

Sensitivity-Aware Visual Parameter-Efficient Fine-Tuning

Haoyu He¹ Jianfei Cai¹ Jing Zhang² Dacheng Tao² Bohan Zhuang^{1†}

¹ ZIP Lab, Monash University ² The University of Sydney

Abstract

Visual Parameter-Efficient Fine-Tuning (PEFT) has become a powerful alternative for full fine-tuning so as to adapt pre-trained vision models to downstream tasks, which only tunes a small number of parameters while freezing the vast majority ones to ease storage burden and optimization difficulty. However, existing PEFT methods introduce trainable parameters to the same positions across different tasks depending solely on human heuristics and neglect the domain gaps. To this end, we study where to introduce and how to allocate trainable parameters by proposing a novel Sensitivity-aware visual Parameter-efficient fine-Tuning (SPT) scheme, which adaptively allocates trainable parameters to task-specific important positions given a desired tunable parameter budget. Specifically, our SPT first quickly identifies the sensitive parameters that require tuning for a given task in a data-dependent way. Next, our SPT further boosts the representational capability for the weight matrices whose number of sensitive parameters exceeds a pre-defined threshold by utilizing existing structured tuning methods, e.g., LoRA [21] or Adapter [20], to replace directly tuning the selected sensitive parameters (unstructured tuning) under the budget. Extensive experiments on a wide range of downstream recognition tasks show that our SPT is complementary to the existing PEFT methods and largely boosts their performance, e.g., SPT improves Adapter with supervised pre-trained ViT-B/16 backbone by 4.2% and 1.4% mean Top-1 accuracy, reaching SOTA performance on FGVC and VTAB-1k benchmarks, respectively. Source code is at <https://github.com/ziplab/SPT>.

1. Introduction

To effectively adapt the pre-trained representations to the downstream tasks, the de-facto choice is full fine-tuning, which initializes the model with the pre-trained weights and tunes all the parameters. However, vanilla full fine-tuning needs to store a separate instance of parameters for each

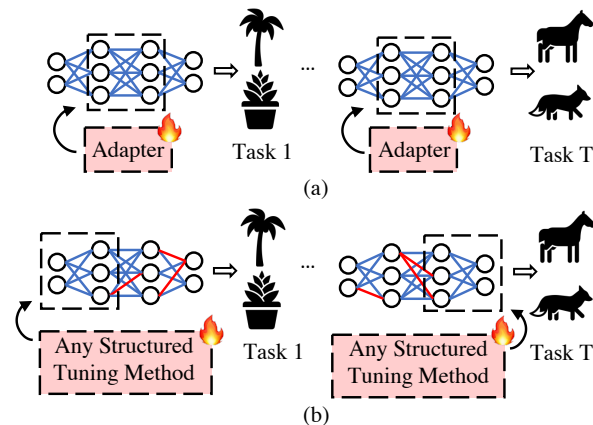


Figure 1: (a) Existing PEFT methods, such as Adapter [20] introduce trainable parameters to the same positions for all downstream tasks. These methods design task-agnostic positions to employ trainable parameters relying on heuristics and neglect consideration of the distinct domain gaps and characteristics for the downstream tasks. (b) Our Sensitivity-aware visual Parameter-efficient fine-Tuning (SPT) introduces trainable parameters to the task-specific important positions and allocates them with both unstructured and structured tuning granularities, simultaneously. For structured tuning, SPT can exploit any existing structured tuning methods, such as LoRA [21] or Adapter [20]. Red lines and blocks represent trainable parameters and modules, while blue lines represent frozen parameters.

task and each deployment scenario. It can be extremely storage-intensive as the storage cost grows linearly with the number of possible cases, considering there are vast varieties of downstream tasks and dynamic deployment environments, especially when deploying the large vision models [12, 32, 56] to mobile systems. For example, even storing a single large pre-trained ViT-H [12] model on a local disk requires at least 2.3GB, while the Top-10 U.S. apps required only collectively 2.2GB in May 2021.¹

Notably, an emerging solution is to replace vanilla fine-tuning with visual Parameter-Efficient Fine-Tuning (PEFT) [22, 9, 63, 23], which only tunes a small num-

[†]Corresponding author. E-mail: bohan.zhuang@gmail.com

¹<https://sensortower.com/blog/ios-app-size-growth-2021>

ber of trainable parameters while freezing the vast majority ones that are shared by multiple tasks. As PEFT approaches exhibit less than 1% of trainable parameters, the storage burden is largely alleviated. Another attractive property of PEFT is that tuning fewer parameters eases the optimization difficulty and mitigates the overfitting issue when adapting large pre-trained models on the target dataset, thereby achieving comparable or even better performance than vanilla fine-tuning [22]. Although promising, the existing PEFT approaches introduce trainable parameters to the same positions for all downstream tasks, relying on human heuristics and neglecting task-specific domain gaps and characteristics, which limits their performance. For instance, in a task-agnostic manner, Prompt Tuning-deep [22] and Adapter [20] respectively add trainable parameters to multi-head self-attention and feed-forward network layers for all distinct tasks as depicted in Figure 1 (a).

To address this fundamental challenge, we explore *where to introduce* and *how to allocate* trainable parameters under a desired parameter budget by presenting a novel Sensitivity-aware visual Parameter-efficient fine-Tuning (SPT) scheme that identifies the *task-specific important positions* to adaptively allocate trainable parameters. Since the pre-trained weights at distinct positions have varying contributions for different downstream tasks [58, 26, 37], we first propose a new criterion to quickly identify the task-specific sensitive parameters that require tuning in a data-dependent way. Inspired by model pruning metrics [44, 35, 2, 3], we propose to measure the parameter sensitivity with the loss reduction when being tuned, which can be approximated by a first-order Taylor expansion derived within a single forward and backward pass ahead of fine-tuning in one-shot. Our sensitivity criterion is simple and effective, which can identify the task-specific important positions to introduce trainable parameters for any backbone quickly. For instance, calculating the sensitivity for ViT-B/16 backbone takes only 5.5 seconds with a single GPU on any of the VTAB-1k datasets.

With our criterion, we empirically observe that the proportions of the sensitivity parameters for each block indeed vary markedly across different tasks in Section 4.4. To allocate the trainable parameters under a desired trainable parameter budget, an intuitive solution is to directly tune the most sensitive weight connections, which we name unstructured tuning. Despite its simplicity and flexibility, unstructured tuning only tunes a few parameters which still lacks representational capability and is challenging to bridge the domain gap. To this end, we propose to further incorporate structured tuning to replace unstructured tuning at the sensitive weight matrices whose numbers of sensitive parameters exceed a pre-defined threshold to improve the representational capability under a similar parameter budget. Structured tuning can be implemented by any parameter-efficient

structured tuning methods [21, 9, 23, 22] that directly adjust the hidden representations, *e.g.*, inserting an adapter module sequentially after the sensitive weight matrices. Therefore, our SPT adaptively combines both unstructured and structured tuning granularity and allocates trainable parameters with high flexibility and representational capability for each distinct downstream task.

This paper has the following key contributions. 1) We make the pioneering exploration to identify the task-specific important positions under the PEFT setting, which is fast, effective, versatile to be applied to various backbones with different pre-training strategies, and orthogonal to the existing PEFT methods. 2) Based on the sensitivity criterion, we propose a trainable parameter allocation strategy that adaptively combines both unstructured and structured tuning under a desired parameter budget to achieve high flexibility, large capacity, and favorable trade-off between parameter efficiency and accuracy. 3) Extensive experiments on a total of 24 downstream recognition tasks with both plain and hierarchical vision Transformer backbones under supervised and self-supervised pre-trainings show that our SPT is complementary to the existing PEFT methods and boosts their performance by large margins. For instance, SPT improves Adapter [20] by 4.2% mean Top-1 accuracy, outperforming the SOTA PEFT methods on the FGVC benchmark.

