

Recurrent Residual Learning for Action Recognition

Ahsan Iqbal, Alexander Richard, Hilde Kuehne, Juergen Gall
{iqbalm, richard, kuehne, gall}@iai.uni-bonn.de

University of Bonn, Germany

Abstract. Action recognition is a fundamental problem in computer vision with a lot of potential applications such as video surveillance, human computer interaction, and robot learning. Given pre-segmented videos, the task is to recognize actions happening within videos. Historically, hand crafted video features were used to address the task of action recognition. With the success of Deep ConvNets as an image analysis method, a lot of extensions of standard ConvNets were purposed to process variable length video data. In this work, we propose a novel recurrent ConvNet architecture called recurrent residual networks to address the task of action recognition. The approach extends ResNet, a state of the art model for image classification. While the original formulation of ResNet aims at learning spatial residuals in its layers, we extend the approach by introducing recurrent connections that allow to learn a spatio-temporal residual. In contrast to fully recurrent networks, our temporal connections only allow a limited range of preceding frames to contribute to the output for the current frame, enabling efficient training and inference as well as limiting the temporal context to a reasonable local range around each frame. On a large-scale action recognition dataset, we show that our model improves over both, the standard ResNet architecture and a ResNet extended by a fully recurrent layer.

1 Introduction

Action recognition in videos is an important research topic [1,4,18] with many potential applications such as video surveillance, human computer interaction, and robotics. Traditionally, action recognition has been addressed by hand crafted video features in combination with classifiers like SVMs as in [21,22]. With the impressive achievements of deep convolutional networks (ConvNets) for image classification, a lot of research was devoted to extend ConvNets to process video data, however, with unsatisfying results. While ConvNets have shown to perform very well for spatial data, they perform poorly for temporal data since they fail to model temporal dependencies. Heuristics were therefore developed for modeling temporal relations. First attempts, which simply stacked the frames and used a standard ConvNet for image classification [10], performed worse than hand crafted features. More successful have been two stream architectures [18] that use two ConvNets. While the first network is applied to the independent

frames, the second network processes the optical flow, which needs to be computed beforehand. While two stream architectures achieve lower classification error rates than hand-crafted features, they are very expensive for training and inference since they need two ConvNets and an additional approach to extract the optical flow.

In this work, we propose a more principled way to integrate temporal dependencies within a ConvNet. Our model is based on the state of the art residual learning framework [6] for image classification, which learns a residual function with respect to the layer’s input. We extend the approach to a sequence of images by having a residual network for each image and connecting them by recurrent connections that model temporal residuals. In contrast to the two stream architecture [4], which proposes residual connections from the motion to the appearance stream, our approach is a single stream architecture that directly models temporal relations within the spatial stream and does not require the additional computation of the optical flow.

We evaluate our approach on the popular UCF-101 [19] benchmark and show that our approach reduces the error of the baseline [6] by 17%. Although two stream architectures, which require the computation of the optical flow, achieve a lower error rate, the proposed approach of temporal residuals could also be integrated into a two stream architecture.

2 Related Work

Due to the difficulty of modeling temporal context with deep neural networks, traditional methods using hand-crafted features have been state of the art in action recognition much longer than in image classification [12, 21, 22, 24]. The most popular approaches are dense trajectories [21] with a bag-of-words and SVM classification as well as improved dense trajectories [22] with Fisher vector encoding. Due to the success of deep architectures, first attempts in action recognition aimed at combining those traditional features with deep models. In [15], for instance, a combination of hand crafted features and recurrent neural networks have been deployed. Peng et. al. [14] proposed Stacked Fisher Vectors, a video representation with multi-layer nested Fisher vector encoding. In the first layer, they densely sample large subvolumes from input videos, extract local features, and encode them using Fisher vectors. The second layer compresses the Fisher vectors of subvolumes obtained in the previous layer, and then encodes them again with Fisher vectors. Compared with standard Fisher vectors, stacked Fisher vectors allow to refine and abstract semantic information in a hierarchical way. Another hierarchical approach has been proposed in [8], who apply HMAX [16] with pre-defined spatio-temporal filters in the first layer. Trajectory pooled deep convolutional descriptors are defined in [23]. CNN features are extracted from a two stream architecture and are combined with improved dense trajectories.

