# Serial Low-rank Adaptation of Vision Transformer

Houqiang Zhong[1*], Shaocheng Shen[2*], Ke Cai[3*], Zhenlong Wu[1], Jiangchao Yao[2], Yuan Cheng[3]
Xuefei Li[3], Xiaoyun Zhang[2], Li Song[1†], Qiang Hu[2†]

[1] School of Information Science and Electronic Engineering, Shanghai Jiao Tong University, Shanghai, China
[2] Cooperative Medianet Innovation Center, Shanghai Jiao Tong University, Shanghai, China
[3] Glodon Company, Shanghai, China

*Abstract*—Fine-tuning large pre-trained vision foundation models in a parameter-efficient manner is critical for downstream vision tasks, considering the practical constraints of computational and storage costs. Low-rank adaptation (LoRA) is a well-established technique in this domain, achieving impressive efficiency by reducing the parameter space to a low-rank form. However, developing more advanced low-rank adaptation methods to reduce parameters and memory requirements remains a significant challenge in resource-constrained application scenarios. In this study, we consider on top of the commonly used vision transformer and propose Serial LoRA, a novel LoRA variant that introduces a shared low-rank matrix serially composite with the attention mechanism. Such a design extracts the underlying commonality of parameters in adaptation, significantly reducing redundancy. Notably, Serial LoRA uses only 1/4 parameters of LoRA but achieves comparable performance in most cases. We conduct extensive experiments on a range of vision foundation models with the transformer structure, and the results confirm consistent superiority of our method.

## I. INTRODUCTION

Vision foundation models, such as SAM [1], CLIP [2], and diffusion models [3], particularly those based on large-scale transformer architectures, have driven rapid advancements in computer vision, achieving breakthroughs in tasks like classification, segmentation and image generation. Despite the impressive performance, deploying these large-scale pre-trained models in real-world applications faces inevitable limitations, especially in resource-constrained environments like mobile devices and edge computing platforms, where high computational and memory demands present challenges. Consequently, parameter-efficient fine-tuning (PEFT) methods are developed to reduce storage and computational costs, making vision foundation models more accessible.

Current PEFT methods encompass different paradigms like prompt tuning [4]–[6], adapter-based methods [7]–[9] and low-rank adaptation (LoRA) [10]–[12], where we place focus on the latter in this study, given its potential in efficiency and parameter size reduction. As a representative approach, LoRA constructs the product of two low-rank matrices to narrow down the optimization space, making it possible to adapt vision foundation models to new tasks with fewer resources.
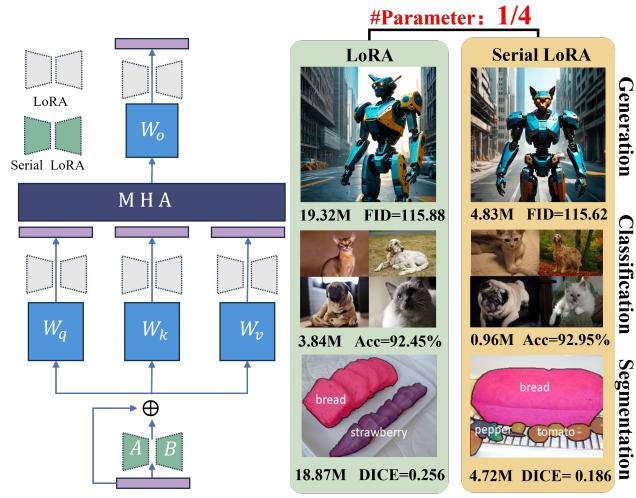
Fig. 1: **Left**: Serial LoRA structure, illustrating the application of Serial LoRA in Transformers. **Right**: Comparison results of Serial LoRA and LoRA across various computer vision tasks, including image generation, image classification and semantic segmentation.

For example, recent studies [13], [14] have shown that LoRA can effectively adapt vision transformer blocks in SAM for precise anatomical segmentation, while significantly reducing the parameter overhead in domain adaptation. In text-to-image generation, recent works [15]–[17] have successfully integrated LoRA into diffusion models, enabling efficient personalization of generative models while maintaining high-fidelity image synthesis capabilities. Despite effectiveness, current LoRA has not well considered the structural characteristic of vision transformer. This makes us rethink how to better align PEFT methods with the architectural nuances of vision transformers, ensuring scalability and efficiency in fine-tuning large-scale vision models.

Generally, in vision transformers, the multi-head attention (MHA) mechanism and the QKV design tend to use larger feature dimensions and more attention heads to enhance the model's ability to capture relationships across multiple feature spaces. However, the fine-tuning process must adjust parameters across the individual feature spaces of Q, K, V, and O for each attention head. this parallel structure leads to an

increase in fine-tuning parameters, making adaptation more computationally expensive. This paper explores a fundamental question: Is there a shared low-rank LoRA space that can efficiently adapt the transformer structure? We hypothesize that such a unified LoRA space could capture underlying parameter commonalities across feature subspaces, reducing redundancy and the number of parameters required for fine-tuning. This concept serves as the foundation that motivates us to propose a more efficient LoRA method to finetuning vision foundation models with the transformer structure.

