



Wed, 17 / 01 - 07

IST-FP6 -027657 / PACO-PLUS

Last saved by:

Confidential

Project no.: IST-FP6-IP-027657

Project full title: Perception, Action & Cognition through Learning of Object-Action Complexes

Project Acronym: PACO-PLUS

Deliverable no.: D3.2.1

Title of the deliverable: Representation of Actions

Contractual Date of Delivery to the CEC:	31st January 2007	
Actual Date of Delivery to the CEC:	31st January 2007	
Organisation name of lead contractor for this deliverable:	Aalborg University (AAU)	
Author(s):	Volker Krüger and Danica Kragic	
Participants(s):	AAU, KTH, UniKarl, BCCN, JSI	
Work package contributing to the deliverable:	WP3.2	
Nature:	R	
Version:	1.0	
Total number of pages:	100	
Start date of project:	1 st Feb. 2006	Duration: 48 month

**Projectco-funded by the European Commission within the Sixth Framework Programme (2002-2006)
Dissemination Level**

PU	Public	X
PP	Restricted to other programme participants (including the Commission Services)	
RE	Restricted to a group specified by the consortium (including the Commission Services)	
CO	Confidential, only for members of the consortium (including the Commission Services)	

Abstract:

This technical report summarizes the research efforts undertaken within the first report period of PACO-Plus. The key points of investigation included representation of action in the space of human joint settings and how to recognize actions. For action recognition, two subproblems have been investigated: 1) Given a set of example actions, how can novel actions be recognized that are of the same type as the example actions. 2) Given a set of example actions, how can actions be recognized that are a *composition* of novel actions of the example action type. The questions we wanted to answer were: 1) What are useful modeling strategies for action representation and recognition. 2) Is it possible to distinguish between very similar actions such as *pick up* or *push*. 3) What role does an object play that is involved in an action? The representations we have investigated are HMMs, PCA and spatio-temporal isomap. For our experiments, we have recorded three different action databases. The actions we have considered in our work are single arm actions of humans, partially including the hand and objects. The actions were performed by different individuals with several repetitions.

Keyword list: action recognition, action representation

Contents

1	Introduction	2
2	Decomposing Complex Actions	3
3	Action Representation with PCA and Isomap	3

1 Introduction

The field of action and activity representation for synthesis and recognition is relatively old, yet still immature. This area is presently subject to intense investigation which is also reflected by the large number of different ideas and approaches. The approaches depend on the goal of the researcher and applications for activity recognition are interesting for surveillance, medical studies and rehabilitation, robotics, video indexing and animation for film and games, see [7] for an extensive review. For example, in applications for scene interpretation the data is often represented statistically and is meant to distinguish “regular” from “irregular” activities.

In scene interpretation, the representations should be independent from the objects causing the activity and thus are usually not meant to distinguish explicitly, e.g. cars from humans. On the other hand, some surveillance applications focus explicitly on human activities and the interactions between humans. Here, one finds both, holistic approaches, that take into account the entire human body without considering particular body parts, and local approaches. Most holistic approaches attempt to identify “holistic” information such as gender, identity or simple actions like walking or running. Researchers using local approaches appear often to be interested in more subtle actions or attempt to model actions by looking for action primitives with which the complex actions can be modeled.

In PACO-Plus, we are particularly interested in understanding and representing action for both, recognition and synthesis. There is strong neurobiological evidence that human actions and activities are directly connected to the motor control of the human body [2; 8; 9]. When viewing other agents performing an action, the human visual system seems to relate the visual input to a sequence of motor primitives. The neurobiological representation for visually perceived, learned and recognized actions appears to be the same as the one used to drive the motor control of the body. These findings have gained considerable attention from the robotics community [1; 10]. Consequently, it is ongoing research in PACO-Plus to identify a set of primitives that allow a) representation of the visually perceived action and b) motor control for imitation. In addition, the biological findings give rise to the idea of interpreting and synthesizing activities through a hierarchy of simple actions primitives, actions and activities. The representations used to describe the primitives vary a lot across the literature and are subject to ongoing research [7].

Within PACO-Plus we have focused our attention on representing actions in the space of joint settings. This way, we hope to be able to recognize as well as synthesize actions.

A number of interesting scientific problems and questions were investigated, e.g.:

1. should an action be represented as a set of (key-) poses or as a sequence of poses with temporal dependency? This question is widely discussed in the computer vision community. A comparison was done in [11]. In [3; 4; 5] good recognition results have been achieved by modeling the temporal dependencies statistically using Hidden Markov Models.
2. How can a set of action primitive be found? This question has not yet been investigated.
3. Once action primitives are found, how can they be composed into more complex actions, and how can complex actions be decomposed into the right sequence of simple actions? In [5], we have developed a maximum a posteriori (MAP) classifier that allows to decompose complex actions into a sequence of simple actions.
4. How can equivalence classes for actions be defined where the classes contain semantically or syntactically related actions, such as a class of grasping actions which is independent of reaching distance and direction? This question is presently under investigation. The results are, however, immature and thus not presented, here.

This report collects the research results made within the first 12 months of PACO-Plus. The first of the following papers reviews the recent advances in vision-based human motion capture

and analysis [7]. This is an evaluation of more than 400 hundred related recent articles. A second article discusses the decomposition of complex actions into simple actions [5] (see Sec. 2). A third article presents an evaluation of using PCA, PCA with temporal dependencies and spatio-temporal Isomap for action representation [11] (see Sec 3).

The research in the report period up to month 18 will be focused on the detection of action primitives (point 2) and the detection of equivalence classes (point4).

This report on Action Representation has been published as a technical report at Aalborg University [6].

2 Decomposing Complex Actions

In order to recognize actions, we have mainly used Hidden Markov Models (HMM) (see also deliverable D8.2.1). As it is not possible to train an HMM for every possible action, one has the necessity to decouple complex actions into very simple ones for which an HMM can be learned. All complex actions are then interpreted as being a combination of these simple actions. Given a complex action, one needs then to recover the sequence of simple action. Based on these simple actions, the final recognition can be done. The classical way of recognition with HMMs is a maximum likelihood classification:

$$\max_i P(O_t|\lambda_i) \quad (1)$$

where λ_i specify different HMMs.

We have formulated the problem as a Bayesian one, where all HMMs are evaluated in parallel:

$$P(S_{t+1}, i_{t+1}|O_{0:t+1}) = P(O_{t+1}|S_{t+1}, i_{t+1})P(S_{t+1}, i_{t+1}|S_t, i_t)P(S_t, i_t|O_{0:t}) . \quad (2)$$

Here, i identifies the HMM, S the present state of the HMM and O the observations. In this formation, the maximum a posteriori probability $P(S_{t+1}, i_{t+1}|O_{0:t+1})$, when marginalized over the states S , degenerates to a dirac for the right identifier and thus recognizes the present action. After convergence, this process is restarted until the end of the observations of the complex action is reached. The evaluation of this approach was carried out on more than 10.000 test sequences with actions composed of up to 100 simple actions. The recognition results were even for noisy data above 86%.

3 Action Representation with PCA and Isomap

We have performed an initial study on recognition of four object manipulation actions: pick up, put down, rotate and push. Training and testing was performed with 20 people where the manipulated object was placed on two different heights and people performing the actions multiple times at three different orientations. This study is important to show how small variations in the training data affects the recognition rate. Most of the current systems that utilize robot imitation learning use a single person to train or teach tasks to the robot. Since the intention for the future is that robots will be able to learn from observing multiple persons that perform same actions, we believe that it is important to study how different methods scale with respect to this.

Currently, we have concentrated on evaluation of dimensionality reduction using linear and nonlinear techniques. We have shown how the number of sensors and different parameters affect the classification rate. We are aware of the fact that PCA and nearest neighbor classification are very simple techniques but we hope that our future work and work of other we evaluate more advanced techniques on the same data and compare it to the results obtained in this work. We also believe that this data and evaluation follows the current trend of designing different benchmarking criteria in robotics in general.

References

- [1] B. Dariush. Human Motion Analysis for Biomechanics and Biomedicine. *Machine Vision and Applications*, 14:202–205, 2003.
- [2] M. Giese and T. Poggio. Neural Mechanisms for the Recognition of Biological Movements. *Nature Reviews*, 4:179–192, 2003.
- [3] V. Krueger. Recognizing action primitives in complex actions using hidden markov models. In *Advances in Visual Computing*, pages 538–547, Second Int. Symp. on Visual Computing, Lake Tahoe, NV, USA, November 6-8, 2006.
- [4] V. Krueger. Recognition of action as a bayesian parameter estimation problem over time. In D. Metaxas, R. Klette, and B. Rosenhahn, editors, *Human Motion*. Springer, 2007. under review.
- [5] V. Krueger and D. Grest. Using hidden markov models for recognizing action primitives in complex actions. In *Scandinavian Conference on Image Analysis*,, 2007. submitted.
- [6] V. Krueger and D. Kragic. Representation of action. Technical Report CVMI-2007:3, Copenhagen Institute of Technology, Aalborg University, SE-100 44 Stockholm, 2007.
- [7] T. Moeslund, A. Hilton, and V. Krueger. A survey of advances in vision-based human motion capture and analysis. *Computer Vision and Image Understanding*, 104(2-3):90–127, 2006.
- [8] G. Rizzolatti, L. Fogassi, and V. Gallese. Parietal Cortex: from Sight to Action. *Current Opinion in Neurobiology*, 7:562–567, 1997.
- [9] G. Rizzolatti, L. Fogassi, and V. Gallese. Neurophysiological Mechanisms Underlying the Understanding and Imitation of Action. *Nature Reviews*, 2:661–670, Sept. 2001.
- [10] S. Schaal. Is Imitation Learning the Route to Humanoid Robots? *Trends in Cognitive Sciences*, 3(6):233–242, 1999.
- [11] I. S. Vicente and D. Kragic. Learning and recognition of object manipulation actions using linear and nonlinear dimensionality reduction. In *The 15th IEEE International Symposium on Robot and Human Interactive Communication, RO-MAN 2007*, 2007. submitted.

A Survey of Advances in Vision-Based Human Motion Capture and Analysis

Thomas B. Moeslund^a, Adrian Hilton^b, and Volker Krüger^c

^a*Laboratory of Computer Vision and Media Technology, Aalborg University, 9220 Aalborg, Denmark*

^b*Centre for Vision, Speech and Signal Processing, University of Surrey, Guildford GU2 7XH, UK*

^c*Aalborg Media Lab, Aalborg University Copenhagen, 2750 Ballerup, Denmark*

Abstract

This survey reviews advances in human motion capture and analysis from 2000 to 2006, following a previous survey of papers up to 2000 [247]. Human motion capture continues to be an increasingly active research area in computer vision with over 350 publications over this period. A number of significant research advances are identified together with novel methodologies for automatic initialization, tracking, pose estimation and movement recognition. Recent research has addressed reliable tracking and pose estimation in natural scenes. Progress has also been made towards automatic understanding of human actions and behavior. This survey reviews recent trends in video based human capture and analysis, as well as discussing open problems for future research to achieve automatic visual analysis of human movement.

Contents

1	Introduction	3
2	Model Initialization	5
2.1	Kinematic Structure Initialization	5
2.2	Shape Initialization	7
2.3	Appearance Initialization	8
2.4	Discussion of Advances in Model Initialization	8

3	Tracking	9
3.1	Background Subtraction	10
3.2	Motion-Based Segmentation	13
3.3	Appearance-Based Segmentation	13
3.4	Shape-Based Segmentation	15
3.5	Depth-Based Segmentation	17
3.6	Temporal Correspondences	18
3.7	Discussion of Advances Human Tracking	21
4	Pose Estimation	21
4.1	Model Free	22
4.2	Indirect Model Use	25
4.3	Direct Model Use	25
5	Discussion of Advances in Human Pose Estimation	30
6	Recognition	31
6.1	Action Hierarchies	31
6.2	Scene Interpretation	32
6.3	Holistic Recognition Approaches	33
6.4	Recognition Based on Body Parts	38
6.5	Action Primitives and Grammars	39
6.6	Discussion of Advances in Human Action Recognition	41
7	Conclusion	42
	References	44

1 Introduction

Automatic capture and analysis of human motion is a highly active research area due both to the number of potential applications and its inherent complexity. The research area contains a number of hard and often ill posed problems such as inferring the pose and motion of a highly articulated and self-occluding non-rigid 3D object from images. This complexity makes the research area challenging from a purely academic point of view. From an application perspective computer vision-based methods often provide the only non-invasive solution making it very attractive.

Applications can roughly be grouped under three titles: Surveillance, control, and analysis. *Surveillance applications* cover some of the more classical types of problems related to automatically monitoring and understanding locations where a large number of people pass through such as airports and subways. Applications could for example be: people counting or crowd flux, flow and congestion analysis. Newer types of surveillance applications - perhaps inspired by the increased awareness of security issues - are analysis of actions, activities and behaviors both for crowds and individuals. For example for queue and shopping behavior analysis, detection of abnormal activities, and person identification.

Control applications where the estimated motion or pose parameters are used to control something. This could be interfaces to games, e.g., as seen in EyeToy [3], Virtual Reality or more generally: Human Computer Interfaces. However, it could also be for the entertainment industry where the generation and control of personalized computer graphic models based on the captured appearance, shape, and motion are making the productions/products more believable.

Analysis applications such as automatic diagnostics of orthopedic patients or analysis and optimization of an athletes' performances. Newer applications are annotation of video as well as content-based retrieval and compression of video for compact data storage or efficient data transmission, e.g., for video conferences and indexing. Another branch of applications is within the car industry where much vision research is currently going on in applications such as automatic control of airbags, sleeping detection, pedestrian detection, lane following, etc.

The number of potential applications, the scientific complexity, the speed and price of current hardware, and the focus on security issues have intensified the effort within the computer vision community towards automatic capture and analysis of human motion. This is evident by looking at the number of publications, special sessions/issues at the major conference/journals as well as the number of workshops directly devoted to such topics. Furthermore, the major funding agencies have also focused on these research fields - especially the surveillance area.

