


# Learning hierarchical video representation for action recognition

Qing Li<sup>1</sup>  · Zhaofan Qiu<sup>1</sup> · Ting Yao<sup>2</sup> · Tao Mei<sup>2</sup> · Yong Rui<sup>2</sup> · Jiebo Luo<sup>3</sup>

Received: 29 September 2016 / Revised: 22 December 2016 / Accepted: 29 December 2016 / Published online: 15 February 2017  
© Springer-Verlag London 2017

**Abstract** Video analysis is an important branch of computer vision due to its wide applications, ranging from video surveillance, video indexing, and retrieval to human computer interaction. All of the applications are based on a good video representation, which encodes video content into a feature vector with fixed length. Most existing methods treat video as a flat image sequence, but from our observations we argue that video is an information-intensive media with intrinsic *hierarchical* structure, which is largely ignored by previous approaches. Therefore, in this work, we represent the hierarchical structure of video with multiple granularities including, from short to long, single *frame*, consecutive frames (*motion*), short *clip*, and the entire *video*. Furthermore, we propose a novel deep learning framework to model each granularity individually. Specifically, we model the *frame* and *motion* granularities with 2D convolutional neural networks and model the *clip* and *video* granularities

with 3D convolutional neural networks. Long Short-Term Memory networks are applied on the *frame*, *motion*, and *clip* to further exploit the long-term temporal clues. Consequently, the whole framework utilizes multi-stream CNNs to learn a hierarchical representation that captures spatial and temporal information of video. To validate its effectiveness in video analysis, we apply this video representation to action recognition task. We adopt a distribution-based fusion strategy to combine the decision scores from all the granularities, which are obtained by using a *softmax* layer on the top of each stream. We conduct extensive experiments on three action benchmarks (UCF101, HMDB51, and CCV) and achieve competitive performance against several state-of-the-art methods.

**Keywords** Video representation learning · Action recognition · Deep learning

---

✉ Qing Li  
sealq@mail.ustc.edu.cn

Zhaofan Qiu  
qiudavy@mail.ustc.edu.cn

Ting Yao  
tiyao@microsoft.com

Tao Mei  
tmei@microsoft.com

Yong Rui  
yongrui@microsoft.com

Jiebo Luo  
jluo@cs.rochester.edu

<sup>1</sup> University of Science and Technology of China, Hefei, People's Republic of China

<sup>2</sup> Microsoft Research, Beijing, People's Republic of China

<sup>3</sup> University of Rochester, New York, USA

## 1 Introduction

Video analysis is attracting more and more research attention in recent years. This is partially due to the explosive increasing amount of video and the wide variety of applications, ranging from video tagging, video summarization, and video action recognition, to video and language [10, 23, 28, 43, 44]. Video analysis heavily relies on a good video representation. However, devising a robust and discriminative video representation is very challenging due to not only the visual variance caused by camera motion, viewpoint changing or illumination conditions, but also the complex temporal structure of video itself. Traditional hand-crafted methods usually start by detecting spatial-temporal interest points and then represent these points with local descriptors. For instance, Wang et al. propose dense trajectory features in

[36], which tracks densely sampled local frame patches over time and extract several traditional descriptors based on the trajectories. The dense trajectory features can achieve good performance on video action recognition by simply training a SVM classifier on them.

In contrast to the hand-crafted features, there is recently a big surge of automatically learning a representation from the raw data using deep neural networks. Among these networks, two-dimensional convolutional neural networks have exhibited state-of-the-art performance in image analysis tasks like classification or detection [12,29,34]. For video analysis, Karpathy et al. [10] extend the 2D CNN into the temporal dimension by stacking frames over time and achieve promising results on action recognition task. Another important work is the two-stream CNN approach proposed by Simonyan et al. [28], which use two different 2D CNNs on individual frame and stack optical flows, respectively, to capture the spatial and motion information.

As we can see, what is in common for these existing methods is that they treat video as a flat image sequence. However, from our observations, video embeds its intensive information in a *hierarchical* structure. Concretely, if we focus on the content of a single frame from the video, we can only obtain information of the objects and scenes. And if we select two consecutive frames and compute the displacement between them, the motion of the objects can be exploited. Furthermore, by leveraging more continuous frames, we can utilize more complex motion pattern of the objects. In short, *hierarchical* means to harness different granularities or modalities in the video, which are complementary and thus have mutual reinforcement for recognition. In the literature, the utilization of multiple granularities in videos has been shown effective for video understanding and retrieval tasks, such as video shot/scene analysis [19,21], video concept detection [30,38], and video search [40,45].

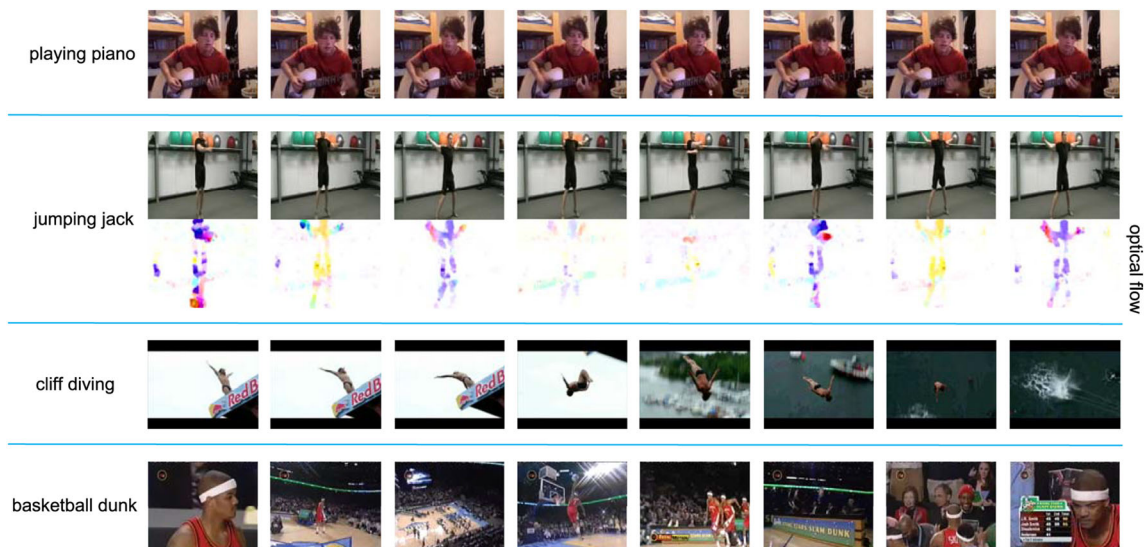
To make full use of the intensive information from video, this paper proposes a multi-stream deep learning framework to learn a hierarchical video representation, which not only harness the spatial–temporal clues, but also consider the multiple granularities of video. To represent the hierarchical structure of video, we define four granularities from short to long, i.e., a single *frame*, consecutive frames (*motion*), a short *clip* and the entire *video*. we model the *frame* and *motion* streams by 2D CNNs, while the *clip* and *video* streams are processed by 3D CNNs. Therefore, the framework can learn both visual appearance and short-term motion information via the multi-stream CNNs. Furthermore, the Long Short-Term Memory networks are utilized on the *frame*, *motion* and *clip* streams for the long-term temporal modeling, while the outputs of 3D CNN on each clip are combined by Mean pooling as the representation of *video* stream.

To verify the power of this hierarchical video representation, we apply it to the action recognition task. First, we equip each stream with a softmax layer to predict the classification scores. And as shown in Fig. 1, an action may span different granularities in a video. Therefore, instead of Mean pooling the classification scores from different streams, we adopt a novel fusion strategy based on the multi-granular score distribution to predict the final probabilities on every action class. We train classifiers on the score distributions to learn the weights of each individual component, which can effectively reflect the importance of each stream and its components to the overall recognition result. And what is worth to mention is that all the granularities are processed in parallel and the whole framework is end-to-end trainable.

The main contributions of this work can be summarized as follows:

- We propose an end-to-end hybrid deep learning framework, which exploits the multiple granularities of video to learn a hierarchical representation for video. The framework can model not only the spatial and short-term motion patterns, but also the long-term temporal clues of the video.
- We adopt the LSTM to model long-term temporal clues on the top of *frame*, *motion* and *clip* streams. We show that all the streams work well with LSTM, which are complementary to the traditional methods without considering the temporal order of frames in video.
- We apply the hierarchical video representation to action recognition task by integrating all the granularities with a novel distribution-based fusion strategy. The fusion scheme not only can reflect the importance of different streams, but also is computationally efficient in training and testing.
- Through an extensive set of experiments, we demonstrate that our proposed framework outperforms several alternative methods with clear margins. On three popular benchmarks (UCF-101, HMDB-51 and CCV), we obtain the state-of-the-art results.

