

Effective Action Recognition with Embedded Key Point Shifts

Haozhi Cao^{a,1}, Yuecong Xu^{a,1,*}, Jianfei Yang^{a,1}, Kezhi Mao^{a,1}, Jianxiong Yin^{b,1}, Simon See^{b,1}

^a50 Nanyang Avenue, 639798, Singapore

^b3 International Business Park Rd, #01-20A Nordic European Centre, 609927, Singapore

Abstract

Temporal feature extraction is an essential technique in video-based action recognition. Key points have been utilized in skeleton-based action recognition methods but they require costly key point annotation. In this paper, we propose a novel temporal feature extraction module, named Key Point Shifts Embedding Module (*KPSEM*), to adaptively extract channel-wise key point shifts across video frames without key point annotation for temporal feature extraction. Key points are adaptively extracted as feature points with maximum feature values at split regions, while key point shifts are the spatial displacements of corresponding key points. The key point shifts are encoded as the overall temporal features via linear embedding layers in a multi-set manner. Our method achieves competitive performance through embedding key point shifts with trivial computational cost, achieving the state-of-the-art performance of 82.05% on Mini-Kinetics and competitive performance on UCF101, Something-Something-v1 and HMDB51 datasets.

Keywords: Action Recognition, Temporal Feature, Key Point Shifts

*Corresponding author

Email addresses: haozhi001@e.ntu.edu.sg (Haozhi Cao), 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)

¹School of Electrical and Electronic Engineering, Nanyang Technological University.

²NVIDIA AI Tech Centre.

1. Introduction

Action recognition has attracted interest in vision and machine learning communities [1, 2] thanks to its applications such as surveillance [3] and smart homes [4]. Besides the spatial information in each video frame, videos contain additional temporal information with different properties. Thus, effective modelling of temporal information in videos is critical for accurate action recognition. Intuitively, we humans could effectively recognize the temporal information of action through the key points of actors while ignoring the unrelated environment. Particularly, we focus on the shifts of an actor’s key points. For instance, in Figure 1a and Figure 1b, the actors in different environments are both climbing. We could still recognize the same action for both actors through the movement of the actors’ feet and elbows according to key point shifts.

Previous methods used in the skeleton-based action recognition task [5, 6] have proven the effectiveness of using key points and their shifts for recognizing human actions, where the key points and their correspondence across frames are provided as inputs. Skeleton data outlines the key points of actors in the video, suppressing the influence of trivial environment information during the feature extraction process. This increases the robustness of networks when encountering actions with complicated backgrounds. However, additional costs, such as depth sensors [7] or estimation algorithms [8, 9], are needed to annotate the raw video data with human key points for supervision.

Currently, most methods for action recognition without the annotation of key points would not utilize key point shifts information at all. Instead of learning the action’s temporal features, many of these methods may learn more about environment knowledge. This is due to the fact that these methods tend to use all pixels in each frame instead of key points and their shifts, while most of the pixels correspond to the environment information which is not as important as action information. Most of these methods fall into three categories: (1) two-stream CNNs [10, 11, 12], (2) 3D CNNs [13, 14, 15, 16] and (3) CNNs with learnable feature correlations [17, 18]. Two-stream CNNs model temporal fea-

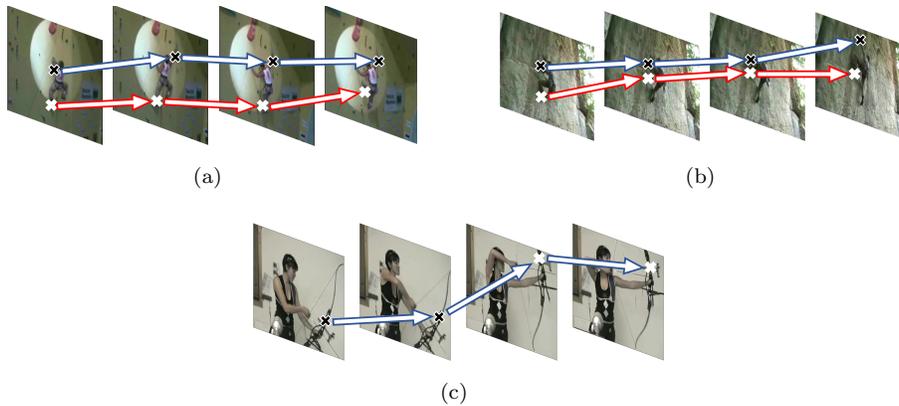


Figure 1: Illustration of using key points and their shifts for action recognition. The key point shifts in Figure 1a and Figure 1b imply the action “Climbing” while the key point shifts on Figure 1c imply the action “Archery”. Figure 1a and Figure 1b show that key points and their shifts could be distributed in different locations and the distributions could be different at each frame. Figure 1c illustrates an example of the cross-region key point shifts. Crosses and arrows in different colours indicate key points in different regions and different key point shifts. The pictures are best viewed in colour and zoomed in.

tures by inputting pre-computed hand-crafted features, such as optical flow, to a CNN. 3D CNNs extract spatiotemporal features jointly by expanding convolution kernels of 2D CNNs to the temporal dimension but they exhibit inferior performances. More recently, learnable correlations of features across frames [17, 18] are used for temporal modelling. The temporal features are extracted by exploiting the multiplicative interactions between each pixel.

To utilize key point shifts for temporal feature extraction without key point annotation, we propose a novel method to extract and embed key point shifts information of each channel from the high-level feature maps. When high-level spatial features of each frame are extracted through CNN layers, the key points related to the action could be viewed as the points with the maximum feature value of each channel in the high-level feature maps. In addition, in many actions, such key points and their shifts would be distributed in different local regions of the video frames and the distributions of key points would also be slightly different for each frame. For example, for the action “Climbing” in

Figure 1b, the shift of the key points corresponding to the elbows and the feet are located at the upper and lower parts, respectively, at the first frame, while both key points are positioned upwards at the last frame. Therefore, to obtain temporal information with respect to the key point shifts at different locations and to cope with the different distributions of key points across all the frames, we propose to split each frame into different regions adaptively with the key points extracted at each region.

Furthermore, for some actions, the key points may not belong to the same region across all the frames. Figure 1c depicts such a case where the key point of the bow locates at the lower region at the first two frames while it shifts to the upper region at the last two frames. To extract the cross-region key point shifts correctly, we compute shift weights to indicate the similarity between any pair of key points across adjacent frames. The key point shifts are then calculated as the spatial displacements between the corresponding key points in adjacent frames according to the shift weights. The proposed temporal feature extractor is termed as Key Point Shifts Embedding Module (*KPSEM*), utilizing key point shifts extracted regionally, termed as Regional Key Point Shifts (*RKPS*). Multiple *RKPS*s are obtained through multiple sets of Adaptive Regional Shift Extractor (*AReSE*) under different region separations. The resulting *RKPS*s are encoded through independent embedding layers to constitute more robust temporal features.

In summary, the main contribution of this work is a novel temporal feature extraction module based on key point shifts: Key Point Shifts Embedding Module (*KPSEM*). First, *KPSEM* is designed to utilize the key point shifts for effective temporal feature extraction without key point annotation. Second, through a multi-set embedding operation, the key point shifts are embedded as effective temporal features of the input video. Third, the extensive experiments on various datasets demonstrate that *KPSEM* can effectively model temporal features, achieving state-of-the-art performance on the Mini-Kinetics dataset without involving high computational cost.

2. Related Work

Temporal Feature Extraction. To extract temporal features effectively, earlier works [11, 10, 12] adopt a two-stream strategy where temporal features are extracted in parallel with the spatial features. The temporal features are extracted by feeding a stack of optical flow frames to CNNs. Typical computation methods of optical flow include LucasKanade [19], HornSchunck [20] and TV-L1 [21]. More recent two-stream CNNs usually apply TV-L1 [21] as the optical flow extraction method due to its robustness and efficiency compared to other methods. Optical flow represents temporal features accurately as it computes pixel-level correlation information across frames. However, the application of optical flow usually forbids end-to-end training of the network, since it requires pre-computation of optical flow before being input to CNNs. Additionally, the process of extracting optical flow is computationally expensive and memory intensive. Therefore, more recent methods try to avoid the need for optical flow.

To address the limitations imposed by utilizing optical flow in two-stream CNNs, later works proposed to extract temporal features jointly with spatial features using 3D CNNs. C3D [13], I3D [15], P3D [22] and 3D-ResNet [14] all belong to this category. C3D [13] is one of the primary works where CNN filters are expanded to the temporal dimension. For faster training, I3D [15] directly inflates 2D CNNs into a 3D structure through endowing filters and pooling kernels with the temporal dimension. Additionally, P3D [22] reduces the computational cost by simulating a 3D convolution filter with a spatial convolution filter and a separate temporal convolution filter. Subsequent networks, such as 3D-ResNet [14], are deeper and larger 3D CNNs trained on the Kinetics [23] to retrace the success of deeper 2D CNNs pretrained on ImageNet [24]. 3D CNNs benefit from end-to-end training and require only RGB input. However, many of these works exhibit inferior results compared to two-stream CNNs. The inferior results could be contributed by the fact that the temporal features are extracted through multiple pooling operations along the temporal dimension. The temporal information which reflects the change in spatial information across time

might be lost during the multiple pooling operations. Therefore, 3D CNNs fail to extract effective temporal features, which results in inferior performance.

