposing additional constraints. For example, a Riemannian solution for the  
extreme eigenpairs search problem was described in Absil et al. (2009); Baker (2008), a matrix

completion task, which is common in collaborative filtering for recommender systems, via optimization on the fixed-rank manifold (Vandereycken, 2013), a Riemannian approach on the manifold of matrices with orthonormal columns (the Stiefel manifold) was used by Wisdom et al. (2016) for diminishing the problem of vanishing and exploding gradients in recurrent neural networks, etc. The book Trendafilov & Gallo (2021) also presents a comprehensive description of useful manifolds for solutions of the data science problems. For the deeper understanding of applied differential geometry techniques see the books Absil et al. (2009) and Boumal (2023).

The idea of using Riemannian optimization has recently started to emerge for the large language models. For example, the fine-tuning of LLMs with the help of the Stiefel manifold was considered in the work Hu et al. (2024). The authors of (Zhang & Pilanci, 2024) introduced the Riemannian inspired modification of Adam. The authors of Mo et al. (2025, LORO) applied the Riemannian optimization techniques for pretraining LLMs on the fixed-rank manifold. Parametrization that is used in our work can potentially help in this setting as well, by additionally avoiding potential overheads and instabilities, arising due the explicit inversion of Gram matrices.

### 3 PRELIMINARIES

#### 3.1 MUON OPTIMIZER

*Muon* is an optimizer designed specifically for matrix-valued parameters in a network’s hidden layers. Empirically, it accelerates training on language and vision workloads while leaving scalar/vector parameters and the input/output layers to a conventional optimizer such as AdamW. At a high level, Muon takes the step that stochastic gradient descent with momentum (SGDM) would make on a weight matrix and *orthogonalizes* that update before applying it. Orthogonalization acts as a per-layer, per-step preconditioner that equalizes singular values of the update, which mitigates the collapse of updates into a few dominant directions (Jordan et al., 2024). More specifically, let  $W \in \mathbb{R}^{n \times m}$  be a hidden-layer weight. With gradient  $G_t = \nabla_W \mathcal{L}(W_t)$  and momentum  $M_t = \beta M_{t-1} + G_t$ , Muon computes

$$\widetilde{M}_t \approx \text{Ortho}(M_t) \quad \text{and} \quad W_{t+1} = W_t - \eta \widetilde{M}_t,$$

where  $\text{Ortho}(\cdot)$  denotes the nearest semi-orthogonal matrix in Frobenius norm, i.e.,

$$\text{Ortho}(G) = \arg \max_O \{ \|O - G\|_F : O^\top O = I \text{ or } O O^\top = I \}. \quad (2)$$

Computing  $\text{Ortho}(G)$  exactly amounts to taking the SVD  $G = U S V^\top$  and returning  $U V^\top$ , which is too slow to do at every iteration. Muon instead applies a Newton–Schulz (NS) iteration that—after normalizing  $G$ —implements a composition of a fixed low-degree polynomial in  $GG^\top$  acting on  $G$  and converges to  $U V^\top$ . Leveraging efficient matrix multiplication operations results in a highly performant iteration. We will write  $\widetilde{M}_t = NS(M_t)$  for short.

**LMO interpretation.** Muon’s step admits a clean linear minimization oracle (LMO) interpretation (Bernstein, 2025). Indeed, consider the operator-norm unit ball  $\mathcal{B}_2 = \{X : \|X\|_2 \leq 1\}$ . The linear minimization oracle (LMO) over  $\mathcal{B}_2$  at matrix  $M_t$  given by its SVD  $M_t = U \Sigma V^\top$  is

$$UV^\top \in \text{Arg max}_{\|S\|_2 \leq 1} \langle M_t, S \rangle. \quad (3)$$

**Applying Muon to LoRA.** In LoRA, a frozen weight  $W_0 \in \mathbb{R}^{n \times m}$  is adapted via a low-rank update  $W = W_0 + \alpha BA$  with  $B \in \mathbb{R}^{n \times r}$ ,  $A \in \mathbb{R}^{r \times m}$  and small  $r$ . Each trainable factor ( $A$  and  $B$ ) is a 2D parameter, so Muon can be applied *per factor*:

$$\begin{aligned} M_t^{(A)} &\leftarrow \beta M_{t-1}^{(A)} + \nabla_A \mathcal{L}, & A_{t+1} &\leftarrow A_t - \eta_A NS(M_t^{(A)}), \\ M_t^{(B)} &\leftarrow \beta M_{t-1}^{(B)} + \nabla_B \mathcal{L}, & B_{t+1} &\leftarrow B_t - \eta_B NS(M_t^{(B)}). \end{aligned}$$

Note that acting on the two factors separately makes Muon non-reparameterization-invariant: its per-factor orthogonalization depends on arbitrary scalings or rotations, skewing the weight-space step and often letting one factor dominate.

---

### 3.2 RIEMANNIAN OPTIMIZATION

Let  $\mathcal{M}_r = \{X \in \mathbb{R}^{m \times n} : \text{rank}(X) = r\} \subseteq \mathbb{R}^{m \times n}$  be a smooth manifold of fixed-rank matrices (Lee, 2003, Example 8.14). Let every point  $X$  of  $\mathcal{M}_r$  be equipped with a *tangent plane*  $\mathcal{T}_X \mathcal{M}_r$ . Thinking geometrically, the tangent plane plays the role of the best local, flat approximation to this curved set: if you "zoom in" at  $X$ , the manifold looks like a plane. We will now discuss how to numerically parametrize points on a manifold and its tangent plane. Every rank- $r$  matrix  $X \in \mathcal{M}_r$  can be represented using matrices  $A_L \in \mathbb{R}^{m \times r}$ ,  $B_R \in \mathbb{R}^{n \times r}$  with orthonormal columns and a square matrix  $T \in \mathbb{R}^{r \times r}$  as

$$X = A_L B^\top = A_L \underbrace{(T B_R^\top)}_B = \underbrace{(A_L T)}_A B_R^\top \quad (4)$$

For example, one may think of a thin SVD, in which case  $T$  becomes the diagonal matrix of singular values, but other representations are also possible and will be convenient for our purposes. **Note** that the last two representations in (4) coincide with the TT parametrization in two-dimensional case (Oseledets, 2011). Combining this representation with the derivations from (Holtz et al., 2012), one can obtain the following form of the tangent vector  $\xi \in \mathcal{T}_X \mathcal{M}_r$ :

$$\xi = [\dot{A} \quad A_L] [B_R \quad \dot{B}]^\top, \quad \dot{A}^\top A_L = 0, \quad (5)$$

so that any tangent vector can be identified in terms of the tuple:  $(\dot{A}, \dot{B})$ . Since the factor matrices contain  $2r$  columns, we immediately have that  $\text{rank } \xi \leq 2r$ . Another remarkable fact is that the point  $X \in \mathcal{M}_r$  itself lies in the  $\mathcal{T}_X \mathcal{M}_r$  with  $\dot{A} = 0, \dot{B} = B$ . Given a matrix  $Z \in \mathbb{R}^{m \times n}$ , its orthogonal projection  $P_{\mathcal{T}_X \mathcal{M}_r} Z$  onto the tangent space  $\mathcal{T}_X \mathcal{M}_r$  (with the parameterization given in (5)) can be computed as follows:

$$P_{\mathcal{T}_X \mathcal{M}_r}(Z) = A_L A_L^\top Z + (I - A_L A_L^\top) Z B_R B_R^\top. \quad (6)$$

and, hence, can be represented in the form of (5) with  $\dot{A} = (I - A_L A_L^\top) Z B_R$ , and  $\dot{B} = Z^\top A_L$ .

Let  $\mathcal{L}: \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$  be a differentiable function with the Euclidean gradient  $\nabla \mathcal{L} \in \mathbb{R}^{m \times n}$ . Within the Riemannian optimization framework, we solve the following task optimization problem:

$$\min_{X \in \mathcal{M}_r} \mathcal{L}(X).$$

When constructing the algorithms, we need to work with three key objects: Riemannian gradient, retraction and vector transport. The Euclidean gradient is a direction of the steepest local increase  $\mathcal{L}$ . Therefore, it is common to use the Riemannian gradient — the direction of the steepest local increase of corresponding smooth function value along the manifold, which lies in the tangent space (Absil et al., 2009, chap. 3.6). Given the Euclidean gradient  $\nabla \mathcal{L}$ , one may endow the tangent space  $\mathcal{T}_X \mathcal{M}_r$  with a natural scalar product and derive a formula for the direction of the local steepest ascent of  $\mathcal{L}$  with respect to the manifold. This unique direction is called the Riemannian gradient and can be computed as follows:

$$\text{grad } \mathcal{L}(X) = P_{\mathcal{T}_X \mathcal{M}_r}(\nabla \mathcal{L}(X)), \quad X \in \mathcal{M}_r. \quad (7)$$

A simple and robust retraction that maps a tangent step  $\xi$  (for example,  $\xi$  is a negative Riemannian gradient) back to the manifold is the truncated SVD:

$$R_X(\xi) \equiv R(X + \xi) = \text{SVD}_r(X + \xi), \quad (8)$$

i.e., the best rank- $r$  approximation of  $X + \xi$  in Frobenius norm. Note that here we do not need to compute the full SVD and can utilize low-rank structure of  $X$  and  $\xi$ , leading to  $\mathcal{O}((m+n)r^2 + r^3)$  operations (Absil & Oseledets, 2015). Finally, because tangent spaces change from an optimization step to step, momentum (an accumulated tangent vector) must be moved between them via a *vector transport*. For embedded manifolds like  $\mathcal{M}_r$ , a standard choice is the *projection transport*

$$\mathcal{T}_{Y \rightarrow X}(\xi) = P_{\mathcal{T}_X \mathcal{M}_r}(\xi), \quad \xi \in \mathcal{T}_Y \mathcal{M}_r, \quad (9)$$

which simply reprojects the same ambient matrix  $\xi$  onto the new tangent space at  $X$ .

---

## 216 4 RIEMANNION

217  
218 In this section, we focus on the setting of parameter-efficient fine-tuning and, hence, to the fixed-rank  
219 manifold. For fine-tuning of one layer the optimization problem becomes:  
220

$$221 \mathcal{L}(W + \Delta W) \rightarrow \min_{\Delta W \in \mathcal{M}_r},$$

222 where  $\mathcal{L}$  is a differentiable loss function. Note that optimizing *on*  $\mathcal{M}_r$  removes the ambiguity  
223 of factorized parameterizations, because all computations are carried out in the intrinsic space of  
224 the product  $X$  rather than in any particular factorization. So the formulas we write below will  
225 naturally be reparameterization-invariant. Let  $G_t = P_{T_W \mathcal{M}_r}(\nabla \mathcal{L}(W_t))$  be the Riemannian gradient.  
226 A Riemannian *heavy-ball* (Polyak, 1964) momentum step reads  
227

$$228 M_t = \beta \widehat{M}_{t-1} + G_t, \quad \widehat{M}_{t-1} = \mathcal{T}_{W_{t-1} \rightarrow W_t}(M_{t-1}), \quad (10)$$

$$229 \Delta W_{t+1} = R(\Delta W_t - \eta M_t), \quad (11)$$

231 with  $M_0 = 0$ , momentum parameter  $\beta \in [0, 1)$ , and stepsize  $\eta > 0$ .

