



A Novel Approach for Human Action Recognition from Silhouette Images

Satyabrata Maity¹, Debotosh Bhattacharjee² and Amlan Chakrabarti¹

1. A.K.Choudhury School of Information Technology, University of Calcutta, 92 A. P. C. Road, Kolkata:700 009, India.

2. Dept of Computer Science and Engineering, Jadavpur University, 188, Raja S. C. Mallick Road, Kolkata, West Bengal 700032.

Abstract

In this paper, a novel human action recognition technique from video is presented. Any action of human is a combination of several micro action sequences performed by one or more body parts of the human. The proposed approach uses spatio-temporal body parts movement (STBPM) features extracted from foreground silhouette of the human objects. The newly proposed STBPM feature estimates the movements of different body parts for any given time segment to classify actions. We also proposed a rule based logic named rule action classifier (RAC), which uses a series of condition action rules based on prior knowledge and hence does not required training to classify any action. Since we don't require training to classify actions, the proposed approach is view independent. The experimental results on publicly available Wizeman and MuHVAi datasets are compared with that of the related research work in terms of accuracy in the human action detection, and proposed technique outperforms the others.

Keywords: Video analysis, Action units, Action recognition, Spatio-temporal body parts movement (STBPM), Rule-Action classifier (RAC).

1. Introduction

The video based applications, which was once restricted for surveillance and entertainment purpose is now extended in every direction like education, health care, social networking etc. This increasing demand directs a large number of researchers towards the domain of video analysis. Human action recognition (HAR) is the thrust to many applications like visual surveillance, video search and retrieval, human computer interaction and many more.

Visual action recognition consists of different sub-topics such as gesture recognition developing human-computer interfaces, facial expression recognition and abnormal activity detection for video surveillance. Conversely, full-body actions generally embrace a number of motions and require an integrated approach for recognition, encompassing facial actions, hand actions and feet actions.

Three main complexity issues, as mentioned in [2], are generally present in any HAR technique viz. i) Environmental complexity: The process of HAR depends on the quality of the video, which differs due to the environmental condition of the scene elements and makes the procedure more complex. This type of complexity includes occlusions, clutter, interaction among multiple objects, changing of illuminations etc. ii) Acquisition complexity: Besides environmental condition, the quality of the video also depends on video acquisition, which varies with respect to view point, movement of the camera etc. iii) Human action complexity: In general sense, human actions are of varied in nature, hence exact determination of human action is a complexed task. Presence of multiple human entities makes any HAR technique more complex. So defining a model to recognize various kinds of human action is extremely difficult. To handle those three cases mentioned above, some constrains have been made in our proposed work viz. i) We have used the videos, which contain human silhouettes only and we didn't consider any silhouette extraction techniques

also. video. So, the environmental and acquisition complicacies are not applicable for the proposed technique. ii) To reduce the complexity, the proposed work considers videos containing only one human object in each of the frame. iii) We also consider that the head should be in the upper portion of the body and the body should not be upside down.

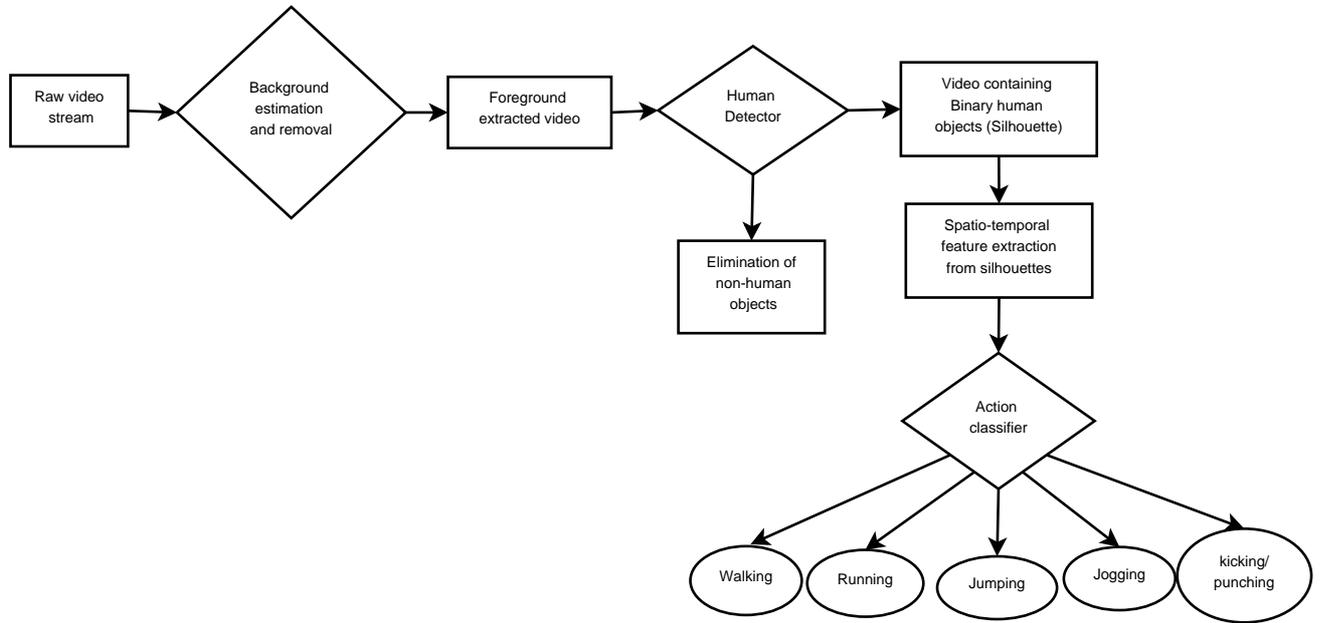


Figure 1. Flow diagram of the general procedure

Fig. 1 illustrates the major components of a generic action recognition system using human silhouettes and their typical arrangement. The general frame-work for human action using silhouette analysis of consecutive frames of a video is mainly divided into three broad sub-steps i) Foreground extraction: Foreground is extracted by eliminating background of the video. This step helps to reduce the searching area of the current frame. ii) Foreground classification: This step determines whether the foreground area contains human or not, as it is unnecessary to analyze nonhuman objects. iii) Feature extraction and action classification: Analysis of the movement of human body parts is performed for the consecutive frames to determine the human action in the successive frames. In our proposed work, we extract spatio-temporal features from human silhouettes and then use the features for action classification.

The actions performed by human are repetitive in nature and those actions are understandable by observing the movement of the body parts. The activity of these body parts determines action, so these body parts are termed as action units (AUs). In the proposed work, we have developed a rule based logic defined over different spatio-temporal values of AUs for determining numerous actions. The proposed approach selects four different parts of the body as AU viz. a) head, b) hand, c) leg and d) remaining body, since these four parts may be visible in a human silhouette. We design the algorithms to extract the location of these AUs in the frame containing human and then track them in the consecutive frames. We have introduced a new feature vector called Spatio-Temporal Body Parts Movement (STBPM), which contains location of a) HEAD, b) LEG(s), c) HAND(s) and d) Body Mass centre (BMC) and the derived features i.e. Head Angle (HA), Stride Angle (SA), Moving Direction (MD), Body Ratio (BR), and Bounding Box (BB). The dimension of the resulting feature vector is $k \times n$, where k is the number of frames to recognize the action, so we can call it as action interval. n is the number of features and the details of k and n will be discussed in the proposed method section. Then Rule Action Classifier (RAC) is introduced to recognize the action. RAC is a classifier, which uses a series of condition action rules to categorize the action instead of any learning mechanism normally used by training based classifiers.

