

Gradient-based Parameter Selection for Efficient Fine-Tuning

Zhi Zhang^{1,2*}, Qizhe Zhang^{2*}, Zijun Gao², Renrui Zhang³, Ekaterina Shutova¹, Shiji Zhou⁴, Shanghang Zhang^{2†}

¹ILLC, University of Amsterdam,

²National Key Laboratory for Multimedia Information Processing,
School of Computer Science, Peking University,

³MMLAB, The Chinese University of Hong Kong,

⁴Department of Automation, Tsinghua University,

{z.zhang, e.shutova}@uva.nl, {theia, shanghang}@pku.edu.cn

Abstract

With the growing size of pre-trained models, full fine-tuning and storing all the parameters for various downstream tasks is costly and infeasible. In this paper, we propose a new parameter-efficient fine-tuning method, **Gradient-based Parameter Selection (GPS)**, demonstrating that only tuning a few selected parameters from the pre-trained model while keeping the remainder of the model frozen can generate similar or better performance compared with the full model fine-tuning method. Different from the existing popular and state-of-the-art parameter-efficient fine-tuning approaches, our method does not introduce any additional parameters and computational costs during both the training and inference stages. Another advantage is the model-agnostic and non-destructive property, which eliminates the need for any other design specific to a particular model. Compared with the full fine-tuning, GPS achieves 3.33% (91.78% vs. 88.45%, FGVC) and 9.61% (73.1% vs. 65.57%, VTAB) improvement of the accuracy with tuning only 0.36% parameters of the pre-trained model on average over 24 image classification tasks; it also demonstrates a significant improvement of 17% and 16.8% in mDice and mIoU, respectively, on medical image segmentation task. Moreover, GPS achieves state-of-the-art performance compared with existing PEFT methods. The code will be available in <https://github.com/FightingFighting/GPS.git>.

1. Introduction

The pre-training and fine-tuning pipeline has become a common paradigm for adapting large models pre-trained on substantial amounts of data to downstream tasks with fewer training samples. However, fine-tuning all the parameters in the model is memory-intensive and data-inefficient, which

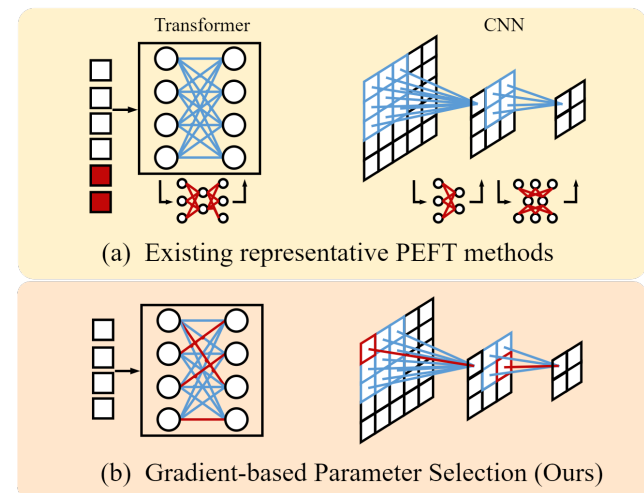


Figure 1. Comparison between our GPS and other PEFT methods. (a) Existing popular methods introduce extra parameters for tuning downstream tasks, which might need a special design for diverse architectures, such as appending prompt into the input token in Transformer or inserting different modules into different layers (b) Our approach avoids the introduction of additional parameters and solely fine-tunes the selected parameters from the model, employing a unified gradient-based parameter selection method across diverse architectural variations, e.g. Transformer and CNN.

is costly and infeasible for multiple downstream tasks given a large-scale model [36, 43, 58]. To tackle this issue, parameter-efficient fine-tuning (PEFT) methods have been proposed with the aim of tuning a minimal number of parameters to fit downstream tasks while keeping most of the parameters frozen. Another benefit of PEFT is that tuning a smaller set of parameters reduces the complexity of optimization and alleviates overfitting concerns when adapting large pre-trained models to downstream tasks with limited data, resulting in comparable or even superior performance

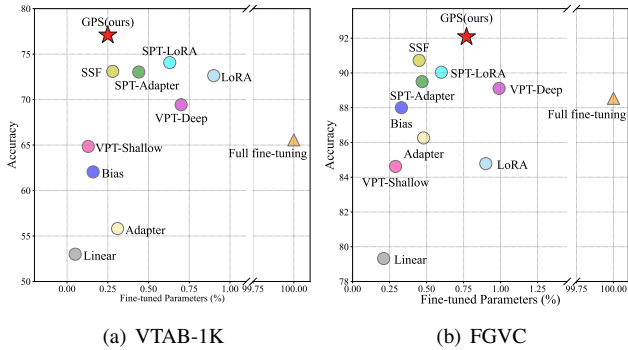


Figure 2. Performance comparisons of 11 fine-tuning methods with a pre-trained ViT-B/16 model on the VTAB-1k (a) and FGVC (b) benchmarks. Our GPS (red stars) achieves state-of-the-art performance on both benchmarks with only 0.25% and 0.77% average trainable parameters respectively.

compared to full fine-tuning [43]. Inspired by the success of PEFT in NLP [13, 31, 38, 40, 57, 79], several methods have been introduced to vision tasks, such as Adapter [36] and Visual Prompt Tuning (VPT) [43] introducing extra learnable parameters into the backbone and the input space of the pre-trained model respectively. SSF, another representative method, transforms features across layers of the pre-trained model using extra learnable layers [58].

However, these methods introduce additional parameters into the pre-trained model and disrupt its original architecture, leading to increased computational costs during training and/or inference stages. Furthermore, these approaches lack generalizability across various model architectures. Specifically, different models are equipped with distinct components (layers), such as MLPs, activation functions, and self-attention layers. These methods need to determine the optimal locations for inserting extra parameters between different layers; moreover, certain transformer-based techniques cannot be directly applied to convolution-based methods like VPT. Therefore, these methods exhibit limited compatibility with diverse architectures.

To tackle the above issues, we propose a non-destructive network architecture and model-agnostic PEFT approach, which introduces no extra parameters during both training and test stages and provides a unified solution for various architectures. We select a small number of essential parameters from the pre-trained model and only fine-tune these parameters for the downstream tasks. To select them, we propose a fine-grained Gradient-based Parameter Selection (GPS) method. For each neuron in the network, we choose top-K of its input connections (weights or parameters) with the highest gradient value, resulting in a small proportion of the parameters in the model being selected.

Such design offers five-fold benefits: i) The pre-trained

Method	Mean Params. Acc.	(%)	Model Agnostic	No extra Train param.	No extra Infer params.	Task Adaptive
Full [43]	70.36	100	✓	✓	✓	✗
Linear [43]	58.48	0.08	✓	✓	✓	✗
Bias [92]	67.54	0.20	✓	✓	✓	✗
Adapter [36]	60.04	0.35	✗	✗	✗	✗
VPT [43]	73.53	0.76	✗	✗	✗	✗
LoRA [38]	75.16	0.90	✗	✗	✓	✗
SSF [58]	76.77	0.32	✗	✗	✓	✗
GPS (ours)	78.64	0.36	✓	✓	✓	✓

Table 1. Comparison between different fine-tuning methods. The ViT-B/16 model accuracy over all 24 tasks in FGVA and VTAB fine-tuned on ViT-B/16 model and the number of tunable parameters are shown in columns Acc. and Params. (%).

model can efficiently tackle downstream tasks because the gradient direction indicates the fastest loss function changes and highest change rate, facilitating efficient gradient descent during model fine-tuning. We also provide a sparse regularized equivalent form for GPS, which indicates better generalization than full fine-tuning; ii) Each neuron within the network possesses the potential to adjust its activation state by fine-tuning selected input connections. Consequently, the pre-trained model exhibits flexibility in modifying features of varying granularities to suit diverse downstream tasks. For instance, when adapting a model pre-trained on ImageNet [10] for CIFAR-100 [75], it is necessary to refine high-level features; whereas for ImageNet-Sketch [85] adaptation, more detailed feature fine-tuning is required. iii) Our approach avoids introducing extra parameters and computational costs and keeps the architecture of the model intact; iv) The selection procedure enables its application across diverse models by adopting a neuron-based rather than a layer-based method, thereby eliminating the necessity for distinct designs for different layers in various models. v) Different from other methods using a pre-defined and consistent strategy for different tasks, our method adaptively selects parameters for each task by our proposed gradient strategy to better fit the domain-specific semantics of different downstream tasks. Please see the difference between our method with others in Fig. 1 and Tab. 1.

