

Human Activity Analysis: A Review

J. K. AGGARWAL, University of Texas at Austin

M. S. RYOO, Electronics and Telecommunications Research Institute, Daejeon and University of Texas at Austin

Human activity recognition is an important area of computer vision research. Its applications include surveillance systems, patient monitoring systems, and a variety of systems that involve interactions between persons and electronic devices such as human-computer interfaces. Most of these applications require an automated recognition of high-level activities, composed of multiple simple (or atomic) actions of persons. This article provides a detailed overview of various state-of-the-art research papers on human activity recognition. We discuss both the methodologies developed for simple human actions and those for high-level activities. An approach-based taxonomy is chosen that compares the advantages and limitations of each approach.

Recognition methodologies for an analysis of the simple actions of a single person are first presented in the article. Space-time volume approaches and sequential approaches that represent and recognize activities directly from input images are discussed. Next, hierarchical recognition methodologies for high-level activities are presented and compared. Statistical approaches, syntactic approaches, and description-based approaches for hierarchical recognition are discussed in the article. In addition, we further discuss the papers on the recognition of human-object interactions and group activities. Public datasets designed for the evaluation of the recognition methodologies are illustrated in our article as well, comparing the methodologies' performances. This review will provide the impetus for future research in more productive areas.

Categories and Subject Descriptors: I.2.10 [Artificial Intelligence]: Vision and Scene Understanding—Motion; I.4.8 [Image Processing and Computer Vision]: Scene Analysis; I.5.4 [Pattern Recognition]: Applications—Computer vision

General Terms: Algorithms

Additional Key Words and Phrases: Computer vision, human activity recognition, event detection, activity analysis, video recognition

ACM Reference Format:

Aggarwal, J. K. and Ryoo, M. S. 2011. Human activity analysis: A review. *ACM Comput. Surv.* 43, 3, Article 16 (April 2011), 43 pages.

DOI = 10.1145/1922649.1922653 <http://doi.acm.org/10.1145/1922649.1922653>

1. INTRODUCTION

Human activity recognition is an important area of computer vision research today. The goal of human activity recognition is to automatically analyze ongoing activities from an unknown video (i.e. a sequence of image frames). In a simple case where a video is segmented to contain only one execution of a human activity, the objective of the system is to correctly classify the video into its activity category. In more general

This work was supported in part by the Texas Higher Education Coordinating Board under award 003658-0140-2007.

Authors' addresses: J. K. Aggarwal, Computer and Vision Research Center, Department of Electrical and Computer Engineering, the University of Texas at Austin, Austin, TX 78705; M. S. Ryoo, Robot Research Department, Electronics and Telecommunications Research Institute, Daejeon 305-700, Korea; email: mryoo@etri.re.kr.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies show this notice on the first page or initial screen of a display along with the full citation. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, to redistribute to lists, or to use any component of this work in other works requires prior specific permission and/or a fee. Permissions may be requested from Publications Dept., ACM, Inc., 2 Penn Plaza, Suite 701, New York, NY 10121-0701 USA, fax +1 (212) 869-0481, or permissions@acm.org.

© 2011 ACM 0360-0300/2011/04-ART16 \$10.00

DOI 10.1145/1922649.1922653 <http://doi.acm.org/10.1145/1922649.1922653>

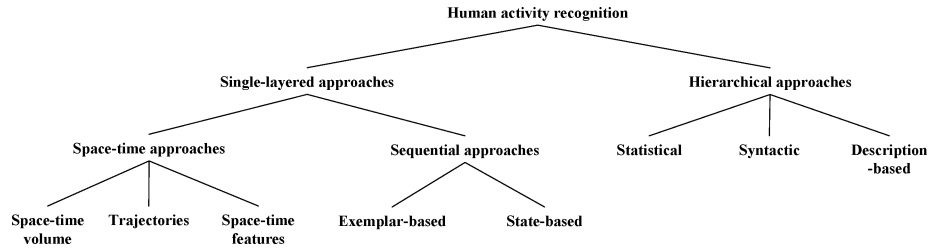


Fig. 1. The hierarchical approach-based taxonomy of this review.

cases, the continuous recognition of human activities must be performed by detecting starting and ending times of all occurring activities from an input video.

The ability to recognize complex human activities from videos enables the construction of several important applications. Automated surveillance systems in public places like airports and subway stations require detection of abnormal and suspicious activities, as opposed to normal activities. For instance, an airport surveillance system must be able to automatically recognize suspicious activities like “a person leaving a bag” or “a person placing his/her bag in a trash bin.” Recognition of human activities also enables the real-time monitoring of patients, children, and elderly persons. The construction of gesture-based human computer interfaces and vision-based intelligent environments becomes possible with an activity recognition system as well.

There are various types of human activities. Depending on their complexity, we conceptually categorize human activities into four different levels: gestures, actions, interactions, and group activities. Gestures are elementary movements of a person’s body part, and are the atomic components describing the meaningful motion of a person. “Stretching an arm” and “raising a leg” are good examples of gestures. Actions are single-person activities that may be composed of multiple gestures organized temporally, such as “walking,” “waving,” and “punching.” Interactions are human activities that involve two or more persons and/or objects. For example, “two persons fighting” is an interaction between two humans and “a person stealing a suitcase from another” is a human-object interaction involving two humans and one object. Finally, group activities are the activities performed by conceptual groups composed of multiple persons and/or objects: “A group of persons marching,” “a group having a meeting,” and “two groups fighting” are typical examples.

The objective of this article is to provide a complete overview of state-of-the-art human activity recognition methodologies. We discuss various types of approaches designed for the recognition of different levels of activities. The previous review written by Aggarwal and Cai [1999] covered several essential low-level components for the understanding of human motion, such as tracking and body posture analysis. However, the motion analysis methodologies themselves were insufficient to describe and annotate ongoing human activities with complex structures, and most of approaches in 1990s focused on the recognition of gestures and simple actions. In this new review, we concentrate on high-level activity recognition methodologies designed for the analysis of human actions, interactions, and group activities, discussing recent research trends in activity recognition.

Figure 1 illustrates an overview of the tree-structured taxonomy that our review follows. We have chosen an approach-based taxonomy. All activity recognition methodologies are first classified into two categories: single-layered approaches and hierarchical approaches. Single-layered approaches are those that represent and recognize human activities directly based on sequences of images. Due to their nature, single-layered approaches are suitable for the recognition of gestures and actions with sequential

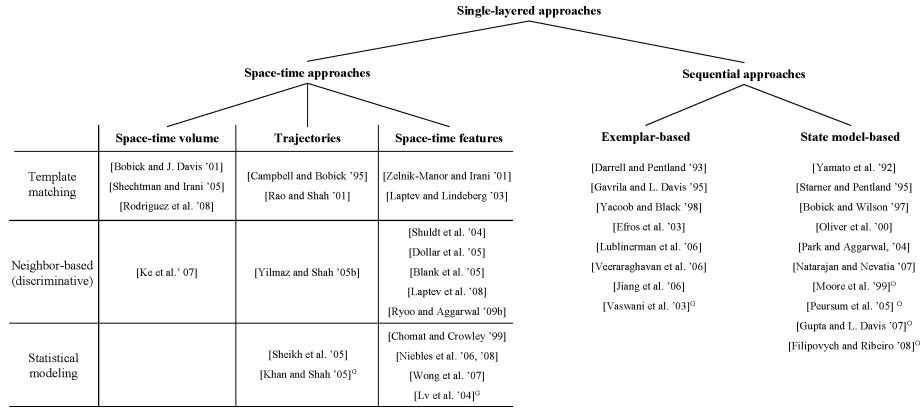


Fig. 2. Detailed taxonomy for single-layered approaches and the lists of selected publications corresponding to each category.

characteristics. On the other hand, hierarchical approaches represent high-level human activities by describing them in terms of other simpler activities, which are generally called *subevents*. Recognition systems composed of multiple layers are constructed, thus making them suitable for the analysis of complex activities.

Single-layered approaches are again classified into two types depending on how they model human activities: that is, space-time approaches and sequential approaches. Space-time approaches view an input video as a 3-D (XYT) volume, while sequential approaches interpret it as a sequence of observations. Space-time approaches are further divided into three categories based on the features they use from the 3-D space-time volumes: volumes themselves, trajectories, or local interest point descriptors. Sequential approaches are classified depending on whether they use exemplar-based recognition methodologies or model-based recognition methodologies. Figure 2 shows a detailed taxonomy used for single-layered approaches covered in the review, together with a number of publications corresponding to each category.

Hierarchical approaches are classified on the basis of the recognition methodologies they use: statistical approaches, syntactic approaches, and description-based approaches. Statistical approaches construct statistical state-based models concatenated hierarchically (e.g., layered hidden Markov models) to represent and recognize high-level human activities. Similarly, syntactic approaches use a grammar syntax such as a stochastic context-free grammar (SCFG) to model sequential activities. Essentially, they model a high-level activity as a string of atomic-level activities. Description-based approaches represent human activities by describing subevents of the activities and their temporal, spatial, and logical structures. Figure 3 presents lists of representative publications corresponding to categories.

In addition, in Figures 2 and 3, we point to previous work that recognizes human-object interactions and group activities by using different colors and by attaching “O” (object) and “G” (group) tags to the right-hand side. The recognition of human-object interactions requires the analysis of interplays between object recognition and activity analysis. This article provides a survey on the methodologies focusing on the analysis of such interplays for the improved recognition of human activities. Similarly, the recognition of groups and the analysis of their structures is necessary for group activity detection, and in this review we cover them as well.

This review is organized as follows: Section 2 covers single-layered approaches. In Section 3 we review hierarchical recognition approaches for the analysis of high-level activities. Section 4.1 discusses recognition methodologies for interactions between

Hierarchical approaches			
	Statistical approaches	Syntactic approaches	Description-based approaches
Human actions	[Nguyen et al. '05]		[Pinhanez and Bobick '98] [Gupta et al. '09]
Human-Human interactions	[Oliver et al. '02]	[Ivanov and Bobick '00] [Joo and Chellappa '06]	[Intille and Bobick '99] [Vu et al. '03] [Ghanem et al. '04] [Ryoo and Aggarwal '06, '09a]
Human-Object interactions	[Shi et al. '04] ^o [Yu and Aggarwal '06] ^o [Damen and Hogg '09] ^o	[Moore and Essa '02] ^o [Minnen et al. '03] ^o	[Siskind '01] ^o [Nevatia et al. '03, '04] ^o [Ryoo and Aggarwal '07] ^o
Group activities	[Cupillard et al. '02] ^o [Gong and Xiang '03] ^o [Zhang et al. '06] ^o [Dai et al. '08] ^o		[Ryoo and Aggarwal '08] ^o

Fig. 3. Detailed taxonomy for hierarchical approaches and the lists of publications corresponding to each category.

humans and objects, while concentrating especially on how previous work handled interplays between object recognition and motion analysis. Section 4.2 presents work on group activity recognition. In Section 5.1 we review available public datasets and compare the systems tested on them. In addition, Section 5.2 covers real-time systems for human activity recognition. Section 6 concludes the article.

1.1. Comparisons with Previous Reviews

There have been other related surveys on human activity recognition. Several previous reviews on human motion analysis [Cedras and Shah 1995; Gavrilu 1999; Aggarwal and Cai 1999] discussed human action recognition approaches as a part of their review. Kruger et al. [2007] reviewed human action recognition approaches while classifying them on the basis of the complexity of features involved in the action recognition process. Their review focused especially on the planning aspect of human action recognitions, considering their potential application to robotics. The Turaga et al. [2008] survey covered human activity recognition approaches, similar to ours. In their paper, approaches are first categorized based on the complexity of the activities that they want to recognize, and are then classified in terms of the recognition methodologies they use.

However, most of the previous reviews have focused on the introduction and summarization of activity recognition methodologies, but do not provide a means to lack of compare different types of human activity recognition approaches. In this review, we present interclass and intraclass comparisons between approaches, while providing an overview of human activity recognition approaches which are categorized on the approach-based taxonomy presented above. To be able to compare the abilities of recognition methodologies is essential for us to take advantage of them. Our goal is to enable a reader (even one from a different field) to understand the context of the development of human activity recognition and comprehend the advantages and disadvantages of the different approach categories.

We use a more elaborate taxonomy and compare and contrast each approach category in detail. For example, differences between single-layered approaches and hierarchical approaches are discussed the highest-level of our review, while space-time approaches are compared with sequential approaches in an intermediate level. We compare the abilities of previous systems within each class as well, pointing out what they are able to recognize and what they are not. Furthermore, our review covers

recognition methodologies for complex human activities, including human-object interactions and group activities, which previous reviews have not focused on. Finally, we discuss the public datasets used by the systems, and compare the performance of the recognition methodologies on the datasets.

2. SINGLE-LAYERED APPROACHES

Single-layered approaches recognize human activities directly from video data. Such approaches consider an activity as a particular class of image sequences, and recognize the activity from an unknown image sequence (i.e., an input) by categorizing it into its class. Various representation methodologies and matching algorithms have been developed to enable the recognition system to make an accurate decision as to whether an image sequence belongs to a certain activity class or not. For recognition from continuous videos, most single-layered approaches have adopted a sliding windows technique that classifies all possible subsequences. Single-layered approaches are most effective when a particular sequential pattern that describes an activity can be captured from training sequences. Due to their nature, the main objective of the single-layered approaches has been to analyze relatively simple (and short) sequential movements of humans, such as walking, jumping, and waving.

In this review, we categorize single-layered approaches into two classes: space-time approaches and sequential approaches. Space-time approaches model a human activity as a particular 3-D volume in a space-time dimension or a set of features extracted from the volume. The video volumes are constructed by concatenating image frames along a time axis, and are compared in order to measure their similarities. On the other hand, sequential approaches treat a human activity as a sequence of particular observations. More specifically, they represent a human activity as a sequence of feature vectors extracted from images and they recognize activities by searching for such a sequence. We discuss space-time approaches in Section 2.1 and compare sequential approaches in Section 2.2.

2.1. Space-Time Approaches

An image is 2-dimensional data formulated by projecting a 3-D real-world scene, and it contains spatial configurations (e.g., shapes and appearances) of humans and objects. A video is a sequence of those 2-D images placed in chronological order. Therefore, a video input containing an execution of an activity can be represented as a particular 3-D XYT space-time volume constructed by concatenating 2-D (XY) images along time (T).

Space-time approaches are those that recognize human activities by analyzing the space-time volumes of activity videos. A typical space-time approach for human activity recognition is as follows. Based on the training videos, the system constructs a model 3-D XYT space-time volume representing each activity. When an unlabeled video is provided, the system constructs a 3-D space-time volume corresponding to the new video. The new 3-D volume is compared with each activity model (i.e., template volume) to measure the similarity in shape and appearance between the two volumes. The system finally deduces that the new video corresponds to the activity that has the highest similarity. This example can be viewed as a typical space-time methodology using the 3-D space-time volume representation and the template-matching algorithm for recognition. Figure 4 shows example 3-D XYT volumes corresponding to the human action of punching.

In addition to the pure 3-D *volume* representation, there are several variations of the space-time representation. First, the system may represent an activity as a *trajectory* (instead of a volume) in a space-time dimension or other dimensions. If the system is able to track feature points such as estimated joint positions of a human,

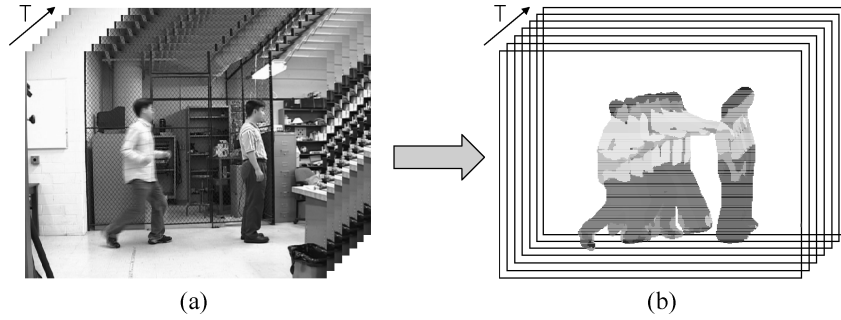


Fig. 4. Example XYT volumes constructed by concatenating (a) entire images and (b) foreground blob images obtained from a punching sequence.

the movements of the person performing an activity can be represented more explicitly as a set of trajectories. Second, instead of representing an activity with a volume or a trajectory, the system may represent an action as a set of *features* extracted from the volume or the trajectory. 3-D volumes can be viewed as rigid objects, and extracting common patterns from them enables their representation.

Researchers have also focused on developing various recognition algorithms using space-time representations to correctly match volumes, trajectories, or their features. We have already seen a typical example of an approach that uses *template-matching*, which constructs a representative model (i.e., a volume) per action using training data. Activity recognition is done by matching the model with the volume constructed from inputs. *Neighbor-based matching* algorithms (i.e., discriminative methods) have also been applied widely. In the case of neighbor-based matching, the system maintains a set of sample volumes (or trajectories) to describe an activity. The recognition is performed by matching the input with all (or a portion) of them. Finally, *statistical modeling* algorithms have been developed that match videos by explicitly modeling a probability distribution of an activity.

