# SHARELORA: LESS TUNING, MORE PERFORMANCE FOR LORA FINE-TUNING OF LLMS

Anonymous authors

Paper under double-blind review

# ABSTRACT

Fine-tuning large language models (LLMs) is prohibitively expensive, prompting the development of various parameter-efficient fine-tuning (PEFT) methods. These methods primarily focus on fine-tuning small, additional modules known as *adapters*, which account for only a small fraction of the total LLM parameters. One such method, low-rank adaptation (LoRA), has shown notable parameter efficiency while maintaining performance comparable to full fine-tuning. However, classical LoRA may still involve tuning more parameters than necessary given the intrinsic rank of pre-trained weights, as highlighted by prior work (Aghajanyan et al., 2020). Recent variants of LoRA aim to enhance fine-tuning performance, but they overlook the layer-wise redundancies that can be leveraged for more efficient weight sharing. In this work, we introduce SHARELORA, a novel approach that further enhances parameter efficiency during LLM fine-tuning by leveraging redundancies in pre-trained model weights to share LoRA modules, thereby significantly reducing the number of trainable parameters. Specifically, SHARELORA automatically identifies redundancies in the pre-trained weights and determines which LoRA adapters can share parameters. This is achieved by measuring the similarity between representations to assess information redundancy and using a greedy algorithm to maximize parameter sharing. We conducted extensive evaluations on the LLMs of the LLaMA family across benchmark tasks. Notably, SHARELORA achieves better parameter efficiency, with up to a 23% reduction in the number of fine-tuned parameters while delivering performance comparable to or better than existing PEFT methods.

031 032

004

010 011

012

013

014

015

016

017

018

019

021

023

025

026

027

028

029

### 033 034

035

# 1 INTRODUCTION

Large language models (LLMs), e.g., GPT-4 and LLaMA2, are at the forefront of advances in the field of machine learning (ML). These large models are pre-trained on vast datasets (e.g., images or text 037 corpora) and are subsequently fine-tuned for specialized tasks, demonstrating proficiency in domains such as natural language, image processing, and fundamental scientific discoveries (Bommasani et al., 2021; Touvron et al., 2023a; Singhal et al., 2022; 2023). These models, often referred as 040 "base model", are pre-trained solely to predict the next token to generate from their entire vocabulary 041 space (Touvron et al., 2023a; Penedo et al., 2023; Team, 2023). To employ the base model for 042 real applications, e.g., building chatbots, then often need further fine-tuning (e.g., on multi-turn 043 human-human or human-chatbot conversations) to follow specific human instructions or align with 044 human preferences (Leike et al., 2018; Ziegler et al., 2019; Chung et al., 2022).

Fine-tuning such large-scale LLMs, however, presents a computational challenge due to the massive number of parameters. For instance, GPU memory must be sufficiently large to handle the fine-tuning process while marinating a reasonably large batch size. Parameter-Efficient Fine-Tuning (PEFT) methods have been proposed to tackle this challenge by allowing fine-tuning of only a small subset of LLM parameters or incorporating small adapter modules on top of the pre-trained model, while leaving the majority of the base LLM parameters frozen (Houlsby et al., 2019; Hu et al., 2021; Zaken et al., 2021; Zhang et al., 2023). These methods democratize LLM fine-tuning by making it feasible on commodity hardware. One popular PEFT method, LoRA, has been shown to effectively reduce the GPU memory requirement during LLM fine-tuning (Hu et al., 2021). LoRA achieves parameter efficiency by adding low-rank adapters in parallel with specific LLM parameters, such as the query,

FT Method	GSM8K $\uparrow$	ARC-Challenge ↑	WinoGrande ↑	Hellaswag $\uparrow$
LoRA(r = 12)	37.98	48.21	64.25	51.96
Naive-shared LoRA	37.23	47.95	62.83	48.38

058 059 060

054

056

Table 1: Accuracy comparison of LoRA and a naive-share LoRA strategy.

key, and value parameter weights in multi-head attention. During fine-tuning, LoRA LoRA optimizesonly the low-rank adapters, while the LLM parameters remain unchanged.

Although LoRA is efficient, it treats all layers uniformly, lacking finer control over which layers are most important or exhibit similar behavior. Recent improvements, such as AdaLoRA (Zhang et al., 2023), which dynamically adjusts the rank based on layer importance, DoRA (Liu et al., 2024a), which decouples weights into direction and magnitude for more nuanced fine-tuning, and LoRA+ (Hayou et al., 2024), which independently adjusts the learning rates of LoRA components, have aimed to enhance LoRA's efficiency. However, these approaches overlook the redundancy present in pre-trained foundation models, where certain layers may exhibit similar behavior and can potentially share parameters, further reducing memory requirements.

071 We observed that sharing the LoRA module's weights across layers does not significantly degrade 072 performance, while effectively reducing the number of trainable parameters. Specifically, we exper-073 imented by sharing the LoRA weights between odd-numbered and even-numbered layers, which 074 halved the number of trainable parameters. In this naive approach, the weights of even-numbered 075 layers were directly mirrored to their adjacent odd-numbered layers. Surprisingly, this straightforward weight-sharing strategy led to only minor performance degradation, as demonstrated in Table 1. This 076 finding suggests that there is significant potential for further optimization through weight sharing 077 in LoRA modules. However, determining which layers should share weights remains an ongoing challenge, largely due to the limited explainability of foundation models (Zhao et al., 2024). The 079 behavior and interaction of different layers within these models are not yet fully understood, and there is no consensus on which layers exhibit sufficiently similar representations to justify weight 081 sharing. This is an active area of research, and more sophisticated methods for identifying redundant 082 layers could lead to even more efficient weight-sharing strategies in the future. 083

In foundation models, certain layer representations often exhibit notable similarities. This redundancy is a result of insufficient training data, which prevents each parameter from learning distinct, unique features. Consequently, this leads to overlapping or redundant representations across layers. Prior research has harnessed these redundancies for model compression (Gromov et al., 2024). Building on this insight, our work leverages this redundancy by sharing LoRA weights across layers with similar representations during fine-tuning. This approach allows us to slightly increase the LoRA rank without increasing the overall number of trainable parameters, thereby enhancing fine-tuning performance.

