

Recognition of human actions using motion history information extracted from the compressed video[☆]

R. Venkatesh Babu^{a,*}, K.R. Ramakrishnan^{b,1}

^aCentre for Quantifiable Quality of Service in Communication Systems, Norwegian University of Science and Technology, O.S. Bragstads plass 2E, N7491 Trondheim, Norway

^bDepartment of Electrical Engineering, Indian Institute of Science, Bangalore 560012, India

Received 27 May 2003; received in revised form 19 November 2003; accepted 20 November 2003

Abstract

Human motion analysis is a recent topic of interest among the computer vision and video processing community. Research in this area is motivated by its wide range of applications such as surveillance and monitoring systems. In this paper we describe a system for recognition of various human actions from compressed video based on motion history information. We introduce the notion of quantifying the motion involved, through what we call Motion Flow History (MFH). The encoded motion information readily available in the compressed MPEG stream is used to construct the coarse Motion History Image (MHI) and the corresponding MFH. The features extracted from the static MHI and MFH compactly characterize the spatio-temporal and motion vector information of the action. Since the features are extracted from the partially decoded sparse motion data, the computational load is minimized to a great extent. The extracted features are used to train the KNN, Neural network, SVM and the Bayes classifiers for recognizing a set of seven human actions. The performance of each feature set with respect to various classifiers are analyzed.

© 2003 Elsevier B.V. All rights reserved.

Keywords: Action recognition; Compressed domain; Content-based retrieval; Feature extraction; Motion history; Video indexing

1. Introduction

Event detection and human action recognition have gained more interest of late, among video processing community because they find various applications in automatic surveillance, monitoring systems [2], video indexing and retrieval, robot motion, human–computer interaction and segmentation [28,30]. One of the important applications of human action recognition is the automatic indexing of video sequences, while most of the multimedia documents available nowadays are in the MPEG [21] compressed form, to facilitate easy storage and transmission, majority of the existing techniques for human action recognition are pixel domain

based [35,8,27,5,32,13,1,26] which are computationally very expensive. Hence, it would be efficient if the classification is performed in the MPEG compressed domain without having to completely decode the bit-stream and subsequently perform classification in the pixel domain. This calls for techniques, which can use information available in the compressed domain such as motion vectors and DCT coefficients.

In the recent past, we reported a technique for recognizing human actions from compressed video using Hidden Markov Model (HMM) [3], where the time-series features used for training the HMM are directly extracted from the motion vectors corresponding to each frame of the video. Though this approach has proven its ability to classify the video sequences, the extracted time series features are not suitable for other efficient classifiers such as *K*-nearest neighbors (KNN), Neural networks, SVM and Bayes.

In this paper we propose a technique for building coarse Motion History Image (MHI) and Motion Flow History (MFH) from the compressed video and extract features from

[☆] An earlier, brief version of this paper has appeared in the Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing 2003 (ICASSP-03) [4].

* Corresponding author. Tel.: +47-7359-2746; fax: +47-7359-6973.

E-mail addresses: venkat@q2s.ntnu.no (R. Venkatesh Babu), krr@ee.isc.ernet.in (K.R. Ramakrishnan).

¹ Tel.: +91-80-293-2441; fax: +91-80-360-0444.

these static motion history information for characterizing human actions. The MHI gives the temporal information of the motion at the image plane, whereas the MFH quantifies the motion at the image plane. The features extracted from MHI and MFH were used to train KNN, Bayes, Neural Network and SVM classifiers for recognizing a set of seven human actions. The encoded motion information available in the MPEG video is exploited for constructing the coarse MHI and MFH. These MHI and MFH represent the human action in a very compact manner. Though the motion information extracted from each frame of the compressed video is very sparse, they are sufficient to construct the MHI and MFH for representing the actions.

This work is motivated by a technique proposed by Davis and Bobick [11] where a view-based approach is used to recognize actions. They have presented a method for recognition of temporal templates. A temporal template is a static image where the value at each point is a function of the motion properties at the corresponding spatial location in an image sequence. The actions were represented by the cumulative motion images called Motion Energy Image (MEI) and MHI. MEI indicates where the motion has occurred in the image plane, whereas MHI indicates the recency of motion using intensity. For recognition, the Hu moments [16], obtained from the templates are known to yield reasonable shape discrimination in a translation and scale invariant manner. Extracted Hu moments are matched using a nearest neighbor approach against the examples of given motions already learned. This work was extended by Rosales [27] using various classification approaches like KNN and Bayes with dimensionality-reduced representation of actions.

This paper is organized as follows: Section 2 gives a brief description of the related work from the literature. Section 3 describes the basics of MPEG video compression and the overview of the proposed work. Section 4 explains about the construction of MHI and MFH. Feature extraction procedures are explained in Section 5. Section 6 presents the classification results. The performance of feature set is analyzed and compared in Section 7. Section 8 concludes the paper.

2. Related work

In this section we will give a brief description of works related to human motion and gesture recognition. The recognition of human motion can be broadly classified into the following two: (i) state-space based and (ii) template matching based approaches.

2.1. State-space based approaches

State-space approach uses time-series features obtained from the image sequences for recognition. The widely used state-space model for activity recognition is HMM due to its

success in the speech community. The first attempt to use HMM for activity recognition is done by Yamato et al. [35], where discrete HMMs are used for recognition of six tennis strokes. In their approach time sequential images expressing human actions are transformed to an image feature vector sequence by extracting mesh [33] feature vector from each image. The mesh features are extracted from a binarized image obtained after subtracting the background image from original image by applying a suitable threshold. The drawbacks of this method are that it is sensitive to position displacement, noise, and also exhibits poor performance if the training and test subjects are different.

The gesture recognition work by Darrell and Pentland [9] uses time-warping technique for recognition which is closely related to HMM. On similar lines, dynamic time warping is used in Ref. [6] to match an input signal to a deterministic sequence of states. Starner and Pentland [31] used HMMs to recognize a limited vocabulary of American Sign Language (ASL) sentences. Here, they used a view based approach with a single camera to extract two-dimensional (2D) features as input to HMMs.