To validate this hypothesis, we propose Serial LoRA, a novel variant specifically designed to address the parameter-intensive parallel structure in Transformers. Serial LoRA employs a shared low-rank matrix that can be serially composite with the attention mechanism, as shown in Fig. 1, enabling the attention heads to efficiently share a unified low-rank space. This design extracts underlying parameter commonalities across heads, significantly reducing redundancy and fine-tuning costs. Our comprehensive evaluations across various transformer-based vision foundation models demonstrate the effectiveness of Serial LoRA. Specifically, we conduct extensive experiments on 24 datasets spanning three major vision tasks: image classification using CLIP, semantic segmentation with SAM, and image generation based on Stable Diffusion 3. The results consistently show that Serial LoRA achieves comparable performance while reducing model parameters by 75% compared to standard LoRA. Furthermore, we successfully integrate learning rate adjustment strategies from LoRA+ into Serial LoRA, demonstrating its compatibility with existing LoRA improvements and highlighting the extensibility of our approach. These systematic evaluations across diverse vision tasks and model architectures validate the efficiency and versatility of Serial LoRA as a general parameter-efficient fine-tuning method.

In summary, our contributions are as follows:

- We propose Serial LoRA, an novel method that leverages a shared low-rank matrix serially composite with the attention mechanisms, which equivalently builds multiple LoRAs respectively adapted with different parameters.
- Serial LoRA can be seamlessly integrated into the adaptation of various vision foundation models with the transformer structure, expanding the scope of LoRA for more efficient fine-tuning in resource-constrained scenarios.
- Extensive experiments on 24 datasets encompassing classification, segmentation and generation, consistently demonstrate the comparable performance of Serial LoRA with 1/4 parameters of LoRA in most cases.

## II. METHODS

### A. preliminary

Without loss of generality, for a pre-trained weight matrix $\mathbf{W}_0 \in \mathbb{R}^{d_1 \times d_2}$ where $d_1$ and $d_2$ are the dimensions of the parameter layer, LoRA models the difference between the pre-

trained and fine-tuned weights as the product of two low-rank matrices:

$$\overline{\mathbf{W}} = \mathbf{W}_0 + \Delta\mathbf{W} = \mathbf{W}_0 + \mathbf{B}\mathbf{A} \tag{1}$$

where $\mathbf{B} \in \mathbb{R}^{d_1 \times r}$, $\mathbf{A} \in \mathbb{R}^{r \times d_2}$, and $r \ll \min\{d_1, \ d_2\}$. During fine-tuning, $\mathbf{W}_0$ remains fixed while only $\Delta\mathbf{W}$ undergoes training. Matrix $\mathbf{A}$ is initialized with random Gaussian noise, ensuring diverse updates, while $\mathbf{B}$ is initialized to zero, setting $\Delta\mathbf{W} = 0$ at the begining, which avoids immediate interference with the pre-trained weights. This low-rank adaptation approach quickly gained traction in large vision models, especially within the context of vision transformers, leading to the development of various extensions such as LoRA+ [12], and DoRA [18]. These methods commonly adapt LoRA to the query $(q)$, key $(k)$, value $(v)$, and output $(out)$ projection matrices in the attention block. Formally, for any projection matrix $\mathbf{W}_{\text{proj}} \in \{\mathbf{W}_q, \mathbf{W}_k, \mathbf{W}_v, \mathbf{W}_{\text{out}}\}$, the fine-tuned weight $\overline{\mathbf{W}_{\text{proj}}}$ is defined by:

$$\overline{\mathbf{W}_{\text{proj}}} = \mathbf{W}_{\text{proj}} + \Delta\mathbf{W}_{\text{proj}} = \mathbf{W}_{\text{proj}} + \mathbf{B}_{\text{proj}}\mathbf{A}_{\text{proj}} \tag{2}$$

By concentrating on essential weight updates, LoRA effectively adapts large-scale vision transformers to new tasks.

### B. motivation

Despite rapid advancements, current LoRA techniques still fall short of meeting real-world demands, especially where the cost of fine-tuning remains high due to substantial parameter requirements. Although several methods have further optimized its parameter redundancies, none of them has considered the inherent network structure, e.g. the vision transformer. This raises us a question: rather than transferring large models with extensive parameters encoding prior knowledge to personalized domains, could we instead learn to adapt the commonality among parameters in vision transformer to further compress the parameter redundancy? To address this, we propose exploring a unified low-rank space within the multi-head attention (MHA) structure of Transformers architectures. By sharing adaptive parameters, we can further reduce the number of trainable parameters to be optimized, enabling more scalable adaption to various resource-constraint tasks and domains.