091 Our SHARELORA method consists of two main components: (i) computing similarity matrices 092 between representations of layers and (ii) sharing the LoRA module parameters among redundancy layers. This method identifies layers with similar representations and shares their weights, reducing 094 the number of trainable parameters while maintaining model performance. By leveraging layer 095 similarities, SHARELORA significantly improves finetuning efficiency. We conduct extensive ex-096 periments on a wide range of tasks and models to demonstrate the effectiveness of SHARELORA. 097 Specifically, we evaluate the performance using LLaMA-7B, LLaMA2-7B, and LLaMA3-8B for 098 natural language commonsense reasoning and LLava-1.5-7B for image-text understanding. The results show that SHARELORA performs similar or better than LoRA, while only using 80% trainable 099 parameter. 100

- <sup>101</sup> The summary of our contributions is as follows:
- 102 103 104

- We propose SHARELORA, a novel parameter-efficient fine-tuning (PEFT) method that leverages the similarity of layer representations to enable weight sharing, achieving this without introducing any additional train or inference latency compared to LoRA.
- 107
- We develop a greedy algorithm to determine which layers should share the weights of the LoRA module based on the similarity of their representations.

• We conduct extensive experiments demonstrating that SHARELORA consistently performs similar or better LoRA across various tasks, while using less trainable parameters.

110 111 112

113

108

#### 2 PRELIMINARY AND MOTIVATION

**Parameter Efficient Fine Tuning.** A primary research trajectory aimed at reducing the fine-tune 114 parameters of pretrained FMs is to model the incremental updates of pretrained weights in a parameter-115 efficient manner. For example, given a pretrained weight matrix W, diff pruning (Guo et al., 2021) 116 initializes  $\Delta$  as the same dimensions as W and performs magnitude-based pruning on  $\Delta$ . Diff pruning 117 characterizes  $\Delta$  as the incremental updates of W, and it can improve the parameter efficiency due to 118 the sparsity of  $\Delta$ . However, it requires specific hardware support to accelerate the computation of 119 unstructured sparse matrices. This hardware-specific dependency underscores a crucial consideration 120 for the practical deployment of such an approach in real-world applications. In addition, it does not 121 significantly reduce computational cost compared to full fine-tuning (Hu et al., 2021), as every entry 122 of  $\Delta$  needs to be updated and then be pruned.

To tackle those limitations, Hu et al. propose LoRA (Hu et al., 2021), which parameterize  $\Delta$  as the product of two low-rank matrices:

$$W' = W + \Delta = W + BA,\tag{1}$$

where  $\Delta \in \mathbb{R}^{d_1 \times d_2}$ ,  $B \in \mathbb{R}^{d_1 \times r}$ , and  $A \in \mathbb{R}^{r \times d_2}$  with  $r \ll \{d_1, d_2\}$ . As the rank r is much smaller than the dimension of W (e.g., r = 8 and  $d_1 = d_2 = 4096$ ), the number of trainable parameters and training overhead dramatically decreases.

However, LoRA has its limitations, as it typically applies separate parameters for each B and A
by default. This approach overlooks the significant variation in the redundancy of weight matrices
across different layers and modules during the fine-tuning of pretrained foundation models. We will
demonstrate this issue later. Consequently, LoRA cannot adaptively share the LoRA modules among
the redundancy layers, which could otherwise achieve comparable performance with fewer trainable
parameters.

137

126

127

Weight Sharing. Although PEFT methods can reduce the number of trainable parameters, thereby 138 decreasing GPU memory footprint during fine-tuning, the number of parameters in LLMs also scales 139 rapidly, making it increasingly challenging to conduct even PEFT on commodity GPUs. In this 140 paper, we explore the possibility of combining LoRA with weight sharing, a method that allows 141 multiple neural network layers to share the same model weights. Weight sharing has been a widely 142 adopted technique to reduce the number of trainable parameters while maintaining performance and 143 sometimes helping mitigate overfitting (Press & Wolf, 2017; de Lhoneux et al., 2018; Lan et al., 2020; 144 Dai et al., 2020; Takase & Kiyono, 2021). 145

146 Similarity Across Layers. Our motivation for weight 147 sharing among LoRAs across layers comes from the observation that the representations among model layers/blocks 148 attain high levels of similarity, especially for bottom layer 149 blocks (as shown in Figure 1). This effect has also been 150 observed by many prior works (Kornblith et al., 2019; Gro-151 mov et al., 2024). This effect is also connected to methods 152 like *layer skipping* and *mixture of depth* for LLMs (Fan 153 et al., 2024; Elhoushi et al., 2024; Raposo et al., 2024). 154 Previous research (He et al., 2024) demonstrates that the 155 distribution of similarity scores remains consistent, even 156 when the calibration dataset sample size varies. Addition-157 ally, prior studies indicate that altering the type of cali-158 bration dataset—whether it be pretraining datasets (e.g., 159 C4 (Raffel et al., 2020)) or instruction-tuning datasets (e.g., , CodeAlpaca-20k, MathInstruct (Xiang Yue, 2023), and 160 LIMA (Zhou et al., 2024))-has minimal effect on the 161 similarity score distribution.



Figure 1: Cosine similarity among layers' representations experimented on LLaMA2-7B model over the GSM8k evaluation dataset.



Figure 2: An overview of our proposed AutoLoRA, which computes the cosine similarity between each layer, and shares the LoRA parameters within redundancy blocks.

