# PEANuT: Parameter-Efficient Adaptation with Weight-aware Neural Tweakers

**Yibo Zhong[1,*]** **Haoxiang Jiang[2,*]** **Lincan Li[3]** **Ryumei Nakada[4]**

**Tianci Liu[5]** **Linjun Zhang[4]** **Huaxiu Yao[6]** **Haoyu Wang[2]**

[1]Independent Researcher  [2]University at Albany  [3]Florida State University  [4]Rutgers University

[5]Purdue University  [6]University of North Carolina at Chapel Hill

Fine-tuning large pre-trained foundation models often yields excellent downstream performance but is prohibitively expensive when updating all parameters. Parameter-efficient fine-tuning (PEFT) methods such as LoRA alleviate this by introducing lightweight update modules, yet they commonly rely on weight-agnostic linear approximations, limiting their expressiveness. In this work, we propose PEANuT, a novel PEFT framework that introduces weight-aware neural tweakers, compact neural modules that generate task-adaptive updates conditioned on frozen pre-trained weights. PEANuT provides a flexible yet efficient way to capture complex update patterns without full model tuning. We theoretically show that PEANuT achieves equivalent or greater expressivity than existing linear PEFT methods with comparable or fewer parameters. Extensive experiments across four benchmarks with over twenty datasets demonstrate that PEANuT consistently outperforms strong baselines in both NLP and vision tasks, while maintaining low computational overhead.

## 1. Introduction

Pre-trained models, trained on large and diverse general-domain corpora, have demonstrated strong generalization capabilities across a variety of tasks, including natural language understanding (Devlin, 2018, Liu, 2019, Howard and Ruder, 2018, Wu and He, 2019, Sun et al., 2023, Wang et al., 2021, 2022), generation (Touvron et al., 2023, AI@Meta, 2024, Xu et al., 2025, Liu et al., 2025, Wang et al., 2024b, Yao et al., 2022, Lewis et al., 2020), and vision tasks such as image classification (Dosovitskiy et al., 2020, Bhojanapalli et al., 2021, Chen et al., 2021). A common strategy for adapting these models to specific downstream tasks is full fine-tuning. However, due to the massive number of parameters involved, full fine-tuning often leads to significant computational and memory costs (Qin et al., 2024).

To mitigate these challenges, various parameter-efficient fine-tuning (PEFT) methods (Ding et al., 2023, Han et al., 2024) have been developed, enabling pre-trained models to be fine-tuned in resource-constrained environments (Lin et al., 2024). These methods retain most of the pre-trained weights in a frozen state and introduce a small set of trainable components, thereby significantly reducing memory and compute overhead (Lin et al., 2024). Among them, Low-Rank Adaptation (LoRA) (Hu

---

[*]These authors contributed equally to this work, order was determined randomly (by rolling a die).

et al., 2021b, Liu et al., 2024, Song et al., 2024, Büyükakyüz, 2024, Zhao et al., 2024) is a popular and widely adopted approach due to its simplicity, strong empirical performance, and compatibility with modern architectures.

Instead of updating pre-trained model weight directly, LoRA introduces two learnable low-rank matrices for it, and approximate weight updates through their product. Since the numbers of parameters of these low-rank matrices are much smaller than that of the original pre-trained weights, LoRA significantly reduces the memory overhead during fine-tuning.

Despite its widespread success, LoRA has inherent limitations, particularly in its ability to model complex weight adaptation behaviors. LoRA approximates the weight change with the product of two low-rank matrices. While recent studies have observed that the cumulative weight updates during fine-tuning often exhibit approximately low-rank structure (Zhao et al., 2024), LoRA itself learns these updates from scratch using randomly initialized parameters, without leveraging any prior knowledge from the pre-trained weights. As a result, the optimization process becomes more challenging, especially under low-rank settings where the parameter space is highly constrained and prone to suboptimal local minima (Pan et al., 2024). Furthermore, due to its linear structure, LoRA may struggle to capture intricate adaptation patterns required by many downstream tasks. To compensate for this limited capacity, LoRA-based methods often resort to increasing the rank of the update matrices, which in turn reduces their parameter efficiency and undermines their original motivation.

To overcome these limitations, we propose a **p**arameter-**e**fficient **a**daptation method with weight-aware **neu**ral **t**weakers, PEANuT, which incorporates a lightweight neural network, which takes the *pre-traiend weight* as the input, into the adaptation process. Unlike LoRA, which approximates weight updates linearly through low-rank decomposition, PEANuT models cumulative weight updates as explicit functions of the pre-trained model's original weights. This enables PEANuT to capture complex, non-linear patterns in the weight space, improving adaptation performance without increasing the number of parameters. The key innovation in PEANuT lies in introducing compact neural networks, *neural tweakers*, that transforms the pre-trained weights, approximating the updates with minimal additional computation. This nonlinear transformation enhances the expressiveness of the parameter updates while maintaining the efficiency. Importantly, this architecture facilitates a more efficient exploration of the optimization landscape, leading to better task adaptation, particularly in cases where linear methods like LoRA would require much larger ranks to achieve competitive results. We theoretically demonstrate that PEANuT can achieve the same or greater expressivity than LoRA with fewer parameters.

The contributions are summarized as follows:

- We propose PEANuT, a new PEFT method that introduces weight-aware neural tweakers to generate adaptive update signals. The method enables efficient and flexible adaptation beyond linear constraints. To the best of our knowledge, this is the first work to introduce nonlinear adaptation for LoRA-based PEFT methods.
- The proposed PEANuT enhances model performance while maintaining the efficiency. We theoretically show that PEANuT can achieve a possibly improved parameter efficiency compared to LoRA.
- We conduct extensive experiments on four benchmarks covering over twenty datasets. The experiments show that the proposed PEANuT can outperform baselines on both vision and text tasks.

## 2. Related Work

In this section, we provide a concise overview of related work on Parameter-Efficient Fine-Tuning (PEFT) methods. PEFT methods aim to reduce the memory overhead of fine-tuning pre-trained models, enabling fine-tuning in resource-constrained environments. According to Han et al. (2024), PEFT methods can be categorized into: 1) **Additive PEFT methods** (Chronopoulou et al., 2023, Edalati et al., 2022, Lester et al., 2021, Wang et al., 2024d, Liu et al., 2022), 2) **Selective PEFT methods** (Guo et al., 2020, Das et al., 2023, Sung et al., 2021, Ansell et al., 2021, Zaken et al., 2021, Vucetic et al., 2022, Chen et al., 2024, Miao et al., 2025, Chen et al., 2025), 3) **Reparameterized PEFT methods** (Hu et al., 2021a, Valipour et al., 2022, Zhang et al., 2023, Karimi Mahabadi et al., 2021, Liu et al., 2024, Kopiczko et al., 2023), and 4) **Hybrid PEFT methods** (Mao et al., 2021, Chen et al., 2023, He et al., 2021, Zhang et al., 2022, Zhou et al., 2024). *Additive PEFT methods* (Chronopoulou et al., 2023, Edalati et al., 2022, Lester et al., 2021, Wang et al., 2024d, Liu et al., 2022) introduces a small set of additional trainable parameters strategically placed within the model. One of the most prominent additive PEFT approaches is Adapter (Chronopoulou et al., 2023, Edalati et al., 2022, Zhao et al., 2022), which involves inserting small adapter layers between pre-trained weight blocks. Prompt Tuning (Wang et al., 2024d, Lester et al., 2021, Vu et al., 2021, Li and Liang, 2021) is another technique, where learnable vectors, or "soft prompts," are prepended to the input sequence without modifying the model's weights. This method is particularly effective for large-scale models and has inspired variants such as Prefix Tuning (Li and Liang, 2021). *Selective PEFT* focuses on optimizing the fine-tuning process by selectively adjusting a subset of the model's parameters rather than introducing additional ones. For instance, Diff Pruning (Guo et al., 2020) uses a learnable binary mask to select parameters for fine-tuning. Similarly, FishMask (Sung et al., 2021) and Fish-Dip (Das et al., 2023) leverage Fisher information to determine parameter importance and identify the most crucial ones for updates. Additionally, BitFit (Zaken et al., 2021) fine-tunes only the bias terms in the model, significantly reducing the number of trainable parameters. *Hybrid PEFT* methods aim to combine the strengths of various existing PEFT techniques to enhance model performance across diverse tasks. UniPELT (Mao et al., 2021) integrates LoRA, prefix-tuning, and adapters within each Transformer block, employing a gating mechanism to determine which module should be active during fine-tuning. S4 (Chen et al., 2023) further explores the design space by partitioning layers into groups and assigning different PEFT methods to each group. Additionally, NOAH (Zhang et al., 2022) and AUTOPEFT (Zhou et al., 2024) leverage neural architecture search (NAS) to automatically discover optimal combinations of PEFT techniques tailored to specific tasks.

*Reparameterized PEFT* methods are most close to our proposed method. Low-Rank Adaptation (LoRA)-based methods, which are representative of reparameterized PEFT approaches, have gained significant attention due to their minimal architectural changes, no additional inference costs, and high efficiency. LoRA (Hu et al., 2021a) introduces two trainable low-rank matrices for each pre-trained model weight to approximate the desired updates of the original model. Extensions of LoRA include DyLoRA (Valipour et al., 2022), which dynamically adjusts the rank of the low-rank matrices during training to optimize for specific tasks; AdaLoRA (Zhang et al., 2023), which adaptively allocates the parameter budget among weight matrices based on their importance scores; and DoRA (Liu et al., 2024), which decomposes the pre-trained weight into magnitude and direction, applying LoRA only for direction updates. Other variants include VeRA (Kopiczko et al., 2023), which introduces shared frozen random matrices across layers to improve efficiency further, and RoseLoRA (Wang et al., 2024c), which employs a row- and column-wise sparse low-rank adaptation mechanism to selectively update the most significant parameters. FourierFT (Gao et al.) replaces the matrix multiplication in LoRA with a Fourier transform,
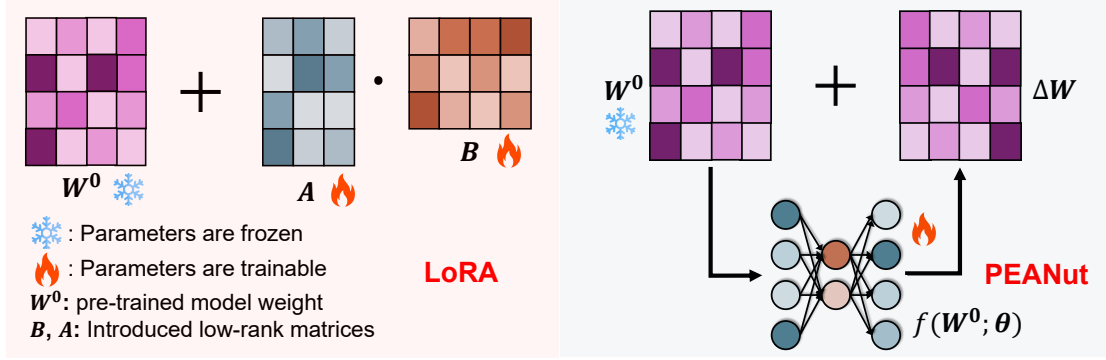
**Figure 1:** Framework of proposed PEANuT.

while PiSSA (Meng et al., 2024) and MiLoRA (Wang et al., 2024a) update the principal and minor singular components of the weight matrix, respectively. However, existing PEFT methods rely on linear transformations to approximate pre-trained weight updates, which struggle to capture the complex relationships inherent in weight updates, leading to a significant performance gap compared to full fine-tuning. Meanwhile, existing research like (Teney et al., 2024) also demonstrates that nonlinear activation is an integral part of the neural network driving its success.

