# Exploiting Inter-Frame Regional Correlation for Efficient Action Recognition

Yuecong Xu[a,][*], Jianfei Yang[a,], Kezhi Mao[a,], Jianxiong Yin[b,], Simon See[b,]

*[a]School of Electrical and Electronic Engineering, Nanyang Technological University,
50 Nanyang Avenue, 639798, Singapore.*
*[b]NVIDIA AI Tech Centre,*
*3 International Business Park Rd, #01-20A Nordic European Centre, 609927, Singapore.*

---

## Abstract

Temporal feature extraction is an important issue in video-based action recognition. Optical flow is a popular method to extract temporal feature, which produces excellent performance thanks to its capacity of capturing pixel-level correlation information between consecutive frames. However, such a pixel-level correlation is extracted at the cost of high computational complexity and large storage resource. In this paper, we propose a novel temporal feature extraction method, named Attentive Correlated Temporal Feature (ACTF), by exploring inter-frame correlation within a certain region. The proposed ACTF exploits both bilinear and linear correlation between successive frames on the regional level. Our method has the advantage of achieving performance comparable to or better than optical flow-based methods while avoiding the introduction of optical flow. Experimental results demonstrate our proposed method achieves the state-of-the-art performances of 96.3% on UCF101 and 76.3% on HMDB51 benchmark datasets.

*Keywords:* Action Recognition, Inter-frame Correlation, Feature Extraction
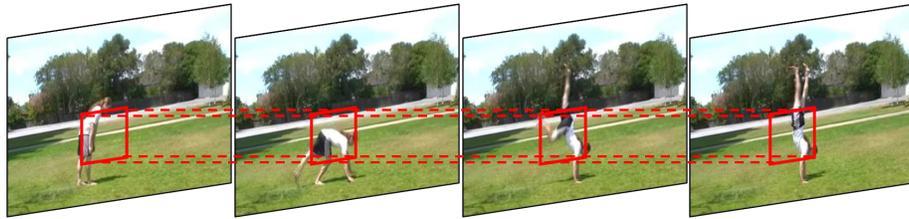
---

*[*]Corresponding author
*Email addresses:* xuyu0014@e.ntu.edu.sg (Yuecong Xu), yang0478@e.ntu.edu.sg
(Jianfei Yang), ekzmao@ntu.edu.sg (Kezhi Mao), jianxiongy@nvidia.com (Jianxiong Yin),
ssee@nvidia.com (Simon See)

## 1. Introduction

Action recognition has received considerable attention from the vision community in recent years (Herath et al., 2017; Yang et al., 2019; Carmona & Climent, 2018; Wang & Wang, 2018; Keeli, 2018; Sahoo & Ari, 2019) thanks to its increasing applications in various fields, such as surveillance (Danafar & Gheissari, 2007; Yang et al., 2018; Xiang & Gong, 2008; Li et al., 2017; Mabrouk & Zagrouba, 2018; Kardas & Cicekli, 2017) and smart homes (Wu et al., 2010; Yang et al., 2018; Ortis et al., 2017; Lundstrm et al., 2016; Lee et al., 2017) etc. Compared with static images, videos contain additional temporal information. Hence, extracting and handling temporal information is very critical in action recognition.

To extract temporal features underlying a video, a few methods have been proposed in the literature. Most of these efforts can be organized into two categories. The first category is the two-stream methods. A typical work in this category is the one proposed in (Simonyan & Zisserman, 2014), which conduct the classification using temporal features and spatial features separately. The two types of features are integrated through classification decision fusion. The second category is the 3D ConvNet methods, which extract spatial and temporal features jointly by expanding the convolution kernel of 2D ConvNets to the temporal dimension. A seminal work in this category is the C3D (Tran et al., 2015) network. A detailed review of the methods in both categories is given in Section 2.

The methods in the two categories have their respective merits and limitations. The two-stream methods often produce the state-of-the-art performance, yet this is achieved at the cost of heavy reliance on accurate temporal feature. Therefore, the two-stream methods usually involve computation or estimation of optical flow, both of which require high computational power and large storage resource. Also, obtaining optical flow needs to be performed prior to the training of the network, thus methods utilizing optical flow cannot be trained end-to-end. On the other hand, the 3D ConvNets-based methods are computationally less demanding, yet their performances are usually inferior to that of the two-stream methods. A possible reason would be the temporal pooling used for dimension reduction towards the complete representation. Temporal

(a) Action "Handstand"



(b) Action "Brushing Teeth"

Figure 1: Illustration of extracting inter-frame corresponding-regional correlation for action recognition. The temporal feature of an action is related to the correlation appearance between frames. Actions that are faster such as "Handstand" in (a) exhibits obvious change within the indicated box. Slower and more static actions such as "Brushing Teeth" in (b) shows little change between frames. To cope with both situations, bilinear operation is employed to extract the inter-frame corresponding-regional correlation

pooling extracts only linear feature along the temporal dimension of the video through pooling operation. With only the linear feature being extracted, we argue that part of the temporal feature is lost during the pooling operation.

In this paper, we present a novel method for temporal feature extraction, which achieves performance comparable to or even better than two-stream methods, yet demands less computational power. Intuitively, temporal feature of an action is related to the correlation of appearance between frames within a certain region. For instance, in Figure 1a, the indicated box across the series of frames shows how a person turns upside down, and is related to the action of "Handstand". Therefore, instead of using optical flow, our proposed method extracts temporal features by extracting the correlation of neighbouring frames with respect to the corresponding regions. The degree of change of appearance varies between different actions. For actions that are slower or more static, neighboring sampled frames could be very similar. One example is the action of "Brushing Teeth" shown in Figure 1b. If linear correlation, such as the differ-

ence in RGB value is employed, the correlation extracted would fail to contain temporal information of the video. To cope with the various type of actions, the inter-frame correlation would thus be computed through bilinear operation. The complete temporal feature, named as Attentive Correlated Temporal Feature (ACTF), is obtained through attentive combination of the inter-frame corresponding-regional correlation feature and the inter-frame mean feature obtained through inter-frame temporal average pooling. Our main contributions are summarized as follows:

* We propose a novel temporal feature extraction method: Attentive Correlated Temporal Feature (ACTF), for action recognition. First, ACTF exploits inter-frame corresponding-regional correlation to implicitly capture temporal information without the use of optical flow. Second, by excluding optical flow estimation or calculation, ACTF can be combined with any spatial feature extraction network under the two-stream structure to implement end-to-end training. Third, ACTF leads to performance comparable to or even better than optical flow-based methods, yet it demands less computation and memory due to the exclusion of optical flow.

* We conduct extensive experiments on two action recognition benchmark datasets: UCF101 (Soomro et al., 2012) and HMDB51 (Kuehne et al., 2011) with a framework utilizing our proposed ACTF. The results demonstrate that our proposed ACTF brings noticeable improvements over baseline methods, achieving state-of-the-art performance for these datasets.

The rest of this paper is organized as follows. Related works for action recognition tasks and the use of regional correlation are discussed in Section 2. In Section 3, we introduce the proposed Attentive Correlated Temporal Feature (ACTF) in detail. After that, we present and analyze the experimental results of our proposed ACTF feature, with a thorough ablation study on the design of ACTF. Finally, we conclude the paper in Section 5.

## 2.  Related Work

Action recognition is one of the core tasks in video understanding. Compared to image understanding tasks, video understanding tasks are more complex due to the additional temporal dimension in videos. The extraction and handling of temporal feature underlying videos is thus the main challenge of the action recognition task.

### 2.1.  Action Recognition with Optical Flow

To extract temporal feature with high quality, previous works (Simonyan & Zisserman, 2014; Feichtenhofer et al., 2016; Wang et al., 2016) adopt a two-stream strategy where temporal feature is extracted in parallel with spatial feature. The temporal feature extraction is performed by feeding a stack of optical flow frames computed by TV-L1 (Zach et al., 2007) to a ConvNet. TSN (Wang et al., 2016) improves performance of the original two-stream network (Simonyan & Zisserman, 2014) through segmenting the video and input the RGB alongside optical flow frames for each segment to ConvNets with shared parameter. The action recognition result is produced through segmental consensus fusion. Meanwhile, ST-ResNet (Feichtenhofer et al., 2016) builds upon ResNet (He et al., 2016), which is a deeper 2D ConvNet. These two-stream methods achieve competitive results in action recognition, but using optical flow for temporal feature extraction has limitations. Optical flow extraction is known to be computationally expensive and memory intensive. In addition, as optical flow requires pre-computation, the use of optical flow as temporal feature prohibits fully end-to-end training of the network.

