AdaptFormer: Adapting Vision Transformers for Scalable Visual Recognition

Anonymous Author(s) Affiliation Address email

Abstract

Although the pre-trained Vision Transformers (ViTs) achieved great success in 1 2 computer vision, adapting a ViT to various image and video tasks is challenging 3 because of its heavy computation and storage burdens, where each model needs to be independently and comprehensively fine-tuned to different tasks, limiting 4 its transferability in different domains. To address this challenge, we propose an 5 effective adaptation approach for Transformer, namely AdaptFormer, which can 6 adapt the pre-trained ViTs into many different image and video tasks efficiently. 7 It possesses several benefits more appealing than prior arts. Firstly, AdaptFormer 8 9 introduces lightweight modules that only add less than 2% extra parameters to a ViT, while it is able to increase the ViT's transferability without updating its 10 original pre-trained parameters, significantly outperforming the existing 100% fully 11 fine-tuned models on action recognition benchmarks. Secondly, it can be plug-and-12 play in different Transformers and scalable to many visual tasks. Thirdly, extensive 13 experiments on five image and video datasets show that AdaptFormer largely 14 improves ViTs in the target domains. For example, when updating just 1.5% extra 15 parameters, it achieves about 10% and 19% relative improvement compared to the 16 fully fine-tuned models on Something-Something v2 and HMDB51, respectively. 17 Project page: anonymous-adaptformer.github.io/. 18

19 1 Introduction

There is a growing interest in adopting a general neural model to tackle a large variety of different 20 tasks since it benefits in reducing the need for task-specific model design and training. Recently, 21 Transformer [75] demonstrates great potential in this goal considering its success in various fields, 22 e.g., natural language processing (NLP) [26, 10, 76, 82], visual recognition [29, 73, 84, 59], dense 23 prediction [77, 11, 92, 90, 80], Generative Adversarial Network (GAN) [49, 45], reinforcement 24 learning (RL) [18, 16, 81], robotics [47, 24], and etc. However, existing literature in computer vision 25 tend to focus on the same network with task-specific weights scenario, where a single network is 26 used to train from scratch or fully fine-tune on a specific dataset, making it infeasible to maintain a 27 separate model weight for every dataset when the number of task grows, especially for the increasing 28 model capacity of state-of-the-art models (e.g., ViT-G/14 [87] with over 1.8 billion parameters). 29

Different from prior arts, we step into the direction of developing *same network with almost same weights* and achieve superior performance than the full-tuning approach by only tuning less than 2% parameters, with the remaining over 98% parameters shared across different tasks. There are two challenges to learning universal representations using a single model. The first one lies in the pre-training stage, which requires algorithms that can learn well-generalized representations that are easy to be applied to many tasks. Recent arts in self-supervised learning [12, 5, 40, 91, 79, 72, 32] can serve as a solution to this challenge. The second one, which is our main concern in this work, is

Submitted to 36th Conference on Neural Information Processing Systems (NeurIPS 2022). Do not distribute.

to build an effective pipeline that can adapt the model obtained at the pre-training stage to various
 downstream tasks by tuning parameters as less as possible and keeping the left parameters frozen.

While fine-tuning pre-trained models has been widely studied in NLP [6, 43, 65, 66, 55, 53, 44, 86, 39 58, 39], this topic is seldomly explored in the vision, where full tine-tuning of model parameters is 40 still the dominant strategy for adapting vision transformers. However, the full fine-tuning cannot 41 satisfy the goal of *universal representation* as it assigns an independent set of weights for every task. 42 Linear probing is a straightforward approach to maintaining the pre-trained model fixed by only 43 tuning a specific lightweight classification head for every task. However, linear probing tends to 44 have an unsatisfactory performance and misses the opportunity of pursuing strong but non-linear 45 features [40], which indeed benefit deep learning. More recently, Bahng et.al., [4] aimed to adapt 46 pre-trained models by modifying raw input pixel space. Jia et.al., [48] proposed Visual Prompt 47 Tuning (VPT) to adapt transformer models for downstream vision tasks, which prepends several 48 learnable parameters (prompts) to the patch embeddings and freezes the whole pre-trained backbone. 49

In this work, we propose a lightweight module, 50 namely AdaptFormer, to adapt vision transform-51 ers by updating the weights of AdaptFormer. We 52 introduce learnable parameters from the model 53 perspective, which is different from VPT, which 54 inserts learnable parameters into the token space. 55 Our AdaptFormer is conceptually simple yet 56 effective. It consists of two fully connected 57 layers, a non-linear activation function, and a 58 scaling factor. This module is set in parallel to 59 the feed-forward network (FFN) of the original 60 ViT model, as shown in Figure 2b. This design 61 is turned out to be effective for model transfer 62 when processing scalable visual tokens for both 63 image and video data (i.e., image data consists 64 of a small scale of visual tokens while video 65 data consists of a large scale). As shown in Fig-66 ure 1, compared with the full-tuning strategy, 67 AdaptFormer achieves comparable performance 68



Figure 1: **Parameter-Accuracy trade-off.** We leverage ViT-Base as backbone and report top-1 accuracy on SSv2 dataset. AdaptFormer can surpass full-tuning with only 0.2% tunable parameters. More detailed results are shown in Table 1.

on video recognition with only about 0.1% tunable parameters. Meanwhile, with less than 2% tunable
 parameters, AdaptFormer surpasses the full-tuning solution by about 10% on top-1 accuracy. Similar
 approaches are also proposed in fine-tuning pre-trained language models (PLMs) [6, 43, 66, 39].

The key contributions of this paper are summarized as follows: (1) We propose a simple yet effective 72 framework, namely AdaptFormer, for adapting vision transformers to a large variety of downstream 73 visual recognition tasks and avoiding catastrophic interference with each other. To the best of 74 our knowledge, this is the first work that explores efficient fine-tuning in video action recognition. 75 (2) We ablate many design choices and demonstrate the superior robustness of AdaptFormer when 76 parameters scale up. (3) Extensive experiments on various downstream tasks demonstrate that 77 78 AdaptFormer outperforms existing fine-tuning approaches significantly. By demonstrating the 79 effectiveness of AdaptFormer on multiple visual benchmarks, we hope our work could inspire the research communities to rethink the fine-tuning mechanism in computer vision and make progress 80 81 toward a flexible yet universal Transformer model for visual recognition.

82 2 Related Works

In the proposed AdaptFormer, we mainly introduce a plug-and-play module for efficiently fine-tuning the current vision Transformer models. In this section, we perform a literature review on related works from two perspectives, *i.e.*, the vision Transformers, and efficient transfer learning for vision Transformers.

87 2.1 Transformer in Vision

The Transformer architecture is first introduced in [75] and has re-energized the natural language processing (NLP) field from then on [26, 10]. Inspired by its huge success, researches in the computer vision filed have also evolved into Transformer era since ViTs [29]. The strong capability of modeling long-range relation has facilitated Transformer in various vision tasks, including image classification [29, 59, 57], object detection [11, 92, 21], semantic/instance segmentation [80], video

⁹³ understanding [8, 2, 31, 54], point cloud modeling [89, 38], 3D Object Recognition [19] and even

low-level processing [17, 56, 78]. Furthermore, transformers have advanced the vision recognition

performance by a large-scale pretraining [20, 63, 13, 33, 40, 72, 67]. In such a situation, given the

⁹⁶ pre-trained Transformer models, which are more larger than the previously prevalent CNN backbones,

one open question is how to fine-tune the big vision models so that they can be adapted into more

specific down-stream tasks. To solve the open question, we propose AdaptFormer to transfer ViTs

⁹⁹ from the pre-trained pre-texts into the target tasks in a more effective and efficient way.

