

LoRA MEETS SECOND-ORDER OPTIMIZATION: TOWARDS OPTIMAL LOW-RANK UPDATES

Anonymous authors

Paper under double-blind review

ABSTRACT

Low-rank fine-tuning is widely applied for the effective adaptation of large models. Most existing methods rely on low-rank matrix factorization, whose performance is limited by the condition number of the associated Jacobi operator. Although these methods are computationally efficient, their performance still falls short compared to full fine-tuning. To address this, we propose SoLoRA, which leverages an adaptive metric to find a low-rank approximation of the full fine-tuning gradient. This low-rank approximation can be viewed as an approximation of Hessian, effectively incorporating second-order information to achieve faster convergence and higher optimization efficiency. Furthermore, the low-rank approximation in SoLoRA is computationally simple and easy to implement, achieving a close approximation to the performance of full fine-tuning with almost no additional computational overhead. We conduct fine-tuning experiments on large language models and diffusion models, and the results consistently demonstrate that SoLoRA achieves superior performance advantages over state-of-the-art low-rank fine-tuning methods.

1 INTRODUCTION

Large language models (LLMs) (Liu et al., 2024a; Yang et al., 2024) and vision-language models (Achiam et al., 2023) have demonstrated outstanding performance in various applications, such as chatbot, image generation, and editing. With their strong generalization capabilities and versatility, they have been widely adopted for a range of downstream tasks. To better adapt LLMs to specific downstream tasks, it is often necessary to fine-tune their parameters. However, full fine-tuning is evidently expensive, incurring significant computational and storage costs. To address this, parameter-efficient fine-tuning (PEFT) has emerged to reduce the overhead of fine-tuning.

Low-Rank Adaptation (LoRA) (Hu et al., 2022) is a representative PEFT method. It assumes that weight updates during fine-tuning exhibit a low “intrinsic rank”. By freezing the pretrained weights and introducing two low-rank matrices, $\mathbf{B} \in \mathbb{R}^{m \times r}$ and $\mathbf{A} \in \mathbb{R}^{r \times n}$, for updates, LoRA reduces the number of trainable parameters. Compared to full fine-tuning, the number of trainable parameters in LoRA is $\mathcal{O}((m+n)r)$, where $r \ll \{m, n\}$, significantly lowering the number of trainable parameters, memory consumption, and fine-tuning costs. Owing to these advantages, LoRA and its numerous variants (Hu et al., 2022; Hayou et al., 2024; Zhang and Pilanci, 2024; Wang et al., 2024; Zhao et al., 2024; Zhu et al., 2024; Wang et al., 2025; Mo et al., 2025; Zhang et al., 2025b) have been widely applied in practical applications.

Although LoRA offers significant advantages, most existing fine-tuning algorithms are based on a factorization framework that updates the two low-rank factors separately. Such factorization-based methods are sensitive to the condition number of the low-rank factors, which can result in slow convergence. ScaledGD (Tong et al., 2021; Zhang and Pilanci, 2024) addresses this issue by introducing two preconditioners, effectively eliminating the dependency on the condition number and making its convergence rate condition-number-independent. However, ScaledGD still suffers from parameter redundancy, and its fine-tuning efficiency falls short of matching that of full-parameter fine-tuning.

LoRA-Pro (Wang et al., 2025) demonstrates that applying gradients \mathbf{G}_A and \mathbf{G}_B to the low-rank factors \mathbf{A} and \mathbf{B} is equivalent to performing full fine-tuning on the weight matrix \mathbf{W} with a low-rank gradient $\tilde{\mathbf{G}}$. Building on this insight, LoRA-Pro reduces the discrepancy between $\tilde{\mathbf{G}}$ and the full

fine-tuning gradient \tilde{G} by solving the optimization problem $\min \|\tilde{G} - G\|_F^2$, thereby bridging the performance gap between LoRA and full fine-tuning. LoRA-Pro employs the standard metric inherited from the Euclidean space of the weight matrices to approximate \tilde{G} . However, approximation is often more effective under a weighted metric rather than the standard metric. For example, AdaGrad (Duchi et al., 2011) and SOAP (Vyas et al., 2025) leverage historical gradient information to adaptively adjust the step size of each gradient component, effectively utilizing weighted metrics in the Euclidean space of the weight matrix. K-FAC (Martens and Grosse, 2015; Eschenhagen et al., 2023) uses a weighted metric based on the Kronecker product to approximate the Hessian, thereby constructing an efficient preconditioner.

Inspired by this, we propose a novel algorithm called **Second-Order Low-Rank Adaption (SoLoRA)**, which aims to further narrow the performance gap between low-rank fine-tuning and full fine-tuning. SoLoRA leverages an adaptive metric derived from AdaGrad (Duchi et al., 2011) and SOAP (Vyas et al., 2025) to identify a low-rank approximation of the full fine-tuning gradient. Notably, this low-rank approximation can also serve as an approximation of the Hessian, enabling SoLoRA to effectively incorporate second-order information from the loss function for faster convergence. Moreover, the optimal low-rank approximation identified by SoLoRA does not directly depend on the full fine-tuning gradient, making SoLoRA simple and easy to implement. Experiments on GPT-2 and diffusion models demonstrate that SoLoRA, by adopting a weighted metric-based approximation, outperforms both standard metric-based approximations and existing low-rank fine-tuning methods, achieving superior performance.

2 LOW-RANK FINE-TUNING OF LARGE LANGUAGE MODELS

In this section, we revisit existing low-rank fine-tuning methods from a fresh theoretical perspective, highlighting their gaps compared to full fine-tuning. Based on this analysis, we discuss the limitations of these low-rank fine-tuning algorithms and elucidate their fundamental distinctions.

2.1 RETHINKING LOW-RANK FINE-TUNING: CONNECTIONS AND LIMITATIONS

As a representative parameter-efficient fine-tuning method, low-rank fine-tuning works by freezing the pretrained weights $\mathbf{W}_0 \in \mathbb{R}^{m \times n}$ and assuming that the weight update \mathbf{W} exhibits a low-rank structure during downstream task adaptation. Consequently, the adaptation process is formulated as a low-rank constrained optimization problem:

$$\min_{\mathbf{W} \in \mathbb{R}^{m \times n}} \mathcal{L}(\mathbf{W}_0 + \mathbf{W}), \quad \text{subject to } \text{rank}(\mathbf{W}) = r,$$

where $\mathcal{L}(\cdot)$ denotes the training loss function and $r \ll \min\{m, n\}$. Proximal gradient descent is a widely adopted method for solving the above low-rank optimization problem by updating the weight matrix via

$$\mathbf{W}_{t+1} = \mathcal{H}_r(\mathbf{W}_t - \alpha_t \nabla_{\mathbf{W}_t} \mathcal{L}(\mathbf{W}_0 + \mathbf{W}_t)),$$

where \mathcal{H}_r represents the r -truncated singular value decomposition (SVD) applied to each weight matrix, α_t is the learning rate of \mathbf{W}_t . This requires performing SVD on every layer at each optimization step, which has a computational complexity $O(m^3)$, leading to [high-computation cost](#).

LoRA and its variants (Hu et al., 2022; Wang et al., 2024; Hayou et al., 2024; Liu et al., 2024b; Wang et al., 2025; Zhang et al., 2025b; Yen et al.; Zhang et al., 2025a) train the network directly via a low-rank factorization, thereby avoiding the expensive SVD computation at each training step. These methods aim to solve the following non-convex optimization problem based on the factorization:

$$\min_{\mathbf{W} \in \mathbb{R}^{m \times n}} \mathcal{L}(\mathbf{W}_0 + \mathbf{W}), \quad \text{subject to } \mathbf{W} = \mathbf{B}\mathbf{A},$$

where $\mathbf{B} \in \mathbb{R}^{m \times r}$, $\mathbf{A} \in \mathbb{R}^{r \times n}$. Here, we define $\mathcal{G}([\mathbf{B}, \mathbf{A}]) = \mathbf{W}$ as a generator that constructs weight matrices from the low-rank factors. Under this definition, the optimization problem can be reformulated as:

$$\min_{\mathbf{B} \in \mathbb{R}^{m \times r}, \mathbf{A} \in \mathbb{R}^{r \times n}} \mathcal{L}(\mathbf{W}_0 + \mathcal{G}([\mathbf{B}, \mathbf{A}])).$$

For factorization-based gradient algorithms, the updates can be expressed as follows, leveraging the chain rule:

$$[\mathbf{B}_{t+1}, \mathbf{A}_{t+1}] = [\mathbf{B}_t, \mathbf{A}_t] - \eta_t J_{\mathcal{G}}^*([\mathbf{B}_t, \mathbf{A}_t]) \nabla_{\mathbf{W}_t} \mathcal{L}(\mathbf{W}_0 + \mathcal{G}([\mathbf{B}_t, \mathbf{A}_t])), \quad (1)$$

where $J_{\mathcal{G}}^*$ is the adjoint of the Jacobian operator of \mathcal{G} and η_t is the learning rate of B_t and A_t . To further analyze the gap between low-rank fine-tuning and full fine-tuning, we return to the update of the weight matrix W . By applying the generator operator \mathcal{G} to both sides of (1), we get

$$\mathcal{G}([B_{t+1}, A_{t+1}]) = \mathcal{G}([B_t, A_t] - \eta_t J_{\mathcal{G}}^*([B_t, A_t]) \nabla_{W_t} \mathcal{L}(W_0 + \mathcal{G}([B_t, A_t])))$$

To facilitate comparison with the gradient descent algorithm based on the weight matrix W , we perform a Taylor expansion around $[B_t, A_t]$,

$$W_{t+1} \approx W_t - \alpha_t J_{\mathcal{G}}([B_t, A_t]) J_{\mathcal{G}}^*([B_t, A_t]) \nabla_{W_t} \mathcal{L}(W_0 + W_t), \quad (2)$$

where $J_{\mathcal{G}}([B_t, A_t])[\cdot, \cdot] : [\mathbb{R}^{m \times r}, \mathbb{R}^{r \times n}] \rightarrow \mathbb{R}^{m \times n}$ is the Jacobian operator. From this update form, it becomes clear that, compared with full fine-tuning, a key limitation of low-rank fine-tuning lies in the explicit dependence of the factor gradients on $J_{\mathcal{G}} J_{\mathcal{G}}^*$, whose condition number is determined by the condition numbers of the low-rank factors B and A (Chen et al., 2019; Chi et al., 2019). This dependency introduces potential instability during training, particularly when fine-tuning complex neural networks or large language models, which often results in performance degradation (Hayou et al., 2024; Zhang and Pilanci, 2024).

2.2 PRECONDITIONED LOW-RANK ADAPTION FINE-TUNING

Under the widely adopted generator form $\mathcal{G}([B, A]) = BA$, the Jacobian operator $J_{\mathcal{G}}([B_t, A_t])[\cdot, \cdot] : [\mathbb{R}^{m \times r}, \mathbb{R}^{r \times n}] \rightarrow \mathbb{R}^{m \times n}$ and its adjoint operator $J_{\mathcal{G}}^*([B_t, A_t])(\cdot) : \mathbb{R}^{m \times n} \rightarrow [\mathbb{R}^{m \times r}, \mathbb{R}^{r \times n}]$ are given by

$$J_{\mathcal{G}}([B_t, A_t])[P, Q] = PA_t + B_t Q,$$

for any factor pairs $[P, Q] \in [\mathbb{R}^{m \times r}, \mathbb{R}^{r \times n}]$, and

$$J_{\mathcal{G}}^*([B_t, A_t])(C) = [CA_t^\top, B_t^\top C],$$

for any matrices $C \in \mathbb{R}^{m \times n}$. For detailed derivations and additional information regarding the Jacobian, please refer to Appendix D.1. Substituting $J_{\mathcal{G}}$ and $J_{\mathcal{G}}^*$ into (2), we can rewrite (2) as

$$\begin{aligned} W_{t+1} &\approx W_t - \alpha_t G_t \cdot A_t^\top A_t - \alpha_t B_t B_t^\top \cdot G_t \\ &\approx (B_t - \eta_t G_t \cdot A_t^\top)(A_t - \eta_t B_t^\top \cdot G_t) \\ &= (B_t - \eta_t G_{B_t})(A_t - \eta_t G_{A_t}), \end{aligned}$$

where $G_t = \nabla_{W_t} \mathcal{L}(W_0 + W_t)$, $G_{B_t} = \nabla_{B_t} \mathcal{L}(W_0 + W_t)$ and $G_{A_t} = \nabla_{A_t} \mathcal{L}(W_0 + W_t)$ are the gradient of the loss function \mathcal{L} with respect to W_t , B_t and A_t . This formulation aligns with the update rule of standard LoRA (Vanilla LoRA) (Hu et al., 2022), in which the factors B and A are updated with the same learning rate. Consequently, the convergence rate of standard LoRA depends on the condition number of $J_{\mathcal{G}}$.

To mitigate this dependence on the condition number of $J_{\mathcal{G}}$, several improvements have been proposed. LoRA+ (Hayou et al., 2024) enhances feature learning efficiency by scaling the update $\eta_t G_{B_t}$ with a factor of 2^4 when training Roberta (Liu et al., 2019) with LeCun initialization (LeCun et al., 2002). This adjustment can be regarded as applying a constant preconditioner on G_{B_t} . However, LoRA+ does not completely eliminate the dependence on the condition number of $J_{\mathcal{G}}$. Imbalance-Regularized LoRA (Zhu et al., 2024) further alleviates the impact of $J_{\mathcal{G}}$ by introducing regularization terms on the low-rank factors B_t and A_t , which effectively reduce parameter redundancy. Going further, Riemannian preconditioned LoRA (Zhang and Pilanci, 2024) applies $r \times r$ preconditioners $(A_t A_t^\top)^{-1}$ and $(B_t^\top B_t)^{-1}$ to G_{B_t} and G_{A_t} respectively, making the update of W_t equivalent to projecting the gradient onto the row space of A_t and the column space of B_t . Specifically,

$$\begin{aligned} B_{t+1} A_{t+1} &= (B_t - \eta_t G_{B_t} \cdot (A_t A_t^\top)^{-1})(A_t - \eta_t (B_t^\top B_t)^{-1} \cdot G_{A_t}) \\ &\approx W_t - \alpha_t G_t \cdot A_t^\top (A_t A_t^\top)^{-1} A_t - \alpha_t B_t (B_t^\top B_t)^{-1} B_t^\top \cdot G_t \\ &= W_t - \alpha_t \text{Proj}_{\text{row}(A_t)}(G_t) - \alpha_t \text{Proj}_{\text{col}(B_t)}(G_t). \end{aligned}$$

Although Riemannian preconditioned LoRA alleviates the influence of condition number of $J_{\mathcal{G}}$ to some extent via two preconditioners, it still has an important limitation: it ignores the projection onto the intersection of the row space of A_t and the column space of B_t . Specifically, the term

$B_t(B_t^\top B_t)^{-1}B_t^\top \cdot G_t \cdot A_t^\top(A_t A_t^\top)^{-1}A_t$ is omitted, which causes the update direction to deviate from the steepest descent direction. To compensate for the missing information in this cross subspace, LoRA-Pro (Wang et al., 2025) proposes solving

$$\min_{\Delta_{B_t}, \Delta_{A_t}} \|G_t - (B_t \Delta_{A_t} + \Delta_{B_t} A_t)\|_F^2,$$

to more accurately approximate the full fine-tuning gradient and obtain an equivalent low-rank gradient for the factors B_t and A_t . Optimizing this objective yields the factor updates

$$\begin{cases} \Delta_{B_t} = [I - B_t(B_t^\top B_t)^{-1}B_t^\top]G_t A_t^\top (A_t A_t^\top)^{-1} - B_t X_t, \\ \Delta_{A_t} = (B_t^\top B_t)^{-1}B_t^\top G_t + X_t A_t, \end{cases}$$

for some $X_t \in \mathbb{R}^{r \times r}$. The corresponding update of the weight matrix is

$$\begin{aligned} & B_t \Delta_{A_t} + \Delta_{B_t} A_t \\ &= B_t(B_t^\top B_t)^{-1}B_t^\top G_t + B_t X_t A_t + [I - B_t(B_t^\top B_t)^{-1}B_t^\top]G_t A_t^\top (A_t A_t^\top)^{-1}A_t - B_t X_t A_t \\ &= (B_t(B_t^\top B_t)^{-1}B_t^\top)G_t + G_t(A_t^\top(A_t A_t^\top)^{-1}A_t) - (B_t(B_t^\top B_t)^{-1}B_t^\top)G_t(A_t^\top(A_t A_t^\top)^{-1}A_t) \\ &= \text{Proj}_{\text{col}(B_t)}(G_t) + \text{Proj}_{\text{row}(A_t)}(G_t) - \text{Proj}_{\text{col}(B_t) \cap \text{row}(A_t)}(G_t) = \mathcal{P}_{\mathbb{T}_t}(G_t). \end{aligned} \tag{3}$$

where \mathcal{M}_r is the Riemannian manifold of all rank r matrices, and \mathbb{T}_t denotes the tangent space of \mathcal{M}_r at the point W_t . By Proposition D.2, $\mathcal{P}_{\mathbb{T}_t}(G_t)$ is the orthogonal projection of G_t onto \mathbb{T}_t .

