

Implicit Temporal Modeling with Learnable Alignment for Video Recognition

Shuyuan Tu^{1,2} Qi Dai³ Zuxuan Wu^{1,2} * Zhi-Qi Cheng⁴ Han Hu³ Yu-Gang Jiang^{1,2}

¹Shanghai Key Lab of Intell. Info. Processing, School of CS, Fudan University

²Shanghai Collaborative Innovation Center of Intelligent Visual Computing

³Microsoft Research Asia ⁴Carnegie Mellon University

Abstract

Contrastive language-image pretraining (CLIP) has demonstrated remarkable success in various image tasks. However, how to extend CLIP with effective temporal modeling is still an open and crucial problem. Existing factorized or joint spatial-temporal modeling trades off between the efficiency and performance. While modeling temporal information within straight through tube is widely adopted in literature, we find that simple frame alignment already provides enough essence without temporal attention. To this end, in this paper, we proposed a novel Implicit Learnable Alignment (ILA) method, which minimizes the temporal modeling effort while achieving incredibly high performance. Specifically, for a frame pair, an interactive point is predicted in each frame, serving as a mutual information rich region. By enhancing the features around the interactive point, two frames are implicitly aligned. The aligned features are then pooled into a single token, which is leveraged in the subsequent spatial self-attention. Our method allows eliminating the costly or insufficient temporal self-attention in video. Extensive experiments on benchmarks demonstrate the superiority and generality of our module. Particularly, the proposed ILA achieves a top-1 accuracy of 88.7% on Kinetics-400 with much fewer FLOPs compared with Swin-L and ViViT-H. Code is released at <https://github.com/Francis-Rings/ILA>.

1. Introduction

Video recognition is rated as one of the most fundamental components of video understanding. Numerous downstream tasks heavily rely on the basic recognition model, e.g., action localization [6, 12, 45, 46, 48], detection [7, 19, 24, 30, 73], and video object tracking [16, 57, 75]. Due to the great potential of video technologies, it has been an active research direction over the past few years. Various approaches have been proposed, including convolution-based

*Corresponding author

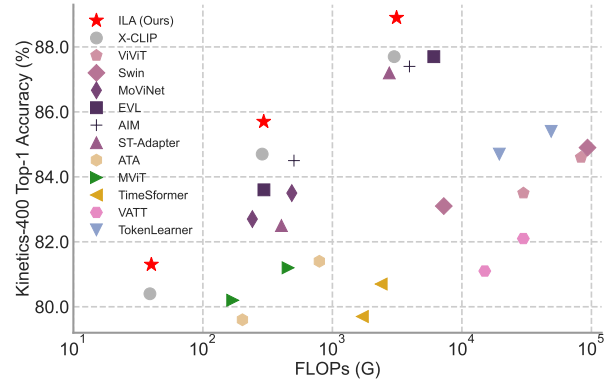


Figure 1. Top-1 accuracy comparison with state-of-the-art methods on Kinetics-400 [28] under different FLOPs. ILA achieves competitive results. Best viewed in color.

methods [49, 55, 53, 10, 60, 18, 17, 76] and transformer-based methods [3, 5, 20, 15, 31, 36, 50, 61]. Recently, Contrastive Language-Image Pretraining (CLIP) [41] has demonstrated strong performance in video domain. Studies [56, 27, 38, 35, 39, 64, 74] attempt to transfer the powerful CLIP model to video tasks, which promote the recognition performance to a brand-new level, showing its general representation ability.

Generally, existing methods devise various temporal modeling schemes to explore the potential of CLIP, including the factorized [64] or frame-level [38, 27] temporal attention, and temporal cross attention [35]. All these tailored methods aim at designing lightweight temporal modules to reuse the CLIP model. Though considerable improvements are achieved, such temporal modeling approaches still depend on the complex self-attention, which we argue is not necessary in CLIP-based framework.

In this paper, we rethink the role of temporal modeling in general CLIP-based video recognition framework. Unlike existing approaches rely on temporal attention, we hypothesize that important motion and action clues can be derived when performing alignment of pairwise frames. As a result, the costly [36, 5] or insufficient [38, 27, 35] temporal at-

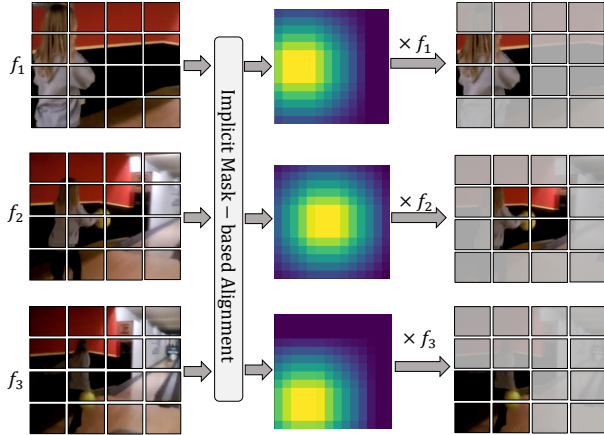


Figure 2. The proposed ILA employs an implicit and coarse mask to align the features, which focus on the active interaction region. We hypothesize the important motion and action clues can be derived from aligned features.

tentions can be avoided, without harming the performance. While explicit patch alignment is time consuming with low efficiency, we prioritize only an implicit and coarse alignment, aiming at involving the vital temporal signals.

In light of this, we present a novel Implicit Learnable Alignment (ILA) method for efficient video recognition. More specifically, ILA employs learnable masks to align features of two adjacent frames. The alignment is achieved with the help of an interactive point that is predicted using a convolution module conditioned on a frame pair. Based on the point, a corresponding region is generated indicating close interactions of adjacent frames. The mask is defined as the map of weights implying which region contains vital information. We then assign higher weights around the interactive point in the mask, while assigning lower weights to other positions, suppressing irrelevant signals among them. By leveraging the generated mask to weight the frame representations, coarsely aligned features are obtained, as shown in Figure 2. Note all above operations are performed in parallel among frame pairs to boost the speed. To efficiently and fully exploit the alignment, the aligned features are pooled into a single mutual information token. The token is subsequently concatenated with other frame tokens to perform the spatial self-attention, which implicitly models the temporal relations between frames. Our method is plugged into each spatial block of vision transformer and forms the Implicit Spatio-Temporal attention (IST) block, which allows temporal modeling without the use of the traditional temporal self-attention.

Our contributions can be summarized as follows: (1) We propose Implicit Learnable Alignment (ILA) for video recognition. Our implicit temporal modeling can be seamlessly plugged into existing vision transformer models. It utilizes the coarse alignment as the key temporal signals,

which enables superior temporal modeling at a low computational cost. (2) We show that such a simple frame alignment already encodes the essence of temporal relations, which allow eliminating the insufficient temporal self-attention. (3) Extensive qualitative and quantitative experiments demonstrate the effectiveness and efficiency of ILA. We achieve 88.7% on Kinetics-400 with low computation overhead. Our method builds a promising bridge for CLIP from image processing to video recognition.

2. Related Work

Visual-language representation learning has demonstrated remarkable success in various tasks [41, 25, 65]. By leveraging contrastive learning between language and image, a joint representation space is learned. Particularly, CLIP [41] has shown its strong power in open domain problems, and dozens of approaches are developed, including few-shot learning [21, 67], point cloud understanding [68, 43], video understanding [62, 56, 27], *etc.*

Recently, several studies extend the existing CLIP model to the video domain. X-CLIP [38] devises the frame-level temporal attention to avoid high computation. EVL [35] employs temporal convolution and cross-attention on top of the CLIP features. ST-Adapter [39] inserts the spatiotemporal adapter into each block, which consists of several 3D convolution layers. AIM [64] reuses the CLIP self-attention as the temporal ones via an additional adapter module. Nevertheless, the above methods explore lightweight adaptations of CLIP using insufficient temporal attention, *e.g.*, frame-level or local temporal attention. In our work, we attempt to perform temporal modeling with signals emerged from a simple alignment process, which involves the comprehensive temporal clues yet remains simplicity.

Video recognition is the key task in video understanding. In the convolution era, two-stream networks [49, 55, 71] and spatiotemporal CNNs [53, 23, 54, 60] are proposed. The former treats spatial representations and optical flow images as two independent modules, and the latter employs (separable) 3D convolution to extract spatiotemporal features. Recently, inspired by vision transformers [14, 52, 69, 50, 22], video transformers [5, 36, 15, 3, 42, 1] have shown promising results compared to CNN methods, due to their much larger receptive fields. TimeSformer [5] adopts factored space time attention as a trade-off between speed and accuracy. ViViT [3] investigates four types of temporal attention, and selected the global spatiotemporal attention as the default. Video Swin [36] uses local spatiotemporal attention to model the temporal information. However, these methods are either computationally intensive or insufficient in modeling the temporal interactions, resulting in high model cost or unsatisfactory performance. In contrast, our method explores how to model the complex temporal

information with minimal effort, demonstrating the redundancy in existing temporal attention models.

