
AdaMSS: Adaptive Multi-Subspace Approach for Parameter-Efficient Fine-Tuning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In this paper, we propose AdaMSS, an adaptive multi-subspace approach for
2 parameter-efficient fine-tuning of large models. Unlike traditional parameter-
3 efficient fine-tuning methods that operate within a large single subspace of the
4 network weights, AdaMSS leverages subspace segmentation to obtain multiple
5 smaller subspaces and adaptively reduces the number of trainable parameters
6 during training, ultimately updating only those associated with a small subset of
7 subspaces most relevant to the target downstream task. By using the lowest-rank
8 representation, AdaMSS achieves more compact expressiveness and finer tuning of
9 the model parameters. Theoretical analyses demonstrate that AdaMSS has better
10 generalization guarantee than LoRA, PiSSA, and other single-subspace low-rank-
11 based methods. Extensive experiments across image classification, natural language
12 understanding, and natural language generation tasks show that AdaMSS achieves
13 comparable performance to full fine-tuning and outperforms other parameter-
14 efficient fine-tuning methods in most cases, all while requiring fewer trainable
15 parameters. Notably, on the ViT-Large model, AdaMSS achieves 4.7% higher
16 average accuracy than LoRA across seven tasks, using just 15.4% of the trainable
17 parameters. On RoBERTa-Large, AdaMSS outperforms PiSSA by 7% in average
18 accuracy across six tasks while reducing the number of trainable parameters by
19 approximately 94.4%. These results demonstrate the effectiveness of AdaMSS in
20 parameter-efficient fine-tuning.

21 1 Introduction

22 With the successful application of large pre-trained models in natural language processing and
23 computer vision, the parameter-efficient fine-tuning (PEFT) methods have gradually become key
24 strategies for efficiently adapting these models. Among these PEFT methods, LoRA (Low-Rank
25 Adaptation) [1] has gained particular attention, as it adopts learnable low-rank structures to represent
26 weight changes, significantly reducing memory requirements. Following LoRA, numerous low-rank
27 adaptation methods have since been developed [2, 3, 4, 5, 6, 7, 8, 9, 10, 11].

28 Despite sharing the common goal of PEFT, these methods primarily differ in how they impose
29 low-rank assumptions. For instance, LoRA [1] and AdaLoRA [3] assume that weight changes
30 $\Delta \mathbf{W} \in \mathbb{R}^{n \times n}$ are low-rank and represent them using the product of two smaller matrices, $\mathbf{A} \in \mathbb{R}^{n \times r}$
31 and $\mathbf{B} \in \mathbb{R}^{r \times n}$, where $r \ll n$. Alternatively, PiSSA [2] assumes that the original weight matrix \mathbf{W}
32 is approximately low-rank and achieves efficient fine-tuning by training its principal components. A
33 further line of research is represented by methods such as LoRA-GA [5], which hypothesize that the
34 gradient of the loss function with respect to the weights—*i.e.*, $\nabla \mathcal{L}(\mathbf{W})$ —exhibits low-rank structure.

35 These approaches are typical single-subspace methods, as they rely on the assumption that weight
36 changes, weights, or gradients lie in, or approximately lie in, a single low-dimensional subspace.

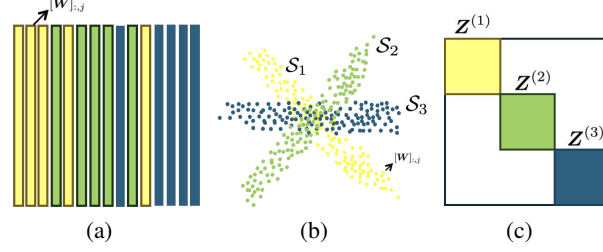


Figure 1: (a) Weight matrix \mathbf{W} ; (b) Multi-subspaces structure of the columns of \mathbf{W} (i.e., $[\mathbf{W}]_{:,j}$); (c) The lowest-rank representation of \mathbf{W} , where each block $\mathbf{Z}^{(k)}$ serves as a new representation for the columns lying in the subspace \mathcal{S}_k .

37 However, this assumption inevitably leads to a fundamental trade-off between limited expressiveness
 38 and parameter efficiency: approximating the weights with a single low-dimensional subspace with
 39 a small r limits the model’s capacity for adaptation, while increasing the subspace dimension r
 40 improves expressiveness but results in a substantial rise in the number of trainable parameters and
 41 memory usage.

42 To address the dilemma between limited expressiveness and parameter efficiency, we propose an
 43 **Adaptive Multi-Subspace** approach (AdaMSS) that enables finer tuning of model parameters and
 44 has stronger expressiveness. As we will illustrate in this work, the columns of the weight matrix are
 45 approximately distributed across multiple linear subspaces (Figure 1 (a)-(b) provides a conceptual
 46 illustration of multi-subspaces structure of the network weights).

47 The key features of this work are summarized as follows:

- 48 • **Compact Expressiveness:** By leveraging the observation that the network weights are
 49 approximately located in multiple subspaces, AdaMSS uses a lowest-rank representation for
 50 the network weights, which exhibit an approximate block-diagonal structure (see Figure 1
 51 (c)). This design leads to a more compact representation while preserving both global and
 52 local low-rank structure from the original weights (see Property 1), which is essential for
 53 ensuring effective adaptation across multiple subspaces.
- 54 • **Generalization Guarantee:** We theoretically show that, under the same rank assumption
 55 r , the proposed AdaMSS achieves a lower Gaussian complexity bound and is therefore
 56 expected to exhibit a stronger generalization capability compared to single-subspace-based
 57 methods. Theorem 1 formalizes an upper bound on the expected loss of AdaMSS, providing
 58 its generalization guarantee.
- 59 • **Multi-Subspace-Based Adaptive Budget Allocation:** Thanks to the multi-subspace struc-
 60 ture of the network weights, we can adaptively freeze the parameters associated with
 61 subspaces of lower importance, and eventually updates only those associated with a small
 62 subset of subspaces relevant to the target downstream task, where the importance of each
 63 subspace is evaluated during training.

64 Experimental results demonstrate the effectiveness of AdaMSS for PEFT across diverse tasks and
 65 show that it outperforms existing methods in most cases while significantly reducing the number of
 66 trainable parameters. For example, compared with LoRA, AdaMSS achieves 4.7% higher average
 67 accuracy on ViT-Large [12] and 2.2% higher accuracy on GSM8K [13] with LLaMA 2-7B [14],
 68 using only 15.35% and 0.28% of the trainable parameters, respectively. On RoBERTa-Large [15],
 69 AdaMSS achieves comparable performance to LoRA while using just 5.62% of the trainable param-
 70 eters. These results suggest that multi-subspace-based fine-tuning offers a promising direction for
 71 improving the trade-off between expressiveness and parameter efficiency.

72 2 Preliminaries and Related Works

73 In this section, we introduce the notations used throughout the paper, and provide preliminaries and
 74 a brief review of existing literature relevant to our study. Unless otherwise specified, we use a , \mathbf{a} ,
 75 and \mathbf{A} to denote scalars, vectors, and matrices, respectively. The transpose of \mathbf{A} is denoted by \mathbf{A}^\top .

We use $[A]_{i,j}$, $[A]_{i,:}$, and $[A]_{:,j}$ to represent the (i, j) -th element, the i -th row vector, and the j -th column vector of A , respectively.

2.1 Preliminaries: Low-Rank Representation for Subspace Segmentation

To get the lowest-rank presentation for n given samples¹ $[M]_{:,1}, [M]_{:,2}, \dots, [M]_{:,n} \in \mathbb{R}^d$, Liu *et al.* proposed the Low-Rank Representation (LRR) model [16] expressed as follows:

$$[Z^*, E^*] = \min_{Z, E} \text{rank}(Z) + \lambda \|E\|_{\diamond} \quad \text{s.t. } M = MZ + E, \quad (1)$$

where Z is known as the coefficient matrix, $\text{rank}(Z)$ denotes the rank of matrix $Z \in \mathbb{R}^{n \times n}$, *i.e.*, the number of singular values of Z , $M \in \mathbb{R}^{d \times n}$ represents the matrix composed of the n samples, E represents the residual component (*i.e.*, noise and redundancy in the samples). The norm $\|\cdot\|_{\diamond}$ is used to quantify the magnitude of the nonzero entries of E , and its specific definition depends on the prior knowledge about E . The problem (1) can be approximated by a convex formulation and solved iteratively [16], and it is computationally expensive. Fortunately, when $E = 0$, problem (1) has a closed-form solution, given by $Z^* = VV^{\top}$, where $V \in \mathbb{R}^{n \times \text{rank}(M)}$ is obtained through the skinny SVD of M , *i.e.*, $M = USV^{\top}$.

Assuming the samples $[M]_{:,1}, [M]_{:,2}, \dots, [M]_{:,n}$ drawn from a union of linear subspaces $\{S_k\}_{k=1}^K$, where the subspace number K and $\{S_k\}_{k=1}^K$ are unknown, the goal of subspace segmentation is to simultaneously assign the samples into their respective subspaces. When the subspaces $\{S_k\}_{k=1}^K$ are independent of each other, samples from different subspaces cannot be expressed as linear combinations of each other. Consequently, Z^* becomes a block-diagonal matrix. To solve the subspace segmentation problem, Liu *et al.* determine the subspace number K and to obtain the segmentation results based on Z^* [16]. This LLR-based subspace segmentation method is widely applied in clustering tasks[16, 17, 18, 19].

2.2 Related Works

2.2.1 Parameter-Efficient Fine-Tuning (PEFT)

In recent years, a variety of PEFT methods have been proposed to adapt large pre-trained models to downstream tasks while minimizing the number of trainable parameters. Broadly, current PEFT methods include prefix-tuning [20, 21], adapter-based methods [22], sparse fine-tuning [23, 24], orthogonal fine-tuning [25, 26], and low-rank-based adaptation [1, 2, 3, 4, 5, 6, 7, 8, 9]. Specifically, prefix-tuning methods introduce a small task-specific trainable vectors, known as prefixes, while keeping the original model parameters frozen [20, 21]. In contrast, adapter-based methods [22] insert small trainable modules within each layer of the model, allowing efficient task adaptation by updating only these modules. Meanwhile, sparse fine-tuning, such as BitFit [23] and FourierFT [24], update only a small subset of parameters in bias terms or transformed weights. On another front, orthogonal fine-tuning introduces a learnable orthogonal transformation applied to the pre-trained weights, with the goal of preserving the model's original capabilities during adaptation [25, 26]. In addition to the above, low-rank-based adaptation methods, for example LoRA [1], PiSSA [2], LoRA-GA [5], FourA [8], and RoSA [7], approximate weight updates using low-rank matrices, significantly reducing the number of trainable parameters while maintaining performance. These techniques provide effective strategies for adapting large models to new tasks with less computational overhead.

