

View-Independent Action Recognition from Temporal Self-Similarities

Imran N. Junejo, *Member, IEEE*, Emilie Dexter, Ivan Laptev, and Patrick Pérez, *Member, IEEE*

Abstract—This paper addresses recognition of human actions under view changes. We explore self-similarities of action sequences over time and observe the striking stability of such measures across views. Building upon this key observation, we develop an action descriptor that captures the structure of temporal similarities and dissimilarities within an action sequence. Despite this temporal self-similarity descriptor not being strictly view-invariant, we provide intuition and experimental validation demonstrating its high stability under view changes. Self-similarity descriptors are also shown to be stable under performance variations within a class of actions when individual speed fluctuations are ignored. If required, such fluctuations between two different instances of the same action class can be explicitly recovered with dynamic time warping, as will be demonstrated, to achieve cross-view action synchronization. More central to the current work, temporal ordering of local self-similarity descriptors can simply be ignored within a bag-of-features type of approach. Sufficient action discrimination is still retained in this way to build a view-independent action recognition system. Interestingly, self-similarities computed from different image features possess similar properties and can be used in a complementary fashion. Our method is simple and requires neither structure recovery nor multiview correspondence estimation. Instead, it relies on weak geometric properties and combines them with machine learning for efficient cross-view action recognition. The method is validated on three public data sets. It has similar or superior performance compared to related methods and it performs well even in extreme conditions, such as when recognizing actions from top views while using side views only for training.

Index Terms—Human action recognition, human action synchronization, view invariance, temporal self-similarities, local temporal descriptors.

1 INTRODUCTION

VISUAL recognition and understanding of human actions have attracted much attention over the past three decades [1], [2] and remain an active research area of computer vision. A good solution to the problem holds a yet unexplored potential for many applications, such as the search for and the structuring of large video archives, video surveillance, human-computer interaction, gesture recognition, and video editing. Recent work has demonstrated the difficulty of the problem associated with the large variation of human action data due to the individual variations of people in expression, posture, motion, and clothing; perspective effects and camera motions; illumination variations; occlusions and disocclusions; and distracting effects of scenes surroundings. Also, actions frequently involve and depend on manipulated objects, which adds another layer of variability. As a consequence, current methods often resort to restricted and simplified scenarios with simple backgrounds,

simpler kinematic action classes, static cameras, or limited view variations.

Various approaches using different constructs have been proposed over the years for action recognition. These approaches can be roughly categorized on the basis of representation used by the authors. Time evolution of human silhouettes was frequently used as an action description. For example, Bobick and Davis [3] proposed capturing the history of shape changes using temporal templates and Weinland et al. [4] extends these 2D templates to 3D action templates. Similarly, the notions of *action cylinders* [5] and *space-time shapes* [6], [7], [8] have been introduced based on silhouettes. Recently, space-time approaches analyzing the structure of local 3D patches in the video have been shown to be promising in [9], [10], [11], [12], [13]. Using space-time or other types of local features, the modeling and recognition of human motion have been addressed with a variety of machine learning techniques, such as Support Vector Machines (SVMs) [14], [15], Hidden Markov Models (HMMs) [16], [17], [18], and Conditional Random Fields (CRFs) [19], [20], [21], [22], [23].

Most of the current methods for action recognition are designed for limited view variations. A reliable and a generic action recognition system, however, has to be robust to camera parameters and different viewpoints while observing an action sequence. View variations originate from the changing and frequently unknown positions of the camera. Similar to the multiview appearance of static objects, the appearance of actions may drastically vary from one viewpoint to another. Differently from the static case, however, the appearance of actions may also be affected by the dynamic view changes of the moving camera.

- I.N. Junejo is with the Department of Computer Sciences, University of Sharjah, PO Box 27272, Sharjah, UAE. E-mail: ijunejo@sharjah.ac.ae.
- E. Dexter is with INRIA Rennes-Bretagne Atlantique, Campus Universitaire de Beaulieu, France. E-mail: emilie.dexter@inria.fr.
- I. Laptev is with INRIA Paris-Rocquencourt/ENS, 23 avenue d'Italie, 75013 Paris, France. E-mail: ivan.laptev@inria.fr.
- P. Pérez is with Thomson R&D, 1 avenue de Belle Fontaine, CS17616, F-35576 Cesson-Sévigné cedex, France. E-mail: Patrick.Perez@thomson.net.

Manuscript received 3 Dec. 2008; revised 16 Nov. 2009; accepted 1 Dec. 2009; published online 2 Mar. 2010.

Recommended for acceptance by B. Schiele.

For information on obtaining reprints of this article, please send e-mail to: tpami@computer.org, and reference IEEECS Log Number TPAMI-2008-12-0831.

Digital Object Identifier no. 10.1109/TPAMI.2010.68.

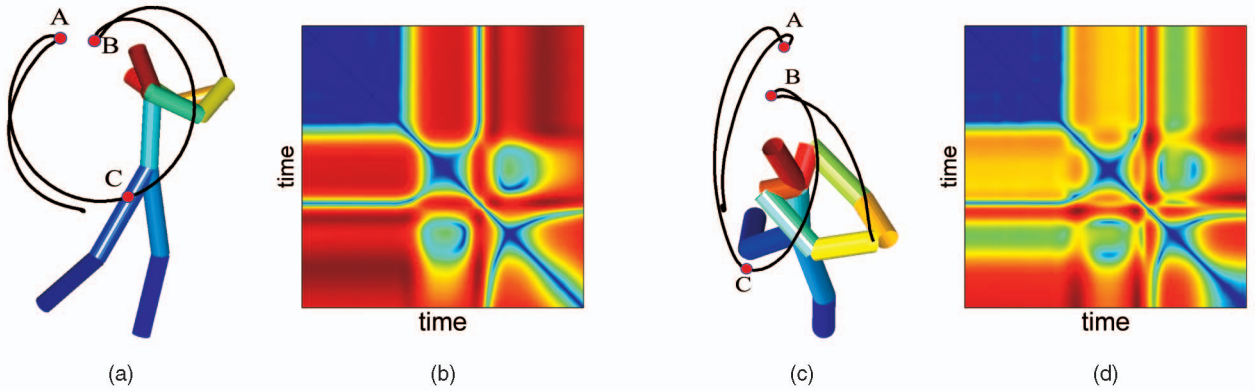


Fig. 1. Cross-view stability of trajectory-based self-similarity matrices on a simple example. (a) and (c) demonstrate, based on motion capture (MOCAP) data, a golf swing action seen from two different views. (b) and (d) represent their respective SSMs for the trajectory of one hand projected in corresponding view. Even though the two views are different, the structures or the patterns of the computed SSMs are very similar.

Multiview variations of actions have been previously addressed using epipolar geometry, such as in [5], [24], [25], [26], [27], [28], by learning poses seen from different viewpoints [29], [30], [31], [32], [33], or by a full 3D reconstruction [34], [35]. Such methods rely either on existing point correspondences between image sequences or/and on many videos representing actions in multiple views. Both of these assumptions, however, are limiting in practice due to 1) the difficulty of estimating nonrigid correspondences in videos and 2) the difficulty of obtaining sufficient video data spanning view variations for many action classes.

In this work, we address multiview action recognition from a different perspective and avoid many assumptions of previous methods. In contrast to the geometry-based methods above, we require neither the identification of body parts nor the estimation of corresponding points between video sequences. Differently from the previous view-based methods, we do not assume multiview action samples either for training or for testing.

Our approach builds upon self-similarities of action sequences over time. For a given action sequence and a given type of low level features, we compute distances between extracted features for all pairs of time frames and store results in a Self-Similarity Matrix (SSM). We claim SSMs to be stable under view changes of an action. Fig. 1 illustrates our idea with an example of a golf swing action seen from two different views. For this example, we compute SSMs as pairwise distances between all 2D points on the projected hand trajectories illustrated in Figs. 1a and 1c. Despite the view variation, close trajectory points A and B remain close in both views, while the distant trajectory points A and C have large distances in both projections. The visualizations of SSMs computed for both sequences in Figs. 1b and 1d have a striking similarity despite the different projections of the action. More generally, if body poses of an action are similar at moments t_1, t_2 , the value of $SSM(t_1, t_2)$, i.e., the distance between some action descriptors at t_1, t_2 will be low for any view of an action. On the contrary, if the body poses are different at t_1, t_2 , the value of $SSM(t_1, t_2)$ is likely to be large for most of the views and nontrivial action descriptors.

In the rest of the paper, we operationalize SSMs for human action sequences and deploy them for view-independent action recognition. In particular, we observe

similar properties of SSMs computed for different image features and use such SSMs in a complementary fashion.

The paper is organized as follows: In the next section, we review related work, with special emphasis on the relationship between SSMs and so-called Recurrence Plots (RPs). SSM can be seen as an extension of RPs. Section 3 gives a formal definition of SSM using alternative image features and reports first experiments on mocap data demonstrating its structural stability across views. Section 4 describes the proposed representation of action sequences based on local temporal SSM descriptors and demonstrate how this representation is at the same time precise, specific to a class of action and largely view-independent, by using it to synchronize (align temporally) different performances of similar actions. In Section 5, we introduce our view-independent action recognition system based on such descriptions and test it on three public data sets. These experiments demonstrate the practicality and the potential of the proposed method. Section 6 concludes the paper.

2 RELATED WORK

This paper concerns view-independent action recognition, a topic which has received considerable attention from researchers recently. To address this problem, epipolar geometry has been employed in [5], [24], [25]. Point correspondences between actions are assumed to be known for imposing fundamental matrix constraints and performing view-invariant action recognition. Rao et al. [26] show that the maxima in space-time curvature of a 3D trajectory are preserved in 2D image trajectories, and are also view-invariant. Parameswaran and Chellappa [28] propose a quasi-view-invariant approach, requiring at least five body points lying on a 3D plane or that the limbs trace a planar area during the course of an action. Recently, Shen and Foroosh [27] showed that for a moving plane, the fundamental ratios, i.e., the ratios among the elements in the upper left 2×2 submatrix of the fundamental matrix, are invariant to the camera parameters as well as its orientation and can be used for action recognition. However, obtaining automatic and reliable point correspondences for daily videos with natural human actions is a very challenging and currently unsolved problem, which limits the application of above-mentioned methods in practice.

One alternative to the geometric approach is to represent actions by samples recorded for different views. A database of poses seen from multiple viewpoints has been created in [29], [30], [31], [32]. Extracted silhouettes from a test action are matched to this database to recognize the action being performed. The drawback of these methods is that each action needs to be represented by many training samples recorded for a large and representative set of views. Other methods [35] and [34] perform a full 3D reconstruction from silhouettes seen from multiple deployed cameras. This approach requires a setup of multiple cameras or training on poses obtained from multiple views, which again restricts the applicability of methods in practice.