Although LoRA-Pro is capable of finding a low-rank approximation of the full fine-tuning gradient under a standard metric, [it incurs higher computational overhead due to solving a Sylvester equation at each step](#). (Lu, 1971; Dmytryshyn et al., 2025). In comparison, gradient approximations based on weighted metrics are often more effective, as they better utilize the second-order information of the loss function. For instance, classical methods such as the Broyden–Fletcher–Goldfarb–Shanno (BFGS) algorithm (Fletcher, 2000) and the Gauss–Newton method (Nocedal and Wright, 2006) both benefit from weighted metrics. Previous research has demonstrated that weighted metrics can significantly enhance algorithmic efficiency across a variety of problems (Duchi et al., 2011; Bian et al., 2024). Motivated by this, we propose designing a novel weighted metric to further improve the approximation of full fine-tuning gradients and to fully exploit the second-order information embedded in the loss function, enabling a more effective low-rank approximation.

3 THE PROPOSED ALGORITHMS

Empirical evidence suggests that the weighted metric are often more effective than the standard metric in deep learning. For instance, AdaGrad (Duchi et al., 2011; Shazeer and Stern, 2018) and SOAP (Gupta et al., 2018; Morwani et al., 2024; Vyas et al., 2025) adaptively adjust the step size of each gradient component based on historical gradient information, which is equivalent to using a weighted metric for the weight matrix. Similarly, K-FAC (Martens and Grosse, 2015; Eschenhagen et al., 2023) employs a Kronecker product-based weighted metric to approximate the Hessian, thereby constructing an efficient preconditioner. In this section, we introduce a novel weighted metric and derive a low-rank approximation of the full fine-tuning gradient G based on this metric. This low-rank approximation can be viewed as an approximation of the Hessian, allowing our algorithm to effectively exploit the second-order information of the loss function, thereby narrowing the gap between the performance of low-rank fine-tuning and full fine-tuning.

3.1 CONSTRUCTION OF THE ADAPTIVE METRIC

The core idea of AdaGrad (Duchi et al., 2011) and Adam (Kingma and Ba, 2014) is to construct a weighted operator h_t through the outer product of gradients, followed by a diagonalization operation. Specifically, by vectorizing the gradient matrix $G_t \in \mathbb{R}^{m \times n}$ into $g_t \in \mathbb{R}^{mn \times 1}$, the linearized weighted operator h_t is denoted as $h_t = (h_{t-1}^2 + \text{diag}(g_t g_t^\top))^{\frac{1}{2}}$, where $\text{diag}(\cdot)$ extracts the diagonal elements of the matrix. This operator h_t is then used to define a new weighted inner product, under which the gradient descent update to linearized weight w_t is derived $w_{t+1} = w_t - g_t/h_t$. AdaGrad employs the elements of h_t to rescale the gradient element-wise. However, when applied to gradient

updates in matrix form, this element-wise rescaling approach ignores the structural information of the matrix, i.e., the relationships between rows and columns. To better utilize matrix structures, Shampoo (Gupta et al., 2018) adopts the Kronecker product to approximate the construction of the weighted matrix. Specifically, Shampoo constructs two matrices, $\mathbf{L}_t = \mathbf{L}_{t-1} + \mathbf{G}_t \mathbf{G}_t^\top$ and $\mathbf{R}_t = \mathbf{R}_{t-1} + \mathbf{G}_t^\top \mathbf{G}_t$, and defines a weighted inner product based on these matrices. Under this inner product, the matrix update is performed by $\mathbf{W}_{t+1} = \mathbf{W}_t - \mathbf{L}_t^{-\frac{1}{4}} \mathbf{G}_t \mathbf{R}_t^{-\frac{1}{4}}$. SOAP (Morwani et al., 2024; Vyas et al., 2025) further improves upon Shampoo by noting that the square root operation in Shampoo is equivalent to running Adafactor (Shazeer and Stern, 2018) in the eigenbasis of the Shampoo preconditioner. To enhance the computational efficiency of Shampoo, SOAP runs Adam in the eigenbasis of the Shampoo preconditioner. However, frequent eigen-decomposition computations result in high computational costs. To balance leveraging matrix structural information and maintaining computational efficiency, we propose a hybrid weighted inner product that is easier to implement, better aligns with matrix structures, and fully utilizes the relationships between matrix rows and columns.

Specifically, we define the weighted factors \mathbf{L}_t and \mathbf{R}_t as follows:

$$\begin{aligned} \mathbf{L}_t &= \text{diag}(\mathbf{l}_t / \sqrt{\|\mathbf{l}_t\|_1}) \quad \text{with} \quad \mathbf{l}_t = \beta_1 \mathbf{l}_{t-1} + (1 - \beta_1) \sum_{j=1}^n (\mathbf{G}_t \odot \mathbf{G}_t)_{i,j}, \\ \mathbf{R}_t &= \text{diag}(\mathbf{r}_t / \sqrt{\|\mathbf{r}_t\|_1}) \quad \text{with} \quad \mathbf{r}_t = \beta_2 \mathbf{r}_{t-1} + (1 - \beta_2) \sum_{i=1}^m (\mathbf{G}_t \odot \mathbf{G}_t)_{i,j}, \end{aligned} \quad (4)$$

where \odot denotes the Hadamard (element-wise) product, $\|\cdot\|_1$ denotes the l_1 -norm, β_1, β_2 are decay factors in the range $[0, 1]$. The term $\sum_{j=1}^n (\mathbf{G}_t \odot \mathbf{G}_t)_{i,j}$ forms a vector of the diagonal elements of the matrix $\mathbf{G}_t \mathbf{G}_t^\top$, and similarly, $\sum_{i=1}^m (\mathbf{G}_t \odot \mathbf{G}_t)_{i,j}$ forms a vector of the diagonal elements of the matrix $\mathbf{G}_t^\top \mathbf{G}_t$. As stated in (Shazeer and Stern, 2018), $\mathbf{l}_t \mathbf{r}_t^\top$ is a rank-1 approximation of $\mathbf{G} \odot \mathbf{G}$, which is optimal with respect to the generalized Kullback-Leibler divergence. In this way, the memory requirement is reduced from $\mathcal{O}(mn)$ to $\mathcal{O}(m + n)$. At the same time, compared to Shampoo, the computational complexity of \mathbf{L}_t and \mathbf{R}_t is reduced to $\mathcal{O}(mn)$.

Based on \mathbf{L}_t and \mathbf{R}_t , we define an adaptive weighted inner product in $\mathbb{R}^{m \times n}$. For any $\mathbf{Y}, \mathbf{Z} \in \mathbb{R}^{m \times n}$, the adaptive weighted inner product is given by:

$$\langle \mathbf{Y}, \mathbf{Z} \rangle_{H_t} = \langle \mathbf{H}_t \mathbf{Y}, \mathbf{Z} \rangle = \langle \mathbf{L}_t^{\frac{1}{2}} \mathbf{Y} \mathbf{R}_t^{\frac{1}{2}}, \mathbf{Z} \rangle. \quad (5)$$

For any matrix $\mathbf{K} \in \mathbb{R}^{m \times n}$, the inverse operation of the operator \mathbf{H}_t is defined as

$$\mathbf{H}_t^{-1} \mathbf{K} = \mathbf{L}_t^{-\frac{1}{2}} \mathbf{K} \mathbf{R}_t^{-\frac{1}{2}}.$$

3.2 SECOND-ORDER LOW-RANK ADAPTION FOR FINE-TUNING

Based on the adaptive weighted inner product, we aim to incorporate second-order information into the update of the weight matrix \mathbf{W} . To achieve this, we consider the update of the weight matrix \mathbf{W} in step t . Let the update to \mathbf{W} at step t be denoted as Δ_t . In this step, we solve the problem:

$$\min_{\Delta_t} \mathcal{L}((\mathbf{W}_0 + \mathbf{W}_t) - \Delta_t).$$

We then expand the loss function $\mathcal{L}(\mathbf{W})$ around the point $\mathbf{W}_0 + \mathbf{W}_t$ using its second-order Taylor expansion. By utilizing the weighted inner product as an approximation of Hessian, the optimization problem can be formulated as:

$$\begin{aligned} & \arg \min_{\Delta_t} \mathcal{L}((\mathbf{W}_0 + \mathbf{W}_t) - \Delta_t) \\ & \approx \arg \min_{\Delta_t} \mathcal{L}(\mathbf{W}_0 + \mathbf{W}_t) - \langle \Delta_t, \mathbf{G}_t \rangle + \frac{1}{2} \langle \mathbf{H}_t \Delta_t, \Delta_t \rangle, \\ & = \arg \min_{\Delta_t} \mathcal{L}(\mathbf{W}_0 + \mathbf{W}_t) - \langle \Delta_t, \mathbf{H}_t^{-1} \mathbf{G}_t \rangle_{H_t} + \frac{1}{2} \langle \Delta_t, \Delta_t \rangle_{H_t} + \frac{1}{2} \langle \mathbf{H}_t^{-1} \mathbf{G}_t, \mathbf{H}_t^{-1} \mathbf{G}_t \rangle_{H_t} \\ & = \arg \min_{\Delta_t} \mathcal{L}(\mathbf{W}_0 + \mathbf{W}_t) + \frac{1}{2} \|\Delta_t - \mathbf{H}_t^{-1} \mathbf{G}_t\|_{H_t}^2. \end{aligned}$$

From this expression, it is evident that the optimization problem is equivalent to finding the optimal Δ_t for the following objective:

$$\min_{\Delta_t} \|\Delta_t - H_t^{-1} G_t\|_{H_t}^2. \quad (6)$$

Let the optimal update be denoted as Δ_t^{opt} . From the form of (6), it becomes clear that Δ_t^{opt} serves as an approximation of the Newton direction

$$-\Delta_t^{\text{opt}} \approx -\nabla^2 \mathcal{L}(W_0 + W_t) \cdot \nabla \mathcal{L}(W_0 + W_t).$$

Thus, the weight matrix is updated as $W_{t+1} = W_t - \Delta_t^{\text{opt}}$. The advantages of this update are evident:

- It completely eliminates the adverse effects of the condition number of the Jacobian operator J_G , thereby improving the stability of the algorithm.
- It effectively incorporates the second-order information of the loss function, enhancing optimization efficiency.

To further reduce memory consumption, we adopt the low-rank factorization strategy of LoRA, representing the update Δ_t in terms of updates to the low-rank factors A_t and B_t , denoted as Δ_{A_t} and Δ_{B_t} , respectively. As noted in (Wang et al., 2025), the changes in the factors A_t and B_t are intrinsically related to the updates in the weight matrix W_t , which can be expressed as

$$\Delta_t = \Delta_{B_t} A_t + B_t \Delta_{A_t}.$$

Therefore, the minimization problem (6) can be equivalently transformed into:

$$\min_{\Delta_{B_t}, \Delta_{A_t}} \|\Delta_{B_t} A_t + B_t \Delta_{A_t} - H_t^{-1} G_t\|_{H_t}^2. \quad (7)$$

To make the optimization process more explicit, we first rewrite (7) as:

$$\begin{aligned} & \arg \min_{\Delta_{B_t}, \Delta_{A_t}} \|\tilde{\mathcal{P}}_{\mathbb{T}_t}(\Delta_{B_t} A_t + B_t \Delta_{A_t} - H_t^{-1} G_t) + \tilde{\mathcal{P}}_{\mathbb{T}_t}^\perp(\Delta_{B_t} A_t + B_t \Delta_{A_t} - H_t^{-1} G_t)\|_{H_t}^2 \\ &= \arg \min_{\Delta_{B_t}, \Delta_{A_t}} \|\Delta_{B_t} A_t + B_t \Delta_{A_t} - \tilde{\mathcal{P}}_{\mathbb{T}_t}(H_t^{-1} G_t)\|_{H_t}^2 + \|\tilde{\mathcal{P}}_{\mathbb{T}_t}^\perp(H_t^{-1} G_t)\|_{H_t}^2, \end{aligned} \quad (8)$$

where $\tilde{\mathcal{P}}_{\mathbb{T}_t}^\perp(\cdot)$ denotes the projection onto the space orthogonal to the tangent space. This equivalence holds because $\Delta_{B_t} A_t + B_t \Delta_{A_t}$ lies in the tangent space \mathbb{T}_t (see Proposition D.4). This implies that, to find the optimal Δ_{B_t} and Δ_{A_t} , we ultimately need to solve the following equivalent problem:

$$\min_{\Delta_{B_t}, \Delta_{A_t}} \|\Delta_{B_t} A_t + B_t \Delta_{A_t} - \tilde{\mathcal{P}}_{\mathbb{T}_t}(H_t^{-1} G_t)\|_{H_t}^2. \quad (9)$$

Here, $\tilde{\mathcal{P}}_{\mathbb{T}_t}(H_t^{-1} G_t)$ represents the projection of $H_t^{-1} G_t$ onto \mathbb{T}_t , with its explicit form given as

$$\tilde{\mathcal{P}}_{\mathbb{T}_t}(L_t^{-\frac{1}{2}} G_t R_t^{-\frac{1}{2}}) = \tilde{P}_{B_t} L_t^{-\frac{1}{2}} G_t R_t^{-\frac{1}{2}} + L_t^{-\frac{1}{2}} G_t R_t^{-\frac{1}{2}} \tilde{Q}_{A_t} - \tilde{P}_{B_t} L_t^{-\frac{1}{2}} G_t R_t^{-\frac{1}{2}} \tilde{Q}_{A_t}, \quad (10)$$

where $\tilde{P}_{B_t} = B_t(B_t^\top L_t^{\frac{1}{2}} B_t)^{-1} B_t^\top L_t^{\frac{1}{2}}$ and $\tilde{Q}_{A_t} = R_t^{\frac{1}{2}} A_t^\top (A_t R_t^{\frac{1}{2}} A_t^\top)^{-1} A_t$. The detailed derivation is provided in Appendix D.3.

For problem (9), we provide its explicit solution in the following Theorem 3.1. For the proof of Theorem 3.1, please refer to Appendix D.3.

Theorem 3.1 (Optimal updates for low-rank factors). *Let $W_t = B_t A_t$ be a rank- r factorization at t -th step, and let $\tilde{\mathcal{P}}_{\mathbb{T}_t}(L_t^{-\frac{1}{2}} G_t R_t^{-\frac{1}{2}})$ denote the projection of the preconditioned gradient $L_t^{-\frac{1}{2}} G_t R_t^{-\frac{1}{2}}$ onto the tangent space \mathbb{T}_t at W_t . Consider the following optimization problem:*

$$\min_{\Delta_{B_t}, \Delta_{A_t}} \frac{1}{2} \|\Delta_{B_t} A_t + B_t \Delta_{A_t} - \tilde{\mathcal{P}}_{\mathbb{T}_t}(L_t^{-\frac{1}{2}} G_t R_t^{-\frac{1}{2}})\|_{H_t}^2, \quad (11)$$

where $\|\cdot\|_{H_t}$ is the norm induced by the operator H_t . Then the optimal solutions for Δ_{B_t} and Δ_{A_t} are given by

$$\begin{aligned} \Delta_{B_t}^{\text{opt}} &= [I - B_t(B_t^\top L_t^{\frac{1}{2}} B_t)^{-1} B_t^\top L_t^{\frac{1}{2}}] L_t^{-\frac{1}{2}} G_{B_t} (A_t R_t^{\frac{1}{2}} A_t^\top)^{-1} - B_t X_t, \\ \Delta_{A_t}^{\text{opt}} &= (B_t^\top L_t^{\frac{1}{2}} B_t)^{-1} G_{A_t} R_t^{-\frac{1}{2}} + X_t A_t, \end{aligned}$$

where $X_t \in \mathbb{R}^{r \times r}$ is an arbitrary matrix.

