# NoRA: Nested Low-Rank Adaptation for Efficient Fine-Tuning Large Models

**Cheng Lin**[1,2†], **Lujun Li**[1†], **Dezhi Li**[1,3], **Jie Zou**[2], **Wei Xue**[1*], **Yike Guo**[1*]

[1]HKUST    [2]UESTC    [3]SEU [*]

## Abstract

In this paper, we introduce Nested Low-Rank Adaptation (NoRA), a novel approach to parameter-efficient fine-tuning that extends the capabilities of Low-Rank Adaptation (LoRA) techniques. Vanilla LoRA overlooks pre-trained weight inheritance and still requires fine-tuning numerous parameters. To addresses these issues, our NoRA adopts a dual-layer nested structure with Singular Value Decomposition (SVD), effectively leveraging original matrix knowledge while reducing tunable parameters. Specifically, NoRA freezes the outer LoRA weights and utilizes an inner LoRA design, providing enhanced control over model optimization. This approach allows the model to more precisely adapt to specific tasks while maintaining a compact parameter space. By freezing outer LoRA weights and using an inner LoRA design, NoRA enables precise task adaptation with a compact parameter space. Evaluations on tasks including commonsense reasoning with large language models, fine-tuning vision-language models, and subject-driven generation demonstrate NoRA's superiority over LoRA and its variants. Code will be released upon acceptance.

## 1 Introduction

In recent years, LLMs [56] have set new performance benchmarks in natural language processing [45] (NLP) and related fields [59, 43, 29], yet their substantial size introduces significant challenges in training and adaptation. The high computational and storage demands of comprehensive fine-tuning make it impractical in resource-limited environments. To address this, Parameter-Efficient Fine-Tuning [3, 11] (PEFT) techniques have been developed, focusing on fine-tuning a reduced subset of model parameters to optimize model adaptability. Among these techniques, LoRA [16] stands out as a prominent method that mitigates the extensive resource requirements of full fine-tuning [36] by introducing trainable low-rank matrices into the model architecture.

LoRA utilizes low-rank matrices to achieve efficient and scalable adaptation to specific downstream tasks [15, 14], allowing the model to focus on essential parameters. This approach not only preserves the model's pre-trained knowledge but also facilitates its specialization with minimal computational overhead. The core hypothesis of LoRA is that the adaptations necessary to tailor a pre-trained large language model (LLM) to a specific task or domain are inherently low-dimensional, which can be effectively achieved through a low-rank decomposition of the weight matrices, minimizing computational overhead while preserving the model's pre-trained knowledge [37].

Despite LoRA's demonstrated utility, various LoRA variants proposed by recent researchers [35, 44], each have their strengths and limitations. For example, while LoRA-FA [53] reduces activation memory demands by freezing part of the weights, it still suffers from the limitation of a fixed rank; VeRA [30] enhances model scalability but remains sensitive to the hidden dimensions of the model;
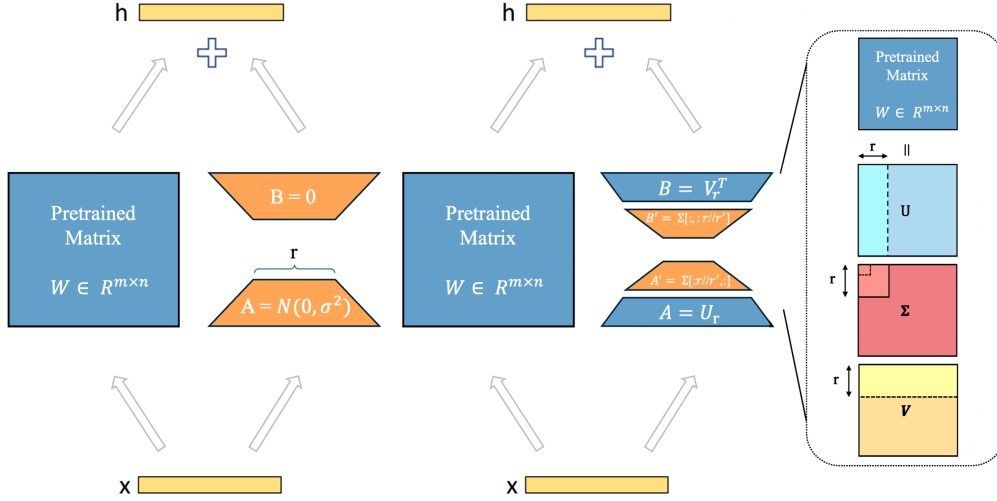
---

Figure 1: In the comparison between training with LoRA and NoRA, the blue modules represent the parts where parameters are frozen during training, while the orange modules indicate the components that need to be updated. Here, $r$ denotes the outer rank, and $r'$ denotes the inner rank.

FLoRA [12] and LoRA-XS [1] have made improvements in real-time performance and memory efficiency, yet they still do not fully address the complexity specific to tasks.

