ROLORA: RANK OPTIMIZATION FOR LOW-RANK ADAPTATION UNDER MEMORY CONSTRAINTS

Anonymous authors

Paper under double-blind review

ABSTRACT

Low-Rank Adaptation (LoRA) has emerged as a prominent technique for finetuning large language models (LLMs) with limited computational resources. However, by injecting low-rank adapters with a rank identical across all layers, standard LoRA overlooks the varying importance of the weight matrices, often leading to suboptimal performance. Therefore, discovering an optimal rank configuration that efficiently utilizes limited training resources remains an open question. Existing solutions typically compromises computational constraints for performance gains, limiting their practical usage in resource-constrained scenarios. To address these issues, in this paper, we propose a novel method named ROLoRA to efficiently discover an effective rank configuration for low-rank adaptation, while strictly adhering to a constrained computational budget during training. In particular, our method iteratively prunes saturated adapters and expands underfitted ones to increase their capacity until they converge to a highly optimized configuration. Our approach is delicately designed within the Frank-Wolfe algorithmic framework, which offers potential theoretical guarantees. Experimentally, we demonstrate that ROLoRA outperforms standard LoRA on common natural language processing tasks, including the GLUE and SQuAD benchmarks. Additionally, we provide a comprehensive analysis to explain why ROLoRA surpasses competing state-of-the-arts.

028 029

031

051 052

004

006

008 009

010 011

012

013

014

015

016

017

018

019

020

021

022

024

025

026

027

1 INTRODUCTION

Fine-tuning pre-trained large language models (LLMs) for downstream tasks has become a common practice to meet customized and domain-specific demands, especially after the recent prevalence of these models. However, the computational cost during training remains a significant concern, particularly due to the increasing number of parameters in these models. Naively fine-tuning the entire models often requires as much computational capacity as training the model from scratch, which limits the feasibility of fine-tuning LLMs on resource-constrained edge devices.

038 To fully exploit the advantages of model fine-tuning, recent works have departed from full fine-039 tuning towards more light-weighted approaches, also known as Parameter-Efficient Fine-Tuning 040 methods (or PEFT) (Houlsby et al., 2019; Karimi Mahabadi et al., 2021; Mao et al., 2022). PEFT 041 methods aim at reducing the number of trainable parameters during the fine-tuning process, which 042 brings about simultaneously two benefits. First, PEFT methods reduce the memory required to store 043 trained parameters, enabling the deployment of fine-tuned models for multiple tasks on the same 044 device. Second, the reduced number of trainable parameters typically lowers the optimization cost during training, particularly by minimizing the memory required to store optimizer's states, which facilitates on-device fine-tuning for these models. 046

Among many recently developed PEFT methods, Low-Rank Adaptation, or LoRA, proposed by Hu
et al. (2022) is among the most important. In principle, LoRA freezes the base pre-trained model
and attach updates to each weight matrix of the base model as separate modules, known as adapters,
in form of a product of two matrices of significantly smaller size. Concretely,

$$W = W_0 + \Delta W = W_0 + BA$$

where $W_0 \in \mathbb{R}^{n \times d}$ is the base pre-trained weight, $\Delta W = BA$ is the update, $A \in \mathbb{R}^{r \times d}$ is called the down projection matrix and $B \in \mathbb{R}^{n \times r}$ the up projection matrix. Typically the size r is chosen to be a small constant significantly smaller than n, d. In this way, LoRA is able to reduce the number of trainable parameters to possibly only 0.5% compared with full fine-tuning while achieves comparable or even better than the latter (Hu et al., 2022; Zhang et al., 2023).

057 As simple and efficient as it may seem, LoRA still has its limitations. In particular, it is unclear how 058 to optimally choose the rank r for each adapter in LoRA. The standard LoRA sets an identical rank r for all adapters. This strategy ignores the fact that different layers in the model carry different 060 amount of knowledge (Chen et al., 2023a) and thus require different level of tuning. In a recent 061 work, Zhang et al. (2023) demonstrate that varying the ranks r according to the importance of the 062 weight matrices can improve the performance of LoRA-style methods. However, existing solutions 063 often require a greater memory budget during training to provide more buffer to explore a high-064 performing adapter configuration. This is especially undesirable for fine-tuning the model on small devices where the memory allowed during training is strictly limited. 065

- In this work, we study the problem of determining the ranks of the adapters that can improve the
 performance of LoRA while imposing a strict memory constraint during training, reflected by a con straint on the number of trainable parameters. Towards answering this question, our contributions
 can be summarized as follows.
- We propose a general framework, named ROLoRA, that iteratively sparsifies and grows the ranks of the adapters. Specifically, at each iteration, our algorithm performs a flexible rank sparsification step to determine the importance of the adapted weight matrices. Our framework allows to use any rank sparsification algorithm as a blackbox. Next, we expand the rank configuration to increase the capacity of the adapters, while strictly adhering to memory budget, before updating the ranks following the direction suggested this configuration. Our approach aims to imitate a first order gradient method (specifically Frank-Wolfe algorithm), which is search-free and highly efficient.
- To demonstrate its effectiveness, we show in experiment that our algorithm improves the performance of LoRA on common NLP benchmarks, including GLUE (Wang et al., 2019) and SQUAD datasets (Rajpurkar et al., 2016; 2018). We show that our algorithm even outperforms AdaLoRA which requires more memory during training. We show at the same time that the output adapters can be smaller in size, reducing the memory storage compared with the standard LoRA.

Finally, to explain the effectiveness of our method, we conduct a comprehensive ablation study. We demonstrate that a certain type of weight matrices in the transformer architecture, such as value matrices, plays a more important role than the others (key and query matrices) in fine-tuning. Mean-while, our approaches could effectively discover the varying saliency pattern and assign proper ranks accordingly. These advantages contribute to the superior performance we achieved.

- 087 088
- 088

2 RELATED WORK

090 091 092

094

095

096

097

098

Parameter-Efficient Fine-Tuning. Together with model compression (Jafari et al., 2021; Chen et al., 2021), PEFT is another category of solutions to improving fine-tuning LLMs. One approach to PEFT is adding adapters between model layers (Houlsby et al., 2019; Rebuffi et al., 2017; Pfeiffer et al., 2021; He et al., 2022). Due to the increase in the number of layers, this approach, however, can introduce latency during inference time. Another direction to PEFT is to directly update the pre-trained weights, under which falls LoRA (Hu et al., 2022). We refer the reader to a recent survey on PEFT methods by Han et al. (2024) for a more comprehensive comparison.