100 2.2 Efficient Transfer learning for Transformers

Transfer learning targets re-adopting a pre-trained model (either via the supervised or the unsupervised 101 manner) as the starting point and further fine-tuning the specific model on a new task. In the NLP 102 field, transferring the large pre-trained language models (PLMs) [26, 10] into downstream tasks has 103 been the popular paradigm for a long time. Conventional arts [26, 10] set all the network parameters 104 105 as learnable ones and adapt them to the target tasks. However, with the growth of model sizes and the complexity of the specific tasks, the conventional paradigm is inevitably limited by the huge 106 computational burden. The NLP community has explored several ways for parameter-efficient transfer 107 learning that only set a few parameters learnable and fine-tune them for efficiency. The pioneer 108 works could be mainly categorized from the token [55, 53] and network perspectives [43, 44, 86, 37]. 109 Basically speaking, the token-related methods [53, 55] typically prepend several learnable prefix 110 vectors/tokens to the projected tokens within the multi-head self-attention layers (MHSA [75]). The 111 philosophy behind it is to assist the pre-trained models in understanding downstream tasks with the 112 guidance of extra token information. On the other hand, network-related methods [43, 44] integrate 113 shallow modules to improve the model transferability. The introduced modules adapt the produced 114 representations into the downstream tasks via features fusion. 115

Recently, with the emergence of a much more large-scale dataset [25, 68, 70, 62, 50], increasing 116 researchers in computer vision have adopted the homologous paradigm, *i.e.*, first pre-training and 117 then fine-tuning, to advance the vision tasks. As for the second stage, traditional methods typically 118 adopt the full-tuning arts in the downstream tasks. Rare attention has been drawn to the field of 119 efficient adaptation, especially in the field of vision Transformers. Inspired by Prompting in NLP, 120 [48] introduced the learnable tokens in exploring the efficient adaptation for ViTs. We empirically 121 found that the performance of prompting is hindered by the scale of tokens. That is to say, for the 122 tasks where the number of tokens is on a small scale, e.g., image classification, Prompting is efficient 123 for improving the model transferability. However, for larger scale tokens, e.g., video understanding, 124 Prompting presents limited potential. This observation motivates us to introduce AdaptFormer, which 125 is effective in the scenarios of scalable visual tokens. 126

127 **3** Approach

We propose AdaptFormer for efficiently transferring large pre-trained vision transformer models to downstream tasks, in both image and video domains. AdaptFormer attains strong transfer learning abilities by only fine-tuning a small number of extra parameters, circumventing catastrophic interference among tasks. We illustrate the overall framework of AdaptFormer in Figure 2b.

132 3.1 Preliminary and Notation

Vision Transformers (ViTs) are first introduced by [29] into vision recognition. A vanilla vision Transformer basically consists of a patch embedding layer and several consecutively connected encoders, as depicted in Figure 2a. Given an image $x \in \mathbb{R}^{H \times W \times 3}$, the patch embedding layer first splits and flatten the sample x into sequential patches $x_p \in \mathbb{R}^{N \times (P^2d)}$, where (H, W) represents the *height* and *width* of the input image, (P, P) is the resolution of each image patch, d denotes the output channel, and $N = HW/P^2$ is the number of image tokens. The overall combination of a prepended [CLS] token and the image tokens x_p are further fed into Transformer encoders for attention calculation.



(a) Full fine-tuning.

(b) AdaptFormer fine-tuning

Figure 2: Comparison of previous *full* and our *AdaptFormer* fine-tuning. AdaptFormer is conceptually simple by replacing the original MLP block with AdaptMLP, which consists of two branches, including the frozen branch (left) and the trainable down \rightarrow up bottleneck module (right).

Each Transformer encoder mainly consists of two types of sub-layers, *i.e.*, a multi-head self-attention layer (MHSA) and a MLP layer. In MHSA, the tokens are linearly projected and further re-formulated into three vectors, namely Q, K and V. The self-attention calculation is performed on Q, K and Vby:

$$x'_{\ell} = \operatorname{Attention}(\boldsymbol{Q}, \boldsymbol{K}, \boldsymbol{V}) = \operatorname{Softmax}(\frac{\boldsymbol{Q}\boldsymbol{K}^{\top}}{\sqrt{d}})\boldsymbol{V}, \tag{1}$$

where x'_{ℓ} are the tokens produced by MHSA at the ℓ -th layer. The output tokens x'_{ℓ} are further sent to

a LayerNorm [3] and a MLP block which is consisted of two fully connected layers with a GELU

147 activation [42] in between. This process is formally formulated as follows,

$$x_{\ell} = \mathrm{MLP}(\mathrm{LN}(x_{\ell}')) + x_{\ell}', \tag{2}$$

where x_{ℓ} is the output of the ℓ -th encoder block. At the last transformer layer, the [CLS] is utilized for the final object recognition. We refer the readers to find more details in [29]. In our work, we

replace the MLP layer with our AdaptMLP module for efficient fine-tuning purposes.

151 3.2 AdaptFormer

We propose a plug-and-play bottleneck module, namely AdaptMLP¹. We denote the vision Transformer equipped with AdaptMLP as AdaptFormer.

Architecture. The design principle of AdaptFormer is simple yet effective, which is illustrated in 154 Figure 2b. Compared to the vanilla full fine-tuning regime, AdaptFormer replaces the MLP block in 155 the transformer encoder with AdaptMLP, which is consisted of two sub-branches. The MLP layer in 156 the left branch is identical to the original network, while the right branch is an additionally introduced 157 lightweight module for task-specific fine-tuning. Specifically, the right branch is designed to be a 158 bottleneck structure for limiting the number of parameters purpose, which includes a down-projection layer with parameters $W_{\text{down}} \in \mathbb{R}^{d \times \hat{d}}$, an up-projection layer with parameters $W_{\text{up}} \in \mathbb{R}^{\hat{d} \times \hat{d}}$, where 159 160 \hat{d} is the bottleneck middle dimension and satisfies $\hat{d} \ll d$. In addition, there is a ReLU layer [1] 161 between these projection layers for non-linear property. This bottleneck module is connected to the 162 original MLP network (left branch) through the residual connection via a scale factor s. For a specific 163 input feature x_{ℓ} , the right branch in AdaptMLP produces the adapted features, \tilde{x}_{ℓ} , formally via: 164

$$\tilde{x}_{\ell} = \text{ReLU}(\text{LN}(x'_{\ell}) \cdot \boldsymbol{W}_{\text{down}}) \cdot \boldsymbol{W}_{\text{up}}.$$
(3)

165 Then both the features \tilde{x}_{ℓ} and x'_{ℓ} are fused with x_{ℓ} by residual connection,

$$x_{\ell} = \mathrm{MLP}(\mathrm{LN}(x_{\ell}')) + s \cdot \tilde{x}_{\ell} + x_{\ell}'.$$
(4)

¹In this paper, we use the term 'AdaptMLP' to denote the designed module and the term 'AdaptFormer' to represent the fine-tuning framework for Vision Transformers. Unless otherwise specified, we apply AdaptFormer to fine-tune the vanilla ViT backbone [29] in this paper.

Fine-tuning. During the fine-tuning phase, the original model parts (blue blocks in Figure 2b) 166 load weights from the pre-trained checkpoint and keeps untouched, avoiding interaction among 167 downstream tasks. The newly added parameters (orange blocks) are updated on the specific data 168 domain with the task-specific losses. 169

Inference. After fine-tuning, we still keep the shared parameters frozen as in the previous fine-170 tuning state, and additionally load the weights of the extra parameters that were fine-tuned in the 171 previous stage. The single overall model is able to be adapted to multiple tasks with the assistance of 172 lightweight introduced modules. 173

3.3 Discussion 174

Tunable parameters analysis. Our AdaptMLP module is lightweight. The total number of param-175 eters introduced to per layer is $2 \times d \times d + d + d$, which includes biases parameters. The middle 176 dimension \hat{d} is a small value compared with d (AdaptFormer still obtains a decent performance even 177 when $\hat{d} = 1$, as discussed in Sec. 4.4). Since most of the shared parameters are fixed and the number 178 of newly introduced parameters is small (< 2% of the pre-trained model parameters), the total model 179 size grows slowly when more downstream tasks are added. 180

Applicability. We note that AdaptMLP is a plug-and-play 181 module that can be adaptively inserted into existing popu-182 lar vision transformer architectures [29, 59, 77, 84, 22, 28] 183 since all of the backbones share the same MLP layers even 184 though they differ in the MHSA architectures (as shown in 185 Figure 2b). Compared to our methods, we notice that recent 186 prompt-related approaches insert trainable parameters into 187 188 the token space, as illustrated in Figure 3. They prepend learnable parameters either into the embedded tokens before 189 linear projection [55] or the key and value tokens after lin-190 ear projection [48]. Therefore, the prompt-related method 191 can not be straightforwardly adapted to special MHSA vari-192 ants, especially for the one that takes the pyramid spatial 193 information into account [59, 77]. Besides, we empirically 194 observe that prompt-related methods perform not well when 195 the number of patch tokens grows up from image to video 196 scale, as shown in Figure 1.



