

OD-LoRA: OVERCOMING THE DILEMMA BETWEEN WEIGHT REPRESENTATION AND GRADIENT APPROXIMATION IN LOW-RANK ADAPTATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Low-Rank Adaptation (LoRA) enables efficient adaptation of large pre-trained models to downstream tasks by representing weight updates with trainable low-rank (LR) matrices. Recent studies have shown a different perspective that learning with LoRA is equivalent to using a low-rank approximation of the full fine-tuning gradient, obtained by mapping it onto low-rank subspaces through the LR matrices. In this paper, we theoretically show that LoRA faces a dilemma between these two perspectives: weight update representation and gradient approximation. We first demonstrate that the quality of gradient approximation is improved if the LR matrices have uniform singular values, since non-uniform singular values anisotropically distort the projection of the full gradient onto the subspaces. However, this condition entails a strict constraint on the weight updates, significantly compromising their representational capacity. To **O**vercome this **D**ilemma, we introduce a new method, named OD-LoRA, which decouples the approximated gradient from the singular values of the LR matrices. Specifically, OD-LoRA ensures that the full gradient is mapped through the orthonormal bases of the low-rank subspaces defined by LR matrices, achieving perfect projection onto the subspaces, while still allowing the singular values to represent the weight updates. Consequently, OD-LoRA achieves both the optimal condition for accurate gradient approximation and unconstrained representation of weight updates simultaneously. The experimental results on natural language and vision benchmarks demonstrate that OD-LoRA improves loss convergence and gradient approximation quality, significantly enhancing the adaptation performance of LoRA.

1 INTRODUCTION

Large-scale pretrained Transformer (Vaswani et al., 2017) models have achieved remarkable success across a wide range of domains, including natural language processing and computer vision. With massive datasets and extensive computational resources, these models are trained to capture rich representations that can be effectively fine-tuned and transferred to diverse downstream tasks. For instance, models such as GPT (Radford et al., 2018), LLaMA (Touvron et al., 2023), BERT (Devlin et al., 2019), and Vision Transformers (ViTs) (Dosovitskiy et al., 2021) have set new benchmarks in language understanding, text generation, and image classification. This pretraining-and-adaptation paradigm has become the standard in the field of artificial intelligence, reducing the need for training models from scratch and driving rapid progress in both research and real-world applications.

However, due to the massive number of parameters, adapting those pretrained models through fine-tuning is often computationally expensive and memory-intensive. For example, GPT-3 (Brown et al., 2020) has about 175 billion parameters, and the parameters of LLaMA-3 models (Grattafiori et al., 2024) range from 8 to 405 billion, making their full fine-tuning impractical. Upon the observation that large pretrained models only require low-rank updates for adaptation (Gooneratne et al., 2020), Low-Rank Adaptation (LoRA) (Hu et al., 2022) addresses the challenge by representing weight updates for adaptation with trainable low-rank matrices (LoRA matrices), reducing the number of trainable parameters to about 1% of the total. Remarkably, LoRA not only demonstrates effective adaptation performance in various natural language and vision benchmarks but also enables the easy merging or detaching of LoRA matrices from the pretrained model, allowing for modular usage.

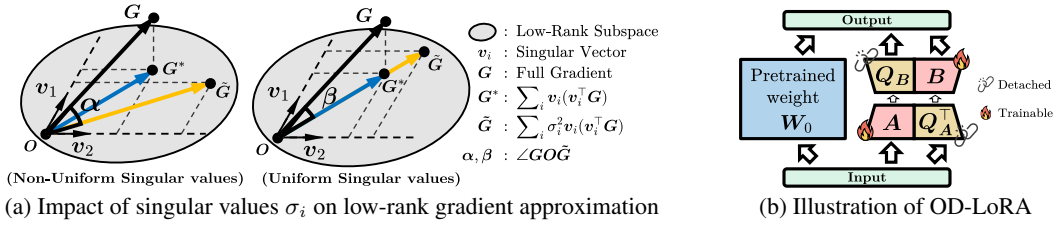


Figure 1: **(a)** Learning dynamics in LoRA are equivalent to learning with a low-rank approximation (\tilde{G}) of a full fine-tuning gradient (G) that is obtained by mapping G onto low-rank subspaces through the LoRA matrices (A and B). Non-uniform singular values of A and B distort the projection of G (G^*) by anisotropically scaling it. By contrast, when the singular values are uniform, the direction of G^* is preserved, resulting in \tilde{G} with the smallest angular distance from G ($\alpha > \beta$). **(b)** OD-LoRA leverages the orthonormal bases (Q_A and Q_B) of the subspaces defined by A and B . The computation process of these bases is not accounted for in the gradient computation. By multiplying A with Q_B^\top and B with Q_A^\top , the gradient approximation is decoupled from the singular values of A and B , while they still represent the weight updates. This allows OD-LoRA to achieve unconstrained weight update representation and accurate gradient approximation simultaneously.

Due to its advantages, LoRA has inspired numerous follow-up studies. Among them, several recent works (Wang et al., 2024; Zhang et al., 2025; Wang et al., 2025) view LoRA from the perspective of low-rank gradient approximation: the learning dynamics in LoRA are equivalent to full fine-tuning with a low-rank approximation of the full gradient. Specifically, the approximated gradient is obtained by mapping the full gradient onto the low-rank subspaces defined by the LoRA matrices. These works propose initialization strategies for the LoRA matrices that leverage the first-step full gradient or adjust the approximated gradient to reduce the approximation error. However, they either overlook the gradient approximation error after the first iteration or incur substantial overhead during training.

In this paper, we provide new insights into LoRA from the perspective of low-rank gradient approximation. Specifically, we prove that, given the subspaces defined by the LoRA matrices, the gradient approximation quality is optimized if the singular values of the LoRA matrices are uniform. As illustrated in Figure 1a, the non-uniform singular values result in anisotropic scaling on the projection of the full gradient onto the low-rank subspaces defined by the LoRA matrices, distorting the projection. By contrast, when the singular values are uniform, the low-rank gradient approximation reduces to a scaled projection, achieving the maximum alignment with the full gradient. However, we show that the uniform singular values of the LoRA matrices compel the non-zero singular values of weight updates to also be uniform, which significantly limits their representational capacity. Therefore, we argue that LoRA faces a dilemma between the representational capacity of weight updates and achieving accurate gradient approximation.

We attribute this dilemma to the dependence of the gradient approximation on the singular values of the LoRA matrices. Thus, to **Overcome the Dilemma in LoRA**, we propose a new formulation of low-rank weight updates, called OD-LoRA, which eliminates this dependence by leveraging the orthonormal bases of the subspaces defined by the LoRA matrices. As illustrated in Figure 1b, OD-LoRA treats these bases as constant in gradient computation and multiplies each LoRA matrix with the bases of the other LoRA matrix. We theoretically show that this formulation makes the low-rank gradient approximation independent of the singular values of the LoRA matrices, preserving the direction of the full gradient projection. Simultaneously, the singular values remain responsible for representing the weight updates, thereby retaining their full representational capacity. Furthermore, because these bases are often unstable due to the small norm of the LoRA matrices early in training, we propose initializing them through a few training iterations. Through experiments on various benchmarks and pretrained models, we demonstrate that OD-LoRA converges faster to a lower loss and outperforms LoRA and its variants with only minor overhead in terms of memory and time.

We summarize the contribution of our work as follows:

- We show that the gradient approximation quality in LoRA is improved if the LoRA matrices have uniform singular values, while leading to a strict constraint on the weight updates.
- To overcome this dilemma, we propose a new method, OD-LoRA, that achieves both accurate gradient approximation and unconstrained weight update representation simultaneously.
- OD-LoRA improves loss convergence and the alignment between the full gradient and its approximation, resulting in better performance than LoRA and existing methods.

2 RELATED WORKS

Low-Rank Adaptation (LoRA). To reduce the computational cost in fine-tuning large pretrained models, several parameter-efficient methods have been proposed, including adapter-based methods (Houlsby et al., 2019; He et al., 2022), prefix- or prompt-tuning (Li & Liang, 2021; Lester et al., 2021), and selective update by searching for salient parameters Guo et al. (2021); Sung et al. (2021). While effective, these methods either introduce additional parameters at inference or rely heavily on carefully designed search criteria. On the other hand, LoRA (Hu et al., 2022) introduces trainable low-rank matrices to approximate weight updates, which can be seamlessly merged with the pre-trained weights at inference. Its advantages and strong performance have motivated numerous follow-up studies Kalajdzievski (2023); Hayou et al. (2024b); Si et al. (2025). Several works propose alternative formulations of low-rank weight updates to further improve the performance–efficiency trade-off (Liu et al., 2024; Kopiczko et al., 2024; Lingam et al., 2024; Albert et al., 2025; Koochpayegani et al., 2024). Other research directions include developing effective initialization strategies for the trainable low-rank matrices (Hayou et al., 2024a; Wang et al., 2024; Meng et al., 2024; Zhang et al., 2025; Li et al., 2025), improving the rank of weight updates (Huang et al., 2025; Lialin et al., 2024), and designing layer-wise adaptive rank search methods (Zhang et al., 2023; Ke et al., 2024).

Low-Rank Gradient Approximation. There have been attempts to reduce the computational cost during training by approximating the original gradient in a low-rank subspace (Gooneratne et al., 2020; Zhao et al., 2024). Several recent studies have examined LoRA from the perspective of low-rank gradient approximation (Wang et al., 2024; Zhang et al., 2025; Wang et al., 2025; Zhang & Pilanci, 2024; Yu et al., 2025). Wang et al. (2024); Zhang et al. (2025) propose initialization strategies for LoRA that minimize the gradient approximation error at the first training step. [Similar to our approach, Wang et al. \(2025\); Zhang & Pilanci \(2024\); Yu et al. \(2025\) aim to optimize the update geometry via preconditioning. However, as discussed in Section D, such methods tend to amplify useless gradients when the singular value distribution of LoRA matrices is skewed, affecting the training dynamics. By contrast, our method structurally decouples the gradient approximation from the spectral properties of LoRA matrices, ensuring more robust training dynamics.](#)

3 PROPOSED METHOD

3.1 PRELIMINARIES

Low-Rank Weight Update Representation. Let $\mathbf{W}_0 \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$ denote the weight matrix of a linear layer in a pretrained model, with input dimension d_{in} and output dimension d_{out} . To adapt the model to a downstream task, we train it to learn task-specific weight updates. Instead of updating all the parameters of \mathbf{W}_0 , Low-Rank Adaptation (LoRA) assumes the low-rank structure of desired weight updates Gooneratne et al. (2020) and learns a low-rank decomposition of the updates:

$$\mathbf{W}_{\text{eff}} = \mathbf{W}_0 + \Delta \mathbf{W}_{\text{LoRA}} = \mathbf{W}_0 + s\mathbf{B}\mathbf{A}, \quad (1)$$

where $\mathbf{A} = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_r]^\top \in \mathbb{R}^{r \times d_{\text{in}}}$ and $\mathbf{B} = [\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_r] \in \mathbb{R}^{d_{\text{out}} \times r}$ are trainable low-rank matrices and s is a scaling factor. For one layer, the total number of trainable parameters in LoRA is $r(d_{\text{out}} + d_{\text{in}})$, which is much smaller than that of full-parameter training case, $d_{\text{out}} \times d_{\text{in}}$, if $r \ll \min(d_{\text{out}}, d_{\text{in}})$. It is common to initialize \mathbf{A} from a random Gaussian or uniform distribution and to initialize \mathbf{B} to zero, ensuring that the model initially behaves identically to the pretrained one. Once \mathbf{A} and \mathbf{B} are trained, they can be easily merged with or detached from the pretrained weights.

Low-Rank Gradient Approximation Perspective. Similar to Wang et al. (2024); Zhang et al. (2025); Wang et al. (2025), we view LoRA from the low-rank gradient approximation perspective. Let $\mathbf{G} = \frac{\partial \mathcal{L}}{\partial \mathbf{W}_{\text{eff}}}$ denote the gradient of loss \mathcal{L} with respect to \mathbf{W}_{eff} . We refer to \mathbf{G} as *full gradient* since it is used in a full-parameter fine-tuning scenario. The gradient with respect to \mathbf{A} and \mathbf{B} is computed as $\frac{\partial \mathcal{L}}{\partial \mathbf{A}} = s\mathbf{B}^\top \mathbf{G}$ and $\frac{\partial \mathcal{L}}{\partial \mathbf{B}} = s\mathbf{G}\mathbf{A}^\top$. Then, the differential of \mathbf{A} and \mathbf{B} in the context of gradient descent with a learning rate η are given by $d\mathbf{A} = -\eta \frac{\partial \mathcal{L}}{\partial \mathbf{A}} = -\eta s\mathbf{B}^\top \mathbf{G}$ and $d\mathbf{B} = -\eta \frac{\partial \mathcal{L}}{\partial \mathbf{B}} = -\eta s\mathbf{G}\mathbf{A}^\top$. Consequently, the differential of \mathbf{W}_{eff} is expressed as

$$d\mathbf{W}_{\text{eff}} = s(d\mathbf{B}\mathbf{A} + \mathbf{B}d\mathbf{A}) = -\eta s^2(\mathbf{G}\mathbf{A}^\top \mathbf{A} + \mathbf{B}\mathbf{B}^\top \mathbf{G}) = -\eta \tilde{\mathbf{G}}. \quad (2)$$

This implies that LoRA is equivalent to full fine-tuning \mathbf{W}_{eff} using the equivalent gradient $\tilde{\mathbf{G}}$. Let $\mathcal{R}(\mathbf{A})$ and $\mathcal{C}(\mathbf{B})$ denote the row space of \mathbf{A} and the column space of \mathbf{B} , respectively. Then, $\mathbf{G}\mathbf{A}^\top \mathbf{A}$ and $\mathbf{B}\mathbf{B}^\top \mathbf{G}$ can be interpreted as the mapping of the rows of \mathbf{G} onto $\mathcal{R}(\mathbf{A})$ and the mapping of columns of \mathbf{G} onto $\mathcal{C}(\mathbf{B})$. Since these subspaces are low rank, $\tilde{\mathbf{G}}$ can be interpreted as a low-rank approximation of \mathbf{G} .

3.2 DILEMMA IN LORA: WEIGHT REPRESENTATION VS. GRADIENT APPROXIMATION

Let $\mathcal{W} = \{\mathbf{W}^l = \mathbf{W}_0^l + s\mathbf{B}^l\mathbf{A}^l\}_{l=1}^L$ denote the set of the effective weights across L layers. Let $\mathcal{G} = \{\mathbf{G}^l\}_{l=1}^L$ and $\tilde{\mathcal{G}} = \{\tilde{\mathbf{G}}^l\}_{l=1}^L$ denote the full gradient of loss \mathcal{L} with respect to \mathcal{W} and its low-rank approximation, respectively. Assuming that \mathcal{L} is β -smooth, we can derive the following inequality on the loss decrease guarantee when \mathcal{W} is updated using $\tilde{\mathcal{G}}$:

$$\mathcal{L}(\mathcal{W} - \eta\tilde{\mathcal{G}}) - \mathcal{L}(\mathcal{W}) \leq -\eta\langle \mathcal{G}, \tilde{\mathcal{G}} \rangle + \frac{\beta}{2}\eta^2\|\tilde{\mathcal{G}}\|^2. \quad (3)$$

Here, $\langle \mathcal{G}, \tilde{\mathcal{G}} \rangle := \sum_{l=1}^L \langle \mathbf{G}^l, \tilde{\mathbf{G}}^l \rangle_F$ and $\|\mathcal{G}\| = \sqrt{\langle \mathcal{G}, \mathcal{G} \rangle}$ where $\langle \cdot, \cdot \rangle_F$ denotes the Frobenius inner product. Details are provided in Appendix E. Equation 3 suggests that for a sufficiently small η , the approximated gradient will result in better loss convergence as $\langle \mathcal{G}, \tilde{\mathcal{G}} \rangle$ increases. Thus, in contrast to prior works that define gradient approximation error as $\|\mathcal{G} - \tilde{\mathcal{G}}\|$ (Wang et al., 2024; Zhang et al., 2025; Wang et al., 2025), we propose using $\langle \mathcal{G}, \tilde{\mathcal{G}} \rangle$ as a measure of the approximation quality of $\tilde{\mathcal{G}}$. We derive the following theorem on the conditions required for the LoRA matrices to improve the approximation quality:

Theorem 1: Optimal Condition for Improving Gradient Approximation Quality in LoRA

Consider the gradient approximation quality in LoRA:

$$\langle \mathcal{G}, \tilde{\mathcal{G}} \rangle = \sum_{l=1}^L \langle \mathbf{G}^l, \tilde{\mathbf{G}}^l \rangle_F = s^2 \sum_{l=1}^L \langle \mathbf{G}^l, \mathbf{G}^l \mathbf{A}^{l\top} \mathbf{A}^l + \mathbf{B}^l \mathbf{B}^{l\top} \mathbf{G}^l \rangle_F.$$

Assume the energy of \mathbf{A}^l and \mathbf{B}^l , the row space of \mathbf{A}^l , and the column space of \mathbf{B}^l are given for all l . Then, $\langle \mathcal{G}, \tilde{\mathcal{G}} \rangle$ is maximized if \mathbf{A}^l and \mathbf{B}^l have uniform singular values for all l .

