# Learning Latent Sub-events in Activity Videos Using Temporal Attention Filters

**AJ Piergiovanni***, **Chenyou Fan***, **and Michael S. Ryoo**
School of Informatics and Computing, Indiana University, Bloomington, IN 47408
{ajpiergi,fan6,mryoo}@indiana.edu

## Abstract

In this paper, we newly introduce the concept of *temporal attention filters*, and describe how they can be used for human activity recognition from videos. Many high-level activities are often composed of multiple temporal parts (e.g., sub-events) with different duration/speed, and our objective is to make the model explicitly learn such temporal structure using multiple attention filters and benefit from them. Our temporal filters are designed to be fully differentiable, allowing end-of-end training of the temporal filters together with the underlying frame-based or segment-based convolutional neural network architectures. This paper presents an approach of learning a set of optimal static temporal attention filters to be shared across different videos, and extends this approach to dynamically adjust attention filters per testing video using recurrent long short-term memory networks (LSTMs). This allows our temporal attention filters to learn latent sub-events specific to each activity. We experimentally confirm that the proposed concept of temporal attention filters benefits the activity recognition, and we visualize the learned latent sub-events.

## Introduction

Human activity recognition is the problem of identifying events performed by humans given a video input. It is formulated as a binary (or multiclass) classification problem of outputting activity class labels, and researchers have been studying better features, representations, and learning algorithms to improve the classification (Aggarwal and Ryoo 2011). Such classification not only allows categorization of videos pre-segmented to contain one single activity, but also enables the 'detection' of activities from streaming videos together with temporal window proposal methods like the sliding window or selective search (Ryoo et al. 2015). Activity recognition is an important problem with many societal applications including smart surveillance, video search/retrieval, intelligent robots, and other monitoring systems.

Particularly, in the past 2-3 years, activity recognition approaches taking advantage of convolutional neural networks (CNNs) have received a great amount of attention. Motivated by the success of image-based object recognition using CNNs, researchers attempted developing CNNs for videos. Some approaches directly took advantage of image-based CNN architectures by applying them to every

---

*These authors contributed equally to the paper.

video frame (Jain, van Gemert, and Snoek 2014; Ng et al. 2015), while some tried to learn 3-D XYT spatio-temporal convolutional filters from short video segments (Tran et al. 2015). In order to represent each video, temporal pooling (e.g., max/average pooling) were often applied on top of multiple (sampled) per-frame or per-video-segment CNNs (Jain, van Gemert, and Snoek 2014; Karpathy et al. 2014; Simonyan and Zisserman 2014; Ng et al. 2015; Tran et al. 2015). Similar to the object recognition, these approaches obtained superior results compared to traditional approaches of using hand-crafted features.

However, in terms of learning and considering activities' temporal structure in videos, previous CNN approaches were limited. Many high-level activities are often composed of multiple temporal parts (i.e., sub-events) with different duration/speed, but approaches to *learn* explicit activity temporal structure together with CNN parameters have not been studied in depth. For instance, the typical strategy of taking max (or average) pooling over sampled per-frame or per-segment CNN responses (Jain, van Gemert, and Snoek 2014; Simonyan and Zisserman 2014; Ng et al. 2015; Tran et al. 2015) completely ignores such temporal structure in longer activity videos. (Ryoo, Rothrock, and Matthies 2015) showed a potential that making the system consider multiple video sub-intervals using a temporal pyramid benefits the recognition, but it was done with predetermined intervals without any learning. LSTM-based recurrent neural network approaches (Ng et al. 2015; Yeung et al. 2016) were able to process per-frame CNN responses sequentially, but no explicit sub-event or interval learning was attempted.

What we need instead is an approach that explicitly 'learns' to focus on important sub-intervals of the activity videos while also optimizing their temporal resolution for the recognition. This is a challenging problem since we want to make this sub-event learning done in an end-to-end fashion together with the training of underlying CNN parameters. Furthermore, it is often the case that ground truth labels of sub-events to be learned are not provided, making them latent variables.