In the past, there have been attempts to address the task of action recognition with deep architectures directly. However, in most of these works, the input to

the model is a stack of consecutive video frames and the model is expected to learn spatio-temporal dependent features in the first few layers, which is a difficult task. In [2, 13, 20], spatio temporal features are learned in unsupervised fashion by using Restricted Boltzmann machines. The approach of [7] combines the information about objects present in the video with the motion in the videos. 3D convolution is used in [9] to extract discriminative spatio temporal features from the stack of video frames. Three different approaches (early fusion, late fusion, and slow fusion) were evaluated to fuse temporal context in [10]. A similar technique as in [9] is used to fuse temporal context early in the network, in late fusion, individual features per frame are extracted and fused in the last convolutional layer. Slow fusion mixes late and early fusion. In contrast to these methods, our method does not rely on temporal convolution but on a recurrent network architecture directly.

More recently, [1] proposed concept of dynamic images. The dynamic image is based on the rank pooling concept [5] and is obtained through the parameters of a ranking machine that encodes the temporal evolution of the frames of the video. Dynamic images are obtained by directly applying rank pooling on the raw image pixels of a video producing a single RGB image per video. And finally, by feeding the dynamic image to any CNN architecture for image analysis, it can be used to classify actions.

The most successful approach to date is the two-stream CNN of [18], where individual frames from the videos are the input to the spatial network, while motion in the form of dense optical flow is the input to the temporal network. The features learned by both networks are concatenated and finally linear SVM is used for classification. Recently, with the success of ResNet [6], [4] proposed a model that combines ResNet and the two stream architecture. They replace both spatial and temporal networks in the two stream architecture by a ResNet with 50 layers. They also introduce a temporal or motion residual, i.e. a residual connection from the temporal network to the spatial network to enable learning of spatio temporal features. In contrast to our method, they incorporate temporal information by extending the convolutions over temporal windows. Note that this leads to a largely increased amount of model parameters, whereas our approach shares the weights among all frames, keeping the network size small. [26] proposed the temporal segment networks, which are mainly based on the two stream architecture. However, rather than densely sampling every other frame in the video, they divide the video in segments of equal length, and then randomly sample snippets from these segments as network input. In this way, the two stream network produces segment level classification scores, which are combined to produce video level output.

Deep recurrent CNN architectures are also explored to model dependencies across the frames. In [3], convolutional features are fed into an LSTM network to model temporal dependencies. [28] considered four networks to address action recognition in videos. The first network is similar to spatial network in the two stream architecture. The second network is a CNN with one recurrent layer, it expects a single optical flow image and in recurrent layer, optical flows over a

range of frames are combined. In the third network, they feed a stack of consecutive frames, the network is also equipped with a recurrent layer to capture the long term dependencies. Similarly, the fourth network expects a stack of optical flow fields as input. However, the network is equipped with a fully connected recurrent layer. Finally, boosting is used to combine the output of all four networks.

Finally, [27] equip a ResNet with recurrent skip connections that are, contrary to ours, purely temporal skip connections, whereas in our framework, we use spatio-temporal skip connections. Note the significant difference in both approaches: while purely temporal skip connections can be interpreted as usual recurrent connections with unit weights, spatio-temporal skip connections are a novel concept that allow for efficient backpropagation and combine both, changes in the temporal domain and changes in the spation domain at the same time.

3 Recurrent Residual Network

In this section, we describe our approach to address the problem of action recognition in videos. Our approach is an extension of ResNet [6], which reformulates a layer as learning the spatial residual function with respect to the layer’s input. State of the art results were achieved in image recognition tasks by learning spatial residual functions. We extend the approach to learn temporal residual functions across the frames to do action recognition in videos. In our formulation, the feature vector at time step t is a residual function with respect to the feature vector at time step $t - 1$. By following the analogy of ResNet, temporal residuals are learned by introducing the temporal skip (recurrent) connections. In the following, we give a brief introduction to ResNet, explain different types of temporal skip connections, and finally describe how to include more temporal context.