# 3 SHARELORA METHOD

As illustrated in Figure 2, our method comprises two key components: (i) layer similarity computation and (ii) LoRA sharing.

#### 190 191 3.1 THE FORMULATION OF SHARELORA

In this subsection, we present the formulation of SHARELORA as an optimization problem. Specifically, we consider an *L*-layer LLM. We denote S as the collection of shared sets of LLM layer indices, e.g.,  $S = \{(2,3), (5,6,7), (11,13)\}$ .

Our objective is to maximize the number of shared LLM layers while ensuring that the similarity measure between any two layers with shared weights exceeds a predefined threshold,  $\epsilon$ . This can be formulated as the following optimization problem:

$$\max \quad \sum_{\forall s \in \mathcal{S}} |s| \tag{2}$$

s.t. 
$$c(\mathbf{x}_i, \mathbf{x}_i) \ge \epsilon$$
 (3)

where |s| is the number of layers in set *s*, and *S* is the collection of sets such that each set *s* contains layers with similarity measure (*i.e.*,  $c(\cdot, \cdot)$ ) between each pair exceeding the threshold  $\epsilon$ .

The objective of SHARELORA is to maximize the total number of layers sharing LoRA. The remaining question is *How to find S*? We will present our algorithm of SHARELORA in section 3.3.

#### 3.2 SIMILARITY BETWEEN REPRESENTATIONS OF LAYERS

To decide which layer's LoRA module could be shared with another, we have to compute the angular distance  $c(\mathbf{x}_i, \mathbf{x}_j)$  between the representations of layer *i* and layer *j*. The similarity of a single sequence of length *T* is given by

214

182

183

185

186 187

188

189

199 200 201

202 203

204

205

208

209

$$c(\mathbf{x}_i, \mathbf{x}_j) \coloneqq \frac{1}{\pi} \cos^{-1} \left( \frac{\mathbf{x}_i^\top \mathbf{x}_j}{\|\mathbf{x}_i\| \|\mathbf{x}_j\|} \right), \tag{4}$$

where the inner product is over the hidden dimension of the model for the final token T of the sequence,  $\|\cdot\|$  denotes the  $\ell_2$ -norm, and the factor of  $1/\pi$  is a convention. This distance should then be summed over a number of examples that is large enough to get a low-fluctuation estimate but overall should be quite small.

220 221 222

# 3.3 SIMILARITY-BASED WEIGHT SHARING

The objective of the algorithm is to identify sets of neural network layers that can share LoRA modules based on the angular distance of their representations. The algorithm employs a greedy strategy to maximize the size of these shared sets while adhering to a predefined similarity threshold.

First, the algorithm processes the similarity matrix to create an upper triangular matrix, excluding the diagonal, which represents the pairwise similarities between layers. It then identifies all eligible pairs of layers (i, j) where the similarity score exceeds the threshold  $\epsilon$ . Formally, this can be expressed as: shared\_pairs = { $(i, j) | c(\mathbf{x}_i, \mathbf{x}_j) \ge \epsilon$  and  $0 \le i < j < L$ }.

Next, the algorithm constructs sets of shared layers by iterating through the identified shared pairs. It maintains a set visited to avoid reprocessing layers. The construction of the shared sets proceeds is shown in Algorithm 1.

234

235		
236	Algo	rithm 1 Construct Shared Sets
237	1: ]	<b>Input:</b> Similarity matrix C, threshold $\epsilon$ ;
238	2: 0	<b>Output:</b> Collection of shared sets $S$ ;
239	3: ]	Initialize shared_pairs $\leftarrow \{(i, j) \mid c(\mathbf{x}_i, \mathbf{x}_j) \ge \epsilon \text{ and } 0 \le i < j < L\};$
240	4: 3	Sort shared_pairs by similarity scores;
241	5:	Initialize $\mathcal{S} \leftarrow \{\}$ , and initialize visited $\leftarrow \text{set}()$ ;
242	6: 1	for all $(i, j)$ in shared_pairs do
2/12	7:	if $i$ not in visited and $j$ not in visited <b>then</b>
044	8:	if $i$ not in $S$ then
244	9:	Initialize $S_i \leftarrow \{i\}$
245	10:	end if
246	11:	Add j to $S_i$ , and add j to visited
247	12:	if All pairs starting with <i>i</i> is done then
248	13:	Add <i>i</i> to visited
249	14:	end if
250	15:	end if
251	16: 0	end for
050		

252 253

254

255

The algorithm ensures that layers are grouped into shared sets only if their pairwise similarities exceed the threshold  $\epsilon$ , thus maximizing the number of layers that can share the LoRA modules while maintaining high similarity within each set.

256 257 258

259

260

261

262

263

264

# Algorithm 2 SHARELORA

- 1: **Input:** Sample Dataset  $\mathcal{D}^*$ ; Train Dataset  $\mathcal{D}$ ; hyper-parameter similarity threshold  $\epsilon$ .
- 2: Inference on  $\mathcal{D}^*$  and save the representations x(l) for each layer l;
- 3: Compute the angular distance  $d(x(l_p), x(l_q))$  between each layer  $l_p$  and  $l_q$ ;
- 4: Compute the share set S and update LoRA rank r;
- 5: Share parameters using S;
- 6: Finetune weight sharing model W with  $\mathcal{D}$ .
- 7: **Output:** The fine-tuned parameters  $W^*$ .

265 266 267

Furthermore, after sharing layers using the set S, the number of layers with shared parameters will be  $\sum_{S_i \in S} |S_i| - |S|$ . This allows us to optionally expand the original rank r to  $\lfloor \frac{L}{L - \sum_{S_i \in S} |S_i| + |S|} r \rfloor$ , where L is the total number of layers. We summarize the algorithm in Algorithm 2.