The main contributions of the proposed research work can be summarized as follows :

- Automatic localization of AUs viz. body, head, hand(s) and leg(s) automatically and the average accuracy rate is about 97%.
- Extraction of newly introduced low-dimensional STBPM features for human action recognition.
- The proposed work is a first of the kind approach, which uses a rule based logic set to determine different human actions.
- The proposed method is view invariant except the top view, it can work in any viewpoint of the camera which makes our methodology more robust.

The paper is organized with seven sections starting with introduction section. Section-2 is related work, which contains a comprehensive study of the related research. The description of the proposed methodology is included in section-3. Section-4 describes our results in details and analyses the efficiency of our results with compare to the results of other techniques of related research. Section-5, section-6 and section-7 include conclusion, acknowledgment and references respectively.

2. Related work

In the domain of computer vision, several communities are working to solve the difficulties of HAR in numerous ways. An outstanding survey of different approaches for HAR are presented in [1]. The HAR tasks can be roughly classified into five different categories which are based on shape models, motion models, geometric human body models, interest-point models, and dynamic models [2]. The shape based models generally exploit the silhouette of moving objects commonly extracted by eliminating the backgrounds of the video frame. The silhouettes of the human object in consecutive frames are changing with respect to time as the movement is due to the displacement of one or more body parts. Some of the methods that apply shape based models are described in [2], [3], [4], [5], [9], [14], [15], [16], [17]. Generally these methods use spatio-temporal features and in an ideal case, these approaches are invariant to luminance, color, and texture of the moving objects (and background) as mentioned in [2]. All the approaches are heavily dependent on proper silhouette extraction of human objects, but accurate silhouette extraction in spite of color, texture and luminance is the main challenges in general.

The second category described in [18], [19], [20], [21] and [22] generally apply flow based (mainly optical flow) techniques to model the motions of the moving objects in the consecutive frames of the video. A series of motion features of a moving object from the consecutive frames of a video provides the spatio-temporal model of any action. The flow-based method computes the optical field between adjacent frames and uses it as a feature for action recognition. This is suitable for recognizing small objects. However, it is computationally expensive and generate coarse features only. Thus different actions may exhibit similar flows over short periods of time.

The geometric human body models [3], [23], [24], [25], [26], attempt to characterize the geometric model of the human body to classify different actions as the geometry of the human body change according to the action he/she performs. This type of approaches construct feature vector, which contains static and dynamic geometric locations of different body parameters responsible to complete the action. Subsequently, a classifier is used to classify the extracted feature. This type of method can be very efficient if we correctly extract the geometric body parameters.

Point of interest models are described in [27], [28], [29], [30] and [31]. The points of interest are the distinguishable area of the image having eminent information. This type of approaches employ the points of interest collected from consecutive frames of the action interval and then process them to construct the feature vector. The classification is taken place exploiting this feature vector. [29] and [27] use 1D Gabor and Gaussian filter respectively to extract point of interest. On the other hand [28] use Harris corner detector to extract salient points, which hold significant local variations in both spatial and temporal direction.

The dynamic models [32], [33], classify human action based on the dynamic variation among the consecutive frames in temporal direction. This type of approaches characterize the static posture of any action as state and then describing the dynamics using this state-space. Thus an action is a set of several states, the state space of which are connected using dynamic probability model (DPN). Hidden Markov Model (HMM) is the most commonly used DPN, which has the capability of directly modeling time variation of the data features.

Our proposed model combine the advantages of shape based model and geometric model to construct the feature vector. In our technique, we use silhouettes to extract spatio-temporal features, which contains the coordinates of different body-parts and some other derived features. Moreover, we designed a rule based logic to classify the human action, which is not required training.

3. Proposed work

To be very precise the proposed work extracts STBPM features from human silhouettes of the consecutive frames of a video and then apply RAC to determine the action. Every shades of color add a dimension in the visual image information which is required to express an image more vividly but these shades of color do not have greater impact in HAR. On the other hand silhouette of a human object has only two levels and it provides us prominent shape based information. More specifically, silhouette of a particular frame of any human action presents the instance of the human shape of that time. Hence, the proposed approach exploits the silhouette towards action recognition, which is divided into two important steps A)Silhouette analysis and feature extraction, and B)Action classification. The block-diagram of the proposed work is shown in Fig. 2. First we analyze silhouettes of the human object in the first frame and compute all spatial feature values and store them in a row vector. The spatial features include the AUs location along with other derived features, which is discussed in the later of this section. The we process all the row vector of the action interval and compute the spatio-temporal features. The rule-action classifier to determine the action of the human in that action interval using spatio-temporal features.

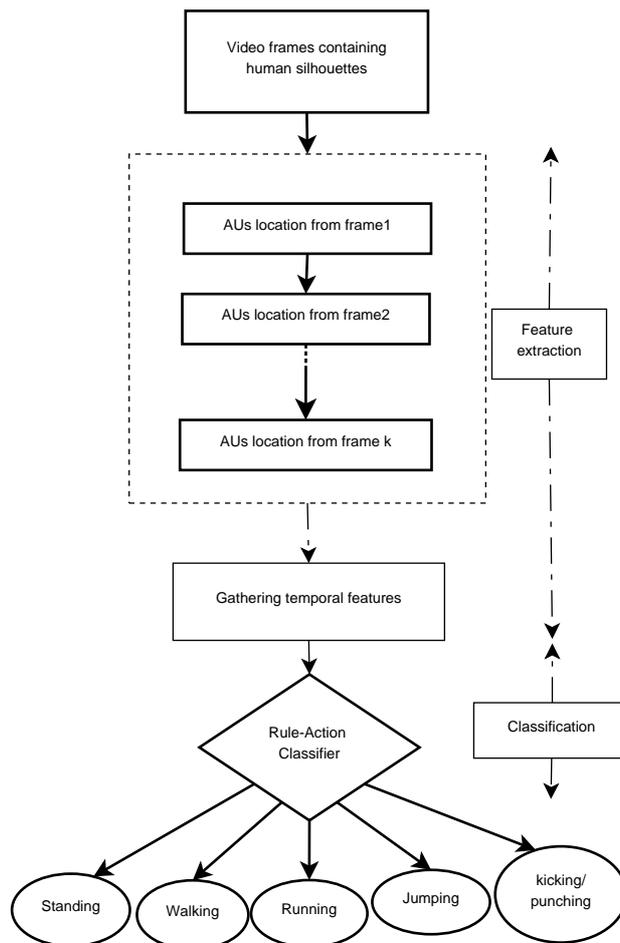


Figure 2. Flow diagram of the proposed work

Different AUs are involved in the various actions performed by human. The three most detectable AUs from human silhouettes, which are taking part of most of the actions, are head, hand(s) and leg(s). Fig. 3 shows an example which depicts these three AUs in case of walking. Every human body in the silhouette is divided with three horizontal and two vertical levels. The change in location or movement of AUs, in case of walking, can determine using these levels and the rules are proposed considering those movements. The following part of this section describes the process of feature extractions and action classification mechanism.