## 3. Methodology

In this section, we start with a brief introduction of LoRA. Motivated by a key limitation in LoRA parameter efficiency that roots from LoRA parameterization form, we propose PEANuT, a novel PEFT method to solve the issue. Notably, PEANuT is able to achieves better parameter efficiency provably.

### 3.1. Preliminary

LoRA (Hu et al., 2021a) assumes that the updates to model weights during the fine-tuning exhibit low-rank properties. Built upon this, LoRA models the *incremental update* of some weight matrix $\mathbf{W}^0 \in \mathbb{R}^{d_1 \times d_2}$ in a pre-trained model approximately by the product of two learnable low-rank matrices

$$\mathbf{W} = \mathbf{W}^0 + \Delta \mathbf{W} = \mathbf{W}^0 + \mathbf{AB},$$

where $\mathbf{A} \in \mathbb{R}^{d_1 \times r}$ and $\mathbf{B} \in \mathbb{R}^{r \times d_2}$ with $r \ll \min(d_1, d_2)$. When conducting fine-tuning, only introduced two low-rank matrices $\mathbf{A}$ and $\mathbf{B}$ will be updated and the pre-trained weight $\mathbf{W}^0$ is frozen, as represented by the following optimization

$$\min_{\mathbf{A}, \mathbf{B}} \ \mathcal{L}(\mathcal{D}_{\text{train}}; \mathbf{W}^0 + \mathbf{AB}), \tag{1}$$

where $\mathcal{D}_{\text{train}}$ is the training set used for fine-tuning and $\mathcal{L}$ is the loss function. Since $\mathbf{A}$ and $\mathbf{B}$ are both low-rank matrices that contain significantly fewer parameters compared with the original $\mathbf{W}^0$, the LoRA costs much less memory space compared to the fully fine-tuning.

## 3.2. Inherent Limitation of LoRA Formulation

While LoRA family have demonstrated remarkable parameter efficiency in fine-tuning pre-trained models for diverse downstream tasks, we argue that their product-based formulation are suboptimal for capturing the full fine-tuning dynamics in an efficient way.

Specifically, when fully fine-tuning a pre-trained model, the update process of weight $\mathbf{W}$ is typically performed through an iterative gradient descent:

$$\mathbf{W}_t^0 = \mathbf{W}_{t-1}^0 - \eta \nabla_{\mathbf{W}_{t-1}^0} \mathcal{L},$$

where $\mathbf{W}_0^0 = \mathbf{W}^0$ is the initial state, $\eta$ is the learning rate, and $\mathbf{W}_t^0$ represents the weights after $t$ iterations. The cumulative change in the weights over time can be represented as:

$$\Delta\mathbf{W} = \mathbf{W}_t^0 - \mathbf{W}_0^0.$$

This weight change $\Delta\mathbf{W}$ can be interpreted as a function of the original pre-trained weights $\mathbf{W}^0$, capturing the model's adaptation to the specific task during fine-tuning.

Nonetheless, LoRA matrices $\mathbf{A}$ and $\mathbf{B}$ are parameterized in a free way without any dependency on $\mathbf{W}^0$. While gradient $\nabla_{\mathbf{A}}\mathcal{L}$ and $\nabla_{\mathbf{B}}\mathcal{L}$ are implicit functions of $\mathbf{W}^0$, making final learned $\mathbf{A}_t, \mathbf{B}_t$ indirectly depends on $\mathbf{W}^0$ as well, as will be proved shortly, the lack of explicit dependency still makes LoRA *inherently suboptimal* for fine-tuning pre-trained models.

## 3.3. Parameter-Efficient Adaptation with Weight-aware Neural Tweakers

Motivated by the above analysis on LoRA's limitation, we propose to approximate $\Delta\mathbf{W}$ using a lightweight neural network that *explicitly* takes pre-trained model weight $\mathbf{W}^0$ as input and outputs the weight update directly. By doing so, our approach captures more complex and richer transformation of the weights in a more efficient manner. We refer to our method as **p**arameter-**e**fficient **a**daptation method with weight-aware **neu**ral **t**weakers (PEANuT).

Following LoRA's updates paradigm, the proposed PEANuT also provides incremental update of pre-trained models. However, PEANuT modifies the forward pass of the model by introducing a dynamic *nonlinear* weight transformation. Specifically, the modified model's forward propagation is formulated as:

$$y = (\mathbf{W}^0 + f(\mathbf{W}^0; \boldsymbol{\theta}))x.$$

Here $x$ and $y$ are the input and output with respect to the current layer, respectively, and $f(\cdot; \boldsymbol{\theta}) : \mathbb{R}^{d_1 \times d_2} \to \mathbb{R}^{d_1 \times d_2}$ is a nonlinear neural network parameterized by learnable parameter $\boldsymbol{\theta}$. The neural network $f(\mathbf{W}^0; \boldsymbol{\theta})$ generates the weight update as a function of $\mathbf{W}^0$.

To ensure the parameter efficiency of our PEANuT, the learnable neural network $f(\mathbf{W}^0; \boldsymbol{\theta})$ should be lightweight, i.e., the number of parameters $\boldsymbol{\theta}$ should be much fewer than that of the original pre-trained weight $\mathbf{W}^0$. Therefore, we parametrize $f(\mathbf{W}^0; \boldsymbol{\theta})$ as a neural network with *bottleneck* layers. For example, a simple case is $f(\mathbf{W}^0; \boldsymbol{\theta}) = \sigma(\mathbf{W}^0 \boldsymbol{\Theta}_1)\boldsymbol{\Theta}_2$, where $\boldsymbol{\theta} = (\boldsymbol{\Theta}_1, \boldsymbol{\Theta}_2) \in \mathbb{R}^{d_2 \times r} \times \mathbb{R}^{r \times d_2}$ with $r \ll \min(d_1, d_2)$, and $\sigma(\cdot)$ is some non-linear activation function such as ReLU (Glorot et al., 2011). We can also increase the layers or add activation function for the output of $f(\mathbf{W}^0; \boldsymbol{\theta})$ to enhance the model expressiveness.

During fine-tuning, the optimization objective is to minimize the task-specific loss function, which can be represented as

$$\min_{\boldsymbol{\theta}} \ \mathcal{L}(\mathcal{D}_{\text{train}}; \mathbf{W}^0 + f(\mathbf{W}^0; \boldsymbol{\theta})),$$

where the original pre-trained weight $\mathbf{W}^0$ is frozen, and only neural network parameters $\theta$ are updated. The overview of PEANuT is shown in Fig. 1.

*Remark* 3.1. The benefit of our new formulation lies in two folds. First, our incremental update $f(\mathbf{W}^0; \boldsymbol{\theta})$ is an explicit function of $\mathbf{W}^0$, allowing it to capture updates in a more effective way. Second, the neural network-based $f(\mathbf{W}^0; \boldsymbol{\theta})$ allows for dynamic, non-linear weight updates that can capture more complex interactions. These two advantages make PEANuT a more effective and efficient PEFT method than existing LoRA-based approaches.

## 3.4. Theoretical Analysis

In this section, we show the theoretical analysis of the sub-optimality of LoRA in terms of parameter efficiency. We prove that PEANuT can achieve equivalent or even superior efficiency under certain conditions. Specifically, suppose PEANuT adopts the following lightweight architecture, as described in Section 3.3:

$$f(\mathbf{W}^0; \boldsymbol{\theta}) = \sigma(\mathbf{W}^0 \boldsymbol{\Theta}_1) \boldsymbol{\Theta}_2.$$

The following proposition demonstrates that PEANuT can match the expressivity of LoRA using fewer parameters under specific conditions. Here, expressivity is measured by the minimum attainable loss.

**Proposition 3.2.** *Given pre-trained weight matrix* $\mathbf{W}^0$. *Let* $\sigma$ *denote ReLU activation function, and* $\boldsymbol{U}^0 \in \mathbb{R}^{d_1 \times \text{rank}(\mathbf{W}^0)}$ *be the left singular vectors of* $\mathbf{W}^0$. *Suppose that the fine-tuning loss* $\mathcal{L}$ *is invariant under the the projection of the weight matrix to the left singular space of* $\mathbf{W}^0$, *i.e.,* $\mathcal{L}(\mathcal{D}_{train}; \mathbf{W}) = \mathcal{L}(\mathcal{D}_{train}; \boldsymbol{U}^0 \boldsymbol{U}^{0\top} \mathbf{W})$ *for any* $\mathbf{W} \in \mathbb{R}^{d_1 \times d_2}$. *Then, for any* $r \geq 1$,

$$\min_{\substack{\boldsymbol{\Theta}_1 \in \mathbb{R}^{d_2 \times 2r}, \\ \boldsymbol{\Theta}_2 \in \mathbb{R}^{2r \times d_2}}} \mathcal{L}(\mathcal{D}_{train}; \mathbf{W}^0 + f(\mathbf{W}^0; (\boldsymbol{\Theta}_1, \boldsymbol{\Theta}_2)))$$

$$\leq \min_{\substack{\mathbf{A} \in \mathbb{R}^{d_1 \times r}, \\ \mathbf{B} \in \mathbb{R}^{r \times d_2}}} \mathcal{L}(\mathcal{D}_{train}; \mathbf{W}^0 + \mathbf{AB})$$

$$\leq \min_{\substack{\boldsymbol{\Theta}_1 \in \mathbb{R}^{d_2 \times r}, \\ \boldsymbol{\Theta}_2 \in \mathbb{R}^{r \times d_2}}} \mathcal{L}(\mathcal{D}_{train}; \mathbf{W}^0 + f(\mathbf{W}^0; (\boldsymbol{\Theta}_1, \boldsymbol{\Theta}_2))).$$

In words, Prop 3.2 demonstrates the (approximate) equivalence of LoRA and PEANuT in terms of their expressivity. Specifically, the minimum attainable loss using rank-$r$ LoRA can be achieved by PEANuT with $2r$ hidden units, and conversely, the minimum attainable loss using PEANuT with $r$ hidden units can be achieved rank-$r$ LoRA, provided the invariance assumption holds. This equivalence further implies that the function classes realized by PEANuT with $O(r)$ hidden dimensions and rank-$r$ LoRA are equivalent in expressivity, as the result holds for any loss functions.

Importantly, *this highlights a potential improvement in parameter efficiency by PEANuT*. Namely, PEANuT with $O(rd_2)$ parameters maintains the expressivity of LoRA with $r(d_1 + d_2)$ parameters. That it to say,

PEANuT offers a significant improvement in parameter efficiency when $d_2 \ll d_1$ ( a condition that widely holds for the down projection matrix of transformers fully-connected layers (Vaswani, 2017, Dosovitskiy et al., 2021) ). In such cases, PEANuT provably achieves better parameter efficiency than LoRA. The added parameter efficiency can also improve sample efficiency by allowing the model to learn representations with the same or fewer data points.

The invariance assumption in Proposition 3.2 pertains to the pre-trained model, and asserts that the later layers of the model depends solely on the task-relevant feature space. Given that we fine-tune a pre-trained model, the later layers are expected to capture this task-relevant feature space, which is described by the left singular space of $\mathbf{W}^0$. In practice, since the later layers primarily rely on this pre-trained feature space, the principal directions of the pre-trained weight matrix, represented by its singular vectors, encode most of the useful features for downstream tasks. This makes the loss largely invariant to changes outside this subspace. The proof is available in Appendix B.1.