Model	PEFT Method	# Params	BoolQ	PIQA	SIQA	HS	WG	ARC-e	ARC-c	OBQA	Avg
	Prefix	10.5M	57.46	50.49	33.62	28.38	64.8	29.84	24.49	26.6	39.4
	AdapterH	20.1M	71.19	74.48	45.45	57.24	59.51	56.90	33.70	39.0	54.0
LLaMA-7B	AdapaterP	20.1M	67.65	73.45	44.27	57.04	58.48	57.62	33.87	37.0	53.
	Parallel	448M	73.36	74.54	73.81	57.08	60.22	56.23	35.58	36.8	54.
	LoRA(r = 8)	14.0M	78.53	78.45	46.98	73.41	70.17	70.24	41.21	42.2	62.
	SHARELORA*	12.1M	73.94	79.38	47.08	72.93	68.11	65.99	37.8	43.0	61.
	SHARELORA(ours)	<b>10.8</b> M	77.55	79.65	47.54	73.04	70.96	69.95	40.70	45.4	63.
LLaMA2-7B	LoRA(r = 8)	14.0M	80.64	79.16	47.85	75.18	69.93	69.70	42.06	44.0	63.
EDuivin 12 7 D	SHARELORA(ours)	11.8 M	80.89	79.22	46.78	74.04	71.27	68.73	41.81	44.0	63.
LLaMA3-8B	LoRA(r = 8)	14.2M	82.17	81.34	49.49	78.38	74.19	76.39	50.77	46.4	67.
	SHARELORA(ours)	11.4M	82.17	81.56	48.77	79.32	73.80	76.85	50.00	44.2	67.

Table 2: Accuracy with LLaMA model family with various PEFT methods on commonsense reasoning tasks.

4 EXPERIMENTS

287 Models and datasets. We implemented SHARELORA to fine-tune LLaMA family LLMs, namely, 288 LLaMA-7B (Touvron et al., 2023a), LLaMA2-7B (Touvron et al., 2023b) and LLaMA3-8B (Meta, 289 2024). We follow the settings of LLM-Adapters (Hu et al., 2023), and evaluate the effectiveness of the several natural language commonsense reasoning task including BoolQ (Clark et al., 2019), 290 PIQA (Bisk et al., 2019), SIQA (Sap et al., 2019), HellaSwag (Zellers et al., 2019), Winogrande (Sak-291 aguchi et al., 2021), ARC-Easy (Clark et al., 2018), ARC-Challenge, and OBQA (Mihaylov et al., 292 2018). Additionally, we implemented SHARELORA to fine-tune LLaVA-1.5-7B (Liu et al., 2023), a 293 popular vision language foundation model (VLM) and used on image-text pair understanding, and evaluated on LLaVA-Bench (in-the-wild) evaluation dataset (Liu et al., 2023). 295

Setup. We use PyTorch (Paszke et al., 2019) to implement all the algorithms. Our fine-tuning algorithm implementation is based on the publicly available Huggingface Transformers (Wolf et al., 2019) and LLM-Adapters code base. All the experiments are conducted on NVIDIA A6000 GPUs.

299 300

301

302

303

305

306

307

308

309

310

311

312

313

314

315

316

283

284 285

- Baselines. We compare SHARELORA with the following baselines.
  - *Prompt learning (Prefix):* (Li & Liang, 2021) Involves fine-tuning a small set of continuous task-specific vectors (prefixes) while keeping the large language model parameters frozen to pre-trained weights.
- Adapter tuning (AdapterH): (Houlsby et al., 2019) Inserts small add-on layers between the multi-head attention modules and FFN modules of the pre-trained model that can be fine-tuned for downstream task learning, while keeping the rest of the model parameters frozen.
- *Pfeiffer adapter (AdapterP):* (Pfeiffer et al., 2020) Unlike adapter tuning it inserts add-on layers after FFN modules and LayerNorm modules, allows them to fine-tune and keeps the rest of the model frozen.
  - *Parallel adapter (Parallel):* (He et al., 2021) Modifies the hidden representations in a transformer model by inserting additional trainable parameters in parallel to the original model's layers.
    - *LoRA:* (Hu et al., 2021) Is the most popular PEFT method that injects trainable low-rank matrices into transformer layers parallel to the frozen main path, to approximate the weight updates.
      - SHARELORA\* is a variation of our SHARELORA method, but instead of using similarity scores to select the set S, it randomly selects adjacent layers for weight sharing.
- 317 4.1 EVALUATIONS ON LLMS

318

We assess the fine-tuning performance of LLaMA-7B, LLaMA2-7B, and LLaMA3-8B using all
baseline methods along with the proposed algorithm. The commonsense reasoning evaluation
includes eight sub-tasks, each with its own predefined training and testing sets. Following the setup
from LLM-Adapters (Hu et al., 2023), we combine the training datasets from all eight tasks to form a
comprehensive training dataset for fine-tuning and performing evaluations on the individual testing
sets for each sub-task.