### C. Serial LoRA for Vision Transformer

Given the aforementioned discussion, we propose **Serial LoRA**, a parameter-efficient fine-tuning approach that introduces a shared low-rank transformation in vision transformers. Serial LoRA learns a pair of low-rank matrices, $\mathbf{A}_s \in \mathbb{R}^{r \times d_2}$ and $\mathbf{B}_s \in \mathbb{R}^{d_1 \times r}$, to directly transform input features for adaptation within the model's pre-trained parameter space. We define the fine-tuning weight matrix as $\Delta\mathbf{W}_s = \mathbf{B}_s\mathbf{A}_s$ and apply it on the input feature $\mathbf{x}$ of MHA block to generate the adapted query, key, and value embeddings. Take query embedding as example, we have:

$$\tilde{\mathbf{q}} = \mathbf{W}_q(\mathbf{I} + \mathbf{B}_s\mathbf{A}_s)\mathbf{x} \tag{3}$$

where $\tilde{\mathbf{q}}$ is the adapted query, The same fine-tuning process applies to the key and value embeddings, with analogous
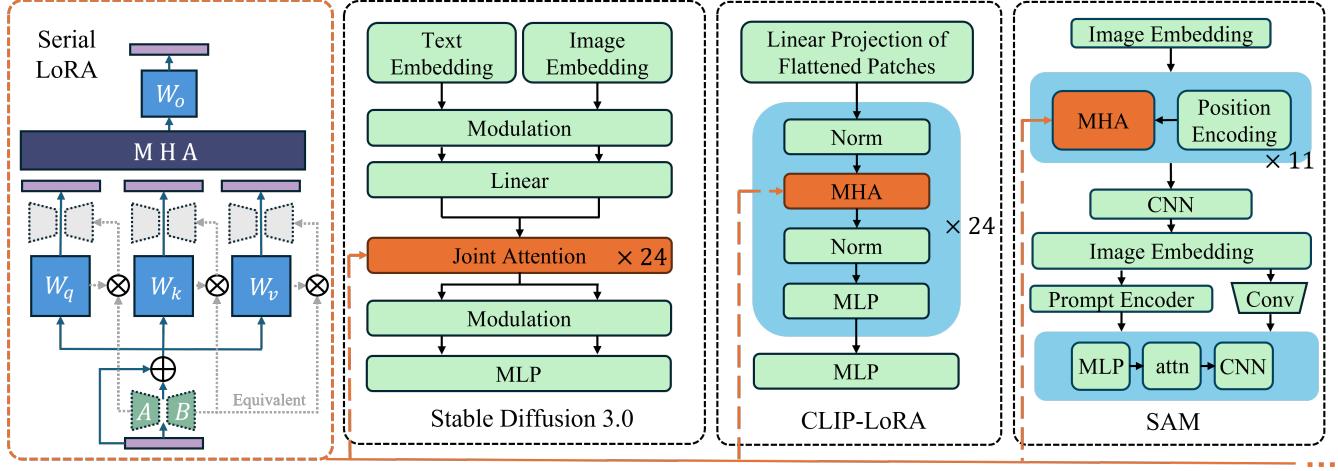
Fig. 2: Instead of learning separate pairs of matrices, Serial LoRA learns a shared pair of low-rank matrices, significantly reducing the training parameter requirements. Its strong scalability allows it to be directly applied to various vision tasks, such as CLIP, Stable Diffusion 3.0 and SAM, enhancing efficiency across diverse applications.
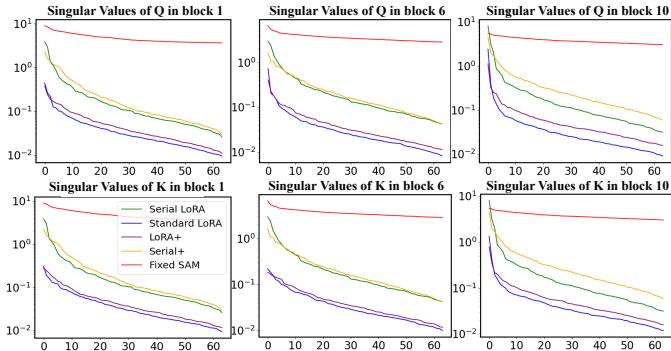


Fig. 3: Singular value analysis varying different transformer blocks in SAM model. The comparison between SAM weights (red), Serial LoRA (green), LoRA (blue), LoRA+ (purple), and Serial LoRA+ (orange) demonstrates that Serial LoRA variants maintain intermediary singular value distributions with gradual decay patterns across different network depths.

expressions for $\tilde{\mathbf{k}}, \tilde{\mathbf{v}}$. These matrices encapsulate the model's pre-trained knowledge and serve as the foundation for our task-specific adaptations. This design differs fundamentally from conventional methods in two aspects. First, while standard LoRA applies separate learnable matrices in parallel for each attention component, Serial LoRA learns a single, shared adaptation matrix that acts as a common adjustment. Second, this transformation is applied before the projection through pre-trained weights, enabling a more uniform adaptation across attention components. The shared nature of $\mathbf{B}_s\mathbf{A}_s$ not only significantly reduces the number of trainable parameters but also aligns the adaptations within a unified parameter space. This serial propagation of a shared transformation achieves parameter efficiency and maintains model adaptability through a more compact and unified design.