3.1 ResNet

ResNet [6] introduces a residual learning framework. In this framework, a stack of convolutional layers fit a residual mapping instead of the desired mapping. Let $H(x)$ denote the desired mapping. The principle of ResNet is to interpret the mapping of the learned function from one layer to another as $H(x) = F(x) + x$, i.e. as the original input x plus a residual function $F(x)$. Introducing the spatial skip connection, the input signal x is directly forwarded and added to the next layer, so it only remains to learn the residual $F(x) = H(x) - x$, see Figure 1b.

3.2 Type of Temporal Skip Connection

There are multiple possibilities to model the temporal skip connection. The standard spatial skip connections in the classical ResNet architecture are either an identity mapping, i.e. they just forward the input signal and add it to the destination layer, or they perform a linear transformation in order to establish

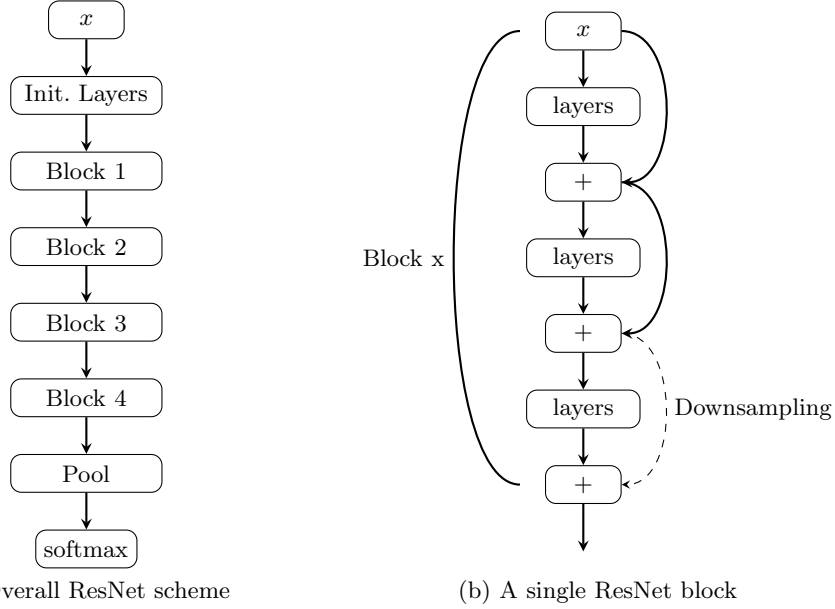


Fig. 1: ResNet architecture, (a) shows the overall ResNet structure with four building blocks and a final classification layer, (b) is the schema of a single block: each block consists of multiple convolutional layers and skip connections to learn the residuals. At the end, the output feature maps are downsampled.

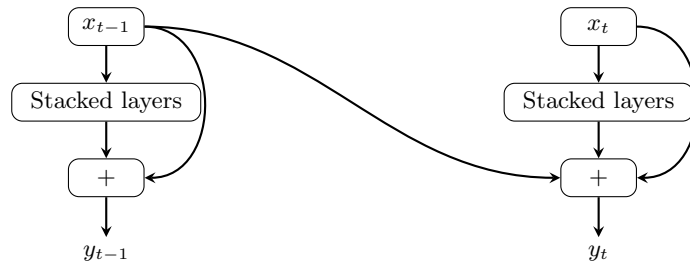


Fig. 2: Building block of our recurrent residual learning, x_{t-1} represents the input to a convolutional layer at time step $t - 1$ and x_t represents the input to the same layer at time step t . While the spatial skip connections within a single time frame allow to learn a spatial residual, the spatio-temporal skip connection from time $t - 1$ to time t adds temporal information to the learned residual.

the downsampling as depicted in Figure 1. The simplest case for the temporal skip connection is to also use an identity mapping. With the notation of Figure 2, the layer output y_t at time t is the residual function

$$y_t = \sigma(x_t * W) + x_t, \quad (1)$$

where σ represents the nonlinear operations performed after the linear transformation. Note that for simplicity of notation, we pretend that the residual block contains a single convolutional layer only and W represents weights for the layer. Extending this for the temporal skip connection, we obtain

$$y_t = \sigma(x_t * W) + x_t + x_{t-1}. \quad (2)$$

In order to allow for a weighting of the temporal skip connection with weights W_s , a linear transformation can be applied to x_{t-1} before adding it to y_t ,