This paper presents a new video classification approach that overcomes such limitations using temporal attention filters. We newly introduce the concept of fully differentiable *temporal attention filters* and describe how they can be learned and used to enable better recognition of human activ-
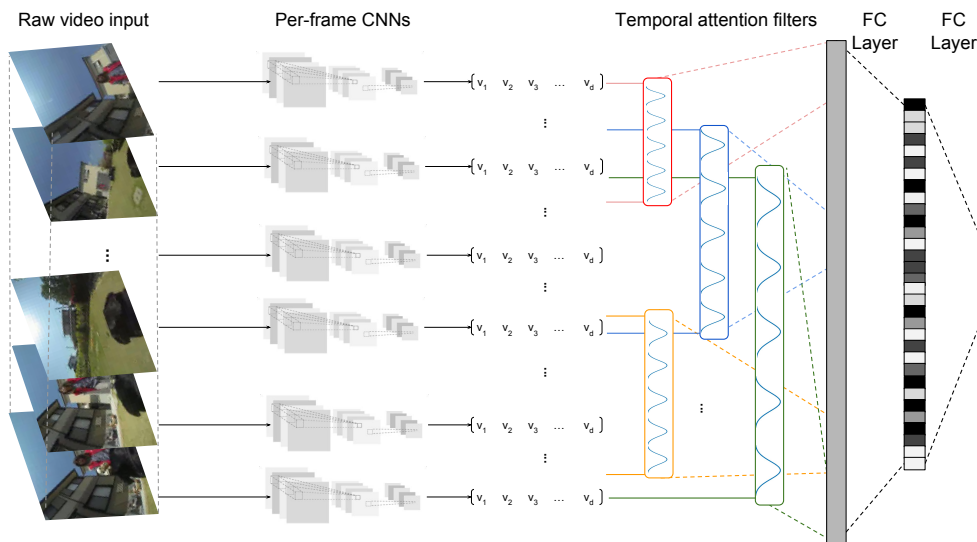
Figure 1: Illustration of our overall recognition architecture with temporal attention filters. $M$ number of temporal filters are learned to focus on different temporal part of video frame features (i.e., latent sub-events). Each filter is composed of a set of Gaussian filters which take a weighted sum of local information. Outputs of the temporal filters are concatenated, and attached with a fully connected layers to perform video classification.

ities from videos. The main idea is to make the network learn and take advantage of multiple temporal attention filters to be applied on top of per-frame CNNs (Figure 1). Each learned attention filter corresponds to a particular sub-interval of the activity the system should focus on (i.e., latent sub-events), and is represented with its center location, duration, and resolution in the relative temporal coordinate. Our approach abstracts per-frame (or per-segment) CNN responses within the sub-interval corresponding to the attention filter, allowing the system to use their results for the classification. Notably, our temporal attention filters are designed to be fully differentiable, motivated by the spatial attention filters for images (Gregor et al. 2015). This allows the *end-to-end* training of the parameters deciding attention filters; the system learns temporal attention filters jointly with the parameters of the underlying per-frame (or per-segment) CNNs. As a result, our temporal filters are trained to be optimized for the recognition tasks, automatically learning its location/scale from the training data without sub-event labels.

The paper not only presents an approach of learning optimal static temporal attention filters to be shared across different videos, but also present an approach of dynamically adjusting attention filters per testing video using recurrent long short-term memory cells (LSTMs). Instead of learning static temporal filters who location/duration/resolution is shared by all videos, our LSTM based approach dynamically and adaptively adjusts its filter parameters depending on the video, by going through multiple iterations. Our proposed approach is able to function in conjunction with any per-frame or per-video-segment CNNs as well as with other types of feature representations (e.g., Fisher vectors), making it very generally applicable for many video understanding scenarios.

## Previous works

As described in the introduction, the direction of using convolutional neural networks for video classification is becoming increasingly popular, since it allows end-to-end training of convolutional filters optimized for the training data. (Simonyan and Zisserman 2014) used optical flows in addition to image feature. (Ng et al. 2015) tested multiple different types of pooling strategies on top of per-frame CNNs, and found that the simple global max pooling of per-frame features over the entire interval performs the best for sports videos. (Karpathy et al. 2014) also tried multiple different (temporal) pooling strategies, gradually combining per-frame CNN responses over a short interval using their 'slow fusion'. (Tran et al. 2015) proposed to do XYT convolution, learning space-time convolutional filters. (Wang, Qiao, and Tang 2015) used local CNN feature maps around tracklets in videos.

However, these prior works focused only on capturing dynamics in very short video intervals without much consideration on long-term temporal structure of activity videos. Optical flow only captures differences between two consecutive frames (Simonyan and Zisserman 2014). Even with the video-based 3-D XYT CNNs (Tran et al. 2015) or trajectory CNNs (Wang, Qiao, and Tang 2015), only the temporal dynamics within short intervals with a fixed duration (e.g., 15 frames) were captured without considering longer-term structure or attempting to learn latent sub-events. (Ryoo, Rothrock, and Matthies 2015) showed a potential that considering temporal structure in terms of sub-intervals (e.g., temporal pyramid) may benefit the recognition, but they did not attempt any learning. Similarly, (Li et al. 2016) considered multiple different temporal scales, but learning of how the system should choose such scales were not attempted.

(Varol, Laptev, and Schmid 2016) also used fixed intervals. Recurrent neural networks such as LSTMs were also used to model sequences (Ng et al. 2015), but they were unable to explicitly consider different temporal sub-intervals and their structure. That is, learning to consider different intervals with different temporal resolution was not possible, and no sub-event learning as involved. (Yeung et al. 2016) proposed the use of LSTM to make the system focus on different frames of videos, but it was unable to represent intervals.