We evaluate our approach on a total of 27 visual tasks (including image classification and semantic segmentation) over 4 different model architectures. Our GPS achieves state-of-the-art performance compared to other PEFT methods and has a good balance between performance and the number of trainable parameters, as illustrated in Fig. 2. Compared with the full fine-tuning, GPS achieves 3.33% (FGVC) and 9.61% (VTAB) improvement of the accuracy while tuning only 0.36% parameters of the pre-trained model on average over 24 tasks; it also demonstrates a significant improvement of 17% and 16.8% in mDice and mIoU, respectively, on medical image segmentation task.

Moreover, we verify the effectiveness of our approach on different network architectures, such as Transformer and Convolutional Neural Networks. Furthermore, we compare GPS with various parameter selection methods and demonstrate its superior properties. GPS provides a new paradigm for PEFT and inspires deeper insights into this field.

2. Related work

Visual parameter efficient fine-tuning In general, there are typically two primary categories of PEFT. Addition-based methods introduce additional parameters to the pre-trained backbone. Adapters [20, 21, 36, 69, 70, 75, 78, 80, 86, 94–96] adopt a residual pathway and learn down and up projection with a nonlinear activation. Others [63] propose a hyper-network to generate model weights or decompose the dense weighted matrix into the low-rank matrix [46]. Prompt methods [12, 22, 39, 45, 57, 59, 60, 97, 101] wrap the input with context. VPT [43] prepend learnable prompts to the input tokens. SSF [58] achieves promising results by scaling and shifting the feature between layers. Selection-based methods select a subset of the parameters for tuning, such as only fine-tuning bias [92], last K layers [36, 43]. While traditionally considered less effective than addition-based methods, our approach of adaptively selecting parameters for each task yielded surprisingly strong results.

Sub-network training Pruning technique [19, 28, 29, 52, 56, 88] uncovers the importance of subnetworks. The lottery ticket hypothesis [17] articulates that subnetworks can reach the accuracy of the original model. Fine-tuning subnetworks are widely studied. SpotTune [27] designs a policy network to make routing decisions for subset networks. Child-tuning [90] iteratively updates a subset of parameters by masking out some gradients during the backward process. However, these methods are not aligned with the PEFT setting. We fix a small number of parameters and only tune them for fitting downstream to achieve PEFT.

3. Approach

Different from the currently popular methods introducing additional parameters to fine-tune the pre-trained model for downstream tasks [36, 43, 58], we select only a small number of parameters from the pre-trained model and then only update these parameters during the fine-tuning stage. Specifically, our method has two stages: parameter selection and masked fine-tuning. For each downstream task t , we first select a small portion of important parameters (task-specific parameters) from the original pre-trained model using a gradient-based method. We then fine-tune the pre-trained model for the task t , keeping all other unimportant parameters frozen and updating only selected parameters

using a sparse binary mask to set the gradient of unimportant parameters to zero (see Fig. 3).

3.1. Gradient-based parameter selection

Relevant studies have indicated that the pre-trained backbone exhibits diverse feature patterns at distinct parameter positions, and the same positions make varying contributions to fine-tuning various tasks [4, 53, 65, 74, 91]. Therefore, we posit that there exists an optimal subset of parameters for fine-tuning a pre-trained model to a downstream task. This subset is essential and necessary for fine-tuning the task, and the different tasks require a distinct subset. Formally, given a downstream task t with the dataset \mathcal{D}_t and a pre-trained model $\Theta = \{w_1, w_2, \dots, w_N\}$, we aim to find a subset of w , i.e. $w = \{w_1, w_2, \dots, w_n\}$ ($n \ll N$). We select parameters following two principles: 1) Important for downstream tasks; 2) Distributed over the whole network.

Importance for downstream tasks We identify the importance of parameters in a pre-trained model for a specific task by selecting those with the highest gradient value, which is obtained by calculating the gradient of a loss function with respect to its parameters. The intuition behind this is that the parameters with the largest gradient value indicate the loss function changes fastest along the gradient direction and has the greatest change rate, which facilitates efficient gradient descent during fine-tuning. Specifically, the gradient of the parameters is calculated by

$$\nabla \mathcal{L}_{\mathcal{D}_t}(\Theta) = \left[\frac{\partial \mathcal{L}}{\partial w_1} \dots \frac{\partial \mathcal{L}}{\partial w_N} \right]^\top \quad (1)$$

where $\mathcal{L}(w)$ is the loss function. Normally, when we fine-tune a pre-trained model on a downstream task, we need a new classification head (*i.e.* MLP) with random initialization. In order to avoid the adverse effects of these randomly initialized parameters on gradient calculation using the cross-entropy loss function, we use Supervised Contrastive Loss (SCL) [48] as the loss function for calculating the gradient during parameter selection, since it does not need to involve the head (We still use cross-entropy loss during fine-tuning stage). SCL is a variant of Contrastive Loss (CL) that aims to bring different augmented samples of the same image closer together in embedding space. In contrast, SCL tries to cluster samples from the same class together, which coincides with our target of the downstream classification tasks. Specifically, given a task with the dataset $\mathcal{D}_t = \{\mathbf{x}_i, y_i\}_{i=1 \dots K}$, SCL is calculated by

$$\begin{aligned} \mathcal{L}^{\text{scl}} &= \sum_{i \in \mathcal{D}_t} \mathcal{L}_i^{\text{scl}} \\ &= \sum_{i \in \mathcal{D}_t} \frac{-1}{|P(i)|} \sum_{p \in P(i)} \log \frac{\exp(\mathbf{z}_i \odot \mathbf{z}_p / \tau)}{\sum_{a \in A(i)} \exp(\mathbf{z}_i \odot \mathbf{z}_a / \tau)} \end{aligned} \quad (2)$$

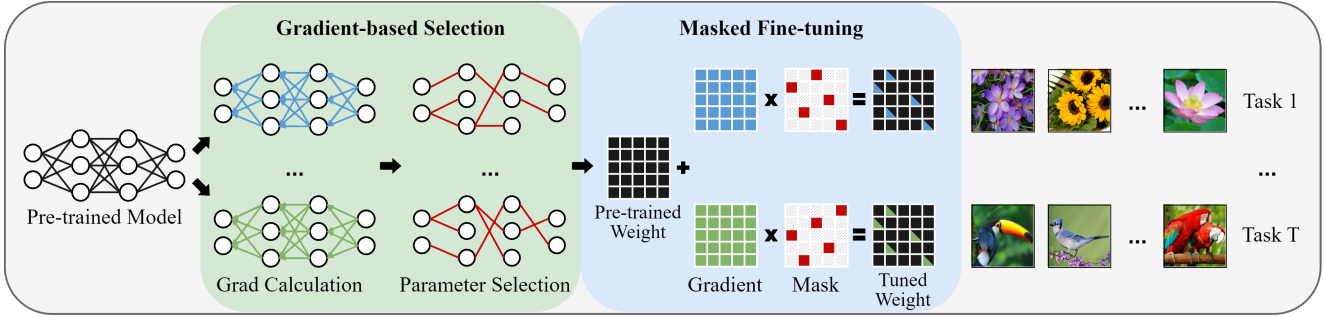


Figure 3. The overall pipeline of GPS. We first select a small portion of important parameters (sub-network) for each task from the original pre-trained model using a gradient-based method. Then only fine-tune the sub-network while keeping other parameters frozen.

where i represents i^{th} sample in \mathcal{D}_t ; $P(i) \equiv \{p \in A(i) : \tilde{y}_p = \tilde{y}_i\}$ is subset of \mathcal{D}_t , in which all samples have the same class with i ; $A(i) \equiv \mathcal{D}_t \setminus \{i\}$; z is the feature extracted from the pre-trained encoder and $\tau \in \mathcal{R}^+$ is a scalar temperature parameter.

