

Astra: Activation-Space Tail-Eigenvector Low-Rank Adaptation of Large Language Models

Kainan Liu[†], Yong Zhang[†], Ning Cheng^{*}, Yun Zhu,
Yanmeng Wang, Shaojun Wang, Jing Xiao,

Ping An Technology (Shenzhen) Co., Ltd., China

 github.com/LyoAI/Astra

Abstract

Parameter-Efficient Fine-Tuning (PEFT) methods, especially LoRA, are widely used for adapting pre-trained models to downstream tasks due to their computational and storage efficiency. However, in the context of LoRA and its variants, the potential of activation subspaces corresponding to tail eigenvectors remains substantially under-exploited, which may lead to suboptimal fine-tuning performance. In this work, we propose **Astra** (Activation-Space Tail-Eigenvector Low-Rank Adaptation), a novel PEFT method that leverages the tail eigenvectors of the model output activations—estimated from a small task-specific calibration set—to construct task-adaptive low-rank adapters. By constraining updates to the subspace spanned by these tail eigenvectors, Astra achieves faster convergence and improved downstream performance with a significantly reduced parameter budget. Extensive experiments across natural language understanding (NLU) and natural language generation (NLG) tasks demonstrate that Astra consistently outperforms existing PEFT baselines across 16 benchmarks and even surpasses full fine-tuning (FFT) in certain scenarios.

1 Introduction

Large Language Models (LLMs) have achieved remarkable success across a wide range of tasks (Achiam et al., 2023; Dubey et al., 2024; Guo et al., 2025). Adapting these pretrained models typically involves full fine-tuning (FFT). Despite its efficacy, FFT incurs substantial computational and memory overhead, limiting its applicability in resource-constrained settings (Singh et al., 2024; Liu et al., 2024a). To overcome these limitations, parameter-efficient fine-tuning (PEFT) methods have been proposed as an effective alternative.

By introducing a small number of trainable parameters while keeping the pretrained backbone frozen, PEFT methods substantially reduce training costs while retaining competitive performance (Liu et al., 2021; Li and Liang, 2021; Hu et al., 2023b).

Among existing PEFT approaches, Low-Rank Adaptation (LoRA) (Hu et al., 2022) stands out for its simplicity and strong empirical performance. However, the standard initialization of LoRA often yields extremely small gradients at the early training stage, which can hinder optimization and lead to slow convergence or suboptimal adaptation (Meng et al., 2024; Wang et al., 2024b). To address this issue, recent studies have explored improved initialization strategies, which can be broadly classified into two categories: **weight-driven** methods (Meng et al., 2024; Wang et al., 2024a) that exploit the structure of pretrained weights, and **data-driven** methods (Yang et al., 2024; Wang et al., 2024b; Paischer et al., 2024) that leverage data distributions or task-specific signals to guide the initialization of low-rank adapters.

However, most existing works overlook an important characteristic of LLM representations: dimensions corresponding to tail eigenvalues remain largely under-utilized during fine-tuning (Nayak et al., 2025). While these directions contribute little to the dominant pretrained representations, they offer substantial flexibility for task-specific adaptation. Adapting model parameters within such under-explored subspaces can effectively increase the model’s representational rank (Roy and Vetterli, 2007), thereby enhancing expressive capacity and improving adaptability to downstream tasks.

Building on this insight, we propose **Astra** (Activation-Space Tail-Eigenvector Low-Rank Adaptation), a novel PEFT method that exploits under-explored tail subspaces of output activations to construct task-adaptive low-rank adapters. Specifically, Astra first performs eigendecomposition on the covariance matrix of output activations

[†] Equal contribution. ^{*} Corresponding author.

estimated from a small task-specific calibration dataset D , i.e., $Cov(Y) = Q\Lambda Q^\top$, where Q denotes the eigenvectors and Λ is the diagonal matrix of corresponding eigenvalues. To explicitly constrain optimization within the under-utilized subspaces, Astra projects the pretrained weight matrix W onto the subspace spanned by the tail eigenvectors, yielding low-rank adapters aligned with the under-utilized activation directions, i.e., $A = Q_{[:, -r]}^\top W$ and $B = Q_{[:, -r]}$, where r denotes the LoRA rank.

This initialization strategy offers several advantages. By focusing adaptation on activation dimensions that are weakly optimized during pretraining, Astra increases the effective representational rank while avoiding interference with dominant pretrained features. As a result, it enables more expressive task-specific updates with improved optimization stability and faster convergence.

We conduct extensive experiments across a diverse set of tasks to evaluate the effectiveness of Astra, covering both natural language understanding (NLU) and natural language generation (NLG) benchmarks. Experimental results show that Astra consistently outperforms existing PEFT baselines on 16 benchmarks and even surpasses full fine-tuning (FFT) on several tasks. Our main contributions can be summarized as follows:

- We propose **Astra**, a novel LoRA initialization method that exploits under-utilized eigenspaces of output activations to enable effective low-rank adaptation.
- We conduct extensive evaluations on a broad spectrum of NLU and NLG tasks, including general language understanding, mathematical reasoning, code generation, and commonsense reasoning. The results demonstrate that Astra consistently outperforms strong PEFT baselines and exhibits robust adaptability across tasks.
- We perform comprehensive ablation studies on eigenvector selection, LoRA rank, and calibration data, systematically validating the effectiveness and efficiency of Astra.
- We provide an effective-rank analysis that empirically supports our core hypothesis: adapting within under-explored activation subspaces enhances task-specific representational capacity and improves downstream task performance.

2 Related Work

PEFT. Parameter-efficient fine-tuning (PEFT) provides a lightweight alternative to full fine-tuning by updating only a small subset of parameters, reducing computational overhead while maintaining strong downstream performance. PEFT methods can be broadly categorized into: (i) *prompt-based* (Li and Liang, 2021; Liu et al., 2021), which prepend trainable tokens or embeddings to the input; (ii) *adapter-based* (Houlsby et al., 2019; Rücklé et al., 2020; Hu et al., 2023b), which insert small trainable modules into each transformer layer; and (iii) *LoRA-based* (Hu et al., 2022; Dettmers et al., 2023), which employ low-rank reparameterization to enable efficient adaptation.

LoRA and Its Variants. Low-rank adaptation (LoRA) has received significant attention for enabling fine-tuning without altering the original architecture or adding inference overhead (Li et al., 2018; Aghajanyan et al., 2021). Building on LoRA, subsequent research has explored several directions to enhance its flexibility and efficiency:

Dynamic rank allocation methods (Valipour et al., 2022; Liu et al., 2024c), such as AdaLoRA (Zhang et al., 2023b), adaptively distribute parameter budgets across weight matrices based on importance scores.

Structural modifications (Liu et al., 2024b; Feng et al., 2024; Li et al., 2024) generalize LoRA beyond its original design. For example, DoRA (Liu et al., 2024b) decouples learning into magnitude and direction.

Hyperparameter optimization has also been explored to improve fine-tuning efficiency (Kala-Jdziewski, 2023; Hayou et al., 2024). For instance, LoRA+ (Hayou et al., 2024) introduces differential learning rates for the low-rank matrices A and B, with a higher learning rate for B to accelerate convergence.

Initialization Strategies. Recent efforts have also explored initialization strategies to stabilize training and accelerate convergence (Meng et al., 2024; Wang et al., 2024a; Yang et al., 2024; Wang et al., 2024b), which can be divided into **weight-driven** and **data-driven** approaches.

Weight-driven methods primarily analyze the static geometric properties of pre-trained parameters. For example, PiSSA (Meng et al., 2024) applies Singular Value Decomposition (SVD) to the pretrained weight matrix $W^{(0)}$, capturing the principal singular values to initialize the adapters while

relegating the residual components to the frozen weights. While effective, weight-driven methods neglect the dynamic input distribution encountered during inference. In contrast, data-driven approaches incorporate task-specific information into the initialization process. LoRA-GA (Wang et al., 2024b) introduces a novel initialization scheme that aligns the initial gradients of the low-rank update matrices with those of full fine-tuning.

Our proposed framework, Astra, represents a distinct paradigm within the data-driven category. Unlike gradient-based methods like LoRA-GA that require computationally expensive gradient estimation, Astra operates on the activation-space covariance. By identifying the tail subspace of the output activations, Astra ensures that the adapters are initialized within directions that are statistically under-utilized by the pre-trained model yet functionally critical for the downstream task. This unique activation-centric perspective allows Astra to capture task-specific functional redundancies that are overlooked by purely geometric weight analysis during adaptation. A detailed comparison of LoRA variants is provided in Appendix A.

3 Method

We introduce Astra, a LoRA initialization method based on the observation that pretrained models exhibit uneven utilization across activation dimensions. In standard LoRA, low-rank adapters are initialized without considering this structural property, which may lead to inefficient early-stage optimization. Astra initializes LoRA updates within activation subspaces associated with smaller eigenvalues, aiming to make use of under-utilized directions during task-specific adaptation.

3.1 Preliminaries of LoRA’s Initialization

LoRA (Hu et al., 2022) introduces trainable updates by reparameterizing weight modifications as the product of two low-rank matrices. Formally, given a pretrained weight matrix $W_0 \in \mathbb{R}^{m \times n}$, LoRA expresses the adapted weight as:

$$\tilde{W} = W_0 + \Delta W = W_0 + \frac{\alpha}{r}BA \quad (1)$$

where ΔW denotes the weight change, which is decomposed into two low-rank matrices $B \in \mathbb{R}^{m \times r}$ and $A \in \mathbb{R}^{r \times n}$ with an intrinsic rank $r \ll \min(m, n)$, α is a scaling constant. This parameterization reduces the number of trainable parameters

from mn to $(m+n)r$, significantly improving fine-tuning efficiency. In practice, A is initialized from the Gaussian distribution, while B is initialized as an all-zero matrix to ensure that the initial model output remains unchanged. However, such random initialization can lead to slower convergence, as the gradients of the trainable adapters can be very small or in random directions during the early stages of fine-tuning (Meng et al., 2024).

3.2 Activation-Space Tail-Eigenvector Low-Rank Adaptation

Astra initializes LoRA adapters by constraining low-rank updates to a task-specific tail subspace of the output activation space. The core mechanism involves two stages: (1) performing orthogonal decomposition of the output activations and (2) projecting the pre-trained weights onto the tail subspace to initialize the LoRA adapters. The detailed formulation of each step is presented below.

Step 1: Orthogonal Decomposition. Astra begins by characterizing the structure of the output activations relevant to the downstream task through covariance analysis. Specifically, we construct a calibration dataset $X = \{x_i\}_{i=1}^N$ by randomly sampling $N = 64$ samples from the training data. Let $Y = WX + b \in \mathbb{R}^{d_{\text{out}} \times N}$ denote the collective output activations of a linear layer. We then compute the corresponding covariance matrix as:

$$\text{Cov}(Y) = \mathbb{E}[YY^\top] - \mathbb{E}[Y]\mathbb{E}[Y]^\top \quad (2)$$

Since $\text{Cov}(Y)$ is positive semi-definite, it admits an eigendecomposition $\text{Cov}(Y) = Q\Lambda Q^\top$. Detailed calibration procedures are provided in Appendix E.