In the work by Bregler [8], this classification problem has been approached from a statistical view point. For each pixel in the image, the spatio-temporal image gradient and the color values are represented as random variables. Then the blob hypothesis is used wherein each blob is represented with a probability distribution over coherent motion, color and spatial support regions.

Recently Ivanov and Bobick [17] proposed a method, which combines statistical techniques used for detecting primitive component of an activity with syntactic recognition of process structure. In this approach the recognition problem is divided into two levels: (i) the lower level detection of primitive components of activity followed by (ii) the syntactic recognition of the primitive features using a stochastic context-free grammar parsing mechanism. Another HMM based human activity recognition method is reported by Psarrou et al. [25]. Here the recognition is based on learning prior and continuous propagation of density models of behavior patterns. Ng et al. [22] proposed a real-time gesture recognition system incorporating hand posture and hand motion. The recognition is done with HMM and recurrent neural networks (RNN).

There are few works reported in literature which use neural networks for gesture recognition [19,7]. Boehm et al. [7] used Kohonen Feature Maps (KFM) [18] for recognizing dynamic gestures. Oliver et al. [23] proposed a system for modeling and recognizing human behaviors in a visual surveillance task. This system segments the moving objects from the background and a kalman filter tracks the object's features such as location, coarse shape, color and velocity. These features are used for modeling the behavior patterns through training HMMs and coupled HMMs (CHMM), which are used for classifying the perceived behaviors. Based on the above-mentioned work Madabhushi

and Aggarwal [20] presented a system for recognition of human action by tracking the head of the subject in an image sequence. The difference in centroids of the head over successive frames form their feature vector. The human actions are modeled based on the mean and covariance of the feature vector. Here detection and segmentation of the head is done manually.

Apart from the above mentioned pixel domain state-space based approaches, recently we have proposed a technique for recognizing human actions using HMM in compressed domain framework [3]. Here the time series features from the MPEG video are extracted from the readily available motion vectors of each inter coded frame. Totally seven actions were considered for recognition (walk, run, jump, bend up, bend down, twist right and twist left). A discrete HMM for each action is trained with the corresponding MPEG video sequences. The recognition of a given action is achieved by feeding the test sequence to all the trained HMMs and employing a likelihood-based measure. The performance of the system for three types of motion-based features were compared.

2.2. Template matching based approaches

One of the earlier works using this approach is found in the work done by Polana and Nelson [24], where the flow information is used as feature. They compute the optical flow fields [15] between consecutive frames and divide each frame into a spatial grid and sum the motion magnitude to get the high dimension feature. Here they assume that the human motion is periodic. The final recognition is performed using nearest neighbor algorithm. Davis and Bobick [11,5] presented a real-time approach for representing human motion using compact MHIs in pixel domain. Here, the recognition of 18 aerobic exercises was achieved by statistically matching the higher order moment based feature extracted from the MHI. The limitations of the above method are related to the ‘global image’ feature calculations and specific label based recognition. To overcome these limitations the author extended the previous approach with a mechanism to compute dense local motion vector field directly from the MHI for describing the movement [10]. For obtaining the dense motion, the MHI is represented at various pyramid levels to tackle multiple speeds of motion. These hierarchical MHIs are not directly created from the original MHI, but through the pyramid representation of the silhouette images. This indirect way of generating MHI pyramid increases the computational load. The resulting motion is characterized by a polar histogram of motion orientation. Rosales [27] use these motion energy and MHIs [11] for obtaining the spatial location and the temporal properties of human actions from raw video sequences. From these motion energy and MHIs, a set of Hu-moment [16] features that are invariant to translation, rotation and scaling are generated. Using principal component analysis, the dimension of

the Hu-moment space is reduced in a statistically optimal way. The recognition performances were evaluated for the following three classifiers namely KNN, Gaussian and mixtures of Gaussian. All the above mentioned techniques process the data in the pixel domain, which is computationally very expensive.

3. System overview

The objective of our work is to rapidly process the video stored in MPEG format, without full-frame decompression, for recognizing human actions. Here we are using the motion vector data, which is easily extractable from the MPEG video bit-stream for our recognition task. Though we have used MPEG-1 video, our algorithm is easily extendable to MPEG-2 or the recent MPEG-4 video streams. To begin with, we briefly describe the relevant parts of MPEG video compression standard.

The MPEG-1 video defines three types of coded pictures namely: intracoded (*I*-frames); predicted (*P*-frames); and bidirectionally predicted (*B*-frames). These pictures are organized into sequences of groups of pictures (GOP) in MPEG video streams. A GOP must start with an *I*-frame, followed by any number of *P*- and *B*-frames. The *I*- and *P*-frames are referred as anchor frames. The *B*-frames appear between each pair of consecutive anchor frames in the GOP and before the *I*-frame of the next GOP. Fig. 1 shows the typical GOP structure that is used in our work with 12 frames in a GOP.

Each frame of the video is divided into non-overlapping macroblocks. For video coding in 4:2:0 format [29], each macroblock consists of six 8×8 pixel blocks: four luminance (*Y*) blocks and two chrominance (*Cb*, *Cr*) blocks. Each macroblock is either intra coded or inter coded. An *I*-frame is completely intra coded. Here every 8×8 pixel block in the macroblock is transformed to frequency domain using the discrete cosine transformation (DCT). The resulting 64 DCT coefficients are then quantized (lossy) and entropy (run length and Huffman, lossless) encoded to achieve compression. Since coding of *I*-frame does not refer to any other video frames, it can be decoded independently and thus provides access points for fast random access to the compressed video.

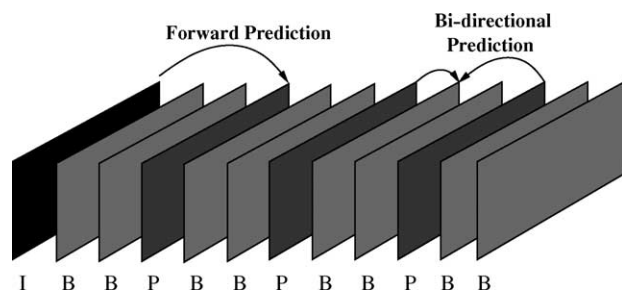


Fig. 1. A typical MPEG GOP structure.