See Section A for the proof. Let $\mathbf{A} = \mathbf{U}_A \Sigma_A \mathbf{V}_A^\top$ and $\mathbf{B} = \mathbf{U}_B \Sigma_B \mathbf{V}_B^\top$ be the compact singular value decomposition (SVD) of \mathbf{A} and \mathbf{B} . To examine how the distribution of singular values of \mathbf{A} and \mathbf{B} affects the gradient approximation quality, we rewrite $\tilde{\mathcal{G}}$ as $\mathbf{G}\mathbf{A}^\top\mathbf{A} + \mathbf{B}\mathbf{B}^\top\mathbf{G} = \mathbf{G}\mathbf{V}_A \Sigma_A^2 \mathbf{V}_A^\top + \mathbf{U}_B \Sigma_B^2 \mathbf{U}_B^\top \mathbf{G}$. Without the singular value terms, $\mathbf{G}\mathbf{V}_A \mathbf{V}_A^\top$ and $\mathbf{U}_B \mathbf{U}_B^\top \mathbf{G}$ are exactly the orthogonal projections of the rows of \mathbf{G} onto $\mathcal{R}(\mathbf{A})$ and the columns of \mathbf{G} onto $\mathcal{C}(\mathbf{B})$, respectively. This implies that the singular values of \mathbf{A} and \mathbf{B} act as weights on the components of the projection, which are aligned with the corresponding singular vectors. As illustrated in Figure 1a, if the singular values are non-uniform, they anisotropically scale the projected components, distorting the projection of \mathbf{G} onto the subspaces. On the other hand, if \mathbf{A} and \mathbf{B} have uniform singular values, the projection is scaled uniformly in all directions, yielding $\tilde{\mathcal{G}}$ that is maximally aligned with \mathcal{G} . The proof shows that given the energy assumption on \mathbf{A} and \mathbf{B} , this maximal alignment result in maximizing $\langle \mathcal{G}, \tilde{\mathcal{G}} \rangle_F$.

However, in the context of the LoRA formulation $\Delta\mathbf{W}_{\text{LoRA}} = s\mathbf{B}\mathbf{A}$, the uniform singular value condition may impose an overly restrictive constraint on the representational capacity. Motivated by the universal approximation theorem (Hornik et al., 1989), we formally define the representational capacity of rank- r weight updates $\Delta\mathbf{W}$ as its capability to represent an arbitrary rank- r matrix:

Definition 1 (Representational capacity). For a rank- r matrix $\mathbf{X} \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$, we define the representational capacity of $\Delta\mathbf{W}$ as the reciprocal of the minimum achievable error between \mathbf{X} and $\Delta\mathbf{W}$:

$$\text{Cap}_{\mathbf{X}}(\Delta\mathbf{W}) := \frac{1}{\min_{\Delta\mathbf{W}} \|\mathbf{X} - \Delta\mathbf{W}\|_F^2}.$$

While $\text{Cap}_{\mathbf{X}}(\Delta\mathbf{W})$ can be further quantified through an optimization or expectation over \mathbf{X} , this is unnecessary in this paper. Using this definition, we prove that LoRA faces a dilemma between the representational capacity $\text{Cap}_{\mathbf{X}}(\Delta\mathbf{W}_{\text{LoRA}})$ and the gradient approximation quality:

Theorem 2: Dilemma in LoRA: Representational Capacity or Gradient Approximation Quality?

The formulation of LoRA, $\Delta\mathbf{W}_{\text{LoRA}} = s\mathbf{B}\mathbf{A}$ where both \mathbf{A} and \mathbf{B} are rank- r , faces the following dilemma between representational capacity and gradient approximation quality:

- **Full representational capacity but suboptimal gradient approximation.** Without any constraints, $\text{Cap}_{\mathbf{X}}(\Delta\mathbf{W}_{\text{LoRA}}) = \infty$, but this requires either \mathbf{A} or \mathbf{B} to have non-uniform singular values, which impedes improving the gradient approximation quality in Theorem 1.
- **Accurate gradient approximation but limited representational capacity.** If both \mathbf{A} and \mathbf{B} have uniform singular values (satisfying the optimal condition of Theorem 1), the non-zero singular values of $\Delta\mathbf{W}_{\text{LoRA}}$ are uniform, significantly limiting the representational capacity: $\text{Cap}_{\mathbf{X}}(\Delta\mathbf{W}_{\text{LoRA}}) = 1 / \sum_{i=1}^r (\sigma_i - \frac{1}{r} \sum_{i=1}^r \sigma_i)^2$, where $\mathbf{X} \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$ is an arbitrary rank- r matrix with non-zero singular values $\{\sigma_1, \sigma_2, \dots, \sigma_r\}$.

See Section B for the proof. Theorem 2 shows that the LoRA formulation is unable to achieve both optimal gradient approximation quality and full representational capacity at the same time.

Notably, under the uniform singular value condition, the representational capacity is reduced to $1 / \sum_{i=1}^r (\sigma_i - \frac{1}{r} \sum_{i=1}^r \sigma_i)^2$ where the denominator is proportional to the variance of the singular values of a target rank- r matrix. This suggests that if the desired weight updates favor singular value spectrum with a high variance, the achievable minimum error between $\Delta\mathbf{W}_{\text{LoRA}}$ and the desired updates is large under the uniform singular value condition. Indeed, Figure 5 shows that the weight updates obtained via full fine-tuning typically exhibit highly skewed singular values, suggesting that the uniform singular value condition may hinder the representation of desired weight updates.

3.3 OD-LoRA: DECOUPLING THE GRADIENT APPROXIMATION FROM SINGULAR VALUES

We argue that this dilemma arises because in the formulation of LoRA, both the weight updates and the approximated gradient depend on the singular values of \mathbf{A} and \mathbf{B} . Thus, to Overcome the Dilemma, we propose a new LoRA formulation, named OD-LoRA, which decouples the low-rank approximation of the full gradient from the singular values of \mathbf{A} and \mathbf{B} , while preserving the maximal representational capacity of the weight updates. To do so, we leverage the orthonormal bases of $\mathcal{R}(\mathbf{A})$ and $\mathcal{C}(\mathbf{B})$ since they are independent of the singular values of \mathbf{A} and \mathbf{B} . Specifically, we compute the orthonormal bases via QR decomposition: $\mathbf{A}^\top = \mathbf{Q}_A \mathbf{R}_A$ and $\mathbf{B} = \mathbf{Q}_B \mathbf{R}_B$ where the columns of \mathbf{Q}_A and \mathbf{Q}_B are the orthonormal bases of $\mathcal{R}(\mathbf{A})$ and $\mathcal{C}(\mathbf{B})$. We then formulate the weight updates so that the gradients for \mathbf{A} and \mathbf{B} depend only on these bases, while the singular values of \mathbf{A} and \mathbf{B} remain responsible for representing the weight updates:

$$\Delta\mathbf{W}_{\text{OD-LoRA}} = s(\mathbf{Q}_B \mathbf{A} + \mathbf{B} \mathbf{Q}_A^\top) = s[\mathbf{Q}_B \ \mathbf{B}] \begin{bmatrix} \mathbf{A} \\ \mathbf{Q}_A^\top \end{bmatrix}, \quad (4)$$

During the backward pass, we treat \mathbf{Q}_A and \mathbf{Q}_B as constants. Hence, the gradient with respect to \mathbf{A} and \mathbf{B} is expressed in terms of \mathbf{Q}_B and \mathbf{Q}_A , respectively: $\frac{\partial \mathcal{L}}{\partial \mathbf{A}} = s \mathbf{Q}_B^\top \mathbf{G}$ and $\frac{\partial \mathcal{L}}{\partial \mathbf{B}} = s \mathbf{G} \mathbf{Q}_A$. Thus, the singular values of \mathbf{A} and \mathbf{B} need not be constrained to satisfy the optimal condition of Theorem 1, allowing unconstrained representation of $\Delta\mathbf{W}_{\text{OD-LoRA}}$. We theoretically show that our new formulation address the dilemma in LoRA:

Theorem 3: OD-LoRA Overcomes the Dilemma in LoRA

The formulation of OD-LoRA in Equation 4 possesses the following properties simultaneously.

1. Define $\mathbf{W}_{\text{eff}} = \mathbf{W}_0 + \Delta\mathbf{W}_{\text{OD-LoRA}}$. Then, the low-rank approximation of full gradient \mathbf{G} is given by

$$\tilde{\mathbf{G}} = s^2(\mathbf{G} \mathbf{Q}_A \mathbf{Q}_A^\top + \mathbf{Q}_B \mathbf{Q}_B^\top \mathbf{G}),$$

where \mathbf{Q}_A and \mathbf{Q}_B have uniform singular values (all equal to 1), satisfying the optimal condition for improving the gradient approximation quality in Theorem 1.

2. For an arbitrary rank- r matrix $\mathbf{X} \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$, $\text{Cap}_{\mathbf{X}}(\Delta\mathbf{W}_{\text{OD-LoRA}}) = \infty$, implying $\Delta\mathbf{W}_{\text{OD-LoRA}}$ can represent any rank- r matrices.

See Section C for the proof. Theorem 3 demonstrates that OD-LoRA is able to achieve both the accurate low-rank approximation of the full gradient and the maximum representational capacity of weight updates simultaneously, overcoming the dilemma in LoRA.

Algorithm 1: Training Process of OD-LoRA

```

270
271
272 Input: pretrained weight  $W_0$ , training dataset  $\mathcal{D} = \{\mathcal{B}_1, \dots, \mathcal{B}_{T_{\text{total}}}\}$ , and  $A$  and  $B$  s.t.  $BA = 0$ 
273 // Initialize  $Q_A$ , and  $Q_B$ 
274 for  $t = 1$  to  $N \times T$  do
275   if  $t \leq T$  then
276      $\Delta W \leftarrow sBA$ 
277   else
278      $\Delta W \leftarrow s[Q_B B] \begin{bmatrix} A \\ Q_A^\top \end{bmatrix}$  // Equation 4
279     Forward with  $\mathcal{B}_t$  using  $W_{\text{eff}} = W_0 + \Delta W$ 
280     Backward and update  $A$  and  $B$ 
281     if  $t \bmod T = 0$  then
282       Compute the rank- $r$  truncated SVD of  $\Delta W$ :  $\Delta W = U_r \Sigma_r V_r^\top$ 
283        $A \leftarrow 0$ ,  $B \leftarrow 0$ ,  $Q_A \leftarrow V_r$ ,  $Q_B \leftarrow U_r$ 
284       Clear all optimizer states
285 // Currently,  $A = 0$ ,  $B = 0$ ,  $Q_A = V_r$ ,  $Q_B = U_r$ 
286 for  $t = 1$  to  $T_{\text{total}}$  do
287    $\Delta W \leftarrow s[Q_B B] \begin{bmatrix} A \\ Q_A^\top \end{bmatrix}$ 
288   Forward with  $\mathcal{B}_t$  using  $W_{\text{eff}} = W_0 + \Delta W$ 
289   Backward and update  $A$  and  $B$ 
290   if  $t = T$  then
291     // Ensure  $Q_A$  and  $Q_B$  form bases for  $\mathcal{R}(A)$  and  $\mathcal{C}(B)$ 
292     Compute the rank- $r$  truncated SVD of  $\Delta W$ :  $\Delta W = U_r \Sigma_r V_r^\top$ 
293      $A \leftarrow \frac{1}{2s} \Sigma_r V_r^\top$ ,  $B \leftarrow \frac{1}{2s} U_r \Sigma_r$ ,  $Q_A \leftarrow V_r$ ,  $Q_B \leftarrow U_r$  // Equation 5
294     Clear all optimizer states
295   else if  $t > T$  and  $t \bmod \lceil \frac{5t}{T_{\text{total}}} \rceil$  then
296     // Adjust the update interval according to  $t$ 
297     Compute the QR decomposition of  $A^\top$  and  $B$ :  $A^\top = Q'_A R'_A$  and  $B = Q'_B R'_B$ 
298      $Q_A \leftarrow Q'_A$ ,  $Q_B \leftarrow Q'_B$ 
299 return  $A, B$ 

```

At the beginning of training, $\Delta W_{\text{OD-LoRA}}$ must be zero to ensure that learning starts from the pretrained weight. Thus, during the initial training phase, representing the weight updates as in Equation 4 may be ineffective, since $\mathcal{R}(A)$ and $\mathcal{C}(B)$ are not yet well-defined and can change rapidly due to small norms. To address this, at the beginning of training, we allow Q_A and Q_B to not coincide with the orthonormal bases of these subspaces and initialize them through a few training iterations. Specifically, we repeat the following processes N times: training A and B for T steps, updating Q_A and Q_B using the top- r right and left singular vectors of the resulting weight updates, and resetting A and B to zero. For the first iteration, we use the LoRA formulation due to the absence of Q_A and Q_B , and switch to the OD-LoRA formulation thereafter. Additionally, at the end of each iteration, we clear all optimizer states (i.e., momentum buffers) to eliminate their impact on the next iteration. After these $N \times T$ steps, we obtain Q_A and Q_B that are effectively initialized for the early training stage. However, they are not yet the bases of $\mathcal{R}(A)$ and $\mathcal{C}(B)$ since $A = B = 0$. To address this mismatch, we train A and B using the fixed Q_A and Q_B for the first T iterations, and then set

$$A = \frac{1}{2s} \Sigma_r V_r^\top, \quad B = \frac{1}{2s} U_r \Sigma_r, \quad Q_A = V_r, \quad Q_B = U_r, \quad (5)$$

where $\Delta W_{\text{OD-LoRA}} = U_r \Sigma_r V_r^\top$ is the rank- r truncated SVD. We also reset all optimizer states after this process. By doing so, we ensure that the columns of Q_A and Q_B form the orthonormal bases of $\mathcal{R}(A)$ and $\mathcal{C}(B)$, respectively. After the initialization phase, we update Q_A and Q_B because $\mathcal{R}(A)$ and $\mathcal{C}(B)$ change during training. See Algorithm 1 for the overall training process.

Training overhead. We set N and T to 5 and 10, respectively, resulting in 6 SVD computations and roughly a 1% increase in training iterations. To reduce the overhead of the QR decomposition, we adjust the update interval for Q_A and Q_B . Specifically, we divide training into five phases, setting the update interval to one training step for the first phase to account for the rapid change of the subspaces during early training, and then increasing it by one after each subsequent phase. As a result, the initialization strategy and the updates of Q_A and Q_B incur minor training overhead, as detailed in Section 4.7.

Table 1: Results on commonsense reasoning tasks.

Model	Method	Rank	BoolQ	PIQA	SIQA	HellaSwag	WinoGrande	ARC-c	ARC-e	OBQA	Avg.
Gemma-2B	Full FT.	-	66.96±0.21	80.31±0.12	75.85±0.47	85.48±0.43	75.61±0.67	65.16±0.70	80.32±0.38	76.27±0.34	75.75
	LoRA		63.94±0.04	76.46±0.28	70.62±0.36	44.50±0.48	65.48±0.74	57.20±0.28	75.06±0.56	65.60±1.41	64.86
	rsLoRA		64.58±0.07	77.68±0.56	72.65±0.60	55.14±1.73	69.48±0.15	58.70±0.00	76.19±0.38	69.40±0.57	67.98
	LoRA+	8	63.35±0.00	76.82±0.03	72.00±0.36	77.18±1.01	70.90±1.49	58.47±0.44	76.06±0.02	70.20±0.28	70.62
	PiSSA		64.94±0.20	77.49±0.13	72.79±0.31	59.89±3.06	69.27±0.04	58.22±0.52	76.18±0.12	71.33±0.19	68.77
	DoRA		64.72±0.16	77.57±0.41	72.77±0.36	56.56±0.61	69.82±0.41	59.02±0.56	75.94±0.08	69.73±0.47	68.27
	OD-LoRA		65.85±0.37	79.22±0.31	74.51±0.43	80.34±0.96	73.11±0.41	62.17±0.20	78.77±0.44	73.13±0.75	73.39
	LoRA		64.24±0.10	76.22±0.15	70.71±0.05	45.05±1.40	65.98±0.11	57.51±0.24	74.96±0.18	65.13±0.47	64.98
	rsLoRA		65.63±0.04	78.89±0.62	73.97±0.17	77.69±0.00	71.03±0.45	61.26±0.72	78.55±0.08	72.80±1.41	72.48
	LoRA+	32	64.62±0.14	78.32±0.10	73.18±0.65	81.73±0.10	71.48±0.45	61.59±0.20	76.91±0.30	73.20±0.28	72.63
PiSSA		65.74±0.32	77.98±0.05	74.14±0.34	72.90±1.05	70.96±0.33	59.47±0.97	76.91±0.28	74.40±0.28	71.56	
DoRA		65.30±0.07	78.89±0.77	74.19±0.10	74.17±1.16	71.56±0.19	61.06±0.68	78.56±0.20	72.60±1.41	72.04	
OD-LoRA		66.40±0.16	80.92±0.28	75.04±0.31	84.89±0.34	74.01±0.15	62.97±0.60	79.52±0.08	75.47±0.38	74.90	
LLaMA3-8B	Full FT.	-	75.24±0.21	89.70±0.22	81.51±0.24	96.11±0.15	87.66±0.17	81.23±0.54	92.07±0.19	88.02±0.31	86.44
	LoRA		72.98±0.06	87.60±0.31	79.67±0.39	94.35±0.14	83.61±0.15	78.95±0.44	90.14±0.02	83.87±1.04	83.89
	rsLoRA		73.40±0.19	88.23±0.26	80.26±0.41	95.19±0.03	84.98±0.19	79.21±0.32	90.70±0.12	84.80±0.28	84.60
	LoRA+	8	73.07±0.61	88.03±0.44	80.22±0.29	94.75±0.01	85.00±0.07	78.95±0.36	90.16±0.06	85.20±0.00	84.42
	PiSSA		73.68±0.01	88.16±0.10	80.37±0.19	95.20±0.05	85.82±0.19	80.12±0.60	90.36±0.06	85.67±0.66	84.92
	DoRA		73.20±0.01	87.85±0.10	80.21±0.27	95.22±0.00	84.37±0.45	79.66±0.16	90.53±0.12	85.53±0.19	84.57
	OD-LoRA		74.02±0.12	88.88±0.36	81.01±0.22	96.06±0.06	87.13±0.22	79.89±0.56	90.88±0.22	86.33±0.09	85.53
	LoRA		73.06±0.13	87.45±0.05	79.68±0.00	94.50±0.08	83.32±0.15	79.58±0.52	90.31±0.02	84.07±0.75	84.00
	rsLoRA		72.56±0.89	88.47±0.08	80.74±0.39	91.80±5.36	85.61±0.52	79.69±0.60	90.90±0.08	85.53±1.04	84.41
	LoRA+	32	73.43±0.30	88.25±0.10	80.03±0.05	94.99±0.01	85.87±0.41	79.92±0.20	90.05±0.32	85.73±1.04	84.78
PiSSA		74.46±0.52	89.03±0.44	80.79±0.41	95.48±0.02	86.98±0.11	80.12±0.36	90.75±0.14	86.33±0.66	85.49	
DoRA		73.86±0.23	88.72±0.33	80.98±0.27	95.67±0.05	85.85±0.41	79.78±0.48	91.08±0.00	85.87±0.09	85.23	
OD-LoRA		74.65±0.09	89.48±0.05	81.01±0.07	96.11±0.00	88.08±0.56	81.06±0.36	92.02±0.22	85.93±1.04	86.04	

Table 2: Results on natural language generation tasks.

