

Matching Mixtures of Trajectories for Human Action Recognition

Michalis Vrigkas¹, Vasileios Karavasilis¹, Christophoros Nikou¹,
and Ioannis A. Kakadiaris²

¹*Department of Computer Science, University of Ioannina, Ioannina, Greece*

²*Computational Biomedicine Lab, Department of Computer Science, University of
Houston, Houston, Texas, USA*

{mvrigkas, vkaravas, cnikou}@cs.uoi.gr, ioannisk@uh.edu

Abstract

A learning-based framework for action representation and recognition relying on the description of an action by time series of optical flow motion features is presented. In the learning step, the motion curves representing each action are clustered using Gaussian mixture modeling (GMM). In the recognition step, the optical flow curves of a probe sequence are also clustered using a GMM, then each probe sequence is projected onto the training space and the probe curves are matched to the learned curves using a non-metric similarity function based on the longest common subsequence, which is robust to noise and provides an intuitive notion of similarity between trajectories. Also, canonical time warping is utilized to find an alignment between the mean trajectories. Finally, the probe sequence is categorized to the learned action with the maximum similarity using a nearest neighbor classification scheme. We also present a variant of the method where the lengths of the time series are reduced by dimensionality reduction in both training and test phases, in order to smooth out the outliers, which are common in these type of sequences. Experimental results on Weizmann, KTH, UCF Sports and UCF YouTube action databases demonstrate the effectiveness of the proposed method.

Keywords:

Human action recognition; Optical flow; Motion curves; Gaussian mixture modeling (GMM); Clustering; Dimensionality reduction; Longest common subsequence.

1. Introduction

Action recognition is a preponderant and difficult task in computer vision. Many applications, including video surveillance systems, human-computer interaction, and robotics for human behavior characterization, require a multiple activity recognition system.

Our goal is to examine human activities from video sequences. However, training an action recognition system with only the knowledge of the motion of the current subject it is on its own a challenging task. The main problem is how we can ensure the continuity of the curves along time as an action occurs uniformly or non-uniformly within a video sequence. Unlike other approaches [1, 2], which use snippets of motion trajectories, our approach uses the full length of motion curves by tracking the optical flow features. Another question concerns the optimal model that one should adopt for recognizing human actions with high accuracy. This is accomplished by a statistical measure based on the data likelihood. The different lengths of the video sequences and therefore the respective lengths of the motion curves is another problem that is addressed. The large variance between benchmark datasets shows how the algorithm may be generalized. All these problems are discussed here and proper solutions are proposed. To this end, we have conducted more experiments on several datasets [3, 4, 5, 6] that would help us to understand how human activity recognition works.

In this paper, we address the problem of human action recognition by representing an action with a set of clustered motion trajectories. Motion curves are generated by optical flow features which are then clustered using a different Gaussian mixture [7] for each distinct action. The optical flow curves of a probe sequence are also clustered using a Gaussian mixture model (GMM) and they are matched to the learned curves using a similarity function [8] relying on the longest common subsequence (LCSS) between trajectories and the canonical time warping (CTW) [9]. Linear [7] and non linear [10] dimensionality reduction methods may also be employed in order to remove outliers from the motion curves and reduce their lengths. The motion curve of a new probe video is projected onto the space of the training sequences, and then the action label of the closest projection is selected according to the learned feature vectors as the identity of the probe sequence. The LCSS is robust to noise and provides an intuitive notion of similarity between trajectories. Since different actors perform the same action in different manners and at different speeds, an advantage of the LCSS similarity is that

it can handle with motion trajectories of varied lengths. On the other hand, CTW, which is based on the dynamic time warping [11], allows the spatio-temporal alignment between two human motion sequences. A preliminary version of this work was presented in [12]. One of the main contributions of this paper is that the training sequences do not need to have the same length. When a new probe sequence comes, it is matched against all the training sequences using the LCSS similarity measure. This measure provides a similarity between motion trajectories without enforcing one-to-one matching. An optimal matching is performed using dynamic programming, which detects similar pairs of curve segments [8].

Tracking of optical flow features along time allows us to collect time series that preserve their continuity along time. It is true that correspondence is missing. However, this is the main assumption in many works [13, 14, 15]. If data association were used the resulting feature trajectories would have short duration and would be incomplete, as the features disappear and reappear due to occlusion, illumination, viewpoint changes and noise. In that case, a combination of sparse approach of clustering trajectories with variant lengths and tracking approaches should be used [16, 17]. This is not the central idea in this paper, as the nature of the feature trajectories drastically changes.

There are three main contributions. First, human motion is represented by a small set of trajectories which are the mean curves of the mixture components along with their covariance matrices. The complexity of the model is considered low, as it is determined by the Bayesian Information Criterion (BIC), but any other model selection technique may be applied. Second, the computational cost is lower since the use of dimensionality reduction allows the algorithm to cope with trajectories of smaller lengths. Finally, the use of the longest common subsequence index allows input curves of different lengths to be compared reliably.

In the rest of the paper, the related work is presented in Section 2, while the extraction of motion trajectories, the clustering and the curve matching are presented in Section 3. In Section 4, we report results on the Weizmann [3], the KTH [4], the UCF Sports [5] and the UCF YouTube [6] action classification datasets. Finally, conclusions are drawn in Section 5.

2. Related Work

The problem of categorizing a human action remains a challenging task that has attracted much research effort in the recent years. The surveys

in [18] and [19] provide a good overview of the numerous papers on action/activity recognition and analyze the semantics of human activity categorization. Several feature extraction methods for describing and recognizing human actions have been proposed [13, 4, 20, 21, 14]. A major family of methods relies on optical flow which has proven to be an important cue. Efros *et al.* [13] recognize human actions from low-resolution sports video sequences using the nearest neighbor classifier, where humans are represented by windows of height of 30 pixels. The approach of Fathi and Mori [14] is based on mid-level motion features, which are also constructed directly from optical flow features. Moreover, Wang and Mori [15] employed motion features as inputs to hidden conditional random fields and support vector machine (SVM) classifiers. Real time classification and prediction of future actions is proposed by Morris and Trivedi [22], where an activity vocabulary is learnt through a three step procedure. Other optical flow-based methods which gained popularity are presented in [23, 24, 25]. The main disadvantage of using a global representation such as optical flow, is the sensitivity to noise and partial occlusions.

The classification of a video sequence using local features in a spatio-temporal environment has also been given much focus. Schuldt *et al.* [4] represent local events in a video using space-time features, while an SVM classifier is used to recognize an action. Gorelick *et al.* [26] consider actions as 3D space time silhouettes of moving humans. They take advantage of the Poisson equation solution to efficiently describe an action by utilizing spectral clustering between sequences of features and applying nearest neighbor classification to characterize an action. Niebles *et al.* [21] address the problem of action recognition by creating a codebook of space-time interest points. A hierarchical approach was followed by Jhuang *et al.* [20], where an input video is analyzed into several feature descriptors depending on their complexity. The final classification is performed by a multi-class SVM classifier. Dollár *et al.* [27] proposed spatio-temporal features based on cuboid descriptors. An action descriptor of histograms of interest points, relying on [4] was presented in [28]. Random forests for action representation have also been attracting widespread interest for action recognition [29, 30]. Furthermore, the key issue of how many frames are required to recognize an action is addressed by Schindler and Van Gool [31]. However, mid-level feature approaches depend on the number of interest points detected.

The problem of identifying multiple persons simultaneously and perform action recognition is presented in [32]. The authors consider that a person

has first been localized by performing background subtraction techniques. Based on the Histograms of Oriented Gaussians [33] they detect a human, whereas classification of actions are made by training a SVM classifier. Action recognition using depth cameras are introduced in [34] and a new feature called “local occupancy pattern” is also proposed. A novel multi-view activity recognition method is presented in [35]. Descriptors from different views are connected together forming a new augmented feature that contains the transition between the different views. A new type of feature called the “Hankelet” is presented in [36]. This type of feature, which is formed by short tracklets, along with a BoW approach is able to recognize actions under different viewpoints, without requiring any camera calibration. Zhou and Wang [37] have also proposed a new representation of local spatio-temporal cuboids for action recognition. Low level features are encoded and classified via a kernelized SVM classifier, whereas a classification score denotes the confidence that a cuboid belongs to an atomic action. The new feature act as complementary material to the low-level feature. The work of Sanchez-Riera *et al.* [38] recognize human actions using stereo cameras. Based on the technique of bag-of-words, each action is presented by a histogram of visual words, whereas their approach is robust to background clutter.

Earlier approaches are based on describing actions by using dense trajectories. The work of Wang *et al.* [39] is focused on tracking dense sample point from video sequences using optical flow. Le *et al.* [40] discover the action label in an unsupervised manner by learning features directly from video data. A high-level representation of video sequences, called Action Bank, is presented by Sadanand and Corso [41]. Each video is represented as a set of action descriptors which are put in correspondence. The final classification is performed by a SVM classifier. Yan and Luo [28] have also proposed a new action descriptor based on spatial temporal interest points (STIP) [42]. In order to avoid overfitting they have also proposed a novel classification technique by combining the Adaboost and sparse representation algorithms. In [43], a visual feature using Gaussian mixture models efficiently represents the spatio-temporal context distributions between the interest point at several space and time scales. An action is represented by a set of features extracted by the interest points over the video sequence. Finally, a vocabulary based approach has been proposed by Kovashka and Grauman [44]. The main idea is to find the neighboring features around the detected interest points quantize them and form a vocabulary. Raptis *et al.* [2] proposed a mid-level approach extracting that spatio-temporal features

construct clusters of trajectories, which can be considered as candidates of an action, and a graphical model is utilized to control these clusters.

Human action recognition using temporal templates has also been proposed by Bobick and Davis [45]. An action was represented by a motion template composed of a binary motion energy image (MEI) and a motion history image (MHI). Recognition was accomplished by matching pairs of MEI and MHI. A variation of the MEI idea was proposed by Ahmad and Lee [46], where the silhouette energy image (SEI) was proposed. The authors have also introduced several variability models to describe an action, and action classification was carried out using a variety of classifiers. Moreover, the proposed model is sensitive to illumination and background changes.