232 Let us now discuss how to introduce a Muon-like variant of this iteration, which we refer to as  
233 Riemannion. A direct projection onto the set of orthogonal matrices  $\text{Ortho}(M_t)$  presents two chal-  
234 lenges. First, such a step does not respect the underlying Riemannian geometry. To address this,  
235 we propose to find the best approximation of  $\text{Ortho}(M_t)$  on the tangent plane  $T_W \mathcal{M}_r$ . Such a so-  
236 lution is given via the projection onto the tangent plane  $P_{T_W \mathcal{M}_r}(\text{Ortho}(M_t))$ . However, a second  
237 issue arises: although this projection is low-rank, its computation remains inefficient because the  
238 input matrix is of full rank. Let us notice that  $M_t \in T_{\Delta W_t} \mathcal{M}_r$  and, hence, is of rank at most  $2r$   
239 (Section 3.2). At the same time, the LMO interpretation (Section 3.1) provides several admissible  
240 low-rank solutions, including one where  $\text{Ortho}_r(\cdot)$  replaces only the first  $2r$  singular values with 1,  
241 while all others are set to 0. Consequently, we obtain the following update rule:

$$242 \widetilde{M}_t = P_{T_{\Delta W_t} \mathcal{M}_r}(\text{Ortho}_r(M_t)). \quad (12)$$

244 Note that  $\text{Ortho}_r(\cdot)$  exactly preserves the column and row spaces of  $M_t \in T_{\Delta W_t} \mathcal{M}_r$ . Although  
245  $P_{T_{\Delta W_t} \mathcal{M}_r}(\text{Ortho}_r(M_t))$  does not yield singular values exactly equal to 1, in practice they remain  
246 in a close proximity. In particular, in our experiments they were always in the interval  $(0.9, 1.1)$ .  
247 This behavior is reminiscent of the Newton–Schulz iteration, which likewise produces approximate  
248 singular values. To eliminate this inexactness and obtain an accurate solution in the intersection  
249 of  $T_{\Delta W_t} \mathcal{M}_r$  with the set of matrices whose first  $2r$  singular values are exactly 1, one could apply  
250 various strategies. For example, one may formally apply the alternating projection method (Cheney  
251 & Goldstein, 1959; Boyd & Dattorro, 2003, Theorem 4, Section 2):

$$252 \widetilde{M}_t = P_{T_{\Delta W_t} \mathcal{M}_r}(\text{Ortho}_r(\dots P_{T_{\Delta W_t} \mathcal{M}_r}(\text{Ortho}_r(M_t)))). \quad (13)$$

254 As an alternative, we could solve the LMO problem on the tangent plane:

$$255 \widetilde{M}_t \in \underset{\substack{S \in T_{\Delta W_t} \mathcal{M}_r \\ \|S\|_2 \leq 1}}{\text{Arg max}} \langle M_t, S \rangle, \quad (14)$$

258 where we proposed to apply observations from (Cesista, 2025) to the fixed-rank manifold. In Ap-  
259 pendix I, we show that (12) is an approximate solution to the maximization problem (14). Apply-  
260 ing (13) or accurately solving (14) on the intersection of two convex sets using convex optimization  
261 methods is more computationally expensive, and our experiments indicate that it has little to no  
262 impact on the overall convergence of the optimizer.  
263

264 Let us finally show that  $\widetilde{M}_t$  from (12) can be computed efficiently using  $\mathcal{O}((m+n)r^2 + r^3)$  arith-  
265 metic operations. Indeed, first of all we need to apply  $\text{Ortho}_r(\cdot)$  to a tangent vector  $M_t$ . From (6),  
266 we know that a tangent vector can be represented in a form of a rank- $2r$  matrix. Such a representa-  
267 tion can always be transformed into the compact SVD form with 2 QR decompositions and a single  
268 full SVD of a  $2r \times 2r$  matrix, see Algorithm 1 called `OrthoLR`. As a next step, we need to project  
269 the obtained result (decomposed matrix of rank  $2r$ ) onto the tangent plane. The operation can also  
be done efficiently via (6) and is summarized in Algorithm 2 called `ProjectLR`.

---

270 **Algorithm 1** OrthoLR (efficient computation of  $\text{Ortho}_r(\xi)$  for  $\xi \in T_X \mathcal{M}_r$ ).

---

271 **Require:**  $\xi \in T_X \mathcal{M}_r$  given by  $(\dot{A}, \dot{B})$  from (5);  $A_L, B_R$  such that  $X = A_L T B_R^\top$  as in (4)

272 **Ensure:**  $A \in \mathbb{R}^{m \times 2r}, B \in \mathbb{R}^{m \times 2r}$ ;  $AB^\top = \text{Ortho}(\xi)$ .

- 273 1:  $Q_L, T_L = \text{qr}([A_L, \dot{A}]), \quad Q_R, T_r = \text{qr}([\dot{B}, B_R]^\top). \quad // \mathcal{O}((m+n)r^2)$   
274 2:  $U_L, V_r^\top = \text{SVD}(T_L T_r^\top). \quad // \mathcal{O}(r^3)$   
275 3:  $\dot{A} = Q_L U_L, \quad \dot{B} = Q_R V_r.$   $// \mathcal{O}((m+n)r^2)$
- 

276 **Algorithm 2** ProjectLR (efficient computation of  $P_{T_X \mathcal{M}_r}(Z)$  for a rank- $r'$  matrix  $Z$ ).

---

277 **Require:**  $A \in \mathbb{R}^{m \times r'}, B \in \mathbb{R}^{n \times r'}$  such that  $Z = AB^\top$ ;  $A_L, B_R$  such that  $X = A_L T B_R^\top$  as in (4).

278 **Ensure:**  $\xi = P_{T_X \mathcal{M}_r}(Z)$  given by  $(\dot{A}, \dot{B})$  from (5).

- 279 1:  $\dot{A} := (A - A_R(A_R^\top A))(B^\top B_R), \quad \dot{B} := B(A^\top A_L) \quad // \mathcal{O}((m+n)r'^2)$
- 

280 **Algorithm 3** RetractionLR (efficient computation of  $R(\xi)$  for a tangent vector  $\xi \in T_X \mathcal{M}_r$ ).

---

281 **Require:**  $A_L, \dot{A} \in \mathbb{R}^{m \times r}, B_R, \dot{B} \in \mathbb{R}^{n \times r}$  such that  $\xi = [\dot{A} \quad A_L] [B_R \quad \dot{B}]^\top$ .

282 **Ensure:**  $\tilde{X} = \tilde{A}_L \tilde{B} \in \mathcal{M}_r$  given by  $\tilde{X} = \text{truncSVD}(\xi, r)$ .

- 283 1:  $Q_L, T_L = \text{qr}([A_L, \dot{A}]), \quad Q_R, T_r = \text{qr}([\dot{B}, B_R]^\top). \quad // \mathcal{O}((m+n)r^2)$   
284 2:  $U, \Sigma, V^\top = \text{truncSVD}(T_L T_r^\top, r). \quad // \mathcal{O}(r^3)$   
285 3:  $\tilde{A}_L, \tilde{B} = (Q_L U), (\Sigma V^\top Q_R^\top). \quad // \mathcal{O}((m+n)r^2)$
- 

## 294 5 LOCALLY OPTIMAL INITIALIZATION (LOI)

295 Once the theoretical framework for the Riemannian optimizer is established, it is natural to consider  
296 an initialization scheme that accounts for the underlying Riemannian geometry. Given any  $\Delta W \in$   
297  $\mathcal{M}_r$ , we may write

$$298 \mathcal{L}(W) = \mathcal{L}(\underbrace{W - \Delta W}_{W'} + \Delta W) = \mathcal{L}(W' + \Delta W). \quad (15)$$

299 This raises the question: how should  $\Delta W$  be chosen to ensure the fastest loss decrease along the  
300 manifold? The solution is to consider the following optimization task:

$$301 \Delta W_*^{(0)} \in \text{Arg max}_{\Delta W \in \mathcal{M}_r} \|P_{T_{\Delta W} \mathcal{M}_r} \nabla_W \mathcal{L}(W)\|_F^2. \quad (16)$$

302 Since  $P_{T_{\Delta W} \mathcal{M}_r}$  is an orthogonal projection matrix to the tangent plane, the task (16) essentially  
303 seeks for the point on the fixed-rank manifold, whose tangent space has most alignment with the  
304 Euclidean gradient. In other words, this means that the direction of the steepest local function  
305 decrease alongside the manifold is aligned with the full model tuning direction. The solution to this  
306 task is presented in Theorem 5.1.

307 **Theorem 5.1.** *Let the SVD of  $\nabla_W \mathcal{L}(W)$  be:*

$$308 \nabla_W \mathcal{L}(W) = [U_{1,r} \quad U_{r,2r} \quad U_\perp] \begin{bmatrix} \Sigma_{1,r} & 0 & 0 \\ 0 & \Sigma_{r,2r} & 0 \\ 0 & 0 & \Sigma_\perp \end{bmatrix} [V_{1,r} \quad V_{r,2r} \quad V_\perp]^\top,$$

309 and let also  $\sigma_{2r} \neq \sigma_{2r+1}$ . Then any optimal solution  $\Delta W_*^{(0)}$  to the problem (16) has the form:

$$310 \Delta W_*^{(0)} \in \left\{ [U_{1,r}, U_{r,2r}, \Sigma_{r,2r}] \begin{bmatrix} S_{11} \\ S_{21} \end{bmatrix} [C_{21} \quad C_{22}] \begin{bmatrix} \Sigma_{1,r} V_{1,r}^\top \\ V_{r,2r} \end{bmatrix} \right\} \quad (17)$$

$$311 S = \begin{bmatrix} S_{11} & S_{12} \\ S_{21} & S_{22} \end{bmatrix} \in \text{GL}_{2r}(\mathbb{R}), \quad S^{-1} = \begin{bmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{bmatrix} \left. \right\}.$$

312 *Proof.* See Appendix A. □

In our experiments, we use  $S = \begin{bmatrix} \alpha I_r & 0 \\ 0 & I_r \end{bmatrix}$ , obtaining:  $\Delta W_*^{(0)} = \alpha U_{1,r} V_{r,2r}^\top \in \mathcal{M}_r, \alpha \in \mathbb{R} \setminus \{0\}$ .

Interestingly, Theorem 5.1 relates to the findings of Wang et al. (2024), although their analysis neither adopts a Riemannian framework nor addresses parametrization-free optimization. Our optimizer further differs from Zhang & Pilanci (2024) and Mo et al. (2025, LORO) in that it avoids inversion of the Gram matrix. As a result, the method remains stable as  $\|\Delta W_*^{(0)}\| \rightarrow 0$ . Empirical results indicate that initializing  $\Delta W_*^{(0)}$  with a small norm leads to improved performance. The procedure for selecting the scaling parameter  $\alpha$  is described in Appendix E.

## 6 SINGLE BACKWARD-PASS GRADIENT TRICK

This section introduces an efficient method for computing gradient-times-matrix in a matrix-free way, at a computational cost equivalent to a *single backward pass*. We then apply this technique in both the LOI initialization procedure in Algorithm 4 and the Riemannion optimizer in Algorithm 5.

The calculation of the full fine-tuning loss gradient  $\nabla_W \mathcal{L}(W)$  with pretrained parameters  $W \in \mathbb{R}^{m \times n}$  is computationally expensive. At the same time, for our framework we only need to compute the products  $(\nabla_W \mathcal{L}(W)^\top M)$  or  $(\nabla_W \mathcal{L}(W) N)$  for some  $M \in \mathbb{R}^{m \times r}, N \in \mathbb{R}^{n \times r}$ . Using the trick from (Novikov et al., 2022), we may calculate both quantities simultaneously using a *single forward-backward pass* with a doubled rank representation.

First, we initialize differentiable parameters  $Z_1 = 0 \in \mathbb{R}^{m \times r}$  and  $Z_2 = 0 \in \mathbb{R}^{n \times r}$  and perform a simple forward pass  $L := \mathcal{L}(W + Z_1 N^\top + M Z_2^\top)$ . This step does not violate the pipelines of LoRA framework since it is equivalent to a standard LoRA forward pass with special adapter:

$$Z_1 N^\top + M Z_2^\top = [Z_1 \quad M] [N \quad Z_2]^\top = \tilde{A} \tilde{B}^\top, \quad \tilde{A} \in \mathbb{R}^{m \times 2r}, \tilde{B} \in \mathbb{R}^{n \times 2r}.$$

Then we invoke an autodiff algorithm for value  $L$ , which ensures us both:

$$\begin{aligned} \nabla_{Z_1} L &= \nabla_{Z_1} \mathcal{L}(W + Z_1 N^\top + M Z_2^\top) |_{Z_1=0, Z_2=0} = \nabla_W \mathcal{L}(W) N, \\ \nabla_{Z_2} L &= \nabla_{Z_2} \mathcal{L}(W + Z_1 N^\top + M Z_2^\top) |_{Z_1=0, Z_2=0} = \nabla_W \mathcal{L}(W)^\top M. \end{aligned} \quad (18)$$

Note, that if  $\Delta W = A_L B^\top = A B_R^\top$  from  $W + \Delta W$ , then one may take  $N = B_R, M = A_L, Z_1 = 0, Z_2 = B, Y = W + \Delta W$  and the exact same operations work.

Notably, (18) is crucial for computation of the Riemannian gradient in (6). Therefore the proposed approach is a key building block that allows us to avoid forming the full fine-tune gradient loss and effectively compute matrix-matrix multiplications. This idea is the core for both: LOI initialization (Algorithm 4) and Riemannion optimizer (Algorithm 5).

**Randomized SVD for efficient initialization** The computation of the  $2r$ -truncated SVD of the full loss gradient, as required by Theorem 5.1, has asymptotic complexity  $\mathcal{O}(\min\{m, n\}mn)$ , which may be infeasible for large-scale models. In addition, explicitly forming this gradient is itself computationally expensive, and thus we need to avoid it in practice. To overcome this problem, we propose to use a randomized SVD with power iterations (see (Halko et al., 2011)), which we also enhance with our *one-step gradient trick* as is described in the Algorithm 4. In a nutshell, we need to compute  $(\nabla_W \mathcal{L}(W) \nabla_W \mathcal{L}(W)^\top)^q Y$ , where  $Y = \nabla_W \mathcal{L}(W) \Omega$  and  $\Omega$  is sampled from standard normal distribution. This iteration can be done in a robust manner using QR decompositions. The steps 2, 4, 5, 6 correspond to a trick from (18).

