BA-LoRA: Bias-Alleviating Low-Rank Adaptation to Mitigate Catastrophic Inheritance in Large Language Models

Anonymous ACL submission

Abstract

Large language models (LLMs) have demonstrated remarkable proficiency across various natural language processing (NLP) tasks. However, adapting LLMs to downstream applications requires computationally intensive and memory-demanding fine-tuning procedures. To alleviate these burdens, parameter-efficient finetuning (PEFT) techniques have emerged as a promising approach to tailor LLMs with minimal computational overhead. While PEFT methods offer substantial advantages, they do not fully address the pervasive issue of bias propagation from pre-training data. This work introduces Bias-Alleviating Low-Rank Adaptation (BA-LoRA), a novel PEFT method designed to counteract bias inheritance. BA-LoRA incorporates three distinct regularization terms: (1) a consistency regularizer, (2) a diversity regularizer, and (3) a singular value decomposition regularizer. These regularizers aim to enhance the models' consistency, diversity, and generalization capabilities during fine-tuning. We conduct extensive experiments on natural language understanding (NLU) and natural language generation (NLG) tasks using prominent LLMs such as LLaMA, Mistral, and Gemma. The results demonstrate that BA-LoRA outperforms LoRA and its state-of-the-art variants. Moreover, the extended experiments demonstrate that our method effectively mitigates the adverse effects of pre-training bias, leading to more reliable and robust model outputs.

1 Introduction

The emergence of large language models (LLMs) has marked a new era in natural language processing (NLP). Models such as GPT-4 (OpenAI, 2023), Llama (Touvron et al., 2023), Mistral (Jiang et al., 2023), and Gemma (Team et al., 2024) have demonstrated exceptional performance across a wide array of NLP tasks, including language comprehension, generation, and reasoning (Zhao et al., 2023; Chang et al., 2024). The remarkable advancements of LLMs can be largely attributed to their training on vast datasets (Zhao et al., 2023). As LLMs continue to evolve rapidly, training on extensively scaled web-derived corpora has become standard practice to improve model generalization, thus bypassing the labor-intensive processes of data curation and annotation (Gao et al., 2020; Penedo et al., 2023). However, the corresponding increase in data volume has introduced several challenges, such as the presence of imbalanced, duplicated, and corrupted information (Parashar et al., 2024; Liu and He, 2024; Chen et al., 2024b; Yang et al., 2023).

Recent research has shown that various forms of bias in training data can negatively affect LLM behavior (Dong et al., 2023; Dodge et al., 2021; Longpre et al., 2023; Chen et al., 2024a). For example, noise within the training data can degrade model generalization (Chen et al., 2024a), while the long-tailed distribution of concepts in web-scale data can cause LLMs to overemphasize overrepresented topics (Zhu et al., 2024). Furthermore, biases introduced during pre-training can persist even after fine-tuning, potentially compromising model performance and safety in real-world applications (Qi et al., 2023; Bommasani et al., 2021; Mallen et al., 2022; Carlini et al., 2023).

This phenomenon, termed "Catastrophic Inheritance" by (Chen et al., 2024a), has spurred investigations into mitigation strategies. While constructing less biased datasets and developing more robust model architectures are prominent approaches (Liu and He, 2024), this study explores an alternative: innovations in fine-tuning. Fine-tuning LLMs is a powerful method for enhancing task-specific performance (Han et al., 2024), aligning models with user intent (Ouyang et al., 2022; Xu et al., 2024), and eliciting desired behaviors (Bai et al., 2022; Rafailov et al., 2024). However, fine-tuning largescale models' computational and memory demands are substantial (Hu et al., 2021). For instance, 16bit fine-tuning of a Llama-65B model requires over 780 GB of GPU memory (Dettmers et al., 2024). To address these limitations, parameter-efficient fine-tuning (PEFT) techniques, such as Low-Rank Adaptation (LoRA) (Hu et al., 2021), have gained prominence.

LoRA enables efficient fine-tuning of large pretrained models by approximating parameter updates using low-rank matrices. Instead of directly updating the entire weight matrix W, LoRA introduces a learnable low-rank adapter $\Delta W = AB$, where $A \in \mathbb{R}^{m \times r}$ and $B \in \mathbb{R}^{r \times n}$ are low-rank matrices with a user-defined rank $r \ll \min(m, n)$. This ensures that only A and B are updated during training, while the original weights W remain frozen. By initializing A with scaled random values and B to zero, LoRA preserves the pretrained model's initial output at the start of training. For an input X, the output is computed as Y = X(W + AB), ensuring consistency with the original model. This method significantly reduces computational costs by limiting the number of trainable parameters, making it a practical solution for adapting large models (Hu et al., 2021).

To mitigate the detrimental effects of Catastrophic Inheritance, particularly noise, and imbalance, we propose Bias-Alleviating Low-Rank Adaptation (BA-LoRA). Building upon Principal Singular Values and Singular Vectors Adaptation (PiSSA) (Meng et al., 2024), which addresses convergence issues in standard LoRA, our approach incorporates three distinct regularization terms: a consistency regularizer, a diversity regularizer, and a singular value decomposition (SVD) regularizer. The consistency regularizer preserves valuable pretrained knowledge during fine-tuning, while the diversity regularizer encourages varied model outputs. The SVD regularizer enhances the generalization capabilities of generative models. Recognizing the fundamental differences between Natural Language Understanding (NLU) and Natural Language Generation (NLG), such as determinism in NLU versus diversity in NLG, we tailor our regularization strategies accordingly.

To evaluate the efficacy of BA-LoRA, we conduct comprehensive experiments across diverse benchmarks, including mathematical reasoning (GSM8K (Cobbe et al., 2021) and MATH (Yu et al., 2023)), coding (HumanEval (Chen et al., 2021) and MBPP (Austin et al., 2021)), natural language understanding (GLUE (Wang et al., 2018)), and general language evaluation (MT-Bench (Zheng et al., 2024)). Our experiments utilize prominent LLMs such as LLaMA 2-7B (Touvron et al., 2023), Mistral-7B (Jiang et al., 2023), and Gemma-7B (Team et al., 2024), as well as encoder-only architectures like RoBERTa-large (Liu et al., 2019) and DeBERTa-v3-base (He et al., 2021b). The results unequivocally demonstrate BA-LoRA's superiority over LoRA and its state-of-the-art variants. Moreover, the extended experiments demonstrate that our method effectively attenuates noise inherited from pre-training, leading to more robust and generalizable models.

2 Related Works

Parameter-efficient fine-tuning (PEFT) techniques (Xu et al., 2023b; Han et al., 2024) have gained attention for adapting LLMs to specific tasks with limited hardware resources. They are divided into three main categories. The first category includes adapter-based methods (Houlsby et al., 2019a; Lin et al., 2020; Lei et al., 2023; He et al., 2021a), which add and fine-tune additional layers with fewer parameters to reduce computational costs. The second category is soft prompt tuning (Hambardzumyan et al., 2021; Lester et al., 2021; Li and Liang, 2021a; Liu et al., 2023), which uses learnable soft prompts at the input to adapt the model to tasks. The third category involves low-rank adaptation (LoRA) and its variants (Hu et al., 2021; Zhang et al., 2023a; Dettmers et al., 2024), which incorporate low-rank matrices to approximate weight updates during fine-tuning (Hu et al., 2021).

Variants of LoRA enhance its efficiency and performance in different ways. AdaLoRA adaptively distributes the parameter budget among weight matrices based on their importance, improving efficiency and performance by pruning unimportant updates and minimizing computational overhead (Zhang et al., 2023b). DoRA increases LoRA's learning capacity and stability by decomposing pre-trained weights into magnitude and direction components for fine-tuning (Liu et al., 2024). LoHA enhances LoRA by employing Hamiltonian products (Hyeon-Woo et al., 2021). DyLoRA addresses the fixed size and rank optimization limitations of LoRA by dynamically training LoRA blocks across varying ranks (Valipour et al., 2022). DeltaLoRA improves the representational capacity of LoRA by updating the model's original weights using parameters from adapter layers (Zi et al., 2023). PiSSA initializes adapter matrices A and B to approximate the original matrix W through singular value decomposition, leading to faster convergence and improved performance (Meng et al., 2024). While many LoRA variants focus on accelerating convergence or reducing memory consumption, our BA-LoRA method uniquely addresses the core challenge of Catastrophic Inheritance in LLM fine-tuning.

3 Method

3.1 Principal Singular Values and Singular Vectors Adaptation (PiSSA)

As a variant of LoRA, PiSSA addresses the convergence speed challenge by retaining the core LoRA architecture while innovating in initialization. Specifically, PiSSA leverages the principal components of the original weight matrix, W, to initialize the adapter matrices, A and B. The remaining components are encapsulated within a residual matrix, $W^{res} \in \mathbb{R}^{m \times n}$. The SVD of $W \in \mathbb{R}^{m \times n}$ is expressed as $W = USV^T$, where $U \in \mathbb{R}^{m \times \min(m, \hat{n})}$ and $V \in \mathbb{R}^{n \times \min(m, n)}$ are orthogonal singular vectors, and $S = diag(s) \in$ $\mathbb{R}^{\min(m,n) \times \min(m,n)}$ is a diagonal matrix, where the operation diag(s) transforms s to S and $s \in$ $\mathbb{R}_{\leq 0}^{\min(m,n)}$ represents the singular values arranged in descending order. PiSSA partitions the singular values and vectors into principal and residual components, denoted as $\{U_{[:,:r]},S_{[:r,:r]},V_{[:,:r]}\}$ and $\{U_{[:,r:]},S_{[r:,r:]},V_{[:,r:]}\}$, respectively, where the matrix slicing notations are the same as those in Py-Torch, [: r] denotes the first r dimensions, and r signifies the intrinsic rank of W. The principal components are then employed to initialize the lowrank adapter with $A \in \mathbb{R}^{m \times r}$ and $B \in \mathbb{R}^{r \times n}$:

$$A = U_{[:,:r]} S_{[:r:r]}^{1/2} \in \mathbb{R}^{m \times r}$$
(1)

$$B = S_{[:r,:r]}^{1/2} V_{[:,:r]}^T \in \mathbb{R}^{r \times n}$$
(2)

The residual matrix W^{res} remains frozen during fine-tuning:

$$W^{res} = U_{[:,r:]} S_{[r:,r:]} V_{[:,r:]}^T \in \mathbb{R}^{m \times n}$$
(3)

PiSSA preserves the pre-trained model's full capacity at the start of fine-tuning by using $W = W^{res} + AB$. This approach prioritizes training the most influential parameters, thereby accelerating convergence. Inheriting LoRA's benefits of reduced parameter count and deployment simplicity, PiSSA further leverages efficient SVD computations to expedite the training process.