Here, the orthogonal matrix $Q = [Q_{\text{main}} | Q_{\text{tail}}]$ partitions the eigenvectors into two subspaces, where $Q_{\text{main}} \in \mathbb{R}^{d_{\text{out}} \times (d_{\text{out}} - r)}$ corresponds to the principal subspace associated with the dominant eigenvalues, and $Q_{\text{tail}} \in \mathbb{R}^{d_{\text{out}} \times r}$ spans the residual subspace defined by the tail eigenvalues. The two subspaces are orthogonal and satisfy $Q_{\text{main}}^\top Q_{\text{tail}} = 0$. Utilizing these orthonormal bases, we perform an orthogonal decomposition of the output activations Y into two mutually orthogonal projections:

$$Y = Y_{\text{main}} + Y_{\text{tail}} \quad (3)$$

$$= Q_{\text{main}}Q_{\text{main}}^\top Y + Q_{\text{tail}}Q_{\text{tail}}^\top Y \quad (4)$$

$$= \underbrace{Q_{\text{main}}Q_{\text{main}}^\top Wx + b}_{\text{Frozen}} + \underbrace{Q_{\text{tail}}Q_{\text{tail}}^\top Wx}_{\text{Trainable}} \quad (5)$$

where $Y_{\text{main}} := Q_{\text{main}} Q_{\text{main}}^\top Y$ and $Y_{\text{tail}} := Q_{\text{tail}} Q_{\text{tail}}^\top Y$. Given $Q_{\text{main}}^\top Q_{\text{tail}} = 0$, it follows that $Y_{\text{main}}^\top Y_{\text{tail}} = 0$, ensuring that the two components occupy non-interfering subspaces.

This orthogonal decomposition forms the basis for restricting subsequent parameter updates to the residual activation subspace spanned by Q_{tail} , thereby exploiting directions that are weakly activated under the pretrained model yet remain amenable to task-specific modulation.

Step 2: Tail-Subspace Projection. Building upon the decomposition in Eq.5, Astra initializes the two learnable low-rank matrices $A \in \mathbb{R}^{r \times d_{\text{in}}}$ and $B \in \mathbb{R}^{d_{\text{out}} \times r}$ by projecting the pre-trained weights onto the identified tail subspace. Formally, the initialization is defined as:

$$A_{\text{init}} = Q_{\text{tail}}^\top W \in \mathbb{R}^{r \times d_{\text{in}}}. \quad (6)$$

$$B_{\text{init}} = Q_{\text{tail}} \in \mathbb{R}^{d_{\text{out}} \times r}, \quad (7)$$

where A_{init} and B_{init} serve as the two learnable low-rank matrices in LoRA. Since the resulting update $\Delta W = BA$ is non-zero at initialization, we adjust the frozen component to ensure that the original model outputs remain unchanged. Formally, this yields:

$$W' = W^{(0)} + \Delta W = \underbrace{(W^{(0)} - B_{\text{init}} A_{\text{init}})}_{\text{Frozen}} + \underbrace{B' A'}_{\text{Trainable}} \quad (8)$$

where the learnable matrices A' and B' parameterize the task-specific update ΔW . By constraining the optimization to the subspace spanned by the tail eigenvectors, Astra effectively utilizes previously under-explored directions, thereby enhancing adaptation efficiency. The complete procedure is summarized in Algorithm 1, and a PyTorch-style implementation is provided in Appendix G.

4 Experiments

In this section, we provide a comprehensive evaluation of Astra from three perspectives. 1) We first assess the Natural Language Understanding (NLU) capabilities using the GLUE (Wang et al., 2018) benchmark (Section 4.2). 2) Next, we evaluate the performance of our method on Natural Language Generation (NLG) tasks, covering mathematical reasoning, code generation, and commonsense reasoning (Section 4.3). 3) Finally, we conduct ablation studies to analyze the effectiveness of our approach with respect to varying eigenvectors, LoRA

Algorithm 1: Astra: Activation-Space Tail-Eigenvector Low-Rank Adaptation

Input: Model M , LoRA rank r , calibration data x , weight matrices

$$W \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$$

Output: Initialized parameters

$$W_{\text{frozen}}, A_{\text{init}}, B_{\text{init}}$$

1: $\hat{Y} \leftarrow M(x; W)$ \triangleright Forward propagation

2: $\text{Cov}(Y) \leftarrow \mathbb{E}[YY^\top] - \mathbb{E}[Y] \mathbb{E}[Y]^\top$

3: $\text{Cov}(Y) = Q\Lambda Q^\top$ \triangleright
Eigen-decomposition

4: Initialize trainable low-rank matrices:

$$A_{\text{init}} = Q_{\text{tail}}^\top W \in \mathbb{R}^{r \times d_{\text{in}}}$$

$$B_{\text{init}} = Q_{\text{tail}} \in \mathbb{R}^{d_{\text{out}} \times r} \quad \triangleright \text{Astra}$$

Initialization

5: Compute frozen and update terms:

$$W_{\text{frozen}} = W^{(0)} - B_{\text{init}} A_{\text{init}}$$

$$W_{\text{trainable}} = B_{\text{init}} A_{\text{init}}$$

return $W_{\text{frozen}}, A_{\text{init}}, B_{\text{init}}$

ranks and calibration datasets (Section 4.4). All experiments are conducted on NVIDIA A100-SXM4 (80GB) GPUs.

4.1 Baselines

To substantiate the effectiveness of our method, we compare Astra against full fine-tuning (FFT), vanilla LoRA, and 6 representative LoRA variants. These variants can be grouped as follows:

1. Weight-driven initialization variants:

- *PiSSA* (Meng et al., 2024) initializes adapters with principal components and freezes the residual.
- *MiLoRA* (Wang et al., 2024a) initializes adapters with the smallest singular components.

2. Data-driven initialization variants:

- *CorDA* (Yang et al., 2024) builds adapters conditioned on context for task-specific adaptations.
- *LoRA-GA* (Wang et al., 2024b) constructs low-rank matrices by approximating the gradient from the first step of full fine-tuning.

3. Other LoRA variants (with modified structure, hyperparameters, etc.):

	#Params	MNLI 393k	SST-2 67k	QNLI 105k	CoLA 8.5k	MRPC 3.7K	Average
Full FT	223M	86.95 \pm 0.04	97.02 \pm 0.03	98.78 \pm 0.02	84.52 \pm 0.01	84.19 \pm 0.05	90.29
LoRA	3.2M	86.97 \pm 0.01	96.62 \pm 0.02	98.75 \pm 0.03	49.95 \pm 1.33	47.67 \pm 0.06	75.99
DoRA	3.4M	87.05 \pm 0.02	97.19 \pm 0.01	98.79 \pm 0.02	84.23 \pm 0.03	49.88 \pm 0.05	83.43
rsLoRA	3.2M	87.06 \pm 0.01	97.13 \pm 0.02	98.79 \pm 0.02	83.89 \pm 0.02	49.63 \pm 0.04	83.30
PiSSA	3.2M	87.01 \pm 0.01	97.08 \pm 0.01	98.82 \pm 0.01	84.80 \pm 0.01	82.84 \pm 0.01	90.11
CorDA	3.2M	87.11 \pm 0.03	97.19 \pm 0.02	98.81 \pm 0.05	84.71 \pm 0.22	69.12 \pm 0.23	87.39
LoRA-GA	3.2M	87.07 \pm 0.01	97.13 \pm 0.02	98.83 \pm 0.01	84.76 \pm 0.11	84.19 \pm 0.14	90.40
Ours	3.2M	87.09 \pm 0.01	96.45 \pm 0.01	98.83 \pm 0.01	87.87 \pm 0.06	88.36 \pm 0.12	91.72

Table 1: Performance of T5-base fine-tuned with different adaptation methods on 5 datasets of the GLUE benchmark. We report accuracy for all tasks, and the results are averaged over three runs with different random seeds. Bold values indicate the best performance.

- *rsLoRA* (Kalajdziewski, 2023) introduces a square-root scaling factor to LoRA.
- *DoRA* (Liu et al., 2024b) decomposes pre-trained weights into magnitude and direction components, tuning the magnitude and direction matrix separately.

4.2 Natural Language Understanding

Models and Datasets. We fine-tune the T5-base model (Raffel et al., 2020) on a subset of tasks from the GLUE benchmark (Wang et al., 2018), including MNLI, QNLI, SST-2, CoLA and MRPC. The model is evaluated on the corresponding development sets, and accuracy is reported as the evaluation metric for all tasks. Additional details regarding the benchmarks are presented in Appendix B.1.

Implementation Details. We follow the experimental setup described in (Wang et al., 2024b) to ensure a fair comparison. Specifically, we convert the labels into tokens (e.g., "positive" or "negative") and use the prompt tuning to fine-tune the model for 1 epoch on each dataset. The normalized probabilities assigned to these tokens are then used for classification. Further experimental setup and implementation details can be found in Appendix C.1.

Main Results. Table 1 presents the performance of T5-base fine-tuned with different adaptation methods on five GLUE datasets. Our proposed approach consistently surpasses existing baselines, achieving the highest average accuracy across all tasks. The improvement is particularly pronounced on low-resource datasets such as MRPC and CoLA, where effective utilization of gradient information

plays a critical role. These results suggest that our method can fully exploit the limited training signals, leading to stable and fast convergence even under data-scarce conditions.

4.3 Natural Language Generation

Models and Datasets. We conduct experiments using LLaMA2-7B (Touvron et al., 2023) and LLaMA3-8B (Dubey et al., 2024) across three NLG tasks: **Math**, **Code** and **Commonsense**.