Overall, the LOI search procedure has asymptotic complexity  $\mathcal{O}((m+n)r^2)$ , plus  $2(q+1)$  additional backward passes. Moreover, since LOI search is executed only once before fine-tuning (which typically involves thousands of backward passes), its runtime overhead is negligible, taking merely 0.25% of the total fine-tuning wall-clock time in our experiments. Section 7.1, about 2 thousand optimization steps were carried out and Algorithm 4 accounted for merely 0.25% of the total fine-tuning wall-clock time.

**Efficient Riemannion implementation** The Riemannion optimizer is presented in Algorithm 5. Similar to vanilla Muon, it relies on a chosen `Ortho` procedure. Since the LoRA approach represents the fine-tuning shift  $\Delta W$  with low-rank matrices, we employ an SVD-based `Ortho` procedure.

378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431

---

**Algorithm 4** BackPropRSVD

---

**Require:** Weights  $W \in \mathbb{R}^{m \times n}$ , rank  $r \in \mathbb{N}$ , oversampling parameter  $p$ , power-step parameter  $q$ .  
**Ensure:** Randomized  $r$ -truncated SVD  $(U_r, \Sigma_r, V_r)$  of  $\nabla_W \mathcal{L}(W)$ .

- 1: Choose  $k = r + p$ , Sample  $\Omega \in \mathbb{R}^{n \times k} \sim \mathcal{N}(0, 1)$ . //  $\mathcal{O}(nr)$
- 2:  $Y := \text{qr}(\nabla_A \mathcal{L}(W + A\Omega^\top)|_{A=0}) \cdot \text{Q}$ . // 1 backward pass +  $\mathcal{O}(mr^2)$
- 3: **for**  $i := 1, \dots, q$  **do**
- 4:    $Y := \text{qr}([\nabla_B \mathcal{L}(W + YB^\top)|_{B=0}]^\top) \cdot \text{Q}$ . // 1 backward pass +  $\mathcal{O}(nr^2)$
- 5:    $Y := \text{qr}(\nabla_A \mathcal{L}(W + AY^\top)|_{A=0}) \cdot \text{Q}$ . // 1 backward pass +  $\mathcal{O}(mr^2)$
- 6:  $Y := [\nabla_B \mathcal{L}(W + YB^\top)|_{B=0}]^\top$ . // 1 backward pass
- 7:  $U, \Sigma, V^\top := \text{truncSVD}(Y, r)$ . //  $\mathcal{O}(nr^2)$
- 8:  $YU, \Sigma, V^\top$  //  $\mathcal{O}(mr^2)$

---

For computational efficiency, this procedure is adapted to the fixed-rank manifold, as described in Section 4 and implemented in the `OrthoLR` (Algorithm 1). In detail, in step 1 we calculate the Riemannian gradient components via a *single backward call* (18). Then, in step 2 the algorithm transports the Heavy-Ball tangent direction to the current point via Algorithm 2 that provides a simple but effective implementation of (6). In line 4, we compute the final optimization direction on the tangent space of the given point with a Heavy-Ball momentum coefficient  $\beta$ . In line 5 algorithm performs the retraction step.

The algorithm finalizes with saving the obtained minimization direction for momentum in line 6 and calculating of a new point representation in step 7.

---

**Algorithm 5** One step of Riemannion

---

**Require:** Weight matrix  $W' \in \mathbb{R}^{m \times n}$ , rank  $r \in \mathbb{N}$ , initial point  $A_L, B_R$ , Heavy-ball momentum  $A_{\text{HB}}, B_{\text{HB}}$ , step size  $\eta$ , momentum coefficient  $\beta$ , weight decay coefficient  $\gamma$ .  
**Ensure:** Tuning parameters  $\Delta W^* = A_L^* B_R^* \in \mathcal{M}_r$ .

- 1:  $\dot{A}, \dot{B} := \nabla_{Z_1} \mathcal{L}(W' + Z_1 B_R^\top + A_L Z_2^\top)|_{Z_1=0, Z_2=B}$ ,  
 $\nabla_{Z_2} \mathcal{L}(W' + Z_1 B_R^\top + A_L Z_2^\top)|_{Z_1=0, Z_2=B}$ . // 1 backward pass
- 2:  $\dot{A}_{\text{prev}}, \dot{B}_{\text{prev}} := \text{ProjectLR}((A_{\text{HB}}, B_{\text{HB}}), A_L, B_R)$ . //  $\mathcal{O}((m+n)r^2)$
- 3:  $\dot{A}, \dot{B} := \beta \dot{A}_{\text{prev}} + (I - A_L A_L^\top) \dot{A}$ ,  $\beta \dot{B}_{\text{prev}} + \dot{B}$  //  $\mathcal{O}((n+m)r^2)$
- 4:  $\dot{A}, \dot{B} := \text{ProjectLR}(\text{OrthoLR}(A_L, B_R, \dot{A}, \dot{B}, r), A_L, B_R)$  //  $\mathcal{O}((m+n)r^2 + r^3)$
- 5:  $A_{\text{HB}}, B_{\text{HB}} := [\dot{A}, A_L], [B_R, \dot{B}]$
- 6:  $A_L^*, B_R^* := \text{RetractionLR}([-\eta \dot{A}, A_L], [B_R, -\eta(\dot{B} + \gamma B_R)])$  //  $\mathcal{O}((m+n)r^2 + r^3)$

---

Overall, one iteration of the Riemannion loop has asymptotic complexity  $\mathcal{O}((m+n)r^2 + r^3)$  and additionally the same number of backward passes as vanilla LoRA. For comparison, the Euclidean Muon optimizer for LoRA (Section 3.1) exhibits the same asymptotic complexity. The complete framework and its time performance are summarized in Algorithm 6 (see Appendix B) and Appendix H.

## 7 EXPERIMENTS

In this section, we also employ the term LOI (Locally Optimal Initialization), referring to the proposed initialization scheme described in Section 5. This is done to distinguish between using RiemannLoRA as an optimizer with zero initialization (as in basic LoRA) and RiemannLoRA-LOI, which is the proposed combination of initialization and optimization. We conduct a series of LLM fine-tuning experiments for different tasks. All of the experiments were computed on NVIDIA V100-32Gb GPU and A100-80Gb GPU. We ran all the experiments within  $\sim 2000$  GPU hours.

## 7.1 COMMONSENSE REASONING FINE-TUNING

The results were obtained for the benchmark (Clark et al. (2019, BoolQ), Bisk et al. (2020, PIQA), Sap et al. (2019, SIQA), Zellers et al. (2019, hellaswag), Sakaguchi et al. (2021, winogrande), Clark et al. (2018, ARC), Mihaylov et al. (2018, OBQA)) common reasoning. The structure of the dataset is described in Appendix D. In the following experiments we conduct a fine-tuning procedures for multilayer perceptron (MLP) and attention layers of Llama 3 8b model (Dubey et al. (2024)).

The commonsense reasoning tasks comprise of 8 sub-tasks, each of them contains a predefined training and a testing set. We follow the setting of Hu et al. (2023) and amalgamate the training datasets from all 8 tasks to create the final training dataset and conduct evaluations on the individual testing dataset for each task.

The hyperparameter tuning protocol and the selected hyperparameters are provided in E. Table 1 contains the accuracy of the trained model’s responses on the test dataset. Within the LoRA framework, the proposed Riemannian fine-tuning method delivers clear performance gains. The Riemannian optimizer consistently outperforms LoRA, DoRA, and achieves superior metric results compared to the standard Muon optimizer applied to LoRA factors (see Section 3.1). Furthermore, relative to other Riemannian-geometry-aware approaches such as RPrecAdamW (Zhang & Pilanci, 2024), our method also demonstrates better results. Finally, the variance of outcomes for the proposed method is the smallest among all compared approaches.

Table 1: The average accuracy (in %) among 8 tasks of fine-tuned Llama 3-8b using different approaches, tested on Commonsense Reasoning benchmark. LoRA rank is set to 16.

Task	BoolQ	PIQA	SIQA	hella- swag	wino- grande	ARC-E	ARC-C	OBQA	All
Raw Initialization									
Raw	65.0	76.6	73.0	66.1	61.3	92.5	82.3	79.6	74.5
Adam	74.8 $\pm$ 1.9	89.8 $\pm$ 0.9	82.6 $\pm$ 0.6	96.2 $\pm$ 0.3	87.9 $\pm$ 1.2	92.4 $\pm$ 0.7	84.9 $\pm$ 0.7	88.5 $\pm$ 0.4	87.1 $\pm$ 0.6
DoRA	74.8 $\pm$ 0.8	89.4 $\pm$ 0.5	82.4 $\pm$ 0.7	95.9 $\pm$ 0.1	87.8 $\pm$ 0.4	90.7 $\pm$ 1.2	83.8 $\pm$ 0.7	87.8 $\pm$ 0.6	86.6 $\pm$ 0.3
Muon	72.9 $\pm$ 0.0	86.4 $\pm$ 0.5	80.8 $\pm$ 0.2	94.1 $\pm$ 0.2	84.4 $\pm$ 0.0	84.2 $\pm$ 0.9	77.3 $\pm$ 2.5	83.9 $\pm$ 1.1	83.0 $\pm$ 0.6
LoRA-RITE	72.2 $\pm$ 0.3	88.6 $\pm$ 0.1	82.0 $\pm$ 0.6	95.1 $\pm$ 0.1	85.6 $\pm$ 0.2	87.7 $\pm$ 1.5	79.3 $\pm$ 2.6	85.7 $\pm$ 0.5	84.5 $\pm$ 0.5
RPrecAdamW	75.8 $\pm$ 0.4	89.5 $\pm$ 0.4	82.4 $\pm$ 0.2	96.1 $\pm$ 0.2	87.7 $\pm$ 0.9	90.6 $\pm$ 1.6	84.1 $\pm$ 1.1	87.7 $\pm$ 0.5	86.8 $\pm$ 0.4
Riemannion	75.7 $\pm$ 0.7	<b>91.2<math>\pm</math>0.2</b>	<b>83.5<math>\pm</math>0.6</b>	<b>96.7<math>\pm</math>0.0</b>	<b>88.6<math>\pm</math>0.4</b>	<b>93.6<math>\pm</math>0.3</b>	<b>86.4<math>\pm</math>0.4</b>	<b>89.3<math>\pm</math>0.8</b>	<b>88.1<math>\pm</math>0.2</b>

## 7.2 SUBJECT-DRIVEN GENERATION

Subject-driven generation (Ruiz et al., 2023; Gal et al., 2023) is a task in which the user provides several reference photos of an object, called a concept, and uses a diffusion model to generate this concept with certain conditions (e.g., a textual prompt). One way to solve this task is to fine-tune a pre-trained diffusion model using this small set of reference images. However, this technique leads to a degradation in understanding of the conditions and to a fast overfitting of the concept. Furthermore, it requires a high computational cost due to the large number of trainable parameters. This is why previous works, such as (Qiu et al., 2023; Liu et al., 2024; Hu et al., 2022; Tewel et al., 2023; Han et al., 2023; Gorbunov et al., 2024), propose training only a lightweight parameterization for the base model. In this section, we demonstrate the performance of our parameterization in this task.

In our experiments, we used Stable Diffusion 2 (Rombach et al., 2022) as the base model. We choose LoRA Hu et al. (2022) as a baseline and train both models with ranks of 4, 8 and 16. We predict the parameterization of the  $\mathbf{c}_l$ ,  $\mathbf{k}$ ,  $\mathbf{v}$ , and  $\text{out} \cdot 0$  matrices in all attention layers. The Dreambooth dataset (Ruiz et al., 2023) was used in all our experiments. LoRA was trained using the Adam optimizer. In Figure 2, we present a visual comparison of LoRA with different ranks. As can be seen, even for complex concepts such as ‘robot toy’, our method requires only 600 steps to learn the concept while preserving appropriate text similarity. We found that for subject-driven generation tasks, the lower the rank, the faster our method converges. To evaluate this, we also calculated metrics for different learning

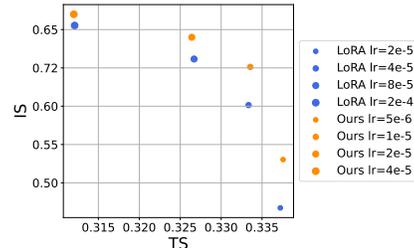


Figure 1: Comparison of text and image similarities for LoRA and our method with rank 4 at different learning rates at 400 step.