### 2.2.  Temporal Feature Extraction without Optical Flow

To address the limitations imposed by utilizing optical flow, subsequent works proposed to extract temporal feature to replace optical flow. One category of methods involves the estimation of optical flow through neural network. FlowNet (Dosovitskiy et al., 2015; Ilg et al., 2017), MotionNet (Zhao et al., 2018a), LMoF (Li et al., 2018), TVNet (Fan et al., 2018) and more recently Representation Flow (Piergiovanni & Ryoo, 2019) all belong to such category. More specifically, FlowNet (Dosovitskiy et al., 2015) learns optical flow from synthetic ground truth data. MotionNet (Zhao

5

et al., 2018a) produces optical flow through next frame prediction. LMoF (Li et al., 2018) further constructs a learnable directional filtering layer to cope with optical flow estimation in blur videos. To further boost the performance of optical flow estimation, TVNet (Fan et al., 2018) unfolds the TV-L1 (Zach et al., 2007) optical flow extraction method and formulates it with neural network. Representation Flow (Piergiovanni & Ryoo, 2019) extends from TVNet (Fan et al., 2018) and constructs fully-differentiable convolutional layers to estimate optical flow. The layers could be stacked on top of each other to obtain flow-of-flow features which could capture longer-term motion representation. Although these optical flow estimation methods render networks to be trained in an end-to-end manner, they are still expensive in computation and intensive in memory, with longer run-time during inference. Another category extracts temporal feature jointly with spatial feature by constructing 3D ConvNets. C3D (Tran et al., 2015), 3D-ResNet(Tran et al., 2018), 3D-ResNext (Hara et al., 2018), I3D (Carreira & Zisserman, 2017) and Asymmetric 3D-CNN (Yang et al., 2019) belong to this category. More specifically, C3D (Tran et al., 2015) is one of the primary works where the CNN network is expanded to the temporal dimension. Subsequent networks such as 3D-ResNet (Tran et al., 2018) and 3D-ResNext (Hara et al., 2018) are deeper and larger 3D ConvNets. To further reduce parameter for even faster training, Carreira *et at.* inflates 2D ConvNets into 3D structure. This simplifies the work of constructing 3D ConvNets by simply convert image classification models to 3D models by endowing filters and pooling kernels with the additional temporal dimension. Whereas Yang *et al.* (Yang et al., 2019) proposed to utilize *MicroNets* to construct asymmetric 3D ConvNets. 3D ConvNets benefit from end-to-end training, and requires only RGB input. Yet the temporal feature is extracted through pooling along the temporal dimension, and thus extracts only linear temporal feature. This causes part of the temporal feature might be lost during feature extraction operation.

*2.3. Correlation Modeling and Bilinear Pooling*

Correlation modeling and bilinear pooling have been used in action recognition and have shown its success in improving temporal feature extraction. More specifically, Diba *et al.* proposed TLE (Diba et al., 2017) which represents the whole video
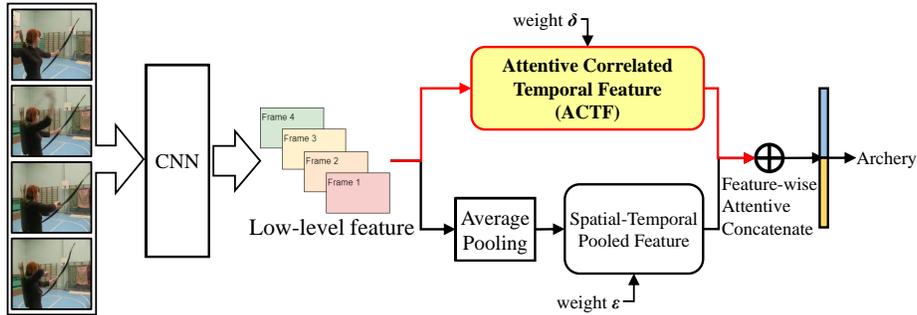
Figure 2: Detailed illustration of applying ACTF for action recognition. The sharp rectangles represent the networks or operations performed, while the rounded rectangles represent the resulting features. The overall framework takes the raw RGB frames as input. The low-level feature of the RGB frames is extracted through a ConvNet (CNN). From the low-level feature, we obtain the spatial-temporal pooled feature of the video through average pooling across both spatial and temporal dimensions. This feature is regarded as the spatial feature of the video. Simultaneously, we obtain the ACTF as the temporal feature of the video. Both features are combined attentively to form the whole representation of the video.

by a bilinear model. The input of TLE is the temporal aggregated feature obtained by aggregating the features of each video segment. Meanwhile, Zhao *et al.* (Zhao et al., 2018b) utilizes correlation modeling for the construction of cost volume, which is an intermediate facility of optical flow estimation. More recently, inspired by the non-local mean operation for image denoising (Buades et al., 2005; Li & Suen, 2016), Wang *et al.* (Wang et al., 2018b) presented non-local operations to capture correlation on a pixel level as the representation of the temporal feature. Unlike the mentioned works above, our work utilizes correlation modeling on a frame-wise regional level, and computes the inter-frame correlation within a certain region through bilinear operation. This guarantees our temporal feature extraction method to be more computation efficient while maintaining improvement in temporal feature extraction.

## 3. Method

The primary goal of our work is to develop an effective video-based action recognition framework with focus on temporal feature extraction. The main idea of the proposed method is to explore correlation of successive frames within a certain region,

which captures temporal information. The extracted correlation feature can work with various low-level feature extraction networks that are normally convolutional neural networks (ConvNets), e.g. C3D and 3D-ResNet. These networks normally adopt a simple temporal pooling operation for obtaining the video representation. We propose an ACTF model to effectively extract the inter-frame corresponding-regional correlation feature and combine it with the feature obtained from simple temporal pooling. Next, we present a general action recognition framework that uses the proposed ACTF for temporal feature extraction, and then describe the details of the ACTF. The attention mechanism employed in ACTF will also be briefly explained.

### 3.1. General Framework for Action Recognition with ACTF

The prominent methods for action recognition employ multiple modality networks, e.g. two-stream convolutional networks (Simonyan & Zisserman, 2014). In these networks, temporal and spatial features are extracted and processed separately. Figure 2 shows the overall framework in our study. Given an input video as a sequence of frames, the low-level feature of each frame is first extracted through a convolutional neural network (ConvNet). The resulted low-level feature is denoted by $\mathbf{F} \in \mathbb{R}^{t \times C_{out} \times H \times W}$, where $t$ denotes the number of frames, $C_{out}$ denotes the number of channels, and $H$, $W$ are the height and width. Subsequently, we obtain two features from this low-level feature, namely the Spatial-Temporal Pooled feature, and the ACTF feature. The Spatial-Temporal Pooled feature $\mathbf{V}_{stpooled}$ is obtained by performing spatial-temporal average pooling over the low-level feature. The ACTF feature $\mathbf{V}_{actf}$ is obtained through an attentive concatenation of features obtained by performing both bilinear and linear operations on successive frames. Each of the two features characterizes a different perspective of the video. Performing average pooling over the low-level feature results in a feature that provides a general appearance pattern of the video. Thus the Spatial-Temporal Pooled feature as shown in Figure 2 is referred to as the spatial feature of the video in this paper. Meanwhile, the ACTF feature captures the correlation pattern of successive frames within a certain region, and is referred to as the temporal feature here. Both features have a dimension of $C_{out}$, *i.e.* $\mathbf{V}_{actf} \in \mathbb{R}^{C_{out}}$, and $\mathbf{V}_{stpooled} \in \mathbb{R}^{C_{out}}$.
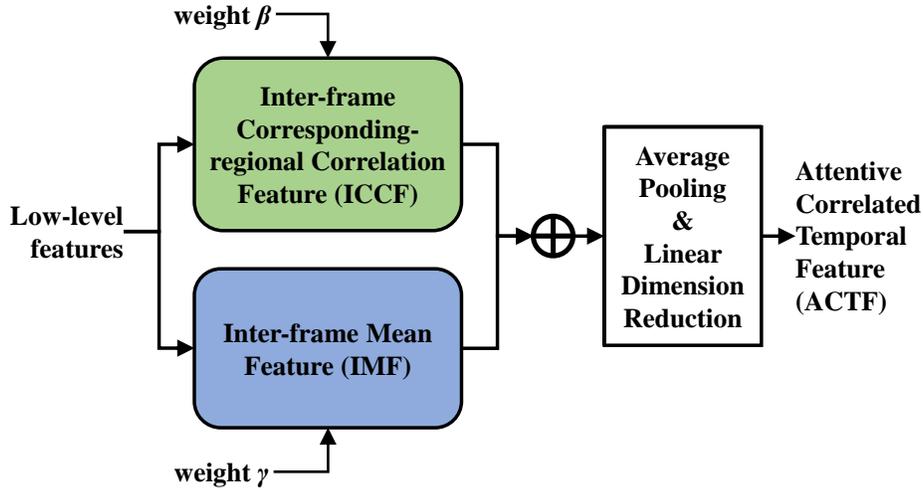
8