2. Related Work

Parameter-efficient fine-tuning. Full fine-tuning is the most predominant approach when adapting a large-scale pre-trained model to downstream tasks, where the model is initialized from the pre-trained weights with all parameters trainable. Yet, when a model becomes larger, parameter-efficient fine-tuning [27, 28] is highly desirable, which tunes only a tiny portion of parameters to alleviate the storage burden. The general PEFT approaches can be categorized into addition-based PEFT methods and reparameterization-based PEFT methods.

Addition-based PEFT attaches additional trainable parameters to the backbone and only tunes these parameters. Apart from Prompt tuning [22] and Adapter [20], recent addition-based methods study connecting or combining existing PEFT methods. For instance, He *et al.* [18] connect Prompt tuning and Adapter and provide a unified view that all PEFT approaches share the same design to adjust the hidden representations. Zhang *et al.* [63] search for the optimal configurations to combine multiple PEFT approaches following once-for-all scheme [5, 54]. Since the additional parameters require extra computations compared to full fine-tuning, a few recent works [46, 48] design specific architectures to avoid storing the intermediate activations, thereby alleviating the fine-tuning memory cost. However, it is noteworthy that enhancing training efficiency is not the primary objective of our work.

Reparameterization-based PEFT aims to avoid extra computational costs by tuning parameters that are inherently in or can be reparameterized into the backbone during inference. Prior works select the parameters that are inherently in the backbone, including the bias terms [59], the last several layers [58, 4], and weight connections [14, 64]. To reparameterize new parameters into the backbone [17, 29], representative work LoRA [21] optimizes two low-rank matrices which can be further merged into the weight matrices. In contrast to the aforementioned works, we argue the importance of tuning parameters at task-specific important positions and quickly identify them with our proposed parameter sensitivity criterion before tuning, which is complementary to and provides valuable guidance for the existing PEFT methods. Moreover, our SPT can also be inference-efficient when implementing structured tuning with any reparameterization-based structured tuning method. Recently, SSF [29] is proposed to introduce trainable scaling and shifting parameters that can be absorbed into the previous linear layers. However, it cannot scale to higher trainable parameter budgets and requires a complex and time-consuming hyper-parameter search for learning rate, weight decay, and drop-path rate on each individual dataset, thus is not directly comparable to our method.

Task-specific transfer learning. The effectiveness of transferring pre-trained models to downstream tasks strongly depends on the relationship between the source and target tasks [43, 53, 26, 40]. This has motivated the community to explore the optimal pre-training data [11, 57], model [47, 38], and weights [16, 55] for the target task. To seek suitable *task-specific pre-training data*, Cui *et al.* [11] select the source domain data from the top-k most similar classes measured by Earth Mover’s Distance; Yoon *et al.* [57] weight each class in the source domain with reinforcement learning; and Puigcerver *et al.* [41] first train a diverse set of experts and then select the most relevant expert for each target task. Another line of work selects a suitable *pre-trained model for the target task* ahead of fine-tuning by measuring the transferability of pre-trained models to the target domain with interclass covariance between the source data and target classes [1] or conditional cross-entropy [47] between the source and target labels. Considering the transferability of the feature representations at distinct layers for the same pre-trained model is different [58, 37], recent works [15, 45] endeavour *transfer task-specific weights* by freezing some pre-trained weights and fine-tuning the rest. For example, the task-specific fine-tuned weights are selected by learning a policy network with Gumbel-Softmax [16], optimizing a sparse mask with L_0 norm [14], and learning binary gates for each parameter [64]. Our SPT also adaptively selects task-specific parameters. In contrast to the previous work, we 1) derive task-specific important positions prior to fine-tuning with only a single

forward and backward pass, which is computationally efficient; 2) mask the gradients for insensitive parameters in unstructured tuning with fixed binary masks, thereby having more affordable fine-tuning memory than optimizing learnable binary masks in [14, 64]. Moreover, we are pioneering work to adaptively allocate task-specific trainable parameters with both fine-grained unstructured and coarse-grained structured tuning granularities to achieve both high flexibility and representational capability.

3. Method

Our sensitivity-aware visual parameter-efficient fine-tuning consists of two stages. In the first stage, SPT measures the task-specific sensitivity for the pre-trained parameters (Section 3.1). Based on the parameter sensitivity and a given parameter budget, SPT then adaptively allocates trainable parameters to task-specific important positions (Section 3.2).

3.1. Task-specific Parameter Sensitivity

Recent research has observed that pre-trained backbone parameters exhibit varying feature patterns [42, 36] and criticality [62, 8] at distinct positions. Moreover, when transferred to downstream tasks, their efficacy varies depending on how much pre-trained features are reused and how well they adapt to the specific domain gap [58, 26, 37]. Motivated by these observations, we argue that not all parameters contribute equally to the performance across different tasks in PEFT and propose a new criterion to measure the sensitivity of the parameters in the pre-trained backbone for a given task.

Specifically, given the training dataset \mathcal{D}_t for the t -th task and the pre-trained model weights $\mathbf{w} = \{w_1, w_2, \dots, w_N\} \in \mathbb{R}^N$ where N is the total number of parameters, the objective for the task is to minimize the empirical risk: $\min_{\mathbf{w}} E(\mathcal{D}_t, \mathbf{w})$. We denote the parameter sensitivity set as $\mathcal{S} = \{s_1, \dots, s_N\}$ and the sensitivity s_n for parameter w_n is measured by the empirical risk difference when tuning it:

$$s_n = E(\mathcal{D}_t, \mathbf{w}) - E(\mathcal{D}_t, \mathbf{w} \mid w_n = w_n^*), \quad (1)$$

where $w_n^* = \underset{w_n}{\operatorname{argmin}}(E(\mathcal{D}_t, \mathbf{w}))$. We can reparameterize the tuned parameters as $w_n^* = w_n + \Delta_{w_n}$, where Δ_{w_n} denotes the update for w_n after tuning. Here we individually measure the sensitivity of each parameter, which is reasonable given that most of the parameters are frozen during fine-tuning in PEFT. However, it is still computationally intensive to compute Eq. (1) for two reasons. Firstly, getting the empirical risk for N parameters requires forwarding the entire network N times, which is time-consuming. Secondly, it is challenging to derive Δ_{w_n} , as we have to tune each individual w_n until convergence.

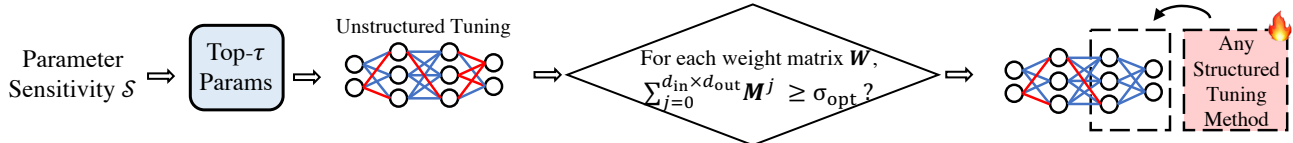


Figure 2: Overview of our trainable parameter allocation strategy. With the parameter sensitivity set \mathcal{S} , we first get the top- τ sensitive parameters. Instead of directly tuning these sensitive parameters, we also boost the representational capability by replacing unstructured tuning with structured tuning at sensitive weight matrices that have a large number of sensitive parameters, which can be implemented by an existing structured tuning method, *e.g.*, LoRA [21] and Adapter [20]. Red lines and blocks represent trainable parameters and modules, while blue lines represent frozen parameters.

