MORE: A MIXTURE OF LOW-RANK EXPERTS FOR ADAPTIVE MULTI-TASK LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

With the rapid development of Large Language Models (LLMs), Parameter-Efficient Fine-Tuning (PEFT) methods have gained significant attention, which aims to achieve efficient fine-tuning of LLMs with fewer parameters. As a representative PEFT method, Low-Rank Adaptation (LoRA) introduces low-rank matrices to approximate the incremental tuning parameters and achieves impressive performance over multiple scenarios. After that, plenty of improvements have been proposed for further improvement. However, these methods either focus on single-task scenarios or separately train multiple LoRA modules for multi-task scenarios, limiting the efficiency and effectiveness of LoRA in multi-task scenarios. To better adapt to multi-task fine-tuning, in this paper, we propose a novel Mixture of Low-Rank Experts (MoRE) for multi-task PEFT. Specifically, instead of using an individual LoRA for each task, we align different ranks of LoRA module with different tasks, which we named *low-rank experts*. Moreover, we design a novel adaptive rank selector to select the appropriate expert for each task. By jointly training low-rank experts, MoRE can enhance the adaptability and efficiency of LoRA in multi-task scenarios. Finally, we conduct extensive experiments over multiple multi-task benchmarks along with different LLMs to verify model performance. Experimental results demonstrate that compared to traditional LoRA and its variants, MoRE significantly improves the performance of LLMs in multi-task scenarios and incurs no additional inference cost. We also release the model and code to facilitate the community¹.

031 032

033 034

045

046 047 048

051 052

004

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

028

029

1 INTRODUCTION

Recent advancements in Large Language Models (LLMs) have revolutionized various domains, of-035 fering unprecedented performance across numerous tasks (Devlin et al., 2019; Raffel et al., 2020; Brown et al., 2020; Touvron et al., 2023). Plenty of tuning strategies are designed to extend the 037 application of LLMs, such as Instruction Tuning (Wei et al., 2022; Zhang et al., 2023b), Continual Pre-Training (Ke et al., 2023), and Parameter-Efficient Fine-Tuning (PEFT) (Houlsby et al., 2019; Liu et al., 2023b; 2022; Lester et al., 2021; Li & Liang, 2021; Hu et al., 2022). Among these strate-040 gies, PEFT has drawn the most attention due to its fewer parameter tuning and lower computational 041 cost. As the representative PEFT method, Low-Rank Adaptation (LoRA) (Hu et al., 2022) intro-042 duces low-rank matrices to approximate the incremental tuning parameters and demonstrate good 043 performance in many scenarios, which has become a standard paradigm for LLM fine-tuning and inspired many improvements (Liu et al., 2024; Valipour et al., 2023; Ding et al., 2023). 044

Table 1: LoRA-based Fine-tuning Performance of T5-base with varying ranks on different tasks

Task/Rank	r=1	r=2	r=4	r=8	r=16	r=32
MRPC	89.7	89.2	88.7	89.2	89.2	89.5
RTE	77.6	78.7	80.5	77.6	80.1	79.1
SST-2	94.4	94.6	94.8	94.5	94.4	94.5
CoLA	60.9	$\overline{60.0}$	61.9	63.3	62.3	60.5

¹https://anonymous.4open.science/r/MoRE-3C37

054 Despite the achieved progress, LoRA relies on a fixed and unalterable intrinsic rank, making it not 055 flexible enough in multi-task scenarios. Taking Table 1 as an example, when dealing with different 056 tasks, LoRA requires different ranks to achieve the best performance (e.g., best ranks for MRPC and 057 CoLA tasks are 1 and 8). Considering the high computational cost and storage cost of LLM fine-058 tuning, training multiple LoRA modules is sub-optimal for applying LLMs to multi-task scenarios. Meanwhile, searching the best rank of LoRA during LLM fine-tuning is also time-consuming and computationally expensive (Valipour et al., 2023), which highlights the limitations of a one-size-060 fits-all approach in LoRA. This phenomenon also emphasizes the need for adaptive mechanisms 061 that dynamically adjust ranks based on task requirements. 062

063 To overcome the limitations of fixed ranks in LoRA, one promising direction is to explore adaptive 064 mechanisms. For example, DyLoRA (Valipour et al., 2023) dynamically trained all ranks during training to avoid separate rank tuning for each task. AdaLoRA (Zhang et al., 2023a) allocated the 065 parameter budget based on the importance scores of the weight matrices and pruned insignificant 066 singular values to exclude unimportant rank spaces. SoRA (Ding et al., 2023) introduced a trainable 067 gating unit and used proximal gradient descent to optimize the sparsity of the update matrices, 068 thereby dynamically adjusting the intrinsic rank size during training. While these improvements 069 enable dynamic adjustment of rank space, they are primarily designed for single-task scenarios. They do not consider the distinctions and connections among different tasks in multi-task scenarios, 071 prohibiting the effectiveness of LoRA in multi-task scenarios. 072

In the meantime, there also exist other strategies that try to exploit the connections among different 073 tasks. However, they are still far from satisfied. For example, HyperFormer (Mahabadi et al., 2021) 074 enhanced adapter-based methods by utilizing a shared hypernetwork to facilitate cross-task knowl-075 edge sharing, while incorporating task-specific adapters to tailor the model for individual tasks. 076 However, they face limitations due to their inherent performance constraints and additional infer-077 ence latency. Prompt Tuning methods (Vu et al., 2022; Asai et al., 2022; Wang et al., 2023b) are proposed to use learned prompts on source tasks to initialize the learning of target tasks. Despite the 079 effectiveness, these approaches typically require a two-stage training process (i.e., first on the source task and then on the target task), which requires higher data quality and results in training efficiency 081 decrease. Meanwhile, parallel LoRA strategies (Wang et al., 2023a; Li et al., 2024; Liu et al., 2023a; Huang et al., 2023) can effectively address the above shortcoming, offering a better adaptability in multi-task scenarios. Nonetheless, the usage of parallel LoRA modules increases the overall pa-083 rameter count and resource consumption, contradicting the original purpose of LoRA to reduce the 084 training parameters. Thus, one important question should be considered: "How to achieve efficient 085 LLM fine-tuning in multi-task scenarios remains challenging". 086