$$y_t = \sigma(x_t * W) + x_t + x_{t-1} * W_s. \quad (3)$$

Moreover, in order to learn a nonlinear spatio-temporal mapping, this can be further extended to

$$y_t = \sigma(x_t * W) + x_t + \sigma(x_{t-1} * W_s). \quad (4)$$

3.3 Temporal Context

While recurrent connections in traditional recurrent neural networks feed the output of a layer at time $t - 1$ to the same layer as input at time t , our proposed spatio-temporal skip connections are different. For an illustration, see Figure 3. Here, we unfold a network with two spatio-temporal skip connections over time. Note that the temporal context that influences the output y_t includes x_{t-2} , x_{t-1} , and x_t as there are paths from y_t leading to all these inputs. If we only used one temporal skip connection instead of two, the accessible temporal context for y_t would only be x_{t-1} and x_t , respectively. In general, if a temporal context over T frames is desired, at least $T - 1$ temporal skip connections are necessary.

In order to use this approach for action recognition, a video is divided into M small sequences each containing T frames. A recurrent residual network with $T - 1$ temporal skip connections is created to capture the dependencies over these T time steps. In training, we optimize the cross-entropy loss of each small video chunk. During inference, for each small sequence $\{x_t^{(i)}\}_{t=1}^T$ within one video, the recurrent residual network computes $P(y = c | \{x_t^{(i)}\}_{t=1}^T)$. In order to obtain an overall classification of a complete video, the individual output probabilities are averaged over the M subsequences of the video. Note that this is similar to existing frame-wise approaches, where an output probability per frame is computed and the overall video action probabilities are obtained by accumulating all single frame probabilities. In our case, instead of frames, we use small subsequences of the original video.

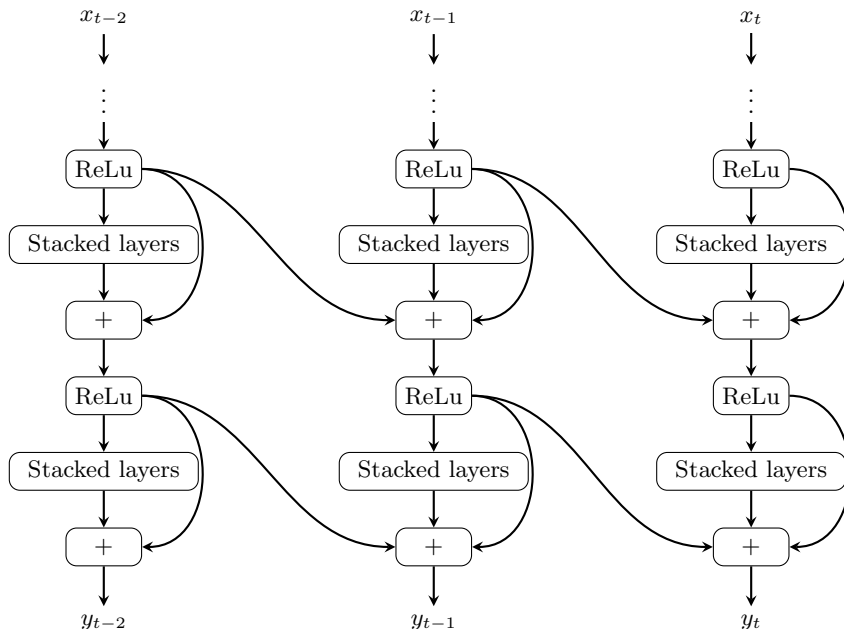


Fig. 3: A network with two temporal skip connections, capable of handling temporal context of three time steps, omitted layers are normal ResNet blocks, i.e. without any temporal skip connections.

4 Experimental Setup

In this section, we describe our experimental setup. We use ADAM [11] as learning algorithm and except for the baseline experiments, we update the model after observing 1% of training data, and every tenth frame from each video is sampled as input to the model. We evaluate our approach on UCF-101 [19], a large-scale action recognition dataset consisting of 13,000 videos from 101 different classes. The dataset comprises about 2.5 million frames in total. All the experimental work was done using our framework squirrel¹. In the following, we describe the baseline experiments and the experiments with our proposed recurrent residual network.