Accordingly, we have classified space-time approaches into several categories. A representation-based taxonomy and a recognition-based taxonomy have been jointly applied for the classification. That is, each of the activity recognition publications with space-time approaches are assigned to a slot corresponding to a specific (representation, recognition) pair. The left part of Figure 2 shows a detailed hierarchy tree of space-time approaches.

2.1.1. Action Recognition with Space-Time Volumes. The core of the recognition using space-time volumes is in the similarity measurement between two volumes. The system must be able to compute how similar human movements described in the two volumes are. In order to calculate the correct similarities, various types of space-time volume representations and recognition methodologies have been developed. Instead of concatenating entire images along time, some approaches only stack the foreground regions of a person (i.e., silhouettes) to track shape changes explicitly [Bobick and Davis 2001]. An approach to compare volumes in terms of their patches has been proposed as well [Shechtman and Irani 2005]. Ke et al. [2007] used over-segmented volumes, automatically calculating a set of 3-D XYT volume segments that corresponds to a moving human. Rodriguez et al. [2008] generated filters capturing characteristics of volumes, in order to match volumes more reliably and efficiently. In this section, we cover each of these approaches while focusing on our taxonomy of “what types of space-time volume they use” and “how they match volumes to recognize activities.”

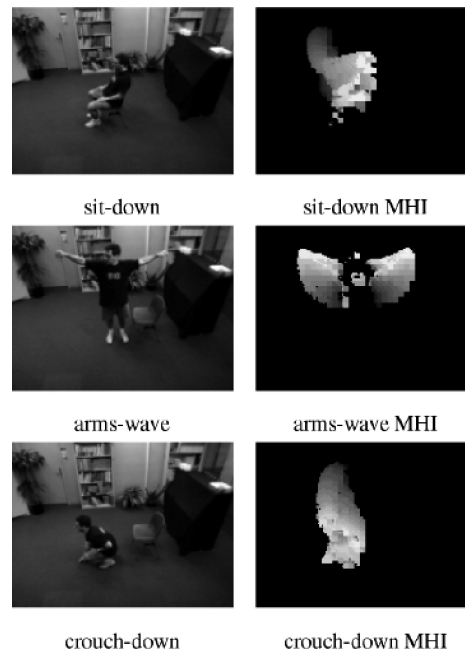


Fig. 5. Examples of space-time action representation: *motion-history images* from Bobick and Davis [2001] (©2001 IEEE). This representation can be viewed as an weighted projection of a 3-D XYT volume into a 2-D XY dimension.

Bobick and Davis [2001] constructed a real-time action recognition system using template matching. Instead of maintaining the 3-dimensional space-time volume of each action, they represented each action with a template composed of two 2-dimensional images: a 2-dimensional binary *motion-energy image* (MEI) and a scalar-valued *motion-history image* (MHI). The two images are constructed from a sequence of foreground images, which essentially are weighted 2-D (XY) projections of the original 3-D XYT space-time volume. By applying a traditional template-matching technique to a pair of (MEI, MHI), their system was able to recognize simple actions like sitting, arm waving, and crouching. Further, their real-time system has been applied to the interactive play environment of children called the Kids-Room. Figure 5 shows example MHIs.

Shechtman and Irani [2005] have estimated motion flows from a 3-D space-time volume to recognize human actions. They have computed a 3-D space-time video-template correlation, measuring the similarity between an observed video volume and maintained template volumes. Their similarity measurement can be viewed as a hierarchical space-time volume correlation. At every location of the volume (i.e., (x, y, t)), they extracted a small space-time patch around the location. Each volume patch captures the flow of a particular local motion, and the correlation between a patch in a template and a patch in video at the same location gives a local match score to the system. By aggregating these scores, the overall correlation between the template volume and a video volume is computed. When an unknown video is given, the system searches for all possible 3-D volume segments centered at every (x, y, t) that best matches the template (i.e., sliding windows). Their system was able to recognize various types of human actions, including ballet movements, pool dives, and waving.

Ke et al. [2007] used segmented spatio-temporal volumes to model human activities. Their system applies a hierarchical meanshift to cluster similarly colored voxels, and obtains several segmented volumes. The motivation is to find the actor volume

segments automatically and to measure their similarity to the action model. Recognition is done by searching for a subset of over-segmented spatio-temporal volumes that best matches the shape of the action model. Support vector machines (SVM) have been applied to recognize human actions while considering both shapes and flows of the volumes. As a result, their system recognized simple actions such as hand waving and boxing from the KTH action database [Schuldt et al. 2004] as well as tennis plays in TV broadcast videos with more complex backgrounds.

Rodriguez et al. [2008] have analyzed 3-D space-time volumes by synthesizing filters: they adopted the maximum average correlation height (MACH) filters that have been used for an analysis of images (e.g., object recognition), to solve the action recognition problem. That is, they have generalized the traditional 2-D MACH filter for 3-D XYT volumes. For each action class, one synthesized filter that fits the observed volume is generated and the action classification is performed by applying the synthesized action MACH filter and analyzing its response on the new observation. They have further extended the MACH filters to analyze vector-valued data using the Clifford Fourier transform. They have not only tested their system on the existing KTH dataset and the Weizmann dataset [Blank et al. 2005], but also on their own dataset constructed by gathering clips from movie scenes. Actions such as kissing and hitting have been recognized.

Table I compares the abilities of the space-time volume-based action recognition approaches. The major disadvantage of space-time volume approaches is the difficulty in recognizing actions when multiple persons are present in the scene. Most of the approaches apply the traditional sliding-window algorithm to solve this problem. However, this requires a large number computations for the accurate localization of actions. Furthermore, they have difficulty recognizing actions that cannot be spatially segmented.

2.1.2. Action Recognition with Space-Time Trajectories. Trajectory-based approaches are recognition approaches that interpret an activity as a set of space-time trajectories. In trajectory-based approaches, a person is generally represented as a set of 2-dimensional (XY) or 3-dimensional (XYZ) points corresponding to his/her joint positions. Human body part estimation methodologies, especially stick figure modeling, have been widely used to extract the joint positions of a person at each image frame. As a human performs an action, changes in his/her joint position are recorded as space-time trajectories, constructing 3-D XYT or 4-D XYZT representations of the action. Figure 6 shows example trajectories. The early work done by Johansson [1975] suggested that the tracking of joint positions is itself sufficient for humans to distinguish actions, and this paradigm has been studied for the recognition of activities in depth [Webb and Aggarwal 1982; Niyogi and Adelson 1994].

Several approaches used the trajectories themselves (i.e., sets of 3-D points) to represent and recognize actions directly [Sheikh et al. 2005; Yilmaz and Shah 2005b]. Sheikh et al. [2005] represented an action as a set of 13 joint trajectories in a 4-D XYZT space. They have used an affine projection to obtain normalized XYT trajectories of an action in order to measure the view-invariant similarity between two sets of trajectories. Yilmaz and Shah [2005b] presented a methodology to compare action videos obtained from moving cameras, also using a set of 4-D XYZT joint trajectories.

Campbell and Bobick [1995] recognized human actions by representing them as curves in low-dimensional *phase spaces*. In order to track joint positions, they took advantage of the 3-D body-part models of a person. Based on the 3-D XYZ models estimated for each frame, they have defined body phase space as a space where each axis represents an independent parameter of the body (e.g., ankle-angle or knee-angle) or its first derivative. In their phase space, a person's static state at each frame corresponds

Table I. A Comparison the Abilities of Important Space-Time Approaches

Approach Type	Authors	Required low-levels	Structural consideration	Scale invariant	Localization	View invariant	Multiple activities
Space-time volume	Bobick and J. Davis '01	Background	Volume-based	Templates needed	✓		
	Shechtman and Irani '05	None	Volume-based	Scaling required	✓		
	Ke et al. '07	None	Volume-based	Templates needed	✓		
Space-time trajectories	Rodriguez et al. '08	None	Volume-based	✓	✓		
	Campbell and Bobick '95	Body-part estimation		✓	✓	✓	
	Rao and Shah '01	Skin detection	Ordering only	✓	✓	✓	
Space-time features	Sheikh et al. '05	Body-part estimation	Ordering only	✓	✓	✓	
	Chomat and Crowley '99	None		✓	✓		
	Zalmik-Manor and Irani '01	None		✓			
	Laptev and Lindeberg '03	None		✓	✓		
	Shuldt et al. '04	None		✓			
	Dollar et al. '05	None		✓			
	Yilmaz and Shah '05a	Background	Ordering only	✓	✓	✓	
	Blank et al. '05	Background		✓	✓	Δ	
	Niebles et al. '06	None		✓	✓		✓
	Wong et al. '07	None	✓	✓	✓		
Space-time features	Savarese et al. '08	None	Proximity-based	✓	✓		✓
	Liu and Shah '08	None	Co-occur only	✓	✓		
	Laptev et al. '08	None	Grid-based	✓			
	Ryoo and Aggarwal '09b	None	✓	✓	✓		✓
		None		✓	✓		✓

The column "required low-levels" specifies the low-level components necessary for the approach to be applicable. "Structural consideration" shows temporal patterns that the approach is able to capture. "Scale invariant" and "view invariant" columns describe whether the approaches are invariant to scale and view changes in videos; "localization" indicates the ability to correctly locate where the activity is occurring spatially and temporally. "Multiple activities" indicates that the system is designed to consider multiple activities in the same scene.

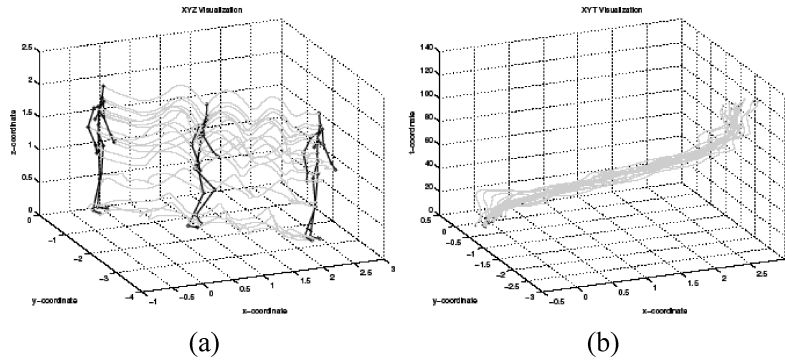


Fig. 6. An example of trajectories of human joint positions when performing the human action of walking [Sheikh et al. 2005] (©2005 IEEE). Figure (a) shows trajectories in XYZ space, and (b) shows those in XYT space.

to a point, and an action corresponds to a set of points (i.e., a curve). The authors have projected the curve in the phase space into multiple 2-D subspaces, and maintained the projected curves to represent the action. Each curve is modeled as a cubic polynomial form, indicating that the authors assume the actions to be relatively simple in the projected subspace. Among all possible curves of 2-D subspaces, their system automatically selects the top k stable and reliable ones to be used for the recognition process.

Once an action representation, that is, a set of projected curves, has been constructed, Campbell and Bobick recognized the action by also converting an unseen video into a set of points in the phase space. Without explicitly analyzing the dynamics of the points from the unseen video, their system simply verifies whether the points are on the maintained curves (i.e., trajectories in the subspaces) when projected. Various types of basic ballet movements have been successfully recognized with markers attached to a subject to track joint positions.

Instead of maintaining trajectories to represent human actions, Rao and Shah's [2001] methodology extracts meaningful curvature patterns from the trajectories. They have tracked the position of a hand in 2-D image space using the skin pixel detection, obtaining a 3-D XYT space-time curve. Their system extracts the positions of peaks of trajectory curves, representing an action as a set of peaks and intervals between them. They verified that these peak features are view-invariant. Automated learning of the human actions is possible in their system by incrementally constructing several action prototypes as representations of human actions. These prototypes can be considered action templates, and the overall recognition process can be regarded as a template-matching process. As a result, by analyzing the peaks of trajectories, their system was able to recognize human actions in an office environment such as "opening a cabinet" and "picking up an object."

Again, Table I compares the trajectory-based approaches. The major advantage of such approaches is their ability to analyze the details of human movements. Furthermore, most of these methods are view-invariant. However, in order to do so, such methods generally require a strong low-level component that is able to correctly estimate the 3-D XYZ joint locations of persons appearing in a scene. The problem of the 3-D body-part detection and tracking is still an unsolved problem and researchers are actively working in this area.

2.1.3. Action Recognition Using Space-Time Local Features. The approaches discussed in this section use local features extracted from 3-D space-time volumes to represent and

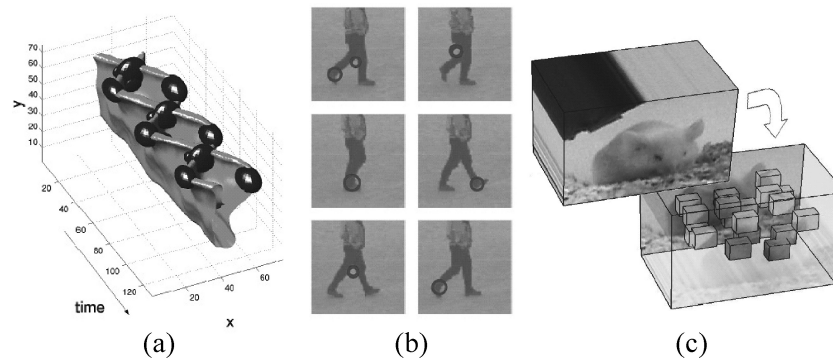


Fig. 7. Example 3-D space-time local features extracted from a video of the human action of walking [Laptev and Lindeberg 2003] (©2003 IEEE), and features from a mouse movement video [Dollar et al. 2005] (©2005 IEEE). Figure (a) shows a concatenated XYT surfaces of a person's legs and detected interest points using Laptev and Lindeberg [2003] methods. Figure (b) shows the same interest points placed on a sequence of original images. Figure (c) shows cuboid features extracted using Dollar et al. [2005] methods.

recognize activities. The motivation behind these approaches is the fact that a 3-D space-time volume essentially is a rigid 3-D object. This implies that if a system is able to extract appropriate features describing characteristics of each action's 3-D volumes, the action can be recognized by solving an object-matching problem.

In this section, we discuss each of the approaches using 3-D space-time features, with a special focus on three aspects: which 3-D local features the approaches extract, how they represent an activity in terms of the extracted features, and what methodology they use to classify activities. In general, we are able to describe the activity recognition approaches using local features by presenting the above three components. Similar to the object recognition process, the system first extracts specific local features that have been designed to capture the local motion information of a person from a 3-D space-time volume. These features are then combined to represent the activities while considering their spatio-temporal relationships or ignoring their relations. Finally, recognition algorithms are applied to classify the activities.

We use the terminology *local features*, *local descriptors*, and *interest points* interchangeably, similar to the case of object recognition problems. Several approaches extract these local features at every frame and concatenate them temporally to describe the overall motion of human activities [Chomat and Crowley 1999; Zelnik-Manor and Irani 2001; Blank et al. 2005]. The other approaches extract sparse spatio-temporal local interest points from 3-D volumes [Laptev and Lindeberg 2003; Dollar et al. 2005; Niebles et al. 2006; Yilmaz and Shah 2005a; Ryoo and Aggarwal 2009b]. Example 3-D local interest points are illustrated in Figure 7. These features have been particularly popular because of their reliability under noise, camera jitter, illumination changes, and background movements.

Chomat and Crowley [1999] proposed the idea of using local appearance descriptors to characterize an action, thereby enabling action classification. Motion energy receptive fields together with Gabor filters are used to capture motion information from a sequence of images. More specifically, local spatio-temporal appearance features that describe motion orientation are detected per frame. Multi-dimensional histograms are constructed based on the detected local features, and the posterior probability of an action occurring given the detected features is calculated by applying the Bayes rule to the histograms. This system first calculates the local probability of an activity occurring at each pixel location and then integrates them for the final recognition of the actions. Even though only simple gestures such as come, go, left, and right are recognized due

to the simplicity of their motion descriptors, they have shown that local appearance detectors may be utilized for the recognition of human activities.

Zelnik-Manor and Irani [2001] proposed an approach utilizing local spatio-temporal features at multiple temporal scales. Multiple temporally scaled video volumes were analyzed to handle execution speed variations of an action. For each point in a 3-D XYT volume, their system estimates a normalized local intensity gradient. Similar to Chomat and Crowley [1999], they computed a histogram of these space-time gradient features per video and presented a histogram-based distance measurement, ignoring the positions of the extracted features. An unsupervised clustering algorithm was applied to these histograms to learn actions, and human activities including outdoor-sport video sequences like basketball and tennis were automatically recognized.