Figure 3: Illustration of the pipeline for extracting ACTF. From the low-level feature extracted, we extract two forms of inter-frame correlation features. A bilinear inter-frame correlation feature, extracted as the Inter-frame Corresponding-regional Correlation Feature (ICCF), as well as a linear inter-frame correlation feature, extracted as the Inter-frame Mean Feature (IMF). The features are combined attentively to form the ACTF.

In certain scenarios, the spatial feature of the video is sufficient to produce satisfactory action recognition result. This occurs when certain action types are associated with certain visual elements. Meanwhile, in other scenarios, temporal feature play a more vital role. This occurs when visual elements of the video may appear in different actions. Thus, we adopt a feature-wise attentive concatenation method to dynamically combine the spatial and temporal features.

### 3.2. Extraction of ACTF feature

Previous works (Simonyan & Zisserman, 2014; Wang et al., 2016; Zhu et al., 2017) show the importance of temporal feature in action recognition. However, most temporal information representations, such as optical flow used in two-stream convolutional network (Simonyan & Zisserman, 2014) or non-local operations in non-local 3D ConvNets (Wang et al., 2018b), are computationally expensive. This is due to the fact that both optical flow and non-local operations compute the correlations between successive frames on a pixel level. Computationally efficient RGB difference is employed

in (Wang et al., 2016) to capture inter-frame relation, but shows inferior performance when used in combination with spatial feature. For slower actions where successive frames are very similar, RGB difference would return zero-valued correlation, and fail to capture the temporal feature, which might explain the inferior performance.

In this work, we propose to explore more sophisticated operations, such as bilinear function, on successive frames for temporal feature extraction. This is inspired by the bilinear operation used for fine-grained image recognition (Gao et al., 2016; Lin et al., 2015), where the within-image bilinear operation is used to learn local pairwise feature correlation through the outer product at every single position of the image. In this paper, the bilinear operation is extended across successive frames to discover inter-frame correlation within a certain region. Figure 3 shows the pipeline for extracting ACTF.

More specifically, given a video sequence, as described in Section 3.1, the low-level feature of the video is extracted through a ConvNet, whose output is $\mathbf{F} \in \mathbb{R}^{t \times C_{out} \times H \times W}$. We then extract a bilinear inter-frame correlation feature, the Inter-frame Corresponding-regional Correlation Feature (ICCF), and a linear inter-frame feature, the Inter-frame Mean Feature (IMF). The extraction function for the ICCF is denoted by $\mathcal{P}_{bilinear}$, while the extraction function for the IMF is denoted by $\mathcal{P}_{mean}$.

We first describe $\mathcal{P}_{bilinear}$, which is the extraction function for the bilinear inter-frame correlation feature denoted as ICCF. Figure 4 shows the details of extracting the ICCF. Denote $\mathbf{f}_i \in \mathbb{R}^{C_{out} \times H \times W}$ as the low-level feature extracted for frame $i$. To extract the bilinear inter-frame correlation feature, $\mathcal{P}_{bilinear}$ computes the pairwise bilinear correlation with respect to two successive frames within a certain region as follows:

$$\mathbf{b}_i = \mathcal{P}_{bilinear}(\mathbf{f}_i, \mathbf{f}_{i+1}) \tag{1}$$

Here $\mathbf{b}_i$ is the bilinear inter-frame correlation feature, and $\mathbf{b}_i \in \mathbb{R}^{C_{bilinear} \times H \times W}$, where $C_{bilinear}$ denotes the number of channels of the ICCF.

More specifically, at the spatial location of $\mathcal{S}$, the feature of the current frame and the next frame is denoted as $\mathbf{f}_{i,\mathcal{S}}$ and $\mathbf{f}_{i+1,\mathcal{S}}$. We denote the bilinear operation function

at location $\mathcal{S}$ to be $\mathcal{B}_\mathcal{S}$, and is formulated by the following equation:

$$\mathcal{B}_{i,\mathcal{S}} = \mathbf{f}_{i,\mathcal{S}}\mathbf{f}_{i+1,\mathcal{S}}{}^T \tag{2}$$

At the spatial location $\mathcal{S}$, the feature for frame $i$ is of size $C_{out} \times 1$. Thus, from Equation 2, the bilinear inter-frame correlation feature at location $\mathcal{S}$ is of size $C_{out} \times C_{out}$. We then reshape it such that the result would be of size $C_{out}{}^2 \times 1$.

Although the bilinear inter-frame correlation feature obtained through Equation 2 is direct, such feature representation is very high-dimensional. In our case where $C_{out}$ is around 750, the dimension of the bilinear inter-frame correlation feature at each spatial location is more than 500,000. Such high dimensional representation is impractical. Therefore, to obtain the desired bilinear correlation, we adopt a compact form of bilinear operation as implemented in (Gao et al., 2016).

The basis of the compact form of the bilinear operation is to find a low dimension projection function of $\mathcal{B}_{i,\mathcal{S}}$, denoted as $\mathcal{C}_{i,\mathcal{S}}$. The two functions are equivalent with respect to a linear kernel machine. Given two pairs of frames: frames $(i, i+1)$ and frames $(j, j+1)$, a linear kernel machine is formulated as:

$$\begin{aligned}
\langle \mathcal{B}_{i,\mathcal{S}}, \mathcal{B}_{j,\mathcal{S}} \rangle &= \langle \mathbf{f}_{i,\mathcal{S}}\mathbf{f}_{i+1,\mathcal{S}}{}^T, \mathbf{f}_{j,\mathcal{S}}\mathbf{f}_{j+1,\mathcal{S}}{}^T \rangle \\
&= \langle \mathbf{f}_{i,\mathcal{S}}, \mathbf{f}_{j,\mathcal{S}} \rangle^2
\end{aligned} \tag{3}$$

We then find a low dimension projection function as $\phi(\mathbf{f}_{i,\mathcal{S}}) \in \mathbb{R}^d$ such that $\langle \phi(\mathbf{f}_{i,\mathcal{S}}), \phi(\mathbf{f}_{j,\mathcal{S}}) \rangle \approx k(\mathbf{f}_{i,\mathcal{S}}, \mathbf{f}_{j,\mathcal{S}})$, where $k$ is a polynomial kernel. Such projection function $\phi(\mathbf{f}_{i,\mathcal{S}})$ would allow us to approximate Equation 3 by:

$$\begin{aligned}
\langle \mathcal{B}_{i,\mathcal{S}}, \mathcal{B}_{j,\mathcal{S}} \rangle &= \langle \mathbf{f}_{i,\mathcal{S}}, \mathbf{f}_{j,\mathcal{S}} \rangle^2 \\
&\approx \langle \phi(\mathbf{f}_{i,\mathcal{S}}), \phi(\mathbf{f}_{j,\mathcal{S}}) \rangle \\
&\equiv \langle \mathcal{C}_{i,\mathcal{S}}, \mathcal{C}_{j,\mathcal{S}} \rangle
\end{aligned} \tag{4}$$

where $\mathcal{C}_{i,\mathcal{S}} = \phi(\mathbf{f}_{i,\mathcal{S}})$ is the compact form of the bilinear operation $\mathcal{B}_{i,\mathcal{S}}$. Hence to obtain the compact form, we need to find the low dimension approximation of the polynomial kernel $k$. Here we utilize the Tensor Sketch approximation method proposed in (Pham & Pagh, 2013). Ultimately, our extraction function $\mathcal{P}_{bilinear}$ computed at each spatial location $\mathcal{S}$ for frame $i$ and the successive frame $i+1$ is equivalent to its compact form $\mathcal{C}_{i,\mathcal{S}}$.