Equivalent with sparse regularization In the above, we implicitly assume that the order of $\left(\left\|\frac{\partial \mathcal{L}}{\partial w_1}\right\| \cdots \left\|\frac{\partial \mathcal{L}}{\partial w_N}\right\|\right)$ is the same as $(\|w'_1 - w_1\| \cdots \|w'_N - w_N\|)$, which means selecting parameters with top- n gradient norm is the same as selecting top- n of the fine-tuning changes. Therefore, GPS captures the top- n important parameters for downstream tasks. The optimization objective can be rewritten as

$$\Theta' = \min \mathcal{L}(\Theta') \quad s.t. \quad \|\Theta' - \Theta\|_0 \leq n \quad (3)$$

where $\|\cdot\|_0$ is the l_0 norm and Θ' is the fine-tuned model. By Lagrangian duality, solving the above problem is equivalent to solving the following problem:

$$\mathcal{L}(\Theta') + \lambda \|\Theta' - \Theta\|_0 \quad (4)$$

with appropriate λ . Hence, GPS can be reviewed as a sparse regularized fine-tuning, which may lead to better generalization. Fu et al. [18] demonstrate that Eq. (4) has smaller generalization bound than pure optimization toward \mathcal{L} with full fine-tuning, resulting in better performance.

Distribution over the whole network A simple idea for parameter selection is to select a certain percentage of parameters with the highest gradient from the entire network. Our experiments have shown that with this idea, the majority of the selected parameters are located in the top layers of the network (see Supplementary for details), which is consistent with the findings reported in [36, 37]. However, solely fine-tuning these top-layer parameters is insufficient to mitigate the impact of the pre-trained model's own inductive bias, particularly when there exist substantial disparities in data distributions between upstream and downstream

tasks, which need to fine-tune more detailed features from shallower layers. Motivated by various studies indicating the distinct roles played by different components of neural networks [3, 15, 72, 77, 87], we posit that when fine-tuning a pre-trained model for downstream tasks, the adjusted parameters should be distributed throughout the entire network. The reason behind this lies in the ability of the model to adapt the information stored in parameters at different levels of granularity to fit downstream tasks. Therefore, our strategy is that for each neuron in the network, we select top- K (at least one) connections (weights) among all the input connections of the neuron, as shown in Fig. 3. By doing so, every neuron within the network possesses the potential to fine-tune its activation state rather than solely adjust high-level information in the top layers. In other words, our approach fine-tunes the detailed information stored in each neuron of the network, which better fits the downstream task during the fine-tuning stage. Our exploratory experiment further substantiates this assertion, as shown in Tab. 6(a).

Combining the two points above, we first calculate the gradient of the loss with respect to all the weights in the models for a specific task. Then for each neuron in the network, we select top- K connections with the highest gradient value (the modulus of gradient) among all input connections to the neuron. Doing so can not only ensure that important parameters for downstream tasks are chosen and allow the model to tune the activation state of all neurons for better fitting of downstream tasks. Another benefit of this selection procedure is its ease of application across various model architectures, such as Transformer and CNN, avoiding any model-specific design. Our experiments also demonstrate the effectiveness of our approach across diverse architectures, as shown in Tab. 2 and Tab. 4.

3.2. Masked fine-tuning

After parameter selection for a specific task, we fine-tune the pre-trained model on the task. During fine-tuning, we only update the selected parameters while keeping the remaining parameters of the pre-trained model frozen. As our

selected parameters are distributed across all neurons in every layer, only a few parameters within a specific weight matrix of the network are chosen, resulting in the updated matrix being sparse. Therefore, we utilize a mask to help with the sparse training. Specifically, for j^{th} weight matrix $\mathbf{W}_j \in \mathbb{R}^{d_{in} \times d_{out}}$ in the network, we build a same size of binary mask $\mathbf{M}_j \in \mathbb{R}^{d_{in} \times d_{out}}$:

$$\mathbf{M}_j = \begin{cases} 1, & w_j^k \in \mathbf{w} \\ 0, & w_j^k \notin \mathbf{w} \end{cases} \quad (5)$$

where w_j^k represents k^{th} element in j^{th} weight matrix. For each element in \mathbf{M}_j , its value is set to 1 if the corresponding parameter in \mathbf{W}_j is selected, and 0 otherwise. Then the weight matrix is updated by

$$\mathbf{W}_j \leftarrow \mathbf{W}_j - \epsilon \nabla \mathcal{L}(\mathbf{W}_j) \odot \mathbf{M}_j \quad (6)$$

where $\nabla \mathcal{L}(\mathbf{W}_j)$ is the gradient of the cross-entropy loss with respect to \mathbf{W}_j . As a result, the gradients of unselected parameters are zeroed out and excluded from updates, while only a small number of our selected parameters are updated during fine-tuning for downstream tasks. Please see Fig. 3 for a visualization of our method.

4. Experiments

We evaluate GPS on various downstream tasks, including image classification tasks and semantic segmentation tasks with different architectures. First, we briefly introduce our experimental settings, including datasets, backbones, and baselines. Then we demonstrate the effectiveness and universality of GPS. Moreover, we systematically study the impacts of different selection schemes and conduct comprehensive ablation experiments.

4.1. Experimental settings

Datasets Following VPT [43] and SSF [58], we evaluate our GPS method on a series of datasets categorized into three groups: i) **FGVC**: Fine-Grained Visual Classification (FGVC) benchmark includes 5 downstream tasks, which are CUB-200-2011 [84], NABirds [81], Oxford Flowers [67], Stanford Dogs [47] and Stanford Cars [23]. ii) **VTAB-1k**: Visual Task Adaptation Benchmark [93] (VTAB) contains 19 visual classification tasks which are grouped into three sets: Natural, Specialized, and Structured. iii) **CIFAR-100** [51] and **ImageNet-1k** [10]: widely use for general image classification task.

Backbones For a fair comparison, we follow VPT and SSF by using ViT-B/16 [14] pre-trained on ImageNet-21K [10] for the main image classification experiments. Moreover, to demonstrate the universality of our GPS, we also explore other backbones, including Swin Transformer [61] and ConvNeXt-B [62] for another variant of

Dataset	CUB-2011	NA-Birds	Oxford Flowers	Stan. Dogs	Stan. Cars	Mean Acc.	Params. (%)
Full [43]	87.3	82.7	98.8	89.4	84.5	88.54	100.00
Linear [43]	85.3	75.9	97.9	86.2	51.3	79.32	0.21
Bias [92]	88.4	84.2	98.8	91.2	79.4	88.40	0.33
Adapter [36]	87.1	84.3	98.5	89.8	68.6	85.66	0.48
LoRA [38]	85.6	79.8	98.9	87.6	72.0	84.78	0.90
VPT-Shallow [43]	86.7	78.8	98.4	90.7	68.7	84.62	0.29
VPT-Deep [43]	88.5	84.2	99.0	90.2	83.6	89.11	0.99
SSF [58]	<u>89.5</u>	<u>85.7</u>	<u>99.6</u>	89.6	<u>89.2</u>	<u>90.72</u>	0.45
SPT-Adapter [30]	89.1	83.3	99.2	91.1	86.2	89.78	0.47
SPT-LoRA [30]	88.6	83.4	99.5	<u>91.4</u>	87.3	90.04	0.60
GPS (Ours)	89.9	86.7	99.7	92.2	90.4	91.78	0.77

Table 2. Performance comparisons on FGVC with ViT-B/16 models pre-trained on ImageNet-21K.

