

# LoRA UNLEASHED: EFFORTLESSLY ADVANCING FROM LOW TO ARBITRARY RANK

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Low-Rank Adaptation (LoRA) has emerged as a prominent technique for fine-tuning large foundation models, facilitating a reduction in trainable parameters through the utilization of low-rank matrices to represent weight changes  $\mathbf{A}$  and  $\mathbf{B}$  (*i.e.*,  $\Delta\mathbf{W} = \mathbf{BA}$ ). Although LoRA has demonstrated considerable success, its expressiveness is inherently limited by the constrained capacity of its low-rank structure. To ameliorate this limitation, we introduce Fourier-based Flexible Rank Adaptation (FoRA), which harnesses the robust expressiveness of the Fourier basis to re-parameterize  $\mathbf{A}$  and  $\mathbf{B}$  from a sparse spectral subspace. Utilizing FoRA, adaptation matrices can overcome conventional rank limitations, achieving up to a 15x reduction in the parameter budget. We illustrate that FoRA achieves an optimal balance of efficiency and performance across various tasks, including natural language understanding, mathematical reasoning, commonsense reasoning, and image classification. Our codes are available at <https://anonymous.4open.science/r/FoRA-0E9C>.

## 1 INTRODUCTION

In recent years, Large Foundation Models (LFMs), have showcased exceptional generalization capabilities, greatly improving performance in a wide array of tasks across natural language processing (NLP) (Brown et al., 2020; Touvron et al., 2023a), computer vision (CV) (Radford et al., 2021; Kirillov et al., 2023), and other fields (Azad et al., 2023). Typically, adapting these general models for specific downstream tasks requires full fine-tuning, which involves retraining all model parameters and can pose significant challenges, particularly in resource-limited environments. To address this issue, Parameter-efficient fine-tuning (PEFT) techniques (Mangrulkar et al., 2022), have been developed, offering more feasible alternatives. Among these, Low-Rank Adaptation (LoRA) (Hu et al., 2021), which decomposes the weight changes into the product of two low-rank matrices  $\mathbf{A}$  and  $\mathbf{B}$ , has stood out for its effectiveness and simplicity.

Despite its success, LoRA’s reliance on low-rank structures can limit its expressive potential. Theoretically, the expressive capacity of LoRA is constrained by the ranks of  $\mathbf{A}$  and  $\mathbf{B}$  (Zeng & Lee, 2023). Therefore, more complex downstream tasks inherently necessitate higher ranks (Hu et al., 2023; Biderman et al., 2024; Gao et al., 2024). To elucidate the significance of rank configurations in practical applications, we delve into their effect on LoRA’s performance across various tasks and present the corresponding observations in Figure 1. As depicted, while different tasks exhibit varying sensitivities to rank, most demonstrate improved performance as the rank increases, with performance peaking at higher ranks (*i.e.*, no less than  $2^6$ ). This pattern aligns with the behavior of LoRA when applied to the LLaMA family, where high ranks yield clear improvements (Biderman et al., 2024). However, adapting LoRA to higher ranks inevitably engenders larger trainable parameter sizes, thereby imposing considerable overhead. Hence, a question is naturally raised:

*How can we unleash the rank-bounded potential of LoRA while still residing in the low-parameter jail?*

This question aligns closely with the principles of sparse learning (Han et al., 2015a), which aim to preserve expressive information while necessitating fewer learnable parameters. Despite the success of its predominant pruning techniques (Han et al., 2015b; Frankle & Carbin, 2018), determining which modules to prune often requires complex strategies (Zhang et al., 2022). In contrast, classical

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

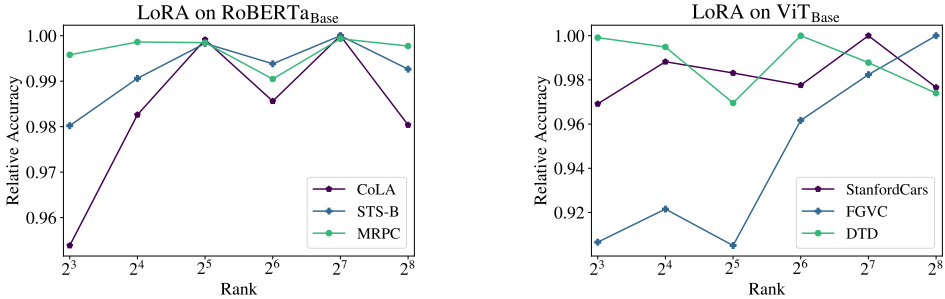


Figure 1: LoRA applied to RoBERTa<sub>BASE</sub> and ViT<sub>BASE</sub> under varying rank configurations. All experiments followed a comprehensive hyperparameter search. The reported relative accuracy, averaged over five distinct random seeds, reflects performance compared to the best results of each task. Notably, performance peaks at higher rank configurations across all tasks.

data compression techniques, such as linear projection (Dony & Haykin, 1995), fractal compression (Cochran et al., 1996), and spectral transformations (Reddy & Murthy, 1986), can be applied directly to weight matrices, providing a simpler yet effective alternative. Among these, the Fourier basis, which enables high-quality data recovery from sparse spectral information (Rudelson & Vershynin, 2006; Duarte & Baraniuk, 2013; Vlaardingerbroek & Boer, 2013), stands out as a promising tool for sparse learning. We refer our readers to Section 4.5 for a more in-depth empirical analysis.

To this end, we propose Fourier-based Flexible Rank Adaptation (FoRA), which leverages Fourier bases to re-parameterize adaptation matrices  $\mathbf{A}$  and  $\mathbf{B}$  as the spatial equivalents of sparse spectral components. Specifically, FoRA learns only  $n$  spectral components at the predefined spectral locations, which are shared among all adaptation matrices. Then, inverse Fast Fourier Transform is applied to derive these adaptation matrices in the spatial space. It is important to note that the use of a fixed quantity of spectral components enables FoRA to facilitate the adjustment of  $\mathbf{A}$  and  $\mathbf{B}$  from lower to potentially unbounded ranks at fixed parameter cost, thus ensuring significant expressiveness within a constrained parameter scope. In summary, our contributions are as follows:

- Given the rank-dependent performance of LoRA, we introduce FoRA, a novel PEFT method that enhances LoRA with Fourier-based compression, maximizing its potential while minimizing the parameter overhead.
- FoRA consistently yields comparable or better performance with up to 15x fewer trainable parameters than LoRA on various tasks, from language to vision domains and across backbones in different scales, including RoBERTa, ViT and LLaMA.
- A thorough analysis is conducted to further substantiate FoRA as a parameter-efficient alternative that replicates LoRA’s potential across different configurations.

## 2 RELATED WORKS

### 2.1 PARAMETER-EFFICIENT FINE-TUNING

Fine-tuning large pre-trained language models is crucial for improving NLP tasks. However, updating all model parameters is computationally intensive and storage-demanding for models like GPT-3 (Brown et al., 2020) and LLaMA (Touvron et al., 2023a). Parameter-efficient fine-tuning (PEFT) methods address these issues by updating fewer parameters or adding lightweight modules.

One prominent approach in PEFT is the use of adapters — small bottleneck layers inserted within each layer of a pre-trained model (Houlsby et al., 2019; Pfeiffer et al., 2020; Karimi Mahabadi et al., 2021; He et al., 2021). Houlsby et al. (2019) introduced adapters that enable task-specific adaptation while keeping the original model weights fixed. Building upon this, Pfeiffer et al. (2020) proposed a modular adapter framework that facilitates multi-task transfer. To further optimize parameter efficiency, Karimi Mahabadi et al. (2021) reduced the number of parameters by employing parameter

108 sharing and low-rank approximations within adapters. Another line of research involves prompt  
 109 tuning, which modifies the input embeddings to guide the model toward specific tasks (Lester et al.,  
 110 2021; Liu et al., 2021; Li & Liang, 2021; Chen et al., 2023a). Lester et al. (2021) optimized con-  
 111 tinuous prompt embeddings while keeping the language model’s parameters fixed, demonstrating  
 112 the effectiveness of prompt tuning for task adaptation. Similarly, Prefix-Tuning (Li & Liang, 2021)  
 113 prepends trainable vectors to the input of each transformer layer without altering the model archi-  
 114 tecture, effectively steering the model toward desired behaviors with minimal parameter updates.

115 While these methods exhibit high efficiency and preserve the originality of the pre-trained model,  
 116 they inevitably introduce higher inference costs due to additional modules or modifications required  
 117 during deployment. In contrast, LoRA (Hu et al., 2021) and its variants (Zhang et al., 2023a; Bałazy  
 118 et al., 2024; Li et al., 2024; Liu et al., 2024; Nikdan et al., 2024; Gao et al., 2024) inject trainable  
 119 low-rank matrix decomposition into transformer layers. This approach not only reduces the number  
 120 of trainable parameters but also allows for merging these decompositions with the original model  
 121 weights, thereby avoiding increased inference burdens. However, the expressiveness of low-rank  
 122 adaptation methods like LoRA is often bounded by the chosen rank (Zeng & Lee, 2023). To ad-  
 123 dress this limitation, Kopiczko et al. (2023) and Jiang et al. (2024) explored high-rank adaptations  
 124 through projection matrices, aiming to enhance expressive capacity without significantly increasing  
 125 parameter counts. Despite these advances, our empirical experiments indicate that while LoRA’s  
 126 performance may peak at certain high-rank configurations, increasing the rank beyond this point  
 127 does not necessarily lead to better results.