Algorithm 1 Computing task-specific parameter sensitivities

Input: Pre-trained model with network parameters \mathbf{w} , training set \mathcal{D}_t for the t -th task, and number of training samples C used to calculate the parameter sensitivities
Output: Sensitivity set $\mathcal{S} = \{s_1, \dots, s_N\}$
Initialize $\mathcal{S} = \{0\}^N$
for $i \in \{1, \dots, C\}$ **do**
 Get the i -th training sample of \mathcal{D}_t
 Compute loss E
 Compute gradients \mathbf{g}
 for $n \in \{1, \dots, N\}$ **do**
 Update sensitivity for the n -th parameter: $s_n = s_n + g_n^2$
 end for
end for

To overcome the first barrier, we simplify the empirical loss by approximating s_n in the vicinity of \mathbf{w} by its first-order Taylor expansion

$$s_n^{(1)} = -g_n \Delta_{w_n}, \quad (2)$$

where the gradients $\mathbf{g} = \partial E / \partial \mathbf{w}$, and g_n is the gradient of the n -th element of \mathbf{g} . To address the second barrier, following [31, 7], we take the one-step unrolled weight as the surrogate for w_n^* and approximate Δ_{w_n} in Eq. (2) with a single step of gradient descent. We can accordingly get $s_n^{(1)} \approx g_n^2 \epsilon$, where ϵ is the learning rate. Since ϵ is the same for all parameters, we can eliminate it when comparing the sensitivity with the other parameters and finally get

$$s_n^{(1)} \approx g_n^2. \quad (3)$$

Therefore, the sensitivity of a parameter can be efficiently measured by its potential to reduce the loss on the target domain. Note that although our criterion draws inspiration from pruning work [35], it is distinct from it. [35] measures the parameter importance by the squared change in loss when removing them, *i.e.*, $(E(\mathcal{D}_t, \mathbf{w}) - E(\mathcal{D}_t, \mathbf{w} | w_n = 0))^2$ and finally derives the parameter importance by $(g_n w_n)^2$, which is different from our formulations in Eqs. (1) and (3).

In practice, we accumulate \mathcal{S} from a total number of C training samples ahead of fine-tuning to generate accu-

rate sensitivity as shown in Algorithm 1, where C is a pre-defined hyper-parameter. In Section 4.3, we show that employing only 400 training samples is sufficient for getting reasonable parameter sensitivity, which requires only 5.5 seconds with a single GPU for any VTAB-1k dataset with ViT-B/16 backbone [12].

3.2. Adaptive Trainable Parameters Allocation

Our next step is to allocate trainable parameters based on the obtained parameter sensitivity set \mathcal{S} and a desired parameter budget τ . A straightforward solution is to directly tune the top- τ most sensitive unstructured connections (parameters) while keeping the rest frozen, which we name unstructured tuning. Specifically, we select the top- τ most sensitive weight connections in \mathcal{S} to form the sensitive weight connection set \mathcal{T} . Then, for a weight matrix $\mathbf{W} \in \mathbb{R}^{d_{in} \times d_{out}}$, we can get a binary mask $\mathbf{M} \in \mathbb{R}^{d_{in} \times d_{out}}$ computed by

$$\mathbf{M}^j = \begin{cases} 1 & \mathbf{W}^j \in \mathcal{T} \\ 0 & \mathbf{W}^j \notin \mathcal{T} \end{cases}, \quad (4)$$

where \mathbf{W}^j and \mathbf{M}^j are the j -th element in \mathbf{W} and \mathbf{M} , respectively. Accordingly, we can train the sensitive parameters by gradient descent and the updated weight matrix can be formulated as $\mathbf{W}' \leftarrow \mathbf{W} - \epsilon \mathbf{g}_{\mathbf{W}} \odot \mathbf{M}$, where $\mathbf{g}_{\mathbf{W}}$ is the gradient for \mathbf{W} .

However, considering PEFT approaches generally limit the proportion of trainable parameters to less than 1%, tuning only a small number of unstructured weight connections might not have enough representational capability to handle the downstream datasets with large domain gaps from the source pre-training data. Therefore, to improve the representational capability, we propose to replace unstructured tuning with structured tuning at the sensitive weight matrices that have a high number of sensitive parameters. To preserve the parameter budget, we can implement structured tuning with an existing efficient structured tuning PEFT method [21, 9, 20, 23] that learns to directly adjust all hidden dimensions at once. We depict an overview of our trainable parameter allocation strategy in Figure 2. For example, we can employ the low-rank reparameterization trick LoRA [21] to the sensitive weight matrices and the one-step update for \mathbf{W} can be formulated as

$$W' = \begin{cases} W + W_{\text{down}}W_{\text{up}} & \text{if } \sum_{j=0}^{d_{\text{in}} \times d_{\text{out}}} M^j \geq \sigma_{\text{opt}} \\ W - \epsilon g_W \odot M & \text{otherwise} \end{cases}, \quad (5)$$

where $W_{\text{down}} \in \mathbb{R}^{d_{\text{in}} \times r}$ and $W_{\text{up}} \in \mathbb{R}^{r \times d_{\text{out}}}$ are two learnable low-rank matrices to approximate the update of W and rank r is a hyper-parameter where $r \ll \min(d_{\text{in}}, d_{\text{out}})$. In this way, we perform structured tuning on W when its number of sensitive parameters exceeds σ_{opt} , whose value depends on the pre-defined type of structured tuning method. For example, since implementing structured tuning with LoRA requires $2 \times d_{\text{in}} \times d_{\text{out}} \times r$ trainable parameters for each sensitive weight matrix, we set $\sigma_{\text{LoRA}} \leftarrow 2 \times d_{\text{in}} \times d_{\text{out}} \times r$ to ensure that the number of trainable parameters introduced by structured tuning is always equal to or lower than the number of sensitive parameters.

In this way, our SPT adaptively incorporates both structured and unstructured tuning granularities to enable higher flexibility and stronger representational power, simultaneously. In Section 4.3, we show that structured tuning is important for the downstream tasks with larger domain gaps and both unstructured and structured tuning contribute clearly to the superior performance of our SPT.

4. Experiments

4.1. Experimental Setup

Datasets and metrics. We evaluate our SPT on total 24 downstream tasks in two groups following [22]. 1) FGVC is a benchmark for fine-grained visual classification, including CUB-200-2011 [52], NABirds [50], Oxford Flowers [39], Stanford Cars [13], and Stanford Dogs [24] datasets. Each FGVC dataset contains between 55 to 200 classes and a few thousand images for train, validation, and test. We follow the validation splits in [22] if the validation set is unavailable. 2) VTAB-1k [60] is a large-scale transfer learning benchmark consisting of a collection of 19 visual classification tasks. VTAB-1k can further be divided into three groups, including Natural tasks with natural images, Specialized tasks with images captured by specialized equipment, e.g., medical images, and Structured tasks with images mostly generated from synthetic environments. Each of the VTAB-1k dataset has only 800 training and 200 validation samples, while the test set sizes vary. We use top-1 accuracy (%) averaged within each group as our main metric following [22].

Pre-trained backbones. We conduct experiments on the plain vision Transformer backbone ViT-B/16 [12] that is pre-trained on ImageNet [25] with different pre-training strategies following [22], including supervised pre-training and self-supervised pre-training with MAE [19] and MoCo v3 [10] following [22]. We also conduct experiments on the representative hierarchical vision Transformer backbone Swin-B [32] under supervised pre-training.

