

Date of publication xxxx 00, 0000, date of current version xxxx 00, 0000.

Digital Object Identifier 10.1109/ACCESS.2017.Doi Number

An LSTM-Based Approach for Understanding Human Interactions Using Hybrid Feature Descriptors over Depth Sensors

MANAHIL WAHEED¹, AHMAD JALAL¹, MOHAMMED ALARFAJ², YAZEED YASIN GHADI³, TAMARA AL SHLOUL⁴, SHAHARYAR KAMAL¹ AND DONG-SEONG KIM⁵

¹Department of Computer Science, Air University, E-9, Islamabad 44000, Pakistan

²Department of Electrical Engineering, King Faisal University, Al-Ahsa, 31982, Saudi Arabia

³Department of Computer Science and Software Engineering, Al Ain University, Al Ain, 15551, UAE

⁴Department of Humanities and Social Science, Al Ain University, Al Ain, 15551, UAE

⁵School of Electronic Engineering, Kumoh National Institute of Technology, Yanghondong, Daehakro 61, 730-701, South Korea

Corresponding author: Dong-Seong Kim (e-mail: dskim@kumoh.ac.kr).

This research work was supported by Priority Research Centers Program through NRF funded by MEST(2018R1A6A1A03024003) and the Grand Information Technology Research Center support program (IITP-2021-2020-0-01612) supervised by the IITP by MSIT, Korea.

ABSTRACT Over the past few years, automatic recognition of human interactions has drawn significant attention from researchers working in the field of Artificial Intelligence (AI). And feature extraction is one of the most critical tasks in developing efficient Human Interaction Recognition (HIR) systems. Moreover, recent researches in computer vision suggest that robust features lead to higher recognition accuracies. Hence, an improved HIR system has been proposed in this paper that combines 2D and 3D features extracted using machine learning and deep learning techniques. These discriminative features result in accurate classification and help avoid misclassification of similar interactions. Ten keyframes have been extracted from each video to reduce computational complexity. Next, these frames have been preprocessed using image normalization and noise removal techniques. The Region Of Interest (ROI), which contains the two humans involved in the interaction, has been extracted using motion detection. Then, the human silhouettes have been segmented using the GrabCut algorithm. Next, the extracted silhouettes have been converted into 3D meshes and their heat kernel signatures (HKS) have been obtained to extract key body points. A Convolutional Neural Network (CNN) has been used to extract full-body features from 2D full-body silhouettes. Then, topological and geometric features have been extracted from the key body points. Finally, the combined feature vector has been fed into Long Short-Term Memory (LSTM) and each interaction has been recognized using a Softmax classifier. The proposed system has been validated via extensive experimentation on three challenging RGB+D datasets. The recognition accuracies of 91.63%, 90.54%, and 90.13% have been achieved with the SBU Kinect Interaction, NTU RGB+D, and ISR-UoL 3D social activity datasets respectively. The results of extensive experiments performed on the proposed system suggest that it can be used effectively for various applications, such as security, surveillance, health monitoring, and assisted living.

INDEX TERMS 3-D mesh, depth videos, geodesic distance, heat kernel signature, human interaction recognition, RGB videos, and topological features.

I. INTRODUCTION

The task of Human-Human Interaction (HHI) recognition involves detecting and understanding social interactions between two humans. These interactions can be everyday activities like talking, passing objects, hugging, and waving. Similarly, these can be assisted living activities such as helping a person standing up, helping another person walk,

or drawing another person's attention. Moreover, suspicious activities including touching someone's pocket, pushing someone, or fighting are also of interest for researchers in this field. HIR has become a trending topic in the field of artificial intelligence because of its wide range of applications, including security [1-3], content-based video retrieval [4-6], healthcare [7-11], and surveillance [12-15].

Even though significant progress has been made in this regard and many efficient HIR systems have been developed for various purposes, detecting human interactions remains challenging because of multiple reasons, such as different viewpoints, change of clothing, poor lighting, different interactions containing similar motions, and unavailability of large datasets. However, low-cost depth sensors, such as Microsoft Kinect [16] are now being used excessively since these are not as affected by lighting conditions as RGB cameras. Moreover, many interactions seem similar and are often misclassified. For example, two humans exchanging a very small object may look very similar to two people shaking hands. On the contrary, the same interaction appears different when viewed from various viewpoints. Hence, it is very important to extract distinctive features from images that can easily differentiate between two interactions that look the same.

This research paper proposes a novel approach for efficient video-based human interaction recognition using both machine learning and deep learning techniques. Human silhouettes have been extracted from both RGB and depth frames using GrabCut. Additional masking has been used to improve the output of the GrabCut algorithm in complex scenarios. Next, the full-body RGB and depth silhouettes have been fed into two separate CNN models and the extracted features have been concatenated. Then 3D meshes have been generated from the full-body silhouettes and their heat kernel signatures have been obtained. These have been used to extract six key body points. These key points have been used to extract topological and geometric features. The three different types of features have been combined and fed into LSTM. Finally, the Softmax classifier has been used for interaction recognition. Three publically available datasets have been used that provide RGB, depth, and skeletal information of human interactions. The major contributions of this research work include:

- Silhouette segmentation from both RGB and depth images using GrabCut algorithm.
- Training and concatenation of two separate CNN models for RGB and depth images.
- 3-D mesh generation from 2-D silhouettes.
- Detection of key points via heat kernel signatures based on geodesic distance.
- Extraction of topological and geometric features using key body points.
- Extensive experimentation on three large and challenging RGBD video datasets.

Section II of the paper describes similar research work and the proposed system architecture has been discussed in Section III. Section IV presents the implementation details and results of the proposed method. Section V contains discussion on various aspects of the designed system and section VI contains the conclusion of this paper and proposes future work of the authors.

II. RELATED WORK

Researchers have been actively contributing to the development of efficient HIR systems. The existing systems have been divided into two categories: marker-based systems and video-based systems. Researches falling into each category have been discussed in detail below:

A. MARKER-BASED HIR SYSTEMS

In marker-based HIR systems, different types of sensors, for example, reflective spheres, light-emitting diodes, and infrared markers, are mounted on the bodies of the humans whose movements are being monitored. These systems are commonly used for rehabilitation treatments [17]. For example, a marker-based motion tracking system is proposed in [18] to analyze the movement of various body parts. The authors have argued that accurate detection of movement of different parts can result in better therapeutic decisions. However, the system was evaluated on a small dataset of only 10 real patients. Similarly, the authors in [19] attached an IR camera and an infrared emitter with a passive hand skateboard training device for conventional upper limb training. The proposed device was used to train eight patients with abnormal upper limb function. After four weeks of training, all the patients were able to move the hand skateboard along the designated figure of eight path.

Capturing body movements is also critical for sports. Hence, researchers have used marker-based sensors for movement detection in walking gait [20], discus [21], dressage [22], and swimming [23] activities. Esfahani et al. [24] developed a trunk motion system (TMS) using printed body-worn sensors (BWS). Twelve BWSs were printed on stretchable clothing to measure the 3D trunk movements and a neural network data fusion algorithm was used to integrate the data from sensors. However, one shortcoming of these marker-based techniques is that they require installation and calibration of multiple cameras. Hence, these systems are quite expensive. Moreover, they can only encode two-dimensional motion information.

B. VIDEO-BASED HIR SYSTEMS

In video-based HIR systems, video cameras are used to record human interactions. In such systems, the first step is to extract important features or interest points [25,26]. Based on these distinctive features, the interaction that has been performed in the video is identified. M. Khan et al. [27] proposed a deformable part-based modeling technique to detect the body parts of a patient and track them in subsequent frames. Their system then performed movement analysis to detect various movement disorders in infants. They captured the data in a local hospital using Microsoft Kinect but it was only RGB data. M. H. Khan et al. [28] proposed a system for analyzing a patient's body movements during Vojta therapy. They proposed the use of color features and pixel locations for segmenting the patient's body in the images. Then they employed a multi-dimensional feature vector to classify the correct movements using multiclass SVM.