3.2 Bias-Alleviating Low-Rank Adaptation (BA-LoRA)

Catastrophic Inheritance encapsulates the challenges posed by biased large-scale training data, which can manifest in LLMs as vulnerabilities and limitations arising from duplicated, noisy, imbalanced, or unethical samples. These inherited flaws can adversely impact downstream tasks, leading to diminished generalization, degraded performance, security breaches, and biased outputs. To address the specific issues caused by noisy and imbalanced data, we introduce BA-LoRA, a method incorporating three distinct regularization terms: (1) consistency regularizer, (2) diversity regularizer, and (3) SVD regularizer. Recognizing the nuanced differences between NLU and NLG, we have tailored specific variants of each regularizer to optimize performance for respective task domains.

3.2.1 Regularizations for NLU Tasks

Consistency Regularization. To safeguard valuable pre-trained knowledge during the fine-tuning process, we introduce a regularization term based on the mean squared error (MSE) loss between normalized output logits produced by the pre-trained model, \mathbf{F}_P , and those generated by the fine-tuned model, \mathbf{F}_F . This loss function incentivizes the fine-tuned model to retain essential pre-trained information while adapting to downstream task requirements.

$$\mathcal{L}_{\text{CR_NLU}} = \left\| \frac{\mathbf{F}_p}{\|\mathbf{F}_p\|_2} - \frac{\mathbf{F}_f}{\|\mathbf{F}_f\|_2} \right\|_2^2 \qquad (4)$$

This objective facilitates the inheritance of critical pre-trained knowledge in \mathbf{F}_{f} after fine-tuning.

Diversity Regularization. To address the detrimental effects of imbalanced data, we introduce a diversity regularizer aimed at eliciting more diverse representational structures within LLMs and preventing the encoding of semantically similar samples during fine-tuning. Inspired by (Bardes et al., 2021), we employ a covariance loss to minimize the off-diagonal elements of the covariance matrix of the fine-tuned outputs \mathbf{F}_f :

$$\mathcal{L}_{\text{DR_NLU}} = \frac{1}{D} \sum_{i \neq j} [C(\mathbf{F}_f)]_{i,j}^2 \qquad (5)$$

316 317

318

28 29 30

33

35

277 278 where D represents the dimensionality of \mathbf{F}_f and $C(\mathbf{F}_f)$ is the covariance matrix of \mathbf{F}_f , which is defined as:

$$C(\mathbf{F}_{f}) = \frac{1}{M-1} \sum_{i=1}^{M} (f_{i} - \bar{f}) (f_{i} - \bar{f})^{T} \quad (6)$$

where M denotes the number of elements involved in \mathbf{F}_f , f_i is the *i*-th element in \mathbf{F}_f , and \bar{f} is the mean value of \mathbf{F}_f .

Singular Value Decomposition Regularization. The SVD regularizer is designed to enhance model generalizability to mitigate the adverse effects of noisy data. Building upon the insight from (Chen et al., 2019) that eigenvectors corresponding to the largest singular values significantly contribute to model generalizability, we propose an SVD regularizer that maximizes the sum of the top k singular values of a batched fine-tuned output matrix:

$$\mathcal{L}_{\text{SVDR_NLU}} = -\frac{\sum_{i=1}^{k} \sigma_i}{\sum_{j=1}^{D} \sigma_j}$$
(7)

where k is a hyperparameter, σ_i denotes the *i*-th singular value of the top k singular values of the output matrix, and $\sum_{j=1}^{D} \sigma_j$ is the sum of all singular values obtained from the SVD of the output matrix. This decomposition represents the matrix as $\mathbf{U}\Sigma\mathbf{V}^{\top}$, where Σ is a diagonal matrix containing singular values $\{\sigma_1, \ldots, \sigma_D\}$. This regularization term emphasizes significant components of the logit matrix, enhancing the model's generalizability across various downstream tasks.

3.2.2 Overall Objective Function for NLU

The overall objective function for NLU tasks is formulated as follows:

$$\mathcal{L}_{\text{NLU}} = \mathcal{L}_{\text{task}_{\text{NLU}}} + \lambda_1 \mathcal{L}_{\text{CR}_{\text{NLU}}} + \lambda_2 \mathcal{L}_{\text{DR}_{\text{NLU}}} + \lambda_3 \mathcal{L}_{\text{SVDR}_{\text{NLU}}}$$
(8)

where \mathcal{L}_{task_NLU} represents the standard crossentropy loss function for the downstream task, and λ_1 , λ_2 , and λ_3 are weighting parameters to balance each regularization term.

3.2.3 Regularizations for NLG Tasks

Consistency Regularization. To ensure that the fine-tuned model retains knowledge from pretraining, we utilize the Kullback-Leibler Divergence (KLD) to measure the divergence between the output distributions of the fine-tuned and pretrained models (Dong et al., 2021). Specifically, we define the consistency regularization loss as:

$$\mathcal{L}_{\text{CR_NLG}} = \frac{1}{T} \sum_{t=1}^{T} \text{KL}(\mathcal{P}_{\text{pt}}(y_t \mid y_{< t}, x) \parallel \mathcal{P}_{\text{ft}}(y_t \mid y_{< t}, x))$$
(9)

where $\mathcal{P}_{pt}(y_t \mid y_{< t}, x)$ and $\mathcal{P}_{ft}(y_t \mid y_{< t}, x)$ represent the conditional probability distributions of the pre-trained model and the fine-tuned model, respectively. For the current token y_t , given the input x and the preceding token sequence $y_{< t}$. KLD encourages the model to continuously retain useful pre-training information during the fine-tuning process, which is crucial for maintaining the style and coherence of the generation task.

Diversity Regularization. To enhance the diversity of the generated text, we introduce an entropybased regularization term, inspired by previous work (Gat et al., 2020). This regularization term aims to increase the entropy of the predicted token distributions during fine-tuning, thus encouraging more varied and diverse outputs.

$$\mathcal{L}_{\text{DR_NLG}} = -\frac{1}{T} \sum_{t=1}^{T} \sum_{i=1}^{N} P_{\text{ft}}(x_i|h_t) \log P_{\text{ft}}(x_i|h_t)$$
(10)

where, $P_{\text{ft}}(x_i|h_t)$ represents the probability assigned to token x_i at time step t, given the model's hidden state h_t . Maximizing entropy at each time step minimizes repetitive outputs and encourages more diverse and enriched text generation.

Singular Value Decomposition Regularization To enhance the generalization of generative models while ensuring computational efficiency, we introduce a regularization technique based on randomized SVD (Pasand and Bashivan). This method efficiently approximates the dominant singular values of a matrix, significantly reducing the computational and memory costs associated with traditional SVD while preserving key information about the data.

$$\mathcal{L}_{\text{SVDR_NLG}} = -\frac{\sum_{i=1}^{k} \tilde{\sigma}_i}{\sum_{j=1}^{D} \tilde{\sigma}_j}$$
(11)

Here, $\tilde{\sigma}_i$ represents the *i*-th largest approximated singular value obtained using randomized SVD, and *D* is the total number of singular values in

the approximation. This regularization maximizes the relative contribution of the top k singular values, encouraging the model to prioritize the most informative patterns in the data.

3.2.4 Overall Objective Function for NLG

The objective function for downstream NLG tasks is formulated as follows:

$$\mathcal{L}_{\text{NLG}} = \mathcal{L}_{\text{task}_{\text{NLG}}} + \lambda_1 \mathcal{L}_{\text{CR}_{\text{NLG}}} + \lambda_2 \mathcal{L}_{\text{DR}_{\text{NLG}}} + \lambda_3 \mathcal{L}_{\text{SVDR}_{\text{NLG}}}$$
(12)

where \mathcal{L}_{task_NLG} denotes the standard loss for the downstream generative task, and λ_1 , λ_2 , and λ_3 are weighting parameters to balance each regularization term.

4 **Experiments**

This section presents a comprehensive evaluation of our proposed BA-LoRA method across a diverse range of NLG and NLU benchmarks. Our results unequivocally demonstrate the superiority of BA-LoRA over existing LoRA variants. Furthermore, through rigorous experimentation, we elucidate BA-LoRA's efficacy in mitigating the adverse impacts of noisy data, thereby enhancing model robustness and generalizability.

4.1 Implementation Details

In our experiments, we adopt the PiSSA (Meng et al., 2024) implementation strategy. We compute the loss using only the responses from the instruction-following dataset, ensuring lora_dropout to 0. We utilize the Float32 computation type for both the base model and the adapter in BA-LoRA. For the NLU tasks, we set the hyperparameters as: k = 5, $\lambda_1 = 1e - 1$, $\lambda_2 = 4e - 2$, and $\lambda_3 = 1e - 2$. We set lora_r = lora_alpha = 128 and use AdamW (Loshchilov and Hutter, 2017) optimizer with a batch size of 128, a learning rate of 2e-5, cosine annealing schedules, and a warmup ratio of 0.03, without any weight decay. For the NLG tasks, the hyperparameters are set as: k = 10, $\lambda_1 = 1e - 4, \lambda_2 = 3e - 1$, and $\lambda_3 = 1e - 1$. We set lora_r as 8 and select lora_alpha in 8, 16. We utilize AdamW with a linear learning rate schedule to optimize and tune the learning rate (LR) from 1e - 4, 2e - 4, 3e - 4, 4e - 4, 5e - 4, 6e - 4,5e-5, 3e-5. Batch sizes (BS) are selected from 6, 8, 16, 32. Appendix Section B presents the detailed hyperparameters we utilized on the GLUE

benchmark. comparison. All experiments were conducted using NVIDIA A40 (48G) GPUs, with results derived from three independent experiments conducted under identical conditions, ensuring consistency and robustness by averaging the results to mitigate variability.

4.2 Results and Analysis

4.2.1 Analysis of the NLG and NLU Performance of BA-LoRA

To evaluate BA-LoRA's effectiveness on NLG tasks, we fine-tuned LLaMA-2-7B, Mistral-7B, and Gemma-7B on the MetaMathQA (Yu et al., 2023) and assessed their mathematical problem-solving capabilities using the GSM8K (Cobbe et al., 2021) and MATH (Yu et al., 2023) validation sets, reporting Accuracy. Similarly, models were fine-tuned on the CodeFeedback (Zheng et al., 2024) and evaluated for coding proficiency via HumanEval (Chen et al., 2021) and MBPP (Austin et al., 2021), with PASS@1 metrics reported. To assess conversational abilities, models were trained on the WizardLM-Evol-Instruct (Xu et al., 2024) and evaluated on MT-Bench (Zheng et al., 2024), with response quality judged by GPT-4 and first turn scores reported. All experiments utilized 100K data points and a single training epoch for efficiency.

