# PAN: Towards Fast Action Recognition via Learning Persistence of Appearance

Can Zhang, Yuexian Zou*, *Senior Member, IEEE,* Guang Chen, and Lei Gan

*Abstract*—Efficiently modeling dynamic motion information in videos is crucial for action recognition task. Most state-of-the-art methods heavily rely on dense optical flow as motion representation. Although combining optical flow with RGB frames as input can achieve excellent recognition performance, the optical flow extraction is very time-consuming. This undoubtably will count against real-time action recognition. In this paper, we shed light on fast action recognition by lifting the reliance on optical flow. Our motivation lies in the observation that small displacements of motion boundaries are the most critical ingredients for distinguishing actions, so we design a novel motion cue called Persistence of Appearance (PA). In contrast to optical flow, our PA focuses more on distilling the motion information at boundaries. Also, it is more efficient by only accumulating pixel-wise differences in feature space, instead of using exhaustive patch-wise search of all the possible motion vectors. Our PA is over $1000\times$ faster (8196fps *vs.* 8fps) than conventional optical flow in terms of motion modeling speed. To further aggregate the short-term dynamics in PA to long-term dynamics, we also devise a global temporal fusion strategy called Various-timescale Aggregation Pooling (VAP) that can adaptively model long-range temporal relationships across various timescales. We finally incorporate the proposed PA and VAP to form a unified framework called Persistent Appearance Network (PAN) with strong temporal modeling ability. Extensive experiments on six challenging action recognition benchmarks verify that our PAN outperforms recent state-of-the-art methods at low FLOPs. *Codes and models are available at:* *https://github.com/zhang-can/PAN-PyTorch*.

*Index Terms*—Fast Action Recognition, Motion Representation, Persistent Appearance Network, Persistence of Appearance
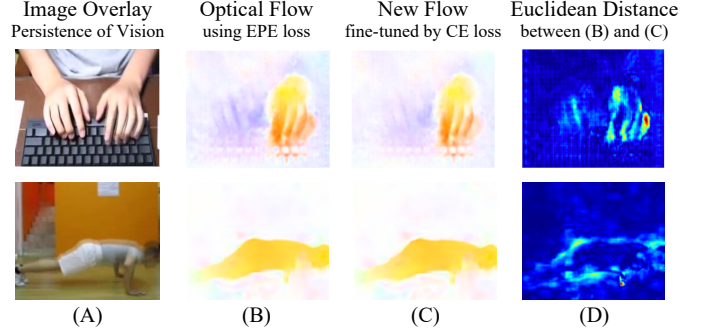


Fig. 1. Comparison of the estimated optical flow using conventional EPE loss (B) and the new flow fine-tuned by Cross-entropy loss (C). The part (D) is obtained by calculating the Euclidean distance at each pixel between (B) and (C). In (D), the optical flow vectors change more around human motion boundaries. *Best viewed in color and zoomed in.*

## I. INTRODUCTION

**V**IDEO action recognition has achieved impressive performance improvements in recent years, mainly due to the aid of deep models [1], [2], [3] and large datasets [4], [5], [6]. Although many prominent 2D CNNs [7], [8], [9] have been designed for image recognition task, these 2D CNNs cannot model effective dynamic motion by naively extending them to video domain. The inherent complexity of temporal evolution in videos makes *motion modeling* in action recognition still a very challenging task.

Optical flow can encode apparent motion of moving objects in visual scenes. When combining optical flow with RGB frames as input, two-stream CNNs variants [1], [10], [11], [12], [13] and 3D CNNs variants [14], [2], [15], [16] largely outperform their counterparts using only RGB frames as input. Thus, optical flow has been commonly used as motion representation for video action recognition task. However, extract-

ing optical flow is time-costly. Aiming at mitigating this inefficiency, several recent works [17], [18], [19], [20] introduce some fast and accurate optical flow estimation methods based on CNNs. But when applied to action recognition task, such solutions are still sub-optimal for two reasons: (1) Computing optical flow in advance makes action recognition a two-stage task. This two-stage paradigm is time-consuming, storage-demanding and not end-to-end trainable. (2) Improving the estimation accuracy of optical flow is not well correlated with boosting the action recognition performance, which has been demonstrated by many works [21], [22], [13].

To address those issues, we design an alternative motion representation to replace optical flow for action recognition task. Ideally, it should be: *efficient* to compute, *effective* in performance, *flexible* to implement and moreover, free of pre-computation and storage. To this end, we need to investigate **which parts of a moving object are most critical for distinguishing actions.**

Generally, most existing action recognition methods follow a two-stage procedure: they first estimate optical flow using the EPE loss[1], and then the estimated optical flow is fed into the subsequent action recognition module. The assumption of these methods is that more accurate optical flow (measured by EPE) will bring superior action recognition performance. However, this correlation is weak as demostrated by Sevilla-Lara *et al.* [21] and Zhu *et al.* [13]. They combined the optical flow estimation and action recognition modules to form a unified network. The optical flow estimation module is fine-

C. Zhang, Y. Zou, G. Chen and L. Gan are with the School of Electrical and Computer Engineering, Peking University, China (e-mail: {zhangcan, zouyx, guangchen, ganlei}@pku.edu.cn).

[1]The evaluation metric of optical flow quality is end-point-error (EPE), the average Euclidean distance between the estimated and the ground-truth flow.

tuned by optimizing the final recognition (CE) loss instead of the EPE loss. Compared with the common two-stage methods, this unified manner achieves better recognition performance. To figure out *what matters most that leads to the performance improvements*, we analyze the visualization result comparisons of the estimated optical flow using the EPE loss, the new flow fine-tuned by the CE loss and their difference computed by the Euclidean distance, as shown in Fig. 1. The visualization results are from [21]. From Fig. 1-D, we can clearly observe that the most salient parts are movement variations occurring at motion boundaries. According to the aforementioned analysis, we can conclude that **small displacements of motion boundaries** play a vital role in action recognition.

Inspired by this observation, we design a new motion cue which derives from optical flow yet focuses more on the small displacements of motion boundaries. From the perspective of human perception, a series of video frames give people a sense of motion when viewed in order at a certain speed. This phenomenon is termed as *Persistence of Vision*, as shown in Fig. 1-A blurred areas. Since we aim to extract motion information directly from RGB frames (*a.k.a*, the appearance information), we name the proposed motion cue as *Persistence of Appearance* (PA). Our PA enjoys high efficiency because we do not use exhaustive search of all the possible motion vectors like optical flow does. Instead, PA only contains pixel-wise operations in feature space. Specifically, given two adjacent RGB frames, our PA first computes pixel-wise intensity variations in feature space and these variations are further accumulated to manifest the motion magnitude. With such a design, PA can model the small displacements at motion boundaries because: (1) *small displacements* are perceived since pixel-wise differences reflect the displacements of a small receptive field in input space; (2) *motion boundaries* are captured since general patterns, *e.g.*, boundaries, texture, etc, can be encoded by the first few convolutional layers [23]. So the differences among low-level feature maps can reflect the variations at boundaries.

In this way, our PA represents instantaneous motion information. However, most human actions last for a while, ranging from seconds to minutes or even longer, so long-term temporal modeling is of great importance for video action recognition. In order to aggregate the short-term dynamics contained in PA to long-term dynamics, we design a global temporal fusion strategy called *Various-timescale Aggregation Pooling* (VAP). It enables the network to model long-range temporal relationships across various timescales. We further incorporate the proposed motion representation PA and the global temporal fusion strategy VAP into a unified ConvNet called *Persistent Appearance Network* (PAN), which achieves fast action recognition with no upfront cost.

The preliminary work is published in ACM MM 2019 [24] and we have extended it in several significant aspects:

- *First*, we add more exploration studies on various network depths and investigate two encoding schemes. These new analyses further demonstrate the efficient motion modeling ability of our PA.
- *Second*, we improve the previous Various-timescale Inference Pooling (VIP) [24] to VAP. In the previous VIP, the

weights of inference function for score fusion are fixed hyper-parameters. As for the enhanced version VAP, we design a more intelligent weight perception scheme that learns to adaptively aggregate the recalibrated features across different timescales. Substantial new analyses are also provided to the improved method VAP.
- *Third*, our experiments are extended from scene-dominant datasets (Kinetics400, UCF101 and HMDB51) to more challenging temporal-dominant datasets (Something-Something-V1 & V2 and Jester). Since temporal modeling is more critical than RGB scene information for recognizing actions in temporal-dominant datasets, the motion modeling ability of our proposed approach can be better demonstrated, *i.e.*, extracting motion information directly from RGB frames.

## II. RELATED WORK

**Network Architecture.** Spatial appearance and temporal motion are two essential ingredients for action recognition. Modeling effective spatiotemporal information still remains a challenging task. Generally, from the architectural perspective, conventional CNN-based methods can be summarized into two categories: (1) Two-stream CNNs [1], [10], [11], [12], [13] that separately process RGB frames (spatial stream) and pre-computed optical flow (temporal stream) using two 2D CNNs and finally apply late fusion strategy to obtain spatiotemporal semantics; (2) 3D CNNs [14], [2], [15], [16] that jointly learn spatiotemporal features from RGB frames using 3D convolutions. However, these two paradigms are both unsatisfactory in terms of efficiency. The two-stream CNNs heavily rely on optical flow as motion representation, while the extraction of optical flow is time-consuming and storage-demanding. And the methods based on 3D CNNs is too expensive to deploy because the 3D convolution kernels require heavy computational cost.

Our PAN follows the 2D two-stream paradigm, but the input modality is only raw RGB video frames. Aiming at fast action recognition, our PAN discards the pre-computed optical flow. Instead, the motion information is distilled by introducing an efficient motion cue PA.

**Motion Representation.** Motion representation serves as an important cue for video-based action recognition. Optical flow is considered as a useful representation of short-term motion, many works [1], [12], [14], [2], [16] have shown that adding optical flow as another input modality can significantly boost the recognition performance. However, conventional optical flow computation approaches [25], [26], [27] computing optical flow in advance and storing that into the disk are absolutely inefficient. In order to alleviate the inefficiency, several recent works speed up the optical flow estimation process by delicately designing some CNN models, such as FlowNet family [17], [18], SpyNet [19] and PWC-Net [20], etc. However, these models aim at accurately estimating optical flow in advance, which is separated from the ultimate action recognition task. Other works [22], [13] employ an encoder-decoder network to reconstruct optical flow and this network can be jointly trained with the subsequent action

recognition network. But the encoder-decoder manner still requires expensive computational cost. Thus, it remains a challenging task to find an efficient and effective motion representation for action recognition. To this end, we decide to replace optical flow with an alternative motion cue for fast action recognition, rather than make optical flow estimation more accurate.

