

## ANALYSIS AND RECOGNITION OF HUMAN BY CONCEPTUAL FEATURES

Shivam Gadhadara

### Abstract

Human action recognition is the process of labeling a video according to human behavior. This process requires a large set of labeled video and analyzing all the frames of a video. The consequence is high computation and memory requirement. This paper solves these problems by focusing on a limited set rather than all the human action and considering the humanobject interaction. This paper employs three randomly selected video frames instead of employing all the frames and, Convolutional Neural Network extracts conceptual features and recognize the video objects. Finally, support vector machine determines the relation between these objects and labels the video. The proposed method have been tested on two popular dataset; UCF Sports Action and Olympic Sports. The results show improvements over state-of-the-art algorithms.

**Index Terms**—Computer Vision, Human Activity Recognition, Convolutional neural networks, Support vector machine



*Scholarly Research Journal's* is licensed Based on a work at [www.srjis.com](http://www.srjis.com)

## INTRODUCTION

Recently, analyzing and understanding human action or activity become an interesting topic because of two factors. First, the advancement of technology and the increase of low cost and powerful imaging equipment result in exponential growth of video creation.

Second, the development of a large number of programs including human robot interaction, human-computer interaction, intelligent video surveillance, face analysis, object tracking, video processing and video annotating, robotics, smart aware-house, rehabilitation center, video games and a variety of systems that involve interactions between human and computer

Human behaviors are analyzed according to gesture, action and activity Gesture or elementary action includes automatic and simple movement such as hand raising or foot forwarding. Action is a series of gesture that temporarily put together and they describe the entire body. Finally, the activity is a series of actions which include interactions and group activities [1]-[2][5]. Also, interactive activities include human-object or human-human interactions (See Fig. 1). This paper focuses on human-object interaction to understand the human activity.



**Fig. 1. Human behaviors based on complexity level.**

Although, human action and activity recognition has started since 1973 [1], unsolved issues have been remained such as view point, clutter, diversity of actions, actor movement variations, high cost computing and memory requirement. One of the main problems in this area relates to the variation of human actions and activity.

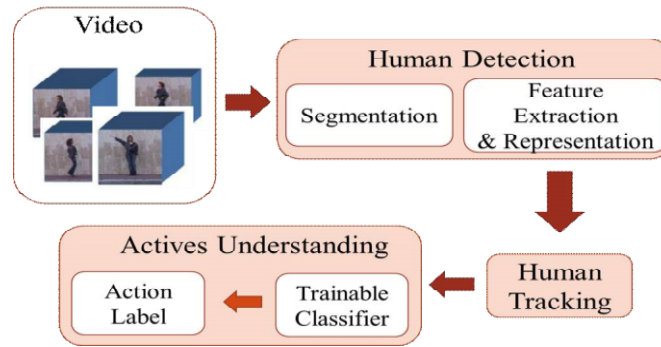
As previously mentioned, we have three categories for human behavior. Also, each person has his own style [2]. Therefore, designing a high performance recognition system for all categories of human behaviors is complicated. A solution is to implement system to recognize a limited number of actions or activities.

This paper focuses on improving the recognition performance by limiting the variation of actions by understanding human-object interaction. It employs only three randomly selected video frame instead of all the frames and, a pre-trained Convolutional Neural Network (CNN) via ImageNet pictures extracts the high level and conceptual features and recognize the video objects and, Support Vector Machine (SVM) understands the action by determining the relation between the objects and labels the video.

The rest of this paper is organized as follows. Section 2 reviews the human action or activity recognition systems. Section 3, explains the theoretical background; convolutional neural networks and support vector machine. Section 4 explains the proposed system. Section 5 present the simulation results and comparison with other along the same systems. The final section concludes the paper and talks about future activities.

## II. LITERATURE REVIEW

Figure 2 shows the general form of a human action or activity recognition systems where it consists of four sections: input, detection, tracking and recognition [1]-[2]-[6] Although action and activity are semantically difference, but in most cases it can be considered that there is no different and they are equal [7]. Next, we discuss each section briefly.



**Fig. 2. Human action recognition system.**

*A. Input*

A recognition system receive input data such as video or still image from visual devices such as cameras. This paper focuses on systems that the location of camera is unimportant and the input is video.

*B. Detection*

The important section in this system is detecting human. This section separates the main object in video frames. The detection is performed by data representation and the features extracted should be robust against small video changes such as human appearance, rotation and displacement, lighting, partial human latency, occlusion, point of view or run operation. The system have to fix its human of interest [1]-[4]-[6]-[8]. Recently, researchers use more hierarchal structure for action recognition.

