

Fine-Grained Action Detection with RGB and Pose Information using Two Stream Convolutional Networks

Leonard Hacker^{1,*}, Finn Bartels^{1,*} and Pierre-Etienne Martin^{2,†}

¹Computer Science Institute, University of Leipzig, Germany

²CCP Department, Max Planck Institute for Evolutionary Anthropology, D-04103 Leipzig, Germany

Abstract

As participants of the MediaEval 2022 Sport Task, we propose a two-stream network approach for the classification and detection of table tennis strokes. Each stream is a succession of 3D Convolutional Neural Network (CNN) blocks using attention mechanisms. Each stream processes different 4D inputs. Our method utilizes raw RGB data and pose information computed from MMPose toolbox. The pose information is treated as an image by applying the pose either on a black background or on the original RGB frame it has been computed from. Best performance is obtained by feeding raw RGB data to one stream, Pose + RGB (PRGB) information to the other stream and applying late fusion on the features. The approaches were evaluated on the provided TTStroke-21 data sets. We can report an improvement in stroke classification, reaching 87.3% of accuracy, while the detection does not outperform the baseline but still reaches an IoU of 0.349 and mAP of 0.110.

1. Introduction

While there have been great advances in detection of coarse-grained action in videos (e.g. the type of sport being performed), fine-grained action detection is inherently more difficult due to its low inter-class variability [1, 2]. The goal of this benchmark task is to provide viable tools to enable analyzing athletes' performance [2]. Table Tennis entails many interesting challenges for fine-grained video detection, e.g. ball trajectory prediction [3] or real-time score and game analysis [4].

In the field of image recognition, a CNN consisting of Convolutional layers, ReLU layers and Max Pooling layers is a conventional practice [5]. The provided baseline model in the competition originates from this [6, 7]. While video data includes an additional dimension (time or frame number), several approaches adapt a plain CNN to determine movement or changes in a range of frames [1]. A popular approach introduced first by Simonyan and Zisserman [8] is the Two-Stream Neural Network. They implemented an architecture consisting of a spatial and a temporal CNN. The aforementioned spatial stream represents a single frame at each time while the temporal stream is built as a multi-frame Optical Flow (OF) recognizing the momentum of movement in multiple RGB frames. Hence, the model obtains the additional benefit of *complementary information* provided by a second stream. Thus, areas within the video covering no movement can be removed easily to reduce noise [1]. A successful implementation built upon this is the Inflated 3D Convolutional Neural Network (I3D) model [9], where multiple images are pushed into a 3D CNN instead of a single image at a time. Feichtenhofer et al. [10] proposed

MediaEval'22: Multimedia Evaluation Workshop, January 13–15, 2023, Bergen, Norway and Online

*These authors contributed equally.

† Corresponding author.

✉ pierre_etienne_martin@eva.mpg.de (P. Martin)

🌐 www.eva.mpg.de/comparative-cultural-psychology/staff/pierre-etienne-martin (P. Martin)

🆔 0000-0002-9593-4580 (P. Martin)

© 2022 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

CEUR Workshop Proceedings (CEUR-WS.org)

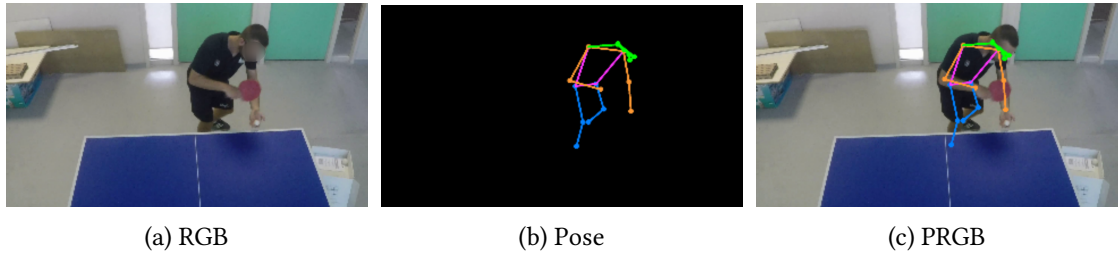


Figure 1: RGB, Pose and PRGB frames of the TTStroke-21 videos. The generated frames serve as input for the network branches.

networks called SlowFast based on a two-stream architecture. The first stream processes the frames with a low frame rate and the second stream with a high frame rate so that the semantic information and motion information are considered.