The approach in [33] exploits transfer learning for constructing view-stable and discriminative features for view-independent action recognition. For a pair of given views, the features are learned from a separate set of actions observed in both views. Given a new action class observed and learned in one view only, transfer learning enables recognition of instances of that class in the second view. While the use of transfer learning in [33] is interesting, the method is limited to a set of predefined views and requires training to be done separately for each pair of views. It also requires (nontarget) actions to be observed and view-tagged for multiple views. Our method avoids these limitations. We compare our results with [33] on the common benchmark in Section 5.3.

The methods most closely related to our approach are those of [36], [37], [38], [39]. For image and video matching, Shechtman and Irani [36] recently explored *local* self-similarity descriptors. The descriptors are constructed by correlating the image (or video) patch centered at a pixel to its surrounding area by the sum of squared differences. The correlation surface is transformed into a binned log-polar representation to form a local descriptor used for image and video matching. Differently from this method, we explore the structure of similarities between *all* pairs of time frames in a sequence. The main focus of our work is on the use of self-similarities for view-invariant action recognition which was not addressed in [36].

Our approach has a close relation to the notion of video self-similarity used by [37], [38]. In the domain of periodic motion detection, Cutler and Davis [38] track moving objects and extract silhouettes (or their bounding boxes). This is followed by building a 2D matrix for the given video sequence, where each entry of the matrix contains the absolute correlation score between the two frames i and j . Their observation is that for a periodic motion, this similarity matrix will also be periodic. To detect and characterize the periodic motion, they resort to Time-Frequency analysis. Following this, Benabdelkader et al. [37] use the same construct of the self-similarity matrix for gait recognition in videos of walking people. The periodicity of the gait creates diagonals in the matrix and the temporal symmetry of the gait cycles are represented by the cross-diagonals. In order to compare sequences of different length, the self-similarity matrix is subdivided into small units. Both of these works focus primarily on videos of walking people for periodic motion detection and gait analysis. The method in [39] also concerns gait recognition using temporal similarities between frames of

different image sequences. None of the methods above explores the notion of self-similarity for view-invariant action recognition.

SSM as a recurrence plot. Recurrence is a fundamental phenomenon of many dynamical systems. The study of such systems is typically based on recorded data time series, $\{\mathbf{x}_t, t = 1 \dots T\}$, from which one wants to learn as much information about observed system as possible. Traditional techniques for understanding a dynamical system involve embedding this time series into an E -dimensional reconstruction phase space using delay coordinates [16], [17], [18]. This process involves estimation of two parameters, i.e., 1) the embedding dimension E and 2) the delay, which is a difficult task.

In order to *visualize* the geometry of a dynamical system's behavior, Eckmann et al. [40] first proposed the *Recurrence Plot*, defined as

$$RP(i, j) = \Theta(\varepsilon - \|\mathbf{x}_i - \mathbf{x}_j\|_2), \quad (1)$$

where $\Theta(\cdot)$ is the Heaviside function. Once a suitable threshold ε is determined, the RP is then a binary image displaying a black dot where the values are within the *threshold corridor*. As we shall see, the proposed self-similarity matrix is a variant of the RP: Instead of capturing the system's behavior using *dots* and *lines* by thresholding, we aim for plots with richer textures, in terms of distinct peaks and valleys, which are hopefully distinctive for different dynamical systems. These patterns on the RPs and a fortiori on SSMs contain a wealth of information about the dynamics of a system and capture specific behaviors of the system. Researchers have attempted to classify dynamic systems into different categories based on these *textures*. Part of this categorization [41] is reproduced in Table 1.

McGuire et al. [42] have shown that RPs not only preserve invariants of a dynamical system (such as the Lyapunov exponents [43]) but are also to some extent independent of the embedding dimension [42], which naturally raises the question of whether embedding is necessary at all for understanding the underlying dynamics of a system [44]. In addition, McGuire et al. [42] have shown that RPs for different systems are identical as long as the transformation is *isometric*. This conclusion also apply to proposed SSMs. As we shall see, SSMs are not strictly invariant under *projective* or *affine* transformations, but are experimentally found stable under 3D view changes.

3 SELF-SIMILARITY MATRIX

Self-similarity matrices have already appeared in the past under various specific forms, including binary recurrence plots associated to time series, as mentioned above. In this section, we define such matrices for different image features, with examples for several action classes, and start investigating their stability across views.

For a sequence of images $\mathcal{I} = \{\mathcal{I}_1, \mathcal{I}_2, \dots, \mathcal{I}_T\}$ in discrete (x, y, t) -space, a SSM of \mathcal{I} is a square symmetric matrix of size $T \times T$,

$$[d_{ij}]_{i,j=1,2,\dots,T} = \begin{bmatrix} 0 & d_{12} & d_{13} & \dots & d_{1T} \\ d_{21} & 0 & d_{23} & \dots & d_{2T} \\ \vdots & \vdots & \vdots & & \vdots \\ d_{T1} & d_{T2} & d_{T3} & \dots & 0 \end{bmatrix}, \quad (2)$$

TABLE 1
Typical Patterns of Recurrence Plots and Their Meaning (Reproduced from [41])

SSM Pattern	Meaning
(1) Homogeneity	The process is stationary
(2) Fading in the corners	Non-stationary data; the process contains a trend or a drift
(3) Periodic/quasi-periodic patterns	Cyclicities in the process; the time distance between periodic patterns (e.g. lines) corresponds to the period
(4) Single isolated points (or structures)	Strong fluctuation in the process; if only single isolated points occur, the process may be an uncorrelated random or even anti-correlated process
(5) Diagonal lines (parallel to the main diagonal)	The evolution of states is similar at different epochs; the process could be deterministic; if these diagonal lines occur beside single isolated points, the process could be chaotic (if these diagonal lines are periodic, unstable periodic orbits can be observed)
(6) Diagonal lines (orthogonal to the main diagonal)	The evolution of states is similar at different times but with reverse time; sometimes this is an indication for an insufficient embedding
(7) Long bowed line structures	The evolution of states is similar at different epochs but with different velocity; the dynamics of the system could be changing

where d_{ij} is the distance between certain low level features extracted in frames \mathcal{I}_i and \mathcal{I}_j , respectively. The diagonal corresponds to comparing a frame to itself (no dissimilarity) and hence is composed of zeros. The exact structures or the patterns of this matrix depend on the features and the distance measure used for computing the entries d_{ij} . For example, after tracking walking people in a video sequence, Benabdelkader et al. [37], and Cutler and Davis [38] compute a particular instance of SSM where d_{ij} is the absolute correlation between two frames, as depicted in Fig. 3. The computed matrix patterns (cf. Fig. 3e) have a significant meaning for their application—the diagonals in the matrix indicate periodicity of the motion.

In this work, we define d_{ij} as the euclidean distance between the different features that we extract from an action sequence. This form of SSM is known in the literature as the Euclidean Distance Matrix (EDM) [45].

To get a first insight into the representation power of SSMs, a comparison with the notion of “dynamic instances” proposed by Rao et al. [26] is illustrated in Fig. 2. The authors of [26] argue that continuities and discontinuities in position, velocity, and acceleration of a 3D trajectory of an object are preserved under 2D projections. For an action of opening a cabinet door, performed by two different actors from considerably different viewpoints, these points are depicted in Fig. 2. Figs. 2c and 2e show the SSMs computed for these two actions based only on one hand trajectory, where red color indicates higher values and dark blue color indicates

lower values. The dynamic instances, red stars in Figs. 2b and 2d, correspond to valleys of different area/spread in our plot of SSM (cf. Figs. 2c and 2e), marked by magenta circles along the diagonal of the matrix. The exact spread of these valleys depends on the width of the peaks in the spatiotemporal curvature of the actions, as shown in Figs. 2b and 2d. However, whereas Rao et al. [26] capture only the local discontinuities in the spatiotemporal curvature, the SSM captures more information about other dynamics of the actions present in the off-diagonal parts of the matrix. Note also that the proposed notion of self-similarity, unlike [5] or [26], does not require estimation of point correspondences or time-alignment between different actions.

3.1 Trajectory-Based Self-Similarities

If a set of M points \mathbf{x}^m , $m = 1 \cdots M$, distributed over a person is “tracked” (in a sense to be specified later) over the duration of an action performance, the mean euclidean distance between each of the k pairs of corresponding points at any two instants i and j of the sequence can be computed as

$$d_{ij} = \frac{1}{M} \sum_{m=1}^M \|\mathbf{x}_i^m - \mathbf{x}_j^m\|_2, \quad (3)$$

where \mathbf{x}_i^k and \mathbf{x}_j^k indicate positions of points on the track k at time instants i and j . We denote the self-similarity matrix computed from (3) by SSM-pos.

In a first set of experiments aimed at investigating SSM properties in a controlled setup, such point trajectories are

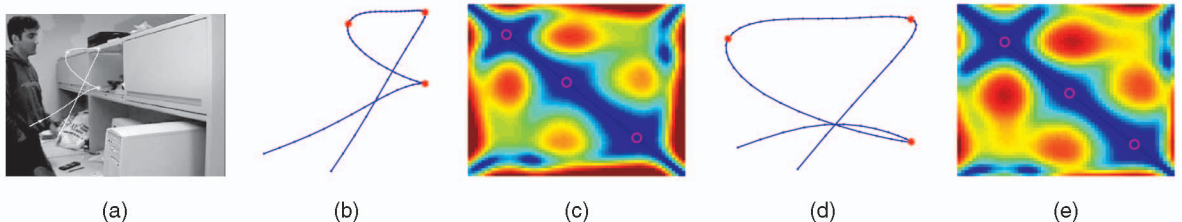


Fig. 2. Relationship between proposed SSM representation and *dynamic instances* introduced in [26]. Two actors perform the action of opening a cabinet door, where the hand trajectory is shown in (b) and (d). The SSMs computed for these two actions based only on one hand trajectory are shown in (c) and (e), respectively. The “dynamic instances” (as proposed by [26]), marked in red stars in (b) and (d), represent valleys in the corresponding SSM, depicted by magenta circles in (c) and (e), respectively. The spread of each valley depends on the peak-width of the corresponding dynamic instance.