## 128 2.2 SPARSE LEARNING

129 Sparse neural networks exploit the fact that many weights in over-parameterized models can be  
 130 pruned with minimal impact on performance (Han et al., 2015b; Lee et al., 2018; Frankle & Carbin,  
 131 2018; Wang et al., 2020; Liu et al., 2022; Frantar & Alistarh, 2023). Techniques such as magnitude  
 132 pruning (Han et al., 2015a) remove weights with small magnitudes, effectively reducing model size.  
 133 Dynamic sparsity methods (Mocanu et al., 2018; Zhang et al., 2022; Chen et al., 2023b) adjust the  
 134 sparsity patterns during training, allowing the network to discover efficient architectures on the fly.  
 135

136 Another innovative approach is learning in transformed domains like the sparse Fourier space. By  
 137 representing weight matrices in the frequency domain using Fourier transforms, neural networks can  
 138 exploit the sparsity inherent in the frequency representation of the data (Yang & Xie, 2016; Chen  
 139 et al., 2016). This allows for efficient compression by retaining only the significant frequency com-  
 140 ponents and discarding the less important ones. Learning in the sparse Fourier space facilitates the  
 141 development of compact models that effectively capture essential patterns with fewer parameters.  
 142

## 143 3 METHODOLOGY

### 144 3.1 BACKGROUND

145 **Low-Rank Adaptation (LoRA)** LoRA (Hu et al., 2021) proposes to use the product of two low-  
 146 rank matrices to update the pre-trained weights incrementally. Let  $\mathbf{W}' \in \mathbb{R}^{d_1 \times d_2}$  denote the fine-  
 147 tuned weight,  $\mathbf{W}_0 \in \mathbb{R}^{d_1 \times d_2}$  the pre-trained weight, and  $\Delta\mathbf{W} \in \mathbb{R}^{d_1 \times d_2}$  the change in weight.  
 148 LoRA models this change  $\Delta\mathbf{W}$  through a low-rank decomposition:  
 149

$$150 \mathbf{W}' = \mathbf{W}_0 + \Delta\mathbf{W} = \mathbf{W}_0 + \mathbf{B}\mathbf{A}, \tag{1}$$

151 where  $\mathbf{W}_0$  is kept unchanged during fine-tuning. The matrices  $\mathbf{A} \in \mathbb{R}^{r \times d_2}$  and  $\mathbf{B} \in \mathbb{R}^{d_1 \times r}$  represents  
 152 the learnable low-rank matrices with the rank  $r \ll \{d_1, d_2\}$ . Typically,  $\mathbf{A}$  adopts Kaiming uniform  
 153 initialization (He et al., 2015) while  $\mathbf{B}$  is initialized to zero at the start of the training process.  
 154

155 In the following parts, we present Fourier-based Flexible Rank Adaptation (FoRA), which re-  
 156 parameterizes adaptation matrices of LoRA by applying inverse Fast Fourier Transform on sparse  
 157 spectral coefficients. The overall framework is presented in Figure 2.  
 158

### 159 3.2 FOURIER-BASED FLEXIBLE RANK ADAPTATION

160 As stated previously, our goal is to re-parameterize  $\mathbf{A}$  and  $\mathbf{B}$  with fewer parameters while maintain-  
 161 ing strong expressiveness, which aligns closely with the foundational principle of sparse learning.

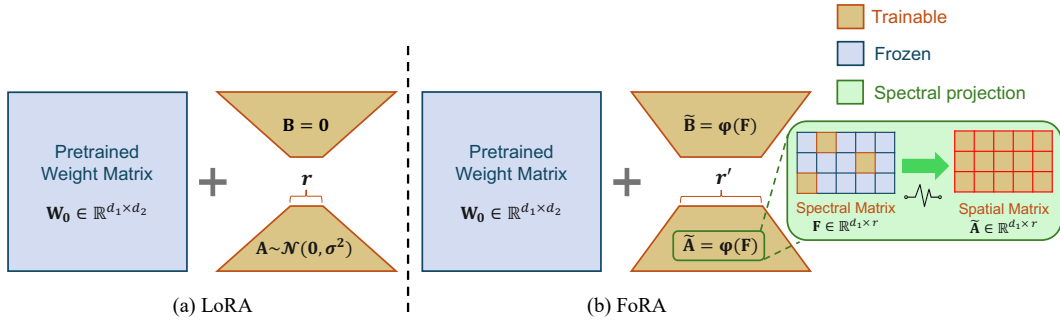


Figure 2: An overview of the schematic comparison between LoRA and our proposed FoRA. While LoRA necessitates training all elements in the low-rank matrices  $\mathbf{A}$  and  $\mathbf{B}$ , FoRA re-parameterizes these matrices from a sparse spectral subspace (highlighted in green). Our approach enables flexible rank adjustment while training fixed and sparse components. In both cases, low-rank matrices can be merged into the original weights matrix  $\mathbf{W}_0$ , ensuring no additional latency is introduced.

Upon revisiting prior successes, we resort to the Fourier basis, known for its robust expressiveness (Candès et al., 2006; Baraniuk, 2007).

Essentially, our approach centers on re-parameterizing the adaptation matrices, termed  $\tilde{\mathbf{A}} \in \mathbb{R}^{r \times d_2}$  and  $\tilde{\mathbf{B}} \in \mathbb{R}^{d_1 \times r}$ , as the spatial recovery of sparse spectral coefficients, while retaining LoRA’s update schema:

$$\mathbf{W}' = \mathbf{W}_0 + \Delta\mathbf{W} = \mathbf{W}_0 + \tilde{\mathbf{B}}\tilde{\mathbf{A}}. \quad (2)$$

To accomplish this, we start by randomly initializing a 2D index matrix  $\mathbf{L} \in \mathbb{R}^{2 \times n}$  to specify spectral locations for all low-rank matrices. To derive  $\tilde{\mathbf{A}}$ , we then define  $n$  learnable spectral coefficient  $\mathbf{s} \in \mathbb{R}^n$ . Using these indices and coefficients, we construct the sparse spectral matrix  $\mathbf{F} \in \mathbb{R}^{r_1 \times d}$  and compute its spatial counterpart  $\mathbf{S} \in \mathbb{R}^{r_1 \times d}$  via the inverse Fast Fourier Transform:

$$\mathbf{S}_{p,q} = \frac{1}{rd_2} \sum_{j=0}^{r-1} \sum_{k=0}^{d_2-1} \mathbf{F}_{j,k} e^{i2\pi(\frac{j}{r}p + \frac{k}{d_2}q)}, \quad (3)$$

where  $i$  denotes the imaginary unit. In particular,  $\mathbf{F}_{j,k} = \mathbf{s}_p$  if  $(j, k) = \mathbf{L}_{:,p}$  and  $\mathbf{F}_{j,k} = 0$  otherwise. The Fourier-based re-parameterized matrix  $\tilde{\mathbf{A}}$  is then defined as the real part of the complex matrix  $\mathbf{S}$  as

$$\tilde{\mathbf{A}} = \text{Re}[\mathbf{S}]. \quad (4)$$

The adaptation matrix  $\tilde{\mathbf{B}}$  is obtained by applying the identical procedure as above.

In this setup, FoRA can be easily integrated as a plug-in by replacing the LoRA linear module with the FoRA linear module in a single line of code, requiring no additional modifications, as outlined in Algorithm 1 in the Appendix. Moreover, despite learning only a fixed number of spectral components, the high expressiveness of the Fourier basis allows FoRA to represent informative matrices with ranks that range from low to very high values. This flexibility enables FoRA to replicate LoRA’s potential, even within a constrained parameter space.

### 3.3 DISCUSSION

**Initialization strategies.** Matrix initialization with consistent variance (Glorot & Bengio, 2010) is crucial for maintaining numerical stability and accelerating convergence. However, unlike LoRA, directly initializing the spectral space in FoRA can lead to suboptimal variance in spatial space due to the involvement of the Fourier transform. To facilitate efficient training, for matrix  $\tilde{\mathbf{A}}$ , we first employ Xavier (Glorot & Bengio, 2010) or Kaiming initialization (He et al., 2015) to its spectral coefficients  $\mathbf{s}$  and a spatial auxiliary matrix  $\mathbf{A}' \in \mathbb{R}^{r \times d_2}$ . Next, we scale  $\mathbf{s}$  by  $\text{Var}(\mathbf{A}')/\text{Var}(\tilde{\mathbf{A}})$  to approximate consistent variance. In contrast, matrix  $\tilde{\mathbf{B}}$  is initialized to zeros following the standard practice of LoRA (Hu et al., 2021). We employ Kaiming initialization by default unless specially stated.