To improve the effectiveness of temporal features extracted through CNNs while avoiding the use of optical flow, multiple methods have been proposed. A prominent category is to utilize learnable correlations of features across frames. Inspired by the non-local mean operation for image denoising [25, 26], Wang *et al.* [17] presented the non-local operation to capture correlations on the pixel level as the representation of the temporal features. Similarly, ARTNet [18] is designed such that its relation branch captures the multiplicative interactions between pixels across multiple frames. Recently, the correlation network [27] utilizes correlation operators to model frame-to-frame correlation in feature maps. The above methods model correlations between all pixels across frames as the temporal features with end-to-end training. However, it is computationally expensive to extract correlations between every single pixel. Moreover, for pixels representing unrelated background information, their correlations might contribute trivial to the overall video features and therefore the correlations computed between these pixels are redundant.

In addition to using learnable feature correlations to extract temporal features, another proposed solution is to enrich the temporal information by extracting temporal features under different frame rates. Specifically, the Slowfast network [28] applies an additional “Fast” pathway, which is a 3D CNN stream with a higher frame rate to capture temporal features in fine temporal resolution. Similarly, the TPN [29] extracts temporal features by aggregating features extracted under different tempos. These methods aim to extract more temporal information under multiple frame rates. However, they usually include multiple sub-streams for the different frame rates, which might introduce extra network structures and result in computation overhead.

Different from two-stream methods, our work proposed an end-to-end method to extract frame-wise temporal information from pure RGB input. In contrast to previous correlation extraction methods, our proposed method aims to extract the temporal information effectively by extracting key correlations instead

of correlations of all pixels or features. Our proposed method utilizes only regional key points and computes their respective shift across successive frames. This implies that our temporal feature extraction method is more computation efficient with less redundant information included. Yet it still brings consistent improvement in temporal feature extraction, supported by stable improvement in action recognition accuracy.

Skeleton-based Action Recognition. Key points and their displacements have been used mainly in skeleton-based action recognition [5, 6]. Most of the skeleton-based methods take skeleton data as the input, which is generated by devices or pose estimation algorithms in the form of 2D or 3D coordinates. These skeleton data already possesses temporal correspondent relationship. The utilization of skeleton-based data excludes the effect of unrelated pixels so that the network can concentrate on the key points and their temporal correlation. Compared to extracting correlations of all pixels or features, the introduction of skeleton data is more effective and efficient in modelling temporal information.

However, such annotations are not available in the first place for most public videos. The generation of skeleton data requires extra cost, such as extra computation resources and additional recording devices. In comparison, we propose a novel method to extract key points in high-level feature maps without the need for collecting and annotating the skeleton data for videos. Empirical results show that our method can bring consistent improvement while resulting in only a trivial amount of extra computational complexity without the need for skeleton data or other key point annotation.

3. Proposed Method

The ultimate goal of our work is to extract temporal information that can represent the frame-wise movement of spatial features. We believe that key points represent the dominant spatial features, such as parts of human body and objects related to the actions as shown in Figure 1. Therefore, their shifts can be viewed as effective temporal features. We propose a novel module, Key

Point Shifts Extraction Module (*KPSEM*), which utilizes Regional Key Point Shifts (*RKPS*) as a new modality to represent the temporal features of videos. The *RKPS*s are extracted by multiple sets of Adaptive Regional Shift Extractor (*AReSE*) and computed as the relative coordinate displacements between key points. In this section, we first illustrate how the *KPSEM* is implemented with the CNN backbone as well as its detailed structure. Subsequently, we expound the details of the *AReSE* with the process of extracting *RKPS*, which is the core of our proposed *KPSEM*.

3.1. Overall Structure

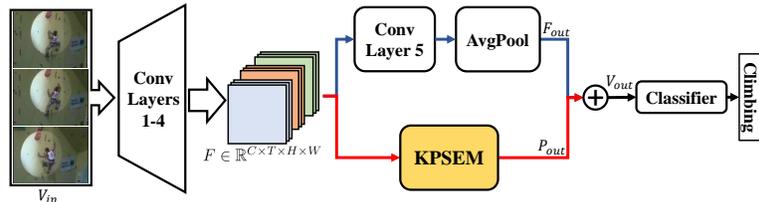


Figure 2: The overall structure of our proposed method. *KPSEM* extracts the temporal features of the video through high-level spatial features F extracted from the CNN backbone. The spatial feature output F_{out} and temporal feature output P_{out} are concatenated and the resulting overall video features V_{out} are passed through a linear classifier. *KPSEM* can be inserted not only at the location shown in this figure but also after any convolution layer.

The overall structure of our network is as shown in Figure 2. Given an input video V_{in} , we utilize CNN to extract its high-level feature maps, denoted as $F \in \mathbb{R}^{C \times T \times H \times W}$ where C , T , H , W are the number of channels, the length of frames, the spatial height and the width of the feature maps, respectively. Our proposed *KPSEM* extracts temporal features from F by embedding *RKPS*s which are the weighted coordinate differences extracted by *AReSE*. The temporal features from *KPSEM* concatenate with the spatial features from CNN and the overall video features pass through the linear classifier. The spatial feature output from the CNN is denoted as $F_{out} \in \mathbb{R}^{C_s}$ while the temporal feature output of *KPSEM* is denoted as $P_{out} \in \mathbb{R}^{C_p}$. Here C_s , C_p denote the number of channels of the spatial feature output and that of temporal feature output from *KPSEM*,

respectively. Subsequently, the overall video features V_{out} are computed by:

$$V_{out} = F_{out} \oplus P_{out}, \quad (1)$$

where \oplus denotes the concatenation operation along the channel dimension. The combined $V_{out} \in \mathbb{R}^{(C_s+C_p)}$ passes through the final classifier. It is worth noting that while the $KPSEM$ is inserted after Conv4 in Fig 2, the $KPSEM$ is an isolated block which can be inserted at any other location.

3.2. Key Point Shifts Embedding Module ($KPSEM$)

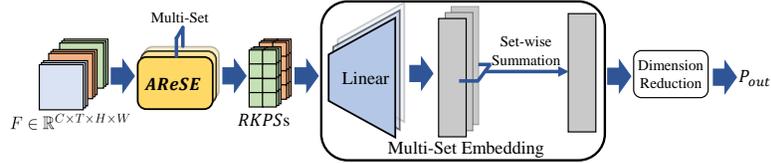


Figure 3: Details of $KPSEM$. From the high-level features F extracted by CNN, G $RKPSs$ are extracted by G sets of $AReSEs$. Each $RKPS$ is linearly projected to a higher dimension vector separately. The resulting multi-set embedding is then processed through a set-wise summation to combine all the information from the different sets. The overall temporal feature output is obtained through a final dimension reduction process.

Our proposed $KPSEM$ module extracts the overall temporal features by embedding the Regional Key Point Shifts ($RKPS$) obtained through $AReSE$ (cf. Section 3.3). The $AReSE$ performs an adaptive region separation on high-level feature maps and subsequently extracts the $RKPS$ based on the region separation. Different region separations generate different $RKPSs$. To enhance the robustness of $KPSEM$, we adopt multiple sets of $AReSEs$, each of which performs region separation and $RKPS$ extraction independently. Formally, the $RKPSs$ are computed as:

$$RKPS_g = AReSE_g(F), \quad g = 1, 2, \dots, G, \quad (2)$$

where $RKPS_g$ is the g^{th} $RKPS$ extracted by the g^{th} set of $AReSE$, denoted as $AReSE_g$. G is the total number of sets of $AReSEs$. The resulting $RKPSs$

are of size $C \times (T - 1) \times G \times 2$. The details about how *AReSE* extracts *RKPS* is illustrated in Section 3.3.

Subsequently, inspired by the multi-head attention introduced in [30], instead of going through a shared linear embedding layer, the G sets of *RKPS*s are linearly projected to dimension d_e separately. This independent embedding operation ensures that our proposed *KPSEM* can gain more abundant representation from different *RKPS*s. A set-wise summation is then utilized to aggregated information across all the *RKPS*s. More specifically, given G sets of *RKPS*s, the multi-set embedding is computed as:

$$\mathcal{M}(RKPSs) = \sum_{g=1}^G \mathcal{L}_g(RKPS_g), \quad (3)$$

where \mathcal{M} denotes the multi-set embedding operation. $RKPS_g$ is the g^{th} *RKPS* computed by *AReSE_g*, while \mathcal{L}_g is the g^{th} linear projection for the corresponding *RKPS_g*. The subsequent set-wise summation across all G *RKPS*s directly merges all embeddings. Ultimately, after a dimension reduction procedure, the overall temporal features are obtained.

We notice that the multi-set mechanism mentioned above plays an important role in the *KPSEM* module. The utilization of multiple sets of *AReSE* provides the overall framework with different adaptive regional feature maps. This allows the network to find better local splits with more representative key points. The separated linear embedding layer for each *RKPS* ensures the abundance of embedded information obtained from the different *RKPS*s. Next, we illustrate how the proposed *AReSE* conducts adaptive region separation and key point shift extraction in details.