The main contribution of this paper is the introduction of the temporal attention filters that allow their end-to-end training together with underlying CNN architectures. We illustrate that our temporal attention filters can be learned to focus on different temporal aspects of videos (i.e., intervals with different temporal resolutions), and experimentally confirm that such learning benefits the activity recognition. The main difference between our approach and previous temporal structure learning methods for activity recognition (e.g., (Niebles, Chen, and Fei-Fei 2010; Ryoo and Matthies 2013)) is that our temporal filters are designed to be fully differentiable, which allows their joint learning and testing with modern CNN architectures.

To our knowledge, this paper is the first paper to enable learning of latent temporal sub-events in an end-to-end fashion using CNN architectures for activity recognition.

## Recognition approach

We design our model as a set of temporal filters, each corresponding to a particular sub-event, placed on top of per-frame (or per-segment) CNN architectures (Figure 1). The idea is to train this fully differential model in an end-to-end fashion, jointly learning latent sub-events composing each activity, underlying CNN parameters, and the activity classifier.

In this section, we introduce the concept of *temporal attention filters*, which extends the spatial attention filter (Gregor et al. 2015) originally designed for digit detection and digit/object synthesis. Next, we present how our proposed model takes advantage of temporal attention filters to learn latent sub-events. The approach learns/mines temporal sub-events optimized for the classification without their ground truth annotations, such as 'stretching an arm' in the activity 'punching'.

### Temporal attention filters

Each temporal attention filter learns three parameters: a center $g$, a stride $\delta$ and a width $\sigma$. These parameters determine where the filter is placed and the size of the interval focused on. A filter consists of $N$ Gaussian filters separated by a stride of $\delta$ frames. The goal of this model is to learn where in the video the most useful features appear. Because the videos are of variable length, $\delta$ and $g$ are relative to the length of the video. Based on the attention model presented in (Gregor et al. 2015), we use the following equations to obtain the mean of the Gaussian filters:

$$g_n = 0.5 \cdot T \cdot (\widetilde{g}_n + 1)$$
$$\delta_n = \frac{T}{N-1} \widetilde{\delta}_n \qquad (1)$$
$$\mu_n^i = g_n + (i - 0.5N + 0.5)\delta_n$$

Using $\mu$ and $\sigma$, the $N$ Gaussian filters are defined by:

$$F_m[i,t] = \frac{1}{Z_m} \exp(-\frac{(t - \mu_m^i)^2}{2\sigma_m^2}) \qquad (2)$$
$$i \in \{0, 1, \ldots, N-1\}, \ t \in \{0, 1, \ldots, T-1\}$$

where $Z_m$ is a normalization constant. Figure 2 shows an illustration of our temporal attention filter.

If each frame has $D$-dimensional features, the filters, $F$, are applied to each dimension, taking the input of size $T \times D$ to $N \times D$ where $T$ is the number of frames in the video. That is, each temporal filter $F$ generates a $N \times D$-dimensional vector as an output for any video, which can be passed to a neural network for the classification. Since this model is fully differentiable, all parameters can be learned using gradient descent.

Let $f_m[i,d]$ be the output of our $m$th temporal attention filter, given the $T \times D$ dimensional input $x$. Each $f_m[i,d]$ describes the response from the $i$th Gaussian filter on the $d$th elements of the input vectors. Then,

$$f_m[i,d] = F_m[i,:] \cdot x[:,d] = \sum_{t=0}^{T-1} F_m[i,t] \cdot x[t,d] \qquad (3)$$
$$i \in \{0, 1, \ldots, N-1\}, \ d \in \{0, 1, \ldots, D-1\}$$

where $x[:,d]$ is the $T$-dimensional vector corresponding to the sequence of the $d$th element values in the underlying CNN feature vectors.

Figure 3 shows how each temporal attention filter is able to capture features from the corresponding sub-interval of the provided video with different temporal resolutions.

### Recognition with temporal attention filters

As described in Figure 1, we take advantage of multiple different temporal attention filters by placing them on top of a sequence of per-frame (or per-segment) CNN models. As a result, our model is able to focus on different sub-intervals of video inputs with different temporal resolutions. Outputs of each temporal filters are concatenated and are connected to fully connected layers performing activity classification. Each of our filters learns a latent sub-event, and concatenating the results allows the later FC-layers to look at the features for each sub-event and classify the activity based on
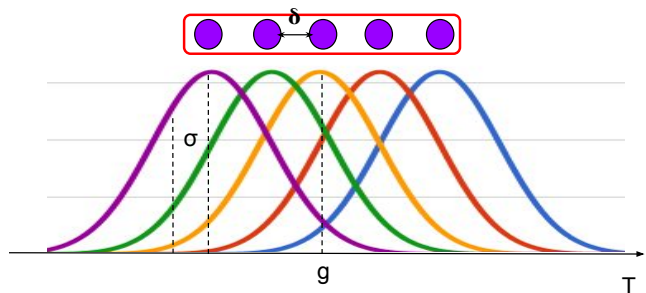