**Extension.** We can easily integrate Serial LoRA into diverse vision foundation models with the transformer structure. The detailed architecture modifications and integration strategy of

Serial LoRA are illustrated in Fig. 2, which demonstrates how our approach can be seamlessly incorporated into various transformer-based models. In the experimental part, we will demonstrate that Serial LoRA learns a unified, parameter-efficient way for adapting vision transformers efficiently across diverse architectures and tasks.

### D. Analysis

**Equivalent Form.** Without loss of generality, we take the query embedding as example. With a pre-trained matrix $\mathbf{W}_q$, both Serial LoRA and standard LoRA adapt to new parameter spaces by learning transformations of input features in the attention mechanism. Due to the different positions and strategies applied in the transformer block, standard LoRA and Serial LoRA exhibit distinct formulations in their transformations, as shown in (4) and (5):

$$\mathbf{q} = (\mathbf{W}_q + \mathbf{B}_q\mathbf{A}_q)\mathbf{x} = \mathbf{W}_q\mathbf{x} + \mathbf{B}_q\mathbf{A}_q\mathbf{x}, \quad (4)$$

$$\tilde{\mathbf{q}} = \mathbf{W}_q(\mathbf{I} + \mathbf{B}_s\mathbf{A}_s)\mathbf{x} = \mathbf{W}_q\mathbf{x} + \mathbf{W}_q\mathbf{B}_s\mathbf{A}_s\mathbf{x}. \quad (5)$$

Standard LoRA directly transform the input by $\mathbf{B}_q\mathbf{A}_q$, but Serial LoRA first forms a composite matrix by $\mathbf{W}_q\mathbf{B}_s\mathbf{A}_s$. It uses the shared matrix to be composite with the weight matrix in attention mechanism, which individually builds different adaptation based on the basis $\mathbf{W}_q$ ($\mathbf{W}_k$ or $\mathbf{W}_v$). From this perspective, our method actually adapt like LoRA by means of the specification of different basis $\mathbf{W}$.

**Non-equivalent Dynamic.** Although there are some potential equivalent form between Serial LoRA and LoRA, we should point out that they perform very different. To be clear, in Fig. 3, we conducted singular value analysis on the Q, K, and V projection matrices across different transformer blocks learned by Serial LoRA and LoRA. As can be seen, Serial LoRA demonstrates two key properties in its singular value distribution: a consistent intermediate positioning between the original pre-trained model and standard LoRA, and a more gradual decay pattern compared to the steep deterioration in standard

LoRA. The intermediate positioning suggests that Serial LoRA achieves an optimal balance between parameter efficiency and representation capacity. Specifically, the higher singular values compared to standard LoRA indicate better preservation of the model's original representation power, while still maintaining significant parameter reduction relative to the original model. The gradual decay pattern further implies that Serial LoRA retains a broader spectrum of feature dimensions, contributing to more robust feature representations. These characteristics persist across both different attention components (Q, K, and V matrices) and various transformer blocks from shallow to deep layers. This consistent behavior throughout the network hierarchy demonstrates that Serial LoRA's enhanced representation capacity is systematically maintained, likely contributing to its superior performance across diverse downstream tasks.

## III. EXPERIMENTS

To validate the effectiveness of our method, we conduct extensive experiments on 24 datasets across three major vision tasks: image generation using Stable Diffusion 3.0 [3], image classification with CLIP [2], and semantic segmentation with SAM [1]. Our method, Serial LoRA, is compared directly with LoRA [10]. Additionally, we extend Serial LoRA into the LoRA+ [12] framework and perform further comparisons with LoRA+ to evaluate its compatibility and performance within advanced fine-tuning structures. Please refer to the supplementary material for more results.

|  | LoRA / LoRA+ | Serial LoRA / Serial LoRA+ |
|---|---|---|
| CLIP [2] | 3.84M | **0.96M** |
| SAM [1] | 18.87M | **4.72M** |
| SD3 [3] | 19.32M | **4.83M** |

TABLE I: Comparison of Parameter Counts Between LoRA and Serial LoRA. Our method achieves a 4-fold reduction in parameters compared to LoRA.