The rest of this paper is organized as follows. Section 2 reviews related work on video representation learning with hand-crafted methods and deep learning architecture. Section 3 describes the proposed multi-granular framework for learning hierarchical video representation in detail, while Sect. 4 formulates the novel fusion scheme based on multi-granular score distribution for action recognition. The experiment settings and implementation details are given in Sect. 5. Section 6 provides experimental results and analyses on three well-known benchmarks (UCF101, HMDB51, and CCV), followed by the conclusions and future work in Sect. 7.



**Fig. 1** An action may span different granularities. For example, the action of “playing piano” can be recognized from individual frames, “jumping jack” may have high correlation with the optical flow images (motion computed from consecutive frames), “cliff diving” should be recognized from a short clip since this action usually lasts for few sec-

onds, while “basketball dunk” can be reliably identified at the video granularity due to the complex nature of this action. Recognizing actions therefore should take the hierarchical multi-granularity and spatial-temporal properties into consideration

## 2 Related work

As aforementioned, video analysis has been an active research topic in multimedia and computer vision. A good video analysis system relies heavily on the extracted video features, so significant efforts have been paid to design discriminative and robust video representations [10, 22, 36]. And here we just focus on the review of recent work in the context of action recognition task of video analysis.

**Hand-crafted representations** There has been numerous work focusing on developing discriminative features by hand that are expected to be able to distinguish different categories and be robust to the large intra-class variances [11, 36, 47]. Designing video representations usually consists of two steps: detecting spatial-temporal interest points and representing these points with local descriptors. For the design of interest points detectors, Laptev and Lindeberg [15] extend the 2D Harris corner detector into 3D space to find the space-time interest points (STIP). To describe the found interest points, we can utilize some image-based descriptors, such as HOG and SIFT, to extract visual appearance information from individual frames of video. In addition to the static appearance information, the motion information is also very crucial for understanding video contents so a lot of efforts are paid to design descriptors taking into account the object movements. A popular way to extract motion descriptors is to extend the frame-based local descriptors into 3D space. For example, Klaser et al. propose HOG3D by extending the idea of integral images for fast descriptor computation [11].

Besides, SIFT-3D [27], Extended SURF [41] and Cuboids [2] are also good choices as the local spatial-temporal descriptors. Recently, Wang et al. propose dense trajectory features, which densely sample local patches from each frame at different scales and then track them in a dense optical flow field [36]. This method has demonstrated very competitive performances on several popular benchmarks. In addition, the further improvements can be achieved by the compensation of camera motion, and the use of advanced feature encoding methods like Fisher vectors. It is worth noting that these spatial-temporal video descriptors can only represent short motion pattern within a very short period, and the popular descriptors encoding methods like bag of words (BoW) just discard the temporal order information of the descriptors totally.

**Deep learning representations** Motivated by the great success of deep neural networks (especially the ConvNets) on image classification tasks [12, 29, 34], there are recently a lot of attempts to devise deep architectures for learning video representations. Qiu et al. perform action recognition using the support vector machine with Mean pooling on the CNN-based representations over frames [25]. Karparthy et al. compare several different fusion architectures for video classification [10] Later in [35], Tran et al. propose to train 3D ConvNets on a large labeled video dataset Sports-1M to learn generic spatial-temporal features which can be computed very efficiently. Xu et al. adopt advanced feature encoding strategies VLAD to make the CNN features generalize better. Zha et al. leverage both spatial and temporal pooling

on the CNN features computed on patches of video frames [49]. Simonyan et al. propose an novel two-stream approach, where two different ConvNets are trained on individual frame and stacked optical flows, respectively, to more explicitly capture the spatial and short-term motion information. Final predictions can be obtained by Mean pooling the decision scores of the two ConvNets or train a SVM classifier on the concatenation of the outputs of the two ConvNets. The late fusion is then exploited to combine spatial–temporal representations. In the work by Wang et al. [37], the local ConvNet responses over the spatiotemporal tubes centered at the trajectories are pooled as the video descriptors. Fisher vector is then used to encode these local descriptors to a global video representation. More recently, local activations of convolutional layer are encoded in a deep generative model for video action recognition in [26]. Similar to the hand-crafted features, the CNN-based representations are also not able to model the long-term temporal information and the proposed several fusion schemes do not consider the temporal order of the different parts of the video.

**Temporal modeling in video** As aforementioned, both the hand-crafted and CNN features cannot model long-term information. There are also extensive works to explore long-term temporal dynamics in video. For instance, Fernando et al. propose to learn a function capable of ordering the frames of a video, which can capture well the evolution of the appearance within the video. Recently, RNN attracts a lot of research attention on many sequential learning tasks such as speech recognition [4] and machine translation [33]. RNN can deal with sequential data with variable length so theoretically it can be utilized to model the long-term temporal dynamics in video. In [20], temporal pooling and LSTM are used to combine frame-level (optical flow images) representation and discover long-term temporal relationships. Srivastava et al. further formulate the video representation learning as an autoencoder model in an unsupervised manner, which consists of the encoder and decoder LSTMs [32].

It can be observed that most existing methods treat video as a flat data sequence while ignoring the aforementioned intrinsic hierarchical structure of the video content deeply. The most closely related work is the two-stream CNN approach proposed by Simonyan et al. [28]. The work applies the CNN separately on individual frame and stacked optical flows. Our method is different from [28] in that we extend two-stream to the multi-granular streams, employ 3D CNN to learn the spatial–temporal representation of video, and further utilize LSTM networks to model long-term temporal cues. Besides, our framework adopts a more principled fusion scheme to integrate each component from all the streams. This paper extends upon a previous conference publication [16]. The extensions include new experiments on more video datasets,

analysis on the power of multi-granular score fusion, and amplified discussions and explanations throughout the paper.

### 3 Hierarchical video representation

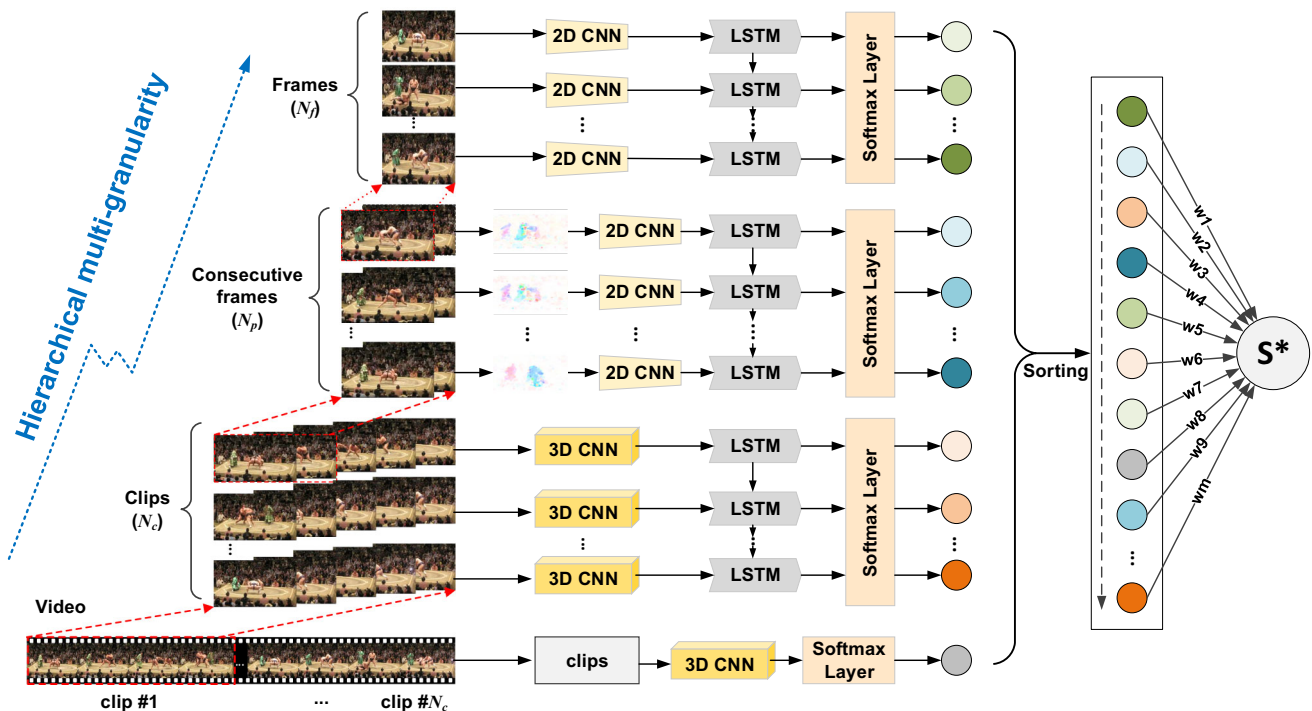
Compared with image, video contains more intensive information, which is essentially embedded in a hierarchical structure. So a good video representation should cover all the aspects of the hierarchical structure. Then how can we define the hierarchical structure of video? In this paper, we represent the hierarchical structure by defining multiple granularities in video from short to long, i.e., single *frame*, consecutive frames (*motion*), short *clip*, and the whole *video*. And we devise a multi-stream deep learning architecture to model each granularity. Figure 2 gives an overview of our framework, and next we will introduce the implementation details of each stream, respectively.