Fig. 3. Earlier example of SSM for motion periodicity analysis. (a)-(d) are frames from a sequence of a walking person [38]. (e) represents the SSM obtained for this sequence by [38] using the absolute correlation score between frames of the sequence. Time-Frequency analysis is performed on this matrix to detect periodicity in a motion sequence.

directly obtained via motion capture, rather than from video sequences. In this case, a “view” corresponds to the projection of 3D point tracks onto a given 2D plane. In these experiments, we track $M = 13$ joints on a person performing different actions [16], as shown in the Fig. 4a. In order to remove the effect of global person translation, without loss of generality, the points are centered to their centroid so that their first moment is zero.

The overall goal of the proposed work being the recognition of actions in videos irrespective to viewpoints, the actual computation of SSM-pos requires that points be extracted and tracked in the input video. We assume that this task is handled automatically by an external module, such as KLT [46] point tracker. Note that our method is not restricted to any particular subset of points as far as the points are distributed over moving body parts. The definition of SSM-pos in (3) needs, however, to be adapted to a set of tracks with arbitrary length and starting time

$$d_{ij} = \frac{1}{|S_{ij}|} \sum_{m \in S_{ij}} \|\mathbf{x}_i^m - \mathbf{x}_j^m\|_2, \quad (4)$$

where $S_{ij} \subset \{1, \dots, M\}$ is the set with indices of point trajectories that are alive between frames i and j .

In addition to the SSM-pos, we also compute similarities based on the first and the second-order derivatives of the 2D positions, i.e., velocities and accelerations. Similarities computed based on these features are denoted by SSM-vel and SSM-acc, respectively.

3.2 Image-Based Self-Similarities

Besides point trajectories, alternative image features can be used to construct other SSMs for the same image sequence. To describe the spatial appearance of a person at each image frame, we compute Histograms of Oriented Gradients (HoG) features [47]. This descriptor, originally used to perform human detection, characterizes the local shape by capturing the gradient structure. In our implementation, we use four bin histograms for each of 5×7 blocks defined on a bounding box around the person in each frame. Feature distance d_{ij} between time instants i and j is then computed as the euclidean distance between two HoG vectors extracted from frames \mathcal{I}_i and \mathcal{I}_j . We denote SSMs computed using HoG features by SSM-hog.

In addition to HoG features, we also test the proposed method by considering optical flow vectors as another input feature. The corresponding SSMs are denoted by SSM-of. More precisely, we assume, as for point trajectories, that optical flow is provided by another module, e.g., Lucas and Kanade algorithm [48] based on two consecutive frames. We

concatenate the components of optical flow vectors computed for all n pixels in a bounding box around a person into a flow vector of size $2n$. Entry d_{ij} of SSM-of matrix then amounts to the euclidean distance between the flow vectors corresponding to the two frames \mathcal{I}_i and \mathcal{I}_j . In practice, we enlarge and resize bounding boxes in order to avoid border effects on the flow computation and to ensure the same size of the flow vectors along an action sequence. We resize the height to a value equal to 150 pixels and the width is set to the greatest value for the considered sequence.

Examples of SSMs computed for different image features are shown in Fig. 4. Fig. 4a contains example actions from the CMU motion capture (mocap) data set projected onto different views. Columns 1 and 5 of Fig. 4a represent two different actors, while columns 2 and 4 represent their computed SSM-pos, respectively. The first two rows represent a bending action performed by two actors and projected onto two considerably different views. The last two rows similarly represent a football kick action for two actors and two different views. Note the similarity of SSMs computed for actions of the same class despite the changes in the actor and the considerable changes of views. Note also the visual difference of SSMs between two action classes. Computing SSMs on real-image features instead of mocap data leads to similar conclusions. Fig. 4b illustrates SSMs obtained for the bending action from the video data set [7]. Row 2 shows SSM-pos computed using point tracks overlaid on images in first row. Rows 3 and 4 show SSM-hog and SSM-of for the same sequences, respectively. For a given type of features, note the similarity of SSMs over the different instances of the same action class. SSMs for different feature types do not look similar since different features capture different properties of the action. This suggests the use of SSMs computed for different features in a complementary manner.

3.3 Structural Stability of SSM across Views

As noted above, the patterns of proposed SSMs are promisingly stable through changes of viewpoints. In order to more thoroughly assess this stability, we conducted the following experiments using the CMU mocap data set. We deployed a total of $K = 684$ synthetic affine cameras (at distinct latitudes and longitudes) on a sphere surrounding the person performing an action, as shown in Fig. 5a. For each of these cameras, we compute the SSM matrix, as described in Section 3.1, and aim to assess qualitatively and quantitatively the stability of the patterns contained in these SSMs. To this end, we consider SSMs as being discrete “images” of size $T \times T$, which allows us to resort to classic tools for image structure analysis. We consider, in particular, orientation of bidimensional gradient as it is known to capture image structures independently of various changes, including changes in the dynamics of intensity levels. We will further rely on this philosophy when building SSM descriptors in the next section. For the time being, we consider a simpler structure analysis based on so-called circular statistics [49].

At each “pixel” (i, j) of the SSM associated to the k th view, we compute the orientation $\theta_{ij}^{(k)}$ of the bidimensional gradient vector. In order to ascertain the effect of different viewing directions on the computed SSMs, we then compute, at this point, the circular mean and standard

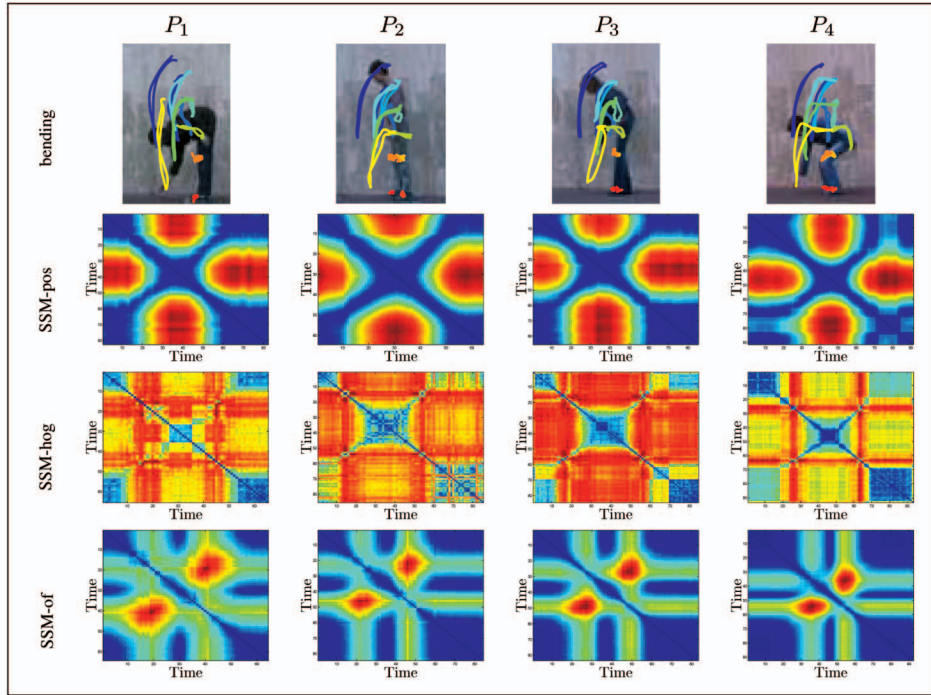
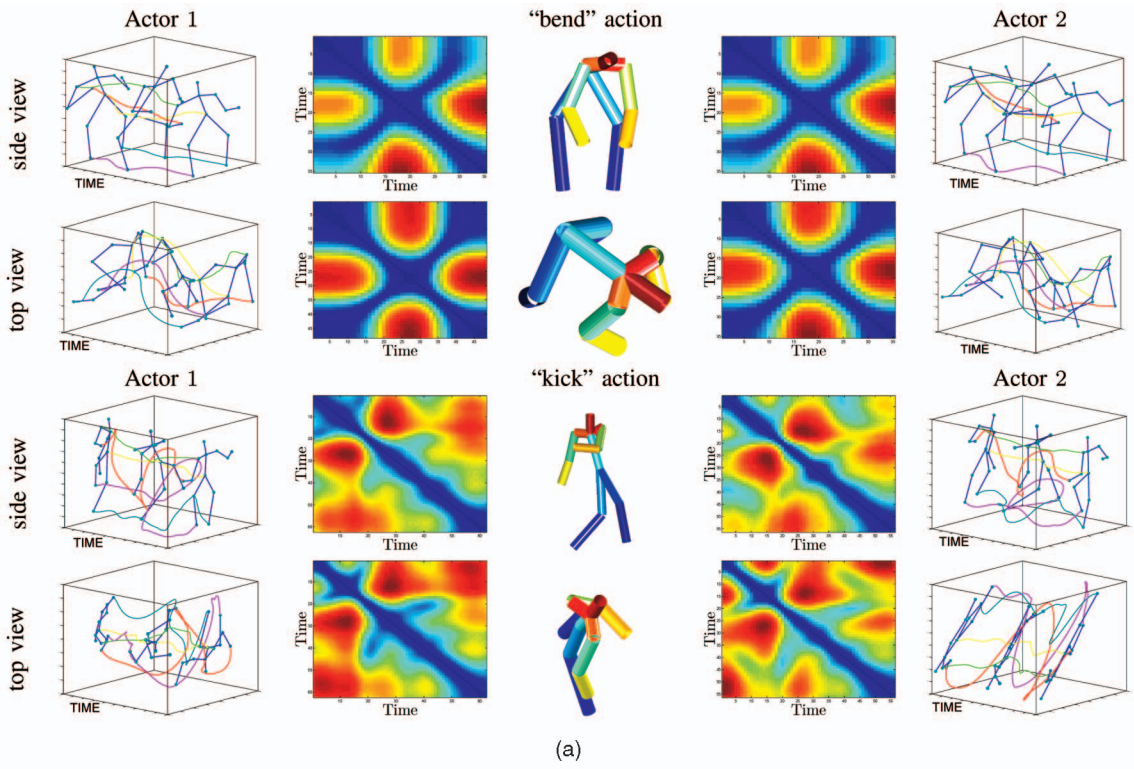


Fig. 4. Examples of SSMs for different types of features and for different actions. (a) Examples from the CMU mocap data set. Columns 1 and 5 represent two actors, while columns 2 and 4 represent corresponding SSM-pos computed with 13 projected point trajectories, respectively. Different rows represent different actions and viewing angles. Note the stability of SSMs over different views and people performing the same action. (b) Examples from the Weizman video data set [7]. Row 1: four bending actions along with manually extracted point trajectories used for computing SSM-pos; rows 2, 3, and 4 represent SSM-pos, SSM-hog, and SSM-of, respectively, for these four bending actions. Note the similarity columnwise.