4.1 Baseline Experiments

As a baseline, we extract imagenet [17] features for individual frames in the video. Averaged individual feature vectors represent the feature vector of the complete video. Feature vectors for individual frames are extracted using a pre-trained ResNet model with 50 layers. A batch normalization layer is added after each layer to normalize the input to a layer. This way the network has in total 106 layers.

¹ <https://github.com/alexanderrichard/squirrel>

| Features | Method | Error Rate (with Z-Norm) | Error Rate (without Z-Norm) |
|----------|-----------|--------------------------|-----------------------------|
| Block4 | Avg. Pool | 0.236 | 0.237 |
| Block4 | GRU | 0.239 | 0.276 |
| Block3 | Avg. Pool | 0.309 | 0.313 |
| Block3 | GRU | 0.403 | 0.325 |
| Block2 | Avg. Pool | 0.431 | 0.434 |
| Block2 | GRU | 0.440 | 0.493 |

Table 1: Results of the baseline experiments.

We extract the imagenet feature vectors for each frame at three different positions of ResNet, i.e. after block4, block3, and block2 respectively, see Figure 1. We performed two sets of experiments on extracted features for each block. In one set, we average the frame level feature vectors, after Z-normalization and without Z-normalization, and train a linear classifier. We call this model the average pooling model. Similarly, in the other set, we use a recurrent neural network with 128 gated recurrent units (GRUs) in order to evaluate the performance of a classical recurrent network. We call this model the GRU.

Table 1 shows the baseline experiments with imagenet features. The average pooling model outperforms the model with gated recurrent units. Also, it is evident from the experiments that with more depth, features become richer. Hence, the depth plays a significant role in getting good classification accuracy.

4.2 Effect of type and position of the recurrent connection

In this set of experiments, we evaluate different types of temporal skip connections along with their position in our proposed model. We evaluate temporal skip connections at four different positions, i.e. at the beginning by making the first skip connection in block1 recurrent (referred to as Block1), in the middle by making last skip connection in block2 recurrent (referred to as Block2), by making last convolutional skip connection in block4 recurrent (referred to as Mid Block4), and finally by making last skip connection in block 4 recurrent (referred to as Block4). Also, we evaluate the type of recurrent connections. In these experiments, we evaluate identity mapping temporal skip connections, and temporal skip connections with convolutional weights having kernels of size 1×1 . Table 2 shows the deeper we place the temporal skip connection in the network, the better is the classification accuracy.

In another set of experiments, we evaluate the effect of the type of temporal skip connection. We change the configuration of the best working setup, i.e. the one with the skip connection in block4. The connection performs a parametrized linear or non linear transformation and identity mapping.

Table 3 shows the results achieved by different type of connections, placed closer to the output layer as our previous analysis shows that works best. Identity mapping connection with non trainable weights performed best, possibly because introducing more weights in the network causes overfitting.

| Position | Type | Error Rate |
|------------|---------------|--------------|
| Block1 | Convolutional | 0.265 |
| Block2 | Convolutional | 0.234 |
| Mid Block4 | Convolutional | 0.231 |
| Block4 | Convolutional | 0.219 |

Table 2: Placing the recurrent connection at different positions in the network.

| Type | Error Rate |
|------------------|--------------|
| Identity Mapping | 0.197 |
| Conv. Linear | 0.219 |
| Conv. Non-Linear | 0.210 |

Table 3: Results achieved by different type of recurrent connections.

4.3 Effect of Temporal Context

In this set of experiments, we explore the effect of temporal context. As discussed earlier, with more recurrent connections, the network is able to include additional temporal dependencies. We already investigated the network with one recurrent connection that is able to include temporal context of two frames. In these experiments, we further explore the temporal context of three frames (by introducing two temporal connections in the network), and the temporal context of five frames (by introducing four temporal skip connections in the network). Figure 3 shows the network architecture to accommodate temporal context of three frames.

As it is evident in Table 4, we do not gain much by including more temporal context. The accuracy improves in case of temporal context three, however it gets worse in case of temporal context five. Hence, considering training time, we consider the model with only one temporal skip connection as the best model. Note that due to the fact that we sample every tenth frame from the video, the overall temporal range is actually ten frames. More precisely, the network learns spatio-temporal residuals between the frames x_t and x_{t-10} , covering a reasonable amount of local temporal progress within the video.