Table 1 presents the experimental outcomes, clearly demonstrating BA-LoRA's superior performance compared to baseline methods. For instance, BA-LoRA enhanced LLaMA 2-7B, Mistral-7B, and Gemma-7B performance on GSM8K by 1.82%, 1.14%, and 0.55%, respectively, compared to PiSSA. HumanEval improvements were 2.03%, 1.11%, and 1.26%, while MT-Bench enhancements reached 0.24%, 0.18%, and 0.07%. Notably, BA-LoRA achieved a remarkable 6.92% performance uplift over full parameter fine-tuning on Gemma, utilizing only 2.3% of trainable parameters across five tasks.

To assess the effectiveness of BA-LoRA on natural language understanding (NLU) tasks, we conducted experiments on the GLUE benchmark (Wang et al., 2018), which includes two single-sentence classification tasks (CoLA, SST), five paired-text classification tasks (MNLI, RTE, QQP, MRPC, QNLI), and one text similarity prediction task (STS-B). The evaluation metrics comprise the overall matched and mismatched accuracy for MNLI, the Matthews correlation coefficient for CoLA, the Pearson correlation coefficient

Models	Methods	#Params	GSM8K	MATH	HumanEval	MBPP	MT-Bench	Avg
	Full FT	6738M	$49.13_{\pm 0.21}$	$7.29_{\pm 0.22}$	$21.20_{\pm 0.30}$	$35.59_{\pm 0.25}$	$4.91_{\pm 0.01}$	23.62
	LoRA	320M	$42.47_{\pm 0.29}$	$5.60{\scriptstyle \pm 0.35}$	$17.03_{\pm 0.61}$	$\overline{32.77_{\pm 0.46}}$	$\overline{4.62_{\pm 0.11}}$	20.50
LLaMA-2-7B	PiSSA	320M	$52.37_{\pm 0.52}$	$7.76_{\pm 0.19}$	$21.55_{\pm 0.44}$	$33.09_{\pm 0.57}$	$4.87_{\pm 0.06}$	<u>23.93</u>
	BA-LoRA	320M	$\overline{54.53_{\pm 0.41}}$	9.21 $_{\pm 17}$	$\overline{\textbf{23.58}_{\pm 0.25}}$	$\textbf{36.86}_{\pm 0.31}$	$5.11_{\pm0.05}$	25.86
	Full FT	6738M	$69.91_{\pm 0.25}$	$18.64_{\pm 0.35}$	$45.31_{\pm 0.14}$	$51.46_{\pm 0.13}$	$4.95_{\pm 0.05}$	38.05
	LoRA	168M	$67.68_{\pm 0.55}$	$19.90_{\pm 0.25}$	$42.54_{\pm 0.44}$	$56.85_{\pm 0.23}$	$4.92_{\pm 0.07}$	38.38
Mistral-7B	PiSSA	168M	$72.25_{\pm 0.64}$	$21.95_{\pm0.37}$	$45.37_{\pm 0.25}$	$61.57_{\pm 0.44}$	$5.23_{\pm0.05}$	<u>41.27</u>
	BA-LoRA	168M	$73.17_{\pm 0.34}$	$22.79_{\pm 0.56}$	$\overline{\textbf{46.31}_{\pm 0.17}}$	$62.77_{\pm 0.33}$	$5.41_{\pm 0.06}$	42.09
	Full FT	6738M	$72.09_{\pm 0.32}$	$22.71_{\pm 0.34}$	$47.02_{\pm 0.27}$	$55.67_{\pm 0.05}$	$5.40_{\pm 0.12}$	40.58
~ ~~	LoRA	200M	$74.64_{\pm 0.58}$	$31.16_{\pm0.33}$	$51.64_{\pm 0.28}$	$63.52_{\pm 0.65}$	$5.01_{\pm0.03}$	45.19
Gemma-7B	PiSSA	200M	$77.58_{\pm 0.41}$	$31.47_{\pm 44}$	$53.22_{\pm 35}$	$65.49_{\pm 0.18}$	$5.66_{\pm 0.05}$	46.68
	BA-LoRA	200M	$78.13_{\pm 0.25}$	$\overline{\textbf{32.25}_{\pm 25}}$	$54.44_{\pm 0.15}$	$\overline{\textbf{66.25}_{\pm 0.33}}$	$\overline{5.73_{\pm 0.07}}$	47.36
	78	1	LaMA 2 78		Mictral 78 v0 1		Mictral 7R	<i>v</i> 0 1
	·/B	LoRA		1 _ [[Mistral-7B-V0.1	22		/0.1
§ 50 BA-LORA			DRA	َ 70			BA-LORA	
20 40	1 1	conracy		curad		20 gird		
		Р Н		ĕ 65 ₩	Lor	– Ÿ ∧ ±18		
^{NS} 30		4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4		[№] 59 60	PiSs BA- Full	LORA ET		
1 2 4 8 16	32 64 128	1 2 4	8 16 32 64 12	8 1	2 4 8 16 32 6	16 4 128	1 2 4 8 16	32 64 128
Rank			Rank		Rank		Rank	

Table 1: Performance Comparison of Various Models and Methods on NLG Tasks. The best and second-best results are highlighted in **bold** and <u>underline</u>.

Figure 1: Performance Comparison of Full Fine-Tuning, LoRA, PiSSA, and BA-LoRA Across Different Ranks.

for STS-B, and accuracy for the remaining tasks. We used the DeBERTa-v3-base model (He et al., 2021b) and compared BA-LoRA against ten baseline methods, including Full Fine-Tuning (Full FT), BitFit (Zaken et al., 2021), HAdapter (Houlsby et al., 2019b), PAdapter (Pfeiffer et al., 2020), LoRA (Hu et al., 2021), LoHA (Hyeon-Woo et al., 2021), DoRA (Liu et al., 2024), AdaLoRA (Zhang et al., 2023a), and PiSSA (Meng et al., 2024).

Table 2 presents the results of DeBERTa-v3-base across eight tasks, demonstrating the consistent superiority of BA-LoRA over all baselines. On average, BA-LoRA outperforms PiSSA and LoRA by 0.44% and 1.35%, respectively. These results underscore the effectiveness of BA-LoRA in enhancing the performance of NLU models.

A comparative analysis of Tables 1 and 2 reveals BA-LoRA's consistent performance advantages across both NLG and NLU tasks. This indicates BA-LoRA's proficiency in augmenting both generative and comprehension capabilities for language models. By incorporating consistency, diversity, and SVD regularization, BA-LoRA effectively mitigates the adverse effects of Catastrophic Inheritance, fostering consistent, diverse, and generalized model outputs. Furthermore, BA-LoRA's modest computational requirements render it suitable for efficient fine-tuning of LLMs with limited resources.

4.3 Evaluating the Performance of Different Ranks

We compare the performance of BA-LoRA, LoRA, and PiSSA at different ranks using LLaMA-2-7B and Mistral-7B-v0.1 models. Each method is finetuned for one epoch on the MetaMathQA-100K dataset with ranks ranging from 1 to 128 and evaluated on the GSM8K and MATH. As Figure 1 shows, BA-LoRA consistently outperforms LoRA and PiSSA across all rank settings and datasets. As the rank increases, the performance of BA-LoRA and PiSSA surpasses full parameter fine-tuning. However, BA-LoRA performs better, especially on Mistral-7B-v0.1.

4.3.1 Analysis on mitigate noisy data

This study aims to evaluate BA-LoRA's efficacy in mitigating the detrimental effects of noise inherent in large-scale pre-training data on downstream tasks. Given the ubiquitous presence of noise in human-annotated datasets, its influence on

Methods	#Params	MNLI	SST-2	MRPC	CoLA	QNLI	QQP	RTE	SST-B	Avg
Full FT	184M	90.34 _{±0.18}	$\textbf{96.33}_{\pm 0.11}$	$89.95_{\pm 1.07}$	$71.43_{\pm 0.72}$	$94.24_{\pm 0.10}$	$92.11_{\pm 0.28}$	$83.75_{\pm 1.81}$	$91.04_{\pm0.48}$	88.86
BitFit	0.1M	$89.54_{\pm 0.29}$	$94.68_{\pm 0.11}$	$87.95_{\pm 1.33}$	$67.31_{\pm 0.49}$	$92.45_{\pm 0.17}$	88.72 ± 0.45	$79.12_{\pm 0.39}$	$91.63_{\pm 0.37}$	86.43
HAdapter	1.22M	$90.23_{\pm 0.07}$	$95.38_{\pm 0.06}$	$89.97_{\pm 0.27}$	$68.73_{\pm 0.27}$	$94.31_{\pm 0.29}$	$91.99_{\pm 0.28}$	$84.76_{\pm 0.39}$	$91.58_{\pm 0.13}$	88.37
PAdapter	1.18M	$90.42_{\pm 0.36}$	$95.49_{\pm 0.10}$	$89.71_{\pm 0.35}$	$69.04_{\pm 0.10}$	$94.38_{\pm 0.26}$	$92.15_{\pm 0.43}$	$85.53_{\pm 0.18}$	$91.69_{\pm 0.13}$	88.55
LoRA	1.33M	$90.71_{\pm 0.16}$	$94.79_{\pm 0.16}$	$89.85_{\pm 0.21}$	$70.05_{\pm 0.34}$	$93.94_{\pm 0.09}$	$92.07_{\pm 0.48}$	$85.43_{\pm 0.09}$	$91.67_{\pm 0.29}$	88.56
LoHA	1.33M	$90.74_{\pm 0.32}$	$94.92_{\pm 0.47}$	$90.43_{\pm 0.48}$	$70.63_{\pm 0.10}$	$93.95_{\pm 0.28}$	$92.05_{\pm 0.09}$	$86.41_{\pm 0.10}$	$91.72_{\pm 0.28}$	88.86
DoRA	1.27M	$90.48_{\pm 0.10}$	$95.85_{\pm 0.08}$	$91.04_{\pm 0.15}$	$71.03_{\pm 0.18}$	$94.21_{\pm 0.37}$	$92.34_{\pm 0.16}$	$86.19_{\pm 0.25}$	$91.92{\scriptstyle\pm0.38}$	89.13
AdaLoRA	1.27M	$90.87_{\pm 0.08}$	$96.18_{\pm 0.43}$	$90.81_{\pm 0.40}$	$71.64_{\pm 0.12}$	$94.68_{\pm 0.46}$	$92.37_{\pm 0.35}$	$87.78_{\pm 0.36}$	$91.97_{\pm 0.43}$	89.53
PiSSA	1.33M	$\overline{90.47_{\pm 0.44}}$	$95.81_{\pm 0.45}$	$91.48_{\pm 0.49}$	$72.27_{\pm 0.29}$	$\overline{94.41_{\pm 0.41}}$	$\overline{92.21_{\pm 0.26}}$	$\overline{87.14_{\pm 0.08}}$	$91.93_{\pm0.25}$	<u>89.47</u>
BA-LoRA	1.33M	90.92 _{±0.38}	$96.25_{\pm 0.09}$	91.83 ±0.25	$72.79_{\pm 0.42}$	94.84 $_{\pm 0.26}$	92.59 ±0.18	$87.87_{\pm 0.31}$	$92.15_{\pm 0.08}$	89.91

Table 2: Performance comparison of different baseline methods on NLU tasks. The best and second-best results are highlighted in **bold** and <u>underline</u>.

pre-training is unavoidable. To comprehensively assess the impact of noisy pre-training data, we employ both ID and OOD evaluation using the GLUE and GLUE-x benchmarks, respectively. BERT-L (Devlin et al., 2018), pre-trained on BooksCorpus (Zhu et al., 2015) and English Wikipedia, and GPT-2-XL (Radford et al., 2019), pre-trained on the noisy WebText dataset derived from Common Crawl, serve as our models.