Some researchers [29,30] also prefer extracting various features and then combining them since hybrid features have yielded better classification results in the past. For example,

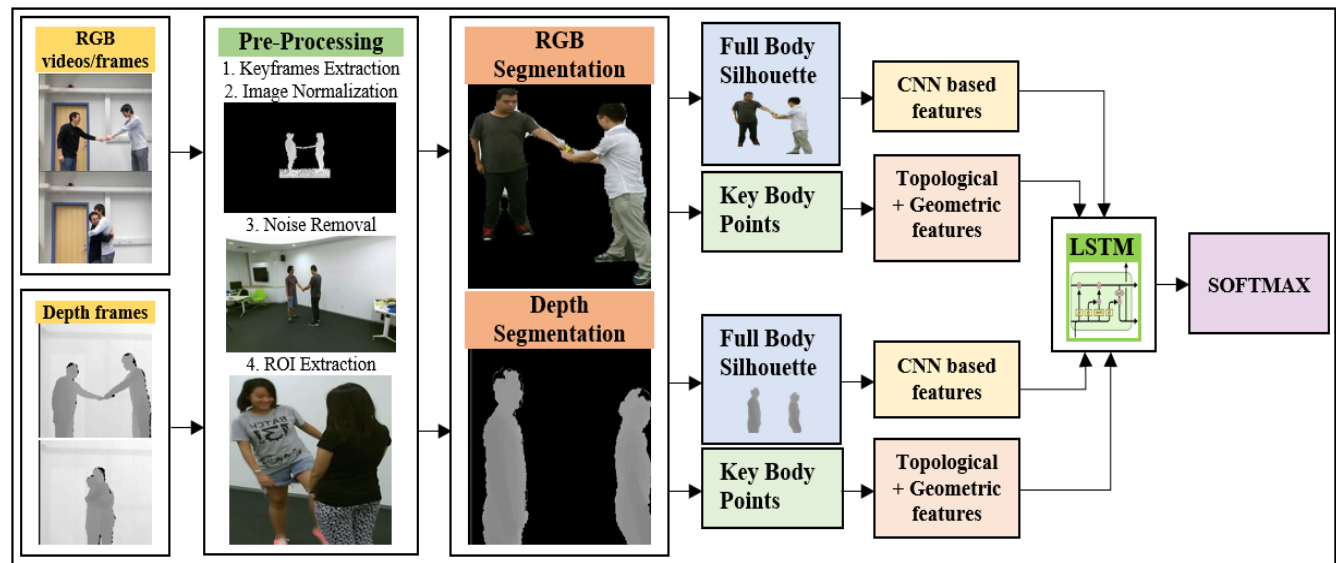


Figure 1. The architecture of the proposed HIR system.

A. Jalal et al. [31] combined four different types of features including blobs, multiple orientations, Fourier transforms, and geometrical points. Similarly, the hybrid features introduced in [32] included energy, sine, distinct body parts movements, and 3D Cartesian views of smoothing gradients. The authors of [33] also used a hybrid of four different local descriptors: spatio-temporal features, energy-based features, shape-based angular and geometric features, and Motion-Orthogonal Histograms of Oriented Gradients (MO-HOG). However, all these approaches only employed 2D features.

III. THE PROPOSED APPROACH

The proposed system can be divided into four main sections: image preprocessing, image segmentation, feature extraction, and interaction recognition. The used methodologies and results of each section are discussed in detail below. Fig. 1 shows a flow chart of the proposed system architecture.

A. IMAGE PREPROCESSING

The RGB videos taken from the NTU RGB+D dataset have been converted into image frames at the rate of 31 frames per second. The image frames of the other two datasets were already available. Since there are multiple videos for each interaction class and each video consists of a large number of image frames, 10 keyframes have been extracted from each video to reduce complexity. The extracted frames have been normalized and noise has been removed from them. Finally, regions of interest (ROI) have been extracted from each frame. These four subsections are explained in detail below. Moreover, Algorithm 1 explains each step of the preprocessing stage.

1) KEYFRAME EXTRACTION

The number of frames varies from video to video. So, to get a fixed number of frames, 10 keyframes have been extracted from each video of every dataset. To extract the keyframes of a video, the histograms of all the image frames have been obtained. The histogram of an image x can be computed as;

$$P_x(i) = \frac{n}{N}, i = 0,1,2 \dots 256 \quad (1)$$

where n is the number of pixels with intensity i and N is the total number of pixels in the input image. Then the histograms of every two consecutive frames have been compared and their differences have been stored in a sorted array. The indices corresponding to the top ten differences have been fetched and the images at those indices are referred to as keyframes. In other words, these frames are the ones with the highest differences in their histograms.

2) IMAGE NORMALIZATION

The purpose behind image normalization is to change the pixel values of an image to a common scale so that the image appears more normal to the senses. The depth images in two out of the three datasets used are too dark to be seen by the naked eye. Moreover, features on drastically different scales can be problematic for an HIR system. In other words, features with a larger scale will dominate others and cause the system to make inaccurate assumptions. Hence, all images have been normalized. Each pixel i in the normalized image $I(x)^{norm}$ after applying min-max normalization technique to the original image $I(x)^{org}$, is defined as;

$$I(x_i)^{norm} = \frac{I(x_i)^{org} - I(x)_{min}^{org}}{I(x)_{max}^{org} - I(x)_{min}^{org}} \quad (2)$$

3) NOISE REMOVAL

The technique of “non-local means denoising” has been used to remove noise from the images. Local-means filters replace

the value of a pixel with the mean of a group of pixels surrounding it. However, a non-local means filter takes the weighted mean of all the pixels in the image. The weight of each pixel depends on how similar it is to the target pixel. A pixel in the denoised image $u(p)$ at point p after applying non-local means denoising technique on a pixel at point q in the original image $v(q)$, is defined as;

$$u(p) = \frac{1}{C(p)} \int v(q) f(p, q) dq \quad (3)$$

where $f(p, q)$ is the weight and $C(p)$ is a normalization factor defined as:

$$C(p) = \int f(p, q) dq \quad (4)$$

4) ROI EXTRACTION

The regions of interest have been extracted from images through motion detection. Frame differencing technique has been used to detect motion in subsequent frames and rectangular boxes have been drawn over the points where motion has been detected. Since different body parts show different movements, rectangular regions of various sizes for different body parts have been obtained. Each region has a starting point x , y , width w , and height h . A minimum area condition has also been set for these rectangular regions to be considered valid. The minimum and the maximum values of x and y and the maximum values of w and h have been obtained. Finally, one rectangular region comprising all the smaller regions has been extracted as the region of interest. Its starting position is the minimum value of x and y obtained from all the smaller regions and its width is equal to the maximum value of x added to the maximum value of w . Similarly, its height is equal to the maximum value of y added to the maximum value of h .

Algorithm 1 Preprocessing

```

Input: raw frames
Output: ROI coordinates (x,y,w,h) in preprocessed frames
           %key frame extraction%
for  $i$  in range(total frames)
   $diff(i) \leftarrow hist(frame(i)) - hist(frame(i+1))$ 
   $indices \leftarrow nlargestindex(10, range(len(difference)))$ 
   $key\_frame(j) \leftarrow frame(indices[j])$ 
end
           %normalization and noise removal%
   $img \leftarrow key\_frames(i)$ 
   $norm\_img \leftarrow zscore\_norm(img)$ 
   $denoised\_img \leftarrow nonlocalmeans\_denoising(norm\_img)$ 
           %ROI extraction%
   $diff\_img \leftarrow absdiff(key\_frame(i), key\_frame(i+1))$ 
   $contours = FindContours(diff\_img)$ 
  for  $contour$  in contours:
     $(x, y, w, h) \leftarrow boundingRect(contour)$ 
    if  $contourArea(contour) > min\_area$ :
       $draw\_rectangle(img1, (x, y), (x+w, y+h))$ 
       $coordinates.append(x, y, w, h)$ 
  end
return  $coordinates$ 

```

B. IMAGE SEGMENTATION

Image segmentation is the process of segmenting the image into two parts: foreground and background. The GrabCut algorithm proposed by C. Rother et al. [34] has proven to be an efficient foreground extraction technique. It takes a rectangular region as input and assumes that all the pixels outside that region belong to the background. Then it uses a Gaussian Mixture Model (GMM) to define the area inside the rectangle by labeling each pixel as probable background and probable foreground depending upon their relation to the provided data.

Using this pixel distribution, a weighted graph is created. All pixels are treated as nodes in the graph. Then two additional nodes are added: the Source node and the Sink node. Every foreground pixel is connected to the Source node and every background pixel is connected to the Sink node. The weights of edges connecting pixels to the Source node depend on the probability of a pixel of belonging to the foreground or background. The weights between the pixels depend on pixel similarity, that is, if there is a large difference in pixel color, the edge between them will get a low weight and vice versa. Next, the graph is segmented using a Min-Cut algorithm. The graph is cut separating the Source node and the Sink node with a minimum cost function. The cost function is the sum of all weights of the edges that are cut. After cutting the graph, all the pixels connected to the Source node are labeled foreground and those connected to the Sink node are labeled background. The process continues until convergence.