099 LoRA and related works. We provide backgrounds on LoRA in Section 3. This seminal work 100 by Hu et al. (2022) opens up a plethora of works on PEFT via low-rank adaptation, focusing on ad-101 dressing shortcomings of LoRA. For example, QLoRA (Dettmers et al., 2024) employs quantization 102 technique to improve the memory consumption of LoRA. Liu et al. (2024) uses a weight decom-103 position of the pre-trained weight matrices to improve the performance of LoRA. Notably, Hayou 104 et al. (2024) show that using different step sizes for different adapters in LoRA can lead to better 105 performances. This latter work also highlight the need for algorithmic frameworks that exploit the layer-wise varying importance which our work targets to develop. However, these works are orthog-106 onal and can be used jointly for further improvements, which we leave for future investigation. Our 107 work is inspired by Xia et al. (2024), which uses multiple runs of LoRA to improve this algorithm.

108 Most related to our work is DyLoRA (Valipour et al., 2023), AdaLoRA Zhang et al. (2023) and 109 SoRA (Ding et al., 2023). DyLoRA focuses on determining an optimal constant for setting the rank 110 of *all* adapters in LoRA. This is in contrast to our work which aims at finding ranks for adapters that 111 can vary across layers when given a budget on the number of trainable parameters. AdaLoRA and 112 SoRA use different techniques to reduce (sparsify) the ranks initialized by LoRA, which in effect, gives us a rank configuration that reflects the importances of the adapted weight matrices. However, 113 both methods violate the budget constraint when starting with a rank higher than the one initialized 114 by LoRA. Our solution to the problem uses a sparsification method such as AdaLoRA and SoRA as 115 a subroutine in an iterative framework, while strictly guarantees the budget constraint is satisfied. 116

117

120

121

3 BACKGROUND AND PROBLEM STATEMENT

3.1 LOW-RANK ADAPTATION

Low-rank adaptation, or LoRA, was introduced by Hu et al. (2022), as a Parameter-Efficient Fine-Tuning method, based on the assumption that the incremental update to the pretrained model has a low intrinsic rank. Specifically, given the pretrained weight $W_0 \in \mathbb{R}^{n \times d}$, LoRA considers only updates in the form of a product of two small matrices $\Delta W = BA$, where $B \in \mathbb{R}^{n \times r}$ and $A \in \mathbb{R}^{r \times d}$ with the rank $r \ll \min\{n, d\}$. We will refer to r as the rank of the adapter, although r is only an upper bound on the rank of A or B. For $h = W_0 x$, the forward pass becomes

$$h = (W_0 + \Delta W)x = W_0x + BAx.$$

In this way, the number of trainable parameters is significantly reduced from nd to r(n+d). For example, when the model is trained with commonly used adaptive optimizers such as Adam (Kingma, 2014), the memory required to store the optimizer's states, *e.g.*, momentums, is proportional to the number of trainable parameters. By reducing the number of trainable parameters, LoRA saves on the memory required to fine-tune the model.

134 In this work, we study a more generalized variation of LoRA by introducing a diagonal matrix to the 135 update process. Let $g \in \mathbb{R}^r$ be an r-dimensional vector and $\operatorname{diag}(g) \in \mathbb{R}^{r \times r}$ be a diagonal matrix 136 formed by q. We consider an update of the form $\Delta W = B \operatorname{diag}(q) A$. Remark here that when q is 137 set to 1, this variant becomes numerically equivalent to the standard LoRA. When g is made train-138 able, the number of trainable parameters increases marginally by r, which is negligible compared to 139 the total number of trainable variables in LoRA, *i.e.*, r(n+d). This generic adjustment offers greater flexibility during training and can often lead to a more favorable optimization landscape, potentially 140 resulting in improved performance. 141

142 143

144

145

146 147

148

149

154 155

158 159

3.2 PROBLEM STATEMENT

Suppose that we aim to use K LoRA adapters for K pretrained weight matrices, with corresponding ranks r_1, \ldots, r_K . Our research question is as follows:

Given a strict budget constraint on the memory required during training, how can we determine the optimal assignment of ranks for these adapters?

Formally, let us denote the rank configuration $\mathcal{R} := (r_1, r_2, \cdots, r_K)$, and the parameters to be finetuned carried out by \mathcal{R} as $\theta_{\mathcal{R}}$. Let $L_{W_0}(\theta_{\mathcal{R}}; \mathcal{D})$ be the loss function we optimize during the training with regard to parameters $\theta_{\mathcal{R}}$, given the pretrained weight W_0 , and training data \mathcal{D} . We let $\theta_{\mathcal{R}}^*$ be a minimizer to $L_{W_0}(\theta_{\mathcal{R}}; \mathcal{D})$, *i.e.*,

$$\theta_{\mathcal{R}}^* \in \arg\min_{\theta_{\mathcal{P}}} L_{W_0}(\theta_{\mathcal{R}}; \mathcal{D}).$$
(1)

We measure the performance of the rank configuration $\mathcal{R} = (r_1, r_2, \cdots, r_K)$ by evaluating the objective loss values of $\theta_{\mathcal{R}}^*$ measured over the dataset \mathcal{D} as follows:

$$f(\mathcal{R}) := L_{W_0}(\theta_{\mathcal{R}}^*; \mathcal{D}). \tag{2}$$

To express the memory constraint, for simplicity, we assume that all pretrained weight matrices used for adaptation have the same dimensions $n \times d$. This assumption typically holds in standard LoRA, where fine-tuning is applied to the key, query, and value matrices, *i.e.*, W_k , W_q , W_v of the



Figure 1: Visualization of an iteration in ROLoRA. Starting the iteration with an initial rank configuration (top left), the algorithm does the following: (a) Sparsify by pruning saturated adapters (grayed out in top right picture); (b) Extend the sparsified configuration; (c) Interpolate the ranks at the beginning of the iteration (top left) and the extended ranks (bottom right) to obtain the new configuration (bottom left); (d) Carry the trained weight during sparification (orange) to the new configuration and proceed to the next iteration.

self-attention modules. However, we can generalize this assumption to matrices of any size. Under this assumption, we deduce that the memory required to maintain the optimizer's states depends linearly on $\sum_{i=1}^{K} r_i$. Therefore, we express the memory constraint as:

$$\sum_{i=1}^{K} r_i \le K \cdot R,\tag{3}$$

where R is the average rank we can afford with the memory budget for each adaptor. We obtain the following main optimization problem to seek the optimal adapter configuration under strict budget:

$$\underset{\mathcal{R}}{\text{minimize } f(\mathcal{R}), \text{ subject to } \sum_{i=1}^{K} r_i \leq K \cdot R.$$
(4)

Remark 1. Standard LoRA typically uses a small constant (*e.g.*, 4, 8, 16) as the value of R and applies it uniformly across all adapters, *i.e.*, $r_1 = \cdots = r_K = R$. However, it neglects the varying importance of different weight matrices and has been shown to result in inferior performance compared to methods like AdaLoRA (Zhang et al., 2023), which adaptively vary the rank assignments.