3.3. Adaptive Regional Shift Extractor (*AReSE*)

As mentioned in Section 1, key point shifts can represent the temporal movement of dominant features of each channel. We propose a novel module named Adaptive Regional Shift Extractor (*AReSE*) to extract key point shifts from high-level feature maps as our temporal features. To cope with the different key point distributions illustrated as Figure 1b and preserve local information,

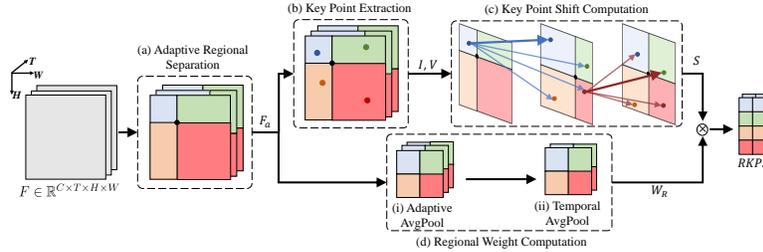


Figure 4: Details of *AReSE* applied to a single channel. Here each frame is split into $K = 4$ regions. The high-level features F extracted by CNN is first separated adaptively, forming the adaptive regional feature maps F_a . The key point coordinates I and feature values V are then extracted from each region at each channel. The key point shifts S are computed as the weighted coordinate shift of corresponding regional key points across adjacent frames. The regional weights W_R are generated by performing adaptive average pooling and temporal average pooling operations to indicate the importance of a region. The *RKPS* is computed as the weighted key point shifts weighted by the regional weights.

we first adaptively separate the feature maps into multiple regions. The key points and their shifts are subsequently extracted from each region to generate Regional Key Point Shifts (*RKPS*). In this section, we illustrate how the proposed *AReSE* extracts *RKPS* step by step, including (a) Adaptive Separation of Feature Maps, (b) Key Point Extraction, (c) Key Point Shift Computation and (d) Regional Weight Computation as shown in Figure 4.

Adaptive Separation of Feature Maps. Taking the high-level features F as the input, the *AReSE* first separates each frame of feature maps into K regions. Instead of manually setting the boundaries between the regions, the splits of the different regions at each frame are trainable and adaptive in order to adjust to the different key point distributions in each frame. Figure 4 shows an example of adaptive separation of $K = 4$ regions, where regions are separated by a single separation centre. For a sequence of T frames, we denote the stack of geometric centres of the feature maps F as $O \in \mathbb{R}^{T \times 2}$, with O_i being the geometric center of the i^{th} frame located at (x_i, y_i) . O splits the original feature maps F into K equal regions. The corresponding stack of adaptive separation centres, denoted as $O_a \in \mathbb{R}^{T \times 2}$, is computed by adding an adaptive bias, which is obtained from

each frame through a Multilayer Perceptron (MLP). The adaptive separation centres O_a are computed by:

$$B_i = \mathcal{I}_i(F_i) = (\Delta x_i, \Delta y_i) \quad (4)$$

$$O_{a_i} = O_i + B_i, \quad (5)$$

where O_i and O_{a_i} are the geometric centre and the adaptive separation point of the i^{th} frame $F_i \in \mathbb{R}^{C \times H \times W}$, respectively. B_i is the adaptive bias of the i^{th} frame F_i . \mathcal{I}_i is the MLP which generates the adaptive bias B_i . The MLPs are trained jointly with the network. Given the adaptive separation centers O_a , the high-level features are then separated into $K \times T$ regions, resulting in the adaptive regional feature maps F_a consisting of $\{F_a^{(1,1)}, \dots, F_a^{(T,K)}\}$ of sizes $\{C \times H^{(1,1)} \times W^{(1,1)}, \dots, C \times H^{(T,K)} \times W^{(T,K)}\}$ as in Figure 4(a).

The separation of feature maps provides regional key points with local characteristics. If the split of each frame is fixed, different key points that are spatially close to each other may locate within the same region. In such cases, key points except the one with the maximum feature value would be ignored. Since the adaptive separation directly affects the resulting key points, we utilize multi-set of *AReSE* as mentioned in Section 3.2 to generate more key points under different region separations. The multi-set operation improves the diversity of feature map separation and therefore improves the robustness of extracted temporal features.

Key Point Extraction. Given the high-level feature map of region r located at channel c of frame f denoted as F_r , the maximum point is extracted as the key point since the maximum feature value point represents the area of raw pixels with key spatial information. Given $F_r \in \mathbb{R}^{H_r \times W_r}$, the coordinate as well as the feature value of the key point are extracted as:

$$I_r = \arg \max_{(h,w)} F_r(h, w) = (x_r^{max}, y_r^{max}) \quad (6)$$

$$V_r = \max(F_r) = F_r(x_r^{max}, y_r^{max}), \quad (7)$$

where $I_r \in \mathbb{R}^2$ is the coordinate of the key point at the region r of the feature map. V_r is the feature value of the key point.

Note that the key point extraction is operated at each channel in practice. In another word, given a region $F_R \in \mathbb{R}^{C \times H_r \times W_r}$ with multiple channels, the key point extraction generates C key points in total, each of which represents the spatial location of the key feature in its own channel and therefore we can compute the movement of these key features independently in subsequent procedures. Given the adaptive regional feature maps F_a consisting of $\{F_a^{(1,1)}, \dots, F_a^{(T,K)}\}$, the key point extraction results in key point coordinates $I \in \mathbb{R}^{K \times C \times T \times 2}$ and key point values $V \in \mathbb{R}^{K \times C \times T}$ indicated as the points with darker color in Figure 4(b).

Key Point Shift Computation. Given the key point 2D coordinates I and their respective feature values V , we compute the regional key point shifts across adjacent frames by two steps, including spatial location difference computation and shift weight computation.

There are many cases where corresponding key points may not belong to the same region across adjacent frames (e.g. Figure 1c). To cope with these situations, we first find the spatial location difference between any pair of regional key points in adjacent frames. The key point shifts are then computed as the weighted sum of the location differences based on the correlation of the key points. This ensures that the computed shift are indeed extracted between two corresponding key points. Given any two key points $I_{i,\alpha}, I_{i+1,\beta}$ where α, β are two regions located at frames $i, i + 1$, respectively, the spatial location difference of these two key points is computed as:

$$\Delta_{i,\alpha,\beta} = I_{i+1,\beta} - I_{i,\alpha}, \quad (8)$$

where $\Delta_{i,\alpha,\beta}$ is the spatial location difference between key points $I_{i+1,\beta}$ and $I_{i,\alpha}$. In practice, the spatial location computation is operated across any adjacent frames at each channel. Therefore, given the key point coordinates $I \in \mathbb{R}^{K \times C \times T \times 2}$, the resulting stack of spatial location differences Δ of size $K_r \times K_n \times C \times (T - 1) \times 2$ contains all channel-wise spatial location differences between any two key points in adjacent frames. Here we use K_r and K_n to re-

fer the region dimension in the current and the next frame for spatial location differences Δ , respectively, with their values both equal to K .

The key point shifts are then obtained through attending to the spatial location differences of corresponding key points, which should have the strongest correlation. Formally, given the spatial location differences between any two adjacent frames i and $i + 1$ denoted as $\Delta_i \in \mathbb{R}^{K_r \times K_n \times 2}$, the key point shifts S_i are weighted sums of Δ_i across K_n . Here the shift weight $W_{i,\alpha,\beta}$ is related to the correlation between the key points located at region α at the recent frame i and the key point located at region β at the next frame $i + 1$. The shift weight indicates the probability of the key points at $I_{i,\alpha}$ of the recent frame falling at the location $I_{i+1,\beta}$ at the next frame. Given the stack of spatial location differences $\Delta_i \in \mathbb{R}^{K_r \times K_n \times 2}$, the shift weight and the key point shifts at the region α of the recent frame is computed by:

$$W_{i,\alpha,\beta} = \mathcal{G}_n \left(\frac{1}{|V_{i,\alpha} - V_{i+1,\beta}| + 0.1} \right) \quad (9)$$

$$S_{i,\alpha} = \sum_{\beta \in K_n} W_{i,\alpha,\beta} \cdot \Delta_{i,\alpha,\beta}, \quad (10)$$

where $V_{i,\alpha}$, $V_{i+1,\beta}$ are the feature values of key points in the region α at the recent frame i and the region β at the next frame $i + 1$, respectively. The \mathcal{G}_n is the softmax function along the K_n dimension. The above equations show that the correlation between any two key points across adjacent frames is computed by the reciprocal of the difference between their feature values. Similar to the spatial location difference computation mentioned in Equation 8, the shift weight computation is also applied to all channels and all frames. Given the spatial location differences $\Delta \in \mathbb{R}^{K_r \times K_n \times C \times (T-1) \times 2}$, the resulting key point shifts $S \in \mathbb{R}^{K \times C \times (T-1) \times 2}$ are obtained from the weighted stack of spatial location differences Δ summed across the K_n dimension per channel.