Temporal correspondences reflect the motions in video and can be used in several video understanding tasks [40, 26, 29, 42, 59, 32]. For example, in video super resolution, alignment-based methods [47, 9, 51, 58, 11, 34] have been proposed to keep frames coherent. PSRT-recurrent [47] points out that patch level alignment can reduce memory cost when computing optical flow. While in video recognition, the recent ATA [70] adopts Hungarian matching to align the patches between frames, and then performs temporal attention within aligned patches, followed by the de-alignment. However, the model is significantly encumbered with the slow serial alignment, followed by computationally expensive operations to calculate temporal attention. In contrast, our approach employs learnable masks to align frames in parallel with an aim to involve important motion and action clues, thus benefiting the video understanding. Therefore, the alignment in our method is implicit and coarse.

3. Method

In this section, we elaborate our proposed architecture in detail. First, we introduce the overview of our ILA in Section 3.1. Second, we depict the concrete implicit mask-based alignment in Section 3.2. Finally, we describe the loss functions of our dedicated framework.

3.1. Architecture Overview

The proposed ILA model consists of several Implicit Spatio-Temporal attention (IST) blocks. The model is built upon a pretrained image vision transformer (ViT) [14]. While previous methods [3, 5] mostly rely on the ImageNet initialized models, recent approaches [38, 35, 39, 64] have revealed the powerful representation ability of large-scale visual-language pretrained models [41, 65]. Our method follows the predecessors and is initialized from the CLIP model [41]. Given an input video clip $\mathbf{x} = [\mathbf{x}_1, \dots, \mathbf{x}_t, \dots, \mathbf{x}_T], \mathbf{x}_t \in \mathbb{R}^{H \times W \times 3}$, we decompose each frame into $\frac{H}{P} \times \frac{W}{P}$ non-overlapping patches $\{\mathbf{x}_{t,i}\}_{i=1}^{hw}$, where T, H, W are the number of frames, height and width, $h = \frac{H}{P}, w = \frac{W}{P}$. P is the patch size. The patches are linearly mapped to embedding vectors $\mathbf{z}_t^{(0)} = [\mathbf{z}_{t,1}^{(0)}, \dots, \mathbf{z}_{t,i}^{(0)}, \dots, \mathbf{z}_{t,hw}^{(0)}], \mathbf{z}_{t,i}^{(0)} \in \mathbb{R}^d$:

$$\mathbf{z}_{t,i}^{(0)} = \mathbf{E}\mathbf{x}_{t,i} + \mathbf{e}_{t,i}^{pos}, \quad (1)$$

where $\mathbf{E} \in \mathbb{R}^{d \times 3P^2}$ is the projection matrix, and $\mathbf{e}_{t,i}^{pos}$ is the spatial positional embedding. We also add a classification token $\mathbf{z}_{t,cls}^{(0)}$ for each frame.

The structure of the IST block is illustrated in Figure 3. At each IST block ℓ , we align the semantic features of

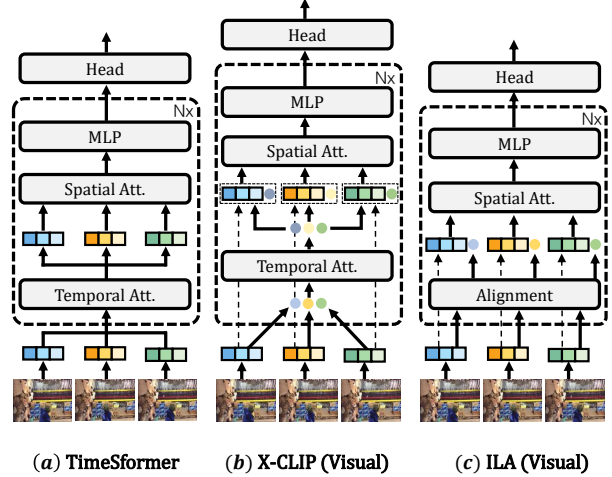


Figure 3. The structures of three different models. (a) The divided spatiotemporal attention in TimeSformer [5]. (b) The frame-level temporal attention in X-CLIP [38]. (c) The alignment-based temporal modeling in our ILA.

each consecutive frame pair $(\mathbf{z}_t^{(\ell-1)}, \mathbf{z}_{t-1}^{(\ell-1)})$ via finding an interactive position (as will be introduced in Section 3.2) per frame, which serves as a mutual information (MI) rich region. By simply weighting the feature map with higher weights surrounding the interactive position, the aligned features $\mathbf{a}_t^{(\ell)}, \mathbf{a}_{t-1}^{(\ell)}$ are obtained:

$$\mathbf{a}_t^{(\ell)}, \mathbf{a}_{t-1}^{(\ell)} = \text{Align}(\mathbf{z}_t^{(\ell-1)}, \mathbf{z}_{t-1}^{(\ell-1)}). \quad (2)$$

The aligned features are subsequently mean-pooled into a single mutual information token $\hat{\mathbf{z}}_{t,mult}^{(\ell)}$, which is further concatenated with corresponding frame to perform the spatial Multi-head Self Attention (MSA):

$$\hat{\mathbf{z}}_{t,mult}^{(\ell)} = \text{Avg}(\mathbf{a}_t^{(\ell)}), \quad (3a)$$

$$[\tilde{\mathbf{z}}_t^{(\ell)}, \tilde{\mathbf{z}}_{t,mult}^{(\ell)}] = \text{MSA}(\text{LN}([\mathbf{z}_t^{(\ell-1)}, \hat{\mathbf{z}}_{t,mult}^{(\ell)}])) + [\mathbf{z}_t^{(\ell-1)}, \hat{\mathbf{z}}_{t,mult}^{(\ell)}], \quad (3b)$$

where $\text{LN}(\cdot)$ indicates layer normalization [4]. $\tilde{\mathbf{z}}_{t,mult}^{(\ell)}$ is then dropped before feeding to the MLP, and the output of block ℓ is formulated as:

$$\mathbf{z}_t^{(\ell)} = \text{MLP}(\text{LN}(\tilde{\mathbf{z}}_t^{(\ell)})) + \tilde{\mathbf{z}}_t^{(\ell)}. \quad (4)$$

Unlike common supervised frameworks that use one-hot labels as the target, to fully leverage the pretrained visual-language model, we follow [38] to optimize the similarity loss supervised by textual information of categories. Formally, the text representation \mathbf{c} is computed by inputting the category name to the text encoder $f_t(\cdot)$. Then a video-specific prompt is obtained by querying \mathbf{c} among video representation $\{\mathbf{z}_t^{(L)}\}_{t=1}^T$ (L is the number of IST blocks), which is further used to enhance \mathbf{c} . Finally, the model maximizes the cosine similarity between the video and text representations if they are matched, otherwise minimizes it.

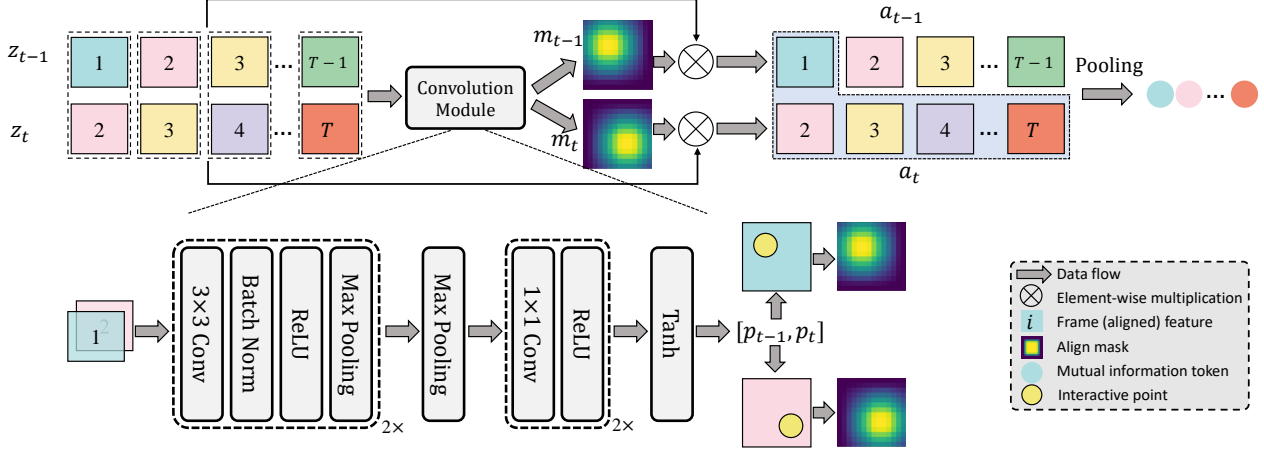


Figure 4. Details of the proposed alignment method. For each adjacent frame pair, a convolution module is leveraged to predict one interactive point per frame, which refers to region with close interactions between frames. A mask is generated by assigning higher weights around the interactive point, while assigning lower weights to other positions. The mask is then adopted to weight the frame features, obtaining aligned features. Finally, the aligned features are pooled into a single mutual information token. Best viewed in color.