Another line of works make efforts to find auxiliary representations of motion in an end-to-end manner [28], [29]. In this way, various input modalities are carefully designed, such as RGBdiff [12], EMV [11], dynamic image [30] and Displacement Map [31], which can be computed on-the-fly. These works still cannot perform on par with optical flow in terms of action recognition accuracy. Lee *et al.* [32] proposed MFNet, which exploits five fixed directions searching strategy to encode temporal features in a unified manner. Recently, Sun *et al.* [33] introduced Optical Flow Guided Feature (OFF), which obtains spatial and temporal features utilizing Sobel operator and element-wise subtraction respectively with ground-truth optical flow as supervision. It should be noted that our work does not require optical flow during both the training and testing phases.

In this paper, we distill the motion cue PA at the bottom of the network with concise pixel-wise computation. Our PA can be viewed as feature-level pixel-wise variation accumulation, where we encode the output motion as a single channel saliency map reflecting small displacement at movement boundaries, which is of the same spatial resolution as the input RGB frames.

**Temporal Modeling.** 3D CNNs are naturally expert in temporal modeling, as 3D convolutional operators are designed to fuse both spatial and temporal information within local receptive fields. Non-local Networks [3] use self-attention mechanism to capture long-range temporal correlations. Slow-fast Networks [34] contain two pathways that capture categorical semantics and motion semantics at slow and fast frame rate respectively, and lateral connections are employed to fuse the two pathways. As for 2D CNNs, 2D convolutional operators only focus on local spatial regions within single frame, without exchanging information among neighboring frames. Such frame-level features are prone to cause partial observations, thus temporal modeling is necessary. Several works [10], [12], [35], [36] use 2D CNNs to process video frames independently and then obtain video-level features by late fusion strategies. Temporal Segment Networks (TSN) [12] aggregate the frame-level features through average pooling consensus function to obtain video-level representations. But average operation can not infer the temporal order or more complicated temporal relationships. Temporal Relation Networks (TRN) [36] further improve the TSN to utilize temporal relations in videos. Another noteworthy work is Temporal Shift Module (TSM) [37], it enables local temporal fusion among neighboring frames with temporal shift operation.

In this paper, we devise a global temporal fusion strategy VAP with negligible parameters for long-term temporal modeling at multiple timescales.

## III. PERSISTENCE OF APPEARANCE (PA)

To lift the reliance on optical flow, we devise a novel motion cue called Persistence of Appearance (PA). In this section, we first present the theoretical derivation of PA (Sec. III-A). Then we design an efficient PA module to speed up the motion modeling procedure (Sec. III-B) in action recognition task. Afterwards, we investigate the function of PA in motion modeling by exploring two meaningful encoding schemes (Sec. III-C).

### A. Theoretical Derivation of PA

As discussed in Sec. I, the small displacements at motion boundaries matter most for action recognition. Thus, given an adjacent two-frame pair, we make efforts to obtain a saliency map that highlights such small motion variations at boundaries.

For traditional optical flow, the *brightness constancy constraint* is defined as follows:

$$I(x, y, t) \approx I(x + \Delta x, y + \Delta y, t + \Delta t) \qquad (1)$$

where $I(x, y, t)$ denotes the pixel value at the location $(x, y)$ of a video frame at time $t$. As time varies from $t$ to $(t + \Delta t)$, the spatial displacements in horizontal and vertical axis are $\Delta x$ and $\Delta y$ respectively. This constraint formulation assumes that the brightness of a point remains unchanged if it moves from $(x, y)$ at time $t$ to $(x + \Delta x, y + \Delta y)$ at time $(t + \Delta t)$. The optical flow can be estimated by finding the optimal solution $(\Delta x^*, \Delta y^*)$ through optimization methods, and additional constraints, *e.g.*, local smoothness assumption, are also considered for estimating the actual flow.

We extend Eq. 1 to the feature space by replacing the image $I(x, y, t)$ with its $i$-th feature map $F_i(x, y, t)$ after a specific layer:

$$F_i(x, y, t) \approx F_i(x + \Delta x, y + \Delta y, t + \Delta t) \qquad (2)$$

The difference map $D$ between $i$-th feature maps is given as:

$$D_i(x, y, \Delta t) = F_i(x + \Delta x, y + \Delta y, t + \Delta t) - F_i(x, y, t) \quad (3)$$

If we apply the optical flow constraint in feature space, $D$ tends to have lower absolute value. However, searching the neighboring areas to find the optimal solution $(\Delta x^*, \Delta y^*)$ in each location is time-consuming, so we do not use such a complex searching strategy. In contrast to optical flow, we only capture the motion variation at a certain point in feature space without considering the direction of the movement, which perfectly aligns with our idea of modeling the small displacements at motion boundaries because: (1) *small displacements* are perceived since one pixel in low-level feature map contains information of a small receptive field in input space; (2) *motion boundaries* are captured since the first few convolutional layers tend to capture general patterns, *e.g.*, boundaries, texture, etc [23]. So the differences among low-level feature maps will pay more attention to the variations at boundaries. In summary, differences in low-level feature maps
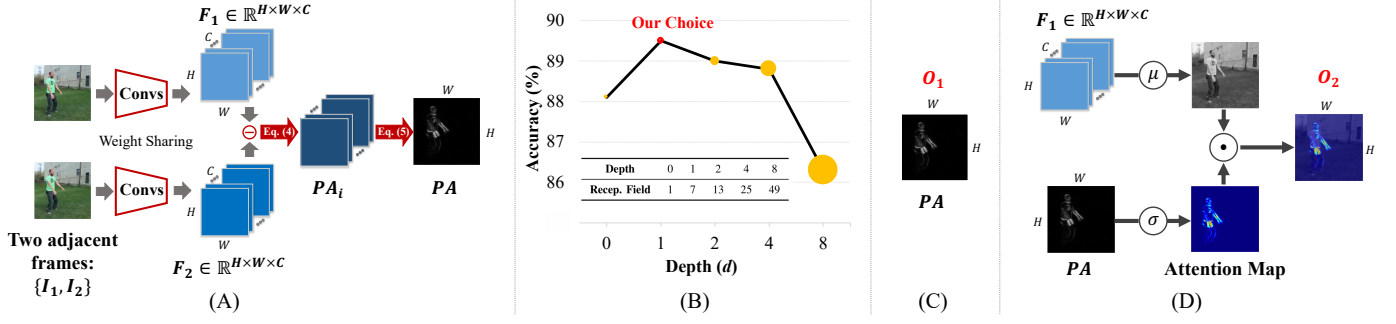
Fig. 2. (A) Illustrations of Persistence of Appearance (PA) design. (B) Depth of conv-layers in "PA Module" *vs.* accuracy. (C) Encoding scheme $e_1$: PA as motion modality. (D) Encoding scheme $e_2$: PA as attention. Here we only provide exemplars that processing two adjacent frames (*i.e.*, $m = 2$) for clarity.

can reflect small displacements of motion boundaries due to convolutional operations.

Therefore, we define the $i$-th PA component as follows:

$$PA_i(p, \Delta t) = F_i(p, t + \Delta t) - F_i(p, t) \qquad (4)$$

where $p = (x, y)$ and $i = 1, \ldots, C$, and $C$ is the channel number. Thus, we can conclude that our PA is highly correlated with optical flow. This definitely provides theoretical support for its effectiveness in modeling motion information.

All the computed $PA_i$ can be further accumulated to 1 channel to manifest the motion magnitude, which can reflect the motion variations at boundaries.

$$PA = \sqrt{\sum_{i=1}^{C} (PA_i(p, \Delta t))^2} \qquad (5)$$

### B. PA Module Design

Since our PA operates in feature space, we need to search for the best depth choice of convolutional layers (conv-layers) to generate feature maps. We define the basic conv-layer as eight $7{\times}7$ convolutions with stride=1 and padding=3, so that the spatial resolutions of the obtained feature maps are not reduced. Assume that $d$ basic conv-layers are sequentially stacked to form the $d$-depth network, we experiment with 5 networks having depth of $d$ equal to 0, 1, 2, 4 and 8. The experimental results on UCF101 split 1 dataset are depicted in Fig. 2-B. The area of the circles indicates the computational cost (FLOPs). We find that directly applying the pixel-wise differences accumulation in input space ($d = 0$) does not perform best. The best performance is achieved when $d = 1$, *i.e.*, only one basic conv-layer is adopted. As the network goes deeper, FLOPs significantly increases and the performance degrades. This is mainly because high-level features with large receptive fields have been highly abstracted and thus may not be able to reflect small motion variations in input images. The experimental results are consistent with our claim that differences in low-level feature maps can reflect small displacements of motion boundaries, which are the most critical ingredients for recognizing actions.

As $d = 1$ performs best, we design a light-weighted "PA module" which only contains single basic conv-layer (eight $7{\times}7$ convolutions) to obtain low-level features and several

computing operations based on Eq. 4 and Eq. 5. This module performs low-level representations comparison pixel by pixel between two adjacent frames, and outputs one saliency map (PA) reflecting small displacements of motion boundaries for further processing. This module is located at the bottom of our network, as shown in Fig. 3-B&3-C (detailed architectural information will be given in Sec. IV-A).

Formally, as shown in Fig. 2-A, given two adjacent frames $\in \mathbb{R}^{H \times W \times 3}$ with $H$, $W$ and 3 being their height, width and channel number. First, low-level feature maps $F_1, F_2 \in \mathbb{R}^{H \times W \times C}$ are obtained without spatial resolution reduction. Then, the pixel-wise value difference is computed between the two feature maps with the same index $i$ (See in Eq. 4). Finally, all the computed $PA_i$ are accumulated to 1 channel based on Eq. 5, so the result PA $\in \mathbb{R}^{H \times W}$ is two-dimensional. Therefore, in "PA module", a mapping $\mathbb{R}^{H \times W \times 3} \to \mathbb{R}^{H \times W}$ is established from the appearance to the dynamic motion.

### C. Encoding Schemes

As elaborated above, PA is a concise motion cue focusing on the small displacements of motion boundaries between two adjacent frames. We would also like to understand the practical function of PA in motion modeling. Intuitively, PA can serve as either auxiliary input modality or spatial attention map. So here we explore two meaningful encoding schemes: PA as motion modality *vs.* attention map. Given $m$ adjacent frames set $\{I^{(i)}\}_{i=1}^{m}$, the corresponding low-level feature maps in PA module are defined as $\{F^{(i)}\}_{i=1}^{m}$, and each two adjacent frames are processed to obtain total $(m-1)$ PA: $\{PA^{(i)}\}_{i=1}^{m-1}$. Assuming that the input modality to the subsequent backbone network is $O$, so in this subsection, we will discuss two encoding schemes $e_1$, $e_2$ that carry out the mapping procedure $e_i : PA \to O$, *i.e.*, aggregating $PA$ to $O$.