Regional Weight Computation. In many videos, the temporal features of the action would be located in a certain region of the video frames. It is therefore reasonable to attend to a certain region while extracting key point shifts. Here an attention mechanism is designed to represent the relative significance of re-

gions across each pair of adjacent frames based on its regional average value. Formally, given the adaptive regional feature maps F_a , the regional weight W_R is computed by:

$$W_R = \mathcal{G}_k(\mathcal{P}_T(\mathcal{P}_R(F_a))). \quad (11)$$

This implies that the regional weight $W_R \in \mathbb{R}^{K \times C \times (T-1)}$ is computed as the average values of F_a over the same region for two adjacent frames. The \mathcal{P}_R is an adaptive average pooling operation along the spatial dimension to generate the mean values of each adaptively split region. Whereas the \mathcal{P}_T is an average pooling operation across adjacent frames with the kernel size set to 2 along the temporal dimension. The \mathcal{G}_k is the softmax function along the K dimension. The computed weight W_R represents the attention level to the key point shifts by considering regional feature values across adjacent frames. The resulting *RKPS* is generated as the weighted key point shifts S , weighted by the regional weight W_R . Formally, given the key point shifts S and the regional weight W_R , the *RKPS* is computed as:

$$RKPS = W_R \odot S, \quad (12)$$

where \odot is the element-wise multiplication. The overall *RKPS* $\in \mathbb{R}^{K \times C \times (T-1)}$ is therefore the overall weighted key point shifts of corresponding key points across two adjacent frames of the input feature map F .

As mentioned in Section 3.2, there are multi-set of *AReSEs* in *KPSEM*, each of which generates *RKPS* based on its own adaptive region separation. Our experiments as shown in Section 4.4 indicate that while the multi-set *AReSEs* can further increase the performance compared to a single *AReSE*, using a single *AReSE* in *KPSEM* can also improve the performance compared to the baseline model.

4. Experiments

In this section, we present our evaluation results of the proposed work. The evaluation is conducted through action recognition experiments on four public

benchmark datasets, namely **Mini-Kinetics** [31], **UCF101** [32], **Something-Something v1** [33] and **HMDB51** [34]. We present state-of-the-art results on Mini-Kinetics dataset, and competitive performances on UCF10, Something-Something v1 and HMDB51 datasets. We also present detailed ablation study performed on HMDB51 [34] dataset to verify our design. We further provide heat maps as well as key point shifts visualization of our proposed framework to justify the effectiveness of our proposed work.

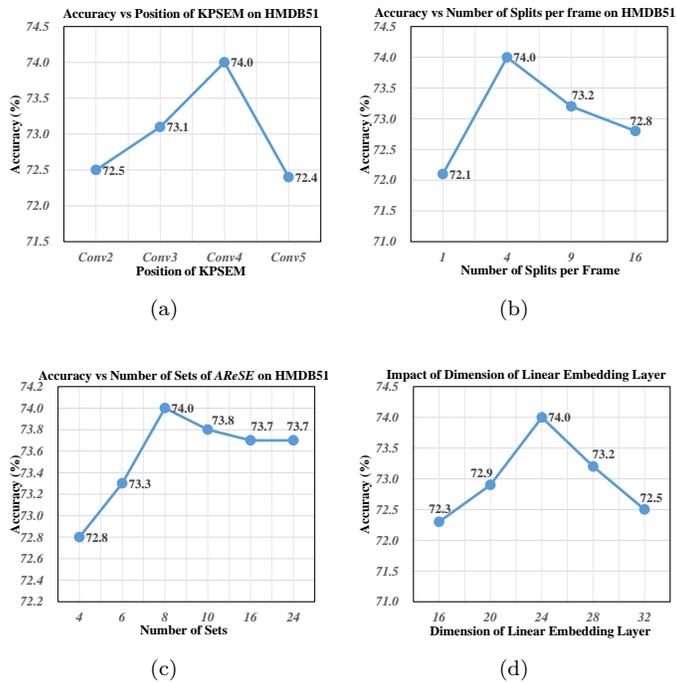


Figure 5: Experiments for best settings of *KPSEM* with MFNet [35]. The experiments are conducted on HMDB51 split-1 with a small batch size of 16. (a) Accuracy *vs.* the position of *KPSEM* on HMDB51, with *KPSEM* placed at the end of each Conv layer. (b) Accuracy *vs.* the number of splits per frame on HMDB51. (c) Accuracy *vs.* the number of sets of *AReSE* on HMDB51. (d) Accuracy *vs.* dimension of the linear embedding layer. Our implementation obtains a 70.8% accuracy with MFNet, which is much lower than that reported in [35] mainly due to the much smaller batch size. All settings surpass the baseline by at least 1.3%.

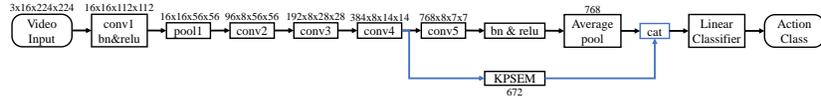
4.1. Experimental Settings

We conduct experiments on four benchmark datasets of action recognition: Mini-Kinetics, Something-Something v1, UCF101 and HMDB51. **Mini-Kinetics** is a subset of the Kinetics [23] dataset, with 200 of its categories. It contains 80K training data and 5K validation data. **Something-Something v1** [33] contains 108,499 videos from 174 human-object-interaction action classes, consisted with 86,017 training, 11,522 validation and 10,960 test videos. **UCF101** [32] contains 13,320 videos from 101 action categories. **HMDB51** [34] contains 51 action categories including a total of 7,000 videos. For UCF101 and HMDB51 datasets, we follow the experiment settings as in [35, 13, 14] that adopt the three train/test splits for evaluation. We report the average top-1 accuracy over the three splits. Our proposed module for temporal feature extraction can be used with any CNN networks. To obtain the state-of-the-art result on Mini-Kinetics and competitive results on UCF101, Something-Something v1 and HMDB51, we instantiate MFNet [35] thanks to its superior performance on Kinetics. The variant of MFNet combined with *KPSEM* is referred to as MF-KPSEM.

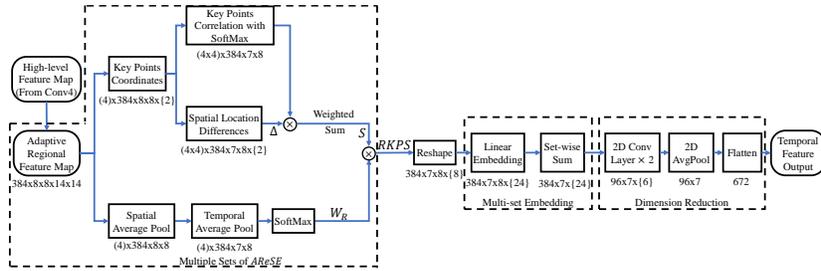
Our experiments are implemented using PyTorch [36]. Following the implementation in [35], the input is a frame sequence with each frame of size 224×224 . Our *KPSEM* extracts the temporal features from the output of the Conv4 layer of MFNet. We choose to split each frame into $K = 4$ splits and adopt $G = 8$ sets based on empirical results in Figure 5. The output dimension of the linear embedding layer is set as 24. This setting exhibits the ultimate results as shown in Figure 5. To accelerate training, we utilize the pretrained model of MFNet [35] trained on Kinetics [23]. The stochastic gradient descent algorithm [37] is used for optimization, with the weight decay set to 0.0001 and the momentum set to 0.9. Our initial learning rate is set to 0.005. A more detailed settings analysis is illustrated in Section 4.4.

4.2. Detailed Implementation of *KPSEM*

As mentioned in Section 3, our proposed *KPSEM* utilizes MFNet as the backbone CNN network. Here we present a more detailed implementation of



(a)



(b)

Figure 6: Detailed implementation of (a) our action recognition framework MF-KPSEM, utilizing MFNet as the backbone for spatial feature extraction, with our proposed *KPSEM* for temporal feature extraction; and (b) more detailed implementation of *KPSEM* with 8 sets of *AReSEs*. For (a), the figures above or below the blocks are the output size of the respective blocks expressed in $C \times T \times H \times W$. Here C is the number of channels, T , H and W are the length of time, height and width of the features. For (b), the figures below the blocks are the output size of the respective blocks. The number within the parentheses corresponds to the dimension reflecting the split of frames K , while the number within the braces corresponds to the key point shifts or its embedding. The other figures are expressed in the order of $C \times T \times G \times H \times W$. Here C is the number of channels, T , G , H and W are the length of time, number of sets for *AReSEs* (*cf.* Section 3.3), height and width of the features. The black path is the path to extract spatial features with MFNet. Whereas the blue path is the path to extract temporal features with *KPSEM*. Figures are best viewed zoomed in.

the overall framework of MF-KPSEM, shown in Figure 6a, and the detailed implementation of our proposed *KPSEM* as shown in Figure 6b.

