

Rethinking Efficient Tuning Methods from a Unified Perspective

Zeyinzi Jiang¹ Chaojie Mao¹ Ziyuan Huang² Yiliang Lv¹ Deli Zhao¹ Jingren Zhou¹

Abstract

Parameter-efficient transfer learning (PETL) based on large-scale pre-trained foundation models has achieved great success in various downstream applications. Existing tuning methods, such as prompt, prefix, and adapter, perform task-specific lightweight adjustments to different parts of the original architecture. However, they take effect on only some parts of the pre-trained models, *i.e.*, only the feed-forward layers or the self-attention layers, which leaves the remaining frozen structures unable to adapt to the data distributions of downstream tasks. Further, the existing structures are strongly coupled with the Transformers, hindering parameter-efficient deployment as well as the design flexibility for new approaches. In this paper, we revisit the design paradigm of PETL and derive a unified framework U-Tuning for parameter-efficient transfer learning, which is composed of an operation with frozen parameters and a unified tuner that adapts the operation for downstream applications. The U-Tuning framework can simultaneously encompass existing methods and derive new approaches for parameter-efficient transfer learning, which prove to achieve on-par or better performances on CIFAR-100 and FGVC datasets when compared with existing PETL methods.

1. Introduction

The increasingly large amount of data and powerful computing resources have driven the emergence of foundation models with large capacities (Devlin et al., 2018; Brown et al., 2020; Radford et al., 2021; Jia et al., 2021), which demonstrate strong generalization ability across multiple downstream tasks in visual (He et al., 2022b; Bao et al., 2021), language (Liu et al., 2019; Raffel et al., 2020) and

¹Alibaba Group ²National University of Singapore. Emails: Zeyinzi Jiang, Chaojie Mao, Yiliang Lv <zeyinzi.jyz, chaojie.mcj, yiliang.lyl@alibaba-inc.com>, Ziyuan Huang <ziyuan.huang@u.nus.edu>, Deli Zhao <zhaodeli@gmail.com>, Jingren Zhou <jingren.zhou@alibaba-inc.com>.

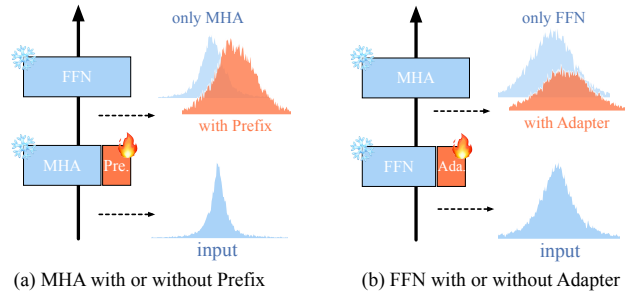


Figure 1. Visualization of input and output data distributions (a) in the first Transformer block and (b) between the first and the second Transformer block. Existing PETL methods only adjust the distribution of some parts of the pre-trained Transformer models, which brings difficulty for the remaining frozen model parts to adapt to the new distribution.

multi-modal paradigms (Li et al., 2019; Yu et al., 2022; Alayrac et al., 2022). On the flip side, the increasing size of pre-trained models has also made the training, storage, and deployment cost incredibly high for downstream adaptation.

To address this issue, the recent efforts are made in parameter-efficient transfer learning (PETL) (Houlsby et al., 2019), which proves to perform competitively against the fully-fine-tuning counterparts (Karimi Mahabadi et al., 2021a; Jia et al., 2022). Existing approaches for PETL either tune a few existing parameters (Zaken et al., 2021) or insert additional trainable structures (Lester et al., 2021; Li & Liang, 2021; Hu et al., 2021) with most or all pre-trained parameters frozen, maintaining a low training and deployment cost for downstream adaptation.

Despite the strong performance, the current practice has two shortcomings. (i) Most existing approaches insert learnable structures to only some parts of the model, leaving the remaining frozen parts struggling to adapt to the new distribution generated by the tuned blocks. For example, introducing prefix tokens (Li & Liang, 2021) essentially changes the output distribution of multi-head attention (MHA) in pre-trained Transformers, while the subsequent feed-forward network (FFN) is trained for the original distribution and will thus encounter difficulty in adapting to the new distribution. This goes for inserting trainable structures in FFNs as well, as can be observed in Figure 1. (ii) The structure of existing tuning methods and the main branch of the pre-trained

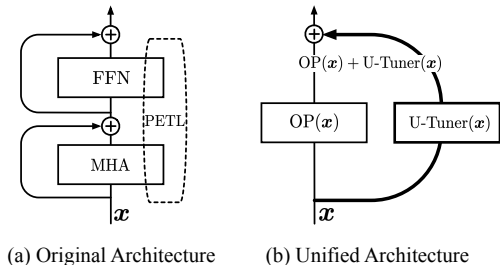


Figure 2. The diagram of existing methods and our unified framework for PETL. (a) Existing PETL methods are deeply embedded into original structures. (b) Our U-Tuning framework is composed of a frozen operation OP and a unified tuner U-Tuner .

models are strongly coupled, which means most existing structures are deeply embedded in the Transformers. On one hand, this hinders parameter-efficient deployment since the pre-trained model and the tuning structures need to be deployed entirely for every downstream application. On the other hand, the strong degree of coupling also limits the design flexibility for new approaches for PETL.

In this work, we rethink parameter-efficient transfer learning from a unified perspective towards existing approaches. Specifically, we revisit existing tuning paradigms and found a parallel form for mainstream tuning methods, such as adapters (Houlsby et al., 2019), prefix tuning (Li & Liang, 2021), and prompt tuning (Lester et al., 2021), which reduces the degree of coupling for tuning structures. This ability to transform to a parallel form for most mainstream methods also allows us to reveal a unified formulation. Based on this, we provide a unified framework for parameter-efficient transfer learning, which we call U-Tuning . Composed of an operation with frozen parameters and a lightweight unified trainable structure (as in Figure 2), the U-Tuning framework allows for flexible insertion of tuning structures. Hence, it can not only encompass most existing works but also allows for the derivation of new tuning structures. Extensive experiments on transfer learning show the generality of the U-Tuning framework, and the derived new PETL models prove to achieve on-par or better performances on various downstream tasks.

Our contributions can be summarized as follows: (i) We derive a parallel form for mainstream PETL methods, which reduces the degree of coupling and facilitates parameter-efficient deployment of large pre-trained models. (ii) We propose a unified tuning framework U-Tuning that encompasses existing PETL methods and allows for the derivation of new ones. (iii) Comprehensive studies on transfer learning prove the generality and powerfulness of U-Tuning .

To the best of our knowledge, this is the first work unifying PETL methods in a backbone-independent form, providing a different perspective for efficient tuning.

2. Related Work

Transformers in computer vision. Transformers (Vaswani et al., 2017) have achieved great success in various fields (Brown et al., 2020; Ramesh et al., 2022). In computer vision, Vision Transformers (ViT) (Dosovitskiy et al., 2020) are widely applied in various vision tasks, *e.g.*, visual classification (Liu et al., 2021d; Li et al., 2022), object detection (Song et al., 2022; Carion et al., 2020) and segmentation (Zheng et al., 2021; Wang et al., 2021), and have demonstrated strong generalization ability when pre-trained on a large corpus of data. A typical ViT consists of a cascade of Transformer blocks in which each block is constructed with a multi-head attention (MHA) and a feed-forward network (FFN). The existing PETL methods mainly perform lightweight adjustments to particular parts in the MHA, FFN, or the whole block.