From Theorem 3.1, we observe that although G_t appears in the solution, it does not directly appear in the closed-form expression. Instead, the solution depends on the low-rank gradients G_{A_t} and G_{B_t} , ensuring low memory overhead. This efficient representation allows for straightforward gradient updates: first, compute the gradients using standard backpropagation, and then adjust Δ_{B_t} and Δ_{A_t} according to the closed-form solution. While Δ_{B_t} and Δ_{A_t} depend on X_t , the choice of X_t is critical for balancing the updates. Next, we minimize the weighted norm of the difference between the two update components, $\Delta_{B_t}A_t$ and $B_t\Delta_{A_t}$. This yields the optimal X_t in Theorem 3.2 (proof provided in Appendix D.3).

Once the matrix X_t is computed, Δ_{B_t} and Δ_{A_t} can be derived. Using the updates Δ_{B_t} and Δ_{A_t} , we propose Second-order Low-Rank Adaption (SoLoRA), summarized in Algorithm 1. The computational complexity is analyzed in Appendix C.

Theorem 3.2 (Optimal Solution for Balancing Matrix X_t). *Let $X_t \in \mathbb{R}^{r \times r}$. Consider the following optimization problem with respect to X_t ,*

$$\min_{X_t \in \mathbb{R}^{r \times r}} \frac{1}{2} \|\Delta_{B_t}A_t - B_t\Delta_{A_t}\|_{H_t}^2, \quad (12)$$

where Δ_{B_t} and Δ_{A_t} are functions of X_t given in Theorem 3.1. Then the optimal solution for X_t is given by

$$X_t^{opt} = -\frac{1}{2} (B_t^\top L_t^{\frac{1}{2}} B_t)^{-1} B_t^\top G_t A_t^\top (A_t R_t^{\frac{1}{2}} A_t^\top)^{-1}.$$

Algorithm 1 Second-order Low-Rank Adaption (SoLoRA) with SGD for Fine-tuning.

- 1: Initialize $B_1 = \mathbf{0}_{m \times r}$, $A_1 = \text{Kaiming uniform}_{r \times n}$, $l_0 = \mathbf{0}_m$, $r_0 = \mathbf{0}_n$, $\epsilon = 1e-6$.
 - 2: **for** $t = 1, \dots, T$ **do**
 - 3: $l_t = \beta_1 l_{t-1} + (1 - \beta_1) \sum_{j=1}^n (G_t \odot G_t)_{i,j}$, $L_t = \text{diag}(l_t / \sqrt{\|l_t\|_1})$.
 - 4: $r_t = \beta_2 r_{t-1} + (1 - \beta_2) \sum_{i=1}^m (G_t \odot G_t)_{i,j}$, $R_t = \text{diag}(r_t / \sqrt{\|r_t\|_1})$.
 - 5: $\Delta_{B_t} = \left[I - \frac{1}{2} B_t (B_t^\top L_t^{\frac{1}{2}} B_t)^{-1} B_t^\top L_t^{\frac{1}{2}} \right] L_t^{-\frac{1}{2}} G_{B_t} (A_t R_t^{\frac{1}{2}} A_t^\top)^{-1}$.
 - 6: $\Delta_{A_t} = (B_t^\top L_t^{\frac{1}{2}} B_t)^{-1} G_{A_t} R_t^{-\frac{1}{2}} \left[I - \frac{1}{2} R_t^{\frac{1}{2}} A_t^\top (A_t R_t^{\frac{1}{2}} A_t^\top)^{-1} A_t \right]$.
 - 7: $B_{t+1} = B_t - \eta_t \Delta_{B_t}$, $A_{t+1} = A_t - \eta_t \Delta_{A_t}$.
 - 8: **end for**
 - 9: **Note:** Add ϵI to matrix $B_t^\top L_t^{\frac{1}{2}} B_t$ if it is not invertible.
-

3.3 SECOND-ORDER LOW-RANK ADAPTION WITH MOMENTUM FOR FINE-TUNING.

First-order momentum methods, such as Adam and AdamW (Kingma and Ba, 2014; Loshchilov and Hutter, 2017), have been shown to be highly effective in stochastic optimization. By maintaining an exponential moving average of both the per-coordinate gradient statistics and the raw gradients, Adam stabilizes updates, reduces gradient variance, and minimizes sensitivity to manual learning rate tuning. To incorporate these advantages into our second-order low-rank adaptation framework, we integrate the exponential moving average of the gradients into SoLoRA. The enhanced method preserves the curvature-aware geometric properties of SoLoRA while inheriting the stability and adaptivity of Adam, resulting in more reliable and efficient fine-tuning. The pseudocode is present in Algorithm 2.

4 EXPERIMENTAL RESULTS

To evaluate the performance of our SoLoRA algorithm, we apply it to fine-tuning tasks for the large language model GPT-2 (see Section 4.1 and Appendix A) and diffusion models (see Appendix B). In the experiments, we compare two kinds of optimization algorithms: SGD-based algorithms and AdamW-based algorithms. The SGD-based algorithms include: LoRA with SGD optimizer (referred to as SGD) (Hu et al., 2022), Scaled GD (Zhang and Pilanci, 2024; Tong et al., 2021), LoRA-Pro with SGD optimizer (Wang et al., 2025), and our SoLoRA with SGD optimizer (Algorithm 1). The AdamW-based algorithms include: LoRA with AdamW optimizer (referred to as AdamW) (Hu et al.,

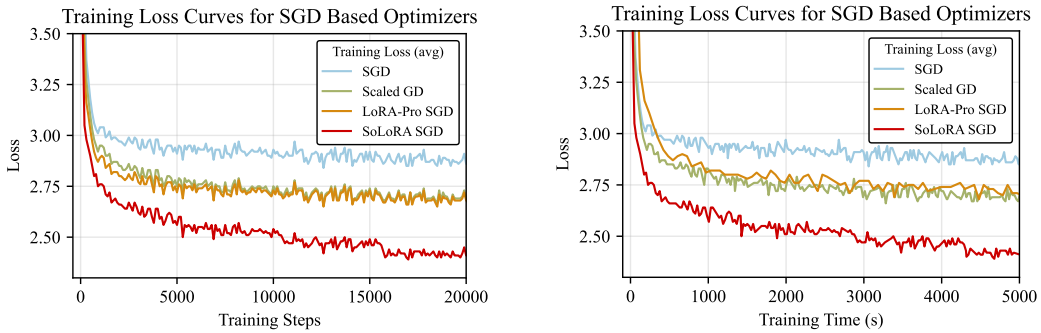
Algorithm 2 Second-order Low-Rank Adaption (SoLoRA) with Momentum for Fine-tuning.

- 1: Initialize moment $\mathbf{M}_0 = \mathbf{0}_{m \times n}$, $\mathbf{B}_1 = \mathbf{0}_{m \times r}$, $\mathbf{A}_1 = \text{Kaiming uniform}_{r \times n}$; $\mathbf{l}_0 = \mathbf{0}_m$, $\mathbf{r}_0 = \mathbf{0}_n$, weight decay λ , coefficients $\beta_1 = \beta_2$, and β_3 , $\epsilon = 1e - 6$.
- 2: **for** $t = 1, \dots, T$ **do**
- 3: $\mathbf{l}_t = \beta_1 \mathbf{l}_{t-1} + (1 - \beta_1) \sum_{j=1}^n (\mathbf{G}_t \odot \mathbf{G}_t)_{i,j}$, $\mathbf{L}_t = \text{diag}(\mathbf{l}_t / \sqrt{\|\mathbf{l}_t\|_1})$.
- 4: $\mathbf{r}_t = \beta_2 \mathbf{r}_{t-1} + (1 - \beta_2) \sum_{i=1}^m (\mathbf{G}_t \odot \mathbf{G}_t)_{i,j}$, $\mathbf{R}_t = \text{diag}(\mathbf{r}_t / \sqrt{\|\mathbf{r}_t\|_1})$.
- 5: $\mathbf{M}_t = \beta_3 \mathbf{M}_{t-1} + (1 - \beta_3) \mathbf{G}_t$.
- 6: $\Delta_{\mathbf{B}_t} = \left[\mathbf{I} - \frac{1}{2} \mathbf{B}_t \left(\mathbf{B}_t^\top \mathbf{L}_t^{\frac{1}{2}} \mathbf{B}_t \right)^{-1} \mathbf{B}_t^\top \mathbf{L}_t^{\frac{1}{2}} \right] \mathbf{L}_t^{-\frac{1}{2}} \mathbf{M}_t \mathbf{A}_t^\top \left(\mathbf{A}_t \mathbf{R}_t^{\frac{1}{2}} \mathbf{A}_t^\top \right)^{-1}$.
- 7: $\Delta_{\mathbf{A}_t} = \left(\mathbf{B}_t^\top \mathbf{L}_t^{\frac{1}{2}} \mathbf{B}_t \right)^{-1} \mathbf{B}_t^\top \mathbf{M}_t \mathbf{R}_t^{-\frac{1}{2}} \left[\mathbf{I} - \frac{1}{2} \mathbf{R}_t^{\frac{1}{2}} \mathbf{A}_t^\top \left(\mathbf{A}_t \mathbf{R}_t^{\frac{1}{2}} \mathbf{A}_t^\top \right)^{-1} \mathbf{A}_t \right]$.
- 8: $\mathbf{B}_{t+1} = (1 - \lambda \eta_t) \mathbf{B}_t - \eta_t \frac{\sqrt{1 - \beta_1^t}}{1 - \beta_3^t} \Delta_{\mathbf{B}_t}$, $\mathbf{A}_{t+1} = (1 - \lambda \eta_t) \mathbf{A}_t - \eta_t \frac{\sqrt{1 - \beta_1^t}}{1 - \beta_3^t} \Delta_{\mathbf{A}_t}$.
- 9: **end for**
- 10: **Note:** Add $\epsilon \mathbf{I}$ to matrix $\mathbf{B}_t^\top \mathbf{L}_t^{\frac{1}{2}} \mathbf{B}_t$ if it is not invertible.

2022), Scaled AdamW (Zhang and Pilanci, 2024), LoRA-Pro with AdamW optimizer (Wang et al., 2025), and our SoLoRA with AdamW optimizer (Algorithm 2). All experiments are implemented using PyTorch (Paszke et al., 2019) and conducted on NVIDIA GeForce RTX 4090 or 3090 GPUs.

4.1 GPT-2 FINE-TUNING

In this section, we conduct fine-tuning experiments on the GPT-2 model (Radford et al., 2019) using SoLoRA. First, we perform fine-tuning on the GPT-2 small model with ranks 16 and 64, evaluated on the E2E natural language generation challenge (Novikova et al., 2017). The results are shown in Table 1. The experimental setup follows (Zhang and Pilanci, 2024), but we independently tune the learning rate for each optimizer using grid search. As shown in Table 1, the model trained with SoLoRA outperforms all other methods across all evaluation metrics, regardless of whether the SGD or AdamW optimizer is used. To further validate the efficiency of SoLoRA, we compare the loss reduction trends when employing different optimizers under the same runtime and the same number of iterations. These results are illustrated in Figures 1 and 2. The findings demonstrate that SoLoRA achieves significantly faster loss reduction than other algorithms within the same runtime, thanks to its effective utilization of second-order information of the loss function.



(a) Training loss curve over training step when fine-tuning using SGD-based methods.

(b) Training loss curve over training time when fine-tuning using SGD-based methods.

Figure 1: Training loss GPT-2 small model ($r = 64$) fine-tuned using different SGD-based optimizers. Evaluation is conducted on E2E Natural Language Generation Challenge.

Optimizing low-rank factorization matrices presents inherent challenges, particularly when the weight matrix contains small singular values — a scenario that often arises with larger ranks, such as ranks 16 and 64 in this experiment. Under these conditions, the curvature of Hessian becomes very large, resulting in a high condition number and making the optimization problem ill-conditioned. Despite

these challenges, SoLoRA demonstrates superior performance in both computational efficiency and final evaluation metrics. This highlights the ability of SoLoRA to effectively mitigate the impact of J_G 's condition number while leveraging the second-order information from the loss function. To further evaluate SoLoRA, we conducted additional experiments on GPT-2 models of varying sizes with rank 4. The results are presented in Table 2 (see Appendix A), reaffirm the advantages of SoLoRA. Finally, we test the stability of SoLoRA under different learning rates, with the results shown in Figure 3 (see Appendix A). The experiments reveal that, compared to other algorithms, SoLoRA exhibits greater stability across varying ranks and learning rates.

Table 1: Scores of GPT-2 small model fine-tuned using different optimizers. Evaluation is conducted on E2E Natural Language Generation challenge.

rank	Method	E2E				
		BLEU	NIST	MET	ROUGE-L	CIDEr
16	SGD	65.4	8.07	40.7	67.0	2.07
	Scaled GD	68.8	8.75	45.0	69.2	2.39
	LoRA-Pro SGD	68.3	8.67	45.1	69.3	2.37
	SoLoRA SGD (ours)	70.0	8.82	46.6	71.6	2.53
	AdamW	69.5	8.77	46.4	71.2	2.48
	Scaled AdamW	69.8	8.79	46.5	71.7	2.51
	LoRA-Pro AdamW	69.7	8.73	46.8	71.7	2.51
	SoLoRA AdamW (ours)	70.2	8.85	46.6	71.9	2.52
	SGD	64.7	8.08	40.8	66.7	2.04
	Scaled GD	68.5	8.68	45.0	69.4	2.38
64	LoRA-Pro SGD	68.6	8.71	45.4	69.7	2.38
	SoLoRA SGD (ours)	70.1	8.85	46.7	71.8	2.53
	AdamW	69.6	8.76	46.7	71.5	2.50
	Scaled AdamW	70.0	8.83	46.4	71.5	2.50
	LoRA-Pro AdamW	70.0	8.82	46.6	71.5	2.51
	SoLoRA AdamW (ours)	70.2	8.84	46.8	72.1	2.52

5 CONCLUSION

This paper addresses the performance limitations of low-rank fine-tuning in efficiently adapting large models by proposing the second-order low-rank adaptation algorithm, **SoLoRA**. SoLoRA leverages an adaptive metric inspired by AdaGrad (Duchi et al., 2011) and SOAP (Vyas et al., 2025) to efficiently compute a low-rank approximation of the full fine-tuning gradient. This approximation, which can be viewed as an approximation of Hessian, effectively incorporates second-order information, accelerating convergence and improving optimization efficiency. Compared to existing low-rank fine-tuning methods, SoLoRA not only exploits second-order information but also completely eliminates the impact of the condition number of Jacobian operator. Moreover, as its low-rank approximation does not directly depend on the full gradient, SoLoRA is simpler and more efficient to implement. Experiments on GPT-2 and diffusion models consistently demonstrate that SoLoRA outperforms state-of-the-art low-rank fine-tuning methods. It achieves performance close to full fine-tuning while incurring almost no additional computational cost. This strongly demonstrates that second-order low-rank approximations based on our adaptive weighted metric provide a practical path to bridging the gap between parameter efficiency and optimal performance, paving the way for efficient and robust task transfer and personalized customization in large models.

Ethics statement This paper conforms with the ICLR Code of Ethics.

Reproducibility statement We are committed to the reproducibility of our research. To this end, we have made all source code, environmental configurations, and data access instructions available in the supplementary material. Furthermore, the key parameters for our experiments are provided in Table 3, Table 4, and Table 6 to facilitate the replication of our findings.