2.2.2 Low-Rank-Based Adaptation Methods

Inspired by the success of LoRA [1], a large number of low-rank-based adaptation methods have emerged over the past three years [2, 3, 4, 5, 6, 7, 8, 9, 10, 11]. LoRA assumes a fixed rank for all incremental matrices across different layers, thereby ignoring the diverse importance of different weight parameters. To address this limitation, AdaLoRA [3] was proposed to allocates parameter budgets by adaptively adjusting the rank during training. Following AdaLoRA, research attention has shifted toward improving the initialization strategies of A and B in LoRA, as opposed to relying on random initialization. Representative works include PiSSA [2] and LoRA-GA[5]. In addition, by incorporating both low-rankness and sparsity constraints, RoSA [7] enhances low-rank adaptation

¹In this work, the *samples* here does not refer to layer input x .

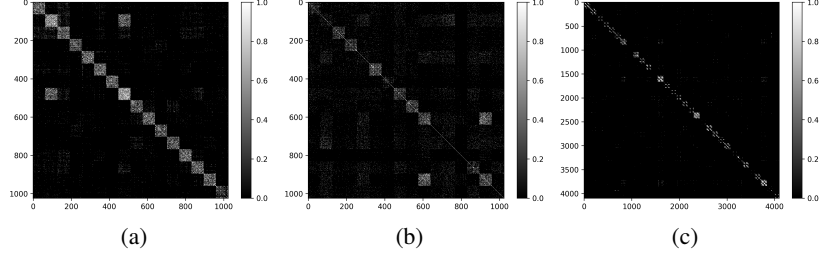


Figure 2: Illustration of multiple subspaces structure in pretrained network weights: an approximate block-diagonal structure of Z_0^* from the first layer of the pretrained (a) ViT-Large (query), (b) Roberta-Large (query), and (c) LLaMA 2-7B (query).

to enable more efficient parameter utilization. In another line of work, FouRA [8] extends LoRA into the frequency domain, yielding disentangled feature spaces that enable fine-grained control and editing. Motivated by the success of tensor decomposition in data compression [27], LoTR [11] and LoRETTA [10] leverage Tucker decomposition and Tensor-train decomposition, respectively, to achieve more compact representations of the low-rank update matrices.

3 Multi-Subspace-Based Adaptation

Let $W_0 \in \mathbb{R}^{d \times n}$ denote the pretrained weight matrix of a network layer, where d and n represent the input and output dimensions of the layer, respectively. The full fine-tuning (FF) of the network is formalized as

$$\Delta W^* = \arg \min_{\Delta W \in \mathbb{R}^{d \times n}} \mathcal{L}(W_0 + \Delta W), \quad (2)$$

where \mathcal{L} denotes the loss function, $\Delta W \in \mathbb{R}^{d \times n}$ is the update on W_0 , and $W = W_0 + \Delta W$ is the corresponding updated weights. In this section, we aim to find a compact representation of W_0 via subspace segmentation [16], which has both low-rank and strictly block-diagonal structure. By doing so, instead of directly updating W_0 as FF, we perform incremental updates within this compact representation.

3.1 Motivation: Multi-Subspace Structure of the Network Weights

Motivated by PiSSA [2], we focus on updating the principle components of W_0 (denoted as \hat{W}_0), and decompose the pretrained weight W_0 as $W_0 = \hat{W}_0 + W_{res}$, where $W_{res} := W_0 - [U_0]_{:,1:R}[S_0]_{1:R,1:R}[V_0]_{:,1:R}^\top$ for a given R , and $\hat{W}_0 = [U_0]_{:,1:R}[S_0]_{1:R,1:R}[V_0]_{:,1:R}^\top$ is the truncated SVD of W_0 . We apply the LRR to obtain a lowest-rank representation of \hat{W}_0 . Considering computational complexity, we adapt the LRR with $E = 0$:

$$Z_0^* = \arg \min_{Z_0 \in \mathbb{R}^{n \times n}} \text{rank}(Z_0) \quad s.t. \quad \hat{W}_0 = \hat{W}_0 Z_0. \quad (3)$$

This formulation uses the coefficient matrix Z_0 to linearly represent the columns of $\hat{W}_0 \in \mathbb{R}^{d \times n}$. By analyzing Z_0^* , we observe that it exhibits an approximate block-diagonal structure (see Figure 2 for $R = 100$), indicating that the columns of $\hat{W}_0 = \hat{W}_0 Z_0^*$ are grouped into different clusters, where columns within the same cluster are approximately linear combinations of each other. In other words, the row space of W_0 can be approximated by multiple smaller low-rank subspaces. By leveraging the approximate block-diagonal structure and low-rank properties of Z_0^* , we drive a more compact representation of W_0 than \hat{W}_0 .

3.2 A Compact Representation By Subspace Segmentation

In the previous subsection, we analyzed the multi-subspace structure of W_0 and obtained a compact representation Z_0^* , which exhibits an approximate block-diagonal structure. We now derive a well-organized and strictly block-diagonal representation of W_0 by subspace segmentation.

Leveraging subspace segmentation and the multi-subspace structure, the row space of \mathbf{W}_0 can be approximately segmented into K smaller low-dimensional subspaces. We assume that the columns in $\mathbf{W}_0 = [\mathbf{W}_0^{(1)}, \mathbf{W}_0^{(2)}, \dots, \mathbf{W}_0^{(K)}]$ are well-organized and partitioned based on the segmentation results, where each block $\mathbf{W}_0^{(k)} \in \mathbb{R}^{d \times n_k}$ contains the columns assigned to the k -th subspace. Accordingly, $\hat{\mathbf{W}}_0$, \mathbf{Z}_0^* , and \mathbf{W}_{res} are also divided into K blocks, leading to the following decomposition

$$\begin{aligned} & [\mathbf{W}_0^{(1)}, \mathbf{W}_0^{(2)}, \dots, \mathbf{W}_0^{(K)}] \\ &= [\hat{\mathbf{W}}_0^{(1)}, \hat{\mathbf{W}}_0^{(2)}, \dots, \hat{\mathbf{W}}_0^{(K)}] \left(\text{diag} \left((\mathbf{Z}_0^*)^{(1)}, (\mathbf{Z}_0^*)^{(2)}, \dots, (\mathbf{Z}_0^*)^{(K)} \right) + [\mathbf{E}_0^{(1)}, \mathbf{E}_0^{(2)}, \dots, \mathbf{E}_0^{(K)}] \right) \\ &+ [\mathbf{W}_{res}^{(1)}, \mathbf{W}_{res}^{(2)}, \dots, \mathbf{W}_{res}^{(K)}], \end{aligned} \quad (4)$$

where $\mathbf{Z}_0^* = [[\mathbf{Z}_0^*]^{(1)}, [\mathbf{Z}_0^*]^{(2)}, \dots, [\mathbf{Z}_0^*]^{(K)}]$ is decomposed into a sum of a strictly block-diagonal matrix $\text{diag} \left((\mathbf{Z}_0^*)^{(1)}, (\mathbf{Z}_0^*)^{(2)}, \dots, (\mathbf{Z}_0^*)^{(K)} \right)$ and a noise matrix $[\mathbf{E}_0^{(1)}, \mathbf{E}_0^{(2)}, \dots, \mathbf{E}_0^{(K)}]$, and each block $(\mathbf{Z}_0^*)^{(k)} \in \mathbb{R}^{n_k \times n_k}$ is the k -th diagonal block of \mathbf{Z}_0^* .

Define $\bar{\mathbf{W}}_{res}^{(k)} := \hat{\mathbf{W}}_0^{(k)} \mathbf{E}_0^{(k)} + \mathbf{W}_{res}^{(k)}$, we have

$$\mathbf{W}_0^{(k)} = \hat{\mathbf{W}}_0^{(k)} (\mathbf{Z}_0^*)^{(k)} + \bar{\mathbf{W}}_{res}^{(k)} \quad (5)$$

for each k , where $(\mathbf{Z}_0^*)^{(k)} \in \mathbb{R}^{n_k \times n_k}$ is used to approximately represents $\mathbf{W}_0^{(k)}$. We thus obtain

$$\begin{aligned} [\mathbf{W}_0^{(1)}, \mathbf{W}_0^{(2)}, \dots, \mathbf{W}_0^{(K)}] &= [\hat{\mathbf{W}}_0^{(1)}, \hat{\mathbf{W}}_0^{(2)}, \dots, \hat{\mathbf{W}}_0^{(K)}] \text{diag} \left((\mathbf{Z}_0^*)^{(1)}, (\mathbf{Z}_0^*)^{(2)}, \dots, (\mathbf{Z}_0^*)^{(K)} \right) \\ &+ [\bar{\mathbf{W}}_{res}^{(1)}, \bar{\mathbf{W}}_{res}^{(2)}, \dots, \bar{\mathbf{W}}_{res}^{(K)}], \end{aligned} \quad (6)$$

and a new representation of \mathbf{W}_0 :

$$\text{diag} \left((\mathbf{Z}_0^*)^{(1)}, (\mathbf{Z}_0^*)^{(2)}, \dots, (\mathbf{Z}_0^*)^{(K)} \right). \quad (7)$$

The following properties hold for this representation.

Property 1. If the subspaces are independent (i.e., $\text{rank}(\mathbf{Z}_0^*) = \sum_{k=1}^K \text{rank}([\mathbf{Z}_0^*]^{(k)})$), then

(1) (Global structure-preserving) $\sum_{k=1}^K \text{rank}(\hat{\mathbf{W}}_0^{(k)}) = \sum_{k=1}^K \text{rank}((\mathbf{Z}_0^*)^{(k)}) = \text{rank}(\hat{\mathbf{W}}_0)$;

(2) (Local structure-preserving) $\text{rank}((\mathbf{Z}_0^*)^{(k)}) = \text{rank}(\hat{\mathbf{W}}_0^{(k)})$ for $k = 1, 2, \dots, K$.

Proof. The corresponding proof is given in the Appendix A. \square

Property 1 ensures both global and local low-rank structure preservation by the representation of \mathbf{W}_0 . Specifically, the total rank across all subspaces remains consistent with that of the original weight matrix (Property 1 (1)), while each subspace individually preserves its intrinsic low-rank structure (Property 1 (2)). These properties are critical for maintaining original model abilities and enabling reliable adaptation across multiple subspaces.