**Comparison to LoRA’s variants.** Recent parameter-efficient variants of LoRA (Kopiczko et al., 2023; Renduchintala et al., 2023; Li et al., 2024) have demonstrated competitive performance by adapting at higher ranks through the use of simple linear projections. However, their strategies for sparse learning, which essentially involve a collection of learnable scaling transformations, suffer from limited expressiveness. To remedy this issue, FoRA leverages the more efficient and expressive Fourier transform for matrix re-parameterization, striking a balance between performance and efficiency. Compared with them, FoRA consistently provides enhanced representational expressiveness while allowing flexible rank adaptation with fixed cost. Further details of the empirical analysis are provided in Section 4.5.

## 4 EXPERIMENTS

In this section, we present a series of experiments to demonstrate the effectiveness of FoRA across diverse tasks, including language and image domains. We begin by evaluating FoRA through fine-tuning RoBERTa on the GLUE benchmark. Next, we focus on instruction tuning within the LLaMA family. Following this, we assess FoRA’s performance by fine-tuning Vision Transformers for image classification. Finally, we provide an in-depth analysis of FoRA’s capabilities.

**Baselines.** We evaluate FoRA against three groups of baselines. The first group follows the classical fine-tuning paradigm, which includes **Full Fine-tuning (FF)** and **BitFit** (Zaken et al., 2021) where only bias vectors are fine-tuned. The second group is adapter-tuning, covering **Adpt<sup>H</sup>** (Houlsby et al., 2019), **Adpt<sup>P</sup>** (Pfeiffer et al., 2020), **Adpt<sup>R</sup>** (He et al., 2021). The third group is the most prevalent low-rank adaptation and its variants, including **LoRA** (Hu et al., 2021), **VeRA** (Kopiczko et al., 2023), **FourierFT** (Gao et al., 2024), **DoRA** (Liu et al., 2024).

### 4.1 GLUE BENCHMARK

We evaluate FoRA on the General Language Understanding Evaluation (GLUE) benchmark (Wang, 2018), a sequence classification benchmark for natural language understanding (NLU) which covers domains such as sentiment classification and natural language inference. We employ the pre-trained RoBERTa<sub>BASE</sub> and RoBERTa<sub>LARGE</sub> (Liu, 2019) as the foundation models for fine-tuning.

Our experimental setup closely follows (Hu et al., 2021), involving fine-tuning only the query and value weights in each transformer block and fully fine-tuning the classification head. For our method, we randomly sample  $n = \{250, 500\}$  trainable spectral coefficients per low-rank matrix, which we denote as FoRA<sup>†</sup> and FoRA, respectively. We adopt the baseline hyperparameters from their original papers. For our approaches, we apply random search (Bergstra et al., 2013) to optimize the learning rates and matrix rank. For comprehensiveness, we report the median performance across 5 random seed trials, selecting the best epoch for each run. Additionally, we report the number of trainable parameters in the fine-tuned layers, excluding the classification head, as suggested by (Hu et al., 2021; Kopiczko et al., 2023). Further specifics are provided in Table 6 in the Appendix.

**Results.** As highlighted in Table 1, FoRA generally delivers better or on-par performance compared with baseline methods, while adapting at higher ranks with extremely lower budget. Notably, under the same parameter constraints, FoRA demonstrates improved performance over FourierFT. The performance gains are more pronounced with the RoBERTa<sub>LARGE</sub> model. Specifically, FoRA<sup>†</sup> not only surpasses adapter tuning by a clear margin but also matches the performance of LoRA, despite requiring 30 times fewer trainable parameters. These results demonstrate that FoRA strikes an effective balance between unleashing LoRA’s rank-bounded potential and parameter efficiency.

### 4.2 MATHEMATICAL REASONING

Instruction tuning involves fine-tuning a language model on a collection of paired prompts and responses (Ouyang et al., 2022). To evaluate the effectiveness of FoRA, we first apply it to LLaMA2<sub>7B/13B</sub> (Touvron et al., 2023b) and LLaMA3<sub>8B</sub> (Dubey et al., 2024) for mathematical reasoning tasks.

This evaluation uses two challenging benchmarks: GSM8K (Cobbe et al., 2021) and MATH (Hendrycks et al., 2020). Both datasets consist of multi-step problems that require chain-



Table 1: Fine-tuning performance of the pre-trained RoBERTa<sub>BASE</sub> and RoBERTa<sub>LARGE</sub> models with different methods on the GLUE benchmark. We report Matthew’s correlation coefficient for CoLA, Pearson correlation coefficient for STS-B, and accuracy for all the remaining tasks. The best results for each dataset are highlighted in **bold**. FoRA<sup>†</sup>: the lightweight version of FoRA.

	Methods	# Trainable Parameters	SST-2	MRPC	CoLA	QNLI	RTE	STS-B	Avg.
BASE	FF	125M	94.8	90.2	63.6	92.8	78.7	91.2	85.2
	BitFit	0.1M	93.7	<b>92.7</b>	62.0	91.8	<b>81.5</b>	90.8	<b>85.4</b>
	LoRA	0.3M	<b>95.1</b> <sub>±0.2</sub>	89.7 <sub>±0.7</sub>	63.4 <sub>±1.2</sub>	<b>93.3</b> <sub>±0.3</sub>	78.8 <sub>±0.5</sub>	<b>91.5</b> <sub>±0.2</sub>	85.3
	VeRA	0.043M	94.6 <sub>±0.1</sub>	89.5 <sub>±0.5</sub>	<b>65.6</b> <sub>±0.8</sub>	91.8 <sub>±0.2</sub>	78.7 <sub>±0.7</sub>	90.7 <sub>±0.2</sub>	85.2
	FourierFT	0.024M	94.2 <sub>±0.3</sub>	90.0 <sub>±0.8</sub>	63.8 <sub>±1.6</sub>	92.2 <sub>±0.1</sub>	79.1 <sub>±0.5</sub>	90.8 <sub>±0.2</sub>	85.0
	<b>FoRA<sup>†</sup></b>	<b>0.012M</b>	94.3 <sub>±0.3</sub>	89.7 <sub>±0.2</sub>	62.6 <sub>±1.6</sub>	92.4 <sub>±0.4</sub>	78.7 <sub>±2.6</sub>	90.0 <sub>±0.3</sub>	84.6
	<b>FoRA</b>	<b>0.024M</b>	94.7 <sub>±0.3</sub>	90.4 <sub>±0.5</sub>	64.6 <sub>±1.0</sub>	92.3 <sub>±0.1</sub>	79.4 <sub>±1.9</sub>	90.7 <sub>±0.2</sub>	<b>85.4</b>
LARGE	Adpt <sup>P</sup>	0.8M	<b>96.6</b> <sub>±0.2</sub>	89.7 <sub>±1.2</sub>	67.8 <sub>±2.5</sub>	<b>94.8</b> <sub>±0.3</sub>	80.1 <sub>±2.9</sub>	91.9 <sub>±0.4</sub>	86.8
	Adpt <sup>H</sup>	0.8M	96.3 <sub>±0.5</sub>	87.7 <sub>±1.7</sub>	66.3 <sub>±2.0</sub>	94.7 <sub>±0.2</sub>	72.9 <sub>±2.9</sub>	91.5 <sub>±0.5</sub>	84.9
	LoRA	0.8M	96.2 <sub>±0.5</sub>	90.2 <sub>±1.0</sub>	<b>68.2</b> <sub>±1.9</sub>	<b>94.8</b> <sub>±0.3</sub>	85.2 <sub>±1.1</sub>	<b>92.3</b> <sub>±0.5</sub>	87.8
	VeRA	0.061M	96.1 <sub>±0.1</sub>	90.9 <sub>±0.7</sub>	68.0 <sub>±0.8</sub>	94.4 <sub>±0.2</sub>	85.9 <sub>±0.7</sub>	91.7 <sub>±0.8</sub>	87.8
	FourierFT	0.048M	96.0 <sub>±0.2</sub>	90.9 <sub>±0.3</sub>	67.1 <sub>±1.4</sub>	94.4 <sub>±0.4</sub>	<b>87.4</b> <sub>±1.6</sub>	91.9 <sub>±0.4</sub>	88.0
	<b>FoRA<sup>†</sup></b>	<b>0.024M</b>	96.1 <sub>±0.2</sub>	91.2 <sub>±1.0</sub>	66.5 <sub>±0.9</sub>	94.2 <sub>±0.5</sub>	86.6 <sub>±1.1</sub>	91.4 <sub>±0.2</sub>	87.7
	<b>FoRA</b>	<b>0.048M</b>	96.3 <sub>±0.1</sub>	<b>91.4</b> <sub>±1.0</sub>	68.0 <sub>±2.0</sub>	94.4 <sub>±0.3</sub>	87.0 <sub>±2.0</sub>	91.9 <sub>±0.4</sub>	<b>88.2</b>

of-thought reasoning (Wei et al., 2022) to reach the final answer, and they are framed as question-answering tasks using the same prompt template as in (Zhang et al., 2023b). Each method is fine-tuned on the respective training sets and evaluated on the testing sets, where we only evaluate the correctness of the final numeric answer.