087 To this end, in this paper, we design a novel Mixture of Low-Rank Experts (MoRE) for efficient LLM 088 fine-tuning in multi-task scenarios. Since different tasks require different ranks of LoRA, we propose to build connections between the ranks and the tasks in a Mixture-of-Expert (MoE) manner. Specif-089 ically, we propose to treat each rank in the LoRA module as an expert and design a novel Adaptive 090 Rank Selector. Thus, the different experts corresponding to different tasks can share common infor-091 mation and maintain distinctive information simultaneously (i.e., the ranks r_i and r_j can share some 092 common parameters). Meanwhile, our proposed selector uses a gating mechanism to select the appropriate rank expert for each task. Moreover, to fully exploit the distinctions and connections among 094 different tasks for accurate rank selection, we develop a novel *CL-based Task Embedding* module, 095 which assigns a task embedding to each task and uses a Contrastive Learning (CL) optimization to 096 ensure the quality of learned task embeddings. Furthermore, we incorporate the Balanced Dataset Sampling strategy to address the severe dataset imbalance in multi-task scenarios. Along this line, 098 MoRE can fully exploit the potential of LoRA and realize efficient LLM fine-tuning in multi-task scenarios. Finally, we conduct extensive experiments on multi-task benchmarks to validate the ef-099 fectiveness of MoRE. Experimental results demonstrate the efficiency of MoRE in multi-task and 100 low-resource transfer scenarios. 101

- 102
- 2 RELATED WORK
- 103 104 105
- 2.1 PARAMETER-EFFICIENT FINE-TUNING (PEFT)
- 107 PEFT methods are designed to adapt LLMs to new tasks with minimal additional parameters. Representative works include BitFit (Zaken et al., 2021), Adapters (Houlsby et al., 2019), Prompt Tun-

108 ing (Liu et al., 2023b; 2022; Lester et al., 2021), Prefix Tuning (Li & Liang, 2021) and Low-Rank Adaptation (LoRA) (Hu et al., 2022). Among these methods, LoRA is the most representative one. 110 It introduces trainable low-rank matrices to approximate weight updates, realizing highly efficient 111 fine-tuning with low cost, which has led to various extensions (Kopiczko et al., 2024; Liu et al., 2024; 112 Valipour et al., 2023; Zhang et al., 2023a; Ding et al., 2023). For example, VeRA (Kopiczko et al., 2024) further reduced the number of trainable parameters in LoRA by employing shared low-rank 113 matrices and trainable scaling vectors. DoRA (Liu et al., 2024) enhanced fine-tuning performance 114 and stability by decomposing the pre-trained weights into magnitude and direction components. For 115 greater flexibility in LoRA's rank, DyLoRA (Valipour et al., 2023) dynamically trained all ranks dur-116 ing training to avoid separate rank tuning for each task. AdaLoRA (Zhang et al., 2023a) allocated the 117 parameter budget based on the importance scores of the weight matrices and pruned insignificant 118 singular values to exclude unimportant rank spaces. SoRA (Ding et al., 2023) introduced a train-119 able gating unit and used proximal gradient descent to optimize the sparsity of the update matrices, 120 dynamically adjusting the intrinsic rank size during training. 121

However, LoRA's fixed-rank constraint limits its flexibility. Although recent works (Valipour et al., 2023; Zhang et al., 2023a) have enhanced LoRA's adaptability, they predominantly address single-task training scenarios. These approaches do not consider multi-task scenarios, where selecting the most suitable rank for different tasks remains an open challenge. This gap underscores the need for more flexible and adaptive methods capable of efficiently handling diverse and concurrent tasks in multi-task learning scenarios.

127 128 129

2.2 MULTI-TASK LEARNING

130 Multi-task learning (MTL) focuses on simultaneously solving multiple related tasks with a single 131 model, which has been studied from multiple perspectives and offers several advantages (Zhang & 132 Yang, 2021; Vandenhende et al., 2022). As large language models advance, multi-task learning has 133 become an essential skill for them. However, it still faces several challenges, such as conflicts be-134 tween different tasks, balancing task weights, and the demands on training resources (Chen et al., 135 2021; Kollias et al., 2024). Optimizing PEFT methods for MTL scenarios is a highly valuable direc-136 tion. For instance, HyperFormer (Mahabadi et al., 2021) improved Adapter-based methods by using a shared hypernetwork for cross-task knowledge sharing while integrating task-specific Adapters. 137 In Prompt Tuning, SPoT (Vu et al., 2022) learned prompts from source tasks and adapted them for 138 target tasks, enhancing model performance. ATTEMPT (Asai et al., 2022) used an attention mecha-139 nism to merge source and target prompts for effective knowledge transfer. MPT (Wang et al., 2023b) 140 employed prompt decomposition and knowledge distillation to create a transferable prompt, which 141 was then fine-tuned with low-rank modifications for specific tasks. 142

Additionally, LoRA-based improvements have shown significance in multi-task scenarios. Multi-143 LoRA (Wang et al., 2023a) employed multiple parallel LoRA modules during training, ensuring 144 that the rank space closely approximates fine-tuning. MixLoRA (Li et al., 2024) used multiple 145 parallel LoRA experts with a gating mechanism to select the appropriate expert for each token. 146 MOELoRA (Liu et al., 2023a) learned task-shared and specific knowledge with multiple experts 147 and adjusted their contributions for each task using a gating function. However, these methods also 148 have limitations. The parallel LoRA modules increase the number of trainable parameters, signifi-149 cantly reducing training efficiency. Besides, they do not account for the varying rank requirements 150 of different tasks.

151 152

3 PRELIMINARY

153 154

156

3.1 PROBLEM DEFINITION

In multi-task learning scenarios, the objective is to concurrently learn multiple tasks, each characterized by potentially diverse data distributions and goals. Formally, we consider a set of tasks $T = \{\mathcal{T}_1, \mathcal{T}_2, \ldots, \mathcal{T}_T\}$, where each task \mathcal{T}_t is associated with a dataset $\mathcal{D}_t = \{(x_i^t, y_i^t)\}_{i=1}^{N_t}$ comprising N_t input-output pairs. Here, x_i^t denotes the input data and y_i^t denotes the corresponding label or output for task \mathcal{T}_t . The target is to learn a shared model F with parameters θ to satisfy the requirements of different simultaneously.





3.2 LORA: LOW-RANK ADAPTATION

178 (LoRA) (Hu et al., 2022) is designed to reduce the computational cost and memory footprint of 179 adapting LLMs to new tasks. Instead of fine-tuning all parameters, LoRA adapts the model by in-180 troducing low-rank updates to the existing weight matrices. Formally, let $W_0 \in \mathbb{R}^{m \times d}$ be a weight 181 matrix of a specific LLM, where m and d denote the input and output dimensions, respectively. 182 LoRA approximates the weight update ΔW using a low-rank decomposition:

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

where $A \in \mathbb{R}^{r \times d}$ and $B \in \mathbb{R}^{m \times r}$ are the learned low-rank matrices, and $r \ll min(m, d)$ is the pre-defined rank. This decomposition significantly reduces the number of trainable parameters from $m \times d$ to $r \times (m + d)$. For original $h = W_0 x$, the modified forward pass yields:

$$\boldsymbol{h} = \boldsymbol{W}_0 \boldsymbol{x} + \Delta \boldsymbol{W} \boldsymbol{x} = \boldsymbol{W} \boldsymbol{x} + \boldsymbol{B} \boldsymbol{A} \boldsymbol{x}. \tag{2}$$

However, this process highly depends on the pre-defined rank r, which is time-consuming and computationally expensive to search. And this problem will be amplified in multi-task scenarios, limiting the potential of LoRA. Thus, *How to use LoRA to achieve efficient LLM fine-tuning in multi-task scenarios* is the main focus of our paper.