*1) PA as motion modality.* This is the most straightforward scheme to directly exploit motion information contained in PA. Generally, for action recognition methods, taking the stacked optical flow as input to capture motion information can significantly boost the performance. Since PA also has the capability of describing the pixel-level apparent motion information between two continuous RGB frames, we use stacked PA as input modality as shown in Fig. 2-C. This scheme can be represented as:

$$O_1 = e_1(PA) = \overset{m-1}{\underset{i=1}{\Upsilon}} \left( PA^{(i)} \right) \tag{6}$$

Here, we define $\overset{m-1}{\underset{i=1}{\Upsilon}} (\cdot)$ as the *cumulative channel concatenation function* that chronologically concatenate the input tensor along the channel dimension. Thus, if the input tensor $PA^{(i)} \in \mathbb{R}^{H \times W \times 1}$, then the output tensor $O_1 \in \mathbb{R}^{H \times W \times (m-1)}$.

*2) PA as attention.* Human perception researches [38], [39] suggest that instantaneous motion can attract attention. Recent video analysis works benefit a lot under the attention guidance of motion captured by optical flow, such as video salient object detection [40], video captioning [41], etc. Motivated by this, we attempt to exploit motion information in PA to emphasize some important regions in appearance feature maps, as shown in Fig. 2-D. This PA-guided spatial attention scheme is defined as follows, we employ PA with a sigmoid activation to attend the corresponding mean feature map:

$$O_2 = e_2(PA) = \overset{m-1}{\underset{i=1}{\Upsilon}} \left( \sigma(PA^{(i)}) \odot \mu(F^{(i)}) \right) \tag{7}$$

where $\sigma(\cdot)$ is a sigmoid function and $\mu(\cdot)$ returns the mean value of the input feature maps along the channel dimension. $\odot$ denotes element-wise multiplication, so if the input tensor $PA^{(i)} \in \mathbb{R}^{H \times W \times 1}$ and $F^{(i)} \in \mathbb{R}^{H \times W \times C}$, then $\mu(F^{(i)}) \in \mathbb{R}^{H \times W \times 1}$ and $O_2 \in \mathbb{R}^{H \times W \times (m-1)}$.

***Which encoding scheme is better?*** We compare the performance of the PA module using these two encoding schemes in the aspect of their runtime efficiency and action recognition accuracy on UCF101 split 1 dataset. The results are shown in Table I. To measure the efficiency, we consider computational cost (FLOPs), the number of parameters (#Param) and inference speed (Speed) of the PA module. To evaluate the performance of these two encoding schemes on action recognition task, we follow the TSN manner: firstly frames are sampled from evenly divided video segments, then these frames are fed into the PA module and backbone CNN (ResNet-50) sequentially, finally the output activations are averaged as the final prediction scores. More implementation details are in the supplementary material.

The results in Table I clearly indicate that encoding scheme $e_1$, directly exploiting the motion information contained in PA, performs better. It has not only fewer FLOPs but also higher runtime speed and superior recognition performance. The number of parameters of the two encoding schemes are the same, because the 1.184K parameters are completely from PA module, and the subsequent encoding procedure does not introduce any extra learnable parameters. The more FLOPs and lower speed of $e_2$ is mainly caused by the sigmoid function and element-wise multiplication. Notably, $e_2$ also degrades the accuracy by 1.5%. We hypothesize that when using appearance-dominant features (*i.e.*, appearance feature maps) as input, the features must secure integral regions of appearance to represent the semantic information for the video category. However, for $e_2$, attending appearance feature maps with PA will highlight the motion boundaries, leading to the imbalanced appearance responses both inside and at the

TABLE I
PA AS MOTION CUE *vs.* PA AS ATTENTION. ACCURACIES ARE EVALUATED ON UCF101 SPLIT 1 WITH THE SAME NETWORK SETTINGS.

| Encoding Schemes | Efficiency Metrics | | | Accuracy |
| --- | --- | --- | --- | --- |
| | FLOPs | #Param | Speed | |
| $e_1$: PA as motion cue | **2.868G** | **1.184K** | **8196fps** | **89.5%** |
| $e_2$: PA as attention | 2.884G | 1.184K | 6752fps | 88.0% |

boundaries of the moving objects, thus $e_2$ is limited in terms of such integrity. Encoding scheme $e_1$, on the contrary, only utilizes motion-dominant features (*i.e.*, PA), so there is no need to consider the appearance integrity.

Therefore, based on the above observations, we employ $e_1$ (*i.e.*, directly exploit PA as input motion modality) as the default encoding scheme in our paper.

## IV. PERSISTENT APPEARANCE NETWORK (PAN)

Our primary objective in this paper is to realize fast action recognition in real-time scenarios, so we propose the Persistent Appearance Network (PAN). In this section, we first give an architectural overview of our PAN framework (Sec. IV-A). Then we introduce our VAP method, a temporal feature aggregation strategy, applied within PAN framework. It assists learning long-term video-level representations by integrating information across various timescales (Sec. IV-B).

### A. $PAN_{Full}$ and $PAN_{Lite}$

The two network variants of our proposed Persistent Appearance Network (PAN) is shown in Fig. 3-B&3-C, namely $PAN_{Full}$ and $PAN_{Lite}$ respectively. They have identical sampling strategy but model the spatiotemporal features in different ways. Their sampling strategy is shown in Fig. 3-A: the input video sequence $V$ is firstly divided into $N$ segments with equal length $\{S_t\}_{t=1}^N$. And $m$ adjacent frames are randomly chosen from each segment as a "$m$-frame stack": $\{I_{t[m]}\}_{t=1}^N$, the first frames of each "$m$-frame stack" are denoted as $\{I_t\}_{t=1}^N$. The main difference between both variants are analyzed below.

*1) $PAN_{Full}$: separate and accurate.* As illustrated in Fig. 3-B, this network is composed of dual branches: *RGB branch* and *PA branch*, capturing the spatial and temporal features separately. The RGB branch encodes the spatial appearance information. It takes the selected $N$ frames $\{I_t\}_{t=1}^N$ as input and processes them to obtain frame-level features through the backbone network $\mathcal{H}_B$ (blue blocks). Then these obtained features are further aggregated as video-level features using VAP module $\mathcal{H}_{VAP}$. Mathematically, the output spatial features $y_s$ of RGB branch can be written as:

$$y_s = \mathcal{H}_{VAP}(\mathcal{H}_B(\{I_t\}_{t=1}^N)) \tag{8}$$

The other PA branch distills apparent motion information purely from adjacent RGB frames (*i.e.*, appearance information). Firstly, $N$ stacks of $m$ adjacent frames $\{I_{t[m]}\}_{t=1}^N$ are transformed to motion cue PA after the PA module $\mathcal{H}_{PA}$. Then
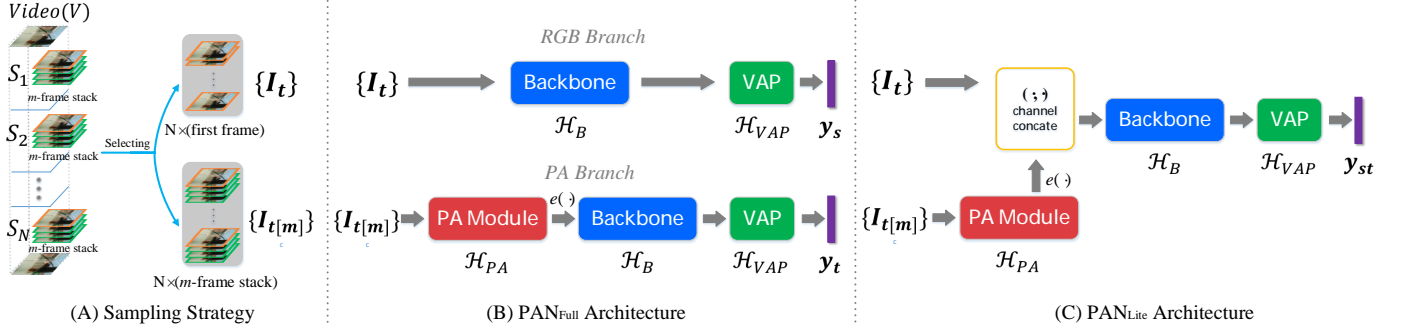
Fig. 3. The overall architecture of Persistent Appearance Network (PAN). It has two network variants: (B) PAN$_{\text{Full}}$ - "*divide and conquer*", *i.e.*, capturing the spatial and temporal semantics separately; (C) PAN$_{\text{Lite}}$ - "*unified and efficient*", *i.e.*, extracting the spatial and temporal semantics simultaneously.

the temporal features $\boldsymbol{y_t}$ are obtained through the backbone network $\mathcal{H}_B$ and VAP module $\mathcal{H}_{VAP}$:

$$\boldsymbol{y_t} = \mathcal{H}_{VAP}(\mathcal{H}_B(e(\mathcal{H}_{PA}(\{\boldsymbol{I_{t[m]}}\}_{t=1}^N)))) \quad (9)$$

where $e(\cdot)$ is the encoding scheme described in Sec. III-C. Following the common practice [12], [36], [37], we merge the branch-level scores to obtain final prediction of the whole video through score fusion strategy, *i.e.*, calculating the weighted average scores from the two branches.

*2) PAN$_{Lite}$: unified and light-weighted.* As illustrated in Fig. 3-C, our unified network PAN$_{Lite}$ stacks RGB and PA together and feeds them through the backbone network, allowing the network to decide itself how to extract the spatiotemporal information. Given $N$ sampled $m$-frame stacks $\{\boldsymbol{I_{t[m]}}\}_{t=1}^N$ and their first frames $\{\boldsymbol{I_t}\}_{t=1}^N$, the output $\boldsymbol{y_{st}}$ can be written as:

$$\boldsymbol{y_{st}} = \mathcal{H}_{VAP}(\mathcal{H}_B(\boldsymbol{I_t}, e(\mathcal{H}_{PA}(\{\boldsymbol{I_{t[m]}}\}_{t=1}^N)))) \quad (10)$$

where $(\cdot, \cdot)$ is the channel concatenation operation and $\mathcal{H}_{PA}$, $\mathcal{H}_B$ and $\mathcal{H}_{VAP}$ indicate PA module, backbone network and VAP module, respectively.