A major current focus in action recognition from still images or videos has been made in the context of scene appearance [47, 48, 49]. More specifically, Thureau and Hlavac [47] represented actions by histograms of pose primitives and n-gram expressions are used for action classification. Also, Yang *et al.* [48] combined actions and human poses together, treating poses as latent variables, to deduce the action label of a still image in order to recognize an action, while Maji *et al.* [49] introduced a representation of human poses called the poselet activation vector, which is defined by the 3D orientation of the head and torso and provides a robust representation of human pose and appearance. Moreover, action categorization based on modeling the motion of parts of the human body was presented by Tran *et al.* [50], where sparse representation was used to model and recognize complex actions. In the sense of template matching techniques Rodriguez *et al.* [5] introduced the Maximum Average Correlation Height (MACH) filter which is a method for capturing intra-class variability by synthesizing a single action MACH filter for a given action class. However, these approaches are limited by the fact that a human action is a continuous act in time and space and therefore, the estimation of still human poses may lead to incorrectly inferences.

Social interactions are an important part of humans daily life. A fundamental component of human behavior is the ability to interact with other people via their actions. Fathi *et al.* [51] models social interactions by estimating the location and the orientation of the faces of the persons taking part in a social event, computing thus a line of sight for each face. This information is used to infer the location where an individual person attend. The type of interaction is recognized by assigning social roles in each person. The authors are able to recognize three types of social interactions: dialogue, discussion and monologue. Human behavior on sport datasets are

introduced by Lan *et al.* [52]. The idea of social roles in conjunction with low-level actions and high-level events model the behavior of humans in a scene. Burgos-Artizzu *et al.* [53] discussed the social behavior of mice. Each video sequence is segmented into periods of activities by constructing a temporal context that combines spatio-temporal features. Kong *et al.* [54] proposed a new high-level descriptor called “interactive phrases” in order to recognize human interactions. This descriptor is a binary motion relationship descriptor for recognizing complex human interactions. Interactive phrases are treated as latent variables, while the recognition is performed in a SVM framework.

3. Action Representation and Recognition

Our goal is to analyze and interpret different classes of actions to build a model for human activity categorization. Given a collection of figure-centric sequences, we represent motion templates using optical flow [55] at each frame. Assuming that a bounding box can be automatically obtained from the image data, we define a rectangle region of interest (ROI) around the human. A brief overview of our approach is depicted in Figure 1. In the training mode, we assume that the video sequences contain only one actor performing only one action per frame. However, in the recognition mode, we allow more than one action per video frame. The optical flow vectors as well as the motion descriptors [13] for each sequence are computed. These motion descriptors are collected together to construct motion curves, which are clustered using a mixture model to describe a unique action. Then, the motion curves are clustered and each action is modeled by a set of clustered motion curves. Action recognition is performed by matching the clusters of motion curves of the probe sequence and the clustered curves in each training sequence.

3.1. Motion Representation

The proposed approach employs optical flow features [55]. These motion descriptors are commonly used in many recognition problems and they are shown to be quite reliable despite the existence of noisy features. Within a figure-centric scene, any human motion may be decomposed to the motion of different body parts (e.g., head and limbs). We can easily localize the motion by computing the optical flow vectors for the regions around the human torso.

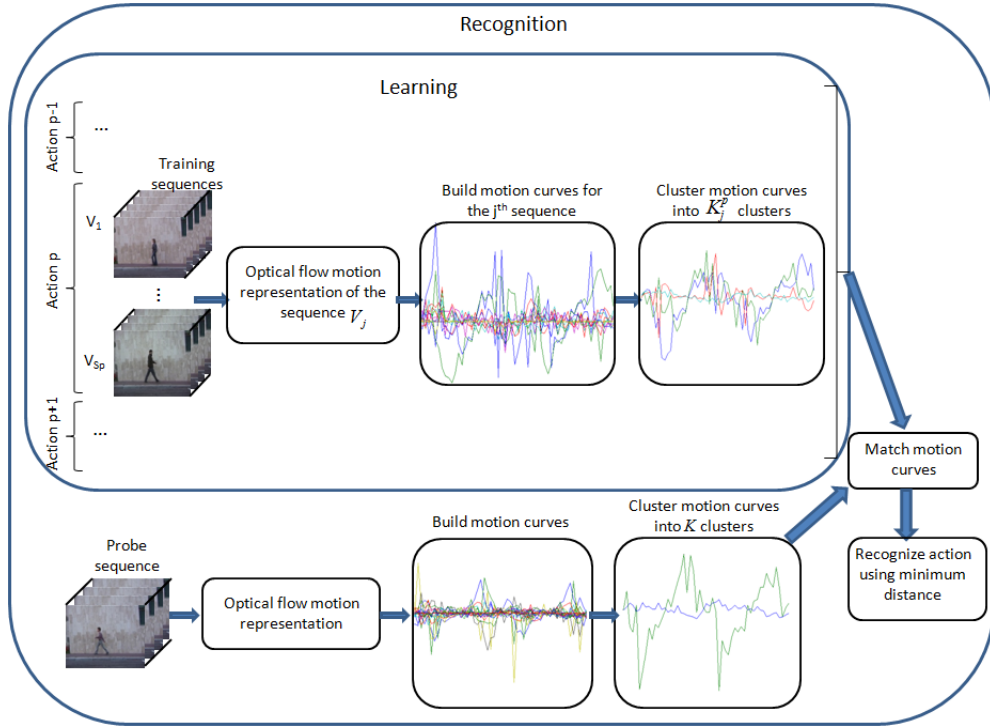


Figure 1: Overview of our approach.

Following the work of Efron *et al.* [13], we compute the motion descriptor for the ROI as a four-dimensional vector $\mathbf{F}_i = (F_{x_i}^+, F_{x_i}^-, F_{y_i}^+, F_{y_i}^-) \in \mathbb{R}^4$, where $i = 1, \dots, N$, with N being the number of pixels in the ROI. Also, the matrix \mathbf{F} refers to the blurred, motion compensated optical flow. We compute the optical flow \mathbf{F} , which has two components, the horizontal \mathbf{F}_x , and the vertical \mathbf{F}_y , at each pixel. It is worth noting that the horizontal and vertical components of the optical flow \mathbf{F}_x and \mathbf{F}_y are half-wave rectified into four non-negative channels $F_x^+, F_x^-, F_y^+, F_y^-$, so that $\mathbf{F}_x = F_x^+ - F_x^-$ and $\mathbf{F}_y = F_y^+ - F_y^-$. In the general case, optical flow is suffering from noisy measurements and analyzing data under these circumstances will lead to unstable results. To handle any motion artifacts due to camera movements, each half-wave motion compensated flow is blurred with a Gaussian kernel. In this way, the substantive motion information is preserved, while minor variations are discarded. Thus, any incorrectly computed flows are removed. Since all curves are considered normally distributed there is an intrinsic smoothing

of the optical flow curves. Moreover, at a preprocessing step, we discard flows whose amplitude is over 20% of the standard deviation of the mean amplitude of all curves for each video.

3.2. Extraction of Motion Curves

A human action is represented by a set of primitive motion curves which are constructed directly from the optical flow motion descriptors. The main idea is to extract the salient features, which describe a relative motion from each frame and associate them with the corresponding feature in the next frame.

Consider T to be the number of image frames and $C = \{c_i(t)\}, t \in [0, T]$, is a set of motion curves for the set of pixels $i = 1, \dots, N$ of the ROI. Each motion curve is described as a set of points corresponding to the optical flow vector extracted in the ROI. Specifically, we describe the motion at each pixel by the optical flow vector $\mathbf{F}_i = (F_{x_i}^+, F_{x_i}^-, F_{y_i}^+, F_{y_i}^-)$. A set of motion curves for a specific action is depicted in Figure 1. Given the set of motion descriptors for all frames, we construct the motion curves by following their optical flow components in consecutive frames. If there is no pixel displacement we consider a zero optical flow vector displacement for this pixel.

The set of motion curves describes completely the motion in the ROI. Once the motion curves are created, pixels and therefore curves that belong to the background are eliminated. We assume that the motion are normally distributed, thus, we keep flows whose values are inside 6 standard deviations of the amplitude distributions. In order to establish a correspondence between the motion curves and the actual motion, we perform clustering of the motion curves using a Gaussian mixture model. We estimate the characteristic motion which is represented by the mean trajectory of each cluster.

3.3. Motion Curves Clustering

A motion curve is considered to be a 2D time signal:

$$c_{ji}(t) = (F_{x_{ji}}(t), F_{y_{ji}}(t)), t \in [0, T], \quad (1)$$

where the index $i = 1, \dots, N$ represents the i^{th} pixel, for the j^{th} video sequence in the training set. To efficiently learn human action categories, each action is represented by a GMM by clustering the motion curves in every sequence of the training set. The p^{th} action ($p = 1, \dots, A$), in the j^{th} video sequence ($j = 1, \dots, S_p$), is modeled by a set of K_j^p mean curves learned by a

GMM. The likelihood of the i^{th} curve $c_{ji}^p(t)$ of the p^{th} action in the j^{th} video is given by:

$$p(c_{ji}^p; \pi_j^p, \mu_j^p, \Sigma_j^p) = \sum_{k=1}^{K_j^p} \pi_{jk}^p \mathcal{N}(c_{ji}^p(t); \mu_{jk}^p, \Sigma_{jk}^p), t \in [0, T], \quad (2)$$

where $\pi_j^p = \{\pi_{jk}^p\}_{k=1}^{K_j^p}$ are the mixing coefficients, $\mu_j^p = \{\mu_{jk}^p\}_{k=1}^{K_j^p}$ is the set of the mean curves and $\Sigma_j^p = \{\Sigma_{jk}^p\}_{k=1}^{K_j^p}$ is the set of covariance matrices. The covariance matrix in equation (2) is a diagonal $\Sigma_{jk}^p = \text{diag}(\sigma_{jk,1}^{2p}, \dots, \sigma_{jk,T}^{2p})$. Therefore, the log-likelihood of the p^{th} action in the j^{th} video can be written as:

$$L(c_j^p) = \prod_{i=1}^{N_j^p} \ln \sum_{k=1}^{K_j^p} \pi_{jk}^p \mathcal{N}(c_{ji}^p(t); \mu_{jk}^p, \Sigma_{jk}^p), t \in [0, T], \quad (3)$$

where N_j^p is the number of motion curves in the training set describing the p^{th} action in the j^{th} video.