Method	Rank	Gemma-2B			LLaMA3-8B		
		MATH	GSM8K	HumanEval	MATH	GSM8K	HumanEval
Full FT.	-	19.17±0.22	56.23±0.23	33.35±0.65	26.47±0.41	77.04±0.18	49.11±0.44
LoRA		16.12±0.28	45.59±0.95	27.24±0.57	23.68±0.28	73.44±0.22	42.53±1.80
rsLoRA		17.10±0.24	49.10±1.44	29.27±1.79	24.52±0.09	75.71±0.23	42.23±1.63
LoRA+	8	17.09±0.46	50.30±0.81	27.64±0.76	24.65±0.63	75.27±0.59	45.27±2.84
PiSSA		16.44±0.38	50.97±0.74	28.86±0.58	24.43±0.23	74.96±1.01	46.19±2.64
DoRA		17.19±0.28	50.61±0.85	28.35±1.52	24.63±0.21	75.41±1.32	43.29±1.14
OD-LoRA		17.79±0.18	52.08±0.53	31.30±0.76	25.89±0.30	76.62±0.34	47.15±1.44
LoRA		16.31±0.13	45.67±0.91	28.66±0.50	23.66±0.25	73.49±0.59	42.53±1.09
rsLoRA		18.09±0.04	51.93±0.53	31.10±1.32	25.97±0.17	76.28±0.96	44.36±1.39
LoRA+	32	17.45±0.39	52.74±0.75	31.30±1.25	25.01±0.11	76.11±0.20	46.34±1.29
PiSSA		17.23±0.24	53.38±0.77	32.11±0.76	24.80±0.41	76.39±0.75	47.10±2.38
DoRA		18.27±0.14	52.41±0.24	30.79±1.52	25.87±0.42	76.42±0.17	44.97±1.00
OD-LoRA		18.47±0.14	55.34±1.10	34.15±0.50	26.45±0.39	77.22±0.19	47.76±1.25

Table 3: Results on image classification tasks with ViT-Base. See Section G for details and more results.

Method	Rank	Cars	CUB200	SUN397
Full FT.	-	84.26±0.19	86.35±0.18	74.76±0.26
LoRA		78.66±0.23	85.65±0.05	74.35±0.28
rsLoRA		78.48±0.07	85.67±0.41	74.63±0.13
LoRA+	8	79.05±0.25	85.28±0.11	72.87±0.28
PiSSA		78.03±0.32	85.17±0.24	72.82±0.10
DoRA		78.93±0.30	85.74±0.35	74.46±0.09
OD-LoRA		80.47±0.14	86.08±0.16	74.92±0.16
LoRA		81.48±0.20	85.75±0.28	70.82±0.07
rsLoRA		80.57±0.22	85.42±0.41	73.92±0.36
LoRA+	32	79.53±0.60	85.68±0.31	67.66±0.06
PiSSA		80.43±0.17	84.90±0.28	69.73±0.12
DoRA		81.07±0.03	85.67±0.22	74.05±0.21
OD-LoRA		83.04±0.32	86.06±0.24	74.92±0.21

4 EXPERIMENTS

4.1 EXPERIMENTAL SETUP

Models and Datasets. For natural language processing tasks, we use the pretrained Gemma-2B (Team et al., 2024) and LLaMA3-8B (Grattafiori et al., 2024). For commonsense reasoning tasks, we fine-tune on the Commonsense-170K dataset (Hu et al., 2023) and evaluate on eight standard benchmarks including BoolQ (Clark et al., 2019), PIQA (Bisk et al., 2020), SIQA (Sap et al., 2019), HellaSwag (Zellers et al., 2019), WinoGrande (Sakaguchi et al., 2020), ARC-c/e (Clark et al., 2018), and OBQA (Mihaylov et al., 2018). For natural language generation tasks, we select 100K examples from the MetaMathQA (Yu et al., 2024) and Code-Feedback (Zheng et al., 2024) datasets for fine-tuning. To assess performance, we use MATH (Hendrycks et al., 2021) and GSM8K (Cobbe et al., 2021) benchmarks for the models fine-tuned on MetaMathQA, and HumanEval (Chen et al., 2021) benchmark for the models fine-tuned on Code-Feedback. For the details on image classification tasks, see Section G.

Implementation Details. For optimization, we adopt the AdamW optimizer (Loshchilov & Hutter, 2019) following the standard setting. We set $\beta_1 = 0.9$, $\beta_2 = 0.999$, and apply zero weight decay for the optimizer. For all experiments, we set the hyperparameters of the proposed initialization to $N = 5$ and $T = 10$. For the scaling factor s , we fix it to 2 for LoRA and $\frac{\alpha}{\sqrt{r}}$ for the other methods, following rsLoRA (Kalajdzievski, 2023), where α is set to 4. We compare our method with the original LoRA and its recent variants, including rsLoRA (Kalajdzievski, 2023), LoRA+ (Hayou et al., 2024b) with scaling ratio 4, PiSSA (Meng et al., 2024), and DoRA (Liu et al., 2024). We reproduce the results of these methods under the same setting. We report the mean and standard deviation across three trials. Further details are provided in Table 9.

Table 4: **Ablation studies.** ‘Init.’ indicates the proposed initialization method, and $\text{Cap}_X(\Delta W)$ and ‘Optimal \tilde{G} ’ denote the representational capacity of weight updates and the equivalent gradient achieving the optimal condition in Theorem 1.

Method	$\text{Cap}_X(\Delta W) = \infty$	Optimal \tilde{G}	Init.	Gemma-2B			LLaMA3-8B		
				MATH	GSM8K	HumanEval	MATH	GSM8K	HumanEval
Full FT.	-	-	-	19.17±0.22	56.23±0.23	33.35±0.65	26.47±0.41	77.04±0.18	49.11±0.44
rsLoRA	✓	✗	✗	18.09±0.04	51.93±0.53	31.10±1.32	25.97±0.17	76.28±0.96	44.36±1.39
rsLoRA w/ Optimal \tilde{G}	✗	✓	✗	17.22±0.24	51.18±0.29	31.30±0.99	24.14±0.27	75.68±0.35	45.04±1.02
OD-LoRA w/ $\text{Cap}_X(\Delta W) < \infty$	✗	✓	✓	17.28±0.31	51.25±0.29	31.71±0.49	25.64±0.22	76.27±0.41	45.73±0.86
OD-LoRA w/ Suboptimal \tilde{G}	✓	✗	✓	18.18±0.11	53.73±0.19	32.32±1.73	25.56±0.40	76.02±0.32	44.92±0.29
OD-LoRA	✓	✓	✓	18.47±0.14	55.34±1.10	34.15±0.50	26.45±0.39	77.22±0.19	47.76±1.25

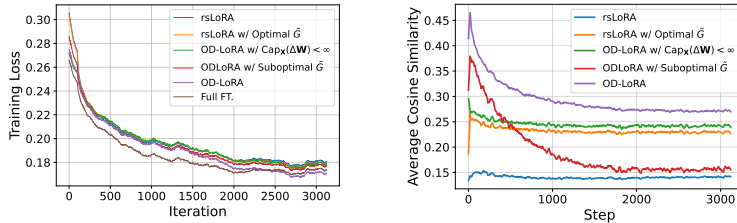


Figure 2: **Ablation studies using training loss curve (left) and alignment between the full gradient (G) and the equivalent gradient (\tilde{G}) (right).** Experiments are conducted using LLaMA3-8B and MetaMATHQA. We average the cosine similarity across all target layers.

4.2 COMMONSENSE REASONING

Table 1 shows the experimental results on the commonsense reasoning tasks. OD-LoRA outperforms LoRA and its variants in nearly all configurations, leading to the highest adaptation performance on average across all models and ranks. Specifically, OD-LoRA yields improvements of approximately 9 and 2 percentage points over LoRA on Gemma-2B and LLaMA3-8B, respectively, and surpasses recent variants including PiSSA and DoRA by 3–5 and 0.5–1 percentage points on Gemma-2B and LLaMA3-8B, respectively. In particular, Gemma-2B with OD-LoRA shows a substantial improvement over the others on the challenging HellaSwag benchmark, achieving about a 40 percentage point improvement over LoRA. Compared to the full fine-tuning (Full FT.), OD-LoRA with rank-32 achieves minimal performance gaps, demonstrating the effectiveness of OD-LoRA.

4.3 NATURAL LANGUAGE GENERATION

Table 2 shows the experimental results on the natural language generation tasks. OD-LoRA achieves the best adaptation performance on all benchmarks across all models and ranks. Specifically, OD-LoRA achieves improvements of approximately 1–2, 1–10, and 2–6 percentage points over existing methods on MATH, GSM8K, and HumanEval, respectively. Consistent with the commonsense reasoning results, OD-LoRA demonstrates substantial performance gains with Gemma-2B, highlighting its effectiveness, especially in challenging adaptation contexts. Furthermore, OD-LoRA with rank-32 achieves results comparable to or exceeding those of full fine-tuning. These results demonstrate that OD-LoRA significantly enhances the performance of LoRA.

4.4 IMAGE CLASSIFICATION

In addition to natural language processing tasks, we provide the results on image classification tasks. See Section G for details and more results. The results in Table 3 demonstrate that our method improves the performance of LoRA substantially, surpassing that of the other methods by significant margins. In particular, we observe that OD-LoRA shows about 2-3 percentage points higher accuracy in the Cars dataset, achieving the smallest performance gap compared to the full fine-tuning case. These results demonstrate the effectiveness of OD-LoRA in the vision domain.

4.5 ABLATION STUDIES AND CONVERGENCE ANALYSIS

To validate our claim on the dilemma in LoRA, we conduct ablation studies. Specifically, we propose 3 variants of rsLoRA and OD-LoRA: ‘rsLoRA w/ Optimal \tilde{G} ’, ‘OD-LoRA w/ $\text{Cap}_X(\Delta W) < \infty$ ’ and ‘OD-LoRA w/ Suboptimal \tilde{G} ’. For the first and second methods, we enforce uniform singular values on the weight updates to satisfy the condition for accurate gradient approximation, which restricts their representational capacity. The second method, in addition, incorporates the proposed initialization phase. For the third method, we modify OD-LoRA by intentionally violating the uniform singular value condition for the equivalent gradient \tilde{G} in Theorem 3. See Section I for more details.

Table 5: **Comparison with existing gradient-based methods.** We reproduce the results of existing methods. Experiments are conducted with rank 8. See Table 7 for more results and analyses.

Method	Gemma-2B			LLaMA3-8B		
	MATH	GSM8K	HumanEval	MATH	GSM8K	HumanEval
LoRA-GA	17.02±0.28	48.60±0.74	30.69±0.76	24.74±0.23	75.21±1.01	45.12±1.22
LoRA-Pro	16.52 ±0.14	44.58±0.55	28.66±0.99	23.90±0.44	72.93±0.17	40.24±0.99
rsLoRA + ScaledAdamW	16.64±0.31	45.56±0.53	26.83±0.99	24.00±0.17	73.39±0.20	43.29±1.15
AltLoRA	16.62±0.28	45.34±0.81	28.66±0.91	23.66±0.68	73.69±0.23	42.07±0.82
OD-LoRA	17.79±0.18	52.08±0.53	31.30±0.76	25.89±0.30	76.62±0.34	47.15±1.44

Table 6: **Analysis of training overhead.** We train the LLaMA3-8B model on a single NVIDIA H200 GPU for one epoch. For OD-LoRA, we separately report the time overhead incurred by the proposed initialization method (first term) and the subsequent training (second term).

Method	MetaMATHQA		Code-Feedback	
	Training Time (hours)	GPU-Memory (GB)	Training Time (hours)	GPU-Memory (GB)
LoRA	1.61	128.54	1.72	128.80
OD-LoRA	0.03 + 1.67 (+5.59%)	129.16 (+0.48%)	0.03 + 1.80 (+5.17%)	129.77 (+0.75%)

Table 4 presents the adaptation performance, and Figure 2 presents the training loss curves and gradient approximation quality. We observe that satisfying the uniform singular value condition leads to better gradient approximation quality, substantiating Theorem 1. However, we find that satisfying only one of the two requirements (i.e., full representational capacity of weight updates or accurate gradient approximation) fails to improve performance and loss convergence, highlighting the negative impact of the dilemma between them. In contrast, we observe that OD-LoRA, which satisfies both requirements simultaneously, achieves better loss convergence and gradient approximation quality, resulting in improved adaptation performance.

4.6 COMPARISON WITH EXISTING GRADIENT-BASED METHODS

We compare OD-LoRA with existing methods that view LoRA from the low-rank gradient approximation perspective, including LoRA-GA (Wang et al., 2024), LoRA-Pro (Wang et al., 2025), ScaledAdamW (Zhang & Pilanci, 2024), and AltLoRA (Yu et al., 2025). We reproduce the results of these methods with the same experimental setup as detailed in Table 9. For a fair comparison, we implement LoRA-Pro without tracking the momentum of full fine-tuning gradients. Table 5 shows that OD-LoRA outperforms existing methods by a substantial margin. Particularly, we observe that adjusting gradients for LoRA matrices (LoRA-Pro, ScaledAdamW, AltLoRA) typically results in substantial performance degradation. As discussed in Section D, we argue that although these methods minimize the gap between the full fine-tuning gradient and its low-rank approximation, modifying the original gradients for LoRA matrices results in ineffective subspace learning, degrading the gradient approximation quality. Indeed, we observe that these methods exhibit worse loss convergence (Figure 3) and lower alignment between the full gradient and its low-rank approximation (Figure 4b). These results further demonstrate that OD-LoRA effectively improves gradient approximation quality, which explains its performance gain.

4.7 TRAINING OVERHEAD

To demonstrate the efficiency of OD-LoRA, we present the time and memory overhead incurred by OD-LoRA in Table 6. The results show that OD-LoRA introduces approximately 5% overhead in time and less than 1% overhead in GPU-memory. As noted previously, LoRA-Pro, which reduces the gradient approximation error via gradient adjustment, incurs approximately 10% time and 24% memory overhead. Compared to LoRA-Pro, our method achieves substantially higher efficiency, particularly in memory consumption. These results further highlight the effectiveness of OD-LoRA.

5 CONCLUSION

In this paper, we demonstrate that the low-rank adaptation (LoRA) faces a dilemma between the representational capacity of weight updates and accurate gradient approximation. The proposed method, OD-LoRA, overcomes this dilemma by decoupling the singular values of weight updates from the equivalent low-rank gradient, achieving the two requirements at the same time. The experimental results on various benchmarks and pretrained models demonstrate that OD-LoRA significantly improves the gradient approximation quality, leading to better loss convergence and adaptation performance than existing methods.