324 **Implementation details.** All of LLaMA models have 32 hidden layers. Initially, we calculate the 325 similarity  $c(\mathbf{x}_i, \mathbf{x}_i)$  between each layer's representation using 256 randomly sampled C4 validation 326 dataset (Raffel et al., 2020). We set the similarity threshold  $\epsilon$  to 0.85 for LLaMA-7B and LLaMA2-7B, 327 and 0.80 for LLaMA3-8B. Utilizing our similarity-based weight-sharing algorithm, the share set 328 collection S for LLaMA-7B is  $\{16, 17\}, \{18, 19\}, \{20, 21, 22\}, \{23, 24, 25, 26, 27\}, \{28, 29, 30\}\},\$ enabling the sharing the LoRA modules of ten similar layers. The share set collection S for 329 LLaMA2-7B and LLaMA3-8B are  $\{\{17, 18\}, \{19, 20\}, \{21, 22, 23\}, \{24, 25, 26, 27\}, \{28, 29\}\}$ 330 and  $\{\{19, 20\}, \{21, 22\}, \{23, 24, 25\}, \{26, 27, 28\}\}$  separate. Note that we expand the original 331 LoRA rank from 8 to 9 for these three models. The hyperparameters for fine-tuning across LLaMA-332 7B, LLaMA2-7B, and LLaMA3-8B models are consistent for most settings. A dropout rate of 0.05, 333 AdamW optimizer, a linear learning rate scheduler, a batch size of 16, and 3 epochs are applied for 334 all models. The learning rate (LR) is set to 2e-4 for LLaMA-7B and LLaMA2-7B, while it is reduced 335 to 1e-4 for LLaMA3-8B. The fine-tuning targets the Q, K, V, Down, and Up layers for each model. 336 The share set collection S for SHARELORA\* is randomly generated as  $\{\{2,3\}, \{4,5,6\}, \{9,10\}, \}$ 337  $\{11, 12, 13\}, \{14, 15\}, \{19, 20\}, \{22, 23\}, \{29, 30\}\}.$ 338

Table 2 shows experimental results on the eight commonsense reasoning tasks. SHARELORA achieves similar or better performance across all datasets for all the models. Notably, our proposed approach saves up to 23% of the trainable parameters yet achieves a 1.5% improvement in performance with LLaMA-7B, while maintaining similar performance on LLaMA2-7B and LLaMA3-8B.

- 342
- 343 344

345

4.2 EVALUATIONS ON MULTI-MODAL VLMS

346 We now present the results of SHARELORAON 347 vision-language models. We used the LLaVA-1.5-7B (Liu et al., 2023), which consists of a 348 language model, a visual encoder and a projec-349 tion layer for feature alignment. Specifically, 350 the language model and visual encoder were ini-351 tialized with Vicuna-1.5-7B (Zheng et al., 2024) 352 and CLIP ViT-L/336px (Radford et al., 2021), 353 respectively. In addition, we employed the pre-354 trained projection layer<sup>1</sup> and directly performed 355 visual instruction fine-tuning stage. In contrast 356 to the setup of LLaVA, we chose a subset of 357 their dataset, i. e. LLaVA-Instruct-80K. It con-358 sists of 80k image-instruction pairs filtered out from LLaVA-Instruct-150K<sup>2</sup>. We set the rank to 359

	LoRA(rank=64)	SHARELORA
# Params	159.9M	<b>139.4</b> M
Description	$45.57\pm0.60$	$\textbf{45.80} \pm 0.40$
Conversation	$46.87 \pm 1.27$	$\textbf{54.10} \pm 0.00$
Reasoning	$\textbf{66.90} \pm 0.85$	$65.37\pm0.72$
All	$55.47 \pm 0.55$	<b>57.00</b> $\pm$ 0.44

Table 3: Instruction-following capability comparison between LoRA and SHARELORA (ours). We conduct three repeated evaluations and report the average scores. The results are reported in the format of  $mean \pm std$ .

64 and performed instruction tuning for one epoch of the LLM component using SHARELORA and
 LoRA, respectively. Note, we keep the feature transformation module trainable for both.

362 **Implementation details.** For SHARELORA we compute the similarity between the inputs and 363 outputs of each layer based on 128 samples, with each sample comprising of visual-text inputs. The similarity threshold was set to 0.9 and the share set collection S for language model is  $\{2, 3\}, \{4, 5\}, \{4,$ 364  $\{6, 7\}, \{17, 18\}, \{19, 20\}, \{21, 22, 23\}, \{24, 25, 26\}, \{27, 28, 29\}\}$ . The hyperparameters used for fine-tuning both LoRA and SHARELORA on the vision-language model are as follows. For the 366 similarity threshold, only SHARELORAutilizes a value of 0.9, while LoRA does not incorporate this 367 parameter. Both methods share the same rank r = 64, scaling factor  $\alpha = 128$ , and a dropout rate of 368 0.05. The target layers include O, K, V, O, Up, Down, and Gate for both methods. The training runs 369 for one epoch with a learning rate of 2e-4, and the learning rate is adjusted using a cosine scheduler. 370 Both approaches employ the AdamW optimizer, with a batch size of 64 and a warmup ratio of 371 0.03. These shared hyperparameters ensure consistent conditions for comparing the performance of 372 LoRA and SHARELORA. We use GPT-4 generated answers as the golden answers. Subsequently, 373 we employed GPT-40 to evaluate the instruction-following capabilities of the fine-tuned model and 374 selected the challenging LLaVA-Bench (in-the-wild).

<sup>&</sup>lt;sup>1</sup>https://huggingface.co/liuhaotian/llava-v1.5-mlp2x-336px-pretrain-vicuna-7b-v1.5 <sup>2</sup>https://huggingface.co/datasets/liuhaotian/LLaVA-Instruct-150K

It comprises 24 images with a total of 60 questions. In Table 3, we demonstrate the performance of SHARELORA and LoRA fine-tuned LLaVA model. Overall, SHARELORA exhibits superior performance compared to that with LoRA, with a significant advantage in the conversation aspect. Additionally, SHARELORA requires around 13% fewer trainable parameters to yield this improved performance.

383 384

386 387

#### 4.3 Ablations and Discussions

# Robustness towards different rank settings.

Here we investigate the impact of various rank 389 configurations on SHARELORA and LoRA 390 by adjusting the original r within the set 391  $\{2, 4, 8, 16\}$ . We then evaluate the performance 392 of fine-tuned LLaMA-7B on commonsense rea-393 soning tasks as described in §4.1. Fig. 3 de-394 picts the results with different ranks for both 395 LoRA and SHARELORA. Across all four orig-396 inal rank settings, SHARELORA consistently 397 improves performance compared to the baseline 398 LoRA. Despite using the same original rank, SHARELORA employs only 80% of the train-399 able parameters relative to the LoRA. For in-400 stance, SHARELORA achieves an average ac-401 curacy of 63.35% on the eight commonsense 402



Figure 3: Accuracy comparison with different rank values for LoRA and SHARELORA.

reasoning tasks with 5M trainable parameters, whereas LoRA requires 7M parameters for the same rank level.