As detailed in Tables 3 and 4, BA-LoRA consistently outperforms LoRA across all tasks, underscoring its superior generalization capabilities. Specifically, BA-LoRA achieves average performance improvements of 2.03% and 2.47% for BERT-L and GPT-2-XL, respectively, on the GLUE benchmark. Similarly, on GLUE-x, BA-LoRA surpasses LoRA by 1.61% and 1.90% for BERT-L and GPT-2-XL, respectively. These results substantiate the effectiveness of our proposed regularization terms in mitigating the negative impacts of noise in pre-training and enhancing model robustness.

4.3.2 Analysis on Mitigating Imbalanced Data

This experiment evaluates the effectiveness of BA-LoRA in addressing imbalanced data. Specifically, using the MNLI task of the GLUE benchmark, LoRA, and BA-LoRA are applied to fine-tune the BERT-L and GPT-2-XL models, respectively. The hidden layer features of the last training step are extracted and visualized using t-SNE (Van der Maaten and Hinton, 2008) technology for comparison.

As shown in Figure 2, the models fine-tuned with standard LoRA in sub-figures (a) and (c) have low discrimination between categories and obvious category mixing. In contrast, the models finetuned with BA-LoRA in sub-figures (b) and (d) have clearer category separation, especially the results of BERT-L, which have higher intra-category clustering and clearer boundaries. These analyses



Figure 2: t-SNE Visualizations Comparing Last Hidden Layer Features of BERT-L and GPT-2-XL Fine-Tuned with LoRA and BA-LoRA on the MNLI task of GLUE.

show that BA-LoRA can effectively alleviate the impact of imbalanced data in pre-training.

4.4 Ablation Study

This ablation study evaluates the impact of regularization terms in BA-LoRA on model performance for both NLG and NLU tasks, with results visualized in Figure 3. For NLG tasks, three models (LLaMA-2-7B, Mistral-7B, and Gemma-7B) were evaluated on the GSM8K and MATH using the regularization terms \mathcal{L}_{CR_NLG} , \mathcal{L}_{DR_NLG} , and \mathcal{L}_{SVDR_NLG} . The absence of regularization (w/o Reg) consistently resulted in the lowest performance. Introducing individual terms led to varying improvements: \mathcal{L}_{CR_NLG} provided balanced gains, \mathcal{L}_{DR_NLG} showed stronger effects for Gemma-7B, and the combination of all three terms in BA-LoRA achieved the highest performance.

For NLU tasks, DeBERTa-v3-base was evalu-

Table 3: ID Performance Comparison of BERT-L and GPT-2-XL Using LoRA and BA-LoRA Methods on GLUE Benchmark. The best outcome is highlighted in **bold**.

Model	Methods	MNLI	SST-2	MRPC	CoLA	QNLI	QQP	RTE	STS-B	Avg
	LoRA	$87.24_{\pm 0.46}$	$93.19_{\pm0.25}$	$90.10_{\pm0.42}$	$64.73_{\pm 0.16}$	$93.13_{\pm0.55}$	$90.94_{\pm0.11}$	$73.14_{\pm0.51}$	$90.63_{\pm0.18}$	85.39
BERT-L	BA-LoRA	$\textbf{89.72}_{\pm 0.35}$	$\textbf{94.85}_{\pm 0.12}$	$\textbf{92.23}_{\pm 0.18}$	$\textbf{65.49}_{\pm 0.37}$	$\textbf{95.48}_{\pm 0.20}$	$91.72_{\pm0.19}$	$\textbf{75.77}_{\pm 0.34}$	$91.71_{\pm0.10}$	87.12
	LoRA	$85.28_{\pm 0.39}$	$95.38_{\pm0.46}$	$86.17_{\pm 0.31}$	$50.63_{\pm0.54}$	$89.42_{\pm 0.06}$	$88.56_{\pm 0.24}$	$72.29_{\pm 0.36}$	$89.27_{\pm 0.16}$	82.13
GPT-2-XL	BA-LoRA	$88.14_{\pm 0.28}$	$96.52_{\pm 0.38}$	$89.23_{\pm 0.19}$	$52.76_{\pm 0.26}$	$91.26_{\pm 0.11}$	$89.95_{\pm 0.37}$	74.57 $_{\pm 0.22}$	$90.83_{\pm 0.07}$	84.16

Table 4: OOD Performance Comparison of BERT-L and GPT-2-XL Using LoRA and BA-LoRA Methods on GLUE-x Benchmark. The best outcome is highlighted in **bold**.

Model	Methods	MNLI	SST-2	MRPC	CoLA	QNLI	QQP	RTE	STS-B	Avg
	LoRA	$85.19_{\pm0.24}$	$93.49_{\pm0.46}$	$89.93_{\pm0.13}$	$63.49_{\pm 0.22}$	$92.32_{\pm0.21}$	$87.73_{\pm 0.17}$	$73.65_{\pm0.25}$	$90.57_{\pm 0.06}$	84.55
BERT-L	BA-LoRA	$\textbf{87.91}_{\pm 0.17}$	$\textbf{94.18}_{\pm 0.85}$	$\textbf{90.62}_{\pm 0.34}$	$\textbf{65.81}_{\pm 0.86}$	$\textbf{93.04}_{\pm 0.25}$	$\textbf{89.06}_{\pm 0.38}$	$\textbf{75.41}_{\pm 0.18}$	$91.21_{\pm0.47}$	85.91
	LoRA	$87.02_{\pm 0.05}$	$95.11_{\pm 0.44}$	$86.81_{\pm 0.11}$	$60.95_{\pm0.38}$	$91.77_{\pm0.04}$	$87.59_{\pm 0.07}$	$78.76_{\pm 0.03}$	$89.25_{\pm0.13}$	84.66
GPT-2-XL	BA-LoRA	$89.58_{\pm 0.08}$	$96.40_{\pm 0.13}$	$\textbf{88.18}_{\pm 0.04}$	$63.11_{\pm 0.14}$	$92.68_{\pm 0.08}$	$88.62_{\pm 0.19}$	$81.21_{\pm 0.48}$	$90.37_{\pm0.27}$	86.27



Figure 3: Ablation Study of BA-LoRA Regularizations on GSM8K, MATH, and NLU Tasks. Here, "Reg" denotes "Regularization", and "w/o Reg" indicates "without regularization". \mathcal{L}_{CR} , \mathcal{L}_{DR} , and \mathcal{L}_{SVDR} represent the application of only the corresponding regularization, while "BA-LoRA" refers to the baseline using all regularizations.

ated on the GLUE benchmark using the corresponding regularization terms \mathcal{L}_{CR_NLU} , \mathcal{L}_{DR_NLU} , and \mathcal{L}_{SVDR_NLU} . The reported results represent the average performance across eight tasks in the GLUE benchmark: MNLI, SST-2, MRPC, CoLA, QNLI, QQP, RTE, and STS-B. As shown in Figure 3, the model without regularization achieved the lowest average score of 89.47. Each of the three NLUspecific regularization terms consistently improved performance, with \mathcal{L}_{CR_NLU} yielding a significant gain (average score of 89.78). The full BA-LoRA model, incorporating all three regularization terms, achieved the highest average score of 89.91, strongly demonstrating its effectiveness.

In summary, the ablation results demonstrate that the proposed regularization terms significantly enhance model performance for both NLG and NLU tasks. Each term contributes unique benefits, and their combination in BA-LoRA consistently achieves optimal results. This underscores the importance of integrating multiple regularization strategies to improve generalization and taskspecific performance.

5 Conclusion

This paper introduces BA-LoRA, a novel parameter-efficient fine-tuning method designed to mitigate catastrophic inheritance in pre-trained language models. BA-LoRA incorporates three key components: consistency regularization, diversity regularization, and singular value decomposition regularization. These components work in concert to preserve pre-training knowledge, enhance output diversity, and improve model generalization. Extensive experiments demonstrate that BA-LoRA consistently outperforms existing baselines on various NLG and NLU tasks while robust to noisy and imbalanced pre-training data. Furthermore, our ablation studies confirm the effectiveness of the three regularization terms both individually and in combination. These results highlight the potential of BA-LoRA as a general-purpose fine-tuning method for pre-trained language models and effectively address the key challenges of deploying these models in real applications.

5.1 Limitations

While BA-LoRA demonstrates significant improvements in mitigating catastrophic inheritance and enhancing model performance, several limitations warrant further investigation. Firstly, our evaluations primarily focus on English language tasks, which may limit the generalizability of our findings to other languages and specialized domains. Additionally, the computational overhead introduced by the consistency, diversity, and SVD regularizers adds complexity to the training process, potentially impacting efficiency. Furthermore, the impact of BA-LoRA on other forms of bias, such as fairness and societal stereotypes, remains unexplored. Lastly, the selection and weighting of regularization terms in BA-LoRA are fixed across different tasks, which may not be optimal for all scenarios.

References

AI@Meta. 2024. Llama 3 model card.

