Frozen CLIP Models are Efficient Video Learners

Ziyi Lin¹, Shijie Geng³, Renrui Zhang², Peng Gao^{*2}, Gerard de Melo⁴, Xiaogang Wang¹, Jifeng Dai⁵, Yu Qiao², and Hongsheng Li^{1,6}

¹ Multimedia Laboratory, The Chinese University of Hong Kong

² Shanghai AI Laboratory
 ³ Rutgers University
 ⁴ Hasso Plattner Institute
 ⁵ SenseTime Research
 ⁶ Centre for Perceptual and Interactive Intelligence Limited

⁵ SenseTime Research zylin@link.cuhk.edu.hk, gaopeng@pjlab.org.cn, hsli@ee.cuhk.edu.hk

Abstract. Video recognition has been dominated by the end-to-end learning paradigm – first initializing a video recognition model with weights of a pretrained image model and then conducting end-to-end training on videos. This enables the video network to benefit from the pretrained image model. However, this requires substantial computation and memory resources for finetuning on videos and the alternative of directly using pretrained image features without finetuning the image backbone leads to subpar results. Fortunately, recent advances in Contrastive Vision-Language Pre-training (CLIP) pave the way for a new route for visual recognition tasks. Pretrained on large open-vocabulary image-text pair data, these models learn powerful visual representations with rich semantics. In this paper, we present Efficient Video Learning (EVL) – an efficient framework for directly training high-quality video recognition models with frozen CLIP features. Specifically, we employ a lightweight Transformer decoder and learn a query token to dynamically collect frame-level spatial features from the CLIP image encoder. Furthermore, we adopt a local temporal module in each decoder layer to discover temporal clues from adjacent frames and their attention maps. We show that despite being efficient to train with a frozen backbone, our models learn high quality video representations on a variety of video recognition datasets. Code is available at https://github.com/OpenGVLab/efficient-video-recognition.

Keywords: Video recognition; Efficient learning; Vision-language model; Spatiotemporal Fusion

1 Introduction

As a fundamental component of video understanding, learning spatiotemporal representations remains an active research area in recent years. Since the beginning of the deep learning era, numerous architectures have been proposed to learn spatiotemporal semantics, such as traditional two-stream networks [35,40,53], 3D convolutional neural networks [36,5,17,30,38,44,42,11,10], and spatiotemporal Transformers [3,28,27,8,1,24,46]. As videos are high-dimensional and exhibit substantial spatiotemporal redundancy, training video recognition models from

^{*} Corresponding author.



Fig. 1: Left: illustration of the difference between our EVL training pipeline and other video recognition methods. Right: despite that EVL targets efficient training, our models set new accuracy vs. *inference FLOPS* Pareto frontiers. On Kinetics-400, the 8-frame ViT-B/16 model achieves 82.9% top-1 accuracy with only 60 V100 GPU-hours of training.

scratch is highly inefficient and may lead to inferior performance. Intuitively, the semantic meaning of a video snippet is highly correlated with each of its individual frames. Previous studies [5,3,1,46] have shown that the datasets and methodologies for image recognition can benefit video recognition as well. Owing to the close relationship between image and video recognition, as a routine practice, most existing video recognition models take advantage of pretrained image models by using them for initialization and then re-training all parameters for video understanding in an end-to-end manner.

However, the end-to-end finetuning regime has two major drawbacks. The first is *efficiency*. Video recognition models are required to process multiple frames simultaneously and are several times larger than their image counterparts in terms of model size. Finetuning the entire image backbone inevitably incurs an enormous computation and memory consumption cost. As a result, this issue limits the adoption and scalability of some of the largest image architectures for video recognition under restricted computational resources. The second issue is known as *catastrophic forgetting* [29] in the context of transfer learning. When conducting end-to-end finetuning on downstream video tasks, we risk destroying the powerful visual features learned from image pretraining and obtaining subpar results if the downstream videos are insufficiently informative. Both concerns suggest that end-to-end finetuning from pre-trained image models is not always an ideal choice, which calls for a more efficient learning strategy to transfer knowledge from images to videos.

Considerable efforts have been made on learning high-quality and general visual representations through contrastive learning [31,21], masked vision modeling [18,45,2], and traditional supervised learning [48,32]. Masked vision modeling approaches such as MAE [18] train an encoder-decoder architecture to reconstruct the original image from the latent representation and mask tokens. Supervised learning-based methods train image backbones with a fixed set of predefined category labels. Since they are usually trained uni-modally, they both lack the ability to represent rich semantics. In contrast, contrastive vision–language models such as CLIP [31] are pretrained with large-scale open-vocabulary image–text pairs. They can learn more powerful visual representations aligned with much richer language semantics. Another advantage of CLIP is its promising feature transferability, which forms a strong foundation for a series of transfer learning methods on various downstream tasks [13,51,54,52,34,22].

The above reasons inspire us to rethink the relationship between image and video features and devise efficient transfer learning methods to make use of frozen CLIP image features for video recognition. To this end, we propose an Efficient Video Learning (**EVL**) framework based on a lightweight Transformer decoder [39]. The difference between EVL and other video recognition models is illustrated in Fig. 1 Left. Specifically, EVL learns a query token to dynamically gather frame-level spatial features from each layer of the CLIP image encoder. On top of that, we introduce a local temporal module to collect temporal cues with the help of temporal convolution, temporal positional embeddings, and cross-frame attention. Finally, a fully-connected layer is used to predict scores of video categories. We conduct extensive experiments to show the effectiveness of our method and find EVL to be a simple and effective pipeline with higher accuracy but lower training and inference costs, as shown in Fig. 1 **Right**. Our contributions are as follows:

- We point out the shortcomings of the current end-to-end learning paradigm for video understanding and propose to leverage frozen CLIP image features to facilitate video recognition tasks.
- We develop EVL an efficient transfer learning pipeline from image to video recognition, in which we train a lightweight Transformer decoder module on top of *fixed* transferable image features to perform spatiotemporal fusion.
- Extensive experiments demonstrate the effectiveness and efficiency of EVL. It incurs much shorter training time than end-to-end finetuning, yet achieves competitive performance. This makes video recognition accessible to a broader community with average computation resources.