#### 405 Similarity threshold.

406 Table 4 presents the experimental results of fine-tuning LLaMA-407 7B with different similarity thresholds  $\{0.75, 0.80, 0.85, 0.90\}$ . 408 The best performance is achieved when the similarity threshold 409  $\epsilon$  is set to 0.85. Higher similarity thresholds result in more el-410 igible shared layers. As the number of shared layers increases, 411 the rank of the LoRA modules can be adjusted upward within the constraints of the trainable parameter budget, which is de-412 termined by the original rank. Thus, there is a tradeoff between 413 having more independent layers and a larger rank. Conse-414 quently, selecting an optimal similarity threshold  $\epsilon$  is crucial 415 for balancing the trade-off between the number of shared layers 416 and the rank of LoRA modules. Higher thresholds allow more 417

Similarity threshold	Accuracy
0.75	61.80
0.80	62.81
0.85	63.10
0.90	61.80

Table 4: Accuracy comparison of different similarity thresholds evaluated with LLaMA-7B.

layers to share weights, permitting an increase in the rank within the limits of the trainable parameter budget. Therefore, the choice of  $\epsilon$  significantly influences the method's performance.

In practice, we typically choose  $\epsilon$  between 0.80 and 0.90, resulting in about one-third of the layers sharing weights with others. We can also run Algorithm 1 independently to obtain the shared dictionary, helping select an appropriate  $\epsilon$  before fine-tuning.

#### 425 Adapter sensitivity to different layer types.

426 In subsection 4.1, we adapt our proposed method to the 427 Q, K, V, Down, and Up weights. Table 5 presents the 428 performance of fine-tuning the weights associated with 429 each module type separately. We use r = 8 for the LoRA 430 modules. The results indicate that adapting our proposed 431 method to the Q and K weights yields the most significant

Target Modules	Accuracy
Q	64.19
Κ	64.31
V	63.71
Down	62.63
Up	62.91
Q,K,V,Down,Up	63.10

Table 5: Accuracy comparison of placing LoRA at different target modules.

benefits in fine-tuning. This improvement can be attributed to the critical role these weights play

in attention mechanisms, highlighting the importance of carefully selecting target modules for optimization.

#### 434 435 436

# 5 RELATED WORKS

437 Generative foundation models. Generative deep learning models pre-trained on large datasets are 438 called generative foundation models (Bommasani et al., 2021). These foundation models can be 439 applied to downstream tasks by fine-tuning. Advanced generative foundation models in natural lan-440 guage processing (NLP) such as GPT (Brown et al., 2020; Ouyang et al., 2022) and LLaMA (Touvron 441 et al., 2023a) model have shown great success in assisting and generating human-like text across 442 a wide range of topics. These generative language models can also be applied to many practical downstream tasks, such as education (Kasneci et al., 2023) and healthcare (Thirunavukarasu et al., 443 2023). Another kind of generative foundation model that has developed maturely is the diffusion 444 model (Ho et al., 2020; Rombach et al., 2022). The diffusion model works well in various computer 445 vision tasks such as text-to-image generation (Everaert et al., 2023) and image editing (Kawar et al., 446 2023). 447

448 Efficient fine-tuning methods. Efficient fine-tuning methods aim to reduce the number of trainable 449 parameters to save the GPU memory and training time during fine-tuning large-scale models. Some PEFT methods freeze most of the parameters in the model and only fine-tune specific modules, e.g., 450 BitFit (Zaken et al., 2021) fine-tunes only the bias of the model, which significantly saves the GPU 451 memory. However, it cannot be executed on models without bias parameters. (Houlsby et al., 2019) 452 and (Pfeiffer et al., 2020) add adapter layers between transformer blocks. These methods accelerate 453 fine-tuning by transferring knowledge from adapter layers pre-trained on general tasks. LoRA (Hu 454 et al., 2021; Liu et al., 2024b) is the most popular adapter for fine-tuning large foundation models. 455 It adopts the product of two small matrices to represent the full gradient during fine-tuning. It can 456 reduce the number of trainable parameters by 10,000 times and the GPU memory requirement by 3 457 times. Some adaptive algorithms work together with LoRA that can dynamically adjust the number 458 of trainable parameters to fit specific needs. For instance, AdaLoRA (Zhang et al., 2023) adaptively 459 allocates the trainable parameters to fit the GPU memory budget. These adaptive algorithms on LoRA 460 require heterogeneous LoRA configuration when implemented in federated fine-tuning.

# 6 FUTURE WORK

463 Our proposed SHARELORA method is orthogonal to other LoRA variants, suggesting that future work 464 could explore combining SHARELORA with existing PEFT methods. For example, LoRA+ (Hayou 465 et al., 2024) sets different learning rates for the adapter matrices A and B. By adjusting the 466 learning rate ratio between these two matrices, the efficiency and performance of fine-tuning can 467 be significantly improved. Specifically, LoRA+ sets the learning rate of B to be  $\lambda$  times that of A. This technique could easily be integrated with SHARELORA. Additionally, DoRA (Liu et al., 2024a) 468 is a novel PEFT method that enhances LoRA by decomposing pre-trained weights into magnitude 469 and direction components for more precise fine-tuning. This weight decomposition allows DoRA to 470 optimize the magnitude and direction of weights separately. SHARELORA could also be combined 471 with DoRA to share the magnitude and direction components among redundant layers. In summary, 472 SHARELORA not only provides a robust solution for fine-tuning but also offers the flexibility to 473 integrate with advanced PEFT methods, opening new avenues for research and potentially leading to 474 significant gains in PEFT performance and efficiency. 475

476 477

478

461

462

# 7 CONCLUSION

In this paper, we present SHARELORA, a parameter-efficient fine-tuning method that determines
which layers share the LoRA weights based on layer redundancy. SHARELORA leverages the
cosine similarity between each layer's representations to ascertain redundancy. Utilizing a greedy
algorithm, we maximize the sharing of LoRA weights while adhering to a predefined similarity
threshold. This approach effectively reduces the number of trainable parameters. We conduct
extensive experiments on large language models and multi-modal vision-language foundation models.
The results demonstrate that SHARELORA achieves comparable or superior performance to existing
PEFT methods, while using only 80% of the trainable parameter budget.