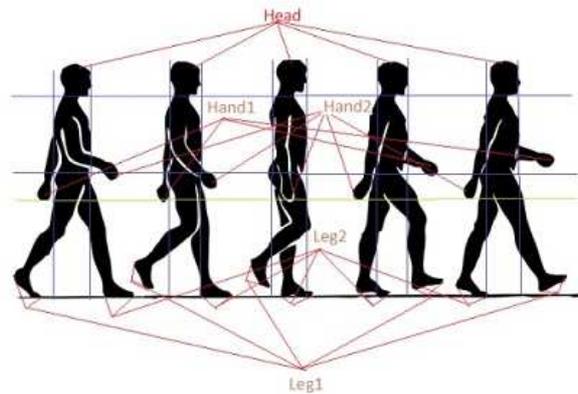


Figure 3. Action unit of human body parts

3.1. Silhouette analysis and feature extraction

In the proposed work, we define actions in terms of AUs and their movement, which is needed to perform the action. The movement of AUs are either i) local movement (LM): movement with the BBOX, ii) global movement (GM): movement with in the frame. An idea of local and global position of AU are shown in Fig. 4, Where the point P has two coordinates i) P_local to represent the local coordinate with respect to BBOX, and ii) P_global to represent the global coordinate with respect to frame. This hierarchy of human action as shown in Fig. 5 is characterized in rule base logic for HAR.

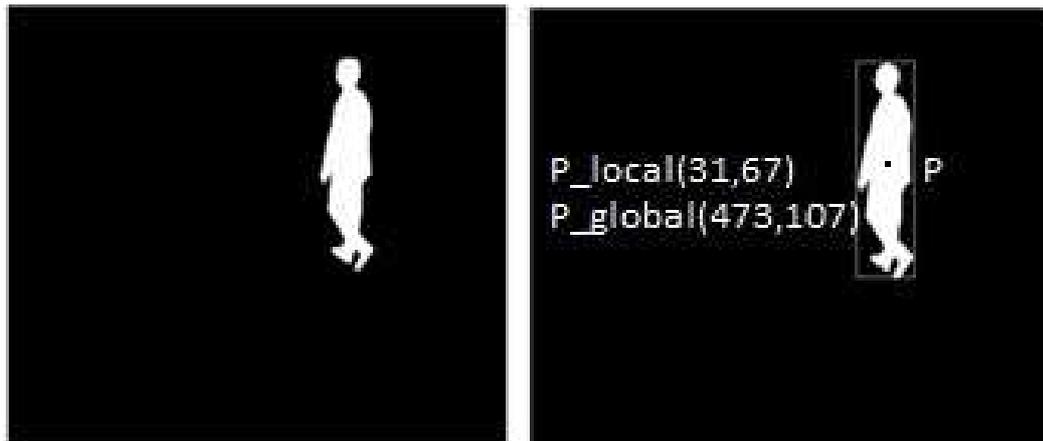


Figure 4. Algorithm of BMC computation

The objective of this step is generation of STBPM feature vector, which holds different state of AUs. The feature vector contains location and other relevant information from the silhouettes of the k consecutive frames of any video,

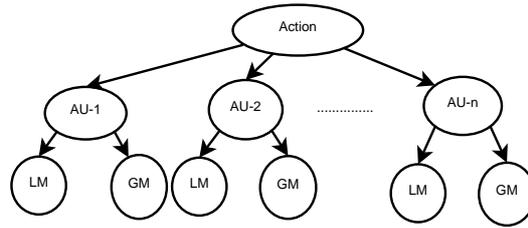


Figure 5. Hierarchy of human an action

which are considered for taking the decision of any human action. In our experiments, we have set the action interval $k =$ number of frames/sec i.e. the frame rate of the video. The STBPM feature vector is formed with 28 different values viz.

- Bounding box (*BBOX*) with four values (x_{st}, y_{st}, dx and dy).
- Body mass center (*BMC*) with two values (x_{cnt}, y_{cnt}).
- Global body mass center (*BMC_G*) with two values (x_{gcnt}, y_{gcnt}).
- Body Ratio (*BR*) with one value.
- Extreme left and extreme right sided locations of head with four values ($(HD1_x, HD1_y)$ & $(HD2_x, HD2_y)$).
- Angle of head (*HA*) with one value measured in form of degree, which is discuss in the later part of this section.
- *HAND1* and *HAND2* with four values ($(xh1, yh1)$ & $(xh2, yh2)$).
- *HEEL1, TOE1, HEEL2 & TOE2* with eight values ($(xh1, yh1), (xt1, yt1), (xh2, yh2)$ & $(xt2, yt2)$).
- Moving direction (*MD*) with one value.

These 28 values extracted each from foreground silhouettes of k consecutive frames make the size of the feature vector equal to $k \times 28$ for each action interval. Thus two levels of features are extracted for classification 1) Spatial features: The Location of AUs and other derived feature extracted from each human silhouette. 2) Generation of spatio-temporal features: The collection of spatial features from action interval of consecutive frames and then reduce the dimensionality of the huge number of features. In the feature vector, all the locations are represented using 2-D spatial co-ordinate in the form of (x, y) .

3.1.1. Spatial features

The spatial features include the spatial location of AUs and their derived features. So, we divided the process of spatial feature extraction into four phases viz.

- Extraction of BBOX, BMC, BMC_G & BR.
- Extraction of head locations and HA.
- Extraction of hand(s) location(s).
- Extraction of leg(s) location(s) and SA.

Before describing the details of extraction of spatial features in details, we explain four vectors viz. X_1, X_2, Y_1 and Y_2 , which are used to extract the location of AUs and other derived features. Silhouette image has only two levels viz. '0' and '1'. So there is no information in the inner region of the silhouette. Thus the analysis is done on the boundary region of the silhouette. The proposed methodology extract four extreme bounding vectors, two of them along the row viz. X_1 and X_2 and other two are along the column viz. Y_1, Y_2 . We named the 2-D matrix containing the silhouette of a human object as SIL.

X_1 vector characterizes the extreme left-sided boundary curve of the human silhouette. The left-sided curve contains the first nonzero point of each row of the human silhouette as computed in Eq-1.

$$X_1(i) = j, IF \quad SIL(i, j) = 1, \& \quad \forall \quad SIL(i, l) = 0. \quad (1)$$

Where, $l < j$.

X_2 vector characterizes the extreme right-sided curve of the human silhouette. The right-sided curve contains the last nonzero point of each row of the human silhouette as shown in Eq-2.

$$X_2(i) = j, IF \quad SIL(i, j) = 1, \& \quad \forall \quad SIL(i, l) = 0. \quad (2)$$

Where, $l > j$.

Y_1 vector characterizes the peak curve of the human silhouette. The peak curve contains the starting nonzero point of each column of the human silhouette as computed in Eq-3.

$$Y_1(i) = j, IF \quad SIL(j, i) = 1, \& \quad \forall \quad SIL(l, i) = 0. \quad (3)$$

Where, $l < j$.

Y_2 vector characterizes the base curve of the human silhouette. The base curve contains the ending nonzero point of each column of the human silhouette as computed in Eq-4.

$$Y_2(i) = j, IF \quad SIL(j, i) = 1, \& \quad \forall \quad SIL(l, i) = 0. \quad (4)$$

Where, $l > j$.