2 Related Work

Video Recognition. Recent advances in video recognition can be divided into two major directions – improving model architectures and proposing new training strategies. Following the success of Transformers in image recognition, video recognition has as well seen a transition from 3D-CNN [5,11,10] to Transformer-based architectures [3,8,27,26]. Uniformer [24] is a custom fused CNN-Transformer architecture achieving good speed–accuracy trade-off. Yan et al. [46] propose a multi-stream Transformer operating on different resolutions with lateral connections. Prior work [5,3,46] has shown the benefit of image pretraining for



Fig. 2: Model architecture overview. (a) Top-level architecture: multiple intermediate feature maps from a massively pretrained image backbone are fed into a Transformer decoder to gather information from them. (b) Motion-enhanced Transformer decoder block: temporal modeling is added on top of raw frame features X_i to retain structural information of the spatiotemporal features.

video recognition tasks. However, the end-to-end finetuning remains expensive, especially due to the large memory footprint. In terms of new training strategies, pretext task design for self-supervised learning [12,43] and multi-task co-training [15,47,50] are two mainstream directions. However, both are even more expensive than regular supervised training. Unlike previous efforts, we leverage fixed CLIP image features and directly learn an efficient video recognition model with an additional Transformer encoder.

Large-scale Image Representation Learning. With the availability of webscale weakly labeled data, we have witnessed a surge of new models for general visual representation learning. Image models built with regular supervised learning have grown dramatically in size. For example, Zhai et al. [48] train a ViT-G model on the large JFT-3B dataset. Riquelme et al. [32] create a Mixture-of-Experts vision model that scales to over 10 billion parameters. To further boost the visual representation power, efforts began to focus on large-scale contrastive learning and self-supervised learning. The success of BERT [7] sparked an emerging direction of building large-scale vision models with masked vision modeling [18,45,2]. Meanwhile, CLIP [31] and ALIGN [21] pretrain vision–language models with a contrastive loss on large-scale datasets consisting of open-vocabulary image–text pairs. The multimodal pretraining environment makes them suitable for downstream tasks requiring rich semantics.

Efficient Transfer Learning. This set of work is most related to our method. Most previous works on efficient transfer learning is for natural language processing and image recognition. Some methods learn parameter-efficient weight difference vectors during finetuning, exploiting sparsity [16] or low rank decomposition [20]. A collection of approaches [19,13,51] train *adapters*, which are additional fully-connected layers with residual connections, keeping the origi-

nal weights in the pretrained model fixed. Another line of methods [25,23,54] learn *prompts*, which are additional learnable tokens appended to the input or intermediate feature sequence for task-specific adaption. While these category of methods share the same motivation as ours, we employ Transformer decoders, which is more flexible and also efficient to train, as we will analyze in the Methods section. In terms of video recognition, the exploration in efficient transfer learning is still limited. Ju at el. [22] transfer CLIP models to video recognition by learning prompts and temporal modeling. Wang et al. [41] utilize CLIP models for video recognition by traditional end-to-end finetuning. We will compare with them in the Experiments section. There are also several works utilizing transferable image features for video-text tasks [6,9,14], but these works focus more on cross-modality modeling. In contrast, our work aims to improve the single-modal video representations, which should be complementary to most of the video-text learning methods.

3 Our Method

The three primary goals of our image to video transfer learning pipeline are (1) capability to summarize multi-frame features and infer video-level predictions; (2) capability to capture motion information across multiple frames; and (3) efficiency. We thus propose the Efficient Video Learning (**EVL**) framework, which we detail in the following.

3.1 Overall Structure

The overall structure of EVL, as illustrated in Fig. 2, is a multi-layer spatiotemporal Transformer decoder on top of a fixed CLIP backbone. The CLIP backbone extracts features from each frame independently. The frame features are then stacked to form a spatiotemporal feature volume, modulated with temporal information, and fed into the Transformer decoder. The Transformer decoder performs global aggregation of multi-layer features: a video-level classification token [CLS] is learned to act as query, and multiple feature volumes from different backbone blocks are fed to the decoder blocks as key and value. A linear layer projects the output of the last decoder block to class predictions. Formally, the operations of the Transformer decoder can be expressed as follows:

$$\mathbf{Y}_{i} = \operatorname{Temp}_{i}\left(\left[\mathbf{X}_{N-M+i,1}, \mathbf{X}_{N-M+i,2}, \dots, \mathbf{X}_{N-M+i,T}\right]\right),$$
(1)

$$\tilde{\mathbf{q}}_{i} = \mathbf{q}_{i-1} + \mathrm{MHA}_{i} \left(\mathbf{q}_{i-1}, \mathbf{Y}_{i}, \mathbf{Y}_{i} \right), \tag{2}$$

$$\mathbf{q}_{i} = \tilde{\mathbf{q}}_{i} + \mathrm{MLP}_{i}\left(\tilde{\mathbf{q}}_{i}\right),\tag{3}$$

$$\mathbf{p} = \mathrm{FC}\left(\mathbf{q}_{M}\right),\tag{4}$$

where $\mathbf{X}_{n,t}$ denotes the frame features of the *t*-th frame extracted from the *n*-th layer of the CLIP backbone, \mathbf{Y}_i denotes the temporal modulated feature volume fed into the *i*-th layer of the Transformer decoder, \mathbf{q}_i is the progressively refined query token with \mathbf{q}_0 as learnable parameters and \mathbf{p} is the final prediction.