Remark 2. To determine the optimal rank configuration, existing methods like AdaLoRA (Zhang et al., 2023) iteratively prune redundant adapter variables until convergence. However, these methods require starting with larger adapter ranks, which directly violates the memory constraint outlined in (3) during training. In contrast, our method adheres to the memory constraint from the outset and effectively discovers the optimal rank configuration, provides significant practical advantages under limited resource situations and distinguishes it from existing approaches.

Efficiently solving problem (4) is challenging for several reasons. First, accurately evaluating the performance measure $f(\mathcal{R})$ is computationally expensive. Second, standard optimization tools, such as gradient methods, struggle with discrete variables, while naive search is prohibitively expensive due to the size of the search space. In the next section, we will introduce a simple but effective method that is delicately designed to solve problem (4) efficiently.

215 4 Algorithm

1:	Input : Target average rank R
2:	Initialize $r_1 = \dots r_K = R, \gamma = 0$
3:	Initialize the adapters, initialize μ the parameter of SPARSIFY algorithm according to
	SCHEDULE algorithm
4:	for iteration $t = 1 \dots T$ do
5:	Let $(s_1, \ldots, s_K), \theta \leftarrow \text{SPARSIFY}((r_1, \ldots, r_K); \theta; \mu)$ (e.g.via AdaLoRA)
6:	Let $S = \frac{1}{K}(s_1 + \dots + s_K); (s_1, \dots, s_K) \leftarrow \text{round} ((s_1, \dots, s_K) \times \frac{R}{S})$
7:	Update $r_i = (1 - \gamma)r_i + \gamma s_i$, for all $i \in [K]$ for $\gamma = 1/t$
8:	Update current model by EXTEND $((r_1, \ldots, r_K); \theta)$
9:	Update $\mu \leftarrow \text{SCHEDULE}(\mu; t)$
10:	end for
11:	return (r_1,\ldots,r_K)

229

In this section, we describe our algorithm ROLoRA to solve Problem (4) shown in Algorithm 1.
ROLoRA is an iterative method to automatically discover an effective rank configuration. After
proper initialization, ROLoRA begins with a given adapter rank configuration, then iteratively identifies saturated and underfitting weights. It adjusts the adapter rank assignment by performing pruning
and growth operations, and is delicately structured within the Frank-Wolfe algorithmic framework
to leverage potential theoretical advantages. Each main step is elaborated as follows, which is also
demonstrated in Figure 1.

238 **Sparsifying / pruning the ranks** (Step (a) in Figure 1). Given a rank configuration \mathcal{R} = 239 (r_1,\ldots,r_K) , certain weight matrices may become saturated during further fine-tuning, meaning 240 they might not require as many trainable parameters. To address this, we introduce a sparsification 241 or pruning operator to identify redundant ranks across the adapter sets. In Algorithm 1, we flexi-242 bly employ any rank sparsification algorithm for the SPARSIFY operator, such as AdaLoRA (Zhang et al., 2023) and SoRA (Ding et al., 2023) or pruning approaches based on saliency scores (Chen 243 et al., 2023b; 2024). Without loss of generality, we assume that SPARSIFY takes in a pruning aggres-244 siveness hyper-parameter μ to produce a new trial rank configuration $\mathcal{S} = (s_1, \dots, s_K)$. \mathcal{S} reflects 245 more accurately than \mathcal{R} the (relative) importance of the corresponding weight matrices (Line 5). 246

247 **Reassign and grow the ranks** (Step (b) in Figure 1). In addition to certain adapters requiring fewer 248 ranks due to saturation, some adapters may need their capacity increased instead to learn more effectively. Therefore, it is essential to identify which adapters should grow to meet this demand. 249 To achieve this, we use the signals from the pruned rank configuration. Typically, the ranks that 250 remain unchanged in S indicate that these adapters require more parameters, as no redundancy is 251 detected. However, because we operate under a strict memory budget, the amount of growth must 252 strictly adhere to this constraint. Taking these factors into account, we scale up S to align with the 253 target average rank R (Line 6). 254

Intuitive connection to gradient estimation to f. The sparsification and growth steps in our algorithm effectively identify adapters that require more capacity or exhibit redundancy. Consequently, it builds on the intuition that the new rank configuration output S suggests a favorable search direction that can enhance the performance measure $f(\mathcal{R})$ of the current rank configuration \mathcal{R} .

259 Update the ranks and adapters (Step (c-d) in Figure 1). Once obtaining the new trial configuration 260 S, we update the current configuration \mathcal{R} by moving it towards S, *i.e.*, taking an interpolation 261 step with step size γ (Step (c)). This step is a reminiscence of the update step in the Frank-Wolfe algorithm (Frank et al., 1956) for solving constrained optimization problems. To update the adapter 262 weights, we can carry over the trained adapters obtained from the SPARSIFY algorithm (Step (d)). 263 Since the some of trial ranks in S are increased during the rank growth step, additional rows and 264 columns will be added to the up and down projection matrices, as well as additional entries to the 265 diagonal matrix, requiring only the initialization of these extended parts. These extended parts are 266 initialized in the same way as LoRA adapters so that their product is 0 to ensure computational 267 invariance after the adapter growth. 268

Scheduling the sparsification. After each sparsification-growth-update step, we update the sparsification parameter of the SPARSIFY operator via SCHEDULE. One approach is to keep the sparsifi-

cation parameter as a constant. Alternatively, we can employ a sparsification schedule that gradually
 reduces the pruning aggressiveness over time. Intuitively, as the model converges to a better rank
 configuration over time, it requires fewer adjustments to determine an update direction.

273 274

275

5 OUTLINE FOR CONVERGENCE ANALYSIS

In this section, we present the outline for the convergence analysis of ROLoRA to solve the target problem (4). Towards this end, we first show that under certain assumptions on the SPARSIFY operator, Algorithm 1 will improve the performance measure f over iterations.