Figure 3: Prompt tuning illustration.

In summary, we present a strategy for tuning a pre-trained vision Transformer on a set of scalable 198 vision recognition tasks (e.g. image domain and video domain). It adds limited learnable parameters 199 for tuning while achieving comparable or even better performance than the full-tuning strategy. 200 Moreover, AdaptFormer could serve as a generic module for a large variety of recognition tasks. 201

Experiments 4 202

197

We evaluate the effectiveness of AdaptFormer by conducting extensive visual recognition experiments 203 in both the image and video domains. We first describe our experimental settings in Sec. 4.1, covering 204 the pre-trained backbones, baseline methods, downstream tasks and training details. We then compare 205 AdaptFormer with baseline methods and provide a thorough analysis in Sec. 4.2. In addition, we also 206 conduct ablation studies to explore different experimental configurations and explain what makes for 207 the superiority of AdaptFormer in Sec 4.4. 208

4.1 Experimental Settings 209

Pre-trained backbone. We adopt the plain Vision Transformer (ViT) [29], *i.e.*, ViT-Base (ViT-B/16) 210 as our backbone model and pre-train the model with both supervised and self-supervised approaches. 211

Specifically, for **image**, we directly use the ImageNet-21k [25] supervised pre-trained model² and 212

MAE [40] self-supervised model³. For video, we take both supervised and self-supervised pre-trained 213

²https://github.com/rwightman/pytorch-image-models/releases/download/v0.1-vitjx/jx_vit_base_patch16_ 224_in21k-e5005f0a.pth

Table 1: **Fine-tuning with self-supervised pre-trained model.** For tunable parameters, we also report the parameter percentage in the brackets. Besides, we report the top-1 accuracy on different dataset with the absolute value and the gap value relative to the *full-tuning* regime. [†] denotes $0.1 \times$ learning rate due to unstable training.

Mathod	Avg.		Image		Vic	leo
Wiethou	Params (M)	CIFAR-100	SVHN	Food-101	SSv2	HMDB51
Full-tuning	86.04 (100%)	85.90	97.67 [†]	90.09 [†]	53.97	46.41
Linear	0.07 (0.08%)	69.83 (-16.07)	66.91 (-30.76)	69.74 (-20.35)	29.23 (-24.74)	49.84 (+3.43)
VPT [48]	0.08 (0.09%)	82.44 (-3.46)	94.02 (-3.65)	82.98 (-7.11)	43.73 (-10.24)	52.67 (+6.26)
AdaptFormer-1	0.10 (0.12%)	83.52 (-2.38)	93.04 (-4.63)	83.64 (-6.45)	50.03 (-3.94)	51.68 (+5.27)
AdaptFormer-4	0.15 (0.17%)	84.83 (-1.07)	96.19 (-1.48)	85.42 (-4.67)	54.70 (+0.73)	51.81 (+5.40)
AdaptFormer-64	1.26 (1.46%)	85.90 (0.00)	96.89 (-0.78)	87.61 (-2.48)	59.02 (+5.05)	55.69 (+9.28)

models from VideoMAE [72]. More details about pre-training approaches and datasets can be found in Appendix.

Initialization of AdaptFormer. For the original networks, we directly load the weights pre-trained on the upstream tasks and keep them frozen/untouched during the fine-tuning process. For the newly added modules, the weights of down-projection layers are initialized with Kaiming Normal [41], while the biases of the additional networks and the weights of the up-projection layers are configured with zero initialization.

Baseline methods. We compare AdaptFormer with three commonly used fine-tuning approaches, including (1)*Linear probing:* adding an extra linear layer on top of the backbone and tuning the added parameters for evaluation. (2) *Full Fine-tuning:* setting all the parameters learnable and tuning them together. (3) *Visual Prompt Tuning (VPT):* [48] fine-tuning the extra token parameters as shown in Figure 3.

Downstream tasks. We evaluate our AdaptFormer on both image and video recognition tasks to verify its effectiveness. The specific datasets leveraged in this work are presented in the following.

• **Image domain :** CIFAR-100 [51] contains 50,000 training images and 10,000 validation images of resolution 32×32 with 100 labels. Street View House Numbers (SVHN) [34] is a digit classification benchmark dataset. In total, the dataset comprises over 600,000 labeled images, containing 73,257 training samples, 26,032 testing samples and 531,131 extra training data. The Food-101 [9] dataset consists of 101 food categories with a total of 101k images, including 750 training and 250 testing samples per category.

• Video domain : Something-Something V2 (SSv2) [36] is a large collection of video clips showing the people perform several normal actions in the daily life (*e.g.*, moving stuff and opening the door). It consists of 168,913 training samples, 24,777 validation samples and 27,157 testing samples, making a total of 220,847 videos with 174 labels. HMDB51 [52] is composed of 6,849 videos with 51 categories, making a split of 3.5k/1.5k train/val videos.

Implementation details. In this work, we use PyTorch toolkit [64] to conduct all experiments on NVIDIA V100 GPUs. Unless otherwise stated, we use 8×8 GPUs for video experiments and 1×8 GPUs for image experiments. Our default configurations follow the *linear probing* settings in [20, 40], which do *not* utilize many common regularization strategies, such as mixup [88], cutmix [85], color jittering and so on. More details can be found in Appendix.

244 **4.2 Main Properties and Analysis**

We compare the performance of different fine-tuning approaches in Table 1 with the backbones pretrained via the self-supervised paradigms. The results show that AdaptFormer consistently surpasses linear probing and Visual Prompt tuning (VPT) methods. Specifically, AdaptFormer outperforms VPT on image benchmark CIFAR-100, SVHN, and Food-101, by 9.50%, 5.95%, and 10.04% respectively. On the more challenging video action recognition dataset Something-Something V2, the superiority becomes even more significant, *i.e.*, about 34.96%. Note that even compared with the full fine-tuning

³https://dl.fbaipublicfiles.com/mae/pretrain/mae_pretrain_vit_base.pth





Figure 4: The trend of performance as the number of tunable parameters grows up. The accuracy of VPT drops dramatically when the parameter number exceeds task-specific value, while AdaptFormer is robust to the increasing parameters.

Figure 5: **Test accuracy of VPT [48] with different number of introduced tokens.** The optimization procedure becomes unstable when the token number is equal or larger than eight on HMDB51 dataset [52].

strategy, our AdaptFormer still outperforms by about 10% Top-1 accuracy on SSv2 dataset. To summarize, our AdaptFormer is highly parameter-efficient, as well as yielding good performance with parameter size at most 2% times than the full fine-tuning manner.

235 with parameter size at most 270 times than the run me tuning m

254 4.3 Scaling Tunable Parameters Up

Even though there are only limited parameters introduced, one might also argue that more tunable parameters of AdaptFormer contribute to its higher accuracy compared with VPT [48]. We conduct experiments to make a comprehensive discussion on this aspect.