3.2. Implicit Mask-based Alignment

The IST block employs an implicit mask-based alignment component to align the semantic features between two frames. A previous study [70] had explored patch-level alignment through Hungarian matching [8], which however suffered from limited performance and low efficiency. On one hand, the explicit patch alignment focuses on patch coherence across frames, which can eliminate possible beneficial temporal interactions. On the other hand, such alignment must be operated frame by frame with cubic time complexity, incurring significant computational overhead. In contrast, our implicit alignment attempts to enhance favorable mutual information and in turn suppress irrelevant information with learned masks. As such, the key temporal clues are preserved while allowing flexible and efficient computation.

Figure 4 illustrates the details of our alignment method, which is concretely described as follows. In the ℓ -th block, we duplicate each input clip $\{\mathbf{z}_t^{(\ell-1)}\}_{t=1}^T$ to form an adjacent input pair $\{(\mathbf{z}_t^{(\ell-1)}, \mathbf{z}_{t-1}^{(\ell-1)})\}_{t=2}^T$. Each pair of representations are then concatenated along the channel dimension, which are further fed into a dedicated lightweight convolution module for predicting two interactive points:

$$\mathbf{p}_t^{(\ell)}, \mathbf{p}_{t-1}^{(\ell)} = \text{Conv}(\text{Concat}(\mathbf{z}_t^{(\ell-1)}, \mathbf{z}_{t-1}^{(\ell-1)})), \quad (5)$$

where the convolution module $\text{Conv}(\cdot)$ consists of a sequence of convolution, normalization and pooling layers. The interactive points $\mathbf{p}_t^{(\ell)}, \mathbf{p}_{t-1}^{(\ell)} \in \mathbb{R}^2$ represent the most semantically similar positions in two frames, indicating the region with favorable mutual information. We assume the closer the position is to the interactive point, the more temporal information it involves. On the contrary, a position that is far away from the interactive point can contain re-

dundant and irrelevant information, which should be suppressed. To this end, two align masks $\mathbf{m}_t^{(\ell)}, \mathbf{m}_{t-1}^{(\ell)} \in \mathbb{R}^{h \times w}$ are generated by endowing positions closer to the interactive points with higher weights. Formally, for a spatial position \mathbf{u} in $\mathbf{m}_t^{(\ell)}$, its weight $w_{\mathbf{u}}$ is computed by:

$$s = \text{dist}(\mathbf{u}, \mathbf{p}_t^{(\ell)}),$$

$$w_{\mathbf{u}} = \begin{cases} \eta, & \text{if } s \leq \delta, \\ \max(0, \eta - \beta(s - \delta)), & \text{if } s > \delta, \end{cases} \quad (6)$$

where $\text{dist}(\cdot)$ is the distance function, and η, δ, β are the parameters. The weights of $\mathbf{m}_{t-1}^{(\ell)}$ are obtained by similar calculation with $\mathbf{p}_{t-1}^{(\ell)}$. Note that all the coordinates of positions are scaled to the range $[-1, 1]$ to facilitate the mask calculation. The aligned feature representations $\mathbf{a}_t^{(\ell)}, \mathbf{a}_{t-1}^{(\ell)}$ are produced by weighting the frame features with the align masks:

$$\mathbf{a}_t^{(\ell)} = \mathbf{m}_t^{(\ell)} \mathbf{z}_t^{(\ell-1)}, \quad (7a)$$

$$\mathbf{a}_{t-1}^{(\ell)} = \mathbf{m}_{t-1}^{(\ell)} \mathbf{z}_{t-1}^{(\ell-1)}. \quad (7b)$$

We hypothesize that the aligned feature can implicitly preserve the mutual information and already encodes essential temporal information, which could be leveraged to model the temporal relations across frames. Nevertheless, directly replacing $\mathbf{z}_t^{(\ell-1)}$ with the aligned feature $\mathbf{a}_t^{(\ell)}$ would prejudice the performance, since $\mathbf{a}_t^{(\ell)}$ focuses more on the interaction region while ignoring the spatial correlations. Instead, we consider $\mathbf{a}_t^{(\ell)}$ as a specific temporal signal. Thus, we averagely pool the feature into a single mutual information token $\hat{\mathbf{z}}_{t,mut}^{(\ell)}$ (Eq. (3a)), which is further utilized in spatial multi-head self attention (Eq. (3b)). Note that since we duplicate the input clip to form frame pairs, there are two aligned features for frame $\mathbf{z}_t^{(\ell-1)}$,

$2 \leq t \leq T - 1$. For example, $\mathbf{a}_t^{(\ell)}$ can be computed from both pairs $(\mathbf{z}_t^{(\ell-1)}, \mathbf{z}_{t-1}^{(\ell-1)})$ and $(\mathbf{z}_{t+1}^{(\ell-1)}, \mathbf{z}_t^{(\ell-1)})$. In our implementation, only $\mathbf{a}_t^{(\ell)}$ computed from $(\mathbf{z}_t^{(\ell-1)}, \mathbf{z}_{t-1}^{(\ell-1)})$ is exploited for pooling to the mutual information token.

Our simple alignment implicitly introduces cross-frame cross-location interactions to the model, thus capturing semantically rich actions. We reveal that the primitive pairwise interaction already contains sufficient information for modeling the complex temporal relations, which allows eliminating the costly temporal self-attention in video. Therefore, there is no additional temporal modeling design in IST block.

3.3. Training

The loss function of our framework consists of two parts. The first part is the supervised prompt-enhanced similarity loss, where the cosine similarity between video representation \mathbf{v} and text representation \mathbf{c} is computed by:

$$\begin{aligned} \mathbf{v} &= \text{Avg}(\text{MSA}([\mathbf{z}_{1,cls}^{(L)}, \dots, \mathbf{z}_{T,cls}^{(L)}])), \\ \cos(\mathbf{v}, \mathbf{c}) &= \frac{\langle \mathbf{v}, \mathbf{c} \rangle}{\|\mathbf{v}\| \cdot \|\mathbf{c}\|}. \end{aligned} \quad (8)$$

Here $\text{Avg}(\cdot)$ is the average pooling. The model maximizes $\cos(\mathbf{v}, \mathbf{c})$ if \mathbf{v} and \mathbf{c} are matched, otherwise minimizes it.

The second part is the alignment loss for aligning pairwise frames in each IST block. Particularly, we align the mean-pooled feature, *i.e.* the mutual information token $\hat{\mathbf{z}}_{t,mut}^{(\ell)}$ as in Eq. (3a), using the cosine similarity:

$$\cos_t^{(\ell)} = \frac{\langle \hat{\mathbf{z}}_{t,mut}^{(\ell)}, \hat{\mathbf{z}}_{t-1,mut}^{(\ell)} \rangle}{\|\hat{\mathbf{z}}_{t,mut}^{(\ell)}\| \cdot \|\hat{\mathbf{z}}_{t-1,mut}^{(\ell)}\|}, \quad (9)$$

where $\cos_t^{(\ell)}$ is the similarity score for t -th frame pair in block ℓ . The loss function l_a is formulated by summing up the similarity scores:

$$l_a = - \sum_{\ell=1}^L \sum_{t=2}^T \cos_t^{(\ell)}. \quad (10)$$

Finally, we optimize Eq. (8) and Eq. (10) simultaneously with a loss weight parameter γ .

4. Experiments

We evaluate our method on two datasets: Kinetics-400 [28] and Something-Something-V2 [37]. Four variants are considered, namely the ILA model based on ViT-B/32, ViT-B/16, ViT-L/14, and ViT-L/14@336, respectively. We sparsely sample 8 or 16 frames to form a video clip, both in training and inference. Additional implementation, hyperparameter details, and more experiments are provided in the supplementary materials.

4.1. Main Results

Kinetics-400. In Table 1, we report the performance of our proposed method on Kinetics-400. Comparisons with recent state-of-the-art approaches are listed, including methods with random initialization, pretrained on ImageNet-1k/21k pretraining, and pretrained with web-scale data.

Compared to methods pretrained on ImageNet [13], ILA-ViT-L with 8 frames outperforms the best competitor MViTv2-L [33] by 1.9% in accuracy with $4\times$ fewer FLOPs. We also observe ILA surpasses other baselines with large margins, *e.g.*, Swin [36] and TimeSformer [5]. It indicates the strong representations of the CLIP model, showing the great potential of large-scale visual-language pretraining.

In comparison with methods pretrained on web-scale images, *e.g.* JFT-300M/3B, ILA exhibits significant advantages. Our ILA-ViT-L exceeds ViViT-H by 3.2% with $12\times$ less computation, and exceeds CoVeR by 0.8%. Note that CoVeR uses much more training data (3B images) compared to CLIP (400M image-text pairs).