### B. Various-timescale Aggregation Pooling (VAP)

Long-term temporal modeling is of great importance for the video understanding task as discussed in Sec. II. In this paper, we devise a temporal fusion strategy called Various-timescale Aggregation Pooling (VAP), which adaptively emphasizes expressive features and suppresses less informative ones by observing global information across various timescales. As shown in Fig. 3-B&3-C, VAP module is adopted at the top of each network. Overall, VAP contains two main steps as illustrated below.

**(A) Specific-timescale Pooling**. As shown in the top part of Fig. 4, assume that the original video is divided into $N$ segments, and the sampled video frames are fed into the backbone 2D CNN to generate $d$-dim feature vectors $\boldsymbol{f} : \{f_1, f_2, \ldots, f_N\}$, where $f_i \in \mathbb{R}^d$. In order to temporally integrate these $N$ features without overlapped scope, we adopt dilated max pooling over the time dimension, where the dilation rate controls the spacing between the kernel points. For brevity, this pooling function can be expressed
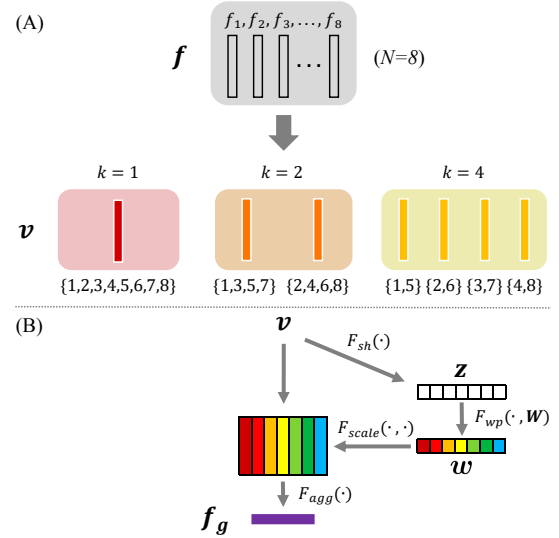


Fig. 4. Various-timescale Aggregation Pooling (VAP). Top part: (A) Specific-timescale Pooling. Bottom part: (B) Various-timescale Aggregation. The numbers in curly brackets indicate which $\boldsymbol{f}$ are involved in generating the timescale-wise features $\boldsymbol{v}$ (Eq. 11). *Best viewed in color and zoomed in.*

as $maxpool_{(ks,st,dr)}\{activations\}$, where $ks$, $st$, $dr$ refer to kernel size, stride and dilation rate respectively. Accordingly, the $k$-timescale pooling is defined as follows:

$$\boldsymbol{v_k} = maxpool_{(\frac{N}{k},1,k)}\{f_1, f_2, \ldots, f_N\} \quad (11)$$

where $k$ is a positive integer *s.t.* $\frac{N}{k} \in Z$. By convention, we use pyramidal timescale settings (*i.e.*, $k = 2^0, 2^1, \ldots, 2^{log_2(N)-1}$). After Eq.11, the time span changes from $N$ to $k$. So $2^0 + 2^1 + \ldots + 2^{log_2(N)-1} = N-1$ timescale-level features are obtained in total (*i.e.*, $\boldsymbol{v} \in \mathbb{R}^{(N-1) \times d}$).

**(B) Various-timescale Aggregation**. Our basic idea is to fuse temporal information at each timescale by weighted timescale-wise aggregation. As shown in the bottom part of Fig. 4, given that $\boldsymbol{v} \in \mathbb{R}^{T \times d}$ represents the total $T$ pooled features at different timescales, we first *shrink* global spatial semantics in each feature into a temporal descriptor reflecting the corresponding timescale-wise statistics. Thus, the output $\boldsymbol{z} \in \mathbb{R}^{T \times 1}$ can be expressed by:

$$z = F_{sh}(\boldsymbol{v}) = \frac{1}{d}\sum_{i=1}^{d}\boldsymbol{v}(i) \qquad (12)$$

To capture the cross-timescale interdependencies, we adopt a concise nonlinearity learning mechanism with a softmax activation for *weight perception*:

$$\boldsymbol{w} = F_{wp}(\boldsymbol{z},\boldsymbol{W}) = softmax(\boldsymbol{W_2}(\delta(\boldsymbol{W_1 z}))) \qquad (13)$$

where $\delta$ denotes the ReLU function, $\boldsymbol{W_1} \in \mathbb{R}^{\alpha T \times T}$ and $\boldsymbol{W_2} \in \mathbb{R}^{T \times \alpha T}$ ($\alpha$ is the expansion ratio) are the learnable parameters of two fully-connected (FC) layers. The output of $F_{wp}$ function is the weight vector $\boldsymbol{w} \in \mathbb{R}^{T \times 1}$.

The final global video-level representation $\boldsymbol{f_g} \in \mathbb{R}^{d}$ of the proposed VAP is obtained by firstly rescaling $\boldsymbol{v}$ with the weight vector $\boldsymbol{w}$ and then aggregate the recalibrated features along $T$ dimension.

$$\boldsymbol{f_g} = F_{agg}(F_{scale}(\boldsymbol{w},\boldsymbol{v})) = sum(\boldsymbol{wv}) \qquad (14)$$

In particular, for action recognition task, the prediction scores $\boldsymbol{s}$ are obtained by:

$$\boldsymbol{s} = F_{pred}(\boldsymbol{f_g},\boldsymbol{W}) = \boldsymbol{W_3 f_g} \qquad (15)$$

where $\boldsymbol{W_3} \in \mathbb{R}^{c \times d}$, $c$ is the number of classes.

## V. Experiments

In this section, we first introduce the evaluation datasets and implementation details. Then in ablation studies, we investigate the importance of our proposed motion cue PA and various-timescale aggregation strategy VAP for real-time action recognition. During this investigation, we also explore some basic settings. Extensive results show the superior performance achieved by PAN compared with baselines and other state-of-the-art methods, on both temporal-dominant datasets and scene-dominant datasets. Finally, we visualize the proposed motion cue PA to qualitatively justify its superiority in motion modeling and discuss the future work.

### A. Experimental Settings

**Datasets.** We evaluate our approach on six challenging benchmarks for action recognition. They can be grouped into two categories: (a) ***Temporal-Dominant Datasets***, including Something-Something-V1 & V2 [4] and Jester [6]. Recognizing actions in these datasets requires strong temporal modeling ability, as many action classes are symmetrical, *e.g.*, "Moving something up" and "Moving something down". Something-Something datasets (2 released versions) include 174 categories with 86,017(V1)/168,913(V2) training videos, 11,522(V1)/24,777(V2) validation videos, and 10,960(V1)/27,157(V2) test videos. Jester contains 27 human hand gestures with 148,092 videos. (b) ***Scene-Dominant Datasets***, including Kinetics400 [5], UCF101 [42] and HMDB51 [43]. RGB scene information in these datasets is more critical than temporal relations for action recognition. Kinetics400 is a large-scale action recognition benchmark

including ∼300k videos with 400 human action classes. The performances are evaluated with the top-1 and top-5 accuracies. The UCF101 dataset includes 13,320 video clips with 101 action classes and the HMDB51 dataset contains 6,766 videos with 51 action categories. For these two datasets, the mean class accuracy over the three official splits is calculated as the final result.

**Input & Backbone.** The input modality of our PAN is only raw RGB frames. We set $N = 8$ and $m = 4$ as default, thus 8 "4-frame stack" are sampled from a video. For PAN$_{Full}$, the first frame in each "4-frame stack" will be selected and fed into RGB branch, meanwhile, all the frames will be fed into the PA branch. For PAN$_{Lite}$, we take all the 8 sampled "4-frame stack" as input. For comparative purposes, we use ResNet-50 as the default backbone like other state-of-the-art methods [44], [37], [34]. Since TSM [37] can exchange information between neighboring frames at zero cost, we opt to inject TSM into our backbone to enable *local temporal fusion*. Note that TSM module is orthogonal to our PA (efficient motion representation) and VAP (*global temporal fusion* strategy).

**Training & Testing.** For relatively large-scale datasets (Something-Something-V1 & V2, Jester and Kinetics400), the training parameters of PAN$_{Lite}$ and the PA branch of PAN$_{Full}$ are: initial learning rate 0.01 (total 80 epochs, divided by 10 after every 30 epochs). As for RGB branch of PAN$_{Full}$, we also set initial learning rate as 0.01 (but total 50 epochs, decreases at epoch 20 & 40). We train our PAN using SGD algorithm, with weight decay 1e-4, mini-batch size 64. ImageNet pre-trained weights are employed for initialization. While for small-scale datasets (UCF101 and HMDB51) that are easy to over-fit, we use Kinetics400 for pre-training and scale the training epochs by half. During testing procedure, we report "view" (spatial crops × temporal clips) used in [34]. For temporal-dominant datasets, we only use 1 view (1 center spatial '224×224' crop × 1 temporal clip) unless otherwise specified. While for scene-dominant datasets, we use 2 views (full resolution image with shorter side 256 pixels × 2 temporal clips), which is much more efficient than the common practice (30 views) in [3], [34], [45]. Especially, for PAN$_{Full}$, we average the class prediction scores of the RGB and PA branches.

### B. Ablation Studies for PA

Following the common practice [12], [10], [46], all experiments in this ablation study are conducted on UCF101 split 1. In this subsection, we prove that our designed PA is a significant motion representation by performing in-depth studies to answer the following two questions:

***Q1: Is PA efficient, effective and flexible enough?*** As mentioned in Sec. I, our main target in this paper is to design a motion representation alternative to optical flow with the merits of efficient, effective and flexible. We compare the PA with other mainstream optical flow extraction methods from the perspective of efficiency and effectiveness for action recognition task, including conventional optical flow [26] and CNN-based estimated optical flow [17], [18] methods. The comparison results are listed in Table II. Note that all the

TABLE II
COMPARISON RESULTS OF PA WITH MAINSTREAM OPTICAL FLOW
COMPUTATION METHODS. ACCURACIES ARE EVALUATED ON UCF101
SPLIT 1 WITH THE SAME NETWORK SETTINGS EXCEPT THE INPUT
MODALITY.

| Motion Rep. Method | Efficiency Metrics | | | Accuracy |
| | FLOPs | #Param | Speed | |
|---|---|---|---|---|
| TV-L1 [26] | - | - | 8fps | 88.2% |
| FlowNetS [17] | 356G | 38.7M | 204fps | 86.8% |
| FlowNetC [17] | 444G | 39.2M | 151fps | 87.3% |
| FlowNet2.0 [18] | 2019G | 162.5M | 25fps | 87.7% |
| **PA (Ours)** | **2.868G** | **1.184K** | **8196fps** | **89.5%** |

TABLE III
FLEXIBILITY VALIDATION OF PA. OUR MOTION CUE PA CONSIDERABLY
IMPROVES ALL THE BASELINES ON UCF101 SPLIT 1. ALL THE
EXPERIMENTS ARE CONDUCTED USING THE SAME NETWORK SETTINGS.