The interest in this area has led to a large body of research which has been digested in a number of surveys, see table 1.

Even though some of these surveys are recent, it should be noted that the number of papers reviewed after 2000 is limited as seen in the table. In the relatively short period since 2000 a massive number of papers have been published advancing state of the art. This indicates increased activity in this research area compared to the number of papers identified in previous surveys.

Recent contributions have among other things addressed the limiting assumptions identified in previous approaches [247]. For example, many systems now address natural outdoor scenes and operate on long sequences of video containing multiple (occluded) people. This is possible, especially, due to more advanced segmentation algorithms. Other examples are model-based pose estimation where the introduction of learnt motion models and stochastic sampling methods have helped to achieved much faster and more precise results. Also within the recognition area there have been significant advances in both the representation and interpretation of actions and behavior.

Due to the significance of recent advances within this field we present the current survey. The survey is based on 351¹ recent papers (2000 - 2006) and structured using the functional taxonomy presented in the 2001 survey by Moeslund and

¹ Note that this number is different from the one listed in table 1 (331). The reason being that we also include papers from the last half of 2000 since this is where the previous survey [247] ends.

Year	Author	#Papers	Focus
1994	Aggarwal <i>et al.</i> [10]	69/0	Articulated and elastic nonrigid motion
1994	Cedras and Shah [54]	76/0	Motion extraction
1995	Aggarwal <i>et al.</i> [11]	104/0	Articulated and elastic nonrigid motion
1995	Ju [180]	91/0	Motion estimation and recognition
1997	Aggarwal and Cai [9]	51/0	Motion extraction
1997	Gavrila [113]	87/0	Motion estimation and recognition
2000	Moeslund and Granum [247]	155/0	Initialization, tracking, pose estimation and recognition
2001	Buxton [48]	88/6	Recognition
2001	Wang <i>et al.</i> [388]	164/14	Detection, tracking and recognition
2003	Hu <i>et al.</i> [156]	185/54	Surveillance
2004	Aggarwal and Park [12]	58/10	Recognition
2006	This survey	424/331	Initialization, tracking, pose estimation and recognition

Table 1

Previous surveys. Note that the *Year* is not necessary the publication year but rather the year of the most recent paper in a survey. The two numbers in the *#Papers* column state the total number of publications and the publications after 2000.

Granum [247]:

Initialization Ensuring that a system commences its operation with a correct interpretation of the current scene.

Tracking Segmenting and tracking humans in one or more frames.

Pose estimation Estimating the pose of a human in one or more frames.

Recognition Recognizing the identity of individuals as well as the actions, activities and behaviors performed by one or more humans in one or more frames.

The different papers are further divided into sub-taxonomies as seen in the table of contents. Inspired by [247] we also provide a visual overview of all the recent referenced papers, see table 2. For readers new to this field it is recommended to read [247] before preceding with the survey at hand. In fact this survey can be seen as a sequel to [247].

2 Model Initialization

Initialization of vision-based human motion capture and analysis often requires the definition of a humanoid model approximating the shape, appearance, kinematic structure and initial pose of the subject to be tracked. The majority of algorithms for 3D pose estimation continue to use a manually initialized generic model with limb lengths and shape which approximate the individual. To automate the initialization and improve the quality of tracking a limited number of authors have investigated the recovery of more accurate reconstructions of the subject from single or multiple view images.

Initialization captures prior knowledge of a specific person which can be used to constrain tracking and pose estimation. A priori knowledge used in human motion capture can be broken into a number of sources: kinematic structure; 3D shape; color appearance; pose; and motion type. In this section we review recent research which advances estimation of kinematic structure, 3D shape and appearance. Initialization of appearance is commonly an integral part of the tracking and pose estimation and is therefore also considered in conjunction with specific approaches in sections 3 and 4.

2.1 Kinematic Structure Initialization

The majority of vision-based tracking systems assume a priori a humanoid kinematic structure comprising a fixed number of joints with specified degrees-of-freedom. The kinematic initialization is then limited to estimation of limb lengths. Commercial marker-based motion capture systems typically require a fixed se-

quence of movements which isolate individual degrees-of-freedom. The known correspondence between markers and limbs together with reconstructed 3D marker trajectories during movement are then used to accurately estimate limb lengths. Hard constraints on left-right skeletal symmetry are commonly imposed during estimation. A number of approaches [26,28,278,361] have addressed initialization of body pose and limb lengths from manually identified joint locations in monocular images. Anthropometric constraints between ratios of limb lengths are imposed to allow estimation of the kinematic structure up to an unknown scale factor.

Direct estimation of the kinematic structure from sequences of a moving person has also been investigated. Krahnstover *et al.* [199,200] present a method for automatically initializing the upper-body kinematic structure based on motion segmentation of a sequence of monocular video images. Song *et al.* [350] introduce an unsupervised learning algorithm which uses point feature tracks from cluttered monocular video sequences to automatically construct triangulated models of whole-body kinematics. Learnt models are then used for tracking of walking motions from lateral views. These approaches provide more general solutions to the problem of initializing a kinematic model by deriving the structure directly from the scene.

Methods that derive the kinematic structure from 3D shape sequences reconstructed from multiple views have also been proposed. Cheung *et al.* [59] initialize the kinematic structure from the visual-hull of a person moving each joint independently. A full skeleton together with the shape of each body part is obtained by alignment of the segmented moving body parts with the visual-hull model in a fixed pose. Menier *et al.* [233] present an automated approach to 3D human pose estimation from the medial axis of the visual-hull. The kinematic structure is initialized independently at each frame enabling robust tracking. More general frameworks are presented in [44,65] to estimate the underlying skeletal spine structure from a temporal sequence of the 3D shape. The spine is estimated from the shape at each frame and common temporal structures identified to estimate the underlying structure. This work demonstrates reconstruction of approximate kinematic structures for babies, adults and animals.

Initializing the joint angle limits on the human kinematic structure is an important problem to constrain motion estimation to valid postures. Manual specification of joint angle limits has been common in many motion estimation algorithms using anthropometric data. This does not take into account the complex nature of joint limits and coupling between limits for different degrees-of-freedom. To overcome these limitations recent research has investigated learning models of joint limits and their correlations. Anthropometric models for the relationship between arm joint angles (shoulder, elbow, wrist) have been used to provide constraints in visual tracking and 3D upper-body pose estimation [248,253,262]. Recent research has investigated the modeling of joint limits from measurements of human motion captured using marker based systems [143,144] and from clinical data [252]. This is demonstrated to improve the performance of human pose estimation for complex

upper-body movement.

Increasingly, human motion capture sequences from commercial marker-based systems have been used to learn prior models of human kinematics and specific motions to provide constraints for subsequent tracking. Similarly motion capture databases [1,2,4] have recently been used to synthesize image sequences with known 3D pose correspondence to learn a priori the mapping from image to pose space for reconstruction.

2.2 Shape Initialization

A generic humanoid model is used in many video-based human motion estimation techniques to approximate a subject's shape. Representations have used either simple shape primitives (cylinders, cones, ellipsoids, super-quadratics) or a surface (polygonal mesh, sub-division surface) articulated using the kinematic skeleton [247]. A number of approaches have been proposed to refine the generic model shape to approximate a specific person.

In previous research [146] a generic mesh model was refined based on front and side view silhouettes taken with a single camera. Texture mapping was then applied to approximate detailed surface appearance. Recently simultaneous capture from multiple calibrated views has been used [53,289,352] to achieve more accurate shape and appearance. Plaenkers and Fua [289] initialize upper-body shape by fitting an implicit ellipsoidal metaball representation to stereo point clouds prior to tracking. Carranza *et al.* [53] fit a generic mesh model to multiple view silhouette images of a person in a fixed pose prior to tracking whole-body motion. Starck and Hilton [352] reconstruct whole-body shape and appearance for a person in an arbitrary pose by optimizing a generic mesh model with respect to both silhouette, stereo and feature correspondence constraints in multiple views. These model fitting approaches provide an accurate parameterized approximation of a person provided the assumed shape of the generic model is a reasonable initial approximation. Model fitting methods commonly assume short hair and close fitting clothing which limits their generality.

The availability of sensors for whole-body 3D scans provides accurate measurement of surface shape. Techniques to fit generic humanoid models to the whole-body scans in a specific pose enable a highly detailed representation of a person's shape to be parameterized for animation and tracking [14,351]. Allen *et al.* [14] fit a sub-division surface to multiple scans of a person in different poses to parameterize the change in body surface shape with pose. Databases of 3D scans have also been used to learn statistical models of the inter-person variation in whole-body shape [15,363]. Reconstruction of shape from images can then be constrained by the learnt model to improve performance.

2.3 Appearance Initialization

Due to the large intra and inter person variability in appearance with different clothing, initialization of appearance has commonly been based on the observed image set. Statistical models of color are commonly used for tracking, see section 3.3. Initialization of the detailed surface appearance for model-based pose estimation has also used texture maps derived from multiple view images [53,352]. A cost function evaluating the difference in appearance between the projected model and observed images is then used in pose estimation.

Sidenbladh and Black [332,333] address modeling the likelihood of image observations for different body parts. They learn the statistics of appearance and motion based on filter responses for a set of training examples. In a related approach, Roberts *et al.* [309] learn the likelihood of body part color appearance using multi-modal histograms on a 3D surface model. Results are presented for 2D tracking of upper-body and walking motions in cluttered scenes.

A recent trend has been towards the learning of body part detectors to identify possible locations for body parts which are then combined probabilistically to locate people [235,296,310,314], see section 4.1.1. Initialization of such models requires a large training corpus of both positive and negative training examples for different body parts. Approaches such as AdaBoost have been successfully used to learn body part detectors such as the face [380], hands, arms, legs and torso [235,310]. Alternatively, Ramanan *et al.* [296] detect key-frame poses in walking sequences and initialize a local appearance model to detect body parts at intermediate frames.

Lim *et al.* [220] address the problem of changing appearance due to motion by modeling the dynamics of the appearance for walking humans. This is done by mapping the pixels inside a bounding box to a low dimensionality space (only 3D) using a nonlinear Local Linear Embedding algorithm. In this space the temporal continuity of the appearance is preserved, which allows for learning a dynamic model of the appearance for walking humans. This model can then be used to predict not only the position and 2D shape of a walking human, but also the appearance.

The initialization of models which accurately represent the change in appearance over time due to creases in clothing, hair and change in body shape with movement remains an open problem. Recent introduction of robust local body part detectors provides a potential solution for tracking and pose estimation.

2.4 Discussion of Advances in Model Initialization

Initialization of shape, appearance and pose remains an import step to automate the process of human motion capture and analysis. As illustrated in this review signifi-

cant advances have been made towards automatic solutions. The problem of initializing the kinematic structure and pose from feature tracks for monocular sequences has been addressed [350]. A number of researchers have presented methods for initializing the kinematic structure from multiple view image sequences using an intermediate volumetric reconstruction [59,233]. These approaches provide a solution to the problem of automatic kinematic model initialization for human pose estimation. Learning approaches [144] and anthropometric models [252] have been presented to initialize the joint angle limits on the kinematic structure to constrain tracking and pose estimation.

Over the past five years there has been substantial research in the automatic initialization of model shape from multiple view images [53,59,289,352]. These approaches reconstruct an articulated model which approximates the shape of a specific person providing the basis for improved accuracy in tracking. Recent research has also started to address the modeling of changes in human body shape during movement [14]. Similarly multiple view reconstruction techniques have allowed the automatic initialization of model appearance to that of a specific individual.

Initialization of appearance models for monocular tracking and pose estimation remains an open problem. A number of approaches have been proposed for initialization of appearance based on image patch exemplars or color mixture models. Recent work on body part detectors has exploited supervised learning approaches to discriminate individual body part appearance from background [296,310,314]. Only limited research has addressed the problem of modeling changes in a person's appearance during movement. The problem of fully automatic initialization of model kinematics, shape and appearance for human pose estimation from monocular image sequences remains open for future research.

3 Tracking

Since 2000 tracking algorithms have focused primarily on surveillance applications leading to advances in areas such as outdoor tracking, tracking through occlusion, and detection of humans in still images. In this section we review recent advances in these areas as well as more general tracking problems.

The notion of *tracking* in visual analysis of human motion is used differently throughout the literature. Here we define it as consisting of two processes: 1) *figure-ground segmentation* and 2) *temporal correspondences*. The latter, temporal correspondences, is the process of associating the detected humans in the current frame with those in the previous frames, providing temporal trajectories through the state space. Recent advances are mainly due to processing more natural scenes where multiple people and occlusions are present.

Figure-ground segmentation is the process of separating the objects of interest (humans) from the rest of the image (the background). Methods for figure-ground segmentation are often applied as the first step in many systems and therefore a crucial process. Recent advances are mostly a result of expanding existing methods. We categorize these methods in accordance with the type of image measurements the segmentation is based on: motion, appearance, shape, or depth data. Before describing these we first review recent advances in background subtraction as this has become the initial step in many tracking algorithms.

3.1 Background Subtraction

Up until the late 90s background subtraction was known as a powerful preprocessing step but only in controlled indoor environments. In 1998 Stauffer and Grimson [354] presented the idea of representing each pixel by a mixture of Gaussians (MoG) and updating each pixel with new Gaussians during run-time. This allows background subtraction to be used in outdoor environments. Normally the updating was done recursively, which can model slow changes in a scene, but not rapid changes like clouds. The method by Stauffer and Grimson has today become the standard of background subtraction. However, since 1998 a number of advances have been seen which can be divided into *background representation*, *classification*, *background updating*, and *background initialization*.

3.1.1 Background Representation

The MoG representation can be in RGB space, but also other color spaces can be applied, see [201] for an overview. Often a representation where the color and intensities are separated is applied, e.g., YUV [394], HSV [69] and normalized RGB [232], since this allows for detecting shadow-pixels wrongly classified as object-pixels [293]. Using a MoG in a 3D color space corresponds to ellipsoids or spheres (depending on the assumptions on the covariance matrix) of the Gaussian representations [232,354,421]. Other geometric representations are truncated cylinders [195] and truncated cones [18].