In addition, when compared with the recent CLIP-based methods, ILA achieves the best performance. ILA-ViT-B with 16 frames surpasses the typical CLIP-based model ActionCLIP-B by 1.9% with $2\times$ fewer FLOPs. Moreover, our largest model outperforms the best competitors X-CLIP-L and EVL-L by 1% with comparable or much less computation. Though MTV-H performs a little higher (89.1%) than ILA (88.7%), it employs the WTS dataset that contains 70M video-text pairs with about 17B images, which are much larger than that in CLIP. The observations show that our alignment-based temporal modeling could capture more comprehensive motion clues than the insufficient temporal attention of X-CLIP and EVL, without increasing the computational burden.

Something-Something-V2. Table 2 reports the comparisons on SSv2. This dataset focuses on the human object action recognition, in which the open domain semantics are limited. We assume the rich textual representation of CLIP language branch can help less. Therefore, we use the cross-entropy loss with one-hot labels, instead of the visual-text similarity loss in Eq. (8). We also increase the number of convolution layers for better alignment. Moreover, we freeze the weights of CLIP for stability.

SSv2 is a motion-heavy dataset and largely depends on temporal modeling. Methods pretrained on CLIP usually produce weaker results compared to those pretrained on Kinetics-400. For example, X-CLIP-B only achieves 57.8% in accuracy, while MViTv1-B produces much higher results (64.7%) with similar computation. Similarly, the result of EVL-ViT-B is also unsatisfactory (61.7%). This phenomenon can be attributed to three factors. (1) The temporal modeling in X-CLIP and EVL is insufficient. In pursuit of high efficiency, X-CLIP and EVL adopt frame-level or local

Table 1. Comparison with the state-of-the-arts on Kinetics-400. The FLOPs per view of each method is reported. We categorize methods by different pretraining data.

Model	Pretrain	Frames	Top-1	Top-5	Views	FLOPs (G)
<i>Random initialization</i>						
MViTv1-B [15]	-	64	81.2	95.1	3x3	455
<i>ImageNet pretraining</i>						
Uniformer-B [31]	IN-1K	32	83.0	95.4	4x3	259
TimeSformer-L [5]	IN-21K	96	80.7	94.7	1x3	2380
ATA [70]	IN-21K	32	81.9	95.5	4x3	793
Mformer-HR [40]	IN-21K	16	81.1	95.2	10x3	959
Swin-L (@384px) [36]	IN-21K	32	84.9	96.7	10x5	2107
MViTv2-L (@312px) [33]	IN-21K	40	86.1	97.0	5x3	2828
<i>Web-scale image pretraining</i>						
ViViT-H/16x2 [3]	JFT-300M	32	84.8	95.8	4x3	8316
TokenLearner-L/10 [44]	JFT-300M	-	85.4	96.3	4x3	4076
CoVeR [66]	JFT-3B	-	87.2	-	1x3	-
<i>Web-scale language-image pretraining</i>						
ActionCLIP-B/16 [56]	CLIP-400M	32	83.8	96.2	10x3	563
A6 [27]	CLIP-400M	16	76.9	93.5	-	-
EVL-ViT-B/16 [35]	CLIP-400M	16	83.6	-	1x3	296
EVL-ViT-L/14 [35]	CLIP-400M	16	87.0	-	1x3	1350
EVL-ViT-L/14@336px [35]	CLIP-400M	32	87.7	-	1x3	6068
X-CLIP-B/16 [38]	CLIP-400M	16	84.7	96.8	4x3	287
X-CLIP-L/14 (@336px) [38]	CLIP-400M	16	87.7	97.4	4x3	3086
AIM-ViT-L/14 [64]	CLIP-400M	16	87.3	97.6	1x3	1868
ST-Adapter-ViT-L/14 [39]	CLIP-400M	16	86.9	97.6	1x3	1375
MTV-H [63]	WTS	32	89.1	98.2	4x3	3705
ILA-ViT-B/32	CLIP-400M	8	81.3	95.0	4x3	40
ILA-ViT-B/32	CLIP-400M	16	82.4	95.8	4x3	75
ILA-ViT-B/16	CLIP-400M	8	84.0	96.6	4x3	149
ILA-ViT-B/16	CLIP-400M	16	85.7	97.2	4x3	295
ILA-ViT-L/14	CLIP-400M	8	88.0	98.1	4x3	673
ILA-ViT-L/14@336px	CLIP-400M	16	88.7	97.8	4x3	3130

temporal attention on top of the CLIP features, which inevitably harms the results. (2) Tuning the weights of CLIP is very challenging, where small perturbations can easily prejudice the primal CLIP. We assume the reason is that SSv2 is a dataset with relatively small semantics. Even assigning a very small learning rate to CLIP weights and a large one to other weights, the model is still prone to encounter exploding gradients. This phenomenon reduces the flexibility of parameter tuning, which leads to the insufficient training of the model. (3) The pretraining on Kinetics can bring significant advantages compared to pretraining on CLIP data.

As shown in the table, ILA-ViT-B (8 frames) achieves a comparable 65.0% with MViTv1-B, which is much higher than X-CLIP and EVL. Moreover, ILA-ViT-L/14@336px obtains promising performance referring to 70.2% on top-1 and 91.8% on top-5. It outperforms EVL-ViT-L/14@336px by 2.2% on top-1 with 2x fewer frames and over 2x fewer FLOPs. It indicates that the proposed implicit alignment

can comprehensively model the temporal information with a low computational cost.

4.2. Ablation Study

Generalization to different backbones. To demonstrate ILA is a versatile module and can be plugged into various backbones, we experiment with a CLIP-based model (EVL-ViT-B/16, 8frames [35]) as well as an ImageNet-based architecture (TimeSformer-ViT-B/16, 8frames [5]). For EVL, we insert our alignment into the CLIP backbone, while keep others unchanged. For TimeSformer, we replace the temporal attention with the proposed alignment module. The results are summarized in Table 3. The utilization of ILA results in a 0.6% and 1.8% performance gain for the CLIP-based and ImageNet-based backbones, respectively, demonstrating ILA is compatible with modern networks.

Effectiveness of implicit alignment. We compare ILA with ATA [70], an alternative of patch alignment, and other temporal modeling approaches, *i.e.* X-CLIP [38], Divided

Table 2. Performance comparison with the state-of-the-arts on Something-Something-V2. The FLOPs per view of each method is reported.

Model	Pretrain	Frames	Top-1 Acc.	Top-5 Acc.	Views	FLOPs (G)
ViViT-L [3]	IN-21K+K400	16	65.4	89.8	1×3	903
TimeSformer-L [5]	IN-21K	96	62.4	81.0	1×3	2380
TimeSformer-HR [5]	IN-21K	16	62.2	78.0	1×3	1703
ATA [70]	IN-21K	32	67.1	90.8	4×3	793
MViTv1-B [15]	K400	16	64.7	89.2	1×3	70.5
MViTv1-B [15]	K400	32	67.1	90.8	1×3	170
Mformer-B [40]	IN-21K+K400	16	66.5	90.1	1×3	370
Mformer-L [40]	IN-21K+K400	32	68.1	91.2	1×3	1185
Mformer-HR [40]	IN-21K+K400	64	67.1	90.6	1×3	959
X-CLIP-B/16 [38]	CLIP-400M	8	57.8	84.5	4×3	145
AIM-ViT-B/16 [64]	CLIP-400M	8	66.4	90.5	1×3	208
AIM-ViT-L/14 [64]	CLIP-400M	32	69.4	92.3	1×3	3836
EVL-ViT-B/16 [35]	CLIP-400M	16	61.7	-	1×3	345
EVL-ViT-L/14 [35]	CLIP-400M	32	66.7	-	1×3	3216
EVL-ViT-L/14@336px [35]	CLIP-400M	32	68.0	-	1×3	8090
ILA-ViT-B/16	CLIP-400M	8	65.0	89.2	4×3	214
ILA-ViT-B/16	CLIP-400M	16	66.8	90.3	4×3	438
ILA-ViT-L/14	CLIP-400M	8	67.8	90.5	4×3	907
ILA-ViT-L/14@336px	CLIP-400M	16	70.2	91.8	4×3	3723

Table 3. Generalization ability of ILA on various visual backbones for Kinetics-400.

Model	Pre-training	Acc. (%)	FLOPs
EVL [35]	CLIP-400M	82.9	150G
EVL + ILA	CLIP-400M	83.5	162G
TimeSformer [5]	IN-21K	78.0	196G
TimeSformer + ILA	IN-21K	79.8	164G

Table 4. Effectiveness of implicit alignment on Kinetics-400. Average Pooling indicates forming the mutual information token in Eq. (3a) without alignment.