However, in some cases, the extracted foreground contains portions that belong to the background. This problem came up while segmenting images from the NTU RGB+D dataset. In all those images, a major portion of the region of interest is the floor. Hence, a floor mask has been created by extracting a certain range of the intensity values from the original image and the GrabCut output is masked to get accurate results. The results of the segmentation process are shown in Fig. 2.



Figure 2. Segmentation results on RGB images showing (a) the original RGB image, (b) GrabCut output, and (c) segmented humans after applying floor mask.

C. KEY BODY POINTS SELECTION

First, the full-body silhouettes have been converted into 3d meshes [35] as shown in Fig. 3. The center points of the 3d meshes have been considered the source and then geodesic distance-based heat kernel signatures (HKS) of the 3d meshes have been achieved as shown in Fig. 4. HKS, as

introduced by J. Sun et al. [36], is based on a heat kernel, which is a fundamental solution to the heat equation. The heat equation describes the variation of heat distribution with time. HKS is one of the many recent shape descriptors which are based on the Laplace–Beltrami operator associated with the shape [37]. The thermal diffusion process can be described by the heat equation as given;

$$(\Delta - \partial/\partial t)u(x, t) = 0 \quad (7)$$

where Δ is the Laplace–Beltrami operator and $u(x, t)$ is the heat distribution at any point x at a given time t . The solution to this heat equation can be expressed as follows;

$$u(x, t) = \int h_t(x, y)u_o(y)dy \quad (8)$$

where $h_t(x, y)$ is called heat kernel function. The heat kernel equation is the fundamental solution to the heat equation. The Eigenvalue decomposition of the heat kernel is expressed as follows;

$$h_t u(x, t) = \sum_{i=0}^{\infty} \exp(\lambda_i t) \phi_i(x) \phi_i(y) \quad (9)$$

where λ_i and ϕ_i are the i^{th} eigenvalue and Eigen function of Δ . For a concise feature descriptor, HKS restricts the heat kernel only to the temporal domain.

$$h_t(x, x) = \exp(-\lambda_i t) \phi_i^2(x) \quad (10)$$

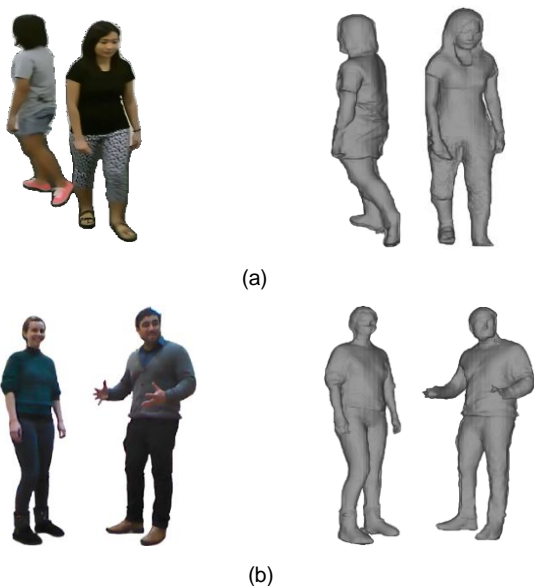


Figure 3. 2D images and the respective 3D meshes of both humans involved in interactions: (a) walking apart; (b) talk.

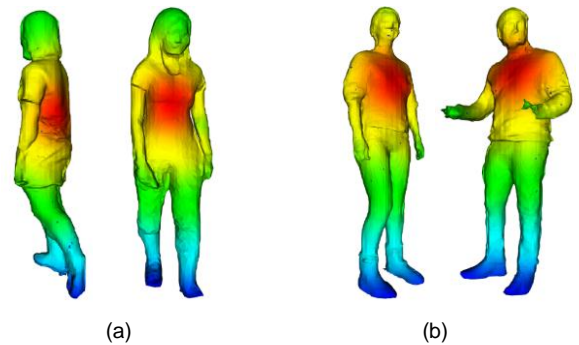


Figure 4. Heat kernel signatures of humans involved in interactions: (a) walking apart; (b) talk.

After obtaining the heat kernel signatures, all vertices in a mesh are grouped into multiple clusters based on their color or intensity value. Moreover, the centroid of each cluster is detected and is stored as a key body point. In this way, six key body points are obtained for each silhouette. When a geodesic path is drawn from the source vertex to the other five target vertices, a 2D leaf skeleton model is obtained as shown in Fig. 5. Algorithm 2 explains the process of extraction of key body points from full-body silhouettes.

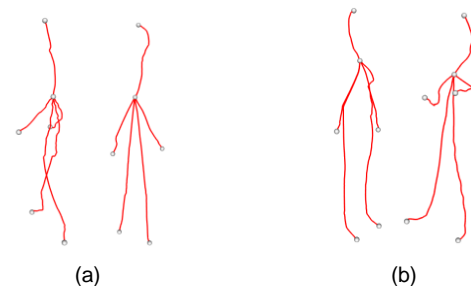


Figure 5. 2D leaf skeleton models using key body points of humans involved in interactions: (a) walking apart; (b) talk.

Algorithm 2 Key Body Points Extraction

Input: segmented silhouettes

Output: key body points ($p1, p2, p3 \dots pn$)

$mesh \leftarrow Get3Dmesh(segmentedsilhouette)$

$HKS \leftarrow GetHeatKernelSignature(mesh)$

$Clusters \leftarrow GetIntensityBasedClusters(HKS)$

for cluster in Clusters:

$KeyPoint \leftarrow GetClusterCentroid(Cluster)$

$KeyBodyPoints.append(KeyPoint)$

end

return KeyBodyPoints

D. FEATURE EXTRACTION

This section can be divided into two phases. In the first phase, a Convolutional Neural Network (CNN) has been used to extract features from full-body silhouettes. Full-body silhouettes have been extracted from the segmented images by removing the black background from the images and making them transparent. In the second phase, topological and geometric features have been extracted using key body

points. Both these phases are described in the following sub-sections.

1) FULL BODY SILHOUETTES: CNN-BASED FEATURES

For extraction of features from images, a convolutional neural network has been used. The transfer learning approach has been employed, which includes using VGG16 as the base model and then fine-tuning its weights according to the used datasets. Visual Geometry Group-16 layers deep (VGG16) [38] is a CNN model that achieved 92.7% on the ImageNet dataset which has 1000 classes. Fig. 6 shows all the layers in the VGG16 model. All images have been reshaped to 224x224x3 to match the desired input size of the VGG16 model. After training on the VGG16 base model, input images are resized to 7x7x512. These are then trained on the proposed CNN model.

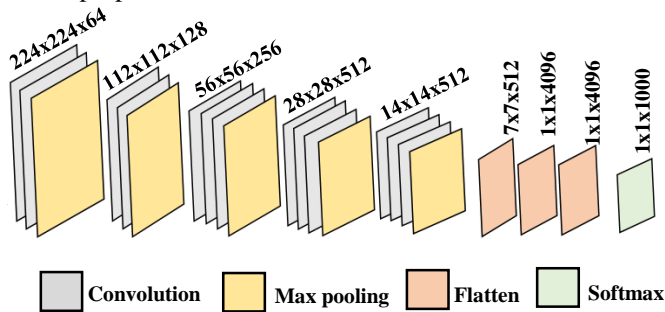


Figure 6. Different layers of the VGG16 architecture with configurations.

There are three convolutional layers in the proposed model with 128, 64, and 32 filters respectively. The size of each filter is 3x3. The convolutional layers compute the output of neurons that are connected to local regions in the input. Convolution is similar to sliding a filter over an image, computing the dot product of filter weights and image pixels. Rectified Linear Unit (RELU) is used as the activation function for all three convolutional layers. It simply rounds up all the negative values to zero as shown;

$$y_k = \max(0, x_k) \tag{5}$$

The convolutional layers are followed by a batch normalization layer. The pixels x_k of input images of each batch are normalized as follows;

$$\widehat{x}_k = \frac{x_k - E(x_k)}{\sqrt{Var(x_k)}} \tag{6}$$

where $E(x_k)$ is the mean and $Var(x_k)$ is the variance of pixel values.

