# ConvLSTM for Table Tennis Stroke Classification

Jansi Rani Sella **Veluswami**[1], Ananth Narayanan P[1], Bhuvan S[1] and Shobith Kumar R[1]

[1]*Sri Sivasubramaniya Nadar College of Engineering, Tamil Nadu, India*

#### Abstract
Our study concentrates on sports video analytics, particularly stroke classification. We utilize a model that combines Convolutional Neural Networks (CNN) and Long Short-Term Memory (LSTM) trained on the MediaEval Fine-Grained Action Classification of the Table Tennis Strokes dataset. With an accuracy of 81.4%, our model effectively classifies table tennis moves, providing insights for post-match commentary and playstyle analysis. This effectiveness is demonstrated in the context of the MediaEval 2023 benchmark.

## 1. Introduction

The field of action recognition involves associating a predefined set of actions with video content to meet the increasing demand for automated action analysis in videos. This paper presents a method that specifically targets the classification of strokes within a dataset of various table tennis strokes performed in match and practice settings. The action recognition process involves localizing objects, identifying them, and then classifying the detected actions. The ability to detect and classify actions is crucial for making strategic decisions, particularly in the context of athletic performance analysis.

The Overview paper [1] describes the dataset TTStroke-21 used in this study which includes 21 different classes of strokes, where two annotated sets are provided: a training and a validation set. Utilizing machine learning in this domain has the potential to enhance athletic performance through computer-aided analysis of moves. In this study, we developed a model implemented using TensorFlow, using a combination of Convolutional Neural Network (CNN) and Long Short-Term Memory (LSTM) architecture. Our approach aims to contribute to the improvement of athletic performance by automating the analysis of various strokes. We discuss the results obtained using our model on the given dataset, highlighting the significance of effective action recognition in sports analytics.

## 2. Related Work

The provided baseline methodology [2] proposes two types of 3D-CNN architectures to solve the subtask. Both the methods are 3D-CNN architectures using Spatio-temporal convolutions and attention mechanisms. The predominant strategies have centered around the utilization of CNN and LSTM-based methodologies. For example, in the paper by Kaustubh Milind Kulkarni et al. [3], an LSTM model, a TCN model, and a combined TCN + LSTM model were presented. They used Pose Estimation and a Savitzky-Golay filter for feature extraction. Kadir Aktas et al. [4], present another approach where RGB images were used as the input data without any prior feature extraction. They used an LSTM model to achieve about 79.8% accuracy in validation data.

CEUR Workshop Proceedings (CEUR-WS.org)

We were inspired by the idea of using the RGB images directly without any feature extraction and executed the same in our work.

## 3. Approach

A Convolutional Neural Network (CNN or ConvNet) is a type of deep neural network specifically designed for processing image data. This network excels in analyzing images and making predictions based on them. It utilizes kernels, known as filters, to examine the image and generate feature maps, which represent the presence of specific features at various locations within the image. Initially, the network produces a limited number of feature maps, which are augmented and refined through subsequent layers using pooling operations, while retaining critical information without loss.

On the other hand, a Long Short-Term Memory (LSTM) network is specifically designed to handle sequential data, taking into account all previous inputs to generate an output. LSTMs are a type of Recurrent Neural Network (RNN) that addresses the vanishing gradient problem, a limitation of traditional RNNs in handling long-term dependencies in input sequences. This enables LSTM cells to maintain context for extended periods, making them better suited for tasks such as time series prediction, speech recognition, language translation, and music composition.

In the context of action recognition, we will employ a CNN + LSTM network to leverage the spatial-temporal aspects of videos. This combination will enable the network to effectively analyze and recognize actions within video sequences.

### 3.1. Data Preprocessing

The data preparation process involves class identification and two pivotal functions: one for frame extraction, ensuring resizing and normalization, and another for dataset construction, incorporating features, labels, and video paths. Notably, the dataset creation rigorously filters videos to align with the specified sequence length. The execution of the dataset creation function on a specified directory results in the generation of dataset objects, including features, labels, and video file paths. These components include features, representing extracted frames from videos, and labels, serving as identifiers for subsequent machine learning model training. The third component consists of paths associated with videos in the dataset, functioning as references to the physical location of each video.