The fusion step is an essential part of multiple stream networks. Within the models, there are different approaches differing from where the fusion is performed. Late fusion is the easiest and one of the most efficient ways to proceed [11, 1]. Fusion can be performed after a fully connected ReLU layer before a final Softmax function [12, 13]. Early fusion is usually performed at an earlier stage of the network, meaning that the information of the second stream is pushed into the first stream [10, 14]. The following sections describe the architecture of two streams, the results of stroke classification and stroke detection and a discussion on preliminary factors and approaches.

2. Method

As mentioned, the two-stream architecture is one of the state-of-the-art methods. For this contribution, the baseline provided by MediaEval is extended into a two-stream network utilising raw RGB images and pose information. The baseline itself is a single stream 3D CNN with an attention mechanism. Recent papers suggest that the utilization of pose information can achieve better results for fine-grained action detection than OF [15, 16, 17]. When considering actions that are performed by people, the pose holds significant information. We superimpose this information by drawing the pose information on top of RGB images and feeding that into a second stream. The implementation of our method is available online¹.

2.1. Pose Estimation

The pose information is the traced human pose in each frame as depicted in Figure 1. Since OF is already well researched and the use of pose information yields promising results, this work focuses on using pose information to the two-stream CNN framework. The pose information is extracted using the MMPose [18] toolbox from OpenMMLab: an open-source package for detailed video understanding. Each frame of the input video is analyzed. First, a person detector is used to draw bounding boxes around every person in the frame, then a pose estimator is deployed to extract the poses out of the bounding boxes. For person detection we utilized a faster region-based CNN [19] and for pose estimation deep high-resolution representation learning [20]. A top-down classifier is utilized for keypoint extraction, as it performs better than bottom-up classifiers [15, 21]. Both models were pre-trained on the COCO data set [21] by OpenMMLab. The individual keypoints are then connected by differently colored segments representing body parts superimposed over an image. Since different body parts contribute

¹<https://github.com/fidsinn/SportTaskME22>

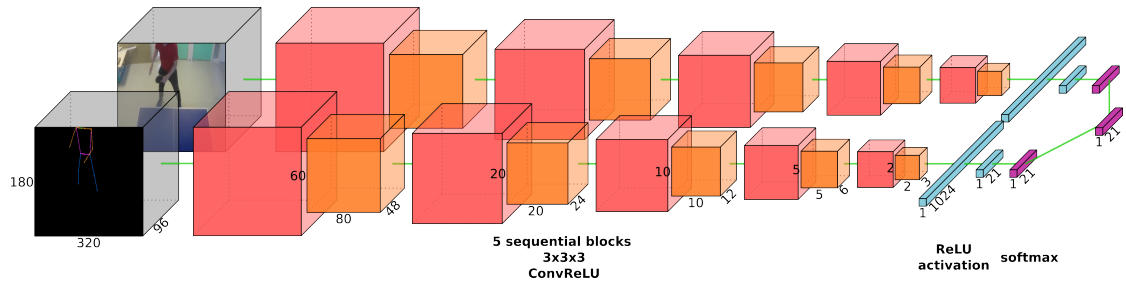


Figure 2: Two Stream Pose Convolutional Neural Network with RGB images (top branch) and Pose information on a black background (lower branch). A branch consists of five successive convolutional layers using ReLU activation function with a growing number of filters (32, 64, 128, 256, 512), each followed by an attention block (both represented in red) and a pooling layer (in orange) of size $(4 \times 3 \times 2)$ for the two first ones and $(2 \times 2 \times 2)$ for the others.

to the strokes in different ways, we assume that the model can distinguish them based on the coloring of each body part. In this work, two methods utilizing pose information are investigated: Pose and PRGB. The Pose variant contains the computed pose over a black background, while the PRGB uses the original RGB frame as a background. By comparing pose, PRGB and RGB performance, it is possible to determine how well the network utilizes the pose information.

2.2. Architecture

Our model is a variant of the two-stream architecture [22] and an extension of the provided baseline [7], which also uses 3D convolutional blocks and an attention mechanism. As depicted in Figure 2, the Two Stream Pose Convolutional Neural Network (TSPCNN) consists of two identical streams, each with five convolutional layers and pooling layers with an increasing number of filters leading to a linear layer with ReLU activation. The latest feeds a second linear layer followed by a Softmax function to convert the output into a 21-dimensional probabilistic vector for classification (21 different stroke types including non-stroke class) or a two-dimensional vector for detection (stroke and non-stroke class). The output of the two branches is then summed and processed by the last Softmax function to have a probabilistic output for classification and detection. The first Softmax function normalizes the output of each individual stream before fusion to minimize vanishing gradients.