N, M denote the number of blocks in the backbone image encoder and the spatiotemporal decoder, respectively. MHA stands for multi-head attention, and the three arguments are the query, key, and value, respectively. Temp is the temporal modelling, which produces feature tokens modulated by more fine-grained temporal information, as is elaborated in the next section.

The network is optimized as a standard classification model by cross-entropy loss with ground-truth labels, except that the back-propagation stops at image features \mathbf{X} and no weight in the image encoder is updated.

3.2 Learning Temporal Cues from Spatial Features

While CLIP models generate powerful spatial features, they entirely lack temporal information. Despite the Transformer decoder being capable of weighted feature aggregation, which is a form of global temporal information, fine-grained and local temporal signals may also be valuable for video recognition. Hence, we introduce the following temporal modules to encode such information before features are fed into the Transformer decoder.

Temporal Convolution. Temporal depthwise convolutions are capable of capturing local feature variations along the temporal dimension, and in known to be efficient and effective [37,10]. Formally the feature encoded by this convolution is written as \mathbf{Y}_{conv} , and

$$\mathbf{Y}_{\text{conv}}\left(t,h,w,c\right) = \sum_{\Delta t \in \{-1,0,1\}} \mathbf{W}_{\text{conv}}\left(\Delta t,c\right) \mathbf{X}\left(t + \Delta t,h,w,c\right) + \mathbf{b}_{\text{conv}}\left(c\right).$$
 (5)

Temporal Positional Embeddings. We learn a set of T vectors of dimension C, denoted as $\mathbf{P} \in \mathbb{R}^{T \times C}$, to serve as temporal positional embedding. Image features are added with one of the vectors according to their temporal position t, or formally

$$\mathbf{Y}_{\text{pos}}\left(t,h,w,c\right) = \mathbf{P}\left(t,c\right).$$
(6)

While temporal convolutions may also capture temporal position information implicitly, positional embeddings are more explicit by making similar features at different time distinguishable. Positional embeddings are also more powerful for long-range temporal modelling, for which multiple convolutional blocks have to be stacked to achieve a large receptive field.

Temporal Cross Attention. Another interesting but often overlooked source of temporal information lies in the attention maps. As attention maps reflect feature correspondence, calculating attention maps between two frames naturally reveals object movement information. More specifically, we first construct attention maps between adjacent frames using the original query and key projections in CLIP:

$$\mathbf{A}_{\text{prev}}(t) = \text{Softmax}\left(\left(\mathbf{QX}(t)\right)^{T}\left(\mathbf{KX}(t-1)\right)\right),$$

$$\mathbf{A}_{\text{next}}(t) = \text{Softmax}\left(\left(\mathbf{QX}(t)\right)^{T}\left(\mathbf{KX}(t+1)\right)\right).$$
(7)

We omitted the attention heads for simplicity, and average across all heads in our implementation. Then we linearly project it into the feature dimension:

$$\mathbf{Y}_{\text{attn}}(t, h, w, c) = \sum_{h'=1}^{H} \sum_{w'=1}^{W} \mathbf{W}_{\text{prev}}(h - h', w - w', c) \,\mathbf{A}_{\text{prev}}(t, h', w') + \\ \mathbf{W}_{\text{next}}(h - h', w - w', c) \,\mathbf{A}_{\text{next}}(t, h', w') \,.$$
(8)

Experiments have shown that, despite the query, key, and input features all being learned from pure 2D image data, such attention maps still provide useful signals.

The final modulated features are obtained by blending the temporal features with the original spatial features in a residual manner, *i.e.* $\mathbf{Y} = \mathbf{X} + \mathbf{Y}_{conv} + \mathbf{Y}_{pos} + \mathbf{Y}_{attn}$.

3.3 Complexity Analysis

Inference The additional Transformer decoder introduces only a negligible amount of computational overhead given that only one query token is used. To show this, we consider ViT-B/16 as our image backbone, and write out the FLOPS for a Transformer block as follows:

$$FLOPS = 2qC^2 + 2kC^2 + 2qkC + 2\alpha qC^2 \tag{9}$$

Here, q, k, C, α stand for the number of query tokens, number of key (value) tokens, number of embedding dimensions, and MLP expansion factor. With this formula, we can roughly compare the FLOPS of an encoder block and decoder block (h, w, t is the feature size along the height, width, temporal dimensions, and we adopt a common choices $\alpha = 4$, h = w = 14, C = 768 for estimation):

$$\frac{\text{FLOPS}_{\text{dec}}}{\text{FLOPS}_{\text{enc}}} \approx \frac{2hwtC^2}{t(12hwC^2 + 2h^2w^2C)} \approx \frac{1}{6}$$
(10)

From this, we can see that a decoder block is much more lightweight compared to an encoder block. Even with a full configuration (one decoder block on every encoder output, no channel reduction and all temporal modules enabled), the FLOPS increase is within 20% of the backbone.

Training As we use a fixed backbone and a non-intrusive Transformer decoder head (i.e., our inserted module does not change the *input* of any backbone layer), we can completely avoid back-propagation through the backbone. This vastly reduces both the memory consumption and the time per training iteration.