Transformer-based and CNN-based architecture, respectively. In addition, we finetune semantic segmentation tasks on SAM [50], a strong segmentation foundation model.

Baselines We compare our GPS with a variety of fine-tuning protocols that can be mainly categorized into three types: i) **Full**: Full fine-tuning is the most commonly used protocol updating all parameters of the whole model during tuning. ii) **Selection-based**: This kind of method selects a subset of parameters in the original model for fine-tuning, including linear probing and Bias [92]. Such methods are easy to implement and require no extra computations but have not worked well. Our method belongs to this group and achieves the best performance while ensuring convenience and universality. iii) **Addition-based**: This kind of method adds new trainable parameters to the backbone, including Adapter [36], VPT [43] and SPT-Adapter [30]. Such methods require extra computations in both the training and inference stages. Other methods like LoRA [38], SSF [58], and SPT-LoRA [30] also add new tunable parameters during the training stage, but these parameters can be reparameterized into the backbone during testing.

Implementation details We follow SSF to process the images in all the FGVC, VTAB-1k and CIFAR-100 datasets. We employ the Adam [49] optimizer with cosine learning rate decay to fine-tune models for 100 epochs, and the linear warm-up is used in the first 10 epochs. All experiments are conducted on the NVIDIA A100 GPU.

4.2. Performance on image classification

We present a comprehensive evaluation of the effectiveness of our GPS by comparing it against multiple baselines on 3 benchmarks, comprising a total of 26 datasets. In addition to common benchmarks (FGVC and VTAB-1k), we also compare our method with others on different architectures. We evaluate the performance and effectiveness by Top-1 accuracy (%) and the number of fine-tuned parameters.

Method	Dataset	Natural							Specialized				Structured						VTAB			
		CIFAR-100	Caltech101	DTD	Flowers102	Pets	SVHN	Sun397	Patch Camelyon	EuroSAT	Resisc45	Retinopathy	Clevr/count	Clevr/distance	DMLab	KITTI/distance	dSprites/loc	dSprites/ori	SmallNORB/azi	SmallNORB/ele	Mean Acc.	Mean Params. (%)
Full [43]		68.9	87.7	64.3	97.2	86.9	87.4	38.8	79.7	95.7	84.2	73.9	56.3	58.6	41.7	65.5	57.5	46.7	25.7	29.1	65.57	100.00
Linear [43]		63.4	85.0	64.3	97.0	86.3	36.6	51.0	78.5	87.5	68.6	74.0	34.3	30.6	33.2	55.4	12.5	20.0	9.6	19.2	53.00	0.05
Bias [92]		72.8	87.0	59.2	97.5	85.3	59.9	51.4	78.7	91.6	72.9	69.8	61.5	55.6	32.4	55.9	66.6	40.0	15.7	25.1	62.05	0.16
Adapter [36]		74.1	86.1	63.2	97.7	87.0	34.6	50.8	76.3	88.0	73.1	70.5	45.7	37.4	31.2	53.2	30.3	25.4	13.8	22.1	55.82	0.31
LoRA [38]		68.1	91.4	69.8	99.0	90.5	86.4	53.1	85.1	95.8	84.7	74.2	83.0	66.9	50.4	81.4	80.2	46.6	32.2	41.1	72.63	0.90
VPT-Shallow [43]		77.7	86.9	62.6	97.5	87.3	74.5	51.2	78.2	92.0	75.6	72.9	50.5	58.6	40.5	67.1	68.7	36.1	20.2	34.1	64.85	0.13
VPT-Deep [43]		78.8	90.8	65.8	98.0	88.3	78.1	49.6	81.8	96.1	83.4	68.4	68.5	60.0	46.5	72.8	73.6	47.9	32.9	37.8	69.43	0.70
SSF [58]		69.0	92.6	75.1	99.4	91.8	90.2	52.9	87.4	95.9	87.4	75.5	75.9	62.3	53.3	80.6	77.3	54.9	29.5	37.9	73.10	0.28
SPT-ADAPTER [30]		72.9	93.2	72.5	99.3	91.4	88.8	55.8	86.2	96.1	85.5	75.5	83.0	68.0	51.9	81.2	51.9	31.7	41.2	61.4	73.03	0.44
SPT-LoRA [30]		73.5	93.3	72.5	99.3	91.5	87.9	55.5	85.7	96.2	85.9	75.9	84.4	67.6	52.5	82.0	81.0	51.1	30.2	41.3	74.07	0.63
GPS (Ours)		81.1	94.2	75.8	99.4	91.7	91.6	52.4	87.9	96.2	86.5	76.5	79.9	62.6	55.0	82.4	84.0	55.4	29.7	46.1	75.18	0.25

Table 3. Performance comparisons on VTAB-1k with ViT-B/16 models pre-trained on ImageNet-21K.

Image classification performance As shown in Tab. 2 and Tab. 3, our GPS outperforms all other fine-tuning methods by a large margin on both FGVC and VTAB benchmarks, sufficiently demonstrating that our method of parameter selection is a simple yet effective way for model tuning. On FGVC, GPS outperforms all other fine-tuning methods, including full fine-tuning, on all 5 tasks. It obtains 1.02% and 3.24% accuracy improvement of the mean accuracy compared to the previous SOAT method SSF [58] and full fine-tuning, while it only uses 0.77% of trainable parameters. On VTAB, GPS also outperforms all other fine-tuning methods. Specifically, it obtains 1.11% and 9.61% improvement of the mean accuracy on 19 VTAB tasks compared to the previous SOAT method SPT-LoRA [30] and full fine-tuning. GPS beats the previous SOTA by 1.75%, 0.23%, and 0.63% in the Natural, Specialized and Structured subsets, respectively. Meanwhile, GPS also uses fewer trainable parameters compared to VPT-Deep, SSF, and SPT-LoRA (0.25% vs. 0.70%, 0.28% and 0.63%), which further illustrates the high efficiency of our approach. For most tasks, we exclusively select the top 1 input connection for each neuron; however, for more challenging tasks, multiple connections are chosen (see Supplement for details). The percentage of learnable parameters in our GPS can be explicitly controlled by adjusting the number of connections selected, allowing for a balance between parameter count and performance in tasks.

Generalization on different architectures Since our method only selects a subset of parameters from the pre-trained model for fine-tuning, it is naturally model-agnostic. We compare GPS with other representative methods across ViT-B/16 (Tab. 2), Swin-B and ConvNeXt-B architectures on the FGVC dataset (Tab. 4), CIFAR-100 and ImageNet-1k (Please see full results in Supplementary). Among all three architectures, GPS consistently outperforms all other

Architecture	Swin-B		ConvNeXt-B	
	Ave. Acc.	Params.(%)	Ave. Acc.	Params.(%)
Full [43]	92.42	100.00	93.04	100.00
Linear [43]	87.90	0.28	88.00	0.28
SSF [58]	91.54	0.56	92.48	0.56
GPS (Ours)	92.56	0.95	93.32	0.90

Table 4. Performance comparisons on FGVC benchmark (Average accuracy over 5 tasks) with different model architectures.

baselines, demonstrating its model-agnostic advantage. The Swin and Convnext have more complex designs than ViT, enabling them to acquire comprehensive and high-quality features during pre-training. Consequently, even the simplest linear probing method yields commendable results on these two architectures, reducing the effectiveness of the PEFT method and causing the previous SOTA SSF to underperform Full. However, in this scenario, our GPS still maintains a lead over Full with gains of 0.12% and 0.28%, respectively, further showing the effectiveness of our method.