#### 3.1 Modeling frame stream

For representing video, individual frames can provide some static useful clues like particular scenes and objects. In recent years, convolution neural network has proved its surprising power in image analysis, so CNN is very good choice to make full use of the static frame appearance. Among the proposed CNN architectures in recent work, we choose the VGG\_19 [29], a widely adopted CNN architecture for image classification, to extract the high-level visual features from each sampled individual frame. VGG\_19 is a deep convolutional network with up to 19 weight layers (16 convolutional layers and 3 fully connected layers). Because there are not enough samples in current labeled video dataset, training the network from scratch will cause a very heavy overfitting problem. In view of the great similarity between the video frames and the images from ImageNet, it is reasonable to pre-train the network on ImageNet, which is a much larger dataset with labeled images. Then we fine-tune the network on the frames from the video dataset to get the final model. Thanks to the power of this transfer learning, we alleviate the overfitting problem to a great extent. Finally, we utilize the outputs of the fully connected layer of the fine-tuned VGG\_19 model as the representation of individual frame, which contains the scenes and objects information of the video.

#### 3.2 Modeling motion stream

Complementary to frame, motion is another important clue for video representation. The crucial difference between video and static image is the movement of the objects in video, so it is very necessary to consider the motion information when representing video. Following [28], we compute



**Fig. 2** Multi-granular spatial–temporal architecture for video action recognition. A video is represented by the hierarchical structure with multiple granularities including, from short to long, *frame*, consecutive frames (*motion*), *clip*, and *video*. Each granularity is modeled as a single stream. 2D CNNs are used to model the *frame* and *motion* (optical flow images) streams, while 3D CNNs are used to model the *clip* and *video* streams. LSTMs are used to further model the temporal information in the *frame*, *motion*, and *clip* streams. A *softmax* layer is built on the

top of each stream to obtain the prediction from each component. Suppose we have  $N_c$  clips,  $N_p$  motions (consecutive frame pairs), and  $N_f$  frames, then we have  $N_c + N_p + N_f + 1$  components. The final action recognition result of the input video is obtained by linearly fusing the prediction scores from all the components with the weights learnt on the score distribution. Note that this deep architecture is trainable in an end-to-end fashion

the optical flow [1] to explicitly measure the displacement between two consecutive frames. Furthermore, to alleviate the effect of camera motion in video, we subtract from every optical flow its Mean vector. This preprocessing can be viewed as a very rough estimation for the camera motion, and more advanced techniques can be explored for compensating the camera motion, but that topic is just out of scope of our work. After getting the measurements of the objects’ motion in video, we need to encode these measurements into a fixed length representation. A smart way proposed in [3] is to convert the optical flow into “image” by centering horizontal (x) and vertical (y) flow values around 128 and multiplying by a scalar such that flow values fall between 0 and 255. By this transformation, we can obtain two channels of optical flow “image,” while the third channel is created by calculating the flow magnitude. Having converted flow into “image,” What we do next is very similar to the frame stream, i.e., to fine-tune the pre-trained VGG\_19 on the extracted optical flow “images.” And then we compute the outputs of the fully connected layer of the fine-tuned VGG\_19 model as the representation of motion in video.

### 3.3 Modeling clip and video streams

As we can see above, the frame only contains the scenes and objects information, while the motion is considered between only two consecutive frames. Then we try to figure out an approach than can explore the motion pattern of the objects in multiple consecutive frames. 3D CNN is just a very good option, which takes a video clip (multiple continuous frames) as the input and conduct 3D convolutions in both spatial and temporal dimensions. In our work, we adopt the superior 3D CNN architecture proposed in [35], named C3D, which takes 3D convolution and 3D pooling alternatively. C3D is pre-trained on a large-scale labeled video dataset, i.e., Sports-1M, to learn a generic spatial–temporal video representation. Following the spirit of frame and motion streams, we also fine-tune the pre-trained C3D model on clips extracted from the video dataset to get the final model. Then we take the outputs of fully connected layer of C3D as the representations for the sampled clips, and we think these representations can capture the appearance and motion pattern information simultaneously. To represent the video globally, we just adopt

Mean pooling of the features of all the clips as the video-level representations. We have to admit that this is simple as well as rough, and a better way to get the video representation can be explored in the future.

### 3.4 Temporal modeling with LSTM

During the modeling of the frame, motion, and clip streams, each CNN architecture just takes one component (i.e., one frame, optical flow “image,” or clip) in video, and the temporal order of the components is totally discarded. To learn the long-term dependencies between different components of video, we employ the Long Short-Term Memory (LSTM) on frame, motion, and clip streams. LSTM is a type of RNN with special memory cells and controllable gate mechanism and has achieved great success in many long-range sequential modeling task, like speech recognition [4] and machine translation [33]. When dealing with long sequential data, LSTM does not suffer the “vanish gradients” issue like the traditional RNN. In general, LSTM recursively maps the input representations at the current step to the output representation based on the current hidden state, and thus the training process of LSTM should be in a sequential manner. At last, we can compute a decision score at each time step with a softmax layer using the hidden states from the LSTM layer.

More formally, given a sequence of representations  $(\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^T)$ , LSTM maps the input sequence to an output sequence of hidden states  $(\mathbf{h}^1, \mathbf{h}^2, \dots, \mathbf{h}^T)$  by updating the hidden state in the network with following formula recursively from  $t = 1$  to  $t = T$ :

$$\begin{aligned} \mathbf{g}^t &= \phi(\mathbf{T}_g \mathbf{x}^t + \mathbf{R}_g \mathbf{h}^{t-1} + \mathbf{b}_g) && \text{cell input} \\ \mathbf{i}^t &= \sigma(\mathbf{T}_i \mathbf{x}^t + \mathbf{R}_i \mathbf{h}^{t-1} + \mathbf{b}_i) && \text{input gate} \\ \mathbf{f}^t &= \sigma(\mathbf{T}_f \mathbf{x}^t + \mathbf{R}_f \mathbf{h}^{t-1} + \mathbf{b}_f) && \text{forget gate} \\ \mathbf{c}^t &= \mathbf{g}^t \odot \mathbf{i}^t + \mathbf{c}^{t-1} \odot \mathbf{f}^t && \text{cell state} \\ \mathbf{o}^t &= \sigma(\mathbf{T}_o \mathbf{x}^t + \mathbf{R}_o \mathbf{h}^{t-1} + \mathbf{b}_o) && \text{output gate} \\ \mathbf{h}^t &= \phi(\mathbf{c}^t) \odot \mathbf{o}^t && \text{cell output} \end{aligned}$$

where  $\mathbf{x}^t$  and  $\mathbf{h}^t$  are the input and output hidden state vector, respectively, with the superscript  $t$  denoting the  $t$ th time step,  $\mathbf{i}^t, \mathbf{f}^t, \mathbf{c}^t, \mathbf{o}^t$  are the activation vectors of input gate, forget gate, memory cell and output gate, respectively.  $\mathbf{T}$  are input weights matrices,  $\mathbf{R}$  are recurrent weight matrices, and  $\mathbf{b}$  are bias vectors. *Logic sigmoid*  $\sigma(x) = \frac{1}{1+e^{-x}}$  and *hyperbolic tangent*  $\phi(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$  are element-wise nonlinear activation functions, mapping real values to  $(0, 1)$  and  $(-1, 1)$  separately. The dot product and sum of two vectors are denoted with  $\odot$  and  $\oplus$ , respectively. Figure 3 is an illustration of an LSTM unit.