**Low-level feature** weight $\alpha_1$

weight $\alpha_2$

weight $\alpha_3$

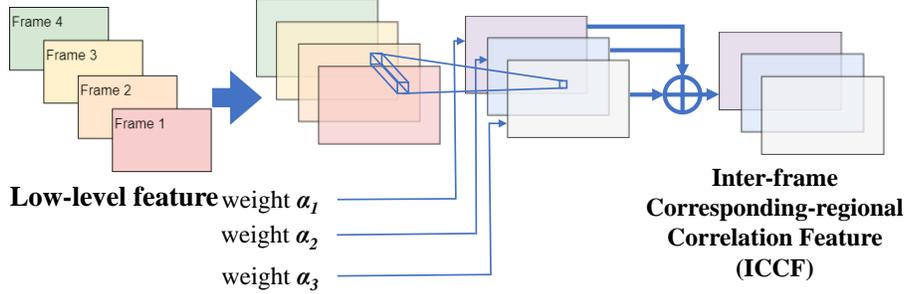**Inter-frame Corresponding-regional Correlation Feature (ICCF)**

Figure 4: Illustration of the details for extracting the bilinear inter-frame correlation feature which is the ICCF. The pairwise bilinear correlation with respect to two successive frames within a certain region is computed for each pair of successive frames. The complete bilinear feature is extracted through temporal-wise attentive concatenation of each inter-frame correlation feature.

For image recognition tasks, features extracted through bilinear function go through a sum pooling operation to extract the complete representation of the image. However, if such a pooling method is used in videos, the temporal information may be lost. This conflicts with our goal of extracting temporal information through the bilinear inter-frame correlation feature. To dynamically combine all bilinear inter-frame correlation features temporally, we apply a temporal-wise attentive concatenation to each pair of successive frames. A learnable weight parameter $\alpha_i$ is assigned to each inter-frame correlation feature $\mathbf{b}_i$. Such attentive concatenation allows the extracted ICCF to focus on the pair of frames where the action most likely takes place. The result is a feature $\mathbf{B} \in \mathbb{R}^{(t-1) \times C_{bilinear} \times H \times W}$. For each pair of successive frames, $\mathbf{B}_i = \alpha_i \mathbf{b}_i$.

To extract the temporal feature of the video more accurately, besides the bilinear inter-frame correlation feature, we would also need the linear inter-frame correlation feature denoted as IMF. The IMF provides a baseline for the bilinear inter-frame correlation feature, and is important when actions are similar temporally but very different in appearance. Following this idea, we feed the low-level feature $\mathbf{F}$ to extract the IMF in parallel with the ICCF. Unlike common temporal pooling layers where the average pooling is performed across the whole temporal dimension, the purpose of the extraction function $\mathcal{P}_{mean}$ is to capture the average of two successive frames. More specifi-

cally, the extraction function for IMF is an average pooling function with a kernel size of $(k_i, k_h, k_w)$. Here $k_h$ and $k_w$ are the kernel size corresponding to the spatial dimensions. As we need to preserve all information along the spatial dimensions, hence $k_h, k_w = 1$. To obtain the average pooling along successive temporal features $\mathbf{f}_i$ and $\mathbf{f}_{i+1}$, the kernel size along the temporal dimension $k_i$ is set to 2 instead of the whole temporal dimension length. The IMF $\mathbf{L} \in \mathbb{R}^{(t-1) \times C_{out} \times H \times W}$ is computed by:

$$\mathbf{L} = \mathcal{P}_{mean}(\mathbf{F}) \tag{5}$$

Similar to the temporal-wise attentive concatenation for bilinear inter-frame correlation features, we adopt a feature-wise attentive concatenation approach to combine the bilinear inter-frame correlation feature with the linear inter-frame correlation feature. Each of the two types of features is assigned a separate weight parameter, denoted as $\beta, \gamma$ respectively. This allows the network to dynamically focus on either feature for different actions. The result of the attentive concatenation $\mathbf{H} \in \mathbb{R}^{(t-1) \times C_{concat} \times H \times W}$ is obtained as follows:

$$\mathbf{H} = \beta\mathbf{B} \oplus \gamma\mathbf{L} \tag{6}$$

where $\oplus$ denotes the concatenation operation along the feature channel dimension. $C_{concat}$ is the total number of feature channels, which is the sum of $C_{out}$ and $C_{bilinear}$.

The complete ACTF feature $\mathbf{V}_{actf}$ is obtained as follows:

$$\mathbf{V}_{actf} = \mathcal{L}(\mathcal{P}_{average}(\mathbf{H})) \tag{7}$$

Where $\mathcal{P}_{average}$ is an average pooling operation with a kernel size of $((t-1), H, W)$ corresponding to the temporal and spatial dimensions respectively. ACTF feature summarizes both the ICCF and the IMF. $\mathcal{L}$ is a linear dimension reduction function, constructed as a multi-layer linear neural network. This allows $\mathcal{L}$ to be learnable and the overall system to be trainable in an end-to-end manner. The resulting feature is thus the ACTF feature $\mathbf{V}_{actf} \in \mathbb{R}^{C_{out}}$.

### 3.3. Attentive Concatenation of Features

Our network is designed to focus on the pairs of time steps which are more relevant to the action. Meanwhile, it is also designed to focus on the more important type of

13

feature, *i.e.* the spatial or temporal feature . To achieve both goals, we adopt attentive concatenation at each location where different features are combined. In this section, we describe how to extract the ICCF by the temporal-wise attentive concatenation of bilinear inter-frame correlation features. The combination of features $\mathbf{V}_{stpooled}$ and $\mathbf{V}_{actf}$ mentioned in Section 3.1 as well as the combination of features $\mathbf{B}$ and $\mathbf{L}$ mentioned in Section 3.2 follow similar implementations.

The attentive concatenation of all $(t-1)$ bilinear correlation features is achieved by assigning each feature with a weight $\alpha_i$ for the $i^{th}$ correlation feature $\mathbf{b}_i$. Inspired by the cascade attention network proposed in (Wang et al., 2018a), we adopt an attentive concatenation approach for the computation of weight $\alpha_i$. Formally, $\alpha_i$ is computed by:

$$\alpha_i = g(h((\mathcal{P}_{spatial}(\mathbf{b}_i))W)) \tag{8}$$

More specifically, given the $i^{th}$ bilinear correlation feature $\mathbf{b}_i \in \mathbb{R}^{C_{bilinear} \times H \times W}$, $\mathcal{P}_{spatial}$ is a spatial average pooling function with kernel size $(H, W)$. The output of $\mathcal{P}_{spatial}$ is a pooled feature vector $\mathbf{b}_{pooled,i} \in \mathbb{R}^{C_{bilinear}}$. $W \in \mathbb{R}^{C_{bilinear} \times 1}$ denotes a trainable parameter matrix, shared among all $(t-1)$ bilinear correlation features. The result of this matrix multiplication is a primitive weight parameter denoted as $\alpha_{prime,i}$.