If we consider a sinusoid activation function $\sigma_{\mathrm{p}}(x) = \sin(2\pi x)$, then stronger result that **PEANuT has expressivity (almost) greater than or equal to a LoRA with possibly more parameters can be established without the invariance assumption**. We defer the result to the Appendix B.2.

## 4. Complexity Analysis

In this section, we compare the computational and space complexity of PEANuT and LoRA.

**Space Complexity.** Because we set the introduced parameters of LoRA and PEANuT to be the same, we only discuss the space complexity of the training in this section. Both LoRA and PEANuT require storing the added parameters and their gradients. PEANuT may incur a slightly higher activation memory during backpropagation due to the extra nonlinearity, but our empirical results (see Sec. 5.4) show that this overhead is minimal and does not affect scalability in practice.

**Computational complexity.** In terms of per-step computation cost, LoRA computes the residual update as $\mathbf{AB}x$, which costs $\mathcal{O}(d_1 r + r d_2)$ per input vector $x$. PEANuT requires computing $f(\mathbf{W}_0; \theta)x$. When using the aforementioned example with one-hidden layer and having the same latent dimension as LoRA, the main cost is $\mathcal{O}(d_1 d_2 r)$. Although this complexity is higher than LoRA's, both methods benefit from highly matrix-friendly implementations. In practice, our experiments (see Sec. 5.4) show that the empirical training time per step is comparable. Importantly, during inference, both PEANuT and LoRA allow their update modules to be merged into the original weight matrix $\mathbf{W}_0$, ensuring that no additional forward-pass cost is incurred in deployment.

## 5. Experiment

In the experiments, we evaluate the proposed PEANuT and answer the following questions:

**RQ1** How does PEANuT compare to state-of-the-art PEFT methods on NLP and vision tasks?
**RQ2** What is the role of nonlinear approximation in the proposed PEANuT?
**RQ3** What is the real runtime and memory consumption of proposed PEANuT?
**RQ4** How does the performance of PEANuT vary with different fine-tuned modules, depths of the lightweight neural network, or non-linear activation functions?

## 5.1. Benchmarks and Experiment Setups

We experiment PEANuT on datasets from four representative benchmarks: 1) **Commonsense Reasoning** covers diverse multi-choice problems from BoolQ (Clark et al., 2019), PIQA (Bisk et al., 2020), SIQA (Sap et al., 2019), HellaSwag (Zellers et al., 2019), WinoGrande (Sakaguchi et al., 2019), ARC-e and ARC-c (Clark et al., 2018), and OpenBookQA (Mihaylov et al., 2018) datasets. Following Wang et al. (2024a), we finetune LLaMA2-7B (Touvron et al., 2023), LLaMA3-8B (AI@Meta, 2024) and Qwen3-8B (Yang et al., 2025) on Commonsense170K (Hu et al., 2023) benchmark which combines all previous training sets, and evaluate the accuracy on their testing sets separately. 2) **Arithmetic Understanding** consists of two math reasoning datasets: GSM8K (Cobbe et al., 2021) and MATH (Hendrycks et al., 2021). We finetune LLaMA2-7B (Touvron et al., 2023) and Qwen3-8B (Yang et al., 2025) on MetaMath (Yu et al., 2023) dataset following Wang et al. (2024a). Models need to generate correct answers, and accuracy is used as the evaluation metric. 3) **Natural Language Understanding** consists of eight datasets from the GLUE benchmark (Wang et al., 2018). We follow the evaluation metrics and setups from Gao et al. (2024), Wu et al. (2024b). 4) **Image Classification** consists of Oxford-Pets (Parkhi et al., 2012), CIFAR10 (Krizhevsky, 2009), DTD (Cimpoi et al., 2014), EuroSAT (Helber et al., 2019), RESISC45 (Cheng et al., 2017), StanfordCars (Krause et al., 2013), FGVC (Maji et al., 2013), and CIFAR100 (Krizhevsky, 2009) following Gao et al. (2024). The first five datasets have small label spaces, while the last three have large label spaces.

**Baselines methods** are constructed on a task basis. Specifically, for each task, the proposed PEANuT is compared with representative baselines from the corresponding domain. For both Commonsense Reasoning and Arithmetic Understanding, following Wang et al. (2024a), LoRA (Hu et al., 2021b), PiSSA (Meng et al., 2024) and MiLoRA (Wang et al., 2024a) are employed as baselines. PEANuT is applied to `query`, `key`, `value`, `MLP up` and `MLP down` layers. For Natural Language Understanding, we follow the setup from prior works (Gao et al., 2024, Wu et al., 2024b) that evaluate various representative PEFT methods, including LoRA (Hu et al., 2021b), Adapter Houlsby et al. (2019), BitFit (Zaken et al., 2021), RED (Wu et al., 2024a), DoRA (Liu et al., 2024), ReFT Wu et al. (2024b), and FourierFT (Gao et al., 2024). For Image Classification, we follow the setting of Gao et al. (2024) and take linear probing (LP), LoRA (Hu et al., 2021b) and FourierFT (Gao et al., 2024) as baselines. PEANuT is applied to the query and value layers. See our appendix for details about the datasets (App D) and hyper-parameters (App C).

## 5.2. Performance Comparison

We showcase PEANuT performance on different tasks.

**Commonsense Reasoning.** We experiment PEANuT with eight commonsense reasoning datasets to address RQ1, results are shown in Tab 1. We compare the performance of three state-of-the-art baselines with the proposed PEANuT, and PEANuT consistently outperforms all of them, achieving the highest accuracy on all tasks. Specifically, PEANuT surpasses LoRA, PiSSA, and MiLoRA in terms of average accuracy by 4.6%, 10%, and 2.5%, respectively, when using LLaMA2-7B as the backbone. On LLaMA3-8B as the backbone, PEANuT demonstrates average improvements of 4.9%, 11.8%, and 2.9% over LoRA, PiSSA, and MiLoRA, respectively. With Qwen3-8B as the backbone, PEANuT improves average accuracy over LoRA and MiLoRA by 3.9% and 2.4%, respectively. These results highlight the effectiveness and superiority of PEANuT as a PEFT method.

**Arithmetic Reasoning.** In this section, we present results on two arithmetic reasoning tasks in Tab 4 to help address RQ1. From the table, while full fine-tuning (FFT) achieves highest accuracy across the two datasets, the performance gap between the proposed PEANuT and FFT is very small, despite that PEANuT relies on significantly fewer trainable parameters. Moreover, compared to state-of-the-art PEFT baselines, PEANuT achieves remarkable performance improvements. In terms of average accuracy, PEANuT demonstrates improvements of 7.5%, 12.4%, and 2.4% over LoRA, PiSSA, and MiLoRA, respectively, when using LLaMA2-7B as the backbone. With Qwen3-8B as the backbone, PEANuT improves average accuracy over LoRA and MiLoRA by 7.1% and 3.7%. These results on clearly confirm that PEANuT is highly effective and efficient for complex reasoning tasks.

**Natural Language Understanding.** We further conduct experiments on the GLUE to answer RQ1, results are shown in Tab 3. From the table, PEANuT significantly outperforms state-of-the-art PEFT methods. Specifically, PEANuT-S, which uses a similar number of trainable parameters as FourierFT (Gao et al., 2024), DiReFT (Wu et al., 2024b), and LoReFT (Wu et al., 2024b), surpasses all PEFT baselines and experiences only a small performance drop (0.2%) compared to FFT. Additionally, PEANuT-L exceeds the performance of all baselines, including FFT, with roughly the same number of trainable parameters as in LoRA. These results demonstrate that PEANuT exhibits excellent generalization ability while maintaining great parameter efficiency.

**Image Classification.** In this section, we conduct experiments on image classification tasks to address RQ2, PEANuT uses depth of 6, and results are shown in Tab 2. From the table, PEANuT significantly outperforms LoRA and FourierFT using the same number of trainable parameters. Specifically, PEANuT achieves performance improvements of 11.05%, 7.30%, and 26.02% compared to LoRA, FourierFT, and LP, respectively. Furthermore, compared to FFT, the proposed PEANuT shows negligible performance drop (86.49% v.s. 86.34%), while using only 0.3% of the trainable parameters required by FFT. This demonstrates that PEANuT exhibits exceptional adaptation capability not only on NLP tasks, but also on vision tasks as well. Additionally, it verifies the effectiveness of the nonlinear adaptation used in PEANuT.

## 5.3. Ablation Study

In this section, in order to answer RQ2, we present an ablation study with two variants of LoRA to validate the effectiveness of our proposed framework: 1) nonlinear LoRA $y = (W_0 + \sigma(\mathbf{A})\mathbf{B})x$, and 2) multiplicative LoRA $y = (W_0 + W_0\mathbf{A}\mathbf{B})x$. Experiments are conducted on image classification benchmarks, and results are reported in Tab 5. According to the table, both nonlinear LoRA and multiplicative LoRA perform worse than PEANuT. This highlights the effectiveness of incorporating nonlinear approximations and explicitly using model weights as input to the nonlinear function in PEANuT.

## 5.4. Runtime and Memory Cost

To answer RQ3, we evaluate the computational efficiency of our proposed method, PEANuT, by measuring its runtime and memory consumption across three representative datasets: MRPC, SST-2, and Commonsense Reasoning. Table 7 summarizes the results, comparing PEANuT against the LoRA approach under identical settings. For a fair comparison, we ensure that the number of trainable parameters is matched between PEANuT and LoRA. All experiments are conducted on the same hardware setup using a single NVIDIA A100 GPU. As shown in Table 7, PEANuT exhibits comparable

**Table 1:** Common Reasoning performance of PEANuT and PEFT baselines on LLaMA 2-7B, LLaMA 3-8B and Qwen 3-8B. Results marked with "+" are taken from Liu et al. (2024), and those marked with "∗" are taken from Wang et al. (2024a). Best results are in **bold**. "AVG" means the average accuracy of all datasets.

| Model | PEFT | Accuracy (↑) | | | | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | BoolQ | PIQA | SIQA | HellaSwag | WinoGrande | ARC-e | ARC-c | OBQA | AVG |
| LLaMA2-7B | LoRA[+] | 69.8 | 79.9 | 79.5 | 83.6 | 82.6 | 79.8 | 64.7 | 81.0 | 77.6 |
| | PiSSA[∗] | 67.6 | 78.1 | 78.4 | 76.6 | 78.0 | 75.8 | 60.2 | 75.6 | 73.8 |
| | MiLoRA[∗] | 67.6 | 83.8 | 80.1 | 88.2 | 82.0 | 82.8 | 68.8 | 80.6 | 79.2 |
| | PEANuT | **71.9** | **84.0** | **80.4** | **88.9** | **84.6** | **86.5** | **71.6** | **83.0** | **81.4** |
| LLaMA3-8B | LoRA[+] | 70.8 | 85.2 | 79.9 | 91.7 | 84.3 | 84.2 | 71.2 | 79.0 | 80.8 |
| | PiSSA[∗] | 67.1 | 81.1 | 77.2 | 83.6 | 78.9 | 77.7 | 63.2 | 74.6 | 75.4 |
| | MiLoRA[∗] | 68.8 | 86.7 | 77.2 | 92.9 | 85.6 | 86.8 | 75.5 | 81.8 | 81.9 |
| | PEANuT | **72.1** | **87.0** | **80.9** | **94.3** | **86.7** | **91.4** | **78.9** | **84.8** | **84.5** |
| Qwen3-8B | LoRA | 86.3 | 87.2 | 84.1 | 92.5 | 81.5 | 89.6 | 78.8 | 89.5 | 86.2 |
| | MiLoRA | 85.2 | 89.3 | 84.2 | 94.6 | 82.2 | 92.3 | 82.7 | 89.5 | 87.5 |
| | PEANuT | **89.4** | **90.2** | **87.4** | **95.7** | **85.5** | **92.7** | **82.6** | **93.6** | **89.6** |