These limitations prompted us to explore a new Nested Low-Rank Adaptation structure called NoRA, which enhances parameter efficiency and task adaptability through a dual-layer nested design combined with SVD. In the NoRA structure, the outer layer provides a stable low-rank foundation by keeping its weights fixed, while the inner layer introduces fine-grained adjustments tailored to different tasks. This innovative dual-layer nested structure effectively leverages the knowledge within the original matrix and further enhances model adaptability through precise optimization of certain parameters. NoRA not only retains the parameter efficiency advantages of LoRA but also improves the accuracy of low-rank matrix approximation through the decomposition properties of SVD, thereby enhancing the model's generalization capability and adaptability across different tasks. Moreover, the hierarchical separation between the outer and inner LoRA layers provides finer optimization control, enabling the model to more flexibly adjust the contribution of each layer's parameters, thus enhancing adaptability and flexibility.

We conducted experiments on multiple downstream tasks, including fine-tuning the LLaMA [45] model for commonsense reasoning, few-shot tuning on the CLIP [41] model, and subject-driven generation on the Stable Diffusion XL [40] model. NoRA effectively reduced the required parameters to as low as 10.2M for LLaMA3 8B while enhancing performance, achieving an average score of 85.0%, surpassing LoRA's 82.8%. Additionally, in visual few-shot tasks using ViT-B/16, NoRA achieved the highest average accuracies of 81.8% (4 shots) and 85.4% (16 shots), demonstrating its superior efficiency and effectiveness over existing methods. We summarize our contribution as below:

- We propose NoRA, which uses a dual-layer nested structure and SVD. The outer LoRA provides a stable low-rank foundation, while the inner LoRA allows for precise task adjustments.

- NoRA leverages the original matrix's knowledge through its nested structure and SVD integration, enhancing low-rank matrix approximation accuracy. This design improves the model's adaptability across diverse tasks, while the hierarchical separation between layers allows for finer control over optimization.

- Our extensive experimental validation on multiple downstream tasks has shown that NoRA outperforms traditional LoRA and other recent variants such as VeRA and FLoRA.

## 2   Related Work

**Parameter-Efficient Fine-Tuning.** Parameter Efficient Fine-Tuning [39, 27, 51, 31, 16, 18] (PEFT) is a practical solution for adapting large pre-trained models to downstream tasks [51]. As the number of parameters in large-scale language models continues to grow, traditional fine-tuning methods face significant challenges, such as high computational resource demands and training costs [36]. PEFT optimizes the parameter adjustment process, effectively reducing the number of additional parameters introduced and the computational resources required for specific tasks or domains. The core idea is to enhance the model's task adaptability by streamlining parameter updates and introducing auxiliary modules while maintaining the structure and performance of the pre-trained model, without the need for complete retraining. This approach is particularly suitable for large-scale language models, as it reduces the computational burden and significantly improves the model's applicability across diverse downstream tasks. As a result, PEFT has become the mainstream trend for fine-tuning large pre-trained models, greatly promoting their adoption and application across different scenarios. Early PEFT research primarily focused on selective update strategies, such as Bias-Free Fine-Tuning [51] and Partial Network Training [10], which achieve task-specific model fine-tuning by modifying only the most critical parameters of the pre-trained model. As research progressed, adapter modules became another direction in PEFT development, embedding additional modules, such as Prompt Tuning and Prefix Tuning, into the pre-trained architecture to further enhance the model's adaptability to specific tasks. Delta-weight techniques represent the latest innovation in PEFT, with LoRA and OFT [34] being notable examples. These strategies involve fine-tuning pre-trained parameters using trainable delta weights, providing fine-grained model adjustments that improve task-specific performance while minimizing computational overhead.

**Low-rank Adaptation.** LoRA has proven to be an efficient fine-tuning strategy in various task scenarios, leveraging low-rank decomposition to enhance model adaptation while minimizing computational overhead. However, LoRA's fixed rank limitation can restrict its flexibility in handling diverse tasks. To address these limitations, researchers have proposed various LoRA variants, each showcasing unique strengths in enhancing model adaptability and performance. For instance, LoRA-FA [53] reduces activation memory demands by freezing part of the weights, though it remains limited by the fixed rank; VeRA [30] excels in enhancing model scalability but is still sensitive to the model's hidden dimensions; FLoRA [12] and LoRA-XS [53] improve real-time performance and memory efficiency, yet they do not fully address the complexity specific to certain tasks. Other variants, such as VB-LoRA [28] and Tied-LoRA [42], introduce new adaptation mechanisms that enhance model adaptability but also bring additional computational overhead and complexity, posing challenges in practical applications. Moreover, as research has progressed, techniques such as VeRA and DoRA [32] have further optimized the LoRA method, making innovative improvements in parameter efficiency, as well as in the distribution and structure of update matrices. Advanced methods like AdaLoRA [54] and PiSSA [38] push the limits of parameter update efficiency by selectively adjusting matrix ranks and distributions, significantly improving the applicability of large-scale pre-trained models in complex tasks. However, despite the significant advancements achieved by these LoRA variants, they still face certain limitations and challenges in practical applications, particularly in meeting diverse task requirements and optimizing computational resource usage, which necessitate further research and solutions.