In addition, FoRA only re-parameterizes the adaptation matrix with Fourier transform, thus allowing it to be adapted to other LoRA variants. To test the adaptability, we select DoRA, where the directional component of the decomposed weight is learnable, and apply FoRA to the directional update, resulting in a combination called DFoRA. We use  $n = 30000$  learnable spectral coefficients for LLaMA2<sub>13B</sub> and  $n = 20000$  for the rest. To ensure a fair comparison, we fine-tuned the models following the setup suggested in (Hu et al., 2023; Liu et al., 2024), keeping the baseline models at a fixed rank of  $r = 32$  while experimenting with different learning rates. In contrast, for our approaches, we optimize both the learning rates and matrix ranks. For comprehensiveness, we consider two scenarios: (1) a standard single training pass and (2) extended training over three epochs, reporting the best results for each (Nikdan et al., 2024). A more detailed setup is provided in Table 7 in the Appendix.

Table 2: Comparison of LLaMA2<sub>7B</sub>, LLaMA2<sub>13B</sub> and LLaMA3<sub>8B</sub> fine-tuned on mathematical benchmark datasets. Avg. denotes the average accuracy. The best results for each dataset are highlighted in **bold**.

	Methods	# Parameters	GSM8K MATH Avg.			GSM8K MATH Avg.		
			1 Epoch			Extended		
LLaMA2 <sub>7B</sub>	LoRA	16.8M	27.07	4.35	15.71	<b>38.53</b>	5.70	<b>22.12</b>
	DoRA	17.0M	<b>28.20</b>	<b>4.55</b>	<b>16.38</b>	38.06	<b>6.05</b>	22.06
	<b>FoRA</b>	<b>2.56M</b>	26.99	4.15	15.57	37.63	5.70	21.67
	<b>DFoRA</b>	<b>2.82M</b>	27.77	4.30	16.04	37.76	5.90	21.83
LLaMA2 <sub>13B</sub>	LoRA	26.2M	38.51	5.30	21.90	49.20	8.45	28.83
	DoRA	26.6M	38.82	5.85	22.34	50.34	<b>9.00</b>	29.67
	<b>FoRA</b>	<b>4.80M</b>	37.54	<b>6.20</b>	21.87	48.98	8.65	28.81
	<b>DFoRA</b>	<b>5.21M</b>	<b>39.58</b>	5.55	<b>22.56</b>	<b>50.49</b>	8.90	<b>29.70</b>
LLaMA3 <sub>8B</sub>	LoRA	13.6M	53.16	18.95	36.06	62.45	21.25	41.85
	DoRA	13.8M	54.28	<b>20.55</b>	37.42	62.55	22.20	42.38
	<b>FoRA</b>	<b>2.56M</b>	54.13	19.55	36.84	<b>63.00</b>	21.35	42.18
	<b>DFoRA</b>	<b>2.72M</b>	<b>55.65</b>	19.40	<b>37.53</b>	62.77	<b>22.45</b>	<b>42.61</b>

**Results.** The results in Table 2 show that FoRA and DFoRA achieve accuracy that closely matches or slightly surpasses baseline methods, even with over 5 times fewer trainable parameters, in both single-pass and extended training scenarios. Notably, DFoRA shows significant improvements over FoRA, highlighting the flexible adaptability of the FoRA framework. Our approaches are particularly effective with the more advanced LLaMA3<sub>8B</sub> model, indicating that FoRA is especially well-suited to the sophisticated post-training techniques used in the latest LLaMA family. Overall, these empirical observations underscore the effectiveness and strong compatibility of FoRA.

### 4.3 COMMONSENSE REASONING

For a comprehensive evaluation of instruction tuning, we further compare our methods with LoRA and DoRA on LLaMA<sub>7B/13B</sub> (Touvron et al., 2023a), LLaMA2<sub>7B</sub> (Touvron et al., 2023b), and LLaMA3<sub>8B</sub> (Dubey et al., 2024) for commonsense reasoning tasks.

These commonsense reasoning tasks are framed as multiple-choice questions across eight distinct datasets, including BoolQ (Clark et al., 2019), PIQA (Bisk et al., 2020), SIQA (Sap et al., 2019), HellaSwag (Zellers et al., 2019), WinoGrande (Sakaguchi et al., 2021), ARC-e, ARC-c (Clark et al., 2018), and OBQA (Mihaylov et al., 2018). Consistent with the approach in (Hu et al., 2023), we use the Commonsense170K dataset for training, which integrates the training sets of all eight datasets, while evaluations are conducted on the test sets of the individual datasets.

In our experiments, we set rank  $r = 32$  for all models as suggested by (Liu et al., 2024). Given the complexity of the tasks, we use  $n = 40000$  learnable spectral coefficients for LLaMA<sub>13B</sub> and  $n = 30000$  for the rest. A detailed configuration setup is provided in Table 8 in the Appendix.

Table 3: Comparison of LLaMA<sub>7B</sub>, LLaMA<sub>13B</sub>, LLaMA2<sub>7B</sub> and LLaMA3<sub>8B</sub> against various methods on eight commonsense datasets. Results of all baseline methods are taken from (Liu et al., 2024). The best and runner-up models for each dataset are highlighted in **bold** and underline.

	Methods	# Parameters	BoolQ	PIQA	SIQA	HellaS.	WinoG.	ARC-e	ARC-c	OBQA	Avg.
ChatGPT	—	—	73.1	85.4	68.5	78.5	66.1	89.8	79.9	74.8	77.0
LLaMA <sub>7B</sub>	Adpt <sup>H</sup>	132M	63.0	79.2	76.3	67.9	75.7	74.5	57.1	72.4	70.8
	Adpt <sup>R</sup>	239M	67.9	76.4	<b>78.8</b>	69.8	78.9	73.7	57.3	75.2	72.2
	LoRA	55.7M	<b>68.9</b>	<u>80.7</u>	77.4	78.1	78.8	77.8	61.3	74.8	74.7
	DoRA	56.5M	68.0	<u>80.6</u>	77.9	<b>83.9</b>	<b>80.8</b>	<b>81.4</b>	<b>63.4</b>	<u>77.6</u>	<b>76.7</b>
	<b>FoRA</b>	<b>9.60M</b>	67.8	80.1	77.5	76.6	<u>79.8</u>	<u>80.3</u>	62.8	75.2	75.0
	<b>DFoRA</b>	<b>10.5M</b>	<u>68.8</u>	<b>81.2</b>	<u>78.0</u>	<u>81.3</u>	79.2	78.9	<u>63.1</u>	<b>79.6</b>	<u>76.3</u>
LLaMA <sub>13B</sub>	Adpt <sup>H</sup>	206M	71.8	83.0	79.2	88.1	82.4	82.5	67.3	81.8	79.5
	Adpt <sup>R</sup>	377M	<b>72.5</b>	84.8	79.8	<u>92.1</u>	<b>84.7</b>	84.1	<b>71.2</b>	82.2	<b>81.5</b>
	LoRA	87.2M	72.1	83.5	80.5	90.5	83.7	82.8	68.3	82.4	80.5
	DoRA	88.6M	<u>72.4</u>	<b>84.9</b>	<b>81.5</b>	<b>92.4</b>	<u>84.2</u>	84.2	69.6	<b>82.8</b>	<b>81.5</b>
	<b>FoRA</b>	<b>16.0M</b>	72.0	84.5	80.0	91.5	83.8	83.6	<u>70.8</u>	<u>82.0</u>	81.0
	<b>DFoRA</b>	<b>17.4M</b>	71.8	84.4	<u>81.0</u>	91.8	<u>84.5</u>	<b>84.4</b>	70.1	81.8	<u>81.2</u>
LLaMA2 <sub>7B</sub>	LoRA	55.7M	69.8	79.9	<u>79.5</u>	83.6	<u>82.5</u>	79.8	64.7	<u>81.0</u>	77.6
	DoRA	56.6M	<b>71.8</b>	<b>83.7</b>	76.0	<u>89.1</u>	<b>82.6</b>	<b>83.7</b>	<u>68.2</u>	<b>82.4</b>	<b>79.7</b>
	<b>FoRA</b>	<b>9.60M</b>	<u>71.6</u>	81.5	<b>80.0</b>	<b>90.5</b>	81.9	<u>83.6</u>	68.0	80.0	<u>79.6</u>
	<b>DFoRA</b>	<b>10.5M</b>	71.5	<u>82.4</u>	<u>79.5</u>	88.2	<b>82.6</b>	83.5	<b>68.5</b>	<u>81.0</u>	<b>79.7</b>
LLaMA3 <sub>8B</sub>	LoRA	56.2M	70.8	85.2	79.9	91.7	84.3	84.2	71.2	79.0	80.8
	DoRA	57.0M	<b>74.6</b>	<b>89.3</b>	79.9	<b>95.5</b>	85.6	<u>90.5</u>	<b>80.4</b>	<u>85.8</u>	<u>85.2</u>
	<b>FoRA</b>	<b>9.60M</b>	74.0	88.7	<u>80.0</u>	<u>95.2</u>	<b>86.2</b>	90.4	77.8	85.0	84.7
	<b>DFoRA</b>	<b>10.4M</b>	<u>74.5</u>	<u>89.1</u>	<b>80.4</b>	95.1	<u>85.8</u>	<b>90.6</b>	<u>79.7</u>	<b>86.8</b>	<b>85.3</b>

**Results.** Table 3 presents an overview of general performance across different backbone models. Our findings indicate that FoRA consistently outperforms LoRA at the same rank while requiring less than 1/5 parameter count. Furthermore, despite the greater complexity of generalized reasoning tasks, DFoRA either closely matches or even exceeds the performance of DoRA on more advanced LLaMA models, mirroring trends observed in mathematical reasoning. Overall, there is significant variability in the results for commonsense reasoning, with no single method emerging as a definitive leader across all datasets.