As described in Sec. 3.3, the number of tunable parameters can be adjusted by changing the number of 258 introduced tokens for VPT, or the hidden feature dimension for AdaptFormer. As shown in Figure 4, 259 we conduct experiments with a wide range of tunable parameters on both SSv2 and HMDB-51 260 datasets. Since AdaptFormer and VPT share the same number of parameters of classification head 261 on a specific dataset, we only report the tunable parameters on the x-axis, which comes from the 262 visual prompts (VPT) or weight/bias of the down-up fully-connected layers (AdaptFormer), without 263 calculating the parameters of classification head. For VPT, the number of introduced tokens is chosen 264 from {1, 2, 4, 8, 16, 32, 48, 64}. Similarly, the number of hidden dimensions in AdaptFormer is in {1, 265 2, 4, 8, 16, 32}. AdaptFormer has a slight performance gain or maintains the accuracy stably when 266 the parameters scale up. On the contrary, the performance of VPT decreases dramatically when the 267 parameters exceed the task-specific value. Moreover, choosing the most suitable number of token 268 number becomes laborious since it might be task-specific (*i.e.* varying from one dataset to the other 269 one). For example, the accuracy of VPT keeps going up when the number of tunable parameters 270 increases up to 300K on SSv2, whereas it begins to drop when the number of tunable parameters 271 exceeds 50K on HMDB-51. 272

We further study the optimization procedures of VPT by monitoring the test accuracy of the training 273 stage. As shown in Figure 5, we gradually increase the number of tokens in VPT and plot the Top-1 274 accuracy of each epoch. The training stages are stable when the number of tokens is less than or equal 275 to 4, e.g., $\{1, 2, 4\}$. However, when the number becomes 8 or larger, e.g., $\{8, 16, 32\}$, the training 276 procedure collapses at about the tenth epoch and achieves poor performance at the end of the training 277 stage. On the contrary, the optimization procedures of AdaptFormer are stable when the number of 278 parameters varies across a large range, as shown in Table 2a. The top-1 accuracy fluctuates within 279 1.5% when the number of parameters increases from 0.44M (dim=16) to 4.87M (dim=256). 280

281 4.4 Ablation Studies

We ablate our AdaptFormer to study what properties make for a good AdaptFormer and observe several intriguing properties. The ablation studies conducted in this work are all performed on the SSv2 validation set [36]. Middle dimension. The middle dimension controls the number of introduced parameters by Adapt-Former. Lower middle dimensions introduce fewer parameters with a possible performance cost. We ablate AdaptFormer on the middle feature dimension to study this effects. As shown in Table 2a, the accuracy consistently improves when the middle dimension increases up to 64 and reaches the saturation point when the middle dimension is about 64. We note that our AdaptFormer can achieve a decent performance when the middle dimension reduces even to one, about 50.03% top-1 accuracy.

Table 2: AdaptFormer ablation experiments with ViT-B/16 on SSv2. We report the top-1 accuracy on the val set. Most suitable settings are marked in color.

(a) Midd	le dimensi	ion d.	(b)) AdaptM	LP inserted	layers	and form.	(c) Scalin	g factor s.
mid dim	#params	top-1		layers	form	#param	stop-1		factor	top-1
1	0.16M	50.03		$1 \rightarrow 6$	parallel	0.73	50.48		0.01	53.44
16	0.44M	57.62		$7 \rightarrow 12$	parallel	0.73	57.99		0.05	58.85
52 64	0.73M	58.27 59.02		$1 \rightarrow 12$	parallel	1.32	59.02		0.10	59.02
256	4.87M	58.87		$1 \rightarrow 12$	sequential	1.32	58.17		0.20	58.89

Scaling factor. The scaling factor *s* is introduced to balance the *task-agnostic* features (generated by the original frozen branch) and the *task-specific* features (generated by the tunable bottleneck branch). We evaluate AdaptFormer with multiple *s* values and the results are summarized in Table 2c. Different from the scaling factor in NLP field which prefer *s* larger than 1 (*e.g.*, s = 4 in [39]), we empirically found that the *s* should be < 1 for vision tasks, otherwise the fine-tuning would become unstable. Besides, we found that AdaptFormer achieves optimal performance with s = 0.1. A larger or smaller *s* would bring slight performance drop. Thus, we choose s = 0.10 as a default setting.

AdaptFormer position. As shown in Table 2b, we further ablate on the specific position to introduce the AdaptMLP block. We gradually increase the number of AdaptMLP layers with a step of three (start \rightarrow end, both included). We observe that the performance of AdaptFormer has a positive correlation with the number of added layers. In addition, AdaptFormer prefers the top part (the one far away from the input image) of the network to the bottom part when introducing the same number of layers, *e.g.*, AdaptFormer with 7 \rightarrow 12 obtains over 14.5% higher accuracy than 1 \rightarrow 6, though both equipped with six AdaptMLP layers.

Insertion form. We study the insertion formulation by comparing the *parallel* and *sequential* instances which are illustrated in Figure 6. As shown in Table 2b, the parallel AdaptFormer is able to outperform the sequential one by 0.85% top-1 accuracy. The reason might be: (1) the parallel design maintains the original feature using an independent branch and aggregating updated context by element-wise scaled sum; (2) the sequential design is equivalent to adding more layers, which might cause optimization difficulty. Therefore, we adopt the parallel design as our default setting due to its superiority.





Figure 6: **Illustration of the** *parallel* **and** *sequential* **insertion form**. Comparison results are shown in Table 2b.

Figure 7: **Performance with video frames number.** AdaptFormer outperforms VPT and linear fine-tuning.

Number of frames. The number of embedded patch tokens increases linearly with the number of video frames for the plain ViT [29]. We conduct experiments with the different number of frames, *i.e.*, {2, 4, 8} and the results are shown in Figure 7. We observe that increasing the number of frames is beneficial for all these three fine-tuning methods. However, AdaptFormer consistently outperforms

the linear manner (*e.g.*, +30% top-1 accuracy on 8 input frames) and VPT method(*e.g.*, +14% top-1 accuracy on 8 input frames).

318 4.5 Towards Visual Recognition Generalist Agent

In the above experiments, we typically utilize a modality-specific pre-trained checkpoint for the corresponding downstream tasks. For example, we use Kinetics-400 (video domain) pretrained model for downstream video action recognition on Something-Something V2 and HMDB-51 benchmarks. Besides, we use ImageNet-21K (image domain) pre-rained model for downstream image classification on CIFAR-100, SVHN and Food-101 benchmarks. Our AdaptFormer achieves superior performances in this *same network with modality-specific weights* scenario. Next, we take a further step to ask what would Table 3: Fine-tuning on video data with image

Next, we take a further step to ask what would happen if using *the same network with the modality-agnostic weights* for multiple tasks in the multi-modalities downstream tasks? Table 3: Fine-tuning on video data with image pre-trained model.

We use the model pre-trained on ImagNet-21k to do action recognition on SSv2. As shown in Table 3, AdaptFormer is robust to domain shift caused by modality. The experimental results show that the linear probe approach obtains a very poor accuracy (*i.e.*, 6.56% top-1 accuracy) when fine-tuning on SSv2. Meanwhile,

Mathod	Avg.	Fine-tuning	
Method	Params (M)	SSv2	
Full-tuning	86.36	41.50	
Linear	0.15	6.56	
VPT [48]	0.16	16.94	
AdaptFormer	1.33	46.06	

VPT [48] achieves a better performance than linear probe but it is not decent (*i.e.*, 16.94% top-1
 accuracy). Our AdaptFormer, compared to the above two methods, attains a promising 46.06% top-1
 accuracy, which is even higher than the full-tuning schedule (+4.56%).

339 4.6 Visualization



Figure 8: **t-SNE visualizations on SSv2 val dataset.** We extract the final classification features from the top linear layer for t-SNE visualizations. The top-1 accuracy is reported in red, while the relative parameter (compared to the full fine-tuning strategy) is reported in blue.

To evaluate the quality of the produced features, we conduct t-SNE [74] visualizations on Adapt-Former and other baseline methods. The features are extracted from the SSv2 validation set via the ViT-Base backbone. Figure 8 shows that the linear fine-tuning and the VPT methods tend to output mixed features as shown in Figure 8(a)-(b). Compared with the above two methods, the full fine-tuning strategy performs well in projecting features. However, it consumes huge computational sources to tune the whole network parameters. Figure 8(d) validates that our AdaptFormer facilitates ViT-Base in generating more separable representations with fewer learnable parameters.