- Extraction of BBOX, BMC, BMC_G & BR features values: *BBOX* is represented by (x_{st}, y_{st}, dx, dy) where (x_l, y_l) is the left top point and dx & dy are the width and height of the BBOX and all the values are computed using four vectors as described in Eq.5 to Eq.12. Now we have two windows viz. BBOX and the image frame. The co-ordinate values with respect to BBOX are the local values and the same with respect to image frame are the global values. The conversion of any local co-ordinate to its corresponding global value can be easily done by adding x_{st} and y_{st} values with the corresponding x and y values of the co-ordinate and the opposite mechanism is done to convert the global co-ordinate to its local one. This conversion is needed to estimate global movements of any AU.

$$x_{min} = MIN(X_1). \quad (5)$$

$$y_{min} = MIN(Y_1). \quad (6)$$

$$x_{max} = MAX(X_2). \quad (7)$$

$$y_{max} = MAX(Y_2). \quad (8)$$

$$x_{st} = x_{min}. \quad (9)$$

$$y_{st} = y_{min}. \quad (10)$$

$$dx = x_{max} - x_{min}. \quad (11)$$

$$dy = y_{max} - y_{min}. \quad (12)$$

Body mass center (*BMC*) is a key parameter for action detection. The location of *BMC* in consecutive frames determines several local movements. *BMC* is computed using weight of rows and columns. The weight of each row or column is the total number of nonzero values along the respective rows or columns. $(x_{cnt}$ and $y_{cnt})$ are the weighted average of all rows and all columns of the BBOX respectively. The proposed work use Eq-13 to Eq-16 to extract *BMC*. Fig. 6 shows the *BMC* and the contour of corresponding human silhouette.

$$ROW_{wt}(i) = \sum_{j=1}^{col} SIL(i, j) \quad (13)$$

$$COL_{wt}(i) = \sum_{j=1}^{row} SIL(j, i) \quad (14)$$

$$x_{cnt} = \left(\sum_{i=1}^{row} ROW_{wt}(i) \times i \right) / \sum_{i=1}^{row} row_{wt}(i) \quad (15)$$

$$y_{cnt} = \left(\sum_{i=1}^{col} COL_{wt}(i) \times i \right) / \sum_{i=1}^{col} COL_{wt}(i) \quad (16)$$

The global body mass center(BMC_G) is represented by (xg_{cnt}, yg_{cnt}) is computed using Eq.13 and Eq.14.

$$xg_{cnt} = x_{cnt} + x_{st}. \quad (17)$$

$$yg_{cnt} = y_{cnt} + y_{st}. \quad (18)$$

Most of the feature values are computed in the proposed methodology in their respective local windows i.e BBOX as we need to tract the movement of any AU. The movement of human objects should reflect in the BMC values of consecutive frames, but it is not always reflected by its local coordinate values. In some cases, the values of BMC remain nearly unchanged inspite of huge movement in that action interval. For instance, the location value of BMC and HEAD change a little with respect to BBOX in case of "running". The value of BMC_G has a big role in such cases as the values of BMC_G of the corresponding frames are changed appreciably if there is any global movement.

Human body ratio is the ratio of the height and the width of the human body as computed in Eq.15.

$$BR = dx/dy \quad (19)$$

The anatomical studies done in [36] proposed the vertical position of head and hip are $0.87H$ and $0.53H$ respectively, where H is the height(from head to toe) of the human body. The ending point of hand can go little beyond hip if we simply keep our hand down and straight and hence, we consider this as $0.47H$. We divide the BBOX in smaller regions using three horizontal and two vertical lines. i) Upper level (UL): The line, which separates upper region of the body from in the BBOX. In Fig. 7, the horizontal solid line in the top area of the BBOX is the UL and the region above this line is upper region. Upper region contains head and hand(s) (when someone uplift his/her hand) and the line is at $0.87H$. ii) Middle level (ML): This is a line, which defines the middle region in the BBOX. The middle region, which contains body and hands, is started from UL and end at ML. ML is the dashed horizontal line at $0.47H$ as shown in Fig. 7. iii) Lower level (LL): Human legs are included in the region below LL and the line is at $0.53H$ as shown in Fig. 7 lower level in dotted line. Some times hand(s) and leg(s) share a common portion in BBOX, a common area between lower and middle regions is there for that reason. iv) Left side level (LSL): This is a vertical line with a distance of $v1$ from the BMC. $v1$ is the average of all values of $X1$ with in ML. v) Right sided level (RSL): This is computed in same manner with LSL except taking the values from $X2$. In Fig. 7, left sided vertical line is LSL and right sided vertical line is RSL.

- Extraction of head locations and HA: Head, located above the UL, is the highest part of the body unless one raises hand(s) in most of the cases of common human action like walking, running; jumping, jogging etc. On the other hand the vector $Y1$ contains the starting points from the top of BBOX. Head, hand or other body parts have certain shape. So the $Y1$ curve contains several numbers of curvatures, which differ depending upon different shape of the different body parts.

We have from Eq.20, Y'_1 contains the column having $Y1$ value greater than UL mark. From have from Eq.21, the component array $COMP$ consists the length of the components present in the upper level. The highest length component is the HEAD and the most left sided point of the HEAD is starting point H_{st} and the most right sided point of the HEAD ending point H_{en} .

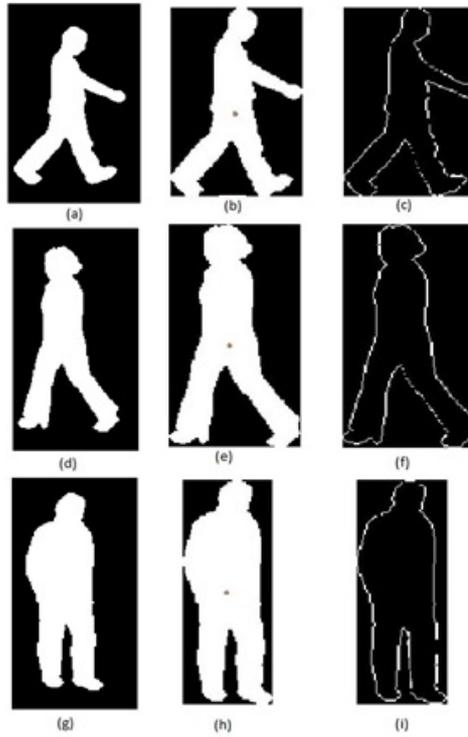


Figure 6. Flow diagram of the proposed method

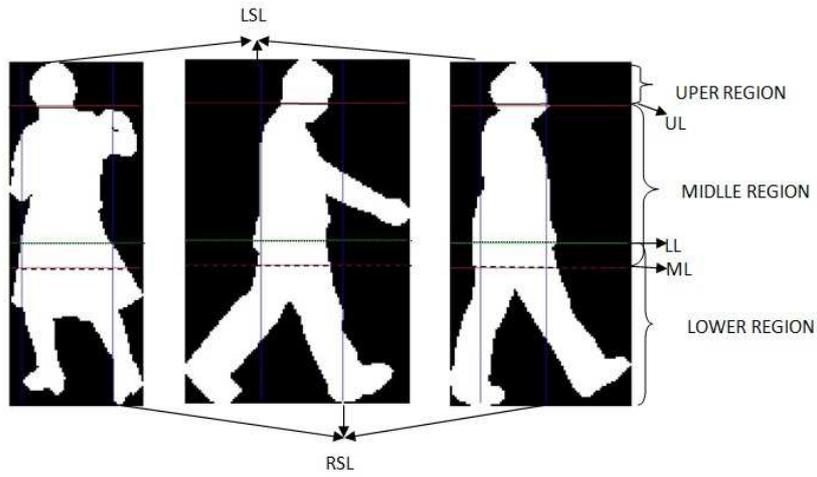


Figure 7. Three logical horizontal level of human silhouette

$$Y'_1(i) = \begin{cases} 1 & \text{if } Y'_1(i) > UL; \\ 0 & \text{Otherwise.} \end{cases} \quad (20)$$

$$COMP(i) = \sum_{j=k1}^{k2} Y'_1(j) \text{ where } \forall Y'_1(j) > 0 \& k1 \leq j \leq k2 \quad (21)$$

Fig. 8 shows some results after applying HEAD extraction formula using Eq.20 and Eq.21. In the figure two end points H_{st} and H_{en} are connected with the BMC of the corresponding human silhouette.