#### 4.4 IMAGE CLASSIFICATION

This section concentrates on image classification to evaluate whether FoRA can remain competitive. We adopt Vision Transformer (ViT) (Dosovitskiy et al., 2020), which is pre-trained on the vast ImageNet-21K dataset (Ridnik et al., 2021), as the foundation model. Specifically, we fine-tune ViT<sub>BASE</sub> and ViT<sub>LARGE</sub> on a variety of datasets, including OxfordPets (Parkhi et al., 2012), StanfordCars (Krause et al., 2013), DTD (Cimpoi et al., 2014), EuroSAT (Helber et al., 2019), FGVC (Maji et al., 2013), and RESISC45 (Cheng et al., 2017). Notably, RESISC45 and EuroSAT offer rich labeled data, while the other datasets serve as few-shot adaptations with sparse training samples.

We follow the same fine-tuning protocols as in the GLUE benchmark, reporting the number of trainable parameters excluding the classification head. For LoRA, we set the rank to  $r = 16$ . To maintain the same parameter constraints, we use  $n = 16000$  learnable spectral entries for FourierFT and  $n = 8000$  for FoRA. Learning rates are tuned over a maximum of 10 training epochs, and we report average results across 5 random trials. Detailed hyperparameters are provided in Table 9 in the Appendix.

Table 4: Fine-tuning results with ViT<sub>BASE</sub> and ViT<sub>LARGE</sub> models on different image classification datasets. Linear Probing (LP) represents fine-tuning only the classification head. Results are averaged across 5 runs with different random seeds. The best performance is shown in **bold**.

	Methods	# Trainable Parameters	OxfordPets	StanfordCars	DTD	EuroSAT	FGVC	RESISC45	Avg.
BASE	LP	-	90.28 $\pm$ 0.43	25.76 $\pm$ 0.28	69.77 $\pm$ 0.67	88.72 $\pm$ 0.13	17.44 $\pm$ 0.43	74.22 $\pm$ 0.10	61.03
	FF	85.8M	92.82 $\pm$ 0.54	<b>85.10</b> $\pm$ 0.21	80.11 $\pm$ 0.56	<b>99.11</b> $\pm$ 0.07	<b>61.60</b> $\pm$ 1.00	<b>96.00</b> $\pm$ 0.23	<b>85.79</b>
	LoRA	0.59M	93.76 $\pm$ 0.44	78.04 $\pm$ 0.33	78.56 $\pm$ 0.62	98.84 $\pm$ 0.08	56.64 $\pm$ 0.55	94.66 $\pm$ 0.17	83.42
	FourierFT	0.384M	93.37 $\pm$ 0.30	81.22 $\pm$ 0.48	78.90 $\pm$ 0.75	98.92 $\pm$ 0.09	58.82 $\pm$ 0.37	94.91 $\pm$ 0.24	84.36
	<b>FoRA</b>	<b>0.384M</b>	<b>94.05</b> $\pm$ 0.37	81.46 $\pm$ 0.78	<b>80.34</b> $\pm$ 1.03	98.85 $\pm$ 0.10	58.67 $\pm$ 0.37	94.89 $\pm$ 0.15	84.71
LARGE	LP	-	91.11 $\pm$ 0.30	37.91 $\pm$ 0.27	73.33 $\pm$ 0.26	92.64 $\pm$ 0.08	24.62 $\pm$ 0.24	82.02 $\pm$ 0.11	66.94
	FF	303M	94.30 $\pm$ 0.31	<b>88.15</b> $\pm$ 0.50	80.18 $\pm$ 0.66	<b>99.06</b> $\pm$ 0.10	<b>67.38</b> $\pm$ 1.06	<b>96.08</b> $\pm$ 0.20	<b>87.53</b>
	LoRA	1.57M	94.62 $\pm$ 0.47	86.11 $\pm$ 0.42	80.09 $\pm$ 0.42	98.99 $\pm$ 0.03	63.64 $\pm$ 0.83	95.94 $\pm$ 0.21	86.56
	FourierFT	0.768M	<b>94.91</b> $\pm$ 0.33	85.93 $\pm$ 0.58	81.17 $\pm$ 0.71	99.04 $\pm$ 0.07	62.48 $\pm$ 0.45	95.59 $\pm$ 0.23	86.52
	<b>FoRA</b>	<b>0.768M</b>	94.90 $\pm$ 0.20	86.23 $\pm$ 0.29	<b>81.91</b> $\pm$ 0.82	<b>99.06</b> $\pm$ 0.09	65.61 $\pm$ 0.72	95.81 $\pm$ 0.13	87.25

**Results.** Table 4 presents a comprehensive overview across 6 distinct image classification datasets using ViT<sub>BASE</sub> and ViT<sub>LARGE</sub>. FoRA consistently outperforms LoRA by a significant margin while using only half the number of trainable parameters. Additionally, FoRA demonstrates superior performance compared to FourierFT under the same parameter constraints. Notably, FoRA even achieves results on par with full fine-tuning, despite utilizing substantially fewer parameters. These findings, along with the insights from Figure 3, highlight the importance of enabling flexible rank adaptation with reduced overhead to enhance representational power.

#### 4.5 ANALYSIS

**Sparse Learning Strategy.** To explore the impact of various sparse learning strategies **applied to LoRA**, we compare FoRA with two prominent strategies, random masking (Masking) and linear projection (VeRA) (Kopiczko et al., 2023), assessing their performance compared to LoRA across different tasks and ranks. We fine-tune RoBERTa<sub>BASE</sub> and ViT<sub>BASE</sub> on three representative datasets respectively, following the setup in Section 4.1 and 4.4. To ensure fairness, the number of retained parameters for random masking matches the learnable coefficients in FoRA.

The average accuracies across different ranks are depicted in Figure 3, with the corresponding parameter counts detailed in Table 10 in the Appendix. FoRA demonstrates a performance pattern akin to LoRA, closely matching its results across various ranks, particularly at higher ranks, while maintaining a more flexible and reduced parameter count that can be adjusted based on task complexity. In contrast, random masking shows degraded performance compared to FoRA in the GLUE, likely due to the limited expressiveness of trivial masking with extremely sparse parameters. Surprisingly, despite the decent performance in GLUE, VeRA shows a notable drop in more challenging image classification tasks, even when using high-rank matrices. This drop may stem from its inflexible parameter count constrained by the size of the adaptation matrices. Overall, these findings suggest



432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485

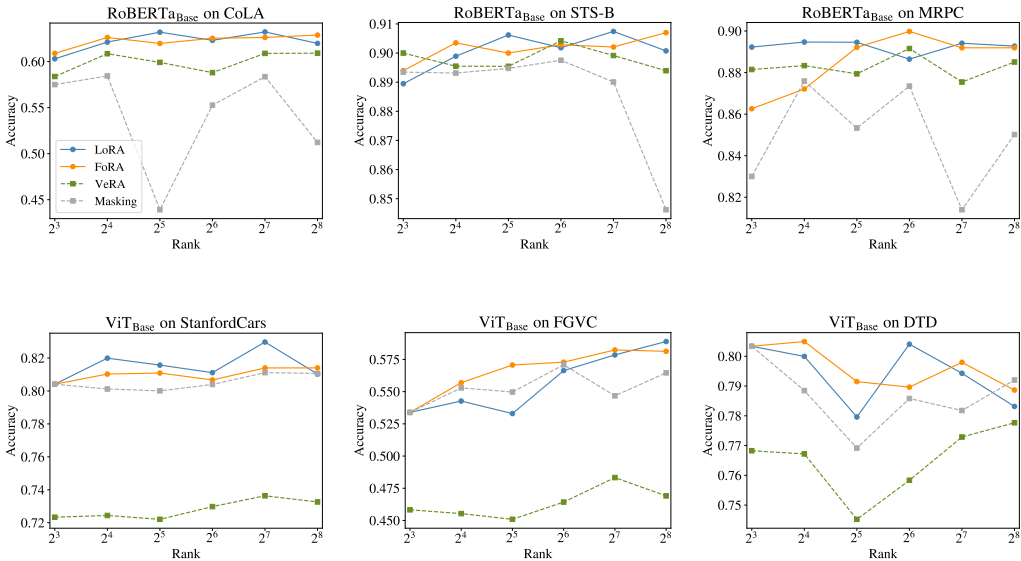


Figure 3: Performance comparison of LoRA variants with different parameter-reduction strategies applied to  $\text{RoBERTa}_{\text{BASE}}$  and  $\text{ViT}_{\text{BASE}}$  across various rank configurations. FoRA consistently matches LoRA’s performance, while other variants show varying levels of degradation.

that the stronger expressive power of the Fourier basis, combined with the flexible adjustment of trainable parameters, positions FoRA as a promising and parameter-efficient alternative to LoRA.