- Görkem Algan and Ilkay Ulusoy. 2021. Image classification with deep learning in the presence of noisy labels: A survey. *Knowledge-Based Systems*, 215:106771.
- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. 2021. Program synthesis with large language models. arXiv preprint arXiv:2108.07732.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. 2023. Qwen technical report. arXiv preprint arXiv:2309.16609.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. 2022. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*.
- Adrien Bardes, Jean Ponce, and Yann LeCun. 2021. Vicreg: Variance-invariance-covariance regularization for self-supervised learning. *arXiv preprint arXiv:2105.04906*.

- Solon Barocas and Andrew D Selbst. 2016. Big data's disparate impact. *Calif. L. Rev.*, 104:671.
- Abeba Birhane and Vinay Uday Prabhu. 2021. Large image datasets: A pyrrhic win for computer vision? In 2021 IEEE Winter Conference on Applications of Computer Vision (WACV), pages 1536–1546. IEEE.
- Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. 2021. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*.
- Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramer, and Chiyuan Zhang. 2022. Quantifying memorization across neural language models. *arXiv preprint arXiv:2202.07646*.
- Nicholas Carlini, Matthew Jagielski, Christopher A Choquette-Choo, Daniel Paleka, Will Pearce, Hyrum Anderson, Andreas Terzis, Kurt Thomas, and Florian Tramèr. 2023. Poisoning web-scale training datasets is practical. *arXiv preprint arXiv:2302.10149*.
- Isaac Caswell, Julia Kreutzer, Lisa Wang, Ahsan Wahab, Daan van Esch, Nasanbayar Ulzii-Orshikh, Allahsera Tapo, Nishant Subramani, Artem Sokolov, Claytone Sikasote, et al. 2021. Quality at a glance: An audit of web-crawled multilingual datasets. *arXiv e-prints*, pages arXiv–2103.
- Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, et al. 2024. A survey on evaluation of large language models. ACM Transactions on Intelligent Systems and Technology, 15(3):1–45.
- Hao Chen, Bhiksha Raj, Xing Xie, and Jindong Wang. 2024a. On catastrophic inheritance of large foundation models. arXiv preprint arXiv:2402.01909.
- Hao Chen, Jindong Wang, Ankit Shah, Ran Tao, Hongxin Wei, Xing Xie, Masashi Sugiyama, and Bhiksha Raj. 2024b. Understanding and mitigating the label noise in pre-training on downstream tasks. *Preprint*, arXiv:2309.17002.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles

Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. 2021. Evaluating large language models trained on code. *Preprint*, arXiv:2107.03374.

- Xinyang Chen, Sinan Wang, Mingsheng Long, and Jianmin Wang. 2019. Transferability vs. discriminability: Batch spectral penalization for adversarial domain adaptation. In *International conference on machine learning*, pages 1081–1090. PMLR.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. arXiv preprint arXiv:2110.14168.
- Damai Dai, Chengqi Deng, Chenggang Zhao, RX Xu, Huazuo Gao, Deli Chen, Jiashi Li, Wangding Zeng, Xingkai Yu, Y Wu, et al. 2024. Deepseekmoe: Towards ultimate expert specialization in mixture-of-experts language models. *arXiv preprint arXiv:2401.06066*.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2024. Qlora: Efficient finetuning of quantized llms. *Advances in Neural Information Processing Systems*, 36.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805.
- Ning Ding, Yujia Qin, Guang Yang, Fuchao Wei, Zonghan Yang, Yusheng Su, Shengding Hu, Yulin Chen, Chi-Min Chan, Weize Chen, et al. 2023. Parameter-efficient fine-tuning of large-scale pretrained language models. *Nature Machine Intelligence*, 5(3):220–235.
- Jesse Dodge, Maarten Sap, Ana Marasović, William Agnew, Gabriel Ilharco, Dirk Groeneveld, Margaret Mitchell, and Matt Gardner. 2021. Documenting large webtext corpora: A case study on the colossal clean crawled corpus. *arXiv preprint arXiv:2104.08758*.
- Guanting Dong, Hongyi Yuan, Keming Lu, Chengpeng Li, Mingfeng Xue, Dayiheng Liu, Wei Wang, Zheng Yuan, Chang Zhou, and Jingren Zhou. 2023. How abilities in large language models are affected by supervised fine-tuning data composition. *arXiv preprint arXiv:2310.05492*.
- Xinshuai Dong, Anh Tuan Luu, Min Lin, Shuicheng Yan, and Hanwang Zhang. 2021. How should pretrained language models be fine-tuned towards adversarial robustness? *Advances in Neural Information Processing Systems*, 34:4356–4369.
- Yanai Elazar, Akshita Bhagia, Ian Magnusson, Abhilasha Ravichander, Dustin Schwenk, Alane Suhr, Pete Walsh, Dirk Groeneveld, Luca Soldaini, Sameer

Singh, et al. 2023. What's in my big data? *arXiv* preprint arXiv:2310.20707.

- Lijie Fan, Kaifeng Chen, Dilip Krishnan, Dina Katabi, Phillip Isola, and Yonglong Tian. 2024. Scaling laws of synthetic images for model training... for now. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 7382– 7392.
- Benoît Frénay and Michel Verleysen. 2013. Classification in the presence of label noise: a survey. *IEEE* transactions on neural networks and learning systems, 25(5):845–869.
- Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, et al. 2020. The pile: An 800gb dataset of diverse text for language modeling. *arXiv preprint arXiv:2101.00027*.
- Itai Gat, Idan Schwartz, Alexander Schwing, and Tamir Hazan. 2020. Removing bias in multi-modal classifiers: Regularization by maximizing functional entropies. *Advances in Neural Information Processing Systems*, 33:3197–3208.
- Demi Guo, Alexander M Rush, and Yoon Kim. 2020. Parameter-efficient transfer learning with diff pruning. *arXiv preprint arXiv:2012.07463*.
- Karen Hambardzumyan, Hrant Khachatrian, and Jonathan May. 2021. Warp: Word-level adversarial reprogramming. *arXiv preprint arXiv:2101.00121*.
- Zeyu Han, Chao Gao, Jinyang Liu, Sai Qian Zhang, et al. 2024. Parameter-efficient fine-tuning for large models: A comprehensive survey. *arXiv preprint arXiv:2403.14608*.
- Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. 2021a. Towards a unified view of parameter-efficient transfer learning. *arXiv preprint arXiv:2110.04366*.
- Pengcheng He, Jianfeng Gao, and Weizhu Chen. 2021b. Debertav3: Improving deberta using electra-style pretraining with gradient-disentangled embedding sharing. *Preprint*, arXiv:2111.09543.
- Danny Hernandez, Tom Brown, Tom Conerly, Nova DasSarma, Dawn Drain, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Tom Henighan, Tristan Hume, et al. 2022. Scaling laws and interpretability of learning from repeated data. *arXiv preprint arXiv:2205.10487*.
- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin de Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019a. Parameter-efficient transfer learning for nlp. *Preprint*, arXiv:1902.00751.
- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly.

2019b. 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. 2021. Lora: Low-rank adaptation of large language models. *Preprint*, arXiv:2106.09685.
- Nam Hyeon-Woo, Moon Ye-Bin, and Tae-Hyun Oh. 2021. Fedpara: Low-rank hadamard product for communication-efficient federated learning. *arXiv* preprint arXiv:2108.06098.
- Joel Jang, Seonghyeon Ye, Sohee Yang, Joongbo Shin, Janghoon Han, Gyeonghun Kim, Stanley Jungkyu Choi, and Minjoon Seo. 2021. Towards continual knowledge learning of language models. *arXiv preprint arXiv:2110.03215*.
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. *arXiv preprint arXiv:2310.06825*.
- Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. 2024a. Mixtral of experts. *arXiv preprint arXiv:2401.04088*.
- Minhao Jiang, Ken Ziyu Liu, Ming Zhong, Rylan Schaeffer, Siru Ouyang, Jiawei Han, and Sanmi Koyejo. 2024b. Investigating data contamination for pre-training language models. *arXiv preprint arXiv:2401.06059*.
- Divyansh Kaushik, Eduard Hovy, and Zachary C Lipton. 2019. Learning the difference that makes a difference with counterfactually-augmented data. *arXiv* preprint arXiv:1909.12434.
- Tao Lei, Junwen Bai, Siddhartha Brahma, Joshua Ainslie, Kenton Lee, Yanqi Zhou, Nan Du, Vincent Zhao, Yuexin Wu, Bo Li, et al. 2023. Conditional adapters: Parameter-efficient transfer learning with fast inference. Advances in Neural Information Processing Systems, 36:8152–8172.
- Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning. *Preprint*, arXiv:2104.08691.
- Xiang Lisa Li and Percy Liang. 2021a. Prefix-tuning: Optimizing continuous prompts for generation. *arXiv preprint arXiv:2101.00190*.
- Xiang Lisa Li and Percy Liang. 2021b. Prefixtuning: Optimizing continuous prompts for generation. *Preprint*, arXiv:2101.00190.
- Zhaojiang Lin, Andrea Madotto, and Pascale Fung. 2020. Exploring versatile generative language model via parameter-efficient transfer learning. *arXiv preprint arXiv:2004.03829*.

- Shih-Yang Liu, Chien-Yi Wang, Hongxu Yin, Pavlo Molchanov, Yu-Chiang Frank Wang, Kwang-Ting Cheng, and Min-Hung Chen. 2024. Dora: Weightdecomposed low-rank adaptation. *arXiv preprint arXiv:2402.09353*.
- Xiao Liu, Kaixuan Ji, Yicheng Fu, Weng Lam Tam, Zhengxiao Du, Zhilin Yang, and Jie Tang. 2021. Ptuning v2: Prompt tuning can be comparable to finetuning universally across scales and tasks. *arXiv preprint arXiv:2110.07602*.
- Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding, Yujie Qian, Zhilin Yang, and Jie Tang. 2023. Gpt understands, too. *AI Open*.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. arXiv preprint arXiv:1907.11692.
- Zhuang Liu and Kaiming He. 2024. A decade's battle on dataset bias: Are we there yet? *arXiv preprint arXiv:2403.08632*.
- Shayne Longpre, Gregory Yauney, Emily Reif, Katherine Lee, Adam Roberts, Barret Zoph, Denny Zhou, Jason Wei, Kevin Robinson, David Mimno, et al. 2023. A pretrainer's guide to training data: Measuring the effects of data age, domain coverage, quality, & toxicity. *arXiv preprint arXiv:2305.13169*.
- Ilya Loshchilov and Frank Hutter. 2017. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*.
- Andrew Maas, Raymond E Daly, Peter T Pham, Dan Huang, Andrew Y Ng, and Christopher Potts. 2011. Learning word vectors for sentiment analysis. In Proceedings of the 49th annual meeting of the association for computational linguistics: Human language technologies, pages 142–150.
- Rabeeh Karimi Mahabadi, Sebastian Ruder, Mostafa Dehghani, and James Henderson. 2021. Parameterefficient multi-task fine-tuning for transformers via shared hypernetworks. *arXiv preprint arXiv:2106.04489*.
- Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi. 2022. When not to trust language models: Investigating effectiveness of parametric and non-parametric memories. *arXiv preprint arXiv:2212.10511*.
- Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. 2021. A survey on bias and fairness in machine learning. *ACM computing surveys (CSUR)*, 54(6):1–35.
- Fanxu Meng, Zhaohui Wang, and Muhan Zhang. 2024. Pissa: Principal singular values and singular vectors adaptation of large language models. *Preprint*, arXiv:2404.02948.