192 193 194

195

175

176

177

183

187 188

189

190

191

4 MIXTURE OF LOW-RANK EXPERTS

To tackle the inefficient problem of LoRA in multi-task scenarios, we propose a novel Mixture of Low-Rank Experts (MoRE). The cores lie in *how to learn experts* and *how to select them*. As illustrated in Figure 1, we focus on parameters in attention layer and FFN layer of the Transformer block. We first assign a task embedding for each task to describe the abstract task characteristics. Then, based on the task embedding, we design a novel *adaptive rank selector* to select the appropriate rank for each task, term as the rank expert. Finally, we incorporate contrastive learning to ensure the quality of learned task embedding and design a *Balanced Data Sampling strategy* to stabilize the learning process for better multi-task learning. Next, we will introduce each part in detail.

4.1 TASK EMBEDDING

Existing multi-task learning methods focus on mining useful information from task data and transferring knowledge from one task to another. Despite the progress, they are still weak at sharing common information among tasks and distinguishing specific information aligning with each task. This shortcoming will prohibit the efficiency of PEFT methods when using them to tune LLMs in multi-task scenarios. Therefore, we propose using task embeddings to represent different tasks so that task characteristics can be summarized comprehensively. This operation is also the precondition of our designed rank expert for measuring the connections and distinctions among different tasks.

213 Specifically, we use matrix $E = \{e_1, e_2, ..., e_l\}$ to denote all tasks, where e_i represents the i^{th} task 214 in the multi-task scenarios. Then, we leverage Kaiming Initialization to initialize them and learn 215 precise E during model training. Since there is no supervised signal for E, we design a Contrastive

Learning (CL) based optimization target to learn them, which will be introduced in Section 4.3.

2164.2ADAPTIVE RANK SELECTOR217

234

239 240 241

246

250

256

257

260 261

262

As illustrated in Section 1, vanilla LoRA and its typical variances usually have a fixed rank *r*, which is pre-defined by experts. However, different tasks may benefit from different ranks depending on their complexity and data distributions (Valipour et al., 2023; Ding et al., 2023). Searching the best rank is time-consuming and computationally expensive. Meanwhile, training parallel LoRA modules or multiple LoRAs when applying LLMs to multi-task scenarios will amplify the problem and prohibit the effectiveness, causing high computational and storage costs. Therefore, we employ Mixture-of-Experts (MoE) framework and design a novel *Adaptive Rank Selector*.

Different from previous work that treated the entire LoRA module as an expert, we propose to treat 225 the rank r as the expert and use one LoRA to realize LLM fine-tuning in multi-task scenarios. As-226 suming the selected rank of LoRA is r, the rank expert can be selected within the range [1, r]. Along 227 this line, different experts can share common information at the overlap part in the learned metrics 228 (i.e., A and B) and align specific information corresponding to each task at the non-overlap part. 229 Formally, we use the learned task embedding e_t to select the appropriate rank from the LoRA mod-230 ule and leverage a gating network $G(\cdot)$ to guarantee the quality of the selection. Let $\{1, 2, \ldots, r\}$ 231 be the set of experts' ranks. For task $\mathcal{T}_t, G(\cdot)$ takes \mathbf{e}_t as input and outputs a probability distribution 232 over rank experts as follows: 233

$$\mathbf{p}_t = G(\mathbf{e}_t) = \operatorname{softmax}(\mathbf{W}_q \mathbf{e}_t + \mathbf{b}_q), \tag{3}$$

where $\{\mathbf{W}_g, \mathbf{b}_g\}$ are learnable parameters. The probability distribution $\mathbf{p}_t \in \mathbb{R}^r$ indicates the relevance of each rank to task \mathcal{T}_t . During the forward pass, we select the rank with the highest probability and use the selected rank to truncate the LoRA module for rank expert construction. Then, MoRE uses LoRA paradigm to realize the fine-tuning as follows:

$$r_t = \arg \max \mathbf{p}_t,$$

$$h = \mathbf{W}_0 x + \mathbf{B}_t \mathbf{A}_t x, \quad \mathbf{A}_t = \mathbf{A}[:r_t,:], \quad \mathbf{B}_t = \mathbf{B}[:,:r_t].$$
(4)

One step further, during the backward pass, the arg max in Eq.(3) is non-differentiable, causing $G(\cdot)$ unable to be learned. Thus, we incorporate Straight-Through Estimator (STE) (Bengio et al., 2013) technique to address this issue. Specifically, we use STE to calculate the approximate gradient to allow the gradient to propagate back to $G(\cdot)$ correctly:

$$Ste(\mathbf{p}_t) = \mathbf{p}_t + sg[one_hot(\mathbf{p}_t) - \mathbf{p}_t],$$
(5)

where $one_hot(\cdot)$ is a function that converts a vector into its one-hot version. $sg(\cdot)$ stands for stop gradient. Then, we modify the forward process in Eq.(4) as:

$$\mathbf{h} = \mathbf{W}_{\mathbf{0}} x + \operatorname{Ste}(\mathbf{p}_{t})[r_{t}] \cdot \mathbf{B}_{t} \mathbf{A}_{t} x.$$
(6)

Thus, Adaptive Rank Selector module can realize a precise selection of rank experts. Furthermore,
since MoRE uses the overlap part among LoRA metrics to share the common information across
different tasks, the lower part will be updated more frequently during fine-tuning. Thus, its learning
rate should be small for a slow and stable updating. To realize this goal, we perform a linear scaling
on its weights for the balance:

$$\mathbf{h} = \mathbf{W}_{\mathbf{0}}x + \operatorname{Ste}(\mathbf{p}_{t})[r_{t}] \cdot \frac{r_{t}}{|T|} \mathbf{B}_{\mathbf{t}} \mathbf{A}_{\mathbf{t}} x,$$
(7)

where |T| represents the total number of tasks. To verify the effectiveness of this design, we also conducted an ablation study on this operation in Section 5.4.

4.3 BALANCED DATA SAMPLING AND CL-BASED OPTIMIZATION

Balanced Data Sampling. In multi-task scenarios, data distributions of different tasks are also
essential for LLM fine-tuning. For instance, in GLUE benchmark (Wang et al., 2018), MNLI and
RTE datasets have proportionally disparate data distributions (i.e., 392, 000 v.s. 2, 500 examples).
If this attribute is not considered when fine-tuning LLMs in multi-task scenarios, it is obvious that
fine-tuned LLMs will underfit the task with smaller datasets.