To scale the primitive weight parameter $\alpha_{prime,i}$ to a range of $[0, 1]$, we apply a sigmoid function denoted as $h(\alpha_{prime,i})$, which is computed by:

$$h(\alpha_{prime,i}) = \frac{1}{1 + e^{-\alpha_{prime,i}}} \tag{9}$$

The weight $\alpha_i$ is then further processed from $h(\alpha_{prime,i})$ to satisfy $\Sigma \alpha_i = 1$. This is achieved by applying a softmax function denoted by $g(\cdot)$, and the weight $\alpha_i$ is calculated as follows:

$$\begin{aligned} \alpha_i &= g(h(\alpha_{prime,i})) \\ &= \frac{e^{h(\alpha_{prime,i})}}{\sum_{i=1}^{t-1} e^{h(\alpha_{prime,i})}} \end{aligned} \tag{10}$$

The weight $\alpha_i$ indicates the importance of the $i^{th}$ bilinear inter-frame correlation feature, $\mathbf{b}_i$.

14

## 4. Experiments

In this section, we present our evaluation results of the proposed work. The evaluation is conducted through action recognition experiments on two public benchmark datasets. We present state-of-the-art results on a competitive architecture, and prove the novelties on another similar baseline. We also present detailed ablation study of the components of our proposed framework to verify our design.

### 4.1. Experimental Settings

We conduct experiments on two benchmark datasets of action recognition: UCF101 (Soomro et al., 2012) and HMDB51 (Kuehne et al., 2011). The UCF101 dataset contains 13,320 videos from 101 action categories while the HMDB51 dataset contains 6,766 videos from 51 action categories. We follow the experiment settings as in (Chen et al., 2018; Tran et al., 2015, 2018) that adopt the three training/testing splits for evaluation. We report the average top-1 accuracy of the three splits. Our proposed framework for temporal feature extraction can be used in any ConvNet based networks. To obtain the state-of-the-art result, we instantiate MFNet (Chen et al., 2018).

Our experiments are implemented using PyTorch (Paszke et al., 2017). Following the implementation in (Chen et al., 2018), the input is a frame sequence with each frame of size $224 \times 224$. The output from MFNet (Chen et al., 2018) is a low-level feature of size $8 \times 768 \times 7 \times 7$, where the number of output channels is 768 . Each frame is represented by a feature of size $7 \times 7$. We set the number of channels of our ICCF to 3,840. Thus, the size of $\mathbf{H}$, described in Section 3.2 of the paper, is $7 \times 4608 \times 7 \times 7$. We design the linear dimension reduction function in Equation 7 as a three-layer linear neural network with RELU activation. For training, we utilize the pretrained model of MFNet (Chen et al., 2018) trained on Kinetics (Kay et al., 2017), a large-scale human action dataset. To accelerate our training, the pretrained model is used for the initialization of the network which includes our framework for temporal feature extraction. We use stochastic gradient descent algorithm (Bottou, 2010) for optimization, setting the weight decay to 0.0001 and the momentum to 0.9. For both datasets, our initial learning rate is set to 0.005. For UCF101 (Soomro et al., 2012)

dataset, the learning rate is decreased for four times, while for HMDB51 (Kuehne et al., 2011) dataset, the learning rate is decreased for three times. The learning rate is decreased with a factor of 0.1.

To prove that our approach can be applied to other 3D ConvNet approaches, we also apply our proposed ACTF in another 3D ConvNet. We instantiate C3D (Tran et al., 2015), a classical 3D ConvNet baseline for action recognition. Our proposed ACTF is extracted after conv5 layer of the C3D network, in parallel with the spatial-temporal pooling layer pool5, as well as the linear layer that follows. We follow the setup as in (Tran et al., 2015), using stochastic gradient descent (Bottou, 2010) with initial learning rate of 0.001. We compare the results of the C3D network with and without the temporal feature extracted by our proposed framework on HMDB51 dataset.

*4.2. Results and Comparison*

Table 1 shows the comparison of top-1 accuracy on UCF101 and HMDB51 datasets with other state-of-the-art methods including:

1. **Two-stream methods:** the original two-stream method (original TS) (Simonyan & Zisserman, 2014), Hidden Two-Stream (Hidden TS) (Zhu et al., 2017), Long-term Temporal Convolutions (LTC) (Varol et al., 2018), ActionVLAD (Girdhar et al., 2017) and Temporal Segment Network (TSN) (Wang et al., 2016)

2. **3D ConvNets-based methods:** C3D (Tran et al., 2015), TSN with RGB input (Wang et al., 2016), Res3D (Tran et al., 2017), ST-ResNet (Feichtenhofer et al., 2016), 3D-ResNext (Hara et al., 2018), R(2+1)D with RGB input (Tran et al., 2018), I3D with RGB input (Carreira & Zisserman, 2017), TVNet (Fan et al., 2018), MFNet (Chen et al., 2018) and T-C3D (Liu et al., 2018)

Our state-of-the-art performance is achieved by instantiating MFNet, denoted as MFNet-ACTF. For this experiment, we set our batch size to 80 and conduct the experiment using four NVIDIA Tesla P100 GPUs.

The performance results in Table 1 show that our network achieves the best results on both benchmark datasets. More specifically, our MFNet-ACTF network achieves a 1.7% improvement on HMDB51 dataset over the networks whose input are solely

| | Method | UCF101 | HMDB51 | FPS |
|---|---|---|---|---|
| | original TS | 88.0% | 59.4% | 14 |
| Two-stream | Hidden TS | 90.3% | 58.9% | < 14 |
| | LTC | 91.7% | 64.8% | < 14 |
| | TSN | 94.2% | 69.4% | 5 |
| | C3D | 85.2% | 65.5% | 314 |
| | TSN (RGB) | 86.2% | - | N/A |
| | Res3D | 85.8% | 54.9% | N/A |
| | T-C3D | 91.8% | 62.8% | 969 |
| | ST-ResNet | 93.5% | 66.4% | N/A |
| 3D ConvNets | 3D-ResNext | 94.5% | 70.2% | < 314 |
| | R(2+1)D (RGB) | 93.6% | 66.6% | N/A |
| | I3D (RGB) | 95.6% | 74.8% | N/A |
| | TVNet | 95.4% | 72.5% | N/A |
| | MFNet | 96.0% | 74.6% | N/A |
| Ours | **MFNet-ACTF** | **96.3%** | **76.3%** | **478** |

Table 1: Comparison of top-1 accuracy and speed with state-of-the-art methods on UCF101 and HMDB51 datasets.

| Method | Top-1 HMDB51 |
|---|---|
| C3D | 65.5% |
| C3D-single-ACTF | 67.9% |
| C3D-ACTF | 69.2% |

Table 2: Top-1 accuracy of C3D network on HMDB51 dataset with and without our proposed framework.

RGB frames. Our method even surpasses several networks with both RGB and optical flow as input. For UCF101 dataset, our MFNet-ACTF also produces the best result. It is noted that the improvement is not as significant as that on HMDB51 dataset, mainly due to the fact that there is little room for improvement.

The speed results in Table 1 show that our proposed method balances between high accuracy and relatively high inference speed. Despite achieving high accuracy on both dataset, two-stream methods such as TSN could not achieve real-time requirements, reaching only 5 FPS. Compared with two-stream methods, our proposed method is much faster in inference speed, reaching a speed of 478 FPS, which is well above real-time requirements. Our speed is even faster than that achieved by C3D network. Note that our speed is slower than that achieved by T-C3D network, but we achieved a much higher accuracy compared to theirs, with a $13.5\%$ increase in top-1 accuracy on HMDB51 dataset.

We suggest in Section 3 that our proposed framework which includes ACTF feature can be used in combination with any ConvNet-based low-level feature extraction networks, such as C3D network. To verify this, we conducted experiments on the baseline network C3D with and without ACTF. We first perform action recognition with only the temporal feature extracted through our proposed ACTF. The low-level feature is extracted through conv5 layer of the C3D network. We denote this network as C3D-single-ACTF. We then perform action recognition by attentively combining the ACTF feature as well as the Spatial-Temporal Pooled feature which is extracted from *pool5* layer of C3D, similar to the implementation of MFNet-ACTF. We denote this modification as C3D-ACTF. The top-1 accuracy of the networks are shown in Table 2.