deviation, $\bar{\theta}_{ij}$ and $\bar{\sigma}_{ij}$ of the orientation over the $K = 684$ SSMs. Let $\bar{\mathbf{r}}_{ij} = [\bar{c}_{ij} \ \bar{s}_{ij}]'$, where

$$\bar{c}_{ij} = \sum_{k=1}^K \cos \theta_{ij}^{(k)} / K \quad \bar{s}_{ij} = \sum_{k=1}^K \sin \theta_{ij}^{(k)} / K. \quad (5)$$

$$\bar{\theta}_{ij} = \begin{cases} \arctan(\bar{s}_{ij}/\bar{c}_{ij}), & \text{if } \bar{c}_{ij} \geq 0, \\ \arctan(\bar{s}_{ij}/\bar{c}_{ij}) + \pi \text{sign}(\bar{s}_{ij}), & \text{if } \bar{c}_{ij} < 0. \end{cases}$$

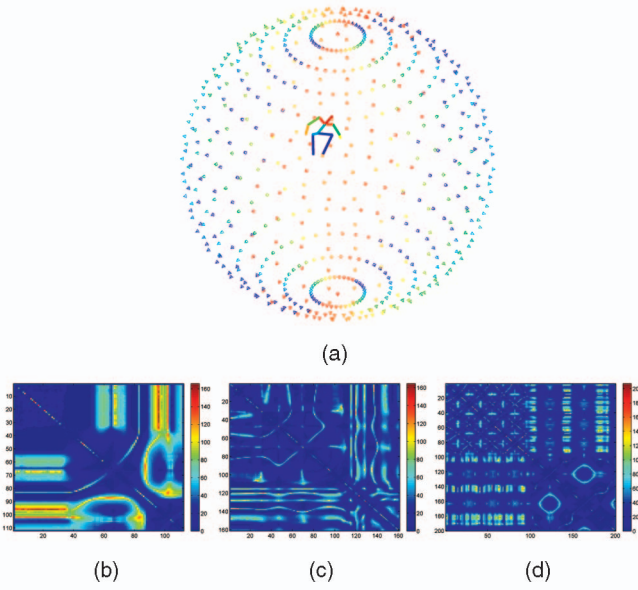


Fig. 5. Stability of SSM-pos structures across viewpoints for the mocap data sequence. (a) Synthetic cameras around a person performing an action. SSMs are generated for each of these synthetic cameras and, for each of these computed SSMs, a gradient angle is computed at each matrix point. From these orientations, circular standard deviations are computed for (b) a *golf swing*, (c) a *kick*, and (d) a *jumping jack* action sequence (code provided by [27] for (a)).

The mean resultant length, $\bar{r}_{ij} = \sqrt{\bar{c}_{ij}^2 + \bar{s}_{ij}^2}$, is used to compute the *circular standard deviation* as $\bar{\sigma}_{ij} = \sqrt{-2 \ln \bar{r}_{ij}}$. We computed $[\bar{\sigma}_{ij}]_{i,j=1 \dots T}$ for some sample action sequences from the *golf swing*, *kick*, and *jumping jack* action classes, as shown in Figs. 5b, 5c, and 5d, respectively. One can notice that, for each action, the standard deviations are low over most parts of the SSM support, which is a good indicator of SSM structure stability across views. The highest values delineate what can be seen as the strong contours of the average SSM structure for the action concerned.

4 SSM-BASED ACTION DESCRIPTION AND ALIGNMENT

As discussed in the previous section, SSMs have view-stable and action-specific structure. Here, we aim to capture this structure and to construct SSM-based descriptors for subsequent view-independent action analysis, such as alignment and recognition. We note the following properties of SSM: 1) Absolute values of SSM may depend on the varying properties of the data, such as the projected size of a person in the case of SSM-pos; 2) fluctuations in the individual performances of a type of actions and temporal desynchronization of the views may effect the global structure of SSM; and 3) the uncertainty of values in SSM increases with the distance from the diagonal due to the increasing difficulty of measuring self-similarity over long time intervals. These properties led us to the SSM description that follows.

As already mentioned in the previous section, we avoid dependency on varying absolute SSM values by resorting to gradient orientations computed from neighboring elements of the matrix seen as an image. We also avoid global descriptors and, in a manner reminiscent of popular local

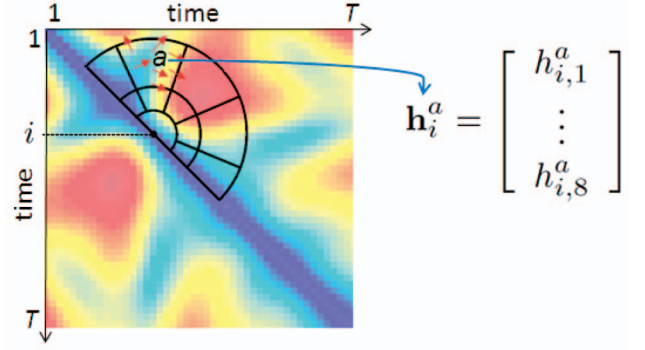


Fig. 6. Local descriptors for SSM: Each individual descriptor is centered at a diagonal point $i \in \{1 \dots T\}$ of the SSM and has a log-polar block structure. Histograms \mathbf{h}_i^a of eight gradient directions are computed separately for each of the 11 blocks of the analysis support and are concatenated into a descriptor vector \mathbf{h}_i .

image descriptors used for object detection and recognition, we accumulate histograms of gradient orientations in local patches. These patches, however, are only centered on the diagonal of SSM. Our patch descriptor has a log-polar block structure, as illustrated in Fig. 6. The diameter of the circular regions under consideration should be seen as temporal window extent. For log-polar block a at time i , we compute the normalized 8-bin histogram $\mathbf{h}_i^a = [\mathbf{h}_{i,b}^a]_{b=1:8}'$ of SSM gradient orientations within the block. We then concatenate the histograms of the 11 blocks of the analysis support into a descriptor vector $\mathbf{h}_i = [\mathbf{h}_{i,a}^a]_{a=1:11}'$. For descriptors at boundaries with blocks falling outside SSM, we set \mathbf{h}_i^a to a zero vector.

Choosing a temporal extent of the descriptor involves a trade-off between the amount of captured temporal information and its variability, which is delicate to tune. In addition, using a single descriptor size may be suboptimal when representing events of varying lengths and with irregularly changing speed. We address this issue by considering local SSM descriptors of multiple sizes and demonstrate the impact of this approach on action recognition in Section 5.3.

When constructing a joint local descriptor for multiple SSMs computed for F different features, we concatenate F corresponding local descriptors \mathbf{h}_i^f from each SSM into a single vector $\mathbf{h}_i = [\mathbf{h}_i^f]_{f=1:F}'$. In such a way, we obtain, for instance, SSM-hog-of descriptors by concatenating image-based SSM-hog and SSM-of descriptors. When temporal ordering is required, the representation for a video sequence can finally be defined by the sequence of local descriptors $H(\mathcal{I}) = (\mathbf{h}_1, \dots, \mathbf{h}_T)$ computed over all diagonal elements of SSMs associated to all feature types.

Temporal cross-view action synchronization. Before moving to action recognition based on representation previously defined, we first test this representation on the problem of temporal alignment, or synchronization, of video sequences representing the *same* action from different view-points. The problem amounts to finding the monotonic mapping between the timeline of the first sequence and the timeline of the second one. Consider, for instance, two videos, \mathcal{I}^1 and \mathcal{I}^2 , recorded simultaneously for the side and the top views of a person in action, as shown in Fig. 7a. To further challenge the alignment, we apply a nonlinear time transformation to one of the sequences. To solve the

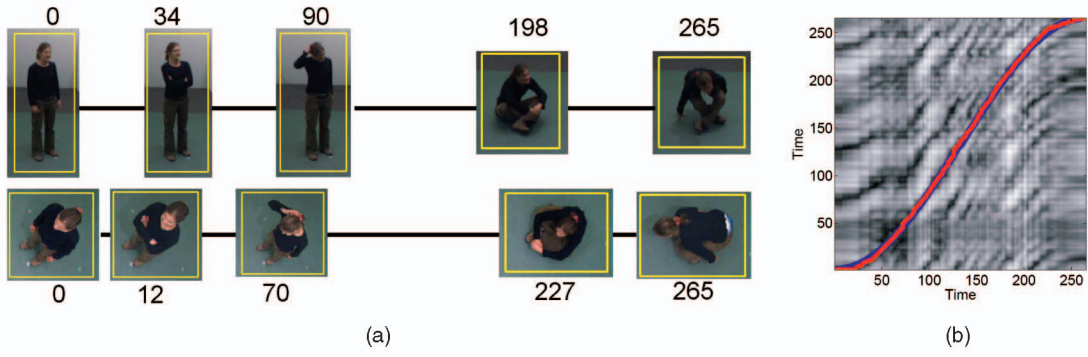


Fig. 7. Temporal alignment of same action performances in videos from different viewpoints and with synthetic desynchronization. (a) Two desynchronized sequences with the side and the top views of the same action are represented with a set of matching key-frames. The second sequence has been time warped according to $t' = a \cos(bt)$ transformation. (b) Distance matrix between sequences $H(\mathcal{T}^1)$ and $H(\mathcal{T}^2)$ of SSM-pos descriptors (bright colors represent large distance values). Dynamic Programming (red curve) finds the minimum cost monotonic path from $(0,0)$ to (T,T) in this matrix. This path coincides almost perfectly with the original warping (blue curve), despite drastic view variations.

alignment, we 1) compute optical flow-based SSM-of for both image sequences, 2) represent both videos by a sequence of local SSM descriptors, $H(\mathcal{T}^1)$ and $H(\mathcal{T}^2)$, respectively, computed for a single temporal scale as described above, and 3) align the two descriptor sequences using Dynamic Programming (DP). The estimated time transformation is

illustrated by the red curve in Fig. 7b, which closely follows the ground truth transformation (blue curve), despite the drastic change of viewpoint between sequences.