The batch normalization layer is followed by a flatten layer that remaps the output of the batch normalization layer to a column vector. Lastly, a drop-out layer of 0.2 has been used to avoid overfitting. Two such models have been trained: one for RGB images and one for depth images. The two CNN models are then concatenated. Table I shows a summary of the CNN model.

TABLE I. A BRIEF SUMMARY OF THE CNN MODEL

Layer	Output Shape	Parameters
-------	--------------	------------

Conv:128	(None,7,7,128)	65664
Conv:64	(None,7,7,64)	8256
Con:32	(None,7,7,32)	2080
BatchNorm	(None,7,7,32)	128
Flatten	(None,7,7,1568)	0
Dropout	(None,7,7,1568)	0

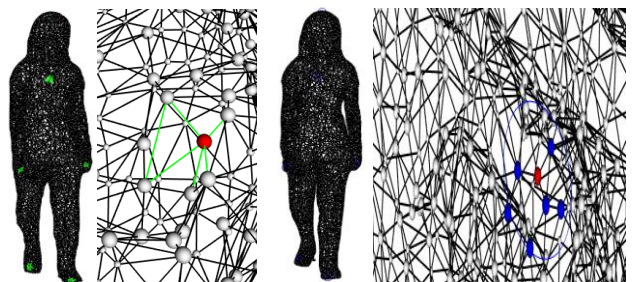
2) KEY BODY POINTS: TOPOLOGICAL FEATURES

Topology can be defined as the spatial relationship between adjacent or neighboring features. Topological features are the properties of a geometric object that are preserved under continuous deformations. In the proposed architecture, four types of topological features have been extracted using the key body points:

1. Geodesic distance from the source.
2. Geodesic path.
3. Connected faces.
4. Nearest neighbors.

A mesh is a collection of vertices, edges, and faces that describe the shape of a 3D object. Every single point in a mesh is a vertex, a line connecting two vertices is an edge, and a flat surface enclosed by edges is called a face. In the proposed approach, the 3d meshes have been converted into graph models and these four topological features have been extracted for each key point.

First, the geodesic distance gd_i between the source vertex and each key point or target vertex has been obtained. Geodesic distance gives the distance between two vertices in a graph along the shortest path between the vertices. Hence, unlike Euclidean distance, geodesic distance considers the shape of the object while computing the distance between two points. If any two vertices are not connected in a graph, the geodesic distance between them will be infinite. After storing the value of geodesic distance, an array of all the vertices lying on this shortest path from source to target vertex have been stored as the geodesic path gp_i . For finding the connected faces, each key point or target vertex has been compared with the three vertices in each face of the mesh. In this way, the faces containing one or more of these target vertices have been found and stored as connected faces cf_i . Finally, the distance of each target vertex from all other vertices in the graph has been computed and stored in a sorted array. Then the top 128 vertices corresponding to the shortest 128 distances have been acquired. These have been stored as the nearest 128 neighbors nn_i . These features are shown in Fig. 7. Hence, for each key point, a topological feature vector $[gd_i, gp_i, cf_i, nn_i]$ has been obtained.



(a) (b)

Figure 7. Topological features including: (a) connected edges (full mesh + zoomed in on one face) and (b) nearest neighbors (full mesh + zoomed in on vertex)

3) KEY BODY POINTS: GEOMETRIC FEATURES

Similar to topological features, some geometric features have also been obtained using the key body points. Ten triangular shapes have been drawn by joining different combinations of key points as shown in Fig. 8. These key points are labeled as left hand (LH), right hand (RH), left foot (LF), right foot (RF), head (H), and torso (T). Finally, the feature vector is also updated as geometric features are added to it. Algorithm 3 explains how these topological and geometric features have been extracted and concatenated in the proposed system.

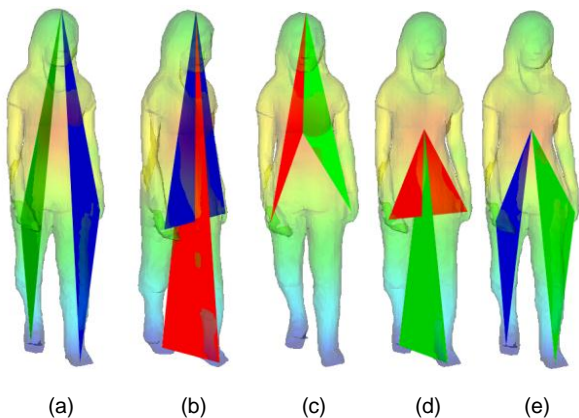


Figure 8. Geometric features including (a) H+LH+LF, H+RH+RF, (b) H+LH+RH, H+LF+RF, (c) H+LH+T, H+RH+T, (d) T+LH+RH, T+LF+RF, and (e) T+LH+LF, T+RH+RF.

Algorithm 3 Topological and Geometric Features

```

Input: key body points
Output: combined feature vectors ( $f_1, f_2, f_3, \dots, f_n$ )
           % Graph Model %
           input ← mesh
points, faces ← getpointsandcellsfrompolydata(input)
for i in range(len(points)):
    actor1 ← createSphere(points[i], radius=0.003)
end
for j in range(len(faces)):
    actor2 ← createLine(points[faces[j][0]], points[faces[j][1]])
    actor2 ← createLine(points[faces[j][0]], points[faces[j][2]])
    actor2 ← createLine(points[faces[j][1]], points[faces[j][2]])
end
           % Feature Extraction %
for i in range(len(target)):
    Distance ← GetGeodesicDistance(source, target[i])
    Path ← GetGeodesicPath(source, target[i])
    Connected ← GetConnectedEdges(target[i])
    Neighbors ← GetNeighbors(target[i])
    Geometricfeatures ← GetGeometricShape(target[i])
    FeatureVector.append(Path, Distance, Connectedfaces, Neighbors,
    Geometricfeatures)
    
```

end

return FeatureVector

E. INTERACTION RECOGNITION

At this stage of the proposed model, the interaction that has been performed in the input video has been recognized. After concatenating the different features extracted using full-body silhouettes and key body points, the feature vector has been fed into an LSTM model which is followed by a dense layer and a Softmax classifier. Hence, this section is subdivided into two sections: LSTM and Softmax classifier.

1) LSTM

Long Short-Term Memory (LSTM) [39] is a special type of Recurrent Neural Network (RNN) that is capable of learning long-term dependencies. The cell structure of LSTM is shown in Fig. 9. The working of LSTM has been described as follows:

1. The output value at a previous time h_{t-1} and the input value at the current time x_t are entered into the forget gate, and the output value of the forget gate f_t is obtained using (11).

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t]) \quad (11)$$

2. The output value at a previous time h_{t-1} and the input value at the current time x_t are also entered into the input gate. The output value i_t and the candidate cell state \check{c}_t of the input gate are obtained using (12) and (13).

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t]) \quad (12)$$

$$\check{c}_t = \tanh(W_c \cdot [h_{t-1}, x_t]) \quad (13)$$

3. The current cell state c_t is updated using (14).

$$c_t = f_t * c_{t-1} + i_t * \check{c}_t \quad (14)$$

4. The output and input are received as input values at the output gate at time t , and the output of the output gate o_t is obtained using (15).

$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t]) \quad (15)$$

5. Finally, the output value of LSTM is calculated by using the output of the output gate o_t and the state of the cell c_t , as shown;

$$h_t = o_t * \tanh(c_t) \quad (16)$$

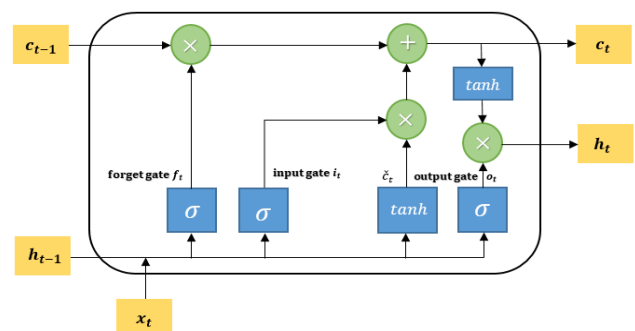


Figure 9. LSTM Cell Structure

2) SOFTMAX

The Softmax classifier has been used to recognize human interactions. The Softmax function is a popular choice for multiclass classification [40]. It is an activation function that computes the probabilities of all the classes based on the output of the fully connected layer. The probabilities are between the values of 0 and 1 and the normalized sum of these probabilities is always equal to 1. It uses cross-entropy loss. The Softmax output for each class is computed using (17).

$$SM(z_j) = \frac{e^{z_j}}{\sum_{i=1}^n e^{z_i}} \quad (17)$$

where z is the probability of each class, i is a vector of the inputs to the output layer, j is the set of the output units, and n is the total number of classes.