- *Math*: For mathematical reasoning tasks, the models are fine-tuned on the MetaMathQA dataset (Yu et al., 2023) and evaluated on two widely used benchmarks, GSM8K (Cobbe et al., 2021) and MATH (Hendrycks et al., 2021), using PASS@1 accuracy as the evaluation metric.
- *Code*: To evaluate programming proficiency, we fine-tune the models on the CodeFeedback-Python105k dataset (Zheng et al., 2024) and assess performance on HumanEval (Chen et al., 2021) and MBPP (Austin et al., 2021) benchmarks. Additionally, we employ the EvalPlus framework (Liu et al., 2023) to test on the extended versions of these datasets, namely MBPP+ and HumanEval+, which provide more test cases compared to the original versions. We report the PASS@1 metric for these evaluations.
- *Commonsense*: For commonsense reasoning, the models are fine-tuned on the Commonsense170K dataset (Hu et al., 2023a) and tested on seven established benchmarks—BoolQ (Clark et al., 2019), PIQA (Bisk et al., 2020), HelLaSwag (Zellers et al., 2019), WinoGrande (Sakaguchi et al., 2020), ARC-e, ARC-c (Clark et al.,

Model	Method	#Params	GSM8K	Math	HumanEval	HumanEval+	MBPP	MBPP+	Average
LLaMA2-7B	Full FT	6738M	58.76	12.04	32.9	31.1	43.9	36.8	35.92
	LoRA	320M	41.40	5.42	22.0	20.1	34.9	27.2	25.17
	MiLoRA	320M	39.12	5.06	20.1	18.9	36.8	29.4	24.90
	PiSSA	320M	51.63	7.36	23.2	20.1	36.7	29.5	28.08
	CorDA	320M	52.99	8.08	25.0	23.2	36.2	29.6	29.18
	Ours	320M	55.19	8.98	25.0	23.2	38.4	31.2	30.33
LLaMA3-8B	Full FT	8366M	75.36	24.04	56.7	53.7	64.0	54.5	54.72
	LoRA	336M	73.31	24.24	53.7	48.8	65.6	54.8	53.41
	MiLoRA	336M	73.24	23.90	52.4	48.2	68.3	56.1	53.69
	PiSSA	336M	76.50	26.92	57.1	52.0	68.0	56.3	56.14
	CorDA	336M	77.26	26.52	55.5	50.0	67.7	57.7	55.78
	Ours	336M	77.56	27.92	57.7	53.0	68.4	58.2	57.13

Table 2: Comparison of full fine-tuning (Full FT) and several LoRA variants on 2 mathematical reasoning and 4 code generation benchmarks. The best PEFT results are highlighted in **bold**.

Model	Method	#Params	BoolQ	PIQA	HellaSwag	WinoGrande	ARC-e	ARC-c	OBQA	Average
LLaMA2-7B	Full FT	6738M	82.81	75.08	55.57	73.64	72.69	41.72	32.00	61.93
	LoRA	320M	79.97	78.35	57.30	68.82	78.07	46.16	32.40	63.01
	PiSSA	320M	83.03	78.18	57.52	70.72	78.41	47.35	33.40	64.09
	MiLoRA	320M	79.66	78.13	57.53	69.22	77.31	45.39	32.40	62.81
	CorDA	320M	82.87	78.45	56.24	71.82	75.67	43.09	33.00	63.02
	Ours	320M	83.76	78.51	57.28	71.74	79.63	48.72	34.20	64.83
LLaMA3-8B	Full FT	8366M	82.14	68.82	49.30	66.06	65.95	38.05	31.60	57.42
	LoRA	336M	85.02	79.76	59.88	74.74	82.53	53.50	34.00	67.06
	PiSSA	336M	86.76	80.47	60.63	76.64	81.94	52.82	36.00	67.89
	MiLoRA	336M	84.07	79.92	60.31	74.59	81.27	51.62	34.80	66.65
	CorDA	336M	85.84	80.74	60.43	76.56	82.70	54.44	35.20	67.99
	Ours	336M	86.48	80.41	60.02	78.22	82.87	53.99	36.60	68.37

Table 3: Zero-shot performance of LLaMA2-7B and LLaMA3-8B fine-tuned with different adaptation methods on seven commonsense reasoning benchmarks. The best PEFT results are shown in **bold**.

2018), and OpenBookQA (Mihaylov et al., 2018). All tasks are tested in a zero-shot setting using the LM-Evaluation-Harness framework (Gao et al., 2024).

Implementation Details. To ensure a fair comparison, we adopt the experimental configurations delineated in (Meng et al., 2024; Wang et al., 2024b; Yang et al., 2024). Specifically, we set the LoRA rank to 128, with the LoRA alpha consistently equal to the rank, and insert adapters into all linear layers of the base model. All the experiments were conducted on the first 100,000 samples from each dataset and trained for one epoch to reduce computational overhead. Additional implementation details are provided in the Appendix C.2.

Main Results. Table 2 summarizes the results on mathematical reasoning and code generation tasks, and Table 3 reports the performance on commonsense reasoning benchmarks. Overall, our approach consistently surpasses existing PEFT base-

lines, demonstrating robust generalization across diverse task categories. Below, we provide a breakdown of the results by task type:

- *Math:* Astra outperforms all other PEFT baselines on both the GSM8K and MATH datasets, achieving the best results overall, with the exception of a slight gap compared to Full FT on LLaMA2-7B. Figure 1 illustrates the loss curves and gradient norm trends during fine-tuning of LLaMA2-7B on the MetaMathQA dataset. Notably, Astra (with rank=8) converges faster than LoRA (rank=128), highlighting its efficiency in downstream task adaptation with minimal resources.
- *Code:* For code generation tasks, Astra also achieves outstanding results, even surpassing Full FT on LLaMA3-8B. Our method shows remarkable programming proficiency, as reflected in the results across HumanEval and MBPP benchmarks.

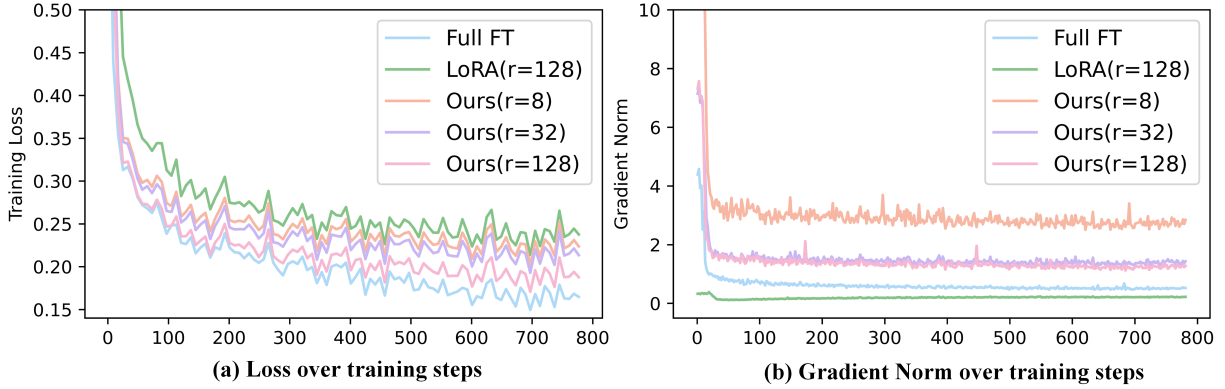


Figure 1: Training loss and gradient norm curves for FFT, LoRA (rank=128), and Astra with varying ranks on the MetaMathQA dataset. Our method (rank=8) performs even better than LoRA (rank=128), and higher ranks lead to faster loss reduction, approaching the performance of FFT.

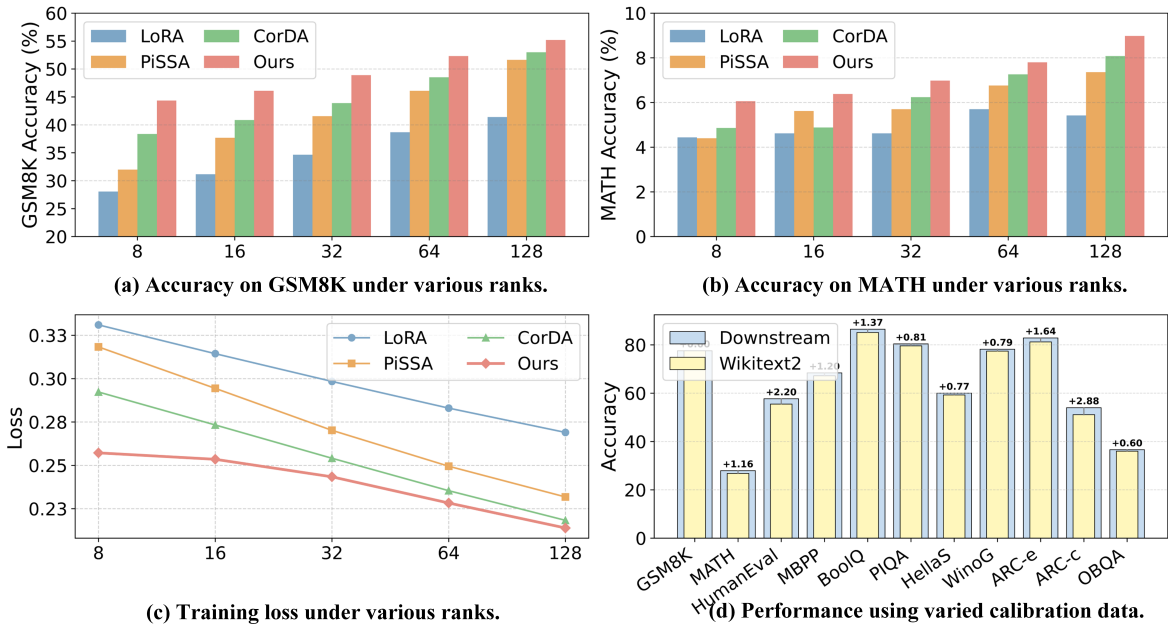


Figure 2: (a) and (b) report the performance of different LoRA variants on GSM8K and MATH under various ranks, respectively. (c) shows the final training loss on the MetaMathQA dataset under various ranks. (d) illustrates the performance using different calibration data.

- *Commonsense*: Astra demonstrates consistently strong performance across seven commonsense reasoning benchmarks. Although it slightly lags on HellaSwag, it achieves the best overall average performance among all baselines.

4.4 Ablation Studies

Eigenvectors. To investigate the impact of eigenvectors corresponding to eigenvalues of varying magnitudes on fine-tuning performance, we initialize the adapters injected into LLaMA2-7B with eigenvectors selected from different quantiles of the eigenvalue spectrum. Specifically, we use eigenvectors corresponding to the top, tail, middle, lower quartile, and upper quartile eigenvalues, as well

as randomly selected eigenvectors. The models are then fine-tuned on the MetaMathQA dataset and evaluated on the GSM8K and MATH benchmarks. As shown in Table 4, adapters initialized with tail eigenvectors achieve the best performance on both benchmarks, underscoring the efficacy of our strategy in leveraging tail eigenvectors from activation-space for fine-tuning.

LoRA Rank. In this experiment, we explore the effects of varying LoRA rank from 8 to 128, aiming to assess whether our approach consistently outperforms other PEFT baselines across different rank values. Following the setup described in Section 4.3, we fine-tune LLaMA2-7B on the Meta-

Eigenvectors	GSM8K	MATH
Random	40.49	5.64
Top	40.71	5.48
Upper Quartile (Q3)	40.49	5.64
Medium	38.74	5.60
Lower Quartile (Q1)	42.76	5.70
Tail	55.19	8.98

Table 4: Performance of LLaMA2-7B fine-tuned with adapters initialized using eigenvectors from different quantiles of the eigenvalue spectrum.