Similarly, Blank et al. [2005] also calculated local features at each frame. Instead of utilizing optical flows for the calculation of local features, they calculated appearance-based local features at each pixel by constructing a space-time volume whose pixel values are solutions to the Poisson equation. The solution to the Poisson equation has proved to be able to extract a wide variety of useful local shape properties, and their system has extracted local features capturing space-time saliency and space-time orientation using the equation. Each sequence of an action is represented as a set of global features, which are the weighted moments of the local features. The authors have applied a simple nearest neighbor classification with a Euclidean distance to recognize the actions. Simple actions such as walking, jumping, and bending in their Weizmann dataset as well as basic ballet movements were recognized successfully.

On the other hand, there are approaches that extract sparse local features from video volumes to represent activities. Laptev and Lindeberg [2003] recognized human actions by extracting sparse spatio-temporal interest points from videos. They extended the previous local feature detectors [Harris and Stephens 1988] commonly used for object recognition, in order to detect interest points in a space-time volume. This scale-invariant interest-point detector searches for spatio-temporal corners in a 3-D space (XYT), which captures various types of non-constant motion patterns. Motion patterns such as a change in direction of an object; splitting and merging of an image structure; and/or collision and bouncing of objects are detected as a result (Figure 7(a) and (b)). In their work, these features were used to distinguish a walking person from complex backgrounds. Furthermore, Schuldt et al. [2004] classified multiple actions by applying SVMs to Laptev and Lindeberg's [2003] features, illustrating their applicability for activity recognition. A new database called the KTH actions dataset, containing action videos (e.g., jogging and hand waving) was introduced, and has been widely adopted. We discuss this dataset further in Section 5.1.1.

The paradigm of recognizing actions by extracting sparse local interest points from a 3-D space-time volume was adopted by several researchers. They have focused on the fact that sparse local features characterizing local motion are sufficient to represent actions, as Laptev and Lindeberg [2003] have suggested. These approaches are particularly motivated by the success of the object recognition methodologies using sparse local appearance features, such as SIFT descriptors [Lowe 1999]. Instead of extracting features at every frame, these approaches extract features only when there exists a salient appearance or shape change in 3-D space-time volume. Most of these features have been verified to be invariant to scale, rotation, and translations, similar to object recognition descriptors.

Dollar et al. [2005] proposed a new spatio-temporal feature detector for the recognition of human (and animal) actions. Their detector is especially designed to extract space-time points with local periodic motions, obtaining a sparse distribution of interest points from a video. Once detected, their system associates a small 3-D volume called *cuboid* to each interest point (Figure 7(c)). Each cuboid captures pixel appearance

values of the interest point's neighborhoods. They have tested various transformations to be applied to cuboids to extract final local features, and have chosen the flattened vector of brightness gradients that shows the best performance. A library of cuboid prototypes is constructed per each dataset by clustering cuboid appearances with k-means. As a result, each action is modeled as a histogram of cuboid types detected in 3-D space-time volume while ignoring their locations (i.e., bag-of-words paradigm). They have recognized facial expressions, mouse behaviors, and human activities (i.e., the KTH dataset) via their method.

Niebles et al. [2006, 2008] presented an unsupervised learning and classification method for human actions using the feature extractor above [Dollar et al. 2005]. Their recognition method is a generative approach, modeling an action class as a collection of spatio-temporal feature appearances. A probabilistic Latent Semantic Analysis (pLSA) commonly used in the field of text mining has been applied to recognize actions statistically. Each feature in the scene is categorized into an action class by calculating its posterior probability of being generated by the action. As a result, Niebles et al. were able to recognize simple actions from public datasets [Schuldt et al. 2004; Blank et al. 2005] as well as figure-skating actions.

In this context, various spatio-temporal feature extractors have been developed recently. Yilmaz and Shah [2005a] proposed an action recognition approach to extract sparse features called *action sketches* from a 3-D contour concatenation, which were confirmed to be view-invariant. Scovanner et al. [2007] designed the 3-D version of the SIFT descriptor, similar to the cuboid features [Dollar et al. 2005]. Liu et al. [2009] presented a methodology to prune cuboid features so as to choose important and meaningful features. Breconzio et al. [2009] proposed an improved detector for extracting cuboid features, and presented a feature selection method similar to Liu et al. [2009]. Rapantzikos et al. [2009] extended the cuboid features to color and motion information as well, in contrast to previous features using intensities only (e.g., Laptev and Lindeberg [2003]; Dollar et al. [2005]).

In most approaches using sparse local features, spatial and temporal relationships among detected interest points are ignored. The approaches that we discussed above show that simple actions can be recognized successfully, even without any spatial and temporal information among features. This is similar to the success of object recognition techniques that ignore the local features' spatial relationships, typically called *bag-of-words*. The bag-of-words approaches were particularly successful for simple periodic actions.

Recently, action recognition approaches that consider spatial configurations among the local features are getting an increasing amount of interest. Unlike the approaches following the bag-of-words paradigm, these approaches attempt to model spatio-temporal distribution of the extracted features for better recognition of actions. Wong et al. [2007] extended the basic pLSA by constructing a pLSA with an implicit shape model (pLSA-ISM). In contrast to the pLSA used by Niebles et al. [2006], the Wong et al. pLSA-ISM captures the relative spatio-temporal location information of the features from the activity center, successfully recognizing and localizing activities in the KTH dataset.

Savarese et al. [2008] proposed a methodology to capture spatio-temporal proximity information among features. For each action video, they measured feature co-occurrence patterns in a local 3-D region, constructing histograms called *ST-correlograms*. Liu and Shah [2008] also considered correlations among features. Similarly, Laptev et al. [2008] constructed spatio-temporal histograms by dividing an entire space-time volume into several grids. The method roughly measures how local descriptors are distributed in the 3-D XYT space by analyzing which feature falls into which grid. Both methods were tested on the KTH dataset as well, obtaining successful

recognition results. Notably, similar to Rodriguez et al. [2008], Laptev et al. [2008] was tested on realistic videos obtained from various movie scenes.

Ryoo and Aggarwal [2009b] introduced the *spatio-temporal relationship match* (STR match), which explicitly considers spatial and temporal relationships among detected features so as to recognize activities. Their method measures structural similarity between two videos by computing pairwise spatio-temporal relations among local features (e.g., *before* and *during*), enabling the detection and localization of complex-structured activities. Their system not only classified simple actions (i.e., those from the KTH datasets), but also recognized interaction-level activities (e.g., hand-shaking and pushing) from continuous videos.

The space-time approaches that extract local descriptors have several advantages. By its nature, background subtraction or other low-level components are generally not required, and the local features are scale-, rotation-, and translation-invariant in most cases. They were particularly suitable for recognizing simple periodic actions such as walking and waving, since periodic actions will repeatedly generate feature patterns.

2.1.4. Comparison. Table I compares the abilities of the space-time approaches reviewed in this article. Space-time approaches are suitable for the recognition of periodic actions and gestures, and many have been tested on public datasets (e.g., the KTH dataset [Schuldt et al. 2004] and the Weizmann dataset [Blank et al. 2005]). Basic approaches using space-time volumes provide a straightforward solution, but often have inherent difficulties in handling speed and motion variations. Recognition approaches using space-time trajectories are able to perform detailed-level analysis and are view-invariant in most cases. However, 3-D modeling of body parts from videos, which is still an unsolved problem, is required for applying a trajectory-based approach.

The spatio-temporal local feature-based approaches are getting an increasing amount of attention due to their reliability under noise and illumination changes. Furthermore, some approaches [Niebles et al. 2006; Ryoo and Aggarwal 2009b] are able to recognize multiple activities without background subtraction or body-part modeling. The major limitation of the space-time feature-based approaches is that they are not suitable for modeling more complex activities. The relations among features are important for a nonperiodic activity that takes a certain amount of time, which most of the previous approaches ignored. Several researchers have worked on approaches to overcome such limitations [Wong et al. 2007; Savarese et al. 2008; Laptev et al. 2008; Ryoo and Aggarwal 2009b]. Viewpoint invariance is another issue that space-time local feature-based approaches must handle.

2.2. Sequential Approaches

Sequential approaches are the single-layered approaches that recognize human activities by analyzing sequences of features. Such approaches consider an input video as a sequence of observations (i.e., feature vectors), and deduce that an activity has occurred in the video if they are able to observe a particular sequence characterizing the activity. Sequential approaches first convert a sequence of images into a sequence of feature vectors by extracting features (e.g., degrees of joint angles) that describe the status of a person per image frame. Once feature vectors have been extracted, sequential approaches analyze the sequence to measure how likely it is that the feature vectors were produced by the person performing the activity. If the likelihood between the sequence and the activity class (or the posterior probability of the sequence belonging to the activity class) is high enough, the system decides that the activity has occurred.

We classify the sequential approaches into two categories by using a methodology-based taxonomy: exemplar-based recognition approaches and state model-based

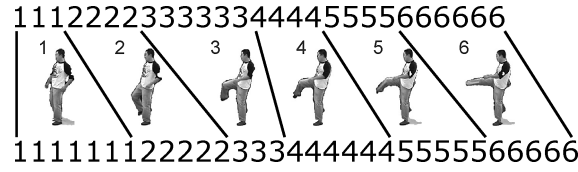


Fig. 8. An example matching between two “stretching a leg” sequences with different nonlinear execution rates. Each number represents a particular status (i.e., pose) of the person.

recognition approaches. Exemplar-based sequential approaches describe classes of human actions by using training samples directly. They maintain either a representative sequence per class or a set of training sequences per activity, and match them with a new sequence to recognize its activity. On the other hand, state model-based sequential approaches represent a human action by constructing a model that is trained to generate sequences of feature vectors corresponding to the activity. By calculating the likelihood (or posterior probability) that a given sequence is generated by each activity model, the state model-based approaches are able to recognize the activities.

2.2.1. Exemplar-Based Approaches. Exemplar-based approaches represent human activities by maintaining a template sequence or a set of sample sequences of action executions. When a new input video is given, the exemplar-based approaches compare the sequence of feature vectors extracted from the video with the template sequence (or sample sequences). If their similarity is high enough, the system is able to deduce that the given input contains an execution of the activity. Humans may perform an identical activity in different styles and/or different rates, and the similarity must be measured considering such variations. The dynamic time warping (DTW) algorithm, originally developed for speech processing, has been widely adopted for matching two sequences with variations [Darrell and Pentland 1993; Gavrilu and Davis 1995; Veeraraghavan et al. 2006]. The DTW algorithm finds an optimal nonlinear match between two sequences with a polynomial amount of computations. Figure 8 shows a conceptual matching between two sequences (i.e., strings) with different execution rates.

Darrell and Pentland [1993] proposed a DTW-based gesture recognition methodology using *view* models to represent the dynamics of articulated objects. Their system maintains multiple models (i.e., template images) of an object in different conditions, which they called views. Each view-model abstracts a particular status (e.g., rotation and scale) of an articulated object such as a hand. Given a video, the correlation scores between image frames and each view are modeled as a function of time. Means and variations of these scores of training videos are used as a gesture template. The templates are matched with a new observation using the DTW algorithm, so that speed variations of action executions are handled. Their system successfully recognized hello and good-bye gestures, and was able to distinguish them from other gestures such as a come closer gesture.

Gavrilu and Davis [1995] also developed the DTW algorithm to recognize human actions, utilizing a 3-D (XYZ) model-based body-part tracking methodology. The motivation is to estimate a 3-D skeleton model at each image frame and to analyze his/her movement by tracking them. Multiple cameras were used to obtain 3-D body-part models of a human, which is composed of a collection of segments and their joint angles (i.e., the stick figure). This stick figure model with 17 degree-of-freedom (DOF) is tracked throughout the frames, recording the values of joint angles. These angle values are treated as features characterizing human movement at each frame. The sequences of

angle values are analyzed using the DTW algorithm to compare them with a reference sequence pretrained per action, similar to Darrell and Pentland [1993]. Gestures including waving hello, waving-to-come, and twisting have been recognized with their system.

Yacoob and Black [1998] have treated an input as a set of signals (instead of discrete sequences) describing sequential changes of feature values. Instead of directly matching the sequences (e.g., DTW), they have decomposed signals using singular value decompositions (SVD). That is, they used principle component analysis (PCA)-based modeling to represent an activity as a linear combination of a set of *activity basis* that is essentially a set of eigen vectors. When a new input is provided to the system, their system calculates the coefficients of the activity basis while considering transformation parameters such as scale and speed variations. The similarity between the input and an action template is measured by comparing the coefficients of the two. Their approach showed successful recognition results for walking-related actions and lip movements by utilizing different types of features.

Efros et al. [2003] presented a methodology for recognizing actions at a distance, where each human is around 30 pixels tall. In order to recognize actions in such environments where the detailed motions of humans are unclear, they used motion descriptors based on optical flows obtained per frame. Their system first computes the space-time volume of each person being tracked, and then calculates 2-D (XY) optical flows at each frame by tracking humans via a temporal difference image similar to Yacoob and Black [1998]. They used blurry motion channels as a motion descriptor, converting optical flows into a spatio-temporal motion descriptor per frame. That is, they are interpreting a video of a human action as a sequence of motion descriptors obtained from the optical flows of a human. The basic nearest neighbor classification method was applied to a sequence of motion descriptors for the recognition of actions. First, frame-to-frame similarities between all possible pairs of frames from two sequences (i.e., a frame-to-frame similarity matrix) are calculated. The recognition is done by detecting diagonal patterns in the frame-to-frame similarity matrix. Their system was able to classify ballet movements, tennis plays, and soccer plays, even from moving cameras.

Lublinerman et al. [2006] presented a methodology that recognizes human activities by modeling them as linear-time-invariant (LTI) systems. Their system converts a sequence of images into a sequence of silhouettes, extracting two types of contour representations: silhouette width and Fourier descriptors. An activity is represented as a LTI system capturing the dynamics of changes in silhouette features. SVMs have been applied to classify a new input which has been converted to the parameters of a LTI model. Four types of simple actions, slow walk, fast walk, walk on an incline and walk with a ball have been correctly recognized as a consequence.

Veeraraghavan et al. [2006] described an activity as a function of time that describes parameter changes similar to Yacoob and Black [1998]. The main contribution of Veeraraghavan et al.'s system is in the explicit modeling of inter- and intra-personal speed variations of activity executions and the consideration of such methods for matching activity sequences. Focusing on the fact that humans may be able to change the speed of an execution of a part of the activity while it may not be possible for other parts, they learn nonlinear characteristics of activity speed variations. More specifically, their system learns the nature of time-warping transformation per activity. They are modeling an action execution with two functions: (i) a function of feature changes over time and (ii) a function space of possible time warping. They have developed an extension of a DTW matching algorithm to take the time-warping function into account when matching two sequences. Human actions including picking up an object, throwing, pushing, and waving have been recognized with high recognition accuracy.

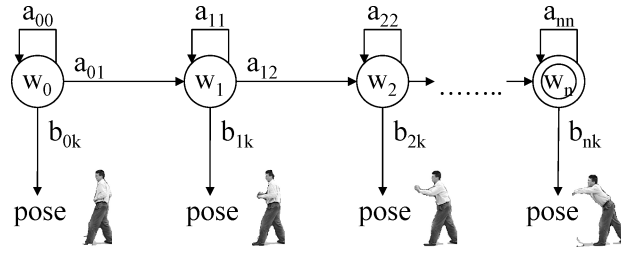


Fig. 9. An example hidden Markov model for the action *stretching an arm*. The model is one of the simplest cases among HMMs, which is designed to be strictly sequential. Each actor image in the figure represents a pose with the highest observation probability b_{jk} for its state w_j .

2.2.2. State Model-Based Approaches. State model-based approaches are the sequential approaches that represent a human activity as a model composed of a set of states. The model is statistically trained so that it corresponds to sequences of feature vectors belonging to its activity class. More specifically, the statistical model is designed to generate a sequence with a certain probability. Generally, one statistical model is constructed for each activity. For each model, the probability of the model generating an observed sequence of feature vectors is calculated to measure the likelihood between the action model and the input image sequence. Either the maximum likelihood estimation (MLE) or the maximum a posteriori probability (MAP) classifier is constructed as a result, in order to recognize activities.

Hidden Markov models (HMMs) and dynamic Bayesian networks (DBNs) have been widely used for state model-based approaches. In both cases, an activity is represented in terms of a set of hidden states. A human is assumed to be in one state at each time frame, and each state generates an observation (i.e., a feature vector). In the next frame, the system transitions to another state considering the transition probability between states. Once transition and observation probabilities are trained for the models, activities are commonly recognized by solving the evaluation problem. The evaluation problem is the problem of calculating the probability of a given sequence (i.e., new input) generated by a particular state-model. If the calculated probability is high enough, the state model-based approaches are able to decide that the activity corresponding to the model occurred in the given input. Figure 9 shows an example of a sequential HMM.