The GMM is trained using the Expectation-Maximization (EM) algorithm [7], which provides a solution to the problem of estimating the model's parameters. The initialization of the EM algorithm is performed by the K-means algorithm. We have examined several configurations for the initialization of K-means and we decided to employ K-means with 50 different random initializations which were consistent and had no significant impact on the final classification. However, the number of mixture components should be determined. To select the number of the Gaussians K_j^p , for the j^{th} training video sequence, representing the p^{th} action, the Bayesian Information criterion (BIC) [7] is used:

$$BIC(c_j^p) = L(c_j^p(t)) - \frac{1}{2}MN_j^p, t \in [0, T], \quad (4)$$

where M is the number of parameters of the GMM to be inferred. Thus, when EM converges the cluster labels of the motion curves are obtained. This is schematically depicted in Figure 1, where a set of motion trajectories, representing a certain action (e.g., p), in a video sequence (e.g., labeled by j) is clustered by a GMM into $K_j^p = 2$ curves for action representation. Note that, a given action is generally represented by a varying number of mean trajectories as the BIC criterion may result in a different number of components in different sequences.

Apart from the BIC criterion, there are other techniques for determining the appropriateness of a model such as the Akaike Information Criterion (AIC) [7].

$$AIC(c_j^p) = L(c_j^p(t)) - M, t \in [0, T], \quad (5)$$

where M is the number of parameters of the GMM to be inferred. BIC is independent of the prior, it can measure the efficiency of the parameterized model in terms of predicting the data and it penalizes the complexity of the model, where complexity refers to the number of parameters in the model. It is also approximately equal to the minimum description length criterion [7] but with negative sign, it can be used to choose the number of clusters according to the intrinsic complexity present in a particular dataset and it is closely related to other penalized likelihood criteria such as the AIC. BIC tends to select highly parsimonious models, while AIC tends to include more parameters [56, 57]. Complexity measures such as BIC and AIC have the virtue of being easy to evaluate, but can also give misleading results.

3.4. Matching of Motion Curves

Once a new probe video is presented, where we must recognize the action depicted, the optical flow is computed, motion trajectories are created and clustered, and they are compared with the learned mean trajectories of the training set. Recall that human actions are not uniform sequences in time, since different individuals perform the same action in different manner and at different speeds. This means that motion curves have varied lengths. An optimal matching may be performed using dynamic programming which detects similar pairs of curve segments. The longest common subsequence (LCSS) [8] is robust to noise and provides a similarity measure between motion trajectories since not all points need to be matched.

Let $\mu(t)$, $t \in [0, T]$ and $\nu(\tau)$, $\tau \in [0, T']$ be two curves of different lengths. Then, we define the affinity between the two curves as:

$$\alpha(\mu(t), \nu(\tau)) = \frac{LCSS(\mu(t), \nu(\tau))}{\min(T, T')}, \quad (6)$$

where the $LCSS(\mu(t), \nu(\tau))$ (Eq. (7)) indicates the quality of the matching between the curves $\mu(t)$ and $\nu(\tau)$ and measures the number of the matching

points between two curves of different lengths.

$$\begin{aligned}
 LCSS(\mu(t), \nu(\tau)) = & \\
 & \begin{cases} 0, & \text{if } T = 0 \text{ or } T' = 0, \\
 1 + LCSS(\mu(t)^{T_t-1}, \nu(\tau)^{T'_\tau-1}), & \text{if } |\mu(t) - \nu(\tau)| < \varepsilon \text{ and } |T - T'| < \delta \\
 \max\{LCSS(\mu(t)^{T_t-1}, \nu(\tau)^{T'_\tau}), LCSS(\mu(t)^{T_t}, \nu(\tau)^{T'_\tau-1})\}, & \text{otherwise} \end{cases}
 \end{aligned} \tag{7}$$

Note that the LCSS is a modification of the edit distance [11] and its value is computed within a constant time window δ and a constant amplitude ε , that control the matching thresholds. The terms $\mu(t)^{T_t}$ and $\nu(\tau)^{T'_\tau}$ denote the number of curve points up to time t and τ , accordingly. The idea is to match segments of curves by performing time stretching so that segments that lie close to each other (their temporal coordinates are within δ) can be matched if their amplitudes differ at most by ε (Fig. 2). A characteristic example of how two motion curves are matched is depicted in Figure 2.

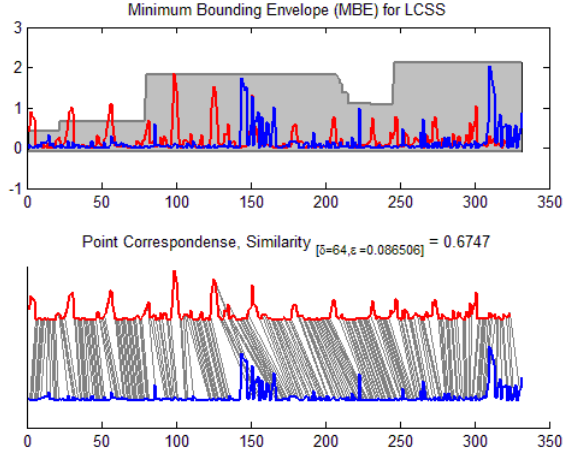


Figure 2: The LCSS matching between two motions considering that they should be within $\delta = 64$ time steps in the horizontal axis and their amplitudes should differ at most by $\varepsilon = 0.086$.

When a probe video sequence is presented, its motion trajectories $z = \{z\}_{i=1}^N$ are clustered using GMMs of various numbers of components using

the EM algorithm. The BIC criterion is employed to determine the optimal value of the number of Gaussians K , which represent the action in the probe sequence. Thus, we have a set of K mean trajectories ν_k , $k = 1, \dots, K$ modeling the probe action, whose likelihood is given by:

$$L(z) = \prod_{i=1}^N \ln \sum_{k=1}^K \pi_k \mathcal{N}(z_i; \nu_k, \Sigma_k), \quad (8)$$

where Σ_k is the covariance matrix for the k^{th} component.

Recognition of the action present in the probe video sequence is performed by assigning the probe action to the action of the labeled sequence which is most similar. As both the probe sequence and the j^{th} labeled video sequence of the p^{th} action in the training set are represented by a number of mean curves $\nu = \{\nu_i\}_{i=1}^K$ and $\mu_j^p = \{\mu_{jk}^p\}_{k=1}^{K_j^p}$ respectively, the overall distance between them is computed by:

$$d(\mu_j^p, \nu) = \sum_{k=1}^{K_j^p} \sum_{\ell=1}^K \pi_{jk}^p \pi_\ell [1 - \alpha(\mu_{jk}^p(t), \nu_\ell(\tau))], \quad (9)$$

where π_{jk}^p and π_ℓ are the GMM mixing proportions for the labeled and probe sequence, respectively, that is $\sum_k \pi_{jk}^p = 1$ and $\sum_\ell \pi_\ell = 1$. The probe sequence ν is categorized with respect to its minimum distance from an already learned action:

$$[j^*, p^*] = \arg \min_{j,p} d(\mu_j^p, \nu). \quad (10)$$

The canonical time warping (CTW) [9] solves the problem of spatio-temporal alignment of human motion between two time series. Based on dynamic time warping, the algorithm in [11] finds the temporal alignment of two subjects maximizing the spatial correlation between them. Given two time series $\mathcal{C}_1 = [c_1(0), \dots, c_1(T)]$ and $\mathcal{C}_2 = [c_2(0), \dots, c_2(T')]$ canonical time warping minimizes the following energy function:

$$J_{ctw}(\mathbf{W}_{\mathcal{C}_1}, \mathbf{W}_{\mathcal{C}_2}, \mathbf{V}_{\mathcal{C}_1}, \mathbf{V}_{\mathcal{C}_2}) = \|\mathbf{V}_{\mathcal{C}_1}^T \mathcal{C}_1 \mathbf{W}_{\mathcal{C}_1}^T - \mathbf{V}_{\mathcal{C}_2}^T \mathcal{C}_2 \mathbf{W}_{\mathcal{C}_2}^T\|_F^2, \quad (11)$$

where $\mathbf{W}_{\mathcal{C}_1}$ and $\mathbf{W}_{\mathcal{C}_2}$ are binary selection matrices that need to be inferred to align \mathcal{C}_1 and \mathcal{C}_2 , and $\mathbf{V}_{\mathcal{C}_1}$, $\mathbf{V}_{\mathcal{C}_2}$ parameterize the spatial warping by projecting sequences into the same coordinate system.

Dimensionality reduction methods [10] may be employed in order to reduce the dimension of the motion curves and to enforce them to be of equal length. In the experiments, Principal Components Analysis (PCA) [58] was chosen as a simple linear method but any other non-linear technique [10] could also be applied. When PCA is employed the time ordering is suppressed and trajectories are then transformed into feature vectors. In that case, the Bhattacharyya distance [11] is (among others) an appropriate matching measure.

Let v_1 and v_2 , be two feature vectors following Gaussian distributions, with means μ_1 and μ_2 and covariance matrices Σ_1 and Σ_2 , respectively. The Bhattacharyya distance has the form :

$$d_B(v_{1j}^p, v_2) = \frac{1}{8}(\mu_1 - \mu_2)^T \left(\frac{\Sigma_1 - \Sigma_2}{2} \right)^{-1} (\mu_1 - \mu_2) + \frac{1}{2} \ln \left(\frac{\frac{\Sigma_1 - \Sigma_2}{2}}{2\sqrt{|\Sigma_1||\Sigma_2|}} \right). \quad (12)$$

To perform the match, one can project a probe video feature vector $v_2 = \{v_{2i}\}_{i=1}^K$ onto the subspace of the training feature vectors $v_{1j}^p = \{v_{1jk}^p\}_{k=1}^{K_j^p}$ and assume the label of the closest projection of the training feature vector. For Gaussian mixture models, we define the Bhattacharyya distance as:

$$d_{GMM}(v_{1j}^p, v_2) = \sum_{k=1}^{K_j^p} \sum_{\ell=1}^K \pi_{jk}^p \pi_\ell d_B(v_{1jk}^p, v_{2\ell}), \quad (13)$$

where π_{jk}^p and π_ℓ are the GMM mixing proportions for the labeled and probe sequence, respectively. This is common in GMM modeling [59]. The probe feature vector v_2 is categorized with respect to the minimum distance from an already learned action:

$$[j^*, p^*] = \arg \min_{j,p} d_{GMM}(v_{1j}^p, v_2). \quad (14)$$

The overall approach for learning an action and categorizing a probe are summarized in Algorithm 1 and Algorithm 2, respectively. The steps inside the parenthesis indicate the extra steps when PCA is employed.

