

# Action Recognition in videos using Deep Learning

Thesis submitted as a partial fulfillment of the requirements for the degree of Master of Science in Computer and Information Sciences

# By Bassel Safwat Chawky Hakim

Teaching Assistant at Scientific Computing Department, Faculty of Computer and Information Sciences, Ain Shams University

#### Supervised by:

#### Prof. Dr. Howida Abdel Fattah Shedeed

Professor in Scientific Computing Department, Faculty of Computer and Information Sciences, Ain Shams University

#### Dr. Mohammed Abd El-Rahman Marey

Assistant Professor in Scientific Computing Department, Faculty of Computer and Information Sciences, Ain Shams University

#### **Dr. Ahmed Samir**

Ph. D. in Scientific Computing, Faculty of Computer and Information Sciences, Ain Shams University

### Acknowledgment

First of all, I would to like to thank GOD for his endless blessings, for giving me the power and strength to complete this work and for surrounding me supportive people.

I would like also to thank my Professor Dr. Howida. Her office door was always open when I ran into trouble or had a question about my research or writing. She consistently allowed our researches to be my own work but steered me in the right direction whenever she thought I needed it. Also, I would like to express my deep appreciations for the scientific vision of Dr. Mohammed Marey and his support through all the meetings that lasted till late evenings besides his continuous pushes to reach my extremes and to achieve several successes. Next, I would like to thank Dr. Ahmed Samir Elons not only for his scientific vision and guidance but also for his patience during my early years in my masters.

Special thanks for my family who always believed in me and supported me through my toughest times. I would like to explicitly thank my mother for her permanent trust and acceptance for the way I am, pushing me to always do my best and to seek quality in my work. Thank you again for all the efforts you put to allow me achieving my goals and targets and your endless support.

My friends who always look at me as a man with potential, thank you for your trust and for existing in my life. I appreciate those time when you kept asking me to pursuit my dreams and to finish my masters successfully.

Special thanks for the theatre team for their efforts to record the dataset and support for the research. Last but not least, I would like to thank my professors, colleagues and students who kept encouraging me. Thank You!

**Bassel Safwat** 

#### **Abstract**

Human action recognition for a given video is a difficult problem containing many challenges ranging from partial occlusion to variations in the action speeds and viewpoints. However, this problem is at the very core of various systems like abnormal behavior detection, action localization and/or online video analysis, and that is why it is a very important problem which got the attention from many researchers in the past decades and till now.

Despite the number of studies in the literature, the action recognition problem remains a difficult problem. Traditional approaches require a lot of efforts to find the best combination of features not only to represent the action in a compact form but also to handling the so many existing challenges. Recent methodologies relied on deep learning models to learn and extract good representation from the datasets. Although the deep learning-based models requires a large training datasets and costly training time, this type of models illustrated advances on several action recognition dataset. Moreover, several techniques like transfer learning allowed faster convergence of the accuracy by pretraining the model.

At first, a novel ranking and listing for 14 action recognition datasets is set. The ranking is based on the number of challenges each dataset covers. Therefore, the higher the dataset's rank, the more realistic the dataset is, indicating its ability to provide a realistic measurement for the models. Based on these datasets, a comparison between the advances in traditional approaches and deep-learning based models is illustrated. Based on this deep survey, deep learning baselines models, namely two stream convolutional neural network and 3D convolutional neural networks, are described. These

baseline models are almost used by all the other studies including state of the art models.

Secondly, a novel action recognition benchmark has been set and used to study several action recognition challenges like the change in viewpoints and shaky videos effects using the baseline models described. Moreover, it is used to study the overfitting problem of the deep learning models with potential solutions.

Finally, two main techniques are suggested where both can be integrated by the several existing action recognition models in order to improve the accuracy. These techniques leverage the video temporal dimension to learn several features representing varying temporal lengths. The first technique is called Single Temporal Resolution Single Model (STR-SM) which suggests training the desired model on one specific temporal resolution of the video. The temporal resolution is defined as the number of frames for a specific amount of time. Therefore, a low temporal resolution means that a small number of frames is used to represent the action while for a high temporal resolution, a large number of frames is used. Therefore, a good model that uses the STR-SM technique uses a temporal resolution that is low enough to represent long temporal duration but also, high enough to capture the motion details. Such technique is faster when compared to traditional approaches as it tackles long temporal range at once with better accuracy as it covers more information. On The second technique is called Multi Temporal Resolution Multi Model (MTR-MM) which tackles the problem of varying action speeds in a novel way. Applying the MTR-MM technique on the desired model requires building several STR model versions, each trained on a specific temporal resolution with a late fusion. This leverage the different existing information in each temporal resolution leading to a better and improved accuracy. Additionally, the STR-SM and the MTR-MM techniques are applied on 3D Convolutional Neural Network model and have improvements over the traditional training approach of 3.63% and 6% videowise accuracy respectively.