**Efficiency Comparison.** To assess the computational efficiency, we compare the training time and GPU overhead of FoRA against LoRA for fine-tuning LLaMA2<sub>7B</sub> on MATH and Comonsense170K, adhering to the setup in Section 4.2 and 4.3. Our evaluation covers both low-rank ( $r = 32$ ) and high-rank ( $r = 256$ ) scenarios to ensure a comprehensive comparison. As shown in Table 5, despite

Table 5: Comparison of GPU memory and training time.

Dataset	Methods	$r = 32$		$r = 256$	
		Memory	Time	Memory	Time
MATH	LoRA	34.9 GB	<b>37 min</b>	37.3 GB	<b>38 min</b>
	FoRA	<b>34.4 GB</b>	37.5 min	<b>35.3 GB</b>	38.5 min
Common	LoRA	42.4 GB	<b>442 min</b>	45.3 GB	<b>466 min</b>
	<b>FoRA</b>	<b>41.9 GB</b>	454 min	<b>43.4 GB</b>	485 min

the additional operations introduced by the Fourier transform in FoRA’s forward pass, the impact on training time remains modest, with an increase of up to only 4%, even when fine-tuning high-rank, large-scale datasets. Moreover, FoRA demonstrates improved GPU memory efficiency, particularly in high-rank scenarios, reducing memory usage by up to 5.3%. These findings highlight that FoRA also strikes a great balance between memory efficiency and training time

## 5 CONCLUSION

In this work, we aim to unlock the rank-bounded potential of LoRA while minimizing and controlling parameter overhead. We present FoRA, a fine-tuning method that re-parameterizes adaptation matrices from spectral subspace and is compatible with LoRA and its variants. With Fourier basis, FoRA allows for the representation of informative adaptation matrices from lower to potentially unbounded ranks at fixed parameter cost. Empirically, FoRA consistently matches or surpasses LoRA’s performance across various fine-tuning tasks and backbone models, requiring up to 15x fewer trainable parameters. Moreover, a comprehensive analysis further substantiates FoRA as a parameter-efficient alternative to LoRA. Our work demonstrates the potential for efficiently replicating LoRA’s capabilities, with opportunities for further exploration in future research.

## REFERENCES

- 486  
487  
488 Bobby Azad, Reza Azad, Sania Eskandari, Afshin Bozorgpour, Amirhossein Kazerouni, Islem  
489 Rekik, and Dorit Merhof. Foundational models in medical imaging: A comprehensive survey  
490 and future vision. *arXiv preprint arXiv:2310.18689*, 2023.
- 491 Klaudia Bałazy, Mohammadreza Banaei, Karl Aberer, and Jacek Tabor. Lora-xs: Low-rank adapta-  
492 tion with extremely small number of parameters. *arXiv preprint arXiv:2405.17604*, 2024.
- 493  
494 Richard G Baraniuk. Compressive sensing [lecture notes]. *IEEE signal processing magazine*, 24  
495 (4):118–121, 2007.
- 496 James Bergstra, Daniel Yamins, and David Cox. Making a science of model search: Hyperparameter  
497 optimization in hundreds of dimensions for vision architectures. In *International conference on*  
498 *machine learning*, pp. 115–123. PMLR, 2013.
- 499  
500 Dan Biderman, Jose Gonzalez Ortiz, Jacob Portes, Mansheej Paul, Philip Greengard, Connor Jen-  
501 nings, Daniel King, Sam Havens, Vitaliy Chiley, Jonathan Frankle, et al. Lora learns less and  
502 forgets less. *arXiv preprint arXiv:2405.09673*, 2024.
- 503 Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. Piqa: Reasoning about physical com-  
504 monsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*,  
505 volume 34, pp. 7432–7439, 2020.
- 506  
507 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
508 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are  
509 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 510 Emmanuel J Candès, Justin Romberg, and Terence Tao. Robust uncertainty principles: Exact signal  
511 reconstruction from highly incomplete frequency information. *IEEE Transactions on information*  
512 *theory*, 52(2):489–509, 2006.
- 513  
514 Aochuan Chen, Yuguang Yao, Pin-Yu Chen, Yihua Zhang, and Sijia Liu. Understanding and improv-  
515 ing visual prompting: A label-mapping perspective. In *Proceedings of the IEEE/CVF Conference*  
516 *on Computer Vision and Pattern Recognition*, pp. 19133–19143, 2023a.
- 517 Aochuan Chen, Yimeng Zhang, Jinghan Jia, James Diffenderfer, Jiancheng Liu, Konstantinos  
518 Parasyris, Yihua Zhang, Zheng Zhang, Bhavya Kailkhura, and Sijia Liu. Deepzero: Scaling  
519 up zeroth-order optimization for deep model training. *arXiv preprint arXiv:2310.02025*, 2023b.
- 520  
521 Wenlin Chen, James Wilson, Stephen Tyree, Kilian Q Weinberger, and Yixin Chen. Compress-  
522 ing convolutional neural networks in the frequency domain. In *Proceedings of the 22nd ACM*  
523 *SIGKDD international conference on knowledge discovery and data mining*, pp. 1475–1484,  
524 2016.
- 525 Gong Cheng, Junwei Han, and Xiaoqiang Lu. Remote sensing image scene classification: Bench-  
526 mark and state of the art. *Proceedings of the IEEE*, 105(10):1865–1883, 2017.
- 527  
528 Mircea Cimpoi, Subhansu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. De-  
529 scribing textures in the wild. In *Proceedings of the IEEE conference on computer vision and*  
530 *pattern recognition*, pp. 3606–3613, 2014.
- 531 Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina  
532 Toutanova. Boolq: Exploring the surprising difficulty of natural yes/no questions. *arXiv preprint*  
533 *arXiv:1905.10044*, 2019.
- 534  
535 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and  
536 Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge.  
537 *arXiv preprint arXiv:1803.05457*, 2018.
- 538  
539 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,  
Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to  
solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.

- 540 Wayne O Cochran, John C Hart, and Patrick J Flynn. Fractal volume compression. *IEEE transactions on visualization and computer graphics*, 2(4):313–322, 1996.
- 541
- 542
- 543 Robert D Dony and Simon Haykin. Neural network approaches to image compression. *Proceedings of the IEEE*, 83(2):288–303, 1995.
- 544
- 545 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- 546
- 547
- 548
- 549 Marco F Duarte and Richard G Baraniuk. Spectral compressive sensing. *Applied and Computational Harmonic Analysis*, 35(1):111–129, 2013.
- 550
- 551
- 552 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- 553
- 554
- 555 Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. *arXiv preprint arXiv:1803.03635*, 2018.
- 556
- 557
- 558 Elias Frantar and Dan Alistarh. Sparsegpt: Massive language models can be accurately pruned in one-shot. In *International Conference on Machine Learning*, pp. 10323–10337. PMLR, 2023.
- 559
- 560 Ziqi Gao, Qichao Wang, Aochuan Chen, Zijing Liu, Bingzhe Wu, Liang Chen, and Jia Li. Parameter-efficient fine-tuning with discrete fourier transform. *arXiv preprint arXiv:2405.03003*, 2024.
- 561
- 562
- 563 Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pp. 249–256. JMLR Workshop and Conference Proceedings, 2010.
- 564
- 565
- 566
- 567 Song Han, Huizi Mao, and William J Dally. Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding. *arXiv preprint arXiv:1510.00149*, 2015a.
- 568
- 569
- 570 Song Han, Jeff Pool, John Tran, and William Dally. Learning both weights and connections for efficient neural network. *Advances in neural information processing systems*, 28, 2015b.
- 571
- 572 Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. Towards a unified view of parameter-efficient transfer learning. *arXiv preprint arXiv:2110.04366*, 2021.
- 573
- 574
- 575 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pp. 1026–1034, 2015.
- 576
- 577
- 578 Patrick Helber, Benjamin Bischke, Andreas Dengel, and Damian Borth. Eurosat: A novel dataset and deep learning benchmark for land use and land cover classification. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 12(7):2217–2226, 2019.
- 579
- 580
- 581 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*, 2020.
- 582
- 583
- 584
- 585 Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for nlp. In *International conference on machine learning*, pp. 2790–2799. PMLR, 2019.
- 586
- 587
- 588 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- 589
- 590
- 591
- 592 Zhiqiang Hu, Lei Wang, Yihuai Lan, Wanyu Xu, Ee-Peng Lim, Lidong Bing, Xing Xu, Soujanya Poria, and Roy Ka-Wei Lee. Llm-adapters: An adapter family for parameter-efficient fine-tuning of large language models. *arXiv preprint arXiv:2304.01933*, 2023.
- 593