3.3 Multi-Subspace-Based Incremental Update

In the previous subsection, we presented a representation of \mathbf{W}_0 based on subspace segmentation, as shown in (7). In this subsection, we propose a multi-subspace-based adaptation strategy. Rather than updating \mathbf{W}_0 directly, the proposed adaptation performs incremental updates for each $(\mathbf{Z}_0^*)^{(k)}$.

We introduce an incremental update $\Delta \mathbf{Z}^{(k)}$ (with rank r_k) for each $(\mathbf{Z}_0^*)^{(k)}$ in (6), and compute the updated weight block $\mathbf{W}^{(k)}$ as

$$\mathbf{W}^{(k)} = \hat{\mathbf{W}}_0^{(k)} ((\mathbf{Z}_0^*)^{(k)} + \Delta \mathbf{Z}^{(k)}) + \bar{\mathbf{W}}_{res}^{(k)} = \mathbf{W}_0^{(k)} + \hat{\mathbf{W}}_0^{(k)} \Delta \mathbf{Z}^{(k)}, \quad (8)$$

where the second equality follows from (5). Since the computation involves multiplication with the large matrix $\hat{\mathbf{W}}_0^{(k)}$, calculating the gradient with respect to $\Delta \mathbf{Z}^{(k)}$ can be computationally expensive.

184 To address this, we decompose $\hat{\mathbf{W}}_0^{(k)}$ by the skinny SVD: $\hat{\mathbf{W}}_0^{(k)} = \mathbf{U}_{\hat{\mathbf{W}}_0}^{(k)} \mathbf{S}_{\hat{\mathbf{W}}_0}^{(k)} (\mathbf{V}_{\hat{\mathbf{W}}_0}^{(k)})^\top$. Let
 185 $R_k := \text{rank}(\hat{\mathbf{W}}_0^{(k)})$ and $\mathbf{A}^{(k)} := \mathbf{U}_{\hat{\mathbf{W}}_0}^{(k)} \in \mathbb{R}^{d \times R_k}$. Then, the update becomes:

$$\mathbf{W}^{(k)} = \mathbf{W}_0^{(k)} + \mathbf{A}^{(k)} \mathbf{S}_{\hat{\mathbf{W}}_0}^{(k)} (\mathbf{V}_{\hat{\mathbf{W}}_0}^{(k)})^\top \Delta \mathbf{Z}^{(k)}. \quad (9)$$

186 Since $\text{rank}(\mathbf{S}_{\hat{\mathbf{W}}_0}^{(k)} (\mathbf{V}_{\hat{\mathbf{W}}_0}^{(k)})^\top \Delta \mathbf{Z}^{(k)}) \leq r_k$, we introduce two trainable matrices $\mathbf{B}^{(k)} \in \mathbb{R}^{R_k \times r_k}$ and
 187 $\mathbf{C}^{(k)} \in \mathbb{R}^{r_k \times n_k}$, and reparameterize the k -th updated weight block as:

$$\mathbf{W}^{(k)} = \mathbf{W}_0^{(k)} + \mathbf{A}^{(k)} \mathbf{B}^{(k)} \mathbf{C}^{(k)}. \quad (10)$$

188 This formulation yields a low-rank update (*i.e.*, $\mathbf{A}^{(k)} \mathbf{B}^{(k)} \mathbf{C}^{(k)}$) within each subspace.

189 To ensure that the initial value of $\mathbf{W}^{(k)}$ remains consistent with $\mathbf{W}_0^{(k)}$, we initialize $\mathbf{C}^{(k)}$ in (10) as
 190 $\mathbf{0}$. The matrix $\mathbf{B}^{(k)} \in \mathbb{R}^{R_k \times r_k}$ can be initialized either as an orthogonal matrix, as described in Steps
 191 7-8 of the Algorithm 1 (see Appendix C), or randomly (following the LoRA initialization strategy).
 192 The complete initialization procedure of AdaMSS is presented in Algorithm 1, and the estimation
 193 of the subspace number K is detailed in Algorithm 2, both of which are included in Appendix C.
 194 The parameter count and computational complexity analysis of the proposed method, along with
 195 comparisons to LoRA and PiSSA, are provided in the Appendices D and E, respectively.

196 4 Multi-Subspace-Based Adaptive Budget Allocation

197 In this section, we introduce an adaptive multi-subspace-based budget allocation mechanism. By
 198 integrating this strategy with the Multi-Subspace-Based Adaptation framework proposed in the
 199 previous section, we derive the final method, AdaMSS.

200 Thanks to the multi-subspace structure of the weights, we can perform adaptive budget allocation
 201 during training by calculating the importance score [28] for each $\mathbf{H}^{(k)} = [\mathbf{B}^{(k)}; (\mathbf{C}^{(k)})^\top] \in$
 202 $\mathbb{R}^{(R_k + n_k) \times r_k}$ directly. Following [28], the importance score $s^{(t)}(\cdot)$ is defined as

$$s^{(t)}(\mathbf{H}^{(k)}) = \frac{1}{r_k(R + n_k)} \sum_{i,j} \bar{I}^{(t)}([\mathbf{H}^{(k)}]_{i,j}) \cdot \bar{U}^{(t)}([\mathbf{H}^{(k)}]_{i,j}), \quad (11)$$

203 where

$$\begin{aligned} \bar{I}^{(t)}([\mathbf{H}^{(k)}]_{i,j}) &= \beta_1 \bar{I}^{(t-1)}([\mathbf{H}^{(k)}]_{i,j}) + (1 - \beta_1) I^{(t)}([\mathbf{H}^{(k)}]_{i,j}), \\ \bar{U}^{(t)}([\mathbf{H}^{(k)}]_{i,j}) &= \beta_2 \bar{U}^{(t-1)}([\mathbf{H}^{(k)}]_{i,j}) + (1 - \beta_2) |I^{(t)}([\mathbf{H}^{(k)}]_{i,j}) - \bar{I}^{(t)}([\mathbf{H}^{(k)}]_{i,j})|, \\ I([\mathbf{H}^{(k)}]_{i,j}) &= |[\mathbf{H}^{(k)}]_{i,j} \nabla_{[\mathbf{H}^{(k)}]_{i,j}} \mathcal{L}|, \end{aligned} \quad (12)$$

206 and t denotes the training step. This importance score measures how consistently and significantly
 207 each $\mathbf{H}^{(k)}$ contributes to loss reduction during training.

208 At each update step, we rank all $\mathbf{H}^{(k)}$ by their importance scores $s^{(t)}(\mathbf{H}^{(k)})$, retain the top- K_t for
 209 training, and freeze the rest. The value of K_t is gradually reduced from K to K_{target} following
 210 a smooth cubic decay schedule over training steps. The adaptive budget allocation mechanism is
 211 stopped when the number of trainable $\mathbf{H}^{(k)}$ for $k = 1, 2, \dots, K$ falls below K_{target} . Figure 3
 212 presents examples on the StanfordCars², FGVC², and CIFAR100² tasks, illustrating the proposed
 213 adaptive budget allocation strategy, which gradually reduces the number of trainable parameters
 214 during the training process.

215 5 Analysis of AdaMSS

216 In the previous section, we introduced AdaMSS that is based on the following reparametrization:

$$\mathbf{W} = \mathbf{W}_0 + \mathbf{ABC}, \quad (13)$$

²<https://huggingface.co/datasets/Multimodal-Fatima>

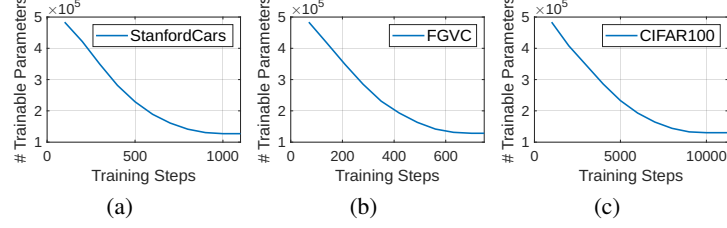


Figure 3: The adaptive change in the number of trainable parameters of AdaMSS during the training process.

where $\mathbf{A} = [\mathbf{A}^{(1)}, \mathbf{A}^{(2)}, \dots, \mathbf{A}^{(K)}] \in \mathbb{R}^{d \times R}$, $\mathbf{B} = \text{diag}(\mathbf{B}^{(1)}, \mathbf{B}^{(2)}, \dots, \mathbf{B}^{(K)}) \in \mathbb{R}^{R \times \sum_{k=1}^K r_k}$,
 $\mathbf{C} = \text{diag}(\mathbf{C}^{(1)}, \mathbf{C}^{(2)}, \dots, \mathbf{C}^{(K)}) \in \mathbb{R}^{\sum_{k=1}^K r_k \times n}$, $(\mathbf{A}^{(k)})^\top \mathbf{A}^{(k)} = \mathbf{I}$ and $\text{rank}(\mathbf{C}^{(k)}) \leq r_k$ for
 $k = 1, 2, \dots, K$.

We establish the upper bound of expected loss by AdaMSS in Lemma 1 and Theorem 1.

As our discussions given in the Appendix B.2, the upbound of the Gaussian complexity for LoRA and its single subspace variants (such as AdaLoRA, PiSSA and LoRA-GA) is given by $\mathcal{L}(\phi) \hat{R} B \sqrt{\frac{\text{rank}(\Delta \mathbf{W}) n}{m}}$ for $\|\Delta \mathbf{W}\|_2 \leq B$. If $\sum_{k=1}^K r_k = \text{rank}(\Delta \mathbf{W})$, applying Cauchy–Schwarz inequality and using $\sum_{k=1}^K n_k = n$, we have $\sum_{k=1}^K \sqrt{r_k n_k} \leq \sqrt{\text{rank}(\Delta \mathbf{W}) n}$. Therefore, AdaMSS attains a lower Gaussian complexity bound when $B_k \leq B$, and can be expected to have a stronger generalization capability than single subspace-based fine-tuning methods, as indicated by Theorem 1. These results can be further extended to deep Lipschitz neural networks by applying Maurer’s Gaussian complexity chain rule [29].