The core idea of LSTM is to introduce a new structure called *memory cell*, which can store information over time to

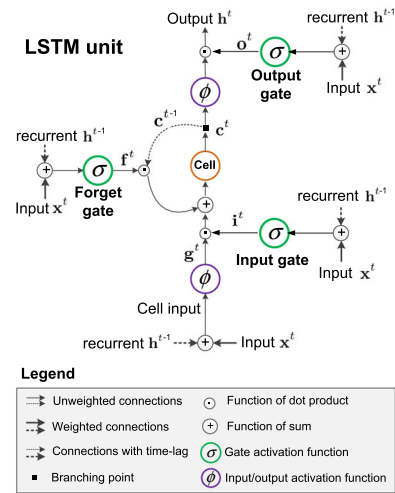


Fig. 3 Diagram of a LSTM memory cell

explore long-range dynamics as well as alleviate the “vanish gradients” effect. As we can see in Fig. 3, a memory cell is composed of four main elements: an input gate, a neuron with a self-recurrent connection, a forget gate, and an output gate. The self-recurrent connection has a weight of 1.0 and ensures that, barring any outside interference, the state of a memory cell can remain constant from one time step to another which is crucial for controlling the gradient flow in backpropagation through time to avoid the gradients vanish. The gates serve to modulate the interactions between the memory cell itself and its environment. The input gate can allow incoming signal to alter the state of the memory cell or block it. On the other hand, the output gate can allow the state of the memory cell to have an effect on other neurons or prevent it. Finally, the forget gate can modulate the memory cell’s self-recurrent connection, allowing the cell to remember or forget its previous state, as needed. Many improvements have been made to the LSTM architecture since its original formulation [6], and we adopt the LSTM architecture as described in [48].

In order to obtain the decision scores over  $C$  classes at the time step  $t$ , we apply a softmax layer on the hidden state of time step  $t$  of LSTM to compute the probabilities as:

$$s_c = \frac{\exp(\mathbf{w}_c^T \mathbf{h}^t + b_c)}{\sum_{c' \in C} \exp(\mathbf{w}_{c'}^T \mathbf{h}^t + b_{c'})} \quad (1)$$

where  $s_c, \mathbf{w}_c,$  and  $b_c$  are the prediction score, the corresponding weight vector and bias term of the  $c$ th class, respectively. The LSTM is trained with the backpropagation through time (BPTT) algorithm [39], which unrolls the model along the time dimension into a feed forward neural nets and backpropagate the gradients.

### 4 Multi-granular score distribution fusion for action recognition (MSD)

In the Sect. 3, we have learnt a hierarchical video representation from the multiple granularities of video, and then through a softmax layer, we can compute the decision score for each component from each stream over each action class. Actually we can simply use the *Mean* or *Max* pooling, but as we have explained in the introduction section, different actions may span different granularities. Therefore, we try to figure out a more adaptive fusion scheme to fuse all the decision scores, which is expected to be able to attach different importance to different component. Inspired by the idea of addressing the temporal ambiguity of actions by learning score distribution in [5], we devise a novel fusion scheme based on the multi-granular score fusion. Specifically, an improved action recognition score will be obtained by automatically aligning the relative importance to each component from all the streams based on the score distribution.

#### 4.1 Formulation

Firstly, we put together the decision scores of all components from the multiple streams to form a distribution matrix as:

$$\mathbf{S} = (\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_C) \in \mathfrak{R}^{L \times C}, \tag{2}$$

where  $\mathbf{s}_c \in \mathfrak{R}^L$  denotes the score vector of  $L$  components from all streams on the  $c^{th}$  action class. Next, we *sort* the score distribution matrix by column like:

$$sort(\mathbf{S}) = (sort(\mathbf{s}_1), sort(\mathbf{s}_2), \dots, sort(\mathbf{s}_C)) \in \mathfrak{R}^{L \times C}, \tag{3}$$

where  $sort(\mathbf{s}_c) \in \mathfrak{R}^L$  is to reorder all elements of the vector  $\mathbf{s}_c$  in descending order. With  $L$  large enough, we can safely presume that  $sort(\mathbf{s}_c)$  can approximate the score distribution of this video over the  $c^{th}$  action class.

Then what we are going to do next is to train a binary classifier for each action class to distinguish whether a score distribution belongs to this class. The merit of training such a distribution classifier is that we can make full use of the information from all the components, rather than a summary statistic (*Mean*) or a extreme one (*Max*). And as we observe, the score distributions vary from class to class. When training the distribution classifier for a certain class, the positive samples consist of the score distributions from this class, while the negative ones come from the other classes. Figure 4 gives an example on the “Surfing” class to illustrate how multi-granular score fusion works.

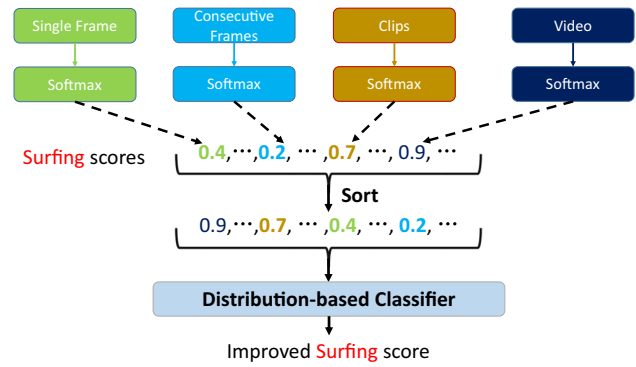


Fig. 4 An example of Multi-granular score fusion over class “Surfing”

More formally, we can formulate the training process as the following optimization problem:

$$\min_{\mathbf{w}_c, b_c} \sum_{i=1}^N \sum_{c=1}^C \max\{1 - y_i^c (\mathbf{w}_c \cdot sort(\mathbf{s}_c) + b_c), 0\} \tag{4}$$

$$\text{s.t.} \quad \sum_{l=1}^L w_c^l = 1, \quad c = 1, \dots, C, \tag{5}$$

$$w_c^1 \geq w_c^2 \geq \dots \geq w_c^L \geq 0, \quad c = 1, \dots, C. \tag{6}$$

For the  $c^{th}$  action class, the weight vector  $\mathbf{w}_c$  and the bias item  $b_c$  form a classifier to compute a confidence score over this action class for every score distribution. The objective function (4) is to minimize the sum of hinge loss. The constraint (5) requires the weights to have unit sum because we are assigning weights to each component. The reason for the constraint (6), which requires the weights to be monotonic and nonnegative, is that  $sort(\mathbf{s}_c)$  is actually classification scores in descending order and naturally we want to attach more importance to the components with high classification scores. By inspecting the two constraints, we can infer that the feasible set of  $\mathbf{w}_c$  contains two special cases:

1.  $w_c^1 = 1, w_c^2 = \dots = w_c^L = 0$ : *Max* pooling.
2.  $w_c^1 = w_c^2 = \dots = w_c^L = \frac{1}{L}$ : *Mean* pooling.

#### 4.2 Solution

Although we can train the distribution classifier in a standalone mode, we try to integrate the training process with the former video representation learning to form an “end-to-end” framework. But the two constraints make it difficult to train the framework in an “end-to-end” manner. To do this, we relax the two constraints into the objective function  $J$  by appending two penalty terms as

$$J = \mathcal{L} + \alpha \sum_{c=1}^C \|\mathbf{w}_c\|^2 + \beta \sum_{c=1}^C \left(1 - \sum_{l=1}^L w_c^l\right)^2 + \gamma \sum_{c=1}^C \sum_{l=1}^L m_c^l, \quad (7)$$

$$m_c^l = \begin{cases} w_c^{l+1} - w_c^l, & \text{if } w_c^{l+1} > w_c^l, \\ 0, & \text{if } w_c^{l+1} \leq w_c^l, \end{cases} \quad l = 1, \dots, L \text{ and } w_c^{L+1} = 0, \quad (8)$$

where the first part  $\mathcal{L}$  is the hinge loss in Eq. (4), the second is a regularization term preventing overfitting, followed by two penalty terms.  $\alpha, \beta, \gamma$  are the tunable hyperparameters.