486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539

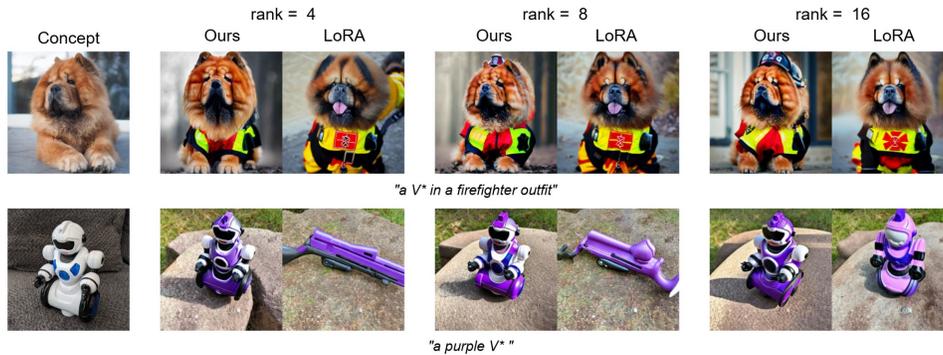


Figure 2: Visual results for Subject-driven generation on 600 training step.

rates on a subset of the Dreambooth dataset. We use CLIP to measure text similarity and DINO to measure image similarity. Figure 1 shows that, even with different learning rates, our method achieves more accurate results in concept preservation. Further details and a visual comparison can be found in the Appendix F.

## 8 CONCLUSION

In this work, we propose a novel fully Riemannian framework that integrates a new muon-based optimization method, locally-optimal initialization, and an efficient implementation. This integrated approach yields a reliable reparametrization-invariant method that outperforms competing approaches on fine-tuning large language models (LLMs) and exhibits additional favorable properties for low-rank approximations in diffusion models. Given the promising empirical results, a natural direction for future research is to investigate the theoretical properties of the proposed method.

## 9 REPRODUCIBILITY STATEMENT

The hyperparameter selection procedure is described in Appendix E. The datasets used in the experiments and the corresponding preprocessing steps are detailed in Appendix D. The proof of the stated assumptions is provided in Appendix A. The hyperparameters for Subject-driven generation task are provided in the Appendix F.

## REFERENCES

P-A Absil and Ivan V Oseledets. Low-rank retractions: a survey and new results. *Computational Optimization and Applications*, 62(1):5–29, 2015.

P-A Absil, Robert Mahony, and Rodolphe Sepulchre. Optimization algorithms on matrix manifolds. In *Optimization Algorithms on Matrix Manifolds*. Princeton University Press, 2009.

Christopher G Baker. *Riemannian manifold trust-region methods with applications to eigenproblems*. The Florida State University, 2008.

Jeremy Bernstein. Deriving muon, 2025. URL <https://jeremybernste.in/writing/deriving-muon>.

Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. PIQA: Reasoning about Physical Commonsense in Natural Language. In *Proceedings of the AAAI conference on artificial intelligence*, number 05 in 34, pp. 7432–7439, 2020.

Nicolas Boumal. *An introduction to optimization on smooth manifolds*. Cambridge University Press, 2023.

- 
- 540 Stephen Boyd and Jon Dattorro. Alternating projections. 2003. URL [https://web.stanford.edu/class/ee392o/alt\\_proj.pdf](https://web.stanford.edu/class/ee392o/alt_proj.pdf).
- 541
- 542
- 543 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
- 544 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
- 545 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 546
- 547 Franz Louis Cesista. Muon and a selective survey on Steepest Descent in Riemannian and
- 548 non-Riemannian Manifolds, April 2025. URL [http://leloykun.github.io/ponder/](http://leloykun.github.io/ponder/steepest-descent-non-riemannian/)
- 549 [steepest-descent-non-riemannian/](http://leloykun.github.io/ponder/steepest-descent-non-riemannian/).
- 550 Ward Cheney and Allen A Goldstein. Proximity maps for convex sets. *Proceedings of the American*
- 551 *Mathematical Society*, 10(3):448–450, 1959.
- 552
- 553 Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina
- 554 Toutanova. BoolQ: Exploring the surprising difficulty of natural yes/no questions. In Jill Burstein,
- 555 Christy Doran, and Thamar Solorio (eds.), *Proceedings of the 2019 Conference of the North Amer-*
- 556 *ican Chapter of the Association for Computational Linguistics: Human Language Technologies,*
- 557 *NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*,
- 558 pp. 2924–2936. Association for Computational Linguistics, 2019. doi: 10.18653/V1/N19-1300.
- 559 URL <https://doi.org/10.18653/v1/n19-1300>.
- 560 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and
- 561 Oyvind Tafjord. Think you have solved question answering? try ARC, the AI2 reasoning chal-
- 562 lenge. *CoRR*, abs/1803.05457, 2018. URL <http://arxiv.org/abs/1803.05457>.
- 563 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
- 564 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony
- 565 Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark,
- 566 Arun Rao, Aston Zhang, Aurélien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Rozière,
- 567 Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris
- 568 Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong,
- 569 Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny
- 570 Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino,
- 571 Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael
- 572 Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Ander-
- 573 son, Graeme Nail, Grégoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Ko-
- 574 revaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel M. Kloumann, Ishan
- 575 Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Ma-
- 576 hadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy
- 577 Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak,
- 578 Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Al-
- 579 wala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, and et al. The
- 580 Llama 3 herd of models. *CoRR*, abs/2407.21783, 2024. doi: 10.48550/ARXIV.2407.21783. URL
- <https://doi.org/10.48550/arXiv.2407.21783>.
- 581 John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and
- 582 stochastic optimization. *Journal of machine learning research*, 12(7), 2011.
- 583
- 584 Carl Eckart and Gale Young. The approximation of one matrix by another of lower rank. *Psychome-*
- 585 *trika*, 1(3):211–218, 1936.
- 586
- 587 Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit Haim Bermano, Gal Chechik, and
- 588 Daniel Cohen-Or. An image is worth one word: Personalizing text-to-image generation using
- 589 textual inversion. In *The Eleventh International Conference on Learning Representations, ICLR*
- 590 *2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL <https://openreview.net/forum?id=NAQvF08TcyG>.
- 591
- 592 Mikhail Gorbunov, Nikolay Yudin, Vera Soboleva, Aibek Alanov, Alexey Naumov, and Maxim
- 593 Rakhuba. Group and shuffle: Efficient structured orthogonal parametrization. In Amir Globersons, Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and

- 
- 594 Cheng Zhang (eds.), *Advances in Neural Information Processing Systems 38: Annual Confer-*  
595 *ence on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada,*  
596 *December 10 - 15, 2024, 2024.*
- 597 N. Halko, P. G. Martinsson, and J. A. Tropp. Finding structure with randomness: Probabilistic  
598 algorithms for constructing approximate matrix decompositions. *SIAM Review*, 53(2):217–288,  
599 2011. doi: 10.1137/090771806. URL <https://doi.org/10.1137/090771806>.
- 600 Ligong Han, Yinxiao Li, Han Zhang, Peyman Milanfar, Dimitris Metaxas, and Feng Yang. Svdiff:  
601 Compact parameter space for diffusion fine-tuning. In *Proceedings of the IEEE/CVF Interna-*  
602 *tional Conference on Computer Vision*, pp. 7323–7334, 2023.
- 603 Soufiane Hayou, Nikhil Ghosh, and Bin Yu. LoRA+: Efficient Low Rank Adaptation of Large Mod-  
604 els. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria,*  
605 *July 21-27, 2024*. OpenReview.net, 2024. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=NEv8YqBROO)  
606 [NEv8YqBROO](https://openreview.net/forum?id=NEv8YqBROO).
- 607 Sebastian Holtz, Thorsten Rohwedder, and Reinhold Schneider. On manifolds of tensors of fixed  
608 tt-rank. *Numerische Mathematik*, 120(4):701–731, 2012.
- 609 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,  
610 Weizhu Chen, et al. LoRA: Low-Rank Adaptation of Large Language Models. *ICLR*, 1(2):3,  
611 2022.
- 612 Jiang Hu, Jiayi Cui, Lin Lin, Zaiwen Wen, Quanzheng Li, et al. Retraction-free optimization over  
613 the Stiefel manifold with application to the loRA fine-tuning, 2024.
- 614 Zhiqiang Hu, Lei Wang, Yihuai Lan, Wanyu Xu, Ee-Peng Lim, Lidong Bing, Xing Xu, Soujanya  
615 Poria, and Roy Ka-Wei Lee. LLM-adapters: An adapter family for parameter-efficient fine-tuning  
616 of large language models. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings*  
617 *of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023,*  
618 *Singapore, December 6-10, 2023*, pp. 5254–5276. Association for Computational Linguistics,  
619 2023. doi: 10.18653/v1/2023.EMNLP-MAIN.319. URL [https://doi.org/10.18653/](https://doi.org/10.18653/v1/2023.emnlp-main.319)  
620 [v1/2023.emnlp-main.319](https://doi.org/10.18653/v1/2023.emnlp-main.319).
- 621 Keller Jordan, Yuchen Jin, Vlado Boza, Jiacheng You, Franz Cesista, Laker Newhouse, and Jeremy  
622 Bernstein. Muon: An optimizer for hidden layers in neural networks, 2024. URL [https://](https://kellerjordan.github.io/posts/muon/)  
623 [kellerjordan.github.io/posts/muon/](https://kellerjordan.github.io/posts/muon/).
- 624 Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *International*  
625 *Conference on Learning Representations*, 12 2014.
- 626 John M Lee. *Smooth manifolds*. Springer, 2003.
- 627 Weiyang Liu, Zeju Qiu, Yao Feng, Yuliang Xiu, Yuxuan Xue, Longhui Yu, Haiwen Feng, Zhen  
628 Liu, Juyeon Heo, Songyou Peng, Yandong Wen, Michael J. Black, Adrian Weller, and Bernhard  
629 Schölkopf. Parameter-efficient orthogonal finetuning via butterfly factorization. In *The Twelfth*  
630 *International Conference on Learning Representations*, 2024. URL [https://openreview.](https://openreview.net/forum?id=7NzgkEdGyr)  
631 [net/forum?id=7NzgkEdGyr](https://openreview.net/forum?id=7NzgkEdGyr).
- 632 Ilya Loshchilov and Frank Hutter. Decoupled Weight Decay Regularization. In *7th International*  
633 *Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*.  
634 OpenReview.net, 2019. URL <https://openreview.net/forum?id=Bkg6RiCqY7>.
- 635 Fanxu Meng, Zhaohui Wang, and Muhan Zhang. PiSSA: Principal Singular Values and Singular  
636 Vectors Adaptation of Large Language Models. *Advances in Neural Information Processing*  
637 *Systems*, 37:121038–121072, 2024.
- 638 Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct  
639 electricity? A new dataset for open book question answering. In Ellen Riloff, David Chiang,  
640 Julia Hockenmaier, and Jun’ichi Tsujii (eds.), *Proceedings of the 2018 Conference on Empirical*  
641 *Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018*,  
642 pp. 2381–2391. Association for Computational Linguistics, 2018. doi: 10.18653/v1/D18-1260.  
643 URL <https://doi.org/10.18653/v1/d18-1260>.

- 
- 648 Zhanfeng Mo, Long-Kai Huang, and Sinno Jialin Pan. Parameter and Memory Efficient Pretraining  
649 via Low-rank Riemannian Optimization. In *The Thirteenth International Conference on Learning*  
650 *Representations*, 2025.
- 651 Alexander Novikov, Maxim Rakhuba, and Ivan Oseledets. Automatic differentiation for Riemannian  
652 optimization on low-rank matrix and tensor-train manifolds. *SIAM Journal on Scientific*  
653 *Computing*, 44(2):A843–A869, 2022.
- 654 Ivan V Oseledets. Tensor-train decomposition. *SIAM Journal on Scientific Computing*, 33(5):2295–  
655 2317, 2011.
- 656 Uliana Parkina and Maxim Rakhuba. Coala: Numerically stable and efficient framework for context-  
657 aware low-rank approximation, 2025. URL <https://arxiv.org/abs/2507.07580>.
- 658 Boris T Polyak. Some methods of speeding up the convergence of iteration methods. *Ussr compu-*  
659 *tational mathematics and mathematical physics*, 4(5):1–17, 1964.
- 660 Zeju Qiu, Weiyang Liu, Haiwen Feng, Yuxuan Xue, Yao Feng, Zhen Liu, Dan Zhang, Adrian  
661 Weller, and Bernhard Schölkopf. Controlling text-to-image diffusion by orthogonal finetun-  
662 ing. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL  
663 <https://openreview.net/forum?id=K30wTdIIYc>.
- 664 Herbert Robbins and Sutton Monro. A stochastic approximation method. *The annals of mathemati-*  
665 *cal statistics*, pp. 400–407, 1951.
- 666 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
667 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Con-*  
668 *ference on Computer Vision and Pattern Recognition (CVPR)*, pp. 10684–10695, June 2022.
- 669 Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman.  
670 Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Pro-*  
671 *ceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 22500–  
672 22510, 2023.
- 673 Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. WinoGrande: An Ad-  
674 versarial Winograd Schema Challenge at Scale. *Communications of the ACM*, 64(9):99–106,  
675 2021.
- 676 Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. SocialIQa: Com-  
677 mon-sense Reasoning about Social Interactions. *CoRR*, abs/1904.09728, 2019. URL <http://arxiv.org/abs/1904.09728>.
- 678 Yoad Tewel, Rinon Gal, Gal Chechik, and Yuval Atzmon. Key-locked rank one editing for text-to-  
679 image personalization. In *ACM SIGGRAPH 2023 Conference Proceedings*, pp. 1–11, 2023.
- 680 Tijmen Tieleman. Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent  
681 magnitude. *COURSERA: Neural networks for machine learning*, 4(2):26, 2012.
- 682 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée  
683 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Ar-  
684 mand Joulin, Edouard Grave, and Guillaume Lample. LLaMA: Open and efficient foundation  
685 language models. *CoRR*, abs/2302.13971, 2023a. doi: 10.48550/ARXIV.2302.13971. URL  
686 <https://doi.org/10.48550/arXiv.2302.13971>.
- 687 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-  
688 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrutu Bhosale, Dan Bikel, Lukas Blecher,  
689 Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy  
690 Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn,  
691 Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel  
692 Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya  
693 Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar  
694 Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan  
695 Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen

702 Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan  
703 Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurélien Rodriguez,  
704 Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-  
705 tuned chat models. *CoRR*, abs/2307.09288, 2023b. doi: 10.48550/ARXIV.2307.09288. URL  
706 <https://doi.org/10.48550/arXiv.2307.09288>.