Lemma 1. [Gaussian complexity of the AdaMSS for a shallow Lipschitz neural network] For the class of spectrally bounded shallow Lipschitz network

$$\mathcal{F}_{\text{AdaMSS}} = \{f_{\mathbf{W}}(\mathbf{x}) = \phi(\mathbf{x}\mathbf{W}) \mid \mathbf{W} = \mathbf{W}_0 + \mathbf{A}\mathbf{B}\mathbf{C},$$

$$(\mathbf{A}^{(k)})^\top \mathbf{A}^{(k)} = \mathbf{I}, \|\mathbf{B}^{(k)} \mathbf{C}^{(k)}\|_2 \leq B_k\}, \quad (14)$$

the Gaussian complexity of this function class is upper bounded as follows:

$$\hat{G}_S(\mathcal{F}_{\text{AdaMSS}}) \leq \mathcal{L}(\phi) \sum_{k=1}^K \hat{R} B_k \sqrt{\frac{r_k n_k}{m}},$$

where $n = \sum_{k=1}^K n_k$, n_k denotes the width of the weight matrix $\mathbf{C}^{(k)}$, $\mathcal{L}(\phi)$ is Lipschitz constant for function ϕ , $\mathbf{X} = [\mathbf{x}_1^\top, \dots, \mathbf{x}_m^\top] \in \mathbb{R}^{d \times m}$ for the samples $\{\mathbf{x}_i\}_{i=1}^m$, and $\max_{i=1,2,\dots,m} \|\mathbf{x}_i\| \leq \hat{R}$.

Proof. The corresponding proof is given in the Appendix B.1. \square

Theorem 1. Let $g(\cdot)$ be a $\mathcal{L}(g)$ -Lipschitz loss function from $(f_{\mathbf{w}}(\mathbf{x}), \mathbf{y})$ to $[0, 1]$, where $f_{\mathbf{w}} \in \mathcal{F}_{\text{AdaMSS}}$ and $(\mathbf{x}, \mathbf{y}) \in \mathbb{X} \times \mathbb{Y}$, $\mathbb{X} \subseteq \mathbb{R}^d$ and \mathbb{Y} are feature space and output space, respectively. For any $\delta > 0$, the following holds with probability at least $1 - \delta$ for a randomly chosen i.i.d. samples $\mathbb{S} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^m$:

$$\mathbb{E}[g(f_{\mathbf{w}}(\mathbf{x}), \mathbf{y})] \leq \frac{1}{m} \sum_{i=1}^m g(f_{\mathbf{w}}(\mathbf{x}_i), \mathbf{y}_i) + \sqrt{\mathcal{L}(g) \pi} \hat{G}_S(\mathcal{F}_{\text{AdaMSS}}) + \sqrt{\frac{9 \log \frac{2}{\delta}}{2m}}.$$

Proof. Combining Lemma 1 with Vector-valued Gaussian complexity Generalization Bound Theorem [30], we obtain the generalization bound for AdaMSS. \square

6 Experiments

In this section, we present comparative experiments to evaluate the effectiveness of AdaMSS across tasks in image classification (IC), natural language understanding (NLU), and natural language

generation (NLG). To better understand the contribution of each component, we define a variant of AdaMSS without the proposed adaptive budget allocation mechanism, referred to as AdaMSS_{base}, for clarity of comparison. The hyperparameter R is set to $R = 100$ for both AdaMSS and AdaMSS_{base}. We report the average number of trainable parameters per epoch in our experimental results, as the number of trainable parameters in AdaMSS is adaptively reduced during the training. Since our focus is on low-rank-based adaptation, we compare AdaMSS not only with baseline methods such as full fine-tuning (FF) and linear probing (LP), which only tune the classification head, but also with several low-rank-based PEFT methods, including LoRA [1], DyLoRA [31], AdaLoRA [3], PiSSA [2], LoRETTA [10], and LoRA-PRO [6]. For each method, we prioritize using hyperparameters (e.g., rank r) recommended by the authors. For a fair comparison, we replicate the experimental setups in [1, 2, 24]. A complete list of all hyperparameters and settings is provided in the Appendix G. All experiments were conducted using an RTX A6000 GPU and Python version 3.12.3. The best results are highlighted in **bold**, and the top three results are underlined. **The code for AdaMSS has been included in the submitted supplementary materials. (Ablation studies are given in Appendix F, and additional experimental results are also available in the submitted supplementary materials.)**

6.1 Image Classification

Table 1: The performance of different fine-tuning methods on various image classification datasets for the ViT-Base model and ViT-Large model, averaged over 5 random seeds, where the hyperparameter r_k in AdaMSS_{base} and AdaMSS corresponds to the assumed rank of $\Delta \mathbf{Z}^{(k)}$ and is set to the same value for all k .

Model	Method	# Trainable Parameters	Accuracy (%)							
			OxfordPets	StanfordCars	CIFAR10	EuroSAT	FGVC	RESISC45	CIFAR100	Avg.
ViT-Base	LP	-	90.28 \pm 0.4	25.76 \pm 0.3	96.41 \pm 0.0	88.72 \pm 0.1	17.44 \pm 0.4	74.22 \pm 0.1	84.28 \pm 0.1	68.15
	FF	85.8M	93.14 \pm 0.4	79.78 \pm 1.2	98.92 \pm 0.1	99.05 \pm 0.1	54.84 \pm 1.2	96.13 \pm 0.1	92.38 \pm 0.1	87.74
	LoRA ($r = 16$)	581K	93.19 \pm 0.4	45.38 \pm 0.4	<u>98.78</u> \pm 0.1	98.44 \pm 0.2	25.16 \pm 0.2	92.70 \pm 0.2	92.02 \pm 0.1	77.95
	PiSSA ($r = 8$)	313K	93.84 \pm 0.3	78.43 \pm 0.5	98.74 \pm 0.0	98.67 \pm 0.1	51.56 \pm 1.8	93.81 \pm 1.8	93.31 \pm 0.2	86.90
	PiSSA ($r = 1$)	55K	93.83 \pm 0.1	60.29 \pm 0.3	98.7 \pm 0.0	98.47 \pm 0.1	29.61 \pm 0.2	92.87 \pm 0.2	91.98 \pm 0.2	80.82
	LoRA-PRO	313K	94.03 \pm 0.1	72.12 \pm 0.4	98.77 \pm 0.1	98.65 \pm 0.1	43.39 \pm 0.7	93.66 \pm 0.2	92.54 \pm 0.1	84.74
	LoRETTA ($r = 5$)	57K	93.39 \pm 0.4	74.15 \pm 0.8	98.73 \pm 0.1	<u>98.67</u> \pm 0.1	48.86 \pm 0.5	93.36 \pm 0.1	91.87 \pm 0.1	85.57
	AdaMSS _{base} ($r_k = 3$)	125K	94.02 \pm 0.3	81.55 \pm 0.4	98.81 \pm 0.05	98.7 \pm 0.06	56.95 \pm 0.6	94.18 \pm 0.1	92.13 \pm 0.1	88.05
	AdaMSS ($r_k = 3$)	59K	94.23 \pm 0.1	<u>80.44</u> \pm 0.2	98.69 \pm 0.04	98.59 \pm 0.09	<u>54.45</u> \pm 0.3	<u>94.03</u> \pm 0.2	91.91 \pm 0.1	<u>87.47</u>
	AdaMSS _{base} ($r_k = 1$)	42K	93.91 \pm 0.2	78.98 \pm 0.2	98.71 \pm 0.07	98.64 \pm 0.07	53.2 \pm 0.4	93.62 \pm 0.08	91.90 \pm 0.1	86.99
ViT-Large	LP	-	91.11 \pm 0.3	37.91 \pm 0.3	97.78 \pm 0.0	92.64 \pm 0.1	24.62 \pm 0.2	82.02 \pm 0.1	84.28 \pm 0.1	72.91
	FF	303.3M	94.43 \pm 0.6	88.90 \pm 0.3	99.15 \pm 0.1	99.04 \pm 0.1	68.25 \pm 1.6	96.43 \pm 0.1	93.58 \pm 0.2	91.40
	LoRA ($r = 16$)	1.57M	<u>94.82</u> \pm 0.1	73.25 \pm 0.4	<u>99.13</u> \pm 0.0	98.63 \pm 0.1	42.32 \pm 1.0	94.71 \pm 0.3	94.87 \pm 0.1	85.40
	PiSSA ($r = 8$)	835K	94.04 \pm 0.4	84.19 \pm 0.7	99.13 \pm 0.0	98.79 \pm 0.0	59.81 \pm 0.6	94.99 \pm 0.2	92.42 \pm 0.1	89.05
	PiSSA ($r = 1$)	147K	93.98 \pm 0.3	83.04 \pm 0.3	<u>99.04</u> \pm 0.0	98.71 \pm 0.0	56.72 \pm 0.6	94.64 \pm 0.2	93.25 \pm 0.1	88.48
	LoRA-PRO	835K	94.67 \pm 0.1	83.57 \pm 0.3	99.20 \pm 0.1	98.81 \pm 0.1	57.71 \pm 0.5	95.12 \pm 0.1	93.53 \pm 0.1	88.94
	LoRETTA ($r = 5$)	132K	78.28 \pm 0.3	68.44 \pm 0.3	88.80 \pm 0.2	98.68 \pm 0.1	58.04 \pm 0.9	94.53 \pm 0.1	93.28 \pm 0.1	82.86
	AdaMSS _{base} ($r_k = 3$)	483K	<u>94.74</u> \pm 0.1	<u>85.40</u> \pm 0.3	99.11 \pm 0.0	98.93 \pm 0.06	65.30 \pm 0.6	95.32 \pm 0.1	93.51 \pm 0.1	90.33
	AdaMSS ($r_k = 3$)	241K	94.87 \pm 0.1	85.24 \pm 0.3	99.12 \pm 0.0	98.93 \pm 0.1	<u>64.31</u> \pm 0.4	<u>95.2</u> \pm 0.2	93.22 \pm 0.1	<u>90.13</u>
	AdaMSS _{base} ($r_k = 1$)	178K	94.58 \pm 0.1	83.71 \pm 0.2	99.08 \pm 0.05	98.85 \pm 0.1	59.27 \pm 0.8	94.68 \pm 0.3	93.43 \pm 0.2	89.09

We evaluate all methods on IC using the widely adopted Vision Transformer (ViT) [12], a prevalent foundation model in computer vision, across seven public datasets: OxfordPets³, StanfordCars², CIFAR10², EuroSAT⁴, FGVC², RESISC45⁵, and CIFAR100². The number of training epochs is set as 10. As shown in Table 1, AdaMSS achieves higher average accuracy than other PEFT methods, including LP, LoRA, PiSSA, LoRA-PRO, and LoRETTA, and performs comparable accuracy with FF on the seven image classification datasets while using significantly fewer trainable parameters. More specifically, with the ViT-Base model, AdaMSS (59K) achieves average accuracy comparable to PiSSA ($r = 8$), using only 19.5% of its trainable parameters, and achieves 6.5% higher accuracy than PiSSA ($r = 1$) at a similar parameter budget. This is because that AdaMSS’s multi-subspace design allows it to capture richer features with fewer parameters. On the ViT-large model, AdaMSS achieves 4.7% higher average accuracy than LoRA, while requiring only with 15.4% of the trainable parameters, demonstrating AdaMSS’s strong parameter efficiency and transfer performance in image classification tasks.