*C. Tracking*

After finding the human’s location in each frame, the system locates a few points of human body. We use these points to track the human body and extracting the movement pattern. This part distinguishes human’s movements from other objects’ movements [1]-[6] Morris et all [9] divided the object tracking algorithms into 6 groups based on region, contour, feature, model, hybrid and optical flow.

*D. Understanding*

The role of understanding is analyzing the motion pattern and find the best description of actions and activities. Understanding human behavior is classification the human action and activity. The classification matches an unknown data with a group of sample reference and labeling human action and activity. The learning methods involves supervised, unsupervised and semi-supervised model [6]-[10]. Table 1 present the several popular algorithms on human action recognition which used deep neural networks for action recognition.

### III. BACKGROUND THEORY

In This section introduces the sub-systems of the proposed system including CNN and SVM. This introduction makes understanding this paper more easily. The role of CNN and SVM are extracting suitable features and classification of action and activity respectively.

#### A. Convolutional Neural Networks

Alex Krizhevsky and colleagues used CNN to classify 2.1 million high-resolution images into 1000 classes. CNN improves the classification performance compared to the other methods significantly and it overcomes overfitting which is observed in most learning issues. Figure 3 shows the structure of CNN employed by Krizhevsky and colleagues. This network can run in parallel on GPUs [14].

CNN composes of 8 layers, 5 convolutional layers and 3 fully connected layers. The small local parts of the input are captured by the convolutional layer with a set of local filters, and the pooling layer preserves the invariant features. Top fully connected layers combine inputs of all features to perform the classification [14]-[15]. This hierarchical organization generates good results in image processing and it had been applied to large-scale visual recognition tasks because low-level edge information leads to complex representations such as mid-level cues like edge intersections or high-level representation like object parts. Also, this method is convenient for extracting dependent\_local features

**Table I Presents the Results of Methods**

First author	Algorithm
Ji, 2013 [11]	Develop a 3D convolutional neural network. It is an extension of 2D CNN. Extracts features from both spatial and temporal dimensions.
Simonyan, 2014 [12]	Represent a two-stream architecture of CNN. <i>Spatial stream</i> captures appearance information from still frames and <i>Temporal stream</i> extracts

optical flow to capture motion between frames.

Presented model called trajectory pooled deep convolutional descriptor (TTD). He used hand-crafted features to extract trajectories and a two stream CNN for extracting feature maps.

Wang, 2015 [13]

---

and scale invariant feature [16]-[14].

### B. Support Vector Machine

Vapnyk introduced primary SVM in 1963 [17] SVM is a supervised learning algorithm for classification and regression and, it is widely used in object detection & recognition, content-based image retrieval, text recognition, biometrics, speech recognition [18]. The advantages of SVM are the training simplicity, high interoperability and, it works well with low samples and high dimensional input space. Experimental results have shown that this method often outperforms competing clustering methods [18]-[19]-[20].

SVM classifies data by finding the best hyperplane that separates all data points into two classes. The training data is a set of data points  $x \in \mathbb{R}^d$  along with their categories  $y_i = \pm 1$ ,  $i$

where each value corresponds to a data class. The equation of a hyperplane is:

$$g(x) = w^T X - b \tag{1}$$

Where  $g(x)$  is a linear function,  $W \in \mathbb{R}^d$  and  $b$  is a real number. A positive or negative value of  $g(x)$  means that the data will have  $y_i = +1$  or  $y_i = -1$  respectively.

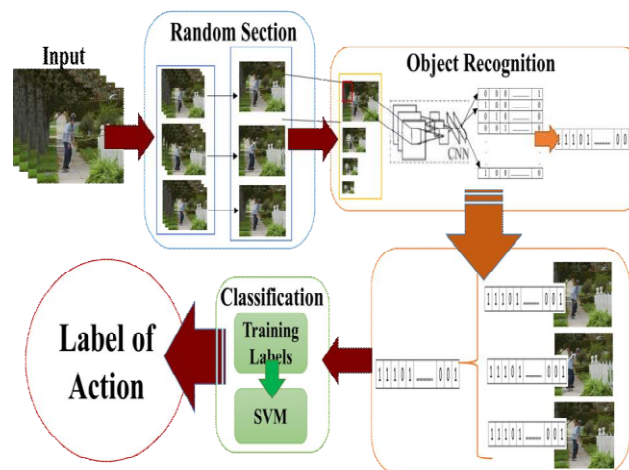
In cases where it is difficult to find a simple hyperplane to separating data, we simply increase the data dimension by a kernel. The commonly-used kernel functions are Linear, Polynomial, Gaussian and Sigmoid [18]-[19]-[20].