REFERENCES

- P-A Absil, Robert Mahony, and Rodolphe Sepulchre. Optimization algorithms on matrix manifolds. In *Optimization Algorithms on Matrix Manifolds*. Princeton University Press, 2009.
- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Fengmiao Bian, Jian-Feng Cai, and Rui Zhang. A preconditioned riemannian gradient descent algorithm for low-rank matrix recovery. *SIAM Journal on Matrix Analysis and Applications*, 45(4):2075–2103, 2024.
- Yuxin Chen, Yuejie Chi, Jianqing Fan, and Cong Ma. Gradient descent with random initialization: Fast global convergence for nonconvex phase retrieval. *Mathematical Programming*, 176(1):5–37, 2019.
- Yuejie Chi, Yue M Lu, and Yuxin Chen. Nonconvex optimization meets low-rank matrix factorization: An overview. *IEEE Transactions on Signal Processing*, 67(20):5239–5269, 2019.
- Andrii Dmytryshyn, Massimiliano Fasi, Nicholas J Higham, and Xiaobo Liu. Mixed-precision algorithms for solving the sylvester matrix equation. *arXiv preprint arXiv:2503.03456*, 2025.
- John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12(7), 2011.
- Runa Eschenhagen, Alexander Immer, Richard Turner, Frank Schneider, and Philipp Hennig. Kronecker-factored approximate curvature for modern neural network architectures. In *Advances in Neural Information Processing Systems (NIPS)*, volume 36, pages 33624–33655, 2023.
- Roger Fletcher. *Practical methods of optimization*. John Wiley & Sons, 2000.
- Yuchao Gu, Xintao Wang, Jay Zhangjie Wu, Yujun Shi, Yunpeng Chen, Zihan Fan, Wuyou Xiao, Rui Zhao, Shuning Chang, Weijia Wu, et al. Mix-of-show: Decentralized low-rank adaptation for multi-concept customization of diffusion models. In *Advances in Neural Information Processing Systems (NIPS)*, volume 36, pages 15890–15902, 2023.
- Vineet Gupta, Tomer Koren, and Yoram Singer. Shampoo: Preconditioned stochastic tensor optimization. In *International Conference on Machine Learning (ICML)*, pages 1842–1850. PMLR, 2018.
- Soufiane Hayou, Nikhil Ghosh, and Bin Yu. Lora+: Efficient low rank adaptation of large models. In *International Conference on Machine Learning (ICML)*, pages 17783–17806. PMLR, 2024.
- Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-free evaluation metric for image captioning. In *EMNLP (1)*, 2021.
- Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Advances in neural information processing systems (NIPS)*, volume 30, 2017.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. In *International Conference on Learning Representations (ICLR)*, volume 1, page 3, 2022.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Yann LeCun, Léon Bottou, Genevieve B Orr, and Klaus-Robert Müller. Efficient backprop. In *Neural networks: Tricks of the trade*, pages 9–50. Springer, 2002.
- Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024a.
- Shih-Yang Liu, Chien-Yi Wang, Hongxu Yin, Pavlo Molchanov, Yu-Chiang Frank Wang, Kwang-Ting Cheng, and Min-Hung Chen. Dora: Weight-decomposed low-rank adaptation. In *International Conference on Machine Learning (ICML)*, 2024b.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.

- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- CS Lu. Solution of the matrix equation $ax + xb = c$. *Electronics Letters*, 7(8):185–186, 1971.
- James Martens and Roger Grosse. Optimizing neural networks with kronecker-factored approximate curvature. In *International conference on machine learning (ICML)*, pages 2408–2417. PMLR, 2015.
- Zhanfeng Mo, Long-Kai Huang, and Sinno Jialin Pan. Parameter and memory efficient pretraining via low-rank riemannian optimization. In *International Conference on Learning Representations (ICLR)*, 2025.
- Depen Morwani, Itai Shapira, Nikhil Vyas, Sham M Kakade, Lucas Janson, et al. A new perspective on shampoo’s preconditioner. In *International Conference on Learning Representations (ICLR)*, 2024.
- Linyong Nan, Dragomir Radev, Rui Zhang, Amrit Rau, Abhinand Sivaprasad, Chiachun Hsieh, Xiangru Tang, Aadit Vyas, Neha Verma, Pranav Krishna, et al. Dart: Open-domain structured data record to text generation. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 432–447, 2021.
- Jorge Nocedal and Stephen J Wright. *Numerical optimization*. Springer, 2006.
- Jekaterina Novikova, Ondřej Dušek, and Verena Rieser. The e2e dataset: New challenges for end-to-end generation. *arXiv preprint arXiv:1706.09254*, 2017.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. In *Advances in neural information processing systems (NIPS)*, volume 32, 2019.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning (ICML)*, pages 8748–8763. PMLR, 2021.
- Noam Shazeer and Mitchell Stern. Adafactor: Adaptive learning rates with sublinear memory cost. In *International Conference on Machine Learning (ICML)*, pages 4596–4604. PMLR, 2018.
- Tian Tong, Cong Ma, and Yuejie Chi. Accelerating ill-conditioned low-rank matrix estimation via scaled gradient descent. *Journal of Machine Learning Research*, 22(150):1–63, 2021.
- Nikhil Vyas, Depen Morwani, Rosie Zhao, Itai Shapira, David Brandfonbrener, Lucas Janson, and Sham M Kakade. Soap: Improving and stabilizing shampoo using adam for language modeling. In *International Conference on Learning Representations (ICLR)*, 2025.
- Shaowen Wang, Linxi Yu, and Jian Li. Lora-ga: Low-rank adaptation with gradient approximation. In *Advances in Neural Information Processing Systems (NIPS)*, volume 37, pages 54905–54931, 2024.
- Zhengbo Wang, Jian Liang, Ran He, Zilei Wang, and Tieniu Tan. Lora-pro: Are low-rank adapters properly optimized? In *International Conference on Learning Representations (ICLR)*, 2025.
- Ke Wei, Jian-Feng Cai, Tony F Chan, and Shingyu Leung. Guarantees of riemannian optimization for low rank matrix recovery. *SIAM Journal on Matrix Analysis and Applications*, 37(3):1198–1222, 2016.
- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.
- Jui-Nan Yen, Si Si, Zhao Meng, Felix Yu, Sai Surya Duvvuri, Inderjit S Dhillon, Cho-Jui Hsieh, and Sanjiv Kumar. Lora done rite: Robust invariant transformation equilibration for lora optimization. In *The Thirteenth International Conference on Learning Representations*.
- Fangzhao Zhang and Mert Pilanci. Riemannian preconditioned lora for fine-tuning foundation models. In *International Conference on Machine Learning (ICML)*, 2024.
- Yilang Zhang, Bingcong Li, and Georgios B Giannakis. Reflora: Refactored low-rank adaptation for efficient fine-tuning of large models. *arXiv preprint arXiv:2505.18877*, 2025a.
- Yuanhe Zhang, Fanghui Liu, and Yudong Chen. Lora-one: One-step full gradient could suffice for fine-tuning large language models, provably and efficiently. In *International Conference on Machine Learning (ICML)*, 2025b.

Jiawei Zhao, Zhenyu Zhang, Beidi Chen, Zhangyang Wang, Anima Anandkumar, and Yuandong Tian. Galore: Memory-efficient llm training by gradient low-rank projection. In *International Conference on Machine Learning (ICML)*, pages 61121–61143. PMLR, 2024.

Zhenyu Zhu, Yongtao Wu, Quanquan Gu, and Volkan Cevher. Imbalance-regularized lora: A plug-and-play method for improving fine-tuning of foundation models. In *Adaptive Foundation Models: Evolving AI for Personalized and Efficient Learning*, 2024.

CONTENTS

A	Supplementary Experiments of GPT-2 Fine-tuning	13
A.1	Experimental Results For Different Datasets	13
A.2	Experimental Results For Different Model Size	13
A.3	Training Loss Curve Using Different Optimizers	14
A.4	Training Efficiency Comparison	14
A.5	Parameter Settings	15
B	Supplementary Experiments of Diffusion Model Fine-tuning	17
B.1	Evaluation Metrics of Diffusion Models	18
B.2	Experimental Results for Different LoRA Scaling Factors	18
B.3	Experimental Results for Different Learning Rates	19
C	Computational and Memory Complexity Analysis of SoLoRA	24
D	Proof of Theoretical Results	24
D.1	Computation of Jacobian	24
D.2	Orthogonal Projection to Tangent Space	25
D.3	Proofs of Theorem 3.1 and Theorem 3.2	27
D.4	Proof of the loss function exhibits a decreasing trend	28
E	Additional Experiments	30
F	The Use of Large Language Models (LLMs)	30

A SUPPLEMENTARY EXPERIMENTS OF GPT-2 FINE-TUNING

A.1 EXPERIMENTAL RESULTS FOR DIFFERENT DATASETS

To further validate the effectiveness of SoLoRA, we also conducted the experiments of GPT-2 fine-tuning on the DART (Nan et al., 2021) dataset. with the results provided in the following table.

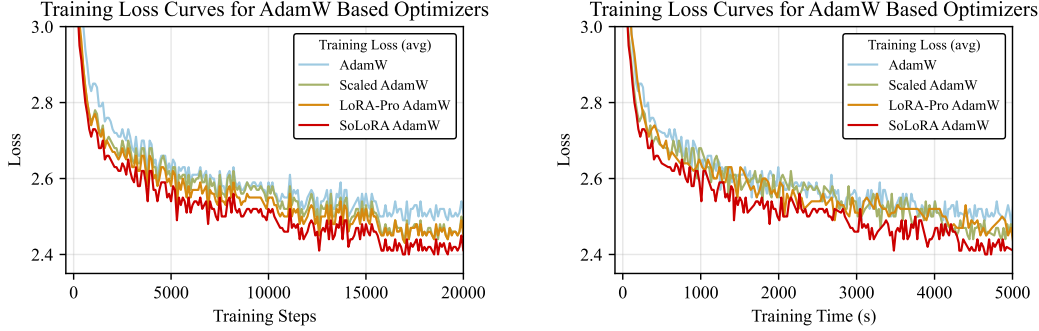
Scores of GPT-2 small model (rank=4) fine-tuned using different optimizers. Evaluation is conducted on DART dataset.

Methods	BLEU \uparrow	METEOR \uparrow	chrF++ \uparrow	TER \downarrow	BLEURT \uparrow
SGD	41.2	0.63	0.59	0.52	0.33
Scaled GD	43.8	0.66	0.61	0.50	0.38
LoRA-Pro SGD	44.1	0.66	0.61	0.50	0.38
SoLoRA SGD (ours)	44.6	0.66	0.62	0.49	0.39
AdamW	43.9	0.66	0.60	0.50	0.38
Scaled AdamW	44.8	0.67	0.62	0.49	0.40
LoRA-Pro AdamW	44.9	0.66	0.62	0.50	0.39
SoLoRA AdamW (ours)	45.4	0.67	0.60	0.49	0.40

A.2 EXPERIMENTAL RESULTS FOR DIFFERENT MODEL SIZE

To more comprehensively validate the advantages of SoLoRA, we conduct experiments not only on the small GPT-2 model but also on GPT-2 models of varying sizes for broader evaluation. All models are fine-tuned with rank 4, and the evaluation results are presented in Table 2. The specific parameter settings can be found in Table 3 and Table 4. By testing on models of different sizes, the experimental

results clearly demonstrate that SoLoRA significantly outperforms other algorithms, regardless of whether the SGD optimizer or the AdamW optimizer is used. This further confirms the effectiveness and stability of the SoLoRA algorithm, enabling it to maintain excellent performance across models of different sizes and under different optimizers.



(a) Training loss curve over training steps when fine-tuning using AdamW-based method.

(b) Training loss curve over training time when fine-tuning using AdamW-based method.

Figure 2: Training loss of GPT-2 small model ($r = 64$) fine-tuned using different AdamW-based optimizers. Evaluation is conducted on E2E Natural Language Generation Challenge.

A.3 TRAINING LOSS CURVE USING DIFFERENT OPTIMIZERS

To further explore the performance advantages of SoLoRA, we compare the runtime of different optimizers when fine-tuning large language models, with the results shown in Figures 1 and 2. These results strongly demonstrate the significant efficiency improvements achieved by the SoLoRA method in fine-tuning tasks. Additionally, Figure 3 illustrates the stability of SoLoRA under different learning rates. The experimental results show that SoLoRA maintains stable performance across a wide range of learning rates, which is crucial for parameter tuning in practical applications. To more comprehensively evaluate the stability of SoLoRA, we also compare it with Scaled AdamW and LoRA-Pro AdamW, under varying learning rates. The results are presented in Figure 3. The comparison reveals that SoLoRA exhibits superior stability across different ranks and learning rates. This indicates that SoLoRA is not only insensitive to changes in learning rates but also robust across varying LoRA ranks. As a result, it reduces the difficulty of hyperparameter tuning and enhances its practicality in fine-tuning.

A.4 TRAINING EFFICIENCY COMPARISON

To validate the training and inference efficiency of SoLoRA, we report in the table below the total training time required for all algorithms on the GPT-2 small model (rank 64). In addition, we recorded the relationship between training time and the number of steps in Figure 4.

Training and Inference Time of GPT-2 small model (rank=64) fine-tuned using different optimizers. Evaluation is conducted on E2E dataset.

Methods	SGD	Scaled GD	LoRA-Pro SGD	SoLoRA SGD
Total Training Time (Hours)	1.79	1.92	2.78	2.04
Total Inference Time (Hours)	1.86	1.87	1.58	1.89
Methods	AdamW	Scaled AdamW	LoRA-Pro AdamW	SoLoRA AdamW
Total Training Time (Hours)	1.79	1.94	2.93	2.04
Total Inference Time (Hours)	1.87	1.88	1.89	1.86

To further validate this, we record GPU memory consumption when the optimizer is called and after the backward is called (fine-tune GPT-2 small model with rank as 4). The results are summarized below.

GPU Memory occupied of GPT-2 small model (rank=4) fine-tuned using different optimizers. Evaluation is conducted on E2E dataset.

Methods	SGD	Scaled GD	LoRA-Pro SGD	SoLoRA SGD
During optimizer computation (MB)	1395.48	1395.57	1395.62	1401.01
After backward (MB)	1395.48	1395.48	1395.48	1395.62
Memory Complexity	0	0	0	m+n
Methods	AdamW	Scaled AdamW	LoRA-Pro AdamW	SoLoRA AdamW
During optimizer computation (MB)	1396.63	1396.72	1529.97	1462.51
After backward (MB)	1396.60	1396.60	1510.98	1457.12
Memory Complexity	(m+n)r	(m+n)r	2mn	mn+m+n

The results confirm that the memory usage of Algorithm 1 is comparable to other algorithms. Specifically, SoLoRA SGD (Algorithm 1) increases memory usage by only $(1401.01-1395.62)/1395.62 = 0.386\%$ compared to LoRA-SGD. However, with this slight increase in memory, Algorithm 1 demonstrates an effective improvement, as shown in Table 2.

Table 2: Scores of GPT-2 small and medium models ($r = 4$) fine-tuned using different optimizers. Evaluation is conducted on E2E Natural Language Generation challenge. See Appendix A.2 for experimental details.