Figure 2: An illustration of our temporal attention filter. The filter is differentiable and represented with three parameters.
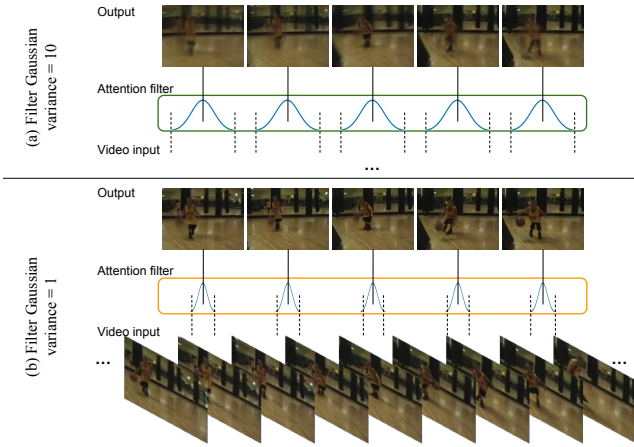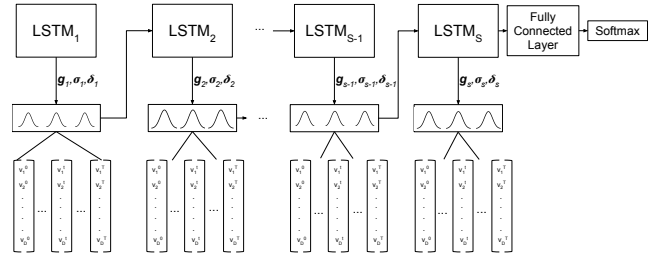
Figure 4: An illustration of our temporal attention filter using an LSTM. The LSTM provides the parameters for the filters at each time step. The filters provide input to the next iteration allowing the LSTM to adjust the location of the filters. After $S$ iterations, a fully-connected layer and softmax layer classify the video.

Figure 3: Examples illustrating how our temporal attention filters work once learned. They are shown with two different Gaussian filter variances. For this visualization, raw video frames are used as inputs directly. In our actual implementation, the input to temporal filters are not raw video frames but CNN feature responses obtained from frames.

them. If we denote the per-frame feature size as $D$ and we have $M$ number of temporal attention filters, each temporal filter generates the output of size $N \times D$, resulting the total dimensionality to be $M \times N \times D$. We used 2 fully connected layers (i.e., one hidden layer and one soft-max layer) for the classification.

Because of the property that our temporal filters are designed to be differentiable, we are able to backpropagate the errors through temporal attention filters reaching the underlying per-frame convolutional layers. This makes end-to-end training of the proposed model possible with video training data. Per-frame CNNs were assumed to share all parameters.

### Recurrent neural networks with temporal filters

Although the model presented in the previous subsection allows us to learn temporal attention filters from the training data, it was assumed that the temporal filters are static once learned and are shared across all videos. However, such assumption that relative locations of sub-events are exactly identical across all activity videos can be dangerous. A particular sub-event of the activity (e.g., a person stretching an arm in the case of 'shake hands') may occur earlier or faster in one video than those in the other videos, due to human action style variations. In such cases, using static temporal filters will fail to capture exact sub-events. Rather, the recognition system must learn how to dynamically and adaptively adjust locations of temporal filters depending on the video content.

Thus, in this subsection, we propose an alternative approach of using a recurrent neural network, the long short-term memory (LSTM). Figure 4 describes our overall LSTM architecture. At each iteration, our LSTM takes the entire video frames as an input and applies per-frame CNNs identical to the previous subsection. Next, instead of using the learned static temporal filters, previous LSTM outputs are used to decided the temporal attention filter parameters in an adaptive fashion. Our approach learns weights that models how previous LSTM iteration outputs (i.e., the abstraction of video information in the previous round) can lead to the better temporal filters in the next iteration

More specifically, our attention filter parameters become the function of previous iteration LSTM outputs:

$$(g_t, \delta_t, \sigma_t) = W_n(h_{h-1}) = \sum_i w_i \cdot h_{t-1}(i) \quad (4)$$

where $W_n$ is the function we need to learn modeled as a weighted sum, $h_{t-1}$ is the LSTM hidden state vector at iteration $t-1$ and $h_{t-1}(i)$ are its elements. These weights are initialized such that the initial iteration places $g$ at the center of the video, and $\delta$ spans the entire duration, allowing the LSTM to get input from the entire sequence of frames. Because of the nature that our temporal attention filters are differentiable, we learn the function $W_n$ through the backpropagation.

Notice that our usage of LSTMs is different from previous approaches (Ng et al. 2015) using LSTMs to capture sequential per-frame feature changes. In our case, the goal of each LSTM iteration is to adjust the temporal filter locations to match the video input.

## Experiments

In order to evaluate the effectiveness of our proposed latent sub-event learning approach with fully differentiable temporal attention filters, we conducted a set of experiments comparing our approach using temporal filters against the previous conventional approaches without them while making the approaches use the same features and classifiers. These approaches were evaluated with multiple different features and multiple different datasets.