Contenders. We categorize the baseline methods into

addition-based and reparameterization-based PEFT methods as introduced in Section 2. Unless specified, all baseline methods keep the backbone frozen. Addition-based methods require extra computations during inference, including MLP- k , PROMPT-SHALLOW [22], PROMPT-DEEP [22], ADAPTER- k [20], ADAPTFORMER [9], and NOAH [63]. Reparameterization-based methods have no additional computational overhead during inference, including LINEAR, PARTIAL- k , BIAS [59], and LORA- k [21]. Here k represents the number of bottleneck dimension in ADAPTER- k and LORA- k . We also compare with full fine-tuning which is denoted by FULL. We introduce the details of these methods in the supplementary material.

We also introduce two variants of our SPT: addition-based SPT-ADAPTER and reparameterization-based SPT-LORA. SPT-ADAPTER directly adjusts the hidden representations that are computed by sensitive weight matrices following [20], while SPT-LORA approximates updating the sensitive weight matrices following [21]. For the two variants, we follow the exact weight initializations that are described in [21] and follow [63] to set the bottleneck dimension as 8.

Implementation details. Following [63], we use the AdamW optimizer [34] with cosine learning rate decay and set the batch size, learning rate, and weight decay as 64, 1×10^{-3} , and 1×10^{-4} , respectively. We also follow [63] for the standard data augmentation pipeline. We set the number of training samples C used to calculate our parameter sensitivities in Algorithm 1 universally to be 800 for the main experiments.

4.2. Main Results

We compare our method with the baseline methods under different backbones, pre-training strategies, and tasks.

Visual recognition with ViT backbone. First, *our proposed SPT-ADAPTER and SPT-LORA achieve the best performance under different trainable parameter budgets with supervised pre-trained ViT-B/16 backbone*, as shown in Table 1 and Figure 3 (a). For instance, SPT-ADAPTER outperforms the SOTA method NOAH by a clear margin of 0.9% mean top-1 accuracy over the 19 VTAB-1k datasets with fewer trainable parameters. We speculate that our SPT variants allocate trainable parameters at task-specific positions compared to the heuristically selected positions in the baseline methods, which contributes to our superior performance. We also observe that our SPT-ADAPTER and SPT-LORA achieve large performance gains over ADAPTER and LORA variants, respectively. For example, SPT-ADAPTER and SPT-LORA with 0.41% trainable parameters respectively improve ADAPTER-8 and LORA-8 significantly by 4.0% and 3.3% mean accuracy on the FGVC benchmark. This suggests that identifying task-specific important positions and combining both unstructured and structured tuning granularities with SPT are complementary to the exist-

ViT-B/16 (85.8M)	Total params	FGVC		VTAB-1k				
		Tuned / Total	Mean Acc.	Tuned / Total	Natural	Specialized	Structured	Mean Acc.
FULL	24.02×	100%	88.5	100%	75.9	83.4	47.6	69.0
Addition-based methods								
MLP-3	1.35×	1.50%	79.8	1.42%	67.8	72.8	30.6	57.1
PROMPT-SHALLOW	1.04×	0.31%	84.6	0.13%	76.8	79.7	47.0	67.8
PROMPT-DEEP	1.18×	0.98%	89.1	1.14%	78.5	82.4	55.0	72.0
ADAPTER-8	1.06×	0.39%	85.5	0.23%	79.0	84.1	58.5	73.9
ADAPTER-32	1.19×	0.95%	85.6	0.71%	79.6	84.0	58.3	74.0
ADAPTFORMER	1.09×	0.44%	85.1	0.36%	80.6	<u>85.4</u>	58.5	74.8
NOAH	-	-	-	0.52%	80.2	84.9	<u>61.3</u>	75.5
SPT-ADAPTER (Ours)	1.08×	0.41%	<u>89.5</u>	0.30%	<u>81.3</u>	85.3	60.8	<u>75.8</u>
SPT-ADAPTER (Ours)	1.10×	0.47%	89.8	0.44%	82.0	85.8	61.4	76.4
Reparameterization-based methods								
LINEAR	1.02×	0.12%	79.3	0.04%	68.9	77.2	26.8	57.6
PARTIAL-1	3.00×	8.38%	82.6	8.30%	69.4	78.5	34.2	60.7
BIAS	1.05×	0.13%	88.4	0.13%	73.3	78.3	44.1	65.2
LoRA-8	1.07×	0.55%	86.0	0.23%	79.5	84.6	60.5	74.9
LoRA-16	1.18×	0.90%	84.8	0.69%	79.8	84.9	60.2	75.0
SPT-LoRA (Ours)	1.08×	0.41%	<u>89.3</u>	0.31%	<u>81.5</u>	<u>85.6</u>	<u>60.7</u>	<u>75.9</u>
SPT-LoRA (Ours)	1.15×	0.60%	90.1	0.63%	81.9	85.9	61.3	76.4

Table 1: Comparisons on FGVC and VTAB-1k [60] benchmarks using supervised pre-trained ViT-B/16 backbone pre-trained on ImageNet-21k. “Total params” denotes the ratio of the total number of parameters needed for all downstream tasks relative to the one for the pre-trained backbone, and “Tuned/Total” denotes the fraction of trainable parameters. Top-1 accuracy (%) is reported. The best result is in **bold**, and the second-best result is underlined.

ViT-B/16 (85.8M)	Total Params	VTAB-1k MAE					VTAB-1k MoCo v3				
		Tuned / Total	Natural	Specialized	Structured	Mean Acc.	Tuned / Total	Natural	Specialized	Structured	Mean Acc.
FULL	38.02×	100%	59.3	79.7	53.8	64.3	100%	72.0	84.7	42.0	69.6
Addition-based methods											
ADAPTER-8	1.08×	0.23%	57.2	78.4	54.7	63.4	0.23%	27.6	70.9	48.4	49.0
ADAPTER-32	1.28×	0.95%	55.3	78.8	53.3	62.5	0.99%	74.2	82.7	47.7	68.2
PROMPT-SHALLOW	1.02×	0.12%	40.0	69.7	27.5	45.7	0.12%	67.3	82.3	37.6	62.4
PROMPT-DEEP	1.05×	0.23%	36.0	60.6	26.6	41.1	0.07%	70.3	83.0	42.4	65.2
SPT-ADAPTER (Ours)	1.07×	0.26%	<u>64.8</u>	<u>82.4</u>	<u>60.4</u>	<u>69.2</u>	0.08%	<u>76.1</u>	<u>84.9</u>	<u>60.1</u>	<u>73.7</u>
SPT-ADAPTER (Ours)	1.13×	0.41%	65.6	82.7	60.7	69.7	0.30%	76.6	85.0	61.7	74.4
Reparameterization-based methods											
LINEAR	1.02×	0.04%	18.9	52.7	23.7	32.1	0.04%	67.5	81.1	30.3	59.6
PARTIAL-1	4.16×	8.30%	58.4	78.3	47.6	61.5	8.30%	72.3	84.6	47.9	68.3
BIAS	1.06×	0.13%	54.6	75.7	47.7	59.3	0.13%	72.9	81.1	53.4	69.2
LoRA-8	1.08×	0.23%	57.5	77.7	57.7	64.3	0.23%	21.2	66.7	45.1	44.3
LoRA-16	1.28×	0.69%	57.3	77.1	59.9	64.8	0.69%	16.0	64.0	48.7	42.9
SPT-LoRA (Ours)	1.11×	0.29%	<u>63.8</u>	<u>81.6</u>	<u>60.0</u>	<u>68.5</u>	0.30%	76.5	<u>85.4</u>	<u>63.0</u>	<u>75.0</u>
SPT-LoRA (Ours)	1.23×	0.69%	65.4	82.4	61.5	69.8	0.50%	76.5	86.0	63.6	75.3

Table 2: Comparisons on VTAB-1k [60] benchmark using self-supervised ViT-B/16 backbone pre-trained by MAE [19] and MoCo v3 [10]. “Total params” denotes the ratio of the total number of parameters needed for all downstream tasks relative to the one for the pre-trained backbone, and “Tuned/Total” denotes the fraction of trainable parameters. Top-1 accuracy (%) is reported. The best result is in **bold**, and the second-best result is underlined.

ing PEFT methods and boost their performance.