MHA-based tuning. MHA-based tuning embeds trainable parameters in MHA. LoRA (Hu et al., 2021) constructs an additional layer with low-rank decomposition matrices of the weights in the network. Derived from specific textual templates (Brown et al., 2020; Liu et al., 2021a), prompt tuning (Lester et al., 2021) prepends extra trainable tokens (Lester et al., 2021; Liu et al., 2021c;b) to the input. Prefix tuning (Li & Liang, 2021) optimizes the task-specific vector in the multi-head attention layer. In computer vision, visual prompt tuning (VPT) (Jia et al., 2022) is proposed to initialize tunable prompt tokens and prepend to the original tokens in the first layer or multiple layers.

FFN-based tuning. For FFNs, the adaptation is generally made by adapter (Houlsby et al., 2019) and its generalized versions (Pfeiffer et al., 2020; Karimi Mahabadi et al., 2021b;a; He et al., 2022a), which usually insert a bottleneck layer into each FFN layer. In the video paradigm, adapters can also be used to introduce temporal information in FFNs (Chen et al., 2022).

Block-based tuning. Block-based tuning updates the common term of the block. A simple solution is BitFit (Zaken et al., 2021), which only tunes the bias terms of the model. Diff pruning (Guo et al., 2020) learns a diff vector, which is adaptively pruned during training process.

Other tuning. In addition, some of the efforts are made in finding out the optimal design paradigm of tuning modules. NOAH (Zhang et al., 2022) attempts to find out the optimal design of tuning modules through a neural architecture search algorithm. He et al. design a variety of forms about previous methods in a unified view and validate the best choices of variant adapters empirically.

In this paper, we rethink and analyze the previous methods systematically and propose a unified tuning framework, which presents to be a flexible and extensible paradigm for PETL on various downstream tasks.

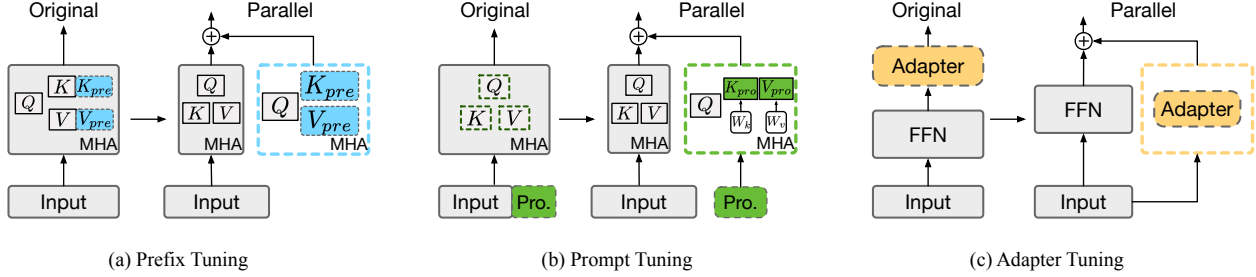


Figure 3. Illustration of (a) prefix tuning (Li & Liang, 2021), (b) prompt tuning (Lester et al., 2021), and (c) adapter tuning (Houlsby et al., 2019) for parameter-efficient transfer learning and their parallel form.

3. Rethinking PETL

In this section, we take a closer look into several existing approaches for PETL. Interestingly, we found that all of them can be equivalently formulated as the parallel combination of (i) existing operations (*e.g.*, multi-head attention (MHA) or feed-forward networks (FFN)) in Transformer architecture, and (ii) newly introduced structures performing a similar operation. Our investigation in these approaches is divided into two groups, MHA-based tuning methods, and FFN-based tuning methods.

3.1. MHA-based Tuning

Multi-head attention (MHA) based tuning usually takes effect on the self-attention part of the Transformer blocks. The original self-attention (Vaswani et al., 2017) is expressed as:

$$\text{Attn}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d}}\right)\mathbf{V}, \quad (1)$$

where \mathbf{Q} , \mathbf{K} and \mathbf{V} denote the query, key and value respectively. Usually, given input tokens \mathbf{x} , the query, key, and value are obtained through a linear projection $\mathbf{Q} = \mathbf{x}\mathbf{W}_q$, $\mathbf{K} = \mathbf{x}\mathbf{W}_k$, and $\mathbf{V} = \mathbf{x}\mathbf{W}_v$, where \mathbf{W}_q , \mathbf{W}_k and \mathbf{W}_v are learnable projection weights.

Prefix tuning. To adapt the output distribution of MHA to downstream tasks, prefix tuning (Li & Liang, 2021) keeps the projection weights in self-attention unchanged and prepends learnable parameters \mathbf{K}_{pre} and \mathbf{V}_{pre} , *i.e.*, prefix tokens, to the projected key and value respectively:

$$\text{MHA}_{pre} = \text{Attn}(\mathbf{x}\mathbf{W}_q, [\mathbf{K}_{pre}; \mathbf{x}\mathbf{W}_k], [\mathbf{V}_{pre}; \mathbf{x}\mathbf{W}_v]). \quad (2)$$

Essentially, when we view the MHA_{pre} as performing MHA separately between the query and the original keys and values, and between the query and the prefix tokens, we can obtain an equivalent form as:

$$\text{MHA}_{pre} = (1 - \lambda) \underbrace{\text{Attn}(\mathbf{Q}, \mathbf{K}, \mathbf{V})}_{\text{original attention}} + \lambda \underbrace{\text{Attn}(\mathbf{Q}, \mathbf{K}_{pre}, \mathbf{V}_{pre})}_{\text{prefix attention in parallel}}, \quad (3)$$

where λ weighs between the original attention and prefix attention in parallel. Detailed value for λ and the derivation process are included in Appendix A.1.

In this way, the original MHA, $\text{Attn}(\mathbf{Q}, \mathbf{K}, \mathbf{V})$ and the MHA for prefixes, $\text{Attn}(\mathbf{Q}, \mathbf{K}_{pre}, \mathbf{V}_{pre})$ can be computed in parallel. Hence, we call Equation (3) the parallel form of prefix tuning. Both the original form and the parallel form of prefix tuning can be seen in Figure 3 (a).

Prompt tuning. Instead of introducing learnable parameters at the same level as projected keys \mathbf{K} and values \mathbf{V} , prompt tuning (Lester et al., 2021) introduces learnable latent tokens \mathbf{x}_{pro} , *i.e.*, prompts, at the same level as the input tokens:

$$\text{MHA}_{pro} = \text{Attn}([\mathbf{x}; \mathbf{x}_{pro}]\mathbf{W}_q, [\mathbf{x}; \mathbf{x}_{pro}]\mathbf{W}_k, [\mathbf{x}; \mathbf{x}_{pro}]\mathbf{W}_v), \quad (4)$$

where the prompts \mathbf{x}_{pro} are concatenated to the input token \mathbf{x} in the first layer or multilayer.