4. Experimental Results

In what follows, we refer to our mixtures of trajectories action recognition method by the acronym TMAR. We evaluated the proposed method

Algorithm 1 Action learning

Input: Training video sequences

Output: GMMs summarizing each action in each sequence

- **For** each action
 - **For** each video sequence representing the action
 - * Compute the optical flow at each pixel and generate half-wave rectified features [13].
 - * Construct the motion curves by concatenating the optical flow features.
 - * (Perform dimensionality reduction of the motion curves.)
 - * Cluster the motion curves by training GMMs with varying number of components and select the model maximizing the BIC criterion [7].

on action recognition by conducting a set of experiments over publicly available datasets. First, we applied the algorithm to the Weizmann human action dataset [3]. The Weizmann dataset is a collection of 90 low-resolution videos, which consists of 10 different actions (i.e., run, walk, skip, jumping jack, jump forward, jump in place, gallop sideways, wave with two hands, wave with one hand, and bend), performed by nine different people. The videos were acquired with a static camera and contain uncluttered background. Nevertheless, the dataset provides a good evaluation context for testing the performance of the proposed algorithm, due to the periodicity of the actions. Figure 3 illustrates some sample frames from the Weizmann dataset.

To test the proposed method on action recognition we adopted the leave-one-out scheme. We learned the model parameters from the videos of eight subjects and tested the recognition results on the remaining video sequences. The procedure was repeated for all sets of video sequences and the final result is the average of the individual results. The optimal number of mixture components K_j^p for the j^{th} video sequence, $j = 1, \dots, S_p$ of the p^{th} action $p = 1, \dots, A$ is found by employing the BIC criterion. The value of BIC was computed for $K_j^p = 1$ to the square root of the maximum number of motion

Algorithm 2 Action categorization

Input: A probe video sequence to be categorized and the GMMs summarizing the actions in the training sequences

Output: Action label

- Compute the optical flow at each pixel of the probe sequence and generate half-wave rectified features [13].
- Construct the motion curves by concatenating the optical flow features.
- (Project the motion curves to the reduced training curve space.)
- Cluster the motion curves by training GMMs with varying number of components and select the model maximizing the BIC criterion [7].
- Compute the distances between the GMM of the probe sequence and each GMM of the learnt actions.
- Classify the probe sequence using a nearest neighbor classifier.



Figure 3: Sample frames from video sequences of the Weizmann dataset [3].

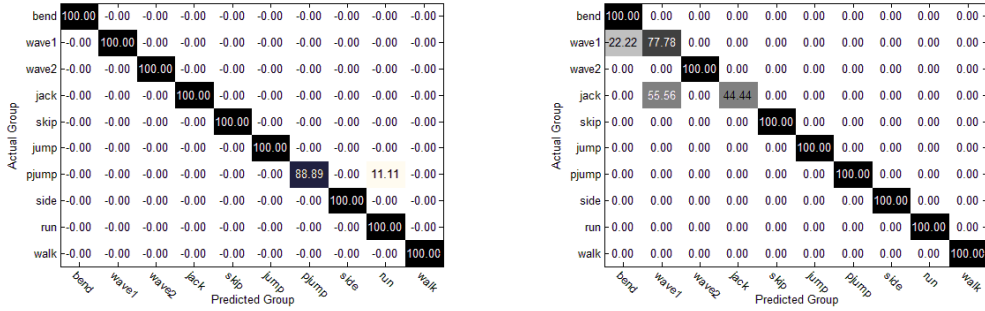
curves.

As shown in Table 1, the average correct classification of the algorithm on this dataset is 98.8%, while it reaches 100% when the proposed method with PCA is utilized. However, the average correct classification falls to 92.2%, when the CTW is utilized. In the same table, the results of [3, 60, 21, 25, 23,

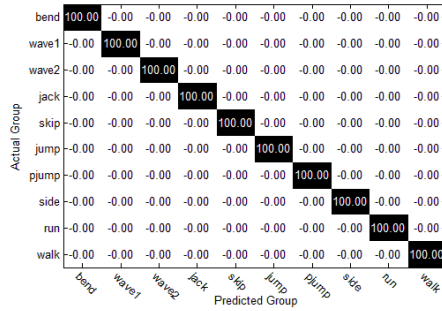
20, 14] are taken from the original papers. Notice that all motion curves are reduced to a length that explains the 90% of the eigenvalue sum, which results in a reduced curve length of 50 time instances with respect to the original 3.000 time instances. Note that better results are achieved with respect to four out of seven state-of-the-art methods for the standard method, whereas for the TMAR(PCA) the highest performance on this dataset is achieved. The proposed method provided only one erroneous categorization as one *jump-in-place* (pjump) action was incorrectly categorized as *run*. It appears that in this case the number of Gaussian components K_j^p computed by the BIC criterion was not optimal. Figure 4 depicts the confusion matrices for the TMAR(LCSS), TMAR(CTW) and TMAR(PCA) approaches.

In order to examine the behavior and the consistency of the method to the BIC criterion, we have also applied the algorithm without using BIC but having a predetermined number of Gaussian components for both the training and the test steps. Therefore, we fixed the number of Gaussians K_j^p to values varying from one to the square root of the maximum number of the motion curves and executed the algorithm. More specifically, for the proposed method, when the LCSS metric is employed, for $K_j^p = 1$, $K_j^p = 2$ and $K_j^p = 3$ recognition rates of 100% are attained and performance begins to decrease for $K_j^p \geq 4$. This is not surprising since the majority of the mixture components provided by the BIC criterion is equal to two. In the case where CTW alignment is employed, the average recognition accuracy begins to fall for $K_j^p \geq 2$. When PCA is employed the recognition is perfect and begins to decrease for $K_j^p \geq 6$. In Figure 5, the recognition accuracy for this dataset with respect to the number of Gaussian components is depicted. As the number of curves representing each action is relatively small (30–60 curves per action), a large number of Gaussian components may lead to model overfitting.

We have further assessed the performance rate of our method by conducting experiments on the KTH dataset [4]. This dataset consists of 2.391 sequences and contains six types of human actions such as walking, jogging, running, boxing, hand waving, and hand clapping. These actions are repeatedly performed by 25 different people in four different environments: outdoors (s1), outdoors with scale variation (s2), outdoors with different clothes (s3), and indoors (s4). The video sequences were acquired using a static camera and include a uniform background. The average length of the video sequences is four seconds, while they were downsampled to a spatial resolution of 160×120 pixels. Figure 6 depicts sample snapshots from the



(a) TMAR(LCSS), accuracy = 98.8% (b) TMAR(CTW), accuracy = 92.2%



(c) TMAR(PCA), accuracy = 100%

Figure 4: Confusion matrices of the classification results for the Weizmann dataset for (a) the proposed method denoted by TMAR(LCSS), (b) the the proposed method using the CTW alignment, denoted by TMAR(CTW), and (c) the proposed method using PCA, denoted by TMAR(PCA), for the estimation of the number of components using the BIC criterion.

KTH dataset.

We tested the action recognition performance of the proposed method by using a leave-one-out cross validation approach. Accordingly, the model from the videos of 24 subjects was learned while the algorithm was tested on the remaining subjects and averaged the recognition results. The confusion matrices over the KTH dataset for this leave-one-out approach are shown in Figure 7. A recognition rate of 96.7% was achieved when only the BIC criterion was employed in conjunction with the LCSS metric, 93.8% when the CTW alignment is employed, and 98.3% using PCA. In addition, comparison

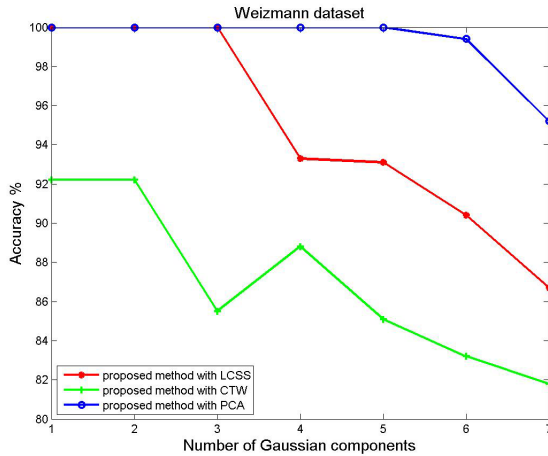


Figure 5: The recognition accuracy with respect to the number of Gaussian components for the Weizmann dataset.

Table 1: Recognition accuracy over the Weizmann dataset.