The results in Table 2 clearly show that applying our proposed framework in the
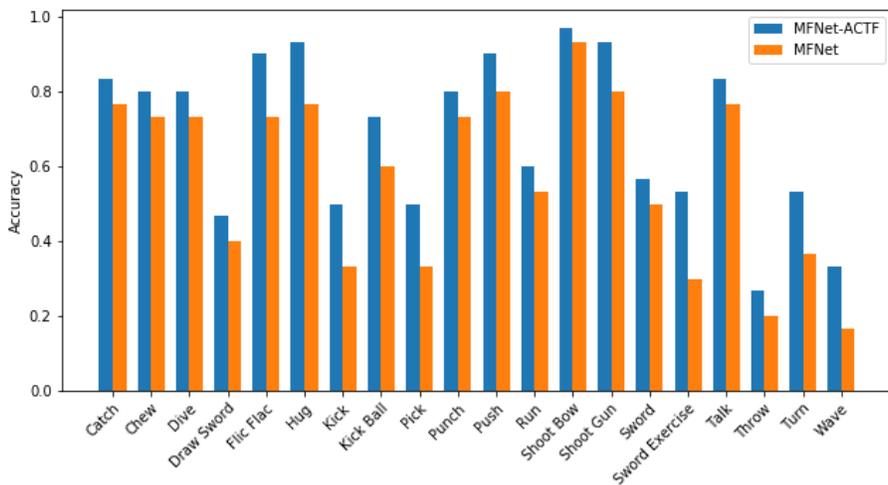
Figure 5: Accuracy comparisons of 20 classes on split 1 of the HMDB-51 between our proposed MFNet-ACTF network and the original MFNet network.

C3D network improves the accuracy of the baseline C3D network. Even by only utilizing the extracted ACTF feature as temporal feature, we obtain an improvement of 2.4%. This shows that our ACTF feature effectively represent the temporal pattern in the video and thus lead to better results. A larger gain of 3.7% is achieved when the extracted ACTF feature is used with the Spatial-Temporal Pooled feature. The results are consistent with that shown in Table 1, where using ACTF feature improves the accuracy of MFNet. This suggests that our proposed framework is generic, and can be used with other baselines.

We further investigate the improvement of performance over different actions and present the comparison of performance between our proposed MFNet-ACTF network and the original MFNet network. Figure 5 shows the accuracy of 20 classes from split-1 of the HMDB51 dataset, where our network outperforms the original network by a noticeable margin. It is worth noticing that for actions with similar spatial appearance but different actions, e.g. "Sword" and "Sword Exercise", our network performs significantly better than the original network. Our network obtains a 23.3% performance gain over the original MFNet on the action class "Sword Exercise". The large perfor-

**Ground Truth: Sword Exercise** | **Ground Truth: Wave** | **Ground Truth: Turn** | **Ground Truth: Kick Ball**

| Network 1: MFNet-ACTF (Ours) | Network 1: MFNet-ACTF (Ours) | Network 1: MFNet-ACTF (Ours) | Network 1: MFNet-ACTF (Ours) |
|---|---|---|---|
| **Sword Exercise  0.3703** | **Wave  0.4595** | **Turn  0.7219** | **Kick Ball  0.3703** |
| Draw Sword  0.2566 | Pick  0.1001 | Sit  0.1849 | Jump  0.2566 |
| Sword  0.2198 | Fall Floor  0.0895 | Walk  0.0237 | Cartwheel  0.2198 |
| Fencing  0.0506 | Throw  0.0605 | Stand  0.0230 | Flic Flac  0.0506 |
| Hit  0.0220 | Run  0.0447 | Kiss  0.0064 | Run  0.0220 |
| Network 2: MFNet | Network 2: MFNet | Network 2: MFNet | Network 2: MFNet |
| Draw Sword  0.5181 | Climb  0.5401 | Drink  0.3571 | Dive  0.5181 |
| Sword  0.2437 | Wave  0.1367 | Turn  0.2485 | Jump  0.2437 |
| Handstand  0.1190 | Pick  0.1296 | Sit  0.1320 | Kick Ball  0.1190 |
| Sword Exercise  0.0466 | Run  0.0778 | Walk  0.1169 | Cartwheel  0.0466 |
| Hit  0.0170 | Talk  0.0572 | Brush Hair  0.0361 | Flic Flac  0.0170 |

Figure 6: Examples from HMDB51 dataset where our proposed MFNet-ACTF succeeds in recognizing the action while the original MFNet fails.

mance gain proves the effectiveness of the additional temporal feature extracted as the ACTF feature in improving the complete video representation. Several examples from HMDB51 dataset is presented in Figure6 where our proposed MFNet-ACTF could accurately recognize the respective actions while the original MFNet network could not. It could be observed that the spatial features of the given examples, or more intuitively the appearance of the given examples, could not provide effective representation for accurate action recognition. For example, for the first video, the scenario as shown could be present in action classes "Sword", in which most videos present people fighting with a sword, and "Draw Sword", in which videos present the action of a sword drawn out. The difference between these action classes could only be determined through the temporal feature instead of the spatial feature. Thus the original network which can only extract the spatial feature of the video cannot distinguish the actions correctly while our proposed framework succeeds in recognizing the different actions.

*4.3. Ablation Study*

In this section, we justify our proposed design of the ACTF feature through ablation study. Specifically, we examine the performance of our proposed ICCF and the ACTF

| Method | Top-1 HMDB51 split-1 |
| --- | --- |
| MFNet | 70.8% |
| MFNet-ICCF | 72.6% |
| MFNet-single-ACTF | 72.9% |

Table 3: Comparison of the network architectures that use only temporal feature for action recognition.

| Method | Top-1 HMDB51 split-1 |
| --- | --- |
| MFNet-ACTF-no-attn | 72.5% |
| MFNet-attn@ACTF | 73.3% |
| MFNet-attn@final | 73.0% |
| MFNet-ACTF | 73.6% |

Table 4: Comparison of the network architectures that use all or partial attentive concatenation.

feature separately. We then examine the performance of the attention mechanisms used to combine the different modules of our proposed generic action recognition framework as discussed in Section 3. All experiments conducted in our ablation study are performed on split 1 of the HMDB51 dataset. We set our batch size to 16 and conduct the experiment using one NVIDIA TITAN Xp GPU. The much smaller batch size is a key reason of the lower accuracy reported than that in Table 1.

We instantiate MFNet to justify our proposed ICCF and the ACTF feature introduced in Section 3.2 and utilize only temporal feature for action recognition. First the proposed ICCF is extracted as our temporal feature. The network that utilizes only ICCF is denoted as MFNet-ICCF. We then employ the ACTF feature as our temporal feature. Similar to the previous denotation, the network that utilizes only the ACTF feature is denoted as MFNet-single-ACTF. The comparison of the performances of these two networks with the baseline MFNet is shown in Table 3.

The result in Table 3 shows that by utilizing only temporal feature, even with only bilinear inter-frame correlation, the performance of the network is improved by a margin of $1.8\%$, indicating that utilizing inter-frame correlation information helps to extract high-quality temporal feature of the video. The improvement achieved by using

high quality temporal feature over feature obtained from spatial-temporal pooling coincides with findings in preceding works (Wang et al., 2016; Carreira & Zisserman, 2017; Feichtenhofer et al., 2016). However, our temporal feature is obtained from RGB input through inter-frame correlation rather than using optical flow.

Table 3 shows further improvement when we employ the ACTF feature. As described in Section 3.2, the ACTF feature is a weighted combination of ICCF, which is a bilinear inter-frame correlation feature, and IMF, which is a linear inter-frame correlation feature. This result proves that the bilinear inter-frame correlation feature and linear inter-frame feature complements each other.