Each *P*-frame is predictively encoded with reference to its previous anchor frame (i.e. previous *I*- or *P*-frame). For each macroblock (16×16 pixel block) in *P*-frame, a local region in the anchor frame is searched for a good match in terms of difference in intensity. If a good match is found, the macroblock is represented by a motion vector to the position of the match together with the DCT encoding of the residue (i.e. difference) between the macroblock and its match. The DCT coefficients of the residue are quantized and entropy coded while the motion vector is differentially and entropy coded with respect to its neighboring motion vector. This is known as *forward motion compensation*, and such macroblocks are referred as *intercoded macroblocks*. If good match cannot be found, the macroblock is intracoded like the macroblocks of *I*-frame. Since the residue of an intercoded macroblock can be coded with fewer bits, it has better compression gain compared to the intracoded macroblock. In our work the motion vectors extracted from these *P*-frames are used for recognizing human actions. The coding of *P*-frame is illustrated in Fig. 2. To achieve further compression, *B*-frames are bidirectionally predictively encoded with forward and/or backward motion compensation referenced to its closest past and/or future *I*- and/or *P*-frames. Since *B*-frames are not used as reference frames, they can accommodate more distortion, and thus, higher compression gain compared to *I*- or *P*-frames.

The overview of the proposed system is shown in Fig. 3. First the motion vectors are extracted from the compressed video by partially decoding the MPEG video bit-stream. This partial decoding is far less expensive compared to the full decoding. Since the sampling rate of the video is normally very high (typically 25 frames/s) compared to human motion dynamics, it is not necessary to extract the motion vectors from all the frames. Hence we have used

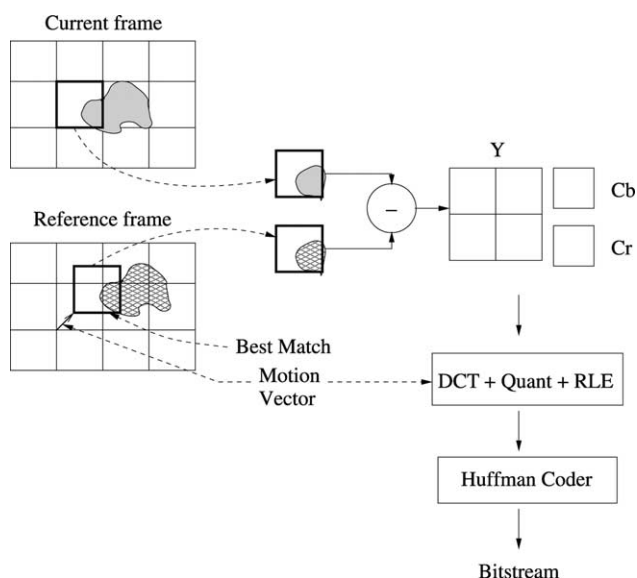


Fig. 2. Coding of *P*-frame.

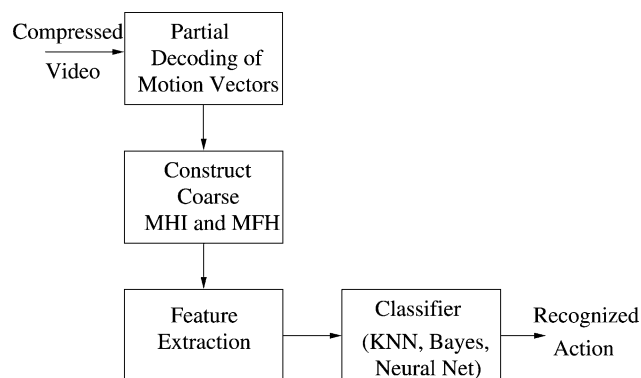


Fig. 3. Overview of the proposed system.

only the motion vectors obtained from the predictive (*P*) frames for constructing the coarse MHI and MFH. As motion vectors are usually noisy, the coarse MHI and MFH are constructed after removing the noisy motion vectors. The constructed coarse MHI and MFH are at macroblock resolution and not at pixel resolution. Hence the size of the MHI and MFH are sixteen times smaller than the original frame size i.e. 16^2 times smaller in terms of number of pixels. In feature extraction phase, various features are extracted from the constructed coarse MHI and MFH, which hold the temporal and motion information of the video sequence. The features based on projection profiles and centroids are extracted from MHI. Affine features and motion vector histogram based features are obtained from the MFH. These features are finally fed to the classifiers such as KNN, Neural network and Bayes for recognizing the action.

4. Representation of action using MHI and MFH

Since we are interested in analyzing the motion occurring in a given window of time, we need a method that allows us to capture and represent motion directly from the video sequence. Such static representations are called MEIs, MHIs and MFH. They are functions of the observed motion parameters at the corresponding spatial image location in the video sequence.

MEI is basically a cumulative binary image with only spatial, and no temporal details of the motion involved. It answers the question ‘where did the motion occur?’. MEI can be obtained by binarizing the MHI. The MHI is a cumulative gray scale image incorporating the spatial as well as the temporal information of the motion [11]. MHI points to, ‘where and when did the motion occur?’. It does not convey any information about the direction and magnitude of the motion. MFH gives the information about the extent of the motion at each macroblock (‘where and how much did the motion occur?’). In case of occlusion, the old motion information is over-written by the new reliable motion information.

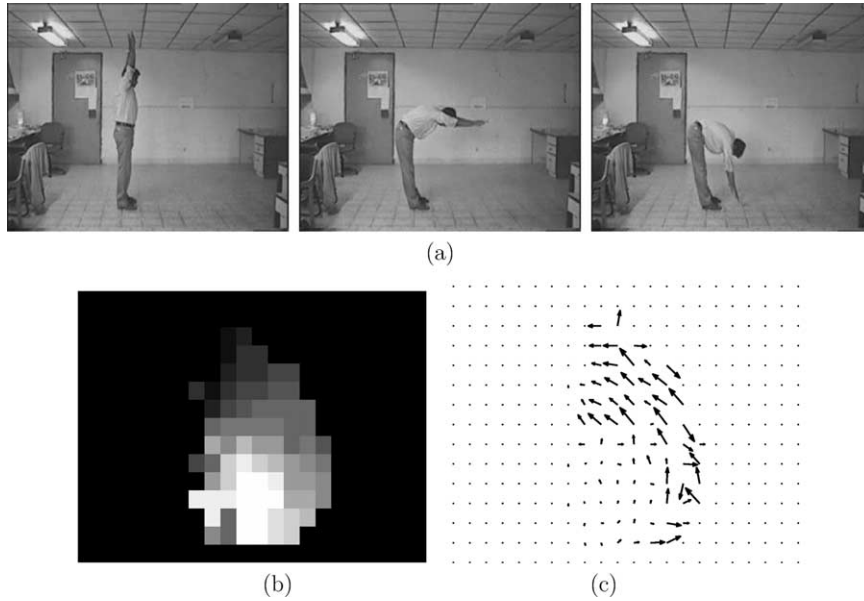


Fig. 4. (a) Key-frames of bend-down sequence and corresponding coarse (b) MHI (c) MFH.