Yamato et al.'s [1992] work is the first in applying standard HMMs to recognize activities. They adopted HMMs which were originally widely used for speech recognition. At each frame, their system converts a binary foreground image into an array of meshes. The number of pixels in each mesh is considered a feature, thereby extracting a feature vector per frame. These feature vectors are treated as a sequence of observations generated by the activity model. Each activity is represented by constructing one HMM that corresponds probabilistically to particular sequences of feature vectors (i.e., meshes). More specifically, the parameters of HMMs (transition probabilities and observation probabilities) are trained with a labeled dataset with the standard learning algorithm for HMMs. Once each of the HMMs is trained, they are used to recognize activities by measuring the likelihoods between a new input and the HMMs by solving the evaluation problem. As a result, various types of tennis plays, such as backhand stroke, forehand stroke, smash, and serve, have been recognized with Yamato et al.'s system. They have shown that the HMMs are able to reliably model feature changes during human activities, encouraging other researchers to pursue further investigations.

Starner and Pentland [1995] also used standard HMMs, in order to recognize American Sign Language (ASL). Their method tracks the location of hands and extracts features that describes shapes and positions of the hands. Each word of ASL is modeled

as one HMM, generating a sequence of features describing hand shapes and positions, similar to the case of Yamato et al. [1992]. Their method uses the Viterbi algorithm for each HMM to estimate the probability that the HMM generated the observations. The Viterbi algorithm provides an efficient approximation of the likelihood distance, enabling an unknown observation sequence to be classified into the most suitable word.

Bobick and Wilson [1997] also recognized gestures by using state models. They represented a gesture as a 2-D XY trajectory describing the location changes of a hand. Each curve is decomposed into sequential vectors, which can be interpreted as a sequence of states computed from a training example. Furthermore, each state is made fuzzy in order to consider speed and motion variance in executions of the same gesture. This is similar to a fuzzy version of a sequential Markov model (MM). Transition costs between states, which correspond to the transition probabilities in the case of HMMs, are also defined in their system. For the recognition of gestures with their model, a dynamic programming algorithm is designed. Their system measures an optimal matching cost between the given observation (i.e., motion trajectory) and each prototype using the dynamic programming algorithm. Applying their framework, they have successfully recognized two different types of gestures: *wave* and *point*.

In addition, approaches using variants of HMMs have also been developed for human activity recognition [Oliver et al. 2000; Park and Aggarwal 2004; Natarajan and Nevatia 2007]. Similar to previous frameworks for action recognition using HMMs [Yamato et al. 1992; Starner and Pentland 1995; Bobick and Wilson 1997], they constructed one model (HMM) for each activity they wanted to recognize, and used visual features from the scene as observations directly generated by the model. The methods with extended HMMs are designed to handle more complex activities (usually combinations of multiple simple actions) by extending the structure of the basic HMM.

Oliver et al. [2000] constructed a variant of the basic HMM, the coupled HMM (CHMM), to model human-human interactions. The major limitation of the basic HMM is its inability to represent activities composed of motions of two or more agents. A HMM is a sequential model and only one state is activated at a time, preventing it from modeling the activities of multiple agents. Oliver et al. introduced the concept of the CHMM to model complex interactions between two persons. Basically, a CHMM is constructed by coupling multiple HMMs, where each HMM models the motion of one agent. They have coupled two HMMs to model human-human interactions. More specifically, they coupled the hidden states of two different HMMs by specifying their dependencies. As a result, their system was able to recognize complex interactions between two persons, such as a concatenation of “two persons approaching, meeting, and continuing together.”

Park and Aggarwal [2004] used a DBN to recognize gestures of two interacting persons. They recognized gestures such as “stretching an arm” and “turning a head left,” by constructing a tree-structured DBN to take advantage of the dependent nature among the motions of body parts. A DBN is an extension of an HMM, composed of multiple conditionally independent hidden nodes that generate observations at each time frame directly or indirectly. In Park and Aggarwal’s work, a gesture is modeled as state transitions of hidden nodes (i.e., body-part poses) in one time point to the next time point. Each pose is designed to generate a set of features associated with the corresponding body part. Features including locations of skin regions, maximum curvature points, and the ratio and orientation of each body-part have been used to recognize gestures.

Natarajan and Nevatia [2007] developed an efficient recognition algorithm using coupled hidden semi-Markov models (CHSMMs), which extend previous CHMMs by explicitly modeling the duration of an activity staying in each state. In the case of

Table II. Comparing Sequential Approaches

Type	Approaches	Required low-levels	Execution variations	Probabilistic	Target activities
Exemplar-based	Darrell and Pentland '93	None	✓		Gesture-level
	Gavrila and L. Davis '95	Body-part estimation	✓		Gesture-level
	Yacoob and Black '98	Body-part estimation	✓		Gesture-level
	Efros et al. '03	Tracking	Linear only		Action-level
	Lubliner et al. '06	Background subtraction	Linear only		Action-level
	Veeraraghavan et al. '06	Background subtraction	✓		Action-level
State model-based	Yamato et al. '92	Background subtraction	Model-based	✓	Action-level
	Starner and Pentland '95	Tracking	Model-based	✓	Gesture-level
	Bobick and Wilson '97	Tracking	Model-based		Gesture-level
	Oliver et al. '00	Background subtraction	Model-based	✓	Interaction-level
	Park and Aggarwal '04	Background subtraction	Model-based	✓	Gesture-level
	Natarajan and Nevatia '07	Action recognition	Model-based	✓	Interaction-level
	Lv and Nevatia '07	3-D pose model	Model-based	✓	Action-level

The column “required low-levels” specifies the low-level components necessary for the approach to be applicable. “Execution variations” shows whether the system is able to handle variations in the execution of human activities (e.g., speed variations). ‘Probabilistic’ indicates that the system makes a probabilistic inference, and ‘target activity’ shows the type of human activities the system aims to recognize. Notably, [Lv and Nevatia 2007]’s system is view-invariant.

basic HMMs and CHMMs, the probability of a person staying in an identical state decays exponentially as time increases. In contrast, each state in a CHSMM has its own duration that best models the activity that the CHSMM represents. As a result, they were able to construct a statistical model that captures the characteristics of activities that the system wants to recognize better compared to HMMs and CHMMs. Similar to Oliver et al. [2000], they tested their system for the recognition of human-human interactions. Due to the CHSMMs’ ability to model the duration of the activity, the recognition accuracy using CHSMMs was better than other simpler statistical models. Lv and Nevatia [2007] also designed a CHMM-like structure called the *action net* to construct a view-invariant recognition system using synthetic 3-D human poses.

2.2.3. Comparison. In general, sequential approaches consider sequential relationships among features in contrast to most of the space-time approaches, thereby enabling detection of more complex activities (i.e., nonperiodic activities such as sign languages). In particular, the recognition of the interactions of two persons, whose sequential structure is important, was attempted in Oliver et al. [2000] and Natarajan and Nevatia [2007].

Compared to the state model-based sequential approaches, exemplar-based approaches provide more flexibility for the recognition system, in the sense that multiple sample sequences (which may be completely different) can be maintained by the system. Further, the dynamic time-warping algorithm generally used for the exemplar-based approaches provides a nonlinear matching methodology that considers execution rate variations. In addition, exemplar-based approaches are able to cope with less training data than the state model-based approaches.

On the other hand, state-based approaches are able to make a probabilistic analysis of the activity. A state-based approach calculates a posterior probability of an activity occurring, enabling it to be easily incorporated with other decisions. One of the limitations of the state-based approaches is that they tend to require a large number of training videos as the activity they want to recognize becomes more complex. Table II is provided for comparing the systems.

3. HIERARCHICAL APPROACHES

The main idea behind hierarchical approaches is to enable the recognition of high-level activities based on the recognition results of other simpler activities. The motivation is to let the simpler subactivities (also called subevents) that can be modeled relatively easily to be recognized first, and then to use them for the recognizing higher-level activities. For example, a high-level interaction like “fighting” may be recognized by detecting a sequence of several punching and kicking interactions. Therefore, in hierarchical approaches, a high-level human activity (e.g., fighting) that the system aims to recognize is represented in terms of its subevents (e.g., punching), which may themselves be decomposable until atomicity is obtained. That is, subevents serve as observations generated by a higher-level activity. The paradigm of hierarchical representation not only makes the recognition process computationally tractable and conceptually understandable, but also reduces redundancy in the recognition process by reusing recognized subevents multiple times.

In general, common activity patterns of motion that appear frequently during high-level human activities are modeled as atomic-level (or primitive-level) actions, and high-level activities are represented and recognized by concatenating them hierarchically. In most hierarchical approaches, these atomic actions are recognized by adopting single-layered recognition methodologies which we presented in the previous section. For example, the gestures “stretching hand” and “withdrawing hand” occur often in human activity, implying that they can become good atomic actions for representing human activities such as shaking hands or punching. Single-layered approaches such as sequential approaches using HMMs can be safely adopted for recognition of those gestures.

The major advantage of hierarchical approaches over nonhierarchical approaches (i.e., single-layered approaches) is their ability to recognize high-level activities with more complex structures. Hierarchical approaches are especially suitable for a semantic-level analysis of interactions between humans and/or objects as well as complex group activities. This advantage is a result of two abilities of hierarchical approaches: the ability to cope with less training data and the ability to incorporate prior knowledge into the representation.

First, the amount of training data required for recognizing activities with hierarchical models is significantly less than that for single-layered models. Even though in some cases it may also be possible for nonhierarchical approaches to model complex human activities, they generally require a large amount of training data. For example, single-layered HMMs need to learn a large number of transition and observation probabilities, since the number of hidden states increases as the activities get more complex. By encapsulating structurally redundant subevents shared by multiple high-level activities, hierarchical approaches model the activities with a lesser amount of training and recognize them more efficiently.

In addition, the hierarchical modeling of high-level activities makes recognition systems to incorporate human knowledge (i.e., prior knowledge of the activity) much easier. Human knowledge can be included in the system by listing semantically meaningful sub-activities composing a high-level activity and/or by specifying their relationships. As mentioned above, when modeling high-level activities, non-hierarchical techniques tend to have complex structures and observation features which are not easily interpretable, preventing a user from imposing prior knowledge. On the other hand, hierarchical approaches model a high-level activity as an organization of semantically interpretable subevents, making the incorporation of prior knowledge much easier.

Using our approach-based taxonomy, we categorize hierarchical approaches into three groups: statistical approaches, syntactic approaches, and description-based

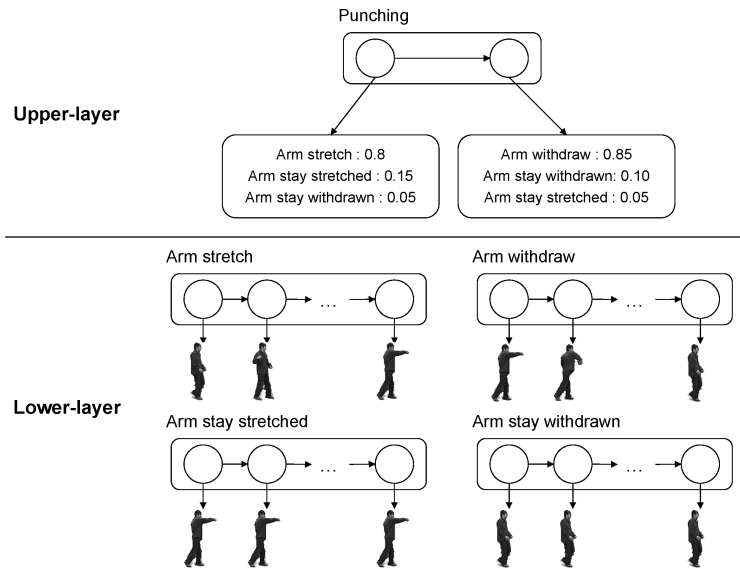


Fig. 10. An example hierarchical hidden Markov model (HHMM) for recognizing the activity of punching. The model is composed of two layers. In the lower layer, HMMs are used to recognize various atomic-level activities, such as stretching and withdrawing. The upper layer HMM treats recognition results of the lower layer HMMs as an input, recognizing that punching stretching, and withdrawing occurred in a sequence.

approaches. Figure 3 illustrates our taxonomy tree, as well as the lists of selected previous work corresponding to the categories.

3.1. Statistical Approaches

Statistical approaches use statistical state-based models to recognize activities. In the case of hierarchical statistical approaches, multiple layers of state-based models (usually two layers) such as HMMs and DBNs are used to recognize activities with sequential structures. At the bottom layer, atomic actions are recognized from sequences of feature vectors, just as in single-layered sequential approaches. As a result, a sequence of feature vectors is converted to a sequence of atomic actions. The second-level models treat this sequence of atomic actions as observations generated by the second-level models. For each model, a probability of the model generating a sequence of observations (i.e., atomic-level actions) is calculated to measure the likelihood between the activity and the input image sequence. Either the maximum likelihood estimation (MLE) or the maximum a posteriori probability (MAP) classifier is constructed as a result. Figure 10 shows an example model of a statistical hierarchical approach, which is designed to recognize “punching.”

Oliver et al. [2002] presented layered hidden Markov models (LHMMs), which is one of the most fundamental forms of the hierarchical statistical approaches (e.g., Figure 10). In this approach, the bottom layer HMMs recognize atomic actions of a single person by matching the models with the sequence of feature vectors extracted from videos. The upper layer HMMs treat recognized atomic actions as observations generated by the upper layer HMMs. That is, essentially they are representing a high-level activity as a sequence of atomic actions by making each state in the upper layer HMM to probabilistically correspond to one atomic action. By its nature, all subevents of an activity are required to be strictly sequential in each LHMM. Human-human

interactions in a conference room environment including “a person giving a presentation” and “face-to-face conversation” have been recognized based on the detection of atomic-level actions (e.g., “nobody”, “one active person”, and “multiple persons present”). Each layer of the HMM is designed to be trained separately with fully labeled data, enabling flexible retraining.

The paradigm of multilayered HMMs has been explored by various researchers. Nguyen et al. [2005] also constructed hierarchical HMMs of two layers to recognize complex sequential activities. Similar to Oliver et al. [2002], they have constructed two levels of HMMs to recognize human activities such as ‘a person having a meal’ and ‘a person having a snack’. Zhang et al. [2006] constructed multi-layered HMMs to recognize group activities occurring in a meeting room. Their framework is also composed of two-layered HMMs. Their system recognized atomic actions of speaking, writing, and idling using the lower-layer HMMs. With the upper layer HMMs, group activities such as monologue, discussion, and presentation have been represented and recognized in terms of the atomic actions. Yu and Aggarwal [2006] used a block-based HMM for the recognition of a person climbing a fence. This block-based HMM can also be interpreted as a 2-layered HMM.

In addition, hierarchical approaches using DBNs have been studied for the recognition of complex activities. DBNs may contain multiple levels of hidden states, suggesting that they can be formulated to represent hierarchical human activities. Gong and Xiang [2003] have extended traditional HMMs to construct dynamic probabilistic networks (DPNs) to represent activities of multiple participants. Their method was able to recognize the group activities of trucks loading and unloading cargo. Dai et al. [2008] constructed DBNs to recognize group activities in a conference room environment similar to Zhang et al. [2006]. High-level activities such as break, presentation, and discussion were recognized based on the atomic actions of talking, asking, and so on. Damen and Hogg [2009] constructed Bayesian networks using a Markov chain Monte Carlo (MCMC) for hierarchical analysis of bicycle-related activities (e.g., drop-and-pick). They used Bayesian networks to model relations between atomic-level actions; and these Bayesian networks were iteratively updated using the MCMC to search for the structure that best explains ongoing observations.

Shi et al. [2004] proposed a hierarchical approach using a propagation network (*P-net*). The structure of a *P-net* is similar to that of a HMM: an activity is represented in terms of multiple state nodes, their transition probabilities, and the observation probabilities. Their work also decomposes actions into several atomic actions, and constructs a network describing the temporal order needed among them. The main difference between a *P-net* and a HMM is that the *P-net* allows activation of multiple state nodes simultaneously. This implies that a *P-net* is able to model a high-level activity composed of concurrent as well as sequential sub-events. If the subevents are activated in a particular temporal order specified through the graph, the system is able to deduce that the activity occurred. Shi et al. have represented an activity of a person performing a chemical experiment using a *P-net*, and have successfully recognized it.

Statistical approaches are especially suitable when recognizing sequential activities. With enough training data, statistical models are able to reliably recognize corresponding activities even in the case of noisy inputs. The major limitation of statistical approaches is their inherent inability to recognize activities with complex temporal structures, such as an activity composed of concurrent subevents. For example, HMMs and DBNs have difficulty modeling the relationship of an activity *A* occurred during, started with, or finished with an activity *B*. The edges of HMMs or DBNs specify the sequential order between two nodes, suggesting that they are suitable for modeling sequential relationships, not concurrent relationships.