To better combine the features extracted from different modules, we introduced attentive concatenation of features as mentioned in Section 3.3. Here we justify the need for utilizing attentive concatenation of features. Table 4 presents the comparison between the networks that utilize attentive concatenation at every step and the networks that partially or do not utilize attentive concatenation for feature combination. Here MFNet-ACTF-no-attn denotes the network where all feature combination utilizes direct concatenation instead of attentive concatenation. Meanwhile, MFNet-ACTF denotes the network utilizing our proposed temporal feature extraction framework with attentive concatenation at every step of feature combination. MFNet-attn@ACTF denotes the network that performs temporal-wise attentive concatenation when constructing ICCF, and feature-wise attentive concatenation of ICCF and IMF as shown in Figure 3. The concatenation of the temporal feature and spatial feature is by direct concatenation. Similarly, MFNet-attn@final denotes that attentive concatenation is adopted only for ACTF and Spatial-Temporal Pooled feature combination while direct concatenation is adopted at other stages.

The result given in Table 4 clearly shows the advantage of adopting attentive concatenation for feature combination. We note that if the network combines features with only direct concatenation, its performance would be even worse than that of MFNet-single-ACTF, whose ACTF feature is constructed with attentive concatenation of ICCF and IMF features. The performance is improved even when attentive concatenation is used in some stages of feature combination only. It can be observed that applying attentive concatenation at different stages complements each other, with over $1\%$ im-
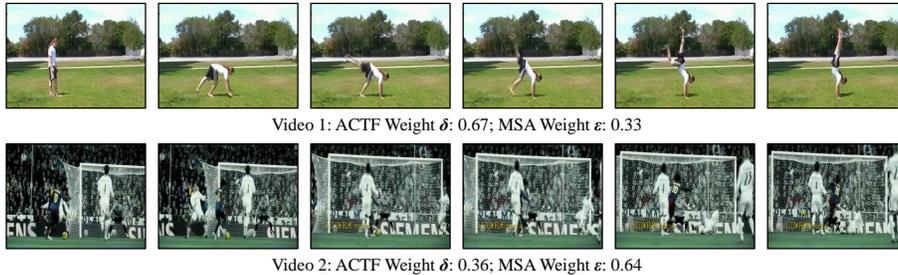
Video 1: ACTF Weight $\delta$: 0.67; MSA Weight $\varepsilon$: 0.33



Video 2: ACTF Weight $\delta$: 0.36; MSA Weight $\varepsilon$: 0.64

Figure 7: The weights of ACTF feature $\delta$ and the weights of Spatial-Temporal Pooled feature $\epsilon$ for two videos. Attentive concatenation learns these weights dynamically.

provement made when all stages adopt attentive concatenation.

We also investigate the weights $\delta$ and $\epsilon$ on the ACTF feature and the Spatial-Temporal Pooled feature for different videos. Figure 7 shows examples where either temporal feature or spatial feature dominates the feature combination process. Video 2 shows a video where the Spatial-Temporal Pooled feature, or the spatial feature, dominates the feature combination. These videos tend to have clear visual characteristics, such as the soccer goal that appears in most videos describing the sport soccer. The appearance of these videos are therefore sufficient for action recognition, and dominate the feature combination process. By contrast, feature combination in Video 1 is dominated by ACTF feature, which is the temporal feature. We observe that similar videos tend to have actions that would mix up with other categories. In this case, the handstand action is similar to actions that may occur in diving or in somersault, where a person would also go upside down. Also, there is no iconic background items in Video 1. For these videos, the temporal features dominate the feature combination, thus has a larger weight $\delta$. The different weights with respect to the different videos could prove that adopting attentive concatenation could attend to the more important feature which is related to the characteristic of the video itself.

## 5. Conclusion

In this work, we propose a new method for extracting the temporal feature of a video while avoiding the use of optical flow. The new temporal feature namely At-

tentive Correlated Temporal Feature (ACTF) is an attentive combination of both bilinear inter-frame correlation and linear inter-frame correlation features. The bilinear inter-frame correlation feature is extracted through a bilinear operation with respect to successive frames within a certain region, while the linear inter-frame feature is extracted through inter-frame temporal pooling. For overall evaluation on UCF101 and HMDB51, our method obtains state-of-the-art results when instantiating MFNet combined with our ACTF feature. We verify our design through thorough ablation study, and then further demonstrate that the proposed feature can be introduced to other similar action recognition networks instead of using optical flow.

## References

## References

Bottou, L. (2010). Large-scale machine learning with stochastic gradient descent. In *Proceedings of COMPSTAT'2010* (pp. 177–186). Springer.

Buades, A., Coll, B., & Morel, J.-M. (2005). A non-local algorithm for image denoising. In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)* (pp. 60–65). IEEE volume 2.

Carmona, J. M., & Climent, J. (2018). Human action recognition by means of subtensor projections and dense trajectories. *Pattern Recognition*, *81*, 443–455.

Carreira, J., & Zisserman, A. (2017). Quo vadis, action recognition? a new model and the kinetics dataset. In *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 6299–6308).

Chen, Y., Kalantidis, Y., Li, J., Yan, S., & Feng, J. (2018). Multi-fiber networks for video recognition. In V. Ferrari, M. Hebert, C. Sminchisescu, & Y. Weiss (Eds.), *Computer Vision – ECCV 2018* (pp. 364–380). Cham: Springer International Publishing.

Danafar, S., & Gheissari, N. (2007). Action recognition for surveillance applications using optic flow and svm. In *Asian Conference on Computer Vision* (pp. 457–466). Springer.

Diba, A., Sharma, V., & Van Gool, L. (2017). Deep temporal linear encoding networks. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition* (pp. 2329–2338).

Dosovitskiy, A., Fischer, P., Ilg, E., Husser, P., Hazirbas, C., Golkov, V., v. d. Smagt, P., Cremers, D., & Brox, T. (2015). Flownet: Learning optical flow with convolutional networks. In *2015 IEEE International Conference on Computer Vision (ICCV)* (pp. 2758–2766). doi:`10.1109/ICCV.2015.316`.

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

Feichtenhofer, C., Pinz, A., & Wildes, R. (2016). Spatiotemporal residual networks for video action recognition. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, & R. Garnett (Eds.), *Advances in Neural Information Processing Systems 29* (pp. 3468–3476). Curran Associates, Inc.

Feichtenhofer, C., Pinz, A., & Zisserman, A. (2016). Convolutional two-stream network fusion for video action recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (pp. 1933–1941). doi:`10.1109/CVPR.2016.213`.

Gao, Y., Beijbom, O., Zhang, N., & Darrell, T. (2016). Compact bilinear pooling. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (pp. 317–326). doi:`10.1109/CVPR.2016.41`.

Girdhar, R., Ramanan, D., Gupta, A., Sivic, J., & Russell, B. (2017). Actionvlad: Learning spatio-temporal aggregation for action classification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 971–980).

Hara, K., Kataoka, H., & Satoh, Y. (2018). Can spatiotemporal 3d cnns retrace the history of 2d cnns and imagenet? In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition* (pp. 6546–6555).

He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (pp. 770–778). doi:`10.1109/CVPR.2016.90`.

Herath, S., Harandi, M., & Porikli, F. (2017). Going deeper into action recognition: A survey. *Image and Vision Computing*, *60*, 4 – 21. URL: `http://www.sciencedirect.com/science/article/pii/S0262885617300343`. doi:`https://doi.org/10.1016/j.imavis.2017.01.010`. Regularization Techniques for High-Dimensional Data Analysis.

Ilg, E., Mayer, N., Saikia, T., Keuper, M., Dosovitskiy, A., & Brox, T. (2017). Flownet 2.0: Evolution of optical flow estimation with deep networks. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (pp. 1647–1655). doi:`10.1109/CVPR.2017.179`.

Kardas, K., & Cicekli, N. K. (2017). Svas: Surveillance video analysis system. *Expert Systems with Applications*, *89*, 343 – 361. URL: `http://www.sciencedirect.com/science/article/pii/S0957417417305286`. doi:`https://doi.org/10.1016/j.eswa.2017.07.051`.

Kay, W., Carreira, J., Simonyan, K., Zhang, B., Hillier, C., Vijayanarasimhan, S., Viola, F., Green, T., Back, T., Natsev, P. et al. (2017). The kinetics human action video dataset. *arXiv preprint arXiv:1705.06950*, .