Computational cost In Fig. 4, we compare the computational cost of GPS with other fine-tuning methods to demonstrate the efficiency of our approach. Following SSF [58], we reimplement VPT [43] with 200 and 50 prompts for the shallow and deep versions, respectively. A batch size of 32 is used in both the training and inference stages. For a fair comparison, for all experiments, we do not use mixed precision training, which was used in SSF. All metrics are measured on a single NVIDIA A100 GPU. In the training stage, GPS has less time and memory consumption than both VPT and SSF. Compared with full fine-tuning, GPS has a much lower time overhead and a similar memory overhead, but it leads to an increased performance by a large margin. Since GPS is a selection-based method, it does not introduce any additional parameters, so it can achieve the same minimum

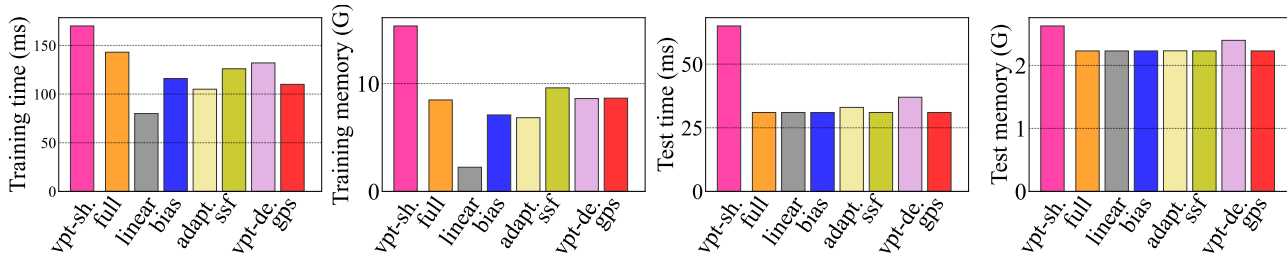


Figure 4. Computational cost of different tuning methods. From left to right: training time, training memory, test time, and test memory. Training/Test time is the time consumed by a mini-batch.

Method	mDice (\uparrow)	mIoU (\uparrow)	Params. (M)
Full [43]	71.1	55.7	93.8
Linear [43]	71.6	46.6	4.06
Bias [92]	86.5	69.1	4.16
Adapter [6]	84.8	66.7	4.12
SSF [58]	87.3	71.7	4.26
GPS (Ours)	88.1	72.5	4.22

Table 5. Quantitative Result for Polyp Segmentation

time and memory overhead as full fine-tuning during inference without any reparameterization operation, which is much lower than the addition-based Adapter and VPT.

4.3. Semantic segmentation

In addition to visual classification tasks, we also explore our method for the task of semantic segmentation. Segment Anything Model (SAM) [50] is a strong foundation model for segmentation. It is pre-trained on a large-scale dataset enabling powerful generalization. However, several studies, e.g. [6], have reported poor performance of SAM on medical segmentation tasks such as polyp segmentation [41]. To address this limitation, they proposed employing Adapter to effectively fine-tune SAM for downstream medical segmentation tasks. Following their experimental setup, we applied our method to SAM and conducted a comparative analysis with other PEFT approaches. Our GPS yielded exceptional results, as shown in Tab. 5 and visually depicted in Fig. 5 (See Supplementary for more case visualization).

4.4. Impacts of different selection schemes

Different selection levels Our GPS selects trainable parameters at the neuron level, i.e. selecting top-k input connection per neuron. We also investigate parameter selection methods at different levels. As shown in Tab. 6 (a), *Net* and *Layer* represent selecting a certain proportion of the parameters with the highest gradient based on the entire network and each layer, respectively. For a fair comparison, we keep the same number of parameters selected over these levels.

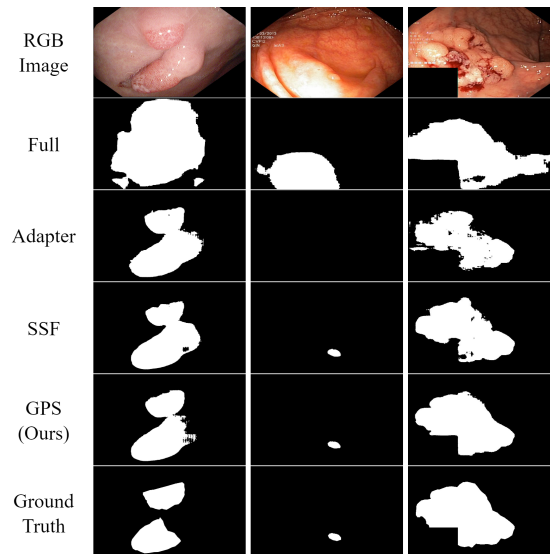


Figure 5. The Visualization of Polyp segmentation task. Our GPS can still handle difficult segmentation cases compared with others.

We can see that the finer the granularity of selection, the better the performance. For example, the accuracy on CUB increases by 0.44% and 0.77% when selection level changes from network to layer, and from layer to neuron.

Different selection criteria We further study the effectiveness of our gradient-based selection method by comparing different selection criteria. As shown in Tab. 6 (b), *Net Random* and *Neuron Random* means randomly selecting top-K the input connection for each neuron and selecting the same number of parameters based on the whole network respectively. *Magnitude* means selecting top-K input connections with the largest weight per neuron. As we can see, the increase in the randomness of parameter selection causes a decrease in performance (*Net Random* < *Neuron Random*). The result of *Magnitude* is similar to *Neuron Random*, demonstrating neuron-level selection is crucial.

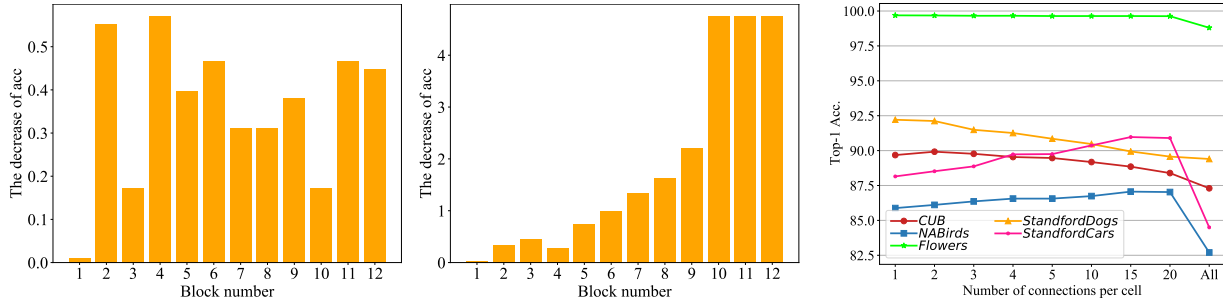


Figure 6. Impacts of different selection locations and quantities. From left to right: (a) Performance drop caused by not selecting parameters from k-th blocks. (b) And by not selecting from the top k blocks. (c) Impacts of different numbers of selected connections on performance.