We follow the implementation in [35] where the input is a frame sequence of 16 frames, with each frame of size 224×224 . As mentioned in Section 3.1 of our paper, the input of *KPSEM* is the output of Conv4 layer of MFNet [35], which is of size $384 \times 8 \times 14 \times 14$. To obtain the temporal feature output, *KPSEM* first computes the 8 *RKPS*s obtained through 8 sets of *AReSEs*. Each set of *AReSEs* performs different region separations, splitting the high-level feature maps from the output of Conv4 layer adaptively. All the *RKPS*s

go through a reshape process before going through the multi-set embedding operation. Each *RKPS* is linearly projected separately to a dimension size of 24 as mentioned in Section 4.1. To obtain the overall feature output, a series of dimension reduction operation is utilized after the set-wise summation of the linear embeddings. Specifically, two 2D convolution layers with a kernel size of $(1, 3)$ and a stride $(1, 2)$ are first utilized to reduce the embedding from the size of 24-d to 6-d. Subsequently, the result passes through a 2D average pooling layer and then flattened to produce the final temporal feature output. It is worth noticing that in addition to reducing embedding size, we also reduce the channels from the size of 384-d to 96-d to discard the non-salient embeddings. The size of the temporal feature output is of 672-d.

4.3. Results and Comparison

Table 1 shows the comparison of top-1 accuracy on Mini-Kinetics, UCF101, Something-Something v1 and HMDB51 datasets with other state-of-the-art methods including:

1. *Two-stream CNN methods*: MARS [38], Residual Frame with two-stream input (ResFrame TS) [39] and I3D with two-stream input (I3D TS) [15].
2. *2D CNN & 3D CNN methods*: C3D [13], I3D with RGB input [15], (2+C1)D [40], S3D [31], MFNet [35], ECO [41], ECO_{Lite} [41], TSM [42] and TSN [12].
3. *CNN with learnable feature correlations*: TBN [43], Res50-NL [17], Res50-CGD [44], Res50-CGNL [45], I3D-NL [17] and I3D-NL-GCN [46].

Our state-of-the-art performance is achieved by instantiating MFNet, denoted as MF-KPSEM. For the experiments as presented in Table 1, the batch size is set to 64 for Mini-Kinetics as well as Something-Something v1 datasets, 80 for UCF101 dataset and 128 for HMDB51 dataset, respectively. The experiments are conducted using two NVIDIA Quadro RTX8000 GPUs.

The performance results in Table 1 show that our network achieves the state-of-the-art result on the Mini-Kinetics with only a minor increase in the

	Method	Mini-Kinetics	UCF101	STH-STH v1	HMDB51	# Params	FLOPs
Two-stream CNNs	MARS [38]	73.5%	98.1%	53.0%	80.9%	-	-
	ResFrame TS [39]	73.9%	90.6%	-	55.4%	-	-
	I3D (TS) [15]	78.7%	97.9%	-	80.2%	25.0M	>107.9G
2D CNNs & 3D CNNs	C3D [13]	66.2%	85.2%	-	-	33.3M	-
	I3D (RGB) [15]	74.1%	95.4%	45.8%	74.5%	12.06M	107.9G
	(2+C1)D [40]	75.74%	96.9%	-	75.2%	7.3M	31.9G
	S3D [31]	78.0%	96.8%	48.2%	75.9%	8.77M	43.47G
	MFNet [35]	78.35%	96.0%	43.0%	74.6%	8.0M	11.1G
	ECO [41]	-	94.8%	41.4%	72.4%	47.5M	64G
	ECO _{Lite} [41]	-	-	46.4%	-	150M	267G
	TSM [42]	-	95.9%	49.7%	73.5%	48.6M	98G
TSN [12]	-	94.2%	19.5%	69.4%	10.7M	16G	
CNN with learnable feature correlations	TBN [43]	69.5%	93.6%	-	69.4	11.4M	-
	Res50-NL [17]	77.53%	82.88%	-	-	27.66M	19.67G
	Res50-CGD [44]	77.56%	84.06%	-	-	25.58M	17.88G
	Res50-CGNL [45]	77.76%	83.38%	-	-	27.2M	19.16G
	I3D-NL [17]	-	-	44.4%	-	27.2M	19.16G
I3D-NL-GCN [46]	-	-	46.1%	-	27.2M	19.16G	
Ours	MF-KPSEM	82.05%	97.4%	48.1%	77.7%	8.11M	11.21G

Table 1: Comparison of top-1 accuracy, number of parameters and computational cost in FLOPs with state-of-the-art methods on Mini-Kinetics, UCF101, Something-Something v1 and HMDB51 datasets. MF-KPSEM instantiates MFNet as the backbone CNN.

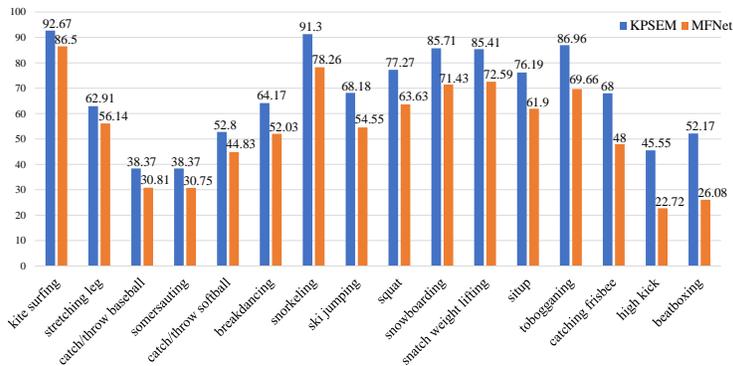


Figure 7: Detailed comparison of accuracy per class on Mini-Kinetics. Here we present the accuracies of 16 classes where MF-KPSEM outperforms by a margin of at least 6%.

number of parameters and required computational cost. Our proposed MF-KPSEM achieves a relative 4.73% increase in recognition accuracy over our baseline method, at a cost of a mere 1.38% increase in parameters and 1% extra FLOPs. Our method surpasses the previous state-of-the-art method which

					
Situp	Stretching Leg	Somersaulting	Snowboarding	Tobogganing	Ski Jumping
KPSEM (Ours)	KPSEM (Ours)				
Situp 45.58	Stretching Leg 83.05	Somersaulting 84.07	Snowboarding 53.11	Tobogganing 51.87	Ski Jumping 71.17
Stretching Leg 23.07	Yoga 16.86	Crawling Baby 14.23	Ski Jumping 30.09	Ski Jumping 39.41	Skiing 25.15
Somersaulting 22.29	Situp 0.0385	Pushing Cart 0.98	Kite Surfing 3.32	Snowboarding 7.43	Snowboarding 1.68
MFNet	MFNet	MFNet	MFNet	MFNet	MFNet
Somersaulting 66.21	Yoga 62.87	Crawling Baby 50.67	Ski Jumping 30.29	Ski Jumping 73.64	Skiing 49.5
Stretching Leg 14.00	Situp 35.66	Somersaulting 38.81	Windsurfing 18.73	Snowboarding 19.07	Ski Jumping 46.97
Situp 11.25	Stretching Leg 0.59	Breakdancing 10.24	Snowboarding 6.85	Tobogganing 3.98	Snowboarding 2.23
(a)	(b)	(c)	(d)	(e)	(f)
					
Catch/Throw Baseball	Catch/Throw Softball	High Kick	Kitesurfing	Squat	Snatch Weight Lifting
KPSEM (Ours)	KPSEM (Ours)				
Catch/Throw Baseball 23.64	Catch/Throw Softball 91.44	High Kick 60.12	Kite Surfing 49.03	Squat 90.67	Snatch Weight Lifting 53.16
Hitting Baseball 18.29	Javelin Throw 3.77	Capoeira 25.65	Jet Skiing 25.58	Deadlifting 5.97	Clean and Jerk 46.77
Juggling Balls 8.91	High Jump 1.86	Side kick 6.00	Surfing Water 14.09	Clean and Jerk 1.55	Lunge 0.03
MFNet	MFNet	MFNet	MFNet	MFNet	MFNet
Catch/Throw Softball 25.42	Javelin Throw 93.3	Capoeira 54.65	Surfing Water 42.54	Deadlifting 43.86	Clean and Jerk 56.66
Lunge 21.68	Catch/Throw Softball 5.17	Breakdancing 42.44	Jet Skiing 21.26	Clean and Jerk 37.89	Snatch Weight Lifting 43.32
Juggling Balls 10.48	High Jump 0.3	High Kick 0.85	Windsurfing 9.49	Squat 13.60	Squat 0.01
(g)	(h)	(i)	(j)	(k)	(l)

Figure 8: Twelve examples taken from the 16 classes presented in Figure 7. The numbers on the right of each class show the probability of the class from the classifier in percentages. We show three classes with the highest probability. The class with the highest probability is the result of the top-1 classification.

utilizes learnable feature correlations through the CGNL module by 4.29%, with a much lighter network and requires lower computational cost. It is noted that the number of parameters in MF-KPSEM exceeds that of $(2+C1)D$, which is built on top of DenseNet [47] that is characterized by its small parameter size. However, our method exceeds theirs by 6.31% in accuracy with 64.86% reduced FLOPs.