### A. Serial Lora on Diffusion

**Experimental Setup**. For our generative tasks, we evaluate eight stylistic datasets: Barbie, Cyberpunk, Art Nouveau, Impressionism. Barbie and Cyberpunk contain 315 and 439 images, respectively,[1] with prompts drawn from captions, while the remaining datasets each include 1,500 images using artwork titles as prompts.[2] According to [19], we assess quality using Fréchet Inception Distance (FID) and CLIP-Score. We utilize Stable Diffusion 3.0 [3] for image generation, comparing LoRA, Serial LoRA, LoRA+, and Serial LoRA+ to assess compatibility and parameter efficiency. Training is set to 10 epochs per dataset with an initial learning rate of $1.5 \times 10^{-5}$. For LoRA+ and Serial LoRA+, we maintain a learning rate ratio of 1:20 between matrices **A** and **B**. All method is configured with a default rank of 64. Evaluation involves 28 inference steps at a $1024 \times 1024$ resolution, with

[1]https://github.com/sjtuplayer/SaRA.
[2]https://github.com/liaopeiyuan/artbench.

classifier-free guidance weight $\omega = 7.0$ to ensure consistency across methods. All experiments run on a single Nvidia RTX 3090 GPU in bfloat16 precision.

|  |  | Barbie | | Cyberpunk | | Art Nouveau | | Impressionism | |
|---|---|---|---|---|---|---|---|---|---|
|  |  | FID↓ | CLIP Score↑ | FID↓ | CLIP Score↑ | FID↓ | CLIP Score↑ | FID↓ | CLIP Score↑ |
| 20% Shot | LoRA | **189.62** | 0.8223 | 172.43 | **0.8146** | 153.31 | 0.6686 | 153.98 | 0.6963 |
|  | Serial LoRA | 193.66 | **0.8302** | 173.54 | 0.8135 | **152.32** | 0.6683 | **153.30** | 0.7003 |
|  | LoRA+ | 187.21 | 0.8122 | 173.68 | 0.8127 | 152.83 | 0.6609 | 153.81 | 0.6983 |
|  | Serial LoRA+ | **182.31** | 0.8222 | **172.21** | **0.8178** | 153.04 | **0.6693** | 153.79 | 0.6998 |
| 50% Shot | LoRA | 142.16 | 0.8256 | 115.88 | **0.8235** | 94.62 | 0.6795 | 102.92 | 0.7051 |
|  | Serial LoRA | **141.93** | 0.8264 | 115.62 | 0.8227 | 94.87 | 0.6841 | 102.65 | 0.7052 |
|  | LoRA+ | 142.79 | 0.8243 | 116.25 | 0.8231 | 94.63 | 0.6815 | 102.10 | 0.7020 |
|  | Serial LoRA+ | 142.14 | **0.8251** | 117.41 | 0.8232 | **94.57** | **0.6884** | 102.33 | 0.7038 |
| Full Shot | LoRA | 116.35 | 0.8298 | **99.39** | 0.8282 | **67.25** | 0.6786 | 73.50 | 0.7010 |
|  | Serial LoRA | 117.93 | **0.8302** | 99.88 | 0.8285 | 67.29 | 0.6791 | 74.08 | 0.7030 |
|  | LoRA+ | 115.51 | 0.8294 | 99.85 | 0.8279 | 68.04 | 0.6792 | **72.81** | 0.7020 |
|  | Serial LoRA+ | **114.42** | 0.8299 | **99.27** | **0.8290** | **66.87** | **0.6794** | 74.10 | 0.7031 |

TABLE II: Quantitative Comparison of FID and CLIP-Score Between Serial LoRA and LoRA on the Stable Diffusion 3.0, with Model Parameters Reduced from 19.32M to **4.83M**.

**Quantitative Results.** Tab. II presents a quantitative comparison of FID and CLIP-Score between our proposed Serial LoRA and Serial LoRA+ methods against standard LoRA and LoRA+ on the Stable Diffusion 3.0 model. The results demonstrate that Serial LoRA and Serial LoRA+ achieve comparable performance to LoRA and LoRA+ while using only 1/4 of the parameters, significantly reducing the model's memory footprint from 19.32M to 4.83M parameters. Specifically, when comparing Serial LoRA to LoRA, we observe that Serial LoRA maintains similar or even improved FID and CLIP-Score across most styles, such as Barbie and Art Nouveau, indicating strong generative quality with lower parameter costs. Similarly, Serial LoRA+ shows robust performance relative to LoRA+, achieving comparable or superior scores in various styles like Cyberpunk and Impressionism. These findings confirm that our method integrates seamlessly with existing LoRA structures, significantly reducing parameter overhead while retaining high-quality outputs across diverse styles.

**Qualitative Results.** We use a consistent prompt as input to generate multiple images with specific styles. Part of the qualitative comparison results of LoRA, Serial LoRA, LoRA+, and Serial LoRA+ are shown in Fig. 4. It can be seen that our model can learn the style accurately while generating images that align well with the given text prompts across different datasets. These findings further validate that Serial LoRA integrates effectively with existing Transformer-based diffusion models, enabling efficient fine-tuning and the creation of high-quality, distinctively styled megapixel images.