Second, SPT variants *outperform baseline methods and full fine-tuning by significant margins with the self-supervised pre-trained ViT-B/16 backbones*. As shown in Table 2, existing PEFT approaches exhibit inferior results than full fine-tuning with the self-supervised pre-trained

backbones MAE and MoCo v3. It is worth noting that previous PEFT methods yield inconsistent results with the backbones of different pre-training strategies. In contrast, SPT variants consistently outperform full fine-tuning. In particular, SPT-ADAPTER achieves remarkable 5.8% and 5.5% mean top-1 accuracy gains over the best-performing

Method	Tuned / Total	Natural	Specialized	Structured	Mean / Acc.
FULL	100%	79.1	86.2	59.7	75.0
Addition-based methods					
MLP-3	1.60%	73.6	75.2	35.7	61.5
PROMPT-SHALLOW	0.04%	79.9	82.5	37.8	66.7
PROMPT-DEEP	0.23%	76.8	84.5	53.4	71.6
ADAPTER-8	1.18%	81.7	87.3	61.2	76.7
SPT-ADAPTER (ours)	0.33%	83.0	87.3	62.1	77.5
Reparameterization-based methods					
LINEAR	0.04%	73.5	80.8	33.5	62.6
PARTIAL-1	2.15%	73.1	81.7	35.0	63.3
LoRA-8	1.18%	81.7	87.2	60.1	76.3
SPT-LoRA (ours)	0.49%	83.1	87.4	60.4	77.2

Table 3: Comparisons on VTAB-1k [60] benchmark with supervised pre-trained Swin-B [32]. ‘‘Tuned/Total’’ denotes the fraction of trainable parameters. Top-1 accuracy (%) is reported. The best result is in **bold**.

baseline method on VTAB-1k benchmark with only 0.26% and 0.08% trainable parameters for MAE and MoCo v3 pre-trained backbones, respectively. Moreover, our observation in supplementary material suggests that self-supervised pre-trained ViT backbones have more diverse sensitivity distributions and a higher variance in sensitivity across different tasks than the supervised pre-trained one. This leads to the conjecture that baseline methods which assign trainable parameters to the same positions for all tasks may fail to mitigate the distinct domain gaps in individual downstream datasets, whereas our SPT allocates trainable parameters to task-specific positions accurately.

Visual recognition with Swin and ConvNeXt backbones.

From Table 3, we observe that our SPT-LoRA and SPT-ADAPTER also achieve SOTA performance with Swin-B backbone on all dataset groups. We also follow VPT [22] to apply SPT variants to ResNet-alike architecture ConvNeXt-Base [33] and report the results in Table 4. We observe that SPT variants achieve better trade-offs between accuracy and parameter efficiency than the baseline methods. The results further demonstrate the versatility and effectiveness of our SPT.

Semantic segmentation. We follow VPT [22] to conduct semantic segmentation on ADE20k dataset. Following the settings of [22], we apply our SPT to SETR-PUP [65] framework with ViT-L backbone and only allocate trainable parameters to the backbone with the head fully fine-tuned. We report the mIoU results in Table 5. Notably, SPT-LoRA and SPT-ADAPTER outperform the baseline methods by large margins, indicating our SPT can be generalized to the semantic segmentation task. We report the results for conducting only structured tuning for both SPT-LoRA and SPT-ADAPTER as it yields higher mIoU.

4.3. Ablation Study

Effect of the sensitivity criterion. We investigate the effectiveness of our sensitivity criterion on VTAB-1k by em-

Method	Tuned/Total	Natural	Specialized	Structured	Mean Acc.
FULL	100%	78.0	83.7	60.4	74.0
LoRA	0.79%	82.2	84.7	64.1	77.0
Adapter	0.47%	83.1	84.9	64.6	77.5
SPT-LoRA (ours)	0.57%	83.4	86.7	65.9	78.7
SPT-Adapter (ours)	0.36%	83.7	86.2	65.3	78.4

Table 4: Comparisons on VTAB-1k [60] benchmark with supervised pre-trained ConvNeXt-Base [33] backbone. ‘‘Tuned/Total’’ denotes the fraction of trainable parameters. Top-1 accuracy (%) is reported. The best result is in **bold**.

Method	mIoU-s.s. (%)	mIoU-m.s. (%)	Trainable Param.
VPT	44.0	45.6	15.8M
LoRA	43.9	45.9	14.7M
Adapter	44.4	46.6	14.6M
SPT-LoRA (ours)	45.4	47.5	14.6M
SPT-Adapter (ours)	45.2	47.2	14.6M

Table 5: Semantic Segmentation: Comparisons on ADE20k [66] val with SETR [65] on ViT-L backbone. We report both single-scale (s.s.) and multi-scale (m.s.) mIoU results. The best result is in **bold**.

ploying structured tuning methods from [22, 20, 21] to the task-specific sensitive weight matrices. Note that we do not conduct unstructured tuning to ensure fair comparisons. The results are presented in Figure 3 (b). Our criterion brings consistent 1.1%, 1.6%, and 0.8% performance gains for PROMPT-DEEP, ADAPTER-32, and LoRA-16, respectively, which demonstrates the effectiveness and versatility of our sensitivity criterion to identify accurate task-specific important positions.

Effect of structured and unstructured tuning. We investigate the effectiveness of unstructured and structured tuning individually on VTAB-1k. The results are presented in Table 6. We start by applying ADAPTER-8 to the sensitive weight matrices identified by our sensitivity criterion (SPT-ADAPTER w/o unstructured). We observe that our sensitivity criterion boosts the performance of all the dataset groups by clear margins, which again demonstrates the importance of our sensitivity criterion. Next, we observe that allocating the trainable parameters to the unstructured sensitive weight connections also brings accuracy improvement to the Natural and Specialized datasets from ADAPTER-8. However, we find that structured tuning is especially important for achieving good performance on Structured datasets. To further investigate this phenomenon, we observe that Structured datasets have larger domain gaps from the pre-training source domain [25] compared to Natural and Specialized datasets as visualized in Figure 3 (c). We hence conjecture that structured tuning has a higher representational capability than unstructured tuning which facilitates mitigating the large domain gaps during fine-tuning (see the supplementary material for visual examples). Finally, we observe