6.2 Natural Language Understanding

We evaluate all methods on the General Language Understanding Evaluation (GLUE) benchmark [32] using the robustly optimized BERT models, *i.e.*, RoBERTa-Base and RoBERTa-Large [15],

³<https://huggingface.co/datasets/timm/oxford-iiit-pet>

⁴<https://huggingface.co/datasets/timm/eurosat-rgb>

⁵<https://huggingface.co/datasets/timm/resisc45>

Table 2: The performance of different fine-tuning methods on six datasets of the GLUE benchmark for the RoBERTa-Base model and RoBERTa-Large model, averaged over 5 random seeds.

Model	Method	# Trainable Parameters	SST-2	MRPC	CoLA	QNLI	RTE	STS-B	Avg.
			Acc.	Acc.	MCC	Acc.	Acc.	PCC	
RoBERTa-Base	FF	125M	94.8	90.2	63.6	92.8	78.7	91.2	85.2
	LoRA	0.3M	95.1 \pm 0.2	89.7 \pm 0.7	63.4 \pm 1.2	93.3 \pm 0.3	78.4 \pm 0.8	91.5 \pm 0.2	85.2
	AdaLoRA	0.3M	94.5 \pm 0.2	88.7 \pm 0.5	62.0 \pm 0.6	93.1 \pm 0.2	81.0 \pm 0.6	90.5 \pm 0.2	85.0
	DyLoRA	0.3M	94.3 \pm 0.5	89.5 \pm 0.5	61.1 \pm 0.3	92.2 \pm 0.5	78.7 \pm 0.7	91.1 \pm 0.6	84.5
	PiSSA ($r = 8$)	0.3M	93.9 \pm 0.1	89.3 \pm 0.8	62.1 \pm 2.9	91.3 \pm 0.1	77.3 \pm 1.4	90.5 \pm 0.2	84.1
	LoRA-PRO	0.3M	94.2 \pm 0.3	90.1 \pm 0.5	64.3 \pm 0.72	92.0 \pm 0.2	80.2 \pm 1.8	90.9 \pm 0.22	85.3
	LoRETTA	0.057M	94.6 \pm 0.5	88.3 \pm 0.7	61.8 \pm 1.3	92.7 \pm 0.2	75.1 \pm 5.3	90.5 \pm 0.1	83.8
	AdaMSS _{base} ($r_k = 1$)	0.042M	94.6 \pm 0.2	89.2 \pm 1.0	64.3 \pm 0.9	92.4 \pm 0.1	77.2 \pm 0.7	90.6 \pm 0.1	84.7
	AdaMSS ($r_k = 1$)	0.032M	94.6 \pm 0.2	88.8 \pm 1.4	64.5 \pm 1.1	92.4 \pm 0.1	77.3 \pm 0.7	90.4 \pm 0.1	84.7
RoBERTa-Large	FF	356M	96.4	90.9	68	94.7	86.6	92.4	88.2
	LoRA	0.8M	96.2 \pm 0.5	90.2 \pm 1.0	68.2 \pm 1.9	94.8 \pm 0.3	85.2 \pm 1.1	92.3 \pm 0.5	87.8
	PiSSA ($r = 8$)	0.8M	95.5 \pm 0.2	86.9 \pm 2.6	61.1 \pm 3.4	92.1 \pm 1.7	56.8 \pm 8.2	91.8 \pm 0.4	80.7
	LoRA-PRO	0.8M	95.9 \pm 0.2	90.9 \pm 0.4	66.7 \pm 2.0	93.0 \pm 0.5	60.5 \pm 13.5	92.0 \pm 0.1	83.2
	LoRETTA	0.132M	96.2 \pm 0.2	90.5 \pm 0.4	69.5 \pm 0.6	94.1 \pm 0.9	53.0 \pm 0.5	92.0 \pm 0.2	82.6
	AdaMSS _{base} ($r_k = 1$)	0.097M	96.3 \pm 0.2	90.5 \pm 0.3	68.0 \pm 0.9	94.6 \pm 0.1	87.3 \pm 1.0	92.0 \pm 0.0	88.1
	AdaMSS ($r_k = 1$)	0.045M	96.1 \pm 0.0	90.3 \pm 0.5	67.2 \pm 1.2	94.5 \pm 0.1	87.1 \pm 2.1	91.9 \pm 0.0	87.9

for the evaluation. We evaluate the performance of the fine-tuned models using three key metrics: Matthew’s correlation coefficient (MCC) for CoLA, Pearson correlation coefficient (PCC) for STS-B, and accuracy (Acc.) for all other tasks. For all methods, the maximum number of training epochs is set to 100 and select the best epoch for each run. Table 2 presents the results of all methods. As shown in the Table 2, on RoBERTa-Base, AdaMSS achieves average accuracy comparable to other low-rank-based PEFT methods while using the fewest trainable parameters. For RoBERTa-Large, AdaMss outperforms PiSSA, LoRA-PRO, and LoRETTA by around 5% in average accuracy while using fewer parameters, and achieves performance comparable to both LoRA and FF.

6.3 Natural Language Generation

In this subsection, we use LLaMA 2-7B [14] and compare AdaMSS with the traditional low-rank-based methods, including PiSSA, LoRA, LoRETTA, LoRA-PRO, and FF on natural language generation (NLG) tasks. As shown in Table 3, AdaMSS achieves competitive or even superior performance on NLG tasks while requiring significantly fewer trainable parameters. For instance, on both the GSM8K and MATH datasets, AdaMSS outperforms LoRA, PiSSA, and LoRETTA, despite using only 0.9M trainable parameters.

Table 3: Comparing different methods on NLG tasks

Model	Method	Trainable Parameters	GSM8K [13]	MATH [33]
LLaMA 2-7B	Full FT	6738M	49.05	7.22
	LoRA	320M	42.30	5.50
	PiSSA	19M	44.11	5.84
	LoRA-PRO	19M	46.61	6.4
	LoRETTA	0.3M	37.86	4.6
	AdaMSS _{base} ($r_k = 3$)	0.9M	44.50	6.04

7 Conclusions and Future Works

In this work, inspired by our observation of the multi-subspace structure in network weights, we proposed AdaMSS, a novel PEFT approach, to address the limitations of single-subspace methods that often struggle with the dilemma between limited expressiveness and parameter efficiency. Compared to low-rank-based fine-tuning methods, AdaMSS enables more expressive and adaptive parameter updates by leveraging subspace segmentation, thereby achieving both parameter efficiency and strong generalization capability. Comprehensive theoretical analysis and extensive empirical evaluations demonstrate the advantages of AdaMSS. In future work, we aim to extend the multi-subspace perspective beyond fine-tuning to broader areas such as network compression and pruning. We believe that subspace segmentation offers a promising direction for learning compact, disentangled, and efficient representations for network weights, which could benefit model compression and pruning in deep neural networks.

References

- [1] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- [2] Fanxu Meng, Zhaohui Wang, and Muhan Zhang. Pissa: principal singular values and singular vectors adaptation of large language models. In *Neural Information Processing Systems*, 2024.
- [3] Qingru Zhang, Minshuo Chen, Alexander Bukharin, Pengcheng He, Yu Cheng, Weizhu Chen, and Tuo Zhao. Adaptive budget allocation for parameter-efficient fine-tuning. In *International Conference on Learning Representations*, 2023.
- [4] Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: efficient finetuning of quantized llms. *Neural Information Processing Systems*, 36, 2024.
- [5] Shaowen Wang, Linxi Yu, and Jian Li. Lora-ga: Low-rank adaptation with gradient approximation. In *Neural Information Processing Systems*, 2024.
- [6] Zhengbo Wang, Jian Liang, Ran He, Zilei Wang, and Tieniu Tan. Lora-pro: Are low-rank adapters properly optimized? 2025.
- [7] Mahdi Nikdan, Soroush Tabesh, Elvir Crnčević, and Dan Alistarh. Rosa: accurate parameter-efficient fine-tuning via robust adaptation. In *International Conference on Machine Learning*, pages 38187–38206, 2024.
- [8] Shubhankar Borse, Shreya Kadambi, Nilesh Pandey, Kartikeya Bhardwaj, Viswanath Ganapathy, Sweta Priyadarshi, Risheek Garrepalli, Rafael Esteves, Munawar Hayat, and Fatih Porikli. Foura: Fourier low-rank adaptation. In *Neural Information Processing Systems*, 2024.
- [9] Yang Li, Shaobo Han, and Shihao Ji. Vb-lora: Extreme parameter efficient fine-tuning with vector banks. In *Neural Information Processing Systems*, 2024.
- [10] Yifan Yang, Jiajun Zhou, Ngai Wong, and Zheng Zhang. Loretta: Low-rank economic tensor-train adaptation for ultra-low-parameter fine-tuning of large language models. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3161–3176, 2024.
- [11] Daniel Bershtatsky, Daria Cherniuk, Talgat Daulbaev, and Ivan Oseledets. Lotr: low tensor rank weight adaptation. *arXiv preprint arXiv:2402.01376*, 2024.
- [12] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [13] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems, 2021.
- [14] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models, 2023.
- [15] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: a robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- [16] Guangcan Liu, Zhouchen Lin, Shuicheng Yan, Ju Sun, Yong Yu, and Yi Ma. Robust recovery of subspace structures by low-rank representation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(1):171–184, 2012.