### 3.2. Proposed Model

The process of creating a dataset for TensorFlow is seamless and incorporates both Convolutional Neural Network (CNN) and Long Short-Term Memory (LSTM) architecture. This choice is informed by the well-established effectiveness of these architectures in video content analysis tasks.

In the construction phase of our model, the Keras ConvLSTM recurrent layers, a critical architectural decision for video classification tasks is utilized. These layers excel at processing spatiotemporal information within video sequences. We configure the layer with parameters such as the number of filters, kernel size, and activation function to facilitate convolutional operations. The resulting sequences are subsequently processed through various other function layers, reducing frame dimensions to alleviate computational load, and Dropout layers, mitigating overfitting risks. The architecture is intentionally kept simple with a limited number of trainable parameters, commensurate with the scale of the dataset. A vital element is the

incorporation of a final Dense layer with softmax activation, yielding probability distributions across action categories.

The constructed model is then compiled using categorical cross-entropy as the loss function, the Adam optimizer, and accuracy as the metric for evaluation. Training is initiated, incorporating an early stopping callback to prevent overfitting. This structure forms a cohesive and efficient framework for the in-depth analysis of spatiotemporal patterns within table tennis stroke videos. The model's adherence to best practices in architectural design and training strategies enhances its adaptability and potential for robust performance in action recognition tasks.

## 4. Results and Analysis

The accuracy of the model was updated after every layer of training, and the results demonstrate a high level of accuracy. The training data accuracy reached a peak of 97%, while the validation accuracy reached 98.8%. Several factors contributed to these results. Firstly, the volume and distribution of the data had a significant impact on the accuracy. Additionally, the image height, width, and sequence length all had a significant effect on the results, with accuracy ranging from 0.7408 for an image of dimensions 90*80 and a sequence length of 60, to 0.9876 for an image of dimensions 64*64 and a sequence length of 60. It is worth noting that the data distribution of some labels is highly biased towards certain classes, leading to biased learning. Over the course of five runs, the highest global accuracy achieved by the model was 81.4

| RUN | Hand | Serve | Hand & Serve | Global |
|-----|------|-------|--------------|--------|
| Run1 | 91.5 | 92.4 | 90.7 | 75.4 |
| Run2 | **93.2** | 89.8 | 89.0 | 74.6 |
| Run3 | 92.4 | **94.1** | **92.4** | **81.4** |
| Run 4 | 89.0 | 89.8 | 86.4 | 72.0 |
| Run 5 | 89.8 | 88.1 | 87.3 | 72.9 |

**Figure 1:** Accuracy across various runs.

## 5. Discussion and Outlook

Throughout the training and validation phase, we have attained encouraging outcomes that lead us to conclude that overfitting is not present.

Nonetheless, the model's performance on the test data reveals that it has not effectively learned and is unable to generalize. In our opinion, this challenge can be remedied by augmenting the quantity of labeled data utilized in training.