Method	Year	Accuracy (%)
Blank <i>et al.</i> [3]	2005	100.0
Chaudhry <i>et al.</i> [25]	2009	95.7
Fathi and Mori [14]	2008	100.0
Jhuang <i>et al.</i> [20]	2007	98.8
Lin <i>et al.</i> [23]	2009	100.0
Niebles <i>et al.</i> [21]	2008	90.0
Seo and Milanfar [60]	2011	97.5
TMAR(LCSS-BIC)	2013	98.8
TMAR(CTW-BIC)	2013	92.2
TMAR(PCA-BIC)	2013	100.0

of the proposed method with other state-of-the-art methods is reported in Table 2. The results of [4, 20, 14, 21, 23, 60, 15, 43, 40, 50, 28, 41] are obtained from the original papers. Note that, the TMAR approach provides the more accurate recognition rates. The TMAR(LCSS) approach attains high action classification accuracy as the BIC criterion determines the optimal value of Gaussians K_j^p for this dataset. Figure 8 depicts the accuracy rate for the

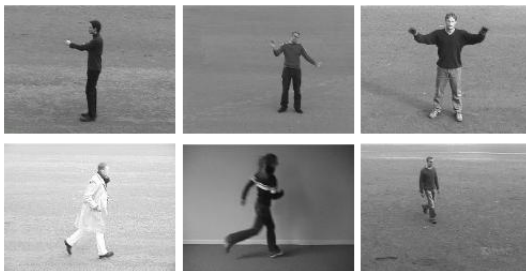


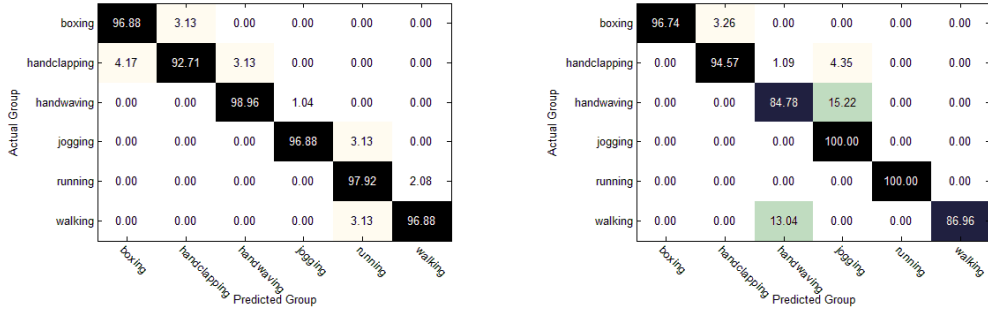
Figure 6: Sample frames from video sequences of the KTH dataset [4].

Table 2: Recognition results over the KTH dataset.

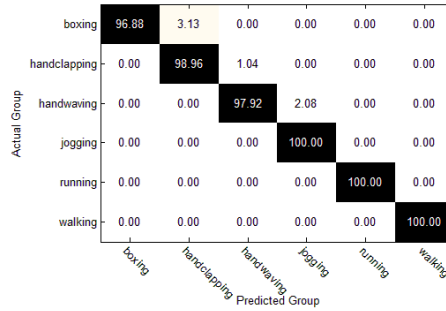
Method	Year	Accuracy (%)
Schuldt <i>et al.</i> [4]	2004	71.7
Jhuang <i>et al.</i> [20]	2007	90.5
Fathi and Mori [14]	2008	90.5
Niebles <i>et al.</i> [21]	2008	83.3
Lin <i>et al.</i> [23]	2009	95.8
Seo and Milanfar [60]	2011	95.1
Wang <i>et al.</i> [15]	2011	94.2
Wu <i>et al.</i> [43]	2011	94.5
Le <i>et al.</i> [40]	2011	93.9
Tran <i>et al.</i> [50]	2012	97.8
Yan and Luo [28]	2012	93.9
Sadanand and Corso [41]	2012	98.2
TMAR(LCSS-BIC)	2013	96.7
TMAR(CTW-BIC)	2013	93.8
TMAR(PCA-BIC)	2013	98.3

TMAR(LCSS), TMAR(CTW) and TMAR(PCA) approaches with respect to the number of mixture components. As the number of Gaussians is $K_j^p \geq 3$ for the TMAR(LCSS), $K_j^p \geq 5$ for the TMAR(CTW) and $K_j^p \geq 4$ for the TMAR(PCA) the accuracy rate drastically falls. This fact indicates the dependency of the recognition accuracy over the number of Gaussian components as an action is represented by few motion curves.

We have also applied our algorithm to the UCF Sports dataset [5]. This



(a) TMAR(LCSS), accuracy = 96.7% (b) TMAR(CTW), accuracy = 93.8%



(c) TMAR(PCA), accuracy = 98.3%

Figure 7: Confusion matrices of the classification results for the KTH dataset for (a) the proposed method denoted by TMAR(LCSS), (b) the the proposed method using the CTW alignment, denoted by TMAR(CTW), and (c) the proposed method using PCA, denoted by TMAR(PCA), for the estimation of the number of components using the BIC criterion.

dataset consists of nine main actions such as diving, golf-swinging, kicking, lifting, horse riding, running, skating, swinging and walking. The dataset contains approximately 200 video sequences at a resolution of 720×480 pixels, which are captured in natural environment with a wide range of scenes and viewpoints. Figure 9 depicts some sample frames from the UCF Sports dataset.

To test the proposed method on action recognition we also adopted the leave-one-out scheme. In Figure 10 are depicted the confusion matrices for the TMAR(LCSS), TMAR(CTW) and the TMAR(PCA) approaches.

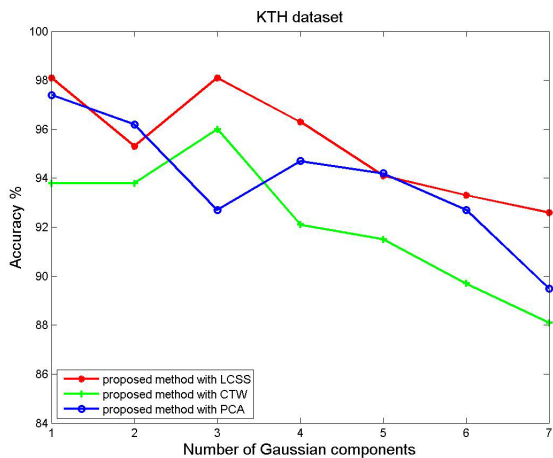


Figure 8: The recognition accuracy with respect to the number of Gaussian components for the KTH dataset.

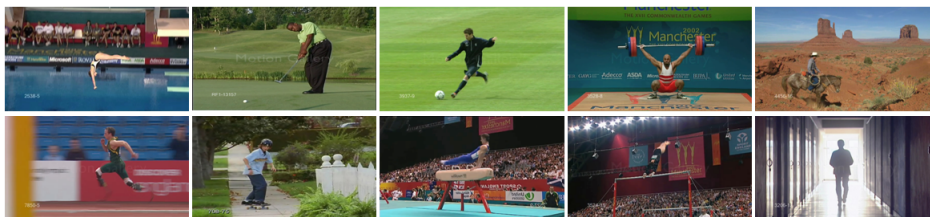
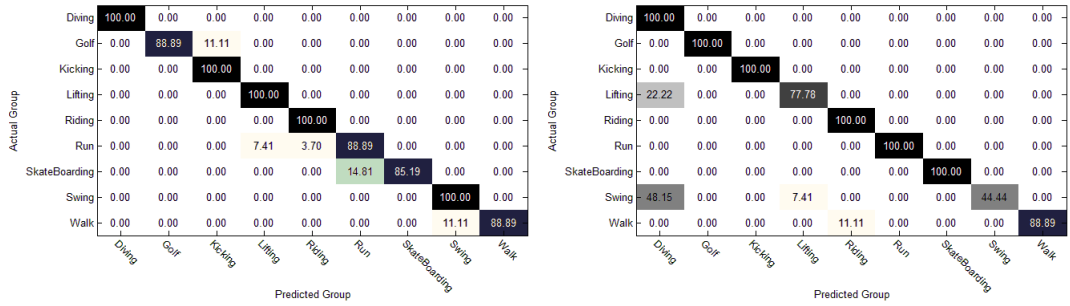


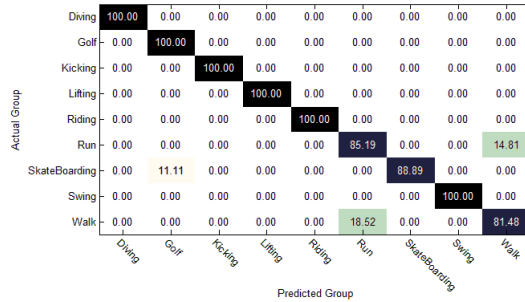
Figure 9: Sample frames from video sequences of the UCF Sports dataset [5].

TMAR(LCSS) achieves 94.6% recognition accuracy with optimal number of components (BIC criterion) and 90.1% when the CTW alignment is employed. We also achieve the highest recognition accuracy of 95.1% when the proposed method uses PCA. In Figure 11, the dependency of the recognition accuracy with respect to the number of the Gaussian components is shown. Note that, for all three approaches as the number of components increases the recognition accuracy decreases, which may occur due to model overfitting. In the case where $K_j^p = 3$ all three approaches reach the highest peak of the graph. For $K_j^p \geq 4$ the recognition accuracy begins to decrease. Table 3, shows the comparison between our TMAR approach, the baseline method using the BIC criterion in conjunction with the LCSS metric and the CTW

alignment, the proposed method with PCA and previous approaches on the UCF Sports dataset. The results of [5, 44, 15, 43, 40, 50, 28, 41] are obtained from the original papers. As it can be observed, the TMAR(PCA) approach preforms better than all the other methods, while TMAR(LCSS) performs better for seven out of eight of the other methods. On the other hand, TMAR(CTW) has the less desirable performance as it outreaches four out of eight of the other methods on the same dataset.



(a) TMAR(LCSS), accuracy = 94.6% (b) TMAR(CTW), accuracy = 90.1%



(c) TMAR(PCA), accuracy = 95.1%

Figure 10: Confusion matrices of the classification results for the UCF Sports dataset for (a) the proposed method denoted by TMAR(LCSS), (b) the the proposed method using the CTW alignment, denoted by TMAR(CTW), and (c) the proposed method using PCA, denoted by TMAR(PCA), for the estimation of the number of components using the BIC criterion.

Finally, we have put our algorithm to test with the UCF YouTube dataset [6]. The UCF YouTube human action data set contains 11 action categories

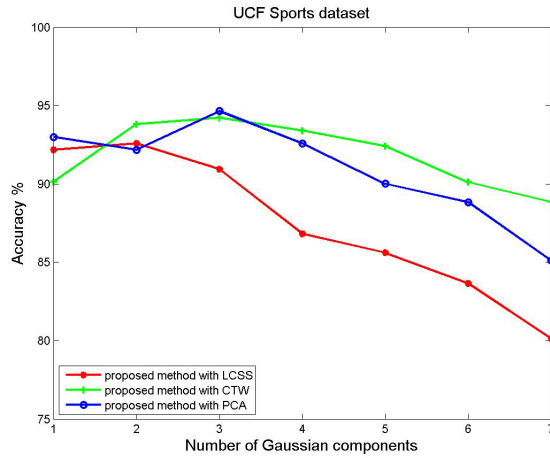


Figure 11: The recognition accuracy with respect to the number of Gaussian components for the UCF Sports dataset.

