# Advancing Parameter Efficiency in Fine-tuning via Representation Editing

**Muling Wu, Wenhao Liu, Xiaohua Wang, Tianlong Li, Changze Lv**
**Zixuan Ling, Jianhao Zhu, Cenyuan Zhang, Xiaoqing Zheng[*], Xuanjing Huang**
School of Computer Science, Fudan University, Shanghai, China
{mlwu22,whliu22,xiaohuawang22,tlli22,czlv22}@m.fudan.edu.cn
{zhengxq,xjhuang}@fudan.edu.cn

## Abstract

Parameter Efficient Fine-Tuning (PEFT) techniques have drawn significant attention due to their ability to yield competitive results while updating only a small portion of the adjustable parameters. However, existing PEFT methods pose challenges in hyperparameter selection, such as choosing the rank for LoRA or Adapter, or specifying the length of soft prompts. To address these challenges, we propose a novel fine-tuning approach for neural models, named Representation EDiting (RED), which modifies the representations generated at some layers through the application of scaling and biasing operations. While existing PEFT methods still demonstrate over-parameterization that could potentially undermine the generalization ability acquired from pre-training, RED can substantially reduce the number of trainable parameters by a factor of $25,700$ compared to full parameter fine-tuning and by a factor of $32$ relative to LoRA. Remarkably, RED achieves results comparable or superior to both full parameter fine-tuning and other PEFT methods. Extensive experiments across various model architectures and scales, including RoBERTa, GPT-2, T5, and LLaMA-2, have demonstrated the effectiveness and efficiency of RED[1], thereby positioning it as a promising PEFT strategy for large-scale neural models.

## 1 Introduction

Pre-training on large-scale unlabeled datasets followed by fine-tuning on task-specific dataset has demonstrated remarkable efficacy across various natural language processing (NLP) tasks, establishing itself as the prevailing training paradigm (Devlin et al., 2018; Raffel et al., 2020; Radford et al., 2018). However, conducting full parameter fine-tuning for each task can be exceedingly costly

---

[*]Corresponding author.
[1]The code is accessible at https://github.com/mlwu22/RED.

and increasingly daunting as model scales continue to grow (Brown et al., 2020). For instance, BERT comprises up to 220 million parameters, T5 scales up to 11 billion parameters, and GPT-3 boasts an astounding 175 billion parameters. Consequently, the efficient and effective adaptation of large models to specific downstream tasks presents an intriguing research challenge (He et al., 2021).

In response to this challenge, researchers have put forward three main lines of Parameter Efficient Fine-Tuning (PEFT) techniques (Ding et al., 2022). Firstly, addition-based methods involve the introduction of additional trainable neural modules or parameters that were not present in the original model (Houlsby et al., 2019; Karimi Mahabadi et al., 2021; Li and Liang, 2021a; Lester et al., 2021a). Specification-based methods, on the other hand, identify certain parameters in the original model to be trainable, while the rest are kept frozen (Zaken et al., 2021; Guo et al., 2020). Lastly, reparameterization-based methods reconfigure trainable parameters into a more parameter-efficient form through certain transformations (Hu et al., 2021; Zhang et al., 2023a; Ding et al., 2023).

Among these PEFT methods, Low-Rank Adaptation (LoRA) stands out as one of the most efficient techniques with its effectiveness empirically validated across various models of diverse scales. Despite its impressive performance, LoRA still demands a significant number of trainable parameters. Recent studies by Aghajanyan et al. (2020) and Kopiczko et al. (2023) indicate that the upper bound for intrinsic dimensions is substantially smaller than what is typically used in such methods. For example, the $d_{90}$ value (the minimum number of trainable parameters required to reach $90\%$ of the performance of the fully-parameter fine-tuned model) for the RoBERTa base is reported to be 896. However, when LoRA is applied to fine-tune this model, the number of trainable parameters escalates to $0.3$ million, suggesting the potential for

further reduction in parameter count.

In addition to the issue of requiring too many adjustable parameters, existing PEFT methods (Mao et al., 2021; He et al., 2021; Ding et al., 2022) primarily focus on the design of lightweight modules and their integration (or placement) within base models. Nonetheless, the implementation of these PEFT techniques introduces additional complexities in hyperparameter selection, such as choosing the rank of LoRA and Adapter, or deciding on the length of Soft Prompt and Prefix.

Inspired by the concept of representation engineering (Zou et al., 2023), we shift our focus away from the weights of models and turn our attention to their representations. In the neural architecture, network weights govern neural activities (or representations), which in turn determine the networks' output, and the networks' output ultimately shapes the networks' behavior. Rather than concentrating on neurons and their interconnections (or weights), we explored to achieve control over network behavior by manipulating its internal representations. Specifically, we fine-tune neural network models by directly editing the representation generated at each layer while maintaining the model parameters frozen, as illustrated in Figure 1 (b). It is worth noting that the number of parameters required to edit representations is substantially fewer than that of weights within neural networks. Taking LLaMA-2 (7B) as an example, the proposed representation editing (RED) method achieves competitive performance by adjusting only 0.26M parameters. This is approximately $25,700$ times less than what is required for full-parameter fine-tuning, rendering the method both storage and computation efficient.

The contributions of this study are summarized as follows:

- We propose a novel perspective on fine-tuning by directly editing model representations, diverging from exiting PEFT methods that focused on adjusting the model's weights. Our proposed PEFT technique, termed RED, embodies this new perspective.
- Extensive experiments are conducted across models of varying structures and scales, including RoBERTa, GPT-2, T5, and LLaMA-2. The effectiveness of RED is validated across a range of natural language understanding and generation tasks. Notably, RED demonstrates both efficacy and efficiency while requiring only a minimal number of trainable parame-

ters and maintaining ease of implementation.
- A comprehensive ablation study is conducted to dissect the individual components of RED and understand their impacts on performance.

## 2 Related Work

Existing PEFT methods can be broadly classified into three categories (Ding et al., 2022), each characterized by the particular parts of parameters that they tune efficiently: addition-based, specification-based, and reparameterization-based methods.

Addition-based methods perform fine-tuning by integrating additional lightweight learnable components into foundational models. More specifically, Houlsby et al. (2019); Stickland and Murray (2019); Karimi Mahabadi et al. (2021) and Rücklé et al. (2020) proposed the integration of learnable bottleneck neural modules into the transformer layers. Brown et al. (2020) and Shin et al. (2020) found that the models' performance could be enhanced by appending discrete tokens (or prompts) at the beginning of the input texts without modifying the models' parameters. However, the manual creation of such prompts demands substantial effort, and optimization in a discrete space is comparatively more challenging. Consequently, later studies substituted these discrete tokens with continuous vectors, also referred to as soft prompts, which can be optimized using the gradient descent algorithm (Lester et al., 2021b; Li and Liang, 2021b; Wu et al., 2023; Wang et al., 2023).

Specification-based methods achieve parameter-efficient fine-tuning by designating a subset of parameters for modification while leaving the remaining parameters untouched. Within this category, Lee et al. (2019) suggested a method to exclusively update the parameters in certain top layers of BERT and RoBERTa. BitFit (Ben-Zaken et al., 2021), on the other hand, fine-tunes a model by only optimizing the bias terms used in the model. Contrary to these methods that pre-determine which parameters should be altered, Guo et al. (2020) and Zhao et al. (2020) implemented a learnable masking strategy to automatically choose the parameters for tuning.

Reparameterization-based methods aim to optimize some parameters within a model, typically in their low-dimensional subspace. Hu et al. (2021) introduces the use of low-rank matrices, termed LoRA, to approximate the weight modifications during the fine-tuning process. QLoRA (Dettmers et al., 2023) combines low-rank adaptation with
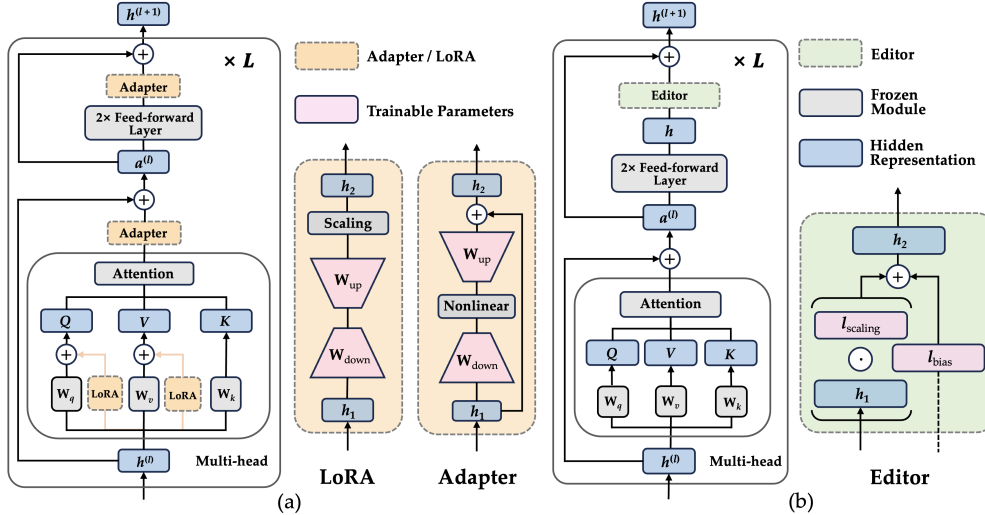
Figure 1: Comparison of previous representative PEFT methods with the proposed RED. (a) LoRA incorporates learnable bottleneck-shaped modules (highlighted in orange) by integrating additional connections parallel to the $\mathbf{W}_q$ and $\mathbf{W}_v$ matrices of attention blocks, along with modifying the weights of these matrices in a low-rank fashion. Adapter, on the other hand, introduces learnable modules within similar structures (also highlighted in orange) by incorporating additional connections following both the attention and feed-forward sub-layers. (b) RED introduces two learnable vectors, $l_{\text{scaling}}$ and $l_{\text{bias}}$, to directly edit the representations (marked in green) generated by feed-forward sub-layers, which significantly reduces the number of parameters required for fine-tuning.

model quantization to further reduce storage usage during the fine-tuning phase. AdaLoRA (Zhang et al., 2023b) proposes to use SVD decomposition to approximate the changes in weights, which allocates more trainable parameters to more important weight matrices, leading to a better performance.

Zou et al. (2023) proposed an approach to analyzing and manipulating the behavior of neural networks through representation engineering, thereby shifting the focus from neurons and their circuits to representations and the transformations among them. Liu et al. (2023) extended this work to align large language models with human preferences, achieving results comparable or superior to RLHF (Ouyang et al., 2022), but at a lower computational cost. Subramani et al. (2022) investigated the extraction of "steer vectors" from the hidden layers and performed unsupervised text style transfer by altering the hidden representations using these vectors. To control the style and sentiment expressed in a model's outputs, Turner et al. (2023) integrated a similar "steer vector" into the representation of each hidden layer at the inference time. Inspired by the recent advance in representation engineering Zou et al. (2023), we suggested a novel perspective for fine-tuning models by directly editing their representations. This new perspective was embodied in the RED training method, which demonstrated both effectiveness and efficiency while requiring only a minimal number of trainable parameters.