Similar to prefix tuning, we can obtain the parallel form of prompt tuning as:

$$\text{MHA}_{pro} = [(1 - \lambda) \underbrace{\text{Attn}(\mathbf{Q}, \mathbf{K}, \mathbf{V})}_{\text{original attention}} + \lambda \underbrace{\text{Attn}(\mathbf{Q}, \mathbf{K}_{pro}, \mathbf{V}_{pro})}_{\text{prompt attention in parallel}}; (1 - \beta) \text{Attn}(\mathbf{Q}_{pro}, \mathbf{K}_{pro}, \mathbf{V}_{pro}) + \beta \text{Attn}(\mathbf{Q}_{pro}, \mathbf{K}, \mathbf{V})], \quad (5)$$

where \mathbf{K}_{pro} and \mathbf{V}_{pro} are the key and value produced by prompt tokens, respectively. λ and β are individual attention weights. More details can be seen in Appendix A.1.

Note that Equation (5) concatenates the prompt tokens with the original input tokens. In practice, prompt tokens are often discarded after the MHA. In this case, the concatenation operation in Equation (5) can be ignored, and thus the MHA_{pro} can have a similar form as Equation (3). The original form and the parallel form of prompt tuning can be seen in Figure 3 (b).

3.2. FFN-based Tuning

In Transformers (Vaswani et al., 2017), feed-forward networks (FFN) is a multi-layer perceptron (MLP) block di-

rectly connected to MHA, which is expressed as:

$$\text{FFN}(\mathbf{x}) = \phi(\mathbf{x}\mathbf{W}_1 + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2, \quad (6)$$

where \mathbf{W}_1 and \mathbf{W}_2 are the projection weights, \mathbf{b}_1 and \mathbf{b}_2 are the bias terms, and ϕ is the non-linear activation function between consecutive fully-connected layers.

Adapter tuning. Adapter (Houlsby et al., 2019) is typically used to adjust the output distribution for FFN in Transformers. Usually, adapters project input tokens \mathbf{x} by an MLP layer, similar to FFN. Differently, instead of expanding the feature dimension of the input tokens, adapters usually first shrink the channel dimension before it is expanded back to the original number:

$$\text{FFN}_{\text{adapter}} = \underbrace{\text{FFN}(\mathbf{x})}_{\text{original module}} + \underbrace{\phi(\text{FFN}(\mathbf{x})\mathbf{W}_{\text{down}})\mathbf{W}_{\text{up}}}_{\text{adapter module in parallel}}, \quad (7)$$

where \mathbf{W}_{down} and \mathbf{W}_{up} denote the weights for the down-projection layer and the up-projection layer, respectively.

Essentially, the FFN module with adapters can be thought of as the parallel connection between the original FFN module and the FFN module with an MLP layer, *i.e.*, an adapter¹. The comparison between the structures of the original form and the parallel form can be seen in Figure 3 (c). Another way to insert adapters is to add a scaling factor and design the adapter explicitly as a parallel module (He et al., 2022a; Chen et al., 2022), which can be similarly viewed as parallel structures.

4. U-Tuning

In this section, we first present a unified perspective on the existing tuning paradigms. Based on this, a Unified Tuning framework, dubbed U-Tuning, is proposed for PETL, which is composed of a frozen operation existing in pre-trained Transformers and a learnable **U-Tuner** for adjusting the output distributions. After that, we instantiate our U-Tuning with various building blocks.

4.1. A Unified Formulation

From Equations (3), (5) and (7), we can make two observations: **(i)** existing tuning methods use a similar structure, *i.e.*, the parallel combination of an existing operation with frozen parameters and a newly-introduced structure with learnable parameters. The frozen parts yield generalized representations learned from large-scale data, while the learnable parts adapt the generalized representations to specific downstream tasks. **(ii)** the learnable parts of existing tuning methods generally have a similar structure to the corresponding

¹For clarity, the residual connection in the outer layer of FFN is not taken into account.

frozen parts, *e.g.*, prefix for MHA essentially introduces self-attention between existing queries and prefixes.

Given the above observations, we generalize existing tuning methods for PETL into a unified formulation:

$$\mathbf{x}' = \text{OP}(\mathbf{x}) + \mathbf{U-Tuner}(\mathbf{x}) \quad (8)$$

where \mathbf{x}/\mathbf{x}' are input/output tokens, OP denotes the existing operations in Transformers with frozen parameters, and **U-Tuner** represents the unified tuner, which is newly-introduced for adjusting the output distributions of OP.

Our unified formulation views each part of Transformer as an operation function OP, while each tuned part as a unified tuner **U-Tuner**. In this way, the formulation encompasses all existing tuning methods when we instantiate OP and **U-Tuner** with similar operations. Meanwhile, it enables the generation of new tuning methods for parameter-efficient transfer learning when we instantiate them with different building blocks.

Further, compared to existing tuning methods where tuning structures are only attached to a subset of operations (*e.g.*, only to MHA or only to FFN), our formulation can attach **U-Tuner** to all operations (MHA and FFN) or even to Transformer blocks, as can be seen in Figure 4.

4.2. Instantiations

The U-Tuning framework can be instantiated in two parts, *i.e.*, the operations OP existing in pre-trained Transformers with frozen parameters and the unified tuner **U-Tuner** with learnable parameters, respectively. To instantiate existing tuning methods, we can instantiate U-Tuning with corresponding OP and **U-Tuner**. On top of that, we can generate various new tuning methods with different combinations of operations and tuners.

Operations. The emergence of parameter-efficient transfer learning is majorly motivated by the strong generalization capabilities of large foundation models pre-trained on a large corpus of data. Hence, on the micro level, we follow the practice of existing tuning methods and align the operations OP term in our unified formulation with the design of pre-trained Transformers, which can be either MHA or FFN. This alignment allows the retainment of the generalization capabilities of the pre-trained Transformers. Different from the existing tuning paradigm for PETL, we additionally introduce a macro level to the operations, where the OP term can also be a Transformer block. This further increases the flexibility of the formulation as well as the U-Tuning framework. The detailed instantiation of the operations can be viewed in Figure 4.

U-Tuner. According to the instantiated operation, our **U-Tuner** can be connected in parallel to either MHA, FFN, or the block, as in Figure 4 (a). The instantiation