For the other three datasets, our MF-KPSEM also performs competitively, gaining a 1.4% increase for UCF101, 5.1% increase for Something-Something v1 and 3.1% for HMDB51 compared to our baseline network. For UCF101 and HMDB51 datasets, our proposed MF-KPSEM performs slightly poorer than the two-stream CNN MARS [38], trailing by 0.7% for UCF101 and 3.2% for HMDB51. Yet our approach does not utilize optical flow as our input during both training and inference, which means a significant reduction in memory and computational cost. For Something-Something v1, while our MF-KPSEM performs slightly inferior compared to S3D [31] and TSM [43], our MF-KPSEM

is a much lighter network with smaller parameter size and requires less computational cost. Specifically, the parameters of our MF-KPSEM are 7.5% and 83.3% less than those of S3D and TSM, respectively. For computational cost, MF-KPSEM requires 74.2% fewer FLOPS compared to S3D and 88.6% fewer FLOPS compared to TSM. Despite the obvious gap of parameters and FLOPS, our MF-KPSEM still performs competitively for Something-Something v1, with a minor 0.1% and 1.8% gap in accuracy compared to S3D and TSM, respectively.

We further investigate the improvement over different actions and present a more detailed comparison of performance between our proposed MF-KPSEM network and the baseline MFNet network. Figure 7 shows the accuracy of 16 classes from the Mini-Kinetics dataset, where our network outperforms the original network by a noticeable margin of over 6%. Many of the examples in these classes are characterized by the fact that the frames in each of these examples would appear similar to other action classes. Therefore, the temporal features showing how the action evolves is the key to correctly classifying these examples.

Figure 8 present 12 examples from the 16 classes mentioned in Figure 7, all of which are better classified through the proposed MF-KPSEM. It could be observed that the spatial features of the given examples, or more intuitively the appearance of the frames in the given examples, could not provide effective representation for accurate action recognition. For example, for Video (a) in Figure 8, the actor is seen rolling up herself. Such a scenario could be present in the action class “Somersaulting”, in which actors would roll themselves up to turn upside down. It could also be presented in the action class “Situp”, where the rolling up is followed by the actor rolling backwards to roll up again. The class of this action could only be determined through the temporal features, which are the change of the actor’s position. In this video, the actor is rolling backwards after this scenario. The actual change of the actor’s position suggests that the action should be classified as “Situp”. If only spatial features are utilized as in the case of MFNet, the video would be instead classified as “Somersaulting” due to the multiple frames showing the actor rolled up. This

clearly shows the importance of temporal features in accurate action recognition and the effectiveness of our *KPSEM* in extracting effective temporal features.

Model	Accuracy
MFNet	70.8%
MF-KPSEM	74.0%
Single Key Point	72.1%

(a) Separation of Feature Maps

Model	Accuracy
MFNet	70.8%
MF-KPSEM	74.0%
Fixed Regions	72.5%

(b) Adaptive Regions

Model	Accuracy
MFNet	70.8%
MF-KPSEM	74.0%
One Set of <i>AReSE</i>	72.4

(c) Multiple Sets of *AReSE*

Model	Accuracy
R3D [16]	62.0%
R3D-KPSEM	65.8%

(d) *KPSEM* with R3D

Table 2: Ablations of *KPSEM* utilizing MFNet on HMDB51 split-1. The ablation is performed with a small batch size of 16. The networks with variants of *KPSEM* that utilize (a) only 1 single key point per frame, (b) only fixed regions per frame and (c) with only 1 set of *AReSE* is compared with the proposed *KPSEM*, and (d) different backbone networks.

4.4. Ablation Study and Visualization

In this section, we justify our proposed design of *KPSEM* through ablation study and visualization of results. Specifically, we examine the performance of our *KPSEM* in four scenarios and justify the need for high-level feature maps as input for *KPSEM*, multiple key points per frame, adaptive regions in each frame and multiple sets of *AReSE*. Additionally, we combine our *KPSEM* with another baseline network R3D to justify the robustness of *KPSEM*. We further examine the effectiveness of *KPSEM* by visualizing heatmaps and the extracted key points with their corresponding shift. The split-1 of HMDB51 dataset is adopted for all ablation studies, trained with a batch size of 16 on a single NVIDIA TITAN Xp GPU.

Position of KPSEM. Our proposed *KPSEM* module utilizes high-level feature maps as the input for extracting key point shifts. Figure 5a compares

the result when *KPSEM* is added to different stages of MFNet. *KPSEM* is added right after the respective layers. Though improvements have been made for all networks utilizing *KPSEM* regardless of the position, the improvement achieved when *KPSEM* is added after Conv2 is 1.5% smaller than that when *KPSEM* is added after Conv4. One possible explanation is that the representation level of the feature maps after Conv2 layer is lower than that of the feature maps after Conv4 layer. This indicates that points with higher feature values from Conv2 layer may be more relevant to pixels with higher values rather than key points. Pixels with higher values may not be key points as they may correspond to the white background pixels. This indicates that the key point shifts extracted from Conv2 feature maps may be semantically less representative than those extracted from Conv4 feature maps, thus resulting in inferior performances. We also observe a sharp drop in performance when *KPSEM* is positioned after Conv5 layer. We believe that this is because of the small spatial size of Conv5 feature maps (7×7). The Conv5 feature maps are therefore insufficient to provide spatial information of the key points, decreasing the accuracy of the extracted key point shifts. This ends up in the inferior performance when *KPSEM* is positioned after Conv5 layer.

Separation of Feature Maps. As mentioned in Section 3.2, frames are split such that localized key points distributed across each frame are preserved. We justify the need for multiple key points per frame by comparing with the variant of *KPSEM* where only a single global key point utilized for each frame at each channel. As indicated in Table 2a, the use of multiple key points (in this case 4 key points) per frame boosts the performance by 1.9%. This demonstrates the effectiveness of extracting multiple regional key points across each frame. The results are consistent with that shown in Figure 5b, where all results are higher than the result using only a single key point.

However, more splits per frame do not guarantee higher accuracy. The result in Figure 5b shows that splitting each frame with $K = 9$ or $K = 16$ regions results in a slight decrease in performance compared to that of splitting $K = 4$

regions per frame. The slightly worse performance for higher K settings could be explained by higher region number results in smaller regions, which may result in regions with only the background being split out. The key points extracted from regions corresponding to the background and their related key point shifts would be redundant. The resulting temporal features would therefore be less effective and result in inferior classification accuracies.

Adaptive Separation of Feature Maps. Also mentioned in Section 3.2, fixed splits for each frame may result in different key points located within the same fixed region. To mitigate this drawback of fixed splits, we split each frame adaptively to obtain an optimal region separation solution for key point extraction. We justify the need for adaptive splits by comparing with the variant of *KPSEM* utilizing fixed splits for each frame. For this variant of *KPSEM*, the split of each set of *AReSE* is randomly initialized, while other network settings including the number of sets and the number of splits remain the same as the default settings of our proposed *KPSEM*. After initialization, the splits of each frame are fixed during the training process. Results in Table 2b demonstrate that a 1.5% increase in accuracy is achieved through splitting frames adaptively. This validates that adaptive splitting could result in better temporal features.

Multiple Sets of AReSE. To extract more key points while avoiding over splitting each frame, we proposed to use multiple sets of *AReSE* in Section 3.2 and Section 3.3. We examine the effect of using multiple sets of *AReSE* with its results as presented in Table 2c. The use of 8 sets of *AReSE* helps improve the network by 1.6% over its variant which uses only a single set of *AReSE*. This matches the results shown in Figure 5c, where the accuracy rises as the number of sets of *AReSE* increases in general. All results using multiple *AReSE* outperform that of using a single set of *AReSE*. The increase justifies the effectiveness of increasing extracted key points without over splitting each frame.

As mentioned in 3.2, the number of *RKPS* increases with the increasing number of sets of *AReSE*. For $G \leq 8$, a significant increase in accuracy with an increase of G can be observed. This suggests that the increase in set number G

results in different *RKPS*, or key point shifts, extracted under different region separations. The different key point shifts extracted by multiple sets of *AReSE* constitutes more robust temporal features. It is worth noticing that the accuracy saturates for $G > 8$. This suggests that a further increase in G may not result in different region separations. The increased *RKPS*s may be a repetition of previous *RKPS*s, which therefore may not increase the robustness of temporal features and the classification accuracy.

KPSEM with Other Backbones. Our proposed *KPSEM* can be applied with any other CNN backbones and improve their performance. To demonstrate the robustness of our *KPSEM*, we have conducted another experiment on *KPSEM* with a different backbone, namely the 3D variant of ResNet50 [16] referred to as R3D [17]. R3D is constructed by inflating 2D convolution kernels directly to 3D convolution kernels, implemented as $1 \times k \times k$ kernels. Therefore, the R3D simply aggregates the input frames and can be directly initialized from weights pretrained on ImageNet. Our *KPSEM* is inserted after the res3 layer of R3D which results in a 3.8% improvement for top-1 accuracy compared to the baseline. The improvement which *KPSEM* achieves with R3D as well as MFNet [35] justifies the robustness of our proposed *KPSEM* module.