Table 3: Recognition results over the UCF Sport dataset.

Method	Year	Accuracy (%)
Rodriguez <i>et al.</i> [5]	2008	69.2
Kovaska and Grauman [44]	2010	87.3
Wang <i>et al.</i> [15]	2011	88.2
Wu <i>et al.</i> [43]	2011	91.3
Le <i>et al.</i> [40]	2011	86.5
Tran <i>et al.</i> [50]	2012	91.6
Yan and Luo [28]	2012	90.7
Sadanand and Corso [41]	2012	95.0
TMAR(LCSS-BIC)	2013	94.6
TMAR(CTW-BIC)	2013	90.1
TMAR(PCA-BIC)	2013	95.1

such as basketball shooting, biking, diving, golf swinging, horse riding, soccer juggling, swinging, tennis swinging, trampoline jumping, volleyball spiking, and walking with a dog. This data set includes actions with large variation in camera motion, object appearance and pose and scale. It also contains viewpoint and illumination changes, and spotty background. The video se-

quences are grouped into 25 groups of at least four actions each for each category, whereas the videos in the same group may share common characteristics such as similar background or actor. Representative frames of this data set are shown in Figure 12.

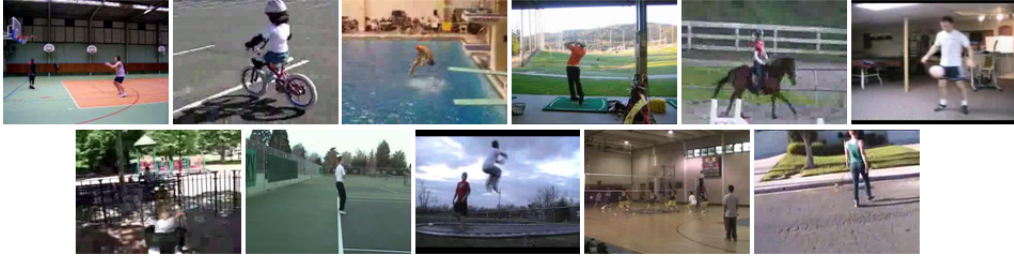
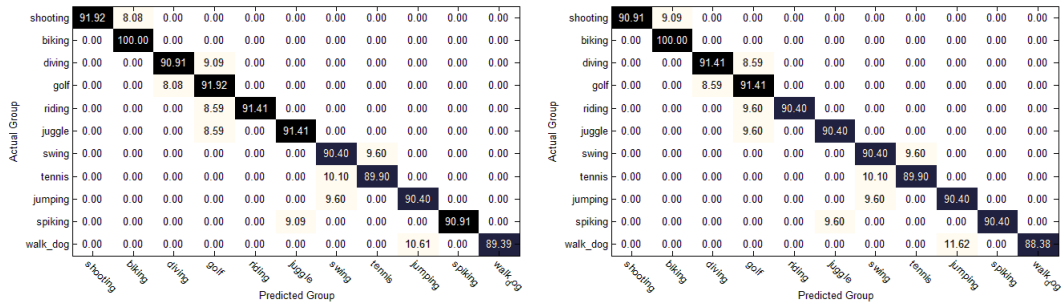


Figure 12: Sample frames from video sequences of the UCF YouTube action dataset [6].

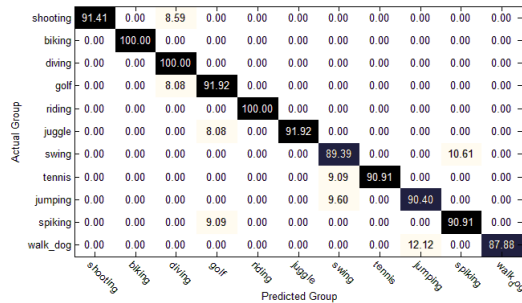
In order to assess our method we have used the leave-one-out cross validation scheme. In Figure 13 the confusion matrices for the TMAR(LCSS), TMAR(CTW) and TMAR(PCA) approaches are shown. We achieve a recognition rate of 91.7% when the LCSS metric is employed and having estimated the Gaussian components using the BIC criterion. We also achieve 91.3% when the CTW alignment is employed and 93.2% when using PCA. In Table 4, comparisons with other state-of-the-art methods for this dataset are reported. The results of [6, 39, 40, 61] were copied from the original papers. As it can be seen, our algorithm achieves the highest recognition accuracy amongst all the others.

The performance of the proposed method with respect to the number of the Gaussian components is depicted in Figure 14. For TMAR(LCSS) the recognition accuracy begins to decrease for $K_j^p \geq 1$ and exhibits the worst performance than the other two approaches. The TMAR(CTW) approach decreases for $K_j^p \geq 2$ while TMAR(PCA) reaches its peak for $K_j^p = 4$ and then it begins to decrease. Note that, the best approach tends to be, attained by TMAR(PCA) which reaches a recognition accuracy of 91%.

In the recognition step, in our implementation of the LCSS (7) the parameters δ and ϵ were optimized using 10-fold cross validation for all four datasets. These parameters need to be determined for each data set separately since each data set perform different types of actions. However, after we have determined the parameters no further action needs to be taken. To classify a new unknown sequence, we have already learned the parameters



(a) TMAR(LCSS), accuracy = 91.7% (b) TMAR(CTW), accuracy = 91.3%



(c) TMAR(PCA), accuracy = 93.2%

Figure 13: Confusion matrices of the classification results for the UCF YouTube dataset for (a) the proposed method denoted by TMAR(LCSS), (b) the the proposed method using the CTW alignment, denoted by TMAR(CTW), and (c) the proposed method using PCA, denoted by TMAR(PCA), for the estimation of the number of components using the BIC criterion.

from the learning step and thus we are able to recognize the new action. For the Weizmann dataset, for all actions, parameter δ was determined to be 1% of the trajectories' lengths, and parameter ε was determined as the smallest standard deviation of the two trajectories to be compared. For the other datasets, Table 5, Table 6 and 7 show the optimal values per action as they have resulted after the cross validation process. Note that, the values in Table 5 for both δ and ε are consistently small. However, the handclapping and walking actions have larger values for ε parameter than the other actions, which may be due to the large vertical movement of the subject

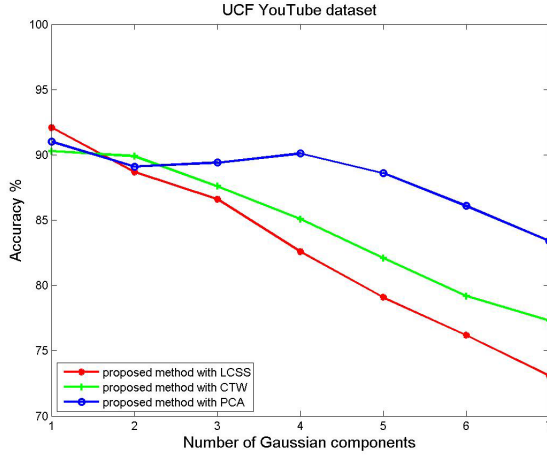


Figure 14: The recognition accuracy with respect to the number of Gaussian components for the UCF YouTube dataset.

Table 4: Recognition results over the UCF YouTube dataset.

Method	Year	Accuracy (%)
Liu <i>et al.</i> [6]	2009	71.2
Ikizler-Cinbis and Sclaroff [61]	2010	75.2
Le <i>et al.</i> [40]	2011	75.8
Wang <i>et al.</i> [39]	2011	84.2
TMAR(LCSS-BIC)	2013	94.6
TMAR(CTW-BIC)	2013	90.1
TMAR(PCA-BIC)	2013	95.1

between consecutive frames. On the other hand, the actions in the UCF Sport dataset holds large movements from one frame to the other for both horizontal and vertical axes, which is the main reason why the actions show large variances between the values of δ and ε (Table 6). Finally, the actions in the UCF YouTube dataset have a uniform distributed representation of the parameters δ and ε , since the parameter δ is determined as the 10% of the mean trajectories length for the most of the actions and the mean of the parameter ε is varies in the 15% of the standard deviation of the two trajectories to be compared.

Table 5: Parameters δ and ε for the KTH dataset estimated using cross validation.

Action	TMAR(LCSS)	
	δ	ε
boxing	10^{-3}	10^{-4}
handclapping	10^{-2}	10^{-2}
handwaving	3×10^{-1}	10
jogging	5×10^{-3}	3×10^{-1}
running	3×10^{-2}	5×10^{-2}
walking	1	12

Table 6: Parameters δ and ε for the UCF Sports dataset estimated using cross validation.

Action	TMAR(LCSS)	
	δ	ε
diving	1	2.1
golf	2.01	6.1
kicking	10	15
lifting	11	10
riding	0.1	15
run	0.1	12
skateboarding	1.4	13
swing	0.6	20
walk	0.1	10

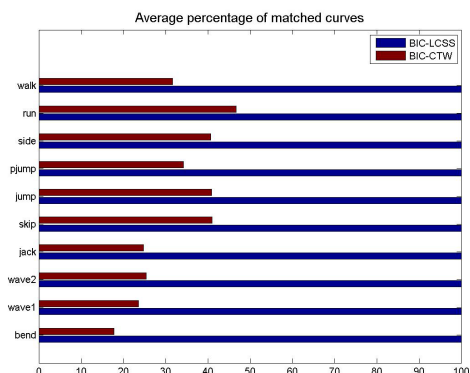
The average percentage of matched curves for the TMAR(LCSS) and TMAR(CTW) approach in the case where the BIC criterion is employed to determine the number of Gaussian components for all four dataset is depicted in Figure 15. As it can be observed, the TMAR(LCSS) method appears to match a larger part of curves for the same dataset than the TMAR(CTW) approach, which is the reason why TMAR(LCSS) performs better than TMAR(CTW).