REFERENCES

- 486
487
488 Paul Albert, Frederic Z. Zhang, Hemanth Saratchandran, Cristian Rodriguez-Opazo, Anton van den
489 Hengel, and Ehsan Abbasnejad. Randlora: Full-rank parameter-efficient fine-tuning of large
490 models. In *ICLR*, 2025.
- 491
492 Yonatan Bisk, Rowan Zellers, Ronan Le bras, Jianfeng Gao, and Yejin Choi. Piqa: Reasoning about
493 physical commonsense in natural language. *AAAI*, 2020.
- 494
495 Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101 – mining discriminative compo-
496 nents with random forests. In *ECCV*, 2014.
- 497
498 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal,
499 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel
500 Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler,
501 Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott
502 Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya
503 Sutskever, and Dario Amodei. Language models are few-shot learners, 2020. URL <https://arxiv.org/abs/2005.14165>.
- 504
505 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared
506 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri,
507 Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan,
508 et al. Evaluating large language models trained on code. 2021.
- 509
510 Mircea Cimpoi, Subhansu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. Describ-
511 ing textures in the wild. In *2014 IEEE Conference on Computer Vision and Pattern Recognition*,
512 2014.
- 513
514 Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina
515 Toutanova. BoolQ: Exploring the surprising difficulty of natural yes/no questions. In *NAACL-HLT*,
516 2019.
- 517
518 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and
519 Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge,
520 2018. URL <https://arxiv.org/abs/1803.05457>.
- 521
522 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
523 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John
524 Schulman. Training verifiers to solve math word problems, 2021. URL <https://arxiv.org/abs/2110.14168>.
- 525
526 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep
527 bidirectional transformers for language understanding. In *NAACL-HLT*, 2019.
- 528
529 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
530 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit,
531 and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale.
532 In *ICLR*, 2021.
- 533
534 Mary Gooneratne, Khe Chai Sim, Petr Zdravzil, Andreas Kabel, Françoise Beaufays, and Giovanni
535 Motta. Low-rank gradient approximation for memory-efficient on-device training of deep neural
536 network. In *ICASSP*, 2020.
- 537
538 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad
539 Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela
540 Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem
541 Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, et al. The llama 3 herd of
542 models, 2024. URL <https://arxiv.org/abs/2407.21783>.
- 543
544 Demi Guo, Alexander Rush, and Yoon Kim. Parameter-efficient transfer learning with diff pruning.
545 In *ACL-IJCNLP*, 2021.

- 540 Soufiane Hayou, Nikhil Ghosh, and Bin Yu. The impact of initialization on lora finetuning dynamics.
541 In *NeurIPS*, 2024a.
- 542
- 543 Soufiane Hayou, Nikhil Ghosh, and Bin Yu. Lora+: Efficient low rank adaptation of large models. In
544 *ICML*, 2024b.
- 545 Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. Towards a
546 unified view of parameter-efficient transfer learning. In *ICLR*, 2022.
- 547
- 548 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,
549 and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. In *NeurIPS*,
550 2021.
- 551 K. Hornik, M. Stinchcombe, and H. White. Multilayer feedforward networks are universal approxi-
552 mators. *Neural networks*, 2(5):359–366, 1989.
- 553
- 554 Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin de Laroussilhe, Andrea
555 Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for nlp. In
556 *ICML*, 2019.
- 557 Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
558 and Weizhu Chen. Lora: Low-rank adaptation of large language models. In *ICLR*, 2022.
- 559
- 560 Zhiqiang Hu, Yihuai Lan, Lei Wang, Wanyu Xu, Ee-Peng Lim, Roy Ka-Wei Lee, Lidong Bing,
561 and Soujanya Poria. Llm-adapters: An adapter family for parameter-efficient fine-tuning of large
562 language models. *arXiv preprint arXiv:2304.01933*, 2023.
- 563 Qiushi Huang, Tom Ko, Zhan Zhuang, Lilian Tang, and Yu Zhang. Hira: Parameter-efficient
564 hadamard high-rank adaptation for large language models. In *ICLR*, 2025.
- 565
- 566 Damjan Kalajdzievski. A rank stabilization scaling factor for fine-tuning with lora, 2023. URL
567 <https://arxiv.org/abs/2312.03732>.
- 568
- 569 Wenjun Ke, Jiahao Wang, Peng Wang, Jiajun Liu, Dong Nie, Guozheng Li, and Yining Li. Unveiling
570 lora intrinsic ranks via salience analysis. In *NeurIPS*, 2024.
- 571
- 572 Soroush Abbasi Koohpayegani, KL Navaneet, Parsa Nooralinejad, Soheil Kolouri, and Hamed
573 Pirsiavash. Nola: Compressing lora using linear combination of random basis. In *ICLR*, 2024.
- 574
- 575 Dawid J. Kopiczko, Tijmen Blankevoort, and Yuki M. Asano. Vera: Vector-based random matrix
576 adaptation. In *ICLR*, 2024.
- 577
- 578 Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained
579 categorization. In *2013 IEEE International Conference on Computer Vision Workshops*, 2013.
- 580
- 581 Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt
582 tuning. In *EMNLP*, 2021.
- 583
- 584 Shiwei Li, Xiandi Luo, Xing Tang, Haozhao Wang, Hao Chen, Weihong Luo, Yuhua Li, Xiuqiang
585 He, and Ruixuan Li. Beyond zero initialization: Investigating the impact of non-zero initialization
586 on lora fine-tuning dynamics. In *ICML*, 2025.
- 587
- 588 Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. In
589 *ACL-IJCNLP*, 2021.
- 590
- 591 Vladislav Lialin, Namrata Shivagunde, Sherin Muckatira, and Anna Rumshisky. Relora: High-rank
592 training through low-rank updates. In *ICLR*, 2024.
- 593
- 594 Vijay Lingam, Atula Tejaswi, Aditya Vavre, Aneesh Shetty, Gautham Krishna Gudur, Joydeep Ghosh,
595 Alex Dimakis, Eunsol Choi, Aleksandar Bojchevski, and Sujay Sanghavi. Svft: Parameter-efficient
596 fine-tuning with singular vectors. In *NeurIPS*, 2024.
- 597
- 598 Shih-Yang Liu, Chien-Yi Wang, Hongxu Yin, Pavlo Molchanov, Yu-Chiang Frank Wang, Kwang-
599 Ting Cheng, and Min-Hung Chen. Dora: Weight-decomposed low-rank adaptation. In *ICML*,
600 2024.

- 594 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *ICLR*, 2019.
595
- 596 Fanxu Meng, Zhaohui Wang, and Muhan Zhang. Pissa: Principal singular values and singular vectors
597 adaptation of large language models. In *NeurIPS*, 2024.
- 598 Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct
599 electricity? a new dataset for open book question answering. In *EMNLP*, 2018.
600
- 601 Yurii Nesterov. *Lectures on Convex Optimization*. Springer, Cham, Switzerland, 2018. ISBN
602 978-3-319-91578-8.
- 603 Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language under-
604 standing by generative pre-training. *Technical report, OpenAI*, 2018.
605
- 606 Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An
607 adversarial winograd schema challenge at scale. *AAAI*, 2020.
- 608 Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. Social IQa: Common-
609 sense reasoning about social interactions. In *EMNLP-IJCNLP*, 2019.
610
- 611 Chongjie Si, Zhiyi Shi, Shifan Zhang, Xiaokang Yang, Hanspeter Pfister, and Wei Shen. Unleashing
612 the power of task-specific directions in parameter efficient fine-tuning. In *ICLR*, 2025.
- 613 Yi-Lin Sung, Varun Nair, and Colin Raffel. Training neural networks with fixed sparse masks. In
614 *NeurIPS*, 2021.
- 615 Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak,
616 Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, Pouya Tafti, Léonard Hussonot,
617 Pier Giuseppe Sessa, Aakanksha Chowdhery, Adam Roberts, Aditya Barua, Alex Botev, Alex
618 Castro-Ros, Ambrose Slone, Amélie Héliou, Andrea Tacchetti, Anna Bulanova, Antonia Paterson,
619 Beth Tsai, Bobak Shahriari, Charline Le Lan, Christopher A. Choquette-Choo, Clément Crepy,
620 Daniel Cer, et al. Gemma: Open models based on gemini research and technology, 2024. URL
621 <https://arxiv.org/abs/2403.08295>.
- 622 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
623 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand
624 Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language
625 models, 2023. URL <https://arxiv.org/abs/2302.13971>.
- 626 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz
627 Kaiser, and Illia Polosukhin. Attention is all you need. In *NIPS*, 2017.
- 628 C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. The caltech-ucsd birds-200-2011 dataset.
629 Technical report, California Institute of Technology, 2011.
630
- 631 Shaowen Wang, Linxi Yu, and Jian Li. Lora-ga: Low-rank adaptation with gradient approximation.
632 In *NeurIPS*, 2024.
633
- 634 Zhengbo Wang, Jian Liang, Ran He, Zilei Wang, and Tieniu Tan. Lora-pro: Are low-rank adapters
635 properly optimized? In *ICLR*, 2025.
636
- 637 Jianxiong Xiao, James Hays, Krista A. Ehinger, Aude Oliva, and Antonio Torralba. Sun database:
638 Large-scale scene recognition from abbey to zoo. In *2010 IEEE Computer Society Conference on
639 Computer Vision and Pattern Recognition*, 2010.
- 640 Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T. Kwok,
641 Zhenguo Li, Adrian Weller, and Weiyang Liu. Metamath: Bootstrap your own mathematical
642 questions for large language models. In *ICLR*, 2024.
- 643 Xin Yu, Yujia Wang, Jinghui Chen, and Lingzhou Xue. Atllora: Towards better gradient approxima-
644 tion in low-rank adaptation with alternating projections. In *NeurIPS*, 2025.
645
- 646 Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. HellaSwag: Can a machine
647 really finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association for
Computational Linguistics*, 2019.

648 Fangzhao Zhang and Mert Pilanci. Riemannian preconditioned lora for fine-tuning foundation models.
649 In *ICML*, 2024.
650

651 Qingru Zhang, Minshuo Chen, Alexander Bukharin, Nikos Karampatziakis, Pengcheng He, Yu Cheng,
652 Weizhu Chen, and Tuo Zhao. Adalora: Adaptive budget allocation for parameter-efficient fine-
653 tuning. In *ICLR*, 2023.

654 Yuanhe Zhang, Fanghui Liu, and Yudong Chen. Lora-one: One-step full gradient could suffice for
655 fine-tuning large language models, provably and efficiently. In *ICML*, 2025.
656

657 Jiawei Zhao, Zhenyu Zhang, Beidi Chen, Zhangyang Wang, Anima Anandkumar, and Yuandong
658 Tian. Galore: Memory-efficient llm training by gradient low-rank projection. In *ICML*, 2024.

659 Tianyu Zheng, Ge Zhang, Tianhao Shen, Xueling Liu, Bill Yuchen Lin, Jie Fu, Wenhui Chen, and
660 Xiang Yue. Opencodeinterpreter: Integrating code generation with execution and refinement. *arXiv*
661 *preprint arXiv:2402.14658*, 2024.
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

APPENDIX

The contents of the Appendix are as follows:

- In Section A, B, and C, we provide the proof of Theorem 1, 2, and 3, respectively.
- In Section D, we provide theoretical and empirical comparisons with related works.
- In Section E, we discuss the importance of improving gradient approximation quality.
- In Section F, we show the singular values spectrum learned via full fine-tuning and LoRA.
- In Section G, we provide the details and experimental results on image classification tasks.
- In Section H, we provide additional ablation studies.
- In Section I, we provide the details on the ablation studies in Section 4.5
- In Table 9, we provide the additional implementation details.

A PROOF OF THEOREM 1

Proof. Let $\langle \cdot, \cdot \rangle_F$ and $\|\cdot\|_F$ denote the Frobenius inner product and Frobenius norm, respectively. Since $\langle \mathcal{G}, \tilde{\mathcal{G}} \rangle = \sum_{l=1}^L \langle \mathcal{G}^l, \tilde{\mathcal{G}}^l \rangle_F$, it is enough to prove that $\langle \mathcal{G}^l, \tilde{\mathcal{G}}^l \rangle_F$ is maximized if \mathbf{A}^l and \mathbf{B}^l have uniform singular values. For simplicity, we omit the subscript l . Using the property of the inner product, we rewrite $\langle \mathcal{G}, \tilde{\mathcal{G}} \rangle_F$ as follows:

$$\begin{aligned} \langle \mathcal{G}, \tilde{\mathcal{G}} \rangle_F &= \langle \mathcal{G}, s^2(\mathbf{G}\mathbf{A}^\top\mathbf{A} + \mathbf{B}\mathbf{B}^\top\mathbf{G}) \rangle_F \\ &= s^2(\langle \mathcal{G}, \mathbf{G}\mathbf{A}^\top\mathbf{A} \rangle_F + \langle \mathcal{G}, \mathbf{B}\mathbf{B}^\top\mathbf{G} \rangle_F). \end{aligned} \quad (6)$$

Since the row space of \mathbf{A} and the column space of \mathbf{B} are given, we can define unique orthogonal projections onto those subspaces: $\mathbf{P}_{\mathcal{R}(\mathbf{A})} = \mathbf{A}^\top(\mathbf{A}\mathbf{A}^\top)^{-1}\mathbf{A}$ and $\mathbf{P}_{\mathcal{C}(\mathbf{B})} = \mathbf{B}(\mathbf{B}^\top\mathbf{B})^{-1}\mathbf{B}^\top$. This allows us to rewrite \mathcal{G} as $\mathcal{G} = \mathbf{G}\mathbf{P}_{\mathcal{R}(\mathbf{A})} + \mathcal{G}(\mathbf{I} - \mathbf{P}_{\mathcal{R}(\mathbf{A})})$ or $\mathcal{G} = \mathbf{P}_{\mathcal{C}(\mathbf{B})}\mathcal{G} + (\mathbf{I} - \mathbf{P}_{\mathcal{C}(\mathbf{B})})\mathcal{G}$. Thus, we can rewrite each term in the numerator of the last equation of Equation 6 as

$$\begin{aligned} \langle \mathcal{G}, \mathbf{G}\mathbf{A}^\top\mathbf{A} \rangle_F &= \langle \mathbf{G}\mathbf{P}_{\mathcal{R}(\mathbf{A})}, \mathbf{G}\mathbf{A}^\top\mathbf{A} \rangle_F + \langle \mathcal{G}(\mathbf{I} - \mathbf{P}_{\mathcal{R}(\mathbf{A})}), \mathbf{G}\mathbf{A}^\top\mathbf{A} \rangle_F, \\ \langle \mathcal{G}, \mathbf{B}\mathbf{B}^\top\mathbf{G} \rangle_F &= \langle \mathbf{P}_{\mathcal{C}(\mathbf{B})}\mathcal{G}, \mathbf{B}\mathbf{B}^\top\mathbf{G} \rangle_F + \langle (\mathbf{I} - \mathbf{P}_{\mathcal{C}(\mathbf{B})})\mathcal{G}, \mathbf{B}\mathbf{B}^\top\mathbf{G} \rangle_F. \end{aligned} \quad (7)$$

The second terms in Equation 7 become zero since

$$\begin{aligned} \langle \mathcal{G}(\mathbf{I} - \mathbf{P}_{\mathcal{R}(\mathbf{A})}), \mathbf{G}\mathbf{A}^\top\mathbf{A} \rangle_F &= \text{Tr}(\mathcal{G}(\mathbf{I} - \mathbf{P}_{\mathcal{R}(\mathbf{A})})\mathbf{A}^\top\mathbf{A}\mathcal{G}^\top) \\ &= \text{Tr}(\mathcal{G}(\mathbf{A}^\top\mathbf{A} - \mathbf{A}^\top\mathbf{A})\mathcal{G}^\top), \\ \langle (\mathbf{I} - \mathbf{P}_{\mathcal{C}(\mathbf{B})})\mathcal{G}, \mathbf{B}\mathbf{B}^\top\mathbf{G} \rangle_F &= \text{Tr}(\mathcal{G}^\top\mathbf{B}\mathbf{B}^\top(\mathbf{I} - \mathbf{P}_{\mathcal{C}(\mathbf{B})})\mathcal{G}) \\ &= \text{Tr}(\mathcal{G}^\top(\mathbf{B}\mathbf{B}^\top - \mathbf{B}\mathbf{B}^\top)\mathcal{G}). \end{aligned}$$

Thus, we rewrite $\langle \mathcal{G}, \tilde{\mathcal{G}} \rangle_F$ as follows:

$$\langle \mathcal{G}, \tilde{\mathcal{G}} \rangle_F = s^2(\langle \mathbf{G}\mathbf{P}_{\mathcal{R}(\mathbf{A})}, \mathbf{G}\mathbf{A}^\top\mathbf{A} \rangle_F + \langle \mathbf{P}_{\mathcal{C}(\mathbf{B})}\mathcal{G}, \mathbf{B}\mathbf{B}^\top\mathbf{G} \rangle_F).$$