Model	Method	E2E				
		BLEU	NIST	MET	ROUGE-L	CIDEr
GPT-2 small	SGD	54.8	4.56	34.0	63.3	1.29
	Scaled GD	68.5	8.72	45.5	69.4	2.40
	LoRA-Pro SGD	68.4	8.72	45.5	69.6	2.43
	SoLoRA SGD (ours)	69.5	8.77	46.5	71.5	2.50
	AdamW	69.1	8.75	46.0	70.5	2.47
	Scaled AdamW	69.5	8.80	46.2	70.9	2.48
	LoRA-Pro AdamW	69.2	8.73	45.9	70.8	2.47
	SoLoRA AdamW (ours)	70.0	8.84	46.3	71.3	2.50
	SGD	66.6	8.54	44.2	68.2	2.32
	Scaled GD	69.2	8.71	46.3	70.9	2.48
GPT-2 medium	LoRA-Pro SGD	69.7	8.77	46.5	70.9	2.50
	SoLoRA SGD (ours)	70.3	8.84	46.9	71.7	2.54
	AdamW	68.9	8.69	46.5	71.3	2.51
	Scaled AdamW	69.6	8.77	46.6	71.8	2.52
	LoRA-Pro AdamW	69.8	8.78	46.5	71.7	2.52
	SoLoRA AdamW (ours)	70.3	8.84	46.7	71.8	2.53

A.5 PARAMETER SETTINGS

To ensure the reproducibility of the experiments described in Section 4 and to facilitate verification and comparison by others, we provide the complete details of the experimental parameter settings. Tables 3 and 4 list the parameters used during the fine-tuning of GPT-2 models and the learning rates corresponding to different optimizers, respectively. Specifically, we conduct experiments with GPT-2 models of various sizes. “Rank 4 (M)” represents a medium-sized model using LoRA with rank 4, while “Rank 4”, “Rank 16”, and “Rank 64” represent small models using LoRA with ranks 4, 16, and 64, respectively. To ensure the fairness of the experimental setup, we follow the parameter settings in LoRA (Hu et al., 2022) and Riemannian Preconditioned LoRA (Zhang and Pilanci, 2024). However,

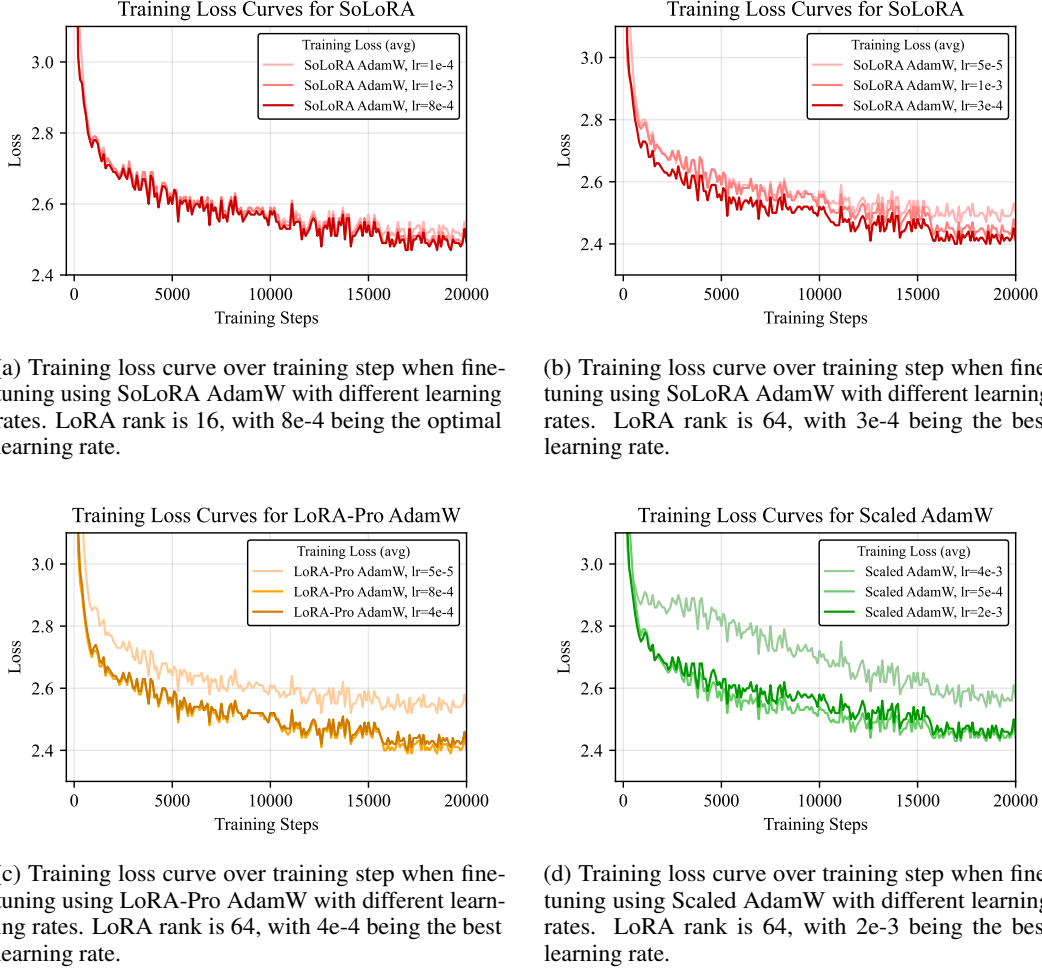


Figure 3: Training loss curve over training step of GPT-2 small model ($r = 16$ and 64) fine-tuned using different learning rates. Evaluation is conducted on E2E Natural Language Generation Challenge. Our optimizer is stable across different learning rates under varying ranks.

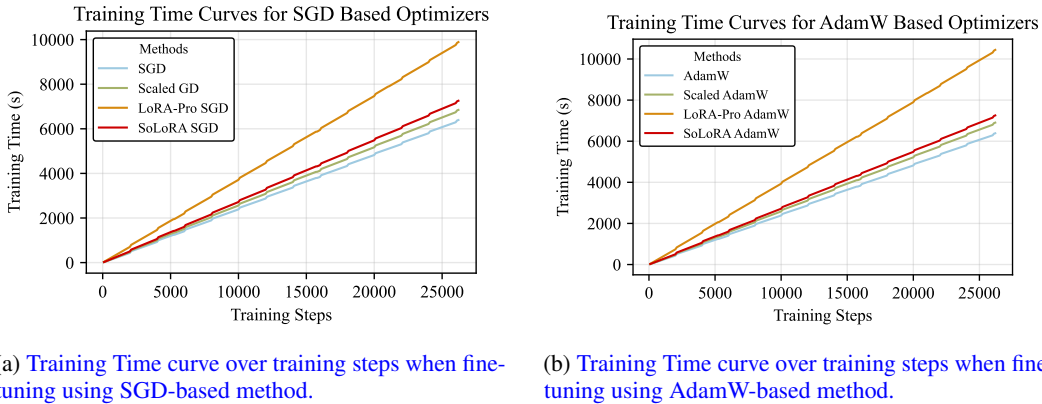


Figure 4: Training Time of GPT-2 small model ($r = 64$) fine-tuned using different optimizers. Evaluation is conducted on E2E Natural Language Generation Challenge.

considering the sensitivity of different optimizers to learning rates, we use a grid search strategy to independently tune the optimal learning rate for each optimizer. This ensures that each optimizer operates under its best-performing configuration, providing more objective and reliable experimental results.

Table 3: Training and Inference Configuration for GPT-2 Fine-tuning.

Training		LoRA α		Inference	
Parameter	Value	Parameter	Value	Parameter	Value
Dropout Probability	0.1				
Batch Size	8				
Number of Epochs	5	α (for Rank 4)	32	Beam Size	10
Warm-up Steps	500	α (for Rank 16)	32	Length Penalty	0.8
Learning Rate Scheduler	Linear	α (for Rank 64)	128	No Repeat Ngram Size	4
Label Smoothing	0.1				
Weight Decay	0.01				

Table 4: Core Optimizer Parameters for GPT-2 fine-tuning.

Methods	Learning Rate ($\times 10^{-3}$)				β_3	$\beta_1 = \beta_2$
	Rank 4	Rank 4 (M)	Rank 16	Rank 64		
SGD	90	90	200	90	/	/
Scaled GD	20	20	40	10	/	/
LoRA-Pro SGD	40	40	40	40	/	/
SoLoRA SGD	0.05	0.05	0.5	0.8	/	0.98
AdamW	0.2	0.2	0.2	0.2	0.9	0.999
Scaled AdamW	0.8	0.8	2	4	0.7	0.8
LoRA-Pro AdamW	0.1	0.1	0.2	0.4	0.9	0.999
SoLoRA AdamW	0.5	0.1	0.8	0.3	0.9	0.98

B SUPPLEMENTARY EXPERIMENTS OF DIFFUSION MODEL FINE-TUNING

As diffusion models increasingly become the mainstream method in image generation, LoRA plays an indispensable role in personalization and style transfer for specific characters. It demonstrates unique advantages, particularly in terms of parameter efficiency, training stability, and rapid convergence. To systematically evaluate the effectiveness of our optimizer SoLoRA in such personalized generation scenarios, we conduct experiments using the Mix-of-Show framework (Gu et al., 2023). This framework integrates Embedding Decomposed LoRA (EDLoRA) into the model, which further reduces the number of trainable parameters while maintaining expressive power. This design better aligns with the dual demands of computational efficiency and generalization stability in real-world applications. To ensure reproducibility and fair comparison, we follow the training and inference settings from (Zhang and Pilanci, 2024; Gu et al., 2023). Specifically, we disable fine-tuning of all embedding vectors and only fine-tune LoRA-related components of the text encoder and U-Net submodules.

Our evaluation encompasses two main aspects: quantitative assessment of the generated images based on objective metrics, as detailed in Section B.1, and qualitative demonstrations of the generated images to visually showcase the effectiveness of the optimizer. For qualitative evaluation, we use examples of Harry Potter and Hermione Granger to visually compare the performance of different optimizers in terms of identity preservation, scene conformity with prompt, and style diversity. As shown in Section B.2 and Section B.3, we conduct image generation under different LoRA scaling factors and compare the performance of various optimizers across multiple learning rates. This design not only evaluates the robustness of the optimizers under multi-scale hyperparameters but also reflects their overall impact on generation quality and consistency in real-world scenarios.

All experimental results consistently demonstrate the advantages of the SoLoRA algorithm. Both the quantitative evaluation metrics and the qualitative image demonstrations highlight the superior performance of SoLoRA. This success can be attributed to the ability of SoLoRA to effectively leverage the second-order information of the loss function, enabling more precise updates to model parameters. Furthermore, the low-rank approximations of gradients derived from our proposed adaptive weighted gradient strategy bring the performance of low-rank fine-tuning closer to that of full-parameter fine-tuning. This allows SoLoRA to achieve comparable performance to full-parameter fine-tuning while significantly reducing computational costs.

B.1 EVALUATION METRICS OF DIFFUSION MODELS

For quantitative evaluation, we employ two metrics: CLIP score (Hessel et al., 2021) and FID (Heusel et al., 2017). The CLIP score, based on the ViT-B/32 variant of the CLIP model (Radford et al., 2021), measures the consistency between the generated images and the input text prompts. The score ranges from 0 to 100, with higher scores indicating better alignment between the generated image and the text prompt. On the other hand, FID assesses the similarity between the distribution of generated images and the reference images. Lower FID values indicate higher similarity and better overall image quality. The experimental results are shown in Table 5.

In terms of FID, regardless of whether the SGD or AdamW optimizer is used, or whether the scaling factor is 0.7 or 1, our algorithm consistently achieve significantly lower FID values compared to all other methods. This strongly indicates that the distribution of images generated by our algorithm closely matches the distribution of the reference images. For the CLIP score, our algorithm achieve the best performance when the scaling factor is set to 1, outperforming all other methods. However, when the scaling factor is 0.7, the CLIP score of our algorithm is comparable to those of LoRA-Pro and AdamW algorithm. It is important to note that this does not imply that the quality of the images generated by our algorithm is inferior to others. On the contrary, this highlights one of the key strengths of our algorithm: by effectively leveraging second-order information from the loss function, the images generated by our method SoLoRA, not only maintain strong relevance to the text prompt but also exhibit richer details and greater diversity. For instance, in Figure 7, the clothing worn by the generated Harry Potter characters is more diverse, incorporating features that go beyond the simple text prompt. Similarly, in Figure 9, in addition to generating Harry Potter wearing a brown hat, our algorithm introduces more varied gestures for Harry Potter. These richer and more diverse features, while potentially causing a slight decrease in the CLIP score (as CLIP tends to prioritize strict prompt-image alignment and might not fully reward additional details beyond the prompt), actually enhance the overall quality and creativity of the generated images.

Table 5: CLIP and FID scores of different optimizers with different scaling factors for Mix-of-Show.

Methods	scaling=0.7		scaling=1	
	CLIP↑	FID↓	CLIP↑	FID↓
SGD	27.79	69.90	31.40	40.95
Scaled GD	31.23	35.86	30.60	29.62
LoRA-Pro SGD	31.47	34.30	30.48	29.19
SoLoRA SGD (ours)	31.47	30.17	31.58	28.18
AdamW	31.47	34.15	30.68	27.80
Scaled AdamW	24.21	48.23	24.51	34.18
LoRA-Pro AdamW	31.04	29.18	30.60	28.18
SoLoRA AdamW (ours)	31.47	29.01	30.73	27.13

B.2 EXPERIMENTAL RESULTS FOR DIFFERENT LORA SCALING FACTORS

To validate the effectiveness of the proposed optimizer, we compared the generated images of models trained using each optimizer under different LoRA scaling factors s . For the sake of fairness, we employed the optimal parameters for each optimizer, detailed in Table 6 for ease of replication.

Figure 5 and 6 show the generated results for Harry Potter and Hermione Granger when fine-tuning the model using different AdamW-based optimizers, with the scaling factor set to 1.0. Figure 7 and 8 show the model’s generated results when fine-tuned using different SGD-based optimizers, with scaling factors uniformly set to 1.0. Figure 9 and 10 present the generated results by using different AdamW-based optimizers, employing the scaling factor of 0.7. Experimental results demonstrate that the models trained with our optimizer generate high-quality images, accurately reproducing the identity of Harry Potter and Hermione Granger while demonstrating diverse scene layouts adhering to the input prompts.



Figure 5: Generated results based on the prompt “Harry Potter is walking near Mount Fuji” when fine-tuned using AdamW-based optimizers. All optimizers employed a LoRA scaling factor of 1.0, with the best learning rate. The results indicate that the output of the model trained with our optimizer incorporates the character “Harry Potter”, the action “walking”, and the scene “Mount Fuji”, yielding superior image quality compared to alternative approaches.

B.3 EXPERIMENTAL RESULTS FOR DIFFERENT LEARNING RATES

To illustrate the stability of the proposed optimizer, we fix the scaling factor to 1.0 and conduct experiments for each optimizer when using different learning rates. For AdamW-based optimizers, we set AdamW to employ the “Small LR” learning rate combination of $5e-6$ and $5e-5$ for text-encoder and U-Net, and the “Large LR” learning rate combination of $1e-5$ and $1e-4$. For Scaled AdamW, LoRA-Pro AdamW, and SoLoRA AdamW, we employed the same learning rate combinations, the “Small LR” of $5e-6$ and $5e-6$, and the “Large LR” combination of $1e-5$ and $1e-5$. For SGD-based optimizers, SGD, Scaled GD, and LoRA-Pro SGD, we employ the “Small LR” combination of $1e-2$ and $1e-2$, and the “Large LR” combination of $1e-1$ and $1e-1$, whereas SoLoRA SGD utilized the “Small LR” combination of $5e-6$ and $5e-6$, and the “Large LR” combination of $1e-5$ and $1e-5$.

The experimental results, presented in Figures 11 and 12, illustrate the effectiveness of our proposed optimizer across both small and large learning rates. This consistent performance signifies a higher degree of stability compared to the alternatives. Such stability is paramount when fine-tuning diffusion models, as their training is characterized by a non-stationary loss landscape. Therefore, the optimizer’s ability to remain effective under varying learning rates makes it a robust and advantageous choice for this application.



Figure 6: Generation results from the prompt “A photo of Hermione Granger on the beach, small waves, detailed symmetric face, beautiful composition” using AdamW-based optimizers. All the optimizers apply LoRA scaling factor as 1.0, with the best learning rate. Results demonstrate that the model trained with our optimizer generates higher-quality images than others, especially the face of Hermione Granger and the scene.



Figure 7: Generated results based on the prompt “Harry Potter standing near the lake” when fine-tuned using SGD-based optimizers. All optimizers employed a LoRA scaling factor of 1.0, with the best learning rate. Results demonstrate that the output images of the model trained with our optimizer have higher-quality than others, especially the face of Harry Potter.



Figure 8: Generated results based on the prompt "Hermione Granger wearing a brown shirt" when fine-tuned using SGD-based optimizers. All optimizers employed a LoRA scaling factor of 1.0, with the best learning rate. Results demonstrate that the model trained with SoLoRA generates higher-quality images than others, especially the face of Hermione Granger.