Finally, the objective function  $J$  is minimized with regard to  $\{\mathbf{w}_c\}_{c=1}^C$  and the gradients are calculated by

$$\frac{\partial J}{\partial w_c^l} = \frac{\partial \mathcal{L}}{\partial w_c^l} + 2\alpha w_c^l - 2\beta(1 - w_c^l) + \gamma \left( \frac{\partial m_c^l}{\partial w_c^l} + \frac{\partial m_c^{l-1}}{\partial w_c^l} \right), \quad (9)$$

$$\frac{\partial m_c^l}{\partial w_c^l} = \begin{cases} -1, & \text{if } w_c^{l+1} > w_c^l \\ 0, & \text{if } w_c^{l+1} \leq w_c^l \end{cases} \quad l = 1, \dots, L, \quad (10)$$

$$\frac{\partial m_c^{l-1}}{\partial w_c^l} = \begin{cases} 1, & \text{if } w_c^l > w_c^{l-1} \\ 0, & \text{if } w_c^l \leq w_c^{l-1} \end{cases} \quad l = 2, 3, \dots, L \text{ and } \frac{\partial m_c^0}{\partial w_c^l} = 0. \quad (11)$$

After the optimization of  $J$  in Eq. (7), we can obtain the optimal  $\{\mathbf{w}_c\}_{c=1}^C$ . With this, we compute the final improved action score for the video as

$$p_c = \mathbf{w}_c \times \text{sort}(\mathbf{s}_c) + b_c. \quad (12)$$

The gradient is also backpropagated to the score distribution. As the order of the decision scores have been changed by the *sort* function, we need to store the index of the sorted scores in original vector in the forward process and propagate the gradients to the corresponding element in the original vector when backpropagating. This practice is very similar to that of Max pooling layer. After incorporating the learning of distribution classifier into the previous framework, we can train the whole architecture in an “end-to-end” fashion using the mini-batch SGD and standard backpropagation algorithm, which are implemented in the deep learning framework *Caffe* [8].

## 5 Experimental setup

### 5.1 Datasets

We adopt three popular datasets to evaluate the proposed hierarchical video representation and multi-granular score fusion scheme.

**HMDB51** [13] and **UCF101** [31] The UCF101 dataset is one of the most popular action recognition benchmarks. It consists of 13,320 videos clips from 101 action categories. The action categories are divided into five groups: Human–Object Interaction, Body-Motion Only, Human–Human Interaction, Playing Musical Instruments, and Sports. The HMDB51 dataset contains 6849 video clips divided into 51 action categories, each containing a minimum of 101 clips. The experimental setup is the same for both datasets, and three training/test splits are provided by the dataset organizers. Each split in UCF101 includes about 9.5K training and 3.7K test video, while a HMDB51 split contains 3.5K training and 1.5K test videos. Following [28], we conduct our analyses of different streams on the first split of the UCF101 and HMDB51 datasets. The average accuracy over three splits on both datasets is reported when compared with the state-of-the-art techniques. What needs noting is that the average length of videos in UCF101 and HMDB51 is about 6 seconds, which is much shorter than that of the videos from the real world.

**Columbia consumer videos (CCV)** [9] The CCV dataset contains 9,317 YouTube videos annotated according to 20 classes, which are mainly events like “Soccer,” “Cat,” “Wedding Ceremony,” “Beach.” We follow the convention defined in [9] to use a training set of 4,659 videos and a test set of 4,658 videos. The results are evaluated by average precision (AP) over each class, and we report the Mean AP (mAP) as the overall measure when compared with the baselines. Comparing with UCF101 and HMDB51, the videos in CCV usually last about several minutes and the categories are relatively more complex and high level than the simple actions in UCF101 and HMDB51.

### 5.2 Implementation details

**Frame stream** We uniformly select 25 frames per video and adopt the VGG\_19 [29] to extract frame features. The VGG\_19 is first pre-trained with the ILSVRC-2012 training set of 1.2 million images and then fine-tuned by using the video frames, which is observed to be much better than training from scratch. Following [28], we also use data augmentation like cropping and flipping. The learning rate starts from  $10^{-3}$  and decreases to  $10^{-4}$  after 14,000 iterations, then to  $10^{-5}$  after 20,000 iterations. For temporal modeling, we extract the outputs of 4096-way fc6 layer from VGG\_19 as inputs and adopt one-layer LSTM. We conduct experiments with different number of hidden states in LSTM. The LSTM weights are learnt by using the BPTT algorithm with a mini-batch size of 10. The learning rate starts from  $10^{-2}$  and decreases to  $10^{-3}$  after 100K iterations. The training is stopped after 150,000 iterations.



**Motion stream** We compute the optical flow between consecutive frames using the GPU implementation of [1] in OpenCV toolbox. The optical flow is converted to a flow “image” by linearly rescaling horizontal ( $x$ ) and vertical ( $y$ ) flow values to  $[0, 255]$  range. The transformed  $x$  and  $y$  flows are the first two channels for the flow image, and the third channel is created by calculating the flow magnitude. Moreover, the settings of VGG\_19 and LSTM are the same with frame stream.

**Clip stream** We define a *clip* as consecutive 16 frames, which is the same setting as [35]. The C3D is exploited to model video clip, which is pre-trained on Sports-1M [10] dataset with 1.1 million sports videos and then fine-tuned on UCF101 and HMDB51, respectively. As designed in C3D architecture, the input of C3D model is 16-frame clip and we uniformly sample 20 clips in each video. The learning rate starts from  $10^{-4}$  and decreases to  $10^{-5}$  after 10,000 iterations and then the training is stopped after 20,000 iterations. Again, the LSTM setting is the same with frame stream.

**Video stream.** The settings of video stream are similar to the clip stream. The only difference is that we do not involve LSTM after C3D and simply fuse the features of all video clips by *Mean* pooling to generate the video-level representations.

## 6 Results and analyses

In this section, we show the experiment results and analyses. As we can see, each granularity of the hierarchical representation of video can be used for action recognition solely. So we first evaluate the performance of different streams on

UCF101 and HMDB51, followed by the evaluation of multi-granular score fusion. Then we compare the performance of our whole framework with the state-of-the-art approaches on these two dataset. Next we evaluate the proposed multi-granular framework on the more complex dataset CCV for video classification. Last but not the least, we report the running time of each stream and the whole framework to justify the feasibility of our framework in the real applications.

### 6.1 Evaluation of frame and motion streams

Table 1 shows the results and comparisons of frame and motion streams on UCF101 (split 1). First, we examine the influence of using 2D CNNs with different depths (AlexNet and VGG\_19) on the frame and motion stream. As we can see in the first and third rows in Table 1a, compared with AlexNet [12] (8 weight layers), VGG\_19 [29] (19 weight layers) exhibits significantly better performance (more than 10 percent gain) on frame stream. The boosting performance should be credited to the increased depth on the CNN architecture, which can learn a better representation for the scene and object information in the frames. But this case does not apply to the motion stream, and we can only get very marginal improvement when using the deeper CNN. We speculate that deeper CNN suffers from the overfitting problem on the motion stream, which can be explained by two reasons: (1) optical flow “image” is relatively simple and does not contain as many details as natural images and (2) the big gap between the optical flow “image” and natural images from ImageNet makes the transfer learning work not very well. We also evaluate the LSTM to investigate the significance of leveraging the long-term temporal clues for action recognition. On UCF101, LSTM can improve both the frame

**Table 1** Accuracy of *frame* and *motion* streams on UCF101 (split 1)

(a) The accuracy of different 2D CNN and LSTM used on <i>frame</i> and <i>motion</i> streams. The results are reported for late fusion		
Training setting	Frame (%)	Motion (%)
AlexNet	67.1	68.4
AlexNet + LSTM	69.3	70.3
VGG_19	77.9	70.6
VGG_19 + LSTM	79.3	73.8
VGG_19 + LSTM + Augmentation	80.2	74.6
(b) The effect of hidden layer size in the LSTM (VGG_19)		
Hidden layer size	Frame (%)	Motion (%)
128	78.2	71.2
256	78.8	72.6
512	79.1	73.5
1024	79.3	73.8
2048	78.5	73.1

**Table 2** Accuracy of *clip* and *video* streams on UCF101 (split 1) and HMDB51 (split 1)

(a) The comparisons of using features from different layers of C3D on <i>clip</i> stream				
Dataset	fc6 (%)	fc7 (%)	Prob (%)	fc6+LSTM (%)
HMDB51	50.36	48.65	38.97	51.3
UCF101	83.11	81.23	69.81	83.9
(b) The comparisons of using the features from different layers of C3D on <i>video</i> stream				
Dataset	fc6 (%)	fc7 (%)	Prob (%)	
HMDB51	51.09	48.52	39.10	
UCF101	83.77	80.76	67.01	

**Table 3** Comparisons of the proposed MSD with Mean and Max fusion schemes in terms of accuracy on three splits of HMDB51 and UCF101

Dataset	Split 1			Split 2			Split 3		
	Mean (%)	Max (%)	MSD (%)	Mean (%)	Max (%)	MSD (%)	Mean (%)	Max (%)	MSD (%)
HMDB51	61.5	59.6	63.1	61.8	59.5	63.5	62.1	60.1	64.1
UCF101	89.6	87.6	90.2	89.6	87.4	90.3	91.2	88.1	91.9

and motion streams. Furthermore, when augmenting the test frame (flow “image”) by cropping and flipping four corners and the center of the frame and averaging the scores across the frame and its crops, the performance can achieve 80.2 and 74.6% on frame and optical flow, respectively.