	Dataset	CUB	NABirds	Flowers	Cars	Dogs
(a)	Net	86.86	86.55	99.62	89.65	91.32
	Layer	87.30	86.79	99.64	90.03	91.90
(b)	Net Random	86.60	85.98	99.61	89.10	91.34
	Neuron Random	87.17	86.02	99.62	89.52	91.82
	Magnitude	87.29	85.99	99.62	89.29	91.30
(c)	Head+CE	87.05	86.20	99.64	89.25	91.29
	GPS	88.07 ±0.11	86.64 ±0.03	99.69 ±0.01	90.10 ±0.10	92.30 ±0.10

Table 6. The result on FGVC for investigating impacts of different selection schemes and ablations. (a) Different selection levels. (b) Different selection criteria. (c) Gradient calculating method.

Different selection location To investigate the impact of selected parameters located at different layers within the network, we conducted experiments using the ViT-B/16 model fine-tuning on CUB and evaluated accuracy degradation when applying our GPS method to select parameters from the entire network except for a specific transformer block or previous several transformer blocks. As shown in Fig. 6(a), it is surprising to note that when we do not select parameters from a specific block, the biggest drop in the accuracy comes from the shallow layers (block 2 and block 4). This finding supports our GPS approach that selects parameters from the entire network rather than just the last few layers. When we do not select the parameters from the first specific number of blocks, it is observed that the accuracy drop increases with more blocks removed (Fig. 6(b)).

4.5. Ablation study

Head-free contrastive loss To obtain more accurate gradients for selecting parameters, inspired by the representation learning pre-training methods, Our GPS adopts the supervised contrastive loss to calculate gradient (without random initialization of the classification head). As shown in Tab. 6 (c), when we use the cross-entropy loss (with the head) to calculate the gradient, the average accuracy on FGVC is dropped by 0.67%, illustrating the importance of

obtaining accurate gradients.

Selected connection number As shown in Fig. 6(c) we select top-K input connections per neuron as trainable parameters, ranging from 1 to 15, and conduct experiments on the 5 tasks. We can observe that more trainable parameters do not necessarily lead to better performance, but each data set has a performance peak. In addition, on the dataset with sufficient training data, the addition of trainable parameters can greatly improve the accuracy. Our GPS can easily control the number of trainable parameters and achieve optimal results on each dataset.

Robustness to seeds Addition-based fine-tuning methods like VPT are sensitive to the initialization of additional parameters as well as random seeds, whereas select-based methods are not. All results in Tab. 6 are the average accuracy of three seeds on FGVC datasets (Only shows the std of GPS here. Please see details in supplementary). The results show random seeds have little influence on our method.

5. Conclusion

In this paper, we propose a new paradigm for PEFT, *i.e.* Gradient-based Parameter Selection (GPS). Our approach does not introduce any additional parameters and only fine-tunes a small subset of the pre-trained model’s parameters for downstream tasks, resulting in robust generalization across diverse models and adaptively selecting a subset of parameters for each task. Remarkably, GPS achieves significant improvement on a range of tasks (including image classification and semantic segmentation), compared to the full fine-tuning method. GPS also attains SOTA performance compared to other PEFT methods.

Acknowledgement

Shanghang Zhang is supported by the National Science and Technology Major Project of China (No. 2022ZD0117801).

References

- [1] Ankur Bapna, Naveen Arivazhagan, and Orhan Firat. Simple, scalable adaptation for neural machine translation. *arXiv preprint arXiv:1909.08478*, 2019. 8
- [2] Charles Beattie, Joel Z Leibo, Denis Teplyashin, Tom Ward, Marcus Wainwright, Heinrich Küttler, Andrew Lefrancq, Simon Green, Víctor Valdés, Amir Sadik, et al. Deepmind lab. *arXiv preprint arXiv:1612.03801*, 2016. 7
- [3] Jiezhong Cao, Jincheng Li, Xiping Hu, Xiangmiao Wu, and Mingkui Tan. Towards interpreting deep neural networks via layer behavior understanding. *Machine Learning*, 111: 1159–1179, 2022. 4
- [4] Niladri S Chatterji, Behnam Neyshabur, and Hanie Sedghi. The intriguing role of module criticality in the generalization of deep networks. *arXiv preprint arXiv:1912.00528*, 2019. 3
- [5] Tianlong Chen, Jonathan Frankle, Shiyu Chang, Sijia Liu, Yang Zhang, Zhangyang Wang, and Michael Carbin. The lottery ticket hypothesis for pre-trained bert networks. *ArXiv*, abs/2007.12223, 2020. 9
- [6] Tianrun Chen, Lanyun Zhu, Chaotao Ding, Runlong Cao, Shangzhan Zhang, Yan Wang, Zejian Li, Lingyun Sun, Papa Mao, and Ying Zang. Sam fails to segment anything? – sam-adapter: Adapting sam in underperformed scenes: Camouflage, shadow, and more, 2023. 7
- [7] 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. 7
- [8] Mircea Cimpoi, Subhansu 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. 7
- [9] Y. L. Cun, J. S. Denker, and S. A. Solla. *Optimal brain damage*. Advances in neural information processing systems 2, 1990. 9
- [10] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 2, 5, 7
- [11] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. 8
- [12] Ning Ding, Shengding Hu, Weilin Zhao, Yulin Chen, Zhiyuan Liu, Hai-Tao Zheng, and Maosong Sun. Openprompt: An open-source framework for prompt-learning. *arXiv preprint arXiv:2111.01998*, 2021. 3, 8
- [13] Ning Ding, Yujia Qin, Guang Yang, Fuchao Wei, Zonghan Yang, Yusheng Su, Shengding Hu, Yulin Chen, Chi Min Chan, and Weize and Chen. Parameter-efficient fine-tuning of large-scale pre-trained language models. *Nature Machine Intelligence*, 5(3):220–235, 2023. 2
- [14] 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. *arXiv preprint arXiv:2010.11929*, 2020. 5, 1, 8
- [15] Fenglei Fan, Jinjun Xiong, Mengzhou Li, and Ge Wang. On interpretability of artificial neural networks: A survey. *IEEE Transactions on Radiation and Plasma Medical Sciences*, 5:741–760, 2020. 4
- [16] Li Fei-Fei, Robert Fergus, and Pietro Perona. One-shot learning of object categories. *IEEE transactions on pattern analysis and machine intelligence*, 28(4):594–611, 2006. 7
- [17] Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. *arXiv preprint arXiv:1803.03635*, 2018. 3, 8
- [18] Zihao Fu, Haoran Yang, Anthony Man-Cho So, Wai Lam, Lidong Bing, and Nigel Collier. On the effectiveness of parameter-efficient fine-tuning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 12799–12807, 2023. 4
- [19] Trevor Gale, Erich Elsen, and Sara Hooker. The state of sparsity in deep neural networks. *arXiv preprint arXiv:1902.09574*, 2019. 3, 8
- [20] Peng Gao*, Shijie Geng*, Renrui Zhang*, Teli Ma, Rongyao Fang, Yongfeng Zhang, Hongsheng Li, and Yu Qiao. Clip-adapter: Better vision-language models with feature adapters. *IJCV 2023*, 2021. 3
- [21] Peng Gao, Jiaming Han, Renrui Zhang, Ziyi Lin, Shijie Geng, Aojun Zhou, Wei Zhang, Pan Lu, Conghui He, Xiangyu Yue, et al. Llama-adapter v2: Parameter-efficient visual instruction model. *arXiv preprint arXiv:2304.15010*, 2023. 3
- [22] Tianyu Gao, Adam Fisch, and Danqi Chen. Making pre-trained language models better few-shot learners. *arXiv preprint arXiv:2012.15723*, 2020. 3, 8
- [23] Timnit Gebru, Jonathan Krause, Yilun Wang, Duyun Chen, Jia Deng, and Li Fei-Fei. Fine-grained car detection for visual census estimation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2017. 5, 6, 7
- [24] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The kitti dataset. *The International Journal of Robotics Research*, 32(11):1231–1237, 2013. 7
- [25] Ben Graham. Kaggle diabetic retinopathy detection competition report. *University of Warwick*, pages 24–26, 2015. 7
- [26] Demi Guo, Alexander M Rush, and Yoon Kim. Parameter-efficient transfer learning with diff pruning. *arXiv preprint arXiv:2012.07463*, 2020. 8
- [27] Yunhui Guo, Honghui Shi, Abhishek Kumar, Kristen Grauman, Tajana Rosing, and Rogerio Feris. Spottune: transfer learning through adaptive fine-tuning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4805–4814, 2019. 3, 8
- [28] Song Han, Huizi Mao, and William J Dally. Deep compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding. *arXiv preprint arXiv:1510.00149*, 2015. 3, 8
- [29] Song Han, Jeff Pool, John Tran, and William Dally. Learning both weights and connections for efficient neural net-