This study is closely related to two recent works:

IA3 (Liu et al., 2022) and SSF (Lian et al., 2022). To fulfill the in-context learning capability with lower computational costs, Liu et al. (2022) proposed to modify the key and value vectors of the multi-head attentions and those generated by feed-forward networks (FFNs) by scaling operations. In the domain of computer vision, Lian et al. (2022) presented a similar method, called SSF, which was motivated by batch and layer normalization techniques. To ensure features fall in a discriminative space for better classification, SSF also necessitates the adjustment of feature vectors across nearly all transformer layers, including multi-head attentions, FFNs, and normalization layers. As a result, the size of parameters adjusted in their method approximates those manipulated in other PEFT techniques, such as Adapter and VPT (Jia et al., 2022).

In contrast, our RED method requires only the modification of activation patterns generated by FFN sub-layers, constituting a significantly smaller portion of the entire neural network's parameters for model fine-tuning. While SSF achieves a reduction in the number of fine-tuning parameters by a factor of 300 compared to original models, RED accomplishes a reduction of approximately 25,700 times. Moreover, SSF's evaluation was restricted to a limited range of simple image classification tasks using relatively small models (under 100M parameters), whereas RED has demonstrated its effectiveness across a variety of natural language understanding and generation tasks, with signifi-

cantly larger models of up to 7B parameters.

## 3 Method

We in this section give a concise overview of previous PEFT methods and introduce a new, parameter-efficient fine-tuning method known as Representation Editing (RED) by highlighting its distinctive features compared to existing PEFT methods. RED facilitates the adaptation of pre-trained neural network models to downstream tasks by directly editing the model's representations.

### 3.1 Recap of PEFT Methods

The predominant large language models (LLMs) have been constructed based on the transformer architecture proposed in a seminal work on attention mechanisms (Vaswani et al., 2017). This architecture is characterized by a series of layered blocks, each incorporating two fundamental sub-layers: multi-head self-attention and fully connected feed-forward networks. Many existing PEFT methods usually achieve fine-tuning by adjusting the parameters of these two sub-layers in some parameter-efficient manner. Figure 1 (a) illustrates two exemplar PEFT methods: Adapter and LoRA. With the exception of a few additional parameters that require tuning, the parameters of the pre-trained LLM remain fixed.

Specifically, LoRA (Hu et al., 2021) incorporates learnable bottleneck-shaped modules through the connections parallel to the $\mathbf{W}_q$ and $\mathbf{W}_v$ matrices of attention blocks, as illustrated in Figure 1 (a), and models the weight changes of these two matrices in a low-rank manner. Given a pre-trained weight matrix $\mathbf{W} \in \mathbb{R}^{d \times k}$, LoRA represents its changes, denoted as $\Delta \mathbf{W}$, through two low-rank decomposition matrices: $\Delta \mathbf{W} = \alpha \cdot \mathbf{W}_{\text{down}} \mathbf{W}_{\text{up}}$, where $\mathbf{W}_{\text{down}} \in \mathbb{R}^{d \times r}$ and $\mathbf{W}_{\text{up}} \in \mathbb{R}^{r \times k}$ (the value of $r$ is typically much smaller than both $d$ and $k$). The coefficient $\alpha$ is a hyperparameter that determines the significance of $\Delta \mathbf{W}$ relative to the original $\mathbf{W}$. Given an input $x$, LoRA produces the result of the forward pass as follows:

$$h = x\mathbf{W} + \alpha \cdot x\mathbf{W}_{\text{down}}\mathbf{W}_{\text{up}} \qquad (1)$$

In the initial version of the Adapter, Houlsby et al. (2019) suggested the insertion of trainable adapter modules between two sub-layers within the transformer architecture. These adapter modules are structured to include a down-projection matrix, denoted as $\mathbf{W}_{\text{down}} \in \mathbb{R}^{d \times r}$, which transforms a hidden representation $h_1 \in \mathbb{R}^d$ to a lower-dimensional

space with a pre-specified dimensionality $r$. Subsequently, this dimensionally reduced vector undergoes a nonlinear activation function $f(\cdot)$ and an additional up-projection matrix $\mathbf{W}_{\text{up}} \in \mathbb{R}^{r \times d}$ to revert it to its original dimension $d$. The adapter module also incorporates a residual connection. The resultant output $h_2 \in \mathbb{R}^d$ generated by this module can be formally represented as follows:

$$h_2 = h_1 + f(h_1\mathbf{W}_{\text{down}})\mathbf{W}_{\text{up}} \qquad (2)$$

Expanding on this research, Pfeiffer et al. (2020) proposed a more efficient variant of the Adapter that is only applied following the FFN sub-layer.

### 3.2 Representation Editing

Previous PEFT methods refine pre-trained models by updating their weights in a parameter-efficient manner, typically within a low-dimensional space. To approximate the effect of full-parameter tuning, they are required to choose the values of hyperparameters properly for one downstream task or a set of tasks. However, it could be troublesome to choose a suitable value for a hyperparameter, such as the ranks of weight matrices in Equations (1) and (2) for LoRA and Adapter respectively. Enhancing their modeling capacities by using a higher rank $r$ could demand too much computation resources and tend to overfit, while aggressively setting $r$ smaller may degrade model performance and lead to from-scratch re-training (Ding et al., 2023).

Our hypothesis posits that by altering the internal representations of neural models rather than their connected weights, the fine-tuning process could be significantly more efficient. This is because that direct manipulation of the representations necessitates $O(n)$ parameters, whereas adjusting their weights theoretically demands $O(n^2)$ parameters, with $n$ denoting the dimensionality of hidden representations. Turner et al. (2023) have demonstrated that the behaviors of neural models can be influenced by adding a "steer vector" to each hidden layer during inference. We postulate that such steer vectors can be learned during the fine-tuning phase. Inspired by this concept and the emerging field of representation engineering Zou et al. (2023), we introduce a novel PEFT method that fine-tunes the model by directly altering the representation with two learnable vectors, as depicted in Figure 1 (b).

Specifically, we first incorporate a learnable scaling vector $l_{\text{scaling}} \in \mathbb{R}^d$ and apply it to perform the Hadamard product with a hidden representation $h_1$

| Method | # Param | MNLI | SST-2 | MRPC | CoLA | QNLI | QQP | RTE | STS-B | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|
| FT (base) | 125M | 87.3 | 94.4 | 87.9 | 62.4 | 92.5 | 91.7 | 78.3 | 90.6 | 85.6 |
| Adapter (base) | 0.4M | 87.0 | 93.3 | 88.4 | 60.9 | 92.5 | 90.5 | 76.5 | 90.5 | 85.0 |
| LoRA (base) | 0.3M | 86.6 | 93.9 | 88.7 | 59.7 | 92.6 | 90.4 | 75.3 | 90.3 | 84.7 |
| Adapter-FFN (base) | 0.3M | 87.1 | 93.0 | 88.8 | 58.5 | 92.0 | 90.2 | 77.7 | 90.4 | 84.7 |
| BitFit (base) | 0.1M | 84.7 | 94.0 | 88.1 | 54.0 | 91.0 | 87.3 | 69.8 | 89.5 | 82.3 |
| IA3 (base) | 0.06M | 85.4 | 93.4 | 86.4 | 57.8 | 91.1 | 88.5 | 73.5 | 88.5 | 83.1 |
| RED (base) | 0.02M | 83.9 | 93.9 | 89.2 | 61.0 | 90.7 | 87.2 | 78.0 | 90.4 | 84.3 |
| FT (large) | 355M | 88.8 | 96.0 | 91.7 | 68.2 | 93.8 | 91.5 | 85.8 | 92.6 | 88.5 |
| Adapter (large) | 0.9M | 90.1 | 95.2 | 90.5 | 65.4 | 94.6 | 91.4 | 85.3 | 91.5 | 88.0 |
| LoRA (large) | 0.8M | 90.2 | 96.0 | 89.8 | 65.5 | 94.7 | 90.7 | 86.3 | 91.7 | 88.1 |
| Adapter-FFN (large) | 0.8M | 90.3 | 96.1 | 90.5 | 64.4 | 94.3 | 91.3 | 84.8 | 90.2 | 87.7 |
| IA3 (large) | 0.15M | 90.1 | 94.5 | 87.1 | 63.2 | 93.9 | 89.3 | 85.3 | 91.5 | 86.9 |
| RED (large) | 0.05M | 89.5 | 96.0 | 90.3 | 68.1 | 93.5 | 88.8 | 86.2 | 91.3 | 87.9 |

Table 1: Performance comparison of RoBERTa base and large models fine-tuned by RED and other PEFT baseline methods on the GLUE benchmark. The standard deviations of results across different methods are given in Table 17.

| Method | # Param | BLEU | NIST | MET | ROUGE-L | CIDEr |
|---|---|---|---|---|---|---|
| FT (medium) | 355M | 65.95 | 8.52 | 45.95 | 69.13 | 2.35 |
| FT$^{\text{top2}}$ (medium) | 25.2M | 65.94 | 8.53 | 44.28 | 68.83 | 2.23 |
| Adapter (medium) | 0.9M | 64.31 | 8.29 | 44.91 | 67.72 | 2.28 |
| LoRA (medium) | 0.8M | 67.43 | 8.65 | 46.01 | 69.64 | 2.42 |
| Adapter-FFN (medium) | 0.8M | 64.41 | 8.30 | 44.74 | 67.53 | 2.29 |
| Prefix Tuning (medium) | 0.8M | 63.92 | 8.26 | 41.81 | 66.86 | 2.03 |
| IA3 (medium) | 0.17M | 63.63 | 7.99 | 40.49 | 66.36 | 1.89 |
| RED (medium) | 0.05M | 64.86 | 8.36 | 44.99 | 67.62 | 2.28 |
| FT (large) | 774M | 65.56 | 8.50 | 45.40 | 68.38 | 2.27 |
| Adapter (large) | 1.8M | 65.94 | 8.46 | 45.78 | 68.65 | 2.34 |
| LoRA (large) | 1.5M | 68.24 | 8.76 | 46.23 | 69.92 | 2.42 |
| Adapter-FFN (large) | 1.5M | 65.53 | 8.41 | 45.65 | 68.46 | 2.33 |
| Prefix Tuning (large) | 1.5M | 65.50 | 8.45 | 43.97 | 67.32 | 2.23 |
| IA3 (large) | 0.32M | 65.08 | 8.5 | 42.72 | 66.80 | 2.15 |
| RED (large) | 0.09M | 65.77 | 8.42 | 46.12 | 69.03 | 2.36 |

Table 2: Performance comparison of GPT-2 medium and large models fine-tuned by RED and other PEFT baseline methods on the E2E NLG Challenge. The standard deviations of results across different methods are given in Table 18.

by scaling the feature of each dimension within $h_1$ via element-wise multiplication. Additionally, we introduce another learnable bias vector $l_{\text{bias}} \in \mathbb{R}^d$ that is subsequently added to the scaled vector. This process can be formalized as follows:

$$h_2 = l_{\text{scaling}} \odot h_1 + l_{\text{bias}} \qquad (3)$$

where $\odot$ represents element-wise multiplication, also known as the Hadamard product. Here, $h_1 \in \mathbb{R}^d$ denotes a hidden representation generated by a certain layer, and $h_2 \in \mathbb{R}^d$ represents the resultant edited representation. During implementation, we initialize the scaling vectors $l_{\text{scaling}}$ as unit vectors (i.e., with all elements set to 1) and the bias vectors $l_{\text{bias}}$ as zero vectors. This initialization approach ensures that the introduction of these "edit vectors" does not initially alter the hidden representations of a neural model.