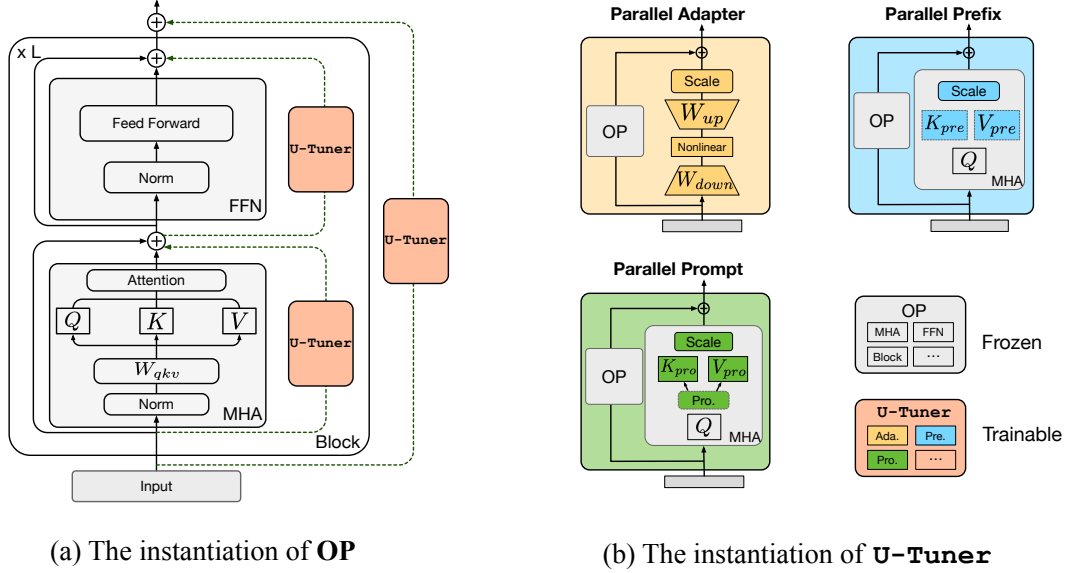


Figure 4. Our U-Tuning framework for parameter-efficient transfer learning (PETL) consists of an operation (OP) with frozen parameters and a unified tuner (**U-Tuner**) that adapts the output distribution to downstream applications. Unlike existing tuning methods, the OP and the **U-Tuner** are instantiated in a decoupled way, which means the instantiation of the OP and **U-Tuner** can have arbitrary combinations. (a) The instantiation of OP can be MHA, FFN, and the whole block in pre-trained Transformers. (b) The instantiation of **U-Tuner** can be a parallel adapter, parallel prefix, and parallel prompt. Note that when new structures or operations emerge, our U-Tuning framework can always be flexibly extended to encompass new operations.

of the **U-Tuner** is decoupled from the instantiation of the operations, which means the structure of the **U-Tuner** does not necessarily need to be consistent with the corresponding operation. This is in contrast to the existing tuning methods, where the tuner performs a similar operation as the existing operations in Transformers, *e.g.*, adapter MLP for FFN.

For the instantiation of **U-Tuner**, we introduce three variants extended from adapter (Houlsby et al., 2019), prefix tuning (Li & Liang, 2021), and prompt tuning (Lester et al., 2021), respectively. Specifically, we use the parallel form of the above methods that we derive in Section 3, *i.e.*, parallel adapter, parallel prefix, and parallel prompt (as shown in Figure 4), such that **U-Tuner** can be easily plugged into or removed from pre-trained Transformers with high flexibility, which facilitates the deployment of such large models in various downstream applications. Such design also allows for easy extension to new tuning methods when new operations or structures emerge.

Another important difference between our U-Tuning framework and existing approaches is that the instantiation of our **U-Tuner** is not limited to only one OP in one block. This means the distributions of all modules in pre-trained models can be adjusted to fit the new data distributions in various downstream tasks. We study the performance of our U-Tuning framework with one, two, and three **U-Tuner** instantiated in the experiments.

4.3. Analysis and Discussion

Scaling. Our U-Tuning framework is able to attach the **U-Tuner** to all modules in pre-trained Transformers. To further improve the flexibility as well as to balance between different modules, we introduce a channel-wise scaling factor to all the parallel tuners, *i.e.*,

$$\mathbf{U-Tuner}(x)' = s \cdot \mathbf{U-Tuner}(x), \quad (9)$$

where s is the channel-wise scaling factor. We also study other forms of scaling factors in experiments.

Parameter analysis. The PETL methods fix most of the pre-trained model parameters and introduce only a few learnable parameters, as is our method. The total number of our tunable parameters basically depends on the current method used in **U-Tuner**. For example, the amount of parameters for VPT-Deep (Jia et al., 2022) is $L \times n \times d$, where L is the total layers, n is the number of prompts, and d is the dimension of prompts. Our following experiments show that U-Tuning can surpass the method used by the full layer (L layers) with a lower number of parameters, thus achieving a fewer parameter count by comparison. Other parameters have similar results, such as the number of prefix tokens, and the hidden layer dimension of the adapter.

Diverse unified perspective. In NLP, He et al. also look at PETL from a unified view, aiming to establish close connections and unify to the adapter method. Although

Table 1. Performance and parameter comparisons on FGVC. The bold font represents the best accuracy, and the underline represents the second best accuracy. “Mean” denotes the average accuracy of datasets. “Params” is the average tunable parameters needed for training.

METHOD	CUB_200_2011	NABIRDS	OXFORDFLOWERS	STANFORDCARS	STANFORDDOGS	MEAN	PARAMS.(M)
FULLY FINE-TUNING	87.3	82.7	98.8	84.5	89.4	88.54	85.98
LINEAR PROBING	85.3	75.9	97.9	51.3	86.2	79.32	0.18
BITFIT	88.4	84.2	98.8	79.4	<u>91.2</u>	88.41	0.28
PREFIX	87.5	82.0	98.0	74.2	<u>90.2</u>	86.37	0.36
ADAPTER	87.1	<u>84.3</u>	98.5	68.6	89.8	85.67	0.41
VPT-SHALLOW	86.7	78.8	98.4	68.7	90.7	84.62	0.25
VPT-DEEP	<u>88.5</u>	84.2	<u>99.0</u>	83.6	90.2	<u>89.11</u>	0.85
U-TUNING	89.16	85.39	99.15	<u>84.14</u>	92.07	89.98	0.36

Table 2. Comparison between the U-Tuning framework and existing tuning approaches for PETL on CIFAR-100. U-Tuning/1 and U-Tuning/12 represent introducing **U-Tuner** in the first layer and in all layers of Vision Transformers respectively. † denotes results from our reimplementation. * denotes tri-**U-Tuner** variant for U-Tuning.

METHOD	CIFAR-100	PARAMS.(M)
FULLY FINE-TUNING	89.12	86.04
LINEAR PROBING	85.95	0.07
PREFIX†	90.95	0.26
VPT-SHALLOW†	86.62	0.08
VPT-DEEP†	91.58	0.17
ADAPTER†	91.80	0.27
ADAPTFORMER	91.86	1.26
U-TUNING/1	91.86	0.11
U-TUNING/12	<u>92.57</u>	0.59
U-TUNING/12*	92.75	0.67

U-Tuning takes inspiration, our approach is different in motivations, design principles, and fields. U-Tuning derives existing methods in a completely equivalent form and proposes a unified paradigm in which any tuning method can be independent of the main branch, whereas He et al. unify PETL methods to find several variant adapters in an approximate manner.

Expansibility and applicability. Our approach proposes a unified paradigm in which the two core operators, OP and **U-Tuner** represent the submodules to be optimized and the tuning methods in parallel. Both operators need not be limited to the approach presented in this paper, but extend to other approaches, *e.g.*, LoRA (Hu et al., 2021), and even newer ones. Furthermore, U-Tuning also allows for plug and play of **U-Tuner** that connects to the original modules through residual connection. In this way, rapid adaption and deployment are possible for all kinds of downstream tasks.

5. Experiments

In this section, we evaluate and analyze the performance and design of our proposed U-Tuning method on several downstream tasks. Specifically, we describe our experi-

ment set up in Section 5.1, compare with SOTA methods in Section 5.2, and present comprehensive ablative analysis in Section 5.3.