Besides, we evaluate the influence of hidden state size of LSTM as seen in Table 1b. In general, increasing the hidden layer size of LSTM can lead to the improvement of the accuracy. When the hidden layer size reaches 1024 in our case, no further improvement can be obtained on both frame and optical flow streams. Note that the performances are reported based on the original frame or optical flow image with only cropping center and no flipping operation in this comparison.

## 6.2 Evaluation of clip and video streams

Next, we turn to measure the performance of the clip and video streams in terms of features extracted from different layers of 3D CNN (C3D) on UCF101 and HMDB51. We extract activations of the C3D layers: fc6, fc7, and prob for each clip. The recognition score is computed by late fusing the predicted score on each clip, and the accuracy comparison by using the outputs from these three different layers is shown in Table 2a. As indicated by our results, the recognition using the C3D feature of fc6 layer leads to a larger performance boost against the C3D features of fc7 and prob layers. Furthermore, the accuracy by using the feature of fc6 can achieve 51.3 and 83.9% on HMDB51 and UCF101 after longer-term temporal modeling with LSTM networks, respectively. The features for video stream are computed by averaging the video clip features separately for each type of feature, and Table 2b reports the comparison of different C3D features on video stream. Similar to the observations

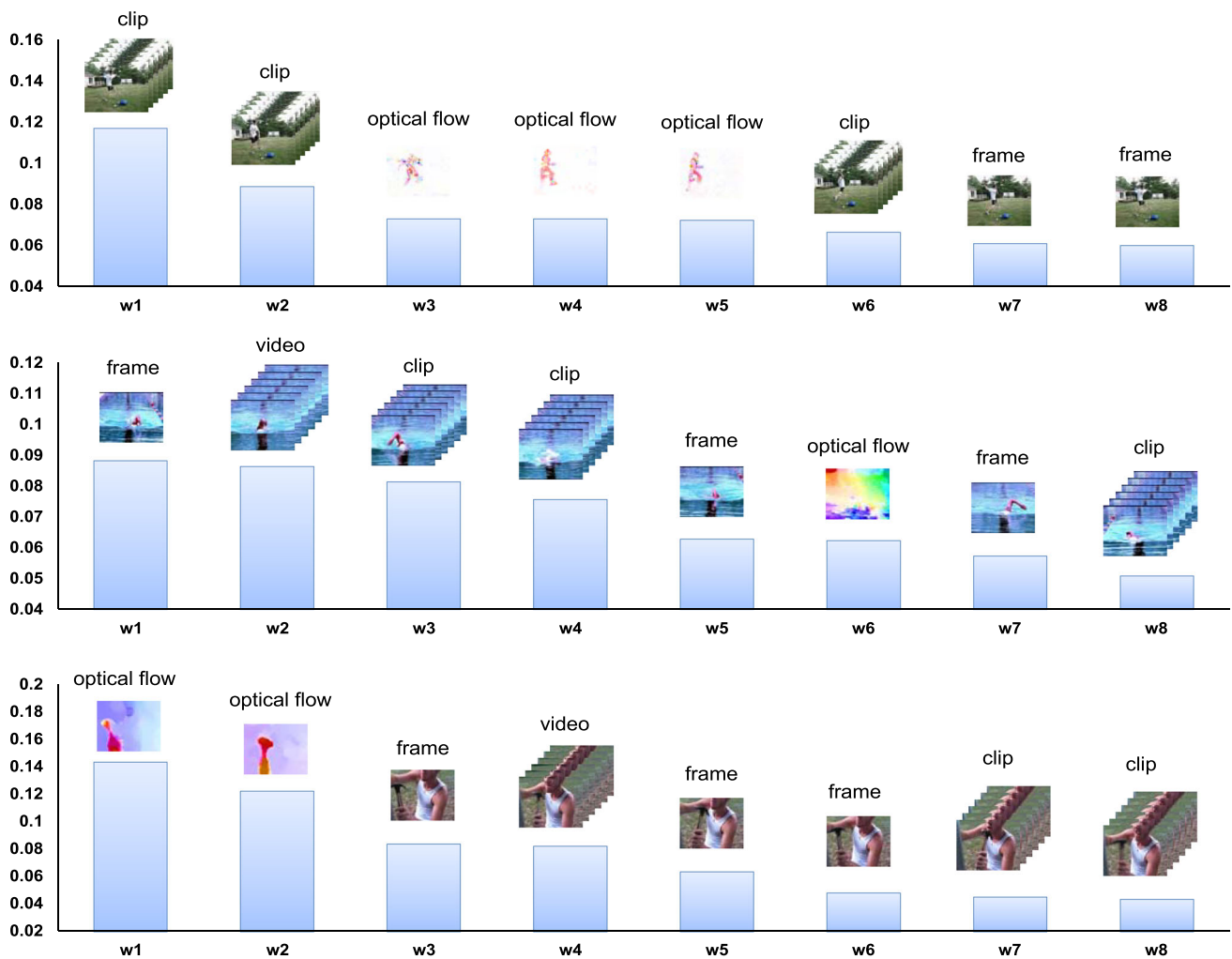
on video clip stream, the features of fc6 layer achieve the best performance among all the three layers with a large margin.

## 6.3 Evaluation of multi-granular score fusion

Here we evaluate the complete multi-granular architecture, which combines the four streams with the MSD fusion method. Table 3 details the accuracy across different fusion strategies on three splits of HMDB51 and UCF101, respectively. MSD consistently outperforms Max and Mean in every split of both two datasets. The improvement is observed in different types of actions. For instance, the actions “playing piano” and “biking” are better fused with Mean as the videos relevant to the two actions are consistent in content. On the other hand, the recognition of actions “cliff diving” and “basketball dunk” shows much better results with Max fusion. In the experiment, MSD boosts the accuracy of these actions. Figure 5 shows the top eight weights learnt by MSD and their corresponding components of three exemplary videos from category “baseball pitch,” “front crawl,” and “hammering.” We can easily see that all the eight components are highly related to each action. More importantly, the top eight components come from four different streams, which validates the effectiveness of MSD on fusing multi-granular information.

## 6.4 Comparisons with the state of the art

We compare with several state-of-the-art techniques on three splits of UCF101 and HMDB51. As shown in Table 4, our multi-granular spatial–temporal architecture exhibits the highest performance on UCF101 dataset. It makes the



**Fig. 5** Examples showing the *top* eight weights learned by the MSD and their corresponding components in a video (*top* baseball pitch, *middle* front crawl, *bottom* hammering). We can see that MSD is able to learn the contributions from different components for particular actions.

For example, two *clip* components play important roles for recognizing “baseball pitch,” while two *motion* (optical flow) components contribute more to the recognition of “hammering”

improvement over [37] by 0.5%, which is generally considered as a significant progress on this dataset. On the HMDB51, the works [14,37] with competitive results are based on the motion trajectory, while our approach fully relies on the deep learning architecture and is trained end-to-end. Compared with the two-stream model [28], our architecture by additionally incorporating more temporal modeling and utilizing a sophisticated fusion strategy leads to a performance boost on both datasets. It is also worth noting that in the training of the HMDB51 dataset, [28] exploit UCF101 as additional training data through multitask learning while our architecture is only trained on the HMDB51 data. In addition, the recent works in [3,20,32] also use the LSTM to exploit the temporal information. Our method achieves more promising results as more dimensions of cues are taken into account.