- [17] Pan Zhou, Canyi Lu, and Jiashi Feng. Tensor low-rank representation for data recovery and clustering. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.
- [18] Shiqiang Du, Baokai Liu, Guangrong Shan, Yuqing Shi, and Weilan Wang. Enhanced tensor low-rank representation for clustering and denoising. *Knowledge-Based Systems*, 243:108468, 2022.
- [19] Hongyang Zhang, Zhouchen Lin, Chao Zhang, and Junbin Gao. Robust latent low rank representation for subspace clustering. *Neurocomputing*, 145:369–373, 2014.
- [20] Xiang Lisa Li and Percy Liang. Prefix-tuning: optimizing continuous prompts for generation. *arXiv preprint arXiv:2101.00190*, 2021.
- [21] Xiao Liu, Kaixuan Ji, Yicheng Fu, Weng Lam Tam, Zhengxiao Du, Zhilin Yang, and Jie Tang. P-tuning v2: prompt tuning can be comparable to fine-tuning universally across scales and tasks, 2022.
- [22] Neil Houlsby, Andrei Giurghi, Stanislaw Jastrzebski, Bruna Morrone, Quentin de Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for nlp, 2019.
- [23] Elad Ben Zaken, Shauli Ravfogel, and Yoav Goldberg. Bitfit: simple parameter-efficient fine-tuning for transformer-based masked language-models, 2022.
- [24] Ziqi Gao, Qichao Wang, Aochuan Chen, Zijing Liu, Bingzhe Wu, Liang Chen, and Jia Li. Parameter-efficient fine-tuning with discrete fourier transform. In *International Conference on Machine Learning*, 2024.
- [25] Zeju Qiu, Weiyang Liu, Haiwen Feng, Yuxuan Xue, Yao Feng, Zhen Liu, Dan Zhang, Adrian Weller, and Bernhard Schölkopf. Controlling text-to-image diffusion by orthogonal finetuning. In *Neural Information Processing Systems*, 2024.
- [26] Massimo Bini, Karsten Roth, Zeynep Akata, and Anna Khoreva. Ether: efficient finetuning of large-scale models with hyperplane reflections. In *International Conference on Machine Learning*, pages 4007–4026, 2024.
- [27] Frank Lauren Hitchcock. The expression of a tensor or a polyadic as a sum of products. *Journal of Mathematical Physics*, (6):164–189, 1927.
- [28] Qingru Zhang, Simiao Zuo, Chen Liang, Alexander Bukharin, Pengcheng He, Weizhu Chen, and Tuo Zhao. Platon: Pruning large transformer models with upper confidence bound of weight importance. In *International conference on machine learning*, pages 26809–26823. PMLR, 2022.
- [29] Andreas Maurer. A chain rule for the expected suprema of gaussian processes. *Theoretical Computer Science*, 650:109–122, 2016.
- [30] Andrea Pinto, Akshay Rangamani, and Tomaso Poggio. On generalization bounds for neural networks with low rank layers. *arXiv preprint arXiv:2411.13733*, 2024.
- [31] Mojtaba Valipour, Mehdi Rezagholizadeh, Ivan Kobyzev, and Ali Ghodsi. Dylora: parameter efficient tuning of pre-trained models using dynamic search-free low-rank adaptation. *arXiv preprint arXiv:2210.07558*, 2022.
- [32] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. Glue: a multi-task benchmark and analysis platform for natural language understanding. *arXiv preprint arXiv:1804.07461*, 2018.
- [33] Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T. Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. Metamath: Bootstrap your own mathematical questions for large language models, 2024.
- [34] Roman Vershynin. High-dimensional probability. 2020.
- [35] Jianbo Shi and J. Malik. Normalized cuts and image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(8):888–905, 2000.

401 A The Proof of Property 1

402 *Proof.* From $\mathbf{Z}_0^* = [\mathbf{V}_0]_{:,1:R} [\mathbf{V}_0]_{:,1:R}^\top$, we know $\text{rank}(\mathbf{Z}_0^*) = \text{rank}(\hat{\mathbf{W}}_0)$.

Since $\hat{\mathbf{W}}_0^{(k)} = \hat{\mathbf{W}}_0^{(k)} (\mathbf{Z}_0^*)^{(k)}$, it follows that $\text{rank}(\hat{\mathbf{W}}_0^{(k)}) \leq \text{rank}((\mathbf{Z}_0^*)^{(k)})$. Thus, we have

$$\sum_{k=1}^K \text{rank}(\hat{\mathbf{W}}_0^{(k)}) \leq \sum_{k=1}^K \text{rank}((\mathbf{Z}_0^*)^{(k)}) \leq \sum_{k=1}^K \text{rank}([\mathbf{Z}_0^*]^{(k)}) = \text{rank}(\mathbf{Z}_0^*) = \text{rank}(\hat{\mathbf{W}}_0).$$

On the other hand, we know that $\text{rank}(\hat{\mathbf{W}}_0) \leq \sum_{k=1}^K \text{rank}(\hat{\mathbf{W}}_0^{(k)})$. Therefore, we can conclude

$$\sum_{k=1}^K \text{rank}(\hat{\mathbf{W}}_0^{(k)}) = \sum_{k=1}^K \text{rank}((\mathbf{Z}_0^*)^{(k)}) = \text{rank}(\hat{\mathbf{W}}_0).$$

Since $\text{rank}(\hat{\mathbf{W}}_0^{(k)}) \leq \text{rank}((\mathbf{Z}_0^*)^{(k)})$ for all k and $\sum_{k=1}^K \text{rank}(\hat{\mathbf{W}}_0^{(k)}) = \sum_{k=1}^K \text{rank}((\mathbf{Z}_0^*)^{(k)})$, it must hold that

$$\text{rank}(\hat{\mathbf{W}}_0^{(k)}) = \text{rank}((\mathbf{Z}_0^*)^{(k)}).$$

403

□

404 B The Proofs of the Results in Section 5

405 B.1 The Proof of Lemma 1

406 *Proof.* By Talagrand's contraction lemma [34], we have

$$\hat{G}_S(\mathcal{F}_{\text{AdaMSS}}) \leq \mathcal{L}(\phi) \frac{1}{m} \mathbb{E}_\gamma \sup_{\mathbf{W}} \langle \mathbf{W}, \mathbf{X}\mathbf{\Gamma} \rangle, \quad (15)$$

407 where each entry in $\mathbf{\Gamma} \in \mathbb{R}^{m \times n}$, i.e., $[\mathbf{\Gamma}]_{ij}$, follows standard Gaussian distribution.

408 Since

$$\begin{aligned} \mathbb{E}_\gamma [\sup_{\mathbf{W}} \langle \mathbf{W}, \mathbf{X}\mathbf{\Gamma} \rangle] &= \mathbb{E}_\gamma [\langle \mathbf{W}_0, \mathbf{X}\mathbf{\Gamma} \rangle] + \mathbb{E}_\gamma [\sup_{\mathbf{B}, \mathbf{C}} \langle \mathbf{A}\mathbf{B}\mathbf{C}, \mathbf{X}\mathbf{\Gamma} \rangle] \\ &= \mathbb{E}_\gamma [\sup_{\mathbf{B}, \mathbf{C}} \langle \mathbf{A}\mathbf{B}\mathbf{C}, \mathbf{X}\mathbf{\Gamma} \rangle], \end{aligned} \quad (16)$$

409 we only need to study the upper bound of $\mathbb{E}_\gamma [\sup_{\mathbf{B}, \mathbf{C}} \langle \mathbf{A}\mathbf{B}\mathbf{C}, \mathbf{X}\mathbf{\Gamma} \rangle]$.

410 We divide $\mathbf{\Gamma}$ into K blocks according $\{n_k\}_{k=1}^K$, i.e., $\mathbf{\Gamma} = [\mathbf{\Gamma}^{(1)}, \mathbf{\Gamma}^{(2)}, \dots, \mathbf{\Gamma}^{(K)}]$ for $\mathbf{\Gamma}^{(k)} \in \mathbb{R}^{m \times n_k}$ ($k = 1, 2, \dots, K$). From the Cauchy-Schwarz inequality, we have

$$\begin{aligned} \sup_{\mathbf{B}, \mathbf{C}} \langle \mathbf{A}\mathbf{B}\mathbf{C}, \mathbf{X}\mathbf{\Gamma} \rangle &= \sum_{k=1}^K \sup_{\mathbf{B}^{(k)}, \mathbf{C}^{(k)}} \langle \mathbf{A}^{(k)} \mathbf{B}^{(k)} \mathbf{C}^{(k)}, \mathbf{X}\mathbf{\Gamma}^{(k)} \rangle \leq \sum_{k=1}^K \|\mathbf{B}^{(k)} \mathbf{C}^{(k)}\|_F \|\mathbf{X}\mathbf{\Gamma}^{(k)}\|_F \\ &\leq \sum_{k=1}^K \sqrt{r_k} B_k \|\mathbf{X}\mathbf{\Gamma}^{(k)}\|_F. \end{aligned} \quad (17)$$

412 Using the fact that $\mathbb{E}[|Y|] \leq \sqrt{\mathbb{E}[Y^2]}$ for any random variable Y , we obtain

$$\mathbb{E}_\gamma [\|\mathbf{X}\mathbf{\Gamma}^{(k)}\|_F] \leq \sqrt{\mathbb{E}_\gamma [\|\mathbf{X}\mathbf{\Gamma}^{(k)}\|_F^2]} = \sqrt{\text{Tr}(\mathbf{X} \mathbf{X}^\top \mathbb{E}_\gamma [(\mathbf{\Gamma}^{(k)})^\top \mathbf{\Gamma}^{(k)}])} \leq \hat{R} \sqrt{n_k m}. \quad (18)$$

413 Combining (17) and (18), we obtain

$$\mathbb{E}_\gamma [\sup_{\mathbf{W}} \langle \mathbf{W}, \mathbf{X}\mathbf{\Gamma} \rangle] = \mathbb{E}_\gamma \sup_{\mathbf{B}, \mathbf{C}} \langle \mathbf{A}\mathbf{B}\mathbf{C}, \mathbf{X}\mathbf{\Gamma} \rangle \leq \sum_{k=1}^K \sqrt{r_k} B_k \mathbb{E}_\gamma [\|\mathbf{X}\mathbf{\Gamma}^{(k)}\|_F] \leq \sum_{k=1}^K \sqrt{r_k} B_k \hat{R} \sqrt{n_k m}. \quad (19)$$

Substituting (19) into (15) leads to

$$\hat{G}_S(\mathcal{F}_{\text{AdaMSS}}) \leq \mathcal{L}(\phi) \frac{1}{m} \mathbb{E}_\gamma \sup_{\mathbf{W}} \langle \mathbf{W}, \mathbf{\Gamma} \mathbf{X} \rangle \leq \mathcal{L}(\phi) \sum_{k=1}^K \hat{R} B_k \sqrt{\frac{r_k n_k}{m}}.$$

414

□

415 B.2 Gaussian Complexity Analysis for LoRA and Its Single Subspace Variants

416 In this subsection, we analyze the Gaussian complexity of LoRA, as presented in Lemma 2.

417 **Lemma 2.** *[Gaussian complexity of the LoRA for a shallow Lipschitz neural network] For the class*
418 *of spectrally bounded shallow Lipschitz network*

$$\mathcal{F}_{\text{LoRA}} = \{f_{\mathbf{W}}(\mathbf{x}) = \phi(\mathbf{x}\mathbf{W}) \mid \mathbf{W} = \mathbf{W}_0 + \Delta\mathbf{W}, \|\Delta\mathbf{W}\|_2 \leq B\} \quad (20)$$

419 the Gaussian complexity of this function class is upper bounded as follows:

$$\hat{G}_S(\mathcal{F}_{\text{LoRA}}) \leq \mathcal{L}(\phi) \hat{R} B \sqrt{\frac{rn}{m}}, \quad (21)$$

420 where $r = \text{rank}(\Delta\mathbf{W})$, n denotes the width of the weight matrix $\Delta\mathbf{W}$, $\mathcal{L}(\phi)$ is Lipschitz constant for
421 function ϕ , $\mathbf{X} = [\mathbf{x}_1^\top, \dots, \mathbf{x}_m^\top] \in \mathbb{R}^{d \times m}$ for the samples $\{\mathbf{x}_i\}_{i=1}^m$, and $\max_{i=1,2,\dots,m} \|\mathbf{x}_i\| \leq \hat{R}$.