We further evaluate our best model on all three splits of UCF-101. On average our best model achieves **0.198** on UCF-101 [19], which is a relative improvement of 17% over the ResNet baseline which has an error rate of 0.236.

4.4 Comparison with the state of the art

In this section, our best model (with one temporal skip connection and with sample rate 10) is compared with state of the art action recognition methods. As motion in the frames and appearance in individual frames are two complementary aspects for action recognition, most of the state of the art methods consider two different neural networks, an appearance stream and a motion stream, to extract and use appearance and motion for action recognition. The output of both the networks is fused, and a simple classifier is trained to classify videos. As our model uses the raw video frames only rather than optical flow fields,

| Temporal Context | Model | Error Rate |
|------------------|-------------------------|--------------|
| 1 | baseline | 0.236 |
| 2 | 1 recurrent connection | 0.197 |
| 3 | 2 recurrent connections | 0.194 |
| 5 | 4 recurrent connections | 0.209 |

Table 4: Results achieved by including more temporal context. For the best setup (context two), the error is reduced by 17% from 0.236 to 0.194.

| Method | Appearance | Motion | App.+Motion |
|------------------------------------------|--------------|--------|-------------|
| Improved Dense Trajectories [22] | - | - | 0.141 |
| Dynamic Image Networks [1] | 0.231 | - | - |
| Two Stream Architecture [18] | 0.270 | 0.163 | 0.120 |
| Two Stream Architecture (GoogleNet) [25] | 0.247 | 0.142 | 0.107 |
| Two Stream Architecture (VGG-Net) [25] | 0.216 | 0.130 | 0.086 |
| Spatiotemporal ResNets [4] | - | - | 0.066 |
| Recurrent Residual Networks | 0.198 | - | - |

Table 5: Classification error rates for UCF-101.

fair comparison of our model and the state of the art is only possible for the appearance stream. For completeness, we also compare our results with results achieved after the outputs of the appearance and the motion streams are fused, see Table 5. Our model achieves better error rates than the state of the art appearance stream models. Only fused models perform better. Note that the works [4, 18, 25] are all two-stream architectures. The dynamic image network [1] is a purely appearance base method that reduces that video to a single frame and uses a ConvNet to classify this frame. For a fair comparison, we provide the result of dynamic image network without the combination with dense trajectories as this would include motion features.

5 Conclusion

We extended the ResNet architecture to include temporal skip connections in order to model both, spatial and temporal information in video. Our model performs well already with a single temporal skip connection, enabling to infer context between two frames. Moreover, we showed that fusing temporal information at a late stage in the network is beneficial and that learning a temporal residual is superior to using a classical recurrent layer. Our method is not limited to appearance based models and can easily be extended to motion networks that have shown to further enhance the performance on action recognition datasets. A comparison to both, a ResNet baseline and state of the art methods showed that our approach outperforms other purely appearance based approaches.

Acknowledgement: The authors have been financially supported by the DFG projects KU 3396/2-1 and GA 1927/4-1 and the ERC Starting Grant ARCA (677650). Further, this work was supported by the AWS Cloud Credits for Research program.