Using the same method, we next address alignment of *different* instances of similar actions. Fig. 8 demonstrates alignment of pairs of videos representing actions *throwing a ball*, *drinking*, and *smoking*

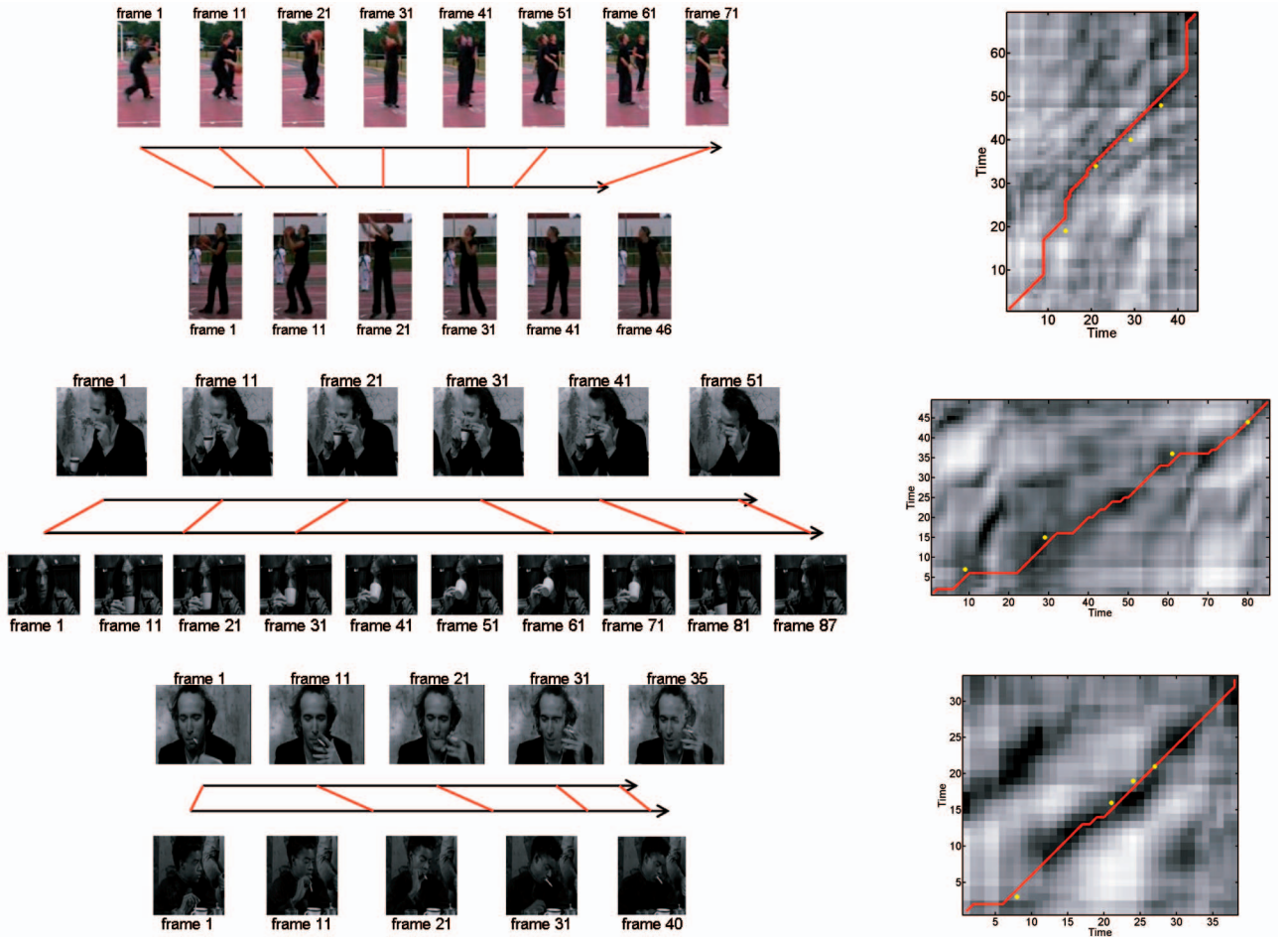


Fig. 8. Temporal alignment of video sequences representing different performances of actions *throwing a ball*, *drinking*, and *smoking*. Left: Pairs of aligned video sequences are illustrated with a few frames and the links between corresponding frames estimated by our algorithm. Right: Distance matrices between sequential descriptors of both videos used as input for aligning video sequences by Dynamic Programming. The estimated temporal alignment is illustrated by red curves. The successful alignment achieved by our method on these sequences is confirmed when comparing red curves with yellow dots illustrating sparse manual alignment for a few key frames of videos (best viewed in color).

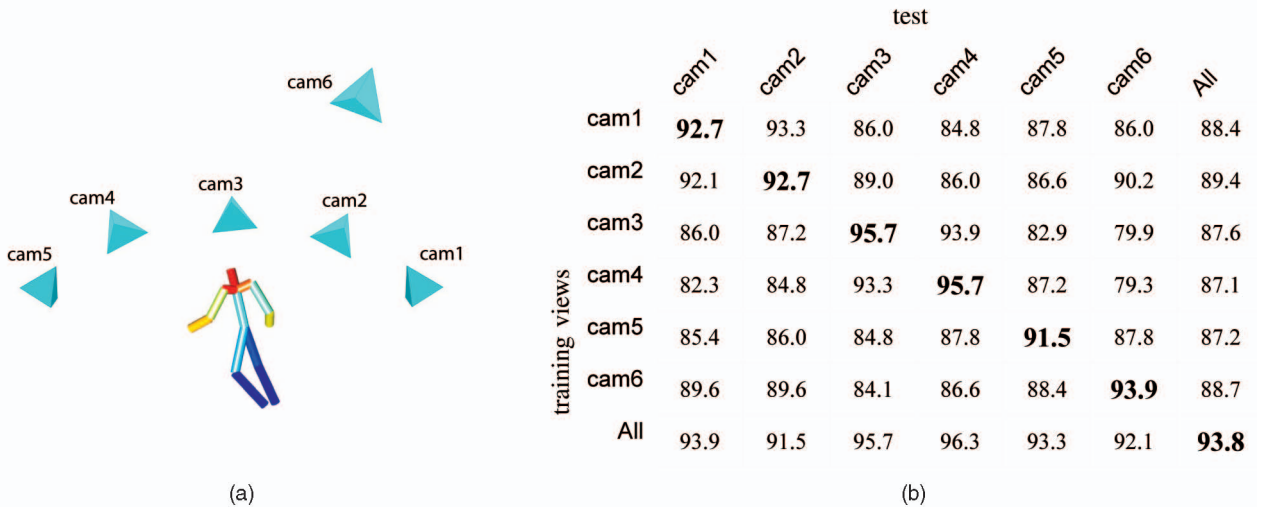


Fig. 9. SSM-based cross-view action recognition on the CMU mocap data. (a) A person figure animated from the motion capture data and six virtual cameras used to simulate projections in our experiments. (b) Accuracy of the cross-view action recognition using SSM-pos-vel-acc descriptors to build the bag-of-features used by nearest-neighbor classifier.

ball, drinking, and smoking performed by different people in varying views. The automatically estimated alignment recovers the manual alignment at key-frames as illustrated in Fig. 8b despite large variations in appearance and viewpoints across videos. Note that in all alignment experiments, we have used known person bounding boxes for computing SSM-of descriptors.

The successful alignment of actions illustrated above indicates the strength of SSM-based descriptors and their ability to cope with video variations in terms of viewpoints, subject appearance, and movement speed. This suggests that SSM-based descriptors can be used for action recognition as will be investigated in the next section.

5 SSM-BASED ACTION RECOGNITION

In this section, we evaluate SSM-based video descriptors for the task of view-invariant action recognition. To recognize action sequences, we follow recently successful bag-of-features (BoF) approaches [12], [50], [51] and represent each video as a set of quantized local SSM descriptors with their temporal positioning in the sequence being discarded. Taking this view that global temporal ordering is not taken into action (as opposed to its use for synchronization where it is a crucial information) permits us to filter out fluctuations between actions from the same class while retaining sufficient action discrimination to build a view-independent action recognition system, as demonstrated below.

As in classic BoF approach, local SSM descriptors are quantized based on a visual “vocabulary” learned offline: By k-means clustering of 10,000 random local SSM descriptors from training sequences, 1,000 clusters are defined, with their centers being the words of this vocabulary. In subsequent classifier training and testing, each feature is then assigned to the closest (we use euclidean distance) vocabulary word. This way, each image sequence \mathcal{I} is now described by a normalized histogram $\mathcal{H}(\mathcal{I})$ of visual words. These histograms are the input data used for recognition.

In the following, we consider two different types of classifiers: the Nearest Neighbor Classifier (NNC) and a

Support Vector Machine classifier. In the case of NNC, we simply assign to test sequence $\mathcal{H}(\mathcal{I})$ the action label of the training sequence \mathcal{I}^* which minimizes distance $D_{NN}(\mathcal{H}(\mathcal{I}), \mathcal{H}(\mathcal{I}^*))$ over all training sequences. The distance D_{NN} is defined by the greedy matching of local descriptors described in [51]. We apply NNC only to data sets with a limited number of samples. For SVM classification, we train nonlinear SVMs using the χ^2 kernel and adopt a one-against-all approach for multiclass classification.

We evaluate SSM-based action recognition on three public data sets. For all recognition experiments, we report results for n -fold cross validation and make sure the actions of the same person do not appear in the training and in the test sets simultaneously. In Section 5.1, we validate the approach in controlled multiview settings using motion capture data. In Section 5.2, we demonstrate and compare the discriminative power of our method on a standard single-view action data set [7]. We finally evaluate the performance of the method on a comprehensive multiview action data set [35] in Section 5.3. We demonstrate the advantage of combining SSM descriptors computed for different types of image features and multiple temporal scales. Multiview recognition results are compared with the results of other methods on the same data sets.