### B. Serial Lora on CLIP

**Experimental Setup** For the image classification task, we evaluate the performance of Serial LoRA on the CLIP model (ViT-B/16) [2] across eight benchmark datasets: Caltech101 (101 classes) [20], Food 101 (101 classes) [21], EuroSAT (10 classes) [22] and Oxford Pets (37 classes) [23]. Average

| Dataset | LoRA | Serial LoRA | LoRA+ | Serial LoRA+ |

Prompt: there is a pink toy tank sitting on the ground, featured on unsplash, resistance, eco-friendly theme

Fig. 4: Qualitative Results about the comparison of Serial LoRA and LoRA.

accuracy is used as the primary metric to assess model performance. Each model is trained for 50 epochs, with a fixed learning rate of 2e-4, and a dropout rate of 0.25. All method is configured with a default rank of 32. For LoRA+ and Serial LoRA+, we maintain a learning rate ratio of 1:20 between matrices $\mathbf{A}$ and $\mathbf{B}$. All experiments are conducted on a single Nvidia RTX 3090 GPU in bfloat16.

| | | Caltech101 | Food 101 | EuroSAT | OxFord Pets |
|---|---|---|---|---|---|
| 1 Shot | LoRA | 93.22 | 83.64 | 69.90 | 87.87 |
| | Serial LoRA | **93.69** | **85.77** | **70.80** | **90.52** |
| | LoRA+ | 93.31 | 84.08 | 59.45 | 88.23 |
| | Serial LoRA+ | **94.44** | **85.88** | **67.36** | **89.59** |
| 4 Shot | LoRA | 94.51 | 84.08 | 86.42 | 90.13 |
| | Serial LoRA | **95.54** | **86.04** | **86.98** | **91.92** |
| | LoRA+ | 92.60 | 79.44 | 76.01 | 85.78 |
| | Serial LoRA+ | **94.48** | **82.79** | **88.25** | **89.51** |
| 16 Shot | LoRA | 96.25 | 84.69 | 92.75 | 92.45 |
| | Serial LoRA | **96.27** | **86.98** | 91.90 | **92.95** |
| | LoRA+ | 95.02 | 77.98 | 88.11 | 87.78 |
| | Serial LoRA+ | **95.94** | **81.92** | **92.12** | **89.07** |

TABLE III: Comparison of Classification Accuracy(%) between Serial LoRA and LoRA on the CLIP Method, with Model Parameters Reduced from 3.84M to **0.96M**.

**Analysis**. Tab. III presents a comparison of classification accuracy between Serial LoRA, Serial LoRA+, and standard LoRA and LoRA+ methods across eight benchmark datasets using the CLIP model. The results demonstrate that our proposed Serial LoRA and Serial LoRA+ methods achieve comparable or even superior accuracy to LoRA and LoRA+ while significantly reducing model parameters from 3.84M to 0.96M. Specifically, in few-shot learning scenarios, particularly in the 1 Shot setting, Serial LoRA consistently exhibits

marginally improved accuracy across almost all datasets. Notably, Serial LoRA+ outperforms LoRA+ in accuracy across all datasets and shot settings, confirming its efficacy in achieving parameter efficiency without compromising performance. This improvement can be attributed to the design of Serial LoRA, which leverages a shared low-rank matrix in a serial configuration. This structure captures parameter commonalities across the attention heads more effectively, reducing redundancy while maintaining expressive power.

### C. Serial LoRA on SAM

**Experimental Setup**

For our segmentation tasks, we construct eight few-shot segmentation datasets across diverse categories: one flood segmentation task[3], human action segmentation task[4] and three food category segmentation tasks[5]. Segmentation performance is measured by the Dice Score [24], where a lower score indicates greater accuracy. For implementation, we use the SAM-ViT-B backbone, fine-tuning only the Vision Transformer (ViT) in the image encoder, while keeping the prompt encoder and mask decoder frozen. Parameter efficiency comparisons are conducted with a rank of 64 for all methods, using an initial learning rate of $1 \times 10^{-4}$ and a 1:20 ratio between matrices $\mathbf{A}$ and $\mathbf{B}$ for LoRA+ variants. Models are trained for 5 epochs across all datasets on an Nvidia RTX 3090 GPU in bfloat16 precision.

**Analysis**. Tab. IV indicates that Serial LoRA and Serial LoRA+ outperform LoRA and LoRA+ in nearly all cases, despite using only **1/4** of the original parameter count. For instance, Serial LoRA shows a slight improvement over LoRA across various shot settings on datasets such as Flood and Strawberry. Additionally, the DICE score indicates that Serial LoRA yields more accurate segmentation masks. Specifically, Serial LoRA+ achieves state-of-the-art results in few-shot learning involving comparatively fewer samples, demonstrating both the compatibility of our approach and the superior performance of Serial LoRA under more limited data conditions. These findings demonstrate that Serial LoRA can be smoothly incorporated into existing LoRA architectures, enabling more accurate segmentation even with significantly reduced trainable parameters.

As shown in Tab. I For the experiments conducted on CLIP, SAM, and SD3 models, our shared parameter design in Serial LoRA ensures a structural parameter reduction of 1/4 compared to the standard LoRA approach.