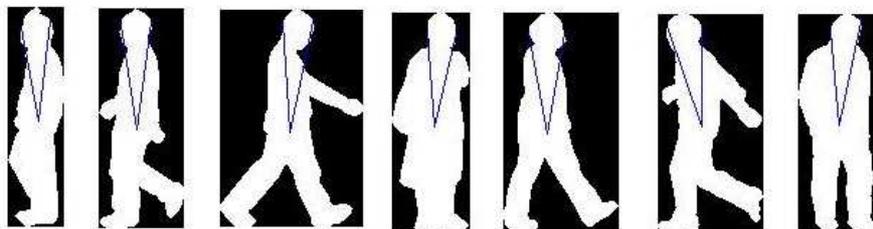


Figure 8. Locating head from BMC

HA is the angle between head and the shoulder. HA is required to determine the body moving direction of the human as the head is angled towards the moving direction in most of the cases. Head angle (HA) is measured by using the middle point (H_{mid}) of the head and the middle point of the shoulder (SH_{mid}). The location of shoulder is computed by standard anatomical height of $0.818 \times H$, where H is the height of the human body (from lowest heel to top point of the head).

$$H_{mid}(x_1, y_1) = \begin{cases} x_1 = (X_1(H_{st}) + X_2(H_{mid}))/2 \\ y_1 = (H_{st} + H_{en})/2. \end{cases} \quad (22)$$

$$SH_{mid}(x_2, y_2) = \begin{cases} x_2 = (X_1(H \times 0.818) + X_2(H \times 0.818))/2 \\ y_2 = (H \times 0.818). \end{cases} \quad (23)$$

$$HA = \tan^{-1} [(y_2 - y_1)/(x_2 - x_1)] \quad (24)$$

- Hand localization: Unlike head, hand(s) can be in two regions above UL and ML. Sometimes it is occluded by different body area. It can be tracked when it is appearing outside the body area. Two vectors X1 and X2, which hold the left and the right sided points of the silhouette must contain the information about the hand.

$$std_{x1} = STD(X'_1(i)) \text{ where } X'_1 = X_1(LL : UL) \quad (25)$$

$$X''_1(i) = X_1(i) \forall abs(X_1(i) - BMC(x)) \geq std_{x1} \quad (26)$$



Figure 9. Locating hand(s) from BMC.

The area beyond two vertical lines contain the information of hand. Thus we analyze those area stored in x''_1 of Eq. 21. If there is any area beyond hand determine threshold Th_{md} and sustain up to $0.03H$, then there is a hand. This approach is applied on x_2 to get second hand if visible.

- Heel and toe localization and stride angle computation: This step includes most important feature of human action recognition as most of the active action are depending on the different movement of legs. The action area is performing below the LL line of the BBOX. Y_2 vector is exploited to find the position of leg and corresponding stride angle (where two legs are detected). The proposed method extract the position of toe and heel of the leg(s) as it is important to understand the direction of human motion. In natural view the human motion is along the direction of the toe and the heel which holds the body weight remaining the back side. Eq-5 is used to extract the bottom point(s) of Y_2 curve. bs_pnt is used instead of pk_pnt and loc are two vectors contain bottom value(s) and corresponding location of the curvatures. One more function leg_find analyzes the shape of the curvature and return the toe and hill vectors.

$$Th_l = mean(Y'_2(i)) \text{ where } Y'_2(i) \in \{x | \forall Y_2(i) > LL\} \quad (27)$$

$$Y'_2(i) = Y_2(i) \forall abs(Y_2(i) - BMC(y)) \geq Th_l \quad (28)$$

If only one leg found, it is assumed that stride angle is zero as one leg is occluded by others; otherwise the angle between two legs is measured and stored in SA . Fig. 10 shows some results of leg finder where the heel(s) and toe(s) are connected with the BMC.

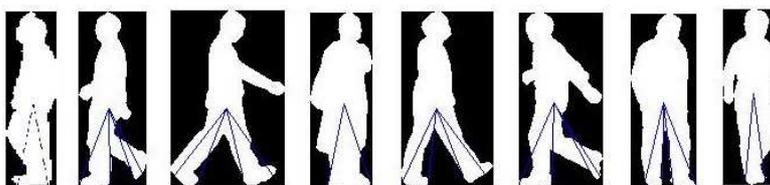


Figure 10. Locating leg(s) from BMC.

3.1.2. Generation of spatio-temporal features

Extraction of moving direction (MD): MD is temporal feature, which is extracted by exploiting the extracted values of BMC_G. BMC_G is of two types i) Static: there is no such change in BMC_G, ii) Dynamic: The value of BMC_G is changing with respect to time. In case of static BMC_G, MD is determined using HA and heel-toe locations of the last frame. If the direction is close to 90° degree refers to non-moving head; otherwise i) Angle greater than 95° refers the movement is towards right. or ii) Angle less than 85° refers to left sided movement. The locations of heel(s) and toe(s) also have contribution to determine the MD. If $TOE(s) > HEEL(s)$ with respect to Y values, then the direction is towards right and vice-versa. Undefined refers to the fact when leg is not correctly defined or HEEL(s)-TOE(s) are not correctly extracted i.e. smaller or bigger than normal or one HEEL is greater than TOE and other TOE is greater than HEEL. In the cases HA and TOE-HEEL provide different values, the value of the other is taken if one is zero, else zero is taken. On the other hand MD is determined by the average value of BMC_G features. The direction is computed by comparing the value of average value of BMC_G with the starting value of BMC_G. The direction has any of the three values '-1' for left direction, '0' for straight and '1' for right direction.

The feature values as depicted in the spatial feature section are extracted from k number of consecutive frames. Thus we have the feature size $k \times 28$. But all AUs may not active for every work. So we can eliminate the non-active features to take the classification decision in faster way. We have now 28 columns and each one contains k number of values. The average(AV), min (MIN), max (MAX) and standard deviation (STD) values from each column are extracted. In case of nonactive feature, the values AV, MIN and MAX are same and STD values are near to zero. Now we have a resultant feature vector of size $k \times 28$ contains AV, MIN, MAX and STD of all feature values. This is used to classify the action.

Table 1. Actions and action locations of AU

Action	AU		
	UL	ML	LL
Stand, Run, Walk, jump, turn back	Head	Hand	Leg
Wave, Punch	Head, Hand	Hand	leg
Kick	Head	Hand, Leg	Leg

3.2. Action Classification

The successful completion of previous step, we got a feature vector, which contains the activity of AUs and other derived features. The activity of any AU is described by its active region (s) i.e. UL, ML, LL. An overview general observation of different actions and action locations of AUs for the corresponding actions is shown in Table. 1.

Table. 2 and Table. 3 give an instance of AUs, their acting regions and some other details of two actions walking and running respectively.