IV. EXPERIMENTAL SETUP AND RESULTS

This section explains the details of the experiments conducted to validate the proposed system. All the processing and experiments have been performed using Python 3.8 with Tensorflow 2.5.0 and Keras 2.4.3. A hardware system with an Intel Core i5 processor and a 64-bit Windows-10 has been used. The system has an 8 GB and 5 (GHz) CPU. The proposed system has been tested on three different datasets and the recognition accuracies for each interaction class have been computed in the form of their confusion matrices along with precision, sensitivity, and F1 scores. For further validation, the accuracies have been compared with those of other State-Of-The-Art (SOTA) methods. This section is further divided into two sections: dataset description and experimental results.

A. DATASETS

The three datasets that are used for experimentation are the SBU Kinect Interaction dataset [41], the NTU RGB+D dataset [42,43], and the ISR-UoL 3D social activity dataset [44]. Details of each dataset are given in the following subsections:

1) THE SBU KINECT INTERACTION DATASET

This dataset consists of RGB, depth, and skeletal information for various interactions performed by two people. The interactions have been recorded using Microsoft Kinect sensors in an indoor environment. It consists of eight interaction classes including *approaching*, *departing*, *kicking*, *punching*, *pushing*, *shaking hands*, *exchanging an object*, and *hugging*. The dataset has a total of 21 folders with subfolders for each interaction class performed by seven different actors. For interactions in which one person is performing and the other is receiving the action, there are two videos. The person performing the action in one video is receiving the action in the second video and vice versa. Videos have been segmented at the rate of 15 frames per second (fps). The sizes of both RGB and depth images are 649x480.

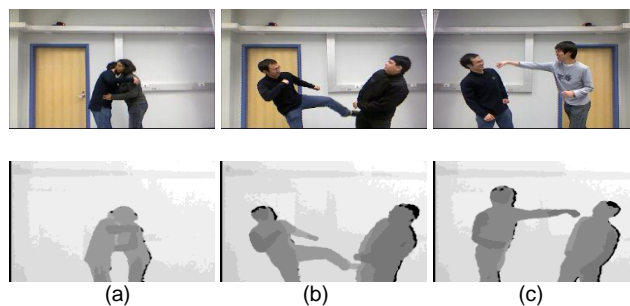


Figure 10. RGB and depth frames from the SBU Kinect interaction dataset. (a) hugging; (b) kicking; (c) punching.

2) THE NTU RGB+D DATASET

This dataset provides RGB, depth, and skeletal information. It consists of 60 classes, 11 of which are two-person interactions including *punching*, *kicking*, *pushing*, *pat on back*, *point finger*, *hugging*, *giving object*, *touch pocket*, *shaking hands*, *walking towards*, and *walking apart*. There are 48 videos for each interaction class. Each session has three sets of videos since each video has been recorded from three different viewpoints.

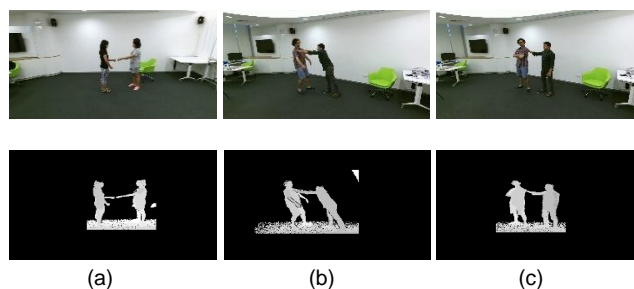
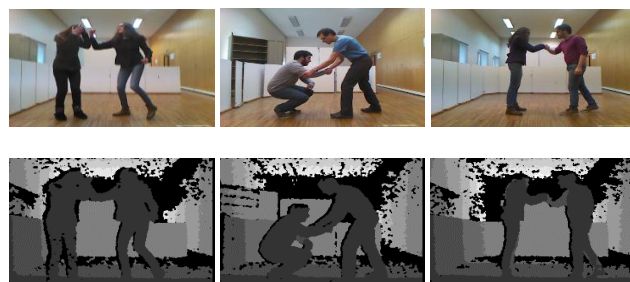


Figure 11. RGB and depth frames from the NTU RGB+D dataset. (a) giving object; (b) pushing; (c) pat on back.

3) THE ISR-UOL 3D SOCIAL ACTIVITY DATASET

This dataset also consists of RGB, depth, and skeletal information recorded using Kinect 2 sensor. In this dataset, some interactions are everyday interactions while others are assisted living interactions. There are a total of eight interactions including *shaking hands*, *hugging*, *help walk*, *help stand-up*, *fight*, *push*, *talk*, and *draw attention*. The actions are performed by four males and two females. There are ten sessions and each session contains all eight interactions. For each interaction, 24-bit RGB images, 8-bit and 16-bit resolution depth images, and the skeletal information of 15 joints are available. Each interaction is repeated over a period of 40–60 repetitions in one video.



(a) (b) (c)
Figure 12. RGB and depth frames from the ISR-UOL 3D dataset.
 (a) fight; (b) help stand; (c) shaking hands.

TABLE 2. A BRIEF SUMMARY OF THE DATASETS

Dataset	# videos	# classes	Modality
SBU Kinect Interaction	231	8	RGB, depth, skeletal
NTU RGB+D	528/2880	11/60	RGB, depth, skeletal
ISR UOL 3D Social Activity	80	8	RGB, depth, skeletal

B. EXPERIMENTS AND RESULTS

For validating the performance of the proposed system, different metrics have been used. The experimentation phase has been divided into two categories: classification accuracy of each class in terms of confusion matrix, precision,

sensitivity, and F1 score, and comparison of the proposed system with other state-of-the-art methods. The results for each stage are given in the following sub-sections.

1) INDIVIDUAL CLASS ACCURACY

The results of the proposed model's performance are given in the form of confusion matrices showing true positives, true negatives, false positives, and false negatives for each class individually. The confusion matrices for SBU Kinect Interaction, NTU RGB+D, and ISR-UoL 3D social activity datasets are given in Tables 3, 4, and 5 respectively. It can be observed, from the above-mentioned tables that the interaction classes of all three datasets achieved higher recognition rates with the mean accuracy rates of 91.63%, 90.54%, and 90.13% with the SBU Kinect Interaction, NTU RGB+D, and ISR-UoL 3D social activity datasets respectively.

TABLE 3. CONFUSION MATRIX OF INDIVIDUAL CLASSES OF THE SBU KINECT INTERACTION DATASET

Classes	Approaching	Departing	Kicking	Pushing	SH	Hugging	EAO	Punching
Approaching	0.90	0.06	0.00	0.00	0.02	0.02	0.00	0.00
Departing	0.05	0.91	0.02	0.00	0.00	0.02	0.00	0.00
Kicking	0.02	0.02	0.92	0.00	0.00	0.00	0.00	0.04
Pushing	0.00	0.00	0.04	0.89	0.00	0.00	0.03	0.04
SH	0.00	0.00	0.00	0.00	0.92	0.02	0.06	0.00
Hugging	0.02	0.02	0.00	0.00	0.00	0.94	0.02	0.00
EAO	0.00	0.00	0.00	0.00	0.03	0.02	0.92	0.03
Punching	0.00	0.00	0.00	0.03	0.00	0.00	0.04	0.93

Mean recognition accuracy = **91.63%**

*SH=Shaking hands, EAO=Exchanging an object.