**Table 2:** Image Classification performance on ViT-base. Best results are in **bold**. "AVG" means the average accuracy of all datasets. Results marked with "∗" are taken from Gao et al. (2024).

| Method | Params (M) | OxfordPets | StanfordCars | CIFAR10 | DTD | EuroSAT | FGVC | RESISC45 | CIFAR100 | AVG |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| FFT[∗] | 85.8M | 93.14 | 79.78 | **98.92** | 77.68 | **99.05** | **54.84** | **96.13** | **92.38** | **86.49** |
| LP[∗] | - | 90.28 | 25.76 | 96.41 | 69.77 | 88.72 | 17.44 | 74.22 | 84.28 | 68.36 |
| LoRA[∗] | 581K | 93.19 | 45.38 | 98.78 | 74.95 | 98.44 | 25.16 | 92.70 | 92.02 | 77.58 |
| FourierFT[∗] | 239K | 93.05 | 56.36 | 98.69 | 77.30 | 98.78 | 32.44 | 94.26 | 91.45 | 80.29 |
| PEANuT | 263K | **93.62** | **80.21** | 98.78 | **79.61** | 98.85 | 52.93 | 94.71 | 92.02 | 86.34 |

runtime and memory usage to LoRA across all tasks. On the MRPC and SST-2 datasets, PEANuT incurs only marginal overhead in training time, with identical memory consumption. For the larger Commonsense Reasoning dataset, PEANuT takes 5.7 hours and 23.8GB of memory, compared to LoRA's 5.6 hours and 22.7GB. The slightly higher memory usage is attributed to the additional nonlinear transformation module in PEANuT, but the increase remains slight. Overall, the results demonstrate that PEANuT achieves improved performance (as discussed in earlier sections) with negligible additional cost in training runtime and memory, highlighting its practicality and scalability for real-world deployment.

## 5.5. Sensitivity w.r.t. Depth

To answer RQ4, we analyze the impact of depth on the performance of PEANuT. Deeper architectures are generally more expressive and can better model the complex, nonlinear relationships involved in ideal weight updates (Raghu et al., 2017). We evaluate PEANuT with varying depth across NLU, vision, commonsense reasoning, and arithmetic reasoning tasks.

**Table 3:** GLUE benchmark performance on RoBERTa-base. Results marked with "∗" are taken from Wu et al. (2024a). Best results are in **bold**. "AVG" means the average accuracy of all datasets. PEANuT-S applies trainable modules to layers starting from the 4th layer, with hidden dimensions set to 1. This matches the parameter numbers of FourierFT. PEANuT-L applies PEANuT to all layers with hidden dimension 8, aligning the parameter budget of LoRA.

| PEFT | Params (%) | Accuracy (↑) | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | MNLI | SST-2 | MRPC | CoLA | QNLI | QQP | RTE | STS-B | AVG |
| FFT | 100% | 87.3 | 94.4 | 87.9 | 62.4 | 92.5 | 91.7 | 78.3 | 90.6 | 85.6 |
| Adapter∗ | 0.318% | 87.0 | 93.3 | 88.4 | 60.9 | 92.5 | **90.5** | 76.5 | 90.5 | 85.0 |
| LoRA∗ | 0.239% | 86.6 | 93.9 | 88.7 | 59.7 | **92.6** | 90.4 | 75.3 | 90.3 | 84.7 |
| Adapter^FNN∗ | 0.239% | **87.1** | 93.0 | 88.8 | 58.5 | 92.0 | 90.2 | 77.7 | 90.4 | 84.7 |
| BitFit∗ | 0.080% | 84.7 | 94.0 | 88.0 | 54.0 | 91.0 | 87.3 | 69.8 | 89.5 | 82.3 |
| RED∗ | 0.016% | 83.9 | 93.9 | 89.2 | 61.0 | 90.7 | 87.2 | 78.0 | 90.4 | 84.3 |
| FourierFT | 0.019% | 84.7 | 94.2 | 90.0 | 63.8 | 92.2 | 88.0 | 79.1 | **90.8** | 85.3 |
| DiReFT∗ | 0.015% | 82.5 | 92.6 | 88.3 | 58.6 | 91.3 | 86.4 | 76.4 | 89.3 | 83.2 |
| LoReFT∗ | 0.015% | 83.1 | 93.4 | 89.2 | 60.4 | 91.2 | 87.4 | 79.0 | 90.0 | 84.2 |
| PEANuT-S | 0.019% | 84.9 | 94.3 | 90.2 | 64.6 | 92.0 | 88.3 | 78.3 | 90.5 | 85.4 |
| PEANuT-L | 0.241% | 86.9 | **95.2** | **90.0** | **64.8** | 92.3 | 90.3 | **82.7** | 90.7 | **86.6** |

**Table 4:** Arithmetic Reasoning performance on LLaMA 2-7B and Qwen 3-8B. Results marked with "+" are taken from Yu et al. (2023), and those marked with "∗" are taken from Wang et al. (2024a). Best results are in **bold**. "AVG" means the average accuracy of all datasets.

| Model | Method | GSM8K | MATH | AVG |
|---|---|---|---|---|
| | FFT + | 66.50 | 19.80 | 43.20 |
| LLaMA2-7B | LoRA∗ | 60.58 | 16.88 | 38.73 |
| | PiSSA∗ | 58.23 | 15.84 | 37.04 |
| | MiLoRA∗ | 63.53 | 17.76 | 40.65 |
| | PEANuT | **65.05** | **18.30** | **41.68** |
| Qwen3-8B | LoRA | 85.22 | 67.26 | 76.24 |
| | MiLoRA | 89.01 | 68.58 | 78.80 |
| | PEANuT | **92.87** | **70.50** | **81.69** |

**Table 5:** Ablation Study on image classification task. The parameters count is the same and "AVG" means the average accuracy of all datasets. For simple and fair comparison, PEANuT uses depth of 2.

| Method | OxfordPets | StanfordCars | CIFAR10 | DTD | EuroSAT | FGVC | RESISC45 | CIFAR100 | AVG |
|---|---|---|---|---|---|---|---|---|---|
| Nonlinear LoRA | 94.11 | 72.84 | 98.68 | 79.16 | 98.61 | 39.33 | 93.79 | 92.38 | 83.31 |
| Multiplicative LoRA | 93.57 | 77.32 | 98.68 | 77.57 | 98.81 | 46.79 | 94.34 | 91.86 | 84.81 |
| PEANuT | 93.77 | 80.03 | 98.70 | 77.57 | 98.79 | 53.60 | 94.27 | 92.47 | 86.15 |

We increase the number of intermediate layers inserted between PEANuT's input and output projections. Each intermediate layer is a small feedforward block of shape $\mathbb{R}^{r \times r}$ with non-linear activations. These

**Table 6:** Accuracy comparison of PEANuT using RoBERTa-base with different depth configurations on the GLUE benchmark. The highest accuracy of methods per category are in **bold**. "AVG" means the average accuracy of all datasets.

| depth | Params (%) | Accuracy ($\uparrow$) | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | MNLI | SST-2 | MRPC | CoLA | QNLI | QQP | RTE | STS-B | AVG |
| 2 | 0.239% | 86.6 | 94.6 | 90.0 | 64.4 | **92.7** | 89.7 | 78.7 | **90.9** | 86.0 |
| 4 | 0.239% | 86.7 | 94.5 | **90.2** | **65.1** | 92.4 | **90.5** | 80.5 | 90.8 | 86.3 |
| 6 | 0.241% | **86.9** | **95.2** | 90.0 | 64.8 | 92.3 | 90.3 | **82.7** | 90.7 | **86.6** |

**Table 7:** Runtime and memory consumption of proposed PEANuT.

| Dataset | Method | Time | Memory |
|---|---|---|---|
| MRPC | LoRA | 77.7s | 6916MB |
| | PEANuT | 78.4s | 6916MB |
| SST-2 | LoRA | 870.7s | 2410MB |
| | PEANuT | 911.4s | 2410MB |
| Commonsense | LoRA | 5.6h | 22.7GB |
| Reasoning | PEANuT | 5.7h | 23.8GB |

layers are lightweight compared to the input/output projections ($\mathbf{A} \in \mathbb{R}^{d_2 \times r}$, $\mathbf{B} \in \mathbb{R}^{r \times d_2}$), and add minimal overhead since $r \ll d_2$. The adaptation starts from the frozen base weight $\mathbf{W}^0$, which is transformed through multiple layers to predict $\Delta\mathbf{W}$. We adopt residual connections for stable optimization and improved convergence. All other hyperparameters are kept fixed during this analysis. The layer structure is illustrated in Fig. 2.

The results in Table 6 (GLUE) and Fig. 3 (RTE, Cars, PIQA, MATH) indicate that increasing depth consistently improves accuracy. For instance, average GLUE accuracy increases from 86.0 to 86.6 when moving from 2 to 6 layers, with no significant change in parameter count. On other benchmarks, deeper configurations yield steady gains up to 6 layers. Beyond this, performance may slightly drop (e.g., at depth 10), likely due to optimization difficulties without fine-grained hyperparameter tuning.

In summary, depth enhances PEANuT's effectiveness across tasks, offering better adaptation capability with negligible cost in memory or parameters. However, very deep settings may require further tuning to maintain stability.

## 5.6. Sensitivity w.r.t. Activations

One key innovation of PEANuT compared to LoRA and other PEFT methods, which rely solely on linear transformations for modeling weight updates, is the introduction of non-linear activations within the adaptation neural network. Since the choice of non-linear activations directly affects the learning process
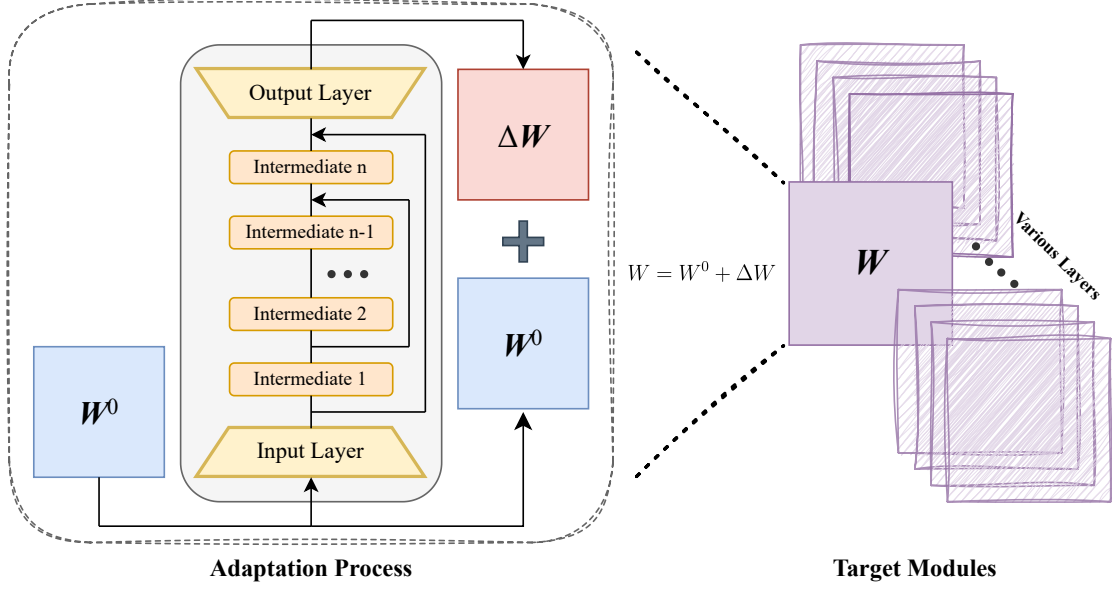
**Figure 2:** Implementation of introducing more depths to PEANuTt. We insert multiple intermediate layers into the layers from vanilla PEANuT, with non-linear activation in between. The depth is described as the number of layers in PEANuT, with vanilla PEANuT having a depth of 2 (i.e. the input and output layers).

and the dynamics of weight updates, we investigate how different non-linear activations affects the adaptation performance to address RQ4. To this end, we perform experiments on the StanfordCars benchmark using various non-linear activations, including ReLU, Leaky ReLU, GELU, Tanh, and sinusoidal activation ($\sigma_{\mathrm{p}}(x) = \sin(2\pi x)$). Corresponding results are presented in Fig 4. To ensure a fair comparison, the number of trainable parameters is fixed. We optimize other hyperparameters such as learning rate for better performance.