References

1. Bilen, H., Fernando, B., Gavves, E., Vedaldi, A., Gould, S.: Dynamic image networks for action recognition. In: *IEEE Conference on Computer Vision and Pattern Recognition* (2016)
2. Chen, B., Ting, J.A., Marlin, B., de Freitas, N.: Deep learning of invariant spatio-temporal features from video. In: *NIPS 2010 Deep Learning and Unsupervised Feature Learning Workshop* (2010)
3. Donahue, J., Anne Hendricks, L., Guadarrama, S., Rohrbach, M., Venugopalan, S., Saenko, K., Darrell, T.: Long-term recurrent convolutional networks for visual recognition and description. In: *IEEE Conference on Computer Vision and Pattern Recognition* (2015)
4. Feichtenhofer, C., Pinz, A., Wildes, R.: Spatiotemporal residual networks for video action recognition. In: *Advances in Neural Information Processing Systems* 29, pp. 3468–3476 (2016)
5. Fernando, B., Gavves, E., M., J.O., Ghodrati, A., Tuytelaars, T.: Rank pooling for action recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 39(4), 773–787 (2017)
6. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: *IEEE Conference on Computer Vision and Pattern Recognition* (2016)
7. Jain, M., van Gemert, J.C., Snoek, C.G.M.: What do 15,000 object categories tell us about classifying and localizing actions? In: *IEEE Conference on Computer Vision and Pattern Recognition* (2015)
8. Jhuang, H., Serre, T., Wolf, L., Poggio, T.: A biologically inspired system for action recognition. In: *IEEE International Conference on Computer Vision* (2007)
9. Ji, S., Xu, W., Yang, M., Yu, K.: 3D convolutional neural networks for human action recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 35(1), 221–231 (2013)
10. Karpathy, A., Toderici, G., Shetty, S., Leung, T., Sukthankar, R., Fei-Fei, L.: Large-scale video classification with convolutional neural networks. In: *IEEE Conference on Computer Vision and Pattern Recognition*. pp. 1725–1732 (2014)
11. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: *International Conference on Learning Representations* (2015)
12. Laptev, I.: On space-time interest points. *International Journal of Computer Vision* 64, 107–123 (2005)
13. Le, Q.V., Zou, W.Y., Yeung, S.Y., Ng, A.Y.: Learning hierarchical invariant spatio-temporal features for action recognition with independent subspace analysis. In: *IEEE Conference on Computer Vision and Pattern Recognition*. pp. 3361–3368 (2011)
14. Peng, X., Zou, C., Qiao, Y., Peng, Q.: Action recognition with stacked fisher vectors. In: *European Conference on Computer Vision*. pp. 581–595 (2014)
15. Richard, A., Gall, J.: A bag-of-words equivalent recurrent neural network for action recognition. *Computer Vision and Image Understanding* 156, 79–91 (2017)
16. Riesenhuber, M., Poggio, T.: Hierarchical models of object recognition in cortex. *Nature neuroscience* 2(11), 1019–1025 (1999)
17. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A.C., Fei-Fei, L.: ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision* 115(3), 211–252 (2015)

18. Simonyan, K., Zisserman, A.: Two-stream convolutional networks for action recognition in videos. In: *Advances in Neural Information Processing Systems 27*, pp. 568–576 (2014)
19. Soomro, K., Zamir, A.R., Shah, M.: UCF101: A dataset of 101 human actions classes from videos in the wild. CoRR abs/1212.0402 (2012)
20. Taylor, G.W., Fergus, R., LeCun, Y., Bregler, C.: Convolutional learning of spatio-temporal features. In: *European Conference on Computer Vision*. pp. 140–153 (2010)
21. Wang, H., Klaser, A., Schmid, C., Liu, C.L.: Action recognition by dense trajectories. In: *IEEE Conference on Computer Vision and Pattern Recognition*. pp. 3169–3176 (2011)
22. Wang, H., Schmid, C.: Action recognition with improved trajectories. In: *IEEE International Conference on Computer Vision*. pp. 3551–3558 (2013)
23. Wang, L., Qiao, Y., Tang, X.: Action recognition with trajectory-pooled deep-convolutional descriptors. In: *IEEE Conference on Computer Vision and Pattern Recognition*. pp. 4305–4314 (2015)
24. Wang, L., Qiao, Y., Tang, X.: Mofap: A multi-level representation for action recognition. *International Journal of Computer Vision* 119(3), 254–271 (2016)
25. Wang, L., Xiong, Y., Wang, Z., Qiao, Y.: Towards good practices for very deep two-stream convnets. CoRR abs/1507.02159 (2015)
26. Wang, L., Xiong, Y., Wang, Z., Qiao, Y., Lin, D., Tang, X., Van Gool, L.: Temporal segment networks: Towards good practices for deep action recognition. In: *European Conference on Computer Vision* (2016)
27. Wang, Y., Tian, F.: Recurrent residual learning for sequence classification. In: *Conference on Empirical Methods on Natural Language Processing*. pp. 938–943 (2016)
28. Yang, X., Molchanov, P., Kautz, J.: Multilayer and multimodal fusion of deep neural networks for video classification. In: *Proceedings of the 2016 ACM on Multimedia Conference*. pp. 978–987 (2016)