Heatmaps Visualization. To investigate where the MF-KPSEM and the baseline MFNet focus on, we visualize the heatmaps that indicate the focus of either network on sampled frames. Specifically, the heatmaps are computed based on Gradient-weighted Class Activation Mapping (Grad-CAM) [48]. The Grad-CAM results are extracted from the last convolution layer of both MF-KPSEM and the baseline MFNet. Figure 9 illustrates Grad-CAM results of four examples from the test set of split 1 of HMDB51, each of which includes four frames sampled identically from the input 16 frames. For each sub-figure, the upper sequence of frames is extracted from the baseline MFNet while the lower sequence is from MF-KPSEM. It can be observed that our proposed MF-KPSEM focuses on the key regions relevant to the action more accurately compared to the baseline model. Whereas the baseline MFNet sometimes focuses on irrelevant regions

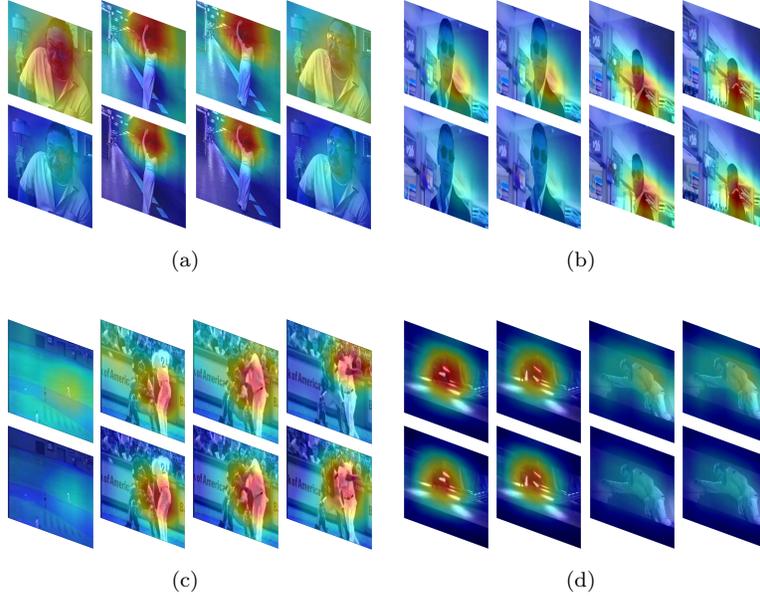


Figure 9: Heatmaps of four samples of HMDB51. For each sub-figure, the upper is the result of baseline MFNet and the lower is the result of our proposed MF-KPSEM. Figure 9a, Figure 9b, Figure 9c and Figure 9d belong to action classes “Wave”, “Turn”, “Swing baseball” and “Sword”, respectively. Compared to the baseline, our proposed MF-KPSEM accurately concentrates on the key regions which are relevant to the ground-truth actions, while the baseline model might focus on irrelevant regions instead.

which bring in redundant temporal information. For example, given the input frames for action “Wave” as demonstrated in Figure 9a, MFNet concentrates not only on the waving hand but also on the face of another actor which is clearly irrelevant to the action “Wave”. In comparison, our MF-KPSEM accurately focuses only on the moving hand of the actress, which are the key regions that imply the ground-truth action.

Visualizing Extracted Key Points and Their Shift. To further investigate the behaviour of our proposed MF-KPSEM, we visualize six examples over the extracted key points and their corresponding shift as shown in Figure 10. Our proposed MF-KPSEM could locate key points and their shifts for each action, such as the elbows and feet of the actor for “Climbing” in Figure 10a. The shift

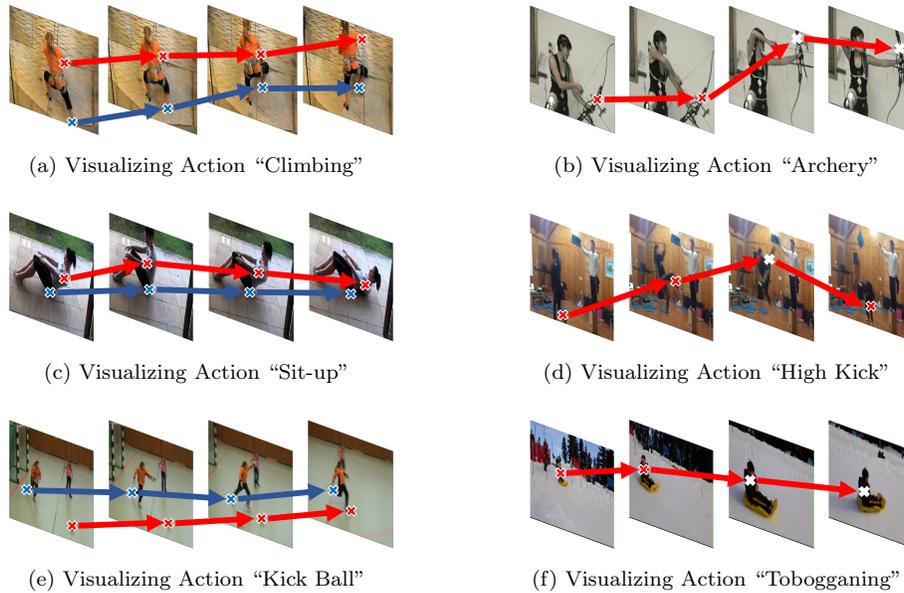


Figure 10: Extracted key points and their corresponding shift from MF-KPSEM for six videos. Crosses coloured differently indicate that they are located in different regions. Arrows in different colours indicate the respective key point shifts. Figure 10a, Figure 10b and Figure 10e are videos from HMDB511 dataset and the others are from Mini-Kinetics. Figure best viewed in colour and zoomed in.

marked in arrows match with the actual moves of the actor for the action, which justify the use of key point shifts for extracting temporal features. In addition, in Figure 10c, the key point shifts of red crosses indicate that the actor is moving back and forth. These key point shifts describe the characteristic of action sit-up and distinguish it from other similar actions such as rolling forward or lying down. Additionally, in Figure 10d, the key point shifts of red and white crosses indicate that the foot shift intensively from the bottom to the top of the video, which is the characteristic of the high kick. It shows that the key point shifts can be applied as the temporal information to differentiate actions from actions.

The MF-KPSEM is also shown to be able to capture corresponding key points and thus computing their shifts even when key points would alter their

location dramatically across regions. Among the samples, Figure 10a, Figure 10c and Figure 10e exhibit the cases where key points remain in the same region across frames. Whereas Figure 10b, Figure 10d and Figure 10f show the possibility that key points could move to another region in the next frame. Take Figure 10c and Figure 10d as examples for these two scenarios. Key points of “Sit-up” in Figure 10c shift within the same regions while the ones of “High kick” in Figure 10d may shift to a different region in a certain frame. This justifies the need for finding corresponding key points across the frame before extracting the key point shifts in *AReSE* as mentioned in Section 3.3.

5. Conclusion and Future Works

In this work, we propose a novel method for extracting the temporal features of a video effectively. The new *KPSEM* exploits key point shifts for temporal feature extraction without additional key point annotation. The overall temporal features encode the key point shifts through linear embedding. Our method obtains state-of-the-art result on Mini-Kinetics when instantiating MFNet, with low additional computational cost compare to other temporal feature extraction methods. We further justify the design and the robustness of our *KPSEM* module through detailed ablation experiments.

In the future, improvements on explicit key point selection without key point annotation could be further explored. Current key point selection method is effective with trivial computational cost. However, in some cases, it may extract key points irrelevant to the actors or objects of the action. More comprehensive key point selection methods can be developed to further improve extracted temporal features without key point annotation.

References

- [1] L. Minh Dang, K. Min, H. Wang, M. Jalil Piran, C. Hee Lee, H. Moon, Sensor-based and vision-based human activity recognition: A comprehen-

- sive survey, *Pattern Recognition* 108 (2020) 107561. doi:10.1016/j.patcog.2020.107561.
- [2] L. Lo Presti, M. La Cascia, 3d skeleton-based human action classification: A survey, *Pattern Recognition* 53 (2016) 130 – 147. doi:10.1016/j.patcog.2015.11.019.
- [3] T. Xiang, S. Gong, Activity based surveillance video content modelling, *Pattern Recognition* 41 (7) (2008) 2309 – 2326. doi:10.1016/j.patcog.2007.11.024.
- [4] J. Yang, H. Zou, H. Jiang, L. Xie, Device-free occupant activity sensing using wifi-enabled iot devices for smart homes, *IEEE Internet of Things Journal* 5 (5) (2018) 3991–4002. doi:10.1109/JIOT.2018.2849655.
- [5] Y. Li, R. Xia, X. Liu, Learning shape and motion representations for view invariant skeleton-based action recognition, *Pattern Recognition* 103 (2020) 107293. doi:10.1016/j.patcog.2020.107293.
- [6] C. Si, Y. Jing, W. Wang, L. Wang, T. Tan, Skeleton-based action recognition with hierarchical spatial reasoning and temporal stack learning network, *Pattern Recognition* 107 (2020) 107511. doi:10.1016/j.patcog.2020.107511.
- [7] A. Shahroudy, J. Liu, T. Ng, G. Wang, Ntu rgb+d: A large scale dataset for 3d human activity analysis, in: 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2016, pp. 1010–1019. doi:10.1109/CVPR.2016.115.
- [8] C. Li, Z. Cui, W. Zheng, C. Xu, J. Yang, Spatio-temporal graph convolution for skeleton based action recognition, 2018.
URL <https://aaai.org/ocs/index.php/AAAI/AAAI18/paper/view/17103>
- [9] Z. Cao, T. Simon, S. Wei, Y. Sheikh, Realtime multi-person 2d pose estimation using part affinity fields, in: 2017 IEEE Conference on Com-

