

MULTI-STREAM SINGLE SHOT SPATIAL-TEMPORAL ACTION DETECTION

Pengfei Zhang, Yu Cao, Benyuan Liu

Department of Computer Science, University of Massachusetts Lowell, USA

ABSTRACT

We present a 3D Convolutional Neural Networks (CNNs) based single shot detector for spatial-temporal action detection tasks. Our model includes: (i) two short-term appearance and motion streams, with single RGB and optical flow image input separately, in order to capture the spatial and temporal information for the current frame; (ii) two long-term 3D ConvNet based stream, working on sequences of continuous RGB and optical flow images to capture the context from past frames. Our model achieves strong performance for action detection in video and can be easily integrated into any current two-stream action detection methods. We report a frame-mAP of 71.30% on the challenging UCF101-24 [1] actions dataset, achieving the state-of-the-art result of the one-stage methods. To the best of our knowledge, our work is the first system that combined 3D CNN and SSD in action detection tasks.

Index Terms— Action Detection, spatial-temporal action localization, 3D convolutional neural networks, SSD

1. INTRODUCTION

The objective of action detection is to recognize and localize all the human action instances in a given video across both space and time. It is a fundamental task for video understanding and important for practical applications such as video surveillance and human-robot interaction. Action detection is a challenging problem due to two main difficulties: (i) it is hard to capture visual representations in the large spatial-temporal search space; (ii) it is difficult to understand the video fast and accurately, while the detection speed is essential for many application scenarios such as fall and violence detection.

To investigate the spatial-temporal video representation, researchers leverage the hand-crafted features such as dense trajectory [2, 3] and optical flow [2, 4] to build a two-stream network [5] to combine the spatial-temporal information together. However, most of these approaches overlook a fundamental issue in action detection, namely, the specific representation of spatial-temporal information for various actions. Many of them only use optical flow, an estimation of motion for each pixel between two images, as the source of temporal information. Conventional approaches estimate optical

flows between adjacent frames, which could only represent temporal information of short time periods but lack of long-term information that is important for human action recognition as well. With the successes that 2D CNN has achieved in the field of visual representation on spatial domain, it is natural to extend it to 3D to capture both spatial and temporal information. In action recognition tasks, even though the 3D CNN (I3D [6]) achieved the best result so far, the improvement brought by 3D CNN, compared to the hand-crafted features based 2D CNN approaches [7], has not reached its full potential. In this work, we revisit the role of optical flow and 3D convolution in temporal reasoning for action detection. To explore the contribution of the 4 streams: 2D RGB, 2D optical flow (OF), 3D RGB and 3D OF in action detection, we propose a multi-stream architecture and examine the performance of different stream combinations for various types of actions. We demonstrate that, for the single-stream framework, 3D CNN based model outperforms 2D CNN based model for RGB and optical flow respectively. However, in a two-stream framework, there are different winners of appearance and motion stream for various actions due to the large intra-class variability. As a result, the best frame level mean average precision (mAP) is achieved by the fusion of all four streams, which adapts to a variety of actions.

As for the second challenge on detection speed, although many conventional action detection methods [4, 8, 9, 10, 11] achieved good results, their two-stage architecture performs region proposal and classification in two steps. While the accuracy is improved, it significantly slows down the detection speed, making it unacceptable for realistic scenarios. To accelerate the detection speed, inspired by [12, 13], our model adopts the one-stage method, Single Shot MultiBox Detector (SSD) [14], as the detection framework. It merges the two stages into a single network, carries out the localization and classification simultaneously, and thus accelerates the entire process.

The key contributions of this paper include: (i) we leverage the single stage object detection architecture SSD to build a time efficient action detector; (ii) we explore different combinations of 2D and 3D streams for the detection task for a variety of action videos; (iii) experiment results show that our model outperforms previous one-stage action detection methods on the challenging untrimmed sports video dataset UCF101-24.