### 6.5 Evaluation on CCV

Now we turn to evaluate our framework on the more challenging dataset CCV with longer video and more complicated categories. Table 5 shows the performance of each stream and the whole framework as well as the comparisons with the state-of-the-art results. As we can see, the performance of motion stream is much lower than that of frame stream, which is inconsistent with that of UCF101 and HMDB51. The reason is twofold. On the one hand, many categories in CCV are of very high level and do not care much about the short-term motion in video, such as “graduation,” “wedding reception,” and “beach.” On the other, since the average duration (around 80 seconds) of the videos of CCV is about 10 times longer than that of UCF101 and HMDB51 and the contents in CCV are more complex and noisy, the optical

**Table 4** Performance in terms of Mean accuracy (over three splits) on UCF101 and HMDB51

Method	UCF101 (%)	HMDB51 (%)
IDT [36]	85.9	57.2
IDT w/encodings [24]	87.9	61.1
MIFS [14]	89.1	65.1
“Slow Fusion” ConvNet [10]	65.4	–
LRCN [3]	82.9	–
C3D [35]	85.2	–
Two-stream model [28]	88.0	59.4
Composite LSTM [32]	84.3	–
CNN + IDT [49]	89.6	–
Temporal pooling + LSTM [20]	88.6	–
TDD [37]	90.3	63.2
Ours	90.8	63.6

Please note that the methods in [14,24,36] are based on traditional dense trajectory which is computationally expensive, while the methods in [37,49] combine dense trajectory and deep learning-based algorithms. Our approach outperforms the deep learning-based methods without combination of dense trajectory [3,20,28] with a large margin. “–” means that the authors did not report their performance on this dataset. IDT improved dense trajectory [36]; MIFS multi-skip feature stacking [14]; LRCN long-term recurrent convolutional networks [3]; C3D: convolutional 3D [35]; TDD trajectory-pooled deep convolutional descriptor [37]

**Table 5** Comparisons with state-of-the-art results on CCV

CCV			
Liu et al. [17]	68.2%	Ours ( <i>frame</i> )	77.5%
Ye et al. [46]	64.0%	Ours ( <i>motion</i> )	59.3%
Jhuo et al. [7]	64.0%	Ours ( <i>clip</i> )	80.9%
Ma et al. [18]	63.4%	Ours ( <i>video</i> )	80.1%
Wu et al. [42]	70.6%	Ours ( <i>fusion</i> )	83.2%

The best performance is obtained by multi-granular score distribution fusion of our framework

**Table 6** Run time of different streams averaged over all test videos in UCF101 dataset (milliseconds)

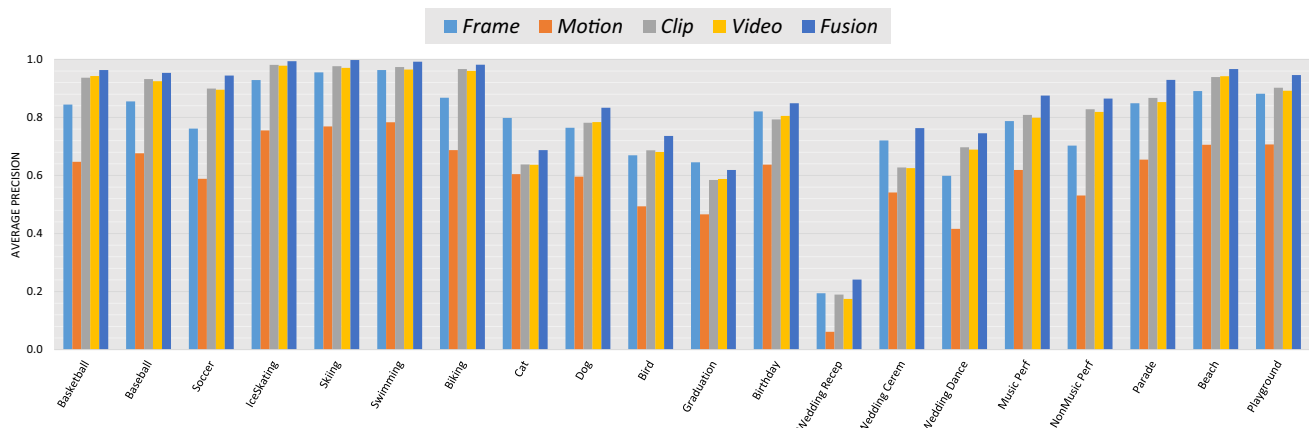
Stream	2D/3D CNN	LSTM	SUM
<i>frame</i>	750	12	762
<i>motion</i>	750	12	762
<i>clip</i>	490	10	500
<i>video</i>	490	–	490

flows are easily disrupted by the big camera motion and useful motion information might be overwhelmed by many disrupted optical flows. When we come to the fusion of all the streams, the performance of the whole framework is clearly better than each single stream, which verifies the effectiveness of the multi-granular score distribution fusion scheme. Table 5 also shows that our multi-granular framework is significantly better than all of the baselines, and we get around 10 percent performance gain over the best baseline [42].

To better understand the contribution of every stream in the multi-granular framework, we further report the performance on each class of CCV in Fig. 6. From the Fig. 6, *frame*, *clip*, and *video* streams, respectively, get the best performance over different classes. For example, *clip* stream achieves the best performance on classes like “Basketball” and “Soccer,” while the results of *frame* stream exceed other streams a big margin on “Cat” class. This proves that different actions may span different granularities. And the multi-granular score fusion of all streams can significantly improve the performance for almost all the classes, although the performance of *motion* stream is relatively low.

## 6.6 Efficiency

In addition to obtaining the superior classification accuracy, our framework also enjoys high computational efficiency. Table 6 lists the details run time of each stream averaged

**Fig. 6** Per-class performance on CCV, using *frame*, *motion*, *clip*, *video* and their *fusion*

over all test videos in UCF101 dataset. The experiments are conducted on a regular server (Intel Xeon 2.40GHz CPU and 256 GB RAM) with a single NVidia K80 GPU. As each stream could be executed in parallel and the fusion with MSD provides instant response, the average prediction time of our multi-granular architecture on each video in UCF101 is about 762 milliseconds, which is very efficient. This is much faster than trajectory-based approaches, e.g., IDT, which requires about seven minutes for each video in UCF101.

## 7 Conclusions and future work

We have proposed a multi-stream deep learning framework to learn a hierarchical video representation, which can exploit the information in video from a multitude of granularities including *frame*, consecutive frames (*motion*), *clip*, and the entire *video*. In the framework, we first apply two types of CNN at each granularity, i.e., 2D CNN on frame and motion streams, and 3D CNN on clip and video streams. The outputs of these CNNs are then used separately as inputs of the LSTM networks to model the long-term temporal dynamics. To verify the effectiveness of the proposed hierarchical video representation, we employ this representation for action recognition by integrating the information from all granularities with a novel score distribution fusion strategy. We conduct extensive experiments on two popular action recognition benchmarks with short videos (UCF101 and HMDB51) and another video classification dataset with complex and long videos (CCV). Our framework has achieved very impressive performance on all the three popular benchmarks. Results not only validate the effectiveness of each single stream, but also demonstrate that the multiple granularities of video are complementary and combining them can significantly boost the performance of action recognition.

There are several future directions. First, video action recognition can be enhanced by further considering audio information. The audio features can be exploited together with the current four streams to more comprehensively characterize the actions in videos. Second, the method for learning the representations of the entire video could be explored by using RNNs in an encoder–decoder framework. In addition, we will continue to conduct more in-depth investigations on how fusion weights of individual streams can be dynamically determined to boost action recognition performance.