From the figure, the best performance achieved by different activation functions is similar, indicating that the adaptation potential of various activations is comparable. This implies that PEANuT can benefit from various type of nonlinearity induced by different activations. However, it is also worth noting that sinusoidal activations encounters a performance drop at large learning rates. Consequently, tuning basic hyperparameters such as learning rate can still be beneficial. In conclusion, we suggest ReLU as a default choice in execution, given its practical simplicity (Teney et al., 2024).

## 6. Sensitivity w.r.t. Fine-tuned Module

We end up this section with a study on applying PEANuT to different modules in a ViT, to help better understand RQ4.

Specifically, given the importance of MLP in Transformer architecture, we compare two settings: 1) Following Hu et al. (2021b), we apply PEANuT to the query and value layers (QV layers) in the multi-head self-attention module (MHSA) in ViT. 2) Besides QV layers, we also apply PEANuT to MLP layers. We tune the hidden dimension $r$ to ensure the same parameter scale for fair comparison, and tune the hyperparameters to maximize performance. Corresponding results are shown in Fig. 5.

From the figure, applying PEANuT to the QV layers yields results comparable to applying PEANuT
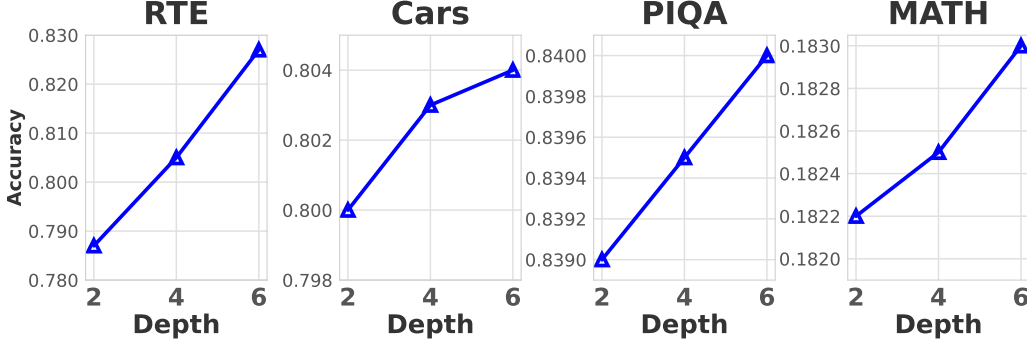
**Figure 3:** Accuracy on the RTE, StanfordCars, PIQA and MATH dataset with varying depths of the neural network used in PEANuT. The depth here represents the total number of layers in the neural network. We choose depth equals to 2, 4 and 6 layers in the figure.
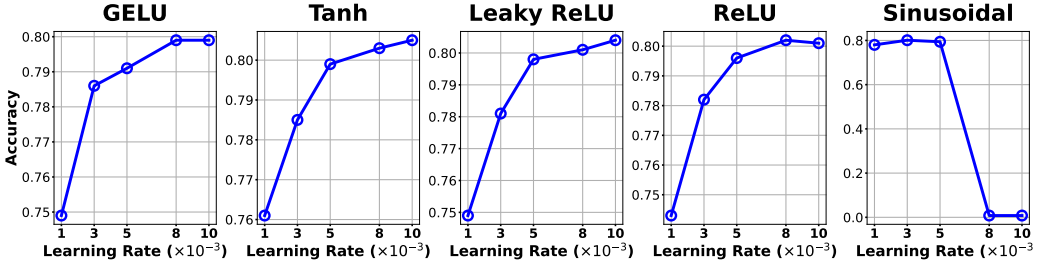


**Figure 4:** Influence of different nonlinear activations choices for PEANuT. Experiments are conducted on StanfordCars, PEANuT depth is fixed to 2. Different activations share a similar pattern of dependency on learning rate.
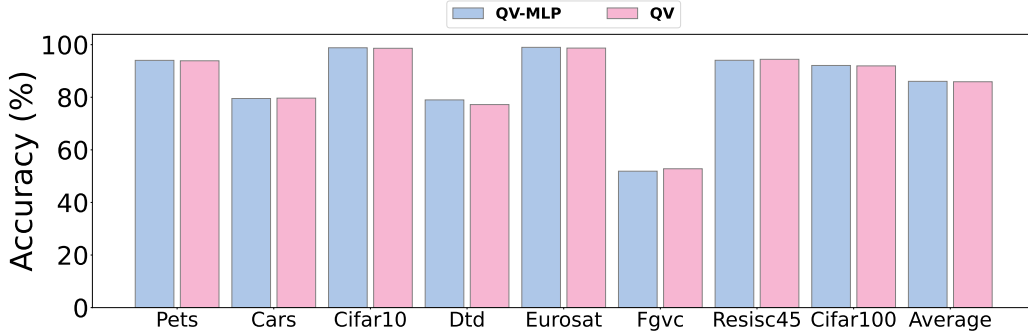


**Figure 5:** Accuracy of PEANuT with different targeted fine-tuning modules, including just QV layers and a combination of QV and MLP layers, on image classification datasets.

to both the QV and MLP layers. This indicates that PEANuT is robust to the selections of fine-tuning different modules. This finding confirms another key advantage of PEANuT: it does not require extensive manual tuning on which parts (modules, layers) of the foundation model PEANuT should be applied. Consequently, PEANuT can be easily incorporated to a wide range of scenarios.

# 7. Conclusion

In this work, we propose PEANuT, a novel parameter-efficient fine-tuning (PEFT) method that introduces nonlinear transformations to enhance model adaptation while maintaining efficiency. By incorporating a lightweight neural network that models cumulative weight updates as functions of the pre-trained weights, PEANuT effectively captures complex, nonlinear structures in the weight space, allowing for more expressive and accurate adaptation to downstream tasks. Our theoretical analysis supports the efficacy of PEANuT, demonstrating that it can achieve greater or equivalent expressiveness compared to existing LoRA, a popular and state-of-the-art PEFT method, with fewer number of parameters. Through extensive experiments on four benchmarks encompassing over twenty datasets with various pre-trained backbones, PEANuT demonstrated superior performance on both NLP and vision tasks compared to existing state-of-the-art methods.

# References

AI@Meta. Llama 3 model card. 2024. URL https://github.com/meta-llama/llama3/blob/main/MODEL_CARD.md.

Alan Ansell, Edoardo Maria Ponti, Anna Korhonen, and Ivan Vulić. Composable sparse fine-tuning for cross-lingual transfer. *arXiv preprint arXiv:2110.07560*, 2021.

Tom M Apostol. Modular functions and dirichlet series in number theory. 1990.

Srinadh Bhojanapalli, Ayan Chakrabarti, Daniel Glasner, Daliang Li, Thomas Unterthiner, and Andreas Veit. Understanding robustness of transformers for image classification. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10231–10241, 2021.

Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. Piqa: Reasoning about physical commonsense in natural language. *arXiv preprint arXiv:1911.11641*, 2020.

Kerim Büyükakyüz. Olora: Orthonormal low-rank adaptation of large language models. *arXiv preprint arXiv:2406.01775*, 2024.

Chun-Fu Richard Chen, Quanfu Fan, and Rameswar Panda. Crossvit: Cross-attention multi-scale vision transformer for image classification. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 357–366, 2021.

Jiaao Chen, Aston Zhang, Xingjian Shi, Mu Li, Alex Smola, and Diyi Yang. Parameter-efficient fine-tuning design spaces. *arXiv preprint arXiv:2301.01821*, 2023.

Wei Chen, Zichen Miao, and Qiang Qiu. Large convolutional model tuning via filter subspace. *arXiv preprint arXiv:2403.00269*, 2024.

Wei Chen, Jingxi Yu, Zichen Miao, and Qiang Qiu. Sparse fine-tuning of transformers for generative tasks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 18703–18713, 2025.

Gong Cheng, Junwei Han, and Xiaoqiang Lu. Remote sensing image scene classification: Benchmark and state of the art. *Proceedings of the IEEE*, 105(10):1865–1883, 2017.

Alexandra Chronopoulou, Matthew E Peters, Alexander Fraser, and Jesse Dodge. Adaptersoup: Weight averaging to improve generalization of pretrained language models. *arXiv preprint arXiv:2302.07027*, 2023.

Mircea Cimpoi, Subhransu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. Describing textures in the wild. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3606–3613, 2014.

Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. Boolq: Exploring the surprising difficulty of natural yes/no questions. *arXiv preprint arXiv:1905.10044*, 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.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.

Sarkar Snigdha Sarathi Das, Ranran Haoran Zhang, Peng Shi, Wenpeng Yin, and Rui Zhang. Unified low-resource sequence labeling by sample-aware dynamic sparse finetuning. *arXiv preprint arXiv:2311.03748*, 2023.

Jacob Devlin. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.

Ning Ding, Yujia Qin, Guang Yang, Fuchao Wei, Zonghan Yang, Yusheng Su, Shengding Hu, Yulin Chen, Chi-Min Chan, Weize Chen, et al. Parameter-efficient fine-tuning of large-scale pre-trained language models. *Nature Machine Intelligence*, 5(3):220–235, 2023.

Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2020.

Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. URL https://openreview.net/forum?id=YicbFdNTTy.

Ali Edalati, Marzieh Tahaei, Ivan Kobyzev, Vahid Partovi Nia, James J Clark, and Mehdi Rezagholizadeh. Krona: Parameter efficient tuning with kronecker adapter. *arXiv preprint arXiv:2212.10650*, 2022.

Ziqi Gao, Qichao Wang, Aochuan Chen, Zijing Liu, Bingzhe Wu, Liang Chen, and Jia Li. Parameter-efficient fine-tuning with discrete fourier transform. In *Forty-first International Conference on Machine Learning*.

Ziqi Gao, Qichao Wang, Aochuan Chen, Zijing Liu, Bingzhe Wu, Liang Chen, and Jia Li. Parameter-efficient fine-tuning with discrete fourier transform. *arXiv preprint arXiv:2405.03003*, 2024.

Michael S Gashler and Stephen C Ashmore. Training deep fourier neural networks to fit time-series data. In *Intelligent Computing in Bioinformatics: 10th International Conference, ICIC 2014, Taiyuan, China, August 3-6, 2014. Proceedings 10*, pages 48–55. Springer, 2014.

Xavier Glorot, Antoine Bordes, and Yoshua Bengio. Deep sparse rectifier neural networks. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages 315–323. JMLR Workshop and Conference Proceedings, 2011.