2. RELATED WORK

Our research builds on previous works in two fields:

Spatial-temporal action localization. Gkioxari and Malik [9] applied a two-stream R-CNN based framework to produce frame level detections, and then linked the result to tubes with a dynamic programming method. Weinzaepfel *et al.* [4] extracted EdgeBoxes as the action proposals and then used a tracking-by-detection method instead of the linking method. Both Saha *et al.* [11] and Peng *et al.* [8] leveraged two-stream Faster R-CNN to do action detection. Singh *et al.* [12] applied a single stage detection method SSD to perform online detection. Kalogeiton *et al.* [13] extend SSD’s anchor boxes to anchor cuboids to perform the temporal-spatial proposal.

3D CNN. Ji *et al.* [15] and Tran *et al.* [16] extended the 2D convolutional kernel to 3 dimensions, and much subsequent studies such as I3D and P3D [7] has gained lots of successes in video related tasks. The most recent state-of-the-art result is achieved by Gu *et al.* [10] based on I3D and faster-RCNN. Hou *et al.* [17] designed a C3D version of one-stage action detection method, however, it is an offline algorithm could not do frame level incremental detection.

3. MODEL DESCRIPTION

Multi-stream model. The architecture of our model is illustrated in Fig.1. Our model consists of 4 streams: 2D and 3D RGB streams, 2D and 3D optical flow streams. The conventional 2D RGB and optical flow streams are employed to capture the short-term spatial-temporal features, meanwhile, 3D streams are added to learn long-term features. The two 2D streams share the same architecture, but are trained individually and have their own parameters. The same applies to the 3D streams.

For the target action instances at time t , the 2D RGB stream’s input is current frame f_t , while the input of the 2D optical flow stream is extracted from the pair of $\{f_{t-1}, f_t\}$ using Brox *et al.*s [18] method. The input dimension of both 2D streams is $C \times H \times W$, where C , H and W denote the number of channels, height and width of the input frame, respectively. To perform spatial-temporal reasoning with 3D CNN, 3D RGB stream’s input is a sequence of continuous N frames $\{f_{t-N+1} \dots f_t\}$. Similarly, 3D optical flow stream’s input is N frames extracted from RGB frame pairs from $\{f_{t-N}, f_{t-N+1}\}$ to $\{f_{t-1}, f_t\}$. The input dimension is $C \times N \times H \times W$. We set $N = 8$ frames in the experiment.

2D SSD network. Each of the 2D networks consists of 3 main parts: backbone network, extra convolutional layers and detection heads. The backbone network is truncated VGG-16 and its last two fully connected layers fc6 and fc7 are converted to convolutional layers. Eight extra layers are added to the end of the backbone network to predict default bounding boxes’ offsets and their confidences for actions. Each of the selected layers has a different spatial output dimension, that

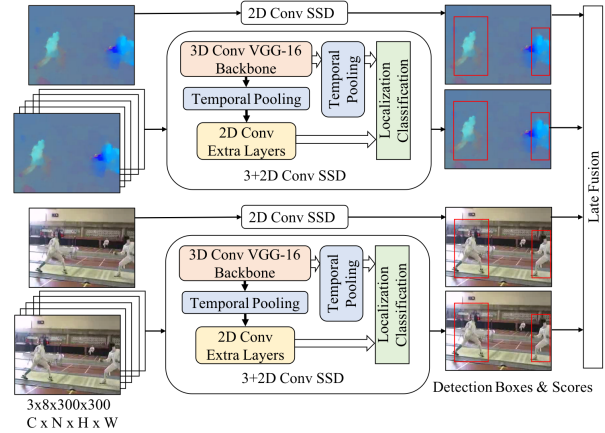


Fig. 1: Illustration of the proposed two-stream architecture. 3D SSD takes consecutive video frames as input and extracts both spatial and temporal information.

represents the action instance in different scales. The final predictions are produced by two detection heads: localization head and classification head, synchronously. We use a VGG-16 model pretrained on ImageNet to initialize our model and fine-tune it on the action dataset.