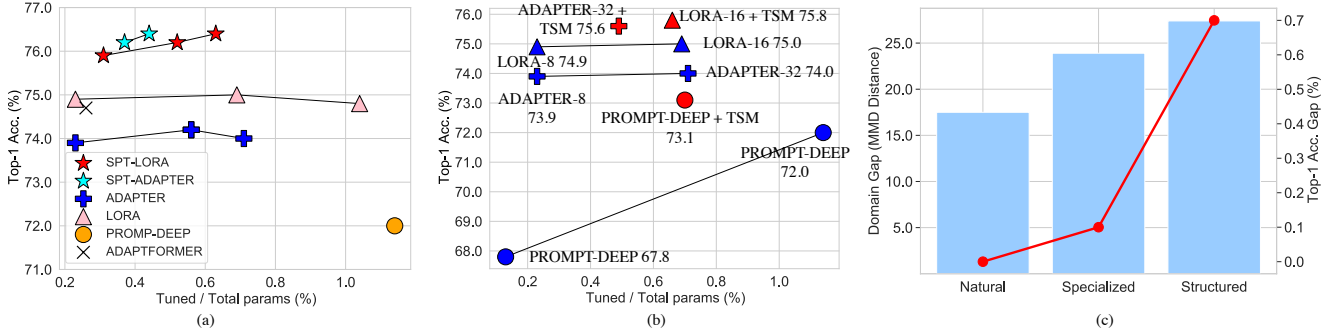


Figure 3: (a) Accuracy vs. parameter efficiency with supervised pre-trained ViT-B/16 backbone on VTAB-1k [60]. SPT variants perform favorably against the other PEFT approaches and are more scalable. (b) Applying other PEFT structured tuning methods [22, 20, 21] to the task-specific sensitive weight matrices (denoted by TSM) identified by our criterion with supervised pre-trained ViT-B/16 backbone on VTAB-1k. Our criterion brings consistent performance gain. (c) Domain vs. performance gaps for different dataset groups in VTAB-1k [60]. The blue bars show the domain gaps between the source domain (ImageNet [25]) and target domains, which are measured by Maximum Mean Discrepancy (MMD) distance [49]. The red line represents the performance gaps between SPT-ADAPTER w/o unstructured and w/o structured, using supervised pre-trained ViT-B/16 backbone. The dataset groups are Natural, Specialized, and Structured. Structured tuning is important for achieving good performance on Structured datasets with larger domain gaps.

Method	Tuned / Total	Natural	Specialized	Structured	Mean / Acc.
ADAPTER-8	0.23%	79.0	84.1	58.5	73.9
SPT-ADAPTER w/o unstructured	0.29%	81.2	85.1	60.3	75.5
SPT-ADAPTER w/o structured	0.34%	81.2	85.0	59.6	75.3
SPT-ADAPTER	0.30%	81.3	85.3	60.8	75.8

Table 6: Ablation study on structured and unstructured tuning only with supervised pre-trained ViT-B/16 backbone. Top-1 accuracy (%) is reported. We set different parameter constraints to align the fractions of trainable parameters for these cases. The best result is in **bold**

that incorporating both structured and unstructured tuning at task-specific important positions achieves the highest performance on all dataset groups.

Effect of number of training samples C to get parameter sensitivity. We investigate the effect of the number of training images C for calculating our parameter sensitivity (Algorithm 1 of the main paper). We randomly sample training samples and report the mean results over three runs in Table 8. We find that our SPT is robust to the number of training samples C and randomly sampling 400 out of a total of 800 training samples is sufficient to obtain accurate task-specific important positions, *e.g.*, calculating the sensitivity for ViT-B/16 backbone takes only 5.5 seconds with a single GPU on any of the VTAB-1k datasets and this computation is required only once.

Computational cost analysis. We investigate the computational cost of SPT-LoRA by comparing with full fine-tuning, addition-based method PROMPT-DEEP, and reparameterization-based method LORA-16. The results are presented in Table 7. We observe that PROMPT-DEEP has higher inference latency and inference GPU memory due to the additional prompts. In contrast, since the up-

Method	Inference Latency (ms/img)	Inference Memory (GB)	Fine-tuning Memory (GB)
FULL	2.8	1.3	11.9
PROMPT-DEEP	3.8	1.9	13.2
LoRA-16	2.8	1.3	8.2
SPT-LoRA w/o unstructured	2.8	1.3	8.3
SPT-LoRA	2.8	1.3	9.8

Table 7: Cost comparisons with ViT-B/16 backbone on the VTAB-1k [60] benchmark with around 0.70% fractions of the trainable parameters. We report the latency (ms/img) and the peak memory usage (GB) with image resolution 224×224 on a single GeForce 3090 GPU. The best result is in **bold**.

C	240	400	560	800
Mean Acc.	76.3	76.4	76.4	76.4

Table 8: Effect of the number of training samples used to get the sensitivity for SPT-LoRA with supervised pre-trained ViT-B/16 backbone on VTAB-1k [60]. Top-1 accuracy (%) is reported. The best result is in **bold**.

dated parameters after fine-tuning can be reparameterized and merged into the pre-trained model, our SPT-LoRA, SPT-LoRA w/o unstructured tuning, and LORA-16 are more efficient than PROMPT-DEEP during inference. However, we observe that our SPT-LoRA has slightly higher fine-tuning memory than the full fine-tuning and LORA-16 which is taken up by updating the unstructurally-tuned parameters with sparse gradients in Eq. (5). Additionally, for memory-intense scenarios, one can employ SPT-LoRA w/o unstructured tuning for improved performance (0.8% Top-1 accuracy on VTAB-1k) and similar fine-tuning memory as LoRA, as shown in Figure 3 (b) and Table 6.

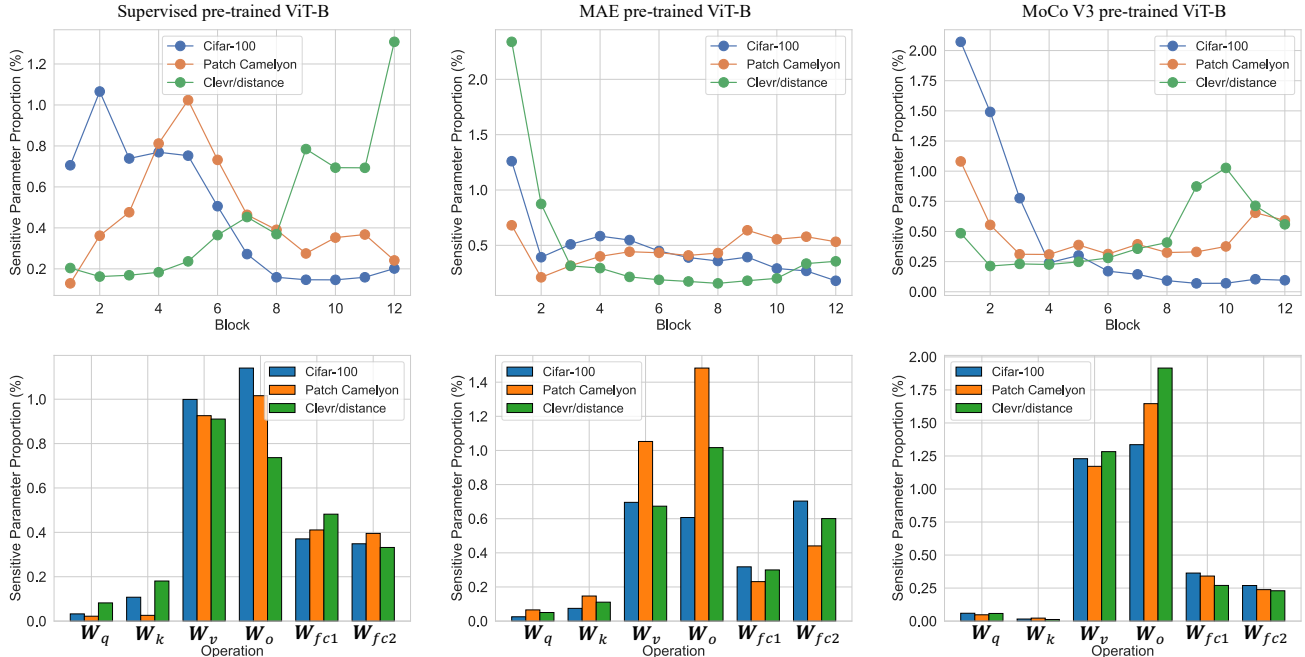


Figure 4: Parameter sensitivity patterns under 0.4M trainable parameter budget for ViT-B/16 backbone with different pre-training strategies on three sample tasks from VTAB-1k [61]. The proportions exhibit task-specific varying patterns in terms of network depth (upper figures) and task-agnostic similar patterns in terms of operations (lower figures).