347 **5** Conclusion

We present a conceptually simple yet effective framework, AdaptFormer, for efficiently adapting a 348 pre-trained Vision Transformer (ViT) backbone to scalable vision recognition tasks. By introducing 349 AdaptMLP, our AdaptFormer is able to fine-tune the lightweight modules for producing features 350 adapted to multiple downstream tasks. The extensive experiments on five datasets, covering both 351 the image and the video domains, validate that our proposed methods are able to increase the 352 ViT's transferability with little computational cost. We hope our work will inspire future research 353 in exploring more efficient fine-tuning methods for large vision models. One limitation is that 354 AdaptFormer is only employed in recognition tasks in this work, it's unclear whether it can work 355 well in tasks beyond recognition, e.g., object detection and semantic segmentation. We leave it for 356 the future exploration. Since our method is specially designed for efficient fine-tuning, we do not 357 foresee obvious undesirable ethical/social impacts at this moment. 358

359 **References**

- [1] Abien Fred Agarap. Deep learning using rectified linear units (relu). *arXiv preprint arXiv:1803.08375*, 2018. 4
- [2] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid. Vivit: A
 video vision transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*,
 pages 6836–6846, 2021. 3, 16
- [3] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. arXiv preprint
 arXiv:1607.06450, 2016. 4
- [4] Hyojin Bahng, Ali Jahanian, Swami Sankaranarayanan, and Phillip Isola. Visual prompting: Modifying
 pixel space to adapt pre-trained models. *arXiv preprint arXiv:2203.17274*, 2022. 2
- [5] Hangbo Bao, Li Dong, and Furu Wei. Beit: Bert pre-training of image transformers. arXiv preprint
 arXiv:2106.08254, 2021. 1
- [6] Ankur Bapna, Naveen Arivazhagan, and Orhan Firat. Simple, scalable adaptation for neural machine translation. *arXiv preprint arXiv:1909.08478*, 2019. 2
- [7] Emanuel Ben-Baruch, Tal Ridnik, Nadav Zamir, Asaf Noy, Itamar Friedman, Matan Protter, and Lihi
 Zelnik-Manor. Asymmetric loss for multi-label classification. *arXiv preprint arXiv:2009.14119*, 2020.
- [8] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is space-time attention all you need for video understanding. *arXiv preprint arXiv:2102.05095*, 2021. 3
- [9] Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101–mining discriminative components
 with random forests. In *European Conference on Computer Vision*, 2014. 6
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
 Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners.
 Advances in Neural Information Processing Systems, 2020. 1, 2, 3
- [11] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey
 Zagoruyko. End-to-end object detection with transformers. In *European Conference on Computer Vision*,
 2020. 1, 3
- [12] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand
 Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2021. 1
- [13] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand
 Joulin. Emerging properties in self-supervised vision transformers. In *IEEE/CVF International Conference on Computer Vision*, 2021. 3
- [14] Joao Carreira, Eric Noland, Andras Banki-Horvath, Chloe Hillier, and Andrew Zisserman. A short note
 about kinetics-600. *arXiv preprint arXiv:1808.01340*, 2018. 17
- [15] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset.
 In proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 6299–6308, 2017. 16
- [16] Chang Chen, Yi-Fu Wu, Jaesik Yoon, and Sungjin Ahn. Transdreamer: Reinforcement learning with
 transformer world models. *arXiv preprint arXiv:2202.09481*, 2022. 1
- [17] Hanting Chen, Yunhe Wang, Tianyu Guo, Chang Xu, Yiping Deng, Zhenhua Liu, Siwei Ma, Chunjing Xu,
 Chao Xu, and Wen Gao. Pre-trained image processing transformer. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021. 3
- [18] Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter Abbeel,
 Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning via sequence
 modeling. Advances in neural information processing systems, 34, 2021. 1
- [19] Shuo Chen, Tan Yu, and Ping Li. Mvt: Multi-view vision transformer for 3d object recognition. *arXiv preprint arXiv:2110.13083*, 2021. 3
- [20] Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised vision transformers. In *IEEE/CVF International Conference on Computer Vision*, 2021. 3, 6

- [21] Cheng Chi, Fangyun Wei, and Han Hu. Relationnet++: Bridging visual representations for object detection
 via transformer decoder. *Advances in Neural Information Processing Systems*, 2020. 3
- [22] Xiangxiang Chu, Zhi Tian, Yuqing Wang, Bo Zhang, Haibing Ren, Xiaolin Wei, Huaxia Xia, and Chunhua
 Shen. Twins: Revisiting the design of spatial attention in vision transformers. In *NeurIPS 2021*, 2021. 5
- Tat-Seng Chua, Jinhui Tang, Richang Hong, Haojie Li, Zhiping Luo, and Yantao Zheng. Nus-wide:
 a real-world web image database from national university of singapore. In *Proceedings of the ACM international conference on image and video retrieval*, pages 1–9, 2009. 18
- [24] Sudeep Dasari and Abhinav Gupta. Transformers for one-shot visual imitation. arXiv preprint
 arXiv:2011.05970, 2020. 1
- [25] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical
 image database. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2009. 3, 5, 16, 17
- [26] 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. 1, 2, 3
- [27] Carl Doersch, Abhinav Gupta, and Alexei A Efros. Unsupervised visual representation learning by context
 prediction. In *Proceedings of the IEEE international conference on computer vision*, pages 1422–1430,
 2015. 16
- [28] Xiaoyi Dong, Jianmin Bao, Dongdong Chen, Weiming Zhang, Nenghai Yu, Lu Yuan, Dong Chen, and
 Baining Guo. Cswin transformer: A general vision transformer backbone with cross-shaped windows,
 2021. 5
- 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. 1, 2, 3,
 4, 5, 8, 16, 17
- [30] Thibaut Durand, Nazanin Mehrasa, and Greg Mori. Learning a deep convnet for multi-label classification
 with partial labels. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*,
 pages 647–657, 2019. 18
- [31] Haoqi Fan, Bo Xiong, Karttikeya Mangalam, Yanghao Li, Zhicheng Yan, Jitendra Malik, and Christoph
 Feichtenhofer. Multiscale vision transformers. In *IEEE/CVF International Conference on Computer Vision*,
 2021. 3
- 437 [32] Christoph Feichtenhofer, Haoqi Fan, Yanghao Li, and Kaiming He. Masked autoencoders as spatiotemporal
 438 learners. *arXiv preprint arXiv:2205.09113*, 2022. 1
- [33] Chongjian Ge, Youwei Liang, Yibing Song, Jianbo Jiao, Jue Wang, and Ping Luo. Revitalizing cnn
 attention via transformers in self-supervised visual representation learning. *Advances in Neural Information Processing Systems*, 2021. 3
- Ian J Goodfellow, Yaroslav Bulatov, Julian Ibarz, Sacha Arnoud, and Vinay Shet. Multi-digit number recog nition from street view imagery using deep convolutional neural networks. *arXiv preprint arXiv:1312.6082*,
 2013. 6
- Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew
 Tulloch, Yangqing Jia, and Kaiming He. Accurate, large minibatch sgd: Training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017. 16
- Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzynska, Susanne Westphal,
 Heuna Kim, Valentin Haenel, Ingo Fruend, Peter Yianilos, Moritz Mueller-Freitag, et al. The" something
 something" video database for learning and evaluating visual common sense. In *IEEE/CVF International Conference on Computer Vision*, 2017. 6, 7
- [37] Demi Guo, Alexander M Rush, and Yoon Kim. Parameter-efficient transfer learning with diff pruning.
 arXiv preprint arXiv:2012.07463, 2020. 3
- [38] Meng-Hao Guo, Jun-Xiong Cai, Zheng-Ning Liu, Tai-Jiang Mu, Ralph R Martin, and Shi-Min Hu. Pct:
 Point cloud transformer. *Computational Visual Media*, 2021. 3
- [39] Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. Towards a unified
 view of parameter-efficient transfer learning. In *International Conference on Learning Representations*,
 2022. 2, 8