Since it is computationally very expensive to decode the full video, we use the readily available encoded motion information in MPEG bit-stream for constructing the coarse MHI and MFH. The motion vectors not only indicate the blocks under motion but also gives the information regarding magnitude and direction of the block with respect to the reference frame. The spurious motion vectors, which do not belong to the moving object, are removed by connected component analysis before constructing MFH and MHI. To remove the spurious motion vectors, first a binary image of the frame is generated from the motion vector magnitude with a threshold of 0.5 to retain the half-pel motion values. Then a simple morphological *clean* operation is employed to remove isolated motion vectors (1's surrounded by 0's).

The MFH is constructed from non-zero P -frame motion vectors according to the following:

$$MFH_d(k, l) = \begin{cases} m_d^{kl}(\tau) & \text{if } E(m_d^{kl}(\tau)) < T_r \\ M(m_d^{kl}(\tau)) & \text{otherwise} \end{cases} \quad (1)$$

where $E(m_d^{kl}(\tau)) = \|m_d^{kl}(\tau) - \text{med}(m_d^{kl}(\tau) \dots m_d^{kl}(\tau - \alpha))\|^2$ and $M(m_d^{kl}(\tau)) = \text{med}(m_d^{kl}(\tau) \dots m_d^{kl}(\tau - \alpha))$.

Here med refers to median filter, $m_d^{kl}(\tau)$ can be horizontal (m_x) component or vertical (m_y) component of motion vector located at k th row and l th column in frame τ and α indicates the number of previous P -frames to be considered for median filtering. Typical range of α is 3–5 for various kinds of noise. Since the correlation of the frames decreases with the temporal distance between them, it is not advisable to increase the α value beyond 5. The function E checks the reliability of the current motion vector with respect to the past non-zero motion vectors at the same location against a predefined threshold T_r . The purpose of this

threshold T_r is to check the reliability of each newly arriving motion vector. Considering the human motion dynamics, the motion vectors of current P -frame cannot change much with respect to the neighboring P -frame motion vectors. At the same time the threshold should not be too tight since most of the recent motion vectors would then be ignored. In our system the threshold T_r is set at 4 for generating MFH. In other words, this threshold T_r makes sure that no reliable motion vector of MFH will be replaced by a recent noisy motion vector. Such spurious motion vectors are replaced by the reliable median value.

The MHI is constructed as given by Eq. (2),

$$MHI(k, l) = \begin{cases} \tau & \text{if } (|m_x^{kl}(\tau)| + |m_y^{kl}(\tau)|) \neq 0 \\ 0 & \text{otherwise} \end{cases} \quad (2)$$

Figs. 4 and 5 show the key frames of the bend-down and twist-left actions and the corresponding coarse MHI and MFH². The coarse MHI and MFH of other actions are shown in Fig. 6. The MHI is a function of the recency of the motion at every macroblock. The brightness of the macroblock is proportional to how recently the motion occurred. The MFH describes the spatial distribution of motion vectors over the video clip. In other words MFH quantifies the motion at spatial locations through horizontal and vertical components of the motion. The MHI, which has spatio-temporal information but no motion vector information, is complemented by the MFH. Thus MHI and MFH together capture the temporal and motion vector (m_x, m_y) information of the entire video sequence. The drawback of this representation is that, self-occlusion or overlapping of motion on the image plane may result in the loss of a part of

² Since the P -frames are predicted from the previous closest I or P -frame, the direction of the motion vectors appear opposite to the actual motion.

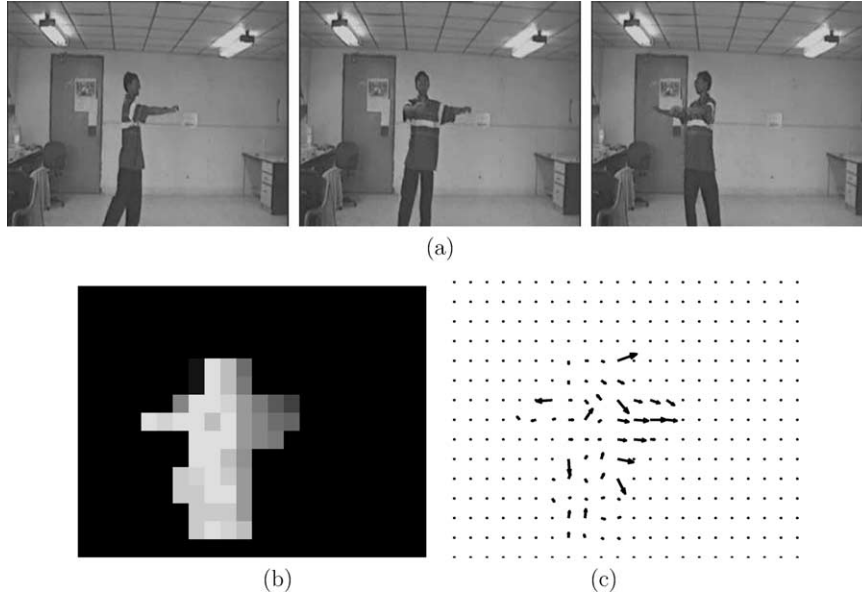


Fig. 5. (a) Key-frames of twist-left sequence and corresponding coarse (b) MHI (c) MFH.

the motion information. However, it might be representative enough for the considered human actions.