5.1 Experiments with the CMU Mocap Data Set

To simulate multiple and controlled view settings, we have used 3D motion capture data from the CMU data set (<http://mocap.cs.cmu.edu>). Trajectories of 13 points on the human body were projected to six cameras with predefined orientations with respect to the human body (see Fig. 9a). We have used 164 sequences in total, corresponding to 12 action classes (*bend, cartwheel, drink, fjump, flystroke, golf, jjack, jump, kick, run, walk, walkturn*). To simulate potential failures of the visual tracker, we distracted trajectories by randomly breaking them into parts with the average length of two seconds. Fig. 9b demonstrates results of NNC action recognition when training and testing on different views using SSM-pos, SSM-vel, and SSM-acc. As observed from the diagonal, the recognition accuracy is the highest when

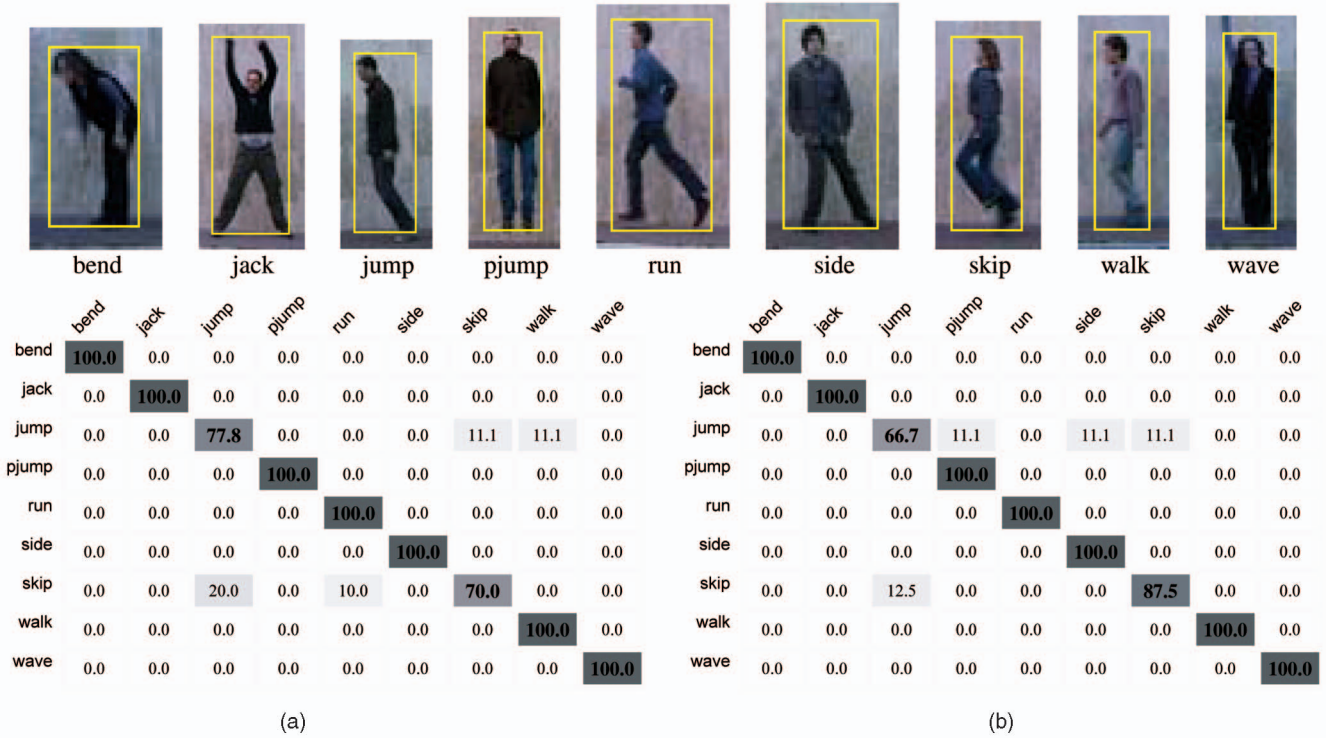


Fig. 10. SSM-based action recognition on the Weizman single-view action data set [7]. (Top) Example frames for nine classes of actions. (Bottom) Confusion matrices corresponding to NNC action recognition using (a) image-based self-similarities SSM-of and (b) trajectory-based self-similarities SSM-pos.

training and testing on the same views, while the best accuracy (95.7 percent) is achieved for cam5 (frontal view). Interestingly, the recognition accuracy degrades only moderately with substantial view changes and still remains high across the top view (camera 6) and side views (cameras 1 to 5). When training and testing on all views, the average accuracy is 90.5 percent.

5.2 Experiments with the Weizman Actions Data Set

To assess the discriminative power of our method on real video sequences, we apply it to a standard single-view video data set with nine classes of human actions performed by nine subjects [7] (see Fig. 10(top)). On this data set, we compute NNC recognition accuracy when using either image-based or trajectory-based self-similarity descriptors according to Section 3. Given the low resolution of image sequences in this data set, the trajectories were acquired by [16] via semi-automatic tracking of body joints. Recognition accuracy achieved by our method for optical flow-based and trajectory-based self-similarities is 94.6 and 95.3 percent, respectively, and the corresponding confusion matrices are illustrated in Figs. 10a and 10b. The recognition results are high for both types of self-similarity descriptors and outperforms the accuracy of 92.6 percent achieved by a recent trajectory-based method [16]. Whereas higher recognition rates on this *single-view* data set have been reported, e.g., in [52], the main strength of our method will be demonstrated for action recognition across *multiple views*, as described in the next section.

5.3 Experiments with the IXMAS Data Set

The IXMAS data set is publicly available and numerous researchers have reported their results on this data set.

Without resorting to engineering a different experimental setup to test view invariance, using this data set allows for a quick and fair comparison of our method to the other methods. Thus, we present results for the IXMAS video data set [35] with 11 classes of actions performed three times by each of 10 actors and recorded simultaneously from five different views. Sample frames for all cameras and four action classes are illustrated in Fig. 11. Here, we use SVM classifier in combination with image-based self-similarity descriptors in terms of SSM-hog, SSM-of and their combination SSM-hog-of. We also consider local SSM descriptors computed at *multiple temporal scales*. For each SSM diagonal point, three local descriptors are computed corresponding to three different diameters for the log-polar domain (respectively, 28, 42, and 56 frames in diameter). The number of descriptors assigned to a given sequence is thus multiplied accordingly. All descriptors are quantized independently of their scale using a single visual vocabulary and are used to compute a single histogram associated to the sequence.

Figs. 12a, 12b, and 12c illustrate recognition accuracy of cross-views action recognition for different combinations of training and test cameras and for different types of SSMs. The results are averaged over all classes and test subjects. Similarly to results on the CMU data set in Section 5.1, here we observe high stability of action recognition over view changes, now using visual data only. The method achieves reasonable recognition accuracy even for extreme variations in views such as for testing on top views (Test Cam4) when using side views only for training. Also, these tables indicate that using jointly HoG-based and optical flow-based SSMs yields better recognition than using either of

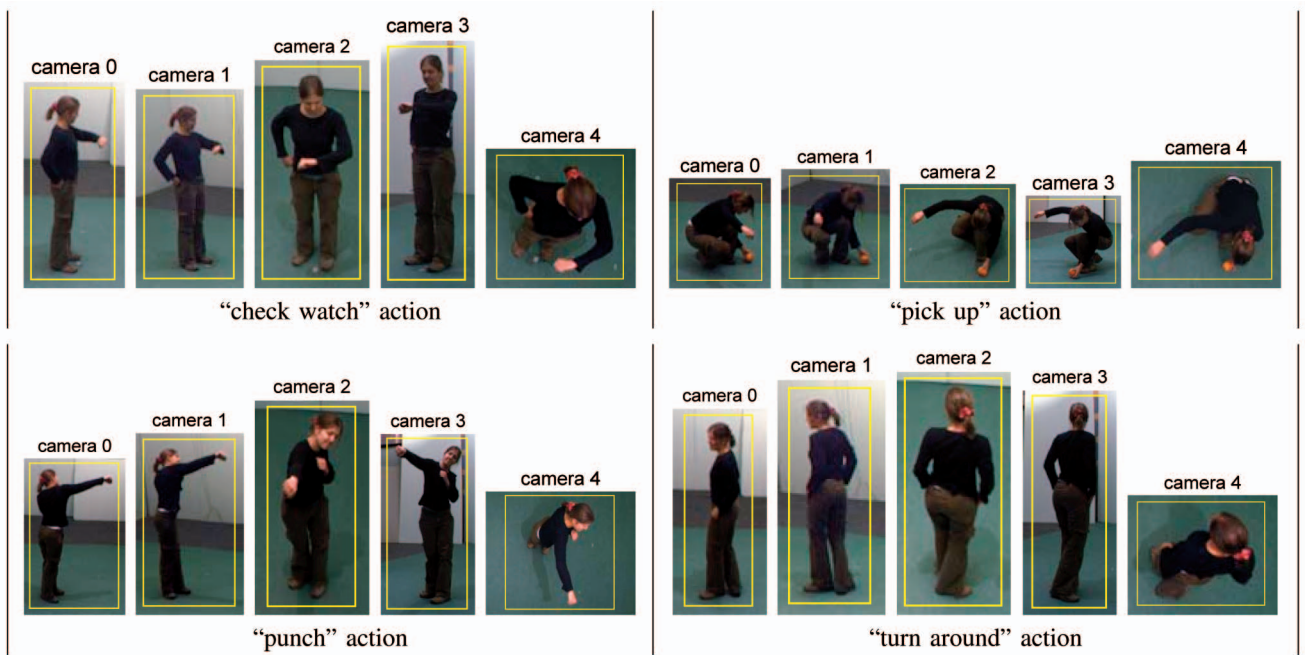


Fig. 11. Example frames from the IXMAS multiview action data set: For four classes of action, the five views at a given instant of one performance of the action are shown.

	Test Cam0	Test Cam1	Test Cam2	Test Cam3	Test Cam4	Test All
Train Cam0	77.0	75.2	69.7	71.8	49.4	68.6
Train Cam1	78.5	77.3	67.9	71.5	48.0	68.6
Train Cam2	70.0	73.0	75.8	68.5	55.2	68.5
Train Cam3	73.6	72.4	67.3	71.2	45.9	66.1
Train Cam4	44.5	41.5	55.2	37.9	68.8	49.6
Train All	77.0	78.8	80.0	73.9	63.3	74.6

■ cross-camera training/testing
 ■ same camera training/testing

(a)

	Test Cam0	Test Cam1	Test Cam2	Test Cam3	Test Cam4	Test All
Train Cam0	66.4	73.3	63.9	60.3	41.5	61.1
Train Cam1	67.3	70.6	62.1	62.4	41.5	60.8
Train Cam2	63.9	65.8	71.8	57.6	52.1	62.2
Train Cam3	62.1	68.2	62.4	61.2	35.9	58.0
Train Cam4	34.8	36.4	55.8	33.3	63.0	44.7
Train All	67.9	73.6	70.6	66.4	59.1	67.5

■ cross-camera training/testing
 ■ same camera training/testing

(b)

	Test Cam0	Test Cam1	Test Cam2	Test Cam3	Test Cam4	Test All
Train Cam0	67.6	66.1	54.8	57.3	44.4	58.0
Train Cam1	73.6	63.6	57.9	59.5	45.3	60.0
Train Cam2	58.2	54.5	63.3	54.2	49.4	55.9
Train Cam3	60.0	58.2	55.8	60.6	42.1	55.3
Train Cam4	46.7	44.2	51.5	43.9	60.0	49.3
Train All	69.7	63.3	64.8	62.7	52.7	62.7