- Curtis G. Northcutt, Lu Jiang, and Isaac L. Chuang. 2021. Confident learning: Estimating uncertainty in dataset labels. Journal of Artificial Intelligence Research (JAIR), 70:1373–1411.
- OpenAI. 2023. Gpt-4 technical report. Preprint, arXiv:2303.08774.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. Advances in neural information processing systems, 35:27730–27744.
- Shubham Parashar, Zhiqiu Lin, Tian Liu, Xiangjue Dong, Yanan Li, Deva Ramanan, James Caverlee, and Shu Kong. 2024. The neglected tails in visionlanguage models. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 12988-12997.
- Ali Saheb Pasand and Pouya Bashivan. Rgp: Achieving memory-efficient model fine-tuning via randomized gradient projection.
- Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Alessandro Cappelli, Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. 2023. The refinedweb dataset for falcon llm: outperforming curated corpora with web data, and web data only. arXiv preprint arXiv:2306.01116.
- Jonas Pfeiffer, Aishwarya Kamath, Andreas Rücklé, Kyunghyun Cho, and Iryna Gurevych. 2020. Adapterfusion: Non-destructive task composition for transfer learning. arXiv preprint arXiv:2005.00247.
- Xiangyu Qi, Yi Zeng, Tinghao Xie, Pin-Yu Chen, Ruoxi Jia, Prateek Mittal, and Peter Henderson. 2023. Finetuning aligned language models compromises safety, even when users do not intend to! arXiv preprint arXiv:2310.03693.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. OpenAI blog, 1(8):9.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2024. Direct preference optimization: Your language model is secretly a reward model. Advances in Neural Information Processing Systems, 36.
- Manley Roberts, Himanshu Thakur, Christine Herlihy, Colin White, and Samuel Dooley. 2023. Data contamination through the lens of time. arXiv preprint arXiv:2310.10628.
- Rylan Schaeffer. 2023. Pretraining on the test set is all you need. arXiv preprint arXiv:2309.08632.

- Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc V. Le, Geoffrey E. Hinton, and Jeff Dean. 2017. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. In 5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings. OpenReview.net.
- Hwanjun Song, Minseok Kim, Dongmin Park, Yooju Shin, and Jae-Gil Lee. 2022. Learning from noisy labels with deep neural networks: A survey. IEEE transactions on neural networks and learning systems, 34(11):8135-8153.
- Lichao Sun, Yue Huang, Haoran Wang, Siyuan Wu, Qihui Zhang, Chujie Gao, Yixin Huang, Wenhan Lyu, Yixuan Zhang, Xiner Li, et al. 2024. Trustllm: Trustworthiness in large language models. arXiv preprint arXiv:2401.05561.
- Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. 2024. Gemma: Open models based on gemini research and technology. arXiv preprint arXiv:2403.08295.
- Antonio Torralba and Alexei A Efros. 2011. Unbiased look at dataset bias. In CVPR 2011, pages 1521-1528. IEEE.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. arXiv preprint arXiv:2307.09288.
- Nirali Vaghani and Mansi Thummar. 2023. Flipkart product reviews with sentiment dataset.
- 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.
- Laurens Van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-sne. Journal of machine learning research, 9(11).
- 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. In International Conference on Learning Representations.
- Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, Qingwei Lin, and Daxin Jiang. 2024. Wizardlm: Empowering large pre-trained language models to follow complex instructions. In The Twelfth International Conference on Learning Representations.

- Hu Xu, Saining Xie, Xiaoqing Ellen Tan, Po-Yao Huang, Russell Howes, Vasu Sharma, Shang-Wen Li, Gargi Ghosh, Luke Zettlemoyer, and Christoph Feichtenhofer. 2023a. Demystifying clip data. *arXiv preprint arXiv:2309.16671*.
- Lingling Xu, Haoran Xie, Si-Zhao Joe Qin, Xiaohui Tao, and Fu Lee Wang. 2023b. Parameter-efficient fine-tuning methods for pretrained language models: A critical review and assessment. *arXiv preprint arXiv:2312.12148*.
- Linyi Yang, Shuibai Zhang, Libo Qin, Yafu Li, Yidong Wang, Hanmeng Liu, Jindong Wang, Xing Xie, and Yue Zhang. 2022. Glue-x: Evaluating natural language understanding models from an outof-distribution generalization perspective. *arXiv preprint arXiv:2211.08073*.
- Yu Yang, Aaditya K Singh, Mostafa Elhoushi, Anas Mahmoud, Kushal Tirumala, Fabian Gloeckle, Baptiste Rozière, Carole-Jean Wu, Ari S Morcos, and Newsha Ardalani. 2023. Decoding data quality via synthetic corruptions: Embedding-guided pruning of code data. *arXiv preprint arXiv:2312.02418*.
- Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, et al. 2024. Yi: Open foundation models by 01. ai. *arXiv preprint arXiv:2403.04652*.
- Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. 2023. Metamath: Bootstrap your own mathematical questions for large language models. *arXiv preprint arXiv:2309.12284*.
- Elad Ben Zaken, Shauli Ravfogel, and Yoav Goldberg. 2021. Bitfit: Simple parameter-efficient fine-tuning for transformer-based masked languagemodels. *arXiv preprint arXiv:2106.10199*.
- Qingru Zhang, Minshuo Chen, Alexander Bukharin, Nikos Karampatziakis, Pengcheng He, Yu Cheng, Weizhu Chen, and Tuo Zhao. 2023a. Adalora: Adaptive budget allocation for parameter-efficient finetuning. *arXiv preprint arXiv:2303.10512*.
- 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 finetuning. *Preprint*, arXiv:2303.10512.
- Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015. Character-level convolutional networks for text classification. *Advances in neural information processing systems*, 28.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. 2023. A survey of large language models. arXiv preprint arXiv:2303.18223.

- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2024. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36.
- Beier Zhu, Kaihua Tang, Qianru Sun, and Hanwang Zhang. 2024. Generalized logit adjustment: Calibrating fine-tuned models by removing label bias in foundation models. *Advances in Neural Information Processing Systems*, 36.
- Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. 2015 IEEE International Conference on Computer Vision (ICCV).
- Bojia Zi, Xianbiao Qi, Lingzhi Wang, Jianan Wang, Kam-Fai Wong, and Lei Zhang. 2023. Delta-lora: Fine-tuning high-rank parameters with the delta of low-rank matrices. *arXiv preprint arXiv:2309.02411*.
- Andy Zou, Zifan Wang, J Zico Kolter, and Matt Fredrikson. 2023. Universal and transferable adversarial attacks on aligned language models. *arXiv preprint arXiv:2307.15043*.

Appendix

Contonto

5 A	Bacl	kground
	A.1	Challenges of Bias and Noise in
		Pre-training Data
	A.2	Mitigating Bias and Noise through
		Parameter-Efficient Fine-Tuning
		Methods
	A.3	Examples of noise in pre-training
		data
B	Deta	ils of Models and Datasets
	B .1	Details of Models
	B .2	Details of Datasets
	B .3	Specific Hyperparameter Settings
		of RoBERTa-large and DeBERTa-
		v3-base on GLUE
	B .4	Specific Hyperparameter Settings
		of BERT-L and GPT-2-XL on
		GLUE and GLUE-x
	B.5	Model Evaluation Details
C	Mor	e Experiments
	C .1	Analysis on Different Sizes and
		Types of Models
	C.2	Impact of Regularization on LoRA
		Variants Performance
	C.3	t-SNE Visualizations of Feature
		Evolution during the Fine-tuning
		with LoRA and BA-LoRA
D	Mor	e Discussions
	D.1	Computational Cost Analysis
	D.2	Future works

A.1 Challenges of Bias and Noise in Pre-training Data

Bias and noise within pre-training datasets present significant hurdles in constructing dependable machine-learning models. Mislabeled data and imbalanced distributions can lead to models that not only underperform on downstream tasks but also reinforce existing biases in the data (Torralba and Efros, 2011; Barocas and Selbst, 2016; Mehrabi et al., 2021). This issue is especially problematic in large-scale datasets, where manual curation is impractical. Consequently, reliance on automated data collection methods may introduce various inaccuracies and biases (Northcutt et al., 2021; Birhane and Prabhu, 2021). The challenge becomes even more severe when dealing with realworld, instance-dependent label noise. Models trained on such data may inadvertently learn these inaccuracies, resulting in poor generalization (Frénay and Verleysen, 2013; Song et al., 2022; Algan and Ulusoy, 2021). Addressing these challenges is essential for advancing machine learning and ensuring models are both effective and equitable.

A.2 Mitigating Bias and Noise through Parameter-Efficient Fine-Tuning Methods

To counteract the adverse effects of bias and noise in pre-training data, parameter-efficient fine-tuning methods have emerged as promising solutions. These approaches aim to adapt pre-trained models to new tasks with minimal parameter updates, thereby reducing the risk of overfitting to noisy or biased data (Houlsby et al., 2019b; Zaken et al., 2021; Lester et al., 2021). Techniques such as integrating lightweight adaptation modules (Pfeiffer et al., 2020; Jang et al., 2021), utilizing prefix tuning (Li and Liang, 2021b; Liu et al., 2021), and employing low-rank adaptations (Hu et al., 2021; Ding et al., 2023) enable efficient model refinement while preserving the valuable representations acquired during pre-training. Selectively fine-tuning specific model components can enhance performance on downstream tasks, improve generalization, and reduce the influence of noise and bias (Zaken et al., 2021; Mahabadi et al., 2021; Guo et al., 2020). This strategy not only results in more robust models but also contributes to the development of fairer AI systems by directly addressing fundamental data quality issues.