**Features:** We extracted VGG features (Simonyan and Zisserman 2015), INRIA's improved trajectory features (ITF) (Wang and Schmid 2013), C3D features (Tran et al. 2015), and TDD features (Wang, Qiao, and Tang 2015) which were used as inputs to our temporal filter model. VGG is an image-based convolutional neural network originally designed for object classification tasks. ITF and TDD are the state-of-the-art trajectory-based local video features, each taking advantage of HOG/HOF or CNN feature maps observed around of

trajectories. We used the source codes of all these features provided by the authors of the corresponding papers. For TDD, we used the single-scale version of the TDD feature, since it showed the better performance. We used 3-frame short video segments as our unit observations. In the case of VGG features, we applied its CNN architecture to one image frame in every 3 frames, obtaining 4K-dimensional vectors from the final fully connected layer. In the case of trajectory features, we considered the trajectories ending within the 3-frame segment as inputs corresponding to the segment and took advantage of their Fisher vector representation.

**Datasets:** We conducted experiments with two different public video dataset: DogCentric activity dataset (Iwashita et al. 2014) and HMDB dataset (Kuehne et al. 2011). The Dog-Centric dataset is a first-person video dataset, and it was chosen because that the previous pooled time series (PoT) representation (Ryoo, Rothrock, and Matthies 2015), which illustrated potential that considering multiple sub-intervals of videos benefit the CNN-based recognition, achieved the state-of-the-art performance on it. The first-person videos in this dataset display a significant amount of ego-motion (i.e., camera motion) of the camera wearer, and it is an extremely challenging dataset. HMDB was chosen due to its popularity. Our experiments were conducted by following each dataset's standard evaluation setting. Both the datasets are designed for multiclass activity video classification.

**Implementation and baseline classifiers:** As mentioned above, we confirm the advantage of our approach with four different types of underlying feature models: VGG, ITF, C3D, and TDD. Our temporal filters were applied on top of them. In all our experiments, we used 2 hidden layers on top of our (multiple) temporal filters: one of size 2048 and one of size 10 and softmax classification.

As the basic baselines, we tested (1) max-pooling, (2) sum-pooling, and (3) mean-pooling across the time dimension, resulting in features of size $1 \times D$. These $1 \times D$ features were fed to the hidden layers for the classification, which was a standard practice as described in (Jain, van Gemert, and Snoek 2014; Karpathy et al. 2014; Simonyan and Zisserman 2014; Ng et al. 2015; Tran et al. 2015). In addition, in order to explicitly confirm the power of our approach of 'learning' temporal attention filters with different location/duration/resolution, we also implemented the baselines of using fixed-predetermined temporal filters (i.e., without learning). The main idea was to make the systems take advantage of temporal filters identical to ours while disabling their learning capability. This makes the baselines behave very similar to previous pooled time series method (Ryoo, Rothrock, and Matthies 2015). We tested (4) an approach of using a single fixed filter (i.e., the level 1 temporal pyramid), which essentially computes a weighted sum of the whole video. We then used (5) a temporal pyramid of level 4, having 15 filters: 1 viewing the whole video, 2 viewing half the video, 4 viewing a forth, and 8 viewing an eighth, giving a vector of size $15 \times N \times D$.

We implemented our CNN-based recognition architecture with learnable temporal attention filters as described in this paper. First, we implemented our approach of learning static

temporal filters. We tested our model's ability to learn filters by using 15 filters (i.e., $M = 15$) with $N = 1$ or $N = 3$. Finally, we modified the model to use a LSTM to dynamically choose where to look. At each step, the temporal model took the hidden state, $h$ from the LSTM as input, and did a linear transformation with learned weights and bias $W \cdot h + b$ to obtain $g, \delta$ and $\sigma$. We then used the same equations as above to create the filters. The LSTM ran for 4 steps and either had 1 or 3 filters (i.e., $M = 1$ or 3) with $N = 5$.

**Training the network:** To increase training data, we apply random cropping on each video frame, and randomly skip several video frames at beginning. We use log-scale of stride and variance to ensure positivity as (Gregor et al. 2015). We initialize each filter bank parameters $(\widetilde{g}_m, \log \widetilde{\delta}_m, \log \sigma_m^2)$ with normal distribution for $\widetilde{g}_m$ and 0 for $\log \widetilde{\delta}_m$ and $\log \sigma_m^2$.

For all experiments, the first fully-connected layer had 4096 nodes and used a ReLU activation function. The second had either 10 nodes (for DogCentric) or 51 nodes (for HMDB) and used soft-max. The network was trained for 10000 iterations with a batch size of 100 and stochastic gradient descent with momentum set to 0.9.

## DogCentric dataset

The DogCentric dataset consists of 209 videos (102 training and 107 testing videos) and 10 classes. As mentioned above, it is a very challenging dataset with severe camera motion.