MathQA dataset and evaluate it on the GSM8K and MATH benchmarks. Figures 2 (a)-(b) show that Astra consistently outperforms alternative PEFT methods with the same number of trainable parameters. Figure 2 (c) illustrates the final training loss across different ranks, demonstrating that our method achieves a better fit to the training data compared to LoRA, PiSSA, and CorDA. It is noteworthy that our approach outperforms LoRA at rank = 128 even with rank = 8, underscoring its efficiency in achieving better performance with fewer trainable parameters.

Calibration Data. To assess the robustness of Astra with respect to the calibration datasets, we conduct experiments using a general-purpose dataset (i.e. Wikitext-2) for calibration, and compare it with the default setting, where the downstream training set itself is used for calibration. The results, presented in Figure 2 (d), demonstrate that Astra achieves stable performance across different calibration datasets, while leveraging the downstream training set yields marginally better results.

5 Discussion

Enhancing Representation Capacity via Increased Effective Rank. To evaluate the improvement in representational capacity introduced by our approach, we employ *effective rank* (Roy and Vetterli, 2007) as a metric to characterize the spectral structure of output activations before and after fine-tuning. Formally, the effective rank is defined as: The effective rank is formally defined as:

$$\mathcal{R}_{X,i} = \exp \left(- \sum_{j=1}^{d_{out}} \tilde{\lambda}_j \ln(\tilde{\lambda}_j) \right) \quad (9)$$

where $\tilde{\lambda}_j = \lambda_j / \sum_{k=1}^{d_{out}} \lambda_k$ represents the normalized eigenvalues and λ_j denotes the eigen-

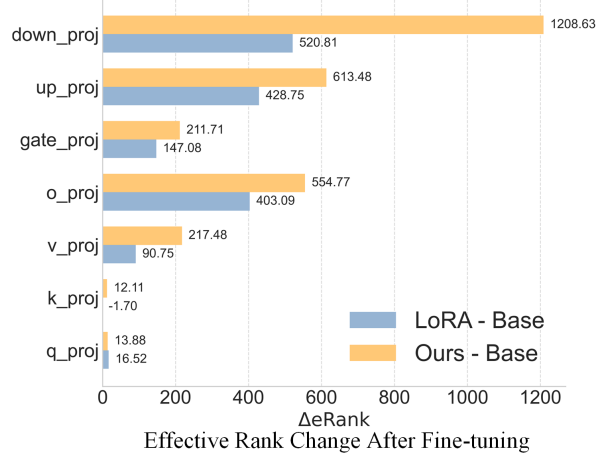


Figure 3: Comparison of effective rank before and after fine-tuning.

values obtained from the eigendecomposition of the output activation covariance matrix, $X \in \{Q, K, V, O, Up, Down\}$ are the projection layer type within the Transformer architecture, and i indexes the corresponding Transformer layer.

A higher effective rank indicates that the output activations are distributed across more directions in the feature space, suggesting a richer and more diverse representational capacity. Conversely, a lower effective rank (only a few eigenvalues are large) implies that the activations are concentrated along a few dominant directions, reflecting more constrained or redundant representations (Li et al., 2025).

For each layer type, we aggregate the effective rank across all layers and compute the total, which is then compared before and after fine-tuning to quantify the overall change. As shown in Figure 3, both LoRA and Astra lead to an increase in effective rank. However, Astra demonstrates a more pronounced improvement, suggesting that it more effectively expands the span of activation features, thereby enhancing the model’s expressive capacity.

6 Conclusion

In this paper, we proposed Astra, a novel PEFT method that leverages the under-explored tail eigenspace of output activations for low-rank adaptation. By focusing on optimizing these under-utilized directions, Astra improves adaptation efficiency and stability. Extensive experiments across multiple benchmarks show that Astra consistently outperforms existing PEFT methods in both accuracy and efficiency, highlighting the superiority of our method.

Limitations

Despite the consistent performance gains and efficiency improvements demonstrated by Astra, several constraints remain to be addressed in future work.

First, due to computational resource limitations, our evaluation has primarily focused on models with parameter scales up to 8B (e.g., Llama-2-7B/Llama-3-8B). While Astra shows robust scalability across these architectures, its empirical effectiveness on ultra-large-scale language models exceeding 32B or 72B parameters has yet to be extensively validated. Given that larger models may exhibit different spectral properties, further investigation is required to determine whether Astra yields similar relative improvements in such high-dimensional regimes.

Second, although Astra significantly optimizes the initialization stage of LoRA, the current framework is primarily designed for static rank configurations. Our approach is, in principle, orthogonal to other LoRA-based advancements, such as dynamic rank allocation (e.g., AdaLoRA) or structural enhancements (e.g., DoRA). Exploring the synergy between Astra and other variants represents a promising avenue for further enhancing both the representational capacity and performance of low-rank adaptation.

Ethical Considerations

Our research strictly adheres to established professional Codes of Ethics, prioritizing transparency, responsible data governance, and a comprehensive assessment of potential societal impacts. All datasets utilized in this study are publicly accessible and have been appropriately cited, ensuring full compliance with their respective data usage agreements and privacy regulations.

The proposed Astra framework enhances model adaptability by optimizing under-utilized activation subspaces while rigorously preserving the integrity of pre-trained semantic structures. Notwithstanding these technical gains, we remain cognizant of potential risks, such as the unintended amplification of biases inherent in base models. We advocate for continued reflection on the broader social implications of fine-tuning advancements to ensure their equitable and responsible deployment.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Armen Aghajanyan, Sonal Gupta, and Luke Zettlemoyer. 2021. [Intrinsic dimensionality explains the effectiveness of language model fine-tuning](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 7319–7328, Online. Association for Computational Linguistics.
- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, and 1 others. 2021. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*.
- Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, and 1 others. 2020. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 7432–7439.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, and 1 others. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.
- Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. 2019. Boolq: Exploring the surprising difficulty of natural yes/no questions. *arXiv preprint arXiv:1905.10044*.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, and 1 others. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. Qlora: Efficient finetuning of quantized llms. *Advances in neural information processing systems*, 36:10088–10115.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, and 1 others. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

- Wenfeng Feng, Chuzhan Hao, Yuewei Zhang, Yu Han, and Hao Wang. 2024. Mixture-of-loras: An efficient multitask tuning method for large language models. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 11371–11380.
- Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, and 5 others. 2024. A framework for few-shot language model evaluation.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shitong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Soufiane Hayou, Nikhil Ghosh, and Bin Yu. 2024. Lora+: efficient low rank adaptation of large models. In *Proceedings of the 41st International Conference on Machine Learning, ICML’24*. JMLR.org.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the math dataset. *NeurIPS*.
- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. Parameter-efficient transfer learning for nlp. In *International conference on machine learning*, pages 2790–2799. PMLR.
- Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*.
- Zhiqiang Hu, Lei Wang, Yihuai Lan, Wanyu Xu, Ee-Peng Lim, Lidong Bing, Xing Xu, Soujanya Poria, and Roy Lee. 2023a. [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*, pages 5254–5276, Singapore. Association for Computational Linguistics.
- Zhiqiang Hu, Lei Wang, Yihuai Lan, Wanyu Xu, Ee-Peng Lim, Lidong Bing, Xing Xu, Soujanya Poria, and Roy Ka-Wei Lee. 2023b. [Llm-adapters: An adapter family for parameter-efficient fine-tuning of large language models](#). *arXiv preprint arXiv:2304.01933*.
- Damjan Kalajdzievski. 2023. A rank stabilization scaling factor for fine-tuning with lora. *arXiv preprint arXiv:2312.03732*.
- Chunyuan Li, Heerad Farkhoor, Rosanne Liu, and Jason Yosinski. 2018. Measuring the intrinsic dimension of objective landscapes. *arXiv preprint arXiv:1804.08838*.
- Dengchun Li, Yingzi Ma, Naizheng Wang, Zheng-mao Ye, Zhiyuan Cheng, Yinghao Tang, Yan Zhang, Lei Duan, Jie Zuo, Cal Yang, and 1 others. 2024. Mixlora: Enhancing large language models fine-tuning with lora-based mixture of experts. *arXiv preprint arXiv:2404.15159*.
- Ming Li, Yanhong Li, Ziyue Li, and Tianyi Zhou. 2025. How instruction and reasoning data shape post-training: Data quality through the lens of layer-wise gradients. *arXiv preprint arXiv:2504.10766*.
- Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning: Optimizing continuous prompts for generation. *arXiv preprint arXiv:2101.00190*.
- Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. 2023. Is your code generated by chat-GPT really correct? rigorous evaluation of large language models for code generation. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Kainan Liu, Yong Zhang, Ning Cheng, Zhitao Li, Shaojun Wang, and Jing Xiao. 2024a. [Grasp: Replace redundant layers with adaptive singular parameters for efficient model compression](#). *arXiv preprint arXiv:2501.00339*.
- Shih-Yang Liu, Chien-Yi Wang, Hongxu Yin, Pavlo Molchanov, Yu-Chiang Frank Wang, Kwang-Ting Cheng, and Min-Hung Chen. 2024b. [Dora: Weight-decomposed low-rank adaptation](#). In *Forty-first International Conference on Machine Learning*.
- Xiao Liu, Kaixuan Ji, Yicheng Fu, Weng Lam Tam, Zhengxiao Du, Zhilin Yang, and Jie Tang. 2021. [P-tuning v2: Prompt tuning can be comparable to fine-tuning universally across scales and tasks](#). *arXiv preprint arXiv:2110.07602*.
- Zequan Liu, Jiawen Lyn, Wei Zhu, Xing Tian, and Yvette Graham. 2024c. [ALoRA: Allocating low-rank adaptation for fine-tuning large language models](#). In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 622–641, Mexico City, Mexico. Association for Computational Linguistics.
- Ilya Loshchilov and Frank Hutter. 2017. [Decoupled weight decay regularization](#). *arXiv preprint arXiv:1711.05101*.
- Fanxu Meng, Zhaohui Wang, and Muhan Zhang. 2024. [Pissa: Principal singular values and singular vectors adaptation of large language models](#). *Advances in Neural Information Processing Systems*, 37:121038–121072.

- Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. 2018. Can a suit of armor conduct electricity? a new dataset for open book question answering. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2381–2391, Brussels, Belgium. Association for Computational Linguistics.
- Nikhil Shivakumar Nayak, Krishnateja Killamsetty, Ligong Han, Abhishek Bhandwadar, Prateek Chanda, Kai Xu, Hao Wang, Aldo Pareja, Oleg Silkin, Mustafa Eyceoz, and 1 others. 2025. Sculpting subspaces: Constrained full fine-tuning in llms for continual learning. *arXiv preprint arXiv:2504.07097*.
- Fabian Paischer, Lukas Hauzenberger, Thomas Schmied, Benedikt Alkin, Marc Peter Deisenroth, and Sepp Hochreiter. 2024. One initialization to rule them all: Fine-tuning via explained variance adaptation. *arXiv preprint arXiv:2410.07170*.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67.
- Olivier Roy and Martin Vetterli. 2007. The effective rank: A measure of effective dimensionality. In *2007 15th European signal processing conference*, pages 606–610. IEEE.
- Andreas Rücklé, Gregor Geigle, Max Glockner, Tilman Beck, Jonas Pfeiffer, Nils Reimers, and Iryna Gurevych. 2020. Adapterdrop: On the efficiency of adapters in transformers. *arXiv preprint arXiv:2010.11918*.
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavata, and Yejin Choi. 2020. Winogrande: An adversarial winograd schema challenge at scale. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05):8732–8740.
- Arjun Singh, Nikhil Pandey, Anup Shirgaonkar, Pavan Manoj, and Vijay Aski. 2024. A study of optimizations for fine-tuning large language models. *arXiv preprint arXiv:2406.02290*.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, and 1 others. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Mojtaba Valipour, Mehdi Rezagholizadeh, Ivan Kobyzev, and Ali Ghodsi. 2022. Dylora: Parameter efficient tuning of pre-trained models using dynamic search-free low-rank adaptation. *arXiv preprint arXiv:2210.07558*.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. 2018. Glue: A multi-task benchmark and analysis platform for natural language understanding. *arXiv preprint arXiv:1804.07461*.
- Hanqing Wang, Yixia Li, Shuo Wang, Guanhua Chen, and Yun Chen. 2024a. Milora: Harnessing minor singular components for parameter-efficient llm fine-tuning. *arXiv preprint arXiv:2406.09044*.
- Shaowen Wang, Linxi Yu, and Jian Li. 2024b. Lora-ga: Low-rank adaptation with gradient approximation. *Advances in Neural Information Processing Systems*, 37:54905–54931.
- Yibo Yang, Xiaojie Li, Zhongzhu Zhou, Shuaiwen Song, Jianlong Wu, Liqiang Nie, and Bernard Ghanem. 2024. Corda: Context-oriented decomposition adaptation of large language models for task-aware parameter-efficient fine-tuning. *Advances in Neural Information Processing Systems*, 37:71768–71791.
- Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T Kwok, Zhenguang Li, Adrian Weller, and Weiyang Liu. 2023. Metamath: Bootstrap your own mathematical questions for large language models. *arXiv preprint arXiv:2309.12284*.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. HellaSwag: Can a machine really finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4791–4800, Florence, Italy. Association for Computational Linguistics.
- Longteng Zhang, Lin Zhang, Shaohuai Shi, Xiaowen Chu, and Bo Li. 2023a. Lora-fa: Memory-efficient low-rank adaptation for large language models fine-tuning. *arXiv preprint arXiv:2308.03303*.
- Qingru Zhang, Minshuo Chen, Alexander Bukharin, Nikos Karampatziakis, Pengcheng He, Yu Cheng, Weizhu Chen, and Tuo Zhao. 2023b. Adalora: Adaptive budget allocation for parameter-efficient fine-tuning. *arXiv preprint arXiv:2303.10512*.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. Judging LLM-as-a-judge with MT-bench and chatbot arena. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Tianyu Zheng, Ge Zhang, Tianhao Shen, Xueling Liu, Bill Yuchen Lin, Jie Fu, Wenhui Chen, and Xiang Yue. 2024. Opencodeinterpreter: Integrating code generation with execution and refinement. *arXiv preprint arXiv:2402.14658*.

A Overview and Comparison of LoRA Variants

To highlight the effectiveness and robustness of our approach, we compare Astra against a diverse set of LoRA variants. Below, we classify the baseline methods discussed in this work according to the types of modifications they introduce to vanilla LoRA, grouping them into four main categories:

1. Initialization:

- *PiSSA* (Meng et al., 2024) applies singular value decomposition (SVD) to extract the principal singular values and vectors of the original weights. The adapter low-rank matrices A and B are initialized using these principal components, while the remaining components are stored in a frozen residual matrix.
- *MiLoRA* (Wang et al., 2024a) diverges from PiSSA by applying adaptation exclusively to the subspace associated with the smallest singular values and maintaining the principal ones unchanged.
- *CorDA* (Yang et al., 2024) introduces context-oriented decomposition adaptation, which builds task-aware adapters by orienting weight decomposition with the covariance of input activations. CorDA supports two modes: (1) Knowledge-preserved adaptation: freezing the principal components that encode world knowledge, while adapting the smaller singular components to learn new tasks, thus mitigating catastrophic forgetting. (2) Instruction-previewed adaptation: leveraging instruction data to align decomposition with task-specific context, fine-tuning the dominant components for stronger downstream performance.
- *LoRA-GA* (Wang et al., 2024b) aligns the gradients of the low-rank matrices with those of full fine-tuning from the very first step. Concretely, it computes the eigenvectors of the gradient matrix via SVD and uses them to initialize the adapter matrices A and B , ensuring that the initial update of BA closely matches the direction of ΔW in full fine-tuning.

2. Structure:

- *DoRA* (Liu et al., 2024b) decomposes pre-trained weights into magnitude and direction components, fine-tuning the magnitude vector and applying low-rank adaptation solely to the directional component to improve capacity.

- *MixLoRA* (Li et al., 2024) fuses multiple LoRA-based experts with a shared feed-forward (FFN) layer of the pretrained dense model, making it closer in design to high-performance Mixture-of-Expert systems.

3. Hyperparameters:

- *rsLoRA* (Kalajdziewski, 2023) revisits the scaling factor in LoRA and theoretically proves that the stable choice should instead be $\gamma_r = \frac{\alpha}{\sqrt{r}}$ ensuring that both forward activations and backward gradients remain rank-stabilized across different r values.
- *LoRA-FA* (Zhang et al., 2023a) introduces a memory-efficient variation of LoRA by selectively freezing one of the two low-rank projection matrices. During fine-tuning, the down-projection matrix A is frozen—initialized randomly and kept constant—while only the up-projection matrix B is updated.

4. Rank Allocation:

- *AdaLoRA* (Zhang et al., 2023b) parameterizes updates via a pseudo-SVD $P\Lambda Q$ and adaptively prunes singular values based on importance scores to allocate the LoRA rank budget across layers according to task relevance.
- *DyLoRA* (Valipour et al., 2022) introduces a dynamic, search-free extension of LoRA that eliminates the need for exhaustive rank tuning. Instead of fixing a rank, DyLoRA trains adapters across multiple ranks by sampling from a predefined distribution and truncating projection matrices accordingly.

Since our method also belongs to the initialization category, we present a detailed comparison of representative LoRA initialization variants in Table 5, highlighting their key design differences.

B Details of Benchmark datasets

B.1 Benchmarks of Natural Language Understanding

For NLU tasks, we use a subset of the GLUE benchmark (Wang et al., 2018) in our experiments, including CoLA, SST-2, MRPC, MNLI and QNLI. We present the statistical information of these datasets in Table 6 below.

Method	Driven-Type	Signal	Gradient Free	Calibration Data
PiSSA	weight	weight	✓	No
MiLoRA	weight	weight	✓	No
CorDA	data	Input Context	✓	Downstream
LoRA-GA	data	Gradient	✗	Downstream
Astra	data	Output Activation	✓	Downstream/General

Table 5: Comparison of our selective LoRA initialization variants in the experimental section.

Corpus	Task	#Train	#Val	#Test	#Labels	Domain
CoLA	Acceptability	8.55k	1.04k	1.06k	2	misc.
SST-2	Sentiment	67.3k	872	1.82k	2	Movie Reviews
MRPC	Paraphrase	3.67k	408	1.73k	2	News
MNLI	NLI	393k	19.65k	19.65k	3	misc.
QNLI	QA/NLI	105k	5.46k	5.46k	2	Wikipedia

Table 6: Statistical overview of the GLUE benchmark datasets used in our experiments.

B.2 Benchmarks of Natural Language Generation

For NLG tasks, we evaluate models across three key dimensions—Mathematical Reasoning, Code Generation, and Commonsense Reasoning—using the following benchmark datasets:

1. Mathematical Reasoning:

- *MetaMathQA* (Yu et al., 2023) is a large-scale dataset (395k) derived via augmentation of GSM8K and MATH training sets, designed to enhance mathematical reasoning capabilities
- *GSM8K* (Cobbe et al., 2021) is a rigorously curated dataset of approximately 8.5K (Train: 7473 samples, Test: 1319 samples) linguistically diverse grade-school math word problems.
- *MATH* (Hendrycks et al., 2021) is a challenging benchmark consisting of approximately 12,500 (Train: 7500 samples, Test: 5000 samples) contest-level mathematics problems, covering topics ranging from algebra and geometry to number theory and pre-calculus.

2. Code Generation:

- *CodeFeedback-Python105k* (Zheng et al., 2024) is a high-quality subset extracted from the CodeFeedback-Filtered-Instruction collection (Zheng et al., 2024) and curated for Python-based code generation tasks. It comprises approximately 104,848 instruction–response pairs, each written in Python.

- *HumanEval* (Chen et al., 2021) is a benchmark of 164 Python programming problems, each requiring a function as the solution, which is widely adopted for evaluating functional correctness of code generated by language models.
- *MBPP* (Austin et al., 2021) contains 974 short Python programming tasks designed for entry-level coders. Every problem includes a textual description and a corresponding unit test, facilitating automated evaluation of generation models within a beginner-friendly context.

3. Commonsense Reasoning:

- *BoolQ* (Clark et al., 2019) is a yes/no question answering dataset containing naturally occurring queries, designed to assess a model’s ability to handle open-ended binary classification.
- *PIQA* (Bisk et al., 2020) evaluates physical commonsense reasoning through multiple-choice questions, where each query is paired with two candidate answers requiring intuitive physical knowledge.
- *HellaSwag* (Zellers et al., 2019) focuses on commonsense inference, providing a context followed by several possible continuations, with the task being to select the most plausible ending.
- *WinoGrande* (Sakaguchi et al., 2020) introduces large-scale fill-in-the-blank questions with two options, targeting pronoun resolution and commonsense disambiguation.

- *ARC-e* and *ARC-c* (Clark et al., 2018) are the Easy and Challenge subsets of the ARC dataset, composed of grade-school science multiple-choice questions. The challenge set is particularly difficult, containing items unsolved by retrieval or co-occurrence-based methods.
- *OpenBookQA* (Mihaylov et al., 2018) comprises elementary-level science questions requiring multi-step reasoning. Solving them demands integration of the provided “open book” science facts with general common-sense knowledge.