#### 486 REFERENCES 487

499

501

505

511

524

488	Armen Aghajanyan, Luke Zettlemoyer, and Sonal Gupta. Intrinsic dimensionality explains the
489	effectiveness of language model fine-tuning. arXiv preprint arXiv:2012.13255, 2020.

- 490 Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. Piqa: Reasoning about 491 physical commonsense in natural language, 2019. 492
- Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, 493 Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportuni-494 ties and risks of foundation models. arXiv preprint arXiv:2108.07258, 2021. 495
- 496 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, 497 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are 498 few-shot learners. Advances in neural information processing systems, 33:1877–1901, 2020.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi 500 Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*, 2022. 502
- Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina 504 Toutanova. Boolq: Exploring the surprising difficulty of natural yes/no questions, 2019.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and 506 Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. 507 arXiv preprint arXiv:1803.05457, 2018. 508
- 509 Dawei Dai, Liping Yu, and Hui Wei. Parameters sharing in residual neural networks. Neural 510 Processing Letters, 51(2):1393–1410, 2020.
- Miryam de Lhoneux, Johannes Bjerva, Isabelle Augenstein, and Anders Søgaard. Parameter sharing 512 between dependency parsers for related languages. arXiv preprint arXiv:1808.09055, 2018. 513
- 514 Mostafa Elhoushi, Akshat Shrivastava, Diana Liskovich, Basil Hosmer, Bram Wasti, Liangzhen Lai, 515 Anas Mahmoud, Bilge Acun, Saurabh Agarwal, Ahmed Roman, et al. Layer skip: Enabling early 516 exit inference and self-speculative decoding. arXiv preprint arXiv:2404.16710, 2024.
- 517 Martin Nicolas Everaert, Marco Bocchio, Sami Arpa, Sabine Süsstrunk, and Radhakrishna Achanta. 518 Diffusion in style. In Proceedings of the IEEE/CVF International Conference on Computer Vision, 519 pp. 2251-2261, 2023. 520
- 521 Siqi Fan, Xin Jiang, Xiang Li, Xuying Meng, Peng Han, Shuo Shang, Aixin Sun, Yequan Wang, and Zhongyuan Wang. Not all layers of llms are necessary during inference. arXiv preprint 522 arXiv:2403.02181, 2024. 523
  - Andrey Gromov, Kushal Tirumala, Hassan Shapourian, Paolo Glorioso, and Daniel A Roberts. The unreasonable ineffectiveness of the deeper layers. arXiv preprint arXiv:2403.17887, 2024.
- Demi Guo, Alexander M Rush, and Yoon Kim. Parameter-efficient transfer learning with diff pruning. 527 In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and 528 the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), 529 pp. 4884–4896, 2021. 530
- 531 Soufiane Hayou, Nikhil Ghosh, and Bin Yu. Lora+: Efficient low rank adaptation of large models. 532 arXiv preprint arXiv:2402.12354, 2024.
- Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. Towards a 534 unified view of parameter-efficient transfer learning. arXiv preprint arXiv:2110.04366, 2021. 535
- Shwai He, Guoheng Sun, Zheyu Shen, and Ang Li. What matters in transformers? not all attention is needed. arXiv preprint arXiv:2406.15786, 2024. 538
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. Advances in neural information processing systems, 33:6840–6851, 2020.

540	Neil Houlsby Andrei Giurgiu Stanislaw Jastrzebski Bruna Morrone Quentin De Laroussilhe
541	Andrea Commundo Mone Attention and Sulvine Cally. Demonster of Scient transfor learning for
540	Andrea Gesmundo, Mona Attariyan, and Sylvani Geny. Parameter-encient transfer learning for
542	nlp. In International Conference on Machine Learning, pp. 2790–2799. PMLR, 2019.
543	

- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- 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.
- Enkelejda Kasneci, Kathrin Seßler, Stefan Küchemann, Maria Bannert, Daryna Dementieva, Frank
   Fischer, Urs Gasser, Georg Groh, Stephan Günnemann, Eyke Hüllermeier, et al. Chatgpt for good?
   on opportunities and challenges of large language models for education. *Learning and individual differences*, 103:102274, 2023.
- Bahjat Kawar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri, and
   Michal Irani. Imagic: Text-based real image editing with diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6007–6017, 2023.
- Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. Similarity of neural network representations revisited. In *International conference on machine learning*, pp. 3519–3529.
   PMLR, 2019.
- <sup>561</sup> Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut.
   <sup>562</sup> Albert: A lite bert for self-supervised learning of language representations. In *International Conference on Learning Representations*, 2020.
- Jan Leike, David Krueger, Tom Everitt, Miljan Martic, Vishal Maini, and Shane Legg. Scalable agent alignment via reward modeling: a research direction. *arXiv preprint arXiv:1811.07871*, 2018.
- Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. *arXiv preprint arXiv:2101.00190*, 2021.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning, 2023.
- Shih-Yang Liu, Chien-Yi Wang, Hongxu Yin, Pavlo Molchanov, Yu-Chiang Frank Wang, Kwang Ting Cheng, and Min-Hung Chen. Dora: Weight-decomposed low-rank adaptation. *arXiv preprint arXiv:2402.09353*, 2024a.
- Zeyu Liu, Souvik Kundu, Anni Li, Junrui Wan, Lianghao Jiang, and Peter Anthony Beerel. Aflora:
   Adaptive freezing of low rank adaptation in parameter efficient fine-tuning of large models. *arXiv* preprint arXiv:2403.13269, 2024b.
- AI Meta. Introducing meta llama 3: The most capable openly available llm to date. *Meta AI*., 2024.
- Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct
  electricity? a new dataset for open book question answering. *arXiv preprint arXiv:1809.02789*, 2018.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35: 27730–27744, 2022.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
   Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
- Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Alessandro Cappelli,
   Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. The refinedweb
   dataset for falcon llm: outperforming curated corpora with web data, and web data only. *arXiv* preprint arXiv:2306.01116, 2023.