## 3   Methodology: Nested Low-Rank Adaptation

In this paper, we present NoRA, a novel approach to parameter-efficient fine-tuning that builds upon traditional LoRA techniques. By incorporating a dual-layer nested structure with SVD, NoRA effectively harnesses the knowledge of the original matrix, further reducing fine-tuning parameters while enhancing the model's adaptability and optimization control.

### 3.1   Review of Low-Rank Adaptation

LoRA is a parameter-efficient method for fine-tuning large-scale pre-trained models. It achieves fine-tuning of the original weights through the introduction of low-rank matrix updates, without altering the stability and performance of the pre-trained models. In LoRA, the weights for each layer are updated using the following formula:

$$h = Wx + \Delta Wx = Wx + BAx, \tag{1}$$

where $\Delta W \in \mathbb{R}^{m \times n}$ is the low-rank weight update, and $A \in \mathbb{R}^{r \times n}$ and $B \in \mathbb{R}^{m \times r}$ are the low-rank matrices with $r \ll \min(m, n)$. During training, $W$ is kept frozen, and $A$ and $B$ are the trainable parameters.

## 3.2 NoRA Structure and Initialization

NoRA employs a dual-layer nested structure, where each layer is initialized using the SVD [9] of the pre-trained weights to enhance the model's learning capabilities. Specifically, NoRA first decomposes the original weights using SVD, then initializes the matrix $A$ with $U$ and the matrix $B$ with $V^T$. The inner LoRA weights are initialized using the intermediate diagonal matrix, thus optimizing the weight updates $\Delta W$.

- **Outer LoRA Layer**: This layer is initialized using the SVD results $U\Sigma V^T$ of the pre-trained weights $W$. The parameters of this outer LoRA layer are frozen during training. Freezing these parameters helps maintain stability and preserve the key features of the pre-trained model, while still allowing for precise adjustments through the inner LoRA layer.

- **Inner LoRA Layer**: Initialized with the diagonal matrix $\Sigma$ from the SVD [9] of the $U\Sigma V^T$. Initializing with $\Sigma$ allows this layer to focus on subtle perturbations in the weight space, enabling finer adjustments without altering the core weights preserved by the outer LoRA layer. This approach ensures that updates are focused on refining and enhancing the model's ability to adapt to new tasks, leveraging minor adjustments that have a targeted impact on the model's performance.

The final weight update formula is:

$$h = Wx + \Delta Wx = Wx + BAx + B'A'x, \tag{2}$$

where $A$ and $B$ are set using the truncated SVD of the original weight matrix $W$, with the outer LoRA weights remaining frozen. The matrices $A$ and $B$ are initialized as $A = U_r$ and $B = V_r^T$ where $U_r$ and $V_r$ are the left and right singular vectors corresponding to the top $r$ singular values, and these vectors are obtained from the decomposition: $W = U\Sigma V^T$ with $U \in \mathbb{R}^{m \times m}$, $\Sigma \in \mathbb{R}^{m \times n}$, and $V \in \mathbb{R}^{n \times n}$. For the inner LoRA layer, the matrices $B'$ and $A'$ are initialized by taking the square roots of selected components from the top $r_{in}$ singular values in the diagonal matrix $\Sigma_r$, where $\Sigma_r \in \mathbb{R}^{r \times r}$. Specifically, $B'$ is initialized using the square roots of the

**Algorithm 1** PyTorch code for NoRA

```
# r_out: rank of the outer LoRA layer.
# r_in: rank of the inner LoRA layer.
# A^ and B^ are the inner LoRA weights.

def init nora param(weight, r_out, r_in):
    U, S, V = torch.svd(weight)

    A = U[:, :r_out]
    B = V.T[:r_out, :]

    S = torch.diag(S)
    S_diag = torch.sqrt(S[:r_out, :r_out])

    A^ = S_diag[:, :r_out // r_in]
    B^ = S_diag[:r_out // r_in, :]
```

first $r_{in}$ columns of $\Sigma_r$, and $A'$ is initialized using the square roots of the first $r_{in}$ rows of $\Sigma_r$. This initialization approach introduces another layer of adaptation within the NoRA framework by decomposing the smaller matrix $W_r$ through an additional LoRA process.

## 3.3 Parameter Efficiency of NoRA

We make the following observation on the parameter efficiency of NoRA compared to LoRA and VeRA methods. Observation: NoRA demonstrates superior parameter efficiency compared to both LoRA and VeRA. For simplicity, let's consider a transformer model with $L$ finetuned layers, each consisting of $q$ number of $W \in \mathbb{R}^{n \times n}$ matrices.