### D. Further Study

We validate the adaptability of the Serial LoRA method under different rank settings and compare its performance with standard LoRA. For image classification with CLIP (Fig. 5), we evaluate Serial LoRA by varying the rank from 8 to 64 across diverse datasets including Food101, DTD, Caltech101, and Oxford Pets. With identical rank settings, Serial LoRA

---

[3]https://datasetninja.com/floodnet.
[4]http://vision.stanford.edu/Datasets/40actions.html.
[5]https://github.com/LARC-CMU-SMU/FoodSeg103-Benchmark-v1.

| | | Flood | Human | Icecream | Pie | Strawberry |
|---|---|---|---|---|---|---|
| 20% Shot | LoRA | 0.2983 | 0.0943 | **0.2079** | **0.2483** | 0.2736 |
| | Serial LoRA | **0.2721** | **0.0790** | 0.2214 | 0.3376 | **0.1675** |
| | LoRA+ | **0.2469** | 0.1575 | 0.2464 | 0.3560 | 0.4512 |
| | Serial LoRA+ | 0.2509 | 0.1574 | 0.2413 | 0.1822 | 0.4490 |
| 50% Shot | LoRA | 0.2487 | 0.0806 | 0.2227 | 0.2910 | 0.1609 |
| | Serial LoRA | **0.2355** | **0.0697** | **0.1887** | **0.2803** | **0.1344** |
| | LoRA+ | 0.2659 | 0.0971 | 0.2747 | 0.3323 | 0.3013 |
| | Serial LoRA+ | 0.2478 | 0.0878 | 0.2370 | 0.1809 | 0.2452 |
| Full Shot | LoRA | 0.2822 | 0.0618 | 0.1673 | 0.2585 | 0.1140 |
| | Serial LoRA | **0.2333** | **0.0439** | 0.1642 | 0.1856 | 0.1116 |
| | LoRA+ | 0.2419 | **0.0692** | 0.1778 | 0.2898 | 0.1453 |
| | Serial LoRA+ | **0.2411** | 0.0727 | **0.1717** | **0.1551** | 0.1383 |

TABLE IV: Dice$\downarrow$ score comparison between Serial LoRA (**4.72MB**) and LoRA (18.87MB) for SAM-based few-shot segmentation tasks.

reduces the parameter count to 1/4 of standard LoRA through shared transformation design.
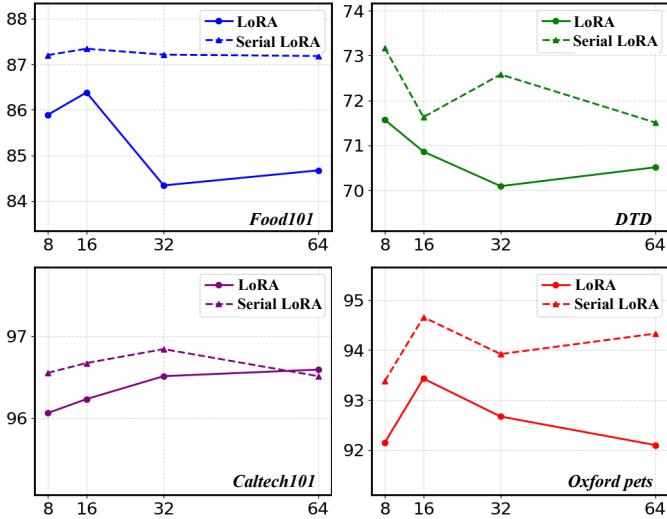


Fig. 5: Comparison of LoRA and Serial LoRA performance in the CLIP model on image classification tasks from rank 8 to 64, measured in accuracy (%). The datasets used include Food101, DTD, Caltech101, and Oxford Pets.

## IV. Conclusion

In this paper, we propose **Serial LoRA**, a novel parameter-efficient fine-tuning method, tailored for various computer vision tasks. By identifying the shared parameter space of the multi-head attention(MHA) block within pretrained Transformers and applying a consistent low-rank adaptation approach, we significantly reduce the number of fine-tuning parameters. In experiments, we applied Serial LoRA to image generation, image classification, and semantic segmentation tasks. With fine-tuning parameters reduced to only **1/4** of the original, Serial LoRA demonstrated the effectiveness of our method while minimally impacting the model's performance. Furthermore, our approach is compatible with existing PEFT methods, underscoring the broad applicability of Serial LoRA.

## References

[1] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollar, and Ross Girshick, "Segment anything," in *Proceedings of the IEEE/CVF ICCV*, 2023, pp. 4015–4026.

[2] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever, "Learning transferable visual models from natural language supervision," in *Proceedings of the 38th ICML*, 2021, vol. 139, pp. 8748–8763.

[3] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, Dustin Podell, Tim Dockhorn, Zion English, and Robin Rombach, "Scaling rectified flow transformers for high-resolution image synthesis," in *Proceedings of the 41th ICML*, 2024.

[4] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed El-hoseiny, "MiniGPT-4: Enhancing vision-language understanding with advanced large language models," in *The Twelfth International Conference on Learning Representations*, 2024.