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

Zeyu Han, Chao Gao, Jinyang Liu, Sai Qian Zhang, et al. Parameter-efficient fine-tuning for large models: A comprehensive survey. *arXiv preprint arXiv:2403.14608*, 2024.

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

Patrick Helber, Benjamin Bischke, Andreas Dengel, and Damian Borth. Eurosat: A novel dataset and deep learning benchmark for land use and land cover classification. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 12(7):2217–2226, 2019.

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

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*, pages 2790–2799. PMLR, 2019.

Jeremy Howard and Sebastian Ruder. Universal language model fine-tuning for text classification. *arXiv preprint arXiv:1801.06146*, 2018.

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

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

Zhiqiang Hu, Lei Wang, Yihuai Lan, Wanyu Xu, Ee-Peng Lim, Lidong Bing, Xing Xu, Soujanya Poria, and Roy Lee. Llm-adapters: An adapter family for parameter-efficient fine-tuning of large language models. *arXiv preprint arXiv:2304.01933*, 2023.

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

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

Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *Proceedings of the IEEE international conference on computer vision workshops*, pages 554–561, 2013.

A Krizhevsky. Learning multiple layers of features from tiny images. *Master's thesis, University of Tront*, 2009.

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

Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in neural information processing systems*, 33:9459–9474, 2020.

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

Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Wei-Ming Chen, Wei-Chen Wang, Guangxuan Xiao, Xingyu Dang, Chuang Gan, and Song Han. Awq: Activation-aware weight quantization for on-device llm compression and acceleration. *Proceedings of Machine Learning and Systems*, 6:87–100, 2024.

Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Mohta, Tenghao Huang, Mohit Bansal, and Colin A Raffel. Few-shot parameter-efficient fine-tuning is better and cheaper than in-context learning. *Advances in Neural Information Processing Systems*, 35:1950–1965, 2022.

Shih-Yang Liu, Chien-Yi Wang, Hongxu Yin, Pavlo Molchanov, Yu-Chiang Frank Wang, Kwang-Ting Cheng, and Min-Hung Chen. DoRA: Weight-decomposed low-rank adaptation. *arXiv:2402.09353*, 2024. URL https://arxiv.org/abs/2402.09353.

Tianci Liu, Haoxiang Jiang, Tianze Wang, Ran Xu, Yue Yu, Linjun Zhang, Tuo Zhao, and Haoyu Wang. Roserag: Robust retrieval-augmented generation with small-scale llms via margin-aware preference optimization. *arXiv preprint arXiv:2502.10993*, 2025.

Yinhan Liu. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.

Subhransu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. Fine-grained visual classification of aircraft. *arXiv preprint arXiv:1306.5151*, 2013.

Yuning Mao, Lambert Mathias, Rui Hou, Amjad Almahairi, Hao Ma, Jiawei Han, Wen-tau Yih, and Madian Khabsa. Unipelt: A unified framework for parameter-efficient language model tuning. *arXiv preprint arXiv:2110.07577*, 2021.

Fanxu Meng, Zhaohui Wang, and Muhan Zhang. Pissa: Principal singular values and singular vectors adaptation of large language models. *arXiv preprint arXiv:2404.02948*, 2024.

Zichen Miao, Wei Chen, and Qiang Qiu. Coeff-tuning: A graph filter subspace view for tuning attention-based large models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 20146–20157, 2025.

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.

Rui Pan, Xiang Liu, Shizhe Diao, Renjie Pi, Jipeng Zhang, Chi Han, and Tong Zhang. Lisa: Layer-wise importance sampling for memory-efficient large language model fine-tuning. *arXiv preprint arXiv:2403.17919*, 2024.

Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, and CV Jawahar. Cats and dogs. In *2012 IEEE conference on computer vision and pattern recognition*, pages 3498–3505. IEEE, 2012.

Ruiyang Qin, Dancheng Liu, Zheyu Yan, Zhaoxuan Tan, Zixuan Pan, Zhenge Jia, Meng Jiang, Ahmed Abbasi, Jinjun Xiong, and Yiyu Shi. Empirical guidelines for deploying llms onto resource-constrained edge devices. *arXiv preprint arXiv:2406.03777*, 2024.

Maithra Raghu, Ben Poole, Jon Kleinberg, Surya Ganguli, and Jascha Sohl-Dickstein. On the expressive power of deep neural networks. In *international conference on machine learning*, pages 2847–2854. PMLR, 2017.

Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale. *arXiv preprint arXiv:1907.10641*, 2019.

Maarten Sap, Hannah Rashkin, Derek Chen, Ronan LeBras, and Yejin Choi. Socialiqa: Commonsense reasoning about social interactions. *arXiv preprint arXiv:1904.09728*, 2019.

Lin Song, Yukang Chen, Shuai Yang, Xiaohan Ding, Yixiao Ge, Ying-Cong Chen, and Ying Shan. Low-rank approximation for sparse attention in multi-modal llms. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13763–13773, 2024.

Xiaofei Sun, Xiaoya Li, Jiwei Li, Fei Wu, Shangwei Guo, Tianwei Zhang, and Guoyin Wang. Text classification via large language models. *arXiv preprint arXiv:2305.08377*, 2023.

Yi-Lin Sung, Varun Nair, and Colin A Raffel. Training neural networks with fixed sparse masks. *Advances in Neural Information Processing Systems*, 34:24193–24205, 2021.

Damien Teney, Armand Mihai Nicolicioiu, Valentin Hartmann, and Ehsan Abbasnejad. Neural redshift: Random networks are not random functions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4786–4796, 2024.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, and et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 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. *arXiv preprint arXiv:2210.07558*, 2022.

A Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017.

Tu Vu, Brian Lester, Noah Constant, Rami Al-Rfou, and Daniel Cer. Spot: Better frozen model adaptation through soft prompt transfer. *arXiv preprint arXiv:2110.07904*, 2021.

Danilo Vucetic, Mohammadreza Tayaranian, Maryam Ziaeefard, James J Clark, Brett H Meyer, and Warren J Gross. Efficient fine-tuning of bert models on the edge. In *2022 IEEE International Symposium on Circuits and Systems (ISCAS)*, pages 1838–1842, 2022.

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. *arXiv preprint arXiv:1804.07461*, 2018.

Hanqing Wang, Zeguan Xiao, Yixia Li, Shuo Wang, Guanhua Chen, and Yun Chen. Milora: Harnessing minor singular components for parameter-efficient llm finetuning. *arXiv preprint arXiv:2406.09044*, 2024a.

Haoyu Wang, Fenglong Ma, Yaqing Wang, and Jing Gao. Knowledge-guided paraphrase identification. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 843–853, 2021.

Haoyu Wang, Handong Zhao, Yaqing Wang, Tong Yu, Jiuxiang Gu, and Jing Gao. Fedkc: Federated knowledge composition for multilingual natural language understanding. In *Proceedings of the ACM Web Conference 2022*, pages 1839–1850, 2022.

Haoyu Wang, Ruirui Li, Haoming Jiang, Jinjin Tian, Zhengyang Wang, Chen Luo, Xianfeng Tang, Monica Xiao Cheng, Tuo Zhao, and Jing Gao. Blendfilter: Advancing retrieval-augmented large language models via query generation blending and knowledge filtering. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 1009–1025, 2024b.

Haoyu Wang, Tianci Liu, Tuo Zhao, and Jing Gao. Roselora: Row and column-wise sparse low-rank adaptation of pre-trained language model for knowledge editing and fine-tuning. *arXiv preprint arXiv:2406.10777*, 2024c.

Yihan Wang, Jatin Chauhan, Wei Wang, and Cho-Jui Hsieh. Universality and limitations of prompt tuning. *Advances in Neural Information Processing Systems*, 36, 2024d.

Muling Wu, Wenhao Liu, Xiaohua Wang, Tianlong Li, Changze Lv, Zixuan Ling, Jianhao Zhu, Cenyuan Zhang, Xiaoqing Zheng, and Xuanjing Huang. Advancing parameter efficiency in fine-tuning via representation editing. *arXiv:2402.15179*, 2024a. URL https://arxiv.org/abs/2402.15179.

Shanchan Wu and Yifan He. Enriching pre-trained language model with entity information for relation classification. In *Proceedings of the 28th ACM international conference on information and knowledge management*, pages 2361–2364, 2019.

Zhengxuan Wu, Aryaman Arora, Zheng Wang, Atticus Geiger, Dan Jurafsky, Christopher D. Manning, and Christopher Potts. ReFT: Representation finetuning for language models. 2024b. URL arxiv.org/abs/2404.03592.

Ran Xu, Wenqi Shi, Yuchen Zhuang, Yue Yu, Joyce C Ho, Haoyu Wang, and Carl Yang. Collab-rag: Boosting retrieval-augmented generation for complex question answering via white-box and black-box llm collaboration. *arXiv preprint arXiv:2504.04915*, 2025.

An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.

Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. In *The eleventh international conference on learning representations*, 2022.

Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T. Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. Metamath: Bootstrap your own mathematical questions for large language models. *arXiv preprint arXiv:2309.12284*, 2023.

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

Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence? *arXiv preprint arXiv:1905.07830*, 2019.

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

Yuanhan Zhang, Kaiyang Zhou, and Ziwei Liu. Neural prompt search. *arXiv preprint arXiv:2206.04673*, 2022.

Hongyu Zhao, Hao Tan, and Hongyuan Mei. Tiny-attention adapter: Contexts are more important than the number of parameters. *arXiv preprint arXiv:2211.01979*, 2022.

Jiawei Zhao, Zhenyu Zhang, Beidi Chen, Zhangyang Wang, Anima Anandkumar, and Yuandong Tian. Galore: Memory-efficient llm training by gradient low-rank projection. *arXiv preprint arXiv:2403.03507*, 2024.

Han Zhou, Xingchen Wan, Ivan Vulić, and Anna Korhonen. Autopeft: Automatic configuration search for parameter-efficient fine-tuning. *Transactions of the Association for Computational Linguistics*, 12: 525–542, 2024.

# A. Appendix

# B. Details of Theoretical Results

In this section, we provide the proof of Proposition 3.2 and introduce additional theoretical results when we assume sinusoid activation.