Figure 9: Generated results based on the prompt "Harry Potter wearing a brown hat" when fine-tuned using AdamW-based optimizers. All optimizers employed a LoRA scaling factor of 0.7, with the best learning rate. The results indicate that the output of the model trained with SoLoRA incorporates the character "Harry Potter", and the "hat", yielding superior image quality compared to alternative approaches.

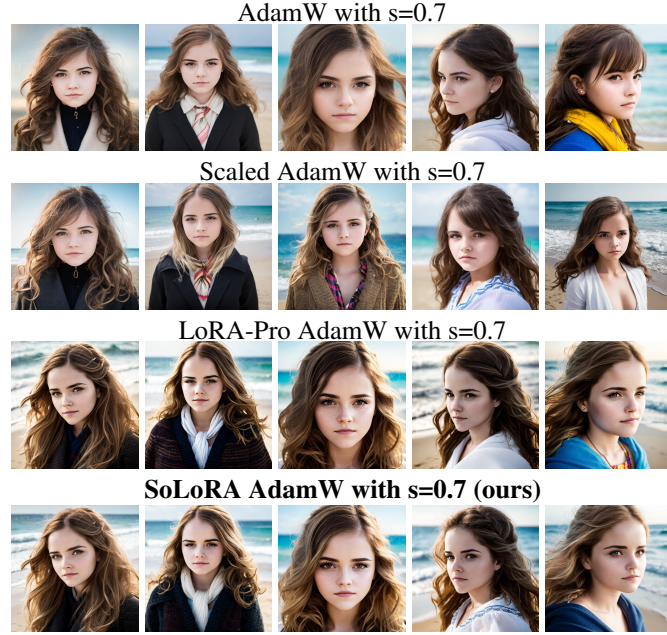


Figure 10: Generation results from the prompt “A photo of Hermione Granger on the beach, small waves, detailed symmetric face, beautiful composition” using AdamW-based optimizers. All the optimizers apply LoRA scaling factor as 0.7. According to the author’s recommendation, the optimizer AdamW and Scaled AdamW utilized a learning rate of $1e-5$ for text-encoder and $1e-4$ for U-Net, whereas LoRA-Pro AdamW and our SoLoRA optimizer adopted $1e-5$ for text-encoder and U-Net. Results demonstrate that SoLoRA generates higher-quality images for both scaling factors than others, including the face of Hermione Granger and the scene.

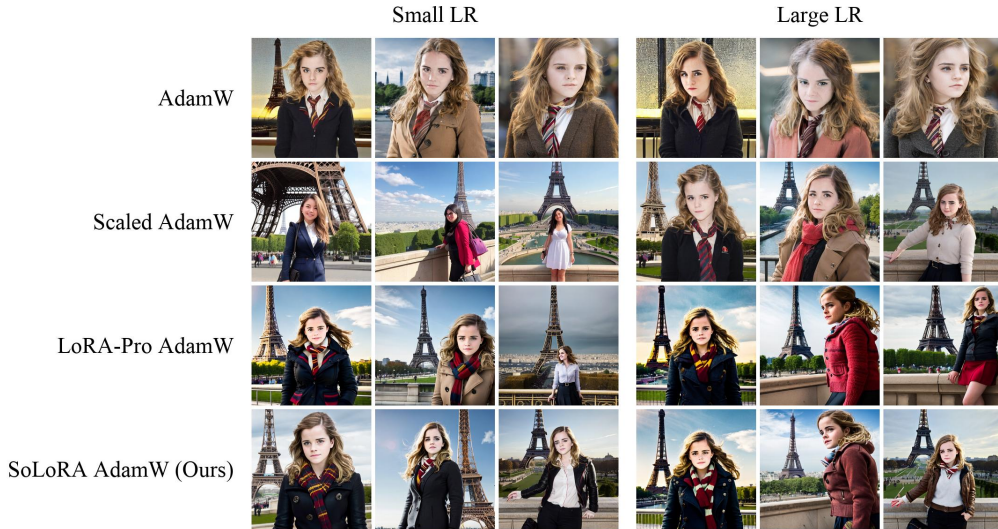


Figure 11: Generated results based on the prompt “Hermione Granger in front of Eiffel Tower” using AdamW-based optimizers. All the optimizers apply LoRA scaling factor as 1.0. “Small LR” and “Large LR” represents using different learning rate, please refer to Appendix B.3 for more details.



Figure 12: Generated results based on the prompt “ Photo of Harry Potter” using SGD-based optimizers. All the optimizers apply LoRA scaling factor as 1.0. “Small LR” and “Large LR” represent using different learning rate, please refer to Appendix B.3 for more details.

Table 6: Optimizer Parameters for fine-tuning the Mix-of-Show Model.

Methods	Learning Rate		β_3	$\beta_1 = \beta_2$
	Text-Encoder	U-Net		
SGD	1e-1	1e-1	/	/
Scaled GD	1e-1	1e-1	/	/
LoRA-Pro SGD	1e-1	1e-1	/	/
SoLoRA SGD	1e-5	1e-5	/	0.98
AdamW	1e-5	1e-4	0.9	0.999
Scaled AdamW	1e-5	1e-4	0.7	0.8
LoRA-Pro AdamW	1e-5	1e-5	0.9	0.999
SoLoRA AdamW	1e-5	1e-5	0.9	0.98

C COMPUTATIONAL AND MEMORY COMPLEXITY ANALYSIS OF SoLoRA

The update rule of SoLoRA is given by

$$\begin{aligned}\Delta_{A_t} &= (B_t^\top L_t^{\frac{1}{2}} B_t)^{-1} \underbrace{B_t^\top G_t R_t^{-\frac{1}{2}}}_{G_{A_t}} \left[I - \frac{1}{2} R_t^{\frac{1}{2}} A_t^\top (A_t R_t^{\frac{1}{2}} A_t^\top)^{-1} A_t \right], \\ \Delta_{B_t} &= \left[I - \frac{1}{2} B_t (B_t^\top L_t^{\frac{1}{2}} B_t)^{-1} B_t^\top L_t^{\frac{1}{2}} \right] L_t^{-\frac{1}{2}} \underbrace{G_t A_t^\top}_{G_{B_t}} (A_t R_t^{\frac{1}{2}} A_t^\top)^{-1}.\end{aligned}$$

We now analyze the computational complexity of computing the updates Δ_{A_t} and Δ_{B_t} . For simplicity, we focus on Δ_{A_t} , as the complexity for Δ_{B_t} is symmetric.

- Compute gradient G_t . The stochastic gradient G_t of W_t is obtained during the backpropagation process.
- Row and column sums for l_t and r_t . Compute l_t and r_t by summing the square of the element of G_t along rows or columns, which is in the computation $\mathcal{O}(mn)$. $L_t^{\frac{1}{2}}$ and $L_t^{-\frac{1}{2}}$ can be computed in $\mathcal{O}(m)$, $R_t^{\frac{1}{2}}$ and $R_t^{-\frac{1}{2}}$ can be computed in $\mathcal{O}(n)$.
- Compute $(B_t^\top L_t^{\frac{1}{2}} B_t)^{-1} G_{A_t} R_t^{-\frac{1}{2}}$. First to compute the inverse matrices $(B_t^\top L_t^{\frac{1}{2}} B_t)^{-1}$ in $\mathcal{O}((m+r)r^2)$. Then multiply the inverse $(B_t^\top L_t^{\frac{1}{2}} B_t)^{-1}$ by G_{A_t} in $\mathcal{O}(nr^2)$, and multiply the diagonal matrix $R_t^{-\frac{1}{2}}$ in $\mathcal{O}(nr)$.
- Compute $(B_t^\top L_t^{\frac{1}{2}} B_t)^{-1} G_{A_t} A_t^\top (A_t R_t^{\frac{1}{2}} A_t^\top)^{-1} A_t$. First to compute the inverse matrices $(A_t R_t^{\frac{1}{2}} A_t^\top)^{-1}$ in $\mathcal{O}((n+r)r^2)$. Use the result from the last step, multiply $(B_t^\top L_t^{\frac{1}{2}} B_t)^{-1} G_{A_t}$ by A_t^\top in computation $\mathcal{O}(nr^2)$, then multiply $(A_t R_t^{\frac{1}{2}} A_t^\top)^{-1}$ in computation $\mathcal{O}(r^3)$, and multiply A_t in $\mathcal{O}(nr^2)$.

The computation complexity of Δ_{A_t} is $\mathcal{O}(mn + (m+n)r^2 + r^3)$. The computation of Δ_{B_t} follows a similar structure, with symmetric terms. Its complexity is also $\mathcal{O}(mn + (m+n)r^2 + r^3)$. Then we have

- **Per Iteration Computational Complexity.** Combining the computations of Δ_{A_t} and Δ_{B_t} , the total computation complexity per iteration is $\mathcal{O}(mn + (m+n)r^2 + r^3)$.
- **Memory Complexity.** The algorithm requires storing the vectors l_t and r_t in each iteration, hence the memory complexity is $\mathcal{O}(m+n)$.

D PROOF OF THEORETICAL RESULTS

D.1 COMPUTATION OF JACOBIAN

Proposition D.1 (Computation of $J_{\mathcal{G}}$ and $J_{\mathcal{G}}^*$). *Let $[B, A]$ be a pair of low-rank factors with $B \in \mathbb{R}^{m \times r}$, $A \in \mathbb{R}^{r \times n}$. Define the generator $\mathcal{G} : [\mathbb{R}^{m \times r}, \mathbb{R}^{r \times n}] \rightarrow \mathbb{R}^{m \times n}$ by $\mathcal{G}([B, A]) = BA$. Denote the Jacobian of \mathcal{G} by $J_{\mathcal{G}}$ and its adjoint by $J_{\mathcal{G}}^*$. Then, for any $[P, Q] \in [\mathbb{R}^{m \times r}, \mathbb{R}^{r \times n}]$ and any $C \in \mathbb{R}^{m \times n}$,*

- $J_{\mathcal{G}}([B, A])[P, Q] = PA + BQ$,
- $J_{\mathcal{G}}^*([B, A])(C) = [CA^\top, B^\top C]$,
- $J_{\mathcal{G}}([B, A])J_{\mathcal{G}}^*([B, A])(C) = CA^\top A + BB^\top C$.

Proof. The Jacobian operator $J_{\mathcal{G}}([B, A])[P, Q] : [\mathbb{R}^{m \times r}, \mathbb{R}^{r \times n}] \rightarrow \mathbb{R}^{m \times n}$ represents the derivative of \mathcal{G} at $[B, A]$ along the direction $[P, Q]$. Similarly, $J_{\mathcal{G}}^*([B, A])(C) : \mathbb{R}^{m \times n} \rightarrow [\mathbb{R}^{m \times r}, \mathbb{R}^{r \times n}]$ is the adjoint of $J_{\mathcal{G}}$ at $[B, A]$ along the direction C . For more details, see (Absil et al., 2009, Section 6.1).

- (i) The computation of $J_{\mathcal{G}}$. Let $\mathbf{B}(t) : \mathbb{R} \rightarrow \mathbb{R}^{m \times r}$ and $\mathbf{A}(t) : \mathbb{R} \rightarrow \mathbb{R}^{r \times n}$ be differentiable curves with $\mathbf{B}(0) = \mathbf{B}$ and $\mathbf{A}(0) = \mathbf{A}$. By the chain rule, the Jacobian of \mathcal{G} at $[\mathbf{B}, \mathbf{A}]$ along these curves is

$$\begin{aligned} J_{\mathcal{G}}([\mathbf{B}(t), \mathbf{A}(t)])[\dot{\mathbf{B}}(t), \dot{\mathbf{A}}(t)] \Big|_{t=0} &= \left[\frac{d\mathcal{G}([\mathbf{B}, \mathbf{A}])}{d\mathbf{B}} \right] \dot{\mathbf{B}}(t) \Big|_{t=0} + \left[\frac{d\mathcal{G}([\mathbf{B}, \mathbf{A}])}{d\mathbf{A}} \right] \dot{\mathbf{A}}(t) \Big|_{t=0} \\ &= \dot{\mathbf{B}}(t)\mathbf{A}(t) \Big|_{t=0} + \mathbf{B}(t)\dot{\mathbf{A}}(t) \Big|_{t=0} \\ &= \dot{\mathbf{B}}(0)\mathbf{A} + \mathbf{B}\dot{\mathbf{A}}(0), \end{aligned}$$

where $\dot{\mathbf{B}}(t)$ and $\dot{\mathbf{A}}(t)$ denote the derivatives of $\mathbf{B}(t)$ and $\mathbf{A}(t)$ with respect to t . The second line follows because $\mathcal{G}([\mathbf{B}, \mathbf{A}]) = \mathbf{B}\mathbf{A}$, hence $\frac{d\mathcal{G}([\mathbf{B}, \mathbf{A}])}{d\mathbf{B}}$ and $\frac{d\mathcal{G}([\mathbf{B}, \mathbf{A}])}{d\mathbf{A}}$ are both linear operators.

Since $\dot{\mathbf{B}}(0)$ and $\dot{\mathbf{A}}(0)$ are arbitrary, for any $[\mathbf{P}, \mathbf{Q}] \in [\mathbb{R}^{m \times r}, \mathbb{R}^{r \times n}]$, we obtain

$$J_{\mathcal{G}}([\mathbf{B}, \mathbf{A}])[\mathbf{P}, \mathbf{Q}] = \mathbf{P}\mathbf{A} + \mathbf{B}\mathbf{Q}.$$

- (ii) The computation of $J_{\mathcal{G}}^*$. For brevity, write $J_{\mathcal{G}}[\mathbf{P}, \mathbf{Q}]$ for $J_{\mathcal{G}}([\mathbf{B}, \mathbf{A}])[\mathbf{P}, \mathbf{Q}]$ and $J_{\mathcal{G}}^*(\mathbf{C})$ for $J_{\mathcal{G}}^*([\mathbf{B}, \mathbf{A}])(\mathbf{C})$. By definition of the adjoint (with respect to the Frobenius inner product), for any $[\mathbf{P}, \mathbf{Q}] \in (\mathbb{R}^{m \times r}, \mathbb{R}^{r \times n})$ and $\mathbf{C} \in \mathbb{R}^{m \times n}$,

$$\langle J_{\mathcal{G}}[\mathbf{P}, \mathbf{Q}], \mathbf{C} \rangle = \langle [\mathbf{P}, \mathbf{Q}], J_{\mathcal{G}}^*(\mathbf{C}) \rangle.$$

For the left-hand side,

$$\begin{aligned} \langle J_{\mathcal{G}}[\mathbf{P}, \mathbf{Q}], \mathbf{C} \rangle &= \langle \mathbf{P}\mathbf{A} + \mathbf{B}\mathbf{Q}, \mathbf{C} \rangle \\ &= \langle \mathbf{P}\mathbf{A}, \mathbf{C} \rangle + \langle \mathbf{B}\mathbf{Q}, \mathbf{C} \rangle \\ &= \langle \mathbf{P}, \mathbf{C}\mathbf{A}^\top \rangle + \langle \mathbf{Q}, \mathbf{B}^\top \mathbf{C} \rangle. \end{aligned}$$

For the right-hand side, writing $J_{\mathcal{G}}^*(\mathbf{C}) = [\mathbf{C}_1, \mathbf{C}_2]$, then

$$\begin{aligned} \langle [\mathbf{P}, \mathbf{Q}], J_{\mathcal{G}}^*(\mathbf{C}) \rangle &= \langle [\mathbf{P}, \mathbf{Q}], [\mathbf{C}_1, \mathbf{C}_2] \rangle \\ &= \langle \mathbf{P}, \mathbf{C}_1 \rangle + \langle \mathbf{Q}, \mathbf{C}_2 \rangle. \end{aligned}$$

Hence $\mathbf{C}_1 = \mathbf{C}\mathbf{A}^\top$ and $\mathbf{C}_2 = \mathbf{B}^\top \mathbf{C}$, and therefore $J_{\mathcal{G}}^*([\mathbf{B}, \mathbf{A}])(\mathbf{C}) = [\mathbf{C}\mathbf{A}^\top, \mathbf{B}^\top \mathbf{C}]$.