Assumption 1. The SPARSIFY($\mathcal{R}; \theta; \mu$) operator has the following properties: 1. During the course of the algorithm, SPARSIFY satisfies the memory budget, indicated by $\sum_i r_i$. 2. For any rank configuration \mathcal{R} and initial weight θ , there exists a sparsification parameter μ such that if $\mathcal{S}, \theta \leftarrow$ SPARSIFY($\mathcal{R}; \theta; \mu$) then $f(\mathcal{S}) \leq f(\mathcal{R})$.

284 We interpret this assumption as follows. As a rank configuration can contain redundant information, 285 the SPARSIFY operator with a well chosen sparsification parameter can remove this redundancy 286 without sacrificing the model performance. In the trivial case, we can choose μ so that no sparsifica-287 tion happens, thus the loss function is not increased (yet in which case the algorithm will terminate 288 without updating the rank configuration). In practice, if the model has saturated weights, keeping 289 fine-tuning these weights can easily result in overfitting problem, which further degrades the model performance. Hence, yielding proper sparsity over the adapters can promote the model performance 290 instead, *i.e.*, $f(\mathcal{S}) \leq f(\mathcal{R})$. 291

We then show the following Proposition.

Proposition 1. Let $\mathcal{R}^{(0)} = (R, ..., R)$, and $\mathcal{R}^{(t)}$ be the rank configuration output by Algorithm 1 after iteration t. For all $t \ge 1$, $f(\mathcal{R}^{(t)}) \le f(\mathcal{R}^{(0)})$.

296 297

298 299 300

301

Proof. Let $\mathcal{S}^{(t)}$ be the sparsified ranks from $\mathcal{R}^{(t)}$ and $\theta^*_{\mathcal{S}^{(t)}}$ be a collection of optimal adapters weights for $\mathcal{S}^{(t)}$. We apply the EXTEND operator and initialization so that for each triple of adapters $(B^{(t)}, g^{(t)}, A^{(t)}) \in \theta^*_{\mathcal{S}^{(t)}}$ we have $B^{(t+1)} \operatorname{diag}(g^{(t+1)}) A^{(t+1)} = B^{(t)} \operatorname{diag}(g^{(t)}) A^{(t)}$. Therefore, we have $f(\mathcal{R}^{(t+1)}) \leq f(\mathcal{S}^{(t)}) \leq f(\mathcal{R}^{(t)})$, where the last inequality is by the assumption.

302 Connection to Frank-Wolfe algorithm and Convergence rate. We now show the connection 303 between Algorithm 1 and Frank-Wolfe algorithm Frank et al. (1956). To minimize a function f over 304 a convex domain \mathcal{X} , in iteration t, Frank-Wolfe algorithm finds a solution in $s \in \mathcal{X}$ that minimizes $\langle s, \nabla f(x^{(t)}) \rangle$, with $x^{(t)}$ being the current solution, and update $x^{(t+1)} \leftarrow (1-\gamma)x^{(t)} + \gamma s$ with some step size γ . Jaggi (2013) shows that Frank-Wolfe with appropriate step sizes converges with rate $O(\frac{1}{T})$ for convex smooth functions; later Lacoste-Julien (2016) shows a $O(\frac{1}{\sqrt{T}})$ convergence rate of 305 306 307 308 the duality gap for nonconvex smooth functions. Our algorithm is driven by this Frank-Wolfe type 309 of update. Our problem (4) while being a discrete problem can be relaxed to a continuous convex 310 domain¹. Although we do not have access to the exact gradient $\nabla f(\mathcal{R})$, we approximate the update 311 direction s using the dedicated designed joint sparsification and growth step as the proxy described in Section 4. Consequently, we could incorporate with the existing theoretical framework of Lacoste-312 Julien (2016) to indicate that within a finite number of iterations, the proposed ROLoRA could find 313 a high-performing rank configuration, which is further numerically demonstrated in Section 6. 314

315 316

317

6 EXPERIMENTS

We implement ROLoRA (Algorithm 1) for fine-tuning RoBERTa-base (Liu et al., 2019) and DeBERTa-v3-base (He et al., 2021) on the GLUE Benchmark (Wang et al., 2019) and SQUAD datasets (SQUADv1 (Rajpurkar et al., 2016) and SQUADv2 Rajpurkar et al. (2018)). We describe the Algorithm implementation and baselines below.

322 323

¹For any two feasible \mathcal{R}^1 and \mathcal{R}^2 satisfying the budget constraint, their convex combination still satisfies the budget constraint since $\sum_{i=1}^{K} (\lambda r_i^1 + (1-\lambda)r_i^2) \leq \lambda K \cdot R + (1-\lambda)K \cdot R = K \cdot R$ for any $\lambda \in [0,1]$.

324 Table 1: Results for RoBERTa-base and GLUE benchmark. The target average rank is R = 8. For 325 AdaLoRA, the initial average rank is set to 12. \dagger represents the vanilla version of the algorithm. \star 326 means the algorithm is executed in T = 3 iterations; in each iteration, the model is initialized to the best checkpoint. We report Pearson correlation for STSB, Matthews correlation for CoLA and 327 Accuracy for the remaining datasets. We report the mean and standard deviation across 5 runs. 328

Method	MRPC	RTE	CoLA	QNLI	MNLI	QQP	SST2	S
$LoRA^{\dagger}$	$89.90_{\pm.73}$	$78.70_{\pm 1.21}$	$64.34_{\pm.78}$	$93.08_{\pm.20}$	$87.57_{\pm.11}$	$90.92_{\pm.10}$	$95.34_{\pm.27}$	90.
LoRA*	$89.36_{\pm.80}$	$79.21_{\pm 1.30}$	$64.91_{\pm.95}$	$93.22_{\pm.18}$	$87.60_{\pm .09}$	$91.21_{\pm .09}$	$95.62_{\pm.21}$	90.9
AdaLoRA [†]	$88.53_{\pm.90}$	$76.10_{\pm.42}$	$58.87_{\pm.46}$	$93.56_{\pm.07}$	$87.57_{\pm.06}$	$90.03_{\pm .04}$	$94.84_{\pm.07}$	90.
AdaLoRA*	$89.01_{\pm .36}$	$79.42_{\pm .65}$	$62.62_{\pm.81}$	$\textbf{93.72}_{\pm.02}$	$87.66_{\pm.06}$	$90.42_{\pm .05}$	$95.18_{\pm.21}$	90.
ROLoRA	$90.15_{\pm.33}$	$81.73_{\pm.37}$	64.98 ±1.20	93.71 ±.07	87.76 _{±.08}	$91.25_{\pm.04}$	$95.30_{\pm.15}$	91.