Labeling a new video sequence does not require any exhaustive search in order to determine the number of the Gaussian components. Exhaustive search is performed once in the learning step while in the testing step classification is performed using the predetermined number of components. In Figure 16, the execution times using the BIC criterion are depicted in

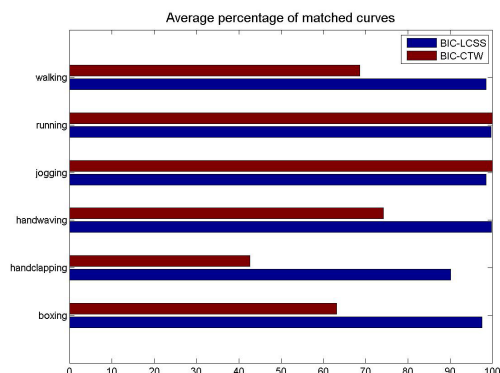
Table 7: Parameters δ and ε for the UCF Youtube dataset estimated using cross validation.

Action	TMAR(LCSS)	
	δ	ε
shooting	20	20
biking	10	10
diving	10	15
golf	20	10
riding	10	5
juggle	10	15
swing	10	5
tennis	10	10
jumping	10	5
spiking	10	30
walk dog	10	20

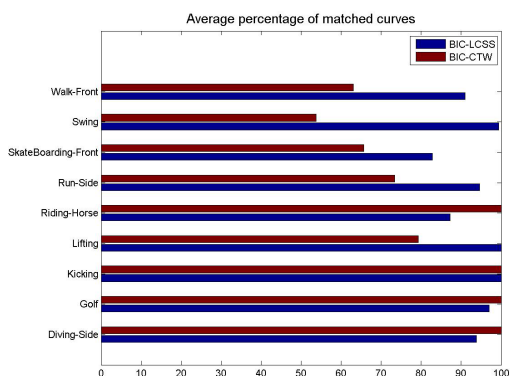
order to determine the number of the Gaussian components, for all three cases, when using the LCSS metric, the CTW alignment and PCA, for all four datasets. For the Weizmann dataset, when PCA is used, the execution time drastically falls below one second per action. On the other hand, TMAR(LCSS) requires the highest execution time, which needs six seconds to recognize the action pjump. Note that, the use of PCA speeds up the execution time for recognizing a single action in all datasets since feature vectors of smaller lengths are being used. In the KTH dataset, TMAR(CTW) requires the highest execution time (needs nine seconds to recognize two out of six actions) and while TMAR(LCSS) takes less than six seconds for one action. However, in UCF Sport dataset TMAR(LCSS) and TMAR(CTW) both have the same upper bound of eight seconds to recognize an action. Finally, in UCF YouTube dataset, the average execution time to recognize an action ranges from two to nine seconds when TMAR(LCSS) approach is used. In the case where TMAR(PCA) is used the upper bound to recognize an action is five seconds in UCF Sports and UCF YouTube datasets, while in Weizmann and in KTH is less than a second. This makes the algorithm capable to adapt to any real video sequence and recognize an action really fast.



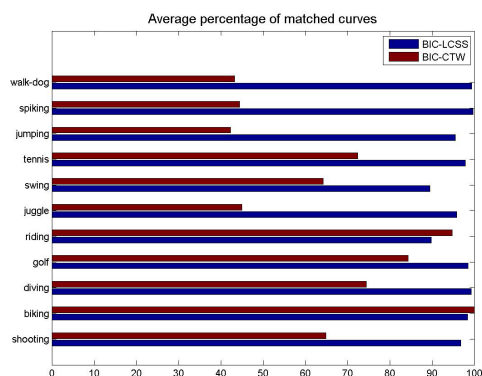
(a) Weizmann dataset



(b) KTH dataset



(c) UCF Sports dataset

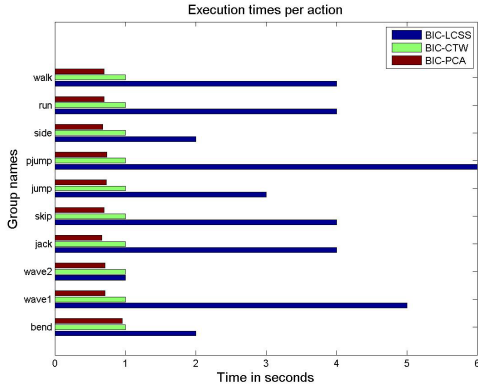


(d) UCF YouTube dataset

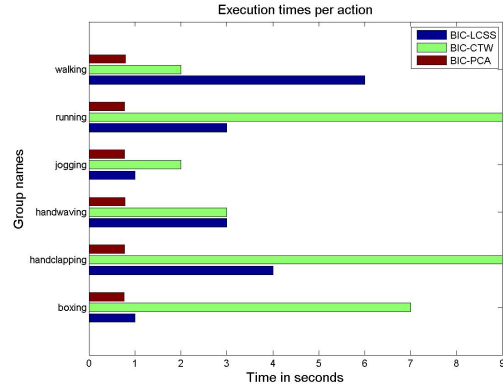
Figure 15: Average percentage of matched curves for TMAR(LCSS) and TMAR(CTW), when the BIC criterion is used, for (a) Weizmann, (b) KTH, (c) UCF Sports and (d) UCF YouTube datasets, respectively.

5. Conclusion

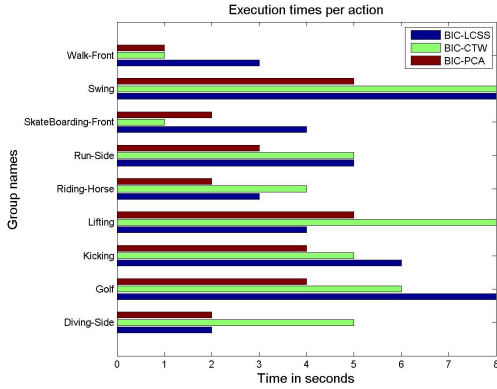
In this paper, we presented an action recognition approach, where actions are represented by a set of motion curves generated by a probabilistic model. The performance of the extracted motion curves is interpreted by computing similarities between the motion trajectories, followed by a classification scheme. The large size of motion curves was reduced via PCA and after noise removal a reference database of feature vectors is obtained. Although a perfect recognition performance is accomplished with a fixed number of



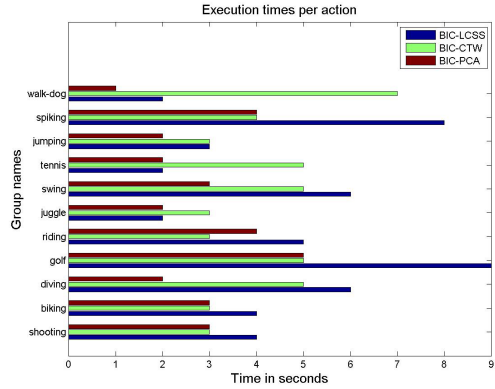
(a) Weizmann dataset



(b) KTH dataset



(c) UCF Sports dataset



(d) UCF YouTube dataset

Figure 16: Execution times per action in seconds for TMAR(LCSS), TMAR(CTW) and TMAR(PCA), when the BIC criterion is used, for (a) Weizmann, (b) KTH, (c) UCF Sports and (d) UCF YouTube datasets, respectively.

Gaussian mixtures, there are still some open issues in feature representation.

Our results classifying activities in four publicly available datasets show that the use of PCA have a significant impact on the performance of the recognition process, while it semantically speeds up the behavior of the proposed algorithm. We demonstrated the effectiveness of the optimal recognition model by using the BIC criterion to determine the number of the Gaussian components. Finally, our algorithm is free of any constraints in the trajectories lengths. Although the proposed method yielded encouraging

results in standard action recognition datasets, it is requirement of a challenging task of performing motion detection, background subtraction, and action recognition in natural and cluttered environments.

References

- [1] P. Matikainen, M. Hebert, R. Sukthankar, Trajectons: Action recognition through the motion analysis of tracked features, in: Workshop on Video-Oriented Object and Event Classification, 2009, pp. 514–521.
- [2] M. Raptis, I. Kokkinos, S. Soatto, Discovering discriminative action parts from mid-level video representations, in: Proc. IEEE Computer Society Conference on Computer Vision and Pattern Recognition, Providence, Rhode Island, 2012, pp. 1242–1249.
- [3] M. Blank, L. Gorelick, E. Shechtman, M. Irani, R. Basri, Actions as space-time shapes, in: Proc. 10th IEEE International Conference on Computer Vision, Beijing, China, 2005, pp. 1395–1402.
- [4] C. Schuldt, I. Laptev, B. Caputo, Recognizing human actions: a local SVM approach, in: Proc. 17th International Conference on Pattern Recognition, Cambridge, UK, 2004, pp. 32–36.
- [5] M. D. Rodriguez, J. Ahmed, M. Shah, Action MACH a spatio-temporal maximum average correlation height filter for action recognition., in: Proc. IEEE Computer Society Conference on Computer Vision and Pattern Recognition, Anchorage, Alaska, 2008, pp. 1–8.
- [6] J. Liu, J. Luo, M. Shah, Recognizing realistic actions from videos ”in the wild”, Proc. IEEE Computer Society Conference on Computer Vision and Pattern Recognition (2009) 1–8.
- [7] C. M. Bishop, Pattern Recognition and Machine Learning, Springer, 2006.
- [8] M. Vlachos, D. Gunopoulos, G. Kollios, Discovering similar multidimensional trajectories, in: Proc. 18th International Conference on Data Engineering, San Jose, California, USA, 2002, pp. 673–682.