For LoRA, the number of trainable parameters is given by:

$$P_{\text{LoRA}} = L \times q \times r \times 2n. \tag{3}$$

For NoRA, the number of trainable parameters is given by:

$$P_{\text{NoRA}} = L \times q \times r_{out} \times r_{in} \times 2. \tag{4}$$

4

To compare the parameter efficiency, we compute the ratios of the number of trainable parameters between the methods. The ratio of trainable parameters for LoRA to NoRA is:

$$\frac{P_{\text{LoRA}}}{P_{\text{NoRA}}} = \frac{L \times q \times r \times 2n}{L \times q \times r_{out} \times r_{in} \times 2} = \frac{r \times n}{r_{out} \times r_{in}}. \tag{5}$$

Since the inner rank ($r_{\text{in}}$) in the NoRA model is generally set to be the same as the rank ($r$) used in the LoRA model, NoRA demonstrates more pronounced advantages over LoRA under certain conditions. Particularly when the outer rank ($r_{\text{out}}$) is less than the maximum allowable rank ($r_{\text{max}}$), the NoRA model exhibits significantly enhanced parameter optimization. Specifically, the ratio of parameter efficiency between LoRA and NoRA can be approximated as follows:

$$\frac{P_{\text{LoRA}}}{P_{\text{NoRA}}} \approx \frac{n}{r_{out}}. \tag{6}$$

## 4 Experiment

### 4.1 Fine-tuning of Large Language Models

**Implementation Details.** We employed the LLaMA3 8B [45] and LLaMA 7B models, each fine-tuned using a variety of parameter-efficient methods to enhance their commonsense reasoning capabilities. We utilized the Commonsense170K dataset [26] for targeted fine-tuning, which is designed to enhance the models' understanding of commonsense knowledge across different contexts. The main objective was to assess the effectiveness of each fine-tuning approach in improving the model's performance on a range of commonsense reasoning tasks. Post-fine-tuning, model performance was evaluated using a suite of eight benchmark tests focused on commonsense reasoning, including ARC-e, OBQA, SIQA, and more.

**Comparison Results.** The results from the experimental evaluations are detailed in Table 1. The fine-tuning methods demonstrated varying levels of success in enhancing the reasoning capabilities of the LLaMA models.

The LLaMA 7B model, when fine-tuned with the NoRA method, exhibited the highest average score of 75.8 across most tasks, demonstrating its superior ability to generalize across different question sets. Remarkably, it achieved top scores in HellaSwag (80.6), WinoGrande (79.6), and ARC-e (80.5), underscoring its robust understanding and reasoning capabilities. The NoRA method not only showed the highest scores but also significantly reduced parameter utilization, emphasizing its efficiency and scalability as a fine-tuning strategy. In contrast, other methods like LoKr and AdaLoRA, while effective in certain areas, required more parameters or resulted in higher computational costs. Notably, the NoRA method outperformed others in terms of parameter efficiency and computational overhead, making it a promising approach for practical applications where resource constraints are a concern.

These findings highlight the potential of NoRA as a scalable and effective fine-tuning strategy that not only achieves high performance across diverse datasets but also maintains a low parameter footprint, enhancing the practical usability of large pre-trained models in varied commonsense reasoning applications.

### 4.2 Fine-tuning of Vision-Language Models

**Implementation Details.** The performance of various adaptation techniques on the Vision Transformer [7] (ViT-B/16) was evaluated across three distinct datasets: Food, Pets, and DTD. Each dataset was chosen to test the robustness and adaptability of these methods across different visual domains. The primary performance metric used was the Top-1 accuracy, which was computed as an average over three random seeds to mitigate randomness and ensure reliability in the results. The experiments were conducted under two shot settings, 4 and 16 shots, to determine the effectiveness of each adaptation technique under limited data conditions.

The adaptation methods tested included various enhancements of existing techniques such as CoOp [58], [58] PLOT++ [2], MaPLe [19], and several variants of LoRA, such as CLIP-LoRA [52], AdaLoRA [54], DyLoRA [46], and others. This comprehensive approach allowed for a detailed assessment of each method's ability to improve model performance under constrained training scenarios.

Table 1: Averaged accuracies (%) for 8 zero-shot tasks. Param denotes the number of trained parameters, Time for the training time on H800 GPU, and Mem for the GPU Memory usage.