By the Cauchy-Schwarz inequality, we obtain

$$\langle \mathcal{G}, \tilde{\mathcal{G}} \rangle_F \leq s^2(\|\mathbf{G}\mathbf{P}_{\mathcal{R}(\mathbf{A})}\|_F \cdot \|\mathbf{G}\mathbf{A}^\top\mathbf{A}\|_F + \|\mathbf{P}_{\mathcal{C}(\mathbf{B})}\mathcal{G}\|_F \cdot \|\mathbf{B}\mathbf{B}^\top\mathbf{G}\|_F),$$

where the equality holds when $\mathbf{G}\mathbf{A}^\top\mathbf{A} = c\mathbf{G}\mathbf{P}_{\mathcal{R}(\mathbf{A})}$ and $\mathbf{B}\mathbf{B}^\top\mathbf{G} = d\mathbf{P}_{\mathcal{C}(\mathbf{B})}\mathcal{G}$ for some $c > 0$ and $d > 0$. The sufficient condition for arbitrary \mathcal{G} is given by: $\mathbf{A}^\top\mathbf{A} = c\mathbf{P}_{\mathcal{R}(\mathbf{A})}$ and $\mathbf{B}\mathbf{B}^\top = d\mathbf{P}_{\mathcal{C}(\mathbf{B})}$. This leads to

$$\langle \mathcal{G}, \tilde{\mathcal{G}} \rangle_F \leq s^2(c\|\mathbf{G}\mathbf{P}_{\mathcal{R}(\mathbf{A})}\|_F^2 + d\|\mathbf{P}_{\mathcal{C}(\mathbf{B})}\mathcal{G}\|_F^2). \quad (8)$$

From the fixed energy assumption, we set $\|\mathbf{A}\|_F^2 = E_A$ and $\|\mathbf{B}\|_F^2 = E_B$. Since \mathbf{A} and \mathbf{B} have rank r , $\text{Tr}(\mathbf{P}_{\mathcal{R}(\mathbf{A})}) = \text{Tr}(\mathbf{P}_{\mathcal{C}(\mathbf{B})}) = r$. This gives $\|\mathbf{A}\|_F^2 = \text{Tr}(\mathbf{A}^\top\mathbf{A}) = \text{Tr}(c\mathbf{P}_{\mathcal{R}(\mathbf{A})}) = cr$ and $\|\mathbf{B}\|_F^2 = \text{Tr}(\mathbf{B}\mathbf{B}^\top) = \text{Tr}(d\mathbf{P}_{\mathcal{C}(\mathbf{B})}) = dr$, leading to $c = \frac{E_A}{r}$ and $d = \frac{E_B}{r}$. Therefore, we can rewrite Equation 8 as

$$\langle \mathcal{G}, \tilde{\mathcal{G}} \rangle_F \leq \frac{s^2}{r}(E_A\|\mathbf{G}\mathbf{P}_{\mathcal{R}(\mathbf{A})}\|_F^2 + E_B\|\mathbf{P}_{\mathcal{C}(\mathbf{B})}\mathcal{G}\|_F^2). \quad (9)$$

Since all the terms except for G in Equation 9 are fixed, we can conclude that $\langle G, \tilde{G} \rangle_F$ is maximized for arbitrary G if $A^\top A = \frac{E_A}{r} P_{\mathcal{R}(A)}$ and $BB^\top = \frac{E_B}{r} P_{\mathcal{C}(B)}$. Let $A = U_A \Sigma_A V_A^\top$ and $B = U_B \Sigma_B V_B^\top$ be the compact singular value decomposition of A and B , respectively. By substituting A and B with their SVD, we obtain $A^\top A = V_A \Sigma_A^2 V_A^\top$, $BB^\top = U_B \Sigma_B^2 U_B^\top$, $P_{\mathcal{R}(A)} = V_A V_A^\top$, and $P_{\mathcal{C}(B)} = U_B U_B^\top$. From $A^\top A = \frac{E_A}{r} P_{\mathcal{R}(A)}$ and $BB^\top = \frac{E_B}{r} P_{\mathcal{C}(B)}$, we obtain

$$\begin{aligned}\Sigma_A^2 &= \frac{E_A}{r} I_r, \\ \Sigma_B^2 &= \frac{E_B}{r} I_r,\end{aligned}$$

where I_r is the $r \times r$ identity matrix. This implies that A and B have uniform singular values. \square

B PROOF OF THEOREM 2

Let $X \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$ be an arbitrary rank- r matrix with non-zero singular values $\{\sigma_1, \sigma_2, \dots, \sigma_r\}$ and $\Delta W_{\text{LoRA}} = sBA$ where both A and B are rank- r . We prove the propositions in Theorem 2.

Proof of "Without any constraints, $\text{Cap}_X(\Delta W_{\text{LoRA}}) = \infty$."

Proof. Let $X = U_r \Sigma_r V_r^\top$ be the rank- r truncated SVD of X . Then, we can always choose A and B such that $A = \Sigma_r V_r^\top$ and $B = \frac{1}{s} U_r$, resulting in $\min_{\Delta W_{\text{LoRA}}} \|X - \Delta W_{\text{LoRA}}\|_F^2 = 0$. Thus, $\text{Cap}_X(\Delta W_{\text{LoRA}}) = \frac{1}{\min_{\Delta W_{\text{LoRA}}} \|\Delta W_{\text{LoRA}}\|_F^2} = \infty$. \square

Proof of "If A and B have uniform singular values, the non-zero singular values of ΔW_{LoRA} are uniform."

Proof. We will examine the eigenvalues of $\Delta W_{\text{LoRA}}^\top \Delta W_{\text{LoRA}}$ since they are identical to the singular values of ΔW_{LoRA} . The compact SVD of B gives us $B = U_B \Sigma_B V_B^\top = \sigma_B U_B V_B^\top$. Then, we obtain $B^\top B = \sigma_B^2 V_B U_B^\top U_B V_B^\top = \sigma_B^2 I_r$ where I_r denotes the $r \times r$ identity matrix. This gives us $\Delta W_{\text{LoRA}}^\top \Delta W_{\text{LoRA}} = s^2 A^\top B^\top B A = s^2 \sigma_B^2 A^\top A$, implying that the eigenvalues of $\Delta W_{\text{LoRA}}^\top \Delta W_{\text{LoRA}}$ are identical to those of $s^2 \sigma_B^2 A^\top A$. Since A has uniform singular values, i.e., σ_A , the non-zero eigenvalues of $A^\top A$ are σ_A^2 , which implies that the non-zero eigenvalues of $\Delta W_{\text{LoRA}}^\top \Delta W_{\text{LoRA}}$ are $s^2 \sigma_A^2 \sigma_B^2$. Therefore, the non-zero singular values of ΔW_{LoRA} are uniform, with value $s \sigma_A \sigma_B$. \square

Proof of "If the non-zero singular values of ΔW_{LoRA} are uniform, $\text{Cap}_X(\Delta W_{\text{LoRA}}) = 1 / \sum_{i=1}^r (\sigma_i - \frac{1}{r} \sum_{i=1}^r \sigma_i)^2$."

Proof. Let W denote ΔW_{LoRA} for simplicity. Recall $\text{Cap}_X(W) = 1 / \min_W \|X - W\|_F^2$. We can rewrite the Frobenius term in the denominator as

$$\begin{aligned}\|X - W\|_F^2 &= \text{Tr}((X - W)^\top (X - W)) \\ &= \text{Tr}(X^\top X) - \text{Tr}(X^\top W) - \text{Tr}(W^\top X) + \text{Tr}(W^\top W) \\ &= \text{Tr}(X^\top X) + \text{Tr}(W^\top W) - 2 \text{Tr}(X^\top W).\end{aligned}$$

Let σ_W be the non-zero singular value of W . Then, we can evaluate the first and the second term: $\text{Tr}(X^\top X) = \|X\|_F^2 = \sum_{i=1}^r \sigma_i^2$ and $\text{Tr}(W^\top W) = \|W\|_F^2 = \sum_{i=1}^r \sigma_W^2 = r \sigma_W^2$. By the Von Neumann's Trace inequality, we obtain the following inequality for the third term:

$$\text{Tr}(X^\top W) \leq \sum_{i=1}^r \sigma_i \cdot \sigma_W,$$

where the equality holds if X and W share the same singular vectors. Thus, we obtain

$$\|X - W\|_F^2 \geq \sum_{i=1}^r \sigma_i^2 + r \sigma_W^2 - 2 \sigma_W \sum_{i=1}^r \sigma_i.$$

As shown in this section, \mathbf{W} can represent any rank- r matrix, implying that there exists \mathbf{W} that shares the same singular vectors with \mathbf{X} . Thus, the minimization of $\|\mathbf{X} - \mathbf{W}\|_F^2$ is converted to the following optimization problem:

$$\min_{\sigma_{\mathbf{W}}} \left(\sum_{i=1}^r \sigma_i^2 + r\sigma_{\mathbf{W}}^2 - 2\sigma_{\mathbf{W}} \sum_{i=1}^r \sigma_i \right),$$

which is a quadratic function of $\sigma_{\mathbf{W}}$. Define $f(\sigma_{\mathbf{W}}) = \sum_{i=1}^r \sigma_i^2 + r\sigma_{\mathbf{W}}^2 - 2\sigma_{\mathbf{W}} \sum_{i=1}^r \sigma_i$. By solving $\frac{df(\sigma_{\mathbf{W}})}{d\sigma_{\mathbf{W}}} = 0$, we obtain $\sigma_{\mathbf{W}} = \frac{1}{r} \sum_{i=1}^r \sigma_i = \bar{\sigma}$. Finally, this leads to $\text{Cap}_{\mathbf{X}}(\mathbf{W}) = 1/\min_{\mathbf{W}} \|\mathbf{X} - \mathbf{W}\|_F^2 = 1/f(\bar{\sigma}) = 1/\sum_{i=1}^r (\sigma_i - \bar{\sigma})^2$. \square

Proof of "Cap $_{\mathbf{X}}(\Delta\mathbf{W}_{\text{LoRA}}) = \infty$ requires either \mathbf{A} or \mathbf{B} to have non-uniform singular values."

Proof. Consider the contrapositive: "If both \mathbf{A} and \mathbf{B} have uniform singular values, $\text{Cap}_{\mathbf{X}}(\Delta\mathbf{W}_{\text{LoRA}}) < \infty$ ". As shown in this section, if \mathbf{A} and \mathbf{B} have uniform singular values, the non-zero singular values of $\Delta\mathbf{W}_{\text{LoRA}}$ are uniform, and if the non-zero singular values of $\Delta\mathbf{W}_{\text{LoRA}}$ are uniform, $\text{Cap}_{\mathbf{X}}(\Delta\mathbf{W}_{\text{LoRA}}) = 1/\sum_{i=1}^r (\sigma_i - \frac{1}{r} \sum_{i=1}^r \sigma_i)^2$. Therefore, the contrapositive holds. \square

C PROOF OF THEOREM 3

C.1 PROOF OF THE FIRST PROPERTY

Proof. Since the gradient does not flow through \mathbf{Q}_A and \mathbf{Q}_B , the gradient with respect to \mathbf{A} and \mathbf{B} is given by $\frac{\partial \mathcal{L}}{\partial \mathbf{A}} = s\mathbf{Q}_B^\top \mathbf{G}$ and $\frac{\partial \mathcal{L}}{\partial \mathbf{B}} = s\mathbf{G}\mathbf{Q}_A$. Then, the differential of \mathbf{A} and \mathbf{B} is expressed as

$$\begin{aligned} d\mathbf{A} &= -\eta \frac{\partial \mathcal{L}}{\partial \mathbf{A}} = -\eta s \mathbf{Q}_B^\top \mathbf{G}, \\ d\mathbf{B} &= -\eta \frac{\partial \mathcal{L}}{\partial \mathbf{B}} = -\eta s \mathbf{G}\mathbf{Q}_A. \end{aligned}$$

Thus, we obtain

$$\begin{aligned} d\mathbf{W}' &= s(\mathbf{Q}_B d\mathbf{A} + d\mathbf{B} \mathbf{Q}_A^\top) \\ &= -\eta s^2 (\mathbf{G}\mathbf{Q}_A \mathbf{Q}_A^\top + \mathbf{Q}_B \mathbf{Q}_B^\top \mathbf{G}) \\ \tilde{\mathbf{G}} &= s^2 (\mathbf{G}\mathbf{Q}_A \mathbf{Q}_A^\top + \mathbf{Q}_B \mathbf{Q}_B^\top \mathbf{G}). \end{aligned}$$

\square

C.2 PROOF OF THE SECOND PROPERTY

Proof. We have $\mathbf{A}^\top = \mathbf{Q}_A \mathbf{R}_A$ and $\mathbf{B} = \mathbf{Q}_B \mathbf{R}_B$. Then, we can rewrite $\mathbf{Q}_B \mathbf{A} + \mathbf{B} \mathbf{Q}_A^\top$ as follows:

$$\begin{aligned} \mathbf{Q}_B \mathbf{A} + \mathbf{B} \mathbf{Q}_A^\top &= \mathbf{Q}_B \mathbf{R}_A^\top \mathbf{Q}_A^\top + \mathbf{Q}_B \mathbf{R}_B \mathbf{Q}_A^\top \\ &= \mathbf{Q}_B (\mathbf{R}_A^\top + \mathbf{R}_B) \mathbf{Q}_A^\top. \end{aligned}$$

For an arbitrary rank- r matrix $\mathbf{X} \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$, let $\mathbf{X} = \mathbf{U}_r \Sigma_r \mathbf{V}_r^\top$ be the rank- r truncated SVD of \mathbf{X} . Then, it is enough to prove that there exist \mathbf{Q}_A , \mathbf{Q}_B , \mathbf{R}_A , and \mathbf{R}_B such that

$$\mathbf{Q}_B (\mathbf{R}_A^\top + \mathbf{R}_B) \mathbf{Q}_A^\top = \mathbf{U}_r \Sigma_r \mathbf{V}_r^\top. \quad (10)$$

By the property of QR decomposition, both \mathbf{Q}_A and \mathbf{Q}_B are matrices with orthonormal columns and both \mathbf{R}_A and \mathbf{R}_B are $r \times r$ invertible upper-triangle matrices. Thus, we can set $\mathbf{Q}_A = \mathbf{V}_r$, $\mathbf{Q}_B = \mathbf{U}_r$, $\mathbf{R}_A = k \Sigma_r$, and $\mathbf{R}_B = (1 - k) \Sigma_r$ for some $0 < k < 1$, achieving the equality in Equation 10. Therefore, $\Delta\mathbf{W}_{\text{OD-LoRA}} = s(\mathbf{Q}_B \mathbf{A} + \mathbf{B} \mathbf{Q}_A^\top)$ can represent any rank- r matrix \mathbf{X} , leading to $\text{Cap}_{\mathbf{X}}(\Delta\mathbf{W}_{\text{OD-LoRA}}) = \infty$. \square

D COMPARISONS WITH RELATED WORKS

In this section, we provide a comprehensive comparison with existing methods that view LoRA from the perspective of low-rank gradient approximation (Zhang & Pilanci, 2024; Wang et al., 2025; Yu et al., 2025). While these methods share our goal of improving the alignment between the full fine-tuning gradient and its low-rank approximation, they rely on preconditioning the gradients of the LoRA matrices. In Section D.1, we theoretically analyze the difference between such preconditioning-based methods and OD-LoRA. Our analysis focuses on two key aspects: the effectiveness in subspace learning (Section D.1.1) and the resulting convergence rates under standard optimization assumptions (Section D.1.2). Finally, in Section D.2, we present empirical comparisons regarding loss convergence and gradient approximation quality to validate our theoretical findings.

D.1 THEORETICAL ANALYSIS

In this subsection, we provide theoretical analyses of OD-LoRA and preconditioning-based methods under the SGD training setting. We adopt ScaledGD (Zhang & Pilanci, 2024) as a representative preconditioning-based method. Note that our analysis can be applied to the other methods (Wang et al., 2025; Yu et al., 2025) since they have similar update rules. For notational brevity, we present the analysis using single-layer notation because the generalization to the multi-layer setting is straightforward due to the additivity of the Frobenius norm and inner product. Let $\theta = (\mathbf{A}, \mathbf{B})$ and $\mathbf{W}_{\text{eff}} = \mathbf{W}_0 + \Delta\mathbf{W}(\mathbf{A}, \mathbf{B})$ denote the LoRA parameters and the effective weight, respectively. For simplicity, we ignore the scaling factor s in our analysis. Also, let $\mathbf{G} = \nabla_{\mathbf{W}_{\text{eff}}}\mathcal{L}$ denote the full fine-tuning gradient for \mathbf{W}_{eff} . Let $\mathbf{P}_{\mathcal{R}(\mathbf{A})}$ and $\mathbf{P}_{\mathcal{C}(\mathbf{B})}$ denote the projection matrix onto the row space of \mathbf{A} and the column space of \mathbf{B} , respectively. For a learning rate η , the update rule for OD-LoRA is given by

$$\begin{aligned}\Delta\mathbf{A} &= -\eta\mathbf{Q}_B\mathbf{G} \\ \Delta\mathbf{B} &= -\eta\mathbf{G}\mathbf{Q}_A \\ \Delta\mathbf{W}_{\text{eff}} &\approx -\eta(\mathbf{P}_{\mathcal{C}(\mathbf{B})}\mathbf{G} + \mathbf{G}\mathbf{P}_{\mathcal{R}(\mathbf{A})}).\end{aligned}\tag{11}$$

The update rule for ScaledGD is given by:

$$\begin{aligned}\Delta\mathbf{A} &= -\eta(\mathbf{B}^\top\mathbf{B})^{-1}\mathbf{B}^\top\mathbf{G} \\ \Delta\mathbf{B} &= -\eta\mathbf{G}\mathbf{A}^\top(\mathbf{A}\mathbf{A}^\top)^{-1} \\ \Delta\mathbf{W}_{\text{eff}} &\approx -\eta(\mathbf{P}_{\mathcal{C}(\mathbf{B})}\mathbf{G} + \mathbf{G}\mathbf{P}_{\mathcal{R}(\mathbf{A})}).\end{aligned}\tag{12}$$

Note that the update rule for \mathbf{W}_{eff} is identical in both methods. However, we demonstrate that the different update rules for \mathbf{A} and \mathbf{B} lead to different abilities to learn desirable subspaces (Section D.1.1), which in turn result in different convergence rates (Section D.1.2).