In response, we propose a simple but effective Balanced Dataset Sampling strategy to ensure each dataset contributes proportionally during the fine-tuning process, regardless of its size. Specifically, we assign a sampling weight ϕ_t to each dataset \mathcal{D}_t , which is inversely proportional to its size: 270 074

273

$$\Phi = [\phi_1, \phi_2, ..., \phi_T], \quad \phi_t = \exp\left(\frac{|\mathcal{D}_t|}{\sum_{i=1}^T |\mathcal{D}_i|}\right),$$

$$D_t = Sampling(D, \Phi),$$
(8)

275 where $|\mathcal{D}_t|$ is the size of dataset \mathcal{D}_t . Sampling (D, Φ) denotes sampling a subset from all datasets D 276 with the probability distribution Φ . This dynamic sampling strategy helps to balance the contribu-277 tions of different datasets, thereby reducing the risk of underfitting smaller datasets and improving 278 the overall performance of the multi-task training.

279 **CL-based Optimization.** As mentioned in Section 4.1, there is no supervised signal for task em-280 bedding learning. Thus, one important question should be considered: "How to ensure the task characteristics and task distinguishability of the learned task embedding without annotation require-282 ments?" In response, we propose to leverage CL to ensure the quality of learned task embeddings. 283 Consider a batch B of samples, where all samples in B belong to the same task \mathcal{T}_t . Let $\{\mathbf{x}_i\}_{i=1}^N$ be the 284 set of N samples in \mathcal{B} , and let \mathbf{h}_i be the representation of sample \mathbf{x}_i obtained from the model. The 285 task embedding for task \mathcal{T}_t is denoted as \mathbf{e}_t . The optimization target can be formulated as follows:

289

300

301 302

305 306 307

281

 $\mathcal{L}_{con} = \frac{1}{N} \sum_{i=1}^{N} \left[\log \frac{\exp\left(\frac{sim(\mathbf{h}_{i}, \mathbf{e}_{t})}{\tau}\right)}{\sum_{k=1}^{T} \exp\left(\frac{sim(\mathbf{h}_{i}, \mathbf{e}_{k})}{\tau}\right)} \right],$ (9)

290 where $sim(\cdot, \cdot)$ denotes a similarity measure, such as the dot product or cosine similarity, and T is 291 the total number of tasks. τ is the temperature. e_t and e_k are the t^{th} and k^{th} tasks where $t \neq k$. 292 By using Eq.(9), we can measure the connection between task embedding e_t and its data samples 293 $\{\mathbf{x}_i\}_{i=1}^N$. Since each data sample is close to the corresponding task embedding, we can conclude the learned task embeddings can be used to describe task characteristics, which is also supported by experimental results in Section 5.3. 295

296 Besides using contrastive loss to learn task embeddings, we also select generation loss \mathcal{L}_{gen} to mea-297 sure the discrepancy between the sequences generated by the model and the target sequences. Let y298 and \hat{y} be target sequence and generation, \mathcal{L}_{gen} can be formulated with the cross-entropy loss: 299

> $\mathcal{L}_{gen} = -\sum_{t=1}^{T} y_t \log \hat{y}_t.$ (10)

303 Then, we leverage a hyperparameter λ to balance the contributions of the generation loss and the 304 contrastive loss, and formulate the overall optimization target of MoRE as follows:

$$\mathcal{L} = \mathcal{L}_{gen} + \lambda \mathcal{L}_{con}.$$
 (11)

Discussion. Compared with existing LoRA-based PEFT methods and MoE-based fine-tuning 308 methods, our proposed MoRE has the following properties. 1) We propose to treat different rank values in one LoRA module as experts, and design an adaptive rank selector to select appropriate 310 rank experts for different tasks, which can effectively share the common information among tasks 311 and emphasize the specific information aligned to each task; 2) We leverage task embeddings to ac-312 curately describe the abstract task characteristics with a CL optimization; 3) We also consider task 313 data distributions and design a simple but effective Balanced Data Sampling strategy to ensure the capability of fine-tuned LLMs on different tasks. 314

- 315 316
- 5 EXPERTMENTS
- 317 318

319

5.1 EXPERIMENTAL SETUP

320 **Datasets.** We utilized GLUE benchmark (Wang et al., 2018) to evaluate the model performance. 321 GLUE covers multiple tasks of paraphrase detection (MRPC, QQP), sentiment classification (SST-2), natural language inference (MNLI, RTE, QNLI), and linguistic acceptability (CoLA). Following 322 previous work (Zhang et al., 2021), for those datasets with fewer than 10,000 samples (i.e., RTE, 323 MRPC, STS-B, CoLA), we split the original validation set into new validation and test sets equally.

Methods	narams/task	MNLI	OOP	ONLI	SST-2	STS-B	MRPC	RTE	CoLA	AVG
E' ('	2014	05.7	× ×-	02.0	02.5	00.0	00.2	75 4	54.0	02.0
Adaptara	28M 1.8M	85.7	91.1	92.0	92.5	88.8	$\frac{90.2}{00.2}$	/5.4 70.2	54.9	83.8
PT	0.6k	85.6	90.5	$\frac{93.2}{03.2}$	95.0	89.9 80.0	<u>90.2</u> 86.3	70.5 67.6	55.3	04.4 82.8
$LoRA_{m-8}$	0.39M	85.8	89.2	$\frac{93.2}{93.1}$	93.2	90.4	89.9	76.3	62.8	85.1
$LoRA_{r=16}$	0.78M	84.9	89.6	93.0	93.7	90.4	88.7	80.6	63.9	85.6
HyperFomer	638K	85.7	90.0	93.0	94.0	89.7	87.2	75.4	63.7	84.8
MPT	10.5K	84.3	90.0	93.0	93.3	90.4	89.2	82.7	63.5	85.8
MultiLoRA	1.56M	85.9	89.7	92.8	94.5	89.8	88.2	80.6	66.9	86.0
MixLoRA	1.49M	85.8	90.0	92.9	93.7	90.3	89.2	78.4	67.2	85.9
MOELoRA	0.78M	86.3	90.1	<u>93.2</u>	<u>94.2</u>	90.0	89.7	81.3	<u>68.4</u>	<u>86.7</u>
MoRE	0.78M	86.2	90.0	93.4	93.7	90.7	91.2	83.5	69.9	87.3
LLaMA2-LoRA	2.5M	86.9	88.6	93.5	96.2	90.2	92.6	89.2	65.0	87.8
LLaMA2-MultiLoRA	10 M	87.6	85.0	93.4	96.7	92.2	88.7	87.8	72.4	88.0
LLaMA2-MixLoRA	12.2M	86.8	88.1	93.6	96.0	91.3	88.2	87.1	73.2	88.0
LLaMA2-MOELoRA	5M	87.0	87.6	91.4	96.3	92.4	<u>91.2</u>	87.8	64.4	87.3
LLaMA2-MoRE	5M	89.4	89.0	94.4	96.9	92.2	89.2	92.1	66.9	88.8

Table 2: Performance on GLUE benchmark. For STS-B, we report Pearson correlation coefficients. For CoLA, we report Matthews correlation. For all other tasks, we report accuracy. **Bold** and <u>underlined</u> fonts indicate the best and the second-best results.