| Method | Param | Time | Mem | BoolQ | PIQA | SIQA | HellaSwag | WinoGrande | ARC-e | ARC-c | OBQA | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| *Fine-tuning on LLaMA-1 7B* | | | | | | | | | | | | |
| LoRA$_{r=16}$ [17] | 8.4M | 5.7h | 21G | 68.9 | 80.7 | 77.4 | 78.1 | 78.8 | 77.8 | 61.3 | 74.8 | 74.4 |
| LoRA$_{r=32}$ [17] | 16.8M | 6.5h | 27G | 68.5 | 81.0 | 77.4 | 77.1 | 79.0 | 77.8 | 63.3 | 77.9 | 75.3 |
| **NoRA** | 8.2M | 5.0h | 19G | 68.1 | 80.3 | 76.8 | 80.6 | 79.6 | 80.5 | 62.6 | 77.8 | 75.8 |
| *Fine-tuning on LLaMA-3 8B* | | | | | | | | | | | | |
| LoRA [17] | 28.3M | 8.0h | 29G | 72.3 | 86.7 | 79.3 | 93.5 | 84.8 | 87.7 | 75.7 | 82.8 | 82.8 |
| LoKr [49] | 0.9M | 26.3h | 66G | 65.1 | 81.6 | 78.7 | 92.0 | 82.1 | 89.2 | 76.7 | 80.9 | 80.9 |
| AdaLoRA [54] | 28.3M | 12.5h | 58G | 75.1 | 86.4 | 76.7 | 75.4 | 83.3 | 90.4 | 79.1 | 81.4 | 81.4 |
| **NoRA** | 7.2M | 6.2h | 28G | 73.3 | 86.4 | 79.1 | 94.1 | 84.3 | 88.2 | 77.5 | 85.0 | 83.1 |

Table 2: Detailed results for 3 datasets with the ViT-B/16 as visual backbone. Top-1 accuracy averaged over 3 random seeds is reported. Highest value is highlighted in bold, and the second highest is underlined.

| Shots 4 | | | | | Shots 16 | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Method | Food | Pets | DTD | Average | Method | Food | Pets | DTD | Average |
| CoOp [58] (4) | 83.5 | 92.3 | 58.5 | 78.1 | CoOp [58] (4) | 85.1 | 92.4 | 81.2 | 86.2 |
| CoOp [58] (16) | 84.5 | 92.5 | 59.5 | 78.8 | CoOp [58] (16) | 84.2 | 92.0 | 69.7 | 81.9 |
| CoCoOp [57] | 86.3 | 92.7 | 55.7 | 78.2 | CoCoOp [57] | 87.4 | 93.4 | 63.7 | 81.5 |
| TIP-Adapter-F [55] | 86.5 | 91.9 | 59.8 | 79.4 | TIP-Adapter-F [55] | 86.8 | 92.6 | 70.8 | 83.4 |
| CLIP-Adapter [8] | 86.5 | 90.8 | 46.1 | 74.5 | CLIP-Adapter [55] | 87.1 | 92.3 | 59.4 | 79.6 |
| PLOT++ [2] | 86.5 | 92.6 | 62.4 | 80.5 | PLOT++ [2] | 87.1 | 93.6 | 71.4 | 84.0 |
| KgCoOp [48] | 86.9 | 92.6 | 58.7 | 79.4 | KgCoOp [48] | 87.2 | 93.2 | 68.7 | 83.0 |
| TaskRes [50] | 86.0 | 91.9 | 60.1 | 79.3 | TaskRes [50] | 86.9 | 92.4 | 71.5 | 83.6 |
| MaPLe [19] | 86.7 | **93.3** | 59.0 | 79.7 | MaPLe [19] | 87.4 | 93.2 | 68.4 | 83.0 |
| ProGrad [60] | 85.4 | 92.1 | 59.7 | 79.1 | ProGrad [60] | 85.8 | 92.8 | 68.8 | 82.5 |
| CLIP-LoRA [52] | 82.7 | 91.0 | 63.8 | 79.2 | CLIP-LoRA [52] | 84.2 | 92.4 | 72.0 | 82.9 |
| LoRA+ [13] | 84.4 | 92.8 | 64.1 | 81.4 | LoRA+ [13] | 85.1 | 93.6 | 72.1 | 83.6 |
| AdaLoRA [54] | 85.6 | 92.8 | **66.2** | <u>81.5</u> | AdaLoRA [54] | 85.9 | 93.7 | <u>72.8</u> | 84.1 |
| DyLoRA [46] | <u>87.0</u> | 92.4 | 64.9 | 81.4 | DyLoRA [46] | <u>87.6</u> | 93.0 | 72.7 | <u>84.4</u> |
| LoRA-FA [53] | 86.7 | 93.0 | 64.4 | 81.4 | LoRA-FA [53] | 87.4 | <u>93.9</u> | 71.9 | <u>84.4</u> |
| VeRA [30] | 84.5 | 92.5 | 65.1 | 80.7 | VeRA [30] | 86.2 | 92.2 | 72.2 | 83.5 |
| **NoRA** | **87.1** | <u>93.1</u> | <u>65.2</u> | **81.8** | **NoRA** | **87.8** | **94.1** | **74.3** | **85.4** |

**Comparison Results.** The detailed results are presented in Table 2, which includes the Top-1 accuracy for each method across the three datasets under both 4-shot and 16-shot settings. Notably, the NoRA model consistently outperformed other adaptation methods in both settings, highlighting its superior adaptability and efficiency. In the 4-shot setting, NoRA achieved an average Top-1 accuracy rate of 81.8, slightly surpassing DyLoRA, the second-best method, by 0.2%. In the more demanding 16-shot setting, NoRA demonstrated even greater efficacy, achieving an average Top-1 accuracy of 85.4 and surpassing DyLoRA's score of 85.0.