Fighting	->	Punching	:	0.3	Punching	->	stretch withdraw	:	0.8
		Punching Fighting	:	0.7			stretch stay_withdrawn	:	0.1
							stay_stretched withdraw	:	0.1

Fig. 11. A simplified example of production rules of an SCFG for representing and recognizing a fighting interaction. *Fighting* is defined as any number of consecutive punching actions which can be decomposed into stretching and withdrawal similar to Figure 10.

3.2. Syntactic Approaches

Syntactic approaches model human activities as a string of symbols, where each symbol corresponds to an atomic-level action. Similar to the case of hierarchical statistical approaches, syntactic approaches also require atomic-level actions to be recognized first, using any of the previous techniques. Human activities are represented as a set of production rules generating a string of atomic actions, and they are recognized by adopting parsing techniques from the field of programming languages. Context-free grammars (CFGs) and stochastic context-free grammars (SCFGs) have been used by previous researchers to recognize high-level activities. The production rules of CFGs naturally lead to a hierarchical representation and recognition of the activities. Figure 11 is an example SCFG.

Ivanov and Bobick [2000] proposed a hierarchical approach for the recognition of high-level activities using SCFGs. They divided the framework into two layers: the lower layer using HMMs for the recognition of simple (i.e., atomic) actions, and the higher layer using stochastic parsing techniques for the recognition of high-level activities. Ivanov and Bobick have encoded a large number of stochastic productions rules which are able to explain all activity possibilities in their environment. The higher layer parses a string of atomic actions generated by the lower layer, recognizing activities probabilistically. The Earley–Stolcke parsing algorithm is extended to handle uncertain observations. Moore and Essa [2002] also used SCFGs for the recognition of activities, focusing on multitask activities. By extending Ivanov and Bobick [2000], they introduced more reliable error detection and recovery techniques for recognition. They were able to recognize human activities happening in a blackjack card game, such as “a dealer dealt a card to a player”, with a high accuracy.

Minnen et al. [2003] adopted SCFGs for activity recognition as well. Their system focuses on the segmentation problem of multiple objects. They have shown that the semantic-level processing of activities using CFGs may help the segmentation and tracking of objects. The concept of the *hallucinations* is introduced to compensate for the failures of atomic-level recognition explicitly. Taking advantage of the CFG parsing techniques while considering hallucinations, they recognized the activity of a person working on the Tower of Hanoi problem, and were able to correctly recognize the activities without any information on the objects’ appearance by depending solely on the motion information of the activities.

Joo and Chellappa [2006] designed an attribute grammar for recognition, which is an extension of the SCFG. Their grammar attaches semantic tags and conditions to the production rules of the SCFG, enabling the recognition of more descriptive activities. That is, their grammar is able to describe feature constraints as well as the temporal constraints of atomic actions. Only when the observations satisfy the syntax of the SCFG (i.e., only when the string can be generated by following the production rules) and the feature constraints are satisfied as well does their system decide that the activity has occurred. Thus events in a parking lot were recognized by tracking cars and humans. Atomic actions including parking, picking up, and walk though are first detected based on changes in the location of cars and humans. By representing the typical activity in a parking lot, normal and abnormal activities are distinguished.

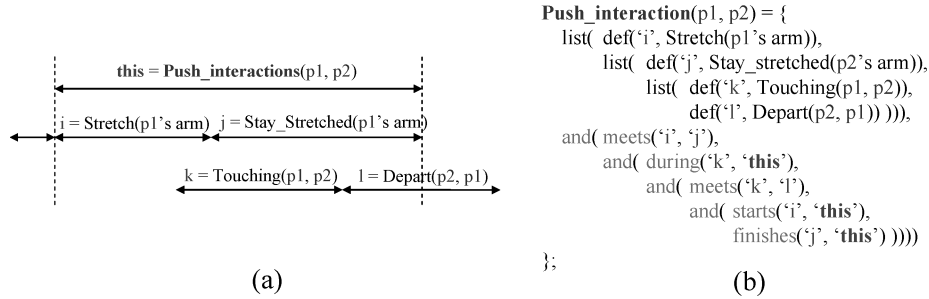


Fig. 12. (a) Time intervals of an interaction *push* and its subevents, and (b) its programming language-like representation following the Ryoo and Aggarwal [2006] syntax (©2009 Springer); (a) is a conceptual illustration describing the activity's temporal structure, whose sub-events are organized sequentially as well as concurrently. Following the CFG, we convert this into a formal representation as shown in (b).

One of the limitations of syntactic approaches is in the recognition of concurrent activities. Syntactic approaches are able to probabilistically recognize hierarchical activities composed of sequential subevents, but are inherently limited for activities composed of concurrent subevents. Since syntactic approaches model a high-level activity as a string of atomic-level activities that compose them, the temporal ordering of atomic-level activities has to be strictly sequential. In addition, syntactic approaches assume that all observations are parsed by applying their production rules. For these systems, a user must provide a set of production rules for all possible events, even for a large domain. Therefore, they tend to have difficulty when an unknown observation (e.g., a pedestrian) interferes with the system. In order to overcome such a limitation, there was an attempt by Kitani et al. [2007] to develop an algorithm to learn grammar rules from observations automatically.

3.3. Description-Based Approaches

A description-based approach is a recognition approach that explicitly maintains human activities' spatio-temporal structures. It represents a high-level human activity in terms of simpler activities that compose the activity (i.e., subevents), describing their temporal, spatial, and logical relationships. That is, description-based approaches model a human activity as an occurrence of its subevent (which might be composed of their own subevents) that satisfies certain relations. Therefore, recognition of the activity is performed by searching the subevents satisfying the relations specified in its representation. All description-based approaches are inherently hierarchical (since they use subevents to represent human activities), and they are able to handle activities with concurrent structures.

In description-based approaches, a *time interval* is usually associated with an occurring subevent to specify necessary temporal relationships among subevents. Allen's temporal predicates [Allen 1983; Allen and Ferguson 1994] have been widely adopted for these approaches, to specify relationships between time intervals [Pinhanez and Bobick 1998; Siskind 2001; Nevatia et al. 2003; Vu et al. 2003; Ryoo and Aggarwal 2006a; Gupta et al. 2009]. Seven basic predicates that Allen has defined are: *before*, *meets*, *overlaps*, *during*, *starts*, *finishes*, and *equals*. Note that the predicates *before* and *meets* describe sequential relationships while the other predicates are used to specify concurrent relationships. Figure 12(a) illustrates a conceptual temporal structure of the human-human interaction *pushing* represented in terms of time intervals.

In a description-based approach, a CFG is often used as a formal syntax for the representation of human activities [Nevatia et al. 2004; Ryoo and Aggarwal 2006a]. Notice

that the description-based approaches' use of CFGs is completely different from that of syntactic approaches: Syntactic approaches directly use CFGs for the recognition, implying that the CFGs themselves describe the semantics of the activities. On the other hand, a description-based approach adopts a CFG as a syntax to formally represent the activities. The activity semantics are usually encoded in a structure similar to that of a programming language (e.g., Figure 12(b)), and the CFG only plays a role to ensure that the activity representation fits its grammar. In general, the recognition is performed by developing an approximation algorithm to solve the constraint satisfaction problem (which is NP-hard).

Pinhanez and Bobick [1998] directly adopted the concept of Allen's interval algebra constraint network (*IA-network*) [Allen 1983] to describe the temporal structure of activities. In an *IA-network*, subevents are specified as nodes and their temporal relationships are described with typed edges between them. Pinhanez and Bobick developed a methodology to convert an *IA-network* into a {past, now, future} network (*PNF-network*). The *PNF-network* that they proposed is able to describe the identical temporal information contained in the *IA-network*, while making it computationally tractable.

Pinhanez and Bobick developed a polynomial time algorithm to process the *PNF-network*. Their system recognizes the top-level activity by checking which subevents have already occurred and which have not. They have shown that their representation is expressive enough to recognize cooking activities occurring in a kitchen environment, such as picking up a bowl. Atomic-level actions were labeled manually from the video in the experiments; the system was able to recognize the activities even when one of the atomic actions was not provided. One of the drawbacks of the system is that a subnetwork corresponding to a subevent has to be specified redundantly if it is used multiple times. Another limitation is that all subevent relations have to be expressed in a network form.

Intille and Bobick [1999] designed a description-based recognition approach to analyze plays in American football. Even though their system was limited to using conjunctions of relatively simple temporal predicates (*before* and *around*), they have shown that complex human activities can be represented by listing the temporal constraints in a format similar to those of programming languages, instead of a network form. Intille and Bobick represented human activities via three levels of hierarchy: atomic-level, individual-level, and team-level activities.

A Bayesian belief network is constructed for the recognition of the activity, based on its temporal structure representation. The root node of the belief network corresponds to the high-level activity that the system aims to recognize. The other nodes correspond to the occurrence of the subevents or describe the temporal relationships between the subevents. The nodes become true if the subevents occur and the relationships are satisfied. Only when all nodes are probabilistically satisfied and propagated to the root node, is the activity said to be detected. Similar to Pinhanez and Bobick [1998], Intille and Bobick have used manually labeled data.

Siskind [2001] also proposed a hierarchical description-based approach for human activity recognition. Notably, it was able to represent and recognize high-level activities with more than three levels. Siskind's methodology uses force dynamics for the recognition of simple actions, and uses the description-based approach called *event logic* to recognize high-level activities. It particularly focused on the recognition of an activity with a *liquid* characteristic, whose occurrences are true for all subintervals of a particular time interval. The approach computes the recognized activity's time interval by calculating the union and intersection of subevents' time intervals, assuming liquidity. This suggests that the recognized activity itself can be used as a subevent of another activity, but is permitted to be used only once.

Nevatia et al. [2003] designed a representation language called VERL to describe human activities. They classified human activities into three categories, similar to Intille and Bobick [1999], enabling the representation of human activities that have three levels of hierarchy: primitive events, single-thread composite events, and multithread composite events. Allen's temporal predicates, spatial predicates, and logical predicates were used to represent human activities by specifying their necessary conditions. Bayesian networks are used for primitive event recognition, and HMMs are used for the recognition of single-thread composite events (i.e., they are strictly sequential). A heuristic algorithm is designed for the constraint satisfaction problem, recognizing interactions between multiple persons. The system is probabilistic, but was not able to overcome the failures of low-level components. The Vu et al. [2003] approach is similar to Nevatia et al. [2003], while extending the representation to describe activities with any levels of hierarchy. However, unlike the Nevatia et al. system, only conjunctive predicates are allowed when concatenating multiple temporal relationships (i.e., only *and* is allowed, not *or*). Hakeem et al. [2004] designed a representation language, CASEE, which also represents an activity as a conjunction of necessary temporal and causal relations.

Several researchers utilized *Petri nets* to represent and recognize human activities [Zaidi 1999; Nam et al. 1999; Ghanem et al. 2004]. Petri nets specify the temporal ordering of an activity's subevents in terms of a graph representation. The recognition is done by sequentially handing tokens in the graph, where each node corresponds to a state before (or after) the completion of particular sub-events. Zaidi [1999] showed that the Petri nets are able to fully represent temporal relationships described by Allen's temporal predicates. Nam et al. [1999] applied the Petri nets for the recognition of hand gestures from videos. Ghanem et al. [2004] took advantage of the Petri nets to represent and recognize interactions between humans and vehicles similar to Ivanov and Bobick [2000]. Due to the Petri net characteristic that the tokens cannot describe multiple possibilities and are nonreversible (i.e., the recognition process is strictly sequential), these deterministic systems have limitations in terms of processing complex scenes.

In order to overcome the limitations of the previous approaches, Ryoo and Aggarwal [2006a] proposed a description-based approach using a CFG as a syntax of their representation language. Their formal grammar enables the representation of human-human interactions with any levels of hierarchy, which are described as logical concatenations (*and*, *or*, and *not*) of complex temporal and spatial relationships among their subevents. As a result, they represent high-level human interactions composed of concurrent subevents (e.g., hand shaking and pushing) in terms of time interval variables and predicates (e.g., Allen's temporal predicates). They have developed a hierarchical semantic matching between the observations and the representations for activity recognition. At the lowest level, Bayesian networks and HMMs are used for the recognition of atomic actions from a sequence of raw image frames. Recognition of represented high-level activities is done by performing a hierarchical matching from the bottom to the top.

In addition, their approach was extended to recognize recursive activities of a continuous nature, such as fighting and greeting [Ryoo and Aggarwal 2006b]. Even though the representation of recursive activities with sequential subevents is possible with syntactic approaches, the recognition of recursive activities with complex concurrent subevents has been limited. The authors have introduced the special time interval *this*, which always corresponds to the activity being represented, and proposed an iterative algorithm to recognize activities described using *this*. With the proposed approach, the recursive activity of fighting was represented as a single negative interaction (e.g., punching and pushing) followed by shorter episode of fighting, and was successfully.

Furthermore, Ryoo and Aggarwal [2009a] proposed a probabilistic extension of their recognition framework that is able to compensate for the failures of its low-level components. One of the limitations of description-based approaches is that they are mostly deterministic, and are fragile when their low-level components are noisy. Ryoo and Aggarwal have overcome such limitations. They have used a logistic regression to model the probability distribution of an activity, and used it to detect the activity even when some of its subevents have been misclassified. In order to compensate for the complete failure of the atomic-level components (i.e., no atomic action detected at all), they took advantage of the concept of the *hallucination* time intervals, similar to the ones used in Minnen et al. [2003].

There have also been attempts to adopt symbolic artificial intelligence techniques to recognize human activities. Tran and Davis [2008] adopted Markov logic networks (MLNs) to probabilistically infer events in a parking lot. This 2-layered approach successfully handled uncertainties in human activities. However, their MLNs relied on the assumption that an identical subevent occurs only once during interactions, thus limiting itself from being applied to dynamically interacting actors.

Gupta et al. [2009] recently presented a description-based approach for a probabilistic analysis as well. Unlike other description-based approaches designed to recognize complex activities, their approach aims to recognize atomic-level actions more reliably by modeling causality among the actions. A tree structured AND-OR graph similar to Hongeng et al. [2004] was used to represent a storyline of a sports game (e.g., baseball), labeling each action (e.g., hitting) that fits the storyline. Their system iteratively searches for the best explanatory AND-OR graph structures and the best video-action associations by taking advantage of captions and video trajectories. That is, a representation-fitting algorithm has been developed.

3.4. Comparison

Hierarchical approaches are suitable for recognizing high-level activities which can be decomposed into simpler subevents. Due to their nature, they can more easily incorporate human knowledge into the systems and require less training data, as pointed out by many researchers [Oliver et al. 2002; Nevatia et al. 2003; Ryoo and Aggarwal 2006a]. Statistical and syntactic approaches provide a probabilistic framework for reliable recognition with noisy inputs. However, they have difficulties representing and recognizing activities with concurrently organized subevents.

Description-based approaches are able to represent and recognize human activities with complex temporal structures. Not only sequentially occurring, but also concurrent organized subevents are handled with description-based approaches. The major drawback of description-based approaches is their inability to compensate for the failures of low-level components (e.g., gesture detection failure). That is, most of the description-based approaches have a deterministic high-level component. Pinhanez and Bobick [1998] showed that the high-level system has the potential to compensate for a single low-level detection failure, and a couple of recent works have proposed probabilistic frameworks for description-based approaches [Ryoo and Aggarwal 2009a; Gupta et al. 2009]. Table III compares the abilities of important hierarchical approaches.

4. HUMAN-OBJECT INTERACTIONS AND GROUP ACTIVITIES

In this section we present and summarize previous papers on the recognition of human-object interactions and those on the recognition of group activities. Such approaches fall into different categories if the approach-based taxonomy of the previous sections is applied as shown in Figures 2 and 3. However, even though they use different methodologies for recognition, they exhibit interesting common properties and characteristics as they share the same objective. In the first section (4.1), we discuss approaches for

Table III. Comparing the Abilities of the Hierarchical Approaches

Type	Approaches	Levels of hierarchy	Complex temporal relations	Complex logical concatenations	Recognition of recursive activities	Handle imperfect low-levels
Statistical	Oliver et al. '02	limited (2-levels)				✓
	Shi et al. '04	limited (2-levels)	one relation: 'before'			✓
	Damen and Hogg '09	limited (2-levels)				✓
Syntactic	Ivanov and Bobick '00	unlimited			✓	✓
	Joo and Chellappa '06	unlimited		conjunction only	✓	✓
	Pinhanez and Bobick '98	limited (redundant nodes required)	network form only	network form only		compensates 1 error
Description-based	Intille and Bobick '99	unlimited	two relations: 'before' and 'around'			✓
	Siskind '01	unlimited	a sub-event participates only once	✓		
	Nevatia et al. '03	limited (3-levels)	✓	✓		
	Vu et al. '03	unlimited	✓	conjunctions only		
	Ghanem et al. '04	unlimited	time intervals of an activity do not overlap	✓		
	Ryoo and Aggarwal '09a	unlimited	✓	✓	✓	✓
	Gupta et al. '09	limited (2-levels)	✓	network form only		✓

The column "levels of hierarchy" describes the possible levels of the activity hierarchy. "Complex temporal relations" suggests that the approach is able to represent and recognize activities with a complex temporal structure. Similarly, "complex logical concatenations" shows whether the system is able to represent activities with complex logical concatenations.