344 345

327

For others, we randomly select 1, 000 examples from training set as the validation set, and use original validation sets as test sets. Additionally, we included the BoolQ (Clark et al., 2019), PIQA (Bisk et al., 2020), OBQA (Mihaylov et al., 2018), and ARC (Clark et al., 2018) datasets to assess the model's performance in commonsense reasoning tasks. These datasets provide a variety of challenges that require understanding of everyday scenarios and logical reasoning. Moreover, we select SciTail (Khot et al., 2018), BoolQ (Clark et al., 2019), and CB (de Marneffe et al., 2019) datasets to evaluate model robustness and generalization capabilities in few-shot learning scenarios. We also provide an early attempt on generation tasks (Nan et al., 2021) and report results in Appendix A.

Baselines. We compare MoRE with the following baselines: (1) Full fine-tuning (FT), (2) Vanilla
Adapter (Houlsby et al., 2019), (3) Vanilla prompt tuning (PT) (Raffel et al., 2020), (4) Vanilla
LoRA (Hu et al., 2022), We also select the following advanced multi-task PEFT baselines: (1) HyperFomer (Mahabadi et al., 2021), (2) MPT (Wang et al., 2023b), (3) MultiLoRA (Wang et al., 2023a), (4) MixLoRA (Li et al., 2024). (5) MOELoRA (Liu et al., 2023a). All methods are tuned
based on reported settings with T5-base and LLaMA2-7B as backbone for a fair comparison.

Training Setup. We select T5-base (Raffel et al., 2020) and LLaMA2-7B (Touvron et al., 2023) as the backbone. The optimizer is AdamW. The learning rate is 3×10^{-4} , with a linear decay schedule and warm-up over the first 500 steps. The batch size is 32, and the training process spanned 5 epochs. The maximum input sequence length is 128 tokens. The λ is set to 0.1, and the temperature τ in Eq.(9) is 0.05. For few-shot domain transfer, We use the best checkpoint trained on the GLUE tasks for initialization. The task embeddings of the most similar tasks will be shared (e.g., MNLI for SciTail and CB). The T5-base is trained on Ubuntu 20.04 platform using two NVIDIA RTX 4090 GPUs, and LLaMA2-7B is trained with four NVIDIA Tesla A100 PCIe GPUs.

Evaluation Setup. For GLUE benchmark and commonsense reasoning tasks, we selected the
 checkpoint with the highest average performance on validation set. For few-shot learning, we per formed training and testing under each shot setting using 5 random seeds. Then, we reported the
 average performance for a fair and robust estimation and comparison.

371

372 5.2 OVERALL PERFORMANCE373

Performance on GLUE Benchmark and Commonsense Reasoning Tables 2 and 3 present model performance on multi-task scenarios. From Table 2, we can observe that MoRE achieves impressive performance over different tasks with a relatively small number of fine-tuned parameters. Moreover, compared with LoRA implementation (r = 16), MoRE achieves significant improvement (e.g., 1.7% average improvement) without extra tuning parameter, demonstrate the superiority of MoRE.

Methods	params/task	BoolQ	PIQA	OBQA	ARC-E	ARC-C	AVG
LLaMA2-LoRA	2.5M	80.9	77.7	79.0	83.7	76.9	79.6
LLaMA2-MultiLoRA	10M	76.5	72.9	68.2	81.6	61.9	72.2
LLaMA2-MixLoRA	12.2M	84.3	79.5	82.6	86.8	76.3	81.9
LLaMA2-MOELoRA	4.5M	84.0	<u>79.9</u>	81.8	86.8	77.3	82.0
LLaMA2-MoRE	4.5M	87.2	82.3	83.0	86.7	74.2	82.7

Table 3: Performance (Accuracy) of all methods on Commonsense Reasoning scenarios.

385 386

384

378

379380381382

By treating ranks as experts and learning accurate task embedding to support the expert selection, MoRE can effectively share common information and specify aligned information across different tasks, and achieve impressive performance without too many fine-tuned parameters. Furthermore, the number of fine-tuned parameters of MoRE is the same as representative $LoRA_{r=16}$, smaller than other LoRA-based improvement methods. When fine-tuning larger models (e.g., LLaMA2-7B), MoRE can achieve much better performance. All these phenomena demonstrate the superiority of MoRE. Additionally, MoRE can be further optimized during inference, allowing the inference-time parameters to be comparable to those of $LoRA_{r=8}$ (see Section 5.4 for details).

For PEFT baselines, we observe that fine-tuned performance over small datasets (i.e., MRPC, RTE, 396 and CoLA) is not so good. One possible reason is that they do not consider the shared knowledge 397 across different tasks, treating each task individually. Moreover, their fine-tuned module requires more training data. Thus, we can observe the sub-optimal multi-task performance on these base-398 lines. For multi-task baselines, though they consider the shared knowledge across different tasks, 399 they do not distinguish the distinctions and connections among different tasks. For example, Hy-400 perFormer just learns task embeddings without any constraints. MPT sacrifices too much on MNLI 401 task. Therefore, their performance is not comparable with MoRE. As for MultiLoRA and MixLoRA, 402 although they improve the performance of LoRA, they do not utilize task-aware modules or consider 403 specific rank allocation for different tasks. As a result, their performance improvements are limited. 404

Meanwhile, we can obtain the similar results from Table 3. MoRE achieves the highest average accuracy when fine-tuning LLaMA2 on commonsense reasoning scenario, surpassing all other approaches with a smaller number of fine-tuned parameters, proving the superiority of MoRE. In contrast, MultiLoRA has worse performance. We speculate the reason is that commonsense reasoning tasks require more fine-grained task information sharing and distinguishing, which cannot be satisfied by a simple LoRA ensemble.

Performance on Few-shot Domain Transfer To further verify the efficiency of MoRE, we conduct few-shot domain transfer experiments and report results in Table 4. We can observe MoRE achieves stable and optimal performance over most datasets with different few-shot settings, This consistency proves the efficiency of MoRE on sharing common information and distinguishing specific information across different tasks, which helps leverage minimal data for effective transfer learning.

For baselines, traditional fine-tuning and multi-task tuning (i.e., HyperFomer and MPT) require more parameters to be tuned. Thus, their performance is subpar with limited training data. Increasing training data will improve their performance. For LoRA-based baselines, multi-task LoRA-based methods have similar performance with vanilla $LoRA_{r=16}$ and do not show performance gains in multi-task learning. We speculate that they may encounter difficulties in efficiently allocating appropriate ranks or adapting parameters to new tasks when only a small number of samples is available. These phenomena highlight the challenges of effectively utilizing few-shot data to achieve good generalization across different domains, demonstrating the superiority of MoRE.