All of the baselines described above as well as 4 different versions of our approach (2 with static filter learning and 2 with LSTMs) were compared while using 3 different types of underlying features. Figure 5 shows the results. We are able to clearly observe that the consideration of multiple sub-intervals improves the recognition performances. The performance increased by using predetermined temporal pyramid, and our proposed approach of learning temporal filters were able to further improve the performance. Additionally, using the LSTM to dynamically choose sub-event locations gives the best performance.

The overall difference between the conventional approach of max/sum/mean pooling and our learned temporal attention filters are around 5% in VGG, 1% in ITF, and 4∼5% in TDD. We believe our approach was more effective with VGG since it's a pure CNN-based feature allowing our differentiable filters to better cope with them. ITF is a completely hand-crafted feature and TDD is a partially hand-crafted feature.

Table 1 compares the performances of our approach with the previously reported state-of-the-arts. The base features we used for this table was TDD. By taking advantage of temporal filter learning and also LSTMs to dynamically adjust the filters matching sub-events, we were able to outperform the state-of-the-arts.

## HMDB

HMDB is a relatively large-scale video dataset with 51 activity classes and more than 5000 videos.

In our first experiment, we made our model to learn general temporal filters to be *shared* by all 51 activity classes, testing whether our approach can learn globally optimal filters.
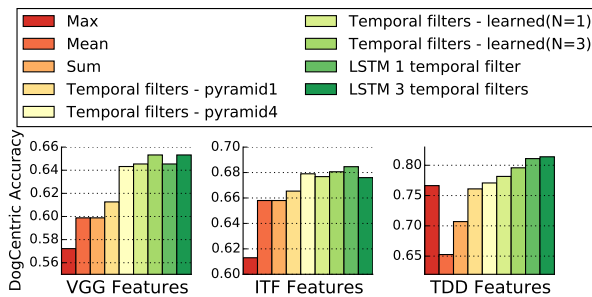
Figure 5: Classification accuracy of baselines and our proposed approach using VGG, ITF and VGG features.

Table 1: Recognition performances of our approach on the DogCentric dataset, compared against previously reported results of state-of-the-art approaches.

| Approach | Accuracy |
|---|---|
| VGG (Simonyan and Zisserman 2015) | 59.9 % |
| Iwashita et al. 2014 | 60.5 % |
| ITF (Wang and Schmid 2013) | 67.7 % |
| ITF + CNN (Jain, van Gemert, and Snoek 2014) | 69.2 % |
| PoT (Ryoo, Rothrock, and Matthies 2015) | 73.0 % |
| TDD (Wang, Qiao, and Tang 2015) | 76.6 % |
| Ours (temporal filters) | 79.6 % |
| Ours (temporal filters + LSTM) | **81.4 %** |

Table 2 shows the results of our approach with the temporal attention filters compared against the approaches without them. We tested this with CNN-based features: VGG, C3D, and TDD. We are able to observe that our approach of learning latent sub-events using temporal attention filters clearly benefits the classification in all cases.

Finally, instead of making our model to learn general filters, we made our models to learn 1-vs-all binary classifier for each activity. This enables each activity classification model to learn class-specific (latent) sub-events tailored for the activity, allowing the model to fully take advantage of our filter learning method. These 1-vs-all results were combined for the final 51-class classification. We were able to get an accuracy of 68.4% (Table 3). This is significant considering that the base feature performance (i.e., TDD) is 57% with max/mean pooling. The setting was $N = 2$ and $M = 3$. The performance increase gained by learning such latent sub-events per activity was significant and consistent: it was 9% increase over fixed temporal pyramid filters with C3D and 10% increase over temporal pyramid with TDD. There are few existing approaches performing comparable to our approach by fusing multiple features (e.g., 69.2 of (Feichtenhofer, Pinz, and Zisserman 2016) + ITF and 67.2 of (Varol, Laptev, and Schmid 2016) + ITF), but our approach using one feature performed superior to them in their original form. Figure 6 shows examples of the learned sub-events.

## Conclusion

We present a new activity recognition approach using *temporal attention filters*. Our approach enables end-to-end learning of the classifier jointly with latent sub-events and underlying CNN architectures. An approach to adaptively update the

Table 2: A table comparing the performance of the approaches using the HMDB Dataset. In this experiment, we made the model to learn general temporal filters *to be shared* across all 51 activities.

| Method | VGG | C3D | TDD |
|---|---|---|---|
| **Baseline** | | | |
| Max Pooling | 37.77 % | 48.45 % | 57.07 % |
| Sum Pooling | 37.00 % | 48.58 % | 55.77 % |
| Mean Pooling | 37.73 % | 49.30 % | 57.17 % |
| **Fixed Temporal Filters** | | | |
| Pyramid 4 | 41.56 % | 49.69 % | 58.87 % |
| **Learned Temporal Filters** | | | |
| $N = 1$ | 41.23 % | 50.35 % | 58.87 % |
| $N = 3$ | 42.50 % | 50.00 % | 59.03 % |
| LSTM 1 filter | 42.72 % | 51.20 % | 58.04 % |
| LSTM 3 filters | 43.03 % | 49.81 % | 58.93 % |