Conceptually different representations have also been developed. Elgammal *et al.* [96] use a kernel-based approach where they represent a background pixel by the individual pixels of the last N frames. Haritaoglu *et al.* [137] represent the minimum and maximum value together with the maximum allowed change of the value in two consecutive frames. Heikkila and Pietikainen [140] represent each background pixel by a bit sequence, where each bit reflects whether the value of a neighboring pixel is above or below the pixel of interest, i.e., a texture operator. This makes the background model invariant to monotonic illumination changes. Oliver *et al.* [270] also use a pixel's neighbors to represent it. They apply an eigenspace representation

of the background and detect new objects by comparing the input image with an image reconstructed via the eigenspace.

Eng *et al.* [101,102] divide a background model learnt over time into a number of non-overlapping blocks. The pixels within each block are grouped into at most three classes according to homogeneity. The means of these classes are then the representation of the background for this block, i.e., a spatio-temporal representation. Heikkila and Pietikainen have also applied their texture operator for a spatio-temporal block-based (overlapping blocks) background segmentation [139]. Other spatio-temporal approaches are [256] and [423] where the background is represented by a predicted region found by an autoregressive process.

The choice of representation is not only dependant on the accuracy but also on the speed of the implementation and the application. This makes sense since the overall accuracy of background subtraction is a combination of representation, classification, updating, and initialization. For example, Cucchiara *et al.* [69] use only one value to represent each background pixel, but still good results (and speed) can be obtained due to advanced classification and updating. It should however be noted that the MoG representation is by far the most widely used method². For scenes with dynamic background the MoG representation does not suffice and methods directly aimed at modeling dynamic background should be applied, see e.g., [256], [327], and [423].

3.1.2 Classification

A number of false positives and negatives will often be present after a background subtraction, for example due to shadows [293]. Using standard filtering techniques based on connected component analysis, size, median filter, morphology, and proximity can improve the result [69,96,128,232,408,420]. Alternatively, the fact that neighboring pixels are likely to be both foreground *or* background can be used in classification. Markov Random Fields have been applied to implement this idea [323,327].

Recent methods have tried to directly identify the incorrect pixels and use classifiers to separate the pixels into a number of sub-classes: unchanged background, changes due to auto iris, shadows, highlights, moving object, cast shadow from moving object, ghost object (false positive), ghost shadow, etc. [57,69,148]. Classifiers have been based on color, gradients [232], flow information [69], and hysteresis thresholding [101].

² See [207,424] for optimizations of the MoG representation.

3.1.3 Background Updating

In outdoor scenes, in particular, the value of a background pixel will change over time and an update mechanism is therefore required. The slow changes in the scene can be updated recursively by including the current pixel value into the model as a weighted combination [69,96,232,354]. A different approach is to measure the overall average change in the scene compared to the expected background and use this to update the model [18,408]. If no real-time requirements are present, both past and future values can be used to update the background [106]. In general, for a good model update only pixels classified as unchanged background should be updated.

Rapid changes in the scene are accommodated by adding a new mode to the model. For the MoG model a new mode is a new Gaussian distribution, which is initiated whenever a non-background pixel is detected. The more pixels (over time) that support this distribution the more weight it will have. A similar approach is seen in [18,195] where the background model, denoted a codebook, for each pixel is represented by a number of codewords (cylinders [195] or cones [18] in RGB-space). During run-time each foreground pixel creates a new codeword. A codeword not having any pixels assigned to it for a certain number of frames is eliminated. A similar idea can be found in [139,140].

3.1.4 Background Initialization

A background model needs to be learned during an initialization phase. Earlier approaches assumed that no moving objects are present in a number of consecutive frames and then learn the model parameters in this period. However, in real scenarios this assumption will be invalid and recent methods have therefore focused on initialization in the presence of moving objects.

In the MoG representation moving objects can to some extent be accepted during initialization since each foreground object will be represented by its own distribution which is likely to have a low weight. However, this erroneous distribution is likely to produce false positives in the classification process. A different approach is to find only pixels that are true background pixels and then only apply these for initialization. This can be done using a temporal median filter if less than 50% of the values belong to foreground objects [101,118,137]. Eng *et al.* [101] combine this with a skin detector to find and remove humans from the training images.

Recent alternatives first divide the pixels in the initialization phase into temporal subintervals with similar values. Second, the "best" subinterval belonging to the background is found as the subinterval with the minimum average motion (measured by optical flow) [129] or the subinterval with the maximum ratio between the number of samples in the subinterval and their variance [385,386]. The codeword method mentioned above uses a temporal filter after the initialization phase to elim-

inate any codeword that has not recurred for a long period of time [195]. A similar approach has used in [139,140].

For comparative studies among some of the different background subtraction methods see [55,61,385,386].

3.2 *Motion-Based Segmentation*

Motion-based figure-ground segmentation is based on the notion that differences in consecutive images arise from moving humans, i.e., by finding the motion you find the human. The motion is measured using either flow or image differencing.

Sidenbladh [331] calculates optical flow for a large number of image windows each containing a walking human. A Support Vector Machine (SVM) is used to detect walking humans in video. Optical flow can be noisy and instead image flow can be measured using higher level entities. For example, Gonzalez *et al.* [121] track KLT-features to obtain flow vectors, Sangi *et al.* [320] extract flow vectors from displacements of pixel-blocks, and Bradski and Davis [40] find flow vectors as gradients in Motion History Images (MHI) [80].

Image differencing adapts quickly to changes in the scene, but pixels from a human that has not moved or are similar to their neighbors are not detected. Therefore, an improved version is to use three consecutive images [66,137,184]. A different type of image differencing is used by Viola *et al.* [382]. They apply the principle of their novel face detector [380], where simple features are combined in a cascade of progressively more advanced classifiers. A rectangle of pixels in the current image is compared to the corresponding rectangle in the previous image. This is done by shifting the rectangle in the current image up, down, left, and right. Image differencing is then performed and the lower the energy in the output the higher the probability that the human has actually moved (shifted) in this direction. The output of these operations is used to build a person detector, which is trained using AdaBoost.

3.3 *Appearance-Based Segmentation*

Segmentation based on the appearance of the human is built on the idea that 1) the appearance of human and background is different and 2) the appearance of individuals are different. The approaches work by building an appearance model of each human and then either building appearance models of the segmented foreground objects in the current image and comparing them with the predicted models, or by directly segmenting the pixels in the current image that belong to each model. Some of these methods are independent on the temporal context, meaning that the meth-

ods apply a general appearance model of a human, as opposed to methods where the appearance model of the human is learned/updated based on previous images in the current sequence.

3.3.1 Temporal context-free

Temporal context-free methods are used to detect humans in a still image [254], to detect humans entering a scene [269], or to index images in databases [275]. Advances are mostly on using massive amount of training data for learning good classifiers. For example, Okuma *et al.* [269] use 6000 images to train an Adaboost-based classifier. Other examples are using DCT coefficients [275], using partial-occlusion handling body-part detectors [254], (see also section 4.1.1), or the block-based method by Utsumi and Tetsutani [374]. In [374] the image is divided into a number of blocks and the mean and covariance matrix of the intensities are calculated for each block. A distance matrix is constructed where an entry represents the generalized Mahalanobis distance between two blocks. The detection is now based on the fact that for non-human images the distances between blocks in the proximity will be larger than for images containing a human.

Common for these methods is that the human is detected as a box (normally a bounding box) and clutter in the background will therefore have an effect on the results. Furthermore, as the methods usually represent the human as one entity, as opposed to a number of sub-entities, occlusion will in general effect the methods strongly. Drastic illumination changes will also effect the methods since the models are general and do not adapt to the current scene.

3.3.2 Temporal context

Temporal context refers to methods where a model which is learned and updated in previous images is used to either detect foreground pixels or to classify foreground pixels to a particular human being tracked. The methods either operate at pixel level or region level. At pixel level the likelihood of each (foreground) pixel belonging to a human model is calculated. The region level is when a region in the image, such as a bounding box, is compared to an appearance model of the humans that are predicted to be present in the current frame, i.e. the probability that a region in an image corresponds to a particular human model. Color-based appearance models have recently received attention leading to advances allowing tracking in outdoor scenes with partial occlusion. This has led to a need for models that can represent the differences between individuals even during partial occlusion.

In many systems the color of a human is represented as either a color histogram [67,154,232,269,401,421] or a MoG [186,193,316,404]³. Color histograms are

³ According to McKenna *et al.* [232] MoG is preferred with small sample sets and many

normally compared using the Bhattacharyya distance, which can be improved by weighting pixels close to the center of the human higher than those close to the border [67,421]. In Zhao [421] the similarity is combined with the dissimilarity with respect to the color histogram of the background. MoG representations are normally compared using the Mahalanobis distance, which can be evaluated efficiently by using only one Gaussian [186] and assuming independence between color channels [70]. Alternatively, only the mean can be used [404].

Representing the entire human by just one color model is often too coarse a representation even though the model contains multiple modes. Recent advances are therefore on including spatial information. For example using a Correlogram, which is a co-occurrence matrix that expresses the probability of two different colored pixels being found at a certain distance from each other [52,161]. Another way of adding spatial information is to divide the human into a number of sub-regions and represent each sub-region with either a color histogram or a MoG [244,269,316,404]. Hu *et al.* [154] use an adaptive approach to obtain three sub-regions representing the head, torso, and legs. A more general approach is to model the human as a number of blobs where each blob is a connected group of pixels having a similar color [193,282]. Grouping the blobs together temporally and spatially into an entire human requires some bookkeeping, but a rough human model can assist as seen in [282].

As mentioned in section 2 - Model Initialization - appearance-based models able to handle changes over time remains an open issue. On one hand a model should adapt quickly to changes, but on the other hand long term temporal consistency is required, e.g., to handle occlusions. The KLT-tracker [329] to some degree handles this dilemma by only updating the model by data from the previous image as long as it is not too different from the initial model. A more general framework is suggested by Jepson *et al.* [177]. They update each pixel in their appearance model by a weighted combination of a slowly changing model, a fast changing model, and a noise model. The weights are updated in accordance with the support of the different models in the current image.

3.4 *Shape-Based Segmentation*

The shape of a human is often very different from the shape of other objects in a scene. Shape-based detection of humans can therefore be a powerful cue. As opposed to the appearance-based models, the shapes of individuals are often very similar. Hence, shape-based methods applied to tracking only involves simple correspondences. The advances are first of all to allow human detection and tracking

possible colors, whereas a color histogram is preferred when many color samples are present in a coarsely quantified color space.

in uncontrolled environments. Due to the recent advances in background subtraction reliable silhouette outlines can describe the shape of the humans in the image sequence. Furthermore, advances in representations and segmentation methods of humans in still images have also been reported. As was done for the appearance-based methods, we divide the shape-based methods into those not using the temporal context and those using the context.

3.4.1 *Temporal context-free*

Zhao and Thorpe [417] use depth data to extract the silhouettes of individuals in the image. A neural network is trained on upright humans and used to verify whether the extracted silhouettes actually originate from humans or not. To make the method more robust the gradients of the outline of a silhouette are used to represent the shape of the human. Leibe *et al.* [214] learn the outlines of walking humans and store them as a number of templates. Each of these are matched with an edge version of the input image over different scales using Chamfer matching. The results are combined with the probability of a person being present, which is measured by comparing small learned image patches of the appearance of humans and their occurrence distribution. Wu and Yu [399] learn a prior shape model for human edges and represent it as a Boltzmann distribution in a Markov Field. The detector searches for different locations, scales, and rotations and is implemented using a Particle Filter. Dalal and Triggs [75] use an SVM to detect humans in a window of pixels. The input is a set of features encoding the shape of a human. The features come from using a spatially arranged set of HOG (Histogram of Oriented Gradients) descriptors. The HOG descriptor operates by dividing an image region into a number of cells. For each cell a 1D histogram of gradient directions over the pixels in the cell is calculated. In [76] the work is extended by including motion histograms. This allows for detecting humans even when the camera and/or background is moving. HOGs are related to Shape Contexts [30] and SIFT (Scale Invariant Feature Transformation) [223]. Zhao and Davis [416] learn a hierarchy of silhouette templates for the upper body. The outline of the silhouettes in the templates is used to detect sitting humans in a frame. This is done using Chamfer matching at different scales together with a color-based detector that is updated iteratively.

3.4.2 *Temporal context*

When the temporal context is taken into consideration shape-based methods can be applied to track individuals over time. In case of temporal smoothness the shape in the previous frame can be used to find the human in the current frame. Haritaoglu *et al.* [137] perform a binary edge correlation between the outlines of the silhouettes in the last frame and the immediate surroundings in the current image. Davis *et al.* [84] use a Point Distribution Model (PDM) to represent the outline of the human.

The most likely configurations of the outline from the previous frame are used to predict the location in the current frame using a particle filter. Predictions are evaluated by comparing the edges of the outline with those in the image. A similar approach is seen in [197] where the active shape model is applied to find a fit in the current frame. Atsushi *et al.* [21] model the pose of the human in the previous frame by an ellipse and predict nine possible poses of the human in to the current frame. Each of these is correlated with the silhouettes in the current image in order to define the current pose of the human. Krüger *et al.* [203] correlate the extracted silhouette with a learned hierarchy of silhouettes of walking persons. At run-time a Bayesian tracking framework concurrently estimates the translation, scale, and type of silhouette.

In situations of partial occlusion the shape-based methods just described often fail due to lack of global shape information. Advances therefore include detection of humans based on only a few parts of the overall shape. In the work by Wu and Nevatia [396] four different (body)parts are detected: full-body, head-shoulder, torso, and legs. For each part a detector is trained using a boosting classifier together with edgelets (small connected chains of edge pixels) which are quantified into different orientations, see also section 4.1.1. When people group together the occlusion often becomes severe and the only reliably shape information is the head or head-shoulder profile. While this work is limited to frontal/rear views, extended work also handles side views [395].