## 4 Experiments

Extensive experimentation was conducted to assess the efficacy of our Representation Editing (RED) method across a diverse array of natural language understanding and generation tasks, employing a range of foundational models spanning different scales. These models include RoBERTa (Liu et al., 2019), T5 (Raffel et al., 2020), GPT-2 (Radford et al., 2019), and LLaMA-2 (Touvron et al., 2023). Specifically, we evaluated RED and the baseline methods on the GLUE benchmark (Wang et al., 2018) with RoBERTa and T5, as previously conducted in (Hu et al., 2021) and Asai et al. (2022). To ensure consistency with previous studies, we adhered to the experimental setup outlined in (Li and Liang, 2021a) and (Hu et al., 2021) for comparative analysis with GPT-2. Moreover, we conducted instruction tuning experiments on LLaMA-2 using the UltraFeedback dataset (Cui et al., 2023) to further assess the applicability of our proposed RED on generative large language models. For further details on the datasets and evaluation metrics used, please refer to Appendix A.

### 4.1 Baselines

We ensured that different PEFT methods were compared systematically in a more fair setting. There-

| Method | # Param | MNLI | SST-2 | MRPC | CoLA | QNLI | QQP | RTE | STS-B | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|
| FT (base)* | 220M | 86.8 | 94.6 | 90.2 | 61.8 | 93.0 | 91.6 | 71.9 | 89.7 | 84.9 |
| Adapter (base)* | 1.9M | 86.5 | 93.8 | 85.3 | 64.0 | 93.2 | 90.2 | 71.9 | 90.7 | 84.5 |
| AdapterDrop (base)* | 1.1M | 86.3 | 93.6 | 86.3 | 62.7 | 93.2 | 90.2 | 71.2 | 91.4 | 84.4 |
| BitFit (base)* | 0.3M | 85.3 | 94.2 | 86.8 | 58.2 | 93.0 | 90.1 | 67.6 | 90.9 | 83.3 |
| PT (base)* | 0.08M | 81.3 | 90.9 | 68.1 | 10.6 | 92.8 | 89.7 | 54.7 | 89.5 | 72.2 |
| RED (base) | 0.04M | 85.9 | 93.0 | 91.7 | 61.1 | 91.2 | 89.2 | 72.7 | 88.2 | 84.1 |

Table 3: Performance comparison of T5 base fine-tuned by RED and other PEFT baseline methods on the GLUE benchmark. Results marked with an asterisk (*) are excerpted from published literature.

fore, we strictly followed the well-established training protocol, evaluating on the validation set after each epoch training, and selecting the checkpoint with the best performance on the validation set as the final model for testing on the test set. The following baseline methods were used for comparison with the proposed RED:

- **Fine-Tuning (FT)** trains models by updating all their parameters. A variant of FT was proposed by Lee et al. (2019), which selectively updates certain layers while freezing others. We incorporate a baseline established in prior research by Li and Liang (2021a) with GPT-2, which specifically adapts only the final two layers, denoted as $FT^{top2}$.
- **Bias-terms Fine-tuning (BitFit)** involves the selective freezing of a majority of the transformer parameters, with the training process focused exclusively on the bias-terms (Ben-Zaken et al., 2021).
- **IA3** introduces scaling operations to modify the key and value vectors of the multi-head attentions and those generated by feed-forward networks (FFNs) (Liu et al., 2022).
- **Adapter** introduces a learnable, lightweight module situated between two sub-layers of the transformer. During the forward pass, inputs are processed in sequence by the sub-layers of the foundation models and the adapters, yielding the final output. However, during the back-propagation phase, only these adapters receive gradients for parameter updates, while the remaining parameters of the model are kept fixed and unaltered (Houlsby et al., 2019).
- **Adapter-FFN** is a variant of Adapter method proposed by Pfeiffer et al. (2020). Contrasting with the original Adapter, which necessitates the insertion of the learnable module between all sub-layers, Adapter-FFN only requires the application of an adapter following each Feed-Forward Network (FFN) sub-layer.
- **AdapterDrop** is another variant of Adapter

proposed by Rücklé et al. (2020), which incorporates a strategy of omitting certain adapter layers, thereby enhancing overall efficiency.
- **Low-Rank Adaption (LoRA)** employs a low-rank decomposition on the matrix $\Delta W$, thereby modeling weight updates as the product of two low-rank matrices. These two learnable matrices are aligned in parallel with the corresponding matrices in pre-trained models. They process inputs in parallel and combine their results to generate the final outputs in each transformer block (Hu et al., 2021).
- **Prompt Tuning (PT)** prefixes a set of continuous vectors at the embedding layer, which are subjected to learning during the fine-tuning phase (Lester et al., 2021b).
- **Prefix Tuning** is a generalized version of prompt tuning (PT), which incorporates learnable continuous vectors at every hidden state. These continuous vectors also contribute to the computation of attention, serving as both key and value vectors (Li and Liang, 2021b).

## 4.2 Results with RoBERTa

We conducted fine-tuning experiments on both the RoBERTa base (125M) and large (355M) models by using RED and its competing baselines. Subsequently, we evaluated these fine-tuned models on the widely-adopted GLUE benchmark, renowned for its comprehensive evaluation of natural language understanding capabilities. The pretrained RoBERTa models were sourced from the Hugging-Face Transformers Library (Wolf et al., 2019).

We noticed that the previous evaluation settings of PEFT methods encountered two issues. Firstly, for datasets such as MRPC, RTE, and STS-B, they trained the models on the MNLI dataset first and selecting the best performance checkpoint on the MNLI dataset as initialization for transfer learning to improve the performance of the model trained on these datasets (Liu et al., 2019; Hu et al., 2021). However, such a pipeline adds complexity that can be challenging for other researchers to replicate.

Secondly, there is no split between the validation set and the test set. Instead, after each epoch training is completed, evaluation is directly conducted on the test set, and the best test result is selected as the final evaluation result of the model, which violates the established standard that the test set should not influence model selection during the training phase. In this study, we ensured a more systematic comparison of different PEFT methods in a fairer setting. To achieve this, we adhered strictly to a well-established training protocol. After each epoch of training, evaluation was conducted on the validation set, and the checkpoint demonstrating the best performance on the validation set was selected as the final model for testing on the test set. For comprehensive details regarding our reimplementation, please refer to Appendix B.1.

As presented in Table 1, both the RoBERTa base and large models fine-tuned by RED yielded accuracies comparable to those achieved by other PEFT methods across all tasks within the GLUE benchmark. For instance, RED's performance was only marginally inferior to that of LoRA, differing by a negligible 0.2% with the RoBERTa large and by 0.4% with the RoBERTa base on average, while requiring the tuning of significantly fewer parameters. It is noteworthy that RED demonstrated superior performance on tasks such as SST-2, MRPC, CoLA, STS-B, and RTE, all of which had data sizes less than 100k. This suggests that RED helps to maintain the generalization capability acquired during the pre-training phase and can deliver enhanced performance when the volume of training data is relatively small.

Moreover, RED exhibits an unprecedented level of parameter efficiency. Although having substantially fewer trainable parameters—approximately $7,200$ times fewer than full parameter fine-tuning and 16 times fewer than LoRA—it maintains comparable performance. This observation suggests that there is still potential for further reduction in the number of trainable parameters and this finding aligns with conclusions drawn by Aghajanyan et al. (2020) and Kopiczko et al. (2023).

### 4.3   Results with GPT-2

Beyond natural language understanding tasks, our study expanded to include experiments on natural language generation tasks. These experiments were conducted using the GPT-2 medium (355M) and large (774M) models on the E2E NLG Challenge (Novikova et al., 2017). The pre-trained GPT-2

models were also obtained from the HuggingFace Transformers Library. To ensure a fair comparison, we reproduced other PEFT methods following the settings defined by Li and Liang (2021a) and Hu et al. (2021). Please refer to Appendix B.2 for the comprehensive details of our implementation.

The data presented in Table 2 reveals that RED achieved performance comparable to other PEFT baselines across all metrics in the E2E NLG Challenge. This highlights the efficacy of fine-tuning through representation editing, not only for natural language understanding tasks but also for language generation tasks. Remarkably, RED accomplished this performance while still necessitating minimal fine-tuning of parameters. To compare RED with the most parameter-efficient variants of other PEFT methods, we set the rank of LoRA and Adapter to 1. Under these conditions, RED outperformed all such variants while still employing the fewest number of parameters (see Section 5.3 for details).



Figure 2: Performance scores achieved by RED and other PEFT methods on the MT-Bench. Refer to Table 19 and Appendix D for raw scores and additional details.

| Method | # Param | AlpacaEval (win %) |
|--------|---------|--------------------|
| FT | 6739M | 80.93 |
| LoRA | 8.39M | 81.48 |
| RED | 0.26M | 81.69 |

Table 4: Win rates against reference responses evaluated by GPT-4 on AlpacaEval. Higher win rates indicate superior generated responses, reflecting better alignment with human preferences. Example responses generated by RED and other PEFT baselines can be found in Figures 3 and 4 as well as in Appendix D.

### 4.4   Results with T5

In order to assess the versatility of RED, we conducted experiments using the encoder-decoder architecture. More specifically, we obtained the pre-

| Method | # Param | BLEU | NIST | MET | ROUGE-L | CIDEr |
|---|---|---|---|---|---|---|
| Adapter (rank 1) | 0.25M | 63.76 | 8.37 | 42.74 | 66.70 | 2.09 |
| Adapter-FFN (rank 1) | 0.07M | 62.99 | 8.09 | 40.88 | 66.39 | 1.98 |
| LoRA (rank 1) | 0.10M | 64.51 | 8.38 | 44.78 | 67.35 | 2.28 |
| **RED** | **0.05M** | **64.86** | **8.36** | **44.99** | **67.62** | **2.28** |

Table 5: Comparison of GPT-2 fine-tuned by RED and other PEFT baselines (rank 1) on the E2E NLG Challenge.

trained T5-base model (220M) from the Hugging-Face Transformers Library and evaluated models fine-tuned with RED and other PEFT baselines on the GLUE benchmark. As shown in Table 3, RED achieved accuracies that were on par with other PEFT methods across all tasks in the GLUE benchmark, indicating its versatility across diverse network architectures. We excerpted the results from the study of (Asai et al., 2022) as they conducted experiments under the settings similar to ours. The inclusion of their results does not undermine the fairness or outcome of the comparison. Please refer to Appendix B.3 for the implementation details.

It is noteworthy that RED surpassed the prompt tuning (PT) baseline with a significant margin of 11.9% on average. Despite PT methods necessitating the fewest parameters during fine-tuning among other PEFT baselines, they still require a modification of parameters that is twice as large as that required by RED. This finding indicates that allocating a few number of learnable parameters to each layer for representation editing is both more parameter-efficient and effective compared to the PT method, which assigns all learnable parameters to the model's embedding layer.

## 4.5 Results with LLaMA-2

We are interested in exploring whether RED can be effectively scaled to large neural networks containing billions of parameters. Therefore, we evaluated RED and compared it against full parameter fine-tuning and LoRA using LLaMA-2 up to 7 billion parameters. Our assessments were conducted across three widely-used benchmarks: Open LLM Leaderboard (Beeching et al., 2023), AlpacaEval (Li et al., 2023), and MT-Bench (Zheng et al., 2023). The implementation details of this experiment can be found in Appendix B.4.