#### IV. THE PROPOSED METHOD

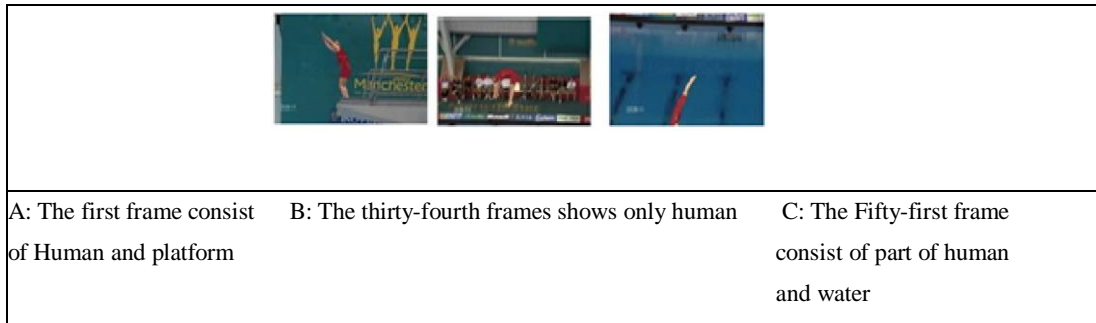
Figure 3 shows the block diagram of the proposed method. This algorithm consists of three steps including: random selection of video frames, object and human recognition and understanding human action.

##### 1) Random selection of video frame

At first we have to find all the video objects. We can use one video frame instead all the frame to find the video objects. However, we may lose video objects in situations where we may face with occlusion or fast video changes. Figure 4 presents such an example in a diving video. Overcoming these problems, we divide a video into three parts and, randomly select a frame of each part as the input of CNN to detect the video objects. The advantages of this random selection are the improvement of computation complexity and reducing the memory requirement.



**Fig. 3. The block diagram of the proposed system including random video frame selection, object recognition by CNN and action recognition according to objects of the frame by SVM.**



**Fig. 4. Three frames of diving movie; fast changes of movement results in different objects in each frame.**

2) *Object and human recognition*

This section employs a pre-trained CNN with ImageNet [21] to recognize all video objects. This CNN processes 400 images with 16x16 dimensions per-second and it detects one object in each process per image. The results shows that CNN classifies objects more accurately than along the line algorithms [16][22]-[23]-[25].

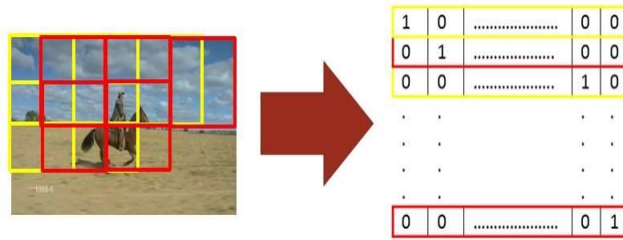
Normally an image consist of different objects with different sizes. Here, we use two methods to find image object; partitioning image into small sizes and image resizing. Partitioning an image to small parts gives us the chance to find an object regardless of its location in the image. Image resizing provides the opportunity to find small and large object and.

An image is partitioned into 16x16 sub-images where each part overlaps with its neighbor with an 8x16 block as it is shown in Figure 5. CNN finds a specific object in a 16 \* 16 sub-image. The final result of CNN for each sub-image is a vector. The vector size is equal to the total number of objects. Each element of the vector corresponds to an object and the existence of this element will be shown by one the vector otherwise by zero (See Fig. 5).

The size of objects are different (some of them are big such as horse and some of them are small such as shoes). For a

better representation, we convert the original image into 4 different sizes in a pyramid form. In this case, big and small objects will be detected at small and large size version of the image respectively. Figure 6 presents an example where an image with original size of 128x128 is

changed into 64x64, 32x32, and 16x16. The tree is a small object and this can be detected in 128x128 size image and the horse is a large object and it can be detected in 16x16 size one.