4 Experiments

We benchmark our method on 2 datasets: Kinetics-400 and Something-Somethingv2. Extra implementation details are provided in the appendix.

4.1 Main Results

In this section we provide a comparison with important baselines from recent work.

Table 1: Comparison with state-of-the-arts on Kinetics-400. We cite a series of models within similar range of accuracy as ours and compare the FLOPS. Frame counts are reported as frames per view \times number of views.

Method	Pretraining	Acc. $(\%)$	#Frames	GFLOPS
Uniformer-B [24]	ImageNet-1k	82.9	32×4	1,036
Swin-B [27]	ImageNet-21k	82.7	32×12	3,384
irCSN-152 [37]	IG-65M	82.6	32×30	2,901
MViT-S [43]	ImageNet-21k	82.6	16×10	710
Omnivore-B [15]	IN1k + SUN	83.3	32×12	3,384
ViViT-L FE [1]	$_{ m JFT}$	83.5	32×3	$11,\!940$
TokenLearner 8at18 $(L/16)$ [33]	$_{ m JFT}$	83.2	32×6	$6,\!630$
MViT-L [43]	MaskFeat, K600	85.1	16×10	3,770
MTV-L [46]	$_{ m JFT}$	84.3	32×12	$18,\!050$
		82.9	8×3	444
EVL ViT-B/16 (Ours)	CLIP	83.6	16×3	888
		84.2	32×3	1,777
		86.3	8×3	2,022
EVL ViT-L/14 (Ours)	CLID	87.0	16×3	4,044
,	OLIP	87.3	32×3	8,088
EVL ViT-L/14 (336px, ours)		87.7	32×3	$18,\!196$

Comparison with State-of-the-art. Comparisons with recent state-of-the-art video recognition models are provided in Table 1. While we aim to build a fast transfer learning pipeline, we find our models achieve competitive accuracy among regular video recognition methods. The models listed in Table 1 achieve similar accuracy as ours but require substantially more computation than our method.

Comparison with CLIP-based Methods. To the best of our knowledge, there are two previous studies that utilize CLIP models for video recognition. As shown in Table 2, we achieve higher accuracy with fewer frames and a smaller number of new parameters, showing a more efficient use of CLIP.

Training Time and Reduced Memory. One of the major advantages of our efficient transfer pipeline is the vastly reduced training time. We cite the training time reported in several previous studies in Table 4 for comparison.¹ In this case,

¹ Training time of Uniformer-B is estimated by halving the value for Kinetics-600 provided in their GitHub repo. Training time of TimeSformer is from our own

Table 2: Comparison with CLIP-based methods on Kinetics-400. All models use ViT-B/16 as backbone. As the paper [22] is vague about the details, we estimate their new parameters to be 3 Transformer blocks with feature size 512 and MLP expansion factor 4. For ActionCLIP [41], we do not count parameters in the text branch.

Method	New Params (M)	$\#{\rm Frames}{\times}\#{\rm Views}$	Acc. (%)
Efficient-Prompting [22] (A5)	9.43*	16×5	76.9
ActionCLIP [41]	105.15	8×1 16×1 32×1 16×3 32×3	81.1 81.7 82.3 82.6 83.8
EVL ViT-B/16 (Ours, 1 Layer) EVL ViT-B/16 (Ours, 4 Layers) EVL ViT-B/16 (Ours, 4 Layers)	$7.41 \\28.70 \\28.78$	8 imes 3 8 imes 3 32 imes 3	81.1 82.9 84.2

Table 3: Inference latency and throughput measured on actual hardware. Both models achieve 82.9% accuracy on Kinetics-400. Results are obtained using V100-32G with PyTorch-builtin mixed precision. Latency is measured using a batch size of 1 and throughput is measured using the largest possible batch size before running out of memory.

Model (# frames)	Acc. (%)	GFLOPS	Latency (ms)	Throughput (V/s)
Uniformer-B (32) [24] EVL ViT-B/16 (Ours, 8)	$\begin{array}{c} 82.9\\ 82.9\end{array}$	$\begin{array}{c} 1036 \ (1.00\times) \\ \textbf{454} \ (0.44\times) \end{array}$	314.58 (1.00×) 102.88 (0.33×)	3.42 (1.00×) 25.53 (7.47×)

powerful pretraining leads to a roughly $10 \times$ training time reduction, and our efficient transfer learning scheme leads to a further reduction of about $8 \times$. We also compare training times in an idealized setting in Table 5: We report single step time (forward + backward + update) using fake data on a single GPU. This bypasses the data loading and distributed communication overhead, which are confounding factors that may be unoptimized and difficult to control.

Inference Latency and Throughput. Despite our method not being specially optimized for inference speed, we show an important advantage of utilizing large-scale pretrained models. Training on small datasets requires injecting hand-crafted inductive biases, which are not necessarily friendly to modern accelerators. On the contrary, ViT models consist almost entirely of standard linear algebra operations. The simplicity of ViT typically enables a higher utilization of hardware resources.

reproduction, which we find to be a few times smaller than the reported number in their paper (reported value is around 400 hours). Training time of ActionCLIP is estimated by doubling the value for 8-frame variant reported in their paper.