422 *Proof.* By Talagrand's contraction lemma [34], we have

$$\hat{G}_S(\mathcal{F}_{\text{LoRA}}) \leq \mathcal{L}(\phi) \frac{1}{m} \mathbb{E}_\gamma \sup_{\Delta\mathbf{W}} \langle \mathbf{W}_0 + \Delta\mathbf{W}, \mathbf{X}\mathbf{\Gamma} \rangle = \mathcal{L}(\phi) \frac{1}{m} \mathbb{E}_\gamma \sup_{\Delta\mathbf{W}} \langle \Delta\mathbf{W}, \mathbf{X}\mathbf{\Gamma} \rangle, \quad (22)$$

423 where each entry in $\mathbf{\Gamma} \in \mathbb{R}^{m \times n}$ is drawn independently from the standard Gaussian distribution.

424 Applying the Cauchy-Schwarz inequality and the inequality $\mathbb{E}[|Y|] \leq \sqrt{\mathbb{E}[Y^2]}$, we obtain

$$\begin{aligned} \mathbb{E}_\gamma \sup_{\Delta\mathbf{W}} \langle \Delta\mathbf{W}, \mathbf{X}\mathbf{\Gamma} \rangle &\leq \sqrt{r} B \mathbb{E}_\gamma [\|\mathbf{X}\mathbf{\Gamma}\|_F] \leq \sqrt{r} B \sqrt{\mathbb{E}_\gamma [\|\mathbf{X}\mathbf{\Gamma}\|_F^2]} \\ &\leq \sqrt{r} B \sqrt{\text{Tr}(\mathbf{X}\mathbf{X}^\top \mathbb{E}_\gamma[(\mathbf{\Gamma})^\top \mathbf{\Gamma}])} \leq \sqrt{r} B \hat{R} \sqrt{nm}. \end{aligned} \quad (23)$$

425 Substituting (23) into (22) yields the desired result: $\hat{G}_S(\mathcal{F}_{\text{LoRA}}) \leq \mathcal{L}(\phi) \hat{R} B \sqrt{\frac{rn}{m}}$.

426

□

427 Using similar proofs, we can obtain the same upper bound in (21) for the Gaussian complexity of
428 PiSSA and other single subspace variants for same samples number.

429 C Initialization of AdaMSS (Algorithm 1) and Subspace Segmentation 430 (Algorithm 2)

431 D Parameter Count

Table 4: Comparison of parameter counts among different low-rank-based methods for $d = n$: r represents the rank of the incremental update used in LoRA and PiSSA.

	LoRA	PiSSA	AdaMSS
# Trainable Parameters	$2rn$	$2rn$	$\sum_{k=1}^K (r_k n_k + r_k R_k)$

432 Given $\mathbf{W}^{(k)} = \mathbf{W}_0^{(k)} + \mathbf{A}^{(k)} \mathbf{B}^{(k)} \mathbf{C}^{(k)}$ for $k = 1, 2, \dots, K$, where $\mathbf{B}^{(k)} \in \mathbb{R}^{R_k \times r_k}$ and $\mathbf{C}^{(k)} \in$
433 $\mathbb{R}^{r_k \times n_k}$ are trainable, the total trainable parameter count in AdaMSS is given by $\sum_{k=1}^K (r_k n_k + r_k R_k)$,
434 which is always less than $2(\sum_{k=1}^K r_k) \max_{k=1,2,\dots,K} n_k$ for $d = n$. Compared to the trainable
435 parameter count in LoRA and PiSSA, as shown in Table 4, AdaMSS introduces significantly fewer
436 parameters for $r = \sum_{k=1}^K r_k$, owing to the fact that $r_k \leq R_k \leq n_k$ and $n = \sum_{k=1}^K n_k$.

Algorithm 1 Initialization of AdaMSS

Input: W_0 , R , and r_k .**Output:** Estimated K , $\{A^{(k)}\}_{k=1}^K$, and the initialized values of $\{B^{(k)}\}_{k=1}^K$ and $\{C^{(k)}\}_{k=1}^K$.**Step 1:** Compute the skinny SVD of W_0 : $W_0 = U_0 S_0 V_0^\top$;**Step 2:** Obtain \hat{W}_0 by $\hat{W}_0 = [U_0]_{:,1:R} [S_0]_{1:R,1:R} [V_0]_{:,1:R}^\top$ for given R ;**Step 3:** Obtain Z_0^* by $Z_0^* = [V_0]_{:,1:R} [V_0]_{:,1:R}^\top$; \triangleright Steps 1-3 are derived from Section 3.1.**Step 4:** Estimate the number of subspaces K and apply clustering to assign the columns of W_0 to K subspaces [16], which are detailed in Algorithm 2;**Step 5:** Obtain the matrix blocks $\{W_0^{(k)}\}_{k=1}^K$, $\{\hat{W}_0^{(k)}\}_{k=1}^K$, and $\{(Z_0^*)^{(k)}\}_{k=1}^K$ by the segmentation results; \triangleright Steps 4-5 are derived from Section 3.2.**Step 6:** For each k , compute the skinny SVD of $\hat{W}_0^{(k)}$: $\hat{W}_0^{(k)} = U_{\hat{W}_0}^{(k)} S_{\hat{W}_0}^{(k)} (V_{\hat{W}_0}^{(k)})^\top$;**Step 7:** Perform the QR decomposition of $S_{\hat{W}_0}^{(k)} (V_{\hat{W}_0}^{(k)})^\top (Z_0^*)^{(k)}$: $S_{\hat{W}_0}^{(k)} (V_{\hat{W}_0}^{(k)})^\top (Z_0^*)^{(k)} = Q^{(k)} R^{(k)}$;**Step 8:** Set $A^{(k)} = U_{\hat{W}_0}^{(k)}$, initialize $B^{(k)} = Q^{(k)}$ and $C^{(k)} = 0$. \triangleright Steps 6-8 are derived from Section 3.3.

Algorithm 2 Estimation of K and Subspace Segmentation [16]

Input: V_0 , R , K_0 , and $\tau > 0$.**Step 1:** Construct the affinity matrix $[F]_{i,j} = ([\tilde{U}\tilde{U}^\top]_{i,j})$, where \tilde{U} is formed by $[V_0]_{:,1:R}$ with normalized rows;**Step 2:** Estimate the number of subspaces K by $K = \max(K_0, n - \text{int}(\sum_{i=1}^n f_\tau(\sigma_i(L)))$ [16], where $\{\sigma_i(L)\}_{i=1}^n$ are the singular values of the Laplacian matrix $L = I - D^{-\frac{1}{2}} F D^{-\frac{1}{2}}$, where $D = \text{diag}(\sum_j [F]_{1j}, \dots, \sum_j [F]_{nj})$, $\text{int}(\cdot)$ is the function of the nearest integer, and

$$f_\tau(\sigma) = \begin{cases} 1, & \text{if } \sigma \geq \tau, \\ \log_2 \left(1 + \frac{\sigma^2}{\tau^2} \right), & \text{otherwise.} \end{cases}$$

Step 3: Construct an undirected graph by using the affinity matrix F ;**Step 4:** Apply the NCut [35] to segment the vertices into K clusters;437 **E Computation Complexity**

438 By the proposed adaptation, we have $Y = X(W_0 + ABC)$, where $X \in \mathbb{R}^{\text{batch size} \times n}$ is the input.
 439 In table 5, we compare the computational complexity of gradient computation for LoRA, PiSSA, and
 440 AdaMSS. The table shows that AdaMSS's computational complexity is on the same order as LoRA
 441 and PiSSA when $\sum_{k=1}^K r_k = r$.

Table 5: Comparison of LoRA, PiSSA and AdaMSS for $d = n$, where underline denotes the trainable parameters, and $A \in \mathbb{R}^{n \times r}$ and $B \in \mathbb{R}^{r \times n}$ for LoRA and PiSSA.

	LoRA	PiSSA	AdaMSS
Forward	$Y = XW_0 + XAB$	$Y = XW_{res} + XAB$	$Y = XW_0 + XABC$
Gradient	$\frac{\partial \mathcal{L}}{\partial A} = X^\top \left(\frac{\partial \mathcal{L}}{\partial Y} \right) B^\top$	$\frac{\partial \mathcal{L}}{\partial A} = X^\top \left(\frac{\partial \mathcal{L}}{\partial Y} \right) B^\top$	$\frac{\partial \mathcal{L}}{\partial B} = (XA)^\top \left(\frac{\partial \mathcal{L}}{\partial Y} \right) C^\top$
	$\frac{\partial \mathcal{L}}{\partial B} = A^\top X^\top \left(\frac{\partial \mathcal{L}}{\partial Y} \right)$	$\frac{\partial \mathcal{L}}{\partial B} = A^\top X^\top \left(\frac{\partial \mathcal{L}}{\partial Y} \right)$	$\frac{\partial \mathcal{L}}{\partial C} = B^\top (XA)^\top \left(\frac{\partial \mathcal{L}}{\partial Y} \right)$
Cost	$\mathcal{O}(rn^2 + \text{batch size} \times n^2)$	$\mathcal{O}(rn^2 + \text{batch size} \times n^2)$	$\mathcal{O}(\sum_{k=1}^K r_k n^2 + \text{batch size} \times n^2)$

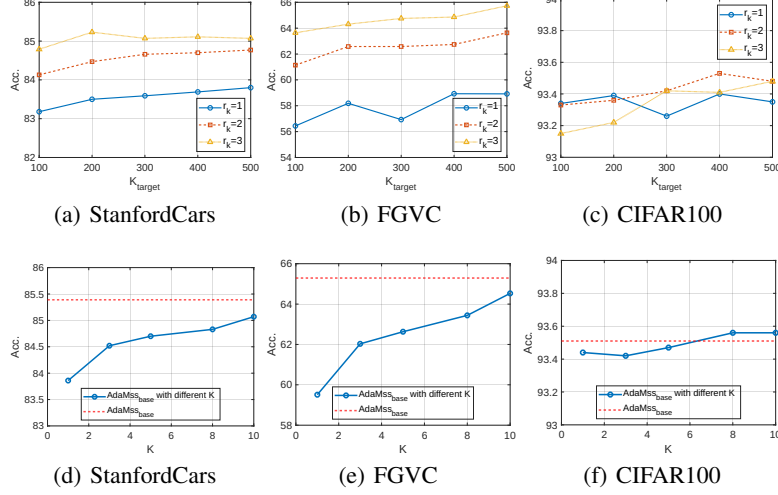


Figure 4: The performance of AdaMSS with different values of r_k , K , and K_{target} .