In [137,155,406] the head candidates are found by analyzing the silhouette boundary and the vertical projected histogram of the silhouette. A similar approach is seen in [419] except that also an edge-based method to find the head-shoulder profile inside silhouettes is applied.

3.5 *Depth-Based Segmentation*

Figure-ground segmentation using depth data is based on the idea that the human stands out in a 3D environment. Methods are either based directly on estimated 3D data for the scene [134,138,170,221,407] or indirectly by combining different camera views after features have been extracted [171,243,244,405]. Advances are mainly due to faster computers allowing for handling multiple camera inputs.

Background subtraction can be sensitive to lighting changes. Therefore a depth-based approach can be taken where the background is modeled as a depth model and compared to estimated depth data for each incoming frame in order to segment the foreground. A real-time dense stereo algorithm is, however, still problematic unless special hardware is applied [221]. An approach to circumvent this is the work by Ivanov *et al.* [170] where an online depth map is not required. Instead the mapping between pixels in two cameras is learnt. This allows for an online

comparison between associated pixels (defined by the mapping) in the two cameras. Detection is now performed based on the assumption that the color and intensity are similar for the pixels if and only if they depict the background. In [221] the merits and drawbacks of this approach are studied in detail.

Other advances in human detection based on depth data include the work by Haritaoglu *et al.* [134] where depth data produced by ceiling-mounted cameras are projected to the ground-plane. Here humans are located by looking for a 3D head-shoulder profile. Similar approaches are seen in [138,407] except for the camera placement and that [138] apply voxels as opposed to 3D points.

Mittal and Davis [243,244] detect humans using an appearance-based method in each camera view. The center of each detected human is combined with those found in another image using region-based stereo constrained by the epipolar geometry. The resulting 3D points are projected to the ground-plane and represented probabilistically using Gaussian kernels and an occlusion likelihood. In Yang *et al.* [405] silhouettes from different cameras are combined into the visual hull. The incorrect interpretations are pruned using a size criterion as well as the temporal history. Iwase and Saito [171] apply multiple cameras to detect and track multiple people. In each camera the feet of each person are detected using background subtraction and knowledge of the environment. For each camera all detected feet are mapped to a virtual ground-plane where an iterative procedure resolves ambiguities. A similar approach can be found in [194].

3.6 Temporal Correspondences

One of the primary tasks of a tracking algorithm is to find the temporal correspondences. That is, given the state of N persons in the previous frame(s) and the current input frame(s), what are the states of the same persons in the current frame(s). Here the state is mainly the image position of a person, but can contain other attributes, e.g., 3D position, color and shape.

Previously tracking algorithms were mostly tested in controlled environments and with only a few people present in the scene. Recently, algorithms have addressed more natural outdoor scenarios where multiple people and occlusions are present. One problem is to have better figure-ground segmentation as discussed above. Another equally important problem is how to handle multiple people that might occlude each other. In this section we discuss advances related to temporal correspondences *before and after occlusion* and temporal correspondences *during occlusion*.

3.6.1 Temporal Correspondences Before and After Occlusion

A model of each individual must be constructed before any tracking can commence. Recent methods are aiming at doing this automatically. One way is to look for (new) large foreground objects possible near the boundaries⁴ [18,19,137,232,316]. Alternatively, a new person can be defined as a foreground object detected far from any predictions [52]. Khan and Shah [193] fit 1D Gaussians to the foreground pixels projected to the horizontal axis. If the number of good fits is higher than the predicted number of people in the scene then a new person has entered the scene.

When the tracking has commenced the problem is to find the temporal correspondences between predicted and measured states. This has recently been approached using a correspondence matrix, which has the predicted objects in one direction and the measured objects in the other direction. For each entry in the matrix a distance between predicted and measured object is calculated. This gives the likelihood that a predicted and measured object are the same. By analyzing the columns and rows the following situations can be hypothesized: new object, object lost, object match, split situation, and merge situation. In case of for example merge and split situations the matrix can not be resolved directly and ad hoc methods are applied. For example by analyzing the motion vectors and the area (change) of each foreground object [52,70,128,232,395,401,408].

Alternatively, global optimizations can also be applied. Polat *et al.* [290] use a Multiple Hypothesis Tracker to construct different hypotheses which each explains all the predictions and measurements, and chooses the hypothesis which is most likely. To prune the combinatorial number of different hypotheses smoothness constraints on the motion trajectories are introduced. If the total number of people in the scene is known in advance the pruning becomes less difficult [29,154]. Another global optimization can be seen in [345,421] where a Particle Filter [168] is applied and where each state is a multi-object configuration (hypothesis). Objects are allowed to enter and exit the scene meaning that the number of elements in the state vector can change. To handle this the particle filter is enhanced by a trans-dimensional Markov chain Monte Carlo approach [125]. This allows new objects to enter and other objects to leave the scene, i.e., the dimensionality of the state space may change. In the work by Li *et al.* [218] a tree-based global optimization for correspondence between multiple objects across multiple views is presented. This approach is used for real-time tracking of hand, head and feet for whole-body pose estimation. Antonini *et al.* [19] learn behavioral models for pedestrians' preferences regarding acceleration and direction. These models are used to find globally coherent trajectories.

⁴ Similar approaches can be used to detect when people are leaving the scene, see e.g., [18,128].

3.6.2 Temporal Correspondences During Occlusion

Tracking during occlusion was not addressed in previous work, instead the track of the group was used to update the states of the individuals. However, this makes it impossible to update the models of the individuals, which can result in unreliable tracking after the group splits up. Furthermore, interactions between humans during occlusions is difficult to analyze when they are represented as one foreground object. Therefore the problem of finding the correspondences during occlusion has been investigated recently.

In some recent systems the first task is to actually detect that an occlusion is present. This can be done using the corresponding matrix mentioned above or as in [52,193,316]. Khan and Shah [193] detect a non-occlusion situation as a situation when the detected foreground objects are far from each other. Capellades *et al.* [52] define a merge as a situation where the total number of foreground objects has decreased and where two or more foreground objects from the previous frame overlap with one foreground object in the current frame. In the work by Roth *et al.* [316] a merge is detected as one of eight different types of occlusion based on the depth ordering and the layout of the bounding boxes. This allows for only using the reliable parts of the bounding box to update the position of the human.

Different approaches for assigning pixels to individuals during occlusion have been reported in recent publications. A local approach is to assign each pixel to the most likely predicted model using a probabilistic method [193,282]. A local approach allows for bypassing the occlusion problem but it is also sensitive to noise and therefore often combined with some post-processing to reassign wrongly classified pixels. Global approaches try to classify pixels based on for example the assumption that people in a group are standing side by side with respect to the camera. This assumption allows for defining vertical dividers between the individuals based on the positions of their heads. Foreground pixels are then assigned to individuals based on these dividers [136,401,406]. When a certain depth ordering is present in the group the assumption fails.

In the work by McKenna *et al.* [232] the depth ordering is found explicitly. During occlusion the likelihood of each pixel in the foreground object belonging to a person is calculated using Bayes rule. The posteriors for each person are added to obtain an overall probability of each person. These probabilities are then used to define the fraction of each person that is visible. This is denoted a visibility index and can be applied to find the depth ordering. In [316] the depth ordering is based on assuming a planar floor. This will result in the closest object to the camera having the highest vertically coordinate. Xu and Puig [401] generalize this idea by using projective geometry to find the line in the image that corresponds to the "horizon line" in the 3D scene. The object closest to the camera is found as the object closest to this horizontal line.

3.7 Discussion of Advances Human Tracking

Advances in figure-ground segmentation have to a large extent been motivated by the increased focus on surveillance applications. For example, in order to have fully autonomous systems operating in uncontrolled environments the segmentation methods have to be adaptive. This has to some extent been achieved within background subtraction where analysis of video sequences of several hours has been reported [18]. However, for 24 hour operation special cameras (and algorithms) are required. Work in this direction has started [66,82] but no one has so far been able to report a truly autonomous system. Furthermore, in most surveillance applications multiple cameras are required to cover the scene of interest at an acceptable resolution. Systems for self-calibrating and tracking across different cameras are being investigated [21,186,192,369], but again, no fully autonomous system has been reported.

Another advance in segmentation is to apply spatial information in the color-based appearance models, for example by dividing each foreground object into a number of regions each having a color representation [154,193,244,269,282,316,404] or by correlograms [52,161]. This has allowed for relatively reliable detection and tracking of people even when multiple people are present with occlusion. Even an accurate appearance model might fail when the lighting changes are significant.

The recent focus on natural scenes has also led to advances within methods for temporal correspondence, especially handling the occlusion problem. Advances are mainly due to the use of probabilistic methods, for example to segment pixels to individuals during occlusion [193,232,282,285] and also to handle multiple hypotheses and uncertainties using stochastic sampling methods [154,269,290,345,404,421]. In fact, concurrent segmentation and tracking can be handled by stochastic sampling methods. It is expected that future work will be based on this framework since it unifies segmentation and tracking *and* the associated uncertainties.

The use of common benchmark data has begun to underpin progress. As has been seen in the speech community for many years and lately in the face recognition community, widely acceptable benchmark data can help to focus research. Within human detection a few recent benchmark data sets have been reported [75,254]. Within tracking in general the PETS and VS-PETS data sets [5] have been applied in many systems.

4 Pose Estimation

Pose estimation refers to the process of estimating the configuration of the underlying kinematic or skeletal articulation structure of a person. This process may be

an integral part of the tracking process as in model-based analysis-by-synthesis approaches or may be performed directly from observations on a per-frame basis. The previous survey [247] separated pose estimation algorithms into three categories based on their use of a prior human model:

Model-Free: This class covers methods where there is no explicit a priori model. Previous methods in this class take a bottom up approach to tracking and labeling of body parts in 2D [394] or direct mapping from 2D sequences of image observations to 3D pose [41].

Indirect Model Use: In this class methods use an a priori model in pose estimation as a reference or look-up table to guide the interpretation of measured data. Previous examples include human body part labeling using aspect ratios between limbs [49] or pose recognition [135].

Direct Model Use: This class uses an explicit 3D geometric representation of human shape and kinematic structure to reconstruct pose. The majority of approaches employ an analysis-by-synthesis methodology to optimize the similarity between the model projection and observed images [147,383].

In this section we identify recent contributions and advances in each category of pose estimation algorithms. A number of trends can be identified from the literature. Three research directions which have each received considerable attention are: the introduction of probabilistic approaches to detect body parts and assemble part configurations in the model-free category; the incorporation of learnt motion models in pose estimation to constrain the recovered 3D human motion; and the use of stochastic sampling techniques in model-based analysis-by-synthesis to improve robustness of 3D pose estimation.

Two important distinctions relating to the difficulty of the pose estimation problem are identified in this analysis: pose estimation from single vs. multiple view images; and 2D pose estimation in the image plane vs. full 3D pose reconstruction. The most difficult and ill-posed problem is the recovery of full 3D pose from single view images towards which initial steps have been made. There has also been substantial research addressing the problems of 2D pose estimation from single view and 3D pose estimation from multiple views. For example recent advances have demonstrated 2D pose estimation in complex natural scenes such as film footage.

4.1 *Model Free*

A recent trend to overcome limitations of tracking over long sequences has been the investigation of direct pose detection on individual image frames. Two approaches have been investigated which fall into this model-free pose estimation category: *probabilistic assemblies of parts* where individual body parts are first detected and then assembled to estimate the 2D pose; and *example-based methods* which directly

learn the mapping from 2D image space to 3D model space.

4.1.1 Probabilistic Assemblies of Parts

Probabilistic assemblies of parts have been introduced for direct bottom-up 2D pose estimation by first detecting likely locations of body parts and then assembling these to obtain the configuration which best matches the observations. A potential advantage of detection over tracking is that the pose can be estimated independently at each frame, allowing pose estimation for rapid movements. Temporal information may be incorporated to estimate consistent pose configurations over sequences. Forsythe and Fleck [109] introduced the notion of body plans to represent people or animals as a structured assembly of parts learnt from images. Following this direction [104,166,167] used pictorial structures to estimate 2D body part configurations from image sequences. Combinations of body part detectors have recently been used to address the related problem of locating multiple people in cluttered scenes with partial occlusion [254,396], see section 3.

Probabilistic assemblies of body part detectors (face, hands, arms, legs, torso) have been investigated for bottom up estimation of whole-body 2D pose in individual frames or sequences [235,296,310,314]. Individual body parts are detected using 2D shape [310], SVM classifiers [314], AdaBoost [235], and locally initialized appearance models [296]. Mikolajczyk *et al.* [240] introduced probabilistic assemblies of robust AdaBoost body part detectors to locate people in images providing a coarse 2D localization. The probabilistic assembly of parts models the joint likelihood of a body part configuration. In [235] this approach is extended to whole-body 2D pose estimation in frontal images using RANSAC to assemble body part configurations with prior pose constraints. Ramanan *et al.* [296] present a related approach where lateral views of a scissor-leg pose for a person walking or running are detected from film footage. Detected poses are then used as key-frames to initialize a local appearance model for body part detection and 2D pose estimation at intermediate frames.

Recent work has also introduced approaches for 2D pose estimation from single images. Ren *et al.* [302] use pairwise constraints between body parts to assemble body part detections into 2D pose configurations. Ramanan *et al.* [297] learn a global body part configuration model based on conditional random fields to simultaneously detect all body parts. Pairwise constraints include aspect ratio, scale, appearance, orientation and connectivity. Hua *et al.* [157] present an approach to 2D pose estimation from a single image using bottom-up feature cues together with a Markov network to model part configurations. Both of these approaches demonstrate impressive results for pose estimation in cluttered scenes such as sports images.

An important contribution of approaches based on the probabilistic assembly of

parts is 2D pose estimation in cluttered natural scenes from a single view. This overcomes limitations of many previous pose estimation methods which require structured scenes, accurate prior models or multiple views.