TT 1 1 4	m • •		•
Table 4:	Training	time	comparison.
100010 10		· · · · · · ·	0011100110011

Method (#Frames per View)	Acc. (#Views)	Pretraining	Training GPU Hours
Uniformer-B [24] (32)	82.9 (4)	ImageNet-1k	$5000 \times V100$
TimeSformer $[3]$ (8)	82.0(3)	CLIP	$100 \times V100$
ActionCLIP $[41]$ (16)	82.6(3)	CLIP	$480 \times \text{RTX}3090$
EVL ViT-B/16 (8)	82.9 (3)	CLIP	$60 \times V100$

Table 5: Idealized training step time. 4 decoder layers are used. All data are measured on a single V100-16G GPU. The step time is measured with 64 training samples.

Backbone	Head	Max Batch Size	Step Time (s)
CLIP (Frozen)	global average pool	inf.	0.57
CLIP (Open)	global average pool	8	3.39
CLIP (Frozen)	\mathbf{EVL}	64	1.03
CLIP (Open)	\mathbf{EVL}	8	4.41

As shown in Table 3, the latency and throughput are even better than the theoretical FLOPS improvement.

4.2 Ablation Studies

We provide detailed ablation studies to clarify the effects of each part of our design. Unless otherwise specified, results are obtained using ViT-B/16 backbone, 8 input frames and 3 testing views on Kinetics-400.

Intermediate Features. We vary the number of features and Transformer decoder layers and present the results in Table 6a and Table 6b. Utilizing multiple decoder blocks improves the accuracy by 1.0%. Feeding each decoder block with multi-layer intermediate features further improves by 0.8%. Another observation is that features in deeper layers provide more effective features for video recognition.

Spatiotemporal Features. We find a crucial design to achieve high transfer performance is to use high-resolution, unpooled feature maps. The results are shown in Table 6c, from which we can see that summarizing along either the temporal or spatial dimension leads to a significant drop in accuracy. We conjecture that this shows the importance of task-specific re-attention, e.g., for human action recognition datasets like Kinetics-400, features relating to the human body are very important, which could be different in the pretraining stage.

Pretraining Quality. One major factor driving the paradigm shift from finetuned to frozen backbone is the improvement in quality of pretrained models. We show that our method outperforms previous methods that fully finetune

number of Transformer decoder blocks. (b) Varying number of feature maps. (c) Varying feature resolution. (a) (b) (c) <u>Depth Acc. (%)</u> Feature Layers Acc. (%) Feature Shape Reduction Acc. (%) <u>1 81 1</u> [-4 -3 -2 -1] 82 9 Temporal only Token 70.8

Table 6: Effects of multi-layer high-resolution feature maps. (a) Varying

Depth	ncc. (70)	reature Layers	mee. (70)	reature shape	ittuttion	ncc. (70)
1	81.1	$\overline{[-4, -3, -2, -1]}$	82.9	Temporal only	Token	79.8
2	82.1	[-2, -2, -1, -1]	82.7	Temporal only	Avg	75.8
3	82.6	[-1, -1, -1, -1]	82.1	Spatial only	Avg	80.1
4	82.9	[-2, -1, -2, -1]	82.4	Spatiotemporal	-	82.9
5	83.0	[-7, -5, -3, -1]	82.0			

the backbone weights given the high quality CLIP backbones in Table 7. All models in the table use the same backbone architecture. While on ImageNet-21k pretrained backbones our method lags behind full-finetuning, on CLIP backbones our method outperforms the competitive full-finetuning baselines.

We also find that, despite being designed for a frozen backbone, our model architecture with a finetuned backbone turns out to be a strong fullfinetuning baseline. However, the tendency of higher training efficiency of frozen backbones given high-quality pretrained models remains the same, as shown in Fig. 3. Full-finetuning with our model architecture yields similar efficiency curve on ViT-B/16, but with the larger ViT-L/14, the gap of the training time to reach the same accuracy becomes clear. We point out that even ViT-L/14 is a relatively small pretrained model by modern standards, with about 300M parameters (for comparison, GPT-3 [4] for natural language processing has 175B parameters, and ViT-G [49] for computer vision has 1.8B parameters). We believe freezing the backbone may potentially



Fig. 3: Training time vs. accuracy with frozen or finetuned backbone. Numbers in the marker are numbers of frames per view. Frozen backbone is more efficient when pretraining quality is higher.

bring further benefits if even larger pretrained models are released in the future.

4.3 Analysis of Temporal Information

An interesting property of our method is to provide a decomposed approach for video recognition: the spatial information is encoded almost entirely in the fixed,

Table 7: **Results of different pretrained image features**. A ViT-B/16 backbone and 8 frames are used unless otherwise specified. We compare with TimeSformer [3] and ActionCLIP [41]. Both of them conduct extensive experiments to determine competitive settings for end-to-end training on video datasets.

Model	Pretraining	Frozen Backbone?	K-400 Acc. (%)
TimeSformer [3] - SOnly	ImageNet-21k	×	76.9
TimeSformer [3] - JointST	ImageNet-21k	X	77.4
TimeSformer [3] - DividedST	ImageNet-21k	×	78.0
EVL (ours)	ImageNet-21k	1	75.4
TimeSformer [3] - DividedST	CLIP	X	82.0
EVL (ours)	CLIP	\checkmark	82.9
ActionCLIP [41] (16 frames)	CLIP	X	82.6
EVL (ours, 16 frames)	CLIP	1	83.3
ActionCLIP [41] (32 frames)	CLIP	X	83.8
EVL (ours, 32 frames)	CLIP	<i>✓</i>	84.2

high quality CLIP backbone, while the temporal information is encoded only in the Transformer decoder head. As shown in Table 8, temporal modelling exhibits vastly different behaviors on the two datasets: On Kinetics-400,temporal modules bring accuracy gains of less than 0.5%, while on Something-Something-v2, adding the temporal module yields a dramatic +13.8% accuracy gain. This shows a clear difference between temporal information required for the two benchmarks. For Kinetics-400, temporal information is primarily captured in the form of global weighted feature aggregation, as shown in Table 6. For Something-Somethingv2, local temporal features (e.g., object motion, feature variations) are also an important source of signals to achieve strong results.