3D SSD network. As for the 3D streams, keeping the extra layers and detection heads unchanged, we inflate all the convolutional and pooling layers in backbone network from 2D to 3D, then apply temporal pooling to bridge the gap between 3D and 2D networks. To initialize the network, we repeat the weights of pretrained model’s 2D kernels T times, where T represents the size of the inflated kernel in temporal dimension. In our model, we convert all 3×3 kernels to $3 \times 3 \times 3$ kernels, set all layers’ temporal padding as 1 and temporal stride for pooling layers as 2.

Temporal pooling. We connect 3D and 2D layers by the temporal pooling layer. This layer performs mean-pooling along the temporal dimension, transforming the input feature map with dimension $C \times N \times H \times W$ to the output with dimension of $C \times H \times W$.

Fusion Method. We adopt late-fusion [5, 2, 19] to merge the spatial and temporal information from each stream together. In this step, we first choose one stream as the appearance stream, such as the 2D or 3D RGB stream, keep its bounding boxes regression result, and then set each box’s confidence score as the average score of the corresponding boxes from all fused streams. In the rest of paper, we will denote $A + M_1 + M_2 + \dots M_n$ as the late fusion of appearance stream A and motion streams from M_1 to M_n , n denotes the number of motion streams.

4. EXPERIMENTS

To evaluate the performance of 3D SSD stream, we examine different stream combinations and their detection accu-

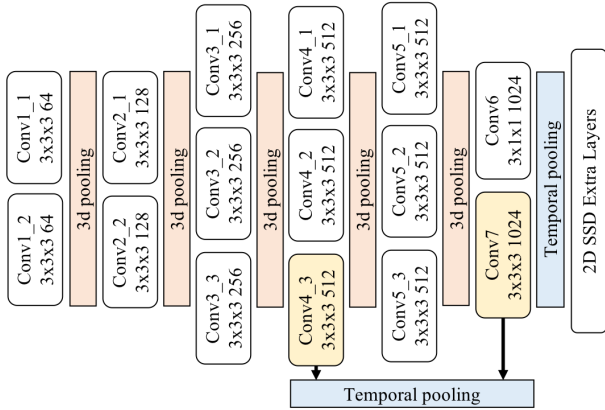


Fig. 2: Details of the inflated 3D backbone.

racy on the UCF101-24 dataset. Singh’s 2D SSD real-time framework [12] is used as a baseline. We keep their fusion and linking methods unchanged and focus on the performance improvement resulted from 3D SSD.

4.1. Settings

Datasets. We choose the first split of UCF101-24 dataset to evaluate our model. It contains 24 sport classes in 3,207 untrimmed videos. Each video is annotated with bounding boxes for each action instance at frame level and each frame may contain multiple actors.

Evaluation metrics. We evaluate the detection accuracy by mean average precision (mAP) for both frame and video levels. At frame level, if the Intersection-over-Union (IoU) between a predicted bounding box and ground truth is greater than a threshold α and this box’s action category is classified correctly, we will mark it as an correct detection. As for video level metric, after we connect the frame level detection into tubes, we can evaluate it with the spatial-temporal overlap between the predicted and annotated tubes. As in [9, 4, 11], we present the performance of our model in Table 1, 2, and 3, for frame-mAP with IoU threshold 0.5 and video-mAP with multiple IoU thresholds, $\alpha = 0.2$, $\alpha = 0.5$, $\alpha = 0.5 : 0.95$.

4.2. Performance

We will first analyse the performance of single streams, then further discuss the contribution of each stream in an ablation study of different two-stream combinations, and show how to get the best result for various of data at last.

Single-Stream. We report the comparison of 2D and 3D streams for RGB and optical flow in Table 1. The 2D streams adopt the same architecture and experiment setup as in Singh et.al [12]. For both frame and video level, each of our 3D streams outperforms the corresponding 2D streams, especially at video level, our 3D RGB network improves the mAP by 5.53% and 4.79% for IoU threshold $\alpha = 0.2$ and $\alpha = 0.5$,

Method	video-mAP			f.-mAP
	IoU	0.2	0.5	0.5
2D RGB [12]	69.8	40.9	18.7	64.96
2D OF [12]	63.7	30.8	11.0	47.26
ours-3D RGB	75.33	45.69	19.15	65.10
ours-3D OF	67.46	35.26	12.51	50.85

Table 1: Comparison of video and frame mAP between 2D and 3D RGB and Optical Flow (OF) streams.