Model	Acc. (%)	FLOPs
Baseline	79.8	37G
X-CLIP [38]	80.4	39G
CLIP + Divided ST Attention [5]	80.6	58G
CLIP + Temporal Shift [56]	80.1	37G
CLIP + ATA [70]	81.0	60G
CLIP + Average Pooling	80.4	39G
CLIP + ILA	81.3	40G

Spatio-Temporal Attention [5], Temporal Shift [56], and Average Pooling. The baseline is employing the loss in Eq. (8) for CLIP without temporal modeling. Average Pooling indicates forming the mutual information token in Eq. (3a) without alignment. Table 4 shows the comparison re-

sults. We have the following observations: (1) ILA outperforms the baseline by 1.5% in top-1 accuracy with minor additional computational cost. It indicates that ILA can promote CLIP for video tasks effectively. (2) Compared to ATA that uses patch-level movement for alignment with a cubic complexity, ILA offers better results with nearly 2× fewer FLOPs through learning an implicit mask with a quadratic complexity. (3) ILA also outperforms other approaches like X-CLIP using temporal attention and temporal shifting, highlighting the effectiveness of ILA. (4) ILA achieves better results compared with average pooling, indicating that the improvement results from our implicit alignment instead of the pooling operation.

Table 5. Comparison of mutual information on Kinetics-400. MI (EMD) refers to the average Wasserstein Distance between neighbouring frames.

Model	Acc. (%)	MI (EMD)
Baseline	79.8	0.56
X-CLIP [38]	80.4	0.51
CLIP + Divided ST Attention [5]	80.6	0.47
CLIP + ATA [70]	81.0	0.30
CLIP + ILA	81.3	0.13

Comparison of mutual information. In our work, we assume that ILA can enhance the mutual information between frames, thereby boosting the recognition performance. Here, we compare the mutual information of ILA

in the last visual layer with other approaches. In particular, we calculate the averaged Wasserstein Distance (*i.e.* Earth Mover’s Distance) [2] between adjacent frames which is negatively correlated to mutual information. Table 5 presents the results. We observe that models with additional alignment have lower Wasserstein Distance and higher performance, suggesting that the alignment can indeed correlate adjacent frames.

Impact of different aligning strategies. ILA aligns two consecutive frames and here we experiment with the following alignment strategies: (1) Align-First: each frame is aligned with the first frame; (2) Align-Middle, each frame is aligned with the middle frame. We can observe in Table 6 that anchor frame based alignments are inferior to adjacent alignment. The reason may be that it is not reliable to align two frames that are too far away.

Table 6. Ablation study of different aligning strategies on K-400.

Aligning Strategy	Top1. (%)	Top5. (%)
Align-First	80.7	94.5
Align-Middle	80.8	94.6
Adjacent frame	81.3	95.0

Impact of inserting locations of alignment. We divide the visual branch of ViT-B/32 (12 blocks) into 4 groups, each containing 3 blocks. We plug our ILA into each group individually for exploring the impact of different inserting locations. Table 7 shows the results. Per-block insertion of ILA outperforms the baseline CLIP by 0.6%, 0.7%, 0.6% and 0.5% in accuracy, respectively. We see that inserting ILA into shallow blocks performs slightly better than inserting it into deep ones, showing that aligning low-level features can encode more temporal clues.

Table 7. Comparisons of different inserting locations.

Configuration	Acc. (%)	FLOPs
None	79.8	37G
Block 1-3	80.4	38G
Block 4-6	80.5	38G
Block 7-9	80.4	38G
Block 10-12	80.3	38G
ILA (Block 1-12)	81.3	40G

Operators in alignment module. To validate the effectiveness and efficiency of the 2D convolution module in ILA, we experiment with an alternative choice of window attention in Eq. (5). Table 8 depicts the comparison results. It demonstrates that window attention requires high computational resources and is difficult to optimize, producing limited results.

Table 8. Evaluation of different operators in Eq. (5). We experiment with an alternative of window attention with size 3×3 , instead of the convolution.

Basic Operators	Acc. (%)	FLOPs
2D Convolution	81.3	40G
Window Attention	80.8	114G

Impact of mutual information token. Here we discuss different approaches of exploiting the aligned features. ILA employs a mutual information (MI) token by pooling and concatenation on aligned features. Another choice is the element-wise addition between the frame and the aligned features. In addition, one can also directly concatenate the tokens of aligned features to frame tokens, resulting $2 \times$ tokens in spatial attention.

The results are shown in Table 9. It can be observed that both element-wise addition and direct concatenation perform inferior to the ILA. Furthermore, their inference latencies are much higher than ILA. The plausible reason is that the aligned features are produced by simple mask weighting of frame features, thus containing much redundant information when performing addition or concatenation. Meanwhile, the pooling operation can effectively remove such irrelevant information and boost the model performance.

Table 9. Ablation study of mutual information (MI) token. ILA employs a MI token by pooling & concatenation with aligned features. Other choices include element-wise addition, or direct concatenation.

Implementation	Acc.	FLOPs	Latency (ms)
Element-wise Addition	80.2	44G	64.029
Direct Concat.	80.6	45G	58.548
Pooling & Concat. (ILA)	81.3	40G	47.075

Impact of CLIP initialization. In order to eliminate the influence of CLIP weights, we initialize the weights of ILA with ViT-B/16 pretrained on IN-21K, as well as removing the text branch and using the one-hot labels instead, which is same to the Swin. The results on K400 are shown on Table 10. ILA outperforms Swin by 0.6% on Top-1 accuracy under the same pretraining setting, indicating the superiority of ILA. The results demonstrate that our proposed model obtains the promising performance due to ILA itself instead of language encoder of CLIP or CLIP pretraining weights.

Table 10. Comparison results between ILA and Swin under the same weight parameters initialization.

Model	Pretraining	R-1	R-5	Views
Swin-B/32f	IN-21K	82.7	95.5	4x3
ILA (ViT-B/16-32f)	IN-21K	83.3	95.8	4x3

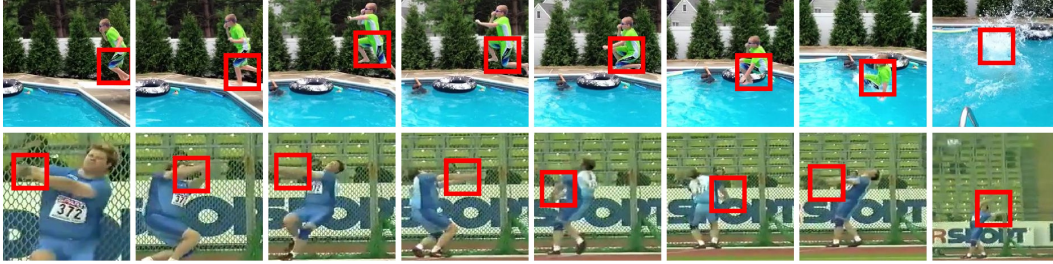


Figure 5. Visualization of mutual information over time. We draw tracking borders around each interactive point in a fixed size, which the tracking borders along temporal dimension depict the temporal corresponding mutual information captured by ILA.

Ablation of text representation tuning. In terms of ILA architecture, We follow X-CLIP to adopt the video-conditioned text representation tuning. For fair comparisons, we run an additional ablation by removing it on K400 with ViT-B/32-8f. In Table 11, we can see that ILA still performs better than X-CLIP even without the video-specific tuning.

Table 11. Ablation results of text representation tuning.

Model	w/o video-specific text	w/ video-specific text
X-CLIP	79.6	80.4
ILA	80.8	81.3

4.3. Additional Comparison Results

Zero-shot results on SSv2. We train ILA and other CLIP-based competitors on K400 with ViT-B/16-8f and evaluate on SSv2 in a zero-shot setting, by training an additional linear classification layer. The results are depicted on Table 12. We see that ILA outperforms two typical competitors X-CLIP and EVL by 5.8% and 8.7% on Top-1 respectively, highlighting the effectiveness of ILA.

Table 12. Zero-shot results on SSv2.

Model	Pretraining	Top1 (%)	Top5 (%)
X-CLIP [38]	CLIP-400M	38.1	68.1
EVL [35]	CLIP-400M	35.2	65.4
AIM [64]	CLIP-400M	39.1	68.7
ILA	CLIP-400M	43.9	71.8

Performance on additional benchmark. We evaluate ILA, X-CLIP [38] and EVL [35] on the temporal understanding benchmark [72] which is without static biases. The results are shown in Table 13. All models are pretrained on CLIP-400M (ViT-B/16-8f). Top-1 (T) and Top-1 (S) refer to the traditional accuracy on Temporal-50 and Static-50, respectively. TS refers to the relative gain of the model on temporal classes compared to static ones (Temporal score). In Table 13, we can see that our ILA outperforms X-CLIP

by 6.0% on Top-1 and by 3.8% on temporal score. ILA also exceeds the best temporal score on RGB modality (5.2%, R3D) in the paper [72]. The result highlights the effectiveness of ILA on the temporal understanding benchmark.