- [40] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. *arXiv preprint arXiv:2111.06377*, 2021. 1, 2, 3, 5, 6, 16
- [41] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing
 human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pages 1026–1034, 2015. 6
- [42] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). arXiv preprint arXiv:1606.08415,
 2016. 4
- [43] Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Ges mundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for nlp. In *International Conference on Machine Learning*, 2019. 2, 3
- 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, 3
- [45] Drew A Hudson and Larry Zitnick. Generative adversarial transformers. In *International Conference on Machine Learning*, pages 4487–4499. PMLR, 2021.
- [46] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing
 internal covariate shift. In *International conference on machine learning*, pages 448–456. PMLR, 2015. 16
- [47] Rishabh Jangir, Nicklas Hansen, Sambaran Ghosal, Mohit Jain, and Xiaolong Wang. Look closer:
 Bridging egocentric and third-person views with transformers for robotic manipulation. *IEEE Robotics and Automation Letters*, 2022. 1
- [48] Menglin Jia, Luming Tang, Bor-Chun Chen, Claire Cardie, Serge Belongie, Bharath Hariharan, and
 Ser-Nam Lim. Visual prompt tuning. *arXiv preprint arXiv:2203.12119*, 2022. 2, 3, 5, 6, 7, 9, 17, 18
- [49] Yifan Jiang, Shiyu Chang, and Zhangyang Wang. Transgan: Two pure transformers can make one strong
 gan, and that can scale up. *Advances in Neural Information Processing Systems*, 34, 2021.
- [50] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan,
 Fabio Viola, Tim Green, Trevor Back, Paul Natsev, et al. The kinetics human action video dataset. *arXiv preprint arXiv:1705.06950*, 2017. 3
- [51] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. *Master's thesis, Department of Computer Science, University of Toronto*, 2009. 6
- 488 [52] Hildegard Kuehne, Hueihan Jhuang, Estíbaliz Garrote, Tomaso Poggio, and Thomas Serre. Hmdb: a large
 video database for human motion recognition. In *IEEE/CVF International Conference on Computer Vision*,
 2011. 6, 7
- [53] Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning.
 arXiv preprint arXiv:2104.08691, 2021. 2, 3
- Kunchang Li, Yali Wang, Gao Peng, Guanglu Song, Yu Liu, Hongsheng Li, and Yu Qiao. Uniformer:
 Unified transformer for efficient spatial-temporal representation learning. In *International Conference on Learning Representations*, 2022. 3
- [55] Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. arXiv
 preprint arXiv:2101.00190, 2021. 2, 3, 5
- [56] Jingyun Liang, Jiezhang Cao, Guolei Sun, Kai Zhang, Luc Van Gool, and Radu Timofte. Swinir: Image
 restoration using swin transformer. In *IEEE/CVF International Conference on Computer Vision*, 2021. 3
- [57] Youwei Liang, Chongjian Ge, Zhan Tong, Yibing Song, Jue Wang, and Pengtao Xie. Not all patches are
 what you need: Expediting vision transformers via token reorganizations. *arXiv preprint arXiv:2202.07800*, 2022. 3
- [58] Xiao Liu, Kaixuan Ji, Yicheng Fu, Zhengxiao Du, Zhilin Yang, and Jie Tang. P-tuning v2: Prompt tuning
 can be comparable to fine-tuning universally across scales and tasks. *arXiv preprint arXiv:2110.07602*,
 2021. 2
- [59] 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*, 2021. 1, 3, 5, 17

- [60] Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang, Stephen Lin, and Han Hu. Video swin transformer.
 arXiv preprint arXiv:2106.13230, 2021. 16, 17
- [61] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016. 16
- [62] Dhruv Mahajan, Ross Girshick, Vignesh Ramanathan, Kaiming He, Manohar Paluri, Yixuan Li, Ashwin
 Bharambe, and Laurens Van Der Maaten. Exploring the limits of weakly supervised pretraining. In
 European Conference on Computer Vision, 2018. 3
- [63] Tian Pan, Yibing Song, Tianyu Yang, Wenhao Jiang, and Wei Liu. Videomoco: Contrastive video
 representation learning with temporally adversarial examples. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021. 3
- [64] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang,
 Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie
 Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In
 H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019. 6, 19
- [65] 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. 2
- [66] 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. 2
- [67] 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*, 2021. 3
- [68] Tal Ridnik, Emanuel Ben-Baruch, Asaf Noy, and Lihi Zelnik-Manor. Imagenet-21k pretraining for the
 masses. *arXiv preprint arXiv:2104.10972*, 2021. 3
- [69] Leslie N Smith. A disciplined approach to neural network hyper-parameters: Part 1–learning rate, batch
 size, momentum, and weight decay. *arXiv preprint arXiv:1803.09820*, 2018.
- [70] Chen Sun, Abhinav Shrivastava, Saurabh Singh, and Abhinav Gupta. Revisiting unreasonable effectiveness
 of data in deep learning era. In *IEEE/CVF International Conference on Computer Vision*, 2017. 3
- [71] Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. On the importance of initialization and
 momentum in deep learning. In *International conference on machine learning*, pages 1139–1147. PMLR,
 2013. 16
- [72] Zhan Tong, Yibing Song, Jue Wang, and Limin Wang. Videomae: Masked autoencoders are data-efficient
 learners for self-supervised video pre-training. *arXiv preprint arXiv:2203.12602*, 2022. 1, 3, 6, 16
- [73] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé
 Jégou. Training data-efficient image transformers & distillation through attention. *arXiv preprint arXiv:2012.12877*, 2020. 1
- [74] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(11), 2008. 9
- [75] 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*,
 2017. 1, 2, 3
- [76] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. GLUE:
 A multi-task benchmark and analysis platform for natural language understanding. In *International Conference on Learning Representations*, 2019. 1
- [77] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and
 Ling Shao. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions.
 In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 568–578, 2021. 1, 5
- [78] Zhendong Wang, Xiaodong Cun, Jianmin Bao, and Jianzhuang Liu. Uformer: A general u-shaped
 transformer for image restoration. *arXiv preprint arXiv:2106.03106*, 2021. 3

- [79] Chen Wei, Haoqi Fan, Saining Xie, Chao-Yuan Wu, Alan Yuille, and Christoph Feichtenhofer. Masked
 feature prediction for self-supervised visual pre-training. *arXiv preprint arXiv:2112.09133*, 2021.
- [80] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer:
 Simple and efficient design for semantic segmentation with transformers. *Advances in Neural Information Processing Systems*, 2021. 1, 3
- [81] Ruihan Yang, Minghao Zhang, Nicklas Hansen, Huazhe Xu, and Xiaolong Wang. Learning vision-guided
 quadrupedal locomotion end-to-end with cross-modal transformers. In *International Conference on Learning Representations*, 2022.
- [82] Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. Xlnet:
 Generalized autoregressive pretraining for language understanding. Advances in neural information
 processing systems, 32, 2019. 1
- [83] Yang You, Igor Gitman, and Boris Ginsburg. Large batch training of convolutional networks. *arXiv preprint arXiv:1708.03888*, 2017. 16
- [84] Li Yuan, Yunpeng Chen, Tao Wang, Weihao Yu, Yujun Shi, Zi-Hang Jiang, Francis EH Tay, Jiashi Feng,
 and Shuicheng Yan. Tokens-to-token vit: Training vision transformers from scratch on imagenet. In
 Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 558–567, 2021. 1, 5
- [85] Sangdoo Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix:
 Regularization strategy to train strong classifiers with localizable features. In *IEEE/CVF International Conference on Computer Vision*, 2019. 6
- [86] 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
- [87] Xiaohua Zhai, Alexander Kolesnikov, Neil Houlsby, and Lucas Beyer. Scaling vision transformers. In
 CVPR, 2022. 1
- [88] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk
 minimization. *arXiv preprint arXiv:1710.09412*, 2017. 6
- [89] Hengshuang Zhao, Li Jiang, Jiaya Jia, Philip HS Torr, and Vladlen Koltun. Point transformer. In
 Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 16259–16268, 2021. 3
- [90] Sixiao Zheng, Jiachen Lu, Hengshuang Zhao, Xiatian Zhu, Zekun Luo, Yabiao Wang, Yanwei Fu, Jianfeng
 Feng, Tao Xiang, Philip Torr, and Li Zhang. Rethinking semantic segmentation from a sequence-to sequence perspective with transformers. In *CVPR*, 2021. 1
- [91] Jinghao Zhou, Chen Wei, Huiyu Wang, Wei Shen, Cihang Xie, Alan Yuille, and Tao Kong. ibot: Image
 bert pre-training with online tokenizer. *International Conference on Learning Representations (ICLR)*,
 2022. 1
- [92] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable
 transformers for end-to-end object detection. *arXiv preprint arXiv:2010.04159*, 2020. 1, 3