A.3 Examples of noise in pre-training data

Pre-training data typically originates from largescale internet sources, which inevitably contain noise and imbalance. Many advanced pre-trained models, such as LLaMA-2-7B/13B (Touvron et al., 2023), Mistral-7B-v0.1 (Jiang et al., 2023), Gemma-7B (Jiang et al., 2023), and GPT-4 (OpenAI, 2023), are trained using large amounts of unlabeled internet text data. These datasets are often not thoroughly cleaned or corrected, leading to training corpora that include irrelevant or inaccurate information. Consequently, during the subsequent fine-tuning phase, models struggle to effectively filter out these undesirable contents, adversely affecting their performance on downstream tasks.

Given the noise and imbalance issues present in pre-training data, understanding the specific types of noise is crucial for improving model performance. Below, we summarize some common examples of noise found in pre-training datasets:

Low quality bias

Duplicity: The presence of identical or similar content in the data can lead to overfitting and privacy leakage risks (Elazar et al., 2023), (Carlini et al., 2022), (Hernandez et al., 2022). *Corruption/noise:* Inconsistent or erroneous inputs in the training data can affect model robustness and downstream task performance (Elazar et al., 2023), (Fan et al., 2024), (Caswell et al., 2021). *Contamination:* Leakage of the test set into the training set may lead to distorted model evaluation results (Roberts et al., 2023), (Schaeffer, 2023), (Jiang et al., 2024b).

Skewed distribution bias

Category imbalance: Too few samples of certain categories cause the model to perform poorly in predicting these categories, producing bias (Xu et al., 2023a), (Zhu et al., 2024), (Parashar et al., 2024).

Unethical content bias

Toxic and harmful content: The training data may contain offensive, biased, or harmful content that may cause the model to generate harmful or inappropriate outputs (Zou et al., 2023), (Sun et al., 2024).

B Details of Models and Datasets

To evaluate the effectiveness of our approach, we conduct experiments using several prominent language models and assess their performance on a diverse array of datasets, covering both NLG and NLU tasks. Specifically, for language generation models, we include LLaMA 2-7B (Touvron et al., 2023), LLaMA 3-8B (AI@Meta, 2024), Mistral-7B (Jiang et al., 2023), Gemma-7B (Team et al., 2024), and GPT-2-XL (Radford et al., 2019). For language understanding models, we use BERT-Large (BERT-L) (Devlin et al., 2018), and DeBERTa-v3-base (He et al., 2021b). For the datasets, we employ a wide range of tasks in both Natural Language Generation (GSM8K (Cobbe et al., 2021), MATH (Yu et al., 2023), HumanEval (Chen et al., 2021), MBPP (Austin et al., 2021), MT-Bench (Zheng et al., 2024)) and Natural Language Understanding. In the latter, we assess in-domain (ID) performance using the GLUE

benchmark (Wang et al., 2018) and out-of-domain (OOD) generalization using the GLUE-X benchmark (Yang et al., 2022). These datasets span a broad range of challenges, allowing for a thorough examination of our method's generalization capabilities.

B.1 Details of Models

We use a variety of pre-trained language models, including Meta AI's LLaMA-2-7B and LLaMA-2-13B (Touvron et al., 2023) and the latest LLaMA-3-8B and LLaMA-3-70B (AI@Meta, 2024), which have good performance in natural language generation tasks. In addition, we also use Mistral AI's Mistral-7B-v0.1 (Jiang et al., 2023) optimized for medium-sized model efficiency, and Google's lightweight open-source model Gemma-7B (Jiang et al., 2023), which performs well in tasks such as question-answering summarization and reasoning. Alibaba Cloud's Qwen-1.5-7B (Bai et al., 2023) model also provides strong language understanding and generation capabilities, while the 34B parameter Yi-1.5-34B (Young et al., 2024) is designed for high-level language tasks. DeepSeek-MoE-16B (Dai et al., 2024) is a model that uses expert routing to increase capacity without significantly increasing computational costs. Mixtral-8x7B-v0.1 (Jiang et al., 2024a) is a Sparse Mixture of Expert models that efficiently utilizes active parameters to outperform larger models like Llama 2 70B and GPT-3.5 across several benchmarks. We also leveraged mature models such as BERT-Large (Devlin et al., 2018), RoBERTa-large (Liu et al., 2019), DeBERTa-v3-base (He et al., 2021b), and GPT-2-XL (Radford et al., 2019), which continue to set standards in natural language processing and text generation tasks.

Table 5 presents an overview of the pre-trained language models used in our study. BERT-Large (BERT-L) and RoBERTa-Large (RoBERTa-L) are pre-trained on the BooksCorpus and English Wikipedia datasets using a masked language modeling objective. In contrast, GPT-2-XL is pre-trained on WebText with an autoregressive language modeling objective. Additionally, DeBERTa-v3-base is trained on a diverse dataset comprising Wikipedia, BooksCorpus, OpenWebText, CC-News, and Stories, utilizing a replaced token detection objective with Gradient Disentangled Embedding Sharing (GDES). These models span a variety of architectures and pre-training strategies, offering a robust basis for evaluating the performance of our proTable 5: Comparison of Pre-trained Data and Methods for Various Language Models.

Model	Pre-trained Data	Pre-training Method
BERT-L (Devlin et al., 2018)	BooksCorpus and English Wikipedia	Masked Language Modeling
RoBERTa-L (Liu et al., 2019)	BooksCorpus and English Wikipedia	Masked Language Modeling
GPT-2-XL (Radford et al., 2019)	WebText	Autoregressive Language Modeling
DeBERTa-v3-base (He et al., 2021b)	Wikipedia, BooksCorpus, OpenWebText, CC-News, and Stories	Replaced Token Detection with GDES

posed approach.

B.2 Details of Datasets

Table 6 provides an overview of the GLUE benchmark datasets and their evaluation metrics. The GLUE benchmark comprises a diverse set of natural language understanding tasks, including grammatical acceptability (CoLA), sentiment analysis (SST-2), paraphrase detection (MRPC and QQP), sentence similarity (STS-B), natural language inference (MNLI, QNLI, and RTE), and coreference resolution (WNLI). The number of training examples varies significantly across datasets, from as few as 634 in WNLI to as many as 393,000 in MNLI. Tasks involve binary or multi-class classification, with up to five classes in STS-B. Evaluation metrics are tailored to each task, employing accuracy, F1 score, Matthews correlation coefficient, and Pearson/Spearman correlation coefficients where appropriate. This comprehensive suite serves as a standard benchmark for assessing and comparing the performance of models across a wide array of linguistic challenges.

Table 7 summarizes the GLUE-X out-of-domain tasks employed for evaluating transfer performance. The datasets cover a broad spectrum of natural language understanding tasks, including natural language inference (SNLI, HANs, SciTail, MNLI mismatched), sentiment analysis (IMDB), question answering (NewsQA), semantic relatedness (SICK), and grammatical error detection (Grammar Test). Each task involves binary classification, with test sizes ranging from 9,832 samples (MNLI mismatched) to 570,152 samples (SNLI). Accuracy is the primary evaluation metric across most datasets, except for the Grammar Test, which uses the Matthews correlation coefficient. These diverse tasks provide a comprehensive benchmark for assessing the models' ability to generalize across different domains and tasks.

Table 8 summarizes the evaluation metrics for the natural language generation (NLG) tasks. Specifically, we use Accuracy for GSM8K and MATH; Pass@1 for HumanEval and MBPP, indicating the percentage of first generated code snippets that pass all unit tests; and GPT-4 Evaluation for MT-Bench, where GPT-4 assesses the quality of the model's responses.

B.3 Specific Hyperparameter Settings of RoBERTa-large and DeBERTa-v3-base on GLUE

We fine-tuned the RoBERTa-large and DeBERTav3-base models on the GLUE benchmark datasets using carefully selected hyperparameters tailored to each task. For RoBERTa-large, we trained on MNLI and SST-2 for 10 epochs with a batch size of 32; MNLI employed a learning rate of 1×10^{-4} , while SST-2 used 2×10^{-4} , both with LoRA_alpha set to 16. Smaller datasets such as MRPC, CoLA, and RTE were trained for 20 epochs with batch sizes of 16, utilizing higher learning rates ranging from 3×10^{-4} to 6×10^{-4} and LoRA_alpha values of 8 or 16. For DeBERTa-v3-base, MNLI was trained for 5 epochs with a batch size of 16, a learning rate of 5×10^{-5} , and LoRA_alpha set to 8. Datasets such as SST-2 and MRPC were trained for 20 epochs with batch sizes of 16 or 32, learning rates between 3×10^{-5} and 2×10^{-4} , and LoRA_alpha of 8. Notably, RTE was trained for 50 epochs with a batch size of 16, a learning rate of 1×10^{-4} , and LoRA_alpha of 8. The LoRA_alpha parameter was set to either 8 or 16, depending on the model and dataset. In all cases, the LoRA_rank was set to 8. These hyperparameters were meticulously chosen to suit the specific requirements of each dataset, ensuring rigorous and optimal training across tasks such as natural language inference, sentiment analysis, paraphrase detection, linguistic acceptability, and semantic textual similarity.

B.4 Specific Hyperparameter Settings of BERT-L and GPT-2-XL on GLUE and GLUE-x

To ensure consistent and reliable performance, the BERT-Large (BERT-L) and GPT-2-XL models were trained on the GLUE benchmark tasks using three different random seeds per task over 10 epochs. A hyperparameter search was conducted over learning rates $\{2 \times 10^{-5}, 3 \times 10^{-5}, 5 \times 10^{-5}\}$,

Dataset	Task Type	Classes	Train Examples	Metric	Description
CoLA	Acceptability	2	8.5k	Matthews Corr.	Grammatical acceptability
SST-2	Sentiment	2	67k	Accuracy	Sentiment analysis
MRPC	Paraphrase	2	3.7k	Accuracy/F1	Paraphrase detection
QQP	Paraphrase	2	364k	Accuracy/F1	Duplicate question detection
STS-B	Similarity	5	7k	Pearson/Spearman Corr.	Sentence similarity
MNLI	NLI	3	393k	Accuracy	Multi-genre NLI
QNLI	NLI/QA	2	108k	Accuracy	QA/NLI converted from SQuAD
RTE	NLI	2	2.5k	Accuracy	Textual entailment
WNLI	Coreference	2	634	Accuracy	Winograd Schema Challenge

Table 6: GLUE Benchmark Datasets and Evaluation Metrics

Table 7: Summary of G	LUE-X Out-of-Domain	Tasks for Transfer	Performance Evaluation
-----------------------	---------------------	--------------------	------------------------

Dataset	Task Type	Classes	Train Examples	Metric	Description
SNLI	NLI	2	570k	Accuracy	Sentence-level inference tasks
IMDB	Sentiment	2	50k	Accuracy	Movie review sentiment analysis
HANs	NLI	2	60k	Accuracy	Adversarial NLI examples to test models
NewsQA	QA	2	119k	Accuracy	QA from news articles
SICK	Semantic Relatedness	2	9.8k	Accuracy	Semantic relatedness and entailment
Grammar Test	Grammar Detection	2	304k	Matthews Corr.	Grammatical error detection
SciTail	NLI	2	26.5k	Accuracy	Science question entailment
MNLI mismatched	NLI	2	9.8k	Accuracy	NLI with mismatched genres

Table 8: Evaluation Metrics for NLG Datasets

Datasets	GSM8K	MATH	HumanEval	MBPP	MT-Bench
Metric	Accuracy	Accuracy	Pass@1	Pass@1	GPT-4 Evaluation

and a batch size of 32 was chosen to balance computational efficiency and memory usage. For finetuning, the training schedule was adjusted to 20 epochs for smaller datasets, while larger datasets such as QNLI, MNLI, and QQP were trained for 5 epochs. Learning rates were explored within $\{2 \times 10^{-4}, 3 \times 10^{-4}, 5 \times 10^{-4}\}$. The parameters were set with LoRA_rank = 8 and LoRA_alpha = 16, with the batch size reduced to 16 due to increased model complexity. All other parameters, including max_length, adhered to Hugging Face Transformers guidelines¹.