### **List of Publications**

- Chawky, Bassel S., A. S. Elons, A. Ali, and Howida A. Shedeed. "A
   Study of Action Recognition Problems: Dataset and Architectures
   Perspectives." In Advances in Soft Computing and Machine
   Learning in Image Processing, pp. 409-442. Springer, Cham, 2018.
- Chawky, Bassel S., A. S. Elons, A. Ali, and Howida A. Shedeed.
   "OA18: Deep Learning for Office Actions Analysis. In 2019 9th IEEE
   International Conference on Intelligent Computing and
   Information Systems (ICICIS), (waiting acceptance)
- Chawky, Bassel S., Mohammed Marey, and Howida A. Shedeed. "Multi-Temporal-Resolution Technique for Action Recognition using C3D: Experimental Study." In 2018 13th International Conference on Computer Engineering and Systems (ICCES), pp. 404-410. IEEE, 2018.

## **Table of Contents**

Abstract		III
Acknowledg	gment	II
List of Publi	ications	VI
Table of Co	ntents	VII
List of Figur	res	IX
List of Table	es	XI
List of Abbr	reviations	XII
Chapter 1.	Introduction	
1.1	Overview of Action Recognition and Deep Learning	2
	1.1.1 Challenges	6
	1.1.2 Traditional approaches	4
	1.1.3 Deep Learning approaches	5
1.2	Problem Definition	10
1.3	Main Contributions of this Thesis	10
1.4	Thesis Outline	11
Chapter 2.	Related Work	14
2.1	Action Recognition Datasets	14
2.2	Action Recognition Models	28
	2.2.1 Shallow Models	28
	2.2.2 Deep Models	35
Chapter 3.	Chapter 3. Scientific Background: Deep Learning based Action	
	Recognition	51
3.1	Convolutional Neural Network	51
3.2	Two Stream Convolutional Neural Network	54
	3.2.1 Optical Flow Field	54
	3.2.2 Spatio-temporal score fusion	57
3.3	C3D: 3D Convolutional Neural Network	57
Chapter 4.	New Dataset Benchmark: OA18	62
4.1	Dataset characteristics	63
4.2	Action Recognition Labeling Software	68
4.3	Experimental Results and Analysis for OA18	69
	4.3.1 Experiment A: Baseline models evaluation	69
	4.3.2 Experiment B: Variation in Viewpoints	70
	4.3.3 Experiment C: Video stabilization	71
	4.3.4 Experiment D: Fighting Overfitting	71
Chapter 5.	Proposed Multi temporal resolution Techniques	78
5.1	Temporal Resolution for Action Recognition	78
5.2	Single Temporal Resolution Technique (STR)	80

	5.2.1 Single Temporal Resolution Single Model (STR-SM) 80	
5.3	Multi Temporal Resolution Techniques (MTR)	
	5.3.1 Multi Temporal Resolution Single Model (MTR-SM) 83	
	5.3.2 Multi Temporal Resolution Multi Model (MTR-SM) 84	
5.4	C3D model with Multi-Temporal-Resolution 85	
5.5	Dataset preparation-UCF56	
5.6	Experimental Results and Discussions 8	
	5.6.1 Experiment A: STR-SM experiments	
	5.6.2 Experiment B: MTR-SM experiments	
	5.6.3 Experiment C: MTR-MM experiments	
Chapter 6.	Conclusion and Future Work	
6.1	Conclusion93	
6.2	Future work	
References	98	

## **List of Figures**

Figure 2-1 Bag-of-Words pipeline
Figure 2-2 Stacked Fisher Vector Architecture
Figure 2-3 Conv-Pooling architecture used in [40]. Temporally stacked
convolution features (denoted by 'C') are max-pooled (blue) and fed to
the fully connected network (yellow) followed by a SoftMax layer
(orange)39
Figure 2-4 Deep Long Short Term Memory Architecture with input as
convolutional features and followed by a softmax layer. [40] 41
Figure 2-5 Encoder-Decoder framework. 3D-CNN is used to encode the video
followed by LSTM to produce video captions
Figure 2-6 Left: A Gated RBM. Right: A Convolutional Gated RBM using
probabilistic max-pooling [43]43
Figure 2-7 Deep Learning – Slow Feature Analysis Architecture 44
Figure 3-1 Tiny VGG: A simple CNN model
Figure 3-2 From top to bottom, (a) and (b) represents two successive frames
and their optical flow. (c) illustrates the color and its magnitude for each
direction used in the optical flow56
Figure 3-3 Comparison between the traditional 2D kernels to the left and 3D
kernels to the right. The difference exists in the depth stride 58
Figure 3-4 C3D model consists of 8 convolutional layers, 5 pooling layers and
2 fully connected layers followed by a SoftMax layer 59
Figure 3-5 Visualization of features embedded for C3D on UCF101 showing
its discriminability potential59
Figure 4-1 Cameras position relative to the scene
Figure 4-2 Histogram showing the number of videos for each action class 65