5.1. Experiment Setup

Datasets. To evaluate the capability of our method on various downstream tasks, we take CIFAR-100 (Krizhevsky et al., 2009) and FGVC datasets, where CIFAR-100 is typically used on general image classification tasks. FGVC is the benchmark of fine-grained visual classification tasks, which consists of CUB-200-2011 (Wah et al., 2011), NABirds (Van Horn et al., 2015), Oxford Flowers (Nilsback & Zisserman, 2008), Stanford Cars (Gebru et al., 2017), and Stanford Dogs (Khosla et al., 2011). We report the Top-1 accuracy on CIFAR-100 and FGVC datasets for the experiments of transfer learning.

Baselines. We first compare our method with baselines in transfer learning and divide them into two categories: (i) Traditional methods: fully fine-tuning, which updates all the parameters of model, and linear probing, which freezes the pre-trained backbone and only tunes the classifier. (ii) PETL methods: BitFit (Zaken et al., 2021), adapter (Houlsby et al., 2019), prefix (Li & Liang, 2021), VPT (Lester et al., 2021), and AdaptFormer (Chen et al., 2022).

Implementation details. If not specified, we experiment with ViT-B/16 (Dosovitskiy et al., 2020) model pre-trained on ImageNet-21K (Deng et al., 2009) as conducted in VPT (Lester et al., 2021). For most datasets, we preprocess the data with a random resized crop and a random horizontal flip to the size of 224×224 . AdamW (Loshchilov & Hutter, 2019) optimizer and cosine annealing learning rate scheduler are used with linear warm-up strategy. We use a combination of **U-Tuner** with a channel-wise scaling strategy instead of a single one to get better performance. More details are introduced in Appendix A.2.

5.2. Comparisons with the SOTA

We compare the transfer ability of our U-Tuning framework with different existing approaches for parameter-

Table 3. Verification of the equivalence between the original and the parallel forms of existing PETL methods.

TYPE	ADAPTER	PREFIX	PROMPT
ORIGINAL	91.80	90.95	91.19
PARALLEL	91.86	91.01	91.26

Table 4. Single **U-Tuner** with different parallel tuners attached to different operations OP.

METHOD	MHA	FFN	BLOCK
VPT	90.97	-	-
P-ADAPTER	92.43	92.19	92.13
P-PREFIX	91.42	90.27	90.26
P-PROMPT	91.54	90.01	90.12

efficient transfer learning on CIFAR-100 and FGVC datasets, which are evaluated in Top-1 accuracy.

FGVC. As shown in Table 1, our **U-Tuning** outperforms VPT (Jia et al., 2022) and other tuning methods (fully fine-tuning, linear probing, bias, prefix, and adapter) on the average accuracy of five FGVC datasets. For all FGVC datasets other than StanfordCars, our **U-Tuning** outperforms existing tuning methods by notable margins. For StanfordCars, the **U-Tuning** falls short of fully fine-tuning but performs stronger than other tuning methods. We suspect that it is because the intra-class difference is too small and hence more parameters are required for this dataset. The used construction form of **U-Tuning** for each dataset is listed in Table 11 in Appendix A.2.3, which is the **U-Tuning** form with the best performance in all possible construction forms. The ablation studies on the performance for different construction forms of **U-Tuning** are discussed in Section 5.3.

CIFAR-100. As in Table 2, our approach goes beyond fully fine-tuning and linear probing by 3.45% and 6.62%, respectively. Compared with existing methods for PETL, **U-Tuning** achieves on-par or better performances. Specifically, **U-Tuning** with the first layer adapted by the unified tuner **U-Tuner** achieves similar performance to the previous best performances achieved by existing approaches (91.86 of both **U-Tuning/1** and **AdaptFormer**), while reducing the parameter requirement by over 10 times (0.11M of **U-Tuning/1** and 1.26M of **AdaptFormer**). When all layers are adapted by the unified tuner, we observe a notable improvement over state-of-the-art accuracy (92.57 of **U-Tuning/12** vs. 91.86 of **AdaptFormer**). Note that, compared with **AdaptFormer**, **U-Tuning/12** only uses less than half of the parameters used by **AdaptFormer** (0.59M of **U-Tuning/12** vs. 1.26M of **AdaptFormer**). Using the **tri-U-Tuner** variant, **U-Tuning/12** achieves the strongest result.

Table 5. Dual-**U-Tuner** with different tuner instantiations attached to MHA and FFN.

	FFN	P-ADAPTER	P-PREFIX	P-PROMPT
MHA				
P-ADAPTER		92.57	92.45	92.30
P-PREFIX		92.10	91.52	91.69
P-PROMPT		92.27	91.60	91.76

Table 6. Tri-**U-Tuner** with different tuner instantiations for the block. The parallel adapter is used for MHA and FFN.

BLOCK	P-ADAPTER	P-PREFIX	P-PROMPT
ACCURACY	92.69	92.75	92.53

5.3. Ablation Study

We present comprehensive ablative analysis here for a deeper understanding of **U-Tuning**.

The equivalency of the original form and the parallel form of existing PETL methods. We first empirically verify the alignment in performance between the existing tuning methods in their original form and their parallel equivalent derivations. As in Table 3, it can be observed that the performance gap is small enough to be neglected. This equivalence also indicates that our unified framework can effectively encompass existing tuning methods for PETL.

Different combinations of OP and U-Tuner. After showing that our **U-Tuning** framework can encompass existing tuning methods for PETL, we further explore various combinations of OP and **U-Tuner** to derive new tuning methods. We focus our explorations here on CIFAR-100 dataset.

(i) **Single U-Tuner.** We start by using one **U-Tuner** per block, which is similar to existing PETL methods. We explore different **U-Tuner** forms on different operations OP. Given the existing paradigms for PETL, we focus our exploration on the aforementioned parallel adapter (P-Adapter), parallel prefix (P-Prefix), and parallel prompt (P-Prompt). As in Table 4, we connect them to MHA, FFN, and the whole block, respectively. It can be observed that connecting P-Adapter to MHA achieves the strongest result, which means that the existing practice is not necessarily the optimal combination. Further, compared to VPT (Jia et al., 2022), our single **U-Tuner** variant can already achieve a notable improvement.

(ii) **Dual-U-Tuner.** We further explore combinations of OP and **U-Tuner** when both MHA and FFN are adapted by the **U-Tuner** in Table 5. Interestingly, the **U-Tuning** framework performs the best when we use parallel adapters for both MHA and FFN. Further, we found the performance

NABirds				StanfordCars					
MHA	P-Adapter	85.13	85.39	85.14	MHA	P-Adapter	84.14	83.75	83.91
	P-Prefix	85.18	83.94	84.36		P-Prefix	82.71	77.37	78.65
	P-Prompt	85.21	84.20	84.22		P-Prompt	82.58	77.86	78.72
		FFN					FFN		

Figure 5. Generalization of our findings in dual-U-Tuner to NABirds and StanfordCars.

Table 7. Ablation studies on the scaling strategies.

TYPE	ACCURACY
DIRECT CONNECTION	92.18
SCALAR SCALING	92.22
CHANNEL-WISE SCALING	92.57
INPUT-DEPENDENT SCALING	92.12

with at least one U-Tuner using the parallel adapter is generally higher. Dual-U-Tuner also outperforms the single U-Tuner variant, indicating that all modules in pre-trained Transformers should be adapted for PETL. Unless otherwise specified, we use dual-U-Tuner as our default version for the remaining of the ablative analysis.