Table 2. Action components and action level for Walking

Action: Walking			
Acting components	Upper level	Middle level	Lower level
Head	Yes	No	No
Hands	No	Yes	Yes
Legs	No	No	Yes
Stride angle: (0-40) degree, Location difference between two legs are very close to or less than $(1/5) \times H$			

Table 3. Action components and action level for running

Action: Running			
Acting components	Upper level	Middle level	Lower level
Head	Yes	No	No
Hands	No	Yes	Yes
Legs	No	No	Yes
Stride angle: (0-65) degree, Location difference between two legs are very close to or less than $(1/5) \times H$			

It is understandable from above observation that the AUs and their action regions are main components to take the decision on actions. The action region is determined from the feature values of the corresponding AUs. Four AUs namely body, head, hands and legs are represented by BMC, head, (hand1, hand2) and (leg1 and leg2) respectively. Three other important derived components are SA, MD and BMC_G. We have studied number of test case of different human actions collected from several publicly available human action datasets and observed that AUs are responsible to perform any action as we have seen values of AUs in different region for different actions. On the other hand the values of AUs are in a similar region when the actions are same however it is done by different actors. Thus we have develop a knowledge base for human action classification. The knowledge based includes the action regions of AUs and other derived features for different actions as shown in Table. 4. Every AU has three properties viz. A) Moving values (MV): this value defines the amount of movement of the corresponding AU and this a temporal feature. The value of MV are any one out of '0', and '1' and '2'; '0' for no movement or static and '1' for slow movement and '2' for rapid movement. B) Moving direction(MDR): the direction of the corresponding AU. MDR can be of three different types X, Y and XY corresponding to horizontal, vertical and diagonal directions. iii) Moving regions: The regions of AUs where action is taken place. Moving regions must have one of the three values viz. UL,ML,LL).

Table 4. Knowledge base for action classification

Action	(MV,MDR)	Values of AU and other derived features (MV,MR)					
	BMC_G	BMC	Head	Hand	Leg	SA	MD
Stand	(0,-)	(0,X)	(0,X)	(0,X)	(0,X)	-	(-1,0,1)
Wave	(0,-)	(1,MID)	(0,UP)	(2,UP/MID)	(0,LL)	-	(-1,0,1)
Punch	(0,-)	(1,MID)	(1,UP)	(2,UP/MID)	(0,LL)	(0-30)	(-1,0,1)
Kick	(0,-)	(1,MID)	(1,UP)	(1,MID)	(0,LL/MID)	(0-130)	(-1,0,1)
Jump	(2,Y/XY)	(1,MID)	(1,UP)	(1,MID)	(1,LL)	-	(-1,0,1)
Run	(2,X/XY)	(1,MID)	(-, -)	(-, -)	(2,LL)	(0-90)	(-1,0,1)
Walk	(2,X/XY)	(1,MID)	(-, -)	(-, -)	(2,LL)	(0-60)	(-1,0,1)
Turn back	(1, X/XY)	(1,MID)	(1,UP)	(1,MID)	(1,LL)	(0-30)	(-1,0,1)

In Table. 4, the symbol '-' used for don't care condition, that is the rule does not depend on that value, which is represented by '-'. The direction of any action is determine by the value of MD.

4. Results and analysis

4.1. Consideration:

The class of actions include the actions performed by a single human object with the body not upside down position . Thus the class of actions includes standing, walking, running, punching, kicking, jumping, jogging, clapping, waving etc. One important concern connected with HAR is that the estimation of the number of frames to determine human action. Most of the human actions are repetitive in nature i.e. the same cycle is continuing over a period of time. The time to complete a cycle is below a second in case of all actions considered in several datasets. We set the value of k as the frame rate for this reason. The ultimate form of the STBPM feature is 28×4 , which does not depend on k . In some of the cases more number of cycles may included in one test case, but that will not create any problem as STBPM includes features, which are extracted from all the sequences. Some of the actions are complicated and perform more than one actions to complete a single one. For example, if there is any action of a running human who is taking rest by standing. This action is a combination of running and standing, but our proposed approach takes the most active action out of all performing actions. Thus for the above example, the proposed method will categorize the action as 'running'. The proposed method extracts features from silhouette of the objects collected from consecutive frames. So to validate the results, we need the help of those datasets having silhouette of the ground truth. The performance evolution of the proposed action recognition framework is done on two publicly available datasets: Weizmann [35] and MuHAVi [34]. All the datasets used here include silhouette sequences.

4.2. Datasets:

4.2.1. Wizeman Dataset [35]

We conducted a series of experiments on the Weizmann Human Action Database available on-line [35]. This is a very common dataset, many state-of-the-art approaches report performance on it thus allowing easy comparison. The database contains 90 low-resolution video and silhouette sequences (180×144 pixels) that show 9 different people each performing 10 different actions, such as jumping, walking, running, skipping, etc. We have applied our algorithm except 'bend', which is violated the consideration upside down.

4.2.2. MuHAVi [34]

To isolate the challenge of object detection, it is assumed that the segmentation problem (i.e. to obtain the silhouettes) has been solved. To address this, the dataset provides a sub-set of data that has been (painstakingly) manually annotated. The whole dataset can still be used in a combined segmentation/action recognition algorithm. In MuHVAi-14 datasets, actions are classified into 14 different classes viz. "CollapseLeft", "CollapseRight", "GuardToKick", "GuardToPunch", "KickRight", "PunchRight", "RunLeftToRight", "RunRightToLeft", "StandUpLeft", "StandUpRight", "TurnBackLeft", "TurnBackRight", "WalkLeftToRight" and "WalkRightToLeft". We can classify the actions irrespective

Table 5. Different action class and actions of Wizeman and MuHVAi datasets

Wizeman Datasets		MuHVAi Datasets	
Action class	Action	Action class	Action
jack	Jumping Jack	AC1	GuardToKick
jump	Jumping Forward	AC2	GuardToPunch
pjump	Jumping in place	AC3	KickRight
run	Running	AC4	PunchRight
side	galloping sideways	AC5	RunLeftToRight
skip	skip one leg while moving	AC6	RunrightToLeft
walk	walking	AC7	TurnBackRight
wave1	waving an hand	AC8	TurnBackright
wave2	waving two hands	AC9	WalkLeftToRight
		AC10	WalkRightToLeft

of direction and view points, but we cannot consider the "Collapse" and "StandUp" actions for our experiment. The extraction of AUs from silhouettes is an important area of our algorithm, but for those two actions when human body get squeezed in a smaller area, and AUs are not properly traceable by silhouette information. Other 10 videos of the datasets are classified into 7 actions in our considerations as the work is direction invariant. A brief of different action class and corresponding actions are shown in Table. 5.

4.3. Evaluation

The proposed methodology implemented a rule-base to classify different human actions. So the approach does not need any training for classification, but uses prior knowledge to develop the rule-base. We consider a few number of actions from a huge variety of human actions, but the effort is to construct the grammar, which can portray different human actions effectively. This technique is totally based on Human silhouette. So the proposed approach used the datasets where silhouette of corresponding video are provided. The following measures are used to measure the efficiency of the proposed method.