Something-Something-v2 also tend to benefit from deep decoders more than Kinetics-400. As shown in Table 8b, Something-Something-v2 benefit from using all 12 decoder blocks, while for Kinetics-400 only around 4 blocks are required (see Table 6a).

Finally we provide our main results on Something-Something-v2 dataset in Table 9. While Something-Something-v2 is a motion-heavy dataset, our lightweight temporal learning module still learns meaningful motion information and reaches mainstream performance (for comparison, a linear probe of CLIP ViT-B/16 achieves only around 20% accuracy). We are also the first CLIP-based method to report results on Something-Something-v2, and we hope this is useful for future reference.

4.4 CLIP-based Models Learn Complementary Knowledge

Another finding is that knowledge learned by our CLIP-based model is highly complementary to that of regular supervised learning. To show this, we consider an ensemble of our model with supervised models and observe the performance

Table 8: Effects of temporal information for video recognition. (a) Local temporal information for both datasets. T-Conv: temporal convolution. T-PE: temporal positional embedding. T-CA: temporal cross attention. (b) Something-Something-v2 needs deeper decoder blocks.

				(b)		
T-Conv	T-PE	T-CA	K-400 Acc. (%)	SSv2 Acc. (%)	Depth	SSv2 Acc. (%)
X	Х	Х	82.5	47.2	4	58.6
1	X	X	82.9	57.1	6	60.1
X	1	X	82.5	58.5	8	60.2
X	X	1	82.6	59.5	10	60.5
1	1	X	82.9	59.4	12	61.0
1	X	1	82.7	60.0		
X	1	1	82.7	60.7		
~	1	1	82.9	61.0		

Table 9: Main results on Something-Something-v2. Ens experiments combine EVL with Uniformer-B (32) pretrained on Kinetics-600.

Method	SSv2 Acc. (%)	#Frames	GFLOPS
EVL ViT-B/16	61.0	8×3	512
EVL ViT-B/16	61.7	16×3	1,023
EVL ViT-B/16	62.4	32×3	$2,\!047$
EVL ViT-L/14	65.1	8×3	2,411
EVL ViT-L/14	66.7	32×3	$9,\!641$
EVL ViT-L/14 (336px) $$	68.0	32×3	$24,\!259$
EVL ViT-B/16 Ens	72.1	$32 \times 3 + 32 \times 3$	2,824

gain. Ensemble is done by weighted averaging the video-level prediction scores and the average weight $\alpha \in [0, 1]$ is searched with a coarse granularity of 0.1 on the validation set. As shown in Table 10 and Table 11, On both Kinetics-400 and Something-Something-v2, we consistently observe more performance gain if CLIP-based models are in the ensemble.

The implications of these ensemble experiments are two-fold. First, they show that, practically, our CLIP-based models can be used in a two-stream fashion [35].Compared to the optical-flow-based second stream in [35], a CLIP-based second stream avoids the expensive optical-flow calculation and is much faster to train. Second, the results suggest that there remains knowledge in the dataset that is not captured by our CLIP-based learning paradigm. This shows the potential of CLIP-based models to further improve once more knowledge from the datasets can be utilized.

Table 10: **Ensemble results of different combinations**. We combine different models with similar accuracy with the same model and measure the accuracy gain.

Model 1	Acc. 1	Model 2	Acc. 2	Model 1 + 2 Acc. (Δ)
Uniformer-B [24] (16)	82.0	Uniformer-B [24] (32) Swin-B [27] EVL ViT-B/16 (8)	82.9 82.7 82.9	83.6 (+1.6) 83.7 (+1.7) 84.5 (+2.5)
Swin-B [27]	82.7	Uniformer-B [24] (32) EVL ViT-B/16 (8)	82.9 82.9	84.7 (+2.0) 85.0 (+2.3)
Uniformer-B [24] (32)	82.9	Swin-B [27] EVL ViT-B/16 (8)	82.7 82.9	84.7 (+1.8) 85.2 (+2.3)

Table 11: Ensemble results on Something-Something-v2. Although EVL (32) has much lower accuracy, it still boosts the performance of a Uniformer-B model. In contrast, a TimeSformer model with slightly higher accuracy brings negligible gains.

Model 1	Acc. 1	Model 2	Acc. 2	Model 1 + 2 Acc. (Δ)
Uniformer-B [24] (32)	71.2	TimeSformer-L [3] EVL ViT-B (32)	$62.4 \\ 62.4$	71.4 (+0.2) 72.1 (+0.9)

5 Conclusion

We present a new form of pipeline for video action recognition: learning an efficient transfer learning head on top of fixed transferable image features. By freezing the image backbone, the training time is vastly reduced. Moreover, the accuracy loss due to the frozen backbone can be largely compensated by leveraging multi-layer high-resolution intermediate feature maps from the backbone. Thus, our method effectively leverage powerful image features for video recognition, while avoiding the heavy or prohibitive full-finetuning of very large image models. We further show that transferable image features learned in an open-world setting harbor knowledge that is highly complementary to that of labeled datasets, which may inspire more efficient ways to build state-of-the-art video models. We believe our observations have the potential to make video recognition accessible to a broader community, and push video models to a new state-of-the-art in a more efficient manner.