Table 3: Final recognition performances of our approach with *per-activity* temporal filters, tested with HMDB. Results are compared against the state-of-the-arts.

| Approach | Accuracy |
|---|---|
| ITF (Wang and Schmid 2013) | 57.2 % |
| 2-stream CNN (Simonyan and Zisserman 2014) | 59.4 % |
| TDD (Wang, Qiao, and Tang 2015) | 63.2 % |
| LTC (Varol, Laptev, and Schmid 2016) | 64.8 % |
| S+T (Feichtenhofer, Pinz, and Zisserman 2016) | 65.4 % |
| Max pooling - C3D | 48.5 % |
| Temporal pyramid - C3D | 49.7 % |
| Ours (temporal filters) - C3D | 57.7 % |
| Max pooling - TDD | 57.1 % |
| Temporal pyramid - TDD | 58.9 % |
| Ours (temporal filters) - TDD | **68.4** % |

temporal filter location/duration/resolution using multiple recurrent LSTM iterations was also proposed, and its potential was experimentally confirmed. The concept of learning latent sub-events using temporal attention clearly benefited the recognition, particularly compared to the baselines using the same features without such learning.

**Discussions:** Unfortunately, we were unable to replicate the TDD's reported recognition performance of 63.2 % with the code provided by the authors. We only obtained 57%. This probably is due to the difference in detailed parameter settings and engineering tricks. If we can replicate the performance of TDD reported in its paper, we would be able to further increase our method's performance using it as a base.

## Acknowledgement

## References

[Aggarwal and Ryoo 2011] Aggarwal, J. K., and Ryoo, M. S. 2011. Human activity analysis: A review. *ACM Computing Surveys* 43:16:1–16:43.
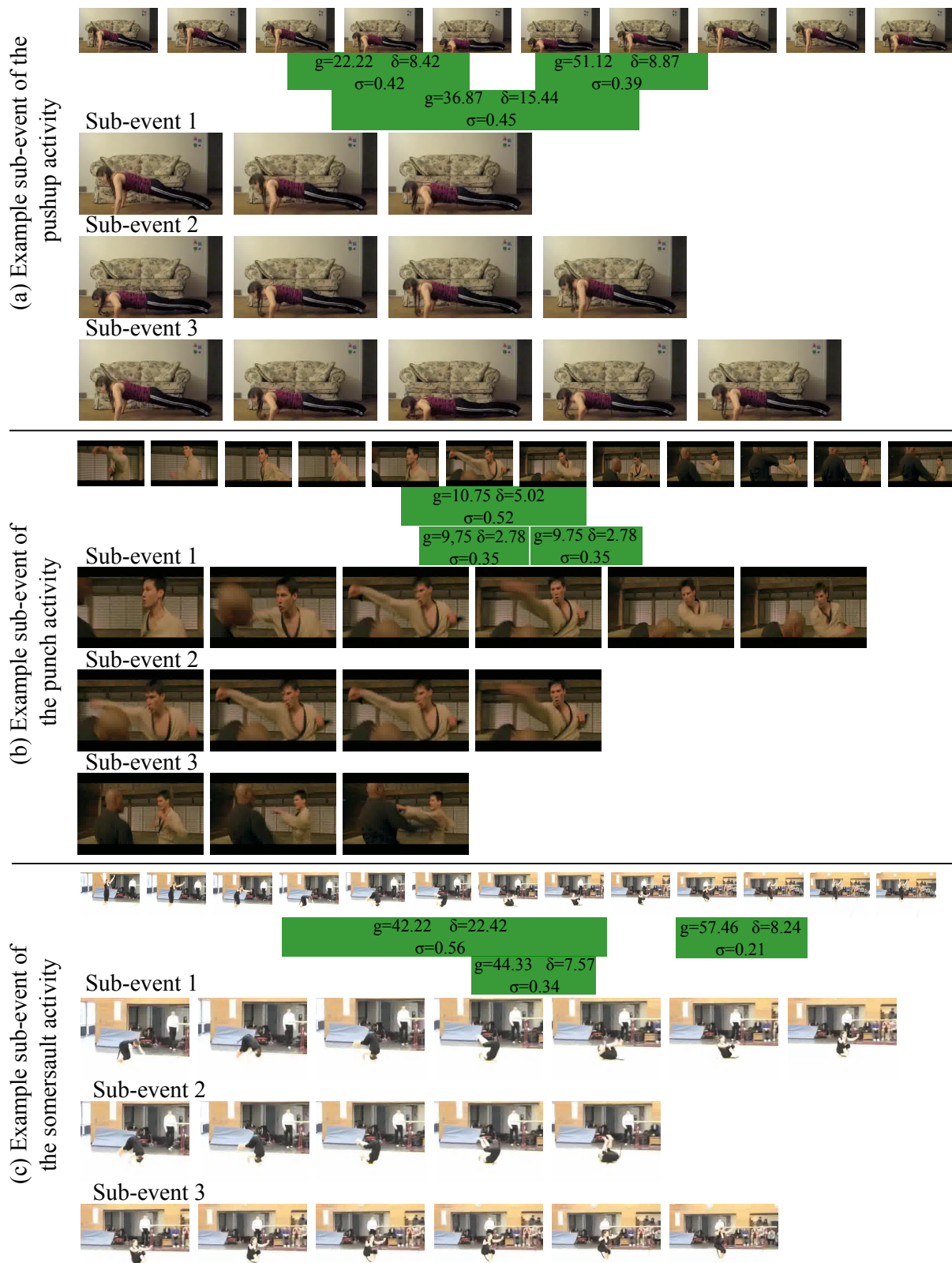
Figure 6: Example frames of the learned latent sub-events. The top row of each sub-figure (a-c) shows an example frame sequence, and the green boxes below shows the locations of the temporal filters learned to capture (latent) sub-events. Actual frames corresponding to each sub-event are also illustrated. We are able to observe that semantic sub-events such as (a) 'doing down' in the pushup activity, (b) 'stretch' in the punch activity and (c) 'roll' and 'stand up' in the somersault activity are captured.