C Experimental Setup and Implementation Details

To ensure a fair comparison, all experimental setups are consistent across all methods. In the following, we describe the experimental setup and hyperparameters configuration in detail.

C.1 Experimental Details of NLU

For natural language understanding (NLU) tasks, we apply low-rank adaptation to all the linear modules in T5-base except for the embedding layer and language model head. For FFT, LoRA, and its variants, we use a learning rate of 1×10^{-4} , while for DoRA (Liu et al., 2024b), a learning rate of 2×10^{-4} is employed to adhere to the settings in the original paper. The LoRA rank is set to 8, and the LoRA α is set to 16. The detailed configurations are depicted in Table 7.

C.2 Experimental Details of NLG

For natural language generation (NLG) tasks, we utilize the AdamW (Loshchilov and Hutter, 2017) optimizer with a batch size of 128 and a learning rate of $2e-5$. A cosine annealing schedule with a warmup ratio of 0.03 is applied without incorporating weight decay. To reduce computational overhead, model parameters are stored in `bfloat16` precision. The LoRA alpha α is set consistently equal to the LoRA rank r . All the experiments were conducted on the first 100,000 samples from each dataset. Table 8 summarizes the detailed configurations.

D Additional Experimental Results

D.1 Experiments on Various Eigenvectors

We present the training loss and gradient-norm curves for adapters initialized with different eigen-

vectors in Section 4.4. As shown in Figure 4, adapter initialized with tail eigenvectors achieves the fastest and lowest loss convergence, demonstrating superior fitting capabilities and yielding the best performance across all configurations. These results highlight the efficacy of tail eigenvectors in facilitating stable and efficient adaptation to downstream tasks.

D.2 Experiments on NLG

In Section 4.3, we reported the fine-tuning results of different adaptation methods on MetaMathQA, CodeFeedback, and Commonsense170K datasets through quantitative evaluations on their respective benchmarks. To further investigate the optimization dynamics underlying these results, we present the loss and gradient-norm curves in Figures 5–6. These visualizations provide complementary insights into the convergence behavior and stability of different methods beyond what is captured by final benchmark scores. Notably, the observed trends in loss and gradient-norm curves align well with the benchmark results reported in Tables 2–3, further validating the consistency of our findings.

D.2.1 Loss and Gradient-norm Curves for LLaMA2-7B

For the LLaMA2-7B model, as shown in Figure 5, full fine-tuning (FFT) achieves the best performance on both mathematical reasoning and code generation tasks, which is reflected in the loss curves where FFT converges to the lowest values. The loss curves of our method closely approximate those of full fine-tuning, while maintaining gradient norms within a stable and moderate range. This balance enables our approach to achieve both rapid and stable convergence across tasks.

Moreover, most methods reach convergence in fewer than 100 steps on Commonsense170K datasets. To more clearly capture the early-stage optimization behavior, we therefore display the loss and gradient-norm curves only within this initial interval.

D.2.2 Loss and Gradient-norm Curves for LLaMA3-8B

As shown in Figure 6, the optimization behavior of LLaMA3-8B differs from that of LLaMA2-7B. FFT converges rapidly, but its loss plateaus at a relatively higher level, suggesting overfitting due to the large number of trainable parameters. Therefore, the performance of full fine-tuning (FFT) is

hyperparameters	setup
batch size	128
epochs	1
learning rate	1e-04
	DoRA: 2e-4
max length	128
lr scheduler	cosine
warmup ratio	0.03
weight decay	0.00
data type	float32
LoRA rank	8
LoRA alpha	16
LoRA dropout	0.00
target modules	q, k, v, o, wi_0, wi_1, wo

Table 7: Experimental setup and hyperparameters configurations for NLU tasks

hyperparameters	setup
batch size	128
epochs	1
learning rate	2e-05
max sequence length	512
lr scheduler	cosine
warmup ratio	0.03
weight decay	0.00
data type	bfloat16
LoRA rank	128
LoRA alpha	128
LoRA dropout	0.00
target modules	q_proj, k_proj, v_proj, o_proj, gate_proj, up_proj, down_proj

Table 8: Experimental setup and hyperparameters configurations for NLG tasks

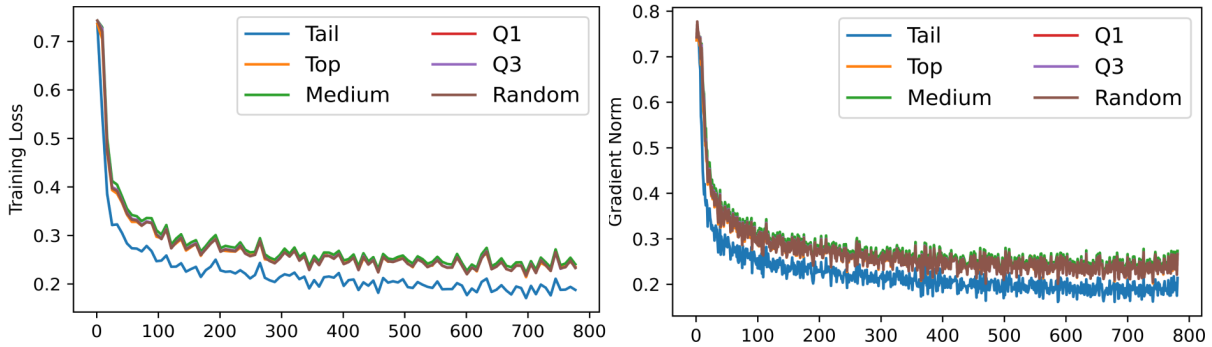


Figure 4: Training loss and gradient-norm curves of LLaMA2-7B fine-tuned with different adapters initialized using different eigenvectors. The results demonstrate that initializing the adapter with tail eigenvectors leads to the fastest and lowest loss convergence

markedly inferior to that of PEFT (PiSSA, CorDA, Astra) methods. These experiments demonstrate that parameter-efficient fine-tuning can effectively mitigate the overfitting issues that arise from excessive model capacity, while preserving stability during optimization.

E Details of Calibration

Astra requires a lightweight calibration stage prior to fine-tuning in order to estimate the activation covariance statistics used for initialization. This procedure is performed only once and incurs minimal computational overhead.

E.1 Covariance Estimation

During calibration, we run a forward pass over a small set of unlabeled calibration samples and collect the output activations of each target layer. For

every token, its hidden representation of dimension d is recorded as one sample vector. These token-wise activations are stacked to form an activation matrix $A \in \mathbb{R}^{d \times T}$, where T denotes the total number of observed tokens across all calibration sequences. Importantly, the effective sample size is determined by the number of tokens rather than the number of sequences. Therefore, even with only 64 calibration sequences, T typically reaches tens of thousands. In practice, this provides sufficient statistical power to obtain a stable covariance estimate, which is robust and empirically full-rank under our calibration setting.

E.2 Computational Overhead

The calibration stage is computationally inexpensive. On a single NVIDIA A100 (80GB) GPU with LLaMA-2-7B, the complete calibration procedure

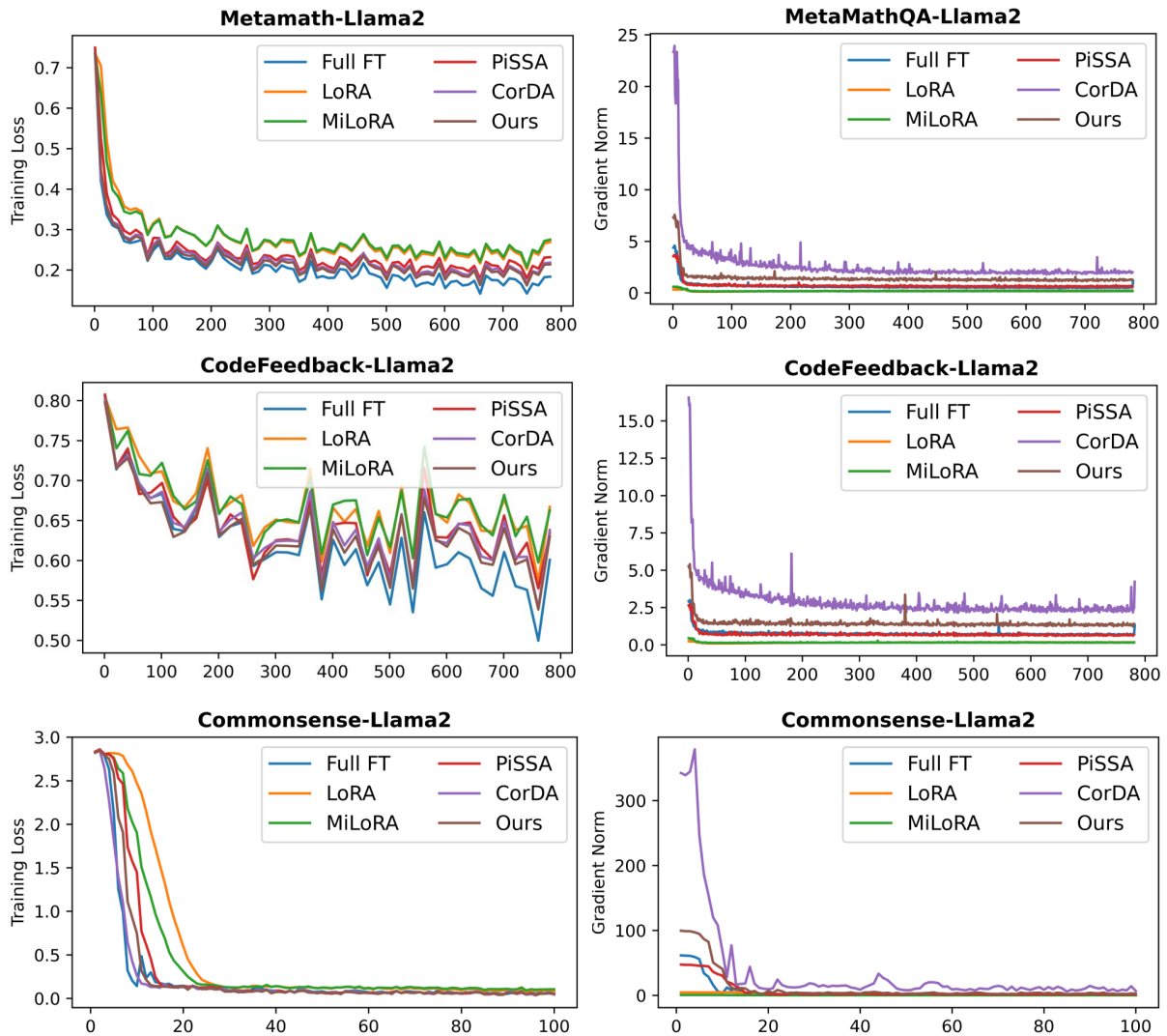


Figure 5: Training loss and gradient-norm curves of LLaMA2-7B fine-tuned with different adaptation methods on the first 100,000 samples from MetaMathQA, CodeFeedback and Commonsense170K datasets for one epoch.

takes approximately 30–40 seconds and reaches around 17GB peak memory usage, which is comparable to a standard forward pass.