- (iii) Finally, $J_{\mathcal{G}}([\mathbf{B}, \mathbf{A}])J_{\mathcal{G}}^*([\mathbf{B}, \mathbf{A}])(\mathbf{C}) = J_{\mathcal{G}}([\mathbf{B}, \mathbf{A}])[\mathbf{C}\mathbf{A}^\top, \mathbf{B}^\top \mathbf{C}] = \mathbf{C}\mathbf{A}^\top \mathbf{A} + \mathbf{B}\mathbf{B}^\top \mathbf{C}$ as claimed. \square

D.2 ORTHOGONAL PROJECTION TO TANGENT SPACE

In this subsection, we derive the orthogonal projection onto the tangent space under both the standard metric and the weighted metric. The specific forms of \mathbf{L}_t and \mathbf{R}_t are presented here and will not be repeated in subsequent propositions and proofs. For the sake of simplicity, the subscript t will be omitted in this subsection.

$$\begin{aligned} \mathbf{L}_t &= \text{diag}(\mathbf{l}_t / \sqrt{\|\mathbf{l}_t\|_1}) \quad \text{with} \quad \mathbf{l}_t = \beta_2 \mathbf{l}_{t-1} + (1 - \beta_2) \sum_{j=1}^n (\mathbf{G}_t \odot \mathbf{G}_t)_{i,j}, \\ \mathbf{R}_t &= \text{diag}(\mathbf{r}_t / \sqrt{\|\mathbf{r}_t\|_1}) \quad \text{with} \quad \mathbf{r}_t = \beta_3 \mathbf{r}_{t-1} + (1 - \beta_3) \sum_{i=1}^m (\mathbf{G}_t \odot \mathbf{G}_t)_{i,j}, \end{aligned} \tag{13}$$

where \odot denotes the Hadamard (elementwise) product and $\mathbf{G}_t = \nabla \mathcal{L}(\mathbf{W}_0 + \mathbf{W}_t)$.

Proposition D.2 (Orthogonal Projection to Tangent Space Under the Standard Metric). *Let $\mathbf{W} \in \mathcal{M}_r$ be a rank- r matrix with a low-rank decomposition $\mathbf{W} = \mathbf{B}\mathbf{A}$, where $\mathbf{B} \in \mathbb{R}^{m \times r}$, $\mathbf{A} \in \mathbb{R}^{r \times n}$. Denote by $\mathbb{T}_{\mathbf{W}}$ the tangent space of the smooth manifold \mathcal{M}_r at the point \mathbf{W} . Then, the orthogonal projection of any matrix $\mathbf{Z} \in \mathbb{R}^{m \times n}$ onto $\mathbb{T}_{\mathbf{W}}$ is given by*

$$\mathcal{P}_{\mathbb{T}_{\mathbf{W}}}(\mathbf{Z}) = \mathbf{B}(\mathbf{B}^\top \mathbf{B})^{-1} \mathbf{B}^\top \mathbf{Z} + \mathbf{Z} \mathbf{A}^\top (\mathbf{A} \mathbf{A}^\top)^{-1} \mathbf{A} - \mathbf{B}(\mathbf{B}^\top \mathbf{B})^{-1} \mathbf{B}^\top \mathbf{Z} \mathbf{A}^\top (\mathbf{A} \mathbf{A}^\top)^{-1} \mathbf{A}.$$

Proof. Suppose \mathbf{W} has a compact singular value decomposition, given by $\mathbf{W} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$, where $\mathbf{U} \in \mathbb{R}^{m \times r}$, $\mathbf{\Sigma} \in \mathbb{R}^{r \times r}$, $\mathbf{V} \in \mathbb{R}^{n \times r}$. Then the tangent space $\mathbb{T}_{\mathbf{W}}$ at \mathbf{W} is characterized as

$$\mathbb{T}_{\mathbf{W}} = \{\mathbf{U}\mathbf{M}^\top + \mathbf{N}\mathbf{V}^\top, \text{ for } \mathbf{M} \in \mathbb{R}^{m \times r}, \mathbf{N} \in \mathbb{R}^{n \times r}\}.$$

Therefore, the orthogonal projection of \mathbf{Z} onto $\mathbb{T}_{\mathbf{W}}$ is known to be (Wei et al., 2016)

$$\mathcal{P}_{\mathbb{T}_{\mathbf{W}}}(\mathbf{Z}) = \mathbf{U}\mathbf{U}^\top \mathbf{Z} + \mathbf{Z}\mathbf{V}^\top \mathbf{V} - \mathbf{U}\mathbf{U}^\top \mathbf{Z}\mathbf{V}^\top \mathbf{V}. \quad (14)$$

Since the columns of \mathbf{B} and \mathbf{U} span the same column space (i.e., the column space of \mathbf{W}), then there exists an invertible matrix $\mathbf{S} \in \mathbb{R}^{r \times r}$ such that $\mathbf{B} = \mathbf{U}\mathbf{S}$ and $\mathbf{U} = \mathbf{B}\mathbf{S}^{-1}$. Using this relation, we have

$$\mathbf{U}^\top \mathbf{U} = (\mathbf{B}\mathbf{S}^{-1})^\top \mathbf{B}\mathbf{S}^{-1} = \mathbf{S}^{-\top} (\mathbf{B}^\top \mathbf{B}) \mathbf{S}^{-1}.$$

Since $\mathbf{U}^\top \mathbf{U} = \mathbf{I}_r$, it follows that

$$\mathbf{S}^{-\top} (\mathbf{B}^\top \mathbf{B}) \mathbf{S}^{-1} = \mathbf{I}_r \implies \mathbf{B}^\top \mathbf{B} = \mathbf{S}^\top \mathbf{S}.$$

Using this, we compute $\mathbf{U}\mathbf{U}^\top$

$$\mathbf{U}\mathbf{U}^\top = \mathbf{B}\mathbf{S}^{-1}\mathbf{S}^{-\top} \mathbf{B} = \mathbf{B}(\mathbf{S}^\top \mathbf{S})^{-1} \mathbf{B}^\top = \mathbf{B}(\mathbf{B}^\top \mathbf{B})^{-1} \mathbf{B}^\top \quad (15)$$

Similarly, since the rows of \mathbf{A} and the columns of \mathbf{V} span the same row space (i.e., the row space of \mathbf{W}), there exists an invertible matrix $\mathbf{Q} \in \mathbb{R}^{r \times r}$ such that $\mathbf{A} = \mathbf{Q}\mathbf{V}^\top$ and $\mathbf{V}^\top = \mathbf{Q}^{-1}\mathbf{A}$. Further, using $\mathbf{V}^\top \mathbf{V} = \mathbf{I}_r$, we obtain

$$\mathbf{V}^\top \mathbf{V} = \mathbf{Q}^{-1}(\mathbf{A}\mathbf{A}^\top)\mathbf{Q}^{-\top} = \mathbf{I}_r,$$

hence $\mathbf{A}\mathbf{A}^\top = \mathbf{Q}\mathbf{Q}^\top$ and

$$\mathbf{V}\mathbf{V}^\top = \mathbf{A}^\top \mathbf{Q}^{-\top} \mathbf{Q}^{-1} \mathbf{A} = \mathbf{A}^\top (\mathbf{Q}\mathbf{Q}^\top)^{-1} \mathbf{A} = \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{A} \quad (16)$$

Substituting (15) and (16) into (14) yields

$$\mathcal{P}_{\mathbb{T}_{\mathbf{W}}}(\mathbf{Z}) = \mathbf{B}(\mathbf{B}^\top \mathbf{B})^{-1} \mathbf{B}^\top \mathbf{Z} + \mathbf{Z}\mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{A} - \mathbf{B}(\mathbf{B}^\top \mathbf{B})^{-1} \mathbf{B}^\top \mathbf{Z}\mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1} \mathbf{A}.$$

□

Proposition D.3 (Orthogonal Projection onto the Tangent Space Under the Weighted Metric). *Let $\mathbf{W} \in \mathcal{M}_r$ has a low-rank decomposition $\mathbf{W} = \mathbf{B}\mathbf{A}$, where $\mathbf{B} \in \mathbb{R}^{m \times r}$, $\mathbf{A} \in \mathbb{R}^{r \times n}$. Denote the tangent space of the Riemannian manifold \mathcal{M}_r at the point \mathbf{W} as $\mathbb{T}_{\mathbf{W}}$. The weighted metric is defined as $\langle \mathbf{Y}, \mathbf{Z} \rangle_{\mathbf{H}} = \langle \mathbf{L}^{\frac{1}{2}} \mathbf{Y} \mathbf{R}^{\frac{1}{2}}, \mathbf{Z} \rangle$ for any $\mathbf{Y}, \mathbf{Z} \in \mathbb{R}^{m \times n}$. Then, the orthogonal projection of any matrix $\mathbf{Z} \in \mathbb{R}^{m \times n}$ onto $\mathbb{T}_{\mathbf{W}}$ under the weighed metric is given by*

$$\begin{aligned} \mathcal{P}_{\mathbb{T}_{\mathbf{W}}}(\mathbf{Z}) &= \mathbf{B}(\mathbf{B}^\top \mathbf{L}^{\frac{1}{2}} \mathbf{B})^{-1} \mathbf{B}^\top \mathbf{L}^{\frac{1}{2}} \mathbf{Z} + \mathbf{Z} \mathbf{R}^{\frac{1}{2}} \mathbf{A}^\top (\mathbf{A} \mathbf{R}^{\frac{1}{2}} \mathbf{A}^\top)^{-1} \mathbf{A} \\ &\quad - \mathbf{B}(\mathbf{B}^\top \mathbf{B})^{-1} \mathbf{B}^\top \mathbf{L}^{\frac{1}{2}} \mathbf{Z} \mathbf{R}^{\frac{1}{2}} \mathbf{A}^\top (\mathbf{A} \mathbf{A}^\top)^{-1} \mathbf{A}. \end{aligned}$$

Proof. This proof is inspired by (Bian et al., 2024). Here, we briefly provide a sketch of the proof.

- (i) *The new orthonormal basis under the weighted metric.* Let $\mathbf{W} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$ be a compact SVD with $\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r] \in \mathbb{R}^{m \times r}$, $\mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r] \in \mathbb{R}^{n \times r}$. Normalize the singular vectors under the weighted vector

$$\langle \mathbf{x}, \mathbf{y} \rangle_{\mathbf{L}^{\frac{1}{2}}} = \langle \mathbf{L}^{\frac{1}{2}} \mathbf{x}, \mathbf{y} \rangle \text{ in } \mathbb{R}^m \quad \text{and} \quad \langle \mathbf{x}, \mathbf{y} \rangle_{\mathbf{R}^{\frac{1}{2}}} = \langle \mathbf{R}^{\frac{1}{2}} \mathbf{x}, \mathbf{y} \rangle \text{ in } \mathbb{R}^n$$

to obtain

$$\begin{aligned} \tilde{\mathbf{U}} &= \mathbf{U}(\mathbf{U}^\top \mathbf{L}^{\frac{1}{2}} \mathbf{U})^{-\frac{1}{2}} := [\tilde{\mathbf{u}}_1, \tilde{\mathbf{u}}_2, \dots, \tilde{\mathbf{u}}_r] \in \mathbb{R}^{m \times r}, \\ \tilde{\mathbf{V}} &= \mathbf{V}(\mathbf{V}^\top \mathbf{R}^{\frac{1}{2}} \mathbf{V})^{-\frac{1}{2}} := [\tilde{\mathbf{v}}_1, \tilde{\mathbf{v}}_2, \dots, \tilde{\mathbf{v}}_r] \in \mathbb{R}^{n \times r}. \end{aligned}$$

Next, we extend $\tilde{\mathbf{U}}$ and $\tilde{\mathbf{V}}$ to full orthonormal basis of $(\mathbb{R}^m, \langle \cdot, \cdot \rangle_{\mathbf{L}^{\frac{1}{2}}})$ and $(\mathbb{R}^n, \langle \cdot, \cdot \rangle_{\mathbf{R}^{\frac{1}{2}}})$ respectively. Then, an orthonormal basis of $\mathbb{T}_{\mathbf{W}}$ with respect to $\langle \cdot, \cdot \rangle_{\mathbf{H}_t}$ is $\{\tilde{\mathbf{u}}_i \tilde{\mathbf{v}}_j^\top\}_{\min\{i,j\} \leq r}$.

(ii) *Orthogonal projection represented by the new orthonormal basis.* Using the orthonormal bases \tilde{U} and \tilde{V} , the projection of Z onto \mathbb{T}_W is expressed as:

$$\begin{aligned}\tilde{\mathcal{P}}_{\mathbb{T}_W}(Z) &= \sum_{(i,j):\min\{i,j\}\leq r} \langle Z, \tilde{u}_i \tilde{v}_j^\top \rangle_{H_t} \cdot \tilde{u}_i \tilde{v}_j^\top = \sum_{(i,j):\min\{i,j\}\leq r} \langle L^{\frac{1}{2}} Z R^{\frac{1}{2}}, \tilde{u}_i \tilde{v}_j^\top \rangle \cdot \tilde{u}_i \tilde{v}_j^\top \\ &= \sum_{(i,j):\min\{i,j\}\leq r} \tilde{u}_i^\top L^{\frac{1}{2}} Z R^{\frac{1}{2}} \tilde{v}_j \cdot \tilde{u}_i \tilde{v}_j^\top \\ &= \tilde{U} \tilde{U}^\top L^{\frac{1}{2}} Z + Z R^{\frac{1}{2}} \tilde{V} \tilde{V}^\top - \tilde{U} \tilde{U}^\top L^{\frac{1}{2}} Z R^{\frac{1}{2}} \tilde{V} \tilde{V}^\top.\end{aligned}$$

(iii) *Express the basis projectors via factors B and A .* Since B and A span the same spaces as U and V , we derive

$$\tilde{U} \tilde{U}^\top = B \left(B^\top L^{\frac{1}{2}} B \right)^{-1} B^\top, \quad \tilde{V} \tilde{V}^\top = A^\top \left(A R^{\frac{1}{2}} A^\top \right)^{-1} A.$$

Substituting these expressions into the formula for $\tilde{\mathcal{P}}_{\mathbb{T}_W}$, we obtain

$$\begin{aligned}\tilde{\mathcal{P}}_{\mathbb{T}_W}(Z) &= B (B^\top L^{\frac{1}{2}} B)^{-1} B^\top L^{\frac{1}{2}} Z + Z R^{\frac{1}{2}} A^\top (A R^{\frac{1}{2}} A^\top)^{-1} A \\ &\quad - B (B^\top L^{\frac{1}{2}} B)^{-1} B^\top L^{\frac{1}{2}} Z R^{\frac{1}{2}} A^\top (A R^{\frac{1}{2}} A^\top)^{-1} A.\end{aligned}$$

□

Proposition D.4. Suppose $W \in \mathcal{M}_r$ has a low-rank decomposition $W = BA$, where $B \in \mathbb{R}^{m \times r}$ and $A \in \mathbb{R}^{r \times n}$. For any matrix $M \in \mathbb{R}^{m \times r}$, $N \in \mathbb{R}^{r \times n}$, the matrix $MA + BN$ lies in the tangent space \mathbb{T}_W at W of \mathcal{M}_r at the point W .