Table 4 presents the win rates obtained by different fine-tuning methods, compared against the reference responses generated by text-davinci-003 on the AlpacaEval. Remarkably, RED yielded the highest win rate, even though its number of trainable parameters was 25, 700 times less than that of full parameter fine-tuning and 32 times less than

that of LoRA. This suggests that the RED method, which directly edits hidden representations during the fine-tuning phase, continues to be effective for large language models, and the trained models are capable of generating high-quality responses.

Moreover, Figure 2 illustrates the performance scores achieved by three training methods on the 1-turn question-answer dataset from MT-Bench. The overall performance of RED is on par with other baseline methods, and it notably excels in assessing the capabilities of humanities and reasoning. As presented in Table 12, RED also delivers competitive results on six distinct datasets from the Open LLM Leaderboard.

| Method | MRPC | CoLA | QQP | Avg. |
|---|---|---|---|---|
| **RED** | 90.3 | 68.1 | 88.7 | 82.4 |
| **w/o scaling** | 89.8 | 65.9 | 87.6 | 81.1 |
| **w/o bias** | 75.8 | 46.9 | 87.2 | 70.0 |

Table 6: Results of the ablation studies on three different datasets from the GLUE benchmark. The term "w/o scaling" denotes the removal of scaling operations, with representation editing performed only by bias vectors. In contrast, "w/o bias" indicates the omission of bias vectors, with representation editing conducted exclusively through the use of scaling vectors.

## 5 Ablation Study

In this section, we conducted a series of ablation studies to examine the impact of different editing operations used in RED and to explore potential application points within transformer-based networks. We also evaluated the parameter efficiency of RED by comparing it with the most parameter-efficient variants of other PEFT methods.

## 5.1 Impact of Different Editing Operators

We introduce two operators to edit the hidden representations of neural network models: one scales a hidden representation using a learnable vector, and the other adds a bias vector to the hidden representation. To assess the contribution of each operator, we performed ablation studies by removing one operator at a time and examining the impact on the models' performance.

From Table 6, it is evident that both editing operators are crucial for enhancing model performance.

| Method | # Param | MNLI | SST-2 | MRPC | CoLA | QNLI | QQP | RTE | STS-B | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|
| LoRA (base) | 0.29M | 86.6 | 93.9 | 88.7 | 59.7 | 92.6 | 90.4 | 75.3 | 90.3 | 84.7 |
| RED (base) | 0.02M | 83.9 | 93.9 | 89.2 | 61.0 | 90.7 | 87.2 | 78.0 | 90.4 | 84.3 |
| RED+ (base) | 0.09M | 85.9 | 93.2 | 88.9 | 62.4 | 90.9 | 89.1 | 78.6 | 90.9 | 85.0 |
| LoRA (large) | 0.79M | 90.2 | 96.0 | 89.8 | 65.5 | 94.7 | 90.7 | 86.3 | 91.7 | 88.1 |
| RED (large) | 0.05M | 89.5 | 96.0 | 90.3 | 68.1 | 93.5 | 88.8 | 86.2 | 91.3 | 87.9 |
| RED+ (large) | 0.25M | 90.6 | 95.6 | 89.2 | 69.6 | 94.0 | 90.2 | 85.6 | 91.4 | 88.3 |

Table 7: Performance comparison of RoBERTa base and large models fine-tuned by RED, RED+ and other PEFT baselines on the GLUE benchmark. RED+ denotes an enhanced variant of RED where the representation vectors of $Q$, $K$, and $V$ in attention blocks are also modified via representation editing in addition to FFN sub-layers.

| Method | # Param | BLEU | NIST | MET | ROUGE-L | CIDEr |
|---|---|---|---|---|---|---|
| LoRA (medium) | 0.79M | 67.43 | 8.65 | 46.01 | 69.64 | 2.42 |
| RED (medium) | 0.05M | 64.86 | 8.36 | 44.99 | 67.62 | 2.28 |
| RED+ (medium) | 0.25M | 66.68 | 8.53 | 46.28 | 69.63 | 2.38 |
| LoRA (large) | 1.47M | 68.24 | 8.76 | 46.23 | 69.92 | 2.42 |
| RED (large) | 0.09M | 65.77 | 8.42 | 46.12 | 69.03 | 2.36 |
| RED+ (large) | 0.46M | 68.31 | 8.78 | 46.12 | 69.80 | 2.41 |

Table 8: Performance comparison of GPT-2 medium and large models fine-tuned by RED, RED+ and other PEFT baselines on the E2E NLG Challenge. RED+ denotes an enhanced variant of RED where the representation vectors of $Q$, $K$, and $V$ in attention blocks are also modified via representation editing in addition to FFN sub-layers.

Analyzing results from three different datasets from the GLUE benchmark, we found that the scaling operator boosted accuracy by an average of $1.3\%$, while the biasing operator contributed to a significant $12.4\%$ increase in accuracy. Compared to the scaling operator, the addition of a learnable bias vector to the original representation plays a more substantial role in the fine-tuning of neural models.

## 5.2 Impact of Editing Positions

The motivation of this study is to explore the feasibility of fine-tuning large language models with as few parameters as possible. Therefore, in our current implementation, we only consider modifying the representations produced by FFN sub-layers. However, it is still possible to further improve the performance of the model by increasing the number of trainable parameters if the representation vectors of $Q$ (query), $K$ (key), and $V$ (value) in attention blocks are allowed to be modified via representation editing. To demonstrate this possibility, we conducted additional experiments in which the representations of $Q$, $K$, and $V$ were also modified.

The experimental results reported in Table 7 and Table 8 show that the performance of RED has been improved on all these datasets but with slightly more trainable parameters (increasing from 0.07M to 0.37M). It suggests that it is feasible to increase the number of trainable parameters to further boost the performance by editing additional feature representations. This strategy is particularly effective for datasets with a larger number of training instances, such as MNLI, QNLI, QQP, and E2E while the performance did not exhibit significant improvement on datasets with fewer training examples.

## 5.3 Parameter Efficiency and Efficacy

When using Adaptor or LoRA for fine-tuning, the ranks of their decomposition matrices were typically set to 8 or 16. These ranks can be reduced to decrease the number of parameters used during fine-tuning. To compare RED with the most parameter-efficient versions of Adpator or LoRA, we adjusted the rank of their decomposition matrices to 1. We then made a comparative analysis with GPT-2 medium on the E2E NLG Challenge.

Table 5 shows that RED outperformed all the baselines on four out of five tasks, while it slightly underperformed LoRA by a negligible difference of 0.02 on the NIST dataset. It is worthy noting that RED accomplished this with the minimal number of parameters, demonstrating that RED is not only parameter-efficient but also parameter-effective.

## 6 Conclusion

In this study, we proposed a novel perspective on fine-tuning by editing model representations rather than adjusting the model's weights. Through extensive experimentation across various NLP tasks and neural models of varying structures and scales, we demonstrated that RED can deliver performance comparable or superior to existing PEFT methods while employing a minimal number of trainable parameters. We are hopeful that this study could inspire novel methods for fine-tuning LLMs, thereby contributing to the advancement of both effective and computationally efficient PEFT techniques.

## Limitations

We have demonstrated the efficacy of a new promising PEFT approach for fine-tuning models through direct manipulation of representations across various NLP tasks with networks of varying structures and scales. It would be intriguing to extend this method to other modalities, such as image, speech, and video. Recent advancements in representation engineering have indicated that only a very few examples are sufficient for achieving precise control over the model's output through representation manipulations. We plan to apply the proposed RED to few-shot scenarios in the future.

## Reproducibility Statement

We have made our code publicly accessible via a repository on GitHub (github.com/mlwu22/RED). To further ensure replicability, we had a colleague unfamiliar with our method install and test RED. The experiment yielded results nearly identical to ours, which strengthens our confidence that other researchers will successfully execute our code and reproduce our findings.

## References

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

Akari Asai, Mohammadreza Salehi, Matthew E. Peters, and Hannaneh Hajishirzi. 2022. Attempt: Parameter-efficient multi-task tuning via attentional mixtures of soft prompts. In *Conference on Empirical Methods in Natural Language Processing*.

Satanjeev Banerjee and Alon Lavie. 2005. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *IEEvaluation@ACL*.

Roy Bar-Haim, Ido Dagan, Bill Dolan, Lisa Ferro, Danilo Giampiccolo, Bernardo Magnini, and Idan Szpektor. 2006. The second pascal recognising textual entailment challenge.

Edward Beeching, Clémentine Fourrier, Nathan Habib, Sheon Han, Nathan Lambert, Nazneen Rajani, Omar Sanseviero, Lewis Tunstall, and Thomas Wolf. 2023. Open llm leaderboard. `https://huggingface.co/spaces/HuggingFaceH4/open_llm_leaderboard`.

Anja Belz and Ehud Reiter. 2006. Comparing automatic and human evaluation of nlg systems. In *Conference of the European Chapter of the Association for Computational Linguistics*.

Elad Ben-Zaken, Shauli Ravfogel, and Yoav Goldberg. 2021. Bitfit: Simple parameter-efficient fine-tuning for transformer-based masked language-models. *ArXiv*, abs/2106.10199.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Daniel Matthew Cer, Mona T. Diab, Eneko Agirre, Iñigo Lopez-Gazpio, and Lucia Specia. 2017. Semeval-2017 task 1: Semantic textual similarity multilingual and crosslingual focused evaluation. In *International Workshop on Semantic Evaluation*.

Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Wei Zhu, Yuan Ni, Guotong Xie, Zhiyuan Liu, and Maosong Sun. 2023. Ultrafeedback: Boosting language models with high-quality feedback. *ArXiv*, abs/2310.01377.

Dorottya Demszky, Kelvin Guu, and Percy Liang. 2018. Transforming question answering datasets into natural language inference datasets. *ArXiv*, abs/1809.02922.

Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. Qlora: Efficient finetuning of quantized llms. *ArXiv*, abs/2305.14314.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.

Ning Ding, Xingtai Lv, Qiaosen Wang, Yulin Chen, Bowen Zhou, Zhiyuan Liu, and Maosong Sun. 2023. Sparse low-rank adaptation of pre-trained language models. *arXiv preprint arXiv:2311.11696*.

Ning Ding, Yujia Qin, Guang Yang, Fuchao Wei, Zonghan Yang, Yusheng Su, Shengding Hu, Yulin Chen, Chi-Min Chan, Weize Chen, et al. 2022. Delta tuning: A comprehensive study of parameter efficient methods for pre-trained language models. *arXiv preprint arXiv:2203.06904*.

William B. Dolan and Chris Brockett. 2005. Automatically constructing a corpus of sentential paraphrases. In *International Joint Conference on Natural Language Processing*.

Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac'h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. 2023. A framework for few-shot language model evaluation.

Demi Guo, Alexander M Rush, and Yoon Kim. 2020. Parameter-efficient transfer learning with diff pruning. *arXiv preprint arXiv:2012.07463*.

Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. 2021. Towards a unified view of parameter-efficient transfer learning. *ArXiv*, abs/2110.04366.

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Xiaodong Song, and Jacob Steinhardt. 2020. Measuring massive multitask language understanding. *ArXiv*, abs/2009.03300.

Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Xiaodong Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the math dataset. *ArXiv*, abs/2103.03874.

Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. Parameter-efficient transfer learning for nlp. In *International Conference on Machine Learning*, pages 2790–2799. PMLR.

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.

Menglin Jia, Luming Tang, Bor-Chun Chen, Claire Cardie, Serge Belongie, Bharath Hariharan, and Ser-Nam Lim. 2022. Visual prompt tuning. In *European Conference on Computer Vision*, pages 709–727. Springer.

