



# HUMAN ACTIVITY RECOGNITION IN REAL TIME USING DEEP LEARNING

AZHAGUMEENATCHI.C<sup>1</sup>, DURGA DEVI.R<sup>2</sup>, KAREESHINI.S<sup>3</sup>, SARANYA.B<sup>4</sup>,  
SANGEETHAPRIYA.J<sup>5</sup>

<sup>1</sup>Information Technology, Saranathan College of Engineering, Trichy, India

<sup>2</sup>Assistant Professor, Information Technology Saranathan College of Engineering, Trichy, India

**Abstract:** In today's world, Human Activity Recognition [HAR] plays a critical role in 'human- to-human' interaction. HAR displays and provides the identification of a human as well as the action done by that human, which is tough to recognize. Due to the high processing time, deep learning techniques such as CNN and LSTM cannot be used, instead we will apply transfer learning to recognize human activities. For many computer vision-based applications, such as video surveillance, criminal investigations, and sports applications, human action recognition is one of the difficult issues. Using the similarities between each pair of frames, each extracted sub-unit is further separated into frames that represent action. We will detect the action by comparing the generated HOG to the existing HOGs in the training phase, which represents all the HOGs of many actions using a dataset, utilising the Histogram of the Oriented Gradient (HOG) of the Temporal Difference Map (TDM) of the frames.

## 1. INTRODUCTION

The goal of Human Action Recognition (HAR) is to comprehend a person's action and assign a label to each and every action. It has a wide range of applications, and as a result, it has gotten a lot of attention in the subject of "computer vision." Human activity recognition [HAR] is a function in Human Computer Interaction [HCI] that identifies various actions to aid computerized operations. HAR is also useful in a variety of security surveillance and CCTV applications, as well as sports forecasting.[3]

Human activities/actions refer to all activities done by people in order to achieve their own objectives. In real life, certain human actions are carried out for monetary benefit, while others are carried out for personal satisfaction. Since Real-Time action has high difficulties in extraction of landmarks due to the speed of the movement/action is not definite all the transfer learning algorithms cannot be used only some of them can be used.[5] This system has made sure to compare all the algorithms and find out which algorithm is best suited and has the best accuracy. Combining two algorithms to make an hybrid system might be useful to disregard the background recognition of the grids.

## 2. EXISTING SYSTEM

In the training phase, most extant action recognition systems use pertained weights of various AI architectures for every visual representation of video frames, which has a significant impact on the feature discrepancy that is identified, such as the separation between visual and temporal indications[1]. The major issues that has not been resolved are the background grids being identified as a voxel which might end up being a wrong prediction. As the training increases the accuracy increases so using a large dataset can be useful but data loading is difficult which is also a disadvantage[13].

**Disadvantages:** System's overall performance is poor. The system's processing speed is also slow[8]. This system is difficult to deploy in real- world circumstances.

## 3. PROPOSED SYSTEM

This is all about a technique that has been utilized to predict using 16 different types of action classes. Each class includes a total of 1000 photos to train with. The MideaPIPE program is a deep learning tool for extracting human body landmarks. This project recognizes human actions such as walking, running, sleeping, fighting, and so on. The entire process entails supplying a video/image input, turning it to frames, passing it through modules that detect the activity, and eventually producing a result in the form of a label, which is then transformed into a speech output[15]. This method is vision-based, meaning it recognizes human activity from video or photos. In this scenario, the system uses a camera to collect data and recognize activities. Smartphones that recognize human activities are also a common data source.

The method of analysing human motion using computer and machine vision technology is known as human activity recognition, or HAR. Sensors record human motion, which can be interpreted as activities, gestures, or behaviours. The movement data is subsequently converted into action commands, which are then executed and analysed by computers[11]. The inputs from the raw sensors of the human activity recognition dataset are fed into HAR models, and the output is a forecast of the user's action/activities.

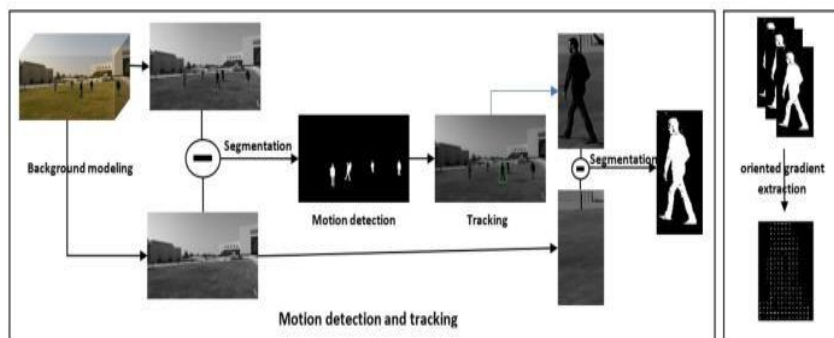
**The step-by-step process of human activity recognition follows the steps below:**

1. At each moment of the video the system identifies the position of the human in space using image triangulation.
2. This spatial position is then matched with an empty 3D grid.
3. Then one voxel is projected onto each of the images.
4. Pixels from the image areas where the human is present is transformed into a voxel[6].
5. The neural network then processes the data from all voxels and the outputs the estimated positions of human's landmark on the original 3D grid.