4.1.2 Example-Based Methods

A number of example-based methods for human pose estimation have been proposed which compare the observed image with a database of samples. Brand [41] used a hidden Markov model (HMM) to represent the mapping from 2D silhouette sequences in image space to skeletal motion in 3D pose space. In this work the mapping for specific motion sequences was learnt using rendered silhouette images of a humanoid model. The HMM was used to estimate the most likely 3D pose sequence from an observed 2D silhouette sequence for a specific view. Similarly, Rosales *et al.* [294,315] learn a mapping from visual features of a segmented person to static pose using neural networks. This representation allows 3D pose estimation invariant to speed and direction of movement. Viewpoint invariant representation of the mapping from image to pose is investigated in [272].

To overcome limitations of tracking researchers have investigated example-based approaches which directly lookup the mapping from silhouettes to 3D pose [6,151,326,340]. Howe [151] uses a direct silhouette lookup using Chamfer distance to select candidate poses together with a Markov chain for temporal propagation for 3D pose estimation of walking and dancing. Shakhnarovich *et al.* [326] present an example-based approach for viewpoint invariant pose estimation of upper-body 3D pose from a single image. Parameter-sensitive hashing is used to represent the mapping between observed segmented images from multiple views and the corresponding 3D pose. Grauman *et al.* [124] learn a probabilistic representation of the mapping from multiple view silhouette contours to whole-body 3D joint locations. Pose reconstruction is demonstrated for close-up images of a walking person from multiple or single views. Similarly, Elgammal and Lee [98] learn multiple view-dependent mapping from silhouettes to 3D pose for walking actions. Agarwal and Triggs [6,8] presented an example-based approach for 3D pose estimation from single view image sequences. Nonlinear regression is used to learn the mapping from silhouette shape descriptors to 3D pose. Results demonstrate reconstruction of long sequences of walking motions with turns from monocular video.

Example-based approaches represent the mapping between image and pose space providing a powerful mechanism for directly estimating 3D pose. Commonly these approaches exploit rendering of motion capture data to provide training examples with known 3D pose. A limitation of current example-based approaches is the restriction to the poses or motions used in training. Extension to a wider vocabulary of movements may introduce ambiguities in the mapping.

4.2 Indirect Model Use

A number of researchers have investigated direct reconstruction of both model shape and motion from the visual-hull [59,237,238] without a prior model. Mikic *et al.* [237,238] present an integrated system for automated recovery of both a human body model and motion from multiple view image sequences. Model acquisition is based on a hierarchical rule-based approach to body part localization and labelling. Prior knowledge of body part shape, relative size and configuration is used to segment the visual-hull. An extended Kalman filter is then used for human motion reconstruction between frames. A voxel labelling procedure is used to allow large inter-frame movements. Cheung *et al.* [59] first reconstruct a model of the kinematic structure, shape and appearance of a person and then use this to estimate the 3D movement. Tracking is performed by hierarchically matching the approximate body model to the visual-hull using color matching along the silhouette boundary edge.

An alternative approach based on full 3D-to-3D non-rigid surface matching using spherical mapping is presented in [353]. Alignment of a skeletal model with the first frame allows the 3D motion to be recovered from the non-rigid surface motion. Results of these approaches demonstrate 3D human pose estimation for rapid movement of subjects wearing tight clothing.

These approaches exploit scene reconstruction from multiple views to directly recover both shape and motion. This approach is suitable for multiple camera studio based systems allowing estimation of complex human movements.

4.3 Direct Model Use

The use of an explicit model of a person's kinematics, shape and appearance in an analysis-by-synthesis framework is the most widely investigated approach to human pose estimation from video. In the previous survey [247] fifty papers (40% of those surveyed) were in this category starting with some of the earliest work in human pose estimation [147]. Model-based analysis-by-synthesis has continued to be a dominant methodology for human pose estimation.

The main novel research directions are: the introduction of stochastic sampling techniques based on sequential Monte Carlo; and the introduction of constraints on the model in particular learnt models of human motion. In this section we review key papers contributing to these advances in multiple and single view model-based pose estimation.

4.3.1 Multiple View 3D Pose Estimation

Up to 2000 the majority of approaches to human pose estimation employed deterministic gradient descent techniques to iteratively estimate changes in pose [86,289]. The extended Kalman filter was widely applied to human tracking with low-order dynamics used to predict change in pose [384]. Recent work using model-based analysis-by-synthesis has extended deterministic gradient descent based approach to more complex motions. For example Plänkers and Fua [289] demonstrated upper body tracking of arm movements with self-occlusion using stereo and silhouette cues. A common limitation of gradient descent approaches is the use of a single pose or state estimate which is updated at each time step. In practice if there is a rapid movement or visual ambiguities pose estimation may fail catastrophically. To achieve more robust tracking, techniques which employ a deterministic or stochastic search of the pose state space have been investigated.

Stochastic tracking techniques, such as the *particle filter*, were introduced for robust visual tracking of objects where sudden changes in movement or cluttered scenes can result in failure. The principal difficulty with their application to human pose estimation is the dimensionality of the state space. The number of samples or particles required increases exponentially with dimensionality. Typically whole-body human models use more than 20 degrees-of-freedom making direct application of particle filters computationally prohibitive. MacCormick and Isard [230] proposed partitioned sampling of the state space for efficient 2D pose estimation of articulated objects such as the hand. However, this approach does not extend directly to the dimensionality required for whole-body pose estimation. Deutscher *et al.* [90] introduced the *annealed particle filter* which combines a deterministic annealing approach with stochastic sampling to reduce the number of samples required. At each time step the particle set is refined through a series of annealing cycles with decreasing temperature to approximate the local maxima in the fitness function. Results [85,90] demonstrate reconstruction of complex motion such as a hand-stand. A hierarchical stochastic sampling scheme to efficiently estimate the 3D pose for complex movements or multiple people is presented in [242]. This approach initially estimates the torso pose for each person and propagates samples with high fitness to estimate the pose of adjacent body parts.

Recent work has combined deterministic or stochastic search with gradient descent for local pose refinement to recover complex whole-body motion. Carranza *et al.* [53] demonstrate whole-body human motion estimation from multiple views combining a deterministic grid search with gradient descent. Pose estimation is performed hierarchically starting with the torso. For each body part a grid search first finds the set of valid poses for which the joint positions project inside the observed silhouettes. A fitness function is then evaluated for all valid poses to determine the best pose estimate. Finally gradient descent optimization is performed to refine the estimated pose. This search procedure is made feasible by the use of graphics hardware to evaluate the fitness function which is based on the overlap between

the projected model and observed silhouette across all views. In related work Kehl *et al.* [190] propose *stochastic meta descent* for whole-body pose estimation with 24 degrees-of-freedom from multiple views. Stochastic meta descent combines a stochastic sampling of the set of model points used at each iteration of a gradient descent algorithm. This introduces a stochastic search element to the optimization which allows the approach to avoid convergence to local minima. The use of a small number of samples (5) per body part together with adaptive step size allows efficient performance. Results of these approaches demonstrate reconstruction of complex movements such as kicking and dancing.

In summary, the introduction of stochastic sampling and search techniques has achieved whole-body pose estimation of complex movements from multiple views. Current approaches are limited to gross-body pose estimation of torso, arms and legs and do not capture detailed movement such as hand-orientation or axial arm rotation. Multiple hypothesis sampling achieves robust tracking but does not provide a single temporally consistent motion estimate resulting in jitter which must be smoothed to obtain visually acceptable results. There remains a substantial gulf between the accuracy of commercial marker-based and marker less video-based human motion reconstruction.

4.3.2 Monocular 3D Pose Estimation

Reconstruction of human pose from a single view image sequence is considerably more difficult than either the problem of 2D pose estimation or 3D pose estimation from multiple views. To resolve the inherent ambiguity in monocular human motion reconstruction additional constraints on kinematics and movement are typically employed [43,384]. Wachter and Nagel [384] used the extended Kalman filter together with kinematic joint constraints to estimate the 3D motion of a person walking parallel to the image plane. As discussed in the previous section the use of a single hypothesis tracking scheme is prone to failure for complex motions. Loy *et al.* [224] employ a manual key-frame approach to 3D pose estimation of complex motion in sports sequences.

Sminchisescu and Triggs [343] have investigated the application of stochastic sampling to estimation of 3D pose from monocular image sequences. They observe that alternative 3D poses which give good correspondence to the observations are most likely to occur in the direction of greatest uncertainty. This motivated the introduction of *covariance scaled sampling* an extension of particle filters which increases the covariance in the direction of maximum uncertainty by approximately an order of magnitude to increase the probability of generating samples close to local minima in the fitness function. Samples are then optimized to find the local minima using a gradient descent approach. Results demonstrate monocular tracking and 3D reconstruction of human movements with moderate complexity including walking with changes in direction. Further research [344] has explicitly enumerated the po-

tential kinematic minima which cause visual ambiguities. Incorporating this in the sampling process increases efficiency and robustness allowing reconstruction of more complex human motion from monocular video sequences.

Probabilistic approaches using assemblies of parts together with higher level knowledge of human kinematics and shape have also been investigated for single view 3D pose estimation. Lee and Cohen [210] combine a probabilistic proposal map representing the estimated likelihood of body parts in different 3D locations with an explicit 3D model to recover the 3D pose from single image frames. A data driven Markov chain Monte Carlo (MCMC) is used to search the high-dimensional pose space. The proposal map for each body part represents the likelihood of the projected 3D pose. Proposal distributions are used to efficiently sample the pose space during MCMC search. Results demonstrate 3D pose estimation from static sports players in a variety of complex poses. Moeslund and Granum [246,252] apply a data driven sequential Monte Carlo approach to pose estimation of a human arm. A part detector provides likely locations of the hand in the image and their uncertainties. This information is applied to correct the prediction lowering the number of particles required.

Navaratnam *et al.* [265] combine a hierarchical kinematic model with a bottom up part detection to recover the 3D upper-body pose. The use of part detection allows individual body parts to be independently located at each frame. Kinematic constraints between body parts are represented hierarchically to recover the 3D pose from a single view. Unlike previous model free probabilistic assembly of parts this approach enables recovery of full 3D pose at each frame. Temporal information is also integrated using a HMM framework to reconstruct temporally coherent movement sequences.

Monocular reconstruction of complex 3D human movement remains an open problem. Recent research has investigated the use of learnt motion models to provide strong priors to constrain the search.

4.3.3 *Learnt Motion Models*

There has been increasing interest in the use of learnt models of human pose and motion to constrain vision-based reconstruction of human movement from single or multiple views. The availability of marker-based human motion capture data [1,2,4] has led to the use of learnt models of human motion for both animation synthesis in computer graphics and vision-based human motion synthesis.

Learnt models have been developed in computer animation to allow synthesis of natural motions with user specified constraints from a motion capture database [20,198,208,255]. This use of learnt models in computer graphics is relevant to the problem of vision-based reconstruction of human movement in developing methods to predict and constrain human pose and motion estimation. Inverse kinematics

of human motion based on learnt models has recently been introduced in computer graphics [127,271]. Ong *et al.* [271] use a learnt model of whole-body configurations to constrain the pose given a set of end effector positions for a motion sequence. Grochow *et al.* [127] use Scaled Gaussian Process Latent Variable Models (SGPLVM) to model the probability distribution over all possible whole-body poses to constrain both character pose in animation and pose reconstruction from images.

Sidenbladh *et al.* [332,334,335] combine stochastic sampling with a strong learned prior of walking motion for tracking. An exemplar based approach is used in [335] similar to work in motion synthesis [20,198,255] where a database of motion capture examples is indexed to obtain possible movement directions. Statistical priors on human appearance and image motion are used [333] to model the likelihood of observing various image cues for a given movement. These are incorporated in an analysis-by-synthesis approach to human motion reconstruction. Similarly, a hierarchical PCA model of human dynamics learnt from motion capture using a Gaussian mixture and HMM to represent dynamics is proposed for monocular tracking in [188]. Agarwal and Triggs [7] use a learned model of local second order dynamics for 2D tracking of more general motions walking and running with transitions and turns in monocular image sequences. Their work demonstrates that strong priors on human dynamics allows 2D pose estimation for fast movements in cluttered scenes.

Subsequent research has investigated the use of learnt motion models for 3D motion reconstruction primarily from monocular image sequences to overcome the inherent visual ambiguity. In [153] learnt models from short motion sequences are used to infer 3D pose from tracked image features of simple movements. Sigal *et al.* [336] combine body part detectors with a learned motion model to infer 3D human pose from monocular images of walking with automatic initialization. Their approach uses belief propagation via stochastic sampling over a loopy graph of loosely attached body parts. Urtasun and Fua [373] introduce the use of temporal motion models learnt from sequences of motion capture data to reconstruct human motion using a deterministic gradient descent optimization. Principal component analysis (PCA) is performed on multiple examples of concatenated joint angle sequences for walking and running to provide a low-dimensional parametrization. The parametric motion model is then used to constrain the movement of a 3D humanoid model for walking and running movements with variable speed from stereo [373] and golf swings from a single view [370]. Urtasun *et al.* [372] advocate an alternative approach to representation of human motion using SGPLVM to learn a low-dimensional embedding of the pose state space for specific movements such as golf-swings or walking from monocular image sequences. Gaussian Process models which incorporate dynamics [259,371] have been introduced to ensure continuous embedding of motion in the latent space for robust tracking. Further research following the methodology of using learnt motion models has addressed the problem of viewpoint invariance in tracking human movement [8,272].

Research introducing the use of learnt statistical models of human motion since 2000 has demonstrated that using strong motion priors facilitates reconstruction of 3D pose sequences from monocular images. To date the generality of these approaches has been limited to specific motion models with relatively small variation in motion and fixed transitions. A challenge for future research is to build more general motion models or methods of transitioning between models, to allow the reconstruction of unconstrained human movement.

5 Discussion of Advances in Human Pose Estimation

As identified in this section research in automatic estimation of human pose has been an active area over the past five years with significant advances being made. A number of novel methodologies have been proposed towards the objective of human pose estimation from monocular image sequences in natural scenes. The introduction of methods based on 2D pose estimation as a probabilistic assembly of parts have achieved significant advances for cluttered natural scenes such as film footage or sports [104,157,167,235,296,302,310,314]. These approaches are based on detection of body parts such as the face, hands or limbs independently for each image frame.