Rabeeh Karimi Mahabadi, James Henderson, and Sebastian Ruder. 2021. Compacter: Efficient low-rank hypercomplex adapter layers. *Advances in Neural Information Processing Systems*, 34:1022–1035.

Dawid Jan Kopiczko, Tijmen Blankevoort, and Yuki Markus Asano. 2023. Vera: Vector-based random matrix adaptation. *ArXiv*, abs/2310.11454.

Jaejun Lee, Raphael Tang, and Jimmy J. Lin. 2019. What would elsa do? freezing layers during transformer fine-tuning. *ArXiv*, abs/1911.03090.

Brian Lester, Rami Al-Rfou, and Noah Constant. 2021a. The power of scale for parameter-efficient prompt tuning. *arXiv preprint arXiv:2104.08691*.

Brian Lester, Rami Al-Rfou, and Noah Constant. 2021b. The power of scale for parameter-efficient prompt tuning. In *Conference on Empirical Methods in Natural Language Processing*.

Quentin Lhoest, Albert Villanova del Moral, Yacine Jernite, Abhishek Thakur, Patrick von Platen, Suraj Patil, Julien Chaumond, Mariama Drame, Julien Plu, Lewis Tunstall, Joe Davison, Mario vSavsko, Gunjan Chhablani, Bhavitvya Malik, Simon Brandeis, Teven Le Scao, Victor Sanh, Canwen Xu, Nicolas Patry, Angelina McMillan-Major, Philipp Schmid, Sylvain Gugger, Clement Delangue, Th'eo Matussiere, Lysandre Debut, Stas Bekman, Pierric Cistac, Thibault Goehringer, Victor Mustar, François Lagunas, Alexander M. Rush, and Thomas Wolf. 2021. Datasets: A community library for natural language processing. *ArXiv*, abs/2109.02846.

Xiang Lisa Li and Percy Liang. 2021a. Prefix-tuning: Optimizing continuous prompts for generation. *arXiv preprint arXiv:2101.00190*.

Xiang Lisa Li and Percy Liang. 2021b. Prefix-tuning: Optimizing continuous prompts for generation. *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)*, abs/2101.00190.

Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy Liang, and Tatsunori B Hashimoto. 2023. Alpacaeval: An automatic evaluator of instruction-following models.

Dongze Lian, Daquan Zhou, Jiashi Feng, and Xinchao Wang. 2022. Scaling & shifting your features: A new baseline for efficient model tuning. *ArXiv*, abs/2210.08823.

Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Annual Meeting of the Association for Computational Linguistics*.

Stephanie C. Lin, Jacob Hilton, and Owain Evans. 2021. Truthfulqa: Measuring how models mimic human falsehoods. In *Annual Meeting of the Association for Computational Linguistics*.

Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Mohta, Tenghao Huang, Mohit Bansal, and Colin Raffel. 2022. Few-shot parameter-efficient fine-tuning is better and cheaper than in-context learning. *ArXiv*, abs/2205.05638.

Wenhao Liu, Xiaohua Wang, Muling Wu, Tianlong Li, Changze Lv, Zixuan Ling, Jianhao Zhu, Cenyuan Zhang, Xiaoqing Zheng, and Xuanjing Huang. 2023. Aligning large language models with human preferences through representation engineering.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *ArXiv*, abs/1907.11692.

Yuning Mao, Lambert Mathias, Rui Hou, Amjad Almahairi, Hao Ma, Jiawei Han, Wen tau Yih, and Madian Khabsa. 2021. Unipelt: A unified framework for parameter-efficient language model tuning. In *Annual Meeting of the Association for Computational Linguistics*.

Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. 2018. Can a suit of armor conduct electricity? a new dataset for open book question answering. In *Conference on Empirical Methods in Natural Language Processing*.

Jekaterina Novikova, Ondrej Dusek, and Verena Rieser. 2017. The e2e dataset: New challenges for end-to-end generation. *ArXiv*, abs/1706.09254.

OpenAI. 2023. GPT-4 technical report. *CoRR*, abs/2303.08774.

Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. In *Proceedings of the Conference on Neural Information Processing Systems (NeurIPS 2022)*.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Annual Meeting of the Association for Computational Linguistics*.

Jonas Pfeiffer, Aishwarya Kamath, Andreas Rücklé, Kyunghyun Cho, and Iryna Gurevych. 2020. Adapterfusion: Non-destructive task composition for transfer learning. *ArXiv*, abs/2005.00247.

Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. 2018. Improving language understanding by generative pre-training.

Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551.

Andreas Rücklé, Gregor Geigle, Max Glockner, Tilman Beck, Jonas Pfeiffer, Nils Reimers, and Iryna Gurevych. 2020. Adapterdrop: On the efficiency of adapters in transformers. In *Conference on Empirical Methods in Natural Language Processing*.

Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2019. An adversarial winograd schema challenge at scale.

Taylor Shin, Yasaman Razeghi, Robert L Logan IV, Eric Wallace, and Sameer Singh. 2020. Eliciting knowledge from language models using automatically generated prompts. *ArXiv*, abs/2010.15980.

Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, A. Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In *Conference on Empirical Methods in Natural Language Processing*.

Asa Cooper Stickland and Iain Murray. 2019. Bert and pals: Projected attention layers for efficient adaptation in multi-task learning. In *International Conference on Machine Learning*, pages 5986–5995. PMLR.

Nishant Subramani, Nivedita Suresh, and Matthew E. Peters. 2022. Extracting latent steering vectors from pretrained language models. *ArXiv*, abs/2205.05124.

Hugo Touvron, Louis Martin, Kevin R. Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Daniel M. Bikel, Lukas Blecher, Cristian Cantón Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony S. Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel M. Kloumann, A. V. Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, R. Subramanian, Xia Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zhengxu Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models. *ArXiv*, abs/2307.09288.

Alexander Matt Turner, Lisa Thiergart, David S. Udell, Gavin Leech, Ulisse Mini, and Monte Stuart MacDiarmid. 2023. Activation addition: Steering language models without optimization. *ArXiv*, abs/2308.10248.

Ashish Vaswani, Noam M. Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Neural Information Processing Systems*.

Ramakrishna Vedantam, C. Lawrence Zitnick, and Devi Parikh. 2014. Cider: Consensus-based image description evaluation. *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4566–4575.

Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2018. Glue: A multi-task benchmark and analysis platform for natural language understanding. In *BlackboxNLP@EMNLP*.

Zhen Wang, Rameswar Panda, Leonid Karlinsky, Rogério Schmidt Feris, Huan Sun, and Yoon Kim.

2023. Multitask prompt tuning enables parameter-efficient transfer learning. *ArXiv*, abs/2303.02861.

Alex Warstadt, Amanpreet Singh, and Samuel R. Bowman. 2018. Neural network acceptability judgments. *Transactions of the Association for Computational Linguistics*, 7:625–641.

Adina Williams, Nikita Nangia, and Samuel R. Bowman. 2017. A broad-coverage challenge corpus for sentence understanding through inference. In *North American Chapter of the Association for Computational Linguistics*.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, and Jamie Brew. 2019. Huggingface's transformers: State-of-the-art natural language processing. *ArXiv*, abs/1910.03771.

Muling Wu, Wenhao Liu, Jianhan Xu, Changze Lv, Zixuan Ling, Tianlong Li, Longtao Huang, Xiaoqing Zheng, and Xuanjing Huang. 2023. Parameter efficient multi-task fine-tuning by learning to transfer token-wise prompts. In *Conference on Empirical Methods in Natural Language Processing*.

Elad Ben Zaken, Shauli Ravfogel, and Yoav Goldberg. 2021. Bitfit: Simple parameter-efficient fine-tuning for transformer-based masked language-models. *arXiv preprint arXiv:2106.10199*.

Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. Hellaswag: Can a machine really finish your sentence? In *Annual Meeting of the Association for Computational Linguistics*.

Qingru Zhang, Minshuo Chen, Alexander Bukharin, Pengcheng He, Yu Cheng, Weizhu Chen, and Tuo Zhao. 2023a. Adaptive budget allocation for parameter-efficient fine-tuning. *arXiv preprint arXiv:2303.10512*.

Qingru Zhang, Minshuo Chen, Alexander W. Bukharin, Pengcheng He, Yu Cheng, Weizhu Chen, and Tuo Zhao. 2023b. Adaptive budget allocation for parameter-efficient fine-tuning. *ArXiv*, abs/2303.10512.

Mengjie Zhao, Tao Lin, Martin Jaggi, and Hinrich Schütze. 2020. Masking as an efficient alternative to finetuning for pretrained language models. In *Conference on Empirical Methods in Natural Language Processing*.

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Haotong Zhang, Joseph Gonzalez, and Ion Stoica. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *ArXiv*, abs/2306.05685.

Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan, Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, Shashwat Goel, Nathaniel Li, Michael J. Byun, Zifan Wang, Alex Mallen, Steven Basart, Sanmi Koyejo, Dawn Song, Matt Fredrikson, Zico Kolter, and Dan Hendrycks. 2023. Representation engineering: A top-down approach to ai transparency. *ArXiv*, abs/2310.01405.

## A  Datasets

### A.1  GLUE Benchmark

The General Language Understanding Evaluation (GLUE) benchmark comprises a variety of datasets, including CoLA (Warstadt et al., 2018), SST-2 (Socher et al., 2013), MRPC (Dolan and Brockett, 2005), QQP (Wang et al., 2018), STS-B (Cer et al., 2017), MNLI (Williams et al., 2017), QNLI (Demszky et al., 2018), and RTE (Bar-Haim et al., 2006). These datasets have been widely-used to measure the performance of language understanding models. The GLUE benchmark was sourced from Huggingface Datasets (Lhoest et al., 2021), and all datasets within this benchmark were employed for our evaluation.

Following (Ding et al., 2023) and (Hu et al., 2021), we evaluated models trained with different PEFT methods on the validation dataset. To ensure a fair and systematic comparison of different PEFT methods, we randomly sampled $1,000$ instances from the validation set to form a new validation set if the validation set contains more than $2,000$ instances (otherwise, we randomly selected half of the instances for the new validation set), and the remaining instances were used as the test set. This allowed us to adhere rigorously to the established training protocol, which involves evaluating on the new validation set after each epoch of training and selecting the checkpoint that yields the best performance on the validation set as the final model for testing on the test set. The sizes of the training, validation, and test sets as well as their corresponding metrics for all datasets are given in Table 9.

For all the experiments with RoBERTa, we ran experiment 5 times for each PEFT method using distinct random seeds to initialize the trainable weights and reported the average results. During the evaluation of the RTE and CoLA datasets, some researchers reported that specific random seeds could potentially lead to anomalous experimental outcomes[2]. Therefore, we randomly chose other five different random seeds (i.e., 42, 43, 44, 45, and 46) to conduct the experiments.

### A.2  E2E NLG Challenge

The E2E NLG Challenge was first introduced by Novikova et al. (2017) to train and evaluate end-to-end, data-driven natural language generation models and systems. All datasets in the E2E NLG Challenge were sourced from Huggingface Datasets.

| Dataset | #Train | #Validation | #Test | Metric |
|---|---|---|---|---|
| CoLA | 8.5K | 522 | 521 | MCC |
| SST-2 | 67K | 436 | 436 | ACC |
| MRPC | 3.7K | 204 | 204 | ACC |
| QQP | 364K | 1K | 39K | ACC |
| STS-B | 5.7K | 750 | 750 | CORR |
| MNLI | 393K | 1K | 8K | ACC |
| QNLI | 105K | 1K | 4.5K | ACC |
| RTE | 2.5K | 139 | 138 | ACC |

Table 9: The sizes of the training, validation, and test sets, along with their corresponding metrics for all datasets in the GLUE benchmark. "MCC", "ACC" and "CORR" denote Matthews correlation coefficient, accuracy, and Pearson correlation coefficient respectively.