Table 2: Results for DeBERTa-base and GLUE benchmark. The target average rank is R = 8. For AdaLoRA, the initial average rank is set to 12.

Method	MRPC	RTE	CoLA	QNLI	MNLI	QQP	SST2	STSB
LoRA [†]	$90.74_{\pm.39}$	$86.43_{\pm.98}$	$71.16_{\pm 1.49}$	$94.24_{\pm.27}$	$90.32_{\pm.06}$	$92.23_{\pm.04}$	$96.54_{\pm.09}$	$90.75_{\pm.16}$
LoRA*	$90.88_{\pm.42}$	$87.65_{\pm 1.03}$	$70.95_{\pm.45}$	$94.40_{\pm.14}$	$90.39_{\pm.04}$	$92.34_{\pm.06}$	$96.51_{\pm.21}$	$90.79_{\pm.19}$
AdaLoRA [†]	$91.08_{\pm.20}$	$87.22_{\pm.37}$	$70.39_{\pm.59}$	$94.70_{\pm.11}$	$90.76_{\pm.04}$	$91.89_{\pm.02}$	$96.24_{\pm.17}$	$91.17_{\pm.15}$
AdaLoRA*	$91.27_{\pm.12}$	$87.80_{\pm.27}$	$\textbf{71.78}_{\pm.73}$	$\textbf{94.72}_{\pm.03}$	$\textbf{90.93}_{\pm.05}$	$92.17_{\pm.03}$	$96.61_{\pm.16}$	$91.44_{\pm.13}$
ROLoRA	$91.72_{\pm.65}$	$\textbf{88.66}_{\pm 1.30}$	$71.17_{\pm.61}$	$\textbf{94.72}_{\pm.07}$	$90.78_{\pm.05}$	$\textbf{92.38}_{\pm.05}$	$96.70_{\pm.22}$	$\textbf{91.48}_{\pm.18}$

Models. We use pre-trained RoBERTa-base (Liu et al., 2019) and DeBERTa-v3-base (He et al., 2021) publicly available on HuggingFace Transformers library (Wolf, 2019). For the former, we fine-tune query and value matrices. For the latter, we fine-tune key, query and value matrices.

Sparsification Algorithm and Scheduling. We use AdaLoRA Zhang et al. (2023) as the SPARSIFY 353 operator. We start with the level of sparsify that reduces the average rank to R/2 and gradually 354 increase the average rank after sparsification by a factor of $\left(1 + \frac{1}{T+t}\right)$ after iteration t. 355

356 **Number of iterations.** We use a small number of iterations to balance the tradeoff between the 357 improvement in the model performance and the increase in the training time. In all experiments, we 358 use three iterations (T = 3). 359

Baselines. We compare our algorithm with two algorithms LoRA (Hu et al., 2022) and 360 AdaLoRA(Zhang et al., 2023). For LoRA, we use the same target average rank R across all adapters. 361 For AdaLoRA, we initialize the rank configuration the same across all adapters to be 1.5R as rec-362 ommended in Zhang et al. (2023) and sparsify it to the target average rank R. However, one should 363 also note that by initializing the total ranks to 1.5RK, this baselines violates the memory constraint. 364 We also compare with two other baselines by repeating LoRA and AdaLoRA T iterations. In each 365 iteration, the model is initialized to the best checkpoint so far. 366

367 368

337 338

350

351

352

6.1 GENERAL LANGUAGE UNDERSTANDING EVALUATION (GLUE) BENCHMARK

369 Datasets. The GLUE Benchmark is a collection of natural language understanding tasks, consisting 370 of two single-sentence classification tasks, three similarity and paraphrase tasks and four natural 371 language inference tasks.

372 **Implementation details.** In our experiment, the target average rank is set to R = 8. The initial aver-373 age rank for AdaLoRA is set to 12. We reuse hyper-parameters such as learning rate, sparsification 374 parameters, etc, as recommended in the original papers by Hu et al. (2022); Zhang et al. (2023). 375

Main results. We report the performances of our algorithm in comparison with the baselines using 376 RoBERTa-base (Liu et al., 2019) in Table 1 and DeBERTa-v3-base (He et al., 2021) in Table 2. 377 For both models, almost across the board, ROLoRA outperforms both vanilla LoRA and AdaLoRA

R = 16R = 4Method MRPC RTE CoLA MRPC RTE CoLA LoRA[†] $88.77_{\pm.57}$ 77.77_{±1.34} 63.21_{±.94} 89.31_{±.59} 79.23_{±1.22} 64.74_{±.53} LoRA* $89.16_{\pm.73}$ 77.83 $_{\pm1.24}$ 63.95 $_{\pm.86}$ 89.85 $_{\pm.65}$ 80.43 $_{\pm.77}$ 65.99 $_{\pm.99}$ AdaLoRA[†] 87.01_{±.31} 75.67_{±.67} 58.50_{±.71} 88.33_{±.74} 77.55_{±2.11} 59.59_{±.44} AdaLoRA* $89.71_{\pm.35}$ $79.49_{\pm1.24}$ $61.51_{\pm.62}$ $89.26_{\pm.39}$ $80.43_{\pm.98}$ $63.07_{\pm.30}$ $90.49_{\pm.72}$ 80.36 $_{\pm1.32}$ 63.98 $_{\pm.83}$ 90.00 $_{\pm.33}$ 82.38 $_{\pm.48}$ 65.57 $_{\pm1.26}$ ROLoRA

Table 3: Results for RoBERTA-base on MRPC, RTE and CoLA when the target average rank R = 16 and R = 4. In the former case, the initial rank of AdaLoRA is set to 24, and in the latter 6.

Table 4: Average rank of the output adapters for Experiment shown in Table 1, using RoBERTAbase, for target rank R = 8. For LoRA and AdaLoRA, the average output rank is fixed to R. For ROLoRA, we report the mean across 5 runs.