4.4. Observations on Sensitivity Patterns

Our sensitivity criterion identifies task-specific important positions, which can reveal the contributions of the pre-trained weights to different downstream tasks during transfer learning. We visualize the proportions of the sensitive parameters for the supervised pre-trained ViT-B/16 backbone under 0.4M trainable parameter budget in Figure 4. First, we investigate the most sensitive blocks, whose numbers of sensitive parameters are summed and normalized over the 12 ViT-B/16 blocks. We observe that the patterns of the sensitive parameter proportions vary markedly across different tasks, which echoes the observations made in [16]. This suggests that we should not introduce trainable parameters to the same positions for each individual task but allocate trainable parameters at task-specific ones as we proposed. Next, we investigate the most insensitive weight matrices within a block. A ViT block consists of a query W_q , a key W_k , a value W_v , and an output W_o weight matrices in the multi-head self-attention layer and two weight matrices W_{fc1} and W_{fc2} in the feed-forward network as elaborated in [51, 12]. We observe that the query W_q and key W_k weight matrices have the lowest proportions of sensitive parameters for all three sample tasks. Since W_q and W_k are responsible for learning the attention scores which indicate the pairwise similarity among the patches, we speculate that although domain changes, the patch relationships learned during pre-training can be efficiently reused when transferred to downstream classification tasks.

5. Conclusion

In this paper, we have explored identifying and allocating trainable parameters to task-specific important positions for visual parameter-efficient tuning. Specifically, we have proposed a novel criterion to quickly measure the sensitivity of the pre-trained parameters for each specific task before fine-tuning. Based on the parameter sensitivity, we have proposed a trainable parameter allocation strategy that adaptively combines both unstructured and structured tuning under a desired trainable parameter budget, enabling high representational capability and flexibility. Finally, we have conducted extensive experiments on a total of 24 downstream recognition tasks with both plain and hierarchical vision Transformer backbones under different pre-training strategies to demonstrate the versatility and effectiveness of our proposed SPT. Notably, we have shown that our approach is complementary to the existing PEFT methods and improves their performance significantly. In the future, we will explore adapting large vision models to more downstream tasks with SPT, *e.g.*, dense prediction and vision-and-language tasks, and improve the training efficiency of SPT for on-device training [6, 30].

Acknowledgement. We thank Jing liu and Ziyi Liu for their helpful discussions. This research is partially supported by Monash FIT Start-up Grant. Dr. Jing Zhang is supported by the Australian Research Council project FL-170100117.

References

- [1] Y. Bao, Y. Li, S.-L. Huang, L. Zhang, L. Zheng, A. Zamir, and L. Guibas. An information-theoretic approach to transferability in task transfer learning. In *ICIP*, pages 2309–2313. IEEE, 2019. 3
- [2] G. Bender, P.-J. Kindermans, B. Zoph, V. Vasudevan, and Q. Le. Understanding and simplifying one-shot architecture search. In *ICML*, pages 550–559. PMLR, 2018. 2
- [3] A. Brock, T. Lim, J. M. Ritchie, and N. Weston. Smash: one-shot model architecture search through hypernetworks. In *ICLR*, 2018. 2
- [4] S. Caelles, K.-K. Maninis, J. Pont-Tuset, L. Leal-Taixé, D. Cremers, and L. Van Gool. One-shot video object segmentation. In *CVPR*, pages 221–230, 2017. 3
- [5] H. Cai, C. Gan, T. Wang, Z. Zhang, and S. Han. Once-for-all: Train one network and specialize it for efficient deployment. In *ICLR*, 2020. 2
- [6] H. Cai, C. Gan, L. Zhu, and S. Han. Tinytl: Reduce memory, not parameters for efficient on-device learning. *NeurIPS*, 33:11285–11297, 2020. 9
- [7] H. Cai, L. Zhu, and S. Han. ProxylessNAS: Direct neural architecture search on target task and hardware. In *ICLR*, 2019. 4
- [8] N. S. Chatterji, B. Neyshabur, and H. Sedghi. The intriguing role of module criticality in the generalization of deep networks. In *ICLR*, 2020. 3
- [9] S. Chen, C. Ge, Z. Tong, J. Wang, Y. Song, J. Wang, and P. Luo. Adaptformer: Adapting vision transformers for scalable visual recognition. *NeurIPS*, 2022. 1, 2, 4, 5
- [10] X. Chen, S. Xie, and K. He. An empirical study of training self-supervised vision transformers. In *ICCV*, pages 9640–9649, 2021. 5, 6
- [11] Y. Cui, Y. Song, C. Sun, A. Howard, and S. Belongie. Large scale fine-grained categorization and domain-specific transfer learning. In *CVPR*, pages 4109–4118, 2018. 3
- [12] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021. 1, 4, 5, 9
- [13] T. Gebru, J. Krause, Y. Wang, D. Chen, J. Deng, and L. Fei-Fei. Fine-grained car detection for visual census estimation. In *AAAI*, 2017. 5
- [14] D. Guo, A. Rush, and Y. Kim. Parameter-efficient transfer learning with diff pruning. In *ACL-IJCNLP*, 2021. 3
- [15] Y. Guo, Y. Li, L. Wang, and T. Rosing. Adafilter: Adaptive filter fine-tuning for deep transfer learning. In *AAAI*, pages 4060–4066, 2020. 3
- [16] Y. Guo, H. Shi, A. Kumar, K. Grauman, T. Rosing, and R. Feris. Spottune: transfer learning through adaptive fine-tuning. In *CVPR*, pages 4805–4814, 2019. 3, 9
- [17] T. Hao, H. Chen, Y. Guo, and G. Ding. Consolidator: Mergable adapter with group connections for vision transformer. In *ICLR*, 2023. 3
- [18] J. He, C. Zhou, X. Ma, T. Berg-Kirkpatrick, and G. Neubig. Towards a unified view of parameter-efficient transfer learning. In *ICLR*, 2022. 2
- [19] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R. Girshick. Masked autoencoders are scalable vision learners. In *CVPR*, pages 16000–16009, 2022. 5, 6
- [20] N. Houlsby, A. Giurigu, S. Jastrzebski, B. Morrone, Q. De Laroussilhe, A. Gesmundo, M. Attariyan, and S. Gelly. Parameter-efficient transfer learning for nlp. In *ICML*, pages 2790–2799, 2019. 1, 2, 4, 5, 7, 8
- [21] E. J. Hu, yelong shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen. LoRA: Low-rank adaptation of large language models. In *ICLR*, 2022. 1, 2, 3, 4, 5, 7, 8
- [22] M. Jia, L. Tang, B.-C. Chen, C. Cardie, S. Belongie, B. Hariharan, and S.-N. Lim. Visual prompt tuning. In *ECCV*, 2022. 1, 2, 5, 7, 8
- [23] S. Jie and Z.-H. Deng. Convolutional bypasses are better vision transformer adapters. *arXiv preprint arXiv:2207.07039*, 2022. 1, 2, 4
- [24] A. Khosla, N. Jayadevaprakash, B. Yao, and L. Fei-Fei. Novel dataset for fine-grained image categorization. In *CVPRW*, 2011. 5
- [25] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. *NeurIPS*, 25, 2012. 5, 7, 8
- [26] A. Kumar, A. Raghunathan, R. M. Jones, T. Ma, and P. Liang. Fine-tuning can distort pretrained features and underperform out-of-distribution. In *ICLR*, 2022. 2, 3
- [27] B. Lester, R. Al-Rfou, and N. Constant. The power of scale for parameter-efficient prompt tuning. In *EMNLP*, 2021. 2
- [28] X. L. Li and P. Liang. Prefix-tuning: Optimizing continuous prompts for generation. In *IJCNLP*, pages 4582–4597. Association for Computational Linguistics, 2021. 2
- [29] D. Lian, D. Zhou, J. Feng, and X. Wang. Scaling & shifting your features: A new baseline for efficient model tuning. *NeurIPS*, 2022. 3
- [30] J. Lin, L. Zhu, W.-M. Chen, W.-C. Wang, C. Gan, and S. Han. On-device training under 256kb memory. *arXiv preprint arXiv:2206.15472*, 2022. 9
- [31] H. Liu, K. Simonyan, and Y. Yang. DARTS: Differentiable architecture search. In *ICLR*, 2019. 4
- [32] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *ICCV*, 2021. 1, 5, 7
- [33] Z. Liu, H. Mao, C.-Y. Wu, C. Feichtenhofer, T. Darrell, and S. Xie. A convnet for the 2020s. In *CVPR*, pages 11976–11986, 2022. 7
- [34] I. Loshchilov and F. Hutter. Fixing weight decay regularization in adam, 2018. 5