4.3.1. Confusion matrix

Confusion matrix gives a clear knowledge about actual and wrong classification of any classifier. Diagonal values of a confusion matrix determines the number true positive out of total number of action of that action class and the values other than diagonal values define the mis-classification. For example, as in table-6 in case of action class 'jump' '8' out of '9' action are successfully classified, and the remaining one is classified as action class 'run'. Confusion matrix for Wizeman and MuHVAi datasets are shown in Table. 6 and Table. 7 respectively. In case of Wizeman dataset we have confusions only for three action classes viz. 'jump', 'side' and 'skip'. All the three actions are confused with the action class 'run', which has a certain similarities with those three actions. On the other hand in case of MuHVAi datasets, the two action classes are confused with each other namely 'GuardToPunch' and 'GuardToKick'. The percentage of successful classification, which proves the efficiency of the proposed method are 95.06 and 93.75 for Wizeman and MuHVAi datasets respectively. This amount of success rate shows the efficiency of our technique.

4.3.2. Misclassified frames (MCF)

Our proposed approach is totally dependent of the correctness of AU localization. So the first thing we emphasize on accuracy of AU localization. The proposed approach gives good accuracy in terms of localization of the AUs for most of the frames of the video. We also apply the moving direction (MD) computation for each of the frames of different videos. MD is computed using two of the AUs viz. leg(s) and head and it is one of the main parameter to recognize the actions. The average mis-classification rate (AMR) of MDs in all videos are mentioned in Table. 5 for the respective datasets and are shown in Table. 8. For MuHVAi and Wizeman datasets, the MCF are 4.46% and 3.20% respectively. This mis-classification rate is very minimal with respect to the frame rate of the video as the decision is taken on the basis of majority response.

Table 6. Confusion matrix for wizeman datasets

	Walk	Run	Jump	Pjump	Side	Skip	Wave1	Wave2	Jack
Walk	9/9								
Run		9/9							
Jump		1/9	8/9						
Pjump				9/9					
Side		1/9			8/9				
skip		2/9				7/9			
Wave1							9/9		
Wave2								9/9	
Jack									8/9

Table 7. Confusion matrix for MuHVAi datasets

	GuardToKick	GuardToPunch	KickRight	PunchRight	RunLeftToRight	RunRightToLeft	TurnBackLeft	TurnBackRight	WalkLeftToRight	WalkRightToLeft
GuardToKick	6/8	2/8								
GuardToPunch	2/8	6/8								
KickRight			8/8							
PunchRight				8/8						
RunLeftToRight					4/4					
RunRightToLeft						4/4				
TurnBackLeft							8/8			
TurnBackRight								8/8		
WalkLeftToRight									4/4	
WalkRightToLeft										4/4

Table 8. Average miss-classification rate of MDs of the corresponding datasets

Datasets	Actions	MAR
MuHAVi [34]	10	4.46
Wiseman [35]	9	3.29

Table 9. Comparative study of several methodologies of the related research with the proposed approach for wizeman datasets

Approach	Input	Actions	Evaluations	Rate
Ikizler and Duyugulu [10]	Silhouette	9	LOSO	100
Tran and Sorokin [11]	Silhouette	10	LOSO	100
Eweiwi et al. [12]	Silhouette	10	LOSO	100
Harnandez et al. [13]	Images	10	LASO	90.3
Cheema et al. [6]	Silhouette	9	LOSO	91.6
Charaoui et al. [8]	Silhouette	9	LOSO	92.8
Proposed	Silhouette	9	No Training	95.06

4.3.3. Leave-one-actor-out (LOAO) and leave-one-sequence-out (LOSO) cross validation

There are several human recognition approaches in related research domain, which are using LOAO and LOSO as the part of estimating the recognition ability of the different classification techniques. Any human action datasets contains several sequences of human actions performed by different actors. In case of LOAO, the classifier is trained by leaving the video of any one actor, that video will be used for testing purpose. On the other hand, in case of LOSO the classifier is trained by leaving any one sequence, which will be used for testing purpose.

LOAO and LOSO can not be applicable for the proposed approach as the approach does not require training to classify human actions. On the other hand, the approaches use LOAO or LOSO for measuring their efficiency involving training and testing with a part of the dataset, which is not used for training, but in our proposed methods all the will be tested with out training. So the rate of successful classification of the proposed methodology can be easily comparable with the efficiency of the related research work. The results of applying our proposed method on Wizeman dataset our method with in comparison with that of the state of the art research works is shown in Table. 9 based on classification efficiency. The success rate of the proposed approach considering Wizeman dataset is 95.06%, which outperforms others. It can be inferred that though our technique doesn't yield hundred percent success rate, but the proposed work is a ready to be implemented one as it doesn't require any learning.

4.3.4. Identical training and testing actors, novel camera

This is used for view invariance test of the algorithm for any dataset having video of the same actors from different angle. So in this case the training is done from the video taking through one camera, which is situated in a certain angle and testing the video which are taking through other camera, which is placed in the different angle. MuHVAi datasets provides multi-view data, which is exploited for novel camera test. Our proposed method does not depend on the camera view point when it is parallel to the object. We have tested all the videos, which are given in the dataset without training a single video, and the actions were recognized irrespective of the camera position. The results of our method and its comparison with the other state of the art is shown in Table. 10 for the MuHVAi dataset.

4.3.5. Identical training and test cameras, novel actors

This testing phase is used for checking the robustness of the algorithm irrespective of actors i.e. the view angle of the actor in the test video differs with that of the actor in the tanning video. MuHVAi dataset provides the videos from different angle of the same actor. We have applied our proposed method on all the videos irrespective of actors. To recognize actions of different actors, the proposed approach doesn't require training. Hence the proposed method qualifies the novel actor test. A comparative study of the results of our approach and that of the other related research works is shown in Table. 10 in context with classification efficiency on MuHVAi dataset for novel actor test.

The accuracy rate of the proposed work over MuHVAi dataset is 93.75%, which is better than any other methods of the related research work. This performance also validate novel actor and novel camera test as we did not use any training mechanism.

5. Conclusion

The present work proposes a new spatio-temporal feature extraction technique coined as STBPM and a rule based logic termed as RAC to classify human actions. Feature extraction is done on silhouette of the foreground for k

Table 10. Comparative study of several methodologies of the related research with the proposed approach for MuHVAi datasets

Approach	Input	Actions	Evaluations	Success Rate
Singh et al. [7]	Silhouette	14	LOSO	82.4
Eweiwi et al. [12]	Silhouette	14	LOSO	91.9
Cheema et al. [6]	Silhouette	14	LOSO	86.0
Chaaroui et al. [8]	Silhouette	14	LOSO	92.8
Proposed	Silhouette	10	No Training	93.75

Table 11. results on MuhVAi datasets for novel actor test

Approach	Input	Actions	Evaluations	Success Rate
Singh et al. [7]	Silhouette	14	LOSO	61.8
Eweiwi et al. [12]	Silhouette	14	LOSO	77.9
Cheema et al. [6]	Silhouette	14	LOSO	73.5
Chaaroui et al. [8]	Silhouette	14	LOSO	82.4
Proposed	Silhouette	10	No Training	93.75

number of consecutive frames. Our proposed technique successfully localizes AUs and determines human actions with high accuracy except the action with head upside down, sitting and lying. Our technique doesn't required any training and is also independent of the camera view angle. Experimental results involving publicly available datasets shows that our technique outperforms the other state of the art techniques in terms of success rate of HAR. In future, we wish to extend our work for detecting action in the sitting and lying conditions making the detection process more extensive.

6. Acknowledgment

This work is funded by DST, Ministry of Science and Technology, Government of India through INSPIRE project.