- [9] F. Zhou, F. D. la Torre, Canonical time warping for alignment of human behavior, in: *Advances in Neural Information Processing Systems Conference*, 2009.
- [10] J. B. Tenenbaum, V. D. Silva, J. C. Langford, A global geometric framework for nonlinear dimensionality reduction, *Science* 290 (2000) 2319–2323.
- [11] S. Theodoridis, K. Koutroumbas, *Pattern Recognition*, 4th Edition, Academic Press, 2008.
- [12] M. Vrigkas, V. Karavasilis, C. Nikou, I. Kakadiaris, Action recognition by matcing clustered trajectories of motion vectors, in: *Proc. 8th International Conference on Computer Vision Theory and Application*, Barcelona, Spain, 2013, pp. 112–117.
- [13] A. A. Efros, A. C. Berg, G. Mori, J. Malik, Recognizing action at a distance, in: *Proc. 9th IEEE International Conference on Computer Vision*, Vol. 2, Nice, France, 2003, pp. 726–733.
- [14] A. Fathi, G. Mori, Action recognition by learning mid-level motion features, in: *Proc. IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, Anchorage, Alaska, 2008, pp. 1–8.
- [15] Y. Wang, G. Mori, Hidden part models for human action recognition: probabilistic versus max margin, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 33 (7) (2011) 1310–1323.
- [16] V. Karavasilis, K. Blekas, C. Nikou, Motion segmentation by model-based clustering of incomplete trajectories, in: *Proc. 2011 European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases (ECML-PKDD)*, 2011, pp. 146–161.
- [17] V. Karavasilis, K. Blekas, C. Nikou, A novel framework for motion segmentation and tracking by clustering incomplete trajectories, *Computer Vision and Image Understanding* 116 (11) (2012) 1135–1148. doi:10.1016/j.cviu.2012.07.004.
- [18] J. K. Aggarwal, M. S. Ryoo, Human activity analysis: a review, *ACM Compututing Surveys* 43 (3) (2011) 1–43.

- [19] R. Poppe, A survey on vision-based human action recognition, *Image and Vision Computing* 28 (6) (2010) 976–990.
- [20] H. Jhuang, T. Serre, L. Wolf, T. Poggio, A biologically inspired system for action recognition, in: *Proc. IEEE International Conference on Computer Vision*, Rio de Janeiro, Brazil, 2007, pp. 1–8.
- [21] J. C. Niebles, H. Wang, L. Fei-Fei, Unsupervised learning of human action categories using spatial-temporal words, *International Journal of Computer Vision* 79 (3) (2008) 299–318.
- [22] B. T. Morris, M. M. Trivedi, Trajectory learning for activity understanding: unsupervised, multilevel, and long-term adaptive approach, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 33 (11) (2011) 2287–2301.
- [23] Z. Lin, Z. Jiang, L. S. Davis, Recognizing actions by shape-motion prototype trees, in: *Proc. IEEE International Conference on Computer Vision*, Miami, Florida, USA, 2009, pp. 444–451.
- [24] N. Dalal, B. Triggs, C. Schmid, Human detection using oriented histograms of flow and appearance, in: *Proc. European Conference on Computer Vision*, Graz, Austria, 2006, pp. 428–441.
- [25] R. Chaudhry, A. Ravichandran, G. D. Hager, R. Vidal, Histograms of oriented optical flow and Binet-Cauchy kernels on nonlinear dynamical systems for the recognition of human actions, in: *Proc. IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, Miami, Florida, USA, 2009, pp. 1932–1939.
- [26] L. Gorelick, M. Blank, E. Shechtman, M. Irani, R. Basri, Actions as space-time shapes, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 29 (12) (2007) 2247–2253.
- [27] P. Dollár, V. Rabaud, G. Cottrell, S. Belongie, Behavior recognition via sparse spatio-temporal features, in: *Proc. 14th International Conference on Computer Communications and Networks*, Beijing, China, 2005, pp. 65–72.

- [28] X. Yan, Y. Luo, Recognizing human actions using a new descriptor based on spatial-temporal interest points and weighted-output classifier, *Neurocomputing* 87 (2012) 51–61.
- [29] A. Yao, J. Gall, L. V. Gool, A Hough transform-based voting framework for action recognition, in: *Proc. IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, San Francisco, CA, USA, 2010, pp. 2061–2068.
- [30] K. Mikolajczyk, H. Uemura, Action recognition with motion-appearance vocabulary forest, in: *Proc. IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, Anchorage, Alaska, USA, 2008.
- [31] K. Schindler, L. V. Gool, Action snippets: How many frames does human action recognition require?, in: *Proc. IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, Anchorage, Alaska, USA, 2008, pp. 1–8.
- [32] S. Khamis, V. I. Morariu, L. S. Davis, A flow model for joint action recognition and identity maintenance, in: *Proc. IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, Providence, Rhode Island, 2012, pp. 1218–1225.
- [33] N. Dalal, B. Triggs, Histograms of oriented gradients for human detection, in: *Proc. IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, Providence, Rhode Island, 2005, pp. 886–893.
- [34] J. Wang, Z. Liu, Y. Wu, J. Yuan, Mining actionlet ensemble for action recognition with depth cameras, in: *Proc. IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, Providence, Rhode Island, 2012, pp. 1290–1297.
- [35] R. Li, T. Zickler, Discriminative virtual views for cross-view action recognition, in: *Proc. IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, Providence, Rhode Island, 2012, pp. 2855–2862.
- [36] B. Li, O. I. Camps, M. Sznaiier, Cross-view activity recognition using hankellets, in: *Proc. IEEE Computer Society Conference on Computer*

Vision and Pattern Recognition, Providence, Rhode Island, 2012, pp. 1362–1369.

- [37] Q. Zhou, G. Wang, Atomic action features: A new feature for action recognition, in: Proc. European Conference on Computer Vision, Firenze, Italy, 2012, pp. 291–300.
- [38] J. Sanchez-Riera, J. Cech, R. Horaud, Action recognition robust to background clutter by using stereo vision, in: Proc. European Conference on Computer Vision, Firenze, Italy, 2012, pp. 332–341.
- [39] H. Wang, A. Kläser, C. Schmid, L. Cheng-Lin, Action recognition by dense trajectories, in: Proc. IEEE Conference on Computer Vision and Pattern Recognition, Colorado Springs, United States, 2011, pp. 3169–3176.
- [40] Q. V. Le, W. Y. Zou, S. Y. Yeung, A. Y. Ng, Learning hierarchical invariant spatio-temporal features for action recognition with independent subspace analysis, in: Proc. IEEE Computer Society Conference on Computer Vision and Pattern Recognition, Washington, DC, USA, 2011, pp. 3361–3368.
- [41] S. Sadanand, J. J. Corso, Action bank: A high-level representation of activity in video, in: Proc. IEEE Computer Society Conference on Computer Vision and Pattern Recognition, Providence, Rhode Island, 2012, pp. 1234–1241.
- [42] I. Laptev, On space-time interest points, *International Journal of Computer Vision* 64 (2-3) (2005) 107–123.
- [43] X. Wu, D. Xu, L. Duan, J. Luo, Action recognition using context and appearance distribution features, in: Proc. IEEE Computer Society Conference on Computer Vision and Pattern Recognition, Washington, DC, USA, 2011, pp. 489–496.
- [44] A. Kovashka, K. Grauman, Learning a hierarchy of discriminative space-time neighborhood features for human action recognition, in: Proc. IEEE Computer Society Conference on Computer Vision and Pattern Recognition, San Francisco, CA, 2010, pp. 2046–2053.

- [45] 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.
- [46] M. Ahmad, S. W. Lee, Variable silhouette energy image representations for recognizing human actions, *Image and Vision Computing* 28 (5) (2010) 814–824.
- [47] C. Thureau, V. Hlavac, Pose primitive based human action recognition in videos or still images, in: *Proc. IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, Anchorage, Alaska, USA, 2008, pp. 1–8.
- [48] W. Yang, Y. Wang, G. Mori, Recognizing human actions from still images with latent poses, in: *Proc. IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, San Francisco, CA, 2010, pp. 2030–2037.
- [49] S. Maji, L. D. Bourdev, J. Malik, Action recognition from a distributed representation of pose and appearance, in: *Proc. IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, Colorado Springs, CO, 2011, pp. 3177–3184.
- [50] K. N. Tran, I. A. Kakadiaris, S. K. Shah, Part-based motion descriptor image for human action recognition, *Pattern Recognition* 45 (7) (2012) 2562–2572.
- [51] A. Fathi, J. K. Hodgins, J. M. Rehg, Social interactions: A first-person perspective, in: *Proc. IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, Providence, Rhode Island, 2012, pp. 1226–1233.
- [52] T. Lan, L. Sigal, G. Mori, Social roles in hierarchical models for human activity recognition, in: *Proc. IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, Providence, Rhode Island, 2012, pp. 1354–1361.
- [53] X. P. Burgos-Artizzu, P. Dollár, D. Lin, D. J. Anderson, P. Perona, Social behavior recognition in continuous video, in: *Proc. IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, Providence, Rhode Island, 2012, pp. 1322–1329.

- [54] Y. Kong, Y. Jia, Y. Fu, Learning human interaction by interactive phrases, in: Proc. European Conference on Computer Vision, Firenze, Italy, 2012, pp. 300–313.
- [55] B. D. Lucas, T. Kanade, An iterative image registration technique with an application to stereo vision, in: Proc. 7th International Joint Conference on Artificial Intelligence, Nice, France, 1981, pp. 674–679.
- [56] R. E. Kass, A. E. Raftery, Bayes factors, *Journal of the American Statistical Association* 90 (1995) 773–795.
- [57] J. Kuha, AIC and BIC: Comparisons of assumptions and performance, *Sociological Methods Research* 33 (2) (2004) 188–229.
- [58] I. T. Jolliffe, *Principal Component Analysis*, Springer Verlag, 1986.
- [59] G. Sfikas, C. Constantinopoulos, A. Likas, N. P. Galatsanos, An analytic distance metric for gaussian mixture models with application in image retrieval, in: Proc of the 15th international conference on Artificial neural networks, Warsaw, Poland, 2005, pp. 835–840.
- [60] H. J. Seo, P. Milanfar, Action recognition from one example, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 33 (5) (2011) 867–882.
- [61] N. Ikizler-Cinbis, S. Sclaroff, Object, scene and actions: combining multiple features for human action recognition, in: Proc 11th European conference on Computer vision: Part I, Heraklion, Crete, Greece, 2010, pp. 494–507.