### B.1. Proof of Proposition 3.2

The intuition behind the proof is that we can always restore an identity function using two ReLU activation functions, i.e., $x = \sigma(x) - \sigma(-x)$ for any $x \in \mathbb{R}$

*Proof of Proposition 3.2.* We first show that

$$\min_{\boldsymbol{\Theta}_1 \in \mathbb{R}^{d_2 \times 2r}, \boldsymbol{\Theta}_2 \in \mathbb{R}^{2r \times d_2}} \mathcal{L}(\mathcal{D}_{\text{train}}; \mathbf{W}^0 + f(\mathbf{W}^0; (\boldsymbol{\Theta}_1, \boldsymbol{\Theta}_2)))$$
$$\leq \min_{\mathbf{A} \in \mathbb{R}^{d_1 \times r}, \mathbf{B} \in \mathbb{R}^{r \times d_2}} \mathcal{L}(\mathcal{D}_{\text{train}}; \mathbf{W}^0 + \mathbf{AB}).$$

Let $(\mathbf{A}^*, \mathbf{B}^*) = \arg\min_{\mathbf{A} \in \mathbb{R}^{d_1 \times r}, \mathbf{B} \in \mathbb{R}^{r \times d_2}} \mathcal{L}(\mathcal{D}_{\text{train}}; \mathbf{W}^0 + \mathbf{AB})$. Take $\boldsymbol{\Theta}_1^{\#} := [(\mathbf{W}^0)^{\dagger} \mathbf{A}^*; -(\mathbf{W}^0)^{\dagger} \mathbf{A}^*] \in \mathbb{R}^{d_2 \times 2r}$ and $\boldsymbol{\Theta}_2^{\#} := [\mathbf{B}^{*\top}; -\mathbf{B}^{*\top}]^{\top} \in \mathbb{R}^{2r \times d_2}$, where $(\mathbf{W}^0)^{\dagger} \in \mathbb{R}^{d_2 \times d_1}$ is the Moore-Penrose inverse of $\mathbf{W}^0$. Then, since $\sigma$ is a ReLU activation function,

$$f(\mathbf{W}^0; (\boldsymbol{\Theta}_1^{\#}, \boldsymbol{\Theta}_2^{\#}))$$
$$= \sigma(\mathbf{W}^0 \boldsymbol{\Theta}_1^{\#}) \boldsymbol{\Theta}_2^{\#}$$
$$= \sigma(\mathbf{W}^0 (\mathbf{W}^0)^{\dagger} \mathbf{A}^*) \mathbf{B}^* - \sigma(-\mathbf{W}^0 (\mathbf{W}^0)^{\dagger} \mathbf{A}^*) \mathbf{B}^*$$
$$= \mathbf{W}^0 (\mathbf{W}^0)^{\dagger} \mathbf{A}^* \mathbf{B}^*.$$

Note that $\mathbf{W}^0 (\mathbf{W}^0)^{\dagger} = \mathbf{U}^0 \mathbf{U}^{0\top}$ is the projection to the left singular space of $\mathbf{W}^0$. Hence

$$\mathcal{L}(\mathcal{D}_{\text{train}}; \mathbf{W}^0 + f(\mathbf{W}^0; (\boldsymbol{\Theta}_1^{\#}, \boldsymbol{\Theta}_2^{\#})))$$
$$= \mathcal{L}(\mathcal{D}_{\text{train}}; \mathbf{U}^0 \mathbf{U}^{0\top} \mathbf{W}^0 + \mathbf{U}^0 \mathbf{U}^{0\top} \mathbf{A}^* \mathbf{B}^*)$$
$$= \mathcal{L}(\mathcal{D}_{\text{train}}; \mathbf{W}^0 + \mathbf{A}^* \mathbf{B}^*),$$

where the last equality follows from the invariance assumption. This gives the first inequality:

$$\min_{\boldsymbol{\Theta}_1 \in \mathbb{R}^{d_2 \times 2r}, \boldsymbol{\Theta}_2 \in \mathbb{R}^{2r \times d_2}} \mathcal{L}(\mathcal{D}_{\text{train}}; \mathbf{W}^0 + f(\mathbf{W}^0; (\boldsymbol{\Theta}_1, \boldsymbol{\Theta}_2)))$$
$$\leq \mathcal{L}(\mathcal{D}_{\text{train}}; \mathbf{W}^0 + f(\mathbf{W}^0; (\boldsymbol{\Theta}_1^{\#}, \boldsymbol{\Theta}_2^{\#})))$$
$$= \mathcal{L}(\mathcal{D}_{\text{train}}; \mathbf{W}^0 + \mathbf{A}^* \mathbf{B}^*)$$
$$= \min_{\mathbf{A} \in \mathbb{R}^{d_1 \times r}, \mathbf{B} \in \mathbb{R}^{r \times d_2}} \mathcal{L}(\mathcal{D}_{\text{train}}; \mathbf{W}^0 + \mathbf{AB}).$$

We next show the following inequality:

$$\min_{\mathbf{A} \in \mathbb{R}^{d_1 \times r}, \mathbf{B} \in \mathbb{R}^{r \times d_2}} \mathcal{L}(\mathcal{D}_{\text{train}}; \mathbf{W}^0 + \mathbf{AB})$$
$$\leq \min_{\boldsymbol{\Theta}_1 \in \mathbb{R}^{d_2 \times r}, \boldsymbol{\Theta}_2 \in \mathbb{R}^{r \times d_2}} \mathcal{L}(\mathcal{D}_{\text{train}}; \mathbf{W}^0 + f(\mathbf{W}^0; (\boldsymbol{\Theta}_1, \boldsymbol{\Theta}_2))).$$

Take $\mathbf{A}^{\#} = \sigma(\mathbf{W}^0 \boldsymbol{\Theta}_1^*) \in \mathbb{R}^{d_1 \times r}$ and $\mathbf{B}^{\#} = \boldsymbol{\Theta}_2^* \in \mathbb{R}^{r \times d_2}$, where $(\boldsymbol{\Theta}_1^*, \boldsymbol{\Theta}_2^*) = \arg\min_{\boldsymbol{\Theta}_1 \in \mathbb{R}^{d_2 \times r}, \boldsymbol{\Theta}_2 \in \mathbb{R}^{r \times d_1}} \mathcal{L}(\mathcal{D}_{\text{train}}; \mathbf{W}^0 + f(\mathbf{W}^0; (\boldsymbol{\Theta}_1, \boldsymbol{\Theta}_2)))$. The conclusion follows from

$$\min_{\mathbf{A} \in \mathbb{R}^{d_1 \times r}, \mathbf{B} \in \mathbb{R}^{r \times d_2}} \mathcal{L}(\mathcal{D}_{\text{train}}; \mathbf{W}^0 + \mathbf{A}\mathbf{B})$$
$$\leq \mathcal{L}(\mathcal{D}_{\text{train}}; \mathbf{W}^0 + \mathbf{A}^{\#}\mathbf{B}^{\#})$$
$$= \mathcal{L}(\mathcal{D}_{\text{train}}; \mathbf{W}^0 + \sigma(\mathbf{W}^0 \boldsymbol{\Theta}_1^*)\boldsymbol{\Theta}_2^*)$$
$$= \min_{\boldsymbol{\Theta}_1 \in \mathbb{R}^{d_2 \times r}, \boldsymbol{\Theta}_2 \in \mathbb{R}^{r \times d_1}} \mathcal{L}(\mathcal{D}_{\text{train}}; \mathbf{W}^0 + f(\mathbf{W}^0; (\boldsymbol{\Theta}_1, \boldsymbol{\Theta}_2))).$$

$\square$

## B.2. Theoretical Analysis of PEANuT under sinusoid activation function

Here we consider a sinusoid activation function $\sigma_{\text{p}}(x) = \sin(2\pi x)$ (Gashler and Ashmore, 2014) and design $f(\mathbf{W}^0; \boldsymbol{\theta}) = \sigma_{\text{p}}(\mathbf{W}^0 \boldsymbol{\Theta}_1)\boldsymbol{\Theta}_2$ with $\boldsymbol{\theta} = (\boldsymbol{\Theta}_1, \boldsymbol{\Theta}_2)$. With this periodic activation function, we can show a stronger result that PEANuT has expressivity (almost) greater than or equal to a LoRA with more parameters when $d_1 \gg d_2$.

**Proposition B.1** (Expressivity of PEANuT with Sine Activation). *Suppose that there exists a row of $\mathbf{W}^0$, whose entries are linearly independent over the rationals. Then, for any $r > 0$, $\mathbf{A} \in \mathbb{R}^{d_1 \times r}$ and $\mathbf{B} \in \mathbb{R}^{r \times d_2}$, and $\epsilon > 0$, there exists some $\boldsymbol{\Theta}_1^* \in \mathbb{R}^{d_2 \times r}$ and $\boldsymbol{\Theta}_2^* \in \mathbb{R}^{r \times d_2}$ such that*

$$\|\mathbf{A}\mathbf{B} - \sigma_{\text{p}}(\mathbf{W}^0 \boldsymbol{\Theta}_1^*)\boldsymbol{\Theta}_2^*\|_{\text{F}} \leq \epsilon.$$

Proposition B.1 shows that the class of updates $\Delta\mathbf{W} = \sigma_{\text{p}}(\mathbf{W}^0 \boldsymbol{\Theta}_1)\boldsymbol{\Theta}_2$ by PEANuT with $2rd_2$ parameters is dense in the class of updates $\Delta\mathbf{W} = \mathbf{A}\mathbf{B}$ by LoRA with $r(d_1 + d_2)$ parameters. When $d_2 \ll d_1$, this shows better parameter efficiency of PEANuT.

Examining the proof of Proposition B.1, it is straightforward to show that the result holds for any continuous and periodic activation function whose range contains an open interval centered at 0.

*Proof.* This proof relies on Kronecker's theorem (Theorem 7.9 in Apostol (1990)) from number theory, which shows that for all $j \in \mathbb{R}^q$, the fractional parts of $(ct_1, ct_2, \ldots, ct_q)^{\top}$ is dense in $[0,1]^q$ over $c \in \mathbb{R}$, as long as $t_1, \ldots, t_q$ are linearly independent over the rationals.

Let $\mathbf{W}_{j^*}$ be the $j^*$-th column of $\mathbf{W}^0$ whose entries are linearly independent over the rationals. Since $\mathbf{A}\mathbf{B}$ has a scale ambiguity, we can assume that $\mathbf{A}$ is a matrix whose entries are bounded by 1 without loss of generality. Write $\mathbf{A} = (\mathbf{A}_1, \mathbf{A}_2, \ldots, \mathbf{A}_r)$.

Take $\epsilon' > 0$ whose value will be determined later. From Kronecker's theorem, for each $\mathbf{A}_j$ there exists some $c_j \in \mathbb{R}$ such that

$$\left| \{c_j \mathbf{W}_{j^*}\} - \frac{\arcsin(\mathbf{A}_j)}{2\pi} \right| \leq \epsilon',$$

where $\{\mathbf{B}\}$ is a vector whose entries are the fractional part of the corresponding entry of $\mathbf{B}$, and arcsin is applied elementwisely.

Let $\boldsymbol{\Theta}_1^* = (c_1 \boldsymbol{e}_{j^*}, c_2 \boldsymbol{e}_{j^*}, \ldots, c_r \boldsymbol{e}_{j^*})$, where $\boldsymbol{e}_{j^*}$ is the $j^*$-th standard basis vector in $\mathbb{R}^{d_2}$. Using the fact that $2\pi \{c_j \boldsymbol{W}_{j^*}\} = 2\pi c_j \boldsymbol{W}_{j^*} \mod 2\pi$, we have

$$
\begin{aligned}
& \left\| \sigma_{\mathrm{p}}(\boldsymbol{W}^0 \boldsymbol{\Theta}_1^*) - \boldsymbol{A} \right\|_{\mathrm{F}}^2 \\
&= \left\| \sigma_{\mathrm{p}}((c_1 \boldsymbol{W}_{j^*}, c_2 \boldsymbol{W}_{j^*}, \ldots c_r \boldsymbol{W}_{j^*})) - \boldsymbol{A} \right\|_{\mathrm{F}}^2 \\
&\leq \sum_j \left\| \sin(2\pi c_j \boldsymbol{W}_{j^*}) - \boldsymbol{A}_j \right\|^2 \leq 4\pi^2 r \epsilon'^2,
\end{aligned}
\tag{2}
$$

where the last inequality follows from equation 2 and the fact that $\sin(x)$ is Lipschitz continuous with Lipschitz constant 1. Hence by choosing $\boldsymbol{\Theta}_2^* \leftarrow \boldsymbol{B}$, we have

$$
\begin{aligned}
& \left\| \boldsymbol{A}\boldsymbol{B} - \sigma_{\mathrm{p}}(\boldsymbol{W}^0 \boldsymbol{\Theta}_1^*) \boldsymbol{\Theta}_2^* \right\|_{\mathrm{F}}^2 \\
&\leq \|\boldsymbol{B}\|^2 \left\| \sigma_{\mathrm{p}}(\boldsymbol{W}^0 \boldsymbol{\Theta}_1^*) - \boldsymbol{A} \right\|_{\mathrm{F}}^2 \\
&\leq 4\pi^2 \|\boldsymbol{B}\|^2 r \epsilon'^2.
\end{aligned}
$$

Choose $\epsilon' = \epsilon / (2\pi \sqrt{r} \|\boldsymbol{B}\|)$, then the proof is complete. $\qquad \square$

## C. Hyperparameters

We provide the specific hyperparameters used in our experiments to ensure reproducibility. For most of our experiments, we use the standard implementation of PEANuT, which we refer to as vanilla PEANuT. The neural network architecture in vanilla PEANuT consists of only two layers: an input layer and an output layer. We select this approach because vanilla PEANuT offers the benefits of simplicity in implementation, a low parameter count, and sufficient adaptation power. Nonetheless, we dedicate Section 5.5 to exploring more complex adaptation networks and their effect on performance.