Method	MRPC RTE CoLA QNLI MNLI QQP SST2 STSB							
LoRA & AdaLoRA				8.	0			
ROLoRA	5.5	5.9	7.3	4.9	8.0	6.4	6.6	6.5

401 (marked with † in Tables 1 and 2). On RTE dataset, for example, for both models, ROLoRA shows 402 more than 1.5% improvement compared with vanilla LoRA and AdaLoRA. This shows that setting identical ranks for all layers as in the standard LoRA is clearly a sub-optimal strategy. Even when 403 we are allowed to violate the budget constraint and go beyond the search space as does AdaLoRA, a 404 single step sparsification algorithm does not result in much improvements. We should also note that, 405 while increasing the number of iterations for LoRA and AdaLoRA (marked with \star in Tables 1 and 406 2) leads to some improvements compared with the vanilla versions, for most instances, ROLoRA 407 still edges out on the model performance. 408

Varying target average rank. In Table 3, we report the performances of our algorithm and the baselines using RoBERTa-base when we use the target average rank R = 4 and R = 16. We discover the same pattern as when R = 8, showing that the ranks discovered by ROLoRA can significantly improve LoRA.

413 **Analysis of the output rank.** In Table 4, we report the average ranks of adapters output by each 414 algorithm. Note that for LoRA and AdaLoRA, the average output rank is fixed to R, which means there is no saving in the memory. On the other hand, ROLoRA can output adapters with significantly 415 lower ranks. For example, for the QNLI dataset, the average rank of the output adapters is reduced by 416 62% while ROLoRA still achieve 0.6% improvement in the performance compared with LoRA. This 417 rank reduction can both benefit memory storage and computational time at inference. This benefit 418 comes inherently from the way the algorithm works by always maintaining the number of trainable 419 parameters within the budget and continually sparsifying and updating the rank configuration. 420

Varying sparsification algorithm. We experiment with SoRA Ding et al. (2023) as the SPARSIFY
algorithm and show the algorithm performance in Table 5. This version of ROLoRA generally
demonstrates improvement over LoRA. However, we should note that since SoRA uses a thresholding based strategy for pruning the ranks as opposed to a ranking based strategy used by AdaLoRA,
it is more difficult to control the sparsity of the solution. For this reason, in this experiment, we
keep the same threshold across all iterations of the algorithm (for the SCHEDULE operator). The
performance of the algorithm could be improved with a better SCHEDULE operator.

428 429

430

378

379

380

382

384

385

386

388

390 391

392

393 394 395

397

399 400

6.2 QUESTION-ANSWERING (SQUAD) BENCHMARK

Datasets. SQUAD consists of two datasets for a Question-Answering task, in which we predict the probability of each token being the start and end of the answer to the question.

า เ ลง

Table 5: Results when varying the sparsification algorithm, using RoBERTA-base model. We compared our algorithms with the vanilla (marked with \dagger) and iterative (marked with \star) versions of LoRA (part of Table 1). Superscript ^A represents ROLoRA with AdaLoRA as the SPARSIFY operator, and ^B represents ROLoRA with SoRA as the SPARSIFY operator.

Method	MRPC	RTE	CoLA	SST2
$LoRA^{\dagger}$	$89.90_{\pm.73}$	$78.70_{\pm 1.21}$	$64.34_{\pm.78}$	$95.34_{\pm.27}$
LoRA*	$89.36_{\pm.80}$	$79.21_{\pm 1.30}$	$64.91_{\pm.95}$	$95.62_{\pm.21}$
ROLoRA	^A 90.15 $_{\pm.33}$	$\textbf{81.73}_{\pm.37}$	$64.98_{\pm 1.20}$	$95.30_{\pm.15}$
ROLoRA	$889.90_{\pm.57}$	$80.0_{\pm 1.52}$	$\textbf{65.14}_{\pm 1.33}$	$95.34_{\pm .22}$

Table 6: Results for DeBERTa-v3-base on SQUADv1 and SQUADv2 datasets when the target average rank is R = 8 and R = 4. We report Exact Match for both datasets.

Rank	R =	= 8	R = 4			
Method	SQUADv1	SQUADv2	SQUADv1	SQUADv2		
LoRA [†]	87.73	84.71	87.71	85.05		
LoRA*	87.85	84.61	87.73	85.13		
$AdaLoRA^{\dagger}$	87.90	85.72	87.74	85.26		
AdaLoRA*	88.25	86.20	87.90	85.35		
ROLoRA	88.29	85.87	87.97	85.60		

458

459

445

Implementation details. We experiment with the target average rank set to R = 8 and R = 4. In the former case, the initial average rank for AdaLoRA is set to 12, and in the latter, it is set to 6.

460 Main results. We report the results using DeBERTa-base model. In both cases R = 8 and R = 4, 461 our algorithm outperforms both vanilla LoRA and AdaLoRA. The improvement over LoRA is more 462 than 1.1% on SQUADv2, when R = 8. Even though in this case, AdaLoRA with T iterations shows 463 a better performance, we should note that AdaLoRA does not adhere to the memory constraint.

464 465

6.3 ABLATION STUDY

466 467 We conduct an in-depth study to show why ROLoRA outperforms LoRA and AdaLoRA.

To understand how ROLoRA assigns ranks, we plot the rank distribution of the output adapters for
DeBERTa-v3-base on CoLA and QNLI datasets in Figure 2(a) (same experimental results in Table
We compare these distributions with the configurations discovered by AdaLoRA in Figure 2(b),
which also serves as the base SPARSIFY operator in our algorithm. We observe the following:

(1) ROLoRA finds a configuration that allows ranks to go beyond the upper limit set by AdaLoRA.
In particular, in the plot for CoLA in Figure 2(a), some adapters have ranks higher than the initial rank 12 set by AdaLoRA in the same experiment in Figure 2(b). This means, ROLoRA can go beyond the search space set by a vanilla sparsification algorithm, and thereby achieves improvements.

476 (2) ROLoRA focuses significantly on the value matrices. The overall ranks assigned to the value matrices are significantly higher than those to the key and query matrices. Compared with AdaLoRA, which shows a similar pattern, ROLoRA further accentuates this difference.

We further investigate this latter point to see the role each type of matrices plays in fine-tuning.

In Figure 3 we plot the gradient norm of the pretrained DeBERTA-v3-base model by making one backward pass over the corresponding datasets (CoLA and QNLI). We can see that in each layer, the value matrices have much larger gradient norms than the others. Intuitively, this implies that these matrices are much more sensitive and thus require higher level of tuning. Therefore, one could expect that the value matrices play a more important role in fine-tuning and by assigning higher ranks to the adapters for them relatively to the other matrices, we could achieve a better performance.