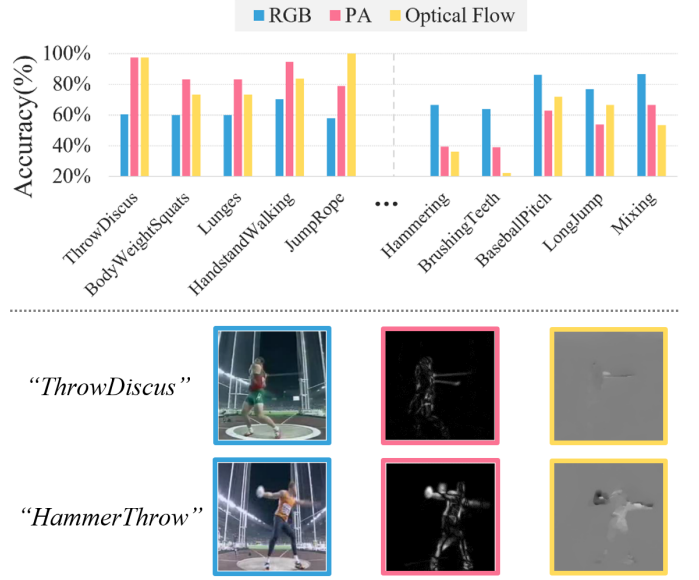| Baseline | Pre-train | Modality | Acc. | ∆Acc. |
|---|---|---|---|---|
| BN-Inception [47] | ImageNet | RGB | 85.6% | **+5.7%** |
| | | RGB+PA | **91.3%** | |
| ECO [48] | Kinetics | RGB | 90.4% | **+4.7%** |
| | | RGB+PA | **95.1%** | |
| 3D ResNet-18 [46] | Kinetics | RGB | 83.9% | **+2.9%** |
| | | RGB+PA | **86.8%** | |



Fig. 5. Top part: The top 5 action classes that show the greatest difference in terms of accuracy between the input modality RGB and PA, and the performance of optical flow in these classes is also plotted for comparison. Bottom part: Input modality visualization of the two classes that are most easily misclassified by RGB. *Best viewed in color and zoomed in.*

network settings are the same and more details can be found in the supplementary material. Surprisingly, compared with conventional optical flow method TV-L1 [26], our proposed PA achieves over $1000\times$ faster speed (8196fps *vs.* 8fps) as well as 1.3% higher action recognition accuracy. Furthermore, PA consistently outperforms all the seminal CNN-based optical flow estimation methods [17], [18], suggesting that PA has fast (8196fps) and effective (89.5%) motion modeling ability.

To demonstrate the flexibility of our proposed PA, we deploy different backbone networks and compare their results. Here, we choose three widely adopted backbone networks for deep action recognition, including BN-Inception [47] (2D CNN), ECO [48] (mixed 2D-3D CNN) and 3D-ResNet [46] (3D CNN). Experimental results on UCF101 split 1 are tabulated in Table III. The results show that the selected backbone networks all significantly benefit from our PA with considerable accuracy improvements, demonstrating the flexibility of our proposed PA.

Therefore, compared with optical flow, our proposed PA is much faster to compute, superior in performance. And it is also flexible to implement.

***Q2: Can PA represents apparent motion?*** A key intuition for designing PA is that it can capture apparent motion by focusing more on moving boundaries without the heavy precomputation of optical flow. In order to investigate whether our PA performs similarly to the optical flow, we compare the performance using three different input modalities: RGB, PA and optical flow. With each modality as input, we train the network first on UCF101 split 1, then we test the trained models to get the prediction scores for all the 101 classes.

Finally, we plot the top 5 classes that show the largest differences in recognition accuracy between RGB and PA. For reference, the performance of optical flow in these classes is also depicted, as shown in the top part of Fig. 5. In this figure, we can clearly see that when PA outperforms RGB, the optical flow is also superior to RGB and vice versa. This indicates that PA can help the network learn patterns different from RGB but similar to optical flow.

After analyzing the recognition results, we find that most "Throw Discus" videos are misclassified as "Hammer Throw" when using just RGB modality. To demonstrate that our PA can represent apparent motion, we visualize the three modalities of these two classes in the bottom part of Fig. 5. In RGB images, it is unclear whether the athletes are holding the discus or the hammer due to the interference of background, illumination, etc, so distinguishing the two classes using only RGB modality may cause confusion. Encouragingly, our PA is able to highlight the moving objects (in this case human and the object in hands) as the optical flow does, which is crucial for differentiating between videos that look much alike.

Therefore, this result is in accordance with our idea that the PA, focusing on capturing small displacements at moving boundaries, can be used as apparent motion representation. See Sec. V-E for more visualization results.

### C. Ablation Studies for VAP

As introduced in Sec. IV-B, our proposed VAP is a temporal aggregation strategy that can learn to exploit global information across different timescales. In this subsection, we conduct ablation experiments to gain more insights about the effect of aggregating semantics across various timescales.
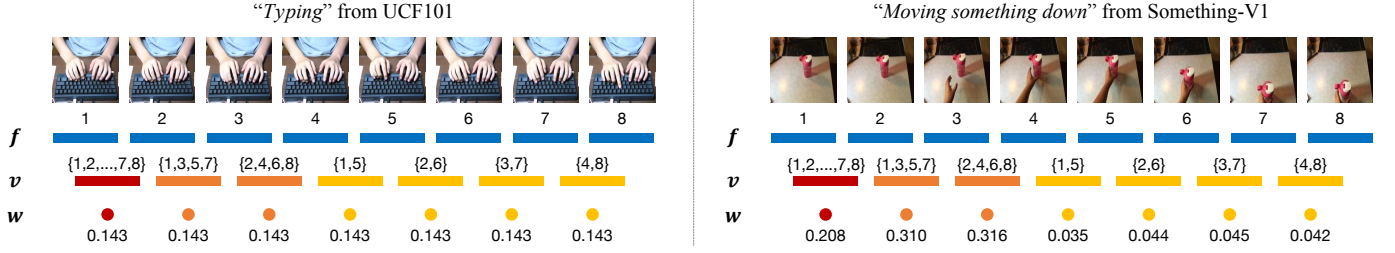
Fig. 6. Visualization of two selected video sequences and their timescale-wise weights $w$ calculated by Eq. 13. The sampled frames are processed by the backbone individually to obtain frame-level features $f$. The numbers in curly brackets indicate which $f$ are involved in generating the timescale-wise features $v$ (Eq. 11). Different colors (red, orange & yellow) are employed to distinguish different timescales. *Best viewed in color and zoomed in.*

TABLE IV
EXPLORATION OF DIFFERENT EXPANSION RATIO $\alpha$ AND COMPARISON AMONG DIFFERENT TEMPORAL AGGREGATION STRATEGIES.

| Aggregation Strategy | Accuracy |
|---|---|
| Avg Pooling (baseline) | 86.5% |
| VIP [24] (Our prev.) | 87.9% |
| VAP($\alpha = 1$) | 88.4% |
| VAP($\alpha = 2$) | 88.3% |
| **VAP($\alpha = 4$)** | **88.5%** |
| VAP($\alpha = 8$) | 88.4% |

***Q1: Is various-timescale aggregation effective?*** Following the common practice [12], [10], [46], we perform ablation studies on UCF101 split 1 dataset using different temporal aggregation strategies. $N = 8$ frames are sampled as input and ResNet-50 is used as backbone network. The results are shown in Table IV. Avg Pooling here is referred to the conventional consensus function introduced in [12] that averages the output logits to get the final prediction. First, we explore the impact of the expansion ratio $\alpha \in \{1, 2, 4, 8\}$. The accuracies between various expansion ratios exhibit stable performances, consistently outperforming the Avg Pooling baseline. This indicates that VAP is not sensitive to the hyperparameter setting $\alpha$. As $\alpha = 4$ achieves the best performance, we choose 4 as the default expansion ratio in VAP. Second, VAP is considered as an enhanced version of our previous work VIP [24]. Different from VIP that exploits preset weights for inference, VAP learns the weights by adaptively using global information at different timescales. We also do comparison between these two variants. VAP achieves higher recognition accuracy than VIP (88.5% *vs.* 87.9%), showing the effectiveness of the enhancement in this paper. Furthermore, VAP and VIP utilizing various-timescale semantics consistently outperform baseline by 2.0% and 1.4% respectively, indicating the significance of long-term temporal modeling through various-timescale aggregation.

***Q2: What does the VAP learn?*** We expect VAP module to let the network capture various-timescale interdependencies. To verify that this is indeed achieved, we output the weight results $w$ learned from the weight perception Eq. 13. Here we consider both scene-dominant dataset (UCF101) and temporal-dominant dataset (Something-Something-V1). $N = 8$ frames are sampled as input, thus $k \in \{1, 2, 3\}$ and total 7 timescale-

level features $v$ are calculated according to Eq. 11. ResNet-50 is used as backbone and VAP module is adopted at the top of the network. We first train this network on the two datasets respectively, then we test the two trained models separately to obtain the corresponding weight values $w$ of different timescales. Fig. 6 shows two selected video sequences and the timescale-wise weights $w$ with respect to the two action classes: "*Typing*" from UCF101 and "*Moving something down*" from Something-Something-V1 dataset. More visualization results can be found in the supplementary material. We can observe that the weight changes for different videos at various timescales, suggesting that VAP can adaptively guide the network to emphasize some expressive features while suppressing other less informative ones. Moreover, we can clearly see that VAP trained on UCF101 dataset tends to learn similar weights for all timescales, while the weights varies widely on Something-Something-V1. We speculate that this has to do with the fact that temporal relationships in Something-Something-V1 is more crucial than that in UCF101 for recognition [36], [37], [49], so VAP module trained on Something-Something-V1 prefers treating all timescales differently. This observation further proves that various-timescale aggregation is an important ingredient for video temporal semantics modeling.

### D. Comparison with the State-of-the-Arts

To validate the efficacy of our approach, we compare PAN with newly published state-of-the-art methods on six datasets. Since our proposed PAN focuses on temporal semantics modeling, we mainly analyze the results on temporal-dominant datasets (Something-Something-V1 & V2 and Jester). We also evaluate PAN on various scene-dominant datasets (Kinetics400, UCF101 and HMDB51) to show its consistent strong performance.