Figure 4-3 Main layout of the developed action labeling software
Figure 4-4. Stabilization Process
Figure 4-5. Training accuracy vs Validation accuracy for the C3D model . 72
Figure 4-6. Samples from the augmented data. From left to right: cropping,
illumination change, moving obstacle, salt & pepper noise
Figure 4-7. Samples from the pretraining dataset
Figure 5-1 Frame selection using both High and Low temporal resolutions.
Figure 5-2 The extraction process of 12 stacks from an input video of 120
frames using steps 1, 3 and 12
Figure 5-3 MTR-SM pipeline. All resolutions are extracted from the video
and then they are all fed to the model for training. For the testing, any
resolution can be used
Figure 5-4 MTR-SM pipeline. The first stage extracts all the different
resolution, and each resolution is used to train a different STR-SM.
Score combination is made using Majority voting at the end of the
pipeline to determine the action class label
Figure 5-5 Illustration of the action classification pipeline using C3D for a
video of 100 frames85

## **List of Tables**

Table 2-1 Summary for action recognition challenges in the 14 datasets
surveilled
Table 2-2 Suggested usage for each dataset
Table 2-3 Summary of the results for both hand-crafted features and deep
learning models
Table 4-1 The min, max, avg and median number of frames for the clips of
each action category67
Table 4-2 The number of clips for each perspective grouped by the offices 67
Table 4-3. C3D vs Two Stream CNN accuracies on OA1870
Table 4-4. Experiment for the side vs front vs both
Table 4-5. Experiment for the stabilization effect
Table 4-6. Accuracy comparison for data augmentation, pretraining, and
pretrained + data augmentation74
Table 5-1 Stack-wise accuracy for testing STR-SM. Each model is trained
on single resolution with the resolution-step stated in the first row
and is tested against different step dataset, with the step stated in
the first column
Table 5-2 Video-wise accuracy for testing STR-SM. Each model is trained
on single resolution with the resolution-step stated in the first row
and is tested against different step dataset, with the step stated in
the first column
Table 5-3 Results of testing a MTR-SM with several datasets extracted from
the UCF-56 dataset. Results are stack wise accuracy90

## **List of Abbreviations**

<u>Abbreviation</u>	Stands for
CNN	Convolutional Neural Network
CRF	Conditional Random Field
FV	Fisher Vectors
HOF	Histogram of Optical Flow
HOG	Histogram of Oriented Gradients
HVS	Human Visual System
LDA	Latent Dirichlet Allocation
LSTM	Long-Short-Term Memory
MBH	Motion Boundary Histogram
MKL	Multiple Kernel Learning
ReLU	Rectified Linear Units
SFA	Slow Feature Analysis
SFV	Stacked Fisher Vectors
SGD	Stochastic Gradient Descent
STR-SM	Single Temporal Resolution Single Model
TDD	Trajectory-pooled Deep-convolution Descriptors
TS	Trajectory Shape

## Chapter 1

# Introduction

Chapter 1.	Introduction	2
-	Action Recognition and Deep Learning Overview	
	1.1.1 Traditional Approaches	4
	1.1.2 Deep Learning Approaches	
1.2	Challenges	
1.3	Problem Definition	10
1.4	Main Contributions of This Thesis	10
	Thesis Outline	

Chapter 1 Introduction

### **Chapter 1. Introduction**

#### 1.1 Action Recognition and Deep Learning Overview

Action recognition is the task of identifying the kind of action presented in the video. The action occurring in the videos vary in complexity, where simple actions are called gestures and complex action is called activity. Therefore, it is important to provide a clear definition for the different action complexities.

Gesture: This is the basic motion performed by the human body parts that have some meaning. This includes hands or legs movements like waving hands or walking, facial expressions like closing/opening eyes and lips movements, and head shaking. This is considered as atomic action as it has the lowest complexity and the shortest amount of time when compared with other action complexities.

**Action:** There is no standard definition for an 'action', however, it is commonly understood as the movements (more than one gesture) performed by a person that involves an interaction with another person or object. In this context, shaking hands, swimming, and kicking a ball are good examples on actions. They usually require orders of few seconds and might last for minutes.

**Activity:** It is a set of actions that occurs either simultaneous or in sequence. Therefore, it is the most complex type of actions and is also called **event**. Examples for different activities: people protesting, a team game or a group meeting. As might be guessed, these activities usually last for long time interval.

Chapter 1 Introduction

This thesis focuses on improving generic deep learning systems for the 'gesture' and 'action'. These actions complexity covers the basic human actions in addition to interaction with another person and or objects.

Following is a discussion of few important applications for action recognition task to illustrate the different areas that are affected by any advances in action recognition.

**A. Video Retrieval [1]:** As the number of videos on the internet is growing, there is a tremendous need to automatically understand the content of the video instead of relying solely on the video tags (i.e. title and description), which might be misleading, to improve the search results quality for the end users.

**B. Automated Surveillance** [2]: The surveillance in large manufactures and governmental institutes requires security cameras that covers the whole place to record the daily activities occurring. There is usually a security person monitoring these videos to notify for any abnormal behaviors and take the required decision. Automating this task would benefit in improving the accuracy and reducing the costs for workers performing a tedious task.

**C. Human-Computer Interaction [3]:** Many modern games provide sensors like web-cameras, kinetic or Virtual Reality (VR) controllers to capture the human action and provide a better entertaining gaming experience. The quality of understanding the human action results in a better user experience.

**D. Video Description and Summarization [4]:** Similar to video retrieval, generating a better video description is another task that relies on understanding the video content. Increasing such system with actions