Acknowledgements. This work is supported in part by Centre for Perceptual and Interactive Intelligence Limited, in part by the General Research Fund through the Research Grants Council of Hong Kong under Grants (Nos. 14204021, 14207319), in part by CUHK Strategic Fund. This work is partially supported by the Shanghai Committee of Science and Technology (Grant No. 21DZ1100100).

15

References

- Arnab, A., Dehghani, M., Heigold, G., Sun, C., Lučić, M., Schmid, C.: Vivit: A video vision transformer. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 6836–6846 (2021)
- Bao, H., Dong, L., Wei, F.: Beit: Bert pre-training of image transformers. arXiv preprint arXiv:2106.08254 (2021)
- Bertasius, G., Wang, H., Torresani, L.: Is space-time attention all you need for video understanding? In: Proceedings of the International Conference on Machine Learning (ICML) (2021)
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J.D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al.: Language models are few-shot learners. Advances in neural information processing systems 33, 1877–1901 (2020)
- Carreira, J., Zisserman, A.: Quo vadis, action recognition? a new model and the kinetics dataset. In: proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 6299–6308 (2017)
- Cheng, X., Lin, H., Wu, X., Yang, F., Shen, D.: Improving video-text retrieval by multi-stream corpus alignment and dual softmax loss. arXiv preprint arXiv:2109.04290 (2021)
- Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. In: NAACL-HLT (2019)
- Fan, H., Xiong, B., Mangalam, K., Li, Y., Yan, Z., Malik, J., Feichtenhofer, C.: Multiscale vision transformers. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 6824–6835 (2021)
- Fang, H., Xiong, P., Xu, L., Chen, Y.: Clip2video: Mastering video-text retrieval via image clip. arXiv preprint arXiv:2106.11097 (2021)
- Feichtenhofer, C.: X3d: Expanding architectures for efficient video recognition. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 203–213 (2020)
- Feichtenhofer, C., Fan, H., Malik, J., He, K.: Slowfast networks for video recognition. In: Proceedings of the IEEE/CVF international conference on computer vision. pp. 6202–6211 (2019)
- Feichtenhofer, C., Fan, H., Xiong, B., Girshick, R., He, K.: A large-scale study on unsupervised spatiotemporal representation learning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 3299–3309 (2021)
- Gao, P., Geng, S., Zhang, R., Ma, T., Fang, R., Zhang, Y., Li, H., Qiao, Y.: Clip-adapter: Better vision-language models with feature adapters. arXiv preprint arXiv:2110.04544 (2021)
- Gao, Z., Liu, J., Chen, S., Chang, D., Zhang, H., Yuan, J.: Clip2tv: An empirical study on transformer-based methods for video-text retrieval. arXiv preprint arXiv:2111.05610 (2021)
- Girdhar, R., Singh, M., Ravi, N., van der Maaten, L., Joulin, A., Misra, I.: Omnivore: A Single Model for Many Visual Modalities. arXiv preprint arXiv:2201.08377 (2022)
- Guo, D., Rush, A.M., Kim, Y.: Parameter-efficient transfer learning with diff pruning. arXiv preprint arXiv:2012.07463 (2020)
- 17. Hara, K., Kataoka, H., Satoh, Y.: Learning spatio-temporal features with 3d residual networks for action recognition. In: Proceedings of the IEEE International Conference on Computer Vision Workshops. pp. 3154–3160 (2017)

- 16 Z. Lin et al.
- He, K., Chen, X., Xie, S., Li, Y., Dollár, P., Girshick, R.: Masked autoencoders are scalable vision learners. arXiv preprint arXiv:2111.06377 (2021)
- Houlsby, N., Giurgiu, A., Jastrzebski, S., Morrone, B., De Laroussilhe, Q., Gesmundo, A., Attariyan, M., Gelly, S.: Parameter-efficient transfer learning for nlp. In: International Conference on Machine Learning. pp. 2790–2799. PMLR (2019)
- Hu, E.J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., Chen, W.: Lora: Low-rank adaptation of large language models. arXiv preprint arXiv:2106.09685 (2021)
- Jia, C., Yang, Y., Xia, Y., Chen, Y.T., Parekh, Z., Pham, H., Le, Q.V., Sung, Y., Li, Z., Duerig, T.: Scaling up visual and vision-language representation learning with noisy text supervision. arXiv preprint arXiv:2102.05918 (2021)
- Ju, C., Han, T., Zheng, K., Zhang, Y., Xie, W.: Prompting visual-language models for efficient video understanding. arXiv preprint arXiv:2112.04478 (2021)
- 23. Lester, B., Al-Rfou, R., Constant, N.: The power of scale for parameter-efficient prompt tuning. arXiv preprint arXiv:2104.08691 (2021)
- 24. Li, K., Wang, Y., Gao, P., Song, G., Liu, Y., Li, H., Qiao, Y.: Uniformer: Unified transformer for efficient spatiotemporal representation learning. In: ICLR (2022)
- Li, X.L., Liang, P.: Prefix-tuning: Optimizing continuous prompts for generation. arXiv preprint arXiv:2101.00190 (2021)
- Li, Y., Wu, C.Y., Fan, H., Mangalam, K., Xiong, B., Malik, J., Feichtenhofer, C.: Improved multiscale vision transformers for classification and detection. arXiv preprint arXiv:2112.01526 (2021)
- Liu, Z., Ning, J., Cao, Y., Wei, Y., Zhang, Z., Lin, S., Hu, H.: Video swin transformer. arXiv preprint arXiv:2106.13230 (2021)
- Patrick, M., Campbell, D., Asano, Y.M., Metze, I.M.F., Feichtenhofer, C., Vedaldi, A., Henriques, J.F.: Keeping your eye on the ball: Trajectory attention in video transformers. In: Advances in Neural Information Processing Systems (NeurIPS) (2021)
- Pfeiffer, J., Kamath, A., Rücklé, A., Cho, K., Gurevych, I.: Adapterfusion: Nondestructive task composition for transfer learning. arXiv preprint arXiv:2005.00247 (2020)
- Qiu, Z., Yao, T., Mei, T.: Learning spatio-temporal representation with pseudo-3d residual networks. In: proceedings of the IEEE International Conference on Computer Vision. pp. 5533–5541 (2017)
- Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. arXiv preprint arXiv:2103.00020 (2021)
- Riquelme, C., Puigcerver, J., Mustafa, B., Neumann, M., Jenatton, R., Pinto, A.S., Keysers, D., Houlsby, N.: Scaling vision with sparse mixture of experts. arXiv preprint arXiv:2106.05974 (2021)
- Ryoo, M.S., Piergiovanni, A., Arnab, A., Dehghani, M., Angelova, A.: Tokenlearner: What can 8 learned tokens do for images and videos? arXiv preprint arXiv:2106.11297 (2021)
- Shridhar, M., Manuelli, L., Fox, D.: Cliport: What and where pathways for robotic manipulation. In: Proceedings of the 5th Conference on Robot Learning (CoRL) (2021)
- Simonyan, K., Zisserman, A.: Two-stream convolutional networks for action recognition in videos. arXiv preprint arXiv:1406.2199 (2014)
- Tran, D., Bourdev, L., Fergus, R., Torresani, L., Paluri, M.: Learning spatiotemporal features with 3d convolutional networks. In: Proceedings of the IEEE international conference on computer vision. pp. 4489–4497 (2015)