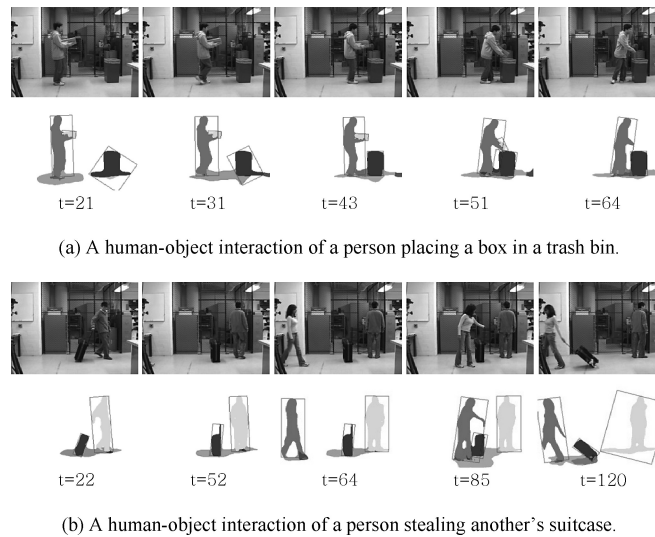


Fig. 13. Example human-object interactions that Ryoo and Aggarwal [2007] have recognized (©2007 IEEE).

analyzing interplays between humans and objects. Next, we compare various recognition approaches for group activities.

4.1. Recognition of Interactions between Humans and Objects

In order to recognize interactions between humans and objects, an integration of multiple components is required. The identification of objects and motion involved in an activity, as well as analysis of their interplays, is essential for the reliable recognition of human activities involving humans and objects. While we provide an overview of general human-object recognition approaches, we particularly focus on the approaches that analyze interplays among object recognition, motion estimation, and activity-level analysis toward robust recognition of human-object interactions. Figure 13 shows an example of human-object interactions.

The most typical human-object interaction recognition approaches are those that ignore interplays between object recognition and motion estimation. In those works, objects are generally recognized first, and activities involving them are recognized by analyzing objects' motion. They have made object recognition and motion estimation independent or made it so that motion estimation is strictly dependent on object recognition. Most of the previous recognition approaches fall into this category, including the approaches that we discussed in previous sections [Siskind 2001; Vu et al. 2003; Nevatia et al. 2003; Shi et al. 2004; Yu and Aggarwal 2006; Damen and Hogg 2009].

On the other hand, several researchers have studied relationships and dependencies among objects, motion, and human activities to improve object recognitions as well as activity recognition [Moore et al. 1999; Gupta and Davis 2007; Ryoo and Aggarwal 2007]. In principle, these components are highly dependent on each other: objects have their own roles, suggesting that the way humans interact with an object depends on the identity of the object. For example, an object, the water bottle, is expected to be involved in a particular type of interaction: *drinking*. Therefore, the motion related to the water bottle must be different from that of spray bottle, even though they appear similar. Several researchers have designed a probabilistic model describing mutual information between objects and the inherent human motions with the objects. The results suggest that the recognition of objects can benefit activity recognition, while

activity recognition helps the classification of objects; we discuss these approaches one by one.

Moore et al. [1999] constructed a system that compensates for the failures of object classification with the recognition results of simple actions. Most of the time, this system performs the object recognition first, and then estimates human activities with objects by depending on the object recognition results, as most other researchers have done. However, when an object component fails to make a concrete decision, their systems uses the objects' action information to compensate for the object recognition. In order to recognize actions, positions of hands and their tracking results are used. HMMs are applied to characterize actions based on the tracking results. Finally, object-based evidence is integrated with action-based evidence by using a Bayesian network to decide the final class of object, making the system recover from the failure of object recognition. Moore et al. have tested their system with various objects in office, kitchen, and automobile environments (such as books, phones, bowls, cups, and steering wheels). They focused on the recognition of simple activities of a single person.

Peursum et al. [2005] proposed a Bayesian framework for the better labeling of objects based on activity context. Similar to Moore et al. [1999], they focused on the fact that humans interact with objects in many different ways, depending on the function of the objects. They pointed out that appearance (i.e., shape) cues of objects are unreliable due to scale and viewpoint variations, and presented an object recognition based solely on the activity information. The system calculates an *interaction signature* per object, which essentially is a concatenation of activity recognition results involving the object. HMMs are adopted for the action recognition: a 3-D pose skeleton of a person as well as the relative locations of objects are analyzed to recognize activities, where each object candidate is computed by region segmentation based on color. Peursum et al. have recognized objects such as floor, chair, and keyboard, by recognizing printing-related activities.

Gupta and Davis [2007] proposed a probabilistic model integrating an objects' appearance, human motion with objects, and reactions of objects. Similar to Moore et al. [1999], a Bayesian network is constructed to combine cues. Two types of motion in which humans interact with objects, that is, "reach motion" and "manipulation motion," are estimated using trajectories as well as HMMs. The reactions of objects, that is, the effect of human activity in relation to their interaction with objects, such as "a light going on after pressing the switch," are considered for classification as well. The Bayesian network integrates all of this information, and makes a final decision to recognize objects and human activities. Human-object interactions involving cups, spray bottles, phones, and flashlights have been recognized in their experiments.

Ryoo and Aggarwal [2007] designed and implemented a recognition system for high-level human-object interactions such as "a person stealing another's suitcase." Similar to the above-mentioned approaches [Moore et al. 1999; Gupta and Davis 2007], the Ryoo and Aggarwal object recognition and motion estimation components were constructed to help each other. Furthermore, their system was designed to compensate for object recognition failures or motion estimation failures using high-level activity recognition results probabilistically. That is, the object recognition and motion estimation components not only help each other, but also get feedback from the high-level activity recognition results for improved recognition. For example, by observing a person pulling an object in an airport environment, their system was able to deduce that it is the activity of "a person carrying a suitcase" and provide feedback that the object in the scene is a suitcase. Via experiments, they have shown that the feedback generated by the high-level activity recognition may benefit object recognition, motion estimation, and low-level tracking of objects.



Fig. 14. Example group activities from the Ryoo and Aggarwal [2008]’s dataset. From left to right, snapshots of group activities: *group carrying*, *group stealing* in an office, *group stealing* in a shop, *group fighting*, and *group arresting*. *Group stealing* is a situation where a thief takes an object while the other thieves distract its owners.

4.2. Recognition of Group Activities

Group activities are those whose actors are one or more conceptual groups. “A group of soldiers marching” and “a group of persons carrying a large object” are examples of simple group activities. In order to recognize group activities, the analysis of activities of individuals as well as their overall relations becomes essential. In this section, we discuss the recognition approaches to group activities, while focusing on the types of activities that they recognized. There are various types of group activities, and most of the work specializes in recognizing a particular type among them. Figure 14 illustrates example snapshots of various group activities.

First of all, researchers have focused on the recognition of group activities where each group member has its own role, different from the others [Cupillard et al. 2002; Gong and Xiang 2003; Lv et al. 2004; Zhang et al. 2006; Dai et al. 2008]. The goal of these approaches is to recognize an activity of a single group with a limited number of members who exhibit nonuniform behaviors. The group activity, *presentation* with a fixed number of participants in a meeting room is an example of this type: the presenter is talking while the other members are taking notes, asking questions, and/or listening. For this type of group activity, the system must recognize the activities of each individual member and then analyze their structures. By nature, most of these approaches are hierarchical, since there exist at least two-levels of activities: activities of the group and activities of individual persons. Statistical hierarchical approaches that use state models we discussed in Section 3.1 are especially popular. Essentially, this type of group activity is equivalent to multiagent interactions recognized by [Intille and Bobick 1999; Ivanov and Bobick 2000; Vu et al. 2003; Nevatia et al. 2003; Joo and Chellappa 2006; Ryoo and Aggarwal 2007].

Cupillard et al. [2002] recognized a group activity using a finite state machine, which is equivalent to a fully observable Markov model. They used multiple cameras, and were able to recognize the activity that “a group is fighting” which is essentially intragroup fighting of a group composed of two members. Similarly, as in Section 3.1, Gong and Xiang [2003] used variations of dynamic Bayesian networks to recognize group activities. With their system, they have recognized “a group of trucks loading (or unloading) baggage on an airplane” which is a group activity of a fixed number of trucks and an airplane. Zhang et al. [2006] recognized a group activity occurring in a meeting room using DBNs, similar to Gong and Xiang [2003]. Sequentially organized group activities including monologues, discussion, presentation, and note-taking have been successfully recognized. Similarly, Dai et al. [2008] recognized break, presentation, and discussion using DBNs with hierarchical structures.

The second type of group activity is that characterized by the overall motion of an entire group. A group of people parading or marching is a typical example of this type. In contrast to the first type of group activity where the individual activities of specific members are important, the analysis of overall motion and formation changes of the entire group is important for the second type of group activity. By their nature, single-layered approaches are appropriate for recognition, since all the motions of group

members must be considered simultaneously [Vaswani et al. 2003; Khan and Shah 2005].

Vaswani et al. [2003] have recognized group activities of people interacting with an airplane. Their approach corresponds to the category of single-layered exemplar-based sequential approaches presented in Section 2.1.3. They represented a group activity as a shape change over time frames. At each frame, they extracted k point objects and constructed a polygon by treating the extracted points as corners. The points are tracked, and the dynamics of shape changes following the statistical shape theory are maintained. Their system was able to distinguish normal and abnormal activities by comparing the activity shape extracted from an input with a maintained model in a tangent space. Similarly, Khan and Shah [2005] recognized a group of people parading by analyzing the overall motion of group members. They use a single-layered space-time approach using trajectory features, discussed in Section 2.1.2. The authors extracted the trajectory of each group member, and analyzed their activities by fitting a 3-D polygon to check the rigidity formation of the group.

Finally, Ryoo and Aggarwal [2008] developed a general representation and recognition methodology that is able to handle various types of group activities. Their approach is a description-based approach (Section 3.3), and various classes of group activities including group actions (e.g., marching), group-group interactions (e.g., group stealing), group-persons interactions (e.g., march by signal), and intra-group interactions (e.g., intragroup fighting) have been represented and recognized with this system. They took advantage of the universal (\forall) and existential (\exists) quantifiers to describe subevents (usually activities of individuals) that need to be performed by any one member of a group or by all members of a group. By attaching the universal and existential quantifiers to the participating group members, their system was able to represent most group activities that previous researchers recognized. The first class of activities that we discussed above is represented by attaching an existential quantifier to each actor of the group activity. The second class of activities is represented by applying the universal quantifier and by posing spatial constraints to the group. In addition, high-level group activities with complex structures that the previous methods had difficulty representing and recognizing, such as “a thief stealing an object while other thieves distract the owners” or “policemen arresting a group of criminals,” have been successfully represented and recognized with their system.

5. DATASETS AND REAL-TIME APPLICATIONS

In this section, we discuss public datasets available for the performance evaluation of the approaches, and review real-time human activity recognition systems.

5.1. Datasets

Public datasets provide a common criterion to measure and compare the accuracy of the proposed approaches. Therefore, a construction of a dataset containing videos of human activities plays a vital role in the advancement in the research of human activity recognition. In this section, we describe the existing human activity datasets that are currently available and discuss the characteristics of the datasets. We also compare the performance of the systems tested on an identical dataset.

Existing datasets that have been made publicly available can be categorized into three groups, as follows. The first type of dataset includes the KTH dataset [Schuldt et al. 2004] and the Weizmann dataset [Blank et al. 2005], which are designed to test general-purpose action recognition systems academically. They contain videos of different participants performing simple actions such as walking and waving, which are taken by the authors in a controlled environment. The second type is a class of more application-oriented datasets obtained from realistic environments (e.g., an airport).

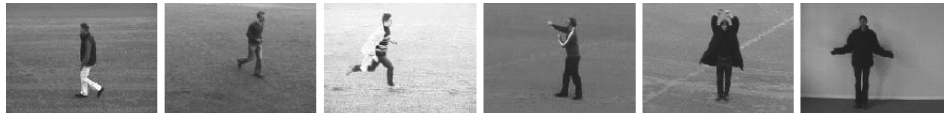


Fig. 15. Example snapshots from the KTH dataset [Schuldt et al. 2004] (©2004 IEEE).

The PETS datasets containing activities like baggage stealing and fighting are typical examples of this type, targeted for surveillance applications. In addition, datasets collected from real video media such as TV broadcasts and movies have recently been constructed and presented.

5.1.1. Action Recognition Datasets. A large number of researchers have tested their system on the KTH dataset [Schuldt et al. 2004; Dollar et al. 2005; Jiang et al. 2006; Niebles et al. 2006; Yeo et al. 2006; Ke et al. 2007; Kim et al. 2007; Jhuang et al. 2007; Savarese et al. 2008; Laptev et al. 2008; Liu and Shah 2008; Bregonzio et al. 2009; Rapantzikos et al. 2009; Ryoo and Aggarwal 2009b] and the Weizmann dataset [Blank et al. 2005; Niebles et al. 2006; Scovanner et al. 2007; Rodriguez et al. 2008; Bregonzio et al. 2009]. The KTH dataset is a large-scale dataset that contains 2391 videos of six actions performed by 25 subjects. Walking, jogging, running, boxing, hand-waving, and hand-clapping are the six actions contained in the dataset (Figure 15). Videos are taken at slightly different scales with various backgrounds in indoor and outdoor environments (but mostly uniform backgrounds). Each video contains repeated executions of a single action in a resolution of 160×120 , 25fps. Similarly, the Weizmann dataset consists of 10 action categories with 9 people, resulting in 90 videos. In the Weizmann dataset, a static and simple background is used throughout the videos. Simple human actions of running, walking, jumping-jack, jumping forward on two legs, jumping in place on two legs, galloping sideways, waving one hand, waving two hands, and bending are performed by the actors. The resolution of the videos is 180×144 , 25fps. Both dataset are composed of relatively simple action-level activities, and only one participant appears in the scene.

What we must note is that these datasets are designed to verify the classification ability of the systems on simple actions. Each video of the datasets contains executions of only one simple action, performed by a single actor. That is, entire motion-related features extracted from each video correspond to a single action, and the goal is to identify the label of the video while knowing that the video belongs to one of a limited number of known action classes. Further, all actions in both datasets except for the “bend” action of the Weizmann dataset are periodic actions (e.g., walking), making the videos suitable for action-level classification systems.

Due to their nature, testing the methodologies utilizing spatio-temporal local features (Section 2.1.3) is popular. As discussed in previous sections, such approaches do not require background subtraction and are robust to scale changes. Further, they are particularly suitable for recognition of periodic actions, since spatio-temporal features will be extracted repeatedly from the periodic actions. Figure 16 compares the classification accuracies of the systems. The X axis corresponds to the time of the publication, while the Y axis shows the classification performance of the systems. Most of the systems tested on the Weizmann dataset have obtained successful results, mainly due to the simplicity of the dataset. Particularly, [Blank et al. 2005; Niebles et al. 2006; Rodriguez et al. 2008; Bregonzio et al. 2009] have obtained more than 0.95 classification accuracy.

5.1.2. Surveillance Datasets. On the other hand, the PETS datasets (i.e., the datasets provided at the PETS workshops on 2004, 2006, 2007) and other similar datasets,

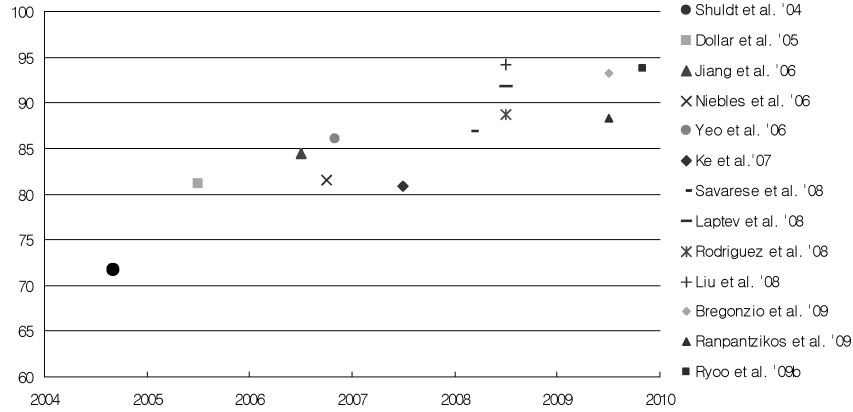


Fig. 16. The classification accuracies of various systems tested on the KTH dataset. The X axis corresponds to the time of publication. Only the results of the systems with the common experimental settings, the original 16 training-9 testing setting [Schuldt et al. 2004; Laptev et al. 2008] or the leave-one actor-out cross-validation setting (i.e., 25-fold cross-validation), are shown. Results of other systems using nontrivial settings, such as Wong et al.'s [2007] system tested with 100-fold cross-validations, Kim et al.'s [2007] system using manually-labeled bounding boxes, and Jhuang et al.'s [2007] system tested with the subsets, are not presented. The results of Niebles et al. [2006] and Savarese et al. [2008] are from the system trained with unlabeled data (i.e., unsupervised learning). [Dollar et al. 2005; Niebles et al. 2006; Savarese et al. 2008; Ryoo and Aggarwal 2009b] used the same *cuboid* features, and most of the other systems developed their own features to recognize the actions.

including the i-Lids dataset, are composed of realistic videos in uncontrolled environments such as crowded subway stations and airports. The camera viewpoints are similar to those of typical CCTVs, and even multiple camera viewpoints are provided in some of the datasets. Cameras are fixed, implying that the backgrounds are static and the scales for persons are mostly constant. Multiple persons and objects appear in the scene simultaneously, and occlusion among them occurs frequently. The goal of these surveillance videos is to test the ability of recognition systems to analyze realistic and specific (e.g., baggage abandonment and baggage theft) activities, which are of practical interest. These datasets are closely related to real-time applications, which we will discuss in the following section.