Table 13. Performance comparisons on additional benchmark without static biases.

Model	Top-1 (T)	Top-5 (T)	Top-1 (S)	TS
X-CLIP	75.9	94.1	73.6	2.3
EVL	70.7	92.5	68.3	2.4
ILA	81.9	97.4	75.8	6.1

4.4. Visualization of mutual information

We visualize the interactive point at each video frame by drawing the bounding boxes centered at the points, as illustrated in the Figure 5. The bounding boxes indicate the regions with rich favorable mutual information, where the moving boxes show the potential of abilities on tracking moving objects.

5. Conclusion

We introduced Implicit Learnable Alignment (ILA), a novel temporal modeling method for video recognition. ILA performs frame alignment so as to encode motion information in lieu of the widely used temporal attention operation. Particularly, ILA employs only an implicit and coarse feature alignment via a weighting mask. By finding the active interaction position in the frame, the mask is generated with higher weights around that position, and lower weights on others. Extensive experiments demonstrates the effectiveness and efficiency of ILA, showcasing its promising adaption ability to CLIP and compatibility with modern visual backbones.

Acknowledgement This project was supported by NSFC under Grant No. 62032006 and No. 62102092.

References

- [1] Hassan Akbari, Liangzhe Yuan, Rui Qian, Wei-Hong Chuang, Shih-Fu Chang, Yin Cui, and Boqing Gong. Vatt: Transformers for multimodal self-supervised learning from raw video, audio and text. In *NeurIPS*, 2021. 2
- [2] Martin Arjovsky, Soumith Chintala, and Leon Bottou. Wasserstein gan. *arXiv preprint arXiv:1701.07875*, 2017. 8
- [3] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid. Vivit: A video vision transformer. In *ICCV*, 2021. 1, 2, 3, 6, 7
- [4] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016. 3
- [5] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is space-time attention all you need for video understanding? In *ICML*, 2021. 1, 2, 3, 5, 6, 7, 12, 13, 14
- [6] Cheng, Zhi-Qi and Wu, Xiao and Liu, Yang and Hua, Xian-Sheng. Video2shop: Exact matching clothes in videos to online shopping images. In *CVPR*, 2017. 1
- [7] Cheng, Zhi-Qi and Wu, Xiao and Liu, Yang and Hua, Xian-Sheng. Video ecommerce++: Toward large scale online video advertising. In *TMM*, 2017. 1
- [8] Dimitri P Bertsekas. A new algorithm for the assignment problem. *Mathematical Programming*, 21(1):152–171, 1981. 4
- [9] Jiezhong Cao, Yawei Li, Kai Zhang, and Luc Van Gool. Video super-resolution transformer. *arXiv preprint arXiv:2106.06847*, 2021. 3
- [10] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *CVPR*, 2017. 1
- [11] Kelvin CK Chan, Shangchen Zhou, Xiangyu Xu, and Chen Change Loy. Basicvsr++: Improving video super-resolution with enhanced propagation and alignment. In *CVPR*, 2022. 3
- [12] Yu-Wei Chao, Sudheendra Vijayanarasimhan, Bryan Seybold, David A Ross, Jia Deng, and Rahul Sukthankar. Rethinking the faster r-cnn architecture for temporal action localization. In *CVPR*, 2018. 1
- [13] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009. 5
- [14] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021. 2, 3
- [15] Haoqi Fan, Bo Xiong, Karttikeya Mangalam, Yanghao Li, Zhicheng Yan, Jitendra Malik, and Christoph Feichtenhofer. Multiscale vision transformers. In *ICCV*, 2021. 1, 2, 6, 7
- [16] He, Jun-Yan and Cheng, Zhi-Qi and Li, Chenyang and Xiang, Wangmeng and Chen, Binghui and Luo, Bin and Geng, Yifeng and Xie, Xuansong. DAMO-StreamNet: Optimizing Streaming Perception in Autonomous Driving. In *IJCAI*, 2023. 1
- [17] Christoph Feichtenhofer. X3d: Expanding architectures for efficient video recognition. In *CVPR*, 2020. 1
- [18] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video recognition. In *ICCV*, 2019. 1
- [19] Cheng, Zhi-Qi and Liu, Yang and Wu, Xiao and Hua, Xian-Sheng. Video ecommerce: Towards online video advertising. In *ACM MM*, 2016. 1
- [20] Cheng, Zhi-Qi and Dai, Qi and Li, Siyao and Mitamura, Teruko and Hauptmann, Alexander. Gsrformer: Grounded situation recognition transformer with alternate semantic attention refinement. In *ACM MM*, 2022. 1
- [21] Peng Gao, Shijie Geng, Renrui Zhang, Teli Ma, Rongyao Fang, Yongfeng Zhang, Hongsheng Li, and Yu Qiao. Clip-adapter: Better vision-language models with feature adapters. *arXiv preprint arXiv:2110.04544*, 2021. 2
- [22] Qi Han, Zejia Fan, Qi Dai, Lei Sun, Ming-Ming Cheng, Jiaying Liu, and Jingdong Wang. On the connection between local attention and dynamic depth-wise convolution. In *ICLR*, 2022. 2
- [23] Kensho Hara, Hirokatsu Kataoka, and Yutaka Satoh. Learning spatio-temporal features with 3d residual networks for action recognition. In *ICCV*, 2017. 2
- [24] Rui Hou, Chen Chen, and Mubarak Shah. Tube convolutional neural network (t-cnn) for action detection in videos. In *ICCV*, 2017. 1
- [25] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *ICML*, 2021. 2
- [26] Shihao Jiang, Dylan Campbell, Yao Lu, Hongdong Li, and Richard Hartley. Learning to estimate hidden motions with global motion aggregation. In *ICCV*, 2021. 3
- [27] Chen Ju, Tengda Han, Kunhao Zheng, Ya Zhang, and Weidi Xie. Prompting visual-language models for efficient video understanding. In *ECCV*, 2022. 1, 2, 6
- [28] 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. In *CVPR*, 2017. 1, 5
- [29] Heeseung Kwon, Manjin Kim, Suha Kwak, and Minsu Cho. Motionsqueeze: Neural motion feature learning for video understanding. In *ECCV*, 2020. 3
- [30] Dong Li, Zhaofan Qiu, Qi Dai, Ting Yao, and Tao Mei. Recurrent tubelet proposal and recognition networks for action detection. In *ECCV*, 2018. 1
- [31] Kunchang Li, Yali Wang, Peng Gao, Guanglu Song, Yu Liu, Hongsheng Li, and Yu Qiao. Uniformer: Unified transformer for efficient spatiotemporal representation learning. In *ICLR*, 2022. 1, 6
- [32] Shuyuan Li, Huabin Liu, Rui Qian, Yuxi Li, John See, Mengjuan Fei, Xiaoyuan Yu, and Weiyao Lin. Ta2n: Two-stage action alignment network for few-shot action recognition. In *AAAI*, 2022. 3
- [33] Yanghao Li, Chao-Yuan Wu, Haoqi Fan, Karttikeya Mangalam, Bo Xiong, Jitendra Malik, and Christoph Feichtenhofer. Improved multiscale vision transformers for classification and detection. In *CVPR*, 2022. 5, 6