NoRA showed exceptional robustness across the various visual domains tested. Specifically, in the 16-shot setting, NoRA achieved the best results in all individual datasets, with scores of 87.8 for Food, 94.1 for Pets, and 74.3 for DTD. These scores underline NoRA's capability to adapt to and perform well across diverse categories, ensuring high reliability for practical applications in fields requiring visual recognition.

Overall, the experimental outcomes underscore the efficacy of NoRA as a scalable and effective fine-tuning strategy, achieving high performance across diverse datasets and maintaining a low parameter footprint, which is crucial for practical deployments where computational efficiency is essential. These results confirm the importance of selecting appropriate adaptation techniques to enhance the performance of vision transformers in limited data environments.

## 4.3 Subject-driven Generation

**Implementation Details.** We explored theme-based image generation by utilizing advanced text-to-image diffusion models, specifically employing a pre-trained text-to-image model. This model was fine-tuned using a combination of images and specific textual prompts (e.g., "[V] photo of a cat"). The fine-tuning involved sophisticated adaptation techniques, namely the LoRA and NoRA methods, to enhance the model's capability to generate images that align closely with given themes.

Figure 2: Comparison of generated images from LoRA and NoRA on the subject-driven generation task.

The model used in our experiments, referred to as SDXL5, was fine-tuned on a 32G V100S GPU. The parameters adjusted during fine-tuning included setting the learning rate to 1e-4 and the batch size to 4. The NoRA method involved initializing the internal and external layers of LoRA with a differential diagonal matrix. The training was conducted over 500 steps, taking approximately 24 minutes to complete.

**Comparison Results.** The results of the image generation process are detailed in Figure 2. During the image generation phase, we performed 50 inference steps for each textual prompt. Our findings indicate that the NoRA method outperformed the traditional LoRA method in several key aspects. Notably, the NoRA method was more effective in capturing complex themes and details within the images. The images generated using the NoRA method showed a higher degree of visual alignment with the specified textual prompts. They better reflected specific environmental contexts, such as forests and snowy scenes, and captured intricate object details, like wet fur.

These outcomes underscore the efficacy of the NoRA method in enhancing thematic consistency and visual expressiveness. The enhanced model demonstrates significant improvements in generating thematic images that are closely aligned with the nuances of the input prompts, offering promising applications in fields requiring detailed and context-specific image generation. This experiment sets a foundation for further exploration into fine-tuning techniques that can more adeptly handle complex thematic prompts in image generation tasks.

## 5 Conclusion

In this paper we present introduces NoRA, a new parameter-efficient fine-tuning method for large models. NoRA employs a dual-layer nested structure combined with SVD to extract latent information from the original matrix while reducing parameter count. The method retains LoRA's parameter efficiency while improving low-rank matrix approximation accuracy through SVD. NoRA's hierarchical structure allows for finer optimization control, enhancing model adaptability across tasks. Experimental results demonstrate NoRA's effectiveness across various tasks, including commonsense reasoning in LLMs, VLMs and subject-driven generation. The significance of NoRA lies in its ability to enhance model adaptability and performance while maintaining excellent parameter efficiency. Our approach addresses the growing need for efficient fine-tuning methods in the era of large models.

Limitations of the study include potential computational overhead due to SVD calculations and the need for further investigation into NoRA's scalability to even larger models. Future work could explore combining NoRA with AutoML [6, 5, 4, 61, 25, 24] and distillation techniques [23, 22, 20, 47, 21], applying it to multimodal models, and investigating its impact on model interpretability and robustness.

# References

[1] Klaudia Bałazy, Mohammadreza Banaei, Karl Aberer, and Jacek Tabor. Lora-xs: Low-rank adaptation with extremely small number of parameters. *arXiv preprint arXiv:2405.17604*, 2024.

[2] Guangyi Chen, Weiran Yao, Xiangchen Song, Xinyue Li, Yongming Rao, and Kun Zhang. Plot: Prompt learning with optimal transport for vision-language models. *arXiv preprint arXiv:2210.01253*, 2022.

[3] Ning Ding, Yujia Qin, Guang Yang, Fuchao Wei, Zonghan Yang, Yusheng Su, Shengding Hu, Yulin Chen, Chi-Min Chan, Weize Chen, et al. Parameter-efficient fine-tuning of large-scale pre-trained language models. *Nature Machine Intelligence*, 5(3):220–235, 2023.

[4] Peijie Dong, Lujun Li, Zhenheng Tang, Xiang Liu, Xinglin Pan, Qiang Wang, and Xiaowen Chu. Pruner-zero: Evolving symbolic pruning metric from scratch for large language models. In *ICML*, 2024.

[5] Peijie Dong, Lujun Li, and Zimian Wei. Diswot: Student architecture search for distillation without training. In *CVPR*, 2023.