595 Checklist

596	1.	For all authors
597 598		(a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
599		(b) Did you describe the limitations of your work? [Yes] Shown in Conclusion Section.
600 601		(c) Did you discuss any potential negative societal impacts of your work? [Yes] Shown in Conclusion Section.
602 603		(d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
604	2.	If you are including theoretical results
605		(a) Did you state the full set of assumptions of all theoretical results? [N/A]
606		(b) Did you include complete proofs of all theoretical results? [N/A]
607	3.	If you ran experiments
608 609 610		(a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] As a URL shown in the abstract.
611 612		(b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] Shown in supplementary materials.
613 614		(c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [N/A]
615 616		(d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] Please see Section 4.1
617	4.	If you are using existing assets (e.g., code, data, models) or curating/releasing new assets
618		(a) If your work uses existing assets, did you cite the creators? [Yes]
619		(b) Did you mention the license of the assets? [Yes] Shown in supplementary materials.
620		(c) Did you include any new assets either in the supplemental material or as a URL? [No]
621 622		(d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [Yes] We used publicly available datasets whose licenses allow research
623		usage.
624		(e) Did you discuss whether the data you are using/curating contains personally identifiable
625		information or offensive content? [No] To the best of our knowledge, the data we used
626		contains no personally identifiable information or offensive content.
627	5.	If you used crowdsourcing or conducted research with human subjects
628 629		(a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
630 631		(b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
632 633		(c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

634 A Appendix

In this supplementary material, we will include the details about the pre-training and fine-tuning processes, the extensive experiments of AdaptFormer on hierarchical vision transformers (*e.g.*, AdaptFormer-Swin), and the pseudo-code of AdaptMLP in a PyTorch-like style.

638 A.1 Experimental Settings

639 A.1.1 Pre-training Approaches

Image. We use MAE [40] as our self-supervised pre-training method in the image domain, a simple yet effective method that first masks nearly 75% patches of the input image and then reconstructs the missing pixels. Specifically, we directly adopt the checkpoint⁴ of ViT-B/16 for convenience, which is pre-trained on ImageNet-1K [25] for 800 epochs.

Video. We use VideoMAE [72] as our self-supervised pre-training method in the video domain, which is an direct extension of MAE to the video domain. VideoMAE utilizes the plain ViT [29] architecture of joint space-time attention mechanism [2, 60] and an extremely high proportion of masking ratio (*i.e.*, 90% to 95%) for pre-training. We also directly use the publicly available checkpoint⁵, which is pre-trained on Kinetics-400 [15].

649 A.1.2 Implementation Details of Fine-tuning

Configuration	Image	Video	
optimizer	SGD		
base learning rate	0.1		
weight decay	0		
optimizer momentum	0.9		
batch size	1024 images	/frames	
learning rate schedule	cosine deca	y [61]	
GPU numbers	8	64	
warmup epochs	20	10	
training epochs	100	90	
augmentation	RandomResizedCrop [40]	MultiScaleCrop [72]	

Table 4: **Fine-tuning settings.** We present the shared configurations, like the optimizer and the base learning rate, the upper part, and show the seperated ones in the lower part.

The implementation details are summarized in Table 4. The video experiments are conducted on 64 Tesla V100 GPUs, while the image experiments are performed on 8 Tesla V100 GPUs. For the optimizer, different from [83] that adopts LARS, we leverage SGD [71] for stable training on small-scale dataset (*e.g.*, CIFAR10). The actual learning rate is calculated by: $lr = base_lr \times batchsize$ / 256 following the linear *lr* scaling rule [35]. More detailed training configurations are presented in Table 4, including the batchsize, learning rate schedule and etc.

The experimental settings of image and video mainly follow the ones utilized in MAE [40] and VideoMAE [72], respectively. We insert an extra BatchNorm layer [46] without affine transformation (*i.e.* affine=False) before the final fully connected layer, following the common practice to normalize the pre-trained features [27, 40]. In addition, there is *no* flip augmentation during the

660 fine-tuning stage for video data.

⁴https://dl.fbaipublicfiles.com/mae/pretrain/mae_pretrain_vit_base.pth

⁵https://drive.google.com/file/d/1tEhLyskjb755TJ65ptsrafUG2llSwQE1/view?usp= sharing

661 A.2 More Supplementary Results

662 A.2.1 AdaptFormer with Supervised Pre-training

In addition to the self-supervised pre-training presented in the main paper, we also evaluate Adapt-Former with the supervised pre-trained model. The results in Table 5 show that AdaptFormer still outperforms linear probe and VPT obviously. In addition, AdaptFormer surpasses full-tuning on four benchmarks (CIFAR100, SVHN, SSv2, HMDB51) with only 1.46% parameters. On the remaining benchmark (Food-101), AdaptFormer achieves an almost comparable performance to full-tuning (90.89% v.s. 90.96%).

Table 5: **Fine-tuning with** *supervised* **pre-trained model.** We report the tunable parameters percentage in the brackets. Besides, we report the top-1 accuracy on different dataset with the absolute value and the gap value relative to the *full-tuning* regime.

Mathad	Avg.		Image		Vi	deo
Method	Params (M)	CIFAR-100	SVHN	Food-101	SSv2	HMDB51
Full-tuning	86.04 (100%)	89.12	95.41	90.96	53.62	59.38
Linear	0.07 (0.08%)	85.95 (-3.17)	55.36 (-40.05)	88.14 (-2.82)	35.49 (-18.13)	70.31 (+10.93)
VPT [48]	0.08 (0.09%)	90.97 (+1.85)	92.77 (-2.64)	90.16 (-0.80)	55.22 (+1.60)	71.56 (+12.18)
AdaptFormer-64	1.26 (1.46%)	91.86 (+2.73)	97.29 (+1.88)	90.89 (-0.07)	60.18 (+6.56)	73.21 (+13.83)

669 A.2.2 AdaptFormer on Swin Transformer

Settings. We further demonstrate the effectiveness of AdaptFormer on hierarchical vision transformers, *e.g.*, Swin [59, 60]. We name AdaptFormer applied to Swin as *AdaptFormer-Swin*, to distinguish plain AdaptFormer (without any suffix) which is applied to the vanilla ViT [29]. It is noted that we can adopt AdaptMLP to Swin easily without any special modification as Swin and ViT share the same MLP architecture. However, VPT [48] needs additional handcraft designs to be suitable for the shifted local windows in the prevalent hierarchical vision transformers, which hinders its general applications.

We utilize Swin-B [59] and the video counterpart [60] for image and video. Similarly, we also directly use the officially provided checkpoints⁶, which are pre-trained on ImageNet-21K [25] and

679 Kinetics-600 [14].

Table 6: **Fine-tuning with** *Swin* **Transformer**. We utilize Swin-B [59] and Video Swin-B [60] for **image** and **video** experiments, respectively. Parameter percentage and performance difference are reported relative to *full-tuning* schedule.

-						
Mathod	Avg.		Image		Vi	deo
Wiethou	Params (M)	CIFAR-100	SVHN	Food-101	SSv2	HMDB51
Full-tuning	87.19 (100%)	89.95	97.03	91.43	52.92	68.73
Linear	0.11 (0.13%)	89.07 (-0.88)	69.06 (-27.97)	90.64 (-0.79)	28.32 (-24.61)	74.00 (+5.27)
AdaptFormer-Swin	1.25 (1.43%)	91.88 (+1.93)	97.31 (+0.28)	91.86 (+0.43)	54.09 (+1.17)	74.65 (+5.92)

Results. Since VPT is not applicable in Swin, we do not report its performance. Table 6 shows

AdaptFormer-Swin performs well compared with other tuning strategies. For image benchmarks, our

method can outperform full-tuning approach with only 1.43% parameters. Moreover, AdaptFormer-

⁶⁸³ Swin surpasses linear probing by a significant margin, especially on the challenging dataset, SSv2.