- puter Vision and Pattern Recognition (CVPR), 2017, pp. 1302–1310. doi:10.1109/CVPR.2017.143.
- [10] C. Feichtenhofer, A. Pinz, R. P. Wildes, Spatiotemporal multiplier networks for video action recognition, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2017, pp. 4768–4777. doi:10.1109/CVPR.2017.787.
- [11] K. Simonyan, A. Zisserman, 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, Curran Associates, Inc., 2014, pp. 568–576.
- [12] L. Wang, Y. Xiong, Z. Wang, Y. Qiao, D. Lin, X. Tang, L. Van Gool, Temporal segment networks for action recognition in videos, IEEE Transactions on Pattern Analysis and Machine Intelligence 41 (11) (2019) 2740–2755. doi:10.1109/TPAMI.2018.2868668.
- [13] D. Tran, L. Bourdev, R. Fergus, L. Torresani, M. Paluri, Learning spatiotemporal features with 3d convolutional networks, in: Proceedings of the IEEE international conference on computer vision, 2015, pp. 4489–4497. doi:10.1109/ICCV.2015.510.
- [14] D. Tran, H. Wang, L. Torresani, J. Ray, Y. LeCun, M. Paluri, A closer look at spatiotemporal convolutions for action recognition, in: Proceedings of the IEEE conference on Computer Vision and Pattern Recognition, 2018, pp. 6450–6459. doi:10.1109/CVPR.2018.00675.
- [15] J. Carreira, 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. doi:10.1109/CVPR.2018.00685.
- [16] K. Hara, H. Kataoka, Y. Satoh, Can spatiotemporal 3d cnns retrace the

- history of 2d cnns and imagenet?, in: Proceedings of the IEEE conference on Computer Vision and Pattern Recognition, 2018, pp. 6546–6555.
- [17] X. Wang, R. Girshick, A. Gupta, K. He, Non-local neural networks, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2018, pp. 7794–7803. doi:10.1109/CVPR.2018.38.
- [18] L. Wang, W. Li, W. Li, 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. doi:10.1109/CVPR.2018.00155.
- [19] B. D. Lucas, T. Kanade, An iterative image registration technique with an application to stereo vision, in: P. J. Hayes (Ed.), Proceedings of the 7th International Joint Conference on Artificial Intelligence, IJCAI '81, Vancouver, BC, Canada, August 24-28, 1981, William Kaufmann, 1981, pp. 674–679.
- [20] B. K. Horn, B. G. Schunck, Determining optical flow, Artificial Intelligence 17 (1) (1981) 185 – 203. doi:10.1016/0004-3702(81)90024-2.
- [21] C. Zach, T. Pock, H. Bischof, A duality based approach for realtime tv-l 1 optical flow, in: Joint pattern recognition symposium, Springer, 2007, pp. 214–223. doi:10.1007/978-3-540-74936-3_22.
- [22] Z. Qiu, T. Yao, 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. doi:10.1109/ICCV.2017.590.
- [23] W. Kay, J. Carreira, K. Simonyan, B. Zhang, C. Hillier, S. Vijayanarasimhan, F. Viola, T. Green, T. Back, P. Natsev, et al., The kinetics human action video dataset, arXiv preprint arXiv:1705.06950.
- [24] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, L. Fei-Fei, Imagenet: A large-scale hierarchical image database, in: 2009 IEEE conference on computer

- vision and pattern recognition, Ieee, 2009, pp. 248–255. doi:10.1109/CVPR.2009.5206848.
- [25] A. Buades, B. Coll, J.-M. Morel, A non-local algorithm for image denoising, in: 2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05), Vol. 2, IEEE, 2005, pp. 60–65.
- [26] H. Li, C. Y. Suen, A novel non-local means image denoising method based on grey theory, Pattern Recognition 49 (2016) 237–248.
- [27] H. Wang, D. Tran, L. Torresani, M. Feiszli, Video modeling with correlation networks, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2020, pp. 352–361. doi:10.1109/CVPR42600.2020.00043.
- [28] C. Feichtenhofer, H. Fan, J. Malik, K. He, Slowfast networks for video recognition, in: Proceedings of the IEEE international conference on computer vision, 2019, pp. 6202–6211. doi:10.1109/ICCV.2019.00630.
- [29] C. Yang, Y. Xu, J. Shi, B. Dai, B. Zhou, Temporal pyramid network for action recognition, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2020, pp. 591–600. doi:10.1109/CVPR42600.2020.00067.
- [30] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, I. Polosukhin, Attention is all you need, in: Advances in neural information processing systems, 2017, pp. 5998–6008.
- [31] S. Xie, C. Sun, J. Huang, Z. Tu, K. Murphy, Rethinking spatiotemporal feature learning: Speed-accuracy trade-offs in video classification, in: Proceedings of the European Conference on Computer Vision (ECCV), 2018, pp. 305–321. doi:10.1007/978-3-030-01267-0\$_19.
- [32] K. Soomro, A. R. Zamir, M. Shah, Ucf101: A dataset of 101 human actions classes from videos in the wild, arXiv preprint arXiv:1212.0402.

- [33] R. Goyal, S. Kahou, V. Michalski, J. Materzynska, S. Westphal, H. Kim, V. Haenel, I. Fründ, P. Yianilos, M. Mueller-Freitag, F. Hoppe, C. Thureau, I. Bax, R. Memisevic, The something something video database for learning and evaluating visual common sense, in: 2017 IEEE International Conference on Computer Vision (ICCV), 2017, pp. 5843–5851. doi:10.1109/ICCV.2017.622.
- [34] H. Kuehne, H. Jhuang, E. Garrote, T. Poggio, T. Serre, Hmdb51: A large video database for human motion recognition, in: Proceedings of the IEEE International Conference on Computer Vision, 2011, pp. 2556–2563. doi:10.1109/ICCV.2011.6126543.
- [35] Y. Chen, Y. Kalantidis, J. Li, S. Yan, J. Feng, Multi-fiber networks for video recognition, in: Proceedings of the european conference on computer vision (ECCV), 2018, pp. 352–367. doi:10.1007/978-3-030-01246-5\$_22\$.
- [36] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, et al., Pytorch: An imperative style, high-performance deep learning library, in: Advances in neural information processing systems, 2019, pp. 8026–8037.
- [37] L. Bottou, Large-scale machine learning with stochastic gradient descent, in: Proceedings of COMPSTAT’2010, Springer, 2010, pp. 177–186. doi:10.1007/978-3-7908-2604-3\$_16\$.
- [38] N. Crasto, P. Weinzaepfel, K. Alahari, C. Schmid, Mars: Motion-augmented rgb stream for action recognition, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2019, pp. 7882–7891. doi:10.1109/CVPR.2019.00807.
- [39] L. Tao, X. Wang, T. Yamasaki, Rethinking motion representation: Residual frames with 3d convnets for better action recognition, arXiv preprint arXiv:2001.05661.

- [40] C. Cheng, C. Zhang, Y. Wei, Y.-G. Jiang, Sparse temporal causal convolution for efficient action modeling, in: Proceedings of the 27th ACM International Conference on Multimedia, 2019, pp. 592–600. doi:10.1145/3343031.3351054.
- [41] M. Zolfaghari, K. Singh, T. Brox, Eco: Efficient convolutional network for online video understanding, in: Proceedings of the European conference on computer vision (ECCV), 2018, pp. 695–712. doi:10.1007/978-3-030-01216-8_43.
- [42] J. Lin, C. Gan, S. Han, Tsm: Temporal shift module for efficient video understanding, in: Proceedings of the IEEE International Conference on Computer Vision, 2019, pp. 7083–7093. doi:10.1109/ICCV.2019.00718.
- [43] Y. Li, S. Song, Y. Li, J. Liu, Temporal bilinear networks for video action recognition, in: Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 33, 2019, pp. 8674–8681. doi:10.1609/aaai.v33i01.33018674.
- [44] X. He, K. Cheng, Q. Chen, Q. Hu, P. Wang, J. Cheng, Compact global descriptor for neural networks, arXiv preprint arXiv:1907.09665.
- [45] K. Yue, M. Sun, Y. Yuan, F. Zhou, E. Ding, F. Xu, Compact generalized non-local network, in: Advances in Neural Information Processing Systems, 2018, pp. 6510–6519.
- [46] X. Wang, A. Gupta, Videos as space-time region graphs, in: Proceedings of the European conference on computer vision (ECCV), 2018, pp. 399–417. doi:10.1007/978-3-030-01228-1_25.
- [47] G. Huang, Z. Liu, L. Van Der Maaten, K. Q. Weinberger, Densely connected convolutional networks, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2017, pp. 4700–4708. doi:10.1109/CVPR.2017.243.

- [48] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, D. Batra, Grad-cam: Visual explanations from deep networks via gradient-based localization, *International Journal of Computer Vision* 128 (2) (2020) 336–359. doi:10.1007/s11263-019-01228-7.