423 424

425

5.3 EXPERT SELECTION ANALYSIS AND TASK VISUALIZATION

Low-Rank Expert Allocation. To investigate the expert distribution after model fine-tuning, we
 analyze the expert allocation across all layers for each task. As shown in Figure 2(a), in most cases,
 all tasks heavily rely on experts with ranks 1, 2, or 3. This indicates a significant amount of parameter redundancy in LoRA, where the higher-rank parameters in many modules do not substantially
 contribute during the fine-tuning process. This is consistent with our design that MoRE leverage
 rank experts to share common information across different tasks with lower-rank parameters. Moreover, to illustrate the dependency of different tasks on various ranks, and to eliminate the influence

Task	k-shot	Finetuning	LoRA	HyperFomer	MPT	MultiLoRA	MixLoRA	MOELoRA	MoRE
	4	50.5	64.2	48.0	62.2	65.2	62.8	64.0	64.6
BoolQ	16	56.5	<u>66.1</u>	50.2	63.3	65.8	64.4	64.8	66.2
	32	58.4	67.4	58.3	68.9	67.6	66.2	65.7	<u>67.9</u>
	4	57.7	84.3	60.7	73.6	85.0	86.6	85.4	85.7
CB	16	77.0	85.7	76.3	78.6	85.7	86.4	86.3	86.4
	32	80.0	87.1	81.4	82.1	86.6	89.3	88.3	88.6
	4	79.6	80.8	82.0	80.2	78.1	77.5	80.4	83.8
SciTail	16	80.0	84.0	86.5	87.3	81.7	82.4	83.1	<u>86.7</u>
	32	81.9	85.3	85.8	<u>86.3</u>	83.6	83.3	84.5	87.4
st = .			Expert 1 Expert 2 Expert 3 Expert 5 Expert 5 Expert 6 Expert 7	8- 24		Espert 1 Espert 2 Espert 3 Espert 4 Espert 6 Espert 6	encoder_query	encoder_key en	coder_value
ion Cou			Expert 8			- CAPRILO	decoder_query	decoder_key de	coder_value
Distribu			distribution of the second					•	•
		և և հետ		·	li II.			•••••	•
cola	mnli mrpc	gnli ggp rte sst	2 stsb	, cola mnli mrpc	qnli qqp	rte sst2 stsb	• cola • mnli • mn	pc • qnli • qqp • rte	sst2 • stsb
		(a)			(b)			(c)	

Table 4: Few-shot domain transfer results (Accuracy) of T5-base models fine-tuned on GLUE averaged across 5 seeds. **Bold** and underlined fonts indicate the best and the second-best results.

Figure 2: (a)-(b) The distribution of expert allocation. (c) Visualization of the task embeddings.

of low-rank experts 1-3, we proportionally scaled the low-rank experts 1-3 for each task. As depicted in Figure 2(b), different tasks indeed exhibit stronger dependencies on different ranked experts. For instance, MRPC relies on expert 4, RTE depends on expert 6, and STSB still relies on expert 1. This phenomenon proves strong evidence to support that MoRE can allocate more suitable ranks for different tasks and use higher-rank parameters to emphasize the specific information aligned to each task. This is also the reason that MoRE can achieve the best performance in multi-task scenarios.

460 **Visualization of Task Embeddings.** In Section 4.1, we mention that task embedding is essential 461 for rank expert selection. Here, we visualize the learned task embeddings to verify the quality and 462 provide some insight for understanding MoRE. Specifically, we use PCA to process the task em-463 beddings from the final layer of the self-attention module and report results in Figure 2(c). We can 464 observe that task embeddings exhibit varying degrees of clustering. Similar tasks (e.g., MRPC and 465 QNLI) have a tendency to cluster, while different tasks have a large distance. Moreover, STSB and CoLA tasks seem relatively independent to all other tasks. It is consistent with our expectations. 466 Since STSB involves similarity computation and other tasks are classification tasks, they indeed 467 have obvious differences in abstract task characteristics. These phenomena also give us some in-468 sight into the efficiency and effectiveness of MoRE. By using task embeddings to extract abstract 469 task information, MoRE can capture and represent the similarities and differences between tasks in 470 diverse ways. This capability helps to select the appropriate rank expert, which in turn ensures the 471 superiority of our proposed MoRE. We also provide more examples in Appendix C. 472

472 473 474

453 454

455

456

457

458

459

5.4 Ablation Study and Parameter Analysis

475 Ablation Study. We perform an ablation study to better verify the contribution of each component 476 in MoRE. As shown in Table 5, we can observe a significant performance decrease (i.e., 0.9% and 477 0.7% average accuracy decrease) when replacing task-specific embeddings with a shared embed-478 ding (w/o Task Embeddings) or removing contrastive optimization (Eq.(9)) (w/o CL optimization), 479 proving the importance of the CL-based task embedding module. This is consistent with our de-480 sign in Section 4.1. Moreover, removing STE and using soft expert selection with Eq.(3) (w/o STE) 481 also have a big impact on the model performance. Since we treat different ranks as experts, using 482 soft expert selection will reduce the discrimination between different experts, leading to a tendency 483 towards vanilla $LoRA_{r=16}$. Furthermore, when using random sampling (w/ Random Sample), we can observe 86.2% average GLUE accuracy, indicating the positive impact of our balanced sample 484 selection strategy. Besides, when removing linear scaling (w/o Linear Scaling), we observed a slight 485 drop in performance, indicating that this adjustment helps mitigate the overfitting of LoRA.

9

433 434

432

Table 5: Ablation study results (Average Results

on GLUE benchmark) of MoRE.

Conditions	GLUE Avg.
MoRE	87.3
w/o Linear Scaling	<u>87.0</u>
w/o Task Embeddings	86.1
w/o CL optimization	86.3
w/o STE	86.4
w/ Random Sample	86.2



Figure 3: (a) Comparison of trainable parameters. (b) Parameter sensitivity experiments.

498 499 500

486

487

488

489

490

491

492

493

494

495

496

497

501 **Parameter Sensitivity Test.** There are two hyper-parameters that affect model performance: (1) 502 λ in Eq.(11); (2) the dimension of task embeddings. Therefore, we conduct additional experiments to verify their impacts and report results in Figure 3(b). From the figure, we have the following 504 observations. First, with the increase of λ , model performance first decreases and then increases. 505 Since different tasks have various data samples, the corresponding contrastive loss would exhibit oscillations, which may have negative impacts on fine-tuning performance. Thus, we set $\lambda = 0.1$ to 506 obtain the best performance. Second, with the increase of the dimension of task embedding, model 507 performance will first increase and then decrease. We attribute this phenomenon to the following 508 reasons. When the dimension is too small, task embedding cannot capture complex task information, 509 which will harm the selection of rank experts. When the dimension is too big, its training requires 510 more data, which cannot always be satisfied. Moreover, since we need to calculate similarity with 511 samples, the sample embedding from LLMs also should be considered. Based on all these reasons, 512 we finally set the dimension of task embeddings as 768. 513