- work. *Advances in neural information processing systems*, 28, 2015. 3, 8
- [30] Haoyu He, Jianfei Cai, Jing Zhang, Dacheng Tao, and Bohan Zhuang. Sensitivity-aware visual parameter-efficient tuning. *arXiv preprint arXiv:2303.08566*, 2023. 5, 6, 2
- [31] Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. Towards a unified view of parameter-efficient transfer learning. *ArXiv*, abs/2110.04366, 2021. 2
- [32] 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. 7
- [33] Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *arXiv preprint arXiv:1903.12261*, 2019. 7, 8
- [34] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A critical analysis of out-of-distribution generalization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8340–8349, 2021. 7, 8
- [35] Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial examples. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15262–15271, 2021. 7, 8
- [36] 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. 1, 2, 3, 4, 5, 6, 8
- [37] Jeremy Howard and Sebastian Ruder. Universal language model fine-tuning for text classification. In *Annual Meeting of the Association for Computational Linguistics*, 2018. 4
- [38] 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*, 2021. 2, 5, 6
- [39] Shengding Hu, Ning Ding, Huadong Wang, Zhiyuan Liu, Jingang Wang, Juanzi Li, Wei Wu, and Maosong Sun. Knowledgeable prompt-tuning: Incorporating knowledge into prompt verbalizer for text classification. *arXiv preprint arXiv:2108.02035*, 2021. 3, 8
- [40] Shengding Hu, Zhen Zhang, Ning Ding, Yadao Wang, Yasheng Wang, Zhiyuan Liu, and Maosong Sun. Sparse structure search for parameter-efficient tuning, 2022. 2
- [41] Debesh Jha, Steven A. Hicks, Krister Emanuelsen, Håvard Johansen, Dag Johansen, Thomas de Lange, Michael A. Riegler, and Pål Halvorsen. Medico multimedia task at mediaeval 2020: Automatic polyp segmentation, 2020. 7, 4, 8
- [42] Debesh Jha, Pia H Smedsrud, Michael A Riegler, Pål Halvorsen, Thomas de Lange, Dag Johansen, and Håvard D Johansen. Kvasir-seg: A segmented polyp dataset. In *International Conference on Multimedia Modeling*, pages 451–462. Springer, 2020. 8
- [43] Menglin Jia, Luming Tang, Bor-Chun Chen, Claire Cardie, Serge Belongie, Bharath Hariharan, and Ser-Nam Lim. Visual prompt tuning. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXIII*, pages 709–727. Springer, 2022. 1, 2, 3, 5, 6, 7, 4, 8
- [44] Justin Johnson, Bharath Hariharan, Laurens Van Der Maaten, Li Fei-Fei, C Lawrence Zitnick, and Ross Girshick. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2901–2910, 2017. 7
- [45] Chen Ju, Tengda Han, Kunhao Zheng, Ya Zhang, and Weidi Xie. Prompting visual-language models for efficient video understanding. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXV*, pages 105–124. Springer, 2022. 3, 8
- [46] 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. 3, 6, 8, 9
- [47] Aditya Khosla, Nityananda Jayadevaprakash, Bangpeng Yao, and Fei-Fei Li. Novel dataset for fine-grained image categorization: Stanford dogs. In *Proc. CVPR workshop on fine-grained visual categorization (FGVC)*. Citeseer, 2011. 5, 6, 7
- [48] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschiot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. *ArXiv*, abs/2004.11362, 2020. 3
- [49] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 5
- [50] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick. Segment anything. *arXiv:2304.02643*, 2023. 5, 7, 4
- [51] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. 5, 2, 7
- [52] John K Kruschke and Javier R Movellan. Benefits of gain: Speeded learning and minimal hidden layers in back-propagation networks. *IEEE Transactions on systems, Man, and Cybernetics*, 21(1):273–280, 1991. 3, 8
- [53] Ananya Kumar, Aditi Raghunathan, Robbie Jones, Tengyu Ma, and Percy Liang. Fine-tuning can distort pretrained features and underperform out-of-distribution. *arXiv preprint arXiv:2202.10054*, 2022. 3
- [54] Yann LeCun, Fu Jie Huang, and Leon Bottou. Learning methods for generic object recognition with invariance to pose and lighting. In *Proceedings of the 2004 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2004. CVPR 2004.*, pages II–104. IEEE, 2004. 7

- [55] Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for efficient convnets. *ArXiv*, abs/1608.08710, 2016. 9
- [56] Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for efficient convnets. *arXiv preprint arXiv:1608.08710*, 2016. 3, 8
- [57] Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. *arXiv preprint arXiv:2101.00190*, 2021. 2, 3, 8
- [58] Dongze Lian, Daquan Zhou, Jiashi Feng, and Xinchao Wang. Scaling & shifting your features: A new baseline for efficient model tuning. *arXiv preprint arXiv:2210.08823*, 2022. 1, 2, 3, 5, 6, 7, 4, 8
- [59] Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Computing Surveys*, 55(9): 1–35, 2023. 3, 8
- [60] Xiao Liu, Kaixuan Ji, Yicheng Fu, Weng Tam, Zhengxiao Du, Zhilin Yang, and Jie Tang. P-tuning: Prompt tuning can be comparable to fine-tuning across scales and tasks. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 61–68, 2022. 3, 8
- [61] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10012–10022, 2021. 5, 8
- [62] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11976–11986, 2022. 5
- [63] Rabeeh Karimi Mahabadi, Sebastian Ruder, Mostafa Dehghani, and James Henderson. Parameter-efficient multi-task fine-tuning for transformers via shared hypernetworks. *arXiv preprint arXiv:2106.04489*, 2021. 3, 8
- [64] Loic Matthey, Irina Higgins, Demis Hassabis, and Alexander Lerchner. dsprites: Disentanglement testing sprites dataset, 2017. 7
- [65] Muhammad Muzammal Naseer, Kanchana Ranasinghe, Salman H Khan, Munawar Hayat, Fahad Shahbaz Khan, and Ming-Hsuan Yang. Intriguing properties of vision transformers. *Advances in Neural Information Processing Systems*, 34:23296–23308, 2021. 3
- [66] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bisacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. 2011. 7
- [67] Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of classes. In *2008 Sixth Indian Conference on Computer Vision, Graphics & Image Processing*, pages 722–729. IEEE, 2008. 5, 6, 7
- [68] 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. 7
- [69] Jonas Pfeiffer, Aishwarya Kamath, Andreas Rücklé, Kyunghyun Cho, and Iryna Gurevych. Adapterfusion: Non-destructive task composition for transfer learning. *arXiv preprint arXiv:2005.00247*, 2020. 3, 6, 8, 9
- [70] Jonas Pfeiffer, Andreas Rücklé, Clifton Poth, Aishwarya Kamath, Ivan Vulić, Sebastian Ruder, Kyunghyun Cho, and Iryna Gurevych. Adapterhub: A framework for adapting transformers. *arXiv preprint arXiv:2007.07779*, 2020. 3, 8
- [71] Sai Prasanna, Anna Rogers, and Anna Rumshisky. When bert plays the lottery, all tickets are winning. In *Conference on Empirical Methods in Natural Language Processing*, 2020. 9
- [72] John G. Proakis and Dimitris G. Manolakis. Digital signal processing: Principles, algorithms, and applications. 1992. 4
- [73] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 8
- [74] Maithra Raghu, Thomas Unterthiner, Simon Kornblith, Chiyuan Zhang, and Alexey Dosovitskiy. Do vision transformers see like convolutional neural networks? *Advances in Neural Information Processing Systems*, 34: 12116–12128, 2021. 3
- [75] Sylvestre-Alvise Rebuffi, Hakan Bilen, and Andrea Vedaldi. Learning multiple visual domains with residual adapters. *Advances in neural information processing systems*, 30, 2017. 2, 3, 8
- [76] Andreas Rücklé, Gregor Geigle, Max Glockner, Tilman Beck, Jonas Pfeiffer, Nils Reimers, and Iryna Gurevych. Adapterdrop: On the efficiency of adapters in transformers. *arXiv preprint arXiv:2010.11918*, 2020. 8
- [77] Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. Learning important features through propagating activation differences. In *International Conference on Machine Learning*, 2017. 4
- [78] Asa Cooper Stickland and Iain Murray. Bert and pals: Projected attention layers for efficient adaptation in multi-task learning. In *International Conference on Machine Learning*, pages 5986–5995. PMLR, 2019. 3, 8
- [79] Yusheng Su, Xiaozhi Wang, Yujia Qin, Chi-Min Chan, Yankai Lin, Zhiyuan Liu, Peng Li, Juanzi Li, Lei Hou, Maosong Sun, et al. On transferability of prompt tuning for natural language understanding. *arXiv preprint arXiv:2111.06719*, 2021. 2
- [80] Yi-Lin Sung, Jaemin Cho, and Mohit Bansal. VI-adapter: Parameter-efficient transfer learning for vision-and-language tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5227–5237, 2022. 3, 6, 8, 9
- [81] Grant Van Horn, Steve Branson, Ryan Farrell, Scott Haber, Jessie Barry, Panos Ipeirotis, Pietro Perona, and Serge Be-longie. Building a bird recognition app and large scale dataset with citizen scientists: The fine print in fine-grained dataset collection. In *Proceedings of the IEEE Conference*