The results validate that AdaptFormer is able to generally boost the transferability of various vision

685 Transformer variants.

⁶Image: https://github.com/SwinTransformer/storage/releases/download/v1.0.4/swin_base_patch244_window877_kinetics600_22k.pth

Video: https://github.com/SwinTransformer/storage/releases/download/v1.0.0/swin_ base_patch4_window7_224_22k.pth

686 A.3 Multi-Label Classification

We further conduct experiments on dataset with larger scale and diversity. Specifically, we evaluate AdaptFormer on NUS-WIDE [23] for multi-label classification. **NUS-WIDE** contains 269,648 images collected from Flicker, which are annotated with 81 visual concepts. Since some images are not available on Flicker, we only use 220,000 images following [7, 30]. We utilize mean average precision (mAP) as performance metric.

Settings: Our training settings mainly follow ASL [7]. Specifically, We trained all models for 40
 epochs using Adam optimize and 1-cycle learning rate policy [69]. The maximal learning rate is
 0.001.

Method	Params (M)	NUS-WIDE [23]
Full-tuning	85.86 (100%)	61.26
Linear	0.06 (0.08%)	51.19 (-27.25)
VPT [48]	0.07 (0.09%)	57.08 (-7.56)
AdaptFormer-1	0.09 (0.12%)	57.51 (-4.08)
AdaptFormer-4	0.15 (0.17%)	58.14 (-2.13)
AdaptFormer-64	1.25 (1.46%)	59.07 (-0.06)

Table 7: AdaptFormer for multi-label classification.

695 A.4 Possible Architectures

We explore other possible architectures utilized in AdaptFormer. Specifically, we further replace the MLP architectures within the AdaptMLP module by the convolution layer, depthwise convolution layer, and LayerNorm layer. For fair comparisons, we carefully design the above modules to meet the comparable number of parameters (1.3M). The experimental results of different adapter modules are shown in Table 8, which validates that the simple MLP modules are simple yet effective compared with the other architectures. For example, our AdaptMLP module surpasses the AdaptConv module by 0.55% Top1 accuracy on SSv2 dataset.

Table 8: **Fine-tuning with different adapter modules**. We use AdaptConv to denote the designed adapter module with convolution layers, while AdaptDepthwise-Conv is utilized to denote the designed adapter module with depthwise convolution layers. Besides, we also replace the MLP architectures with LayerNorm layer as AdaptLayerNorm-In.

Methods	Avg Parameters	SSv2 Top1	NUS-WIDE mAP	CIFAR100 Top1
AdaptMLP	1.28	59.02	59.07	85.93
AdaptConv	1.39	58.47	58.86	85.42
AdaptDepthwise-Conv	1.29	58.15	58.73	85.37
AdaptLayerNorm-In	1.30	57.85	58.51	85.71

703 A.5 Evaluation on ImageNet-1k datasets

To further conduct the evaluation on the ImageNet dataset, we directly load the weights pre-trained on ImageNet-21K, and evaluate the AdaptFormer on ImageNet-1k. The results in Table 9 validate that the performance of AdaptFormer is comparable to the full fine-tuning strategy with only 1.5% tunable parameters. Meanwhile, our AdaptFormer outperforms the linear probe and VPT by 0.81% and 0.14%, respectively.

709 A.6 Extended experiments on middle dimension

We conduct the extended ablation studies on the middle dimension design in this sub-section. We aim to seek for a trade-off between model capacity (i.e., potential) and adaptation efficiency. In fact, the middle dimension has a main influence on the parameter size of adapter. The higher dimension

	1	0
Methods	Parameters (M)	ImageNet-1k Top1 (%)
Full Fine-tuning	86.57	82.26
Liner	0.77	80.95
VPT	0.78	81.68
AdaptFormer	1.96	57.85

Table 9: **Fine-tuning with AdaotFormer on ImageNet-1k dataset**. We load the weights pretrained on ImageNet-21K and evaluate the classification performance on ImageNet-1K.

⁷¹³ brings more parameters while the efficiency and storage are limited. As shown in Table 10, we

evaluate several numbers of middle dimension and found that using 64 is optimal to achieve accuracy,

715 light-weight storage, and efficiency.

Table 10: AdaptFormer ablation experiments with ViT-B/16 on SSv2. The experimental results on middle dimension are investigated.

Middle Dimension	Parameters (M)	SSv2 Top1 (%)	NUS-WIDE mAP (%)
1	0.16	50.03	57.51
4	0.22	54.70	58.14
16	0.44	57.62	59.00
32	0.73	58.27	59.09
64	1.32	59.02	59.07
128	2.51	58.95	59.49
256	4.87	58.87	59.62
512	9.59	58.98	59.82

716 A.7 Analysis on the fine-tuning time and inference latency

To analysis the computational efficiency, we compare the fine-tuning time and inference time on a 717 single NVIDIA A100-40G GPU. We utilize SSv2 video classification for this part. For fine-tuing, 718 we experiment with batchsize of 32. For inference, we test the latency with multiple batch sizes 719 to get a comprehensive comparison under various inference scenarios. All the time is measured in 720 milliseconds averaged over 100 trials. The results are summarized in Table 11 and Table 12. As 721 shown in Table 11, AdaptFormer only costs less than a half of the fine-tuning time compared with the 722 full-tuning. Moreover, AdaptFormer significantly outperforms linear probing in terms of accuracy 723 with a slight longer fine-tuning time. For inference, AdaptFormer introduce a negligible FLOPs and 724 latency compared with the Linear/Full-tuning.

Table 11: Fine-tuning time of a single forward-backward step averaged over 100 trials.

Methods	Latency (B=32)
Full-tuning	355.0 ms
Linear	140.2 ms
VPT	210.3 ms
AdaptFormer	162.2 ms

725

726 A.8 Implementation

The core part of AdaptFormer is replacing the original MLP with AdaptMLP, which consists of the frozen original MLP and newly introduced Down \rightarrow ReLU \rightarrow Up layers, which are tunable at the fine-tuning stage. Algorithms 1 provides the implementation of AdaptMLP written in PyTorch [64].

⁷³⁰ For more implementation details, please refer to the provided source code.

Table 12: Inference time of a single forward step averaged over 100 trials.

Methods	Flops (B=1)	Latency (B=1)	Latency (B=16)	Latency (B=32)
Linear/Full-tuning	78.915G	11.1 ms	22.4 ms	42.3 ms
VPT	79.029G	11.3 ms	22.9 ms	42.4 ms
AdaptFormer	79.840G	11.9 ms	23.2 ms	42.8 ms

Algorithm 1 Implementation of AdaptMLP in PyTorch-like style.

```
class AdaptMLP(nn.Module):
   def __init__(self, original_mlp, in_dim, mid_dim, dropout=0.0, s=0.1):
       super().__init__()
       self.original_mlp = original_mlp # original MLP block
       # down --> non linear --> up
       self.down_proj = nn.Linear(in_dim, mid_dim)
       self.act = nn.ReLU()
       self.up_proj = nn.Linear(mid_dim, in_dim)
       self.dropout = nn.Dropout(dropout)
       self.scale = s # scaling factor
       # initialization
       nn.init.kaiming_uniform_(self.down_proj.weight)
       nn.init.zeros_(self.up_proj.weight)
       nn.init.zeros_(self.down_proj.bias)
       nn.init.zeros_(self.up_proj.bias)
       # freeze original MLP
       for _, p in self.original_mlp.named_parameters():
    p.requires_grad = False
   def forward(self, x):
       down = self.down_proj(x)
       down = self.act(down)
       down = self.dropout(down)
       up = self.up_proj(down)
       output = self.original_mlp(x) + up * self.scale
       return output
```