#### 1) Temporal-Dominant Datasets.

**Something-Something-V1 & V2 and Jester**. Table V shows the comparison results with the state-of-the-arts on Something-Something-V1 & V2 and Jester datasets. We group the listed methods into two sections (separated by solid line) according to their backbone types: 3D CNN based methods [44], [48] and 2D CNN based methods [12], [36], [37]. It is clear that with complex 3D convolutional operations, the FLOPs of I3D [44] and ECO [48] are much higher than the 2D CNN based methods. Our PAN surpasses all the 3D CNN

TABLE V

COMPARISON RESULTS OF PAN WITH OTHER STATE-OF-THE-ART METHODS ON SOMETHING-SOMETHING V1 & V2 AND JESTER DATASETS.

| Methods | Backbone | Flow? | #Frame | FLOPs×views | Something-V1 | | Something-V2 | | Jester | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | Val Top1 | Val Top5 | Val Top1 | Val Top5 | Val Top1 | Val Top5 |
| I3D [44] | 3DResNet-50 | | 32×2 | 153G×2 | 41.6 | 72.2 | - | - | - | - |
| NL I3D [44] | 3DResNet-50 | | | 168G×2 | 44.4 | 76.0 | - | - | - | - |
| NL I3D + GCN [44] | 3DResNet-50+GCN | | | 303G×2 | 46.1 | 76.8 | - | - | - | - |
| ECO$_\text{En}$Lite [48] | BNInception | | 92 | 267G×1 | 46.4 | - | - | - | - | - |
| ECO$_\text{En}$Lite$_\text{RGB+Flow}$ [48] | +3DResNet-18 | ✓ | 92+92×6 † | N/A | 49.5 | - | - | - | - | - |
| TSN$_\text{8F}$ [12] | BNInception | | 8 | 16G×1 | 19.5 | - | - | - | - | - |
| | ResNet-50 | | 8 | 33G×1 | 19.7 | 46.6 | 27.8 | - | 81.0 | 99.0 |
| TRN-Multiscale [36] | BNInception | | 8 | 16G×1 | 34.4 | - | 48.8 | 77.6 | 95.3 | - |
| | ResNet-50 | | 8 | 33G×1 | 38.9 | 68.1 | - | - | - | - |
| TRN$_\text{RGB+Flow}$ [36] | BNInception | ✓ | 8+8×6 † | N/A | 42.0 | - | 55.5 | 83.1 | - | - |
| TSM$_\text{8F}$ [37] | | | 8 | 33G×1 | 45.6 | 74.2 | 59.1 | - | 94.4 | 99.7 |
| TSM$_\text{16F}$ [37] | ResNet-50 | | 16 | 65G×1 | 47.2 | 77.1 | 63.4 | 88.5 | 95.3 | 99.8 |
| TSM$_\text{RGB+Flow}$ [37] | | ✓ | 16+16×6 † | N/A | 52.6 | 81.9 | 66.0 | 90.5 | - | - |
| TEA$_\text{8F}$ [45] | ResNet-50 | | 8 | 35G×30 | 51.7 | 80.5 | - | - | - | - |
| TEA$_\text{16F}$ [45] | | | 16 | 70G×30 | 52.3 | 81.9 | - | - | - | - |
| PAN$_\text{Lite}$ (Ours) | | | 8+8×4 † | 35.7G×1 | 48.0 | 76.1 | 60.8 | 86.7 | 96.2 | 99.8 |
| PAN$_\text{Full}$ (Ours) | ResNet-50+TSM | | 8+8×4 † | 67.7G×1 | 50.5 | 79.2 | 63.8 | 88.6 | 96.6 | 99.8 |
| PAN$_\text{En}$ (Ours) | | | (8+8×4)×2 | (46.6G+88.4G)×2 | 53.4 | 81.1 | 66.2 | 90.1 | 97.2 | 99.9 |
| PAN$_\text{En}$ (Ours) | ResNet-101+TSM | | (8+8×4)×2 | (85.6G+166.1G)×2 | **55.3** | **82.8** | **66.5** | **90.6** | **97.4** | **99.9** |

† The flow streams of [48], [36], [37] take 10-channel (from "6-frame stack") optical flow as input, while our PAN only uses "4-frame stack".

based methods in both efficiency and accuracy aspects. For example, compared with NL I3D+GCN [44], our PAN$_\text{Lite}$ achieves 1.9% higher top-1 accuracy (48.0% *vs.* 46.1% on Something-V1) while with only ∼6% computational cost (35.7G *vs.* 303G×2). As for 2D CNN based methods, the performance of the baseline TSN [12] is relatively inferior than other methods, indicating the significance of temporal modeling for these datasets. Compared with the efficient (33G FLOPs) baselines TSN$_\text{8F}$ [12] and TSM$_\text{8F}$ [37], our PAN$_\text{Lite}$ achieves much higher top-1 accuracy (48.0% *vs.* 19.7%/45.6% on Something-V1, 60.8% *vs.* 27.8%/59.1% on Something-V2, 96.2% *vs.* 81.0%/94.4% on Jester) with only slight extra cost (0.08× FLOPs). With dual-path structure that separately processes RGB and PA modalities, PAN$_\text{Full}$ further bolsters the action recognition performance (50.5% on Something-V1, 63.8% on Something-V2 and 96.6% on Jester). The consistent improvements of our PAN over the other state-of-the-art methods strongly justify the superiority of our proposed PA and VAP for temporal modeling.

Following [48], [37], we also average the class prediction results of PAN$_\text{Lite}$ and PAN$_\text{Full}$ to form our ensemble version PAN$_\text{En}$. Here, two views (1 full-resolution with 256 shorter size × 2 temporal clips) are used for inference. Compared with the optical flow based methods TRN$_\text{RGB+Flow}$ (42.0%), ECO$_\text{En}$Lite$_\text{RGB+Flow}$ (49.5%) and TSM$_\text{RGB+Flow}$ (52.6%), our PAN$_\text{En}$ (53.4%) provides +11.4%, +3.9% and +0.8% top-1 accuracy improvements on Something-V1, respectively. As listed in Table II, the mainstream optical flow extraction methods are computationally intensive, even the most light-weighted FlowNetS model needs 356G FLOPs and extra storage. Taking this cost into account, the optical flow based methods are expensive in both time and space. Surprisingly,

the total cost of our PAN$_\text{En}$ model is only ∼270G. This observation confirms that our proposed PAN, with an efficient motion cue PA, can effectively accelerate the video representation learning process by lifting the reliance on optical flow.

Furthermore, when exploiting a deeper backbone ResNet-101, our PAN brings the current state-of-the-art results to a whole new level (55.3% on Something-V1, 66.5% on Something-V2, 97.4% on Jester). Notably, our method only takes RGB frames as input and thus is free of optical-flow pre-computation.

### 2) Scene-Dominant Datasets.

**Kinetics400, UCF101 and HMDB51**. We also compare the performance of our PAN with other recent state-of-the-art methods on three scene-dominant datasets: Kinetics400, UCF101 and HMDB51. The results are summarized in Table VI. Our method achieves very competitive performance on these datasets, where most actions can be classified by a single frame (*e.g.*, "Typing" action shown in Fig. 6). Our PAN outperforms most of the heavy 3D CNN based architectures (first part in Table VI) [50], [48], [16] using 2D CNN as backbone, indicating that it can enhance the traditional 2D CNN with low cost. For 2D CNN based methods (second part), when pursuing high accuracy, they usually sample many views (spatial crops × temporal clips) per video for inference, thus leading to expensive computational cost. For example, TSM [37] and TEA [45] sample 30 views, their FLOPs are even higher than that of 3D CNN based architectures. In contrast, we only use 2 views for efficiency concerns (detailed in Sec. V-A) and achieve superior performance. For example, compared with baseline TSM [37], our PAN$_\text{En}$ achieves +1.2%, +1.3% and +3.8% top-1 accuracy improvements on Kinetics400, UCF101 and HMDB51 respectively with only
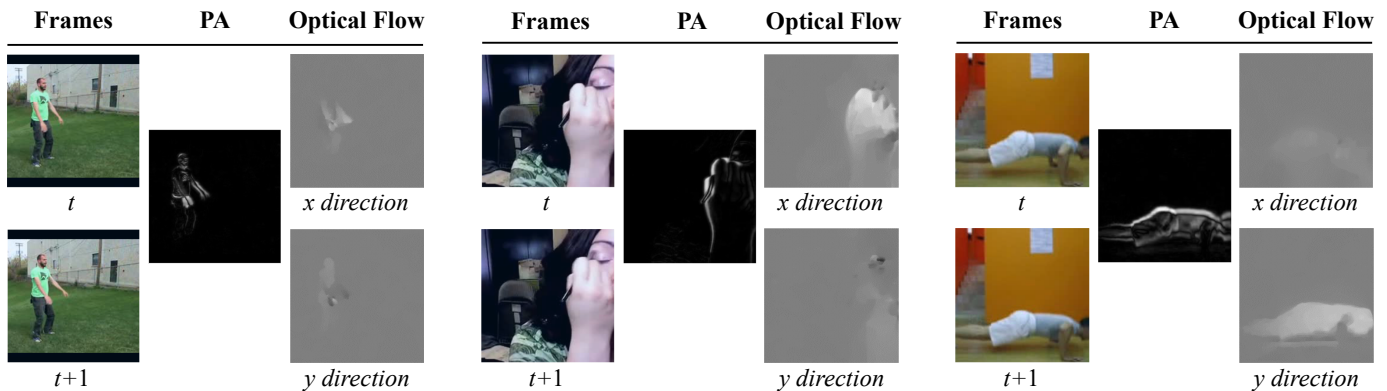
Fig. 7. Visualization of two adjacent frames and their corresponding PA and optical flow. Left: *BodyWeightSquats*. Middle: *ApplyEyeMakeup*. Right: *PushUps*. *Best viewed in color and zoomed in.*

TABLE VI
COMPARISON RESULTS OF PAN WITH OTHER STATE-OF-THE-ART METHODS ON KINETICS400, UCF101 AND HMDB51 DATASETS.

| Method | FLOPs ×views | Kinetics400 top1 (top5) | UCF101 | HMDB51 |
|---|---|---|---|---|
| STC [50] | - | 68.7 (88.5) | 93.7 | 66.8 |
| ECO_En [48] | 368G×1 | 70.0 (89.4) | 94.8 | 72.4 |
| I3D_RGB [16] | 108G×- | 71.1 (89.3) | 95.1 | 74.3 |
| SlowFast-4×16 [34] | 36.1G×30 | **75.6** (92.1) | - | - |
| ARTNet [51] | 23.5×250 | 69.2 (88.3) | 94.3 | 70.9 |
| Zhao *et. al.* [31] | - | 71.5 (89.9) | 95.9 | - |
| TSN_RGB [12] | 53G×10 | 69.1 (88.7) | 91.1 | - |
| OFF_RGB [33] | - | - | 93.3 | - |
| TSM [37] | 43G×30 | 74.1 (91.2) | 95.9 | 73.5 |
| TEA [45] | 35G×30 | 75.0 (91.8) | 96.9 | 73.3 |
| I3D_RGB+Flow [16] | - | 74.2 (91.3) | **98.0** | **80.7** |
| TSN_RGB+Flow [12] | - | 73.9 (91.1) | 97.0 | - |
| OFF_RGB+Flow [33] | - | - | 96.0 | 74.2 |
| PAN_Lite (Ours) | 47G×2 | 73.1 (91.1) | 96.0 | 74.5 |
| PAN_Full (Ours) | 88G×2 | 74.4 (91.6) | 96.5 | 77.0 |
| PAN_En (Ours) | 135G×2 | <u>75.3</u> (**92.4**) | <u>97.2</u> | <u>77.3</u> |

\* For fair comparison, all the results on UCF101 and HMDB51 are obtained with Kinetics400 pre-train.