(iii) **Tri-U-Tuner**. On top of the best dual-U-Tuner, we add an additional block U-Tuner in Table 6. The performances of adding the parallel prefix and the parallel adapter are similar. Compared with dual-U-Tuner, tri-U-Tuner further improves by 0.18%.

Generalization of our findings to fine-grained datasets. We further experiment with dual-U-Tuner on fine-grained datasets to examine our previous findings on CIFAR-100. The results are shown in Figure 5. On fine-grained datasets such as NABirds and StanfordCars, the observations are consistent: models with adapters as the U-Tuner generally achieve stronger performance.

Different scaling strategies. In Table 7, we compare the default channel-wise scaling with other scaling strategies, *i.e.*, direct connection (no scaling), scalar scaling, and input-dependent scaling (SE (Hu et al., 2018)-like). We find channel-wise scaling better than other variants.

Different number of layers with U-Tuner. In Figure 6, we vary the number of layers to insert the U-Tuner from bottom to top on the FGVC dataset. Increasing the number of layers with U-Tuner effectively raises the performance of StanfordCars, while the other datasets are less affected.

Varying the length of dimensions. The length of dimen-

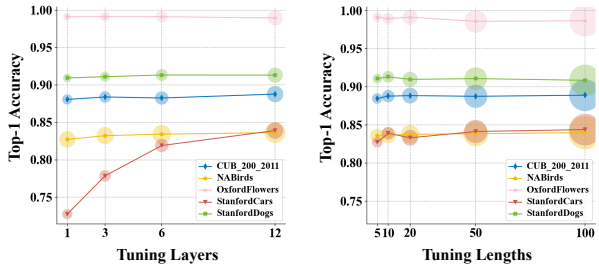


Figure 6. Ablation studies on the tuning layers and the length of dimension for U-Tuning on FGVC.

Table 8. Ablative analysis on different pre-trained Transformers.

ARCHITECTURE PRE-TRAINING	ViT/B-16			ViT/L-14
	IN-21K	MAE	CLIP	CLIP
ADAPTER	91.88	86.11	89.00	92.01
PREFIX	90.95	78.35	86.54	90.50
VPT-SHALLOW	86.62	55.98	85.47	89.51
VPT-DEEP	91.58	82.95	87.77	91.56
U-TUNING	92.57	86.17	89.17	92.38

sions represents the hidden dimension for parallel adapters and the number of prompt/prefix tokens. Altering the length of dimensions affects performances and the number of parameters simultaneously. In Figure 6, we show that increasing the length of dimensions has a relatively great impact when its value is small.

Different pre-training sources and backbones. In Table 8, we compare several tuning methods with our proposed U-Tuning based on different architectures and pre-trained models, *i.e.*, IN-21K (Deng et al., 2009), MAE (He et al., 2022b) and CLIP (Radford et al., 2021) pre-trained ViT/B-16 (Dosovitskiy et al., 2020), as well as CLIP pre-trained ViT/L-14. The U-Tuning framework achieves the best performance for all pre-trained Transformers.

6. Conclusion

In this work, we revisit mainstream methods for parameter-efficient transfer learning (PETL) and establish a unified framework for PETL. Our approach, dubbed U-Tuning, is composed of an operation OP with frozen parameters and a unified tuner U-Tuner for adapting the operation to various downstream tasks. Instantiated with various building blocks, our U-Tuning framework can encompass most existing tuning methods for PETL and meanwhile generate new tuning methods, such that U-Tuning is more flexible and efficient in the training and deployment process. Empirically, we show that the U-Tuning paradigm achieves strong results on transfer learning. We hope that U-Tuning can provide a new perspective on parameter-efficient tuning methods.

References

- Alayrac, J.-B., Donahue, J., Luc, P., Miech, A., Barr, I., Hasson, Y., Lenc, K., Mensch, A., Millican, K., Reynolds, M., et al. Flamingo: a visual language model for few-shot learning. *arXiv preprint arXiv:2204.14198*, 2022.
- Bao, H., Dong, L., Piao, S., and Wei, F. BEiT: BERT pre-training of image Transformers. *arXiv preprint arXiv:2106.08254*, 2021.
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. Language models are few-shot learners. In *Adv. Neural Inform. Process. Syst.*, 2020.
- Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov, A., and Zagoruyko, S. End-to-end object detection with Transformers. In *Eur. Conf. Comput. Vis.*, pp. 213–229, 2020.
- Chen, S., Ge, C., Tong, Z., Wang, J., Song, Y., Wang, J., and Luo, P. AdaptFormer: Adapting vision Transformers for scalable visual recognition. In *Adv. Neural Inform. Process. Syst.*, 2022.
- Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. ImageNet: A large-scale hierarchical image database. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pp. 248–255, 2009.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. BERT: Pre-training of deep bidirectional Transformers for language understanding. *North Am. Chap. Assoc. Comput. Linguist.*, 2018.
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *Int. Conf. Learn. Represent.*, 2020.
- Gebru, T., Krause, J., Wang, Y., Chen, D., Deng, J., and Fei-Fei, L. Fine-grained car detection for visual census estimation. In *Assoc. Adv. Artif. Intell.*, 2017.
- Guo, D., Rush, A. M., and Kim, Y. Parameter-efficient transfer learning with diff pruning. *Assoc. Comput. Linguist.*, 2020.
- He, J., Zhou, C., Ma, X., Berg-Kirkpatrick, T., and Neubig, G. Towards a unified view of parameter-efficient transfer learning. In *Int. Conf. Learn. Represent.*, 2022a.
- He, K., Chen, X., Xie, S., Li, Y., Dollár, P., and Girshick, R. Masked autoencoders are scalable vision learners. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pp. 16000–16009, 2022b.
- Houlsby, N., Giurgiu, A., Jastrzebski, S., Morrone, B., De Laroussilhe, Q., Gesmundo, A., Attariyan, M., and Gelly, S. Parameter-efficient transfer learning for NLP. In *Int. Conf. Mach. Learn.*, pp. 2790–2799, 2019.
- Hu, E., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, L., and Chen, W. LoRA: Low-rank adaptation of large language models. In *Int. Conf. Learn. Represent.*, 2021.
- Hu, J., Shen, L., and Sun, G. Squeeze-and-excitation networks. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pp. 7132–7141, 2018.
- Jia, C., Yang, Y., Xia, Y., Chen, Y.-T., Parekh, Z., Pham, H., Le, Q., Sung, Y.-H., Li, Z., and Duerig, T. Scaling up visual and vision-language representation learning with noisy text supervision. In *Int. Conf. Mach. Learn.*, pp. 4904–4916, 2021.
- Jia, M., Tang, L., Chen, B.-C., Cardie, C., Belongie, S., Hariharan, B., and Lim, S.-N. Visual prompt tuning. In *Eur. Conf. Comput. Vis.*, 2022.
- Karimi Mahabadi, R., Henderson, J., and Ruder, S. Compacter: Efficient low-rank hypercomplex adapter layers. *Adv. Neural Inform. Process. Syst.*, 34:1022–1035, 2021a.
- Karimi Mahabadi, R., Ruder, S., Dehghani, M., and Henderson, J. Parameter-efficient multi-task fine-tuning for Transformers via shared hypernetworks. In *Assoc. Comput. Linguist.*, pp. 565–576, 2021b.
- Khosla, A., Jayadevaprakash, N., Yao, B., and Li, F.-F. Novel dataset for fine-grained image categorization: Stanford dogs. In *IEEE Conf. Comput. Vis. Pattern Recog. Worksh.*, 2011.
- Krizhevsky, A., Hinton, G., et al. Learning multiple layers of features from tiny images. 2009.
- Lester, B., Al-Rfou, R., and Constant, N. The power of scale for parameter-efficient prompt tuning. In *Conf. Empirical Methods NLP*, 2021.
- Li, K., Wang, Y., Gao, P., Song, G., Liu, Y., Li, H., and Qiao, Y. UniFormer: Unified Transformer for efficient spatiotemporal representation learning. *Int. Conf. Learn. Represent.*, 2022.
- Li, L. H., Yatskar, M., Yin, D., Hsieh, C.-J., and Chang, K.-W. VisualBERT: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*, 2019.
- Li, X. L. and Liang, P. Prefix-tuning: Optimizing continuous prompts for generation. In *Assoc. Comput. Linguist.*, pp. 4582–4597, 2021.