5. Feature extraction

Given the MHI and MFH of an action, it is essential to extract some useful features for classification. We have extracted features from MHI based on (i) Projection profiles and (ii) Centroid. The MFH based features are (i) Affine motion model; (ii) projected 1D feature and (iii) 2D polar feature [3].

5.1. MHI features

Projection profile based feature. Let N be the number of rows and M be the number of columns of MHI. Then the vertical profile is given by the vector P_v of size N and defined as $P_v[i] = \sum_{j=1}^M \text{MHI}[i, j]$. The horizontal profile is represented by the vector P_h of size M and define as $P_h[j] = \sum_{i=1}^N \text{MHI}[i, j]$. The features representing the distribution of projection profile with respect to the centroid are computed as

$$\mathbf{F}_{pp} = \begin{bmatrix} \frac{\sum_{i=1}^{h_c} P_h[i]}{\sum_{i=h_c+1}^M P_h[i]} & \frac{\sum_{i=1}^{v_c} P_v[i]}{\sum_{i=v_c+1}^N P_v[i]} \\ \frac{\sum_{i=h_c+1}^M P_h[i]}{\sum_{i=h_c+1}^M P_h[i]} & \frac{\sum_{i=v_c+1}^N P_v[i]}{\sum_{i=v_c+1}^N P_v[i]} \end{bmatrix} \quad (3)$$

where h_c and v_c are the horizontal and vertical centroids of MEI. The above feature (\mathbf{F}_{pp}) indicates the bias of the MHI along horizontal and vertical direction with respect to the centroid of MEI. This indirectly conveys the temporal information of motion along horizontal and vertical direction.

Centroid based feature. This feature is computed as the shift of centroids of MEI and MHI, which is given by the 2D vector

$$\mathbf{F}_c = [\text{MHI}_{xc} - \text{MEI}_{xc} \quad \text{MHI}_{yc} - \text{MEI}_{yc}] \quad (4)$$

The centroid of MHI differs from the centroid of MEI because it is computed using the gray-level time stamp values as weights in the summation. The above vector indicates the approximate direction of the movement of centroid for the corresponding action.

5.2. MFH features

Three types of features are extracted from MFH. Since it holds the entire history of spatial motion information, many useful features are extracted from MFH.

Affine feature. Though it is difficult to capture some complex motion, affine model gives a good approximation to the actual optical flow of the planar surface under orthographic projection [12]. An affine model requires six basic flow fields as shown in Fig. 7. The affine parameters are estimated by standard linear regression techniques. The regression is applied separately on each motion vector component since the x affine parameter depends only on horizontal component of motion vector and y parameter depends only on the vertical component of motion vector. Let $\mathbf{c} = [c_1 \dots c_6]$ be the 6D affine parameter vector. Then the linear least squares estimate of \mathbf{c} is given by:

$$\mathbf{c}^T = \left[\sum \mathbf{\Pi}(\mathbf{p})^T \mathbf{\Pi}(\mathbf{p}) \right]^{-1} \cdot \sum \mathbf{\Pi}^T(\mathbf{p}) \mathbf{v}(\mathbf{p}) \quad (5)$$

where

$$\mathbf{\Pi}(\mathbf{p}) = \begin{bmatrix} 1 & x & y & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & x & y \end{bmatrix}$$

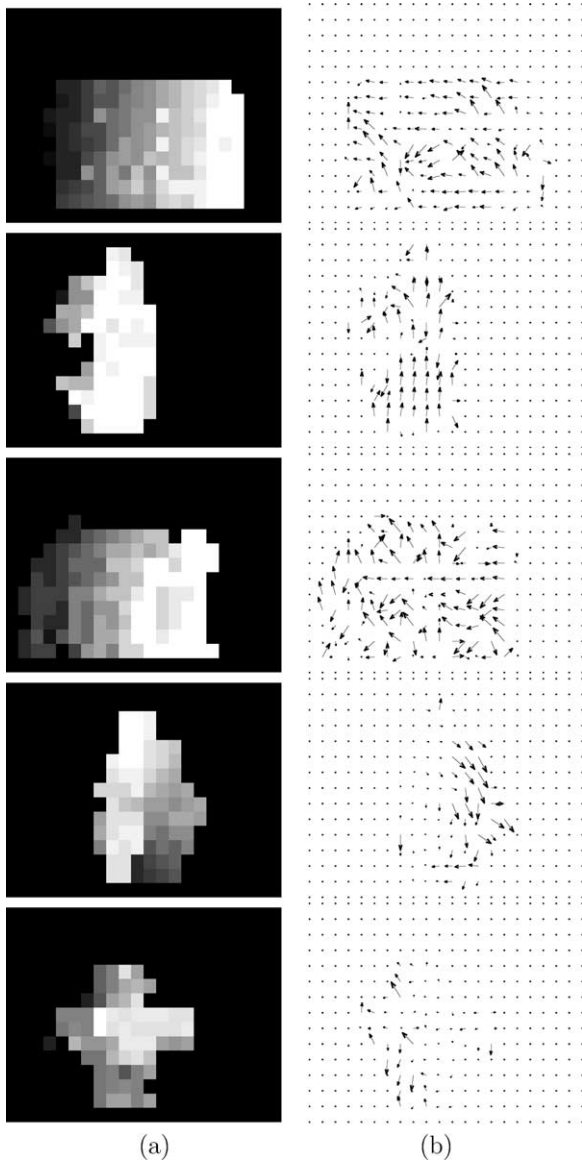


Fig. 6. (a) The coarse MHI and the corresponding (b) MFH of walk, jump, run, bend-up and twist-right action.

is the regressor and $\mathbf{p} = [x \ y]^T$ is the vector representing the position of pixel in the image plane and $\mathbf{v}(\mathbf{p})$ is the motion vector at location \mathbf{p} (here the spatial location of motion vectors are assigned to the center of the corresponding macroblock).

Projected 1D feature. Here horizontal and vertical components of the motion vectors are considered separately. The histogram values are quantized into five bins to cover the entire range in the following intervals: $[\text{Min}, -8], [-8, -3], [-3, 3], (3, 8], (8, \text{Max}]$. The bins are

chosen in such a way so as to capture the low, medium and higher speeds. The distance between the centers of low and medium speeds are set apart by 5 pels approximately. The motion vector magnitude exceeding 8 are considered as high speed.