- [34] Jiayi Lin, Yan Huang, and Liang Wang. Fdan: Flow-guided deformable alignment network for video super-resolution. *arXiv preprint arXiv:2105.05640*, 2021. 3
- [35] Ziyi Lin, Shijie Geng, Renrui Zhang, Peng Gao, Gerard de Melo, Xiaogang Wang, Jifeng Dai, Yu Qiao, and Hongsheng Li. Frozen clip models are efficient video learners. In *ECCV*, 2022. 1, 2, 3, 6, 7, 9
- [36] Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang, Stephen Lin, and Han Hu. Video swin transformer. In *CVPR*, 2022. 1, 2, 5, 6
- [37] Joanna Materzynska, Tete Xiao, Roei Herzig, Huijuan Xu, Xiaolong Wang, and Trevor Darrell. Something-else: Compositional action recognition with spatial-temporal interaction networks. In *CVPR*, 2020. 5
- [38] Bolin Ni, Houwen Peng, Minghao Chen, Songyang Zhang, Gaofeng Meng, Jianlong Fu, Shiming Xiang, and Haibin Ling. Expanding language-image pretrained models for general video recognition. In *ECCV*, 2022. 1, 2, 3, 6, 7, 9, 12, 14
- [39] Junting Pan, Ziyi Lin, Xiatian Zhu, Jing Shao, and Hongsheng Li. St-adapter: Parameter-efficient image-to-video transfer learning for action recognition. In *NeurIPS*, 2022. 1, 2, 3, 6
- [40] Mandela Patrick, Dylan Campbell, Yuki Asano, Ishan Misra, Florian Metze, Christoph Feichtenhofer, Andrea Vedaldi, and João F Henriques. Keeping your eye on the ball: Trajectory attention in video transformers. In *NeurIPS*, 2021. 3, 6, 7
- [41] 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 *ICML*, 2021. 1, 2, 3
- [42] Deva Ramanan, David A Forsyth, and Andrew Zisserman. Strike a pose: Tracking people by finding stylized poses. In *CVPR*, 2005. 2, 3
- [43] Yongming Rao, Wenliang Zhao, Guangyi Chen, Yansong Tang, Zheng Zhu, Guan Huang, Jie Zhou, and Jiwen Lu. Denseclip: Language-guided dense prediction with context-aware prompting. In *CVPR*, 2022. 2
- [44] Michael Ryoo, AJ Piergiovanni, Anurag Arnab, Mostafa Dehghani, and Anelia Angelova. Tokenlearner: Adaptive space-time tokenization for videos. In *NeurIPS*, 2021. 6
- [45] Baifeng Shi, Qi Dai, Judy Hoffman, Kate Saenko, Trevor Darrell, and Huijuan Xu. Temporal action detection with multi-level supervision. In *ICCV*, 2021. 1
- [46] Baifeng Shi, Qi Dai, Yadong Mu, and Jingdong Wang. Weakly-supervised action localization by generative attention modeling. In *CVPR*, 2020. 1
- [47] Shuwei Shi, Jinjin Gu, Liangbin Xie, Xintao Wang, Yujia Yang, and Chao Dong. Rethinking alignment in video super-resolution transformers. In *NeurIPS*, 2022. 3
- [48] Zheng Shou, Dongang Wang, and Shih-Fu Chang. Temporal action localization in untrimmed videos via multi-stage cnns. In *CVPR*, 2016. 1
- [49] Karen Simonyan and Andrew Zisserman. Two-stream convolutional networks for action recognition in videos. In *NeurIPS*, 2014. 1, 2
- [50] Rui Tian, Zuxuan Wu, Qi Dai, Han Hu, Yu Qiao, and Yu-Gang Jiang. Resformer: Scaling vits with multi-resolution training. In *CVPR*, 2023. 1, 2
- [51] Yapeng Tian, Yulun Zhang, Yun Fu, and Chenliang Xu. Tdan: Temporally-deformable alignment network for video super-resolution. In *CVPR*, 2020. 3
- [52] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *ICML*, 2021. 2
- [53] Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. Learning spatiotemporal features with 3d convolutional networks. In *ICCV*, 2015. 1, 2
- [54] Du Tran, Heng Wang, Lorenzo Torresani, Jamie Ray, Yann LeCun, and Manohar Paluri. A closer look at spatiotemporal convolutions for action recognition. In *CVPR*, 2018. 2
- [55] Limin Wang, Yuanjun Xiong, Zhe Wang, Yu Qiao, Dahua Lin, Xiaoou Tang, and Luc Van Gool. Temporal segment networks: Towards good practices for deep action recognition. In *ECCV*, 2016. 1, 2
- [56] Mengmeng Wang, Jiazheng Xing, and Yong Liu. Action-clip: A new paradigm for video action recognition. In *ECCV*, 2022. 1, 2, 6, 7
- [57] Qiang Wang, Li Zhang, Luca Bertinetto, Weiming Hu, and Philip HS Torr. Fast online object tracking and segmentation: A unifying approach. In *CVPR*, 2019. 1
- [58] Xintao Wang, Kelvin CK Chan, Ke Yu, Chao Dong, and Chen Change Loy. Edvr: Video restoration with enhanced deformable convolutional networks. In *CVPR*, 2019. 3
- [59] Mingyu Wu, Boyuan Jiang, Donghao Luo, Junchi Yan, Yabiao Wang, Ying Tai, Chengjie Wang, Jilin Li, Feiyue Huang, and Xiaokang Yang. Learning comprehensive motion representation for action recognition. In *AAAI*, 2021. 3
- [60] Saining Xie, Chen Sun, Jonathan Huang, Zhuowen Tu, and Kevin Murphy. Rethinking spatiotemporal feature learning: Speed-accuracy trade-offs in video classification. In *ECCV*, 2018. 1, 2
- [61] Zhen Xing, Qi Dai, Han Hu, Jingjing Chen, Zuxuan Wu, and Yu-Gang Jiang. Svformer: Semi-supervised video transformer for action recognition. In *CVPR*, 2023. 1
- [62] Hu Xu, Gargi Ghosh, Po-Yao Huang, Dmytro Okhonko, Armen Aghajanyan, Florian Metze, Luke Zettlemoyer, and Christoph Feichtenhofer. Videoclip: Contrastive pre-training for zero-shot video-text understanding. In *EMNLP*, 2021. 2
- [63] Shen Yan, Xuehan Xiong, Anurag Arnab, Zhichao Lu, Mi Zhang, Chen Sun, and Cordelia Schmid. Multiview transformers for video recognition. In *CVPR*, 2022. 6
- [64] Taojiannan Yang, Yi Zhu, Yusheng Xie, Aston Zhang, Chen Chen, and Mu Li. Aim: Adapting image models for efficient video action recognition. In *ICLR*, 2023. 1, 2, 3, 6, 7, 9
- [65] Lu Yuan, Dongdong Chen, Yi-Ling Chen, Noel Codella, Xiyang Dai, Jianfeng Gao, Houdong Hu, Xuedong Huang, Boxin Li, Chunyuan Li, et al. Florence: A new foundation model for computer vision. *arXiv preprint arXiv:2111.11432*, 2021. 2, 3

- [66] Bowen Zhang, Jiahui Yu, Christopher Fifty, Wei Han, Andrew M Dai, Ruoming Pang, and Fei Sha. Co-training transformer with videos and images improves action recognition. *arXiv preprint arXiv:2112.07175*, 2021. 6
- [67] Renrui Zhang, Rongyao Fang, Wei Zhang, Peng Gao, Kunchang Li, Jifeng Dai, Yu Qiao, and Hongsheng Li. Tip-adapter: Training-free clip-adapter for better vision-language modeling. In *ECCV*, 2022. 2
- [68] Renrui Zhang, Ziyu Guo, Wei Zhang, Kunchang Li, Xupeng Miao, Bin Cui, Yu Qiao, Peng Gao, and Hongsheng Li. Pointclip: Point cloud understanding by clip. In *CVPR*, 2022. 2
- [69] Xiaosong Zhang, Yunjie Tian, Lingxi Xie, Wei Huang, Qi Dai, Qixiang Ye, and Qi Tian. Hivit: A simpler and more efficient design of hierarchical vision transformer. In *ICLR*, 2023. 2
- [70] Yizhou Zhao, Zhenyang Li, Xun Guo, and Yan Lu. Alignment-guided temporal attention for video action recognition. In *NeurIPS*, 2022. 3, 4, 6, 7, 12, 13, 14
- [71] Bolei Zhou, Alex Andonian, Aude Oliva, and Antonio Torralba. Temporal relational reasoning in videos. In *ECCV*, 2018. 2
- [72] Sevilla-Lara, Laura and Zha, Shengxin and Yan, Zhicheng and Goswami, Vedanuj and Feiszli, Matt and Torresani, Lorenzo. Only time can tell: Discovering temporal data for temporal modeling. In *WACV*, 2021. 9
- [73] Jiajun Deng, Yingwei Pan, Ting Yao, Wengang Zhou, Houqiang Li, and Tao Mei. Single shot video object detector. In *TMM*, 2020. 1
- [74] Weng, Zejia and Yang, Xitong and Li, Ang and Wu, Zuxuan and Jiang, Yu-Gang. Open-VCLIP: Transforming CLIP to an Open-vocabulary Video Model via Interpolated Weight Optimization. In *ICML*, 2023. 1
- [75] Wang, Junke and Chen, Dongdong and Wu, Zuxuan and Luo, Chong and Zhou, Luwei and Zhao, Yucheng and Xie, Yujia and Liu, Ce and Jiang, Yu-Gang and Yuan, Lu. Omnivl: One foundation model for image-language and video-language tasks. In *NeurIPS*, 2022. 1
- [76] Wu, Zuxuan and Li, Hengduo and Xiong, Caiming and Jiang, Yu-Gang and Davis, Larry Steven. A dynamic frame selection framework for fast video recognition. In *TPAMI*, 2020. 1

Appendix

A. Implementation Details of ILA

Training Details. The experiments are conducted on 8 NVIDIA 32G V100 GPUs. The training configuration is listed in Table 14. It is worth noting that our sampling strategies for Kinetics-400 and Something-Something-V2 are different during the training phase. We implement the sparse sampling strategy on Kinetics-400. For SSv2, we uniformly sample the entire video at predefined temporal intervals without group division. In term of the training on Kinetics-400, the base learning rate indicates the learning rate of the original CLIP parameters. The learning rate

for other additional parameters is $10\times$ larger than the base learning rate. In term of the training on SSv2, we exclude the prompt branch and freeze the weights of CLIP visual branch for training stability. Thus the base learning rate is used for the rest parameters.