⁴⁵⁵ 456 457

To verify this hypothesis, we train separately three sets of adapters: for each type of matrices (key, value, query), we inject a LoRA adapter for each matrix of that type of rank R = 8 and fine-tune the model. We report their performances in Table 7. The result in Table 7 shows that fine-tuning only value matrices can give us better performances than doing so with key and query matrices. This confirms out hypothesis that value matri-ces play a more important role in fine-tuning LLMs.

Table 7: Results when fine-tuning key, query and value matrices separately with LoRA adapters and rank R = 8.

	Value	Key	Query
CoLA	69.7	68.9	67.1
QNLI	94.3	93.7	93.7

Returning to Figure 2, the performance of ROLoRA can by explained by its accentuation on the ranks of value matrices.



Figure 2: Distribution of ranks across layers obtained by (a) ROLoRA and (b) AdaLoRA on CoLA and QNLI using DeBERTA-v3-base. Both algorithms give higher rank to adapters for value matrices, though ROLoRA further accentuates this feature.



Figure 3: Distribution of gradient norms across of the pre-trained model across layers. Value matrices have higher significantly gradient norms than key and query matrices, implying that these matrices are much more sensitive and require more tuning.

7 FUTURE WORK AND CONCLUSION

In this work, we propose a novel iterative algorithm that optimizes rank configurations to enhance
 the performance of the popular fine-tuning algorithm LoRA. Our approach leverages a new insight
 that value matrices in the transformer architecture are more critical for fine-tuning on downstream
 tasks than key and value matrices. We leave further investigation on the applications of the our
 method and its theoretical convergence guarantee to future work.

540 REFERENCES

546

547

548

552

553

554

565

566

578

579

580

581

585

592

- 542 Daniel Cer, Mona Diab, Eneko Agirre, Inigo Lopez-Gazpio, and Lucia Specia. Semeval-2017 task
 543 1: Semantic textual similarity-multilingual and cross-lingual focused evaluation. *arXiv preprint*544 *arXiv:1708.00055*, 2017.
 - Patrick Chen, Hsiang-Fu Yu, Inderjit Dhillon, and Cho-Jui Hsieh. Drone: Data-aware low-rank compression for large nlp models. *Advances in neural information processing systems*, 34:29321–29334, 2021.
- Tianyi Chen, Tianyu Ding, Badal Yadav, Ilya Zharkov, and Luming Liang. Lorashear: Efficient large language model structured pruning and knowledge recovery. *arXiv preprint arXiv:2310.18356*, 2023a.
 - Tianyi Chen, Luming Liang, DING Tianyu, Zhihui Zhu, and Ilya Zharkov. Otov2: Automatic, generic, user-friendly. In *International Conference on Learning Representations*, 2023b.
- Tianyi Chen, Xiaoyi Qu, David Aponte, Colby Banbury, Jongwoo Ko, Tianyu Ding, Yong Ma,
 Vladimir Lyapunov, Ilya Zharkov, and Luming Liang. Hesso: Towards automatic efficient and
 user friendly any neural network training and pruning. *arXiv preprint arXiv:2409.09085*, 2024.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning of quantized llms. *Advances in Neural Information Processing Systems*, 36, 2024.
- Ning Ding, Xingtai Lv, Qiaosen Wang, Yulin Chen, Bowen Zhou, Zhiyuan Liu, and Maosong Sun.
 Sparse low-rank adaptation of pre-trained language models. In *Conference on Empirical Methods in Natural Language Processing*, 2023.
 - Bill Dolan and Chris Brockett. Automatically constructing a corpus of sentential paraphrases. In *Third international workshop on paraphrasing (IWP2005)*, 2005.
- Marguerite Frank, Philip Wolfe, et al. An algorithm for quadratic programming. Naval research logistics quarterly, 3(1-2):95–110, 1956.
- Zeyu Han, Chao Gao, Jinyang Liu, Sai Qian Zhang, et al. Parameter-efficient fine-tuning for large
 models: A comprehensive survey. *arXiv preprint arXiv:2403.14608*, 2024.
- Soufiane Hayou, Nikhil Ghosh, and Bin Yu. Lora+: Efficient low rank adaptation of large models. In *International Conference on Machine Learning*, 2024.
- Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. Towards
 a unified view of parameter-efficient transfer learning. In *International Conference on Learning Representations*, 2022.
 - Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. Deberta: Decoding-enhanced bert with disentangled attention. In *International Conference on Learning Representations*, 2021. URL https://openreview.net/forum?id=XPZIaotutsD.
- 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.
- Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen,
 et al. Lora: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022.
- Aref Jafari, Mehdi Rezagholizadeh, Pranav Sharma, and Ali Ghodsi. Annealing knowledge distillation. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pp. 2493–2504, 2021.
- 593 Martin Jaggi. Revisiting frank-wolfe: Projection-free sparse convex optimization. In *International conference on machine learning*, pp. 427–435. PMLR, 2013.