707 Nickolay Trendafilov and Michele Gallo. *Multivariate data analysis on matrix manifolds*. Springer,  
708 2021.

709 Bart Vandereycken. Low-rank matrix completion by Riemannian optimization. *SIAM Journal on*  
710 *Optimization*, 23(2):1214–1236, 2013.

711 Hanqing Wang, Yixia Li, Shuo Wang, Guanhua Chen, and Yun Chen. MiLoRA: Harnessing minor  
712 singular components for parameter-efficient LLM finetuning. In Luis Chiruzzo, Alan Ritter, and  
713 Lu Wang (eds.), *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of*  
714 *the Association for Computational Linguistics: Human Language Technologies, NAACL 2025 -*  
715 *Volume 1: Long Papers, Albuquerque, New Mexico, USA, April 29 - May 4, 2025*, pp. 4823–4836.  
716 Association for Computational Linguistics, 2025. doi: 10.18653/V1/2025.NAACL-LONG.248.  
717 URL <https://doi.org/10.18653/v1/2025.naacl-long.248>.

718 Shaowen Wang, Linxi Yu, and Jian Li. LoRA-GA: Low-rank adaptation with gradient approxima-  
719 tion. *Advances in Neural Information Processing Systems*, 37:54905–54931, 2024.

720 Scott Wisdom, Thomas Powers, John Hershey, Jonathan Le Roux, and Les Atlas. Full-Capacity  
721 Unitary Recurrent Neural Networks. *Advances in neural information processing systems*, 29,  
722 2016.

723 Yibo Yang, Xiaojie Li, Zhongzhu Zhou, Shuaiwen Song, Jianlong Wu, Liqiang Nie, and Bernard  
724 Ghanem. CorDA: Context-oriented decomposition adaptation of large language models for task-  
725 aware parameter-efficient fine-tuning. *Advances in Neural Information Processing Systems*, 37:  
726 71768–71791, 2024.

727 Jui-Nan Yen, Si Si, Zhao Meng, Felix Yu, Sai Surya Duvvuri, Inderjit S Dhillon, Cho-Jui Hsieh, and  
728 Sanjiv Kumar. Lora done rite: Robust invariant transformation equilibration for lora optimization.  
729 *arXiv preprint arXiv:2410.20625*, 2024.

730 Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. HellaSwag: Can a ma-  
731 chine really finish your sentence? In Anna Korhonen, David R. Traum, and Lluís Màrquez  
732 (eds.), *Proceedings of the 57th Conference of the Association for Computational Linguistics,*  
733 *ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, pp. 4791–  
734 4800. Association for Computational Linguistics, 2019. doi: 10.18653/V1/P19-1472. URL  
735 <https://doi.org/10.18653/v1/p19-1472>.

736 Fangzhao Zhang and Mert Pilanci. Riemannian preconditioned loRA for fine-tuning foundation  
737 models. In *Proceedings of the 41st International Conference on Machine Learning, ICML’24*.  
738 JMLR.org, 2024.

739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

## A INITIAL POINT SEARCH

In this section we introduce the proof of Theorem 5.1

*Proof.* In order to represent the optimization task (16) in a more simple way, we will use a slightly different formula for the projection of the full loss gradient onto the tangent space of the fixed-rank manifold. Let

$$\Delta W = A_L B^\top = A B_R^\top \in \mathcal{M}_r,$$

and the tangent space is parametrized as follows

$$\mathcal{T}_{\Delta W} \mathcal{M}_r = \{\dot{A} B_R^\top + A_L \dot{B}^\top \mid \dot{B} \in \mathbb{R}^{n \times r}, \dot{A} \in \mathbb{R}^{m \times r}, A_L^\top \dot{A} = 0\},$$

then one may derive an orthogonal projection formula for any matrix  $Z$  (see, e.g. (Boumal, 2023, eq. 7.53))

$$P_{\mathcal{T}_{\Delta W} \mathcal{M}_r} Z = Z - (I - A_L A_L^\top) Z (I - B_R B_R^\top) \in \mathcal{T}_{\Delta W} \mathcal{M}_r.$$

As the tangent space is a linear space, the operation of orthogonal projection can be written as an optimization task

$$P_{\mathcal{T}_{\Delta W} \mathcal{M}_r} Z = \arg \min_{\xi \in \mathcal{T}_{\Delta W} \mathcal{M}_r} \|Z - \xi\|_F^2. \quad (19)$$

Since

$$\begin{aligned} \|\nabla \mathcal{L}\|_F^2 &= \|(\nabla \mathcal{L} - P_{\mathcal{T}_{\Delta W} \mathcal{M}_r} \nabla \mathcal{L}) + P_{\mathcal{T}_{\Delta W} \mathcal{M}_r} \nabla \mathcal{L}\|_F^2 = \\ &= \|\nabla \mathcal{L} - P_{\mathcal{T}_{\Delta W} \mathcal{M}_r} \nabla \mathcal{L}\|_F^2 + \|P_{\mathcal{T}_{\Delta W} \mathcal{M}_r} \nabla \mathcal{L}\|_F^2 + 2\langle \nabla \mathcal{L} - P_{\mathcal{T}_{\Delta W} \mathcal{M}_r} \nabla \mathcal{L}, P_{\mathcal{T}_{\Delta W} \mathcal{M}_r} \nabla \mathcal{L} \rangle \\ &= \|\nabla \mathcal{L} - P_{\mathcal{T}_{\Delta W} \mathcal{M}_r} \nabla \mathcal{L}\|_F^2 + \|P_{\mathcal{T}_{\Delta W} \mathcal{M}_r} \nabla \mathcal{L}\|_F^2 + 2\langle P_{\mathcal{T}_{\Delta W} \mathcal{M}_r} (\nabla \mathcal{L} - P_{\mathcal{T}_{\Delta W} \mathcal{M}_r} \nabla \mathcal{L}), \nabla \mathcal{L} \rangle \\ &= \|\nabla \mathcal{L} - P_{\mathcal{T}_{\Delta W} \mathcal{M}_r} \nabla \mathcal{L}\|_F^2 + \|P_{\mathcal{T}_{\Delta W} \mathcal{M}_r} \nabla \mathcal{L}\|_F^2 + 2\underbrace{\langle P_{\mathcal{T}_{\Delta W} \mathcal{M}_r} \nabla \mathcal{L} - P_{\mathcal{T}_{\Delta W} \mathcal{M}_r} \nabla \mathcal{L}, \nabla \mathcal{L} \rangle}_{=0} \\ &= \|\nabla \mathcal{L} - P_{\mathcal{T}_{\Delta W} \mathcal{M}_r} \nabla \mathcal{L}\|_F^2 + \|P_{\mathcal{T}_{\Delta W} \mathcal{M}_r} \nabla \mathcal{L}\|_F^2 \end{aligned}$$

so the optimization task (16) is equivalent to

$$\|\nabla \mathcal{L} - P_{\mathcal{T}_{\Delta W} \mathcal{M}_r} \nabla \mathcal{L}\|_F^2 \rightarrow \min_{\Delta W},$$

which, in turn, using (19) is equivalent to the task

$$\min_{\Delta W} \min_{\xi \in \mathcal{T}_{\Delta W} \mathcal{M}_r} \|\nabla \mathcal{L} - \xi\|_F^2.$$

Due to the fact that every vector of the tangent space is an element of a set of all matrices with  $2r$ -bounded rank  $\mathcal{M}_{\leq 2r}$ , one may use the Eckart-Young-Mirsky theorem (see (Eckart & Young, 1936)) to obtain the following lower bound:

$$\|\nabla \mathcal{L} - P_{\mathcal{T}_{\Delta W} \mathcal{M}_r} \nabla \mathcal{L}\|_F^2 \geq \|\nabla \mathcal{L} - \text{truncSVD}(\nabla \mathcal{L}, 2r)\|_F^2, \forall \Delta W \in \mathcal{M}_r, \xi \in \mathcal{T}_{\Delta W} \mathcal{M}_r. \quad (20)$$

But it is possible to ensure  $\Delta W_*^{(0)} \in \mathcal{M}_r$  and  $\xi_* \in \mathcal{T}_{\Delta W_*^{(0)}} \mathcal{M}_r$  which turn the inequality (20) into the equality. One may take

$$\begin{aligned} \Delta W_*^{(0)} &= A_L (\alpha B)^\top = \alpha U_{1,r} V_{r,2r}^\top = (\alpha A) B_R^\top, \quad \alpha \in \mathbb{R}, \\ \xi_* &= \dot{A} B_R^\top + A_L \dot{B}^\top = (U_{r,2r} \Sigma_{r,2r}) B_R^\top + A_L (\Sigma_{1,r} V_{1,r})^\top = \text{truncSVD}(\nabla \mathcal{L}, 2r). \end{aligned}$$

Note, that  $\dot{A}^\top A_L = 0$ . So the initialization  $\Delta W_*^{(0)}$  ensures that  $\text{truncSVD}(\nabla \mathcal{L}, 2r)$  lies in the tangent space  $\mathcal{T}_{\Delta W_*^{(0)}} \mathcal{M}_r$ , which means that the first step of the RiemannLoRA Algorithm 7 will get the Riemannian gradient equal to  $\text{truncSVD}(\nabla \mathcal{L}, 2r)$ .

To generalize the result, we should change the parametrization of the tangent space. Because of the fact that each of the tangent vectors  $\xi$  lies in  $\mathcal{M}_{\leq 2r}$ , we may use an unconstrained skeleton decomposition (without constraints for  $\dot{A}$  and  $A_L$ ):

$$\begin{aligned}\xi_* &= [A_*, \dot{A}_*] \begin{bmatrix} \dot{B}_*^\top \\ B_*^\top \end{bmatrix} = [A_*, \dot{A}_*] I_{2r} \begin{bmatrix} \dot{B}_*^\top \\ B_*^\top \end{bmatrix} = \\ &= [A_*, \dot{A}_*] S S^{-1} \begin{bmatrix} \dot{B}_*^\top \\ B_*^\top \end{bmatrix}, \quad S \in \mathbb{R}^{2r \times 2r}, \det S \neq 0.\end{aligned}$$

Representing  $S$  and its inverse as block matrices:

$$S = \begin{bmatrix} S_{11} & S_{12} \\ S_{21} & S_{22} \end{bmatrix}, \quad S^{-1} = \begin{bmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{bmatrix},$$

one arrives to (17):

$$\begin{aligned}\xi_* &= [A_*, \dot{A}_*] \begin{bmatrix} S_{11} & S_{12} \\ S_{21} & S_{22} \end{bmatrix} \begin{bmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{bmatrix} [\dot{B}_* \quad B_*]^\top, \\ \Delta W_*^{(0)} &= [A'_*, \dot{A}'_*] \begin{bmatrix} S_{11} \\ S_{21} \end{bmatrix} [C_{21} \quad C_{22}] [\dot{B}'_* \quad B'_*]^\top.\end{aligned}$$

To derive the version that is utilized in practice, one may take  $S$  to be block-diagonal with  $S_{11} = \alpha I_r, S_{22} = I_r, \alpha \in \mathbb{R}$ .  $\square$

## B RIEMANNION WITH LOI

The complete framework is summarized in Algorithm 6, which consists of an LOI initialization step followed by `max_iters` iterations of the Riemannion optimizer.

Specifically, the first 3 lines are dedicated to the computation of LOI via `BackPropRSVD` for given layer weights and oversampling parameter  $p$ , that increases the accuracy of the singular approximation. The 4-th step initializes Heavy-Ball momentum matrices. Therefore, the asymptotical complexity of this stage is determined by the complexity of Algorithm 4.

---

### Algorithm 6 Riemannion with LOI

---

**Require:** Weights  $W \in \mathbb{R}^{m \times n}$ , rank  $r \in \mathbb{N}$ , step size  $\eta$ , momentum coefficient  $\beta$ , weight decay coefficient  $\gamma$ , oversampling parameter  $p$ , power-step parameter  $q$ .