### C.1. Image Classification

Hyperparameters for PEANuT for Fig. 5 are provided in Table 8. We tune the classification head and the backbone separately and provide detailed settings for each dataset. All weight decay values are not tuned and follow the settings from Gao et al. (2024). The scaling factor $s$ is set to 1.0. The hidden layer dimension $r$ for MHSA is set to 7 in the QV-setting, while both hidden layer dimensions for MHSA and MLP are set to 2 in the QV-MLP-setting described in Section 6.

### C.2. Natural Language Understanding

We provide used hyper-parameters for PEANuT in natural language understanding on the GLUE benchmark in Table 11 and Table 12. The reported results are obtained when using a depth of 6 for PEANuT. The learning rates for the head and the backbone are tuned separately. The scaling factor $s$ is searched in $\{0.01, 0.1, 1.0\}$. For reproducibility, we fix the seed as 0. The hidden layer dimension $r$ is set to 8 in PEANuT-L and 1 in PEANuT-S. More specifically, we apply PEANuT to all layers in RoBERTa-base for PEANuT-L, while only applying PEANuT to layers $\{4, 5, 6, 7, 8, 9, 10, 11\}$ for PEANuT-S to reduce the number of trainable parameters. The seed is fixed for reproducibility.

**Table 8:** Hyperparameter of image classification for PEANuT.

| Hyperparameter | OxfordPets | StanfordCars | CIFAR10 | DTD | EuroSAT | FGVC | RESISC45 | CIFAR100 |
|---|---|---|---|---|---|---|---|---|
| Epochs | | | | 10 | | | | |
| Optimizer | | | | AdamW | | | | |
| LR Schedule | | | | Linear | | | | |
| Weight Decay | 8E-4 | 4E-5 | 9E-5 | 7E-5 | 3E-4 | 7E-5 | 3E-4 | 1E-4 |
| QV | | | | | | | | |
| Learning Rate (PEANuT) | 5E-3 | 1E-2 | 5E-3 | 1E-2 | 5E-3 | 1E-2 | 5E-3 | 5E-3 |
| Learning Rate (Head) | 5E-3 | 1E-2 | 5E-3 | 1E-2 | 5E-3 | 1E-2 | 1E-2 | 5E-3 |
| QV-MLP | | | | | | | | |
| Learning Rate (PEANuT) | 5E-3 | 5E-3 | 5E-3 | 1E-2 | 5E-3 | 5E-3 | 1E-2 | 5E-3 |
| Learning Rate (Head) | 5E-3 | 1E-2 | 5E-3 | 1E-2 | 5E-3 | 1E-2 | 1E-2 | 5E-3 |

**Table 9:** Hyperparameter of commonsense reasoning for PEANuT.

| Hyperparameter | Commonsense Reasoning |
|---|---|
| Hidden Layer Dimension | 32 |
| $\alpha$ | 32 |
| Dropout | 0.05 |
| Optimizer | Adam W |
| Learning Rate | 3e-4 |
| Batch Size | 16 |
| Warmup Steps | 100 |
| Epochs | 1 |

### C.3. Commonsense Reasoning

We provide hyperparameters settings of PEANuT for commonsense reasoning task in Table 9. We follow the hyperparameters settings in MiLoRA (Wang et al., 2024a). We limit all samples to a maximum of 256 tokens. For evaluation, we set a maximum token number of 32.

### C.4. Arithmetic Reasoning

We provide hyperparameters settings of PEANuT for arithmetic reasoning task in Table 10. We follow the hyper-parameters settings in MiLoRA (Wang et al., 2024a). We limit all samples to a maximum of 2048 tokens. For evaluation, we set a maximum token number of 256 on GSM8K (Cobbe et al., 2021) dataset. On MATH (Hendrycks et al., 2021), we set the maximum new token to 512.

## D. Datasets

In this section, we provide a detailed description of the datasets used in our experiments.

**Table 10:** Hyperparameter of arithmetic reasoning for PEANuT.

| Hyperparameter | Arithmetic Reasoning |
|---|---|
| Hidden Layer Dimension | 64 |
| $\alpha$ | 64 |
| Dropout | 0.05 |
| Optimizer | Adam W |
| Learning Rate | 3e-4 |
| Batch Size | 16 |
| Warmup Steps | 100 |
| Epochs | 3 |

**Table 11:** Hyperparameter of GLUE benchmark for PEANuT-L.

| Hyperparameter | STS-B | RTE | MRPC | CoLA | SST-2 | QNLI | MNLI | QQP |
|---|---|---|---|---|---|---|---|---|
| Optimizer | AdamW | | | | | | | |
| LR Schedule | Linear | | | | | | | |
| Learning Rate (PEANuT) | 5E-3 | 5E-3 | 5E-3 | 1E-3 | 5E-3 | 1E-3 | 5E-3 | 5E-3 |
| Learning Rate (Head) | 5E-3 | 5E-3 | 5E-3 | 1E-3 | 5E-3 | 1E-3 | 5E-3 | 5E-3 |
| Scaling | 0.1 | 0.01 | 0.01 | 0.1 | 0.01 | 0.01 | 0.01 | 0.01 |
| Max Seq. Len | 512 | 512 | 512 | 512 | 512 | 512 | 512 | 512 |
| Batch Size | 64 | 32 | 64 | 64 | 32 | 32 | 32 | 64 |

### D.1. Image Classification

For image classification, we provide detailed information about the used datasets in Table 8.

### D.2. Natural Language Understanding

The GLUE benchmark comprises 8 NLP datasets: MNLI, SST-2, MRPC, CoLA, QNLI, QQP, RTE, and STS-B, covering tasks such as inference, sentiment analysis, paraphrase detection, linguistic acceptability, question-answering, and textual similarity. We provide detailed information about them in Table 14.

### D.3. Commonsense Reasoning

For commonsense reasoning task, we use 8 datasets, including BoolQ, PIQA, SIQA, HellaSwag, Wino-Grande, ARC-e, ARC-c and OBQA. The detailed information is provided in Table 15.

### D.4. Arithmetic Reasoning

Detailed information for arithmetic reasoning task is provided in Table 16. GSM8K consists of high quality grade school math problems, typically free-form answers. MATH includes classifications from multiple mathematical domains, such as algebra, counting_and_probability, geometry, intermediate_algebra, number_theory, prealgebra and precalculus.

**Table 12:** Hyperparameter of GLUE benchmark for PEANuT-S.

| Hyperparameter | STS-B | RTE | MRPC | CoLA | SST-2 | QNLI | MNLI | QQP |
|---|---|---|---|---|---|---|---|---|
| Optimizer | | | | AdamW | | | | |
| LR Schedule | | | | Linear | | | | |
| Learning Rate (PEANuT) | 5E-3 | 1E-3 | 5E-3 | 5E-3 | 5E-3 | 1E-3 | 5E-3 | 1E-3 |
| Learning Rate (Head) | 1E-3 | 1E-3 | 5E-3 | 1E-3 | 5E-3 | 1E-3 | 5E-3 | 1E-3 |
| Scaling | 0.1 | 1.0 | 0.01 | 0.1 | 0.01 | 0.1 | 0.01 | 1.0 |
| Max Seq. Len | 512 | 512 | 512 | 512 | 512 | 512 | 512 | 512 |
| Batch Size | 64 | 32 | 64 | 64 | 32 | 32 | 32 | 64 |

**Table 13:** Detailed information of image classification tasks.

| Dataset | #Class | #Train | #Val | #Test | Rescaled resolution |
|---|---|---|---|---|---|
| OxfordPets | 37 | 3,312 | 368 | 3,669 | |
| StandfordCars | 196 | 7,329 | 815 | 8,041 | |
| CIFAR10 | 10 | 45,000 | 5,000 | 10,000 | |
| DTD | 47 | 4,060 | 452 | 1,128 | $224 \times 224$ |
| EuroSAT | 10 | 16,200 | 5,400 | 5,400 | |
| FGVC | 100 | 3,000 | 334 | 3,333 | |
| RESISC45 | 45 | 18,900 | 6,300 | 6,300 | |
| CIFAR100 | 100 | 45,000 | 5,000 | 10,000 | |

**Table 14:** Detailed information of the GLUE benchmark. STS-B is a regression task, while all other tasks are either single-sentence or sentence-pair classification tasks.

| Corpus | Task | Metrics | # Train | # Val | # Test | # Labels |
|--------|------|---------|---------|-------|--------|----------|
| Single-Sentence Tasks | | | | | | |
| CoLA | Acceptability | Matthews Corr. | 8.55k | 1.04k | 1.06k | 2 |
| SST-2 | Sentiment | Accuracy | 67.3k | 872 | 1.82k | 2 |
| Similarity and Paraphrase Tasks | | | | | | |
| MRPC | Paraphrase | Accuracy/F1 | 3.67k | 408 | 1.73k | 2 |
| STS-B | Sentence similarity | Pearson/Spearman Corr. | 5.75k | 1.5k | 1.38k | 1 |
| QQP | Paraphrase | Accuracy/F1 | 364k | 40.4k | 391k | 2 |
| Inference Tasks | | | | | | |
| MNLI | NLI | Accuracy | 393k | 19.65k | 19.65k | 3 |
| QNLI | QA/NLI | Accuracy | 105k | 5.46k | 5.46k | 2 |
| RTE | NLI | Accuracy | 2.49k | 277 | 3k | 2 |

**Table 15:** Detailed information of commonsense reasoning task.

| Dataset | #Class | #Train | #Dev | #Test |
|---------|--------|--------|------|-------|
| BoolQ | Binary classification | 9,427 | 3,270 | 3,245 |
| PIQA | Binary classification | 16,113 | 1,838 | 3,000 |
| SIQA | Ternary classification | 33,410 | 1,954 | 2,224 |
| HellaSwag | Quaternary classification | 39,905 | 10,042 | 10,003 |
| WinoGrande | Binary classification | 40,398 | 1,267 | 1,767 |
| ARC-e | Quaternary classification | 2,251 | 570 | 2,376 |
| ARC-c | Quaternary classification | 1,119 | 229 | 1,172 |
| OBQA | Quaternary classification | 4,957 | 500 | 500 |

**Table 16:** Detailed information of arithmetic reasoning task.

| Dataset | #Train | #Dev | #Test |
|---------|--------|------|-------|
| GSM8K | 7,473 | 1,319 | 1,319 |
| MATH | 12,500 | 500 | 5,000 |