- [35] P. Molchanov, A. Mallya, S. Tyree, I. Frosio, and J. Kautz. Importance estimation for neural network pruning. In *CVPR*, pages 11264–11272, 2019. 2, 4
- [36] M. M. Naseer, K. Ranasinghe, S. H. Khan, M. Hayat, F. Shahbaz Khan, and M.-H. Yang. Intriguing properties of vision transformers. *NeurIPS*, 34:23296–23308, 2021. 3
- [37] B. Neyshabur, H. Sedghi, and C. Zhang. What is being transferred in transfer learning? *NeurIPS*, 33:512–523, 2020. 2, 3
- [38] C. Nguyen, T. Hassner, M. Seeger, and C. Archambeau. Leep: A new measure to evaluate transferability of learned representations. In *ICML*, pages 7294–7305. PMLR, 2020. 3
- [39] M.-E. Nilsback and A. Zisserman. Automated flower classification over a large number of classes. In *ICVGIP*, pages 722–729. IEEE, 2008. 5
- [40] J. Plested and T. Gedeon. Deep transfer learning for image classification: a survey. *arXiv preprint arXiv:2205.09904*, 2022. 3
- [41] J. Puigcerver, C. Riquelme, B. Mustafa, C. Renggli, A. S. Pinto, S. Gelly, D. Keysers, and N. Houlsby. Scalable transfer learning with expert models. *arXiv preprint arXiv:2009.13239*, 2020. 3
- [42] M. Raghu, T. Unterthiner, S. Kornblith, C. Zhang, and A. Dosovitskiy. Do vision transformers see like convolutional neural networks? *NeurIPS*, 34:12116–12128, 2021. 3
- [43] M. T. Rosenstein, Z. Marx, L. P. Kaelbling, and T. G. Dietterich. To transfer or not to transfer. In *NeurIPS*, 2005. 3
- [44] R. K. Srivastava, K. Greff, and J. Schmidhuber. Training very deep networks. *NeurIPS*, 28, 2015. 2
- [45] X. Sun, R. Panda, R. Feris, and K. Saenko. Adashare: Learning what to share for efficient deep multi-task learning. *NeurIPS*, 33:8728–8740, 2020. 3
- [46] Y.-L. Sung, J. Cho, and M. Bansal. Lst: Ladder side-tuning for parameter and memory efficient transfer learning. *arXiv preprint arXiv:2206.06522*, 2022. 2
- [47] A. T. Tran, C. V. Nguyen, and T. Hassner. Transferability and hardness of supervised classification tasks. In *ICCV*, pages 1395–1405, 2019. 3
- [48] C.-H. Tu, Z. Mai, and W.-L. Chao. Visual query tuning: Towards effective usage of intermediate representations for parameter and memory efficient transfer learning. *arXiv preprint arXiv:2212.03220*, 2022. 2
- [49] E. Tzeng, J. Hoffman, N. Zhang, K. Saenko, and T. Darrell. Deep domain confusion: Maximizing for domain invariance. *arXiv preprint arXiv:1412.3474*, 2014. 8
- [50] G. Van Horn, S. Branson, R. Farrell, S. Haber, J. Barry, P. Ipeirotis, P. Perona, and S. Belongie. Building a bird recognition app and large scale dataset with citizen scientists: The fine print in fine-grained dataset collection. In *CVPR*, pages 595–604, 2015. 5
- [51] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. Attention is all you need. *NeurIPS*, pages 5998–6008, 2017. 9
- [52] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. The caltech-ucsd birds-200-2011 dataset. *Tech. Rep. CNS-TR-2011-001*, California Institute of Technology, 2011. 5
- [53] Z. Wang, Z. Dai, B. Póczos, and J. Carbonell. Characterizing and avoiding negative transfer. In *CVPR*, pages 11293–11302, 2019. 3
- [54] H. Wu, J. Xu, J. Wang, and M. Long. Autoformer: Decomposition transformers with auto-correlation for long-term series forecasting. *NeurIPS*, 34:22419–22430, 2021. 2
- [55] R. Xu, F. Luo, Z. Zhang, C. Tan, B. Chang, S. Huang, and F. Huang. Raise a child in large language model: Towards effective and generalizable fine-tuning. In *EMNLP*, 2021. 3
- [56] Y. Xu, Q. Zhang, J. Zhang, and D. Tao. Vitae: Vision transformer advanced by exploring intrinsic inductive bias. In *NeurIPS*, 2021. 1
- [57] J. Yoon, S. Arik, and T. Pfister. Data valuation using reinforcement learning. In *ICML*, pages 10842–10851. PMLR, 2020. 3
- [58] J. Yosinski, J. Clune, Y. Bengio, and H. Lipson. How transferable are features in deep neural networks? *NeurIPS*, 27, 2014. 2, 3
- [59] E. B. Zaken, Y. Goldberg, and S. Ravfogel. Bitfit: Simple parameter-efficient fine-tuning for transformer-based masked language-models. In *ACL*, pages 1–9, 2022. 3, 5
- [60] X. Zhai, J. Puigcerver, A. Kolesnikov, P. Ruysen, C. Riquelme, M. Lucic, J. Djolonga, A. S. Pinto, M. Neumann, A. 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, 8
- [61] X. Zhai, J. Puigcerver, A. Kolesnikov, P. Ruysen, C. Riquelme, M. Lucic, J. Djolonga, A. S. Pinto, M. Neumann, A. Dosovitskiy, et al. A large-scale study of representation learning with the visual task adaptation benchmark. *arXiv preprint arXiv:1910.04867*, 2019. 9
- [62] C. Zhang, S. Bengio, and Y. Singer. Are all layers created equal? *arXiv preprint arXiv:1902.01996*, 2019. 3
- [63] Y. Zhang, K. Zhou, and Z. Liu. Neural prompt search. *arXiv preprint arXiv:2206.04673*, 2022. 1, 2, 5
- [64] M. Zhao, T. Lin, F. Mi, M. Jaggi, and H. Schütze. Masking as an efficient alternative to finetuning for pretrained language models. In *EMNLP*, 2020. 3
- [65] S. Zheng, J. Lu, H. Zhao, X. Zhu, Z. Luo, Y. Wang, Y. Fu, J. Feng, T. Xiang, P. H. Torr, et al. Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers. In *CVPR*, pages 6881–6890, 2021. 7
- [66] B. Zhou, H. Zhao, X. Puig, S. Fidler, A. Barriuso, and A. Torralba. Scene parsing through ade20k dataset. In *CVPR*, pages 633–641, 2017. 7