**Ensure:** Tuning parameters  $\Delta W_* \in \mathcal{M}_r$ .

**Function:**

- 1:  $A_L, \_, B^\top := \text{BackPropRSVD}(2r, p, q, \mathcal{L}, W)$   $\ll \mathcal{O}((m+n)r^2)$
  - 2:  $A_L, B^\top := A_L[:, :r], B[:, r:]^\top$ .
  - 3:  $W' := W - A_L B^\top$ .  $\ll \mathcal{O}(mn)$
  - 4:  $A_{\text{HB}}, B_{\text{HB}} := 0, 0$ .
  - 5: **for**  $i := 0, \dots, \text{max\_iters}$  **do**
  - 6:    $B_R := \text{qr}(B)$ .  $\ll \mathcal{O}(nr^2)$
  - 7:    $\dot{A} := \nabla_{Z_1} \mathcal{L}(W' + Z_1 B_R^\top + A_L Z_2^\top) |_{Z_1=0, Z_2=B}$   $\ll \nabla_W \mathcal{L}(W) B_R$
  - 8:    $\dot{B} := \nabla_{Z_2} \mathcal{L}(W' + Z_1 B_R^\top + A_L Z_2^\top) |_{Z_1=0, Z_2=B}$   $\ll \nabla_W \mathcal{L}(W)^\top A_L$
  - 9:    $\dot{A}_{\text{prev}}, \dot{B}_{\text{prev}} := \text{ProjectLR}((A_{\text{HB}}, B_{\text{HB}}), A_L, B_R)$   $\ll \mathcal{O}((m+n)r^2)$
  - 10:    $\dot{A}, \dot{B} := \beta \dot{A}_{\text{prev}} + (I - A_L A_L^\top) \dot{A}, \beta \dot{B}_{\text{prev}} + \dot{B}$   $\ll \mathcal{O}(n+m)r^2$
  - 11:    $\dot{A}, \dot{B} := \text{ProjectLR}(\text{OrthoLR}(A_L, B_R, \dot{A}, \dot{B}, r), A_L, B_R)$   $\ll \mathcal{O}((m+n)r^2 + r^3)$
  - 12:    $A_{\text{HB}}, B_{\text{HB}} := [\dot{A}, A_L], [B_R, \dot{B}]$
  - 13:    $A_L, B := \text{RetractionLR}([\dot{A}, A_L] [B_R, -\eta(\dot{B} + \gamma B_R)])$   $\ll \mathcal{O}((m+n)r^2 + r^3)$
  - 14: **return**  $A_L, B$ .
-

---

## C RIEMANN-SGD ALGORITHM

In Algorithm 7, we present the RiemannSGD version of the fine-tuning algorithm. The first steps are dedicated to the computation of the optimal initialization using `BackPropRSVD` (Algorithm 4). After initialization, the algorithm follows an Adam-like procedure adapted to the Riemannian setting: it maintains exponentially smoothed estimates of momentum and gradient norms ( $4th$  step). The overall cycle encapsulates the computation of the Riemannian gradient at the current point (steps 7–8), the update of the adaptive momentum terms via `ProjectLR` (Algorithm 2), and the retraction step that projects the updated tangent direction back onto the manifold (steps 15–17). This approach has asymptotical complexity of  $\mathcal{O}((m+n)r^2+r^3)$  and  $(2(q+1)+\text{max\_iters})$  amount of backward calls.

---

### Algorithm 7 RiemannSGD with LOI

---

**Require:** Weights  $W \in \mathbb{R}^{m \times n}$ , rank  $r \in \mathbb{N}$ , step size  $\eta$ , momentum coefficient  $\beta$ , oversampling parameter  $p$ , power-step parameter  $q$ , `simulate_Adam`, Adam momentum coefficient  $\gamma$ .

**Ensure:** Tuning parameters  $\Delta W^* \in \mathcal{M}_r$ .

**Function:**

```

1:  $A_L, \_, B^\top := \text{BackPropRSVD}(2r, p, q, \mathcal{L}, W)$  //  $\mathcal{O}((m+n)r^2)$ 
2:  $A_L, B^\top := A_L[:, :r], B[:, r:]^\top$ 
3:  $W' := W - A_L B^\top$  //  $\mathcal{O}(mn)$ 
4:  $A_{\text{HB}}, B_{\text{HB}}, S_A, S_B := 0$ 
5: for  $i := 0, \dots, \text{max\_iters}$  do //  $\mathcal{O}(nr^2)$ 
6:    $B_R := \text{qr}(B) \cdot \mathbb{Q}$ 
7:    $\dot{A} := \nabla_{Z_1} \mathcal{L}(W + Z_1 B_R^\top + A_L Z_2^\top) |_{Z_1=0, Z_2=B}$ 
8:    $\dot{B} := \nabla_{Z_2} \mathcal{L}(W + Z_1 B_R^\top + A_L Z_2^\top) |_{Z_1=0, Z_2=B}$ 
9:    $\dot{A}_{\text{prev}}, \dot{B}_{\text{prev}} := \text{ProjectLR}((A_{\text{HB}}, B_{\text{HB}}), A_L, B_R)$  //  $\mathcal{O}((m+n)r^2)$ 
10:   $\dot{A}, \dot{B} := \beta \dot{A}_{\text{prev}} + (1-\beta)(I - A_L A_L^\top) \dot{A}, \beta \dot{B}_{\text{prev}} + (1-\beta) \dot{B}$  //  $\mathcal{O}(nr^2)$ 
11:  if simulate_Adam then
12:     $S_A, S_B := \gamma \|\dot{A}\|_F + (1-\gamma) S_A, \gamma \|\dot{B}\|_F + (1-\gamma) S_B$  //  $\mathcal{O}((m+n)r)$ 
13:     $\dot{A}, \dot{B} := \dot{A}/S_A, \dot{B}/S_B$  //  $\mathcal{O}((m+n)r)$ 
14:   $U, \Sigma, V^\top := \text{RetractionLR}(\left[ \eta \dot{A}, A_L \right], \left[ B_R, \eta \dot{B} + B \right])$  //  $\mathcal{O}((m+n)r^2+r^3)$ 
15:   $A_{\text{HB}}, B_{\text{HB}} := [\dot{A}, A_L], [B_R, \dot{B}]$ 
16:   $A_L, B := U, \Sigma V^\top$  //  $\mathcal{O}(nr^2)$ 
17: return  $A_L, B$ .
```

---

Table 2: The average accuracy (in %) among 8 tasks of fine-tuned Llama 3.2-1b using different SGD variants, tested on Commonsense Reasoning benchmark. The «R-» prefix stands for Riemannian, the postfix «-LOI» stands for locally optimal initialization. LoRA rank is set to 16. In all experiments except for «-LOI» the  $A$  factor in  $AB^\top$  has orthonormal columns and  $B$  is zero.

Task	BoolQ	PIQA	SIQA	hella-swag	wino-grande	ARC-E	ARC-C	OBQA	All
Initialization									
Raw	40.1	55.4	50.3	25.8	50.0	61.9	41.8	42.8	46.0
LoRA	64.0 $\pm$ 0.2	75.3 $\pm$ 0.3	69.6 $\pm$ 0.4	82.2 $\pm$ 0.1	53.0 $\pm$ 1.0	75.4 $\pm$ 1.7	56.5 $\pm$ 0.6	67.1 $\pm$ 2.3	67.9 $\pm$ 0.4
LoRA-LOI	64.9 $\pm$ 0.4	76.8 $\pm$ 0.3	71.2 $\pm$ 2.2	84.1 $\pm$ 0.1	56.7 $\pm$ 0.6	77.0 $\pm$ 4.2	61.1 $\pm$ 3.6	68.9 $\pm$ 3.3	70.1 $\pm$ 1.6
RSLoRA	64.5 $\pm$ 0.2	77.2 $\pm$ 0.3	68.1 $\pm$ 0.8	86.3 $\pm$ 0.1	58.2 $\pm$ 0.5	76.5 $\pm$ 1.5	58.6 $\pm$ 2.9	70.0 $\pm$ 1.9	70.0 $\pm$ 0.8
RiemannLoRA	<b>65.1</b> $\pm$ 0.4	<b>77.7</b> $\pm$ 0.3	<b>69.1</b> $\pm$ 1.5	<b>85.8</b> $\pm$ 0.2	<b>58.5</b> $\pm$ 0.4	<b>74.4</b> $\pm$ 1.9	<b>56.9</b> $\pm$ 1.3	<b>69.3</b> $\pm$ 2.5	<b>69.6</b> $\pm$ 0.9
RiemannLoRA-LOI	<b>65.2</b> $\pm$ 0.4	<b>79.4</b> $\pm$ 0.4	<b>75.6</b> $\pm$ 0.4	<b>87.3</b> $\pm$ 0.1	<b>62.4</b> $\pm$ 0.4	<b>79.9</b> $\pm$ 0.8	<b>63.6</b> $\pm$ 1.9	<b>73.8</b> $\pm$ 0.5	<b>73.4</b> $\pm$ 0.3

## D DATASET

### D.1 COMMONSENSE REASONING

Following the approach of Hu et al. (2023), we combine the training datasets from all 8 tasks to form the final training set and evaluate performance on each task’s individual test dataset. We adopt

---

918 queries structure from Hu et al. (2023) to Llama 3.2 instruct template. The prompt structure is  
919 demonstrated in Table 3.  
920

## 921 E HYPERPARAMETERS 922

923 For each optimization method, we carefully preselected the learning rate (step size). For the Rie-  
924 mannion method, we set the momentum to 0.9, and for the Muon method, we set the momentum to  
925 0.95. In both cases, we additionally tuned the weight-decay hyperparameter (see Algorithm 5). The  
926 final choices of these hyperparameters are summarized in Table 4.  
927

928 In the Commonsense reasoning benchmark, the hyperparameters are reported in Table 5 for the  
929 SGD-like methods and in Table 4 for the Adam-like methods. For the LOI method, we selected the  
930 following hyperparameters: the backprop RSVD oversampling parameter was set to  $r_k = 16$ , the  
931 backprop powerstep parameter was set to  $q = 1$ , and the multiplier  $\alpha$  was set to  $\frac{-0.01}{\sqrt{r}}$ .  
932

933 The non-tuned hyperparameters used for experiments on the Commonsense Reasoning dataset are  
934 presented in 6.  
935

## 936 F SUBJECT-DRIVEN GENERATION 937

938 **Training details** We used Stable Diffusion-2-base model with a batch size of 4 for all experiments.  
939 For both methods, we set the betas to 0.9 and 0.999 and the weight decay to 0.1. In all variations of  
940 our approach, we set  $q = 15$ ,  $p = \text{rank}$  and  
941

942  $\alpha = -\frac{1}{\sqrt{\text{rank}}}$ . We used a learning rate of 2e-5 to train both our and the LoRA models, which were  
943 used to generate images shown in Figures 2 and 3.  
944

945 **Evaluation details** We used the DreamBooth dataset, which contains 25 different prompts and 30  
946 various concepts. Due to computational costs, we only used half of the proposed concepts to evaluate  
947 metrics: can, candle, cat, cat2, colorful\_sneaker, dog2, dog3, dog5, dog6, dog7,  
948 dog8, fancy\_boot, grey\_sloth\_plushie, pink\_sunglasses, vase. To measure the  
949 similarity between the original concept and the generated images, we used Image Similarity (IS): we  
950 synthesized 30 images for the base prompt «a photo of a V\*» and 10 images for each of 25 editing  
951 prompts (like «a V\* on top of a dirt road»). We then measured the average pairwise cosine similarity  
952 with reference photos of the concept using the DINO model. To check the correspondence between  
953 the generated images and the textual prompts, we calculated Text Similarity (TS): we evaluated each  
954 concept with each of 25 prompts, synthesizing 10 images per prompt, and calculating the average  
955 pairwise cosine similarity with the prompts using the CLIP ViTB/32 model.  
956

957 **Additional results** Figure 3 shows an additional visual comparison of our method and LoRA  
958 after 600 training steps. As can be seen, our method learns the concept much faster than the original  
959 LoRA, while preserving editing capabilities.  
960

## 961 G ADDITIONAL EXPERIMENTAL RESULTS 962

### 963 G.1 ABLATION STUDY ON INITIALIZATION 964

965 The fine-tuning optimization was carried out by AdamW (Loshchilov & Hutter (2019)) First of all,  
966 for every approach tested we preselected a suitable optimization step sizes. The list of all selected  
967 hyper-parameters for each method is specified in the Appendix E. Table 7 contains the accuracy of  
968 the trained model’s responses on the test dataset. The notation «LoRA-A» in the table means the  
969 utilization of the vanilla LoRA with non-zero  $A$ , the notation «LoRA-B» means the same with  
970 non-zero  $B$ , the notations «Stiefel-A» and «Stiefel-B» indicate vanilla LoRA with orthonormal  
971 initialization (for matrix  $A$  and  $B$  respectively). The naming of «Stiefel-both» involves taking the  
factors  $A, B$  with orthonormal columns and with initialization like in (15).