D.1.1 SIGNAL-TO-NOISE RATIO ANALYSIS OF SUBSPACE LEARNING

We first demonstrate that OD-LoRA is effective in learning desirable subspaces for \mathbf{A} and \mathbf{B} . In the context of LoRA, desirable subspaces have higher alignment with the full gradient, leading to a more accurate update of $\Delta\mathbf{W}_{\text{eff}}$, as indicated in Equation 11 and Equation 12. Note that although the following analyses in this subsection focus on the gradient for \mathbf{A} , they are equally applicable to \mathbf{B} . Let $\mathbf{B} = \mathbf{U}\Sigma\mathbf{V}^\top$ be the compact singular value decomposition for \mathbf{B} where $\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r]$, $\mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r]^\top$, and Σ is a diagonal matrix with singular values $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$. To update \mathbf{A} , \mathbf{u}_i is used to evaluate how the column space of \mathbf{B} is aligned with the columns of \mathbf{G} . Then, the row space of \mathbf{A} is updated according to the alignment. To evaluate how effectively the alignment information is used to update \mathbf{A} , we introduce the following definitions.

Definition 2. (Update decomposition using alignment) Express the update for \mathbf{A} as a sum of components corresponding to the alignment between the left singular vectors of \mathbf{B} and the full gradient \mathbf{G} :

$$\Delta\mathbf{A} = \sum_{i=1}^r \mathbf{x}_i(\mathbf{u}_i^\top\mathbf{G}) = \sum_{i=1}^r \mathbf{x}_i\mathbf{h}_i^\top,\tag{13}$$

- **Signal update** ($\mathbf{x}_1\mathbf{h}_1^\top$): The component aligned with the strongest direction \mathbf{u}_1 .

- **Noise update** ($x_r \mathbf{h}_r^\top$): The component aligned with the weakest direction \mathbf{u}_r .

Definition 3. (Signal-to-noise ratio of update) We define the signal-to-noise ratio of the update for \mathbf{A} as

$$\text{SNR}(\Delta \mathbf{A}) := \frac{\|\mathbf{x}_1 \mathbf{h}_1^\top\|_F}{\|\mathbf{x}_r \mathbf{h}_r^\top\|_F}. \quad (14)$$

Definition 4. (Condition number) For a matrix \mathbf{X} , we define the condition number of \mathbf{X} as

$$\kappa(\mathbf{X}) := \frac{\sigma_{\max}(\mathbf{X})}{\sigma_{\min}(\mathbf{X})}, \quad (15)$$

where $\sigma_{\max}(\mathbf{X})$ and $\sigma_{\min}(\mathbf{X})$ represent the largest and smallest singular value of \mathbf{X} .

For example, for the standard LoRA formulation, we can rewrite the update of \mathbf{A} as $\Delta \mathbf{A} = -\eta \mathbf{B}^\top \mathbf{G} = \sum_{i=1}^r -\eta \sigma_i \mathbf{v}_i (\mathbf{u}_i^\top \mathbf{G}) = \sum_{i=1}^r -\eta \sigma_i \mathbf{v}_i \mathbf{h}_i^\top$. Then, the signal update is $-\eta \sigma_1 \mathbf{v}_1 \mathbf{h}_1^\top$ and the noise update is $-\eta \sigma_r \mathbf{v}_r \mathbf{h}_r^\top$. Consequently, the signal-to-noise ratio of the gradient for \mathbf{A} in LoRA is expressed as $\frac{\|-\eta \sigma_1 \mathbf{v}_1 \mathbf{h}_1^\top\|_F}{\|-\eta \sigma_r \mathbf{v}_r \mathbf{h}_r^\top\|_F} = \kappa(\mathbf{B}) \frac{\|\mathbf{h}_1\|}{\|\mathbf{h}_r\|}$. This implies that the alignment information is distorted by $\kappa(\mathbf{B})$, causing strongly aligned directions to be amplified more.

Using these definitions, we derive the following theorem on the distortion of the signal-to-noise ratio:

Theorem 4: Distortion of Signal-to-Noise Ratio

Let $\text{SNR}_{\text{intrinsic}} := \frac{\|\mathbf{h}_1\|}{\|\mathbf{h}_r\|}$ be the intrinsic signal-to-noise ratio of the gradient alignment. The Signal-to-Noise Ratio (SNR) of the parameter update $\Delta \mathbf{A}$ for OD-LoRA and ScaledGD is given by:

- **OD-LoRA:** $\text{SNR}(\Delta \mathbf{A}) = 1 \cdot \text{SNR}_{\text{intrinsic}}$.
- **Preconditioning-based Methods (e.g., ScaledGD):** $\text{SNR}(\Delta \mathbf{A}) = \frac{1}{\kappa(\mathbf{B})} \text{SNR}_{\text{intrinsic}}$.

This implies that existing methods amplify the noise component and suppress the signal component by the inverse of the singular values of \mathbf{B} , diluting the meaningful gradient signal.

Proof. We analyze the update rules for each method using the decomposition in Definition 2.

1. OD-LoRA

The update rule is given by $\Delta \mathbf{A} = -\eta \mathbf{Q}_B^\top \mathbf{G}$. Since \mathbf{Q}_B corresponds to the orthonormal basis \mathbf{U} of the column space $\mathcal{C}(\mathbf{B})$, its singular values are all equal to 1 ($\sigma_i = 1$ for all i). The update can be written as:

$$\Delta \mathbf{A} = \sum_{i=1}^r (-\eta \mathbf{v}_i) (\mathbf{u}_i^\top \mathbf{G}). \quad (16)$$

Here, the signal update is $-\eta \mathbf{v}_1 \mathbf{h}_1^\top$ and the noise update is $-\eta \mathbf{v}_r \mathbf{h}_r^\top$. Consequently, the SNR is calculated as:

$$\text{SNR}(\Delta \mathbf{A}) = \frac{\|-\eta \mathbf{v}_1 \mathbf{h}_1^\top\|_F}{\|-\eta \mathbf{v}_r \mathbf{h}_r^\top\|_F} = \frac{\eta \|\mathbf{h}_1\|}{\eta \|\mathbf{h}_r\|} = 1 \cdot \text{SNR}_{\text{intrinsic}}. \quad (17)$$

2. ScaledGD

The update rule is given by $\Delta \mathbf{A} = -\eta (\mathbf{B}^\top \mathbf{B})^{-1} \mathbf{B}^\top \mathbf{G}$. Using the SVD of $\mathbf{B} = \mathbf{U} \Sigma \mathbf{V}^\top$, the pseudoinverse term is $(\mathbf{B}^\top \mathbf{B})^{-1} \mathbf{B}^\top = \mathbf{V} \Sigma^{-1} \mathbf{U}^\top$. Substituting this into the update rule:

$$\Delta \mathbf{A} = -\eta \mathbf{V} \Sigma^{-1} \mathbf{U}^\top \mathbf{G} = \sum_{i=1}^r \left(-\eta \frac{1}{\sigma_i} \mathbf{v}_i \right) (\mathbf{u}_i^\top \mathbf{G}). \quad (18)$$

Here, the signal update is $-\frac{\eta}{\sigma_1} \mathbf{v}_1 \mathbf{h}_1^\top$, and the noise update is $-\frac{\eta}{\sigma_r} \mathbf{v}_r \mathbf{h}_r^\top$. Consequently, the SNR is calculated as:

$$\text{SNR}(\Delta \mathbf{A}) = \frac{\|-\frac{\eta}{\sigma_1} \mathbf{v}_1 \mathbf{h}_1^\top\|_F}{\|-\frac{\eta}{\sigma_r} \mathbf{v}_r \mathbf{h}_r^\top\|_F} = \frac{\frac{1}{\sigma_1} \|\mathbf{v}_1\| \|\mathbf{h}_1\|}{\frac{1}{\sigma_r} \|\mathbf{v}_r\| \|\mathbf{h}_r\|} = \frac{\sigma_r}{\sigma_1} \frac{\|\mathbf{h}_1\|}{\|\mathbf{h}_r\|} = \frac{1}{\kappa(\mathbf{B})} \text{SNR}_{\text{intrinsic}}. \quad (19)$$

□

Discussion. This SNR analysis is particularly critical in the context of LoRA, where the trainable matrices \mathbf{A} and \mathbf{B} typically exhibit large condition numbers, driven by the intrinsic low-rank nature of task-specific features (see Figure 5 and Figure 6). In such regimes, the smallest singular vectors often encode noise or transient artifacts rather than meaningful data signals. Consequently, scaling the noise components by the inverse of their singular values may dilute the meaningful gradient signal, thereby resulting in ineffective subspace learning. By preserving the intrinsic signal-to-noise ratio, OD-LoRA ensures that the dominant signal component is effectively captured, facilitating effective subspace evolution even when the singular value distribution of the desired weight updates is highly skewed.

D.1.2 CONVERGENCE ANALYSIS

Moreover, we compare OD-LoRA with ScaledGD from the perspective of convergence rate. We first introduce several assumptions used in our analysis.

Assumption 1. (β -smoothness) *The loss function \mathcal{L} is β -smooth with respect to \mathbf{W}_{eff} :*

$$\mathcal{L}(\mathbf{W}_{\text{eff}} + \Delta\mathbf{W}_{\text{eff}}) \leq \mathcal{L}(\mathbf{W}_{\text{eff}}) + \underbrace{\langle \mathbf{G}, \Delta\mathbf{W}_{\text{eff}} \rangle_F}_{\text{Descent}} + \underbrace{\frac{\beta}{2} \|\Delta\mathbf{W}_{\text{eff}}\|_F^2}_{\text{Penalty}}. \quad (20)$$

This assumption is required to ensure the curvature of the optimization landscape is bounded. This allows us to derive a step size regime under which the first-order Taylor approximation remains valid, guaranteeing monotonic descent of the loss function.

Assumption 2. (*PL condition*) *The loss function \mathcal{L} satisfies the μ -Polyak–Łojasiewicz (PL) condition: $\|\nabla_{\mathbf{W}_{\text{eff}}} \mathcal{L}\|_F^2 \geq 2\mu(\mathcal{L}(\mathbf{W}_{\text{eff}}) - \mathcal{L}^*)$ where \mathcal{L}^* denotes the minimal loss, which is assumed to be 0 under the overparameterized regime.*

Unlike strong convexity conditions, the PL condition allows for non-convex landscapes with multiple global minima, which is typical in deep learning. This condition relates the gradient magnitude to the sub-optimality gap, providing the necessary geometry to establish a linear convergence rate.

Assumption 3. (α -alignment) *We assume that the projection of full gradient \mathbf{G} onto $\mathcal{R}(\mathbf{A})$ and $\mathcal{C}(\mathbf{B})$ captures fraction α of the energy of \mathbf{G} where $0 < \alpha \leq 2$:*

$$\|\mathbf{G}\mathbf{P}_{\mathcal{R}(\mathbf{A})}\|_F^2 + \|\mathbf{P}_{\mathcal{C}(\mathbf{B})}\mathbf{G}\|_F^2 = \alpha\|\mathbf{G}\|_F^2. \quad (21)$$

In Equation 20, there are two terms to determine the loss-decrease bound: the descent term and the penalty term. We first derive lemmas on these terms.

Lemma 1. (*Descent term*) *For both OD-LoRA and ScaledGD, whose update rules are formulated as Equation 11 and Equation 12, respectively, the descent term in Equation 20 is given by*

$$\langle \mathbf{G}, \Delta\mathbf{W}_{\text{eff}} \rangle_F = -\alpha\eta\|\mathbf{G}\|_F^2. \quad (22)$$

Proof. For both methods, the effective weight update is a projection of the gradient onto the low-rank subspaces: $\Delta\mathbf{W}_{\text{eff}} \approx -\eta(\mathbf{P}_{\mathcal{C}(\mathbf{B})}\mathbf{G} + \mathbf{G}\mathbf{P}_{\mathcal{R}(\mathbf{A})})$. Substituting this into the inner product:

$$\begin{aligned} \langle \mathbf{G}, \Delta\mathbf{W}_{\text{eff}} \rangle_F &\approx \langle \mathbf{G}, -\eta(\mathbf{P}_{\mathcal{C}(\mathbf{B})}\mathbf{G} + \mathbf{G}\mathbf{P}_{\mathcal{R}(\mathbf{A})}) \rangle_F \\ &= -\eta(\langle \mathbf{G}, \mathbf{P}_{\mathcal{C}(\mathbf{B})}\mathbf{G} \rangle_F + \langle \mathbf{G}, \mathbf{G}\mathbf{P}_{\mathcal{R}(\mathbf{A})} \rangle_F). \end{aligned} \quad (23)$$

For any matrix \mathbf{X} and projection \mathbf{P} , $\langle \mathbf{X}, \mathbf{P}\mathbf{X} \rangle_F = \text{Tr}(\mathbf{X}^\top \mathbf{P}\mathbf{X}) = \text{Tr}(\mathbf{X}^\top \mathbf{P}^\top \mathbf{P}\mathbf{X}) = \|\mathbf{P}\mathbf{X}\|_F^2$. Applying this to our terms:

$$\begin{aligned} \langle \mathbf{G}, \mathbf{P}_{\mathcal{C}(\mathbf{B})}\mathbf{G} \rangle_F &= \|\mathbf{P}_{\mathcal{C}(\mathbf{B})}\mathbf{G}\|_F^2, \\ \langle \mathbf{G}, \mathbf{G}\mathbf{P}_{\mathcal{R}(\mathbf{A})} \rangle_F &= \langle \mathbf{G}^\top, \mathbf{P}_{\mathcal{R}(\mathbf{A})}^\top \mathbf{G}^\top \rangle_F = \|\mathbf{G}\mathbf{P}_{\mathcal{R}(\mathbf{A})}\|_F^2. \end{aligned} \quad (24)$$

Finally, applying the α -alignment assumption in Equation 21, we obtain

$$\begin{aligned} \langle \mathbf{G}, \Delta\mathbf{W}_{\text{eff}} \rangle_F &= -\eta(\|\mathbf{G}\mathbf{P}_{\mathcal{R}(\mathbf{A})}\|_F^2 + \|\mathbf{P}_{\mathcal{C}(\mathbf{B})}\mathbf{G}\|_F^2) \\ &= -\alpha\eta\|\mathbf{G}\|_F^2. \end{aligned} \quad (25)$$

□

1026 **Lemma 2. (Penalty term)** For OD-LoRA, the bound for the penalty term in Equation 20 is given by

$$1027 \text{Penalty}_{\text{OD-LoRA}} \leq \alpha\beta\eta^2\|\mathbf{G}\|_F^2. \quad (26)$$

1028 For ScaledGD, the bound for the penalty term is given by

$$1029 \text{Penalty}_{\text{ScaledGD}} \leq \alpha\beta\kappa^2\eta^2\|\mathbf{G}\|_F^2, \quad (27)$$

1030 where $\kappa = \frac{\max(\sigma_{\max}(\mathbf{A}), \sigma_{\max}(\mathbf{B}))}{\min(\sigma_{\min}(\mathbf{A}), \sigma_{\min}(\mathbf{B}))}$.