2D polar feature. The angular direction and magnitude of motion vectors are considered together to quantize the polar plane into histogram bins. Each bin is defined by the angular range as well as the magnitude (radius) range. Here angular range is quantized into four intervals of length $\pi/2$ from $-\pi$ to $+\pi$. The magnitude range is quantized into the following intervals: $(0, 5], (5, 10], (10, \text{Max}]$. This leads to a feature vector of 12 dimension. Table 1 summarizes the features used in our experiment.

6. Classification results and discussion

The following seven actions were considered for recognition: walk, run, jump, bend up, bend down, twist right and twist left. For collecting the database, each subject was asked to perform each action many times in front of the fixed camera inside the laboratory. The actions were captured at the angle at which the camera could view the motion with minimal occlusion. The subjects are given freedom to perform the actions at their own pace at any distance in front of the camera.

We have used four types of classifiers for recognizing the action, namely Normalized KNN, Bayesian, Neural network: Multi-Layer feed forward Perceptron (MLP) and Support Vector Machines (SVM). As in the previous paper, seven actions (walk, run, jump, bend down, bend up, twist left and twist right) were considered for recognition. In our experimental setup, we trained the system with 10 instances of each action performed by four to five different subjects. For testing, we have used at least five instances per action with the subjects that are not used for training phase. The total number of samples used for training is 70 (10 samples/action) and 51 samples for testing.

6.1. K-nearest neighbors classifier

The KNN algorithm simply selects k -closest samples from the training data to the new instance and the class with the highest number of votes is assigned to the test instance. An advantage of this technique is due to its non-parametric nature, because we do not make any assumptions on the parametric form of the underlying distribution of classes. In higher dimensional spaces these distributions may be often erroneous. Even in situations where second order statistics

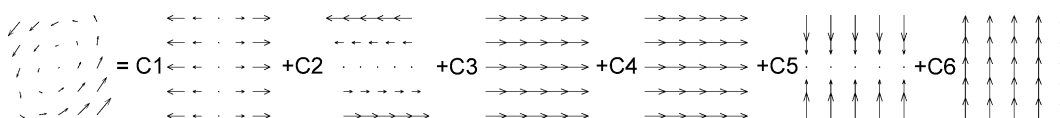


Fig. 7. Affine flow model expressed as a linear sum of basis flows.

Table 1
Features extracted from MHI and MFH

	Feature	Dimension
MHI based	Proj. profile	2
	Centroid	2
MFH based	Affine	6
	1D projected	10
	2D polar	12
	Total	32

Table 2
KNN classification result for $k = 3$

Input class	Result								
	Walk	Run	Jump	BD	BU	TWL	TWR	Error	
Walk	5	0	0	0	0	0	0	0	
Run	0	7	0	0	0	0	0	0	
Jump	0	0	7	0	0	0	0	0	
BD	0	0	0	11	0	0	0	0	
BU	0	0	0	0	8	1	0	1	
TWL	0	0	0	0	0	6	0	0	
TWR	0	0	0	0	0	0	6	0	
Error	0	0	0	0	0	1	0	1	

cannot be reliably computed due to limited training data, KNN performs very well, particularly in high dimensional feature spaces and on atypical samples. Table 2 shows the classification results of KNN classifier with all aforementioned features.

6.2. Bayes classifier

The second classifier used is Bayes—a parametric classifier that assumes normal distribution for class (ω) conditional probability of feature vector \mathbf{x} , $P(\mathbf{x}|\omega_i)$. Though Bayes classifier is optimal, the performance degrades if the models used are erroneous. Since erroneous models degrade classification performance, starting with the first feature, we have added subsequent features only if they improved the classification performance. Table 3 shows the performance

Table 3
Bayes classification result

Input class	Result								
	Walk	Run	Jump	BD	BU	TWL	TWR	Error	
Walk	3	2	0	0	0	0	0	2	
Run	0	7	0	0	0	0	0	0	
Jump	0	0	7	0	0	0	0	0	
BD	0	0	0	11	0	0	0	0	
BU	0	0	0	0	9	0	0	0	
TWL	0	0	0	0	0	6	0	0	
TWR	0	0	0	0	0	1	5	1	
Error	0	2	0	0	0	1	0	3	

Table 4
Neural net classification result

Input class	Result								
	Walk	Run	Jump	BD	BU	TWL	TWR	Error	
Walk	4	1	0	0	0	0	0	1	
Run	0	7	0	0	0	0	0	0	
Jump	0	0	7	0	0	0	0	0	
BD	0	0	0	11	0	0	0	0	
BU	0	0	0	0	9	0	0	0	
TWL	0	0	0	0	0	6	0	0	
TWR	0	0	0	0	0	0	6	0	
Error	0	1	0	0	0	0	0	1	

of Bayes classifier with only four selected features out of total 32 features. The selected feature numbers are 1, 3, 6 and 9, i.e. three from the affine feature and one from MHI centroid-based feature. With all the features Bayes classifier gives a performance of 92.1% (see Table 8).

6.3. Neural network classifier

MLP is a supervised neural network. It can have multiple inputs and outputs and multiple hidden layers with arbitrary number of neurons (nodes). In our network, the commonly used sigmoid function is used as the activation function for nodes in the hidden layer. The MLP utilizes the back-propagation (BP) algorithm for determining suitable weights and biases of the network using supervised training [14]. Table 4 shows the classification results obtained with an MPL trained with two hidden layers with 15 neurons in each layer using all the features.