Proof. Let $W \in \mathcal{M}_r$ has a compact singular value decomposition $W = U \Sigma V^\top$, where $U \in \mathbb{R}^{m \times r}$, $\Sigma \in \mathbb{R}^{r \times r}$, and $V \in \mathbb{R}^{n \times r}$. By definition, the tangent space \mathbb{T}_W at W is given by

$$\mathbb{T}_W = \{U K_1^\top + K_2 V^\top \mid K_1 \in \mathbb{R}^{n \times r}, K_2 \in \mathbb{R}^{m \times r}\}.$$

Since B and A are low-rank factors of W , there exist invertible matrices $S \in \mathbb{R}^{r \times r}$ and $Q \in \mathbb{R}^{r \times r}$ such that

$$B = US, \quad A = QV^\top.$$

Substituting these expressions, the matrix $MA + BN$ can be rewritten as

$$MA + BN = MQV^\top + USN.$$

The first term, MQV^\top , lies in $\text{span}(V^\top)$, and the second term, USN , lies in $\text{span}(U)$. Thus, the sum $MQV^\top + USN$ lies in the tangent space \mathbb{T}_W by the definition of the tangent space. Then, it follows that $MA + BN$ is on the tangent space \mathbb{T}_W . This completes the proof. □

D.3 PROOFS OF THEOREM 3.1 AND THEOREM 3.2

Proof of Theorem 3.1. Define

$$\Gamma(\Delta_{B_t}, \Delta_{A_t}) := \frac{1}{2} \|\Delta_{B_t} A_t + B_t \Delta_{A_t} - \tilde{\mathcal{P}}_{\mathbb{T}_t}(L_t^{-\frac{1}{2}} G_t R_t^{-\frac{1}{2}})\|_{H_t}^2.$$

Differentiating $\Gamma(\Delta_{B_t}, \Delta_{A_t})$ with respect to Δ_{B_t} and Δ_{A_t} yields

$$\nabla_{\Delta_{B_t}} \Gamma(\Delta_{B_t}, \Delta_{A_t}) = L_t^{\frac{1}{2}} \Delta_{B_t} (A_t R_t^{\frac{1}{2}} A_t^\top) + L_t^{\frac{1}{2}} B_t \Delta_{A_t} R_t^{\frac{1}{2}} A_t^\top - G_t A_t^\top, \quad (17)$$

and

$$\nabla_{\Delta_{A_t}} \Gamma(\Delta_{B_t}, \Delta_{A_t}) = B_t^\top L_t^{\frac{1}{2}} \Delta_{B_t} A_t R_t^{\frac{1}{2}} + B_t^\top L_t^{\frac{1}{2}} B_t \Delta_{A_t} R_t^{\frac{1}{2}} - B_t^\top G_t. \quad (18)$$

Setting $\nabla_{\Delta_{B_t}} \Gamma(\Delta_{B_t}, \Delta_{A_t}) = 0$ and using the invertibility of $(A_t R_t^{\frac{1}{2}} A_t^\top)$ and $L_t^{\frac{1}{2}}$ gives

$$\Delta_{B_t} = L_t^{-\frac{1}{2}} G_t A_t^\top (A_t R_t^{\frac{1}{2}} A_t^\top)^{-1} - B_t \Delta_{A_t} R_t^{\frac{1}{2}} A_t^\top (A_t R_t^{\frac{1}{2}} A_t^\top)^{-1}. \quad (19)$$

Substituting (18) into $\nabla_{\Delta_{A_t}} \Gamma(\Delta_{B_t}, \Delta_{A_t}) = \mathbf{0}$ and using the invertibility of $B_t^\top L_t^{\frac{1}{2}} B_t$ and $R_t^{\frac{1}{2}}$ yields

$$\Delta_{A_t}[I - \tilde{Q}_{A_t}] = (B_t^\top L_t^{\frac{1}{2}} B_t)^{-1} B_t^\top G_t R_t^{-\frac{1}{2}} [I - \tilde{Q}_{A_t}],$$

where $\tilde{Q}_{A_t} = R_t^{\frac{1}{2}} A_t^\top (A_t R_t^{\frac{1}{2}} A_t^\top)^{-1} A_t$, which is the projection matrix onto the row space of A_t . Since $I - \tilde{Q}_{A_t}$ is the residual maker matrix, then a general solution is

$$\Delta_{A_t}^{\text{opt}} = (B_t^\top L_t^{\frac{1}{2}} B_t)^{-1} B_t^\top G_t R_t^{-\frac{1}{2}} + X_t A_t,$$

with arbitrary matrix $X_t \in \mathbb{R}^{r \times r}$. Plugging this Δ_{A_t} back into (19) gives

$$\Delta_{B_t}^{\text{opt}} = [I - \tilde{P}_{B_t}] L_t^{-\frac{1}{2}} G_t A_t^\top (A_t R_t^{\frac{1}{2}} A_t^\top)^{-1} - B_t X_t,$$

where $\tilde{P}_{B_t} = B_t (B_t^\top L_t^{\frac{1}{2}} B_t)^{-1} B_t^\top L_t^{\frac{1}{2}}$, which is the projection matrix onto the column space of B_t . \square

Proof of Theorem 3.2. Let the objective function be $\Psi(X_t) = \frac{1}{2} \|\Delta_{B_t} A_t - B_t \Delta_{A_t}\|_{H_t}^2$. To minimize $\Psi(X_t)$, we compute its gradient with respect to X_t ,

$$\nabla_{X_t} \Psi(X_t) = B_t^\top L_t^{\frac{1}{2}} (\Delta_{B_t} A_t - B_t \Delta_{A_t}) R_t^{\frac{1}{2}} A_t^\top.$$

Substituting the expressions for A_t and B_t from Theorem 3.1, we have

$$\begin{aligned} \nabla_{X_t} \Psi(X_t) &= B_t^\top L_t^{\frac{1}{2}} \left([I - B_t (B_t^\top L_t^{\frac{1}{2}} B_t)^{-1} B_t^\top L_t^{\frac{1}{2}}] L_t^{-\frac{1}{2}} G_t A_t^\top (A_t R_t^{\frac{1}{2}} A_t^\top)^{-1} A_t \right. \\ &\quad \left. - B_t (B_t^\top L_t^{\frac{1}{2}} B_t)^{-1} B_t^\top G_t R_t^{-\frac{1}{2}} - 2B_t X_t A_t \right) R_t^{\frac{1}{2}} A_t^\top \\ &= -B_t^\top G_t A_t - 2(B_t^\top L_t^{\frac{1}{2}} B_t) X_t (A_t R_t^{\frac{1}{2}} A_t^\top). \end{aligned}$$

Setting $\nabla_{X_t} \Psi(X_t) = \mathbf{0}$, we obtain

$$-B_t^\top G_t A_t = 2(B_t^\top L_t^{\frac{1}{2}} B_t) X_t (A_t R_t^{\frac{1}{2}} A_t^\top).$$

Since $B_t^\top L_t^{\frac{1}{2}} B_t$ and $A_t R_t^{\frac{1}{2}} A_t^\top$ are invertible, we solve for X_t as

$$X_t^{\text{opt}} = -\frac{1}{2} (B_t^\top L_t^{\frac{1}{2}} B_t)^{-1} B_t^\top G_t A_t^\top (A_t R_t^{\frac{1}{2}} A_t^\top)^{-1}.$$

Thus, the optimal solution for X_t is derived. \square

D.4 PROOF OF THE LOSS FUNCTION EXHIBITS A DECREASING TREND

Under the framework of LoRA, the infinitesimal change in loss is

$$\begin{aligned} d\mathcal{L} &= \langle G_B, dB \rangle_H + \langle G_A, dA \rangle_H \\ &= -\eta (\langle G_B, \Delta_B \rangle_H + \langle G_A, \Delta_A \rangle_H) \\ &= -\eta \left(\langle G_B, [I - B(B^\top L^{\frac{1}{2}} B)^{-1} B^\top L^{\frac{1}{2}}] L^{-\frac{1}{2}} G_B (A R^{\frac{1}{2}} A^\top)^{-1} - B X \rangle_H \right. \\ &\quad \left. + \langle G_A, (B^\top L^{\frac{1}{2}} B)^{-1} G_A R^{-\frac{1}{2}} + X A \rangle_H \right) \\ &= -\eta \left(\langle G_B, [I - B(B^\top L^{\frac{1}{2}} B)^{-1} B^\top L^{\frac{1}{2}}] L^{-\frac{1}{2}} G_B (A R^{\frac{1}{2}} A^\top)^{-1} \rangle_H + \langle G_A, (B^\top L^{\frac{1}{2}} B)^{-1} G_A R^{-\frac{1}{2}} \rangle_H \right. \\ &\quad \left. + \langle G_B, -B X \rangle_H + \langle G_A, X A \rangle_H \right) \\ &= -\eta \left(\langle G_B, [I - B(B^\top L^{\frac{1}{2}} B)^{-1} B^\top L^{\frac{1}{2}}] L^{-\frac{1}{2}} G_B (A R^{\frac{1}{2}} A^\top)^{-1} \rangle_H + \langle G_A, (B^\top L^{\frac{1}{2}} B)^{-1} G_A R^{-\frac{1}{2}} \rangle_H \right), \end{aligned} \tag{20}$$

Where the third line is derived from Theorem 3.1, the forth line is because of the additivity of the inner product, and the last line is because

$$\begin{aligned}\langle G_B, -BX \rangle_H + \langle G_A, XA \rangle_H &= \langle -B^\top G_B + G_A A^\top, X \rangle_H \\ &= \langle -B^\top G_A^\top + B^\top G_A^\top, X \rangle_H \\ &= 0.\end{aligned}$$

Therefore, to prove $d\mathcal{L} \leq 0$, it suffices to prove that

$$\begin{aligned}\langle G_B, [I - B(B^\top L^{\frac{1}{2}} B)^{-1} B^\top L^{\frac{1}{2}} L^{-\frac{1}{2}} G_B (AR^{\frac{1}{2}} A^\top)^{-1}] \rangle_H &\geq 0, \\ \langle G_A, (B^\top L^{\frac{1}{2}} B)^{-1} G_A R^{-\frac{1}{2}} \rangle_H &\geq 0.\end{aligned}$$

1. First to prove $\langle G_B, [I - B(B^\top L^{\frac{1}{2}} B)^{-1} B^\top L^{\frac{1}{2}}] L^{-\frac{1}{2}} G_B (AR^{\frac{1}{2}} A^\top)^{-1} \rangle_H \geq 0$.

- Prove $[I - B(B^\top L^{\frac{1}{2}} B)^{-1} B^\top L^{\frac{1}{2}}]$ is symmetric and positive semi-definite.

From Equation (10), we have $\tilde{P}_B = B(B^\top L^{\frac{1}{2}} B)^{-1} B^\top L^{\frac{1}{2}}$. \tilde{P}_B is symmetric in the weighted space if it is self-adjoint with respect to the weighted inner product. That means we need to prove for all vectors $x, y \in \mathbb{R}^m$

$$\langle \tilde{P}_B x, y \rangle_H = \langle x, \tilde{P}_B y \rangle_H$$

For the left-hand side,

$$\langle \tilde{P}_B x, y \rangle_H = \langle L^{\frac{1}{2}} \tilde{P}_B x, y \rangle = \langle L^{\frac{1}{2}} B (B^\top L^{\frac{1}{2}} B)^{-1} B^\top L^{\frac{1}{2}} x, y \rangle.$$

For the right-hand side,

$$\langle x, \tilde{P}_B y \rangle_H = \langle L^{\frac{1}{2}} x, \tilde{P}_B y \rangle = \langle \tilde{P}_B^\top L^{\frac{1}{2}} x, y \rangle = \langle L^{\frac{1}{2}} B (B^\top L^{\frac{1}{2}} B)^{-1} B^\top L^{\frac{1}{2}} x, y \rangle.$$

Thus, the left-hand side equals to the right-hand side, \tilde{P}_B is symmetric.

By Proposition D.3, \tilde{P}_B is the orthogonal projection under the inner product $\langle \cdot, \cdot \rangle_H$. As a projection matrix, the eigenvalues of \tilde{P}_B are zeros and ones. Thus, \tilde{P}_B^\top and $I - \tilde{P}_B^\top$ are both positive semi-definite matrices.

There exists a Cholesky decomposition of $I - \tilde{P}_B^\top$, denote as $I - \tilde{P}_B^\top = K K^\top$, where $K \in \mathbb{R}^{m \times m}$ is a lower triangular matrix.

- Prove $(AR^{\frac{1}{2}} A^\top)^{-1}$ is symmetric and positive definite.

Since R is a diagonal matrix, then symmetric, and therefore $AR^{\frac{1}{2}} A^\top$ is also symmetric.

For any non-zero vector $x \in \mathbb{R}^r$, $x^\top AR^{\frac{1}{2}} A^\top x = \langle A^\top x, A^\top x \rangle_H = \|A^\top x\|_H^2 > 0$.

Thus, $AR^{\frac{1}{2}} A^\top$ is a positive definite matrix.

Therefore, $AR^{\frac{1}{2}} A^\top$ is invertible and $(AR^{\frac{1}{2}} A^\top)^{-1}$ is symmetric and positive definite.

There exists a Cholesky decomposition of $(AR^{\frac{1}{2}} A^\top)^{-1}$, denote as $(AR^{\frac{1}{2}} A^\top)^{-1} = C C^\top$, where $C \in \mathbb{R}^{r \times r}$ is a lower triangular matrix.

Since $[I - B(B^\top L^{\frac{1}{2}} B)^{-1} B^\top L^{\frac{1}{2}}]$ and $(AR^{\frac{1}{2}} A^\top)^{-1}$ both have a Cholesky decomposition. The inner product can be rewritten as

$$\begin{aligned}\langle G_B, [I - B(B^\top L^{\frac{1}{2}} B)^{-1} B^\top L^{\frac{1}{2}}] L^{-\frac{1}{2}} G_B (AR^{\frac{1}{2}} A^\top)^{-1} \rangle_H &= \langle G_B, K K^\top L^{-\frac{1}{2}} G_B C C^\top \rangle_H \\ &= \langle L^{-\frac{1}{2}} K K^\top G_B, G_B C C^\top \rangle_H \\ &= \langle K K^\top G_B, G_B C C^\top \rangle \\ &= \langle K^\top G_B C, K^\top G_B C \rangle \\ &= \|K^\top G_B C\|_F^2 \\ &\geq 0\end{aligned}$$

2. Then prove $\langle G_A, (B^\top L^{\frac{1}{2}} B)^{-1} G_A R^{-\frac{1}{2}} \rangle_H \geq 0$.

Similar to $(A R^{\frac{1}{2}} A^\top)^{-1}$ is symmetric and positive definite, $(B^\top L^{\frac{1}{2}} B)^{-1}$ is symmetric positive definite, and there exists a Cholesky decomposition of $(B^\top L^{\frac{1}{2}} B)^{-1}$, denote as $(B^\top L^{\frac{1}{2}} B)^{-1} = D D^\top$, where $D \in \mathbb{R}^{r \times r}$ is a lower triangular matrix.

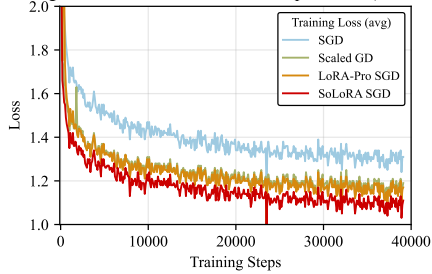
Thus, the inner product can be rewritten as

$$\begin{aligned} & \langle G_A, (B^\top L^{\frac{1}{2}} B)^{-1} G_A R^{-\frac{1}{2}} \rangle_H \\ &= \langle G_A, D D^\top G_A R^{-\frac{1}{2}} \rangle_H \\ &= \langle G_A R^{\frac{1}{2}}, D D^\top G_A R^{-\frac{1}{2}} \rangle \\ &= \langle G_A, D D^\top G_A \rangle \\ &= \langle D^\top G_A, D^\top G_A \rangle \\ &= \|D^\top G_A\|_F^2 \\ &\geq 0 \end{aligned}$$

In the conclusion, we have completed the proof of $d\mathcal{L} \leq 0$. This ensures that the SoLoRA model maintains a decreasing trend in the loss.

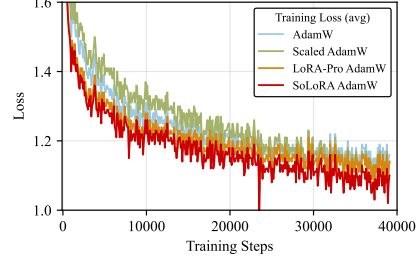
E ADDITIONAL EXPERIMENTS

Training Loss Curves for SGD Based Optimizers (Dart Dataset)



(a) Training loss curve over training steps when fine-tuning using SGD-based method.

Training Loss Curves for AdamW Based Optimizers (Dart Dataset)



(b) Training loss curve over training time when fine-tuning using AdamW-based method.

Figure 13: Training loss of GPT-2 small model ($r = 4$) fine-tuned using different optimizers. Evaluation is conducted on the DART Dataset.

F THE USE OF LARGE LANGUAGE MODELS (LLMs)

We use LLMs to polish writing.