17

- Tran, D., Wang, H., Torresani, L., Feiszli, M.: Video classification with channelseparated convolutional networks. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 5552–5561 (2019)
- Tran, D., Wang, H., Torresani, L., Ray, J., LeCun, Y., Paluri, M.: A closer look at spatiotemporal convolutions for action recognition. In: Proceedings of the IEEE conference on Computer Vision and Pattern Recognition. pp. 6450–6459 (2018)
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. In: Advances in neural information processing systems. pp. 5998–6008 (2017)
- 40. Wang, L., Xiong, Y., Wang, Z., Qiao, Y., Lin, D., Tang, X., Van Gool, L.: Temporal segment networks: Towards good practices for deep action recognition. In: European conference on computer vision. pp. 20–36. Springer (2016)
- Wang, M., Xing, J., Liu, Y.: Actionclip: A new paradigm for video action recognition. arXiv preprint arXiv:2109.08472 (2021)
- Wang, X., Girshick, R., Gupta, A., He, K.: Non-local neural networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 7794–7803 (2018)
- Wei, C., Fan, H., Xie, S., Wu, C.Y., Yuille, A., Feichtenhofer, C.: Masked feature prediction for self-supervised visual pre-training. arXiv preprint arXiv:2112.09133 (2021)
- 44. Xie, S., Sun, C., Huang, J., Tu, Z., Murphy, K.: Rethinking spatiotemporal feature learning: Speed-accuracy trade-offs in video classification. In: Proceedings of the European conference on computer vision (ECCV). pp. 305–321 (2018)
- 45. Xie, Z., Zhang, Z., Cao, Y., Lin, Y., Bao, J., Yao, Z., Dai, Q., Hu, H.: Simmim: A simple framework for masked image modeling. arXiv preprint arXiv:2111.09886 (2021)
- 46. Yan, S., Xiong, X., Arnab, A., Lu, Z., Zhang, M., Sun, C., Schmid, C.: Multiview transformers for video recognition. arXiv preprint arXiv:2201.04288 (2022)
- 47. Yuan, L., Chen, D., Chen, Y.L., Codella, N., Dai, X., Gao, J., Hu, H., Huang, X., Li, B., Li, C., et al.: Florence: A new foundation model for computer vision. arXiv preprint arXiv:2111.11432 (2021)
- Zhai, X., Kolesnikov, A., Houlsby, N., Beyer, L.: Scaling vision transformers. arXiv preprint arXiv:2106.04560 (2021)
- Zhai, X., Kolesnikov, A., Houlsby, N., Beyer, L.: Scaling vision transformers. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12104–12113 (2022)
- Zhang, B., Yu, J., Fifty, C., Han, W., Dai, A.M., Pang, R., Sha, F.: Co-training transformer with videos and images improves action recognition. arXiv preprint arXiv:2112.07175 (2021)
- Zhang, R., Fang, R., Gao, P., Zhang, W., Li, K., Dai, J., Qiao, Y., Li, H.: Tipadapter: Training-free clip-adapter for better vision-language modeling. arXiv preprint arXiv:2111.03930 (2021)
- 52. Zhang, R., Guo, Z., Zhang, W., Li, K., Miao, X., Cui, B., Qiao, Y., Gao, P., Li, H.: Pointclip: Point cloud understanding by clip. arXiv preprint arXiv:2112.02413 (2021)
- Zhou, B., Andonian, A., Oliva, A., Torralba, A.: Temporal relational reasoning in videos. In: Proceedings of the European conference on computer vision (ECCV). pp. 803–818 (2018)
- Zhou, K., Yang, J., Loy, C.C., Liu, Z.: Learning to prompt for vision-language models. arXiv preprint arXiv:2109.01134 (2021)