Keeli, A. S. (2018). Viewpoint projection based deep feature learning for single and dyadic action recognition. *Expert Systems with Applications*, *104*, 235 – 243. URL: `http://www.sciencedirect.com/science/article/pii/S0957417418301933`. doi:`https://doi.org/10.1016/j.eswa.2018.03.047`.

Kuehne, H., Jhuang, H., Garrote, E., Poggio, T., & Serre, T. (2011). Hmdb: a large video database for human motion recognition. In *2011 International Conference on Computer Vision* (pp. 2556–2563). IEEE.

Lee, J. S., Choi, S., & Kwon, O. (2017). Identifying multiuser activity with overlapping acoustic data for mobile decision making in smart home environments. *Expert Systems with Applications*, *81*, 299 – 308. URL: `http://www.sciencedirect.com/science/article/pii/S0957417417302257`. doi:`https://doi.org/10.1016/j.eswa.2017.03.062`.

Li, H., & Suen, C. Y. (2016). A novel non-local means image denoising method based on grey theory. *Pattern Recognition*, *49*, 237–248.

Li, W., Chen, D., Lv, Z., Yan, Y., & Cosker, D. (2018). Learn to model blurry motion via directional similarity and filtering. *Pattern Recognition*, *75*, 327–338.

Li, X., Ye, M., Liu, Y., Zhang, F., Liu, D., & Tang, S. (2017). Accurate object detection using memory-based models in surveillance scenes. *Pattern Recognition*, *67*, 73–84.

Lin, T.-Y., RoyChowdhury, A., & Maji, S. (2015). Bilinear cnn models for fine-grained visual recognition. In *Proceedings of the IEEE international conference on computer vision* (pp. 1449–1457).

Liu, K., Liu, W., Gan, C., Tan, M., & Ma, H. (2018). T-c3d: Temporal convolutional 3d network for real-time action recognition. In *Thirty-second AAAI conference on artificial intelligence*.

Lundstrm, J., Jrpe, E., & Verikas, A. (2016). Detecting and exploring deviating behaviour of smart home residents. *Expert Systems with Applications*, *55*, 429 – 440. URL: `http://www.sciencedirect.com/science/article/pii/S0957417416300616`. doi:`https://doi.org/10.1016/j.eswa.2016.02.030`.

Mabrouk, A. B., & Zagrouba, E. (2018). Abnormal behavior recognition for intelligent video surveillance systems: A review. *Expert Systems with Applications*, *91*, 480

– 491. URL: `http://www.sciencedirect.com/science/article/pii/S0957417417306334`. doi:`https://doi.org/10.1016/j.eswa.2017.09.029`.

Ortis, A., Farinella, G. M., DAmico, V., Addesso, L., Torrisi, G., & Battiato, S. (2017). Organizing egocentric videos of daily living activities. *Pattern Recognition*, *72*, 207–218.

Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., & Lerer, A. (2017). Automatic differentiation in PyTorch. In *NIPS Autodiff Workshop*.

Pham, N., & Pagh, R. (2013). Fast and scalable polynomial kernels via explicit feature maps. In *Proceedings of the 19th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* KDD '13 (pp. 239–247). New York, NY, USA: ACM. URL: `http://doi.acm.org/10.1145/2487575.2487591`. doi:`10.1145/2487575.2487591`.

Piergiovanni, A., & Ryoo, M. S. (2019). Representation flow for action recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 9945–9953).

Sahoo, S. P., & Ari, S. (2019). On an algorithm for human action recognition. *Expert Systems with Applications*, *115*, 524 – 534. URL: `http://www.sciencedirect.com/science/article/pii/S0957417418305220`. doi:`https://doi.org/10.1016/j.eswa.2018.08.014`.

Simonyan, K., & Zisserman, A. (2014). Two-stream convolutional networks for action recognition in videos. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, & K. Q. Weinberger (Eds.), *Advances in Neural Information Processing Systems 27* (pp. 568–576). Curran Associates, Inc.

Soomro, K., Zamir, A. R., & Shah, M. (2012). Ucf101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, .

Tran, D., Bourdev, L., Fergus, R., Torresani, L., & Paluri, M. (2015). Learning spatiotemporal features with 3d convolutional networks. In *Proceedings of the 2015 IEEE International Conference on Computer Vision (ICCV)* ICCV '15 (pp. 4489–4497). Washington, DC, USA: IEEE Computer Society. URL: `http://dx.doi.org/10.1109/ICCV.2015.510`. doi:`10.1109/ICCV.2015.510`.

Tran, D., Ray, J., Shou, Z., Chang, S.-F., & Paluri, M. (2017). Convnet architecture search for spatiotemporal feature learning. *arXiv preprint arXiv:1708.05038*, .

Tran, D., Wang, H., Torresani, L., Ray, J., LeCun, Y., & Paluri, M. (2018). A closer look at spatiotemporal convolutions for action recognition. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition* (pp. 6450–6459).

Varol, G., Laptev, I., & Schmid, C. (2018). Long-term temporal convolutions for action recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, *40*, 1510–1517. doi:`10.1109/TPAMI.2017.2712608`.

Wang, H., & Wang, L. (2018). Learning content and style: Joint action recognition and person identification from human skeletons. *Pattern Recognition*, *81*, 23–35.

Wang, K., Zeng, X., Yang, J., Meng, D., Zhang, K., Peng, X., & Qiao, Y. (2018a). Cascade attention networks for group emotion recognition with face, body and image cues. In *Proceedings of the 20th ACM International Conference on Multimodal Interaction* ICMI '18 (pp. 640–645). New York, NY, USA: ACM. URL: `http://doi.acm.org.ezlibproxy1.ntu.edu.sg/10.1145/3242969.3264991`. doi:`10.1145/3242969.3264991`.

Wang, L., Xiong, Y., Wang, Z., Qiao, Y., Lin, D., Tang, X., & Van Gool, L. (2016). Temporal segment networks: Towards good practices for deep action recognition. In *European conference on computer vision* (pp. 20–36). Springer.

Wang, X., Girshick, R., Gupta, A., & He, K. (2018b). Non-local neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 7794–7803).

Wu, C., Khalili, A. H., & Aghajan, H. (2010). Multiview activity recognition in smart homes with spatio-temporal features. In *Proceedings of the fourth ACM/IEEE international conference on distributed smart cameras* (pp. 142–149). ACM.

Xiang, T., & Gong, S. (2008). Activity based surveillance video content modelling. *Pattern Recognition*, *41*, 2309–2326.

Yang, H., Yuan, C., Li, B., Du, Y., Xing, J., Hu, W., & Maybank, S. J. (2019). Asymmetric 3d convolutional neural networks for action recognition. *Pattern Recognition*, *85*, 1–12.

Yang, J., Zou, H., Jiang, H., & Xie, L. (2018). Carefi: Sedentary behavior monitoring system via commodity wifi infrastructures. *IEEE Transactions on Vehicular Technology*, *67*, 7620–7629.

Yang, J., Zou, H., Jiang, H., & Xie, L. (2018). Device-free occupant activity sensing using wifi-enabled iot devices for smart homes. *IEEE Internet of Things Journal*, *5*, 3991–4002. doi:`10.1109/JIOT.2018.2849655`.

Zach, C., Pock, T., & Bischof, H. (2007). A duality based approach for realtime tv-l1 optical flow. In F. A. Hamprecht, C. Schnörr, & B. Jähne (Eds.), *Pattern Recognition* (pp. 214–223). Berlin, Heidelberg: Springer Berlin Heidelberg.

Zhao, M., Li, T., Abu Alsheikh, M., Tian, Y., Zhao, H., Torralba, A., & Katabi, D. (2018a). Through-wall human pose estimation using radio signals. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 7356–7365).

Zhao, Y., Xiong, Y., & Lin, D. (2018b). Recognize actions by disentangling components of dynamics. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 6566–6575).

Zhu, Y., Lan, Z., Newsam, S., & Hauptmann, A. G. (2017). Hidden two-stream convolutional networks for action recognition. *arXiv preprint arXiv:1704.00389*, .