■ cross-camera training/testing
 ■ same camera training/testing

(c)

	Test Cam0	Test Cam1	Test Cam2	Test Cam3	Test Cam4	Test All
Train Cam0	80.0	75.9	42.3	55.6	21.8	55.6
Train Cam1	74.8	83.9	36.5	58.3	23.6	56.0
Train Cam2	43.6	46.1	80.5	64.7	34.2	53.7
Train Cam3	47.0	50.0	45.8	85.5	18.8	49.5
Train Cam4	19.7	19.4	43.5	26.1	73.3	36.0
Train All	80.3	84.5	79.4	84.8	68.5	79.6

■ cross-camera training/testing
 ■ same camera training/testing

(d)

Fig. 12. Comparative action recognition results for the IXMAS multiview action data set: Results are averaged over 11 action classes and 10 subjects. Results in (a)-(c) are shown for different types of SSMs and the same bag-of-features SVM classification method. Results in (d) are obtained with the same bag-of-features SVM approach, but using quantized descriptors of spatiotemporal interest points (STIPs) instead of quantized local SSM descriptors. Recognition scores are illustrated for different combinations of training and test cameras. (a) Recognition results for SSM-hog-of Multiscale Features. (b) Recognition results for SSM-of multiscale features. (c) Recognition results for SSM-hog multiscale features. (d) Recognition results for STIP-hog-hof multiscale features.

	check-watch	cross-arms	get-up	kick	pick-up	punch	scratch-head	sit-down	turn-around	walk	wave
check-watch	73.7	8.1	0.0	1.1	1.5	1.7	11.3	0.5	0.1	0.0	1.9
cross-arms	3.8	72.6	1.0	1.8	0.4	0.2	15.7	0.6	0.3	0.0	3.6
get-up	0.5	0.6	72.8	3.6	4.1	1.4	0.2	8.8	7.7	0.4	0.0
kick	1.7	1.3	3.9	57.7	1.0	15.4	0.9	1.4	14.4	0.9	1.4
pick-up	0.9	0.1	1.7	0.4	84.5	1.9	0.7	6.5	1.1	2.0	0.3
punch	3.3	1.0	0.9	15.1	2.2	70.5	0.0	1.4	2.6	0.0	3.0
scratch-head	13.5	11.9	1.2	0.6	0.4	0.8	61.1	0.3	0.1	0.7	9.3
sit-down	0.6	0.1	9.6	1.1	2.3	1.2	0.1	81.1	3.5	0.2	0.0
turn-around	0.0	0.1	4.0	5.2	0.8	0.8	0.0	1.8	73.2	14.0	0.2
walk	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	100.0	0.0
wave	3.3	1.6	0.6	1.4	0.1	1.6	8.2	0.0	1.1	0.0	82.1

Fig. 13. Class-confusion matrix for action recognition in the IXMAS data set: This confusion matrix is obtained using SSM-hog-of multiscale SSM local descriptors. It corresponds to the average confusion computed for all *cross-camera* recognition setups in Fig. 12a.

the two types of feature individually. The class confusion matrix in Fig. 13, computed using SSM-hog-of, illustrates good per-class recognition performance for all classes when averaged over all cross-camera setups in Fig. 12a, i.e., using camera- X for training and camera- Y for testing for $X \neq Y$.

Comparison to alternative methods. We compare recognition performance of SSM-based features to space-time interest points (STIPs) [9], [15] representing videos by sets of descriptors computed from local space-time patches. STIP descriptors have recently been demonstrated to achieve competitive performance on several action recognition benchmarks [53]. STIP features, however, are not designed to handle large view variations. The recognition performance of STIP features on the IXMAS data set using the same classification method as for SSM-based features is illustrated in Fig. 12d. It is interesting to observe that STIP features outperform SSM-based features in recognition setups where the same or similar views are used for training and testing. For a large variation between training and test views, however, SSM-based descriptors considerably outperform STIP features, especially when testing on top views after learning on side views, and vice versa. This behavior is consistent with the intuition that SSM-based descriptors gain view independence at the cost of somewhat reduced discriminative power. The comparison of SSM-based and STIP features is summarized in Table 2 for different recognition setups.

TABLE 2
Comparison of Recognition Results on the IXMAS Data Set by Alternative Methods

	cross camera	same camera	any-to-any
SSM-hog-of	61.8	74.0	64.3
SSM-of	55.0	66.6	57.4
SSM-hog	53.9	63.0	55.7
STIP-hog-hof	42.4	80.6	50.0
Farhadi [33]	58.1	68.8	60.3
Weinland [35]	—	57.9	—

Results are presented for different combinations of training camera- X and test camera- Y setups where “cross camera” indicates setups with $X \neq Y$, “same camera” indicates setups with $X = Y$, and “any-to-any” indicates all combinations of X and Y .

TABLE 3
Comparison of the Impact of SSM Descriptor Size on the Recognition Performance in IXMAS Data Set

	cross camera	same camera	any-to-any
multi-scale	61.8	74.0	64.3
56 frames	59.9	70.9	61.8
42 frames	59.3	69.4	61.6
28 frames	54.0	65.3	56.2

We also compare our approach to the two alternative methods in the literature that were evaluated on the same data set. Action recognition in IXMAS data set is addressed by means of 3D reconstruction in [35]. Results of this method reported for the same training/test camera setup are lower compared to our SSM-based recognition scheme, as illustrated in Table 2. Our SSM-based descriptors also outperform results of the transfer-learning approach reported in [33] both for cross-camera and same-camera setups (cf. Table 2). Apart from the superior recognition performance, our method does not require any knowledge about actions in test views which is not the case for [33], [35].

Impact of multiple temporal scales. Table 3 presents recognition results for SSM-hog-of descriptors computed at multiple and single temporal scales. Comparing single-scale descriptors, we observe that accuracy increases with the temporal extent of the descriptor tested for descriptor sizes 28 frames (1.1 sec.), 42 frames (1.9 sec.) and 56 frames (2.2 sec.). Combining different scales, however, results in the considerable increase of performance compared to single-scale results for all camera setups.

6 CONCLUSION

We propose a self-similarity-based descriptor for view-independent video analysis, with human action recognition as a central application. Self-similarity being possibly defined over a variety of image features, either static (histograms of intensity gradient directions) or dynamic (optical flows or point trajectories), these descriptors can take different form and can be combined for increased descriptive power. Experimental validation on action recognition, as well as for the different problem of action synchronization, clearly confirms the stability of this type of description with respect to view variations. Results on public multiview action recognition data sets demonstrate superior performance of our method compared to alternative methods in the literature.

Such encouraging results are simply obtained by exploiting the stability across views of SSM patterns, with no need to rely on the delicate recovery of 3D structures nor on the estimation of correspondences across views. Our method only makes mild assumptions about the rough localization of a person in the frame. This lack of strong assumptions is likely to make this approach applicable to action recognition beyond controlled data sets when combined with modern techniques for person detection and tracking.

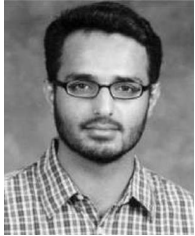
ACKNOWLEDGMENTS

This work was partially funded by the QUAERO project and the MSR/INRIA joint laboratory.