**Fig. 5. Input and output of CNN.**



**Fig. 6. Changing the size of the original image for CNN.**

In this example, the output of CNN will be 337 vectors for one frame by adding the vectors of different image sizes from 12x128 to 16x16. The output vectors of different size of an image are combined into one vector (See Fig. 7).

As we select three frames of three parts of a video the above process is performed for all of the selected frames and we have a final vector for each of the three frame and finally, the three vectors are converted into one vector by union function. Figure 8 shows this process

*3) Understanding the relationship between objects of an image*

Up to now, the objects and human has been recognized. For example if the input video is lifting, the output of the second stage is a vector that includes the objects of video such as human, shoes, barbells or other heavy weights and etc. In this step, we used support vector machine with RBF kernel.

First the classification is trained with 70 samples. After that each vector of the frame is mapped onto the correct class and SVM determines the label of interactions between human and objects.



## V. CONCLUSION

This paper presented an algorithm for human action recognition that suitable in situations where we have humanobject tnteractions in videos. This method collects UCF Sport Action and Olympic Sport videos have several kind of sports that human use the objects, and then we employ a supervised learning algorithm where the underlying network trained by image-net data to train the system for objects recognition. The employing of learning method, significantly improves the performance in vision tasks. The accuracy improvement comes from the using three frames instead of all video frames, the learning algorithm and the approach for objects recognition. Our approach does not require segmentation, tracking of humans, pruning of motion or static features, or stabilization of videos. The results show that, UCF Sport Action(87.5%) and Olympic Sport(93.1%), the proposed method provides better performance in comparison to other state-of-the-arts results.

## REFERENCES

- E.J. Amirbandi and G. Shamsipour, "Exploring methods and systems for vision based human activity recognition," in 2016 1st Conference on Swarm Intelligence and Evolutionary Computation (CSIEC). pp. 160-164, 2016. IEEE*
- R. Poppe, "A survey on vision-based human action recognition. Image and vision computing," 2010. 28(6): pp. 976-990.*
- J. Aggarwal, and L. Xia, "Human activity recognition from 3d data: A review." Pattern Recognition Letters, 2014. 48: pp. 70-80.*
- S. Ke, R., et al., "A review on video-based human activity recognition. Computers," 2013. 2(2): pp. 88-131.*
- T.B. Moeslund, A. Hilton, and V. Krüger, "A survey of advances in vision-based human motion capture and analysis." Computer vision and image understanding, 2006. 104(2): pp. 90-126.*
- S. Vishwakarma, and A. Agrawal, "A survey on activity recognition and behavior understanding in video surveillance." The Visual Computer, 2013. 29(10): pp. 983-1009.*
- J.M. Chaquet, E.J. Carmona, and A. Fernández-Caballero, "A survey of video datasets for human action and activity recognition." Computer Vision and Image Understanding, 2013. 117(6): pp. 633659.*
- K. Grauman, and B. Leibe, "Visual object recognition." Morgan & Claypool Publishers, Synthesis lectures on artificial intelligence and machine learning 5, no. 2, pp.1-181, Apr. 2011.*
- B.T. Morris, and M.M. Trivedi, "A survey of vision-based trajectory learning and analysis for surveillance." Circuits and Systems for Video Technology, IEEE Transactions on, 2008. 18(8): pp. 11141127.*
- M. Fiaz, and B. Ijaz. "Vision based human activity tracking using artificial neural networks." in Intelligent and Advanced Systems (ICIAS), 2010 International Conference on. 2010. IEEE.*
- 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, 2013. 35(1): pp. 221-231.*
- K. Simonyan, and A. Zisserman. "Two-stream convolutional networks for action recognition in videos." in Advances in Neural Information Processing Systems, pp. 568-576. 2014.*