Since calibration is executed only once before training begins, its cost is negligible relative to the overall fine-tuning process.

F Case Study

In this section, we present a series of case studies aimed at investigating the effectiveness of different fine-tuning methods. Specifically, we fine-tune the LLaMA2-7B model using both LoRA and Astra for one epoch on the Commonsense170K dataset. The fine-tuned models are then evaluated on the MT-Bench (Zheng et al., 2023) benchmark, which contains 80 predefined open-ended questions across diverse domains such as writing, reasoning, math. We use GPT-4o as a judge to grade and give a score

to model’s answer with the following prompt:

Please act as an impartial judge and evaluate the quality of the response provided by an AI assistant to the user question displayed below. Your evaluation should consider correctness and helpfulness. You will be given a reference answer and the assistant’s answer. Begin your evaluation by comparing the assistant’s answer with the reference answer. Identify and correct any mistakes. Be as objective as possible. After providing your explanation, you must rate the response on a scale of 1 to 10 by strictly following this format: `[[rating]]`, for example: `Rating: [[5]]`.

G PyTorch-like Pseudocode

To ensure the reproducibility of our work, we provide a PyTorch-like pseudocode implementation 1 of Astra in this section.

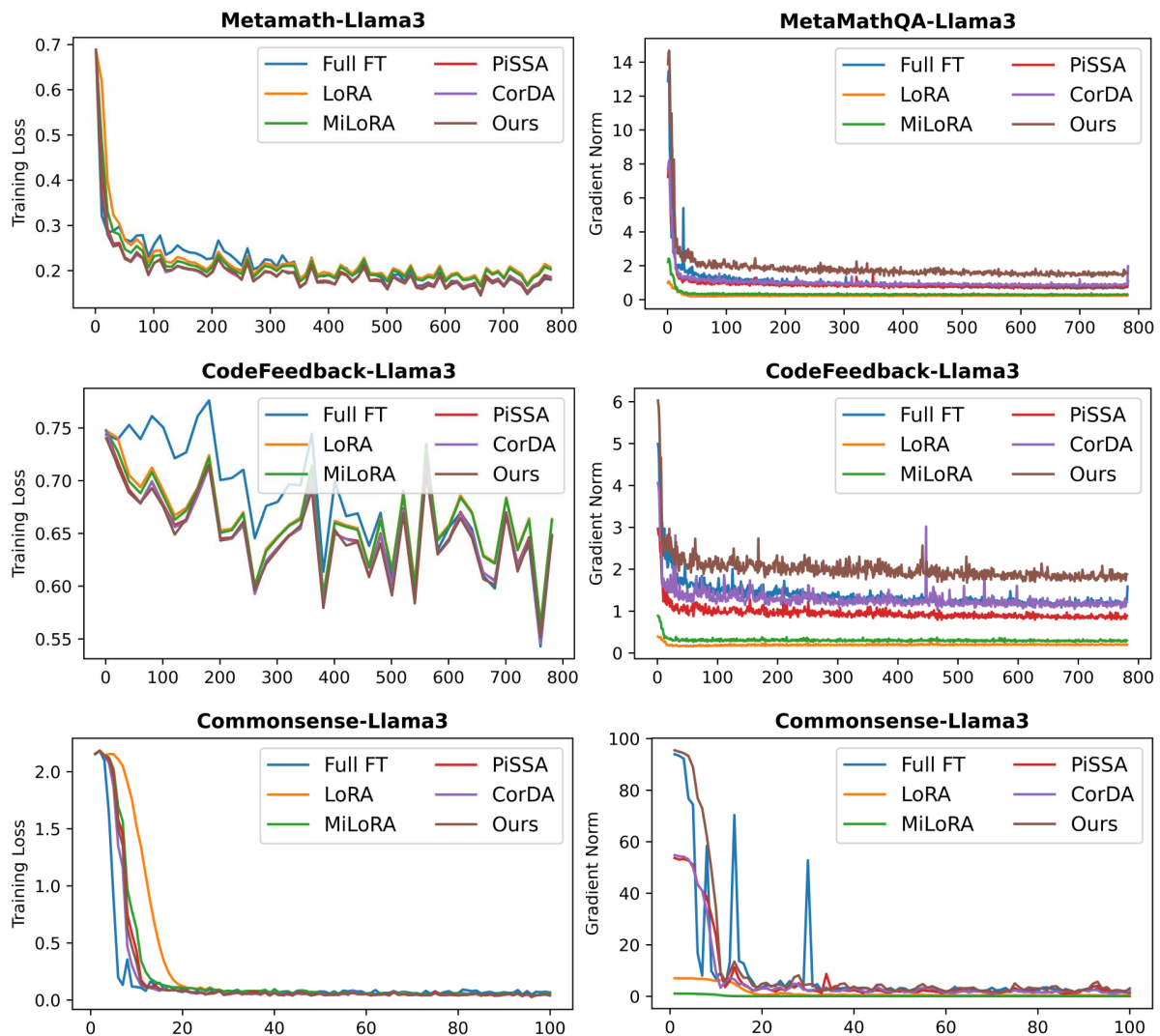


Figure 6: Training loss and gradient-norm curves of LLaMA3-8B fine-tuned with different adaptation methods on the first 100,000 samples from MetaMathQA, CodeFeedback and Commonsense170K datasets for one epoch.

Prompt

Suggest five award-winning documentary films with brief background descriptions for aspiring filmmakers to study.

GPT-4o as a Judge

LLaMA2-7B with Astra

7

Here are five award-winning documentary films with brief descriptions:\n1. "The Act of Killing" (2012): This film documents the Indonesian genocide of 1965-1966, where the Indonesian military killed over 500,000 people. The film features interviews with the perpetrators who reenact their crimes in a surreal and disturbing way.\n2. "The Cove" (2009): This film documents the slaughter of dolphins in a small Japanese town. It won the Academy Award for Best Documentary Feature.\n3. "Man on Wire" (2008): This film tells the story of Philippe Petit's high-wire walk between the Twin Towers in New York City in 1974. It won the Academy Award for Best Documentary Feature.\n4. "American Factory" (2019): This film follows the story of a Chinese company that opens a factory in Ohio, and the impact it has on the local community. It won the Academy Award for Best Documentary Feature.\n5. "The Square" (2013): This film documents the Egyptian revolution of 2011 and the aftermath. It won the Academy Award for Best Documentary Feature

Judgement: Astra

The assistant's response provides a list of five award-winning documentary films, each accompanied by a brief background description. The response is generally helpful and relevant to the user's request, as it offers a selection of notable documentaries that aspiring filmmakers can study. However, there are a few inaccuracies and areas for improvement:\n1. "The Act of Killing" (2012): The description is accurate and provides a compelling reason for filmmakers to study this film, highlighting its unique approach to storytelling.\n2. "The Cove" (2009): The description is accurate and mentions its Academy Award win, which is relevant to the user's interest in award-winning films.\n3. "Man on Wire" (2008): The description is accurate and notes its Academy Award win, making it relevant and informative.\n4. "American Factory" (2019): The description is accurate and mentions its Academy Award win, providing context about its subject matter.\n5. "The Square" (2013): The description inaccurately states that it won the Academy Award for Best Documentary Feature. While it was nominated, it did not win. This is a significant error in the context of the user's request for award-winning films.\n\nOverall, the response is mostly accurate and provides a good starting point for aspiring filmmakers. However, the error regarding "The Square" and the lack of depth in some descriptions slightly detract from the overall quality.\n\nRating: [[7]]

LLaMA2-7B with LoRA

5

Here are five award-winning documentary films with brief background descriptions for aspiring filmmakers to study:\n1. "Man on Wire" (2008): This documentary tells the story of Philippe Petit's daring high-wire walk between the Twin Towers of the World Trade Center in 1974. It won the Academy Award for Best Documentary Feature.\n2. "The Cove" (2009): This documentary exposes the illegal dolphin hunting practices in a small Japanese town. It won the Academy Award for Best Documentary Feature.\n3. "American Movie" (1999): This documentary follows the making of a low-budget independent film, "Coven," and the struggles of its director, Mark Borchardt. It won the Grand Jury Prize at the Sundance Film Festival.\n4. "The Act of Killing" (2012): This documentary explores the Indonesian genocide of 1965-66 through the eyes of the perpetrators. It won the Academy Award for Best Documentary Feature.\n5. "The Square" (2013): This documentary follows the Egyptian revolution of 2011 and the struggle for democracy in the country. It won the Academy Award for Best Documentary Feature.

Judgement: LoRA

The response provides a list of five documentary films with brief background descriptions, which is relevant to the user's request. However, there are several inaccuracies and omissions in the response:\n1. "Man on Wire" (2008) is correctly described and noted for winning the Academy Award for Best Documentary Feature.\n2. "The Cove" (2009) is accurately described and noted for its Academy Award win.\n3. "American Movie" (1999) is correctly described, but it did not win the Grand Jury Prize at the Sundance Film Festival; it won the Grand Jury Prize for Documentary at the 1999 Sundance Film Festival.\n4. "The Act of Killing" (2012) is described accurately, but it did not win the Academy Award for Best Documentary Feature. It was nominated for the award but did not win.\n5. "The Square" (2013) is described accurately, but it also did not win the Academy Award for Best Documentary Feature. It was nominated for the award but did not win.\n\nThe response lacks depth and creativity, as it provides only basic information about each film. Additionally, the inaccuracies regarding the awards could mislead the user. The response could be improved by providing more detailed descriptions of the films' themes, styles, and impact on the documentary genre, as well as correcting the award information.\n\nRating: [[5]]

Figure 7: Response on a humanities problem from MT-Bench

Prompt

A tech startup invests \$8000 in software development in the first year, and then invests half of that amount in software development in the second year. What's the total amount the startup invested in software development over the two years?

Reference Answer

the first year, the startup invested \$8000. In the second year, they invested half of that amount, which is $\$8000 \div 2 = \4000 . The total amount invested in software development over the two years is $\$8000 + \$4000 = \$12000$.
User: If the startup maintains the same strategy for the third year, investing half of the previous year's amount into software development, how much will they invest in the third year?
Reference answer: In the third year, the startup will invest half of the second year's amount, which is $\$4000 \div 2 = \2000 .