[6] Peijie Dong, Lujun Li, Zimian Wei, Xin Niu, Zhiliang Tian, and Hengyue Pan. Emq: Evolving training-free proxies for automated mixed precision quantization. In *ICCV*, 2023.

[7] 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.

[8] Peng Gao, Shijie Geng, Renrui Zhang, Teli Ma, Rongyao Fang, Yongfeng Zhang, Hongsheng Li, and Yu Qiao. Clip-adapter: Better vision-language models with feature adapters. *International Journal of Computer Vision*, 132(2):581–595, 2024.

[9] Gene H Golub and Charles F Van Loan. *Matrix computations*. JHU press, 2013.

[10] Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A Smith. Don't stop pretraining: Adapt language models to domains and tasks. *arXiv preprint arXiv:2004.10964*, 2020.

[11] Zeyu Han, Chao Gao, Jinyang Liu, Sai Qian Zhang, et al. Parameter-efficient fine-tuning for large models: A comprehensive survey. *arXiv preprint arXiv:2403.14608*, 2024.

[12] Yongchang Hao, Yanshuai Cao, and Lili Mou. Flora: Low-rank adapters are secretly gradient compressors. *arXiv preprint arXiv:2402.03293*, 2024.

[13] Soufiane Hayou, Nikhil Ghosh, and Bin Yu. Lora+: Efficient low rank adaptation of large models. *arXiv preprint arXiv:2402.12354*, 2024.

[14] Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. Towards a unified view of parameter-efficient transfer learning. *arXiv preprint arXiv:2110.04366*, 2021.

[15] Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for nlp. In *International conference on machine learning*, pages 2790–2799. PMLR, 2019.

[16] 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.

[17] 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.

[18] Zhiqiang Hu, Lei Wang, Yihuai Lan, Wanyu Xu, Ee-Peng Lim, Lidong Bing, Xing Xu, Soujanya Poria, and Roy Ka-Wei Lee. Llm-adapters: An adapter family for parameter-efficient fine-tuning of large language models. *arXiv preprint arXiv:2304.01933*, 2023.

[19] Muhammad Uzair Khattak, Hanoona Rasheed, Muhammad Maaz, Salman Khan, and Fahad Shahbaz Khan. Maple: Multi-modal prompt learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19113–19122, 2023.

[20] Lujun Li. Self-regulated feature learning via teacher-free feature distillation. In *ECCV*, 2022.

[21] Lujun Li, Yufan Bao, Peijie Dong, Chuanguang Yang, Anggeng Li, Wenhan Luo, Qifeng Liu, Wei Xue, and Yike Guo. Detkds: Knowledge distillation search for object detectors. In *ICML*, 2024.

[22] Lujun Li, Peijie Dong, Anggeng Li, Zimian Wei, and Ya Yang. Kd-zero: Evolving knowledge distiller for any teacher-student pairs. *NeuIPS*, 2024.

[23] Lujun Li and Zhe Jin. Shadow knowledge distillation: Bridging offline and online knowledge transfer. In *NeuIPS*, 2022.

[24] Lujun Li, Haosen Sun, Shiwen Li, Peijie Dong, Wenhan Luo, Wei Xue, Qifeng Liu, and Yike. Guo. Auto-gas: Automated proxy discovery for training-free generative architecture search. In *ECCV*, 2024.

[25] Lujun Li, Zimian Wei, Peijie Dong, Wenhan Luo, Wei Xue, Qifeng Liu, and Yike. Guo. Attnzero: Efficient attention discovery for vision transformers. In *ECCV*, 2024.

[26] Xiang Li, Aynaz Taheri, Lifu Tu, and Kevin Gimpel. Commonsense knowledge base completion. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1445–1455, 2016.

[27] Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. *arXiv preprint arXiv:2101.00190*, 2021.

[28] Yang Li, Shaobo Han, and Shihao Ji. Vb-lora: Extreme parameter efficient fine-tuning with vector banks. *arXiv preprint arXiv:2405.15179*, 2024.

[29] Shanchuan Lin, Anran Wang, and Xiao Yang. Sdxl-lightning: Progressive adversarial diffusion distillation. *arXiv preprint arXiv:2402.13929*, 2024.

[30] Jiacheng Liu, Wenya Wang, Dianzhuo Wang, Noah A Smith, Yejin Choi, and Hannaneh Hajishirzi. Vera: A general-purpose plausibility estimation model for commonsense statements. *arXiv preprint arXiv:2305.03695*, 2023.

[31] Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Computing Surveys*, 55(9):1–35, 2023.

[32] 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. *arXiv preprint arXiv:2402.09353*, 2024.

[33] 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. *arXiv preprint arXiv:2402.09353*, 2024.

[34] Weiyang Liu, Zeju Qiu, Yao Feng, Yuliang Xiu, Yuxuan Xue, Longhui Yu, Haiwen Feng, Zhen Liu, Juyeon Heo, Songyou Peng, et al. Parameter-efficient orthogonal finetuning via butterfly factorization. *arXiv preprint arXiv:2311.06243*, 2023.