[Feichtenhofer, Pinz, and Zisserman 2016] Feichtenhofer, C.; Pinz, A.; and Zisserman, A. 2016. Convolutional two-stream network fusion for video action recognition. In *CVPR*.

[Gregor et al. 2015] Gregor, K.; Danihelka, I.; Graves, A.; and Wierstra, D. 2015. DRAW: A recurrent neural network for image generation. *CoRR* abs/1502.04623.

[Iwashita et al. 2014] Iwashita, Y.; Takamine, A.; Kurazume, R.; and Ryoo, M. S. 2014. First-person animal activity recognition from egocentric videos. In *ICPR*.

[Jain, van Gemert, and Snoek 2014] Jain, M.; van Gemert, J.; and Snoek, C. G. M. 2014. University of amsterdam at THUMOS challenge 2014. In *ECCVW*.

[Karpathy et al. 2014] Karpathy, A.; Toderici, G.; Shetty, S.; Leung, T.; Sukthankar, R.; and Fei-Fei, L. 2014. Large-scale video classification with convolutional neural networks. In *CVPR*.

[Kuehne et al. 2011] Kuehne, H.; Jhuang, H.; Garrote, E.; Poggio, T.; and Serre, T. 2011. HMDB: a large video database for human motion recognition. In *ICCV*.

[Li et al. 2016] Li, Q.; Qiu, Z.; Yao, T.; Mei, T.; Rui, Y.; and Luo, J. 2016. Action recognition by learning deep multi-granular spatio-temporal video representation. In *ICMR*.

[Ng et al. 2015] Ng, J. Y.; Hausknecht, M. J.; Vijaya-narasimhan, S.; Vinyals, O.; Monga, R.; and Toderici, G. 2015. Beyond short snippets: Deep networks for video classification. In *CVPR*.

[Niebles, Chen, and Fei-Fei 2010] Niebles, J.; Chen, C.; and Fei-Fei, L. 2010. Modeling temporal structure of decomposable motion segments for activity classification. In *ECCV*.

[Ryoo and Matthies 2013] Ryoo, M. S., and Matthies, L. 2013. First-person activity recognition: What are they doing to me? In *CVPR*.

[Ryoo et al. 2015] Ryoo, M. S.; Fuchs, T. J.; Xia, L.; Aggarwal, J. K.; and Matthies, L. 2015. Robot-centric activity prediction from first-person videos: What will they do to me? In *HRI*.

[Ryoo, Rothrock, and Matthies 2015] Ryoo, M. S.; Rothrock, B.; and Matthies, L. 2015. Pooled motion features for first-person videos. In *CVPR*.

[Simonyan and Zisserman 2014] Simonyan, K., and Zisserman, A. 2014. Two-stream convolutional networks for action recognition in videos. In *NIPS*.

[Simonyan and Zisserman 2015] Simonyan, K., and Zisserman, A. 2015. Very deep convolutional networks for large-scale image recognition. In *ICLR*.

[Tran et al. 2015] Tran, D.; Bourdev, L.; Fergus, R.; Torresani, L.; and Paluri, M. 2015. Learning spatiotemporal features with 3D convolutional networks. In *ICCV*.

[Varol, Laptev, and Schmid 2016] Varol, G.; Laptev, I.; and Schmid, C. 2016. Long-term Temporal Convolutions for Action Recognition. *arXiv:1604.04494*.

[Wang and Schmid 2013] Wang, H., and Schmid, C. 2013. Action recognition with improved trajectories. In *ICCV*.

[Wang, Qiao, and Tang 2015] Wang, L.; Qiao, Y.; and Tang, X. 2015. Action recognition with trajectory-pooled deep-convolutional descriptors. In *CVPR*.

[Yeung et al. 2016] Yeung, S.; Russakovsky, O.; Mori, G.; and Fei-Fei, L. 2016. End-to-end learning of action detection from frame glimpses in videos. In *CVPR*.