6.4. SVM classifier

SVM [34] are powerful tools for data classification. SVM is based on the idea of hyperplane classifier, that achieves classification by a separating surface (linear or nonlinear) in the input space of the data set. SVMs are modeled as optimization problems with quadratic objective functions and linear constraints. Tables 5 and 6 show the classification

Table 5
Linear SVM classifier result

Input class	Result								
	Walk	Run	Jump	BD	BU	TWL	TWR	Error	
Walk	5	0	0	0	0	0	0	0	
Run	0	6	1	0	0	0	0	1	
Jump	0	0	7	0	0	0	0	0	
BD	0	0	0	11	0	0	0	0	
BU	0	0	0	0	8	1	0	1	
TWL	0	0	0	0	0	6	0	0	
TWR	0	0	0	0	0	0	6	0	
Error	0	0	1	0	0	1	0	2	

Table 6
Classification result using non-linear SVM classifier with a radial based kernel

Input class	Result							
	Walk	Run	Jump	BD	BU	TWL	TWR	Error
Walk	5	0	0	0	0	0	0	0
Run	0	7	0	0	0	0	0	0
Jump	0	0	7	0	0	0	0	0
BD	0	0	0	11	0	0	0	0
BU	0	0	0	0	8	1	0	1
TWL	0	0	0	0	0	6	0	0
TWR	0	0	0	0	0	0	6	0
Error	0	0	0	0	0	1	0	1

Table 7
Comparison of various classifiers

Classifier	No. of features used	Classification accuracy (%)
KNN ($k = 3$)	32	98.0
Neural Net	32	98.0
SVM (RBF-kernel)	32	98.0
Bayes	4	94.1

results of SVM classifier with linear kernel and radial-based kernel.

Comparing the results of the classifiers, the results obtained by KNN, Neural Net and SVM (with RBF-kernel) show excellent performance. Bayes classifier recognizes most of the actions, but is relatively less successful in discriminating between ‘walk’ and ‘run’ actions. This could be due to the parametrization of the underlying feature distribution. Moreover the Bayes result is obtained only with the selected four features, whereas the other classifiers use all features. Table 7 summarizes the recognition results for various classifiers.

Consistency in results obtained using various classifiers, proves the credibility of features used. The uniform results

obtained using SVM, KNN and Neural Nets points to the fact that the system has very few outliers. It must be noted that the performance of the Bayesian classifier is only marginally lesser, in spite of drastically reducing the number of features to only four, compared to 32 used in the other cases. The deterioration in the performance of the Bayesian classifier on using all features may be attributed to ‘curse of dimensionality’.

7. Performance analysis of features

In this section, we present the performance of each feature set for various classifiers. Fig. 8 shows the recognition performance of each feature with test and training samples using the nearest neighbor criterion. The individual performance of the first 10 or 11 features is good on both, the test as well as the training samples. Other features perform slightly better with test samples compared to training samples. Here, the considered test subjects are different from the training ones and the subjects were given freedom to perform the action at their own pace at any location in front of the camera. So the features show invariance to translation, scale and speed of action.

Table 8 shows the performance of each feature set for various classifiers on test samples. For each of the classifiers, different feature sets contribute in different proportions. The 1D projected (11–20) and 2D polar (21–32) features show good performance (around 85–90%) consistently. Affine features (1–6) show better performance than the above two except for KNN classifier with the feature dimension being only six. Other features (projection profile and centroid) improve the overall performance considerably in spite of their low dimension. Though the MHI features: projection profile and centroid individually perform inferior compared to other features, they jointly perform well. The four MHI

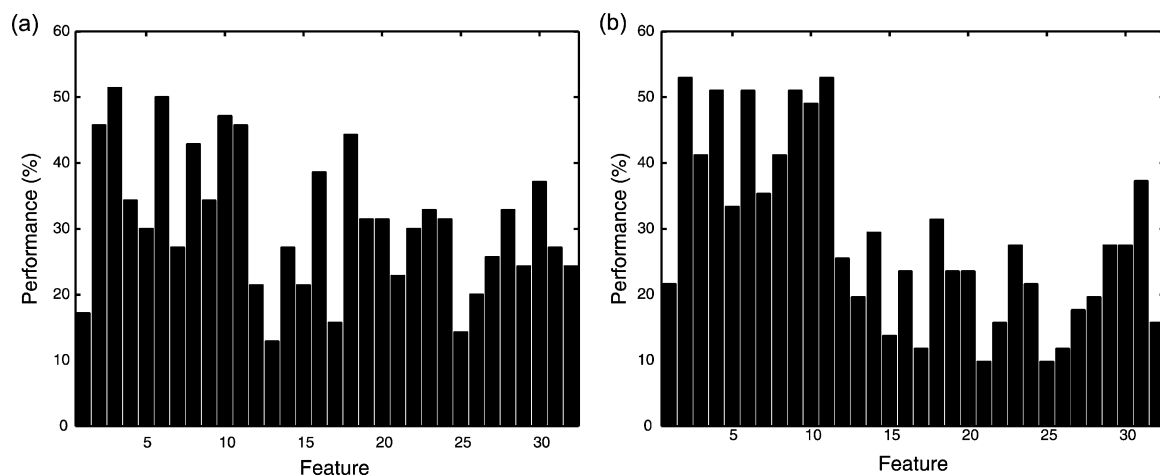


Fig. 8. Feature ranking on individual performance (a) with test samples (b) with training samples.

Table 8
Performance of feature sets (%)

Input features		Classifier results (%)			
Name	Numbers	KNN ($k = 3$)	Bayes	NNet	SVM (RBF)
Affine (6)	1–6	78.4	92.2	94.1	86.3
Proj. profile (2)	7–8	64.7	58.8	70.6	11.8
Centroid (2)	9–10	68.6	68.6	72.6	66.7
1D proj. (10)	11–20	92.2	84.3	86.3	90.2
2D polar (12)	21–32	86.3	86.3	90.2	86.3
Overall (32)	1–32	98.0	92.1	98.0	98.0

features together give recognition accuracy of 82.35% with nearest-neighbor classifier. Hence, MHI features are indeed powerful in recognizing human actions.

The recognition results for the seven Hu moments extracted from the coarse MHI were also tried with all the classifiers. The Hu moment features do not perform well with these four classifiers. With Hu moments the KNN ($k = 3$), Bayes, Neural Net and SVM (RBF-kernel) classifiers give the recognition accuracy of 25.42%, 37.25%, 47.06% and 27.45%, respectively. The reason could be that Hu moments may require the MHIs at pixel resolution to capture the characteristics of the action.