- L. Wang, Y. Qiao, and X. Tang. "Action recognition with trajectory-pooled deep-convolutional descriptors." in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4305-4314. 2015.
- A. Krizhevsky, I. Sutskever, and G.E. Hinton. "Imagenet classification with deep convolutional neural networks." in *Advances in neural information processing systems*, pp. 1097-1105. 2012.
- J. Schmidhuber, "Deep learning in neural networks: An overview. *Neural Networks*," 61: pp. 85-117. 2015.
- R. Girshick, J. Donahue, T. Darrell, and Malik, J. "Rich feature hierarchies for accurate object detection and semantic segmentation." in *Computer Vision and Pattern Recognition (CVPR)*, pp. 580-587. 2014 IEEE.
- V. Vapnik, "The nature of statistical learning theory (Book style)." 2013: Springer Science & Business Media."
- D. Meyer, and F.T. Wien, "Support vector machines." *The Interface to libsvm in package e1071*, 2015 Aug 5.
- M.P. Brown, P., Grundy, W. N., Lin, D., Cristianini, N., Sugnet, C., Ares, M., and Haussler, D "Support vector machine classification of microarray gene expression data." University of California, Santa Cruz, Technical Report UCSC-CRL-99-09, 1999.
- S. Abe, "Support vector machines for pattern classification (Book style)." Vol. 53. 2005: Springer.
- K. Chatfield, K. Simonyan, A. Vedaldi, and A. Zisserman, "Return of the devil in the details: Delving deep into convolutional nets." *arXiv preprint arXiv:1405.3531*, 2014.
- K. Simonyan, and A. Zisserman, "Very deep convolutional networks for large-scale image recognition." *arXiv preprint arXiv:1409.1556*, 2014.
- C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, and A. Rabinovich, "Going deeper with convolutions." in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1-9. 2015
- P. Sermanet, D. Eigen, X. Zhang, M. Mathieu, R. Fergus, and Y. LeCun, "Overfeat: Integrated recognition, localization and detection using convolutional networks." *arXiv preprint arXiv:1312.6229*, 2013.
- K. Soomro, and A.R. Zamir, "Action recognition in realistic sports videos," in *Computer Vision in Sports*. 2014, Springer. pp. 181-208.
- M. Rodriguez, "Spatio-temporal maximum average correlation height templates in action recognition and video summarization." 2010.
- H. Wang, M. M. Ullah, A. Klaser, I. Laptev, and C. Schmid, "Evaluation of local spatio-temporal features for action recognition." in *BMVC 2009-British Machine Vision Conference*, pp. 124-1. BMVA Press, 2009.
- Q.V. Le, Zou. Will Y, Yeung Serena Y, and Ng. Andrew Y. "Learning hierarchical invariant spatio-temporal features for action recognition with independent subspace analysis." in *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*, pp. 3361-3368. IEEE, 2011.
- A. Kovashka, and K. Grauman. "Learning a hierarchy of discriminative space-time neighborhood features for human action recognition." in *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*. 2010. IEEE.
- M. Ravanbakhsh, H. Mousavi, M. Rastegari, V. Murino, and L. S. Davis. "Action Recognition with Image Based CNN Features." *arXiv preprint arXiv:1512.03980*, 2015.

- J.C. Niebles, C.-W. Chen, and L. Fei-Fei, "Modeling temporal structure of decomposable motion segments for activity classification," in *Computer Vision—ECCV 2010*. 2010, Springer. pp. 392-405.
- Ibrahim, M., S. Muralidharan, Z. Deng, A. Vahdat, and G. Mori, "A Hierarchical Deep Temporal Model for Group Activity Recognition." *arXiv preprint arXiv:1511.06040*, 2015.
- M. Jain, H. Jégou, and P. Bouthemy. "Better exploiting motion for better action recognition." in *Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference*, pp. 2555-2562. 2013.
- W. Li, et al. "Dynamic pooling for complex event recognition." in *Computer Vision (ICCV), 2013 IEEE International Conference*, pp. 2728-2735. 2013.
- L. Wang, Y. Qiao, and X. Tang. "Mining motion atoms and phrases for complex action recognition." in *Computer Vision (ICCV), 2013 IEEE International Conference*, pp. 2680-2687. 2013.
- A. Gaidon, Z. Harchaoui, and C. Schmid, "Activity representation with motion hierarchies." *International journal of computer vision*, 2014. 107(3): pp. 219-238.
- H. Wang, D. Oneata, J. Verbeek, and C. Schmid, "A robust and efficient video representation for action recognition." *arXiv preprint arXiv:1504.05524*, pp.1-20, 2015.
- N. Souly, and M. Shah, "Visual Saliency Detection Using Group Lasso Regularization in Videos of Natural Scenes." *International Journal of Computer Vision*, pp. 1-18, 2015.
- H. Wang, A. Kläser, C. Schmid, and C. L. Liu, "Action recognition by dense trajectories." in *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference*, pp. 3169-3176. IEEE, 2011.
- P. Weinzaepfel, Z. Harchaoui, and C. Schmid. "Learning to track for spatio-temporal action localization." in *Proceedings of the IEEE International Conference on Computer Vision*, pp. 3164-3172. 2015.