- on *Computer Vision and Pattern Recognition*, pages 595–604, 2015. [5](#), [6](#), [7](#)
- [82] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. [1](#)
- [83] Bastiaan S Veeling, Jasper Linmans, Jim Winkens, Taco Cohen, and Max Welling. Rotation equivariant cnns for digital pathology. In *Medical Image Computing and Computer Assisted Intervention–MICCAI 2018: 21st International Conference, Granada, Spain, September 16-20, 2018, Proceedings, Part II 11*, pages 210–218. Springer, 2018. [7](#)
- [84] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset. 2011. *California Institute of Technology*. [5](#), [6](#), [7](#)
- [85] Haohan Wang, Songwei Ge, Zachary Lipton, and Eric P Xing. Learning robust global representations by penalizing local predictive power. *Advances in Neural Information Processing Systems*, 32, 2019. [2](#)
- [86] Ruize Wang, Duyu Tang, Nan Duan, Zhongyu Wei, Xuanjing Huang, Guihong Cao, Daxin Jiang, Ming Zhou, et al. K-adapter: Infusing knowledge into pre-trained models with adapters. *arXiv preprint arXiv:2002.01808*, 2020. [3](#), [8](#)
- [87] Zifeng Wang, Shao-Lun Huang, Ercan Engin Kuruoglu, Jimeng Sun, Xi Chen, and Yefeng Zheng. Pac-bayes information bottleneck. *ArXiv*, abs/2109.14509, 2021. [4](#)
- [88] Wei Wen, Chunpeng Wu, Yandan Wang, Yiran Chen, and Hai Li. Learning structured sparsity in deep neural networks. *Advances in neural information processing systems*, 29, 2016. [3](#), [8](#)
- [89] Jianxiong Xiao, James Hays, Krista A Ehinger, Aude Oliva, and Antonio Torralba. Sun database: Large-scale scene recognition from abbey to zoo. In *2010 IEEE computer society conference on computer vision and pattern recognition*, pages 3485–3492. IEEE, 2010. [7](#)
- [90] Runxin Xu, Fuli Luo, Zhiyuan Zhang, Chuanqi Tan, Baobao Chang, Songfang Huang, and Fei Huang. Raise a child in large language model: Towards effective and generalizable fine-tuning. *arXiv preprint arXiv:2109.05687*, 2021. [3](#), [8](#)
- [91] Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. How transferable are features in deep neural networks? *Advances in neural information processing systems*, 27, 2014. [3](#)
- [92] 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. [2](#), [3](#), [5](#), [6](#), [7](#), [4](#), [8](#)
- [93] Xiaohua Zhai, Joan Puigcerver, Alexander Kolesnikov, Pierre Ruysen, Carlos Riquelme, Mario Lucic, Josip Djolonga, Andre Susano Pinto, Maxim Neumann, Alexey Dosovitskiy, et al. A large-scale study of representation learning with the visual task adaptation benchmark. *arXiv preprint arXiv:1910.04867*, 2019. [5](#), [6](#), [7](#)
- [94] Renrui Zhang, Rongyao Fang, Wei Zhang, Peng Gao, Kunchang Li, Jifeng Dai, Yu Qiao, and Hongsheng Li. Tip-adapter: Training-free clip-adapter for better vision-language modeling. *arXiv preprint arXiv:2111.03930*, 2021. [3](#), [8](#)
- [95] Renrui Zhang, Ziyu Guo, Wei Zhang, Kunchang Li, Xupeng Miao, Bin Cui, Yu Qiao, Peng Gao, and Hongsheng Li. Pointclip: Point cloud understanding by clip. In *CVPR 2022*, 2022.
- [96] Renrui Zhang, Jiaming Han, Aojun Zhou, Xiangfei Hu, Shilin Yan, Pan Lu, Hongsheng Li, Peng Gao, and Yu Qiao. Llama-adapter: Efficient fine-tuning of language models with zero-init attention. *ICLR 2024*, 2023. [3](#)
- [97] Renrui Zhang, Xiangfei Hu, Bohao Li, Siyuan Huang, Hanqiu Deng, Hongsheng Li, Yu Qiao, and Peng Gao. Prompt, generate, then cache: Cascade of foundation models makes strong few-shot learners. *CVPR 2023*, 2023. [3](#)
- [98] Mengjie Zhao, Tao Lin, Fei Mi, Martin Jaggi, and Hinrich Schütze. Masking as an efficient alternative to fine-tuning for pretrained language models. *arXiv preprint arXiv:2004.12406*, 2020. [8](#)
- [99] Daquan Zhou, Bingyi Kang, Xiaojie Jin, Linjie Yang, Xiaochen Lian, Zihang Jiang, Qibin Hou, and Jiashi Feng. Deepvit: Towards deeper vision transformer. *arXiv preprint arXiv:2103.11886*, 2021. [8](#)
- [100] Daquan Zhou, Zhiding Yu, Enze Xie, Chaowei Xiao, Animashree Anandkumar, Jiashi Feng, and Jose M Alvarez. Understanding the robustness in vision transformers. In *International Conference on Machine Learning*, pages 27378–27394. PMLR, 2022. [8](#)
- [101] Xiangyang Zhu, Renrui Zhang, Bowei He, Aojun Zhou, Dong Wang, Bin Zhao, and Peng Gao. Not all features matter: Enhancing few-shot clip with adaptive prior refinement. *ICCV 2023*, 2023. [3](#)