Parameter Efficiency. To analyze the model complexity, we give the number of tuning parameters 514 of different LoRA-based methods, which is summarized in Figure 3 (a). The notation explanations 515 are as follows: $\{L, r, (m, d), n, T, h\}$ refer to model layers, LoRA rank, model dimensions, par-516 allel LoRA module number, task number, and task embedding dimension. Compared with tuning 517 parameter size of LoRA, the added parameter number of MoRE is 6Lh(r + T), including the ex-518 tra task embeddings (Th for a single LoRA module) and adaptive rank selector (rh for a single 519 LoRA module). Compared with MultiLoRA and MixLoRA which use parallel module design to 520 tackle multi-task learning, MoRE is more efficient. Moreover, once our task embedding and gate 521 modules are trained, we can construct a mapping from tasks to experts, which allows us to avoid the repeated computation of the task embedding and gate modules during inference, thereby reduc-522 ing the parameter count to be consistent with LoRA. This is also the reason why MoRE achieves 523 impressive performance in multi-task scenarios without too many fine-tuning parameters. We also 524 provide detailed parameter size calculations of baselines in Appendix B. 525

526 527

528

6 CONCLUSION AND FUTURE WORK

529 In this paper, we argued that existing PEFT methods either did not consider multi-task fine-tuning 530 or used parallel structures that added too many tuning parameters, prohibiting the efficiency of LoRA. In response, we proposed a novel MoRE for multi-task PEFT. By treating each low-rank 531 in LoRA module as a specialized expert, MoRE could share common information with lower-rank 532 parameters and emphasize the specific information aligned to each task with higher-rank parameters, 533 which not only fully exploited the potential of LoRA module but also reduced the tuning parame-534 ter size in multi-task scenarios. To ensure the quality of rank experts, we used task embeddings to capture the distinctions and connections among different tasks. We also developed a CL-based 536 optimization target and Balanced Dataset Sampling strategy to ensure the fine-tuning quality. Ex-537 tensive experiments demonstrated that our method achieves significant improvements on the GLUE 538 benchmark and exhibits strong transfer learning performance. In the future, we plan to extend the application of MoRE and design a more efficient module to further improve its capability.

540 REFERENCES

- Akari Asai, Mohammadreza Salehi, Matthew E Peters, and Hannaneh Hajishirzi. Attempt:
 Parameter-efficient multi-task tuning via attentional mixtures of soft prompts. In *Proceedings* of the 2022 Conference on Empirical Methods in Natural Language Processing, pp. 6655–6672, 2022.
- Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients
 through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*, 2013.
- 548
 549
 549
 550
 550
 551
 Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pp. 7432–7439, 2020.
- Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal,
 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
 few-shot learners. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, pp. 1877–1901, 2020.
- Shijie Chen, Yu Zhang, and Qiang Yang. Multi-task learning in natural language processing: An overview. ACM Computing Surveys, 2021.
- Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. Boolq: Exploring the surprising difficulty of natural yes/no questions. In *Proceedings* of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pp. 2924–2936, 2019.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and
 Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge.
 arXiv preprint arXiv:1803.05457, 2018.
- Marie-Catherine de Marneffe, Mandy Simons, and Judith Tonhauser. The commitmentbank: Investigating projection in naturally occurring discourse. In *Proceedings of Sinn und Bedeutung 23*, 2019.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 4171–4186, 2019.
- 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 *Proceedings of the 2023 Confer*-*ence on Empirical Methods in Natural Language Processing*, pp. 4133–4145, 2023.
- ⁵⁷⁸ 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.
- Chengsong Huang, Qian Liu, Bill Yuchen Lin, Chao Du, Tianyu Pang, and Min Lin. Lorahub:
 Efficient cross-task generalization via dynamic lora composition. In *R0-FoMo: Robustness of Few-shot and Zero-shot Learning in Large Foundation Models*, 2023.
- Zixuan Ke, Yijia Shao, Haowei Lin, Tatsuya Konishi, Gyuhak Kim, and Bing Liu. Continual pretraining of language models. In *The Eleventh International Conference on Learning Representations*, 2023.
- Tushar Khot, Ashish Sabharwal, and Peter Clark. Scitail: A textual entailment dataset from science question answering. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.

603

- Dimitrios Kollias, Viktoriia Sharmanska, and Stefanos Zafeiriou. Distribution matching for multi-task learning of classification tasks: a large-scale study on faces & beyond. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pp. 2813–2821, 2024.
- 598 Dawid Jan Kopiczko, Tijmen Blankevoort, and Yuki M Asano. Vera: Vector-based random matrix adaptation. In *The Twelfth International Conference on Learning Representations*, 2024.
- Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 3045–3059, 2021.
- Dengchun Li, Yingzi Ma, Naizheng Wang, Zhiyuan Cheng, Lei Duan, Jie Zuo, Cal Yang, and
 Mingjie Tang. Mixlora: Enhancing large language models fine-tuning with lora based mixture of experts. *arXiv preprint arXiv:2404.15159*, 2024.
- Kiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. In
 Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers),
 pp. 4582–4597, 2021.
- Qidong Liu, Xian Wu, Xiangyu Zhao, Yuanshao Zhu, Derong Xu, Feng Tian, and Yefeng Zheng.
 Moelora: An moe-based parameter efficient fine-tuning method for multi-task medical applications. *CoRR*, 2023a.
- 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. In *Forty-first International Conference on Machine Learning*, 2024.
- Kiao Liu, Kaixuan Ji, Yicheng Fu, Weng Tam, Zhengxiao Du, Zhilin Yang, and Jie Tang. P-tuning: Prompt tuning can be comparable to fine-tuning across scales and tasks. In *Proceedings of the* 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), pp. 61–68, 2022.
- Kiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding, Yujie Qian, Zhilin Yang, and Jie Tang. Gpt understands, too. *AI Open*, 2023b.
- Rabeeh Karimi Mahabadi, Sebastian Ruder, Mostafa Dehghani, and James Henderson. Parameterefficient multi-task fine-tuning for transformers via shared hypernetworks. In *Proceedings of the*59th Annual Meeting of the Association for Computational Linguistics and the 11th International
 Joint Conference on Natural Language Processing (Volume 1: Long Papers), pp. 565–576, 2021.
- 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.
- Linyong Nan, Dragomir Radev, Rui Zhang, Amrit Rau, Abhinand Sivaprasad, Chiachun Hsieh, Xiangru Tang, Aadit Vyas, Neha Verma, Pranav Krishna, et al. Dart: Open-domain structured data record to text generation. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 432–447, 2021.
- 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. *Journal of machine learning research*, 21(140):1–67, 2020.
- 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. *arXiv preprint arXiv:2302.13971*, 2023.
- Mojtaba Valipour, Mehdi Rezagholizadeh, Ivan Kobyzev, and Ali Ghodsi. Dylora: Parameter efficient tuning of pre-trained models using dynamic search-free low-rank adaptation. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pp. 3274–3287, 2023.