442 F Ablation Study

443 In the AdaMSS_{base} and AdaMSS, the key hyperparameters include r_k , K , and K_{target} , where K
 444 represents the number of subspaces for each layer, and K_{target} is used for the proposed adaptive
 445 budget allocation mechanism. This subsection explores how the performance of the proposed methods
 446 (AdaMSS and AdaMSS_{base}) varies with different values of r_k , K and K_{target} using ViT-Large.
 447 Each configuration is evaluated on the StanfordCars², FGVC datasets², and CIFAR100². All results
 448 are presented in Figures 4 (a)–(f).

449 Figures 4 (a)–(c) illustrate the impact of varying r_k and K_{target} on the performance of AdaMSS, where
 450 $r_k \in \{1, 2, 3\}$, $K_{\text{target}} \in \{100, 200, 300, 400, 500\}$, and K is estimated according to Algorithm
 451 2 (see Appendix C). Figures 4 (d)–(f) investigate the effect of varying K on the performance of
 452 AdaMSS_{base} at $r_k = 3$, comparing estimated K (shown as red dashed lines) with given values
 453 of K . As shown in Figures 4, increasing both r_k and K leads to significant improvements in
 454 accuracy. Notably, the improvement associated with increasing K highlights the advantage of the
 455 multi-subspace approach.

Table 6: Hyper-parameter setup for AdaMSS_{base} and AdaMSS for IC.

Model	Hyperparameter	OxfordPets	StanfordCars	CIFAR10	EuroSAT	FGVC	RESISC45	CIFAR100
Both	Optimizer	AdamW						
	Batch Size	10						
	Epochs	10						
	Seeds	{7,77,777,7777,77777}						
ViT-Base	Learning Rate	0.005	0.01	0.01	0.01	0.01	0.01	0.01
	Learning Rate (Head)	0.005	0.05	0.005	0.0005	0.005	0.005	0.005
	Weight Decay	0.0005	0.0	0.05	0.05	0.0005	0.0005	0.05
ViT-Large	Learning Rate	0.001	0.01	0.01	0.01	0.01	0.01	0.01
	Learning Rate (Head)	0.0005	0.005	0.05	0.0005	0.0005	0.0005	0.05
	Weight Decay	0.0005	0.1	0.1	0.01	0.0005	0.1	0.05

Table 7: Hyper-parameter setup for AdaMSS_{base} and AdaMSS for NLU.

Model	Hyperparameter	SST-2	MRPC	CoLA	QNLI	RTE	STS-B
Both	Optimizer	AdamW					
	Batch Size	32					
	Epochs	100					
	Seeds	{0,1111,22222,33333,44444}					
RoBERTa-Base	Learning Rate	0.001	0.01	0.001	0.001	0.0005	0.001
	Learning Rate (Head)	0.005	0.0005	0.005	0.005	0.005	0.005
	Weight Decay	0.0005	0.0	0.005	0.005	0.005	0.005
RoBERTa-Large	Learning Rate	0.001	0.001	0.001	0.0005	0.0005	0.001
	Learning Rate (Head)	0.0005	5e-05	0.05	0.05	0.005	0.0005
	Weight Decay	0.0	0.005	0.0005	0.005	0.0	0.0005

Table 8: The experimental setup

Fine-tuned on	Evaluated Datasets
MetaMathQA [33]	GSM8K [13], MATH [33]

Table 9: Hyper-parameter setup for AdaMSS_{base} for NLG.

Model	Hyperparameter	GSM8K	MATH
LLaMA 2-7B	Optimizer	AdamW	
	Batch Size	4	
	Epochs	1	
	Learning Rate	2e-3	

G Hyper-parameter setup

For a fair comparison, we replicate the experimental setups in [1, 2, 24]:

- IC and NLU: We tune the hyper-parameters for learning rates and weight decay for all methods. A detailed list of all hyper-parameters and settings can be found in the Tables 6-7. The reported results are averaged over five random seeds. All models are fine-tuned by updating only the query and value projection matrices, while the classification head is updated for every method.
- NLG: Following [2], we use the datasets listed in Table 8, and all experiments are conducted on 100K-example subsets and trained for a single epoch. List of all hyper-parameters and settings can be found in the Table 9. The results are averaged over three runs.

If not specified, K_0 in AdaMSS and AdaMSS_{base} is set to 10.

NeurIPS Paper Checklist

The checklist is designed to encourage best practices for responsible machine learning research, addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove the checklist: **The papers not including the checklist will be desk rejected.** The checklist should follow the references and follow the (optional) supplemental material. The checklist does NOT count towards the page limit.

Please read the checklist guidelines carefully for information on how to answer these questions. For each question in the checklist:

- You should answer [Yes], [No], or [NA].
- [NA] means either that the question is Not Applicable for that particular paper or the relevant information is Not Available.
- Please provide a short (1–2 sentence) justification right after your answer (even for NA).

The checklist answers are an integral part of your paper submission. They are visible to the reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it (after eventual revisions) with the final version of your paper, and its final version will be published with the paper.

The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation. While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided a proper justification is given (e.g., "error bars are not reported because it would be too computationally expensive" or "we were unable to find the license for the dataset we used"). In general, answering "[No]" or "[NA]" is not grounds for rejection. While the questions are phrased in a binary way, we acknowledge that the true answer is often more nuanced, so please just use your best judgment and write a justification to elaborate. All supporting evidence can appear either in the main paper or the supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification please point to the section(s) where related material for the question can be found.

IMPORTANT, please:

- **Delete this instruction block, but keep the section heading “NeurIPS Paper Checklist”,**
- **Keep the checklist subsection headings, questions/answers and guidelines below.**
- **Do not modify the questions and only use the provided macros for your answers.**

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope?

Answer: [Yes]

Justification: The abstract and introduction clearly articulate the core contributions of the paper, including the proposed methodology, the scope of experiments, and key findings.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: The authors mention several shortcomings and directions for improvement in the “Future Work” section, reflecting an awareness of the method’s current boundaries.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren’t acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [\[Yes\]](#)

Justification: All theoretical results in the paper are supported by clearly stated assumptions. The proofs are included in the Appendix.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [\[Yes\]](#)

Justification: The details of the experiments are given in the Appendix.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: All datasets employed in the experiments are publicly accessible and appropriately referenced. The code of the proposed method is included in the supplemental material, with sufficient instructions

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.

- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [\[Yes\]](#)

Justification: The details are given in the Appendix.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [\[Yes\]](#)

Justification: The paper reports experimental results with multiple independent runs, and the variance across these runs is provided to capture statistical variability.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [\[Yes\]](#)

Justification: The details are given in the experiment section.

Guidelines:

- The answer NA means that the paper does not include experiments.

- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: The research fully conforms to the NeurIPS Code of Ethics. All data used are publicly available and properly cited, no human subjects or private user information were involved, and no actions were taken that could lead to harm or discrimination.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: The paper presents foundational research. As such, it does not have an immediate societal impact.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [\[Yes\]](#)

Justification: The paper does not release any models or datasets that pose a high risk for misuse. All datasets used are publicly available and widely adopted benchmarks, and no new data is scraped or generated for release.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [\[Yes\]](#)

Justification: All datasets and code packages used in the paper are publicly available and properly cited.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [\[Yes\]](#)

Justification: The paper releases a new code and provides detailed documentation.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

779 **14. Crowdsourcing and research with human subjects**

780 Question: For crowdsourcing experiments and research with human subjects, does the paper
781 include the full text of instructions given to participants and screenshots, if applicable, as
782 well as details about compensation (if any)?

783 Answer: [NA]

784 Justification: The study does not involve any crowdsourcing or human subjects. All data
785 used are collected from publicly available sources or generated synthetically.

786 Guidelines:

- 787 • The answer NA means that the paper does not involve crowdsourcing nor research with
788 human subjects.
- 789 • Including this information in the supplemental material is fine, but if the main contribu-
790 tion of the paper involves human subjects, then as much detail as possible should be
791 included in the main paper.
- 792 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
793 or other labor should be paid at least the minimum wage in the country of the data
794 collector.

795 **15. Institutional review board (IRB) approvals or equivalent for research with human**
796 **subjects**

797 Question: Does the paper describe potential risks incurred by study participants, whether
798 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
799 approvals (or an equivalent approval/review based on the requirements of your country or
800 institution) were obtained?

801 Answer: [NA]

802 Justification: The paper does not involve any experiments with human subjects or crowd-
803 sourcing. Therefore, IRB approval is not applicable.

804 Guidelines:

- 805 • The answer NA means that the paper does not involve crowdsourcing nor research with
806 human subjects.
- 807 • Depending on the country in which research is conducted, IRB approval (or equivalent)
808 may be required for any human subjects research. If you obtained IRB approval, you
809 should clearly state this in the paper.
- 810 • We recognize that the procedures for this may vary significantly between institutions
811 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
812 guidelines for their institution.
- 813 • For initial submissions, do not include any information that would break anonymity (if
814 applicable), such as the institution conducting the review.

815 **16. Declaration of LLM usage**

816 Question: Does the paper describe the usage of LLMs if it is an important, original, or
817 non-standard component of the core methods in this research? Note that if the LLM is used
818 only for writing, editing, or formatting purposes and does not impact the core methodology,
819 scientific rigor, or originality of the research, declaration is not required.

820 Answer: [NA]

821 Justification: LLMs were not used in the development of the core methodology.

822 Guidelines:

- 823 • The answer NA means that the core method development in this research does not
824 involve LLMs as any important, original, or non-standard components.
- 825 • Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>)
826 for what should or should not be described.