Method	video-mAP			f.-mAP
	IoU	0.2	0.5	0.5
2D RGB+2D OF (b)[12]	73.0	44.0	19.2	68.31
2D RGB+2D OF (u)[12]	73.5	46.3	20.4	64.97
2D RGB+2D OF (l)[12]	76.43	45.18	20.08	67.81
ours-3D RGB+2D OF(l)	76.02	47.38	19.35	67.06
ours-2D RGB+3D RGB(l)	76.18	46.52	20.94	68.72
ours-3D RGB+3D OF(l)	76.84	46.38	19.2	68.82
ours-2D RGB+3D OF(l)	77.19	47.75	21.11	69.47

Table 2: Comparison between different combinations of two-stream fusion. (b) boost fusion, (l) late fusion, (u) union fusion.

respectively. Similarly, our 3D OF network improves 3.76% and 4.46%. The result indicates that the temporal information brought in by 3D convolution significantly improves the single-stream model’s performance.

Two-Stream. In this section, we answer the following two questions: (i) which stream is the best appearance stream? (ii) which stream is the best motion stream?

As for appearance stream, because 3D RGB stream contains both spatial and temporal information, our model can choose either 2D or 3D RGB stream as appearance stream. As shown in Table 2, the result of 2D RGB + 2D OF outperform that of 3D RGB + 2D OF by 0.75% with the same fusion method. Meanwhile, when the motion stream is 3D OF, the combination with 2D RGB appearance stream outperforms that of 3D RGB by 0.75%. This can be explained as the 2D RGB stream contains more accurate spatial information for current frame, while the 3D convolution brings in certain noises from the previous frames.

The candidates for motion stream are: 3D RGB, 2D and 3D optical flows. Comparing the 3D RGB stream with 2D optical flow stream, we find that 2D RGB + 3D RGB performs better than 2D RGB + 2D OF in frame-mAP and video-mAP for IoU threshold $\alpha = 0.2$ and $0.5 : 0.95$. The more detailed frame level average precision analysis for each of the 24 action classes is demonstrated in Fig.3. Based on the way how actors and background change with respect to camera, the videos of UCF101-24 can be divided into 3 categories: (i) **active background videos**: videos where the camera moves

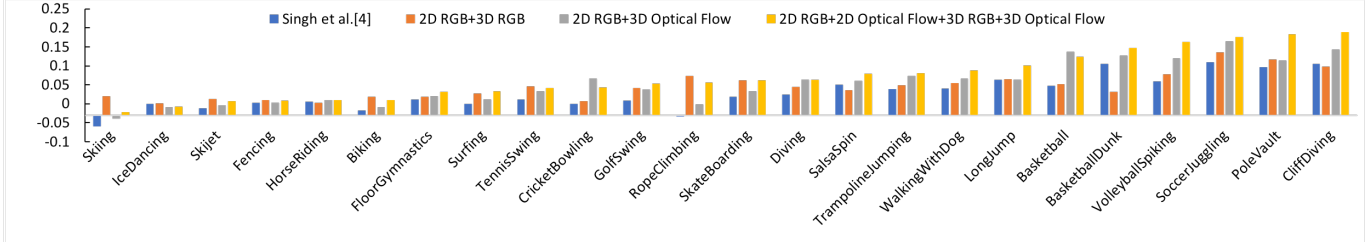


Fig. 3: UCF101-24 frame average precision for each action class compared to 2D RGB with baseline, the value of each class is compute with $AP_{multi-stream} - AP_{2DRGB}$.

along with the actors, meanwhile, the background environment changes sharply, for example, rope climbing, skiing and skateboarding. For these 3 classes of videos, the 2D RGB + 3D RGB combination outperform the 2D RGB + 2D OF combination by 10.75%, 7.89% and 4.35%, respectively. The poor performance of 2D optical flow stream is caused by the noises produced by the fast changing background. (ii) **fixed background videos:** videos where the camera is fixed, the background does not change much and the actors move quickly in short time frame, such as Salsa Spin, Cliff Diving and Basketball Dunk. Because optical flow contains more accurate short-term temporal information than 3D RGB, the performance of 2D RGB + 2D OF is better than that of 2D RGB + 3D RGB. (iii) For other videos that contain more complex circumstance, 3D RGB stream’s contribution is similar to or slightly better than optical flow.