TABLE 4. CONFUSION MATRIX OF INDIVIDUAL CLASS OF THE NTU RGB+D DATASET

Classes	Kicking	Pushing	PB	PF	Hugging	GO	TP	SH	WT	WA	PG
Kicking	0.91	0.03	0.02	0.00	0.00	0.00	0.01	0.00	0.00	0.03	0.00
Pushing	0.02	0.92	0.00	0.00	0.03	0.00	0.00	0.00	0.00	0.00	0.03
PB	0.00	0.04	0.88	0.04	0.00	0.00	0.04	0.00	0.00	0.00	0.00
PF	0.02	0.00	0.06	0.89	0.00	0.00	0.03	0.00	0.00	0.00	0.00
Hugging	0.00	0.04	0.00	0.00	0.92	0.02	0.00	0.02	0.00	0.00	0.01
GO	0.00	0.00	0.00	0.00	0.00	0.94	0.04	0.02	0.00	0.00	0.00
TP	0.00	0.00	0.04	0.02	0.00	0.04	0.90	0.00	0.00	0.00	0.00
SH	0.00	0.00	0.00	0.02	0.02	0.03	0.00	0.93	0.00	0.00	0.00
WT	0.04	0.00	0.00	0.00	0.00	0.02	0.00	0.03	0.90	0.00	0.00
WA	0.04	0.00	0.00	0.00	0.02	0.03	0.00	0.03	0.00	0.88	0.05
PG	0.00	0.04	0.00	0.00	0.02	0.00	0.00	0.00	0.02	0.03	0.89

Mean recognition accuracy = **90.54%**

*PB=Pat on back, PF=Point finger, GO=Giving Object, TP=Touch Pocket, SH=Shaking hands, WT=Walking towards, WA=Walking apart, PG=Punching.

TABLE 5. CONFUSION MATRIX OF INDIVIDUAL CLASS OVER THE ISR-UOL 3D SOCIAL ACTIVITY DATASET

Classes	SH	Hugging	Help walk	HSU	Fight	Push	Talk	DA
SH	0.90	0.02	0.03	0.03	0.02	0.00	0.00	0.00
Hugging	0.05	0.91	0.00	0.04	0.00	0.00	0.00	0.00
Help walk	0.05	0.04	0.89	0.02	0.00	0.00	0.00	0.00
HSU	0.03	0.05	0.04	0.88	0.00	0.00	0.00	0.02

Fight	0.01	0.00	0.00	0.00	0.92	0.04	0.03	0.00
Push	0.00	0.00	0.00	0.03	0.04	0.90	0.03	0.00
Talk	0.00	0.00	0.03	0.00	0.03	0.03	0.90	0.02
DA	0.03	0.02	0.00	0.00	0.00	0.00	0.04	0.91
Mean recognition accuracy = 90.13%								

*SH=Shaking hands, HSU=Help stand up, DA=Draw attention.

However, there is still some confusion between interaction classes that involve similar actions such as the *departing* and *approaching* interactions in the sports dataset. Similarly, shaking hands and exchanging an object interactions of the SBU Kinect Interaction dataset are confused with each other as shown in Table 3. Table 4 shows that the *pat on back* and *point finger* interactions of the NTU RGB+D datasets are confused with each other. As seen in Table 5, there is confusion between the *hugging* and *shaking hands* interaction of the ISR-UOL 3D social activity dataset.

Tables 6, 7, and 8 show the precision, sensitivity, and F1 scores of each class in SBU Kinect Interaction, ISR-UoL 3D social activity, and NTU RGB+D datasets respectively. The precision, sensitivity, and F1 scores of all the interaction classes for each dataset have been calculated as;

$$\text{Precision} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Positive}} \quad (18)$$

$$\text{Sensitivity} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Negative}} \quad (19)$$

$$\text{F1 score} = \frac{2(\text{Precision} \times \text{Recall})}{(\text{Precision} + \text{Recall})} \quad (20)$$

TABLE 6. MEASUREMENTS OF PRECISION, SENSITIVITY, AND F1 SCORES OVER THE SBU KINECT INTERACTION DATASET

Class	Precision	Sensitivity	F1 score
Approaching	0.88	0.90	0.90
Departing	0.90	0.91	0.90
Kicking	0.91	0.92	0.91
Pushing	0.89	0.89	0.90
Shaking hands	0.90	0.92	0.93
Hugging	0.93	0.94	0.93
Exchanging an object	0.91	0.92	0.91
Punching	0.91	0.93	0.92
Mean	0.90	0.91	0.91

TABLE 7. MEASUREMENTS OF PRECISION, SENSITIVITY, AND F1 SCORES OVER THE ISR-UOL 3D SOCIAL ACTIVITY DATASET

Class	Precision	Sensitivity	F1 score
Shaking hands	0.92	0.90	0.90
Hugging	0.92	0.91	0.91
Help walk	0.89	0.89	0.88
Help stand up	0.89	0.88	0.89
Fight	0.92	0.92	0.91
Push	0.90	0.90	0.88
Talk	0.91	0.90	0.91

Draw attention	0.92	0.91	0.92
Mean	0.91	0.90	0.90

TABLE 8. MEASUREMENTS OF PRECISION, SENSITIVITY, AND F1 SCORE OVER THE NTU RGB+D DATASET

Class	Precision	Sensitivity	F1 score
Kicking	0.91	0.91	0.90
Pushing	0.91	0.92	0.91
Pat on back	0.89	0.88	0.88
Point finger	0.91	0.89	0.91
Hugging	0.91	0.92	0.91
Giving object	0.93	0.94	0.94
Touch pocket	0.89	0.90	0.89
Shaking hands	0.95	0.93	0.93
Walking towards	0.90	0.90	0.91
Walking apart	0.87	0.88	0.89
Punching	0.89	0.89	0.90
Mean	0.90	0.90	0.91

2) COMPARISON WITH STATE-OF-THE-ART METHODS

In this section, the proposed method is compared with different methodologies adopted by researchers for HIR recognition from recent years. The action recognition accuracies of each evaluated methodology are used for comparison with the proposed system. Tables 9, 10, and 11 give the comparison of the proposed system with other state-of-the-art (SOTA) systems evaluated on SBU Kinect Interaction, NTU RGB+D, and ISR-UoL 3D social activity datasets respectively.

TABLE 9. COMPARISON WITH OTHER SOTA METHODS OVER THE SBU DATASET

Methods	Accuracy (%)
Joint features [41]	80.30
Body parts contrast mining [45]	86.9
Joint Features [46]	90.3
Deep LSTM [47]	90.41
STA-LSTM [48]	91.51
Proposed Method	91.63 (RGB+D) 89.53 (Depth only) 88.24 (RGB only)

TABLE 10. COMPARISON WITH OTHER SOTA METHODS OVER THE NTU RGB+D DATASET

Methods	Accuracy (%)
---------	--------------

geometric features [49]	70.26
ensemble TS-LSTM v2 [50]	74.60
STA-LSTM [48]	81.2
multitask deep learning [51]	85.5
pair wise features [52]	88.6
	90.54 (RGB+D)
Proposed Method	88.12 (Depth only)
	87.63 (RGB only)

TABLE 11. COMPARISON WITH OTHER SOTA METHODS OVER THE ISR-UOL 3D SOCIAL ACTIVITY DATASET

Methods	Accuracy (%)
probabilistic merging of skeletal features [44]	85.1
multimodal feature level fusion [53]	85.12
statistical and geometrical features [54]	85.56
skeletal data [55]	87
	90.13 (RGB+D only)
Proposed Method	88.14 (Depth only)
	86.83 (RGB only)

V. DISCUSSION

The proposed system is a complete HIR solution that should be applicable to many real-world problems involving the tasks of human behavior monitoring, security, surveillance, and managing smart homes. It is designed for RGB+D datasets but can also be used with RGB only or depth only datasets using only one stream of the proposed CNN model and skipping the model concatenation stage.

Each step from the preprocessing stages to the classification stage contributes to the improved performance achieved by the system. The proposed feature extraction method successfully extracts robust features, which in turn, play a critical role in accurate classification of the interactions. Using two CNN models for training RGB and depth images separately and then concatenating the models gives better results than those obtained by concatenating both RGB and Depth images first and then training the 4-dimensional images using only one CNN model. Moreover, since all three datasets contain video sequences, the LSTM-based classification step gives accurate results.

Despite yielding good results, the proposed system is not without limitations. The proposed 2D leaf skeleton model for the detection of key body points can only extract six key points so far. However, better accuracies can be achieved if more key points are identified and their features are extracted. Moreover, the proposed system is very extensive and computationally expensive. The time complexity of the proposed system is shown in Table 12.

TABLE 12. TIME COMPLEXITY OF THE PROPOSED SYSTEM

Dataset	Execution time (s)
SBU Kinect Interaction	4795.71
NTU RGB+D	6560.05

ISR UOL 3D Social activity 3671.32

VI. CONCLUSION

In this paper, a novel HIR framework had been proposed that uses both machine learning and deep learning techniques for feature extraction from 2D human silhouettes and 3D meshes. Using efficiently segmented silhouettes of the humans from images, multiple features using full-body silhouettes and key body points from their corresponding 3D meshes have been extracted. The features have then been fed to an LSTM and a Softmax-based classifier. The proposed system achieved average accuracies of 91.63%, 90.54%, and 90.13% with the SBU Kinect Interaction, NTU RGB+D, and ISR-UoL 3D social activity datasets respectively.