~20% computational cost (270G *vs.* 1290G). Notably, our PAN is even superior to the optical flow based methods (third part) except for I3D [16]. Since I3D is based on heavy 3D convolutions and takes pre-computed optical flow as input, its computational cost is much more expensive than our PAN.

*E. Visualization and Discussion*

In this section, we present the visualization results of two adjacent frames and their corresponding PA and optical flow in Fig. 7. PA can model motion information as optical flow does, but is more advanced for action recognition task with its special characteristic. Similar to optical flow, the computed PA highlights the moving objects and suppresses the stationary background, which visually demonstrates that our PA well characterizes the instantaneous motion. Differently, our PA focuses more on the motion boundaries instead of the whole moving areas. This aligns with our motivation that modeling

small displacements of motion boundaries matters most for action recognition task. For more visualization results, please refer to our supplementary material.

Visually, our PA contains more noise than optical flow. We speculate that it is because the regularization term, penalizing high variations to obtain smooth displacement fields, is not applied for efficiency concerns, while it is used in conventional optical flow method [26]. Although our PA has been proved superior to several optical flow methods in terms of motion modeling efficiency and action recognition accuracy, we have not fully exploited its potential because of the noise interference. So in the future, we plan to design the noise mitigation method to obtain more smooth PA.

## VI. CONCLUSION

In this paper, we shed light on fast action recognition by lifting the reliance on optical flow. We design a concise motion cue called *Persistence of Appearance* (PA) to capture motion information directly from RGB frames. In contrast to optical flow, our PA is more effective by focusing more on modeling the small displacements of motion boundaries, and it is more efficient by simply calculating the pixel-wise differences between two adjacent frames in feature space. Its efficiency, effectiveness and flexibility have been well elaborated by extensive theoretical support (Sec. III-A), experimental support (Sec. V-B) and visualization support (Sec. V-E). The motion modeling speed of our PA reaches 1000× faster than that of conventional optical flow method (8196fps vs 8fps). To further aggregate the short-term dynamics in PA to long-term dynamics, we also propose a temporal fusion strategy named *Various-timescale Aggregation Pooling* (VAP), which enables the network to capture long-range various-timescale interdependencies. The proposed PA and VAP are finally incorporated to form a unified framework call *Persistent Appearance Network* (PAN). Extensive experiments on six challenging benchmarks demonstrate that our proposed PAN achieves the state-of-the-art recognition performance. Most importantly, it significantly accelerates the inference process of action recognition with the powerful motion cue PA.

## References

[1] K. Simonyan and A. Zisserman, "Two-stream convolutional networks for action recognition in videos," in *Advances in neural information processing systems*, 2014, pp. 568–576.

[2] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri, "Learning spatiotemporal features with 3d convolutional networks," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 4489–4497.

[3] X. Wang, R. Girshick, A. Gupta, and K. He, "Non-local neural networks," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.

[4] R. Goyal, S. E. Kahou, V. Michalski, J. Materzynska, S. Westphal, H. Kim, V. Haenel, I. Fruend, P. Yianilos, M. Mueller-Freitag, and Others, ""The" Something Something" Video Database for Learning and Evaluating Visual Common Sense." in *ICCV*, vol. 2, 2017, p. 8.

[5] W. Kay, J. Carreira, K. Simonyan, B. Zhang, C. Hillier, S. Vijayanarasimhan, F. Viola, T. Green, T. Back, P. Natsev, and Others, "The kinetics human action video dataset," *arXiv preprint arXiv:1705.06950*, 2017.

[6] J. Materzynska, G. Berger, I. Bax, and R. Memisevic, "The jester dataset: A large-scale video dataset of human gestures," in *The IEEE International Conference on Computer Vision (ICCV) Workshops*, Oct 2019.

[7] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.

[8] K. He, X. Zhang, S. Ren, and S. Jian, "Deep Residual Learning for Image Recognition," in *IEEE Conference on Computer Vision & Pattern Recognition*, 2016.

[9] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1–9.

[10] C. Feichtenhofer, A. Pinz, and A. Zisserman, "Convolutional two-stream network fusion for video action recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 1933–1941.

[11] B. Zhang, L. Wang, Z. Wang, Y. Qiao, and H. Wang, "Real-time action recognition with deeply transferred motion vector cnns," *IEEE Transactions on Image Processing*, vol. 27, pp. 2326–2339, 2018.

[12] L. Wang, Y. Xiong, Z. Wang, Y. Qiao, D. Lin, X. Tang, and L. Van Gool, "Temporal segment networks: Towards good practices for deep action recognition," in *European conference on computer vision*. Springer, 2016, pp. 20–36.

[13] Y. Zhu, Z. Lan, S. Newsam, and A. G. Hauptmann, "Hidden two-stream convolutional networks for action recognition," *arXiv preprint arXiv:1704.00389*, 2017.

[14] S. Ji, W. Xu, M. Yang, and K. Yu, "3d convolutional neural networks for human action recognition," *IEEE transactions on pattern analysis and machine intelligence*, vol. 35, no. 1, pp. 221–231, 2013.

[15] Z. Qiu, T. Yao, and T. Mei, "Learning spatio-temporal representation with pseudo-3d residual networks," in *proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 5533–5541.

[16] J. Carreira and A. Zisserman, "Quo vadis, action recognition? a new model and the kinetics dataset," in *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 6299–6308.

[17] A. Dosovitskiy, P. Fischer, E. Ilg, P. Hausser, C. Hazirbas, V. Golkov, P. Van Der Smagt, D. Cremers, and T. Brox, "Flownet: Learning optical flow with convolutional networks," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 2758–2766.

[18] E. Ilg, N. Mayer, T. Saikia, M. Keuper, A. Dosovitskiy, and T. Brox, "Flownet 2.0: Evolution of optical flow estimation with deep networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 2462–2470.

[19] A. Ranjan and M. J. Black, "Optical flow estimation using a spatial pyramid network," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 4161–4170.

[20] D. Sun, X. Yang, M.-Y. Liu, and J. Kautz, "Pwc-net: Cnns for optical flow using pyramid, warping, and cost volume," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 8934–8943.

[21] L. Sevilla-Lara, Y. Liao, F. Guney, V. Jampani, A. Geiger, and M. J. Black, "On the integration of optical flow and action recognition," *arXiv preprint arXiv:1712.08416*, 2017.

[22] J. Y.-H. Ng, J. Choi, J. Neumann, and L. S. Davis, "Actionflownet: Learning motion representation for action recognition," in *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE, 2018, pp. 1616–1624.

[23] M. D. Zeiler and R. Fergus, "Visualizing and understanding convolutional networks," in *European conference on computer vision*. Springer, 2014, pp. 818–833.

[24] C. Zhang, Y.-X. Zou, G. Chen, and L. Gan, "Pan: Persistent appearance network with an efficient motion cue for fast action recognition," *Proceedings of the 27th ACM International Conference on Multimedia*, 2019.

[25] B. K. Horn and B. G. Schunck, "Determining optical flow," *Artificial intelligence*, vol. 17, no. 1-3, pp. 185–203, 1981.

[26] C. Zach, T. Pock, and H. Bischof, "A duality based approach for realtime tv-l 1 optical flow," in *Joint pattern recognition symposium*. Springer, 2007, pp. 214–223.

[27] D. Sun, S. Roth, and M. J. Black, "Secrets of optical flow estimation and their principles," in *2010 IEEE computer society conference on computer vision and pattern recognition*. IEEE, 2010, pp. 2432–2439.

[28] L. Fan, W. Huang, C. Gan, S. Ermon, B. Gong, and J. Huang, "End-to-end learning of motion representation for video understanding," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 6016–6025.

[29] W. Huang, L. Fan, M. Harandi, L. Ma, H. Liu, W. Liu, and C. Gan, "Toward efficient action recognition: Principal backpropagation for training two-stream networks," *IEEE Transactions on Image Processing*, vol. 28, no. 4, pp. 1773–1782, 2018.

[30] H. Bilen, B. Fernando, E. Gavves, and A. Vedaldi, "Action recognition with dynamic image networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, pp. 2799–2813, 2018.

[31] Y. Zhao, Y. Xiong, and D. Lin, "Recognize Actions by Disentangling Components of Dynamics," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 6566–6575.

[32] M. Lee, S. Lee, S. Son, G. Park, and N. Kwak, "Motion feature network: Fixed motion filter for action recognition," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 387–403.

[33] S. Sun, Z. Kuang, L. Sheng, W. Ouyang, and W. Zhang, "Optical flow guided feature: a fast and robust motion representation for video action recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 1390–1399.

[34] C. Feichtenhofer, H. Fan, J. Malik, and K. He, "Slowfast networks for video recognition," *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 6201–6210, 2019.

[35] R. Girdhar, D. Ramanan, A. Gupta, J. Sivic, and B. Russell, "Action-vlad: Learning spatio-temporal aggregation for action classification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 971–980.

[36] B. Zhou, A. Andonian, A. Oliva, and A. Torralba, "Temporal relational reasoning in videos," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 803–818.

[37] J. Lin, C. Gan, and S. Han, "Tsm: Temporal shift module for efficient video understanding," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 7083–7093.

[38] C. J. Howard and A. O. Holcombe, "Unexpected changes in direction of motion attract attention," *Attention, Perception, & Psychophysics*, vol. 72, no. 8, pp. 2087–2095, 2010.

[39] L. Itti and P. F. Baldi, "Bayesian surprise attracts human attention," in *Advances in neural information processing systems*, 2006, pp. 547–554.