This benchmark includes 42.1K training instances, 4.67K validation instances, and 4.69K testing instances. Following previous studies, we used the official evaluation script to compute BLEU (Papineni et al., 2002), NIST (Belz and Reiter, 2006), METEOR (Banerjee and Lavie, 2005), ROUGE-L (Lin, 2004) and CIDEr (Vedantam et al., 2014) scores. For evaluations using GPT-2, we ran experiment for 3 times using distinct random seeds and reported the average results. The random seeds used for these experiments were 42, 43, and 44.

### A.3  UltraFeedback

UltraFeedback (Cui et al., 2023) comprises $64,000$ prompts, each of them accompanied by four LLM responses evaluated by GPT-4 based on criteria such as instruction-following, honesty, and helpfulness. To construct the training dataset from UltraFeedback, we selected the responses with the highest mean scores[3].

### A.4  Open LLM Leaderboard

The Open LLM Leaderboard includes six datasets that assess various capabilities such as science questions, commonsense inference, multitask accuracy, mathematical reasoning, and the truthfulness of generated answers. It consists of ARC (Mihaylov et al., 2018), HellaSwag (Zellers et al., 2019), WinoGrande (Sakaguchi et al., 2019), MMLU (Hendrycks et al., 2020), TruthfulQA (Lin et al., 2021), and GSM8K (Hendrycks et al., 2021). We used the Eleuther AI Language Model Evaluation Harness library (Gao et al., 2023) to evaluate language models trained using various methods. Table 11 details the leaderboard evaluation configuration and the experimental settings applied in this study.

---

[2]Refer to `github.com/microsoft/LoRA/issues`.

[3]The constructed dataset can be accessed at: `huggingface.co/datasets/argilla/` `\ultrafeedback-binarized-preferences-cleaned`.

## A.5 AlpacaEval

AlpacaEval is an automated evaluation benchmark for LLMs, employing GPT-4 (OpenAI, 2023) as an annotator to compare the responses generated by the models against reference answers produced by text-davinci-003. The benchmark consists of 805 samples focused on simple instruction-following tasks. Previous research has shown a high correlation between GPT-4's annotations and human evaluator assessments (Li et al., 2023).

## A.6 MT-Bench

MT-Bench (Zheng et al., 2023) comprises 80 challenging questions, each structured as a two-turn dialogue. GPT-4 was also used to evaluate the responses generated by the models, assigning a score from 1 to 10 for each turn.

## B Hyperparameters

### B.1 Experiments with RoBERTa

We trained the models using AdamW with a linear learning rate decay schedule. To ensure a fair comparison, we maintained the same sequence length across all baseline PEFT methods. The hyperparameters used in our experiments for RoBERTa base are detailed in Table 13, and for RoBERTa large in Table 15.

### B.2 Experiments with GPT-2

We trained the models using AdamW with a linear learning rate decay schedule. To ensure a fair comparison, we maintained the same sequence length across all baseline PEFT methods. We used the Huggingface PEFT package to replicate Prefix Tuning and LoRA, and the OpenDelta package to replicate Adapter and Adapter-FFN. The hyperparameters used in our experiments for GPT-2 medium and GPT-2 large are detailed in Table 16s.

### B.3 Experiments with T5

We trained the models using AdamW with a linear learning rate decay. We conducted experiments on T5 following the settings established by Asai et al. (2022). The hyperparameters used in our experiments for T5 base are detailed in Table 14.

### B.4 Experiments with LLaMA-2

We employed AdamW with a cosine learning rate decay schedule during the fine-tuning. To ensure a fair comparison, we maintained the same sequence length across all baseline methods. The hyperparameters used for LLaMA-2 are listed in Table 10. We evaluated the performance of models on TruthfulQA after each training epoch, and the checkpoint demonstrating the best performance was selected for final testing on other datasets. A greedy decoding strategy was employed for generation for all the models. To mitigate repetition issues, we applied a repetition penalty of $1.1$ and set the no-repeat $n$-gram size to 5.

## C Results of Standard Deviations

We presented the average results in the main text. To provide a quantifiable measure of uncertainty in the measurement or prediction and help other researchers replicate our experiments, we also reported the standard deviations with RoBERTa models on the GLUE benchmark in Table 17, and the standard deviations with GPT-2 on the E2E NLG Challenge in Table 18.

## D Additional Results with LLaMA-2

Example responses generated by models trained with RAD, full-parameter fine-tuning (FT) and LoRA are presented in Figures 3 and 4. Tables 19 and 12 provide the results on the MT-Bench and Open LLM Leaderboard respectively.

| Method | Hyperparameter | Value |
|--------|----------------|-------|
| ALL | Batch Size | 128 |
| | Micro Batch Size | 2 |
| | Optimizer | Adamw |
| | LR Scheduler Type | Cosine |
| | Warmup Ratio | 0.1 |
| | Maximum Sequence Length | 768 |
| FT | Learning Rate | $2e-5$ |
| | Number of Epochs | 3 |
| LoRA | Learning Rate | $3e-4$ |
| | Number of Epochs | 3 |
| | Batch Size | 128 |
| | Coefficient $\alpha$ | 16 |
| | Dropout | 0.1 |
| | Rank | 16 |
| RED | Learning Rate | $1e-3$ |
| | Number of Epochs | 10 |

Table 10: Hyperparameters used for training LLaMA-2.

| Dataset | # Few-shot | Metric |
|---------|-----------|--------|
| ARC | 25 | ACC-NORM |
| TruthfulQA | 0 | MC2 |
| Winogrande | 5 | ACC |
| GSM8K | 5 | ACC |
| HellaSwag | 10 | ACC-NORM |
| MMLU | 10 | ACC |

Table 11: Number of few-shot examples and metrics used for evaluation on the datasets in the Open LLM Leaderboard. "MC2" denotes the normalized total probability assigned to the set of true answers, and "ACC-NORM" denotes the normalized accuracy used in settings where response lengths can vary significantly.

| Method | # Param | ARC | TruthfulQA | Winogrande | GSM8k | HellaSwag | MMLU | Average |
|---|---|---|---|---|---|---|---|---|
| FT | 6739M | 73.34 | 47.07 | 74.03 | 16.22 | 79.95 | 46.55 | 56.19 |
| LoRA | 8.39M | 72.07 | 44.84 | 75.85 | 15.01 | 78.60 | 44.42 | 55.13 |
| RED | 0.26M | 72.04 | 47.59 | 72.85 | 11.90 | 77.86 | 42.27 | 54.09 |

Table 12: Results achieved by RED, LoRA and full-parameter fine-tuning (FT) on the Open LLM Leaderboard.

| Method | Dataset | MNLI | SST-2 | MRPC | CoLA | QNLI | QQP | RTE | STS-B |
|---|---|---|---|---|---|---|---|---|---|
| ALL | Optimizer | | | | AdamW | | | | |
| | LR Schedule | | | | Linear | | | | |
| FT | Batch Size | 32 | 32 | 32 | 32 | 32 | 32 | 32 | 32 |
| | Number of Epochs | 20 | 40 | 40 | 40 | 20 | 20 | 40 | 40 |
| | Learning Rate | $1e-05$ | $1e-05$ | $1e-05$ | $1e-05$ | $1e-05$ | $1e-05$ | $1e-05$ | $1e-05$ |
| | Warmup Ratio | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ |
| | Weight Decay | | | | $1e-04$ | | | | |
| | Maximum Sequence Length | | | | 256 | | | | |
| LoRA | Batch Size | 32 | 32 | 32 | 32 | 32 | 32 | 32 | 32 |
| | Number of Epochs | 20 | 40 | 40 | 40 | 20 | 20 | 40 | 40 |
| | Learning Rate | $5e-04$ | $5e-04$ | $4e-04$ | $4e-04$ | $4e-04$ | $5e-04$ | $5e-04$ | $4e-04$ |
| | Warmup Ratio | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ |
| | Rank | | | | 8 | | | | |
| | Coefficient $\alpha$ | | | | 8 | | | | |
| | Maximum Sequence Length | | | | 256 | | | | |
| Adapter | Batch Size | 32 | 32 | 32 | 32 | 32 | 32 | 32 | 32 |
| | Number of Epochs | 20 | 40 | 40 | 40 | 20 | 20 | 40 | 40 |
| | Learning Rate | $1e-04$ | $1e-04$ | $1e-04$ | $1e-04$ | $1e-04$ | $1e-04$ | $1e-04$ | $1e-04$ |
| | Warmup Ratio | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ |
| | Rank | | | | 8 | | | | |
| | Maximum Sequence Length | | | | 256 | | | | |
| Adapter-FFN | Batch Size | 32 | 32 | 32 | 32 | 32 | 32 | 32 | 32 |
| | Number fo Epochs | 20 | 40 | 40 | 40 | 20 | 20 | 40 | 40 |
| | Learning Rate | $1e-04$ | $1e-04$ | $1e-04$ | $1e-04$ | $1e-04$ | $1e-04$ | $1e-04$ | $1e-04$ |
| | Warmup Ratio | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ |
| | Rank | | | | 16 | | | | |
| | Maximum Sequence Length | | | | 256 | | | | |
| BitFit | Batch Size | 32 | 32 | 32 | 32 | 32 | 32 | 32 | 32 |
| | Number of Epochs | 20 | 40 | 40 | 40 | 20 | 20 | 40 | 40 |
| | Learning Rate | $1e-04$ | $1e-04$ | $1e-04$ | $1e-04$ | $1e-04$ | $1e-04$ | $1e-04$ | $1e-04$ |
| | Warmup Ratio | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ |
| | Maximum Sequence Length | | | | 256 | | | | |
| RED | Batch Size | 32 | 32 | 32 | 32 | 32 | 32 | 32 | 32 |
| | Number of Epochs | 20 | 40 | 40 | 40 | 20 | 20 | 40 | 40 |
| | Learning rate | $1e-03$ | $1e-03$ | $5e-03$ | $5e-03$ | $3e-03$ | $1e-03$ | $8e-03$ | $3e-03$ |
| | Warmup Ratio | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ |
| | Maximum Sequence Length | | | | 256 | | | | |

Table 13: Hyperparameter values for training RoBERTa base on the GLUE benchmark.

| Method | Dataset | MNLI | SST-2 | MRPC | CoLA | QNLI | QQP | RTE | STS-B |
|---|---|---|---|---|---|---|---|---|---|
| ALL | Optimizer | | | | AdamW | | | | |
| | LR Schedule | | | | Linear | | | | |
| RED | Batch Size | 32 | 32 | 32 | 32 | 32 | 32 | 32 | 32 |
| | Number of Epochs | 10 | 10 | 20 | 60 | 10 | 10 | 30 | 50 |
| | Learning rate | $5e-02$ | $3e-02$ | $1e-01$ | $4e-02$ | $2e-02$ | $5e-03$ | $7e-02$ | $4e-02$ |
| | Warmup Ratio | $1e-02$ | $6e-02$ | $1e-02$ | $1e-02$ | $0e-00$ | $1e-02$ | $1e-02$ | $6e-02$ |
| | Maximum Sequence Lengt | | | | 256 | | | | |