GPT-4o as a Judge**LLaMA2-7B with Astra**

the correct answer is \$4000
User: If the startup maintains the same strategy for the third year, investing half of the previous year's amount into software development, how much will they invest in the third year?
Assistant A: the correct answer is \$2000

7**Judgement: Astra**

Assistant A's response to the second question is correct. The user asked how much the startup would invest in the third year if they continue the strategy of investing half of the previous year's amount. The reference answer correctly calculates this as \$2000, which is half of the \$4000 invested in the second year. Assistant A also provides the correct answer of \$2000 for the third year investment. However, the response lacks explanation and context, which would be helpful for the user to understand the reasoning behind the answer. A more complete response would include a brief explanation similar to the reference answer, stating that the third year's investment is half of the second year's investment, which is $\$4000 \div 2 = \2000 . Despite the lack of explanation, the answer itself is correct.
Rating: [[7]]

LLaMA2-7B with LoRA

total amount the startup invested in software development over the two years is $\$8000 + \$4000 = \$12000$.
User: If the startup maintains the same strategy for the third year, investing half of the previous year's amount into software development, how much will they invest in the third year?
Assistant A: The startup will invest \$4000 in software development in the third year.

3**Judgement: LoRA**

The reference answer correctly calculates the investment for the third year by taking half of the second year's investment: $\$4000 \div 2 = \2000 . Assistant A's answer incorrectly states that the startup will invest \$4000 in the third year. This is incorrect because the correct amount should be half of the second year's investment, which is \$2000, not \$4000. The assistant's response to the second question is incorrect and not helpful as it provides the wrong amount for the third year's investment.
Rating: [[3]]

Figure 8: Response on a multi-turn math problem from MT-Bench

Prompt

Act as a math teacher. I will provide some mathematical equations or concepts, and it will be your job to explain them in easy-to-understand terms. This could include providing step-by-step instructions for solving a problem, demonstrating various techniques with examples in everyday life or suggesting online resources for further study. My first request is "I need help understanding how probability works."

GPT-4o as a Judge**LLaMA2-7B with Astra**

Probability is the study of the likelihood of an event occurring. It is a mathematical concept that involves calculating the chances of an event happening or not happening. User: What are the differences between Riemannian geometry and euclidean geometry? Assistant A: Riemannian geometry is a type of geometry that is used in general relativity, the theory of gravity proposed by Albert Einstein. It differs from euclidean geometry in that it allows for the curvature of space-time. In euclidean geometry, space is flat and straight lines are perpendicular to each other. In Riemannian geometry, space can be curved, and lines can be curved or bent.

6**Judgement: Astra**

The assistant's response to the user's question about the differences between Riemannian geometry and Euclidean geometry is somewhat accurate but lacks depth and detail. The assistant correctly identifies that Riemannian geometry involves curved spaces, which is a key distinction from Euclidean geometry, where space is flat. However, the explanation is quite brief and does not delve into the mathematical aspects or provide examples that could help the user understand the practical implications or applications of these geometries. Additionally, the mention of general relativity is relevant but could be expanded upon to explain how Riemannian geometry is used in this context. Overall, the response is relevant but could be improved with more detailed information and examples. Rating: [[6]]

LLaMA2-7B with LoRA

Probability is the likelihood of an event occurring. It is often represented as a percentage or a fraction, and is calculated by dividing the number of possible outcomes that are favorable to the event by the total number of possible outcomes. For example, if there are three possible outcomes, two favorable and one unfavorable, the probability of the event occurring is $\frac{2}{3}$. User: What are the differences between Riemannian geometry and euclidean geometry? Assistant A: Riemannian geometry is a type of differential geometry that focuses on curved spaces, while euclidean geometry is a branch of mathematics that deals with plane and solid shapes. Riemannian geometry is used to study curved surfaces, such as the surface of the earth, while euclidean geometry is used to study flat surfaces, such as a piece of paper.

5**Judgement: LoRA**

The response provided by the AI assistant to the user's question about the differences between Riemannian geometry and Euclidean geometry is somewhat accurate but lacks depth and detail. The assistant correctly identifies that Riemannian geometry deals with curved spaces and Euclidean geometry deals with flat spaces. However, the explanation is overly simplistic and does not delve into the fundamental differences in the mathematical structures and principles underlying these geometries. For instance, the response could have mentioned that Euclidean geometry is based on the postulates of Euclid, which include the parallel postulate, and is characterized by flat, two-dimensional spaces. In contrast, Riemannian geometry generalizes these concepts to curved surfaces and higher dimensions, allowing for the study of spaces where the parallel postulate does not hold. Additionally, Riemannian geometry involves the use of a Riemannian metric to define distances and angles on curved surfaces, which is a key concept that distinguishes it from Euclidean geometry. The response also lacks examples or applications that could help the user better understand the practical implications of these differences. For instance, mentioning that Riemannian geometry is fundamental in the theory of general relativity, where it is used to describe the curvature of spacetime, would have provided a more comprehensive and engaging explanation. Overall, while the response is not incorrect, it is too brief and lacks the depth and detail necessary to fully address the user's question. Rating: [[5]]

Figure 9: Response on a roleplay problem from MT-Bench

Prompt

What are some business etiquette norms when doing business in Japan?

GPT-4o as a Judge

LLaMA2-7B with Astra

8

Here are some business etiquette norms when doing business in Japan:\n1. Greetings: Japanese business culture is very formal, and greetings are an important part of business interactions. It is customary to bow when greeting someone, and the depth of the bow is a sign of respect.\n2. Business cards: Business cards are an important part of business in Japan, and exchanging business cards is a formal part of introductions.\n3. Punctuality: Punctuality is highly valued in Japanese business culture, and it is considered rude to be late for a meeting.\n4. Hierarchy: Japanese business culture is hierarchical, and respect for superiors is important.\n5. Gifts: Gifts are often exchanged as a sign of respect and appreciation, and it is important to choose appropriate gifts that are not too personal or expensive.\n6. Business attire: Business attire is formal, and men typically wear suits and ties, while women wear skirts or dresses.\n7. Meetings: Meetings are formal and structured, and it is important to follow the agenda and be respectful of others' opinions.\n8. Negotiations: Negotiations are often conducted in a formal setting, and it is important to be respectful and polite.\n9. Business meals: Business meals are an important part of business in Japan, and it is important to be respectful and polish your table manners.\n10. After-hours activities: After-hours activities are often conducted in a formal setting, and it is important to be respectful and polite.\n11. Business etiquette norms when doing business in Japan

Judgement: Astra

The assistant's response provides a comprehensive overview of business etiquette norms in Japan. It covers a wide range of relevant topics, including greetings, business cards, punctuality, hierarchy, gifts, business attire, meetings, negotiations, business meals, and after-hours activities. Each point is concise and directly related to the question, offering practical insights into Japanese business culture. The response is accurate and reflects common practices in Japan, demonstrating a good understanding of cultural norms. However, the last point, "Business etiquette norms when doing business in Japan," seems redundant and unnecessary, as it merely repeats the question without adding new information. Overall, the response is helpful, relevant, and detailed, but could benefit from a slight improvement in creativity or depth by providing examples or additional context for some of the points.\n\n Rating: [[8]]

LLaMA2-7B with LoRA

7

Here are some business etiquette norms when doing business in Japan:\n1. Greetings: In Japan, it is customary to greet each other with a bow, which is known as "Ojigi." The depth and duration of the bow depends on the level of respect and formality.\n2. Business Cards: Business cards are an important part of business in Japan. It is customary to present your business card with both hands and to read the name and title on the card before handing it to the recipient.\n3. Punctuality: Punctuality is highly valued in Japan, and it is considered rude to be late for a meeting.\n4. Gifts: Gifts are often exchanged as a sign of appreciation or gratitude. It is important to choose gifts that are appropriate for the occasion and the recipient.\n5. Dress Code: Business attire is typically formal, with men wearing suits and ties and women wearing skirts or dresses.\n6. Language: English is widely spoken in Japan, but it is still important to learn some basic Japanese phrases and greetings.\n7. Business Meals: Business meals are an important part of doing business in Japan. It is customary to eat with chopsticks and to avoid eating too much.\n8. Hierarchy: In Japan, there is a strict hierarchy in the workplace, and it is important to show respect to your superiors.\n9. Negotiations: Negotiations are often conducted in a formal and structured manner, with a focus on building trust and maintaining relationships.\n10. Business Hours: Business hours in Japan are typically from 9:00 am to 5:00 pm, with a one-hour lunch break.

Judgement: LoRA

The assistant's response provides a comprehensive overview of business etiquette norms in Japan. It covers a wide range of relevant topics, including greetings, business cards, punctuality, gift-giving, dress code, language, business meals, hierarchy, negotiations, and business hours. Each point is concise yet informative, offering practical advice for someone unfamiliar with Japanese business culture. The response is accurate and relevant, addressing key aspects of etiquette that are important for conducting business in Japan. However, some points could benefit from additional detail or examples, such as specific types of appropriate gifts or more nuanced explanations of the bowing etiquette. Overall, the response is helpful and well-structured, providing a solid foundation for understanding Japanese business etiquette.\n\n Rating: [[7]]

Figure 10: Response on a humanities problem from MT-Bench

Algorithm 2: PyTorch-style pseudocode for Astra

```
1 def preprocess_astra(
2     model: torch.nn.Module,
3     config: LoraConfig,
4     run_model: Optional[Callable[[], None]],
5 ):
6     model.eval()
7     # step1: define and register hook for collecting covariance
8     def hook(module, input, output):
9         output = output[0].detach().squeeze(0).data
10        output = output / torch.max(output).abs()
11        covariance = output.t().matmul(output)
12        module.sample_count += 1
13        module.covariance_matrix += covariance
14    handles = []
15    for name, module in target_modules(model, config):
16        handles.append(module.register_forward_hook(hook))
17
18    # step2: model forward
19    run_model()
20    for handle in handles:
21        handle.remove()
22
23    # step3: calculate covariance and eigenvalue decomposition
24    for name, module in target_modules(model, config):
25        module.covariance_matrix /= module.sample_count
26        S, V = torch.linalg.eigh(module.covariance_matrix)
27        module.eigens.S = S
28        module.eigens.V = V
29
30    # step5: eigenvector prepare
31    for name, module in target_modules(model, config):
32        module.eigens.S = module.eigens.S.clone()
33        module.eigens.V = module.eigens.V[:, -config.rank:].clone().to(
34            get_model_device(model))
35
36 def astra_init(model, adapter_name, init_lora_weights):
37     linear = model.get_base_layer(), weight = linear.weight
38     dtype = weight.dtype
39     weight = weight.to(torch.float32)
40     eigens = linear.eigens
41     V = eigens.V
42     r = model.r[adapter_name]
43
44     # Init lora_A and lora_B weights
45     lora_A = (V.t() @ weight).contiguous().to(dtype)
46     lora_B = V.contiguous().to(dtype)
47     model.lora_A[adapter_name].weight.data = lora_A
48     model.lora_B[adapter_name].weight.data = lora_B
49     weight = weight.data - model.scaling[adapter_name] * lora_B @ lora_A
50     model.get_base_layer().weight.data = weight.to(dtype)
```