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

Table 3: The structure of the queries for the Commonsense reasoning dataset

Task	Role	Fine-tuning Data Template
BoolQ	system	Please answer the following question with True or False. Follow the answer format, full answer not needed.
	user	Question: [QUESTION] Answer format: True/False
	assistant	The correct answer is [ANSWER]
PIQA	system	Please choose the correct solution to the question. Follow the answer format, full answer not needed.
	user	Question: [QUESTION] Solution1: [SOLUTION_1] Solution2: [SOLUTION_2] Answer format: Solution1/Solution2
	assistant	The correct answer is [ANSWER]
SIQA	system	Please choose the correct answer to the question based on the context provided. Follow the answer format, full answer not needed.
	user	Context: [CONTEXT] Question: [QUESTION] A: [ANSWER_A] B: [ANSWER_B] C: [ANSWER_C] Answer format: A/B/C
	assistant	The correct answer is [ANSWER]
hellaswag	system	Please choose the correct ending to complete the given sentence. Follow the answer format, full answer not needed.
	user	[ACTIVITY_LABEL]: [CONTEXT] Ending1: [ENDING_1] Ending2: [ENDING_2] Ending3: [ENDING_3] Ending4: [ENDING_4] Answer format: Ending1/Ending2/Ending3/Ending4
	assistant	The correct answer is [ANSWER]
winogrande	system	Please choose the correct answer to fill in the blank to complete the given sentence. Follow the answer format, full answer not needed.
	user	Sentence: [SENTENCE] Option1: [OPTION_1] Option2: [OPTION_2] Answer format: Option1/Option2
	assistant	The correct answer is [ANSWER]
ARC-e & ARC-c	system	Please choose the correct answer to the question. Follow the answer format, full answer not needed.
	user	Question: [QUESTION] Answer1: [ANSWER_1] Answer2: [ANSWER_2] Answer3: [ANSWER_3] Answer4: [ANSWER_4] Answer format: Answer1/Answer2/Answer3/Answer4
	assistant	The correct answer is [ANSWER]
OBQA	system	Please choose the correct answer to the question. Follow the answer format, full answer not needed.
	user	Question: [QUESTION] Answer1: [ANSWER_1] Answer2: [ANSWER_2] Answer3: [ANSWER_3] Answer4: [ANSWER_4] Answer format: Answer1/Answer2/Answer3/Answer4
	assistant	The correct answer is [ANSWER]

1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079

Table 4: The parameters for different Adam variants for fine-tuning on the Commonsense reasoning dataset

Optimizer	learning rate	weight decay
Adam	0.0002	0.01
DoRA	0.0003	0.01
Muon	0.0005	0.2
DoneRITE	0.0005	0.0001
RPrecAdamW	0.0005	0.5
Riemannion	0.0001	0.00316

Table 5: The parameters for different SGD variants (RiemannLoRA with `simulate_Adam` flag disabled) for fine-tuning on the Commonsense reasoning dataset

Optimizer	learning rate
LoRA	0.1
LoRA-LOI	0.06
RSLoRA	0.1
RiemannLoRA	0.1
RiemannLoRA-LOI	0.07

Table 6: Other non-tuned hyperparameter configurations for experiments on Commonsense reasoning dataset

Dataset	Commonsense reasoning
Hyperparameters	
Rank $r$	16
Dropout	0.05
LR Scheduler	Linear
Batch size	64
Epochs	2
Warmup ratio	0.1

Table 7: The average accuracy among 8 tasks of fine-tuned Llama 3.2-1b via Adam using different LoRA initialization variants, tested on Commonsense reasoning benchmark. «LOI» stands for locally optimal initialization. LoRA rank is set to 16.

Task Initialization	BoolQ	PIQA	SIQA	hella-swag	wino-grande	ARC-E	ARC-C	OBQA	All
Raw	40.1	55.4	50.3	25.8	50.0	61.9	41.8	42.8	46.0
LoRA-A	66.2	80.2	76.0	87.0	65.1	78.4	63.6	76.3	74.1
LoRA-B	65.9	77.9	74.2	82.9	61.2	73.9	60.6	72.2	71.1
Pissa	66.4	80.1	75.6	87.6	63.1	77.6	64.3	75.0	73.7
Stiefel-A	65.4	79.5	74.9	87.2	62.4	79.3	62.3	75.5	73.3
Stiefel-B	66.4	80.6	76.0	87.5	64.3	78.5	64.9	74.6	74.1
Stiefel-both	66.3	79.7	75.8	86.4	64.0	78.3	62.4	73.9	73.4
LOI	65.6	81.3	75.7	87.7	65.7	77.9	65.0	75.2	74.3
LoRA-GA	65.4	77.1	74.9	84.5	58.3	75.4	61.9	72.2	71.2

1080  
 1081  
 1082  
 1083  
 1084  
 1085  
 1086  
 1087  
 1088  
 1089  
 1090  
 1091  
 1092  
 1093  
 1094  
 1095  
 1096  
 1097  
 1098  
 1099  
 1100  
 1101  
 1102  
 1103  
 1104  
 1105  
 1106  
 1107  
 1108  
 1109  
 1110  
 1111  
 1112  
 1113  
 1114  
 1115  
 1116  
 1117  
 1118  
 1119  
 1120  
 1121  
 1122  
 1123  
 1124  
 1125  
 1126  
 1127  
 1128  
 1129  
 1130  
 1131  
 1132  
 1133

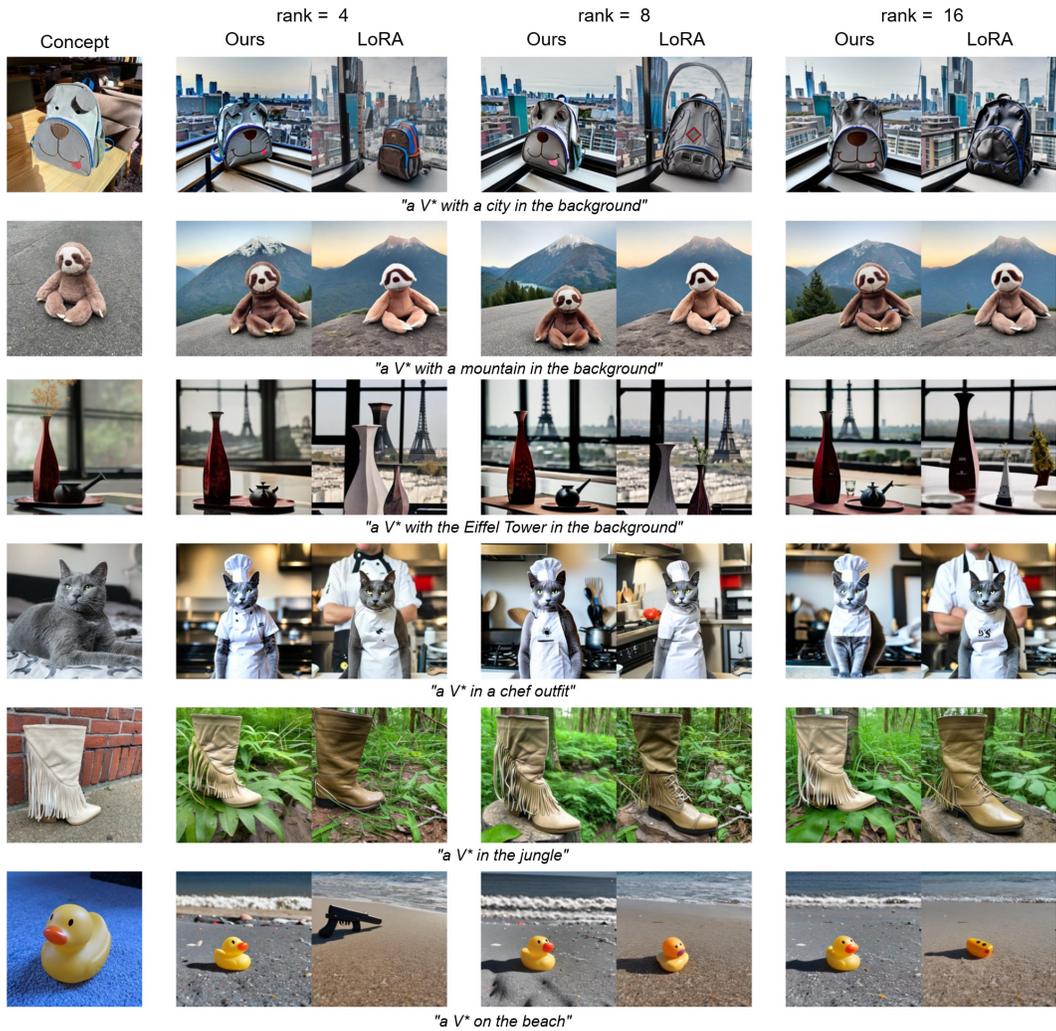


Figure 3: Additional visual comparison of our method and LoRA, checkpoint 600

Table 8: The average accuracy among 8 tasks of fine-tuned Llama 3.2-1b using different LoRA initialization variants, tested on Commonsense reasoning benchmark. «LOI» stands for locally optimal initialization.

Optimizer	Initialization	rank	Accuracy
Raw			46.0
Riemannion	random	32	74.7 $\pm$ 0.2
Riemannion	LOI	32	74.9 $\pm$ 0.3
LoRA	zero	32	73.4 $\pm$ 1.1
Riemannion	random	64	74.6 $\pm$ 0.3
Riemannion	LOI	64	75.1 $\pm$ 0.3
LoRA	zero	64	73.1 $\pm$ 0.1

Table 9: The average accuracy among 8 tasks of fine-tuned Llama 3.2-8b using different LoRA initialization variants, tested on Commonsense reasoning benchmark. «LOI» stands for locally optimal initialization.

Optimizer	Initialization	rank	Accuracy
Raw			74.5
Riemannion	random	32	86.9
Riemannion	LOI	32	88.2
LoRA(Adam)	zero	32	87.2
Riemannion	random	64	86.2
Riemannion	LOI	64	87.9
LoRA(Adam)	zero	64	87.0

## H COMPUTATIONAL COST

To measure execution time, we used a virtual server instantiated on a compute node equipped with an Intel Xeon Gold 6240R CPU and Nvidia Tesla V100 GPU. The virtual machine was provisioned with 8 CPU cores, 16 GB of RAM and a single Tesla V100 GPU. Experiments were executed using the Hugging Face `transformers` library; timing corresponds to the execution time of the `trainer.train()` call. Measurements were conducted according to the following procedure. For each combination of method, LoRA adapter rank, and batch size, four timing runs were performed and the minimum execution time was selected. The measured quantity was the time spent on 16 optimizer steps. Time spent on initialization was excluded. For each combination of rank and batch size, execution times for each of the methods were measured sequentially. In the present work, rank refers to the rank of the LoRA adapter, and method denotes one of two optimization schemes: Adam (baseline) and Riemannion (proposed method).

Figure 4 shows the relative increase in per-step execution time of the proposed Riemannion method compared with Adam, computed as  $((T_{\text{Riemannion}} - T_{\text{Adam}}) / T_{\text{Adam}})$ . Table 10 shows the average time cost for each epoch. Table 11 presents a comparison of both the theoretical time and memory complexity of several methods with our algorithm (more methods provided in Yen et al. (2024))

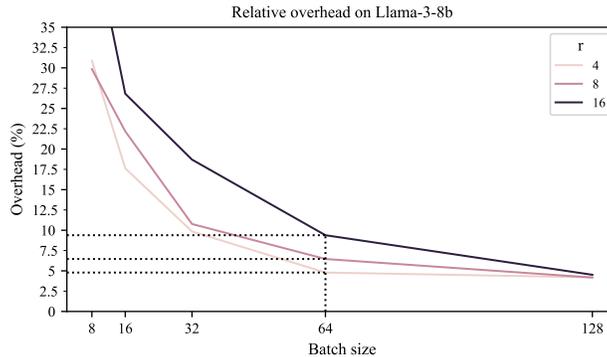


Figure 4: Relative Time Cost (calculated as  $(T_{\text{Riemannion}} - T_{\text{Adam}}) / T_{\text{Adam}}$ ) of Riemannion vs. Adam during Llama 3-8B fine-tuning, as a function of LoRA rank and batch size.

## I LMO BOUNDS

Let  $A \in \mathcal{T}_W \mathcal{M}_r$  be a vector from the tangent space. Consider the problem

$$\begin{aligned} \langle A, X \rangle &\rightarrow \max_X, \\ \text{s.t. } \|X\|_2 &\leq 1, \end{aligned} \tag{21}$$

whose (any) solution will be denoted by  $X_* = \text{Ortho}(A)$ . Denote  $f_* = \langle A, X_* \rangle$  the optimal value.