Similarly there have been significant advances in the use of example-based methods to learn the mapping from 2D image features such as silhouettes to 3D pose [6,41,151,326,340]. These methods commonly exploit databases of human motion capture data to render images of a model in multiple poses providing known 2D image to 3D pose correspondence. Currently example-based methods are limited to the fixed classes of movement and range of viewpoints used in training. A future challenge is to extend these methods to viewpoint invariant 3D pose estimation for general movement. There is also the possibility of combining learnt 2D to 3D mappings with 2D pose detection to achieve 3D pose detection in cluttered scenes from monocular image sequences or single image frames.

Model-based pose estimation using an analysis-by-synthesis methodology to estimate 3D pose from multiple view images has focused on reliable recovery of complex movements [53,90,190]. Significant advances in the complexity of movement that can be reconstructed have been achieved through the use of stochastic sampling and search techniques in pose estimation from multiple views. Similarly research in 3D pose estimation from monocular image sequences using stochastic sampling [343] has achieved reconstruction in cluttered scenes. Monocular reconstruction of complex 3D human movement remains an open problem. Learnt models of human motion have been applied extensively to constrain the monocular reconstruction problem by providing strong priors on motion [7,334,336,372]. Currently learnt motion models are limited to specific classes of motion. The extension of learnt models to reconstruction of general human movement remains an open problem.

Over the past five years there have been significant advances in the range of human motion which can be reconstructed from either monocular or multiple view image sequences. A limitation of existing research which should be addressed in future is the comparison of different approaches on common data sets and performance evaluation of accuracy against ground-truth.

6 Recognition

The field of action and activity representation and recognition is relatively old, yet still immature. This area is presently subject to intense investigation which is also reflected by the large number of different ideas and approaches. The approaches depend on the goal of the researcher and applications for activity recognition are interesting for surveillance, medical studies and rehabilitation, robotics, video indexing and animation for film and games. For example, in scene interpretation the knowledge is often represented statistically and is meant to distinguish “regular” from “irregular” activities.

In scene interpretation, the representations should be independent from the objects causing the activity and thus are usually not meant to distinguish explicitly, e.g, cars from humans. On the other hand, some surveillance applications focus explicitly on human activities and the interactions between humans. Here, one finds both, holistic approaches, that take into account the entire human body without considering particular body parts, and local approaches. Most holistic approaches attempt to identify “holistic” information such as gender, identity or simple actions like walking or running. Researchers using local approaches appear often to be interested in more subtle actions or attempt to model actions by looking for action primitives with which the complex actions can be modeled.

We have structured this review according to a visual abstraction hierarchy yielding the following: *scene interpretation* where the entire image is interpreted without identifying particular objects or humans, *holistic recognition* where either the entire human body or individual body parts are applied for recognition, and *action primitives and grammars* where an action hierarchy gives rise to a semantic description of a scene. Before going into these topics we first look closer at the definition of the action hierarchy used in this survey since it has influence on the remaining categories.

6.1 Action Hierarchies

Terms like *actions*, *activities*, *complex actions*, *simple actions* and *behaviors* are often used interchangeably by the different authors. However, in order to be able

to describe and compare the different publications we see the need for a common terminology. In a pioneering work [264], Nagel suggested to use a hierarchy of *change, event, verb, episode, history*. An alternative hierarchy (reflecting the computational aspects) is proposed by Bobick [37] who suggests to use *movement, activity and action* as different levels of abstraction (see also [12]). Others suggest to also include *situations* [120] or use a hierarchy of *Action primitives* and *Parent Behaviors* [174].

In this survey we will use the following action hierarchy: *action/motor primitives, actions* and *activities*. *Action primitives* or *motor primitives* will be used for atomic entities out of which actions are built. *Actions* are, in turn, composed into *activities*. The granularity of the primitives often depends on the application. For example, in robotics, *motor primitives* are often understood as sets of motor control commands that are used to generate an action by the robot (see section 6.5).

As an example, in tennis *action primitives* could be, e.g., “forehand”, “backhand”, “run left”, “run right”. The term *action* is used for a sequence of action primitives needed to return a ball. The choice of a particular action depends on whether a forehand, backhand, lob or volley etc, is required in order to be able to return the ball successfully. Most of the research discussed below fall into this category. The *activity* then is in this example “playing tennis”. *Activities* are larger scale events that typically depend on the context of the environment, objects or interacting humans.

A good overview of activity recognition is given by Aggarwal and Park [12]. They aim at higher-level understanding of activities and interactions and discuss different aspect such as level of detail, different human models, recognition approaches and high-level recognition schemes. Veeraraghavan *et al.*[379] discuss the structure of an action and activity space.

6.2 Scene Interpretation

Many approaches consider the camera view as a whole and attempt to learn and recognize activities simply by observing the motion of objects without necessarily knowing their identity. This is reasonable in situations where the objects are small enough to be represented as points on a 2D plane.

Stauffer *et al.*[355] present a full scene interpretation system which allows detection of unusual situations. The system extracts features such as 2-D position and speed, size and binary silhouettes. Vector Quantization is applied to generate a codebook of K prototypes. Instead of taking the explicit temporal relationship between the symbols into account, Stauffer and Grimson use co-occurrence statistics. Then, they define a binary tree structure by recursively defining two probability mass functions across the prototypes of the code book that best explain the co-occurrence matrix. The leaf nodes of the binary tree are probability distributions of

co-occurrences across the prototypes and at a higher tree depth define simple scene activities like pedestrian and car movement. These can then be used for scene interpretation. In Eng *et al.* [101] a swimming pool surveillance system is presented. From each of the detected and tracked objects features such as speed, posture, submersion index, an activity index and a splash index, are extracted. These features are fed into a multivariate polynomial network in order to detect water crisis events. Boiman and Irani [39] approach the problem of detection irregularities in a scene as a problem of composing newly observed data using spatio-temporal patches extracted from previously seen visual examples. They extract small image and video patches which are used as local descriptors. In an inference process, they search for patches with a similar geometric configuration and appearance properties, while allowing for small local misalignments in their relative geometric arrangement. This way, they are able to quickly and efficiently infer subtle but important local changes in behavior. Junejo *et al.*[181] describe an approach to focusses on dynamic information for scene interpretation. Their method can distinguish between objects traversing spatially dissimilar paths or objects traversing spatially proximal paths but with different spatio-temporal characteristics. For this, they learn the paths in a training phase where graph-cuts are used for clustering the trajectories. For matching, they use spatial similarity, velocity characteristics and curvature features.

In [64,375] activity trajectories are modeled using non-rigid shapes and a dynamic model that characterizes the variations in the shape structure. Vaswani *et al.* [375] uses Kendall's statistical shape theory [191]. Nonlinear dynamical models are used to characterize the shape variation over time. An activity is recognized if it agrees with the learned parameters of the shape and associated dynamics. Chowdhury *et al.* [63] use a subspace method to model activities as a linear combination of 3D basis shapes. The work is based on the factorization theorem [365]. Deviations from the learned normal activity shapes can be used to identify abnormal ones.

A similar complex task is approached by Xiang and Gong [400]. They present a unified bottom-up and top-down approach to model complex activities of multiple objects in cluttered scenes. Their approach is object-independent and they use a Dynamically Multi-Linked Hidden Markov Models (HMMs) on conjunction with Schwarz's Bayesian Information Criterion [324] to interlink between multiple temporal processes corresponding to multiple event classes. Liu and Chua [222] present an HMM-based approach for recognizing multi-agent activities.

6.3 Holistic Recognition Approaches

The recognition of the identity of a human, based on his/her global body structure and the global body dynamics is discussed in many publications. Of particular interest for identity recognition has been the human gait. Other approaches using global body structure and dynamics are concerned with the recognition of simple actions

such as running and walking. Almost all methods are silhouette or contour based. Subsequent techniques are mostly holistic, e.g., the entire silhouette or contour is being taken into account without detecting individual body parts.

6.3.1 Human Body Based Recognition of Identity

In Wang *et al.* [390] the silhouette of a human is computed and then unwrapped by evenly sampling the contour. Next, the distance between each contour point and its center of gravity is computed. The unwrapped contour is then processed by PCA. To compute the spatio-temporal correlation they compare trajectories in eigenspace by first applying appropriate time warping to minimize the distance between the probe and the gallery trajectories. On outdoor data and in spite of its simplicity, it gives good results while being computationally efficient. BenAbdelkader *et al.* [32] use a variation of co-occurrence techniques. After applying a suitable time-warping and normalization with respect to scale a self-similarity plot is computed where silhouette images of the sequences are pairwise correlated. PCA is applied to reduce the dimensionality of these plots and a k -nearest neighbor classifier is applied in eigenspace for recognition.

Foster *et al.* [110] extract, emboss and normalize silhouettes. Then, a set of binary masks are defined and the area of the silhouette within the mask is computed to give a dynamic signature of the observed person for each mask. A frame rate of 30 fps results in a 30-D vector for each signature giving a $n \times 30$ matrix where n denotes the number of area masks used. To remove the information about the static shape of the silhouette, the average value of each signature can be subtracted. Fisher analysis is applied and the k -nearest neighbor classifier is used for classification. Kale *et al.* [182,183] define a HMM to model the dynamics of individual gait. A HMM is trained for each individual in the database. Five representative binary silhouette are used as the hidden states for which transition probabilities and observation likelihoods are trained. During the recognition phase, the HMM with the largest probability identifies the individual. Yam *et al.* [402] investigate the relationship between walking and running. They define a gait signature based on a frequency analysis of thigh and lower leg rotations. Phase and magnitude of the Fourier descriptions are multiplied to give the phase-weighted magnitude (PWM). It appears that the signatures for walking and running for an individual is related by a phase modulation. The additional individual relationship between walking and running is used to derive improved gait-recognition which can recognize both, walking and running patterns.

6.3.2 Human Body Based Recognition

While a large number of papers recognize individuals based on their dynamics, the dynamics can also be used to recognize *what* the individual is doing. The ap-

proaches discussed in this subsection are again based on holistic body information where no attempt is made to identify individual body parts.

A pioneering work in this context has been presented by Efros *et al.* [94]. They attempt to recognize simple actions of people whose images in the video are only 30 pixels tall and where the video quality is poor. They use a set of features that are based on blurred optic flow (blurred motion channels). First, the person is tracked so that the image is stabilized in the middle of a tracking window. The blurred motion channels are computed on the residual motion that is due to the motion of the body parts. Spatio-temporal cross-correlation is used for matching with a database. Roh *et al.* [312] base their action recognition task on curvature scale space templates of a player's silhouette.

Of further interest is the enhancement where complex actions can be dynamically composed out of the set of simple actions. Robertson and Reid [311] attempt to *understand* actions by building a hierarchical system that is based on reasoning with belief networks and HMMs on the highest level and on the lowest level with features such as position and velocity as action descriptors. Their action descriptor is based on [94]. The system is able to output qualitative information such as *walking – left-to-right – on the sidewalk*.

A large number of publications work with space-time volumes. One of the main approaches is to use spatio-temporal XT -slices from an image volume XYT [304,305] where articulated motions of a human can be associated with a typical trajectory pattern. Ricquebourg and Bouthemy [304] demonstrate how XT -slices can facilitate tracking and reconstruction of 2D motion trajectories. The reconstructed trajectory allows a simple classification between pedestrians and vehicles. Ritscher *et al.* [305] discuss the recognition in more detail by a closer investigation of the XT -slices. Quantifying the braided pattern in the slices of the spatio-temporal cube gives rise to a set of features (one for each slice) and their distribution is used to classify the actions.

Bobick and Davis pioneered the idea of temporal templates [37,38]. They propose a representation and recognition theory [37,38] that is based on *motion energy images* (MEI) and *motion history images* (MHI). The MEI is a binary cumulative motion image. The MHI is an enhancement of the MEI where the pixel intensities are a function of the motion history at that pixel. Matching temporal templates is based on Hu moments. Bradski *et al.* [40] pick up the idea of MHI and develop timed MHI (tMHI) for motion segmentation. tMHI allow determination of the normal optical flow. Motion is segmented relative to object boundaries and the motion orientation. Hu moments are applied to the binary silhouette to recognize the pose. A work conceptually related to [38] is by Masound and Papanikolopoulos [231]. Here, motion information for each video frame is represented by a feature image. However, unlike [38], an action is represented by several feature images. PCA is applied for dimensionality reduction and each action is then represented by a man-

ifold in PCA space.

Yi *et al.* [409] present the idea of a pixel change ratio map (PCRM) which is conceptually similar to the MHI. However, further processing is based on motion histograms which are computed from the PCRM. Weinberg *et al.* [393] suggest replacing the motion history image by a 4D motion history volume. For this, they first compute the visual hull from multiple cameras. Then, they consider the variations around the central vertical axes and use cylindrical coordinates to compute alignments and comparisons. Motion history images can also be used to detect and interpret actions in compressed video data. Babu and Ramakrishnan[23] a motion compute a flow history (MFH) from the motion data available in compressed video. In addition to MFH, they also use motion history images to classify activities.

As the search of activities in large databases gains importance, a full, hierarchical human detection system is presented by Ozer and Wolf [275]. They approach the tracking, pose estimation and action recognition problem in an integrated manner. They apply a number of well-known techniques on (un)compressed video data.

Another approach is that of “Actions Sketches” or “Space-Time Shapes” in the 3D XYT volume. Yilmaz and Shah [410] propose to use spatio-temporal volumes (STV) for action recognition: The 3D contour of a person gives rise to a 2D projection. Considering this projection over time defines the STV. Yilmaz and Shah extract information such as speed, direction and shape by analyzing the differential geometric properties of the STV. They approach action recognition as an object matching task by interpreting the STV as rigid 3D objects. Blank *et al.* [36] also analyze the STV. They generalize techniques for the analysis of 2D shapes [122] for the use on the STV. Blank *et al.* argue that the time domain introduces properties that do not exist in the xy -domain and needs thus a different treatment. For the analysis of the STV they utilize properties of the solution of the Poisson equation [122]. This gives rise to local and global descriptors that are used for recognizing simple actions.