Moreover, we posit that the low variability between the classes and the nature of the task contribute to this issue. Considering that a single class can be sampled in various ways for different players, such as right/left-handed or high/low experienced, we suggest that the dataset could be enhanced by increasing the coverage of the classes and reducing bias among them.
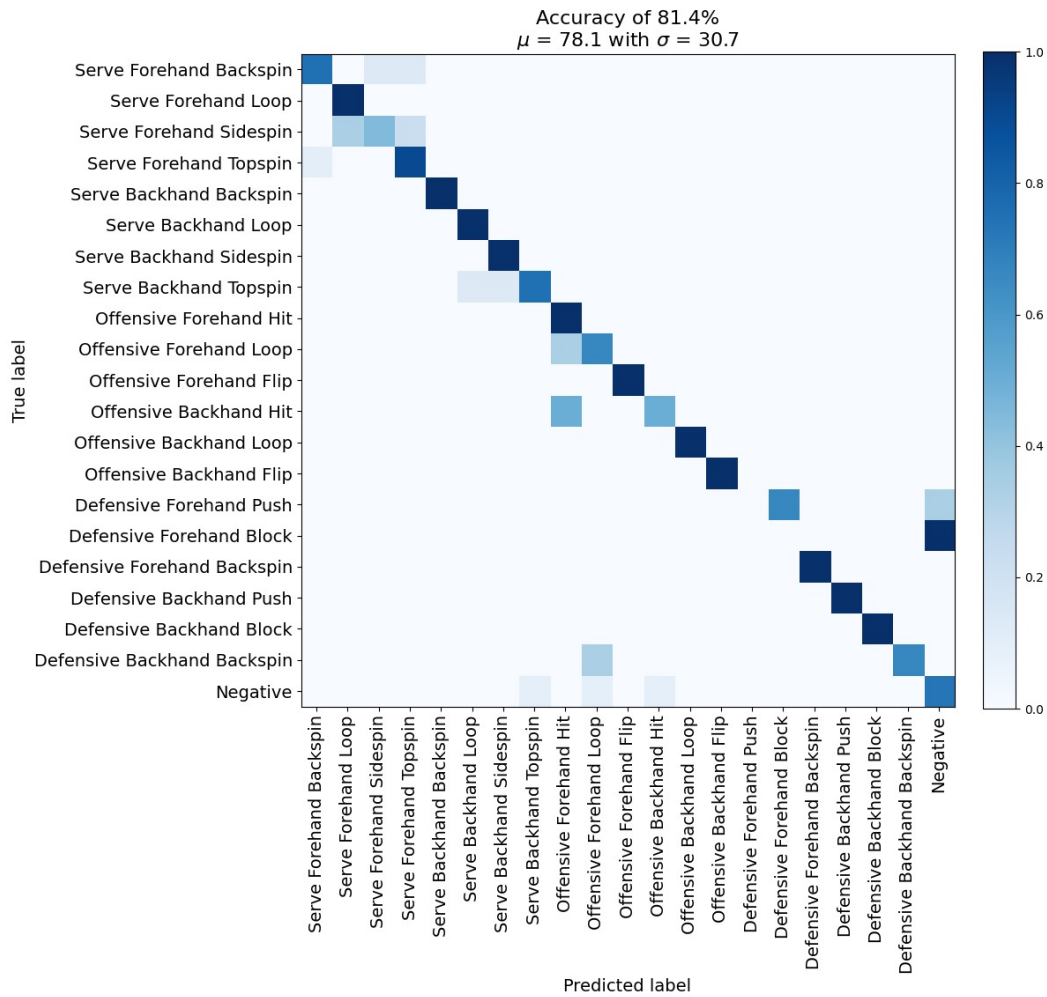
**Figure 2:** Heatmap of the model with highest global accuracy.

# References

[1] P.-E. Martin, Baseline method for the sport task of mediaeval 2023 3d cnns using attention mechanisms for table tennis stoke detection and classification., MediaEval Workshop 2023 (2023).

[2] A. Erades, P.-E. Martin, R. Vuillemot, B. Mansencal, R. Peteri, J. Morlier, S. Duffner, J. Benois-Pineau, SportsVideo: A Multimedia Dataset for Event and Position Detection in Table Tennis and Swimming, MediaEval Workshop 2023 (2023).

[3] K. M. Kulkarni, S. Shenoy, Table tennis stroke recognition using two-dimensional human pose estimation, CVPR Sports Workshop (2021).

[4] K. Aktas, M. Demirel, M. Moor, J. Olesk, G. Anbarjafari, Spatio-temporal based table tennis hit assessment using lstm algorithm, MediaEval (2020).

[5] A. Zahra, P.-E. Martin, Two stream network for stroke detection in table tennis, MediaEval (2021).

[6] HCMUS at MediaEval'20: Ensembles of Temporal Deep Neural Networks for Table Tennis Strokes Classification Task, 2020.

[7] P.-E. Martin, J. B. Pineau, B. Mansencal, R. Péteri, J. Morlier, Siamese spatio-temporal convolutional neural network for stroke classification in table tennis games (2020).