1188 Table 10: Time Cost of Riemannion vs. Adam during Llama 3-8B fine-tuning, as a function of  
 1189 LoRA rank and batch size.

r	batch size	AdamW	Riemannion
8	8	27.32	35.48
8	16	50.41	61.60
8	32	97.46	107.96
8	64	190.14	202.42
8	128	373.89	389.50
16	8	31.59	46.41
16	16	57.82	73.33
16	32	111.42	132.26
16	64	217.11	237.50
16	128	431.31	450.79
32	8	35.08	64.81
32	16	61.68	94.39
32	32	114.21	148.63
32	64	219.34	257.81
32	128	433.32	476.11

1204 Table 11: Time and space complexity comparison for LoRA optimization.

Method	Time	Memory
LoRA (Adam)	$O((m+n)r)$	$O((m+n)r)$
LoRA-RITE	$O((m+n)r^2)$	$O((m+n)r+r^2)$
Riemannion (our proposed)	$O((m+n)r^2)$	$O((m+n)r+r^2)$

1211 Consider the restricted (LMO) problem

$$\begin{aligned}
 1213 & \langle A, X \rangle \rightarrow \max_X, \\
 1214 & \text{s.t. } \|X\|_2 \leq 1, \\
 1215 & X \in \mathcal{T}_W \mathcal{M}_r.
 \end{aligned} \tag{22}$$

1217 Let  $X^*$  denote an optimizer of the task derived in (22) and  $f^* = \langle A, X^* \rangle$  its optimal value.

1218 **Proposition I.1.** Assume  $A \in \mathcal{T}_W \mathcal{M}_r$  and let  $P$  be the orthogonal projector onto  $\mathcal{T}_W \mathcal{M}_r$ . If

$$1220 \quad \|X_* - PX_*\|_2 \leq \varepsilon$$

1221 for some  $\varepsilon \geq 0$ , then

$$1222 \quad \frac{1}{1+\varepsilon} f_* \leq f^* \leq f_*$$

1225 *Proof.* Since  $A \in \mathcal{T}_W \mathcal{M}_r$ , we have  $PA = A$  and  $P = P^\top$ . By assumption  $\|X_* - PX_*\|_2 \leq \varepsilon$   
 1226 and, because  $\|X_*\|_2 = 1$ , the reverse triangle inequality yields

$$1227 \quad \left| \|X_*\|_2 - \|PX_*\|_2 \right| \leq \varepsilon,$$

1228 hence

$$1229 \quad 1 - \varepsilon \leq \|PX_*\|_2 \leq 1 + \varepsilon.$$

1232 Define the normalized matrix

$$1233 \quad Y := \frac{PX_*}{\|PX_*\|_2} \in \mathcal{T}_W \mathcal{M}_r, \quad \|Y\|_2 = 1.$$

1234 By definition of  $f^*$  (the maximum of  $\langle A, \cdot \rangle$  over  $\{X \in \mathcal{T}_W \mathcal{M}_r : \|X\|_2 \leq 1\}$ ) we have

$$1237 \quad f^* \geq \langle A, Y \rangle = \frac{\langle A, PX_* \rangle}{\|PX_*\|_2} = \frac{\langle PA, X_* \rangle}{\|PX_*\|_2} = \frac{\langle A, X_* \rangle}{\|PX_*\|_2} = \frac{f_*}{\|PX_*\|_2}.$$

1239 Using the inequality  $\|PX_*\|_2 \leq 1 + \varepsilon$  we obtain the claimed lower bound

$$1241 \quad f^* \geq \frac{f_*}{1+\varepsilon}.$$

1242 The upper bound  $f^* \leq f_*$  is trivial because the feasible set of the task derived in (22) is a subset of  
 1243 the feasible set of the problem described in (21). Combining the two inequalities gives  
 1244

$$1245 \frac{1}{1+\varepsilon} f_* \leq f^* \leq f_*,$$

1246 which completes the proof.  $\square$

1248 **Proposition I.2.** (*Upper bound on the projection residual.*) *With the notation above, the residual*

$$1249 \varepsilon := \|X_* - PX_*\|_2$$

1250 *satisfies  $\varepsilon \leq 1$ .*

1252 *Proof.* Write the thin SVD of the base point  $W$  as  $W = U_r \Sigma_r V_r^\top$ , where  $U_r \in \text{St}(m, r)$  and  
 1253  $V_r \in \text{St}(n, r)$  have  $r$  orthonormal columns. Let  $U_\perp$  and  $V_\perp$  be orthonormal complements so that  
 1254  $[U_r \ U_\perp] \in \mathbb{R}^{m \times 2r}$  and  $[V_r \ V_\perp] \in \mathbb{R}^{n \times 2r}$  are orthogonal matrices (here we assume  $m, n \geq 2r$ ; if  $m$   
 1255 or  $n$  is smaller the argument is adapted in the obvious way).  
 1256

1257 Because  $A \in \mathcal{T}_W \mathcal{M}_r$ , one can write (using a suitable  $2r \times 2r$  block representation) (Vandereycken,  
 1258 2013)

$$1259 A = [U_r \ U_\perp] \begin{bmatrix} \tilde{\mathcal{A}} & \tilde{\mathcal{B}} \\ \tilde{\mathcal{C}} & 0 \end{bmatrix} \begin{bmatrix} V_r^\top \\ V_\perp^\top \end{bmatrix},$$

1261 where  $\tilde{\mathcal{A}}, \tilde{\mathcal{B}}, \tilde{\mathcal{C}} \in \mathbb{R}^{r \times r}$  (the  $(2, 2)$  block vanishes for tangent-space elements of this form). The rank  
 1262 of  $A$  is therefore at most  $2r$ .  
 1263

1264 Let  $X_* = \text{Ortho}(A)$  be a maximizer of (21); represent  $X_*$  in the same block basis,

$$1265 X_* = [U_r \ U_\perp] \begin{bmatrix} \mathcal{A} & \mathcal{B} \\ \mathcal{C} & \mathcal{D} \end{bmatrix} \begin{bmatrix} V_r^\top \\ V_\perp^\top \end{bmatrix},$$

1266 where the  $2r \times 2r$  block matrix is orthogonal,

$$1269 \begin{bmatrix} \mathcal{A} & \mathcal{B} \\ \mathcal{C} & \mathcal{D} \end{bmatrix} \in \mathcal{O}(2r), \quad (23)$$

1271 where  $\mathcal{O}(2r)$  denotes the orthogonal group in  $\mathbb{R}^{2r \times 2r}$ . By Lemma I.1 below we have

$$1273 PX_* = [U_r \ U_\perp] \begin{bmatrix} \mathcal{A} & \mathcal{B} \\ \mathcal{C} & 0 \end{bmatrix} \begin{bmatrix} V_r^\top \\ V_\perp^\top \end{bmatrix}.$$

1275 Therefore

$$1276 X_* - PX_* = [U_r \ U_\perp] \begin{bmatrix} 0 & 0 \\ 0 & \mathcal{D} \end{bmatrix} \begin{bmatrix} V_r^\top \\ V_\perp^\top \end{bmatrix} = U_\perp \mathcal{D} V_\perp^\top,$$

1278 and hence

$$1279 \|X_* - PX_*\|_2 = \|\mathcal{D}\|_2.$$

1281 Because the block matrix in (23) is orthogonal, its lower block

$$1282 \begin{bmatrix} \mathcal{C} \\ \mathcal{D} \end{bmatrix} \in \mathbb{R}^{2r \times r}$$

1283 has orthonormal columns. Thus  $\|\frac{\mathcal{C}}{\mathcal{D}}\|_2 = 1$ , and in particular  $\|\mathcal{D}\|_2 \leq 1$ . Consequently  $\varepsilon =$   
 1284  $\|X_* - PX_*\|_2 \leq 1$ , as claimed.  $\square$

1287 **Lemma I.1.** *Let*

$$1288 W = U_r \Sigma_r V_r^\top, \quad U_r \in \text{St}(m, r), \quad V_r \in \text{St}(n, r),$$

1289 *and write an arbitrary matrix  $X \in \mathbb{R}^{m \times n}$  in the block basis determined by  $[U_r \ U_\perp]$  and  $[V_r \ V_\perp]$ ,*

$$1291 X = [U_r \ U_\perp] \begin{bmatrix} \mathcal{A} & \mathcal{B} \\ \mathcal{C} & \mathcal{D} \end{bmatrix} \begin{bmatrix} V_r^\top \\ V_\perp^\top \end{bmatrix}.$$

1293 *Then the orthogonal projector  $P$  onto  $\mathcal{T}_W \mathcal{M}_r$  acts as*

$$1294 PX = [U_r \ U_\perp] \begin{bmatrix} \mathcal{A} & \mathcal{B} \\ \mathcal{C} & 0 \end{bmatrix} \begin{bmatrix} V_r^\top \\ V_\perp^\top \end{bmatrix}.$$

1295

1296 *Proof.* One convenient expression for the orthogonal projector onto the tangent space at  $W =$   
 1297  $U_r \Sigma_r V_r^\top$  is

$$P(\cdot) = (\cdot) - (I - U_r U_r^\top) (\cdot) (I - V_r V_r^\top).$$

1299 Applying this to  $X$  written in the chosen block basis yields

$$\begin{aligned} 1300 \quad PX &= X - (I - U_r U_r^\top) X (I - V_r V_r^\top) = \\ 1301 &= [U_r \quad U_\perp] \begin{bmatrix} \mathcal{A} & \mathcal{B} \\ \mathcal{C} & \mathcal{D} \end{bmatrix} \begin{bmatrix} V_r^\top \\ V_\perp^\top \end{bmatrix} - [0 \quad U_\perp] \begin{bmatrix} \mathcal{A} & \mathcal{B} \\ \mathcal{C} & \mathcal{D} \end{bmatrix} \begin{bmatrix} 0 \\ V_\perp^\top \end{bmatrix} \\ 1302 &= [U_r \quad U_\perp] \left( \begin{bmatrix} \mathcal{A} & \mathcal{B} \\ \mathcal{C} & \mathcal{D} \end{bmatrix} - \begin{bmatrix} 0 & 0 \\ 0 & \mathcal{D} \end{bmatrix} \right) \begin{bmatrix} V_r^\top \\ V_\perp^\top \end{bmatrix} \\ 1303 &= [U_r \quad U_\perp] \begin{bmatrix} \mathcal{A} & \mathcal{B} \\ \mathcal{C} & 0 \end{bmatrix} \begin{bmatrix} V_r^\top \\ V_\perp^\top \end{bmatrix}, \end{aligned}$$

1309 which proves the lemma. □

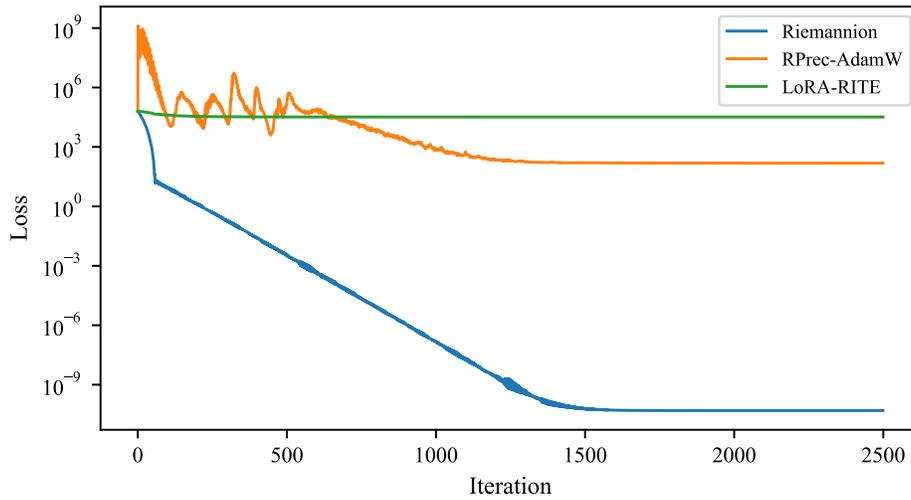
1311 **Remark I.1.** In our experiments we observed that the residual  $\varepsilon$  is typically small (in our reported  
 1312 runs it did not exceed 0.1).

## 1314 J ADDITIONAL EXPERIMENTS ON DIFFERENT TASKS

1316 Let  $M$  be a matrix with  $\text{rank}(M) = k$ . We aim to approximate it with a matrix  $M^* \in \mathcal{M}_r$  such  
 1317 that  $\text{rank}(M^*) = r > k$ . The task can be formulated as

$$1318 \quad \|M^* - M\| \rightarrow \min_{M^*}. \quad (24)$$

1319 Figure 5 shows the dependence of the loss on the epoch number for three different optimizers: Rie-  
 1320 mannion, PrecAdamW, and LoRA-RITE. This experiment shows that our optimizer demonstrates  
 1321 remarkable stability at singular points and achieves significantly lower loss than competing methods.



1322 Figure 5: Comparison of the performance of three optimizers on task (24).

## 1346 K LLM USAGE

1347 We used an LLM only for minor language polishing to improve readability and grammar; it was not  
 1348 involved in research ideation, methodology, analysis, or substantive writing, and all research ideas  
 1349 and arguments were developed entirely by the authors.