[40] H. Li, G. Chen, G. Li, and Y. Yu, "Motion guided attention for video salient object detection," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 7274–7283.

[41] S. Chen and Y.-G. Jiang, "Motion guided spatial attention for video captioning," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, 2019, pp. 8191–8198.

[42] K. Soomro, A. R. Zamir, and M. Shah, "UCF101: A dataset of 101 human actions classes from videos in the wild," *arXiv preprint arXiv:1212.0402*, 2012.

[43] H. Kuehne, H. Jhuang, E. Garrote, T. Poggio, and T. Serre, "Hmdb: a large video database for human motion recognition," in *2011 International Conference on Computer Vision*. IEEE, 2011, pp. 2556–2563.

[44] X. Wang and A. Gupta, "Videos as space-time region graphs," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 399–417.

[45] Y. Li, B. Ji, X. Shi, J. Zhang, B. Kang, and L. Wang, "Tea: Temporal excitation and aggregation for action recognition," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.

[46] D. Tran, J. Ray, Z. Shou, S.-F. Chang, and M. Paluri, "Convnet architecture search for spatiotemporal feature learning," *arXiv preprint arXiv:1708.05038*, 2017.

[47] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," *arXiv preprint arXiv:1502.03167*, 2015.

[48] M. Zolfaghari, K. Singh, and T. Brox, "Eco: Efficient convolutional network for online video understanding," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 695–712.

[49] S. Xie, C. Sun, J. Huang, Z. Tu, and K. Murphy, "Rethinking spatiotemporal feature learning: Speed-accuracy trade-offs in video classification," in *ECCV*, 2018.

[50] A. Diba, M. Fayyaz, V. Sharma, M. M. Arzani, R. Yousefzadeh, J. Gall, and L. V. Gool, "Spatio-temporal channel correlation networks for action classification," in *ECCV*, 2018.

[51] L. Wang, W. Li, W. Li, and L. Van Gool, "Appearance-and-relation networks for video classification," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 1430–1439.

# −Supplementary Material−
## PAN: Towards Fast Action Recognition via Learning Persistence of Appearance

Can Zhang, Yuexian Zou*, *Senior Member, IEEE,* Guang Chen, and Lei Gan

## I. IMPLEMENTATION DETAILS OF TABLE I AND TABLE II

In this section, we describe the implementation details for the comparative experiments of two encoding schemes (see Table I) and the comparative experiments of PA with other mainstream optical flow computation methods (see Table II). To have an apple-to-apple comparison, the implementation details of all the experiments are identical. We set $N = 8$ and $m = 6$, *i.e.*, 8 sampled "6-frame stack" RGB frames are sampled as input.

To measure the efficiency of the motion modeling methods, we use three evaluation metrics including the computational cost (FLOPs), the number of parameters (#Param) and inference speed (Speed). All the efficiency metrics in these two tables are measured on a single NVIDIA TITAN X GPU with 1 mini batch size and 1 CPU thread. As the I/O is most relevant to hardware and operating system, the runtime speeds are reported without considering I/O.

To evaluate the performance of the motion representation (encoded PA or optical flow) on action recognition task, we obtain accuracy results on UCF101 split1 using the motion representation as the input modality. Following the TSN manner, we first sample frames from evenly divided video segments, then these frames are fed into the motion modeling module and backbone CNN (ResNet-50) sequentially. Finally, the output activations are averaged as the final prediction scores. Note that all the accuracy results in both tables are measured with the same network settings: initial learning rate: 0.001; total epochs: 80 (decreases lr by 10 after every 30 epochs); mini batchsize: 16; dropout ratio: 0.7; pretrain: ImageNet.
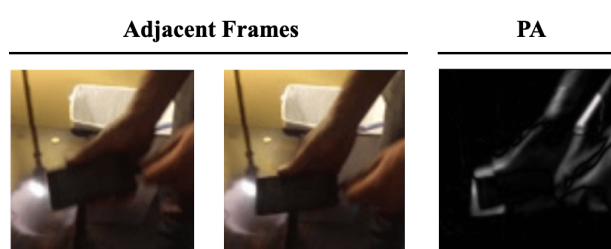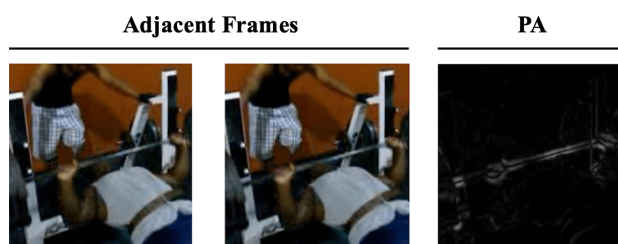
## II. VISUALIZATION RESULTS OF PA.

| Adjacent Frames | PA | | Adjacent Frames | PA |
|:---:|:---:|---|:---:|:---:|



*"ShavingBeard"*
from UCF101 dataset.



*"Pushing something so it spins"*
from Something-V1 dataset.



*"Typing"*
from UCF101 dataset.



*"Showing something to the camera"*
from Something-V1 dataset.

| **Adjacent Frames** | **PA** |
|---|---|



*"CliffDiving"*
from UCF101 dataset.

| **Adjacent Frames** | **PA** |
|---|---|



*"Moving something closer to something"*
from Something-V1 dataset.

| **Adjacent Frames** | **PA** |
|---|---|



*"BaseballPitch"*
from UCF101 dataset.

| **Adjacent Frames** | **PA** |
|---|---|



*"Plugging something into something"*
from Something-V1 dataset.

| **Adjacent Frames** | **PA** |
|---|---|



*"BenchPress"*
from UCF101 dataset.

| **Adjacent Frames** | **PA** |
|---|---|



*"Bending something so that it deforms"*
from Something-V1 dataset.

| **Adjacent Frames** | **PA** |
|---|---|



*"BlowDryHair"*
from UCF101 dataset.

| **Adjacent Frames** | **PA** |
|---|---|



*"Spinning something that quickly stops spinning"*
from Something-V1 dataset.

*"Drumming"*
from UCF101 dataset.



*"Lifting something with something on it"*
from Something-V1 dataset.



*"ParallelBars"*
from UCF101 dataset.



*"Spinning something that quickly stops spinning"*
from Something-V1 dataset.



*"PushUps"*
from UCF101 dataset.



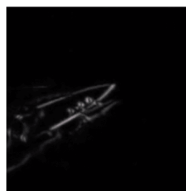*"Turning the camera upwards while filming something"*
from Something-V1 dataset.
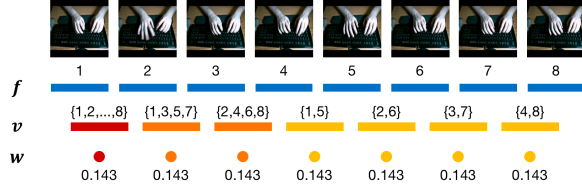


*"CuttingInKitchen"*
from UCF101 dataset.



*"Tipping something over"*
from Something-V1 dataset.

## III. Learned Timescale-wise Weights $w$ of VAP
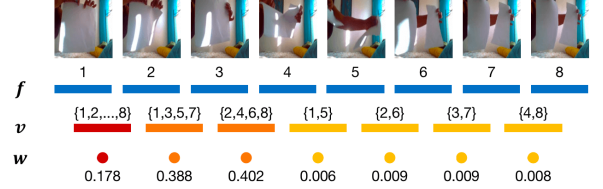


$f$   1   2   3   4   5   6   7   8

$v$   {1,2,...,8}   {1,3,5,7}   {2,4,6,8}   {1,5}   {2,6}   {3,7}   {4,8}

$w$   0.143   0.143   0.143   0.143   0.143   0.143   0.143

*"Typing"*
from UCF101 dataset.



$f$   1   2   3   4   5   6   7   8

$v$   {1,2,...,8}   {1,3,5,7}   {2,4,6,8}   {1,5}   {2,6}   {3,7}   {4,8}

$w$   0.199   0.355   0.366   0.015   0.022   0.022   0.021

*"Dropping something onto something"*
from Something-V1 dataset.



$f$   1   2   3   4   5   6   7   8

$v$   {1,2,...,8}   {1,3,5,7}   {2,4,6,8}   {1,5}   {2,6}   {3,7}   {4,8}

$w$   0.143   0.143   0.143   0.143   0.143   0.143   0.143

*"ShavingBeard"*
from UCF101 dataset.



$f$   1   2   3   4   5   6   7   8

$v$   {1,2,...,8}   {1,3,5,7}   {2,4,6,8}   {1,5}   {2,6}   {3,7}   {4,8}

$w$   0.178   0.388   0.402   0.006   0.009   0.009   0.008

*"Tearing something into two pieces"*
from Something-V1 dataset.



$f$   1   2   3   4   5   6   7   8

$v$   {1,2,...,8}   {1,3,5,7}   {2,4,6,8}   {1,5}   {2,6}   {3,7}   {4,8}

$w$   0.143   0.143   0.143   0.143   0.143   0.143   0.143

*"HulaHoop"*
from UCF101 dataset.



$f$   1   2   3   4   5   6   7   8

$v$   {1,2,...,8}   {1,3,5,7}   {2,4,6,8}   {1,5}   {2,6}   {3,7}   {4,8}

$w$   0.204   0.341   0.350   0.021   0.028   0.029   0.027

*"Holding something behind something"*
from Something-V1 dataset.



$f$   1   2   3   4   5   6   7   8

$v$   {1,2,...,8}   {1,3,5,7}   {2,4,6,8}   {1,5}   {2,6}   {3,7}   {4,8}

$w$   0.143   0.143   0.143   0.143   0.143   0.143   0.143

*"PlayingPiano"*
from UCF101 dataset.



$f$   1   2   3   4   5   6   7   8

$v$   {1,2,...,8}   {1,3,5,7}   {2,4,6,8}   {1,5}   {2,6}   {3,7}   {4,8}

$w$   0.188   0.374   0.388   0.009   0.014   0.014   0.013

*"Throwing something in the air and letting it fall"*
from Something-V1 dataset.



$f$   1   2   3   4   5   6   7   8

$v$   {1,2,...,8}   {1,3,5,7}   {2,4,6,8}   {1,5}   {2,6}   {3,7}   {4,8}

$w$   0.143   0.143   0.143   0.143   0.143   0.143   0.143

*"ApplyLipstick"*
from UCF101 dataset.



$f$   1   2   3   4   5   6   7   8

$v$   {1,2,...,8}   {1,3,5,7}   {2,4,6,8}   {1,5}   {2,6}   {3,7}   {4,8}

$w$   0.203   0.344   0.354   0.019   0.027   0.027   0.026

*"Unfolding something"*
from Something-V1 dataset.