### MOTION DETECTION

A "background-subtraction-based" method is utilized to start the recognition process and ensure proper detection of the human body moving in the scene. One of the most widely utilized techniques in motion detection is background subtraction. This method's major role is background modeling, which entails removing the input video's unmodified pixels, referred to as voxels and regions. Fig. 1 illustrates how detection is carried out.

The computation of a moving object during the time of the input video is represented by a binary image which is created using a selected threshold.



**Fig.1 Detection of human action**

### DATA PREPARATION

The training video is first pre-processed to select the landmark locations of the human bodies. The Histogram of oriented gradient (HOG) has been used to identify features, and the Temporal Difference Map (TDM) is computed between each and every frame (each pair of successive frames). TDM is used to create a motion history of a person's movements across all grids in the regions.

We use HOG to extract information from the oriented gradient representation about the landmarks of each action. Then, for each action, the Histogram of Oriented Gradient (HOG) is compared, and the relevant set of histograms is collected. For each activity, the same procedure is followed.

### ACTION RECOGNITION

**HOG BASED RECOGNITION** - In the case of human interactions, we train with images from the UT-Interaction and INRIA XMAS (IXMAS) datasets. The histogram of the oriented gradient of the temporal difference map between each pair of photos is produced or computed after the detection of human body and extraction of landmarks for each detected individual. The histogram of each observed action is shown in Fig. II

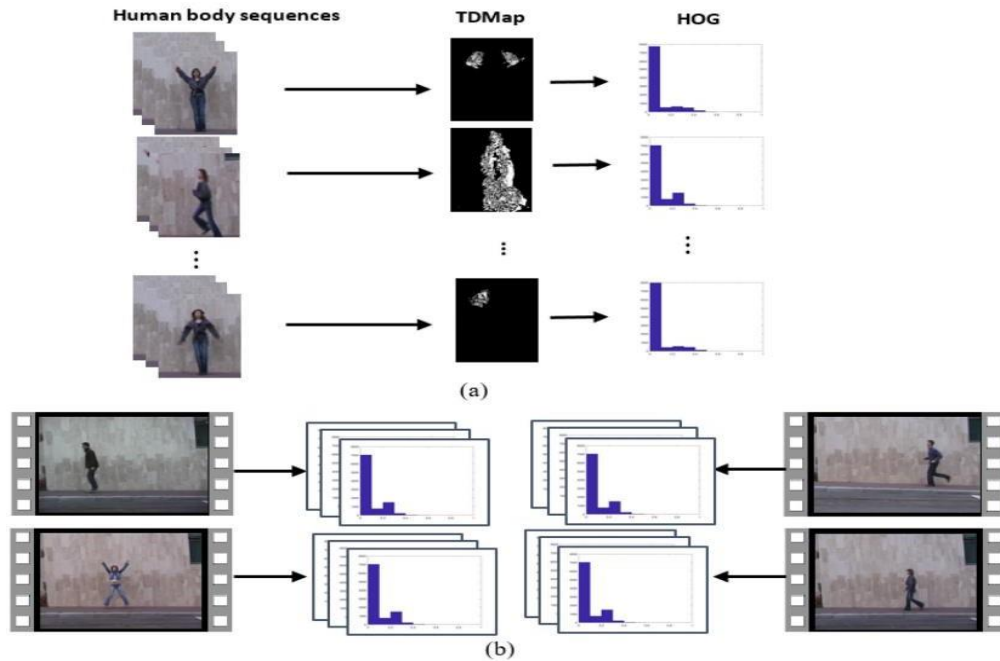


Fig.II HOGs of each action

**SUMMARIZATION OF ACTION**

After completing the recognition of action, the summary is finished by documenting each action taken by a person. The procedure of differentiating the frames of each action from the sequence of actions performed by a person is known as detection. As a result, a sequence from each frame represents a person's summarized actions. Selecting a frame to represent each action can also be used to conduct the summarization. A summary is a collection of action labels from a video in which a person does something specific. The action is summarized in Fig. III