## References

1. Brox T, Bruhn A, Papenbergh N, Weickert J (2004) High accuracy optical flow estimation based on a theory for warping. In: European conference on computer vision
2. Dollár P, Rabaud V, Cottrell G, Belongie S (2005) Behavior recognition via sparse spatio-temporal features. In: 2005 IEEE international workshop on visual surveillance and performance evaluation of tracking and surveillance, IEEE, pp 65–72
3. Donahue J, Hendricks LA, Guadarrama S, Rohrbach M, Venugopalan S, Saenko K, Darrell T (2014) Long-term recurrent convolutional networks for visual recognition and description. arXiv preprint: [arXiv:1411.4389](https://arxiv.org/abs/1411.4389)
4. Graves A, Mohamed A-r, Hinton G (2013) Speech recognition with deep recurrent neural networks. In: 2013 IEEE international conference on acoustics, speech and signal processing, IEEE, pp 6645–6649
5. Hoai M, Zisserman A (2014) Improving human action recognition using score distribution and ranking. In: Asian conference on computer vision
6. Hochreiter S, Schmidhuber J (1997) Long short-term memory. *Neural Comput* 9(8):1735–1780
7. Jhuo I-H, Ye G, Gao S, Liu D, Jiang Y-G, Lee D, Chang S-F (2014) Discovering joint audio-visual codewords for video event detection. *Mach Vis Appl* 25(1):33–47
8. Jia Y, Shelhamer E, Donahue J, Karayev S, Long J, Girshick R, Guadarrama S, Darrell T (2014) Caffe: Convolutional architecture for fast feature embedding. arXiv preprint: [arXiv:1408.5093](https://arxiv.org/abs/1408.5093)
9. Jiang Y-G, Ye G, Chang S-F, Ellis D, Loui AC (2011) Consumer video understanding: a benchmark database and an evaluation of human and machine performance. In: Proceedings of ACM international conference on multimedia retrieval
10. Karpathy A, Toderici G, Shetty S, Leung T, Sukthankar R, Fei-Fei L (2014) Large-scale video classification with convolutional neural networks. In Proceedings of the IEEE conference on computer vision and pattern recognition
11. Klaser A, Marszałek M, Schmid C (2008) A spatio-temporal descriptor based on 3d-gradients. In: BMVC 2008–19th British machine vision conference, British Machine Vision Association, pp 275:1–10
12. Krizhevsky A, Sutskever I, Hinton GE (2012) Imagenet classification with deep convolutional neural networks. In: Advances in neural information processing systems
13. Kuehne H, Jhuang H, Garrote E, Poggio T, Serre T (2011) HMDB: a large video database for human motion recognition. In: Proceedings of the IEEE international conference on computer vision
14. Lan Z, Lin M, Li X, Hauptmann AG, Raj B (2015) Beyond gaussian pyramid: multi-skip feature stacking for action recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition
15. Laptev I, Lindeberg T (2003) Space-time interest points. In: Proceedings of the IEEE international conference on computer vision
16. Li Q, Qiu Z, Yao T, Mei T, Rui Y, Luo J (2016) Action recognition by learning deep multi-granular spatio-temporal video representation. In: Proceedings of ACM international conference on multimedia retrieval
17. Li W, Zhang Z, Liu Z (2008) Expandable data-driven graphical modeling of human actions based on salient postures. *IEEE Trans Circuits Syst Video Technol* 18(11):1499–1510
18. Ma AJ, Yuen PC (2014) Reduced analytic dependency modeling: Robust fusion for visual recognition. *Int J Comput Vis* 109(3):233–251
19. Ma Y-F, Hua X-S, Lu L, Zhang H-J (2005) A generic framework of user attention model and its application in video summarization. *IEEE Trans MM* 7(5):907–919
20. Ng JY-H, Hausknecht M, Vijayanarasimhan S, Vinyals O, Monga R, Toderici G (2015) Beyond short snippets deep networks for video classification. In: Proceedings of the IEEE conference on computer vision and pattern recognition

21. Ngo C-W, Ma Y-F, Zhang H-J (2005) Video summarization and scene detection by graph modeling. *IEEE Trans CSVT* 15(2):296–305
22. Pan Y, Li Y, Yao T, Mei T, Li H, Rui Y (2016) Learning deep intrinsic video representation by exploring temporal coherence and graph structure. In: *International joint conference on artificial intelligence*
23. Pan Y, Mei T, Yao T, Li H, Rui Y (2016) Jointly modeling embedding and translation to bridge video and language. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*
24. Peng X, Wang L, Wang X, Qiao Y (2014) Bag of visual words and fusion methods for action recognition: Comprehensive study and good practice. *arXiv preprint: [arXiv:1405.4506](https://arxiv.org/abs/1405.4506)*
25. Qiu Z, Li Q, Yao T, Mei T, Rui Y (2015) Msr asia msm at thumos challenge 2015. In: *CVPR THUMOS challenge workshop*
26. Qiu Z, Yao T, Mei T (2016) Deep quantization: encoding convolutional activations with deep generative model. *arXiv preprint: [arXiv:1611.09502](https://arxiv.org/abs/1611.09502)*
27. Scovanner P, Ali S, Shah M (2007) A 3-dimensional sift descriptor and its application to action recognition. In: *ACM international conference on multimedia*, ACM, pp 357–360
28. Simonyan K, Zisserman A (2014) Two-stream convolutional networks for action recognition in videos. In: *Advances in neural information processing systems*, pp 568–576
29. Simonyan K, Zisserman A (2015) Very deep convolutional networks for large-scale image recognition. In: *International conference on learning representations*
30. Snoek CGM, van de Sande KEA, de Rooij O et al (2008) The mediamill trecvid 2008 semantic video search engine. In: *NIST TRECVID workshop*
31. Soomro K, Zamir AR, Shah M (2012) UCF101: A dataset of 101 human action classes from videos in the wild. *CRCV-TR-12-01*
32. Srivastava N, Mansimov E, Salakhutdinov R (2015) Unsupervised learning of video representations using LSTMs. In: *Proceedings of international conference on machine learning*
33. Sutskever I, Vinyals O, Le QV (2014) Sequence to sequence learning with neural networks. In: *Advances in neural information processing systems*, pp 3104–3112
34. Szegedy C, Liu W, Jia Y, Sermanet P, Reed S, Anguelov D, Erhan D, Vanhoucke V, Rabinovich A (2015) Going deeper with convolutions. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp 1–9
35. Tran D, Bourdev LD, Fergus R, Torresani L, Paluri M (2014) Learning spatiotemporal features with 3d convolutional networks. *arXiv preprint: [arXiv:1412.0767](https://arxiv.org/abs/1412.0767)*
36. Wang H, Schmid C (2013) Action recognition with improved trajectories. In: *Proceedings of the IEEE international conference on computer vision*, pp 3551–3558
37. Wang L, Qiao Y, Tang X (2015) Action recognition with trajectory-pooled deep-convolutional descriptors. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*
38. Wei X-Y, Jiang Y-G, Ngo C-W (2011) Concept-driven multi-modality fusion for video search. *IEEE Trans CSVT* 21(1):62–73
39. Werbos PJ (1990) Backpropagation through time: what it does and how to do it. *Proc IEEE* 78(10):1550–1560
40. Wilkins P, Ferguson P, Smeaton AF (2006) Using score distributions for query-time fusion in multimedia retrieval. In: *ACM SIGMM international workshop on Multimedia information retrieval*
41. Willems G, Tuytelaars T, Van Gool L (2008) An efficient dense and scale-invariant spatio-temporal interest point detector. In: *European conference on computer vision*, pp 650–663. Springer,
42. Wu Z, Jiang Y-G, Wang J, Pu J, Xue X (2014) Exploring inter-feature and inter-class relationships with deep neural networks for video classification. In: *ACM international conference on multimedia*, pp 167–176. ACM
43. Yao T, Mei T, Ngo C-W, Li S (2013) Annotation for free: video tagging by mining user search behavior. In: *ACM international conference on multimedia*
44. Yao T, Mei T, Rui Y (2016) Highlight detection with pairwise deep ranking for first-person video summarization. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*
45. Yao T, Ngo C-W, Mei T (2013) Circular reranking for visual search. *IEEE Trans Image Process* 22(4):1644–1655
46. Ye G, Liu D, Jhuo I-H, Chang S-F (2012) Robust late fusion with rank minimization. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp 3021–3028. IEEE
47. Yuan X, Lai W, Mei T, Hua X-S, Wu X-Q, Li S (2006) Automatic video genre categorization using hierarchical svm. In: *2006 International conference on image processing*, pp 2905–2908. IEEE
48. Zaremba W, Sutskever I (2014) Learning to execute. *arXiv preprint: [arXiv:1410.4615](https://arxiv.org/abs/1410.4615)*
49. Zha S, Luisier F, Andrews W, Srivastava N, Salakhutdinov R. (2015) Exploiting image-trained CNN architectures for unconstrained video classification. *arXiv preprint: [arXiv:1503.04144](https://arxiv.org/abs/1503.04144)*