[35] Zeyu Liu, Souvik Kundu, Anni Li, Junrui Wan, Lianghao Jiang, and Peter Anthony Beerel. Aflora: Adaptive freezing of low rank adaptation in parameter efficient fine-tuning of large models. *arXiv preprint arXiv:2403.13269*, 2024.

[36] Kai Lv, Yuqing Yang, Tengxiao Liu, Qinghui Gao, Qipeng Guo, and Xipeng Qiu. Full parameter fine-tuning for large language models with limited resources. *arXiv preprint arXiv:2306.09782*, 2023.

[37] Yuren Mao, Yuhang Ge, Yijiang Fan, Wenyi Xu, Yu Mi, Zhonghao Hu, and Yunjun Gao. A survey on lora of large language models. *arXiv preprint arXiv:2407.11046*, 2024.

[38] Fanxu Meng, Zhaohui Wang, and Muhan Zhang. Pissa: Principal singular values and singular vectors adaptation of large language models. *arXiv preprint arXiv:2404.02948*, 2024.

[39] Jonas Pfeiffer, Aishwarya Kamath, Andreas Rücklé, Kyunghyun Cho, and Iryna Gurevych. Adapterfusion: Non-destructive task composition for transfer learning. *arXiv preprint arXiv:2005.00247*, 2020.

[40] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023.

[41] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.

[42] Adithya Renduchintala, Tugrul Konuk, and Oleksii Kuchaiev. Tied-lora: Enhacing parameter efficiency of lora with weight tying. *arXiv preprint arXiv:2311.09578*, 2023.

[43] Maarten Sap, Vered Shwartz, Antoine Bosselut, Yejin Choi, and Dan Roth. Commonsense reasoning for natural language processing. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: Tutorial Abstracts*, pages 27–33, 2020.

[44] Chunlin Tian, Zhan Shi, Zhijiang Guo, Li Li, and Chengzhong Xu. Hydralora: An asymmetric lora architecture for efficient fine-tuning. *arXiv preprint arXiv:2404.19245*, 2024.

[45] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.

[46] 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.

[47] Liu Xiaolong, Li Lujun, Li Chao, and Anbang Yao. Norm: Knowledge distillation via n-to-one representation matching. In *ICLR*, 2023.

[48] Hantao Yao, Rui Zhang, and Changsheng Xu. Visual-language prompt tuning with knowledge-guided context optimization. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6757–6767, 2023.

[49] Shih-Ying Yeh, Yu-Guan Hsieh, Zhidong Gao, Bernard BW Yang, Giyeong Oh, and Yanmin Gong. Navigating text-to-image customization: From lycoris fine-tuning to model evaluation. In *The Twelfth International Conference on Learning Representations*, 2023.

[50] Tao Yu, Zhihe Lu, Xin Jin, Zhibo Chen, and Xinchao Wang. Task residual for tuning vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10899–10909, 2023.

[51] Elad Ben Zaken, Shauli Ravfogel, and Yoav Goldberg. Bitfit: Simple parameter-efficient fine-tuning for transformer-based masked language-models. *arXiv preprint arXiv:2106.10199*, 2021.

[52] Maxime Zanella and Ismail Ben Ayed. Low-rank few-shot adaptation of vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1593–1603, 2024.

[53] Longteng Zhang, Lin Zhang, Shaohuai Shi, Xiaowen Chu, and Bo Li. Lora-fa: Memory-efficient low-rank adaptation for large language models fine-tuning. *arXiv preprint arXiv:2308.03303*, 2023.

[54] Qingru Zhang, Minshuo Chen, Alexander Bukharin, Pengcheng He, Yu Cheng, Weizhu Chen, and Tuo Zhao. Adaptive budget allocation for parameter-efficient fine-tuning. In *The Eleventh International Conference on Learning Representations*, 2023.

[55] Renrui Zhang, Wei Zhang, Rongyao Fang, Peng Gao, Kunchang Li, Jifeng Dai, Yu Qiao, and Hongsheng Li. Tip-adapter: Training-free adaption of clip for few-shot classification. In *European conference on computer vision*, pages 493–510. Springer, 2022.

[56] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 2023.

[57] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Conditional prompt learning for vision-language models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16816–16825, 2022.

[58] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for vision-language models. *International Journal of Computer Vision*, 130(9):2337–2348, 2022.

[59] Xuhui Zhou, Yue Zhang, Leyang Cui, and Dandan Huang. Evaluating commonsense in pre-trained language models. In *AAAI*, 2020.

[60] Beier Zhu, Yulei Niu, Yucheng Han, Yue Wu, and Hanwang Zhang. Prompt-aligned gradient for prompt tuning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15659–15669, 2023.

[61] Zimian Zimian Wei, Lujun Li Li, Peijie Dong, Zheng Hui, Anggeng Li, Menglong Lu, Hengyue Pan, and Dongsheng Li. Auto-prox: Training-free vision transformer architecture search via automatic proxy discovery. In *AAAI*, 2024.