Instead of using spatio-temporal volumes, a large number of papers choose the more classical approach of considering sequences of silhouettes. Yu *et al.*[412] extract silhouettes and their contours are unwrapped and processed by PCA. A three-layer feed forward network is used to distinguish “walking”, “running” and “other” based on the trajectories in eigenspace. In another PCA-based approach, Rahman and Robles-Kelly [295] suggest to use a tuned eigenspace technique. They tuned eigenspaces allow to treat the action problem as a nearest-neighborhood problem in eigenspace. Jiang *et al.*[178] attempt to match a given sequence of poses to a novel video. They treat this problem as an optimal matching problem by changing the usually highly non-convex problem in to a convex one.

Elgammal and Lee [13] use optic flow in addition to the shape features and a HMM is used to model the dynamics. In [98,97], Elgammal and Lee use local linear em-

bedding (LLE) [317,362] in order to find a linear embedding of human silhouettes. In conjunction with a generalized radial basis function interpolation, they are able to separate style and content of the performed actions [97] as well as to infer 3D body pose from 2D silhouettes [98]. Sato and Aggarwal [321] are concerned with the detection of interaction between two individuals. This is done by grouping foreground pixels according to similar velocities. A subsequent tracker tracks the velocity blobs. The distance between two people, the slope of relative distance and the slope of each person's position are the features used for interaction detection and classification. In Cheng *et al.* [58], walking is distinguished from running based on sport event video data. The data comes from real-life programs. They compute a dense motion field and foreground segmentation is performed based on color and motion. Within the foreground region, the mean motion magnitude between frames is computed over time followed by an analysis in frequency space to compute a characteristic frequency. A Gaussian classifier is used for classification. Gao *et al.* [112] consider a smart room application. A dining room activity analysis is performed by combining motion segmentation with tracking. They use motion segmentation based on optical flow and RANSAC. Then, they combine the motion segmentation with a tracking approach which is sensitive to subtle motion. In order to identify activities, they identify predominant directions of relative movements.

In a number of publications, recognition is based on HMMs and dynamic Bayes networks (DBNs). Elgammal *et al.* [99] propose a variant of semi-continuous HMMs for learning gesture dynamics. They represent the observation function of the HMM as non-parametric distributions to be able to relate a large number of exemplars to a small set of states. Luo *et al.* [226] present a scheme for video analysis and interpretation where the higher-level knowledge and the spatio-temporal semantics of objects are encoded with DBNs. The DBNs are based on key-frames and are defined for video objects. Shi *et al.* [330] present an approach for semi-supervised learning of the HMM or DBN states to incorporate prior knowledge. Leo *et al.* [216] attempt to classify actions at an archaeological site. They present a system that uses binary patches and an unsupervised clustering algorithm to detect human body postures. A discrete HMM is used to classify the sequences of poses into a set of four different actions.

Smith *et al.* [347] suggest to use multiple levels of zoom for activity analysis to combine both detailed and coarse views of a scene. They find feature correspondences across different zoom levels using epipolar, spatial, trajectory and appearance constraints.

A totally different approach is presented by Wang *et al.* [392] where the aim is at classifying actions in *still* images. Unsupervised learning is used to generate action classes out of a large training set. These action classes are then used to label test images. The approach uses a technique for deformable matching of edges of image pairs, based on linear programming relaxation techniques.

6.4 Recognition Based on Body Parts

Many authors are concerned with the recognition of actions based on the dynamics and settings of individual body parts. Some approaches, e.g., [83], start out with silhouettes and detect the body parts using a method inspired by the W4-system [137]. Others use 3D-model based body tracking approaches (see section 4) where the recognition of (often periodic) action is used as a loop-back to support pose estimation. Other approaches circumvent the vision problem by using a motion capture system in order to be able to focus on the action issues [81,277,279].

In a work related to [390], Wang *et al.* [389] present an approach where contours are extracted and a mean contour is computed to represent the static contour information. Dynamic information is extracted by using a detailed model composed of 14 rigid body parts, each one represented by a truncated cone. Particle filtering is used to compute the likelihood of a pose given an input image. For classification, a nearest neighbor classifier (NN) was used.

Davis and Taylor [83] present an approach to distinguish walking from non-walking. A method based on the W4-system is used to detect body parts from silhouettes. Based on the feet locations four motion properties are extracted of which three (cycle time, stance/swing ratio, double support time) reflect dynamic features and one (extension angle) reflects a structural feature. The walking category is defined by three pairs of the dynamic features and the structural feature. In a similar approach Ren and Xu [300] use as input a binary silhouette from which they detect the head, torso, hands and elbow angles. Then, a primitive-based coupled HMM is used to recognize natural complex and predefined actions. They extend their work in [301] by introducing primitive-based DBNs. Parameswaran and Chellappa [277,279] consider the problem of view-invariant action recognition based on point-light displays by investigating 2D and 3D invariant theory. As no general, non-trivial 3D-2D invariants exist, Parameswaran and Chellappa employ a convenient 2D invariant representation by decomposing and combining the patches of a 3D scene. For example, key poses can be identified where joints in the different poses are aligned. In the 3D case, six-tuples corresponding to six joints give rise to 3D invariant values and it is suggested to use the progression of these invariants over time for action representation. A similar issue is discussed in the work by Yilmaz and Shah [411] where joint trajectories from several uncalibrated moving cameras are considered. They propose an extension to the standard epipolar geometry based approach by introducing a temporal fundamental matrix that models the effects of the camera motion. The recognition problem is then approached in terms of the quality of the recovered scene geometry. Gritai *et al.* [126] address the invariant recognition of human actions, and investigate the use of anthropometry to provide constraints on matching. Gritai *et al.* use the constraints to measure the similarity between poses and pose sequences. Their work is based on a point-light display like representation where a pose is presented through a set of points in 3D space. Sheikh

et al.[328] pick up these results of [126,411] and discuss that the three most important sources of variability in the task of recognizing actions come from variations in viewpoint, execution rate and anthropometry of the actors. Then, they argue that the variability associated with the execution of an action can be closely approximated by a linear combination of action bases in joint spatio-temporal space. Davis' and Gao's [79,81] aim is to recognize properties from visual target cues, e.g. the sex of an individual or the weight of a carried object is estimated from how the individuals move. Davis and Gao[81] recognize the gender of a person based on the gait. Labeled 2D trajectories from motion capture devices of humans are factored using three-mode PCA into components interpreted as *posture*, *time* and *gender*. An importance weight for each of the trajectories is learned automatically. Davis *et al.*[79] use the three-mode PCA framework to recognize human action efforts. Here, the three modes *pose*, *time* and *effort* are used. In order to detect particular body parts Fanti *et al.* [103] give the structure of a human as model knowledge. To find the most likely model alignment with input data they exploit appearance information which remains approximately invariant within the same setting. Expectation maximization is used for unsupervised learning of the parameters and structure of the model for a particular action and unlabeled input data. Action is then recognized by maximum likelihood estimation. Ning *et al.*[267] use a parabola to model the shoulders of a human. Fisher Discriminant Analysis (FDA) on the parabola parameters are used to detect shrugs.

6.5 Action Primitives and Grammars

There is strong neurobiological evidence that human actions and activities are directly connected to the motor control of the human body [116,307,308]. When viewing other agents performing an action, the human visual system seems to relate the visual input to a sequence of motor primitives. The neurobiological representation for visually perceived, learned and recognized actions appears to be the same as the one used to drive the motor control of the body. These findings have gained considerable attention from the robotics community [77,322]. In *imitation learning* the goal is to develop a robot system that is able to relate perceived actions to its own motor control in order to learn and to later recognize and perform the demonstrated actions. Consequently, it is ongoing research to identify a set of motor primitives that allow a) representation of the visually perceived action and b) motor control for imitation. In addition, this gives rise to the idea of interpreting and recognizing activities in a video scene through a hierarchy of primitives, simple actions and activities. Most of the following researchers attempt to learn the motor or action primitives by defining a "suitable" representation and then learning the primitives from demonstrations. The representations used to describe the primitives vary a lot across the literature and are subject to ongoing research. Most of the subsequently mentioned work is based on motion capture data.

Jenkins *et al.* [175,176] suggest apply a spatio-temporal non-linear dimension reduction technique on manually segmented human motion capture data. Similar segments are clustered into primitive units which are generalized into parameterized primitives by interpolating between them. In the same manner, they define action units (“behavior units”) which can be generalized into actions. Ijspeert *et al.*[164] approach the problem of defining motor primitives from the motor side. They define a set of nonlinear differential equations that form a control policy (CP) and quantify how well different trajectories can be fitted with these CPs. The parameters of a CP for a primitive movement are learned in a training phase. These parameters are also used to compute similarities between movements. Billard and Calinon [34,50,51] use an HMM based approach to learn characteristic features of repetitively demonstrated movements. They suggest to use the HMM to synthesize joint trajectories of a robot. For each joint, one HMM is used. Calinon *et al.*[51] use an additional HMM to model end-effector movement. In these approaches, the HMM structure is heavily constrained to assure convergence to a model that can be used for synthesizing joint trajectories.

A number of publications attempt to decouple actions into action primitives and to interpret actions as a composition on the alphabet of these action primitives, however, without the constraints of having to drive a motor controller with the same representation. Vecchio and Perona [376] employ techniques from the dynamical systems framework to approach segmentation and classification. System identification techniques are used to derive analytical error analysis and performance estimates. Once, the primitives are detected an iterative approach is used to find the sequence of primitives for a novel action. Another approach in this context is presented by Bissacco [35]. They extract some temporal statistics from the images and use them to build a dynamical system that models contact forces explicitly. Then, they explicitly factor out exogenous inputs that are not unique to an individual.

Lu *et al.* [225] also approach the problem from a system theoretic point of view. Their goal is to segment and represent repetitive movements. For this, they model the joint data over time with a second order auto-regressive (AR) model and the segmentation problem is approached by detection significant changes of the dynamical parameters. Then, for each motion segment and for each joint, they model the motion with a damped harmonic model. In order to compare actions, a metric based on the dynamic model parameters is defined. A different problem is studied by Wang *et al.*[387] addressing what kind of cost function should be used to assure smooth transitions between primitives.

While most scientists concentrate on the action representation by circumventing the vision problem, Rao *et al.*[298] take a vision-based approach. They propose a view-invariant representation of action based on *dynamic instants* and *intervals*. Dynamic instants are used as primitives of actions which are computed from discontinuities of 2D hand trajectories. An interval represents the time period between two dynamic instants (key poses). A similar approach of using meaningful instants

in time is proposed by Reng *et al.* [303] where key poses are found based on the curvature and covariance of the normalized trajectories. Cuntoor *et al.*[72] find key poses through evaluation of anti-eigenvalues.

González *et al.*[120] employ the point distribution model [68] to model the variability of joint angle settings of a stick figure model. An action spaces, *aSpace*, is trained by giving a set of joint angle settings coming from different individuals but showing the same action. *aSpaces* are then used for synthesis and recognition of known actions. Modeling of activities on a semantic level has been attempted by Park and Aggarwal[281]. The system they describe has 3 abstraction levels. At the first level, human body parts are detected using a Bayesian network. At the second level, DBNs are used to model the actions of a single person. At the highest level, the results from the second level are used to identify the interactions between individuals. Ivanov and Bobick [169] suggest using stochastic parsing for a semantic representation of an action. They discuss that for some activities, where it comes to semantic or temporal ambiguities or insufficient data, stochastic approaches may be insufficient to model complex actions and activities. They suggest decoupling actions into primitive components and using a stochastic parser for recognition. In [169] they pick up a work by Stolcke [356] on syntactic parsing in speech recognition and enhance this work for activity recognition in video data. Yamamoto *et al.*[403] present an application where a stochastic context free grammar is used for action recognition. A somewhat different approach is taken by Yu and Yang [413]. They use neural networks to find primitives. They apply self-organizing maps (SOMs, Kohonen's feature maps [196]) which cluster the training images based on shape feature data. After training the SOMs generated a label for each input image which converts an input image sequence into a sequence of labels. A subsequent clustering algorithm allows to find repeatedly appearing substructures in these label sequences. These substructures are then interpreted as motion primitives. A very interesting approach is presented by Lv and Nevatia in [229] where the authors are interested in recognizing and segmenting full-body human action. Lv and Nevatia decompose the large joint space into a set feature spaces where each feature corresponds to a single joint or combinations of related joints. They use then HMMs to recognize each action class based on the features and an AdaBoost scheme to detect and recognize the features.

6.6 Discussion of Advances in Human Action Recognition

The field of recognizing human actions has received a considerable increase of attention in the last few years. It is apparent from the published works, that the major interest lies in the field of surveillance and the related action understanding problems. While in some publications, the actions are interpreted without explicitly considering humans, others discuss the dynamics of humans, explicitly. In the latter, a large attention is devoted to rather simple actions such as walking, running,

sitting. Here, only a small body of literature goes beyond these simple actions into motion interpretation where scene context and the interaction with other humans is considered, e.g., [169,281,311,400,403]. Much more work is expected to appear in this context and the approaches will be interesting as they are likely to bridge the traditional vision field with the field of artificial intelligence.

On the other hand, a good understanding of these simple actions is necessary before they can be combined into more complex ones. The issues lie, e.g., in the invariances with respect to viewing angle, speed, and variations between individuals [98,328,379].

Another significant part of the discussed articles draw some of their motivations from neuroscientific studies [116,307,308] and deal explicitly with action primitives, action grammars [169,281,403] and the close relationship between action recognition and action synthesis [34,50,51,77,322]. As these works also build on action primitives a better understanding of action primitives is necessary also in this context, e.g., in order to generalize the HMMs as proposed by Billard and Calinon [34,50,51].