648 649 650	Simon Vandenhende, Stamatios Georgoulis, Wouter Van Gansbeke, Marc Proesmans, Dengxin Dai, and Luc Van Gool. Multi-task learning for dense prediction tasks: A survey. <i>IEEE transactions on pattern analysis and machine intelligence</i> , 44(7):3614–3633, 2022.
652 653 654	Tu Vu, Brian Lester, Noah Constant, Rami Al-Rfou, and Daniel Cer. Spot: Better frozen model adap- tation through soft prompt transfer. In <i>Proceedings of the 60th Annual Meeting of the Association</i> <i>for Computational Linguistics (Volume 1: Long Papers)</i> , pp. 5039–5059, 2022.
655 656 657	Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. Glue: A multi-task benchmark and analysis platform for natural language understanding. In <i>Interna-</i> <i>tional Conference on Learning Representations</i> , 2018.
658 659 660	Yiming Wang, Yu Lin, Xiaodong Zeng, and Guannan Zhang. Multilora: Democratizing lora for better multi-task learning. <i>arXiv e-prints</i> , pp. arXiv–2311, 2023a.
661 662 663	Zhen Wang, Rameswar Panda, Leonid Karlinsky, Rogerio Feris, Huan Sun, and Yoon Kim. Mul- titask prompt tuning enables parameter-efficient transfer learning. In <i>The Eleventh International</i> <i>Conference on Learning Representations</i> , 2023b.
664 665 666	Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, An- drew M Dai, and Quoc V Le. Finetuned language models are zero-shot learners. In <i>International</i> <i>Conference on Learning Representations</i> , 2022.
668 669	Elad Ben Zaken, Shauli Ravfogel, and Yoav Goldberg. Bitfit: Simple parameter-efficient fine-tuning for transformer-based masked language-models. <i>arXiv preprint arXiv:2106.10199</i> , 2021.
670 671 672	Qingru Zhang, Minshuo Chen, Alexander Bukharin, Pengcheng He, Yu Cheng, Weizhu Chen, and Tuo Zhao. Adaptive budget allocation for parameter-efficient fine-tuning. In <i>International Conference on Learning Representations</i> . Openreview, 2023a.
674 675 676	Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang, Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi Hu, Tianwei Zhang, Fei Wu, et al. Instruction tuning for large language models: A survey. <i>arXiv preprint arXiv:2308.10792</i> , 2023b.
677 678	Tianyi Zhang, Felix Wu, Arzoo Katiyar, Kilian Q Weinberger, and Yoav Artzi. Revisiting few- sample bert fine-tuning. In <i>International Conference on Learning Representations</i> , 2021.
679 680 681 682	Yu Zhang and Qiang Yang. A survey on multi-task learning. <i>IEEE transactions on knowledge and data engineering</i> , 34(12):5586–5609, 2021.
683 684	
685	
686	
687	
688	
690	
691	
692	
693	
694	
695	
696	
697	
698	
699	
700	
701	

Method	DART	E2E	WebNLG	AVG
FT	46.1	<u>61.4</u>	44.2	50.6
$LoRA_{r=8}$	43.2	60.6	43.8	49.2
$LoRA_{r=16}$	44.6	60.8	44.3	49.9
MultiLoRA	44.0	61.3	44.9	50.1
MixLoRA	44.3	60.9	45.3	50.2
MoRE	<u>45.0</u>	61.5	<u>45.1</u>	<u>50.5</u>

Table 6: Model performance on NLG tasks

A ADDITIONAL EXPERIMENTS ON NLG

To further validate the effectiveness of our method, we conducted experiments on natural language generation (NLG) tasks using three datasets: DART, E2E, and WebNLG. DART focuses on gener-ating text from structured data, E2E involves generating restaurant descriptions from key attributes, and WebNLG is designed for generating text from knowledge graph triples. As shown in Table 6, none of the methods outperform fine-tuning (FT) on NLG tasks, and LoRA shows a significant per-formance drop. This indicates that using a fixed rank for training all tasks is suboptimal. In contrast, our method achieves performance comparable to FT. This is attributed to our method's ability to allocate an appropriate rank for different tasks efficiently.

B DETAILED CALCULATION OF PARAMETER COUNTS

LoRA parameters: LoRA employs matraix A and B to introduce low-rank adaptations in both the attention layers (q, k, v, o) and the feed-forward network (FFN) layers (w_i , w_o) of the T5-base model. Each LoRA layer has r(m+d) parameters. The total number of parameters for LoRA with L transformer layers is 6Lr(m+d).

MultiLoRA parameters: MultiLoRA employs parallel LoRA models for training, so its parameter count is *n* times that of vanilla LoRA, where *n* is the number of parallel LoRA modules. Additionally, MultiLoRA modifies the scaling factors to be learnable parameters (with parameter count *d*). Therefore, the total number of parameters is 6nLr(m + d) + 6Ld.

MixLoRA parameters: MixLoRA only employs parallel expert LoRA modules in the FFN layers and uses a gating module (with parameter count nm) to select the appropriate LoRA expert. Therefore, the total number of parameters is 2nLr(m + d) + 2Lnm.

741 **MOELoRA Parameters:** MOELoRA utilizes parallel LoRA models with a rank of r/n and in-742 corporates a task embedding module to represent each task (with a parameter count of Th). Addi-743 tionally, it employs a gating module (with a parameter count of nh) to compute the weights for each 100 LoRA. Therefore, the total number of parameters is given by 6Lr(m+d) + 6Lh(n+T).

MoRE parameters: Our proposed MoRE employs the same LoRA modules as vanilla LoRA, but treats LoRA modules with different ranks as experts, thereby introducing an additional gating mod-ule (with parameter count rh). To better adapt to different tasks, we also introduce a task embedding module (with parameter count Th, where h is the hidden dimension). Therefore, the total number of parameters is 6Lr(m+d) + 6Lh(r+T). In the GLUE dataset, T = 8 is consistent with r = 8. If the hidden dimension is set to be the same as d, then the parameter count is 12Lr(m+d), which is exactly the same as the parameter count with $LoRA_{r=16}$. Compared to MultiLoRA and MixLoRA, we do not use a parallel module design, so there is no parameter n that leads to a parameter count far exceeding that of LoRA. Furthermore, once our task embedding and gate modules are trained, we can construct a mapping from tasks to experts. This allows us to avoid the repeated computation of the task embedding and gate modules during inference, thereby reducing the parameter count to be consistent with $LoRA_{r=8}$.

756 C Additional Visualization of Task Embeddings

Further analysis of task embeddings is presented in Figures 4-6. These figures reveal that the patterns observed in other layers and modules of the model are consistent with those reported in the main text. Notably, stronger clustering is observed in the w_i and w_o layers. This enhanced clustering may be attributed to the feed-forward network (FFN) layers' ability to capture shared information underlying different tasks more effectively.



Figure 4: Visualization of Task Embeddings in Layer 1.



Figure 5: Visualization of Task Embeddings in Layer 6.



Figure 6: Visualization of Task Embeddings in Layer 12.