Regarding the GLUE-x tasks, BERT-L and GPT-2-XL models trained on GLUE were evaluated without further fine-tuning. GLUE-x encompasses 13 out-of-distribution (OOD) tasks, introducing domain shifts. For sentiment analysis, models fine-tuned on SST-2 were evaluated on test sets from IMDB (Maas et al., 2011), Yelp (Zhang et al., 2015), Amazon (Kaushik et al., 2019), and Flipkart (Vaghani and Thummar, 2023), offering a broader assessment of domain variability and testing the robustness beyond SST-2.

For t-SNE visualization, we used the MNLI sub-

set from GLUE due to its diverse linguistic styles and label distributions. Training was limited to one epoch to expedite the process, while still providing insights into how well the models differentiate between classes and sentence structures.

B.5 Model Evaluation Details

For evaluation, we employed publicly available frameworks. The model's code generation capabilities were assessed using datasets like HumanEval and MBPP through the BigCode Evaluation Harness². Instruction-following performance was evaluated using MTBench³.

C More Experiments

C.1 Analysis on Different Sizes and Types of Models

This experiment compares LoRA, PiSSA, and BA-LoRA across ten models: LLaMA-2-7/13B (Touvron et al., 2023), LLaMA-3-8B/70B (AI@Meta, 2024), Mistral-7B-v0.1 (Jiang et al., 2023), Gemma-7B (Jiang et al., 2023), Qwen1.5-7B (Bai

²https://github.com/bigcode-project/

bigcode-evaluation-harness

³https://github.com/lm-sys/FastChat

¹https://github.com/huggingface/transformers



Figure 4: Performance Comparison of Different Models on GSM8K and HumanEval Benchmarks.

et al., 2023), Yi-1.5-34B (Young et al., 2024) and Mixture-of-Experts (MoE (Shazeer et al., 2017)) models: DeepSeek-MoE-16B (Dai et al., 2024) and Mixtral-8x7B-v0.1 (Jiang et al., 2024a). These models were fine-tuned on the MetaMathQA-100K and CodeFeedback-100K datasets and evaluated on the GSM8K and HumanEval benchmarks. As depicted in Figure 4, BA-LoRA consistently surpasses LoRA and PiSSA across all models and tasks, underscoring its superior ability to enhance model generalization.

C.2 Impact of Regularization on LoRA Variants Performance

In this experiment, we evaluated the impact of the regularization terms on multiple LoRA variants using the LLaMA-2-7B. Table 9 shows the performance comparison of LoRA, DoRA, PiSSA, and BA-LoRA with and without regularization terms, where "Reg" refers to the three regularization terms designed for each NLG task.

The experimental results indicate that incorporating regularization terms into LoRA and DoRA architectures significantly enhances performance across all evaluated tasks. This finding demonstrates that regularization techniques are broadly effective when applied to different LoRA variants. Furthermore, BA-LoRA, which integrates PiSSA with regularization, achieves the best performance across various tasks and substantially improves the model's generalization capabilities.

C.3 t-SNE Visualizations of Feature Evolution during the Fine-tuning with LoRA and BA-LoRA

This section provides more detailed t-SNE visualization results to compare the feature evolution during fine-tuning of LoRA and BA-LoRA.

Figure 5 shows that during LoRA fine-tuning of BERT-L, the feature separation is slow, with the class distributions remaining scattered and overlapping even towards the end. In contrast, Figure 6 demonstrates that with BA-LoRA fine-tuning, class separation begins earlier and is much clearer, ultimately forming a distinct "Y" shape with well-defined class boundaries.

Similarly, Figure 7 shows that during LoRA finetuning of GPT-2 XL, the feature clusters remain scattered and overlapping throughout the training, with only minimal separation between classes by the final steps. In contrast, Figure 8 demonstrates that BA-LoRA fine-tuning results in much clearer and more distinct class separation, with well-defined boundaries emerging earlier in the training process and becoming more pronounced over time.

D More Discussions

Here, we offer further insights into our work.

D.1 Computational Cost Analysis

Parameter-efficient fine-tuning (PEFT) methods are designed to minimize computational overhead while maintaining competitive performance relative to full fine-tuning. While PEFT methods are computationally efficient overall, specific components can introduce additional demands that warrant careful consideration. For instance, consis-

Methods	#Params	GSM8K	MATH	HumanEval	MBPP	MT-Bench	Avg
Full FT	6738M	$49.13_{\pm 0.21}$	$7.29_{\pm 0.22}$	$21.20_{\pm 0.30}$	$35.59_{\pm0.25}$	$4.91_{\pm 0.01}$	23.62
LoRA	320M	$42.47_{\pm 0.29}$	$5.60_{\pm0.35}$	$17.03_{\pm 0.61}$	$\overline{32.77_{\pm 0.46}}$	$4.62_{\pm 0.11}$	20.50
LoRA + Reg	320M	$51.91_{\pm 0.30}$	$8.63_{\pm0.38}$	$21.11_{\pm 0.03}$	$33.85_{\pm0.16}$	$4.75_{\pm0.35}$	23.83
DoRA	321M	$42.12_{\pm 0.16}$	$\overline{6.28_{\pm 0.19}}$	$16.97_{\pm 0.08}$	$22.07_{\pm 0.38}$	$4.53_{\pm 0.47}$	18.41
PiSSA	320M	$52.37_{\pm 0.52}$	$7.76_{\pm 0.19}$	$21.55_{\pm 0.44}$	$33.09_{\pm 0.57}$	$4.87_{\pm 0.06}$	23.93
DoRA + Reg	321M	$52.80_{\pm 0.38}$	$8.05_{\pm0.07}$	$21.94_{\pm0.14}$	$34.61_{\pm 0.12}$	$5.02_{\pm0.10}$	<u>24.48</u>
BA-LoRA	320M	$54.53_{\pm 0.41}$	$9.21_{\pm 17}$	$\overline{23.58_{\pm 0.25}}$	$\textbf{36.86}_{\pm 0.31}$	$5.11_{\pm 0.05}$	25.86

Table 9: Effect of regularization term in different LoRA variants. The best and second-best results are highlighted in **bold** and <u>underline</u>.



Figure 5: t-SNE Visualization of Feature Evolution during LoRA Fine-Tuning of BERT-L.

tency regularization often necessitates an additional forward pass during training, leading to increased activation storage requirements. Similarly, operations like singular value decomposition (SVD) present scalability challenges for large matrices, as traditional SVD relies on eigenvalue computations that are computationally intensive and inherently CPU-bound. These factors highlight the importance of scalability, especially for high-dimensional models like LLaMA, where even relatively small matrices (e.g., 4096×4096) demand substantial processing time. To address these challenges, we introduce a **randomized SVD** approach, which achieves significant computational savings while preserving sufficient accuracy. Randomized SVD leverages random projections to approximate the dominant singular values and vectors, enabling faster computation and enhanced scalability for large matrices. This optimization ensures that computational demands remain manageable, even for models with high-dimensional weight matrices, thereby supporting scalability across diverse architectures.

To assess the practical implications of our ap-



Figure 6: t-SNE Visualization of Feature Evolution during BA-LoRA Fine-Tuning of BERT-L.

proach, we conducted experiments measuring the average training time and memory usage across multiple runs. Specifically, we fine-tuned the DeBERTa-v3-base model on the full GLUE benchmark using an NVIDIA A40 GPU. The results for three PEFT methods—LoRA, PiSSA, and our proposed BA-LoRA—are summarized below:

- LoRA: Averaged 14.45 hours of training time with an average memory usage of 6.3 GB.
- **PiSSA**: Averaged 12.07 hours of training time with an average memory usage of 6.5 GB.
- **BA-LoRA**: Averaged 15.82 hours of training time with an average memory usage of 7.1 GB.

Although BA-LoRA incurs marginally higher computational costs compared to LoRA and PiSSA, these increases stem primarily from its advanced regularization mechanisms, including the innovative use of randomized SVD. Crucially, these mechanisms effectively mitigate the catastrophic forgetting problem, resulting in substantial performance improvements. Considering the trade-off between computational overhead and model performance, the additional costs of BA-LoRA are both manageable and well-justified, especially in scenarios where achieving superior accuracy is paramount.

D.2 Future works

Future research should extend assessments of BA-LoRA to multilingual settings and specialized domains to ensure broader applicability. Exploring optimization techniques could help reduce the computational overhead introduced by the regularizers, balancing performance gains with efficiency. Investigating the impact of BA-LoRA on other forms of bias, including fairness and societal stereotypes, is crucial for developing more equitable models. Additionally, refining the selection and weighting of regularization terms-possibly through automated or dynamic adjustment methods-could enhance adaptability across different tasks and models. Testing the scalability of BA-LoRA on larger models with hundreds of billions of parameters and exploring its integration with other bias mitigation strategies may yield synergistic effects and further improve model robustness.



Figure 7: t-SNE Visualization of Feature Evolution during LoRA Fine-Tuning of GPT-2-XL.



Figure 8: t-SNE Visualization of Feature Evolution during BA-LoRA Fine-Tuning of GPT-2-XL.