1031 *Proof.* We bound the penalty term $\frac{\beta}{2}\|\Delta\mathbf{W}_{\text{eff}}\|_F^2$ by relating it to the parameter $(\theta = (\mathbf{A}, \mathbf{B}))$
1032 updates. For any factorization $\mathbf{W}_{\text{eff}} \approx \mathbf{W}_0 + \mathbf{B}\mathbf{A}$ (or involving \mathbf{Q}), the triangle inequality and
1033 sub-multiplicative property of norms give us

$$1034 \begin{aligned} \|\Delta\mathbf{W}_{\text{eff}}\|_F &= \|\mathbf{B}\Delta\mathbf{A} + \Delta\mathbf{B}\mathbf{A}\|_F \\ 1035 &\leq \|\mathbf{B}\|_2\|\Delta\mathbf{A}\|_F + \|\Delta\mathbf{B}\|_F\|\mathbf{A}\|_2 \\ 1036 &\leq \sigma_{\max}(\mathbf{A}, \mathbf{B})(\|\Delta\mathbf{A}\|_F + \|\Delta\mathbf{B}\|_F), \end{aligned} \quad (28)$$

1037 where $\sigma_{\max}(\mathbf{A}, \mathbf{B}) = \max(\|\mathbf{A}\|_2, \|\mathbf{B}\|_2)$. Using $(x + y)^2 \leq 2x^2 + 2y^2$, we obtain

$$1038 \|\Delta\mathbf{W}_{\text{eff}}\|_F^2 \leq 2\sigma_{\max}^2(\mathbf{A}, \mathbf{B})(\|\Delta\mathbf{A}\|_F^2 + \|\Delta\mathbf{B}\|_F^2). \quad (29)$$

1044 1. OD-LoRA

1045 \mathbf{Q}_B and \mathbf{Q}_A^\top have orthonormal columns/rows, so their spectral norms are exactly 1. Thus,
1046 $\sigma_{\max}(\mathbf{Q}_A, \mathbf{Q}_B) = 1$. Using the update rules $\Delta\mathbf{A} = -\eta\mathbf{Q}_B^\top\mathbf{G}$ and $\Delta\mathbf{B} = -\eta\mathbf{G}\mathbf{Q}_A$:

$$1047 \begin{aligned} \|\Delta\mathbf{A}\|_F^2 &= \eta^2\|\mathbf{Q}_B^\top\mathbf{G}\|_F^2 = \eta^2\text{Tr}(\mathbf{G}^\top\mathbf{Q}_B\mathbf{Q}_B^\top\mathbf{G}) = \eta^2\|\mathbf{P}_{\mathcal{C}(\mathbf{B})}\mathbf{G}\|_F^2. \\ 1048 \|\Delta\mathbf{B}\|_F^2 &= \eta^2\|\mathbf{G}\mathbf{Q}_A\|_F^2 = \eta^2\text{Tr}(\mathbf{G}\mathbf{Q}_A\mathbf{Q}_A^\top\mathbf{G}^\top) = \eta^2\|\mathbf{G}\mathbf{P}_{\mathcal{R}(\mathbf{A})}\|_F^2. \end{aligned} \quad (30)$$

1049 Substituting these into Equation 29 with $\sigma_{\max} = 1$:

$$1050 \|\Delta\mathbf{W}_{\text{eff}}\|_F^2 \leq 2 \cdot 1^2 \cdot \eta^2(\|\mathbf{P}_{\mathcal{C}(\mathbf{B})}\mathbf{G}\|_F^2 + \|\mathbf{G}\mathbf{P}_{\mathcal{R}(\mathbf{A})}\|_F^2). \quad (31)$$

1051 Using the α -alignment assumption ($\|\mathbf{P}_{\mathcal{C}(\mathbf{B})}\mathbf{G}\|_F^2 + \|\mathbf{G}\mathbf{P}_{\mathcal{R}(\mathbf{A})}\|_F^2 = \alpha\|\mathbf{G}\|_F^2$):

$$1052 \|\Delta\mathbf{W}_{\text{eff}}\|_F^2 \leq 2\alpha\eta^2\|\mathbf{G}\|_F^2. \quad (32)$$

1053 The penalty is therefore $\frac{\beta}{2}\|\Delta\mathbf{W}_{\text{eff}}\|_F^2 \leq \alpha\beta\eta^2\|\mathbf{G}\|_F^2$.

1058 2. ScaledGD

1059 The updates for \mathbf{A} involve pseudoinverse $\mathbf{B}^\dagger = (\mathbf{B}^\top\mathbf{B})^{-1}\mathbf{B}^\top$. From the update rules for \mathbf{A} in
1060 Equation 12, we obtain:

$$1061 \begin{aligned} \|\Delta\mathbf{A}\|_F^2 &= \eta^2\|\mathbf{B}^\dagger\mathbf{G}\|_F^2 \\ 1062 &= \eta^2\|\mathbf{B}^\dagger(\mathbf{P}_{\mathcal{C}(\mathbf{B})}\mathbf{G} + (\mathbf{I} - \mathbf{P}_{\mathcal{C}(\mathbf{B})})\mathbf{G})\|_F^2 \\ 1063 &= \eta^2\|\mathbf{B}^\dagger\mathbf{P}_{\mathcal{C}(\mathbf{B})}\mathbf{G}\|_F^2 \\ 1064 &\leq \eta^2\|\mathbf{B}^\dagger\|_2^2\|\mathbf{P}_{\mathcal{C}(\mathbf{B})}\mathbf{G}\|_F^2. \end{aligned} \quad (33)$$

1065 Since $\|\mathbf{B}^\dagger\|_2 = \frac{1}{\sigma_{\min}(\mathbf{B})}$, we have:

$$1066 \|\Delta\mathbf{A}\|_F^2 \leq \eta^2\frac{1}{\sigma_{\min}^2(\mathbf{B})}\|\mathbf{P}_{\mathcal{C}(\mathbf{B})}\mathbf{G}\|_F^2. \quad (34)$$

1067 Similarly, $\|\Delta\mathbf{B}\|_F^2 \leq \eta^2\frac{1}{\sigma_{\min}^2(\mathbf{A})}\|\mathbf{G}\mathbf{P}_{\mathcal{R}(\mathbf{A})}\|_F^2$.

1068 Substituting these into Equation 29:

$$1069 \begin{aligned} \|\Delta\mathbf{W}_{\text{eff}}\|_F^2 &\leq 2\sigma_{\max}^2(\mathbf{A}, \mathbf{B})\left(\frac{\eta^2}{\sigma_{\min}^2(\mathbf{B})}\|\mathbf{P}_{\mathcal{C}(\mathbf{B})}\mathbf{G}\|_F^2 + \frac{\eta^2}{\sigma_{\min}^2(\mathbf{A})}\|\mathbf{G}\mathbf{P}_{\mathcal{R}(\mathbf{A})}\|_F^2\right) \\ 1070 &\leq 2\kappa^2\eta^2(\|\mathbf{P}_{\mathcal{C}(\mathbf{B})}\mathbf{G}\|_F^2 + \|\mathbf{G}\mathbf{P}_{\mathcal{R}(\mathbf{A})}\|_F^2) \\ 1071 &= 2\alpha\kappa^2\eta^2\|\mathbf{G}\|_F^2. \end{aligned} \quad (35)$$

1072 The penalty is $\frac{\beta}{2}\|\Delta\mathbf{W}_{\text{eff}}\|_F^2 \leq \alpha\beta\kappa^2\eta^2\|\mathbf{G}\|_F^2$. \square

Finally, using Lemma 1 and Lemma 2, we derive the following theorem on the convergence rate for OD-LoRA and ScaledGD:

Theorem 5: Convergence Rate

The convergence rate for OD-LoRA is $\rho_{\text{OD-LoRA}} = \mathcal{O}(\frac{\mu\alpha}{\beta})$:

$$\mathcal{L}_{t+1} \leq (1 - \frac{\mu\alpha}{2\beta})\mathcal{L}_t. \quad (36)$$

The convergence rate for preconditioning-based methods (e.g., ScaledGD) is $\rho_{\text{precond}} = \mathcal{O}(\frac{\mu\alpha}{\beta\kappa^2})$:

$$\mathcal{L}_{t+1} \leq (1 - \frac{\mu\alpha}{2\beta\kappa^2})\mathcal{L}_t, \quad (37)$$

where $\kappa = \frac{\max(\sigma_{\max}(\mathbf{A}), \sigma_{\max}(\mathbf{B}))}{\min(\sigma_{\min}(\mathbf{A}), \sigma_{\min}(\mathbf{B}))}$.

Proof. Let $\mathcal{L}_t = \mathcal{L}(\mathbf{W}_{\text{eff}})$ and $\mathcal{L}_{t+1} = \mathcal{L}(\mathbf{W}_{\text{eff}} + \Delta\mathbf{W}_{\text{eff}})$.

1. OD-LoRA

Substitute the Descent and Penalty lemmas into Equation 20:

$$\mathcal{L}_{t+1} \leq \mathcal{L}_t - \alpha\eta\|\mathbf{G}\|_F^2 + \alpha\beta\eta^2\|\mathbf{G}\|_F^2. \quad (38)$$

To find the optimal step size, we define $f(\eta) = -\alpha\eta + \alpha\beta\eta^2$ and solve $f'(\eta^*) = -\alpha + 2\alpha\beta\eta = 0$. This gives us $\eta^* = \frac{1}{2\beta}$. By substituting η^* back into Equation 38, we obtain:

$$\mathcal{L}_{t+1} \leq \mathcal{L}_t - \frac{\alpha}{2\beta}\|\mathbf{G}\|_F^2 + \frac{\alpha}{4\beta}\|\mathbf{G}\|_F^2 = \mathcal{L}_t - \frac{\alpha}{4\beta}\|\mathbf{G}\|_F^2. \quad (39)$$

Applying the PL condition ($\|\mathbf{G}\|_F^2 \geq 2\mu\mathcal{L}_t$), we have

$$\begin{aligned} \mathcal{L}_{t+1} &\leq \mathcal{L}_t - \frac{\alpha}{4\beta}(2\mu\mathcal{L}_t) \\ &= (1 - \frac{\mu\alpha}{2\beta})\mathcal{L}_t. \end{aligned} \quad (40)$$

2. ScaledGD

Substitute the Descent and Penalty lemmas for ScaledGD:

$$\mathcal{L}_{t+1} \leq \mathcal{L}_t - \alpha\eta\|\mathbf{G}\|_F^2 + \alpha\beta\kappa^2\eta^2\|\mathbf{G}\|_F^2. \quad (41)$$

Similarly, we set $f(\eta) = -\alpha\eta + \alpha\beta\kappa^2\eta^2$ and solve $f'(\eta^*) = -\alpha + 2\alpha\beta\kappa^2\eta = 0$, obtaining $\eta^* = \frac{1}{2\beta\kappa^2}$. By substituting η^* into Equation 41:

$$\mathcal{L}_{t+1} \leq \mathcal{L}_t - \frac{\alpha}{2\beta\kappa^2}\|\mathbf{G}\|_F^2 + \frac{\alpha}{4\beta\kappa^2}\|\mathbf{G}\|_F^2 = \mathcal{L}_t - \frac{\alpha}{4\beta\kappa^2}\|\mathbf{G}\|_F^2. \quad (42)$$

Applying the PL condition:

$$\begin{aligned} \mathcal{L}_{t+1} &\leq \mathcal{L}_t - \frac{\alpha}{4\beta\kappa^2}(2\mu\mathcal{L}_t) \\ &\leq (1 - \frac{\mu\alpha}{2\beta\kappa^2})\mathcal{L}_t. \end{aligned} \quad (43)$$

□

Discussion. Theorem 5 reveals that while both ScaledGD and OD-LoRA achieve geometric optimality by ensuring the effective weight update is an isotropic projection of the full gradient, they differ significantly in their stability and resulting convergence speeds. ScaledGD achieves this alignment through preconditioning, a mechanism that inherently scales parameter updates by the inverse of the smallest singular values. To mitigate the resulting parameter instability, the step size must be dampened by the square of the condition number (κ^2), as evidenced by the derived optimal step size η^* . In contrast, OD-LoRA resolves this issue by restructuring the parameterization via orthonormal bases, effectively enforcing a perfect condition number ($\kappa = 1$) by construction. Consequently, OD-LoRA is the only method that simultaneously secures the optimal descent direction and supports a large, constant learning rate, achieving a convergence rate independent of the weight matrix conditioning.

Table 7: **Comparison with existing gradient-based methods.** We reproduce the results of existing methods. Experiments are conducted with rank 32.

Method	Gemma-2B			LLaMA3-8B		
	MATH	GSM8K	HumanEval	MATH	GSM8K	HumanEval
LoRA-GA	17.92±0.22	52.56±0.44	32.52±0.76	25.64±0.63	76.43±0.92	45.12±1.22
LoRA-Pro	17.32±0.41	52.11±0.25	30.49±0.76	25.02±0.44	76.02±0.37	43.90±0.82
rsLoRA + ScaledAdamW	17.22±0.42	52.37±0.63	31.10±1.22	25.21±0.17	75.49±0.28	44.76±0.86
AltLoRA	17.03±0.19	52.44±0.90	30.46±0.4	25.17±0.48	75.92±0.30	43.97±0.61
OD-LoRA	18.47±0.14	55.34±1.10	34.15±0.50	26.45±0.39	77.22±0.19	47.76±1.25

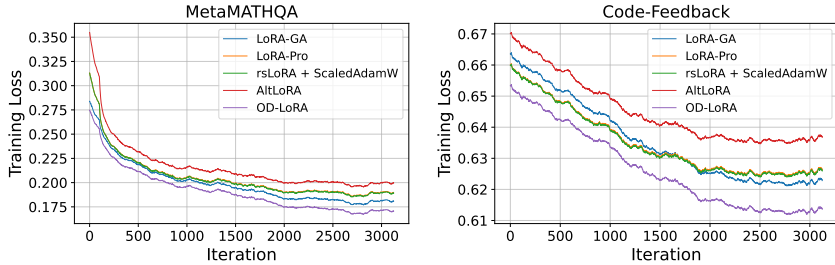


Figure 3: **Comparison with existing gradient-based methods using training loss curve.** Experiments are conducted on LLaMA3-8B with rank 8. We smooth the curves for better visualization.

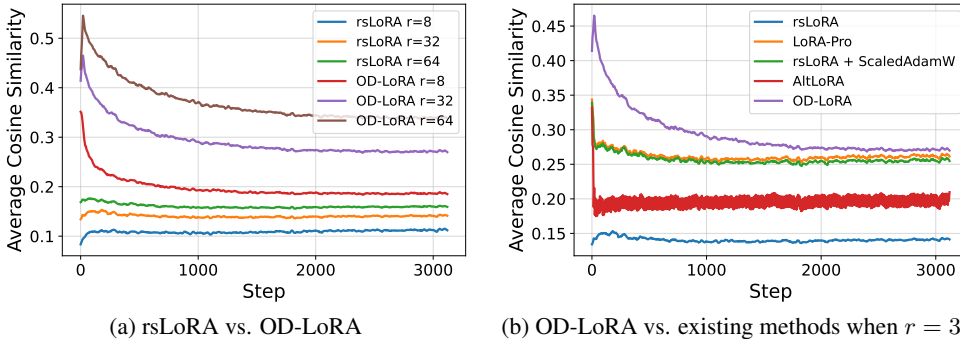


Figure 4: **Cosine similarity between full fine-tuning gradient (G) and the equivalent gradient (\hat{G}).** We measure the average cosine similarity across all target layers and present the rolling window mean over training steps, with a window size of 5.

D.2 EMPIRICAL COMPARISONS

In this subsection, we provide empirical comparisons with existing gradient-based methods (Wang et al., 2024; 2025; Zhang & Pilanci, 2024; Yu et al., 2025). Table 7 shows the adaptation performance comparisons under the rank-32 setting. Similar to the rank-8 results in Table 5, the results demonstrate that OD-LoRA significantly outperforms existing methods. We also compare these methods using training loss curves. Figure 3 shows that OD-LoRA achieves faster convergence toward lower loss. Moreover, we provide the comparisons in terms of alignment between the full fine-tuning gradient and its low-rank approximation in each method. Figure 4b demonstrates that OD-LoRA achieves better gradient alignment. These empirical findings align with our theoretical analyses in Section D.1, which demonstrate that preconditioning in existing methods results in slow convergence due to the condition-number penalty (κ^2) and degraded gradient alignment caused by noise-amplified subspace learning. By decoupling the update magnitude from singular values, OD-LoRA avoids these pitfalls, resulting in the superior convergence speed and higher gradient alignment observed in our experiments.

E IMPORTANCE OF IMPROVING GRADIENT APPROXIMATION QUALITY

In this section, we discuss the reason why improving gradient approximation quality is essential. Using the notations and assumptions introduced in Section 3.2, the descent lemma (Nesterov, 2018) gives us

$$\mathcal{L}(\mathcal{W} + \Delta\mathcal{W}) \leq \mathcal{L}(\mathcal{W}) + \langle \mathcal{G}, \Delta\mathcal{W} \rangle + \frac{\beta}{2} \|\Delta\mathcal{W}\|^2. \quad (44)$$

Equation 3 is simply derived if we set $\Delta\mathcal{W} = -\eta\tilde{\mathcal{G}}$, which implies updating \mathcal{W} using alternative gradient $\tilde{\mathcal{G}}$. Equation 3 suggests that as the alternative gradient exhibits higher alignment with the full fine-tuning gradient, a steeper decrease of loss will be guaranteed. For empirical validation, we measure the cosine similarity between the full fine-tuning gradient and the equivalent gradient during training with rsLoRA or OD-LoRA. We adopt rsLoRA instead of LoRA to compare them under the same scaling factor, and we measure cosine similarity instead of inner product to consider the varying magnitude of the full gradient for each experiment. Figure 4a shows that as the rank of rsLoRA increases, the equivalent gradient exhibits higher alignment with the full gradient. Considering the better performance of LoRA with a higher rank, these results demonstrate the positive correlation between the adaptation performance of rsLoRA and the gradient approximation quality. Moreover, we observe that OD-LoRA significantly improves the alignment, suggesting that the performance gain of OD-LoRA is related to the better gradient alignment.

F SINGULAR VALUE DISTRIBUTION

In the main manuscript, we argue that desired weight updates may not favor uniform singular values. To justify this, we present the distribution of singular values learned by full fine-tuning. Figure 5 shows the singular value distribution in full fine-tuning across various datasets and modules. The results demonstrate that full fine-tuning tends to produce highly skewed singular values. Moreover, in Figure 6, we show the singular value distribution of weight updates learned via LoRA. Similar to the results for full fine-tuning, LoRA typically learns weight updates with non-uniform singular values.

G EXPERIMENTS ON IMAGE CLASSIFICATION TASKS

We use ViT-Base and ViT-Large (Dosovitskiy et al., 2021) trained with 224×224 images and 16×16 patch size. We fine-tune and evaluate them on five datasets, including Cars (Krause et al., 2013), CUB200 (Wah et al., 2011), DTD (Cimpoi et al., 2014), Food101 (Bossard et al., 2014), and SUN397 (Xiao et al., 2010). We use the same setup as for the natural language processing tasks, except for the epochs, learning rate, and batch size, which are detailed in Table 10. We present the complete results in Table 8. The results show that OD-LoRA achieves the best average performance across all models and ranks. In particular, we observe that OD-LoRA significantly reduces the performance gap between LoRA and the full fine-tuning, achieving performance comparable to the full fine-tuning. These results suggest that OD-LoRA generalizes well to vision tasks, highlighting the effectiveness of our method in various domains.