2.3. Fusion

Literature suggests that employing early fusion combined with late fusion boosts the performance of a two-stream model considerably [15]. Adding multiple fusion methods to the TSPCNN showed limited gain in performance. The best performance was achieved using a late fusion approach, i.e. fusing before the last layer. Different fusion styles were i) weighted fusion, where the resulting feature is equal to the weighted sum of the two fused features, ii) summed fusion where the resulting feature is the sum of both features and iii) concatenated fusion, where the resulting feature is a concatenation of both features. Summed fusion performed best, and therefore, only its results are reported in the remaining of this paper.

2.4. Training

All experiments were performed on Tesla V100 GPUs provided by the University of Leipzig. The training took 7 to 8 minutes per epoch for the detection task and 1 to 2 minutes for the classification task over 2000 epochs, with a learning rate of 0.0001 and a momentum of 0.5.

Table 1

Results on classification accuracy (%) and detection metrics (accuracy, IoU, mAP (%)) of baseline, one stream approaches and two stream approaches.

Models	Classification			Detection			
	Train	Validation	Test	Train	Validation	Test IoU	Test mAP
Baseline	-	0.813	0.864 ²	-	-	0.515 ⁵ (0.365 ³)	0.131 ⁵ (0.118 ⁴)
Pose	0.995	0.878	0.847 ¹	0.862	0.591	0.205 ¹	0.046 ¹
PRGB	0.978	0.813	0.864 ¹	0.980	0.834	0.165 ¹	0.036 ¹
RGB and Pose	1	0.830	0.872 ¹	0.987	0.820	0.331 ¹	0.100 ¹
RGB and PRGB	0.998	0.848	0.873 ¹	0.990	0.840	0.349 ¹	0.110 ¹

Decision method: ¹No Window, ²Gaussian, ³Mean, ⁴Vote, ⁵Vote (Sliding Window)

3. Results and Discussion

We evaluated our approach using the TTStroke-21 data set [23] provided by the Sport task organizers of MediaEval [2]. The results are compared with the provided baseline [7]. To compare our approach with the baseline, we evaluated our runs using accuracy for the classification task and IoU and mAP for the detection task. In addition, we added the results for training accuracy and validation accuracy for each model, reported in Table 1. The test results were selected depending on which of several decision methods produced the best results [7]. The baseline and the other single-stream approaches already achieve quite promising results for the classification task. Basic RGB combined with PRGB in a two-stream approach shows the best accuracy in testing. The two-stream approaches slightly improve the classification results compared to the single-stream method by up to .009. In contrast to the improvement regarding the classification task, for detection the two-stream methods using the pose stream lead to a decreasing quality compared to the single-stream baseline. The poor detection performance is likely due to the missing ball and racket information in the pose data. An arm movement without the racket may also be similar to a stroke which can confuse the model. Therefore, we can say that pose information did not improve stroke detection performance.

We have shown that the pose information can be suited for fine-grained action classification but seems to fail to capture discriminant features for detection. While the TSPCNN outperformed the baseline in the classification task, we could not improve detection performance. A contributing factor to the slight classification improvement might be the limited training data, especially since some classes only have a few labelled videos in the training set.

4. Summary and Outlook

As wearable systems can be intrusive and cumbersome to set up while also not being widely available, fine-grained action detection from video is of high interest for athletes and coaches to be able to classify different actions in their game and to improve training efficiency [2]. We have built a TSPCNN on state-of-the-art research. Our main contribution is to use RGB data overlaid with human poses in a two-stream network. Our approach slightly outperforms the baseline in terms of classification accuracy but produces poor performances for stroke detection. To improve the TSPCNN further, some more experiments with different qualities of pose data are needed. Moreover, different representations of pose data can be evaluated such as thicker lines to emphasise poses. The approach should also be validated with different data sets, such as the Finegym data set [24] since it also has low variability between classes but a more even class distribution in the training data.