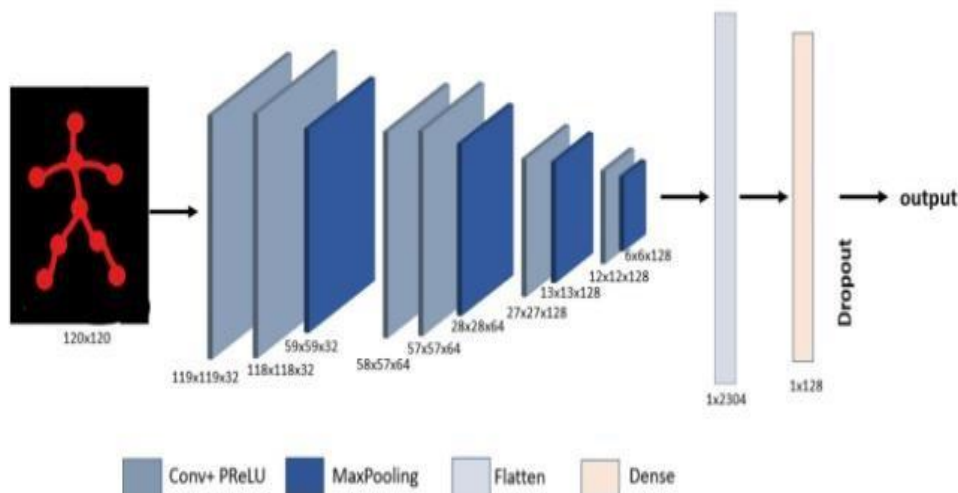


Fig.III Summarization of action

**4. IMPLEMENTATION**

A dataset must be loaded into the system for this process to take place. The training is split into 16 different classes, each with 1000 images. Once the input image/video has been loaded into the Home page, it is transformed into frames and preprocessed to match the training dataset in order to detect the action. When a highly similar activity is identified, it is presented as a main output in the form of a label, which is then transformed into a vocal output.



5. EXPERIMENTAL RESULTS

After the continuous training of epochs, the output is obtained it is attached. Once the training is completed we can start testing the system by providing an input like an image or a video.

```

Epoch 1/100
40/40 [=====] - 109s 2s/step - loss: 3.8389 - accuracy: 0.0875
Epoch 2/100
40/40 [=====] - 92s 2s/step - loss: 3.3989 - accuracy: 0.0891
Epoch 3/100
40/40 [=====] - 98s 2s/step - loss: 3.2572 - accuracy: 0.0961
Epoch 4/100
40/40 [=====] - 94s 2s/step - loss: 3.0922 - accuracy: 0.1250
Epoch 5/100
40/40 [=====] - 105s 3s/step - loss: 3.0154 - accuracy: 0.1281
Epoch 6/100
40/40 [=====] - 98s 2s/step - loss: 3.0091 - accuracy: 0.1289
Epoch 7/100
40/40 [=====] - 98s 2s/step - loss: 2.8683 - accuracy: 0.1594
Epoch 8/100
40/40 [=====] - 112s 3s/step - loss: 2.8190 - accuracy: 0.1398
Epoch 9/100
40/40 [=====] - 103s 3s/step - loss: 2.7138 - accuracy: 0.1773
Epoch 10/100
40/40 [=====] - 93s 2s/step - loss: 2.6968 - accuracy: 0.1648
Epoch 11/100
40/40 [=====] - 98s 2s/step - loss: 2.6389 - accuracy: 0.1844
Epoch 12/100
40/40 [=====] - 92s 2s/step - loss: 2.5535 - accuracy: 0.2078
Epoch 13/100
...
40/40 [=====] - 81s 2s/step - loss: 0.0816 - accuracy: 0.9758
    
```

Fig.IV Training

We use Google Collaborative for training the model. We had trained upto 140 epoches with an accuracy of 96%. Fig.IV shows the training of dataset and the epochs during training session.

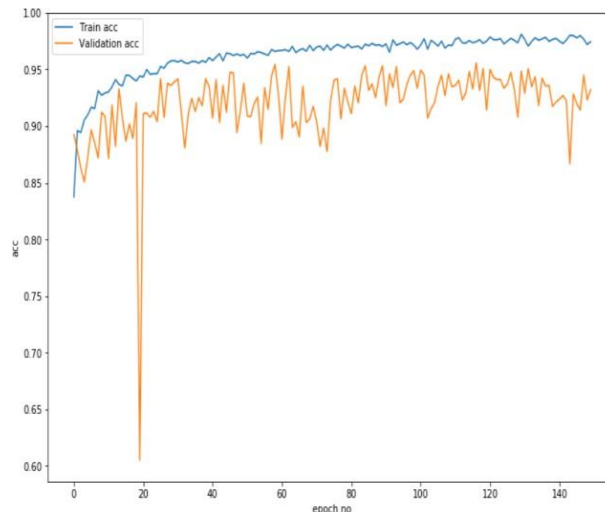
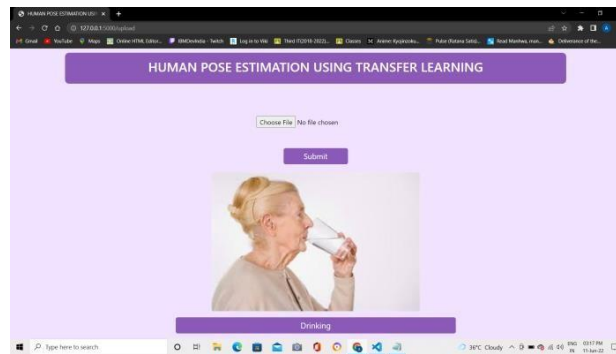


Fig.V Accuracy of Detection

Fig.V shows the accuracy of the detection.