Table 14. Default implementation details of our method.

Training Configuration	Kinetics-400	Something-Something v2
Optimisation		
Optimizer		AdamW
Optimizer betas		(0.9,0.98)
Batch size		256
Learning rate schedule		Cosine
Learning warmup epochs		5
Base learning rate	8e-6	5e-4
Minimal learning rate	8e-8	5e-6
training steps	50000	30000
Data augmentation		
RandomFlip		0.5
MultiScaleCrop		(1, 0.875, 0.75, 0.66)
ColorJitter		0.8
GrayScale		0.2
Label smoothing		0.1
Mixup		0.8
Cutmix		1.0
Other regularisation		
Weight decay	0.003	0.01

Convolution Module in SSv2. In SSv2, we increase the number of convolution layers in alignment. Particularly, two additional 3×3 convolution layers plus batch normalization and ReLU are added. In comparison to the original convolution module, it can bring 0.6% improvement on top-1 accuracy.

B. Complexity of ILA

We analyze various temporal modeling methods (Spatial Attention [5], Joint Attention [5], Divided ST Attention [5], ATA [70], X-CLIP [38] and our proposed ILA) in terms of complexity, as shown in Table 15. The complexity of our alignment process is $O(Thwk^2d)$ due to the 2D convolution-based operations. The complexity of the whole ILA consists of the implicit alignment $O(Thwk^2d)$ and the spatial attention $O(Th^2w^2d)$. In terms of Joint Attention and Divided Spatiotemporal Attention, Joint Attention requires more computational memory since it takes all patches into consideration. Divided ST Attention only considers the temporal attention along the time axis. In terms of ATA, ATA is based on Hungarian Algorithm whose complexity is $O(N^3)$. In practice, the complexity of Hungarian matching is $O(Th^3w^3d)$ in video domain. Moreover, ATA requires additional temporal attention with complexity $O(T^2hwd)$. X-CLIP adopts a frame-level temporal attention with complexity $O(T^2d)$, which however obtains sub-

optimal result. We can observe that our proposed ILA can have better performance in low complexity.

C. Qualitative Analysis

In order to investigate the quality of three temporal modeling approaches (Divided ST Attention [5], ATA [70], and ILA), we visualize their intermediate and last feature maps respectively, as shown in Figure 6 and Figure 7. According to the illustrations, all three approaches capture the static semantic features, such as static flowers on the desk. Moreover, our proposed ILA pays more attention to the action area of arranging flowers (*e.g.* the 5-th frame in the last row of Figure 7) instead of the static flowers on the desk. It indicates that our ILA can leverage the learnable mask to achieve implicit temporal modeling, focusing on the vital motion region. For divided ST attention, the model prefers to focus on static object instead of significant actions. While in ATA, the model attempts to concentrate on discontinuous regions with inaccurate positions. The plausible reason is that ATA utilizes patch movement-based alignment, which may destroys the continuity of semantic distribution.

D. Key differences between ILA and ATA

ATA adopts an explicit patch-level alignment with Hungarian matching, aiming at modeling temporal attention within aligned patches, which has poor efficiency due to the frame-by-frame serial alignment. Our ILA is fundamentally different as we utilize learnable masks to obtain implicit and coarse semantic-level alignment, which attempts to enhance favorable mutual information and can be performed in parallel with high efficiency.

It exits three fundamental different aspects. First, ATA can only align the collection of frames representations in serial frame-by-frame mode due to the limitation of KMA, while our ILA can utilize learnable masks to align semantical correspondences between two neighboring frames in parallel resulting in faster inference. Second, the complexity of ATA is $O(N^3)$ and ATA is unlearnable resulting in difficult optimization. Computational complexity of ILA is $O(N^2)$. Third, the core idea of ATA is to implement KMA algorithm to find out the optimal patch-level movement scheme capturing temporal correspondences, while the core idea of our ILA is to utilize specific masks to suppress irrelevant redundant information and enhance task-related mutual information among frames resulting in implicit alignment. Therefore, ATA still preserve the original irrelevant redundant information, while our ILA has the suppression of irrelevant redundant information due to principle of masks.

Table 15. Complexities of different methods, with results on Kinetics-400. T , h , w , d , and k refer to temporal size, spatial height of input, spatial width of input, channel depth of input, and kernel size of convolution, respectively.

Temporal Modeling	Complexity	Acc.(%)	FLOPs
Spatial Attention [5]	$O(Th^2w^2d)$	79.8	37G
Joint Attention [5]	$O(T^2h^2w^2d)$	80.4	71G
Divided ST Attention [5]	$O(T^2hwd + Th^2w^2d)$	80.6	58G
ATA [70]	$O(Th^3w^3d + T^2hwd + Th^2w^2d)$	81.0	60G
X-CLIP [38]	$O(T^2d + Th^2w^2d)$	80.4	39G
ILA	$O(Thwk^2d + Th^2w^2d)$	81.3	40G

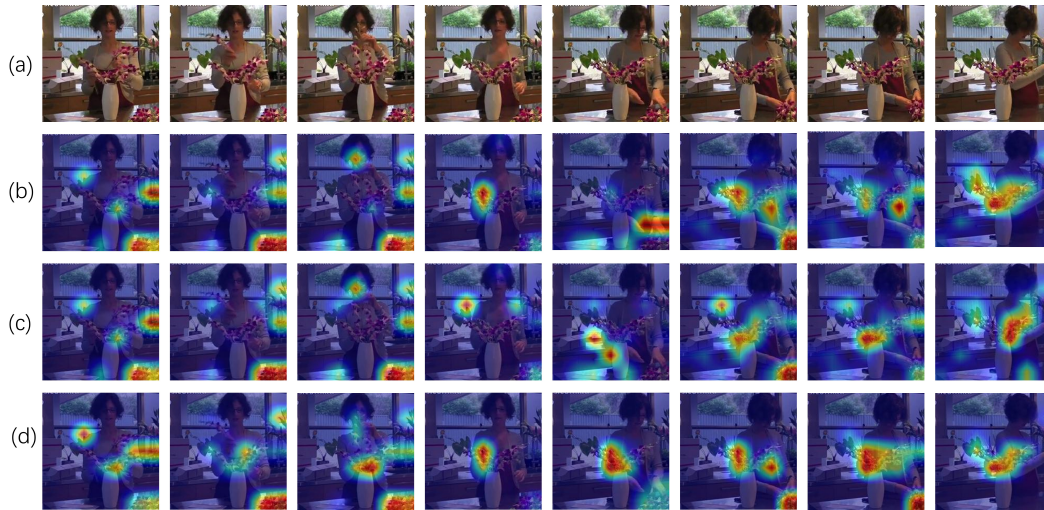


Figure 6. Visualization of intermediate feature map of different temporal modeling approaches on Kinetics-400. (a) refers to raw frames. (b), (c) and (d) refer to Divided ST Attention, ATA and ILA respectively.

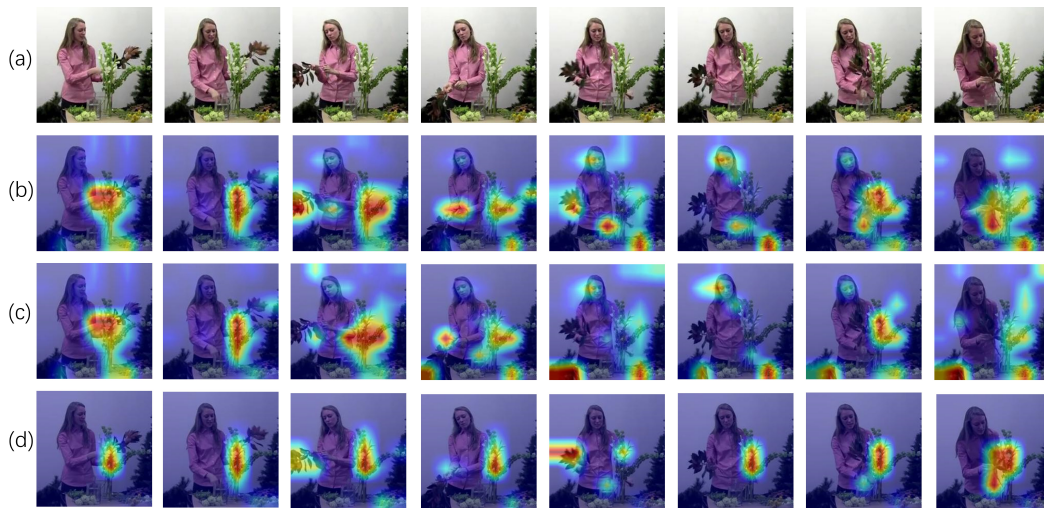


Figure 7. Visualization of the last feature map of different temporal modeling approaches on Kinetics-400. (a) refers to raw frames. (b), (c) and (d) refer to Divided ST Attention, ATA and ILA respectively.