References

- [1] Y. Zhu, X. Li, C. Liu, M. Zolfaghari, Y. Xiong, C. Wu, Z. Zhang, J. Tighe, R. Manmatha, M. Li, A comprehensive study of deep video action recognition, arXiv preprint arXiv:2012.06567 (2020).
- [2] P. Martin, J. Calandre, B. Mansencal, J. Benois-Pineau, R. Péteri, L. Mascarilla, J. Morlier, Sport task: Fine grained action detection and classification of table tennis strokes from videos for mediaeval 2022, in: MediaEval, CEUR Workshop Proceedings, CEUR-WS.org, 2022.
- [3] H. Lin, Z. Yu, Y. Huang, Ball tracking and trajectory prediction for table-tennis robots, *Sensors* 20 (2020) 333.
- [4] R. Voeikov, N. Falaleev, R. Baikulov, Ttnet: Real-time temporal and spatial video analysis of table tennis, in: Proceedings of the IEEE/CVF CVPR Workshops, 2020, pp. 884–885.
- [5] S. Albawi, T. A. Mohammed, S. Al-Zawi, Understanding of a convolutional neural network, in: 2017 International Conference on Engineering and Technology (ICET), 2017, pp. 1–6.
- [6] P. Martin, Spatio-temporal cnn baseline method for the sports video task of mediaeval 2021 benchmark, in: MediaEval, CEUR Workshop Proceedings, CEUR-WS.org, 2021.
- [7] P. Martin, Baseline method for the sport task of mediaeval 2022 benchmark with 3d cnn using attention mechanism, in: MediaEval, CEUR Workshop Proceedings, CEUR-WS.org, 2022.
- [8] K. Simonyan, A. Zisserman, Two-stream convolutional networks for action recognition in videos, *Advances in neural information processing systems* 27 (2014).
- [9] J. Carreira, A. Zisserman, Quo vadis, action recognition? a new model and the kinetics dataset, in: IEEE CVPR, 2017, pp. 6299–6308.
- [10] C. Feichtenhofer, H. Fan, J. Malik, K. He, Slowfast networks for video recognition, in: IEEE/CVF international conference on computer vision, 2019, pp. 6202–6211.
- [11] P. Martin, J. Benois-Pineau, B. Mansencal, R. Péteri, J. Morlier, Classification of strokes in table tennis with a three stream spatio-temporal cnn for mediaeval 2020, in: Proc. of the MediaEval 2020 Workshop, Online, 14-15 December 2020, 2020.
- [12] P.-E. Martin, Fine-Grained Action Detection and Classification from Videos with Spatio-Temporal Convolutional Neural Networks. Application to Table Tennis., Theses, Université de Bordeaux ; Université de la Rochelle, 2020. URL: <https://hal.archives-ouvertes.fr/tel-03099907>.
- [13] P. Martin, J. Benois-Pineau, R. Péteri, J. Morlier, Sport action recognition with siamese spatio-temporal cnns: Application to table tennis, in: CBMI, IEEE, 2018, pp. 1–6.
- [14] C. Feichtenhofer, A. Pinz, A. Zisserman, Convolutional two-stream network fusion for video action recognition, in: IEEE CVPR, 2016, pp. 1933–1941.
- [15] H. Duan, Y. Zhao, K. Chen, D. Shao, D. Lin, B. Dai, Revisiting skeleton-based action recognition, arXiv preprint arXiv:2104.13586 (2021).
- [16] P. Martin, J. Benois-Pineau, R. Péteri, J. Morlier, Optimal choice of motion estimation methods for fine-grained action classification with 3d convolutional networks, in: ICIP, IEEE, 2019, pp. 554–558.
- [17] S. Sato, M. Aono, Leveraging human pose estimation model for stroke classification in table tennis, in: MediaEval, volume 2882 of *CEUR Workshop Proceedings*, CEUR-WS.org, 2020.
- [18] M. Contributors, Openmmlab pose estimation toolbox and benchmark, <https://github.com/open-mmlab/mmpose>, 2020.
- [19] S. Ren, K. He, R. Girshick, J. Sun, Faster r-cnn: Towards real-time object detection with region proposal networks, *Advances in neural information processing systems* 28 (2015).
- [20] K. Sun, B. Xiao, D. Liu, J. Wang, Deep high-resolution representation learning for human pose estimation, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2019, pp. 5693–5703.
- [21] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, C. L. Zitnick, Microsoft coco: Common objects in context, in: ECCV, Springer, 2014, pp. 740–755.
- [22] P. Martin, J. Benois-Pineau, R. Péteri, J. Morlier, 3d attention mechanisms in twin spatio-temporal convolutional neural networks. application to action classification in videos of table tennis games., in: ICPR, IEEE Computer Society, 2021.
- [23] P. Martin, J. Benois-Pineau, R. Péteri, J. Morlier, Fine grained sport action recognition with twin spatio-temporal convolutional neural networks, *Multim. Tools Appl.* 79 (2020) 20429–20447.
- [24] D. Shao, Y. Zhao, B. Dai, D. Lin, Finegym: A hierarchical video dataset for fine-grained action understanding, in: CVPR, IEEE, 2020, pp. 2613–2622.