- Liu, P., Yuan, W., Fu, J., Jiang, Z., Hayashi, H., and Neubig, G. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *arXiv preprint arXiv:2107.13586*, 2021a.
- Liu, X., Ji, K., Fu, Y., Du, Z., Yang, Z., and Tang, J. P-tuning v2: Prompt tuning can be comparable to fine-tuning universally across scales and tasks. *arXiv preprint arXiv:2110.07602*, 2021b.
- Liu, X., Zheng, Y., Du, Z., Ding, M., Qian, Y., Yang, Z., and Tang, J. GPT understands, too. *arXiv preprint arXiv:2103.10385*, 2021c.
- Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., and Stoyanov, V. RoBERTa: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., and Guo, B. Swin Transformer: Hierarchical vision Transformer using shifted windows. In *Int. Conf. Comput. Vis.*, 2021d.
- Loshchilov, I. and Hutter, F. Decoupled weight decay regularization. *Int. Conf. Learn. Represent.*, 2019.
- Nilsback, M.-E. and Zisserman, A. Automated flower classification over a large number of classes. In *Ind. Conf. Comput. Vis. Graph. Image Process.*, pp. 722–729, 2008.
- Pfeiffer, J., Kamath, A., Rücklé, A., Cho, K., and Gurevych, I. AdapterFusion: Non-destructive task composition for transfer learning. In *Eur. Chap. Assoc. Comput. Linguist.*, pp. 487–503, 2020.
- Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al. Learning transferable visual models from natural language supervision. In *Int. Conf. Mach. Learn.*, pp. 8748–8763, 2021.
- Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., and Liu, P. J. Exploring the limits of transfer learning with a unified text-to-text Transformer. *J. Mach. Learn. Res.*, 21(140):1–67, 2020.
- Ramesh, A., Dhariwal, P., Nichol, A., Chu, C., and Chen, M. Hierarchical text-conditional image generation with CLIP latents. *arXiv preprint arXiv:2204.06125*, 2022.
- Song, H., Sun, D., Chun, S., Jampani, V., Han, D., Heo, B., Kim, W., and Yang, M.-H. ViDT: An efficient and effective fully Transformer-based object detector. In *Int. Conf. Learn. Represent.*, 2022.
- Van Horn, G., Branson, S., Farrell, R., Haber, S., Barry, J., Ipeiritos, P., Perona, P., and Belongie, S. Building a bird recognition app and large scale dataset with citizen scientists: The fine print in fine-grained dataset collection. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pp. 595–604, 2015.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. Attention is all you need. *Adv. Neural Inform. Process. Syst.*, 30, 2017.
- Wah, C., Branson, S., Welinder, P., Perona, P., and Belongie, S. The Caltech-UCSD birds-200-2011 dataset. 2011.
- Wang, W., Xie, E., Li, X., Fan, D.-P., Song, K., Liang, D., Lu, T., Luo, P., and Shao, L. Pyramid vision Transformer: A versatile backbone for dense prediction without convolutions. In *Int. Conf. Comput. Vis.*, pp. 568–578, 2021.
- Yu, J., Wang, Z., Vasudevan, V., Yeung, L., Seyedhosseini, M., and Wu, Y. CoCa: Contrastive captioners are image-text foundation models. *arXiv preprint arXiv:2205.01917*, 2022.
- Zaken, E. B., Ravfogel, S., and Goldberg, Y. BitFit: Simple parameter-efficient fine-tuning for Transformer-based masked language-models. *Assoc. Comput. Linguist.*, 2021.
- Zhang, Y., Zhou, K., and Liu, Z. Neural prompt search. *arXiv preprint arXiv:2206.04673*, 2022.
- Zheng, S., Lu, J., Zhao, H., Zhu, X., Luo, Z., Wang, Y., Fu, Y., Feng, J., Xiang, T., Torr, P. H., and Zhang, L. Rethinking semantic segmentation from a sequence-to-sequence perspective with Transformers. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2021.

A. Appendix

A.1. Detailed derivations

- **Prefix tuning:** The following is the detailed derivation of Equation (3).

$$\begin{aligned}
 \text{head} &= \text{Attn}(\mathbf{x}\mathbf{W}_q, \text{concat}(\mathbf{K}_{pre}, \mathbf{x}\mathbf{W}_k), \text{concat}(\mathbf{V}_{pre}, \mathbf{x}\mathbf{W}_v)) \\
 &= \text{softmax}\left(\mathbf{x}\mathbf{W}_q \text{concat}(\mathbf{K}_{pre}, \mathbf{x}\mathbf{W}_k)^\top\right) \begin{bmatrix} \mathbf{V}_{pre} \\ \mathbf{x}\mathbf{W}_v \end{bmatrix} \\
 &= (1 - \lambda(\mathbf{x})) \text{softmax}\left(\mathbf{x}\mathbf{W}_q \mathbf{W}_k^\top \mathbf{x}^\top\right) \mathbf{x}\mathbf{W}_v + \lambda(\mathbf{x}) \text{softmax}\left(\mathbf{x}\mathbf{W}_q \mathbf{K}_{pre}^\top\right) \mathbf{V}_{pre} \\
 &= (1 - \lambda(\mathbf{x})) \text{Attn}(\mathbf{x}\mathbf{W}_q, \mathbf{x}\mathbf{W}_k, \mathbf{x}\mathbf{W}_v) + \lambda(\mathbf{x}) \text{Attn}(\mathbf{x}\mathbf{W}_q, \mathbf{K}_{pre}, \mathbf{V}_{pre}) \\
 &= (1 - \lambda(\mathbf{Q}, \mathbf{K}, \mathbf{K}_{pre})) \underbrace{\text{Attn}(\mathbf{Q}, \mathbf{K}, \mathbf{V})}_{\text{standard attention}} + \lambda(\mathbf{Q}, \mathbf{K}, \mathbf{K}_{pre}) \underbrace{\text{Attn}(\mathbf{Q}, \mathbf{K}_{pre}, \mathbf{V}_{pre})}_{\text{independent of } \mathbf{K}_{pre}, \mathbf{V}_{pre}}
 \end{aligned} \tag{10}$$

where

$$\lambda(\mathbf{Q}, \mathbf{K}, \mathbf{K}_{pre}) = \frac{\sum_i \exp(\mathbf{Q}\mathbf{K}_{pre}^\top)_i}{\sum_i \exp(\mathbf{Q}\mathbf{K}^\top)_i + \sum_j \exp(\mathbf{Q}\mathbf{K}_{pre}^\top)_j}, \tag{11}$$