594 595 596	Rabeeh Karimi Mahabadi, James Henderson, and Sebastian Ruder. Compacter: Efficient low-rank hypercomplex adapter layers. <i>Advances in Neural Information Processing Systems</i> , 34:1022–1035, 2021
597	1055, 2021.
598	Diederik P Kingma. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980,
599	2014.
600	
601	Simon Lacoste-Julien. Convergence rate of frank-wolfe for non-convex objectives. arXiv preprint
602	<i>urxiv.1007.00343</i> , 2010.
603	Shih-yang Liu, Chien-Yi Wang, Hongxu Yin, Pavlo Molchanov, Yu-Chiang Frank Wang, Kwang-
604 605	Ting Cheng, and Min-Hung Chen. Dora: Weight-decomposed low-rank adaptation. In Interna- tional Conference on Machine Learning, 2024.
606	Yinhan Liu Myle Ott Naman Goyal Jingfei Du Mandar Joshi Dangi Chen Omer Levy Mike
607 608	Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized BERT pretraining approach. <i>CoRR</i> , abs/1907.11692, 2019. URL http://arxiv.org/abs/1907.11692.
609	
610	runing Mao, Lambert Mathias, Kui Hou, Amjad Almahairi, Hao Ma, Jiawei Han, Scott Yih, and Madien Khahaa. Uningity A unifold from a substant for a substant for a substant large state of the substant
611	Madian Knabsa. Unipell: A unified framework for parameter-efficient language model tuning. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume
612	1. Long Papers) nn 6253-6264 2022
613	1. Long 1 up (13), pp. 0233-0207, 2022.
614	Jonas Pfeiffer, Aishwarya Kamath, Andreas Rücklé, Kyunghyun Cho, and Iryna Gurevych. Adapter-
615	fusion: Non-destructive task composition for transfer learning. In Proceedings of the 16th Con-
616	ference of the European Chapter of the Association for Computational Linguistics: Main Volume,
617	pp. 487–503, 2021.
618	Pranay Rainurkar Jian Zhang Konstantin Lopyrey and Percy Liang Squad: 100,000+ questions for
619	machine comprehension of text. In Jian Su, Xavier Carreras, and Kevin Duh (eds.), <i>Proceedings</i>
620	of the 2016 Conference on Empirical Methods in Natural Language Processing, EMNLP, pp.
621	2383–2392. The Association for Computational Linguistics, 2016.
622	
623	for squad In Proceedings of the 56th Annual Masting of the Association for Computational
624 625	Linguistics (Volume 2: Short Papers), pp. 784–789, 2018.
626 627	Sylvestre-Alvise Rebuffi, Hakan Bilen, and Andrea Vedaldi. Learning multiple visual domains with residual adapters. <i>Advances in neural information processing systems</i> , 30, 2017.
628	Richard Socher Alex Perelygin Jean Wu Jason Chuang Christopher D Manning Andrew Ng
629	and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment
630	treebank. In Proceedings of the 2013 Conference on Empirical Methods in Natural Language
631	Processing, pp. 1631–1642, Seattle, Washington, USA, October 2013. Association for Computa-
632	tional Linguistics. URL https://www.aclweb.org/anthology/D13-1170.
633	Maitaba Valinaun Mahdi Dagaghaligadah Juan Vahugay and Ali Chadai Dalaga Burayata
634	efficient tuning of pre-trained models using dynamic search free low rank adaptation. In Dre
635	credings of the 17th Conference of the European Chapter of the Association for Computational
636	Linguistics, pp. 3274–3287, 2023.
637	Linguistics, pp. 5271 5267, 2025.
638	Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman.
639	Glue: A multi-task benchmark and analysis platform for natural language understanding. In 7th
640	International Conference on Learning Representations, ICLR 2019, 2019.
641	Alex Warstadt, Amanpreet Singh, and Samuel R Bowman. Neural network acceptability judgments
642	arXiv preprint arXiv:1805.12471, 2018.
643	
644	Adina Williams, Nikita Nangia, and Samuel Bowman. A broad-coverage challenge corpus for sen-
645	tence understanding through interence. In <i>Proceedings of the 2018 Conference of the North Amer-</i>
646	ican Chapter of the Association for Computational Linguistics: Human Language Technologies,
647	http://aclweb.org/anthology/N18-1101.

 T Wolf. Huggingface's transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*, 2019.

Wenhan Xia, Chengwei Qin, and Elad Hazan. Chain of lora: Efficient fine-tuning of language
 models via residual learning. *arXiv preprint arXiv:2401.04151*, 2024.

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

Model	Hyperparams	MRPC	CoLA	RTE	QNLI	MNL	I SST2	QQP	STSB
	# Epochs	30	25	50	25	25	30	25	25
	Warm-up Ratio				0.0)6			
DoDEDTo base	Batch-size	32	32	32	16	16	16	16	32
KODEK1a-Dase	α				8	5			
	Learning rate	1e-3	5e-4	5e-4	5e-4	5e-4	5e-4	5e-4	5e-4
	Max Seq. Length				12	28			
	# Epochs	30	30	50	25	30	25	25	25
	Warm-up Ratio				0.0)6			
DaDEDTa v2 haaa	Batch-size	32	32	32	16	16	16	32	32
DeBER1a-V5-Dase	α				8	5			
	Learning rate				5e	-4			
	Max Seq. Length	l			12	28			

Table 8:	Hyperparameters	are identical	for all methods.
----------	-----------------	---------------	------------------

Table 9: Sparsification parameters for AdaLoRA and ROLoRA.

Model	Hyperparams	MRPC	CoLA	RTE	QNLI	MNLI	SST2	QQP	STSB
	γ	0.1	0.5	0.3	0.1	0.1	0.1	0.1	0.1
DoDEDTo hasa	t_i	600	600	600	2000	8000	600	8000	600
KUDEKIA-Dase	t_f	1800	3500	1800	8000	50000	22000	25000	2000
	Δ_T	1	10	1	100	100	100	100	10
	γ	0.1	0.5	0.3	0.1	0.1	0.1	0.1	0.1
DaDEDTa v2 hasa	t_i	600	600	600	2000	8000	600	8000	600
DeBERTa-V5-Dase	t_f	1800	3500	1800	8000	50000	22000	25000	2000
	Δ_T	1	10	1	100	100	100	100	10

A DATASET DETAILS

GLUE benchmark. GLUE benchmark consists of: MNLI (inference Williams et al. (2018)), SST-2 (sentiment analysis, Socher et al. (2013)), MRPC (paraphrase detection, Dolan & Brockett (2005)), CoLA (linguistic acceptability, Warstadt et al. (2018)), QNLI (inference, Rajpurkar et al. (2018)), QQP8 (question-answering), RTE (inference), and STS-B (textual similarity, Cer et al. (2017)).

SQUAD. SQUAD consists of two Question-Answering datasets: SQUADv1 (Rajpurkar et al., 2016) and SQUADv2 Rajpurkar et al. (2018).

- **B** HYPERPARAMETERS
- 745 B.1 GLUE BENCHMARK

747 Hyper-parameters choices are shown in Tables 8 and 9.

749 B.2 SQUAD BENCHMARK

751 Hyper-parameters choices are shown in Tables 10 and 11.

756			
757			
758			
759			
760			
761			
762			
763			
764			
765 Table 1	0. Uuparparamat	ers are identical for all	methods
766	0. Hyperparamet		methous.
767			-
768	Hyperparams	SQUADVI SQUADV	2 —
769	# Epochs	12	
770	Warm-up	1	
771	Batch-size	16	
772	α	8	
773	Learning rate	1e-3	
774	Max Seq. Lengtl	h 128	
775			_
776			
777			
778			
779			
780			
781			
782			
783			
78/			
785			
786			
700			
788			
780			
700			
790			
702			
792			
Table 11: S	parsification para	meters for AdaLoRA a	nd ROLoRA.
705			
795	Hyperparams S	SOUADv1 SOUADv2	
790		0.01	
709	γ	0.01	
790	t_i	5000	
800	t_f	25000	
801	Δ_T	100	
802			
803			
804			
905			
206			
807			
000			
200			
003			