H ADDITIONAL ABLATION STUDIES

In this section, we provide additional ablation studies on the hyperparameters of the proposed initialization method. In Algorithm 1, N and T denote the number of updates and the update interval during the initialization phase. We investigate how sensitive OD-LoRA’s performance is to the choice of N and T . Figure 7 shows OD-LoRA’s performance across various configurations of N and T . The results indicate that performance generally increases as both N and T increase. However, the overall differences are minor, suggesting that OD-LoRA is not sensitive to these hyperparameters.

I DETAILS ON ABLATION STUDIES

In Section 4.5, we propose three variants of rsLoRA and OD-LoRA for ablation studies: ‘rsLoRA w/ Optimal $\tilde{\mathcal{G}}$ ’, ‘OD-LoRA w/ $\text{Cap}_{\mathcal{X}}(\Delta\mathcal{W}) < \infty$ ’ and ‘OD-LoRA w/ Suboptimal $\tilde{\mathcal{G}}$ ’. The first

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

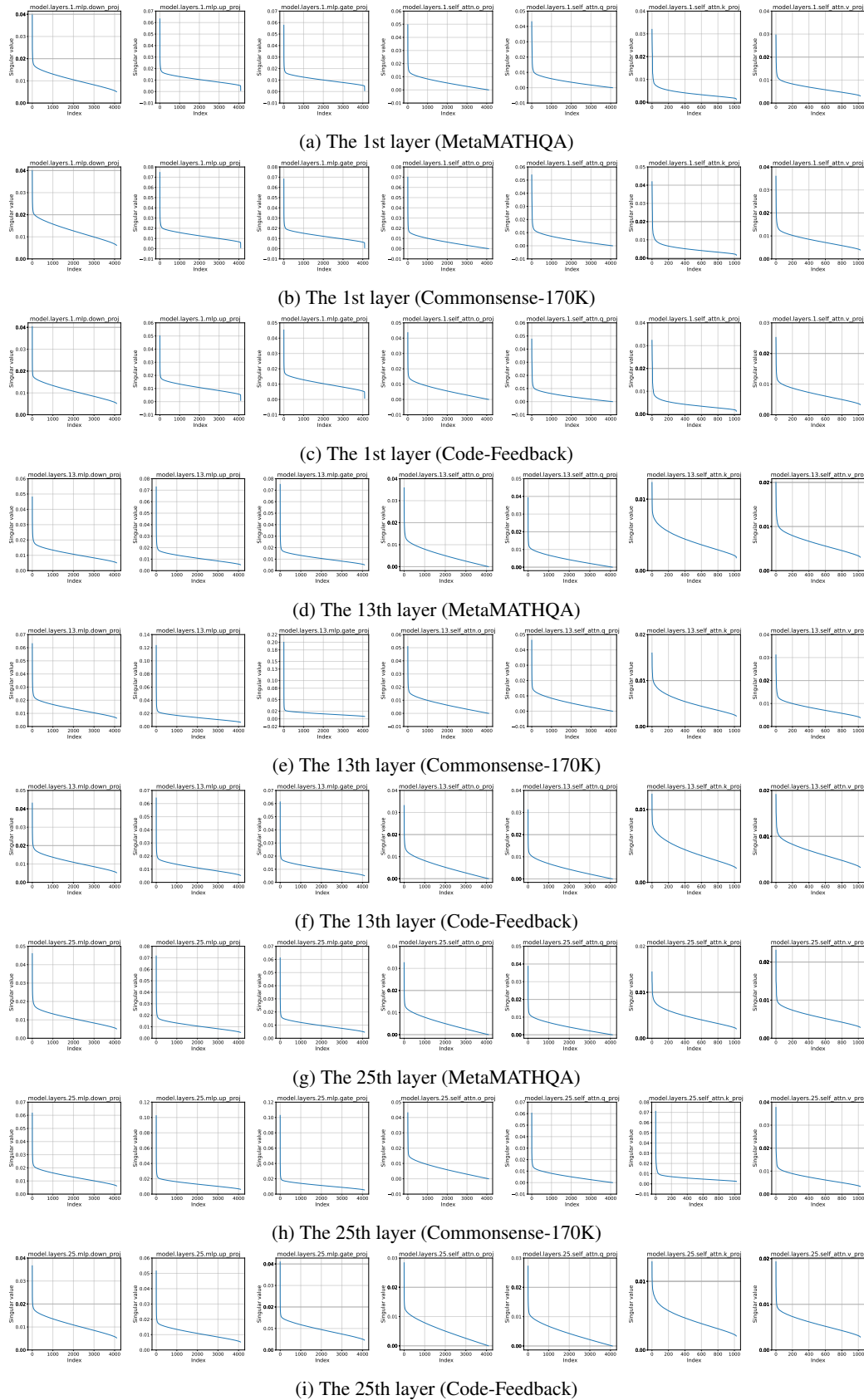


Figure 5: Singular value distributions trained on LLaMA3-8B using full fine-tuning. Best viewed in enlarged form.

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

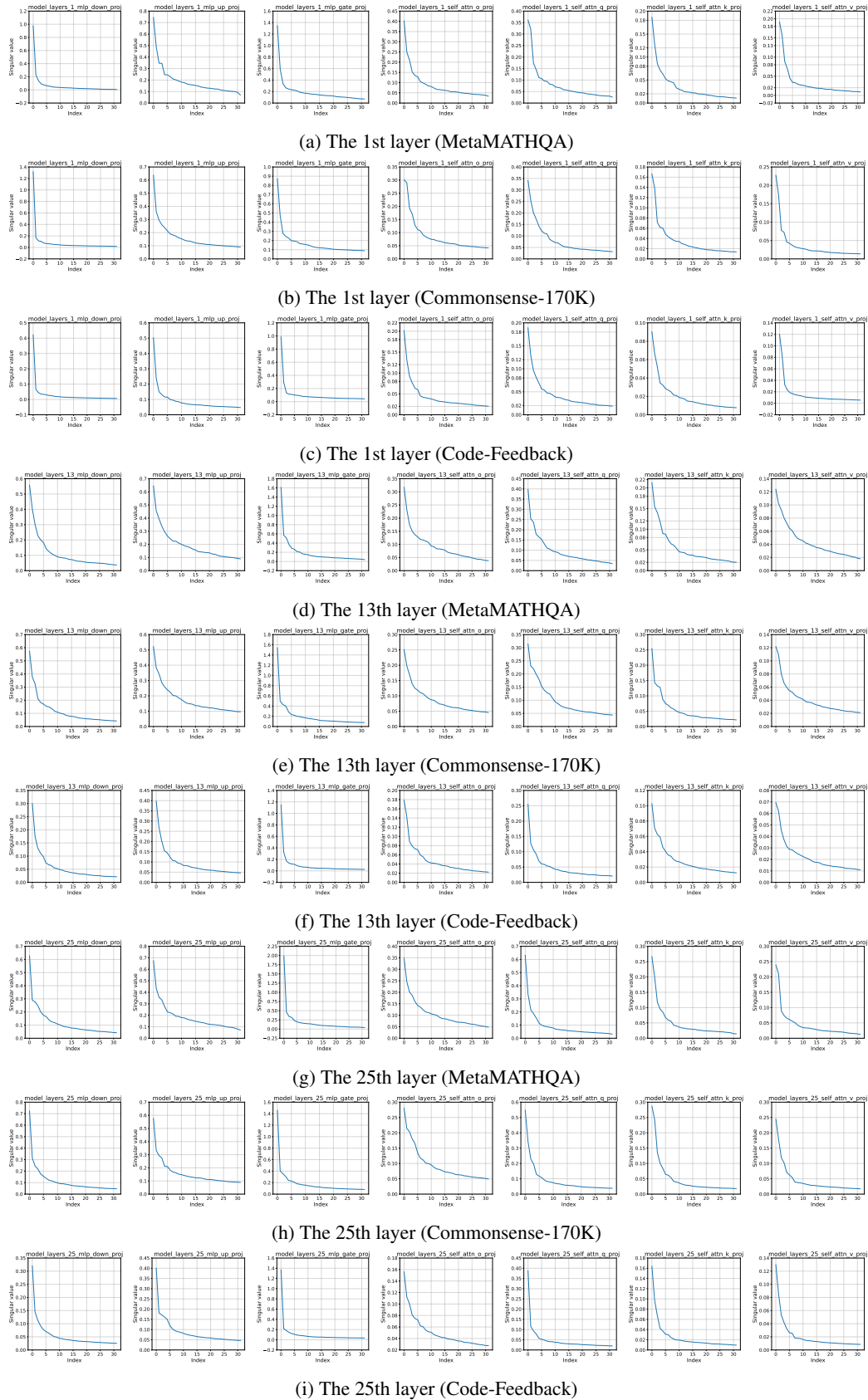


Figure 6: Singular value distributions trained on LLaMA3-8B using LoRA with rank 32. Best viewed in enlarged form.

Table 8: Results on image classification tasks.

Model	Method	Rank	Cars	CUB200	DTD	Food101	SUN397	Avg.
ViT-Base	Full FT.	-	84.26±0.19	86.35±0.18	80.07±0.32	89.27±0.20	74.76±0.26	82.94
	LoRA	8	78.66±0.23	85.65±0.05	78.81±0.95	88.15±0.07	74.35±0.28	81.12
	rsLoRA		78.48±0.07	85.67±0.41	78.79±0.48	88.13±0.16	74.63±0.13	81.14
	LoRA+		79.05±0.25	85.28±0.11	78.32±0.16	88.27±0.05	72.87±0.28	80.76
	PiSSA		78.03±0.32	85.17±0.24	78.33±0.46	87.65±0.07	72.82±0.10	80.40
	DoRA		78.93±0.30	85.74±0.35	78.46±0.16	88.27±0.17	74.46±0.09	81.17
	OD-LoRA	80.47±0.14	86.08±0.16	79.36±0.30	88.35±0.11	74.92±0.16	81.84	
	LoRA	32	81.48±0.20	85.75±0.28	79.34±0.18	88.37±0.12	70.82±0.07	81.15
	rsLoRA		80.57±0.22	85.42±0.41	79.50±0.09	88.81±0.10	73.92±0.26	81.64
	LoRA+		79.53±0.60	85.68±0.31	78.85±0.39	88.05±0.01	67.66±0.06	79.95
PiSSA	80.43±0.17		84.90±0.28	78.90±0.35	87.91±0.08	69.73±0.12	80.37	
DoRA	81.07±0.03		85.67±0.22	79.56±0.43	88.83±0.08	74.05±0.21	81.84	
OD-LoRA	83.04±0.32	86.06±0.24	80.11±0.16	89.05±0.08	74.92±0.21	82.64		
ViT-Large	Full FT.	-	88.19±0.13	88.23±0.19	81.71±0.09	90.90±0.17	76.20±0.22	85.05
	LoRA	8	85.07±0.19	87.53±0.20	80.64±0.23	89.96±0.08	76.30±0.11	83.90
	rsLoRA		84.61±0.11	87.75±0.19	80.87±0.32	89.92±0.05	76.45±0.16	83.92
	LoRA+		81.99±0.97	87.66±0.05	80.62±0.55	90.06±0.03	75.08±0.14	83.08
	PiSSA		85.59±0.17	87.55±0.19	80.34±0.13	89.60±0.17	75.41±0.20	83.70
	DoRA		84.67±0.08	87.58±0.11	80.73±0.33	90.03±0.04	76.50±0.08	83.90
	OD-LoRA	85.80±0.08	87.98±0.08	81.06±0.49	90.05±0.05	76.65±0.13	84.31	
	LoRA	32	86.23±0.05	87.94±0.22	81.15±0.32	90.62±0.03	74.59±0.25	84.11
	rsLoRA		85.93±0.30	87.64±0.06	81.56±0.07	90.58±0.13	76.22±0.13	84.39
	LoRA+		82.40±1.84	87.92±0.22	80.92±0.18	90.48±0.07	75.19±1.23	83.38
PiSSA	86.70±0.15		87.65±0.15	80.82±0.26	90.04±0.09	73.84±0.08	83.81	
DoRA	86.10±0.36		87.69±0.11	81.31±0.13	90.64±0.04	76.23±0.20	84.39	
OD-LoRA	87.07±0.36	88.27±0.10	81.45±0.21	90.75±0.17	76.48±0.06	84.80		

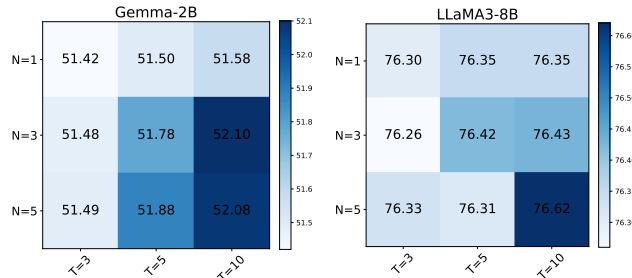
Figure 7: Ablation studies on N and T . Using OD-LoRA, we fine-tune on MetaMATHQA and evaluate on GSM8K.

Table 9: Additional implementation details.

	Full FT.		LoRA-based		
	Gemma-2B	LLaMA3-8B	Gemma-2B	LLaMA2-7B	LLaMA3-8B
Learning Rate	2e-5	1e-5	2e-4	2e-4	1e-4
Learning Rate Scheduler	cosine scheduler				
Epochs	1				
Batch Size	32				
Target Modules	'q_proj', 'k_proj', 'v_proj', 'up_proj', 'down_proj', 'o_proj', 'gate_proj'				

Table 10: Implementation details on image classification tasks.

	ViT-Base					ViT-Large				
	Cars	CUB200	DTD	Food101	SUN397	Cars	CUB200	DTD	Food101	SUN397
Epochs	5	7	7	7	5	5	7	7	7	5
Learning Rate	5e-3	2e-3	2e-3	2e-3	5e-3	2.5e-3	1e-3	1e-3	1e-3	2.5e-3
Batch Size	64									
Target Modules	'query', 'value'									

and second methods formulate the weight updates as $sQ_BQ_A^\top$. The third method formulates the weight updates as $s(Q_B \text{diag}(R_B)A + B(Q_A \text{diag}(R_A))^\top)$. Except for the first method, Q_A and Q_B are initialized the same way as in OD-LoRA for a fair comparison. The update interval for Q_A and Q_B is also the same as in OD-LoRA. In both 'rsLoRA w/ Optimal \tilde{G} ' and 'OD-LoRA w/ $\text{Cap}_X(\Delta W) < \infty$ ', Q_A and Q_B are trainable, but the gradient does not flow through the QR

1404 decomposition process. Thus, the gradient for Q_A and Q_B are expressed as $\frac{\partial \mathcal{L}}{\partial Q_A} = sQ_B^\top G$ and
 1405 $\frac{\partial \mathcal{L}}{\partial Q_B} = sGQ_A$, leading to $\tilde{G} = s^2(GQ_AQ_A^\top + Q_BQ_B^\top G)$. For ‘OD-LoRA w/ Suboptimal \tilde{G} ’,
 1406 the gradient for A and B are expressed as $\frac{\partial \mathcal{L}}{\partial A} = s(Q_B \text{diag}(\mathbf{R}_B))^\top G$ and $\frac{\partial \mathcal{L}}{\partial B} = sGQ_A \text{diag}(\mathbf{R}_A)$,
 1407 leading to $\tilde{G} = s^2(GQ_A \text{diag}(\mathbf{R}_A)^2 Q_A^\top + Q_B \text{diag}(\mathbf{R}_B)^2 Q_B^\top G)$. Hence, ‘rsLoRA w/ Optimal \tilde{G} ’
 1408 and ‘OD-LoRA w/ $\text{Cap}_{\mathcal{X}}(\Delta W) < \infty$ ’ satisfies the uniform singular value condition in Theorem 1,
 1409 whereas ‘OD-LoRA w/ Suboptimal \tilde{G} ’ does not. Moreover, the representational capacity of ‘rsLoRA
 1410 w/ Optimal \tilde{G} ’ and ‘OD-LoRA w/ $\text{Cap}_{\mathcal{X}}(\Delta W) < \infty$ ’ is less than infinity, and that of ‘OD-LoRA w/
 1411 Suboptimal \tilde{G} ’ is infinity as proven in Section B.
 1412
 1413

1414 **Use of Large Language Models.** We use Large Language Models (LLMs) to assist with writing
 1415 and to search for mathematical propositions, such as properties in linear algebra. We acknowledge our
 1416 responsibility for the content of this paper and ensure that the contribution of LLMs is insignificant.
 1417
 1418
 1419
 1420
 1421
 1422
 1423
 1424
 1425
 1426
 1427
 1428
 1429
 1430
 1431
 1432
 1433
 1434
 1435
 1436
 1437
 1438
 1439
 1440
 1441
 1442
 1443
 1444
 1445
 1446
 1447
 1448
 1449
 1450
 1451
 1452
 1453
 1454
 1455
 1456
 1457