In the future, the authors plan to shift their focus to the task of human-object interaction recognition and investigate new features and modeling techniques for better classification results.

REFERENCES

- O. Aran and D. Gatica-Perez, "One of a kind: Inferring personality impressions in meetings," in *Proc. on ICMI (ACM)*, pp. 11–18, 2013.
- A. Jalal, N. Sharif, J. T. Kim, and T.-S. Kim, "Human activity recognition via recognized body parts of human depth silhouettes for residents monitoring services at smart homes," *Indoor and Built Environment*, vol. 22, pp. 271–279, 2013.
- S.U. Khan, T. Hussain, A. Ullah, and S. W. Baik, "Deep-ReID: deep features and autoencoder assisted image patching strategy for person re-identification in smart cities surveillance," *Multimedia Tools and Applications*, 2021.
- G.H. Liu, J.Y. Yang, and Z. Li, "Content-based image retrieval using computational visual attention model," *Pattern Recognition*, vol. 48, pp. 2554–2566, 2015.
- S. Sempena, N.U. Maulidevi, and P.R. Aryan, "Human action recognition using dynamic time warping," in *Proc. on ICEEI (IEEE)*, pp. 1–5, 2011.
- S. U. Khan, I.U. Haq, S. Rho, S.W. Baik, and M.Y. Lee, "Cover the Violence: A Novel Deep-Learning-Based Approach Towards Violence-Detection in Movies," *Applied Sciences*, vol. 9, no. 22, pp. 4963, 2019.
- A. Jalal, M. Batoool, and K. Kim, "Stochastic recognition of physical activity and healthcare using tri-axial inertial wearable sensors," *Applied Sciences*, vol. 10, no. 20, pp. 7122, 2020.
- A. Jalal, M. A. Quaid, S. B. Tahir, and K. Kim, "A study of accelerometer and gyroscope measurements in physical life-log activities detection systems," *Sensors*, vol. 20, no. 22, pp. 6670, 2020.
- A. Jalal, M. Batoool, and K. Kim, "Sustainable Wearable System: Human Behavior Modeling for Life-logging Activities Using K-Ary Tree Hashing Classifier," *Sustainability*, vol. 12, no. 24, pp. 10324, 2020.
- M. Javeed, A. Jalal, and K. Kim, "Wearable sensors based exertion recognition using statistical features and random forest for physical healthcare monitoring," in *Proc. on IEEE IBCAST*, pp. 512–517, 2021.
- S. U. Khan and S.W. Baik, "MPPIF-Net: Identification of Plasmodium Falciparum Parasite Mitochondrial Proteins Using Deep Features with Multilayer Bi-directional LSTM," *Processes*, vol. 8, no. 6, pp. 725, 2020.
- A. Shehzad, A. Jalal, and K. Kim, "Multi-Person Tracking in Smart Surveillance System for Crowd Counting and

- Normal/Abnormal Events Detection,” in *Proc. on Applied and Engineering Mathematics*, pp. 163-168, 2019.
13. P. Mahwish, A. Jalal, and K. Kim, “Hybrid algorithm for multi people counting and tracking for smart surveillance,” in *Proc. on IEEE IBCAST*, pp. 530-535, 2021.
 14. N. Khalid, M. Gochoo, A. Jalal, and K. Kim, “Modeling Two-Person Segmentation and Locomotion for stereoscopic Action Identification: A Sustainable Video Surveillance System,” *Sustainability*, vol. 13, no. 2, pp. 970, 2021.
 15. M. Asadi-Aghbolaghi, A. Clapes, M. Bellantonio, H.J. Escalante, V.P. Lopez, X. Baro, I. Guyon, S. Kasaei, and S. Escalera, “A survey on deep learning based approaches for action and gesture recognition in image sequences,” in *Proc. on Automatic Face & Gesture Recognition*, pp. 476-483, 2017.
 16. A. Jalal, J. T. Kim, and T.-S. Kim, “Human activity recognition using the labeled depth body parts information of depth silhouettes,” in *Proc. on Sustainable Healthy Buildings*, pp. 1-8, 2012.
 17. D. Rado, A. Sankaran, J. Plasek, D. Nuckley, and D.F. Keefe. “A Real-Time Physical Therapy Visualization Strategy to Improve Unsupervised Patient Rehabilitation,” in *Proc. on IEEE Visualization*, 2009.
 18. M.H. Khan, M. Zöllner, M.S. Farid, and M. Grzegorzec. “Marker-Based Movement Analysis of Human Body Parts in Therapeutic Procedure,” *Sensors*, vol. 20, no. 11, pp. 3312, 2020.
 19. C.C. Chen, C.Y. Liu, S.H. Ciou, S.C. Chen, and Y.L. Chen. “Digitized Hand Skateboard Based on IR-Camera for Upper Limb Rehabilitation,” *J. Med. Syst.* vol. 41, pp. 36, 2017.
 20. R.E. Mayagoitia, A.V. Nene, and P.H. Veltink. “Accelerometer and rate gyroscope measurement of kinematics: an inexpensive alternative to optical motion analysis systems,” *J Biomech*, vol. 35, no. 4, pp. 537-542, 2002.
 21. N. Ganter, A. Krüger, M. Gohla, K. Witte, and J. Edelmann-Nusser. “Applicability of a full body inertial measurement system for kinematic analysis of the discus throw,” in *Proc. of the International Society of Biomechanics in Sports*, 2010.
 22. F. Eckardt, A. Münz, and K. Witte. “Application of a full body inertial measurement system in dressage riding.” *J Equine Vet Sci*, vol. 34, pp. 1294-1299, 2014.
 23. F.A. de Magalhaes, G. Vannozzi, G. Gatta, and S. Fantozzi. “Wearable inertial sensors in swimming motion analysis: a systematic review,” *J Sports Sci*, vol. 33, no. 7, pp. 732-745, 2015.
 24. M.I. Mokhlespour Esfahani, O. Zobeiri, B. Moshiri, R. Narimani, M. Mehravar, E. Rashedi, and M. Parnianpour. “Trunk Motion System (TMS) Using Printed BodyWorn Sensor (BWS) via Data Fusion Approach,” *Sensors*, vol. 17, pp. 112, 2017.
 25. Y. Tian, L. Cao, Z. Liu, and Z. Zhang, “Hierarchical filtered motion for action recognition in crowded videos,” *IEEE Transactions on Systems, Man, and Cybernetics*, vol. 42, pp. 313-323, 2012.
 26. S. A. Rizwan, A. Jalal, M. Gochoo, and K. Kim, “Robust active shape model via hierarchical feature extraction with SFS-optimized convolution neural network for invariant human age classification,” *Electronics*, vol. 10, no. 4, pp. 465, 2021.
 27. M. Khan, M. Schneider, M. Farid, and M. Grzegorzec, “Detection of Infantile Movement Disorders in Video Data Using Deformable Part-Based Model,” *Sensor*, vol. 18, no. 10, pp. 3202, 2018.
 28. M. H. Khan, J. Helsen, M.S. Farid, and M. Grzegorzec, “A computer vision-based system for monitoring Vojta therapy,” *Int. J. Med. Inform. Vol. 113*, pp. 85-95, 2018.
 29. M. Javeed, M. Gochoo, A. Jalal, and K. Kim, “HF-SPHR: Hybrid features for sustainable physical healthcare pattern recognition using deep belief networks,” *Sustainability*, vol. 13, no. 4, pp. 1699, 2021.
 30. M. Gochoo, I. Akhter, A. Jalal, and K. Kim, “Stochastic remote sensing event classification over adaptive posture estimation via multifused data and deep belief network,” *Remote Sensing*, vol. 13, no. 5, pp. 912, 2021.
 31. A. Jalal, A. Ahmed, A. Rafique, and K. Kim “Scene Semantic recognition based on modified Fuzzy c-mean and maximum entropy using object-to-object relations,” *IEEE Access*, vol. 9, pp. 27758-27772, 2021.
 32. A. Jalal, I. Akhtar, and K. Kim, “Human Posture Estimation and Sustainable Events Classification via Pseudo-2D Stick Model and K-ary Tree Hashing,” *Sustainability*, vol. 12, no. 23, pp. 9814, 2020.
 33. A. Jalal, N. Khalid, and K. Kim, “Automatic recognition of human interaction via hybrid descriptors and maximum entropy markov model using depth sensors,” *Entropy*, vol. 22, no. 8, pp. 817, 2020.
 34. C. Rother, V. Kolmogorov, and A. Blake, “GrabCut: Interactive Foreground Extraction Using Iterated Graph Cuts,” *ACM transactions on Graphics*, vol. 23, no. 3, pp. 309-314, 2004.
 35. S. Saito, Z. Huang, R. Natsume, S. Morishima, A. Kanazawa, and H. Li, “PIFu: Pixel-Aligned Implicit Function for High-Resolution Clothed Human Digitization,” in *Proc. on ICCV*, pp. 2304-2314, 2019.
 36. J. Sun, M. Ovsjanikov, and L. Guibas (2009). “A Concise and Provably Informative Multi-Scale Signature-Based on Heat Diffusion,” *Computer Graphics Forum*, vol. 28, pp. 1383-1392, 2009.
 37. R. Litman and A. M. Bronstein., “Spectral descriptors for deformable shape Correspondence,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 1, pp. 171-180, 2014.
 38. K. Simonyan and A. Zisserman, “Very Deep Convolutional Networks for Large-Scale Image Recognition,” in *Proc. on Learning Representations*, 2015.
 39. S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural Computation*, vol. 9, no. 8, pp. 1735-1780, 1997.
 40. K. Banerjee, V. Prasad, R. Gupta, K. Vyas, H. Anushree, B. Mishra, “Exploring Alternatives to Softmax Function,” *CoRR abs/2011.11538*, 2020.
 41. K. Yun, J. Honorio, D. Chattopadhyay, T.L. Berg, and D. Samaras. “Two-person interaction detection using body-pose features and multiple instance learning,” in *Proc. on Computer Vision and Pattern Recognition Workshops*, pp. 28-35, 2012.
 42. A. Shahroudy, J. Liu, T. Ng, G. Wang, “NTU RGB+D: A Large Scale Dataset for 3D Human Activity Analysis,” in *Proc. on Computer Vision and Pattern Recognition*, 2016.
 43. J. Liu, A. Shahroudy, M. Perez, G. Wang, L. Duan, and A. Kot, “NTU RGB+D 120: A Large-Scale Benchmark for 3D Human Activity Understanding,” in *Proc. on Pattern Analysis and Machine Intelligence*, 2019.
 44. C. Coppola, D. R. Faria, U. Nunes, and N. Bellotto, “Social Activity Recognition Based on Probabilistic Merging of Skeleton Features with Proximity Priors from RGB-D Data,” in *Proc. on Intelligent Robots and Systems*, 2016.
 45. Y. Ji, G. Ye, and H. Cheng. “Interactive Body Part Contrast Mining for Human Interaction Recognition,” in *Proc. on Multimedia and Expo Workshops*, 2014.
 46. T. Huynh-The, O. Banos, B. Le, D. Bui, S. Lee, Y. Yoon, and, T. Le-Tian, “PAM-Based Flexible Generative Topic Model for 3D Interactive Activity Recognition,” in *Proc. on Advanced Technologies for Communications*, pp. 117-122, 2015.
 47. W. Zhu, C. Lan, J. Xing, W. Zeng, Y. Li, L. Shen, and X. Xie, “Co-occurrence Feature Learning for Skeleton based Action Recognition using Regularized Deep LSTM Networks,” in *Proc. on Artificial Intelligence*, pp. 3697-3703, 2016.
 48. S. Song, C. Lan, J. Xing, W. Zeng, and J. Liu, “An End-to-End Spatio-Temporal Attention Model for Human Action Recognition from Skeleton Data,” in *Proc. on Artificial Intelligence*, 2016.
 49. S. Zhang, X. Liu, and J. Xiao. “On geometric features for skeleton-based action recognition using multilayer lstm net-