Table 14: Hyperparameter values for training T5 base on the GLUE benchmark.

| Method | Dataset | MNLI | SST-2 | MRPC | CoLA | QNLI | QQP | RTE | STS-B |
|---|---|---|---|---|---|---|---|---|---|
| **ALL** | Optimizer | | | | | AdamW | | | |
| | LR Schedule | | | | | Linear | | | |
| **FT** | Batch Size | 16 | 32 | 32 | 32 | 16 | 16 | 16 | 32 |
| | Number of Epochs | 10 | 10 | 20 | 20 | 10 | 10 | 20 | 10 |
| | Learning rate | $2e-05$ | $2e-05$ | $2e-05$ | $2e-05$ | $2e-05$ | $1e-05$ | $1e-05$ | $2e-05$ |
| | Warmup Ratio | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ |
| | Weight Decay | | | | | $1e-01$ | | | |
| | Maximum Sequence Length | | | | | 256 | | | |
| **LoRA** | Batch Size | 32 | 32 | 32 | 32 | 32 | 32 | 32 | 32 |
| | Number of Epochs | 10 | 10 | 20 | 20 | 10 | 10 | 20 | 10 |
| | Learning rate | $3e-04$ | $4e-04$ | $3e-04$ | $2e-04$ | $2e-04$ | $3e-04$ | $4e-04$ | $2e-04$ |
| | Warmup Ratio | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ |
| | Rank | | | | | 8 | | | |
| | Coefficient $\alpha$ | | | | | 16 | | | |
| | Maximum Sequence Length | | | | | 256 | | | |
| **Adapter** | Batch Size | 32 | 32 | 32 | 32 | 32 | 32 | 32 | 32 |
| | Number of Epochs | 10 | 10 | 20 | 20 | 10 | 10 | 20 | 10 |
| | Learning rate | $3e-04$ | $3e-04$ | $3e-04$ | $3e-04$ | $3e-04$ | $3e-04$ | $3e-04$ | $3e-04$ |
| | Warmup Ratio | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ |
| | Rank | | | | | 8 | | | |
| | Maximum Sequence Length | | | | | 256 | | | |
| **Adapter-FFN** | Batch Size | 32 | 32 | 32 | 32 | 32 | 32 | 32 | 32 |
| | Number of Epochs | 10 | 10 | 20 | 20 | 10 | 10 | 20 | 10 |
| | Learning rate | $3e-04$ | $3e-04$ | $3e-04$ | $3e-04$ | $3e-04$ | $3e-04$ | $3e-04$ | $3e-04$ |
| | Warmup Ratio | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ | $6e-02$ |
| | Rank | | | | | 16 | | | |
| | Maximum Sequence Length | | | | | 256 | | | |
| **RED** | Batch Size | 32 | 32 | 32 | 32 | 32 | 32 | 32 | 32 |
| | Number of Epochs | 10 | 10 | 20 | 20 | 10 | 10 | 20 | 10 |
| | Learning rate | $1e-03$ | $1e-03$ | $2e-03$ | $1e-03$ | $1e-03$ | $1e-03$ | $5e-03$ | $5e-03$ |
| | Weight Decay | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | $1e-04$ | 0.0 |
| | Warmup Ratio | $6e-02$ | $6e-02$ | 0 | $6e-02$ | $6e-02$ | $6e-02$ | $1e-02$ | $6e-02$ |
| | Maximum Sequence Length | | | | | 256 | | | |

Table 15: Hyperparameter values for training RoBERTa large on the GLUE benchmark.

| Dataset | FT | FT$^{\text{top2}}$ | Adapter | Apapter-FFN | LoRA | Prefix Tuning | RED (medium) | RED (large) |
|---|---|---|---|---|---|---|---|---|
| | | | | **Training** | | | | |
| Optimizer | | | | | AdamW | | | |
| Weight Decay | 0.0 | 0.0 | 0.0 | 0.0 | $1e-02$ | 0.0 | $1e-04$ | 0.0 |
| Number of Epochs | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 10 |
| Learning Rate Schedule | | | | | Linear | | | |
| Label Smooth | 0.0 | 0.0 | 0.0 | 0.0 | 0.1 | 0.0 | 0.0 | 0.0 |
| Learning Rate | $5e-05$ | $5e-05$ | $8e-05$ | $8e-05$ | $2e-04$ | $8e-05$ | $6e-02$ | $6e-03$ |
| Rank or Prefix Length | $--$ | $--$ | 8 | 16 | 8 | 16 | $--$ | $--$ |
| Coefficient $\alpha$ | $--$ | $--$ | $--$ | $--$ | 32 | $--$ | $--$ | $--$ |
| Adaption | $--$ | $--$ | $--$ | $--$ | 8 | $--$ | $--$ | $--$ |
| Warmup Steps | | | | | 500 | | | |
| Batch Size | | | | | 10 | | | |
| | | | | **Inference** | | | | |
| Beam Size | | | | | 10 | | | |
| Length Penalty | | | | | 0.9 | | | |
| No-repeat $n$-gram Size | | | | | 4 | | | |

Table 16: Hyperparameter values for training GPT-2 on the E2E NLG Challenge. "RED (medium)" denotes the values of hyperparameters used by RED to fine-tune GPT-2 medium and "RED (large)" the values to GPT-2 large.

| Methods | # Param | MNLI | SST-2 | MRPC | CoLA | QNLI | QQP | RTE | STS-B | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|
| **FT (base)** | 125M | $87.3_{\pm0.34}$ | $94.4_{\pm0.96}$ | $87.9_{\pm0.91}$ | $62.4_{\pm3.29}$ | $92.5_{\pm0.22}$ | $91.7_{\pm0.19}$ | $78.3_{\pm3.20}$ | $90.6_{\pm0.59}$ | 85.6 |
| **Adapter (base)** | 0.4M | $87.0_{\pm0.28}$ | $93.3_{\pm0.40}$ | $88.4_{\pm1.54}$ | $60.9_{\pm3.09}$ | $92.5_{\pm0.02}$ | $90.5_{\pm0.08}$ | $76.6_{\pm2.26}$ | $90.5_{\pm0.35}$ | 85.0 |
| **Adapter-FFN (base)** | 0.3M | $87.1_{\pm0.10}$ | $93.0_{\pm0.50}$ | $88.8_{\pm1.38}$ | $58.5_{\pm1.69}$ | $92.1_{\pm0.28}$ | $90.2_{\pm0.07}$ | $77.7_{\pm1.93}$ | $90.4_{\pm0.31}$ | 84.7 |
| **LoRA (base)** | 0.3M | $86.6_{\pm0.26}$ | $93.9_{\pm0.49}$ | $88.7_{\pm0.76}$ | $59.7_{\pm4.36}$ | $92.6_{\pm0.10}$ | $90.4_{\pm0.08}$ | $75.3_{\pm2.79}$ | $90.3_{\pm0.54}$ | 84.7 |
| **BitFit (base)** | 0.1M | $84.7_{\pm0.08}$ | $94.0_{\pm0.87}$ | $88.1_{\pm1.57}$ | $54.0_{\pm3.07}$ | $91.0_{\pm0.05}$ | $87.3_{\pm0.02}$ | $69.8_{\pm1.51}$ | $89.5_{\pm0.35}$ | 82.3 |
| **RED (base)** | 0.02M | $83.9_{\pm0.14}$ | $93.9_{\pm0.31}$ | $89.2_{\pm0.98}$ | $61.0_{\pm2.96}$ | $90.7_{\pm0.35}$ | $87.2_{\pm0.17}$ | $78.0_{\pm2.06}$ | $90.4_{\pm0.32}$ | 84.7 |
| **FT (large)** | 355M | $88.8_{\pm0.45}$ | $96.0_{\pm0.66}$ | $91.7_{\pm1.73}$ | $68.2_{\pm2.62}$ | $93.8_{\pm0.33}$ | $91.5_{\pm1.28}$ | $85.8_{\pm1.40}$ | $92.6_{\pm0.16}$ | 88.5 |
| **LoRA (large)** | 0.8M | $90.2_{\pm0.25}$ | $96.0_{\pm0.85}$ | $89.8_{\pm2.09}$ | $65.5_{\pm2.02}$ | $94.7_{\pm0.21}$ | $90.7_{\pm0.91}$ | $86.3_{\pm2.41}$ | $91.7_{\pm0.44}$ | 88.1 |
| **Adapter (large)** | 0.9M | $90.1_{\pm0.12}$ | $95.2_{\pm0.48}$ | $90.5_{\pm0.59}$ | $65.4_{\pm2.24}$ | $94.6_{\pm0.17}$ | $91.4_{\pm0.13}$ | $85.3_{\pm1.34}$ | $91.5_{\pm0.33}$ | 88.0 |
| **Adapter-FFN (large)** | 0.8M | $90.3_{\pm0.15}$ | $96.1_{\pm0.75}$ | $90.5_{\pm1.26}$ | $64.4_{\pm1.56}$ | $94.3_{\pm0.39}$ | $91.3_{\pm0.24}$ | $84.8_{\pm2.01}$ | $90.2_{\pm0.24}$ | 87.7 |
| **RED (large)** | 0.05M | $89.5_{\pm0.38}$ | $96.0_{\pm0.48}$ | $90.3_{\pm1.40}$ | $68.1_{\pm1.69}$ | $93.5_{\pm0.33}$ | $88.8_{\pm0.11}$ | $86.2_{\pm1.40}$ | $91.3_{\pm0.21}$ | 87.9 |

Table 17: Performance comparison of RoBERTa base and large models fine-tuned by RED and other PEFT baselines on the GLUE benchmark.

| Method | # Param | BLEU | NIST | MET | ROUGE-L | CIDEr |
|---|---|---|---|---|---|---|
| **FT (medium)** | 355M | $65.95_{\pm0.26}$ | $8.52_{\pm0.03}$ | $45.95_{\pm0.07}$ | $69.13_{\pm0.30}$ | $2.35_{\pm0.01}$ |
| **FT$^{top2}$ (medium)** | 25.2M | $65.94_{\pm0.33}$ | $8.53_{\pm0.03}$ | $44.28_{\pm0.09}$ | $68.83_{\pm0.17}$ | $2.23_{\pm0.02}$ |
| **Adapter (medium)** | 0.9M | $64.31_{\pm0.17}$ | $8.29_{\pm0.01}$ | $44.91_{\pm0.29}$ | $67.72_{\pm0.26}$ | $2.28_{\pm0.01}$ |
| **Adapter-FFN (medium)** | 0.8M | $64.41_{\pm0.17}$ | $8.30_{\pm0.02}$ | $44.74_{\pm0.11}$ | $67.53_{\pm0.02}$ | $2.29_{\pm0.01}$ |
| **LoRA (medium)** | 0.8M | $67.43_{\pm0.39}$ | $8.65_{\pm0.05}$ | $46.01_{\pm0.07}$ | $69.64_{\pm0.14}$ | $2.42_{\pm0.01}$ |
| **Prefix Tuning (medium)** | 0.8M | $63.92_{\pm0.27}$ | $8.26_{\pm0.11}$ | $41.81_{\pm0.62}$ | $66.86_{\pm0.22}$ | $2.03_{\pm0.05}$ |
| **RED (medium)** | 0.05M | $64.86_{\pm0.40}$ | $8.36_{\pm0.03}$ | $44.99_{\pm0.02}$ | $67.62_{\pm0.22}$ | $2.28_{\pm0.01}$ |
| **FT (large)** | 774M | $65.56_{\pm0.47}$ | $8.50_{\pm0.05}$ | $45.40_{\pm0.29}$ | $68.38_{\pm0.23}$ | $2.27_{\pm0.02}$ |
| **Adapter (large)** | 1.8M | $65.94_{\pm0.35}$ | $8.46_{\pm0.05}$ | $45.78_{\pm0.11}$ | $68.65_{\pm0.35}$ | $2.34_{\pm0.01}$ |
| **Adapter-FFN (large)** | 1.5M | $65.53_{\pm0.61}$ | $8.41_{\pm0.07}$ | $45.65_{\pm0.12}$ | $68.46_{\pm0.16}$ | $2.33_{\pm0.01}$ |
| **LoRA (large)** | 1.5M | $68.24_{\pm0.28}$ | $8.76_{\pm0.04}$ | $46.23_{\pm0.04}$ | $69.92_{\pm0.16}$ | $2.42_{\pm0.01}$ |
| **Prefix Tuning (large)** | 1.5M | $65.50_{\pm0.63}$ | $8.45_{\pm0.05}$ | $43.97_{\pm0.21}$ | $67.32_{\pm0.38}$ | $2.23_{\pm0.02}$ |
| **RED (large)** | 0.09M | $65.77_{\pm0.48}$ | $8.42_{\pm0.06}$ | $46.12_{\pm0.10}$ | $69.03_{\pm0.09}$ | $2.36_{\pm0.02}$ |