To identify the action, the given input image/video is successively matched with the dataset. The TTS function converts the identified action from label (text) to voice (speech).

In Machine Learning (ML), the TTS [Text To Speech] function is typically used to convert text input into speech format. The primary output is the action's label, and the final output is an audio of the activity. The label, which is the detected output of the action done, is shown in Fig.VI.



**Fig.VI Label of the action detected**

## 6. CONCLUSION

This model uses Random Forest and Grid Search together which has the highest accuracy rate while comparing with a lot of other algorithms like decision tree, k-nearest and 4 more. First and foremost, the confusion matrix for all the algorithms was found then all the algorithms were compared and the best one was found which is being used. Here a lot of graphs were also used to identify features that are suitable for the action detection. This model makes sure of having the highest accuracy by plotting accuracy loss and validation graphs so that it will be more useful to know the accuracy and loss during training period. Since the training has been done with 1000 images for each class, which is 16 in total, which in turn increases the accuracy as the training period increases. The image is loaded as an Input then it is preprocessed into frames and grid to extract landmarks/features. Once extraction is done the Landmark are matched with the trained dataset to identify the highly possible action that resembles/matches the action performed. Once matched the corresponding label is given as a text output which is further converted into voice output. Human Action Recognition is one of the highly interested field which can be incorporated with any computer-vision based problems/systems. This system is to be enhanced further with more features and more accurate models.

## REFERENCES

- [1] A.D. Antar, M. Ahmed, M.A.R. Ahad, Challenges in sensor-based human activity recognition and a comparative analysis of benchmark datasets: A review, in: 2019 Joint 8th International Conference on Informatics, Electronics & Vision (ICIEV) and 2019 3rd International Conference on Imaging, Vision & Pattern Recognition, IcIVPR, IEEE, 2019.
- [2] B. Ali, et al., A volunteer supported fog computing environment for delay-sensitive IoT applications, IEEE Internet Things J. (2020).
- [3] J.K. Aggarwal, M.S. Ryoo, Human activity analysis: A review, ACM Comput. Surv. 43 (3) (2011) 1–43.
- [4] M. Alazab, et al., Intelligent mobile malware detection using permission requests and api calls, Future Gener. Comput. Syst. 107 (2020) 509–521.
- [5] M. Baccouche, et al., Sequential deep learning for human action recognition, in: International Workshop on Human Behavior Understanding, Springer, 2011.
- [6] K.A. da Costa, et al., Internet of things: A survey on machine learning-based intrusion detection approaches, Comput. Netw. 151 (2019) 147–157.
- [7] C. Dai, et al., Human behavior deep recognition architecture for smart city applications in the 5G environment, IEEE Netw. 33 (5) (2019) 206–211.
- [8] C. Dai, X. Liu, J. Lai, Human action recognition using two-stream attention based LSTM networks, Appl. Soft Comput. 86 (2020) 105820.
- [9] M. Elhoseny, et al., A hybrid model of internet of things and cloud computing to manage big data in health system
- [10] R. Girdhar, et al., Actionvlad: Learning spatio-temporal aggregation for action classification. in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2017.
- [11] R. Hou, C. Chen, M. Shah, An end-to-end 3d convolutional neural network for action detection and segmentation in videos, 2017, arXiv preprint arXiv:1712.01111.
- [12] Y.-L. Hsueh, W.-N. Lie, G.-Y. Guo, Human behavior recognition from multiview videos, Inform. Sci. (2020).
- [13] H. Kwon, et al., First person action recognition via two-stream convnet with long-term fusion pooling, Pattern Recognit. Lett. 112 (2018) 161–167.



- [14] R. Khemchandani, S. Sharma, Robust least squares twin support vector machine for human activity recognition, *Appl. Soft Comput.* 47 (2016) 33–46.
- [15] A. Keshavarzian, S. Sharifian, S. Seyedin, Modified deep residual network architecture deployed on serverless framework of IoT platform based on human activity recognition application, *Future Gener. Comput. Syst.* 101 (2019) 14–28.
- [16] Y. Li, et al., Spatiotemporal interest point detector exploiting appearance and motion- variation information, *J. Electron. Imaging* 28 (3) (2019) 033002.