- **Prompt tuning:** The following is the detailed derivation of Equation (5).

$$\begin{aligned}
 \text{head} &= \text{Attn}(\text{concat}(\mathbf{x}, \mathbf{x}_{pro})\mathbf{W}_q, \text{concat}(\mathbf{x}, \mathbf{x}_{pro})\mathbf{W}_k, \text{concat}(\mathbf{x}, \mathbf{x}_{pro})\mathbf{W}_v) \\
 &= \text{concat}\left(\text{softmax}\left(\mathbf{x}\mathbf{W}_q \text{concat}(\mathbf{x}\mathbf{W}_k, \mathbf{x}_{pro}\mathbf{W}_k)^\top\right) \begin{bmatrix} \mathbf{x}\mathbf{W}_v \\ \mathbf{x}_{pro}\mathbf{W}_v \end{bmatrix}, \right. \\
 &\quad \left. \text{softmax}\left(\mathbf{x}_{pro}\mathbf{W}_q \text{concat}(\mathbf{x}\mathbf{W}_k, \mathbf{x}_{pro}\mathbf{W}_k)^\top\right) \begin{bmatrix} \mathbf{x}\mathbf{W}_v \\ \mathbf{x}_{pro}\mathbf{W}_v \end{bmatrix}\right) \\
 &= \text{concat}\left((1 - \lambda(\mathbf{Q}, \mathbf{K}, \mathbf{K}_{pro})) \text{Attn}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) + \lambda(\mathbf{Q}, \mathbf{K}, \mathbf{K}_{pro}) \text{Attn}(\mathbf{Q}, \mathbf{K}_{pro}, \mathbf{V}_{pro}), \right. \\
 &\quad \left. (1 - \beta(\mathbf{Q}_{pro}, \mathbf{K}_{pro}, \mathbf{K})) \text{Attn}(\mathbf{Q}_{pro}, \mathbf{K}_{pro}, \mathbf{V}_{pro}) + \beta(\mathbf{Q}_{pro}, \mathbf{K}_{pro}, \mathbf{K}) \text{Attn}(\mathbf{Q}_{pro}, \mathbf{K}, \mathbf{V})\right)
 \end{aligned} \tag{12}$$

where

$$\lambda(\mathbf{Q}, \mathbf{K}, \mathbf{K}_{pro}) = \frac{\sum_i \exp(\mathbf{Q}\mathbf{K}_{pro}^\top)_i}{\sum_i \exp(\mathbf{Q}\mathbf{K}^\top)_i + \sum_j \exp(\mathbf{Q}\mathbf{K}_{pro}^\top)_j}, \tag{13}$$

$$\beta(\mathbf{Q}_{pro}, \mathbf{K}_{pro}, \mathbf{K}) = \frac{\sum_i \exp(\mathbf{Q}_{pro}\mathbf{K}^\top)_i}{\sum_i \exp(\mathbf{Q}_{pro}\mathbf{K}_{pro}^\top)_i + \sum_j \exp(\mathbf{Q}_{pro}\mathbf{K}^\top)_j}, \tag{14}$$

A.2. Experimental Details

A.2.1. DATASETS

We list the description of each dataset and our experimental settings in Table 9, which includes the number of classes and the amount of images in training set and test set.

Table 9. Datasets and settings in our experiments.

DATASET	CLASSES	TRAIN	TEST
GENERAL IMAGE CLASSIFICATION			
CIFAR-100 (KRIZHEVSKY ET AL., 2009)	100	50000	10000
FINE-GRAINED VISUAL CLASSIFICATION (FGVC)			
CUB_200_2011 (WAH ET AL., 2011)	200	5994	5794
NABIRDS (VAN HORN ET AL., 2015)	555	23929	24633
OXFORDFLOWERS (NILSBACK & ZISSERMAN, 2008)	102	1020	6149
STANFORDCARS (GEBRU ET AL., 2017)	196	8144	8041
STANFORDDOGS (KHOSLA ET AL., 2011)	120	12000	8580

A.2.2. HYPERPARAMETERS

To help reproduce the experiments conducted in this paper, we list all the hyperparameters and the architecture’s units. Concretely, the training hyperparameters are listed in the first section of Table 10. In addition, the choices of the other attributes, such as the number of Transformer layers, the length of dimensions, the used OP of transformer, the basic tuners and so on, are listed in the others sections of Table 10, respectively.

Table 10. Hyperparameters and architecture’s units in detail

CONFIG	VALUE
BATCH SIZE	32
OPTIMIZER	ADAMW (LOSHCHILOV & HUTTER, 2019)
WEIGHT DECAY	0.05
BASE LEARNING RATE RANGE	{0.001, 0.005}
LEARNING RATE SCHEDULE	COSINE DECAY
TRAINING EPOCHS RANGE	{50, 100}
WARMUP EPOCHS	10
AUGMENTATION	RANDOMRESIZEDCROP, RANDOMHORIZONTALFLIP
THE NUMBER OF TRANSFORMER LAYERS	{1, 1-3, 1-6, 1-12}
THE LENGTH OF DIMENSIONS	{5, 10, 20, 50, 100}
OP	MHA (VASWANI ET AL., 2017) FFN (VASWANI ET AL., 2017) BLOCK (VASWANI ET AL., 2017)
U-TUNER	ADAPTER (HOULSBY ET AL., 2019) PREFIX (LI & LIANG, 2021) PROMPT (JIA ET AL., 2022)
SCALING STRATEGIES	DIRECT CONNECTION, SCALAR SCALING, CHANNEL-WISE SCALING, INPUT-DEPENDENT SCALING
ARCHITECTURE	ViT/B-16 (DOSOVITSKIY ET AL., 2020) ViT/B-14 (DOSOVITSKIY ET AL., 2020)
PRE-TRAINED	IMAGENET-21K (DENG ET AL., 2009) MAE (HE ET AL., 2022B) CLIP (RADFORD ET AL., 2021)

A.2.3. THE CONSTRUCTION FORMS OF U-TUNING ON FGVC DATASETS

As described in Appendix A.2.2, we conduct experiments on different construction forms of U-Tuning by attaching the **U-Tuner** units to the OP units, MHA and FFN, respectively. Table 11 lists the construction forms for each FGVC dataset, which achieves the best performance in all the experiments above. We find that there is a discrepancy in the best construction form of U-Tuning for different datasets. For example, CUB_200_2011 benefits from the U-Tuning constructed with

Table 11. The construction forms of U-Tuning on FGVC datasets.

OP	CUB_200_2011	NABIRDS	OXFORDFLOWERS	STANFORDCARS	STANFORDDOGS
MHA	PROMPT	ADAPTER	ADAPTER	ADAPTER	PREFIX
FFN	ADAPTER	PREFIX	PROMPT	ADAPTER	ADAPTER

prompt on MHA and adapter on FFN, while OxfordFlowers prefers the U-Tuning constructed with adapter on MHA and prompt on FFN. This phenomenon further demonstrates that our proposed U-Tuning paradigm offers the shortcuts for the downstream tasks adaption of the pre-trained model.