7. References

References

- [1] J. K. Aggarwal, M. S. Ryoo, Human Activity Analysis: A Review, ACM Computing Surveys (CSUR). Volume 43 Issue 3, April 2011, Article No. 16.
- [2] K. Guo, P. Ishwar, Action Recognition from Video Using Feature Covariance Matrices, IEEE Transaction on Image Processing, Vol. 22, No. 6, June 2013.
- [3] M. A. Mendoza, N. P. Blanca, HMM-Based Action Recognition Using Contour Histograms, Pattern Recognition and Image Analysis, Lecture Notes in Computer Science, 4477, Springer, Berlin/Heidelberg (2007), pp. 394-401.
- [4] C. H. Hsieh, P. S. Huang, and M. D. Tang, Human Action Recognition Using Silhouette Histogram, Proceedings of the Thirty-Fourth Australasian Computer Science Conference (ACSC 2011), pp. 11-15, Perth, Australia, 17-20, 2011.
- [5] L. Gorelick, M. Galun, E. Sharon, A. Brandt, R. Basri, Shape representation and recognition using the poisson equation, in: Computer Vision and Pattern Recognition, Washington DC, USA, June 2003.
- [6] S. Cheema, A. Eweiwi, C. Thureau, C. Bauckhage, Action recognition by learning discriminative key poses, IEEE Internat. Conf. on Computer Vision Workshops (ICCV Workshops), pp. 1302-1309.
- [7] S. Singh, S. Velastin, H. Ragheb, Muhavi: A multicamera human action video dataset for the evaluation of action recognition methods. Seventh IEEE Internat. Conf. on Advanced Video and Signal Based Surveillance (AVSS), pp. 8-15, 2005.
- [8] A. A. Chaaroui, P. Climent-Perez, F. Florez-Revuelta, Silhouette-based human action recognition using sequences of key poses, Pattern Recognition Letters. 34(15), pp. 1799-1807. ISSN (print) 0167-8655.
- [9] L. Wang, T. Tan, Silhouette Analysis-Based Gait Recognition for Human Identification, IEEE Transactions on Pattern Analysis & Machine Intelligence. 2003 vol.25 Issue No.12.
- [10] N. Ikizler, P. Duygulu, Human action recognition using distribution of oriented rectangular patches, Lecture Notes in Computer Science. 4814. Springer, Berlin/Heidelberg, pp. 271-284
- [11] D. Tran, A. Sorokin, Human activity recognition with metric learning. Computer Vision ECCV 2008, Lecture Notes in Computer Science. 5302. Springer, Berlin/Heidelberg, pp. 548-561.
- [12] A. Eweiwi, S. Cheema, C. Thureau, C. Bauckhage, Temporal key poses for human action recognition, IEEE Internat. Conf. on Computer Vision Workshops (ICCV Workshops). 2011, pp. 1310-1317.

- [13] J. Hernández, A. Montemayor, J. Pantrigo, A. Sánchez, Human action recognition based on tracking features, *Foundations on Natural and Artificial Computation, Lecture Notes in Computer Science*. 6686, pp. 471-480.
- [14] M. Ahmad, I. Parvin, and S. W. Lee, Silhouette history and energy image information for human movement recognition, *J. Multimedia*, vol. 5, no. 1, pp. 12-21, Feb. 2010.
- [15] A. Bobick and J. Davis, The recognition of human movement using temporal templates, *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 23, no. 3, pp. 257-267, Mar. 2001.
- [16] Y. Chen, Q. Wu, and X. He, Human action recognition by Radon transform, *IEEE Int. Conf. Data Mining Workshops*, Dec. 2008, pp. 862-868.
- [17] L. Gorelick, M. Blank, E. Shechtman, M. Irani, and R. Basri, Actions as space-time shapes, *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 29, no. 12, pp. 2247-2253, Dec. 2007.
- [18] S. Ali and M. Shah, Human action recognition in videos using kinematic features and multiple instance learning, *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 32, no. 2, pp. 288-303, Feb. 2010.
- [19] S. Danafar and N. Gheissari, Action recognition for surveillance applications using optic flow and SVM, *Asian Conf. Comput. Vis.*, 2007, pp. 457-466.
- [20] A. Fathi and G. Mori, Action recognition by learning mid-level motion features, in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2008, pp. 1-8.
- [21] Y. Ke, R. Sukthankar, and M. Hebert, Efficient visual event detection using volumetric features, in *Proc. IEEE Int. Conf. Comput. Vis.*, vol. 1, Oct. 2005, pp. 166-173.
- [22] J. Liu, S. Ali, and M. Shah, Recognizing human actions using multiple features, *IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2008, pp. 1-8.
- [23] D. Cunado, M. S. Nixon, and J. N. Carter, Automatic extraction and description of human gait models for recognition purposes, *Comput. Vis. Image Understand.* vol. 90, no. 1, pp. 1-41, Apr. 2003.
- [24] L. Goncalves, E. D. Bernardo, E. Ursella, and P. Perona, Monocular tracking of the human arm in 3-D, in *Proc. IEEE Int. Conf. Comput. Vis.*, Jun. 1995, pp. 764-770.
- [25] K. Rohr, Toward model-based recognition of human movements in image sequences, *CVGIP, Image Understand.*, vol. 59, no. 1, pp. 94-115, Jan. 1994.
- [26] L. Wang, H. Ning, T. Tan, and W. Hu, Fusion of static and dynamic body biometrics for gait recognition, *IEEE Trans. Circuits Syst. Video Technol.* vol. 14, no. 2, pp. 149-158, Feb. 2004.
- [27] P. Dollar, V. Rabaud, G. Cottrell, and S. Belongie, Behavior recognition via sparse spatio-temporal features, *2nd IEEE Int. Workshop Vis. Surveill. Perform. Evaluation Tracking Surveill.*, Oct. 2005, pp. 65-72.
- [28] I. Laptev, M. Marszalek, C. Schmid, and B. Rozenfeld, Learning realistic human actions from movies, *IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2008, pp. 1-8.
- [29] J. Niebles, H. Wang, and L. Fei-Fei, Unsupervised learning of human action categories using spatial-temporal words, *Int. J. Comput. Vis.*, vol. 79, no. 3, pp. 299-318, Sep. 2008.
- [30] C. Schuldt, I. Laptev, and B. Caputo, Recognizing human actions: A local SVM approach, in *Proc. Int. Conf. Pattern Recognit.*, vol. 3, Aug. 2004, pp. 32-36.
- [31] S. F. Wong and R. Cipolla, Extracting spatio-temporal interest points using global information, *IEEE Int. Conf. Comput. Vis.*, Oct. 2007, pp. 1-8.
- [32] T. Starner and A. Pentland, Visual recognition of American sign language using hidden Markov model. *IEEE Int. Conf. Autom. Face Gesture Recognit.*, Jan. 1995, pp. 1-52.
- [33] J. Yamato, J. Ohya, and K. Ishii, Recognizing human action in time sequential image using hidden markov model, in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 1992, pp. 379-385.
- [34] <http://dipersec.king.ac.uk/MuHAVi-MAS/>
- [35] <http://www.wisdom.weizmann.ac.il/~vision/SpaceTimeActions.html>
- [36] S. D. Choudhury, T. Tjahjadi, Silhouette-based gait recognition using Procrustes shape analysis and elliptic Fourier descriptor, *Pattern Recognition*. 45(2012). PP. 3414-3426.