8. Conclusion

In this paper, we have proposed a method for constructing coarse MHI and MFH from compressed MPEG video with minimal decoding. Various useful features are extracted from the above mentioned motion representations for human action recognition. We have shown the recognition results for four classification paradigms. The performance of these features is analyzed and compared. Though the test instances are from entirely different subjects other than those used for training the classifiers, the results show excellent recognition accuracy. The KNN, Neural network (MLP) and SVM (RBF-kernel) classifiers give the best classification accuracy of 98% and 1D projected and 2D polar features show consistent performance with all the classifiers. Since the data is handled at macroblock level, the computational cost is extremely less compared to the pixel domain processing.

Acknowledgements

The authors wish to express grateful thanks to the referees for their useful comments and suggestions to improve the presentation of this paper.

References

- [1] J.K. Aggarwal, Q. Cai, Human motion analysis: a review, *Computer Vision and Image Understanding* 73 (3) (1999) 428–440.
- [2] D. Ayers, M. Shah, Monitoring human behavior from video taken in an office environment, *Image and Vision Computing* 19 (12-1) (2001) 833–846.
- [3] R. Venkatesh Babu, B. Anantharaman, K.R. Ramakrishnan, S.H. Srinivasan, Compressed domain action classification using HMM, *Pattern Recognition Letters* 23 (10) (2002) 1203–1213.
- [4] R. Venkatesh Babu, K.R. Ramakrishnan, Compressed Domain Human Motion Recognition Using Motion History Information, *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, April, vol. 3, 2003, pp. 41–44.
- [5] A.F. Bobick, J.W. Davis, The recognition of human movement using temporal templates, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 23 (3) (2001) 257–267.
- [6] A.F. Bobick, A.D. Wilson, A state-based approach to the representation and recognition of gesture, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 19 (12) (1997) 1325–1337.
- [7] K. Boehm, W. Broll, M. Sokolewicz, Dynamic gesture recognition using neural networks; a fundament for advanced interaction construction, *Proceedings of the SPIE—The International Society for Optical Engineering* 2177 (1994) 336–346.
- [8] C. Bregler, Learning and recognizing human dynamics in video sequences, *Proceedings of the IEEE CVPR* (1997) 568–574.
- [9] T.J. Darrell, A.P. Pentland, Space–time gestures, *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition* (1993) 335–340.
- [10] J. Davis, Hierarchical motion history images for recognizing human motion, *IEEE Workshop on Detection and Recognition of Events in Video* (2001) 39–46.
- [11] J. Davis, A.F. Bobick, The representation and recognition of human movements using temporal templates, *Proceedings of the IEEE CVPR* (1997) 928–934.
- [12] D.J. Fleet, M.J. Black, Y. Yacoob, A. Jepson, Design and use of linear models for image motion analysis, *International Journal of Computer Vision* 36 (3) (2000) 171–193.
- [13] D.M. Gavrilu, The visual analysis of human movement: a survey, *Computer Vision and Image Understanding* 73 (1) (1999) 82–98.
- [14] S. Haykin, *Neural Networks: A Comprehensive Foundation*, Prentice Hall, 1998.
- [15] B.K.P. Horn, B.G. Schunck, Determining optical flow, *Artificial Intelligence* 17 (1981) 185–203.
- [16] M.K. Hu, Visual pattern recognition by moment invariants, *IRE Transactions on Information Theory* 8 (2) (1962) 179–187.
- [17] Y.A. Ivanov, A.F. Bobick, Recognition of visual activities and interactions by stochastic parsing, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 22 (8) (2000) 852–872.
- [18] T. Kohonen, The self-organizing map, *Proceedings of the IEEE* 78 (9) (1990) 1464–1480.
- [19] M. Su, H. Huang, C. Lin, C. Huang, C. Lin, Application of neural networks in spatio temporal hand gesture recognition, *Proceedings of the IEEE World Congress on Computational Intelligence*, 1998.
- [20] A. Madabhushi, J. K. Aggarwal, A Bayesian Approach to Human Activity Recognition, *Second IEEE Workshop on Visual Surveillance*, 25–30, 1999.
- [21] J.L. Mitchell, W.B. Pennebaker, C.E. Fogg, D.J. LeGall, *MPEG Video Compression Standard*, International Thomson Publishing, 1996.
- [22] C.W. Ng, S. Ranganath, Real-time gesture recognition system and application, *Image and Vision Computing* 20 (13–14) (2002) 993–1007.
- [23] N.M. Oliver, B. Rosario, A. Pentland, A bayesian computer vision system for modeling human interactions, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 22 (8) (2000) 831–843.

- [24] R. Polana, R. Nelson, Low level recognition of human motion, Workshop on Non-Rigid Motion (1994) 77–82.
- [25] A. Psarrou, S. Gong, M. Walter, Recognition of human gestures and behaviour based on motion trajectories, *Image and Vision Computing* 20 (5–6) (2002) 349–358.
- [26] J. Rittscher, A. Blake, S.J. Roberts, Towards the automatic analysis of complex human body motions, *Image and Vision Computing* 20 (12) (2002) 905–916.
- [27] R. Rosales, Recognition of human action based on moment based features, Technical Report BU 98-020, Boston University, Computer Science, 1998.
- [28] M. Shah, R. Jain, *Motion Based Recognition*, Kluwer, 1997.
- [29] Standard MPEG1: ISO/IEC 11172, Coding of moving pictures and associated audio for digital storage media at upto about 1.5 Mbit/s, 1996.
- [30] T. Starner, A. Pentland, Real-time american sign language recognition from video using hidden Markov models, MIT Media Lab, TR-375, 1995.
- [31] T. Starner, A.P. Pentland, Real-time american sign language recognition using desk and wearable computer based video, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 20 (12) (1998) 1371–1375.
- [32] X. Sun, B.S. Manjunath, Panoramic capturing and recognition of human activity, *IEEE International Conference on Image Processing* 2 (2002) 813–816.
- [33] M. Umeda, Recognition of multi-font printed chinese character, *Proceedings of sixth Computer Vision and Pattern Recognition* (1982) 793–796.
- [34] V.N. Vapnik, *The Nature of Statistical Learning Theory*, Springer, New York, 1995.
- [35] J. Yamato, J. Ohya, K. Ishii, Recognizing human action in time-sequential images using hidden Markov model, *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (1992) 379–385.