The PETS 2004 dataset (also known as the CAVIAR dataset) contains 6 categories of activities where each category is composed of one or more actions: walking, browsing, resting-slumping-fainting, leaving bags behind, people meeting, walking together, and splitting up, and fighting. Each class has 3 to 6 videos, providing a total of 28 videos. Videos have 384*288 spatial resolution, 25fps. Background images are provided, and the videos were taken in a shop environment. It is a single viewpoint dataset (i.e., only one camera was installed).

In the PETS 2006 dataset, 7 long video sequences were provided from every one of the 4 different viewpoints. The PETS 2006 dataset focused on the baggage abandonment problem: each sequence contains an event in which a bag is abandoned in a train station. Either one or two persons participated in the activity, and several other pedestrians were present in the scene. All four cameras had a high spatial resolution of 768*576 with 25fps. The PETS 2007 have a similar setup to the PETS 2006. The videos were taken in an airport hall with four cameras, providing eight executions of human activities. They focused on human-baggage interactions: two sequences of general loitering, four sequences of baggage theft, and two sequences of baggage abandonment similar to those of the PETS 2006. The resolution of the images are identical to the PETS 2006. Actors, objects, and pedestrians were severely occluded in the videos.

Similar to the PETS datasets, the recently introduced i-Lids dataset focuses on the baggage abandonment problem. Videos were taken from a single viewpoint in a London subway station, in a crowded environment. The videos not only contain persons and objects, but also a moving subway train in which people get out and in. Humans and objects were severely occluded by each other, and pedestrians were easily occluded by pillars in the station. Three videos were provided for training and validation purposes, and a lengthy video containing six baggage abandonment activities were given for testing. Videos have a resolution of 720*576 with 25fps. There was a real-time abandoned baggage-detection competition at the AVSS 2007 conference with the dataset [Venetianer et al. 2007; Bhargava et al. 2007].

Several works testing their systems on these surveillance datasets have been published [Lv et al. 2004; Kitani et al. 2005; Ribeiro et al. 2007]. In contrast to the datasets mentioned in Section 5.1.1, the surveillance datasets are motivated by the practical needs for constructing surveillance systems for public safety. They provide more realistic videos in practical environments. However, they lack generality in a certain aspect, since they are highly oriented toward surveillance applications. That is, they are focused on particular types of activities.

5.1.3. Movie Datasets. Movie datasets are composed of challenging videos obtained from real movie scenes (or from TV broadcasts). Unlike the datasets in Section 5.1.1, they are not taken in a controlled environment. They are different from the datasets in Section 5.1.2 as well, since camera viewpoints move frequently, and background information is seldom provided. Most of the movie datasets [Ke et al. 2007; Laptev and Perez 2007; Laptev et al. 2008; Rodriguez et al. 2008] focused on relatively simple actions such as kissing and hitting. Even though the actions are simple, each video of an action exhibits person-dependent, viewpoint-dependent, and situation-dependent variations. Thus, the major challenge is in handling those variations rather than recognizing complex structured activities; space-time feature-based approaches have been applied to solve the problem.

5.2. Real-Time Applications

In this section, we review several computer vision systems designed to recognize activities in real-time. Even though the approaches that we have discussed in the previous sections have shown results on various types of human activities, most of the proposed algorithms are far from being real-time. In order for an activity recognition methodology to be applicable for real-world applications, including surveillance systems, human-computer interfaces, intelligent robots, and autonomous vehicles, this computational gap must be overcome.

Recently, various real-time human activity recognition systems have been proposed, and we review some of them here. The general idea of most of them is to increase the efficiency of the algorithms by simplifying them. These systems focus on the utilization of simple but effective features instead of performing a detailed analysis. Lv et al. [2004] used a traditional Bayesian posterior probability calculation for the recognition of actions. In order to detect activities reliably without spending too much on computational costs, their approach searches for an optimal set of features from a large number of features. They have proposed a dynamic programming algorithm to find a set of features that maximizes the detection accuracy on the training data. The PETS 2004 dataset has been used for the testing.

Yeo et al. [2006] focused on a frame-to-frame similarity measurement based on optical flow calculations. The key to their system is the fact that the modern video compression technology takes advantage of the optical flows to encode the videos. That is, optical flows are naturally embedded in the videos, and are easily extractable. The similarities

between frames are measured based on the optical flow distribution, and are aggregated to measure the similarities between two videos. Their approach can be seen as a sequential exemplar-based approach similar to Efros et al. [2003]; the KTH dataset has been applied to test their system.

The Li et al. [2008] approach is that of a space-time trajectory analysis. They used the principal component analysis (PCA) to compress the trajectories from a high-dimensional space to a low-dimensional one. Several learning algorithms were applied on the low-dimensional trajectories, and the Gaussian mixture model was adopted for the classification. The reduction in dimensionality provided the ability to process videos in real-time. Their system was also tested with the KTH dataset.

Notably, Rofouei et al. [2008] utilized the graphical processing units (GPUs) of computer systems to enable the real-time recognition of human activities. Instead of making the algorithms simpler, they focused on the fact that modern hardware is able to support computationally expensive processing. The state-of-the-art graphic cards are composed of GPUs with many cores, suggesting that they are able to compute repetitive computations in parallel. This implies that they are suitable for the parallel processing of many computer vision algorithms analyzing images and videos (e.g., a GPU-based SIFT feature extraction). Rofouei et al. [2008] designed a GPU-version algorithm of Dollar et al. [2005], which is 50 times faster than the CPU implementation of the algorithm, without sacrificing performance. They illustrated the potential for GPUs (or multi-core CPUs) to greatly improve the speed of computer vision systems, enabling the real-time implementation of existing activity recognition algorithms.

6. CONCLUSION

Computer recognition of human activities is an important area of research in computer vision with applications in many diverse fields. The application to surveillance is natural in today's environment, where the tracking and monitoring people is becoming an integral part of everyday life. Other applications include human-computer interaction, biometrics based on gait or face, and hand and face gesture recognition. We have provided an overview of the current approaches to human activity recognition. The approaches are diverse and yield a spectrum of results. The senior author of this article has been involved in the study of motion since the early 1970s [Aggarwal and Duda 1975] and in the area of human activity recognition since the early 1980s [Webb and Aggarwal 1982]. The impetus for the study of human motion and human activity was provided by Johansson's [1975] pioneering work in the early 1970s; human activity research came to the forefront in the early 1990s.

In this review we have summarized the methodologies previously explored for the recognition of human activity, and discussed the advantages and disadvantages of those approaches. An approach-based taxonomy was designed and applied to categorize previous works. We have discussed nonhierarchical approaches developed for the recognition of gestures and actions as well as hierarchical approaches for the analysis of high-level interactions between multiple humans and objects. Nonhierarchical approaches were again divided into space-time approaches and sequential ones, and we discussed the similarities and differences of the two thoroughly. Previous publications on hierarchical methodologies including statistical, syntactic, and description-based approaches were compared as well.

In 1999, human activity recognition was in its infancy, as Aggarwal and Cai [1999] pointed out. A significant amount of progress on human activity recognition has been made in the past ten years, but it is still far from being an off-the-shelf technology. We are at a stage where experimental systems are deployed at airports and other public places. It is likely that more and more such systems will be deployed. There is a strong interaction between the surveillance authorities and computer vision researchers. For

example, Mubarak Shah of the University of Central Florida and the Orlando Police Department are joining forces to develop a system to monitor downtown Orlando: <http://server.cs.ucf.edu/~vision/projects/Knight/Knight.html>.

Further, today's environment for human activity recognition is significantly different from the scenario at the end of the last decade. The cameras were mostly fixed and without pan-tilt-zoom adjustments. Today's cameras may be mounted on several types of moving platforms, ranging from a moving car or a truck to an unmanned aerial vehicle (UAV). A global positioning system may be attached to the camera system to pin-point its location. The recognition of activity from a moving platform poses many more challenges. Noise, tracking, and segmentation issues arising out of stabilization of video add to the difficulty of the problem for the recognition of activities. Tracking is a difficult problem, although animals and human do it almost effortlessly. If the tracking algorithm does not extract the object of the focus of attention, recognizing the activity being performed becomes enormously more difficult. Designing an activity recognition system that is able to compensate for low-level failures in such environments (i.e., moving platforms) is an extremely challenging task.

The future direction of research is obviously encouraged and dictated by applications. The pressing applications are the surveillance and monitoring of public facilities like train stations, underground subways or airports, monitoring patients in a hospital environment or other health care facilities, monitoring activities in the context of UAV surveillance, and other similar applications. All of these applications are trying to understand the activities of an individual or the activities of a crowd as a whole and as subgroups. These problems will occupy us for a number of years and several generations of graduate students.

As pointed out previously, segmenting and tracking multiple persons in videos is harder than it appears. This difficulty is partly due to poor lighting, crowded environments, noisy images, and camera movements. For example, lighting in subways is almost universally poor. Further, it is difficult to segment individuals or their body parts when occlusion is present in the scene. Alternative approaches to segmenting body parts based on analyzing 3-D XYT volumes by extracting gross features are being developed. In particular, 3-D local patch features described in terms of histogram of gradient (HOG) and/or histogram of optical flow (HOOF), such as cuboids [Dollar et al. 2005] and 3-D SIFT [Scovanner et al. 2007], are gaining popularity. These approaches are motivated by the success of object recognition using 2-D local descriptors (e.g., SIFT [Lowe 1999]).

However, they involve long feature vectors obtained from a large 3-D XYT volume created by concatenating image frames, and are likely to have an impact on real-time analysis. The 3-D search space is much larger than its 2-D versions. Further, the existing local space-time features generally require a nontextured background for reliable recognition, such as the ones in the KTH and Weizmann datasets [Schuldt et al. 2004; Blank et al. 2005]. Also, a limited amount of work has been published on the 3-D feature-based approaches for analysis of complex human activities. What we need is an approach that exploits the easy computation of SIFT, HOG, and HOOF operators and avoids the difficulties of segmentation of body parts and/or combines the two approaches in a meaningful way.

One promising direction for enabling real-time implementation is the study of hardware support. Rofouei et al. [2008] have implemented a GPU-based version of the cuboid feature extractor, utilizing graphical processing units (GPUs) with tens of cores running thousands of threads. The GPU-version turned out to be 50 times faster than its CPU counterpart, while obtaining the same results. Modern CPUs and GPUs are composed of multiple cores, and the number of cores is likely to increase continually for the next few years, suggesting that computer vision researchers explore their utilization.

There are a number of other innovative approaches being explored. One such approach is exploiting the fact that images, high-dimensional signals, may possibly reside in low-dimensional manifolds. Several researchers are pursuing issues relating to characterizing the manifolds and exploring the relationships of the manifolds of different activities of the same person or the same activity of different persons [Veeraraghavan et al. 2006]. The temporal segmentation of activities and gestures is still a difficult issue. The inability to simultaneously register rigid and nonrigid parts of a face (parts of the human body in general) contributes to this difficulty. In certain activities, parts of the body move fairly rigidly whereas other parts undergo nonrigid motions, for example, the movement of the head/face. Shape deformations may be modeled as a linear combination of unknown shape bases [la Torre Frade et al. 2007], providing another approach to the recognition of facial expressions.

Hierarchical recognition approaches are being studied intensively, especially for the recognition of complex multiperson activities. In particular, description-based approaches are gaining an increasing amount of popularity due to their ability to represent and recognize human interactions with complex spatio-temporal structures. Activities with structured scenarios (e.g., most surveillance scenarios) require hierarchical approaches, and the description based approaches show the potential for making a reliable decision probabilistically. In the near future, hierarchical approaches together with strong action-level detectors such as the ones mentioned previously, will be explored for reliable recognition of complex activities. As we have shown in previous sections, hierarchical approaches have their advantages in the recognition of high-level activities performed by multiple persons, hence they must be explored further in the future to support the demands from surveillance systems and other applications.

The preceding areas of research, the space-time feature-based approaches, manifold learning, rigid/nonrigid motion analysis, and hierarchical approaches briefly mentioned are but a small glimpse into the large number of methodologies being pursued today. Hopefully, a review in another ten years will document significant progress in human activity recognition, to the extent that off-the-shelf systems will be readily available.

REFERENCES

- AGGARWAL, J. K. AND CAI, Q. 1999. Human motion analysis: A review. *Comput. Vision Image Understand.* 73, 3, 428–440.
- AGGARWAL, J. K. AND DUDA, R. O. 1975. Computer analysis of moving polygonal images. *IEEE Trans. Comput.* 24, 10, 966–976.
- ALLEN, J. F. 1983. Maintaining knowledge about temporal intervals. *Comm. ACM* 26, 11, 832–843.
- ALLEN, J. F. AND FERGUSON, G. 1994. Actions and events in interval temporal logic. *J. Logic Comput.* 4, 5, 531–579.
- BHARGAVA, M., CHEN, C.-C., RYOO, M. S., AND AGGARWAL, J. K. 2007. Detection of abandoned objects in crowded environments. In *Proceedings of the IEEE Conference on Advanced Video and Signal Based Surveillance (AVSS)*. IEEE, Los Alamitos, CA.
- BLANK, M., GORELICK, L., SHECHTMAN, E., IRANI, M., AND BASRI, R. 2005. Actions as space-time shapes. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. IEEE, Los Alamitos, CA, 1395–1402.
- BOBICK, A. AND DAVIS, J. 2001. The recognition of human movement using temporal templates. *IEEE Trans. Patt. Anal. Mach. Intell.* 23, 3, 257–267.
- BOBICK, A. F. AND WILSON, A. D. 1997. A state-based approach to the representation and recognition of gesture. *IEEE Trans. Patt. Anal. Mach. Intell.* 19, 12, 1325–1337.
- BREGONZIO, M., GONG, S., AND XIANG, T. 2009. Recognising action as clouds of space-time interest points. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE Los Alamitos, CA.