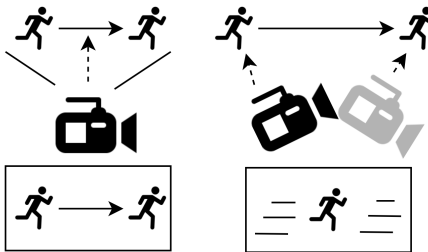


Fig. 4: Fixed background and Active background videos.

While 2D OF and 3D RGB outperforms each other in different scenarios, the best performance of all two-stream combinations is achieved by 2D RGB + 3D OF. It improve the frame-mAP of 2D RGB + 2D OF by 1.66%, and the result of 2D RGB + 3D RGB by 0.75%, which means 3D optical flow is the best choice of motion stream in two-stream framework. However, as shown in Fig.3, we can observe that the 3D optical flow still inherit the drawback of 2D optical flow, resulting in poor performance for active background videos.

Multi-Stream. We present the fusion results of three-stream and four-stream models in Table 3. Compared to Singh *et al.* [12]’s two-stream model, our three-stream model (2D RGB + 3D RGB + 2D OF) obtains 2.23% improvement for the frame-mAP with the 3D RGB stream integrated, and

Method	frame-mAP@0.5
(SSD) Kalogeiton <i>et al.</i> [13]	67.10
Hou <i>et al.</i> [20]	67.3
(SSD) Singh <i>et al.</i> [12]	67.81
ours-2D RGB+3D RGB+2D OF	70.04
ours-2D RGB+3D RGB+3D OF	71.10
ours-3D RGB+3D OF+2D RGB+2D OF	71.28
ours-2D RGB+2D OF+3D RGB+3D OF	71.30

Table 3: Comparison of frame-mAP to the state-of-the-art on UCF101-24 dataset in split1.

3.49% improvement with the fusion of all four streams. To the best of our knowledge, our model outperforms all the one-stage methods with better action localization and classification accuracy. In practice, we also need to consider the time consumption to prepare a stack of optical flows, which is important for developing an online real-time system. For different kinds of action videos and applications, our model is flexible to be reorganized or integrated into other models to meet the requirements.

5. CONCLUSIONS AND FUTURE PLANS

This paper introduced a multi-stream action detector which achieves state-of-the-art results of the one-stage methods on UCF101-24 dataset. We present an empirical study of the properties of the combinations of 2D RGB, 2D OF, 3D RGB and 3D OF streams. Based on the results of those experiments, the following conclusions could be obtained: (i) 2D RGB stream is a better choice for appearance stream comparing to other streams; (ii) for active background videos, 3D RGB motion stream can tolerate more environmental noises; (iii) optical flow, especially the 3D stream, performs well for videos that have fixed background and significant short-term action instances. Future work will be devoted to two directions: (i) optimize the framework with other one-stage methods, such as YOLO [21] series. (ii) Improve the temporal convolutional module with more lightweight 3D kernels to accelerate the whole forward process.