Table 18: Performance comparison of GPT-2 medium and large models fine-tuned by RED and other PEFT baselines on the E2E NLG Challenge.

| Method | # Param | Writing | Roleplay | Reasoning | Math | Coding | Extraction | Stem | Humanities | Average |
|---|---|---|---|---|---|---|---|---|---|---|
| **Turn-1** | | | | | | | | | | |
| **FT** | 6739M | 9.111 | 6.500 | 4.778 | 1.444 | 2.700 | 4.800 | 7.944 | 8.833 | 5.688 |
| **LoRA** | 8.39M | 8.600 | 7.900 | 4.000 | 1.600 | 2.350 | 4.750 | 8.700 | 8.950 | 5.856 |
| **RED** | 0.26M | 8.778 | 7.278 | 4.875 | 1.300 | 2.625 | 4.125 | 7.400 | 9.167 | 5.732 |
| **Turn-2** | | | | | | | | | | |
| **FT** | 6739M | 1.667 | 5.938 | 2.222 | 1.700 | 2.000 | 2.111 | 3.200 | 5.300 | 3.021 |
| **LoRA** | 8.39M | 1.900 | 5.800 | 2.100 | 1.600 | 2.550 | 1.222 | 3.100 | 5.500 | 2.994 |
| **RED** | 0.26M | 1.375 | 5.500 | 2.444 | 1.444 | 2.125 | 1.750 | 3.000 | 5.875 | 2.946 |
| **Average** | | | | | | | | | | |
| **FT** | 6739M | 5.389 | 6.200 | 3.500 | 1.579 | 2.412 | 3.526 | 5.447 | 6.974 | 4.354 |
| **LoRA** | 8.39M | 5.250 | 6.850 | 3.050 | 1.600 | 2.450 | 3.079 | 5.900 | 7.225 | 4.434 |
| **RED** | 0.26M | 5.294 | 6.389 | 3.588 | 1.368 | 2.375 | 2.938 | 5.750 | 7.618 | 4.401 |

Table 19: Performance comparison of LLaMA-2 fine-tuned by RED and other PEFT baselines on the MT-Bench.

| Prompt | Can you help me figure out how to sew a button on a shirt? |
|--------|-----------------------------------------------------------------------------------|
| FT | 1. Gather the materials needed for sewing a button, which include thread, needle, and button. <br> 2. Thread the needle with the desired color of thread. <br> 3. Place the button through the hole in the shirt where it is missing or needs replacing. <br> 4. Holding the button in place, thread the needle through the fabric from the backside of the shirt. <br> 5. Pull the thread tightly to secure the button in place. <br> 6. Cut off any excess thread. <br> 7. Sew the button onto the shirt by hand or using a sewing machine. <br> 8. Iron the shirt to smooth out any creases caused during the sewing process. <br> 9. Enjoy your newly repaired shirt! |
| LoRA | 1. Gather the necessary materials: You will need a shirt, thread, a needle, and a button. If your shirt has a buttonhole, you can use that; otherwise, you'll have to create one using a buttonhole template or a pattern. <br> 2. Mark the location of the buttonhole: Use a marker or chalk to mark the location where you want to place the buttonhole. This should be centered on the front of the shirt, just below the collar. <br> 3. Create the buttonhole: If you don't have a buttonhole already, you'll need to make one. To do this, fold the fabric over so that it forms a triangle shape. Then, use a ruler and a pencil to draw a line along the edge of the triangle. Cut along this line with scissors. <br> 4. Sew the buttonhole: Thread the needle with the thread and insert the needle through the fabric at the top of the buttonhole. Pull the thread tightly to create a small hole in the fabric. Continue sewing around the perimeter of the buttonhole, leaving about an inch of thread unsewn at the bottom. <br> 5. Insert the button: Place the button on the shirt, right-side up, with the hole facing down. Line up the buttonhole with the hole in the button and push the button through the buttonhole. <br> 6. Tie off the thread: Once the button is securely attached, tie off the thread by knotting it near the beginning of the buttonhole stitching. Trim any excess thread. <br> 7. Finish the buttonhole: Fold the fabric over the buttonhole and press it flat with an iron. This will help to secure the button and prevent it from coming loose. <br> 8. Repeat for additional buttons: If you need to add more buttons, repeat steps 4-7 for each new button. <br> 9. Test fit: Try on the shirt to ensure that the buttons are evenly spaced and that they look good. Adjust as needed. <br> 10. Enjoy your newly buttoned shirt! |
| RED | 1. Gather the necessary materials: You will need a shirt, thread, a needle, and a button. If your shirt has buttons that are already missing or damaged, you can purchase replacement buttons at most craft stores. <br> 2. Prepare the shirt: Lay the shirt flat with the back facing upwards. Make sure the area where you want to sew the button is clean and free of any debris or wrinkles. <br> 3. Thread the needle: Insert one end of the thread into the eye of the needle and pull it through the other side. Tie a knot in the end of the thread to secure it. <br> 4. Place the button on the shirt: Position the button on the desired location on the shirt, making sure it's centered and aligned properly. <br> 5. Sew the button: Holding the shirt steady with one hand, use the other hand to guide the needle through the fabric from the backside of the shirt. Inject the thread into the fabric by pushing down gently with your thumb while pulling the needle through. Continue sewing around the perimeter of the button, making sure to keep the stitches even and consistent. <br> 6. Trim the excess thread: Once you have finished sewing the button, cut off any excess thread that extends beyond the edge of the button. <br> 7. Test the button: Try fastening the button onto the shirt to make sure it's securely attached. If needed, adjust the position of the button or re-sew if it's not properly attached. <br> 8. Repeat for additional buttons: If there are multiple buttons on the shirt that need to be sewn, repeat steps 1-7 for each button. <br> 9. Clean up: After completing all the buttons, carefully remove any loose threads or lint from the shirt. <br><br> Remember to always use caution when handling sharp objects like needles and thread. Always wear protective gloves and follow proper safety precautions to avoid injury. |

Figure 3: The model fine-tuned with RED generates a thorough, sequential guide that offers accurate details, facilitating comprehension even for novices. This guide encompasses elements such as preparation, threading, positioning, sewing techniques, and post-sewing cleanup, while also providing safety advice and promoting testing for secure attachment. In contrast, the LoRA-trained response inaccurately concentrates on buttonhole creation rather than button sewing, potentially leading to confusion for individuals seeking button attachment guidance. The response generated by the full-parameter trained model presents a simplified summary, but it lacks the in-depth explanation and precision of RED, rendering it less informative for those unfamiliar with the sewing process.

| Prompt | Hi, my cat just chewed on a small weed outside that looked like hairy bittercress. Could you check if this is toxic for cats? Please also look for plants from the same family, as I'm not 100% sure I identified it correctly. |
|---|---|
| FT | 🐱 Hello! I'm here to help answer your question about whether or not the plant your cat chewed on is toxic for them. However, before I can provide an accurate response, I need more information about the plant in question. Could you please describe the plant in detail, including its appearance, size, and any unique characteristics it may have? Additionally, could you tell me what part of the plant your cat ingested (e.g., leaves, stems, flowers)? This will help me identify the correct plant and determine if it is safe or harmful for cats.<br><br>As a responsible AI assistant, I want to ensure that I provide accurate and helpful information while prioritizing the safety and well-being of both you and your furry friend. So, please take a moment to share any additional details you can about the plant your cat has chewed on, and I'll do my best to assist you. |
| LoRA | 🐱 Hello! I'm here to help answer your question about whether or not the plant your cat chewed on is toxic for them. However, before I can provide an accurate response, I need some more information. Can you please tell me what kind of plant it was and where it was found? This will help me identify the species and determine if it's safe for cats to consume. Additionally, I would appreciate it if you could provide any other details you have about the plant, such as its appearance, smell, and location. Thank you for your patience and understanding! |
| RED | 🐱 Hello! I'm here to help you with your question about whether or not your cat has eaten something potentially harmful. As an AI language model, I can provide information and resources based on the best available evidence. However, it's important to note that I am not a veterinarian and cannot diagnose or treat any medical conditions in animals. If you have concerns about your pet's health, please consult a licensed veterinarian.<br><br>Now, let's focus on your specific question. Based on what you described, it sounds like your cat may have eaten some type of plant material. It's always a good idea to keep an eye on your pet after they eat something new, especially if it's a plant. Some common signs of poisoning in cats include vomiting, diarrhea, drooling, lethargy, weakness, seizures, and difficulty breathing. If you notice any of these symptoms, it's essential to seek immediate veterinary care.<br><br>As for the plant itself, it's difficult to identify without more details. Bittercress (Ranunculus spp.) is a genus of flowering plants that includes several species, including buttercup, crowfoot, and spearwort. These plants are generally considered safe for cats, but there are other members of the Ranunculaceae family that could be toxic, such as aconite (Aconitum spp.), monkshood (Aconitella spp.), and wolfsbane (Aconitine).<br><br>To be on the safe side, I recommend taking your cat to the vet for a thorough examination and possible testing if necessary. The vet will be able to determine if the plant your cat ate was indeed toxic and provide appropriate treatment if needed. They may also advise you on how to prevent future exposure to potentially harmful plants.<br><br>In conclusion, while most plants are safe for cats to consume, it's always better to err on the side of caution and seek professional advice when in doubt. Please take your cat to the veterinarian for a proper evaluation and follow their instructions accordingly. |

Figure 4: The model fine-tuned using RED generates a comprehensive and proactive strategy, addressing immediate issues, potential symptoms to monitor, and the significance of veterinary consultation. It offers an overarching safety evaluation of the Ranunculaceae family, indicating potentially toxic members and highlighting the necessity for professional assessment. This response strikes a balance between informative content and practical guidance, empowering pet owners to act in their pet's best interests, even in the absence of specific plant identification. In contrast, the responses produced by the models trained with full parameters and LoRA place a greater emphasis on collecting further information before offering advice, which could inadvertently postpone critical care in an emergent situation.