- 594 Ting Jiang, Shaohan Huang, Shengyue Luo, Zihan Zhang, Haizhen Huang, Furu Wei, Weiwei Deng,  
595 Feng Sun, Qi Zhang, Deqing Wang, et al. Mora: High-rank updating for parameter-efficient fine-  
596 tuning. *arXiv preprint arXiv:2405.12130*, 2024.
- 597 Rabeeh Karimi Mahabadi, James Henderson, and Sebastian Ruder. Compacter: Efficient low-rank  
598 hypercomplex adapter layers. *Advances in Neural Information Processing Systems*, 34:1022–  
599 1035, 2021.
- 601 Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete  
602 Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceed-  
603 ings of the IEEE/CVF International Conference on Computer Vision*, pp. 4015–4026, 2023.
- 604 Dawid Jan Kopiczko, Tijmen Blankevoort, and Yuki Markus Asano. Vera: Vector-based random  
605 matrix adaptation. *arXiv preprint arXiv:2310.11454*, 2023.
- 607 Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained  
608 categorization. In *Proceedings of the IEEE international conference on computer vision work-  
609 shops*, pp. 554–561, 2013.
- 611 Namhoon Lee, Thalaisyasingam Ajanthan, and Philip HS Torr. Snip: Single-shot network pruning  
612 based on connection sensitivity. *arXiv preprint arXiv:1810.02340*, 2018.
- 613 Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt  
614 tuning, 2021.
- 615 Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. *arXiv  
616 preprint arXiv:2101.00190*, 2021.
- 618 Yang Li, Shaobo Han, and Shihao Ji. Vb-lora: Extreme parameter efficient fine-tuning with vector  
619 banks. *arXiv preprint arXiv:2405.15179*, 2024.
- 621 Shih-Yang Liu, Chien-Yi Wang, Hongxu Yin, Pavlo Molchanov, Yu-Chiang Frank Wang, Kwang-  
622 Ting Cheng, and Min-Hung Chen. Dora: Weight-decomposed low-rank adaptation. *arXiv  
623 preprint arXiv:2402.09353*, 2024.
- 624 Shiwei Liu, Tianlong Chen, Xiaohan Chen, Li Shen, Decebal Constantin Mocanu, Zhangyang  
625 Wang, and Mykola Pechenizkiy. The unreasonable effectiveness of random pruning: Return  
626 of the most naive baseline for sparse training. *arXiv preprint arXiv:2202.02643*, 2022.
- 627 Xiao Liu, Kaixuan Ji, Yicheng Fu, Weng Lam Tam, Zhengxiao Du, Zhilin Yang, and Jie Tang. P-  
628 tuning v2: Prompt tuning can be comparable to fine-tuning universally across scales and tasks.  
629 *arXiv preprint arXiv:2110.07602*, 2021.
- 631 Y Liu. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*,  
632 2019.
- 633 Subhansu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. Fine-grained  
634 visual classification of aircraft. *arXiv preprint arXiv:1306.5151*, 2013.
- 636 Sourab Mangrulkar, Sylvain Gugger, Lysandre Debut, Younes Belkada, Sayak Paul, and Benjamin  
637 Bossan. Pefit: State-of-the-art parameter-efficient fine-tuning methods. [https://github.  
638 com/huggingface/peft](https://github.com/huggingface/peft), 2022.
- 639 Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct  
640 electricity? a new dataset for open book question answering. *arXiv preprint arXiv:1809.02789*,  
641 2018.
- 642 Decebal Constantin Mocanu, Elena Mocanu, Peter Stone, Phuong H Nguyen, Madeleine Gibescu,  
643 and Antonio Liotta. Scalable training of artificial neural networks with adaptive sparse connect-  
644 ivity inspired by network science. *Nature communications*, 9(1):2383, 2018.
- 645 Mahdi Nikdan, Soroush Tabesh, and Dan Alistarh. Rosa: Accurate parameter-efficient fine-tuning  
646 via robust adaptation. *arXiv preprint arXiv:2401.04679*, 2024.

- 648 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong  
649 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to fol-  
650 low instructions with human feedback. *Advances in neural information processing systems*, 35:  
651 27730–27744, 2022.
- 652 Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, and CV Jawahar. Cats and dogs. In *2012*  
653 *IEEE conference on computer vision and pattern recognition*, pp. 3498–3505. IEEE, 2012.
- 654  
655 Jonas Pfeiffer, Aishwarya Kamath, Andreas Rücklé, Kyunghyun Cho, and Iryna Gurevych. Adapter-  
656 fusion: Non-destructive task composition for transfer learning. *arXiv preprint arXiv:2005.00247*,  
657 2020.
- 658 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
659 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
660 models from natural language supervision. In *International conference on machine learning*, pp.  
661 8748–8763. PMLR, 2021.
- 662 BR Shankara Reddy and ISN Murthy. Ecg data compression using fourier descriptors. *IEEE Trans-*  
663 *actions on Biomedical Engineering*, (4):428–434, 1986.
- 664  
665 Adithya Renduchintala, Tugrul Konuk, and Oleksii Kuchaiev. Tied-lora: Enhancing parameter effi-  
666 ciency of lora with weight tying. *arXiv preprint arXiv:2311.09578*, 2023.
- 667 Tal Ridnik, Emanuel Ben-Baruch, Asaf Noy, and Lihi Zelnik-Manor. Imagenet-21k pretraining for  
668 the masses. *arXiv preprint arXiv:2104.10972*, 2021.
- 669  
670 Mark Rudelson and Roman Vershynin. Sparse reconstruction by convex relaxation: Fourier and  
671 gaussian measurements. In *2006 40th Annual Conference on Information Sciences and Systems*,  
672 pp. 207–212. IEEE, 2006.
- 673 Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adver-  
674 sarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106, 2021.
- 675 Maarten Sap, Hannah Rashkin, Derek Chen, Ronan LeBras, and Yejin Choi. Socialliqa: Common-  
676 sense reasoning about social interactions. *arXiv preprint arXiv:1904.09728*, 2019.
- 677  
678 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée  
679 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and  
680 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.
- 681  
682 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-  
683 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda-  
684 tion and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.
- 685 Marinus T Vlaardingerbroek and Jacques A Boer. *Magnetic resonance imaging: theory and prac-*  
686 *tice*. Springer Science & Business Media, 2013.
- 687  
688 A Wang. Glue: A multi-task benchmark and analysis platform for natural language understanding.  
689 *arXiv preprint arXiv:1804.07461*, 2018.
- 690 Chaoqi Wang, Guodong Zhang, and Roger Grosse. Picking winning tickets before training by  
691 preserving gradient flow. *arXiv preprint arXiv:2002.07376*, 2020.
- 692  
693 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny  
694 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in*  
695 *neural information processing systems*, 35:24824–24837, 2022.
- 696 Zai Yang and Lihua Xie. Exact joint sparse frequency recovery via optimization methods. *IEEE*  
697 *Transactions on Signal Processing*, 64(19):5145–5157, 2016.
- 698  
699 Elad Ben Zaken, Shauli Ravfogel, and Yoav Goldberg. Bitfit: Simple parameter-efficient fine-tuning  
700 for transformer-based masked language-models. *arXiv preprint arXiv:2106.10199*, 2021.
- 701  
702 Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a ma-  
chine really finish your sentence? *arXiv preprint arXiv:1905.07830*, 2019.



702 Yuchen Zeng and Kangwook Lee. The expressive power of low-rank adaptation. *arXiv preprint*  
703 *arXiv:2310.17513*, 2023.  
704

705 Longteng Zhang, Lin Zhang, Shaohuai Shi, Xiaowen Chu, and Bo Li. Lora-fa: Memory-efficient  
706 low-rank adaptation for large language models fine-tuning. *arXiv preprint arXiv:2308.03303*,  
707 2023a.

708 Renrui Zhang, Jiaming Han, Chris Liu, Peng Gao, Aojun Zhou, Xiangfei Hu, Shilin Yan, Pan Lu,  
709 Hongsheng Li, and Yu Qiao. Llama-adapter: Efficient fine-tuning of language models with zero-  
710 init attention. *arXiv preprint arXiv:2303.16199*, 2023b.  
711

712 Yihua Zhang, Yuguang Yao, Parikshit Ram, Pu Zhao, Tianlong Chen, Mingyi Hong, Yanzhi Wang,  
713 and Sijia Liu. Advancing model pruning via bi-level optimization. *Advances in Neural Informa-*  
714 *tion Processing Systems*, 35:18309–18326, 2022.  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

## A ADDITIONAL EXPERIMENTAL DETAILS

### A.1 COMPUTATIONAL HARDWARE

All our experiments were carried out on Linux servers equipped with an AMD EPYC 7763 64-Core CPU processor, 512GB RAM, and NVIDIA RTX 6000 ADA 48G / A800 80G GPU.