6. REFERENCES

- [1] K. Soomro, A. R. Zamir, and M. Shah, “Ucf101: A dataset of 101 human actions classes from videos in the wild,” *CoRR*, vol. abs/1212.0402, 2012.
- [2] K. Simonyan and A. Zisserman, “Two-stream convolutional networks for action recognition in videos,” in *Advances in Neural Information Processing Systems 27*, Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, Eds., pp. 568–576. Curran Associates, Inc., 2014.
- [3] H. Wang and C. Schmid, “Action recognition with improved trajectories,” in *2013 IEEE International Conference on Computer Vision*, Dec 2013, pp. 3551–3558.
- [4] P. Weinzaepfel, Z. Harchaoui, and C. Schmid, “Learning to track for spatio-temporal action localization,” in *2015 IEEE International Conference on Computer Vision (ICCV)*, Dec 2015, pp. 3164–3172.
- [5] A. Karpathy, G. Toderici, S. Shetty, T. Leung, R. Sukthankar, and L. Fei-Fei, “Large-scale video classification with convolutional neural networks,” in *2014 IEEE Conference on Computer Vision and Pattern Recognition*, June 2014, pp. 1725–1732.
- [6] J. Carreira and A. Zisserman, “Quo vadis, action recognition? a new model and the kinetics dataset,” in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017, pp. 4724–4733.
- [7] Z. Qiu, T. Yao, and T. Mei, “Learning spatio-temporal representation with pseudo-3d residual networks,” in *2017 IEEE International Conference on Computer Vision (ICCV)*, Oct 2017, pp. 5534–5542.
- [8] X. Peng and C. Schmid, “Multi-region two-stream R-CNN for action detection,” in *ECCV - European Conference on Computer Vision*, Amsterdam, Netherlands, Oct. 2016, vol. 9908 of *Lecture Notes in Computer Science*, pp. 744–759, Springer.
- [9] G. Gkioxari and J. Malik, “Finding action tubes,” in *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015, pp. 759–768.
- [10] C. Gu, C. Sun, D. A. Ross, C. Vondrick, C. Pantofaru, Y. Li, S. Vijayanarasimhan, G. Toderici, S. Ricco, R. Sukthankar, C. Schmid, and J. Malik, “Ava: A video dataset of spatio-temporally localized atomic visual actions,” in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, June 2018, pp. 6047–6056.
- [11] S. Saha, G. Singh, M. Sapienza, P. H. S. Torr, and F. Cuzzolin, “Deep learning for detecting multiple space-time action tubes in videos,” in *Proceedings of the British Machine Vision Conference 2016, BMVC 2016, York, UK, September 19-22, 2016*.
- [12] G. Singh, S. Saha, M. Sapienza, P. Torr, and F. Cuzzolin, “Online real-time multiple spatiotemporal action localisation and prediction,” in *2017 IEEE International Conference on Computer Vision (ICCV)*, Oct 2017, pp. 3657–3666.
- [13] V. Kalogeiton, P. Weinzaepfel, V. Ferrari, and C. Schmid, “Action tubelet detector for spatio-temporal action localization,” in *2017 IEEE International Conference on Computer Vision (ICCV)*, Oct 2017, pp. 4415–4423.
- [14] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. E. Reed, C. Fu, and A. C. Berg, “SSD: single shot multibox detector,” in *ECCV (1)*. 2016, vol. 9905 of *Lecture Notes in Computer Science*, pp. 21–37, Springer.
- [15] S. Ji, W. Xu, M. Yang, and K. Yu, “3d convolutional neural networks for human action recognition,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 1, pp. 221–231, Jan 2013.
- [16] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri, “Learning spatiotemporal features with 3d convolutional networks,” in *2015 IEEE International Conference on Computer Vision (ICCV)*, Dec 2015, pp. 4489–4497.
- [17] R. Hou, C. Chen, and M. Shah, “An end-to-end 3d convolutional neural network for action detection and segmentation in videos,” *arXiv preprint arXiv:1712.01111*, 2017.
- [18] T. Brox, A. Bruhn, N. Papenberg, and J. Weickert, “High accuracy optical flow estimation based on a theory for warping,” in *European Conf. on Computer Vision*, Prague, Czech Republic, 2004, vol. 3024, pp. 25–36.
- [19] C. Feichtenhofer, A. Pinz, and A. Zisserman, “Convolutional two-stream network fusion for video action recognition,” in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016, pp. 1933–1941.
- [20] R. Hou, C. Chen, and M. Shah, “Tube convolutional neural network (t-cnn) for action detection in videos,” in *2017 IEEE International Conference on Computer Vision (ICCV)*, Oct 2017, pp. 5823–5832.
- [21] J. Redmon and A. Farhadi, “YOLO9000: better, faster, stronger,” in *CVPR*. 2017, pp. 6517–6525, IEEE Computer Society.