works,” in *Proc. on Application of Computer Vision WACV*, pp. 148–157, 2017.

50. Lee, D. Kim, S. Kang, and S. Lee, “Ensemble Deep Learning for Skeleton-Based Action Recognition Using Temporal Sliding LSTM Networks,” in *Proc. on Computer Vision*, 2017.
51. D.C. Luvizon, D. Picard, H. Tabia, “2D/3D pose estimation and action recognition using multitask deep learning,” in *Proc. on Computer Vision and Pattern Recognition*, 2018.
52. B. Li, M. He, X. Cheng, Y. Chen, and Y. Dai, “Skeleton based action recognition using translation-scale invariant image mapping and multi-scale deep CNN,” in *Proc. on Multimedia & Expo Workshops*, 2017.
53. M. Ehatisham-Ul-Haq, A. Javed, M.A. Azam, H.M.A. Malik, A. Irtaza, I.H. Lee, and M.T. Mahmood, “Robust human activity recognition using multimodal feature-level fusion,” *IEEE Access*, vol. 7, pp. 60736–60751, 2019.
54. C. Coppola, S. Cosar, D.R. Faria, and N. Bellotto, “Automatic detection of human interactions from RGB-D data for social activity classification,” in *Proc. on Robot and Human Interactive Communication*, 2017.
55. A. Manzi, L. Fiorini, R. Limosani, P. Dario, and F. Cavallo, “Two-person activity recognition using skeleton data,” *IET Computer. Vision*, vol. 12, pp. 27-35, 2018.



MANAHIL WAHEED is currently enrolled in the MS Data Science program at Air University, Islamabad. She received a BS in electronics engineering from International Islamic University in 2018. Her research interests include digital image processing, data science, and artificial intelligence.

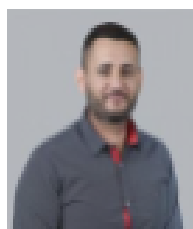


AHMAD JALAL is currently an Associate Professor from the Department of Computer Science and Engineering, Air University, Pakistan. He received his Ph.D. degree in the Department of Biomedical Engineering at Kyung Hee University, Republic of Korea. Now, he was working as a Post-doctoral Research fellowship at POSTECH. His research interest includes multimedia contents and artificial intelligence.



MOHAMMED ALARFAJ is an Assistant Professor of Electrical Engineering at King Faisal University. He was awarded the B.S., M. Eng. and a Ph.D. in Electrical and Computer Engineering from Oregon State University in 2011, 2014, and 2019, respectively. He is now the head of the Electrical Engineering department at King Faisal University. His current research interests include MIMO, mmWave and wireless communications, signal

processing, and applications in wireless communication and sensor networks.



YAZEED YASIN GHADI received his Ph.D. in Electrical and Computer Engineering from Queensland University. His dissertation on developing novel hybrid plasmonic photonic on-chip biochemical sensors received the Sigma Xi best Ph.D. thesis award. He is currently an assistant professor of Software engineering at Al Ain University. He was a postdoc researcher at Queensland University before joining Al Ain. His current research is on developing novel electro-acoustic-optic neural interfaces for large-scale high-resolution electrophysiology and distributed optogenetic stimulation. Yazeed has published more than 25 peer-reviewed journal and conference papers and he holds three pending patents. He is the recipient of several awards.



TAMARA AL SHLOUL is an assistant professor (humanities) at Al Ain University. She has vast experience of teaching education and humanities courses, along with experience in school supervision, thinking skills, and higher education improvement ability. Her research interest includes teacher socialization and professional development.



SHAHARYAR KAMAL received his M.S. degree in Computer Engineering from Mid Sweden University, Sweden and Ph.D. degree in the Department of Radio and Electronics Engineering at Kyung Hee University, Republic of Korea. He is currently an Assistant Professor from the Department of Computer Science and Engineering, Air University, Pakistan. His research interest includes advanced wireless communication, image and signal processing.



DONG-SEONG KIM received his Ph. D degree in EECS, Seoul National University, Seoul, Korea. From 1994 to 1998, he worked as a full-time Researcher in ERC-ACI at Seoul National University. From September 2000 to December 2001, he worked as a part-time Lecturer in department of information and communication at Dong-Guk University. He is currently Professor in Department of IT Convergence Engineering, School of Electronic Engineering, Kumoh National Institute of Technology, Gumi, Korea. He has other responsibilities as Director and Head of Convergence Technology Institute, ICT-CRC, Korea. His main research interests include Industrial networked control system, Fieldbus and Real-time systems and wired/wireless military networks.