- CAMPBELL, L. W. AND BOBICK, A. F. 1995. Recognition of human body motion using phase space constraints. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. IEEE Los Alamitos, CA, 624–630.
- CEDRAS, C. AND SHAH, M. 1995. A motion-based recognition: A survey. *Image Vision Comput.* 13, 2, 129–155.
- CHOMAT, O. AND CROWLEY, J. 1999. Probabilistic recognition of activity using local appearance. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. Vol. 2, IEEE Los Alamitos, CA.
- CUPILLARD, F., BREMOND, F., AND THONNAT, M. 2002. Group behavior recognition with multiple cameras. In *Proceedings of the IEEE Workshop on Applications of Computer Vision (WACV)*. IEEE, Los Alamitos, CA, 177–183.
- DAI, P., DI, H., DONG, L., TAO, L., AND XU, G. 2008. Group interaction analysis in dynamic context. *IEEE Trans. Syst. Man Cybern. Part B* 38, 1, 275–282.
- DAMEN, D. AND HOGG, D. 2009. Recognizing linked events: Searching the space of feasible explanations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, Los Alamitos, CA.
- DARRELL, T. AND PENTLAND, A. 1993. Space-time gestures. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, Los Alamitos, CA, 335–340.
- DOLLAR, P., RABAUD, V., COTTRELL, G., AND BELONGIE, S. 2005. Behavior recognition via sparse spatio-temporal features. In *Proceedings of the 2nd Joint IEEE International Workshop on Visual Surveillance and Performance Evaluation of Tracking and Surveillance (VS-PETS)*. IEEE, Los Alamitos, CA, 65–72.
- EFROS, A., BERG, A., MORI, G., AND MALIK, J. 2003. Recognizing action at a distance. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. Vol. 2, IEEE, Los Alamitos, CA, 726–733.
- GAVRILA, D. AND DAVIS, L. 1995. Towards 3-D model-based tracking and recognition of human movement. In *Proceedings of the International Workshop on Face and Gesture Recognition*. 272–277.
- GAVRILA, D. M. 1999. The visual analysis of human movement: A survey. *Comput. Vision Image Understand.* 73, 1, 82–98.
- GHANEM, N., DEMENTHON, D., DOERMANN, D., AND DAVIS, L. 2004. Representation and recognition of events in surveillance video using Petri nets. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshop (CVPRW)*. IEEE, Los Alamitos, CA.
- GONG, S. AND XIANG, T. 2003. Recognition of group activities using dynamic probabilistic networks. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. IEEE, Los Alamitos, CA, 742.
- GUPTA, A. AND DAVIS, L. S. 2007. Objects in action: An approach for combining action understanding and object perception. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, Los Alamitos, CA.
- GUPTA, A., SRINIVASAN, P., SHI, J., AND DAVIS, L. S. 2009. Understanding videos, constructing plots. Learning a visually grounded storyline model from annotated videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, Los Alamitos, CA.
- HAKHEEM, A., SHEIKH, Y., AND SHAH, M. 2004. CASEE: A hierarchical event representation for the analysis of videos. In *Proceedings of the 20th National Conference on Artificial Intelligence (AAAI)*. 263–268.
- HARRIS, C. AND STEPHENS, M. 1988. A combined corner and edge detector. In *Proceedings of the Alvey Vision Conference*. 147–152.
- HONGENG, S., NEVATIA, R., AND BREMOND, F. 2004. Video-based event recognition: Activity representation and probabilistic recognition methods. *Comput. Vision Image Understand.* 96, 2, 129–162.
- INTILLE, S. S. AND BOBICK, A. F. 1999. A framework for recognizing multi-agent action from visual evidence. In *Proceedings of the AAAI Conference on Innovative Applications of Artificial Intelligence*. AAAI/IAAI. 518–525.
- IVANOV, Y. A. AND BOBICK, A. F. 2000. Recognition of visual activities and interactions by stochastic parsing. *IEEE Trans. Patt. Anal. Mach. Intell.* 22, 8, 852–872.
- JHUANG, H., SERRE, T., WOLF, L., AND POGGIO, T. 2007. A biologically inspired system for action recognition. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. IEEE, Los Alamitos, CA.
- JIANG, H., DREW, M., AND LI, Z. 2006. Successive convex matching for action detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, Los Alamitos, CA.
- JOHANSSON, G. 1975. Visual motion perception. *Sci. Amer.* 232, 6, 76–88.
- JOO, S.-W. AND CHELLAPPA, R. 2006. Attribute grammar-based event recognition and anomaly detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshop (CVPRW)*. IEEE, Los Alamitos, CA, 107.

- KE, Y., SUKTHANKAR, R., AND HEBERT, M. 2007. Spatio-temporal shape and flow correlation for action recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, Los Alamitos, CA.
- KHAN, S. M. AND SHAH, M. 2005. Detecting group activities using rigidity of formation. In *Proceedings of the ACM International Conference on Multimedia (ACM MM)*. ACM, New York, 403–406.
- KIM, T.-K., WONG, S.-F., AND CIPOLLA, R. 2007. Tensor canonical correlation analysis for action classification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, Los Alamitos, CA.
- KITANI, K. M., SATO, Y., AND SUGIMOTO, A. 2005. Deleted interpolation using a hierarchical Bayesian grammar network for recognizing human activity. In *Proceedings of the Second Joint IEEE International Workshop on Visual Surveillance and Performance Evaluation of Tracking and Surveillance (VS-PETS)*. IEEE, Los Alamitos, CA.
- KITANI, K. M., SATO, Y., AND SUGIMOTO, A. 2007. Recovering the basic structure of human activities from a video-based symbol string. In *Proceedings of the IEEE Workshop on Motion and Video Computing (WMVC)*. IEEE, Los Alamitos, CA.
- KRUGER, V., KRAGIC, D., UDE, A., AND GEIB, C. 2007. The meaning of action: A review on action recognition and mapping. *Advanced Robotics* 21, 13, 1473–1501.
- LA TORRE FRADE, F. D., CAMPOY, J., COHN, J., AND KANADE, T. 2007. Simultaneous registration and clustering for temporal segmentation. In *Proceedings of the International Conference on Computer Vision Theory and Applications*. 110–115.
- LAPTEV, I. AND LINDBERG, T. 2003. Space-time interest points. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. IEEE, Los Alamitos, CA, 432.
- LAPTEV, I., MARSZALEK, M., SCHMID, C., AND ROZENFELD, B. 2008. Learning realistic human actions from movies. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, Los Alamitos, CA.
- LAPTEV, I. AND PEREZ, P. 2007. Retrieving actions in movies. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. IEEE, Los Alamitos, CA.
- LI, Z., FU, Y., HUANG, T., AND YAN, S. 2008. Real-time human action recognition by luminance field trajectory analysis. In *Proceedings of the ACM International Conference on Multimedia (ACM MM)*. ACM, New York, 671–676.
- LIU, J., LUO, J., AND SHAH, M. 2009. Recognizing realistic actions from videos “in the wild”. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, Los Alamitos, CA.
- LIU, J. AND SHAH, M. 2008. Learning human actions via information maximization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, Los Alamitos, CA.
- LOWE, D. G. 1999. Object recognition from local scale-invariant features. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. IEEE, Los Alamitos, CA, 1150–1157.
- LUBLINERMAN, R., OZAY, N., ZARPALAS, D., AND CAMPS, O. 2006. Activity recognition from silhouettes using linear systems and model (in)validation techniques. In *Proceedings of the International Conference on Pattern Recognition (ICPR)*. 347–350.
- LV, F., KANG, J., NEVATIA, R., COHEN, I., AND MEDIONI, G. 2004. Automatic tracking and labeling of human activities in a video sequence. In *Proceedings of the IEEE International Workshop on Performance Evaluation of Tracking and Surveillance (PETS)*. IEEE, Los Alamitos, CA.
- LV, F. AND NEVATIA, R. 2007. Single view human action recognition using key pose matching and Viterbi path searching. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, Los Alamitos, CA.
- MINNEN, D., ESSA, I. A., AND STARNER, T. 2003. Expectation grammars: Leveraging high-level expectations for activity recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. Vol. 2, IEEE, Los Alamitos, CA, 626–632.
- MOORE, D. J. AND ESSA, I. A. 2002. Recognizing multitasked activities from video using stochastic context-free grammar. In *Proceedings of the AAAI Conference on Innovative Applications of Artificial Intelligence*. 770–776.
- MOORE, D. J., ESSA, I. A., AND HAYES, M. H. 1999. Exploiting human actions and object context for recognition tasks. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. Vol. 1, IEEE, Los Alamitos, CA, 80–86.
- NAM, Y., WOHN, K., AND LEE-KWANG, H. 1999. Modeling and recognition of hand gesture using colored Petri nets. *IEEE Trans. Syst. Man Cybern.* 29, 5, 514–521.
- NATARAJAN, P. AND NEVATIA, R. 2007. Coupled hidden semi-Markov models for activity recognition. In *Proceedings of the IEEE Workshop on Motion and Video Computing (WMVC)*. IEEE, Los Alamitos, CA.

- NEVATIA, R., HOBBS, J., AND BOLLES, B. 2004. An ontology for video event representation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshop (CVPRW)*. Vol. 7, IEEE, Los Alamitos, CA.
- NEVATIA, R., ZHAO, T., AND HONGENG, S. 2003. Hierarchical language-based representation of events in video streams. In *Proceedings of the IEEE Workshop on Event Mining*. IEEE, Los Alamitos, CA.
- NGUYEN, N. T., PHUNG, D. Q., VENKATESH, S., AND BUI, H. H. 2005. Learning and detecting activities from movement trajectories using the hierarchical hidden Markov models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. Vol. 2, IEEE, Los Alamitos, CA, 955–960.
- NIEBLES, J. C., WANG, H., AND FEI-FEI, L. 2006. Unsupervised learning of human action categories using spatial-temporal words. In *Proceedings of the British Machine Vision Conference (BMVC)*.
- NIEBLES, J. C., WANG, H., AND FEI-FEI, L. 2008. Unsupervised learning of human action categories using spatial-temporal words. *Int. J. Comput. Vision* 79, 3.
- NIYOGI, S. AND ADELSON, E. 1994. Analyzing and recognizing walking figures in XYT. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, Los Alamitos, CA, 469–474.
- OLIVER, N., HORVITZ, E., AND GARG, A. 2002. Layered representations for human activity recognition. In *Proceedings of the IEEE International Conference on Multimodal Interfaces (ICMI)*. IEEE, Los Alamitos, CA, 3–8.
- OLIVER, N. M., ROSARIO, B., AND PENTLAND, A. P. 2000. A Bayesian computer vision system for modeling human interactions. *IEEE Trans. Patt. Anal. Mach. Intell.* 22, 8, 831–843.
- PARK, S. AND AGGARWAL, J. K. 2004. A hierarchical Bayesian network for event recognition of human actions and interactions. *Multimedia Syst.* 10, 2, 164–179.
- PEURSUM, P., WEST, G., AND VENKATESH, S. 2005. Combining image regions and human activity for indirect object recognition in indoor wide-angle views. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. IEEE, Los Alamitos, CA.
- PINHANEZ, C. S. AND BOBICK, A. F. 1998. Human action detection using PNF propagation of temporal constraints. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, Los Alamitos, CA, 898.
- RAO, C. AND SHAH, M. 2001. View-invariance in action recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. Vol. 2, IEEE, Los Alamitos, CA, 316–322.
- RAPANTZIKOS, K., AVRITHIS, Y., AND KOLLIAS, S. 2009. Dense saliency-based spatiotemporal feature points for action recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, Los Alamitos, CA.
- RIBEIRO, P. C., MORENO, P., AND SANTOS-VICTOR, J. 2007. Detecting luggage related behaviors using a new temporal boost algorithm. In *Proceedings of the IEEE International Workshop on Performance Evaluation of Tracking and Surveillance (PETS)*. IEEE, Los Alamitos, CA.
- RODRIGUEZ, M. D., AHMED, J., AND SHAH, M. 2008. Action MACH: A spatio-temporal maximum average correlation height filter for action recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, Los Alamitos, CA.
- ROFOUEI, M., MOAZENI, M., AND SARRAFZADEH, M. 2008. Fast GPU-based space-time correlation for activity recognition in video sequences. In *Proceedings of the IEEE/ACM/IFIP Workshop on Embedded Systems for Real-Time Multimedia (ESTImedia)*. ACM, New York, 33–38.
- RYOO, M. S. AND AGGARWAL, J. K. 2009a. Semantic representation and recognition of continued and recursive human activities. *Int. J. Comput. Vision* 32, 1, 1–24.
- RYOO, M. S. AND AGGARWAL, J. K. 2009b. Spatio-temporal relationship match: Video structure comparison for recognition of complex human activities. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, IEEE, Los Alamitos, CA.
- RYOO, M. S. AND AGGARWAL, J. K. 2008. Recognition of high-level group activities based on activities of individual members. In *Proceedings of the IEEE Workshop on Motion and Video Computing (WMVC)*. IEEE, Los Alamitos, CA.
- RYOO, M. S. AND AGGARWAL, J. K. 2007. Hierarchical recognition of human activities interacting with objects. In *Proceedings of the 2nd International Workshop on Semantic Learning Applications in Multimedia (SLAM)*.
- RYOO, M. S. AND AGGARWAL, J. K. 2006a. Recognition of composite human activities through context-free grammar based representation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, Los Alamitos, CA, 1709–1718.
- RYOO, M. S. AND AGGARWAL, J. K. 2006b. Semantic understanding of continued and recursive human activities. In *Proceedings of the International Conference on Pattern Recognition (ICPR)*. 379–382.

- SAVARESE, S., DELPOZO, A., NIEBLES, J., AND FEI-FEI, L. 2008. Spatial-temporal correlatons for unsupervised action classification. In *Proceedings of the IEEE Workshop on Motion and Video Computing (WMVC)*. IEEE, Los Alamitos, CA.
- SCHULDT, C., LAPTEV, I., AND CAPUTO, B. 2004. Recognizing human actions: A local SVM approach. In *Proceedings of the International Conference on Pattern Recognition (ICPR)*. Vol. 3, 32–36.
- SCOVANNER, P., ALI, S., AND SHAH, M. 2007. A 3-dimensional sift descriptor and its application to action recognition. In *Proceedings of the ACM International Conference on Multimedia (ACM MM)*. ACM, New York, 357–360.
- SHECHTMAN, E. AND IRANI, M. 2005. Space-time behavior based correlation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. Vol. 1, IEEE, Los Alamitos, CA, 405–412.
- SHEIKH, Y., SHEIKH, M., AND SHAH, M. 2005. Exploring the space of a human action. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. Vol. 1, IEEE, Los Alamitos, CA, 144–149.
- SHI, Y., HUANG, Y., MINNEN, D., BOBICK, A. F., AND ESSA, I. A. 2004. Propagation networks for recognition of partially ordered sequential action. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. Vol. 2, IEEE, Los Alamitos, CA, 862–869.
- SISKIND, J. M. 2001. Grounding the lexical semantics of verbs in visual perception using force dynamics and event logic. *J. Artif. Intell. Res.* 15, 31–90.
- STARNER, T. AND PENTLAND, A. 1995. Real-time American Sign Language recognition from video using hidden Markov models. In *Proceedings of the International Symposium on Computer Vision*. 265.
- TRAN, S. D. AND DAVIS, L. S. 2008. Event modeling and recognition using Markov logic networks. In *Proceedings of European Conference on Computer Vision (ECCV)*. 610–623.
- TURAGA, P., CHELLAPPA, R., SUBRAHMANIAN, V. S., AND UDREA, O. 2008. Machine recognition of human activities: A survey. *IEEE Trans. Circuits Syst. Video Technol.* 18, 11 (Nov), 1473–1488.
- VASWANI, N., ROY CHOWDHURY, A., AND CHELLAPPA, R. 2003. Activity recognition using the dynamics of the configuration of interacting objects. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. Vol. 2, IEEE, Los Alamitos, CA.
- VEERARAGHAVAN, A., CHELLAPPA, R., AND ROY-CHOWDHURY, A. 2006. The function space of an activity. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. Vol. 1, IEEE, Los Alamitos, CA, 959–968.
- VENETIANER, P., ZHANG, Z., YIN, W., AND LIPTON, A. 2007. Stationary target detection using the ObjectVideo surveillance system. In *Proceedings of the IEEE Conference on Advanced Video and Signal Based Surveillance (AVSS)*. IEEE, Los Alamitos, CA, 242–247.
- VU, V.-T., BREMOND, F., AND THONNAT, M. 2003. Automatic video interpretation: A novel algorithm for temporal scenario recognition. In *Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI)*. 1295–1302.
- WEBB, J. A. AND AGGARWAL, J. K. 1982. Structure from motion of rigid and jointed objects. *Artif. Intell.* 19, 107–130.
- WONG, S.-F., KIM, T.-K., AND CIPOLLA, R. 2007. Learning motion categories using both semantic and structural information. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- YACOOB, Y. AND BLACK, M. 1998. Parameterized modeling and recognition of activities. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. IEEE, Los Alamitos, CA, 120–127.
- YAMATO, J., OHYA, J., AND ISHII, K. 1992. Recognizing human action in time-sequential images using hidden Markov models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, IEEE, Los Alamitos, CA, 379–385.
- YEO, C., AHAMMAD, P., RAMACHANDRAN, K., AND SHANKAR SASTRY, S. 2006. Compressed domain real-time action recognition. In *Proceedings of the IEEE Workshop on Multimedia Signal Processing*. IEEE, Los Alamitos, CA, 33–36.
- YILMAZ, A. AND SHAH, M. 2005a. Actions sketch: A novel action representation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. Vol. 1, IEEE, Los Alamitos, CA, 984–989.
- YILMAZ, A. AND SHAH, M. 2005b. Recognizing human actions in videos acquired by uncalibrated moving cameras (ICCV). IEEE, Los Alamitos, CA.
- YU, E. AND AGGARWAL, J. K. 2006. Detection of fence climbing from monocular video. In *Proceedings of the International Conference on Pattern Recognition (ICPR)*. 375–378.

- ZAIDI, A. K. 1999. On temporal logic programming using Petri nets. *IEEE Trans. Syst. Man Cybern.* 29, 3, 245–254.
- ZELNIK-MANOR, L. AND IRANI, M. 2001. Event-based analysis of video. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, Los Alamitos, CA.
- ZHANG, D., GATICA-PEREZ, D., BENGIO, S., AND McCOWAN, I. 2006. Modeling individual and group actions in meetings with layered hmms. *IEEE Trans. Multimedia* 8, 3, 509–520.

Received May 2008; revised March 2009; accepted September 2009