### A.2 HYPERPARAMETERS

Table 6: Hyperparameter configurations for GLUE benchmark.

Model	Hyperparameter	SST-2	MRPC	CoLA	QNLI	RTE	STS-B
	Optimizer	AdamW					
	LR Scheduler	Linear					
	Warmup Ratio	0.06					
	Max Seq. Len.	512					
	Spectral Coefficients $n$	{250,500}					
BASE	Rank $r$	32	64	64	8	32	256
	Epochs	50	30	100	40	100	90
	Batch Size	128	32	128	32	32	32
	LR (Head)	6E-4	6E-4	3E-4	6E-5	3E-4	2E-4
	LR (FoRA)	2E-2	4E-2	4E-2	7E-2	3E-2	2E-2
LARGE	Rank $r$	32	32	32	32	32	32
	Epochs	20	50	100	30	70	40
	Batch Size	128	32	128	8	32	32
	LR (Head)	1E-4	2E-4	4E-4	4E-4	3E-4	7E-5
	LR (FoRA)	3E-2	5E-2	4E-2	2E-2	2E-2	3E-2

Table 7: Hyperparameter configurations for mathematical reasoning.

Hyperparameter	LLaMA2 <sub>7B</sub>		LLaMA2 <sub>13B</sub>		LLaMA3 <sub>8B</sub>	
	GSM8k	MATH	GSM8k	MATH	GSM8k	MATH
Optimizer	AdamW					
LR Scheduler	Cosine					
Batch Size	16					
Warmup Ratio	0.05					
Dropout	0.05					
Epochs	3					
Where	Q,V					
Spectral Coefficients $n$	20000		30000		20000	
Rank $r$ (FoRA)	256	128	256	128	256	128
Rank $r$ (DFoRA)	256	128	256	128	128	128
LR (LoRA)	5E-4	5E-4	5E-4	6E-4	5E-4	5E-4
LR (DoRA)	4E-4	5E-4	4E-4	6E-4	6E-4	2E-4
LR (FoRA)	6E-3	5E-3	5E-3	5E-3	1E-3	9E-4
LR (DFoRA)	5E-3	3E-3	6E-3	6E-3	1E-3	9E-4

Table 8: Hyperparameter configurations for commonsense reasoning.

Hyperparameter	LLaMA <sub>7B</sub>		LLaMA <sub>13B</sub>		LLaMA <sub>27B</sub>		LLaMA <sub>38B</sub>	
	FoRA	DFoRA	FoRA	DFoRA	FoRA	DFoRA	FoRA	DFoRA
Optimizer	AdamW							
LR Scheduler	Linear							
Batch Size	16							
Warmup Steps	100							
Dropout	0.05							
Epochs	3							
Rank $r$	32							
Alpha $\alpha$	64							
Where	Q,K,V,Up,Down							
Spectral Coefficients $n$	30000		40000		30000		30000	
LR	1E-3	1.4E-3	9E-4	9E-4	8E-4	8E-4	5E-4	5E-4

Table 9: Hyperparameter configurations for finetuning ViT on the image classification datasets.

Model	Hyperparameter	OxfordPets	StanfordCars	DTD	EuroSAT	FGVC	RESISC
	Optimizer	AdamW					
	Epochs	10					
	Batch Size	64					
	Rank $r$ (LoRA)	16					
	Spectral Coefficients $n$	8000					
BASE	Rank $r$ (FoRA)	32	128	64	64	256	32
	LR (Head)	8E-3	1E-2	1E-2	1E-4	1E-2	1E-2
	LR (FoRA)	4E-3	5E-2	5E-3	2E-2	5E-2	2E-2
	Weight Decay	4E-2	1E-5	2E-4	4E-3	2E-2	9E-2
LARGE	Rank $r$ (FoRA)	64	128	128	64	256	32
	LR (Head)	6E-3	5E-3	1E-2	1E-3	1E-2	1E-2
	LR (FoRA)	5E-3	3E-2	4E-3	3E-2	8E-2	1E-2
	Weight Decay	3E-4	2E-5	3E-5	3E-3	1E-2	1E-3

### 864 A.3 PARAMETER COUNT OF SPARSE LEARNING STRATEGIES

865  
866 As the rank increases, the number of learnable parameters in LoRA grows linearly, leading to a  
867 significant parameter overhead. While VeRA exhibits a minimal increase in parameters, its strong  
868 dependence on the size of its adaptation matrices limits its flexibility in adapting to more complex  
869 tasks. In contrast, both FoRA and random masking maintain a fixed number of learnable parameters  
870 across different ranks, providing greater flexibility by allowing parameter adjustments based on task  
871 complexity.

872 Table 10: Comparison of learnable parameters across different compression strategies.

874	875	876	877 Rank $r$				
			878 $2^3$	879 $2^4$	880 $2^5$	881 $2^6$	882 $2^7$
884 RoBERTa <sub>BASE</sub>	885 LoRA	886 6,144	887 12,288	888 24,576	889 49,152	890 98,304	891 196,608
	892 VeRA	893 776	894 784	895 800	896 832	897 896	898 1024
	899 FoRA/Mask	900 500	901 500	902 500	903 500	904 500	905 500
906 ViT <sub>BASE</sub>	907 LoRA	908 6,144	909 12,288	910 24,576	911 49,152	912 98,304	913 196,608
	914 VeRA	915 776	916 784	917 800	918 832	919 896	920 1024
	921 FoRA/Mask	922 6,144	923 8000	924 8000	925 8000	926 8000	927 8000

## 884 B IMPLEMENTATION

885 Algorithm 1 presents the PyTorch implementation of FoRA. Our approach allows for a straight-  
886 forward plug-in, with the only modification needed being the replacement of the PyTorch linear  
887 module with the FoRA linear module. Additionally, we create a cached empty spectral matrix that  
888 matches the size of the adaptation matrix. This caching strategy not only accelerates GPU compu-  
889 tations but also minimizes GPU overhead, as only the sparse spectral coefficients require backward  
890 propagation.  
891  
892  
893

---

### 894 Algorithm 1 FoRA PyTorch Implementation

---

```

895 class PEFT(nn.Module):
896     def __init__(self, ...):
897         if LoRA:
898             A, B = torch.nn.Linear(d1, r), torch.nn.Linear(r, d2)
899         elif FoRA:
900             A, B = FoRALinear(n, d1, r), FoRALinear(n, r, d2)
901
902 class FoRALinear(nn.Module):
903     def __init__(self, n, din, dout):
904         self.n = n # number of spectral coefficients
905         self.din = din # input dimension
906         self.dout = dout # output dimension
907         # spectral location
908         self.L = torch.randperm(din * dout)[:n]
909         self.L = torch.stack([self.L // din, self.L % dout], dim=0)
910         # spectral coefficients
911         self.s = nn.Parameter(torch.zeros(n))
912         # cached spectral matrix
913         self.cached_F = torch.zeros(self.din, self.dout)
914
915     def weight(self):
916         # spectral assignment
917         F = self.cached_F.clone()
918         F[self.L[0], self.L[1]] = self.s
919         # spatial conversion
920         W = torch.fft.ifft2(F).real()
921         return W

```

---

## C GENERATION EXAMPLES

## Instruction Tuning (MATH)

How many inches are in the circumference of a great circle of a sphere which has a surface area of  $196\pi$  sq in? Express your answer in terms of  $\pi$ .

Note: To measure the circumference, take the largest circle on the surface of the sphere.

Answer the above question. First think step by step and then answer the final number.

## Response

**Model****LLaMA2<sub>7B</sub> (FoRA)****Response**

The surface area of a sphere is  $\pi r^2$ , where  $r$  is the radius.

So, the radius is  $r = \sqrt{196\pi} = 12\sqrt{3}$ .

The circumference of the sphere is  $2\pi r = \boxed{24\pi}$ .

The final answer is:  $24\pi$ .

**LLaMA3<sub>8B</sub> (FoRA)**

The surface area of a sphere is  $4\pi r^2$ , so  $r^2 = 49$  and  $r = 7$ .

The circumference of the great circle is  $2\pi r = \boxed{14\pi}$ .

The final answer is:  $14\pi$ .

## Instruction Tuning (HellaS.)

Please choose the correct ending to complete the given sentence: Roof shingle removal: A man is sitting on a roof. he

Ending1: is using wrap to wrap a pair of skis. Ending2: is ripping level tiles off. Ending3: is holding a rubik's cube. Ending4: starts pulling up roofing on a roof.

Answer format: ending1/ending2/ending3/ending4

## Response

**Model****LLaMA2<sub>7B</sub> (FoRA)****Response**

the correct answer is ending2

**LLaMA3<sub>8B</sub> (FoRA)**

the correct answer is ending4