594 595 596	Jonas Pfeiffer, Aishwarya Kamath, Andreas Rücklé, Kyunghyun Cho, and Iryna Gurevych. Adapter- fusion: Non-destructive task composition for transfer learning. <i>arXiv preprint arXiv:2005.00247</i> , 2020.
597 598 599 600	Ofir Press and Lior Wolf. Using the output embedding to improve language models. In <i>Proceedings</i> of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers, pp. 157–163, 2017.
601 602 603 604	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In <i>International conference on machine learning</i> , pp. 8748–8763. PMLR, 2021.
606 607 608	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>Journal of machine learning research</i> , 21(140):1–67, 2020.
609 610 611 612	David Raposo, Sam Ritter, Blake Richards, Timothy Lillicrap, Peter Conway Humphreys, and Adam Santoro. Mixture-of-depths: Dynamically allocating compute in transformer-based language models. <i>arXiv preprint arXiv:2404.02258</i> , 2024.
613 614 615	Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High- resolution image synthesis with latent diffusion models. In <i>Proceedings of the IEEE/CVF confer-</i> <i>ence on computer vision and pattern recognition</i> , pp. 10684–10695, 2022.
617 618 619	Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale. <i>Communications of the ACM</i> , 64(9):99–106, 2021.
620 621	Maarten Sap, Hannah Rashkin, Derek Chen, Ronan LeBras, and Yejin Choi. Socialiqa: Commonsense reasoning about social interactions. <i>arXiv preprint arXiv:1904.09728</i> , 2019.
623 624 625	Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, et al. Large language models encode clinical knowledge. <i>arXiv preprint arXiv:2212.13138</i> , 2022.
626 627 628 629	Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Le Hou, Kevin Clark, Stephen Pfohl, Heather Cole-Lewis, Darlene Neal, et al. Towards expert-level medical question answering with large language models. <i>arXiv preprint arXiv:2305.09617</i> , 2023.
630 631	Sho Takase and Shun Kiyono. Lessons on parameter sharing across layers in transformers. <i>arXiv</i> preprint arXiv:2104.06022, 2021.
632 633 634	MosaicML NLP Team. Introducing mpt-7b: A new standard for open-source, commercially usable llms, 2023. URL www.mosaicml.com/blog/mpt-7b. Accessed: 2023-05-05.
635 636 637	Arun James Thirunavukarasu, Darren Shu Jeng Ting, Kabilan Elangovan, Laura Gutierrez, Ting Fang Tan, and Daniel Shu Wei Ting. Large language models in medicine. <i>Nature medicine</i> , 29(8): 1930–1940, 2023.
639 640 641	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> , 2023a.
642 643 644	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> , 2023b.
646 647	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Huggingface's transformers: State-of-the-art natural language processing. <i>arXiv preprint arXiv:1910.03771</i> , 2019.

648 649 650	Ge Zhang Yao Fu Wenhao Huang Huan Sun Yu Su Wenhu Chen Xiang Yue, Xingwei Qu. Mammoth: Building math generalist models through hybrid instruction tuning. <i>arXiv preprint</i> arXiv:2309.05653, 2023
651	······································
652	Elad Ben Zaken, Shauli Rayfogel, and Yoay Goldberg, Bitfit: Simple parameter-efficient fine-tuning
653	for transformer-based masked language-models. arXiv preprint arXiv:2106.10199, 2021.
654	Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine
655	really finish your sentence? In Proceedings of the 57th Annual Meeting of the Association for
656	Computational Linguistics, pp. 4791–4800, 2019.
657	
658	Qingru Zhang, Minshuo Chen, Alexander Bukharin, Pengcheng He, Yu Cheng, Weizhu Chen,
659	and Tuo Zhao. Adaptive budget allocation for parameter-efficient fine-tuning. In <i>The Eleventh</i>
660	International Conference on Learning Representations, 2023. URL https://openreview.net/
661	forum/ld=lq6/uwkJjlY.
662	Haiyan Zhao, Hanije Chen, Fan Yang, Ninghao Liu, Huigi Deng, Hengyi Cai, Shuaigiang Wang,
663	Dawei Yin, and Mengnan Du. Explainability for large language models: A survey. ACM
664	Transactions on Intelligent Systems and Technology, 15(2):1–38, 2024.
665	
666	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,
667	Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and
822	chatbot arena. Advances in Neural Information Processing Systems, 36, 2024.
000	Chunting Zhou, Pengfei Liu, Puxin Xu, Sriniyasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia
670	Efrat, Ping Yu, Lili Yu, et al. Lima: Less is more for alignment. Advances in Neural Information
671	Processing Systems, 36, 2024.
670	
672	Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul
674	Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. arXiv
675	preprint arXiv:1909.08593, 2019.
676	
677	
678	
679	
680	
681	
682	
683	
684	
685	
686	
687	
688	
689	
690	
691	
692	
693	
694	
695	
696	
697	
698	
699	