REFERENCES

- [1] T. Moeslund, A. Hilton, and V. Krüger, "A Survey of Advances in Vision-Based Human Motion Capture and Analysis," *Computer Vision and Image Understanding*, vol. 103, nos. 2-3, pp. 90-126, Nov. 2006.
- [2] L. Wang, W. Hu, and T. Tan, "Recent Developments in Human Motion Analysis," *Pattern Recognition*, vol. 36, no. 3, pp. 585-601, Mar. 2003.
- [3] A. Bobick and J. Davis, "The Recognition of Human Movement Using Temporal Templates," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 23, no. 3, pp. 257-267, Mar. 2001.
- [4] D. Weinland, R. Ronfard, and E. Boyer, "Free Viewpoint Action Recognition Using Motion History Volumes," *Computer Vision and Image Understanding*, vol. 103, nos. 2-3, pp. 249-257, Nov. 2006.
- [5] T. Syeda-Mahmood, M. Vasilescu, and S. Sethi, "Recognizing Action Events from Multiple Viewpoints," *Proc. IEEE Workshop Detection and Recognition of Events in Video*, pp. 64-72, 2001.
- [6] A. Yilmaz and M. Shah, "Actions Sketch: A Novel Action Representation," *Proc. IEEE CS Conf. Computer Vision and Pattern Recognition*, vol. 1, pp. 984-989, 2005.
- [7] L. Gorelick, M. Blank, E. Shechtman, M. Irani, and R. Basri, "Actions as Space-Time Shapes," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 29, no. 12, pp. 2247-2253, Dec. 2007.
- [8] M. Grundmann, F. Meier, and I. Essa, "3D Shape Context and Distance Transform for Action Recognition," *Proc. Int'l Conf. Pattern Recognition*, pp. 1-4, 2008.
- [9] I. Laptev, "On Space-Time Interest Points," *Int'l J. Computer Vision*, vol. 64, nos. 2/3, pp. 107-123, 2005.
- [10] E. Shechtman and M. Irani, "Space-Time Behavior Based Correlation," *Proc. IEEE CS Conf. Computer Vision and Pattern Recognition*, vol. 1, pp. 405-412, 2005.
- [11] P. Dollár, V. Rabaud, G. Cottrell, and S. Belongie, "Behavior Recognition via Sparse Spatio-Temporal Features," *Proc. Second Joint IEEE Int'l Workshop Visual Surveillance and Performance Evaluation of Tracking and Surveillance*, pp. 65-72, 2005.
- [12] J. Niebles, H. Wang, and F. Li, "Unsupervised Learning of Human Action Categories Using Spatial-Temporal Words," *Proc. British Machine Vision Conf.*, 2006.
- [13] A. Gilbert, J. Illingworth, and R. Bowden, "Scale Invariant Action Recognition Using Compound Features Mined from Dense Spatio-Temporal Corners," *Proc. European Conf. Computer Vision*, part 1, pp. 222-233, 2008.
- [14] C. Schödl, I. Laptev, and B. Caputo, "Recognizing Human Actions: A Local SVM Approach," *Proc. Int'l Conf. Pattern Recognition*, vol. 3, pp. 32-36, 2004.
- [15] I. Laptev, M. Marszałek, C. Schmid, and B. Rozenfeld, "Learning Realistic Human Actions from Movies," *Proc. IEEE CS Conf. Computer Vision and Pattern Recognition*, 2008.
- [16] S. Ali, A. Basharat, and M. Shah, "Chaotic Invariants for Human Action Recognition," *Proc. IEEE Int'l Conf. Computer Vision*, 2007.
- [17] D. Weinland and E. Boyer, "Action Recognition Using Exemplar-Based Embedding," *Proc. IEEE CS Conf. Computer Vision and Pattern Recognition*, 2008.
- [18] K. Jia and D.-Y. Yeung, "Human Action Recognition Using Local Spatio-Temporal Discriminant Embedding," *Proc. IEEE CS Conf. Computer Vision and Pattern Recognition*, 2008.
- [19] C. Sminchisescu, A. Kanaujia, Z. Li, and D. Metaxas, "Conditional Models for Contextual Human Motion Recognition," *Proc. IEEE CS Conf. Computer Vision and Pattern Recognition*, 2005.
- [20] L.-P. Morency, A. Quattoni, and T. Darrell, "Latent-Dynamic Discriminative Models for Continuous Gesture Recognition," *Proc. IEEE CS Conf. Computer Vision and Pattern Recognition*, 2007.
- [21] S.B. Wang, A. Quattoni, L.-P. Morency, D. Demirdjian, and T. Darrell, "Hidden Conditional Random Fields for Gesture Recognition," *Proc. IEEE CS Conf. Computer Vision and Pattern Recognition*, 2006.
- [22] L. Wang and D. Suter, "Recognizing Human Activities from Silhouettes: Motion Subspace and Factorial Discriminative Graphical Model," *Proc. IEEE CS Conf. Computer Vision and Pattern Recognition*, 2007.
- [23] P. Natarajan and R. Nevatia, "View and Scale Invariant Action Recognition Using Multiview Shape-Flow Models," *Proc. IEEE CS Conf. Computer Vision and Pattern Recognition*, pp. 1-8, 2008.
- [24] A. Yilmaz and M. Shah, "Recognizing Human Actions in Videos Acquired by Uncalibrated Moving Cameras," *Proc. IEEE Int'l Conf. Computer Vision*, vol. 1, pp. 150-157, 2005.
- [25] S. Carlsson, "Recognizing Walking People," *Int'l J. Robotics Research*, vol. 22, no. 6, pp. 359-370, 2003.
- [26] C. Rao, A. Yilmaz, and M. Shah, "View-Invariant Representation and Recognition of Actions," *Int'l J. Computer Vision*, vol. 50, no. 2, pp. 203-226, Nov. 2002.
- [27] Y. Shen and H. Foroosh, "View Invariant Action Recognition Using Fundamental Ratios," *Proc. IEEE CS Conf. Computer Vision and Pattern Recognition*, 2008.
- [28] V. Parameswaran and R. Chellappa, "View Invariance for Human Action Recognition," *Int'l J. Computer Vision*, vol. 66, no. 1, pp. 83-101, Jan. 2006.
- [29] A. Ogale, A. Karapurkar, and Y. Aloimonos, "View-Invariant Modeling and Recognition of Human Actions Using Grammars," *Proc. IEEE Workshop Dynamic Vision*, pp. 115-126, 2006.
- [30] M. Ahmad and S. Lee, "HMM-Based Human Action Recognition Using Multiview Image Sequences," *Proc. Int'l Conf. Pattern Recognition*, vol. 1, pp. 263-266, 2006.
- [31] R. Li, T. Tian, and S. Sclaroff, "Simultaneous Learning of Nonlinear Manifold and Dynamical Models for High-Dimensional Time Series," *Proc. IEEE Int'l Conf. Computer Vision*, 2007.
- [32] F. Lv and R. Nevatia, "Single View Human Action Recognition Using Key Pose Matching and Viterbi Path Searching," *Proc. IEEE CS Conf. Computer Vision and Pattern Recognition*, 2007.
- [33] A. Farhadi and M. Tabrizi, "Learning to Recognize Activities from the Wrong View Point," *Proc. European Conf. Computer Vision*, part 1, pp. 154-166, 2008.
- [34] P. Yan, S.M. Khan, and M. Shah, "Learning 4D Action Feature Models for Arbitrary View Action Recognition," *Proc. IEEE CS Conf. Computer Vision and Pattern Recognition*, 2008.
- [35] D. Weinland, E. Boyer, and R. Ronfard, "Action Recognition from Arbitrary Views Using 3D Exemplars," *Proc. IEEE Int'l Conf. Computer Vision*, 2007.
- [36] E. Shechtman and M. Irani, "Matching Local Self-Similarities Across Images and Videos," *Proc. IEEE CS Conf. Computer Vision and Pattern Recognition*, 2007.
- [37] C. Benabdelkader, R. Cutler, and L. Davis, "Gait Recognition Using Image Self-Similarity," *EURASIP J. Applied Signal Processing*, vol. 2004, no. 1, pp. 572-585, Jan. 2004.
- [38] R. Cutler and L. Davis, "Robust Real-Time Periodic Motion Detection, Analysis, and Applications," *IEEE Trans. Pattern Recognition and Machine Intelligence*, vol. 22, no. 8, pp. 781-796, Aug. 2000.
- [39] S. Carlsson, "Recognizing Walking People," *Proc. European Conf. Computer Vision*, pp. 472-486, 2000.
- [40] J. Eckmann, S. Kamphorst, and D. Ruelle, "Recurrence Plots of Dynamical Systems," *Europhysics Letters*, vol. 4, pp. 973-977, 1987.
- [41] N. Marwan, M.C. Romano, M. Thiel, and J. Kurths, "Recurrence Plots for the Analysis of Complex Systems," *Physics Reports*, vol. 438, nos. 5-6, pp. 237-329, 2007.
- [42] G. McGuire, N.B. Azar, and M. Shelhamer, "Recurrence Matrices and the Preservation of Dynamical Properties," *Physics Letters A*, vol. 237, nos. 1-2, pp. 43-47, 1997.
- [43] E. Bradley and R. Mantilla, "Recurrence Plots and Unstable Periodic Orbits," *Chaos: An Interdisciplinary J. Nonlinear Science*, vol. 12, no. 3, pp. 596-600, 2002.
- [44] J.S. Iwanski and E. Bradley, "Recurrence Plots of Experimental Data: To Embed or Not to Embed?" *Chaos: An Interdisciplinary J. Nonlinear Science*, vol. 8, no. 4, pp. 861-871, 1998.
- [45] S. Lele, "Euclidean Distance Matrix Analysis (EDMA): Estimation of Mean Form and Mean Form Difference," *Math. Geology*, vol. 25, no. 5, pp. 573-602, 1993.
- [46] C. Tomasi and J. Shi, "Good Features to Track," *Proc. IEEE CS Conf. Computer Vision and Pattern Recognition*, 1994.
- [47] N. Dalal and B. Triggs, "Histograms of Oriented Gradients for Human Detection," *Proc. IEEE CS Conf. Computer Vision and Pattern Recognition*, vol. 1, pp. 886-893, 2005.
- [48] B. Lucas and T. Kanade, "An Iterative Image Registration Technique with an Application to Stereo Vision," *Proc. Imaging Understanding Workshop*, pp. 121-130, 1981.
- [49] J.P.M. de Sa, *Applied Statistics Using SPSS, STATISTICA, MATLAB and R*. Springer, 2007.
- [50] M. Marszałek, C. Schmid, H. Harzallah, and J. van de Weijer, "Learning Object Representations for Visual Object Class Recognition," *Proc. PASCAL VOC'07 Challenge Workshop, in conjunction with IEEE Int'l Conf. Computer Vision*, 2007.

- [51] I. Laptev, B. Caputo, C. Schödl, and T. Lindeberg, "Local Velocity-Adapted Motion Events for Spatio-Temporal Recognition," *Computer Vision and Image Understanding*, vol. 108, no. 3, pp. 207-229, 2007.
- [52] N. Ikizler and P. Duygulu, "Human Action Recognition Using Distribution of Oriented Rectangular Patches," *Proc. Workshop Human Motion*, pp. 271-284, 2007.
- [53] H. Wang, M. Ullah, A. Kläser, I. Laptev, and C. Schmid, "Evaluation of Local Spatio-Temporal Features for Action Recognition," *Proc. British Machine Vision Conf.*, 2009.



modeling, video surveillance, scene understanding, and event detection. He is a member of the IEEE.



Imran N. Junejo received the PhD degree in computer science from the University of Central Florida in 2007. After one year as postdoctoral researcher at INRIA-Rennes, he joined the Department of Computer Sciences, University of Sharjah, where he is currently working as an assistant professor. His current focus of research is human action recognition from arbitrary views. Other areas of research interests include camera calibration, metrology, path modeling, video surveillance, scene understanding, and event detection. He is a member of the IEEE.

Emilie Dexter received the ME degree in signal and image processing from the Ecole Nationale Supérieure d'Electronique, Informatique et Radiocommunications of Bordeaux, France, the MSc degree in signal and image processing from the University of Bordeaux 1, France, in 2005, and the PhD degree from the University of Rennes in 2009. Her current research focuses primarily on event recognition/detection and image sequence synchronization.



Ivan Laptev received the master of science degree from the Royal Institute of Technology (KTH) in 1997 and the PhD degree in computer science from the same institute in 2004. He is a full time researcher at the French National Institute for Research in Computer Science and Control, INRIA Paris. From 1997 to 1999, he worked as a research assistant at the Technical University of Munich (TUM). He joined the VISTA research team at INRIA Rennes in

2004 and moved to the WILLOW research team at INRIA/ENS in 2009. His research areas include action, scene, and object recognition from video and still images. He has published more than 30 papers in international conferences and journals of computer vision. He serves as an associate editor for *Image and Vision Computing Journal* and he is a regular member of the program committees of major international conferences on computer vision.



Patrick Pérez received the engineering degree from the École Centrale Paris in 1990 and the PhD degree from the University of Rennes in 1993. After one year as a postdoctoral researcher in the Department of Applied Mathematics at Brown University, he joined INRIA (France) in 1994 as a full time researcher. From March 2000 to February 2004, he was with Microsoft Research (Cambridge, United Kingdom). He then returned to INRIA as a senior researcher and, in 2007, took the direction of Vista research team of the INRIA Rennes Center where the current work was conducted. In October 2009, he joined Thomson Corporate Research (France) as a senior researcher. His research focuses on models and algorithms for understanding, analyzing, and manipulating still and moving images. He is currently an associate editor for the *IEEE Transactions on Pattern Analysis and Machine Intelligence* and a member of the editorial board of the *International Journal of Computer Vision*. He is a member of the IEEE.

► For more information on this or any other computing topic, please visit our Digital Library at www.computer.org/publications/dlib.