[5] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi, "Instruct-BLIP: Towards general-purpose vision-language models with instruction tuning," in *the 37th NeurIPS*, 2023.

[6] Shuo Chen, Jindong Gu, Zhen Han, Yunpu Ma, Philip Torr, and Volker Tresp, "Benchmarking robustness of adaptation methods on pre-trained vision-language models," in *the 37th NeurIPS Datasets and Benchmarks Track*, 2023.

[7] Zhiqiang Hu, Lei Wang, Yihuai Lan, Wanyu Xu, Ee-Peng Lim, Lidong Bing, Xing Xu, Soujanya Poria, and Roy Lee, "LLM-adapters: An adapter family for parameter-efficient fine-tuning of large language models," in *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, 2023, pp. 5254–5276.

[8] Jonas Pfeiffer, Aishwarya Kamath, Andreas Rücklé, Kyunghyun Cho, and Iryna Gurevych, "Adapterfusion: Non-destructive task composition for transfer learning," in *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, 2021, pp. 487–503.

[9] Alexandra Chronopoulou, Matthew E. Peters, Alexander Fraser, and Jesse Dodge, "Adaptersoup: Weight averaging to improve generalization of pretrained language models," in *EACL (Findings)*, 2023, pp. 2009–2018.

[10] 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," in *International Conference on Learning Representations*, 2022.

[11] Yixiao Li, Yifan Yu, Chen Liang, Nikos Karampatziakis, Pengcheng He, Weizhu Chen, and Tuo Zhao, "Loftq: LoRA-fine-tuning-aware quantization for large language models," in *The Twelfth International Conference on Learning Representations*, 2024.

[12] Soufiane Hayou, Nikhil Ghosh, and Bin Yu, "LoRA+: Efficient low rank adaptation of large models," in *Forty-first International Conference on Machine Learning*, 2024.

[13] Zelin Peng, Zhengqin Xu, Zhilin Zeng, Lingxi Xie, Qi Tian, and Wei Shen, "Parameter efficient fine-tuning via cross block orchestration for segment anything model," in *Proceedings of the IEEE/CVF Conference on CVPR*, June 2024, pp. 3743–3752.

[14] Kevin Li and Pranav Rajpurkar, "Adapting segment anything models to medical imaging via fine-tuning without domain pretraining," in *AAAI 2024 Spring Symposium on Clinical Foundation Models*, 2024.

[15] Simian Luo, Yiqin Tan, Suraj Patil, Daniel Gu, Patrick von Platen, Apolinário Passos, Longbo Huang, Jian Li, and Hang Zhao, "Lcm-lora: A universal stable-diffusion acceleration module," *arXiv preprint arXiv:2311.05556*, 2023.

[16] Zane K.J. Hartley, Rob J. Lind, Michael P. Pound, and Andrew P. French, "Domain targeted synthetic plant style transfer using stable diffusion lora and controlnet," in *Proceedings of the IEEE/CVF Conference on CVPR Workshops*, June 2024, pp. 5375–5383.

[17] James Seale Smith, Yen-Chang Hsu, Lingyu Zhang, Ting Hua, Zsolt Kira, Yilin Shen, and Hongxia Jin, "Continual diffusion: Continual customization of text-to-image diffusion with c-loRA," *Transactions on Machine Learning Research*, 2024.

[18] Shih yang Liu, Chien-Yi Wang, Hongxu Yin, Pavlo Molchanov, Yu-Chiang Frank Wang, Kwang-Ting Cheng, and Min-Hung Chen, "DoRA: Weight-decomposed low-rank adaptation," in *Forty-first International Conference on Machine Learning*, 2024.

[19] Haoning Wu, Shaocheng Shen, Qiang Hu, Xiaoyun Zhang, Ya Zhang, and Yanfeng Wang, "Megafusion: Extend diffusion models towards

higher-resolution image generation without further tuning," in *Proceedings of the Winter Conference on Applications of Computer Vision (WACV)*, February 2025, pp. 3944–3954.

[20] Li Fei-Fei, Rob Fergus, and Pietro Perona, "Learning generative visual models from few training examples: An incremental bayesian approach tested on 101 object categories," *Computer vision and Image understanding*, vol. 106, no. 1, pp. 59–70, 2007.

[21] Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool, "Food-101 – mining discriminative components with random forests," in *European Conference on Computer Vision*, 2014.

[22] Patrick Helber, Benjamin Bischke, Andreas Dengel, and Damian Borth, "Eurosat: A novel dataset and deep learning benchmark for land use and land cover classification," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 2019.

[23] Omkar M. Parkhi, Andrea Vedaldi, Andrew Zisserman, and C. V. Jawahar, "Cats and dogs," *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3498–3505, 2012.

[24] Fausto Milletari, Nassir Navab, and Seyed-Ahmad Ahmadi, "V-net: Fully convolutional neural networks for volumetric medical image segmentation," in *2016 Fourth International Conference on 3D Vision (3DV)*, 2016, pp. 565–571.