7 Conclusion

Over the past five years vision-based human motion estimation and analysis has continued to be a thriving area of research. This survey has identified over three-hundred related publications over the period 2000-06 in major conferences and journals. Increased activity in this research area has been driven by both the scientific challenge of automatic scene interpretation and the demands of potential mass-market applications in surveillance, entertainment production and indexing visual media.

During this period there has been substantial progress towards automatic human motion tracking and reconstruction. Recognition of human motion has also become a central focus of research interest. Key advances identified in this review include:

Initialization: Automatic initialization of model shape, appearance and pose has been addressed in recent work [59,238]. A major advance is the introduction of methods for pose detection from static images [157,302,315,326] which potentially provide automatic initialization for human motion reconstruction.

Tracking: Surveillance applications have motivated research advances towards reliable tracking of multiple people in unstructured outdoor scenes. Advances in especially the use of appearance, shape and motion for figure-ground segmentation have increased reliability of detecting and tracking people with partial occlusion [154,193,244,269,280,316,404]. Probabilistic classification methods [193,232,280,285] and stochastic sampling [154,269,290,345,404,421] have

been introduced to improve the reliability of temporal correspondence during occlusion. Systems for self-calibrating and tracking across multiple cameras have been investigated [21,186,192,369]. There remains a gap between the state-of-the-art and robust tracking of people for surveillance in outdoor scenes.

Human motion reconstruction from multiple views: Significant progress has been made towards the goal of automatic reconstruction of human movement from video. The model-based analysis-by-synthesis methodology, pioneered in early work [147], has been extended with the introduction of techniques to efficiently search the space of possible pose configurations for robust reconstruction from multiple view video acquisition [53,90,190,238]. Current approaches capture gross body movement but do not accurately reconstruct fine detail such as hand movements or axial rotations.

Monocular human motion reconstruction: Progress has also been made towards human motion capture from single views with stochastic sampling techniques [210,265,332,343]. An increasing trend in monocular tracking has been the use of learnt motion models to constrain reconstruction based on movement [7,8,332,334,373,372]. Research has demonstrated that the use of strong a priori models enables improved monocular tracking of specific movements.

Pose estimation in natural scenes: A recent trend to overcome limitations of monocular tracking in video of unstructured scenes has been direct pose detection on individual frames. Probabilistic assemblies of parts based on robust body part detection has achieved 2D pose estimation in challenging cluttered scenes such as film footage [157,235,240,296,302,314]. Example based methods which learn a mapping from image to 3D pose space have been presented for reconstruction of specific movements [8,315,326].

Recognition: Understanding behavior and action has recently seen an explosion of research interest. Considerable steps have been made to advance surveillance applications towards automatic detection of unusual activities. Progress can also be seen for the recognition of simple actions and the description of action grammars. Relatively few papers have so far dealt with higher abstraction levels in action grammars which touch the border of semantics and AI. Association of actions and activities with affordances of objects will also bring a new perspective to object recognition.

Future research in visual analysis of human movement must address a number of open problems to satisfy the common requirements of potential applications for reliable automatic tracking, reconstruction and recognition. Body part detectors which are invariant to viewpoint, body shape and clothing are required to achieve reliable tracking and pose estimation in cluttered natural scenes. The use of learnt models of pose and motion are currently restricted to specific movements. More general models are required to provide constraints for capturing a wide range of human movement. Whilst there has been substantial advances in human motion reconstruction the visual understanding of human behavior and action remains immature despite a surge of recent interest. Progress in this area requires fundamental advances in behavior representation for dynamic scenes, viewpoint invariant re-

relationships for movement and higher level reasoning for interpretation of actions [325].

Industrial applications also require specific advances: human motion capture for entertainment production requires accurate multiple view reconstruction; surveillance applications require both reliable detection of people and recognition of movement and behavior from relatively low quality imagery; human-computer interfaces require low-latency real-time recognition of gestures, actions and natural behaviors. The potential of these applications will continue to inspire the advances required to realize reliable visual capture and analysis of moving people.

Acknowledgement

The authors would like to thank the following people for providing valuable comments to the paper: the anonymous reviewers, Prof. Larry S. Davis, Dr. Jordi González, Dr. Hedvig Kjellström (formerly Sidenbladh), Prof. Hans-Hellmut Nagel, Prof. Ramakant Nevatia, and Prof. Mubarak Shah.

Thomas B. Moeslund is supported by the Danish National Research Councils and HERMES (FP6 IST-027110). Adrian Hilton is supported by EPSRC GR/S13576 Visual Media Platform Grant. Volker Krüger is supported by PACO-PLUS (FP6 IST-IP-027657).

References

- [1] <http://www.charactermotion.com/products/powermoves/megamocap/>. Cited on page(s): 7, 28
- [2] <http://mocap.cs.cmu.edu/>. Cited on page(s): 7, 28
- [3] <http://www.eyetoy.com>. Cited on page(s): 3
- [4] <http://www.ict.usc.edu/graphics/animWeb>. Cited on page(s): 7, 28
- [5] <http://www.visualsurveillance.org>. Cited on page(s): 21
- [6] A. Agarwal and B. Triggs. 3D Human Pose from Silhouettes by Relevance Vector Regression. In *Computer Vision and Pattern Recognition*, Washington DC, USA, June, 2004. Cited on page(s): 24, 24, 30, 75
- [7] A. Agarwal and B. Triggs. Tracking Articulated Motion with Piecewise Learned Dynamic Models. In *European Conference on Computer Vision*, Prague, Czech Republic, May 11-14, 2004. Cited on page(s): 29, 30, 43, 75
- [8] A. Agarwal and B. Triggs. Recovering 3D Human Pose from Monocular Images. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 28(1):44–58, 2006. Cited on page(s): 24, 29, 43, 43, 77

- [9] J.K. Aggarwal and Q. Cai. Human Motion Analysis: A Review. *Computer Vision and Image Understanding*, 73(3):428–440, 1999. Cited on page(s): 4
- [10] J.K. Aggarwal, Q. Cai, W. Liao, and B. Sabata. Articulated and Elastic Non-Rigid Motion: A Review. In *Workshop on Motion of Non-Rigid and Articulated Objects*, Austin, Texas, Nov 1994. Cited on page(s): 4
- [11] J.K. Aggarwal, Q. Cai, W. Liao, and B. Sabata. Nonrigid Motion Analysis: Articulated and Elastic Motion. *Computer Vision and Image Understanding*, 70(2):142–156, 1998. Cited on page(s): 4
- [12] J.K. Aggarwal and S. Park. Human Motion: Modeling and Recognition of Actions and Interactions. In *Second International Symposium on 3D Data Processing, Visualization and Transmission*, Thessaloniki, Greece, September 6-9 2004. Cited on page(s): 4, 32, 32, 75
- [13] M. Ahmad and S. Lee. Human Action Recognition Using Multi-view Image Sequence Features. In *International Conference on Automatic Face and Gesture Recognition*, Southampton, UK, April 10-12, 2006. Cited on page(s): 36, 77
- [14] B. Allen, B. Curless, and Z. Popovic. Articulated Body Deformation from Range Scan Data. In *ACM SIGGRAPH*, pages 612—619, 2002. Cited on page(s): 7, 7, 9, 73
- [15] B. Allen, B. Curless, and Z. Popovic. The Space of Human Body Shapes: Reconstruction and Parameterization from Range Images. In *ACM SIGGRAPH*, pages 587—594, 2003. Cited on page(s): 7, 74
- [16] J. Ambrosio, J. Abrantes, and G. Lopes. Spatial Reconstruction of Human Motion by Means of a Single Camera and a Biomechanical Model. *Journal of Human Movement Science*, 20:829–851, 2001. Cited on page(s): 72
- [17] J. Ambrosio, G. Lopes, J. Costa, and J. Abrantes. Spatial Reconstruction of the Human Motion Based on Images of a Single Camera. *Journal of Biomechanics*, 34:1217–1221, 2001. Cited on page(s): 72
- [18] P.F. Andersen and R. Corlin. Tracking of Interacting People and Their Body Parts for Outdoor Surveillance. Master’s thesis, Laboratory of Computer Vision and Media Technology, Aalborg University, Denmark, 2005. Cited on page(s): 10, 12, 12, 12, 19, 19, 21, 76
- [19] G. Antonini, S.V. Martinez, M. Bierlaire, and J.P. Thiran. Behavioral Priors for Detection and Tracking of Pedestrians in Video Sequences. *International Journal of Computer Vision*, 69(2):159–180, 2006. Cited on page(s): 19, 19, 77
- [20] O. Arikan and D.A. Forsyth. Synthesizing Constrained Motions from Examples. In *ACM SIGGRAPH*, pages 483—490, 2002. Cited on page(s): 28, 29
- [21] N. Atsushi, K. Hirokazu, H. Shinsaku, and I. Siji. Tracking Multiple People using Distributed Vision Systems. In *International Conference on Robotics & Automation*, Washington DC, USA, May 2002. Cited on page(s): 17, 21, 43, 73
- [22] Y. Azoz, L. Devi, and M. Yeasin. Tracking the Human Arm using Constraint Fusion and Multiple-Cue Localization. *Machine Vision and Applications*, 13(5-6):286–302, 2003. Cited on page(s): 74
- [23] R.V. Babu and K.R. Ramakrishnan. Compressed Domain Human Motion Recognition using Motion History Information. In *International Conference on Acoustics, Speech and Signal Processing*, Hong Kong, China, April 6-10, 2003. Cited on page(s): 36, 74
- [24] A.O. Balan and M.J. Black. An Adaptive Appearance Model Approach for Model-based Articulated Object Tracking. In *Computer Vision and Pattern Recognition*, New York City, New York, USA, June 17-22, 2006. Cited on page(s): 77

- [25] A.O. Balan, L. Sigal, and M.J. Black. A Quantitative Evaluation of Video-based 3D Person Tracking. In *Workshop on Visual Surveillance and Performance Evaluation of Tracking and Surveillance*, Beijing, China, Oct 15-21, 2005. Cited on page(s): 76
- [26] C. Barron and I.A. Kakadiaris. Estimating Anthropometry and Pose from a Single Image. In *Computer Vision and Pattern Recognition*, Hilton Head Island, South Carolina, June 13-15 2000. Cited on page(s): 6, 71
- [27] C. Barron and I.A. Kakadiaris. Estimating Anthropometry and Pose from a Single Uncalibrated Image. *Computer Vision and Image Understanding*, 81(3):269–284, 2001. Cited on page(s): 72
- [28] C. Barron and I.A. Kakadiaris. On the Improvement of Anthropometry and Pose Estimation from a Single Uncalibrated Image. *Machine Vision and Applications*, 14(4):229–236, 2003. Cited on page(s): 6, 74
- [29] C. Beleznai, B. Fruhstuck, and H. Bischof. Tracking Multiple Humans using Fast Mean Shift Mode Seeking. In *Workshop on Performance Evaluation of Tracking and Surveillance*, Breckenridge, Colorado, Jan 2005. Cited on page(s): 19, 76
- [30] S. Belongie, J. Malik, and J. Puzicha. Matching Shapes. In *International Conference on Computer Vision*, Vancouver, Canada, July 9-12 2001. Cited on page(s): 16
- [31] J. Ben-Arie, Z. Wang, P. Pandit, and S. Rajaram. Human Activity Recognition Using Multidimensional Indexing. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 24(8):1091–1104, 2002. Cited on page(s): 73
- [32] C. BenAbdelkader, R. Cutler, and L. Davis. Motion-based Recognition of People in EigenGait Space. In *International Conference on Automatic Face and Gesture Recognition*, Washington DC, USA, May 20-21 2002. Cited on page(s): 34, 73
- [33] J. Berclaz, F. Fleuret, and P. Fua. Robust People Tracking with Global Trajectory Optimization. In *Computer Vision and Pattern Recognition*, New York City, New York, USA, June 17-22, 2006. Cited on page(s): 77
- [34] A. Billard, Y. Epars, S. Calinon, S. Schaal, and G. Cheng. Discovering Optimal Imitation Strategies. *Robotics and Autonomous Systems*, 47:69–77, 2004. Cited on page(s): 40, 42, 42, 75
- [35] A. Bissacco and S. Soatto. Classifying Human Dynamics Without Contact Forces. In *Computer Vision and Pattern Recognition*, New York City, New York, USA, June 17-22, 2006. Cited on page(s): 40, 77
- [36] B. Blank, L. Gorelick, E. Shechtman, M. Irani, and R. Basri. Actions as Space-Time Shapes. In *International Conference on Computer Vision*, Beijing, China, Oct 15-21, 2005. Cited on page(s): 36, 76
- [37] A. Bobick. Movement, Activity, and Action: The Role of Knowledge in the Perception of Motion. *Philosophical Trans. Royal Soc. London*, 352:1257–1265, 1997. Cited on page(s): 32, 35, 35
- [38] A. Bobick and J. Davis. The Recognition of Human Movement Using Temporal Templates. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 23(3):257–267, 2001. Cited on page(s): 35, 35, 35, 35, 72
- [39] O. Boiman and M. Irani. Detecting Irregularities in Images and in Video. In *International Conference on Computer Vision*, Beijing, China, Oct 15-21, 2005. Cited on page(s): 33, 76
- [40] G.R. Bradski and J.W. Davis. Motion Segmentation and Pose Recognition with Motion History Gradients. *Machine Vision and Applications*, 13(3):174–184, 2002. Cited on page(s): 13, 35, 72, 73

- [41] M. Brand. Shadow Puppetry. In *International Conference on Computer Vision*, Corfu, Greece, September 1999. Cited on page(s): 22, 24, 30
- [42] M. Bray, P. Kohli, and P.H.S. Torr. PoseCut: Simultaneous Segmentation and Pose Estimation of Humans Using Dynamic Graph-Cuts. In *European Conference on Computer Vision*, Graz, Austria, May 7-13, 2006. Cited on page(s): 77
- [43] C. Bregler, J. Malik, and K. Pullen. Twist Based Acquisition and Tracking of Animal and Human Kinematics. *International Journal of Computer Vision*, 56(3):179–194, 2004. Cited on page(s): 27, 75
- [44] G.J. Brostow, I. Essa, D. Steedly, and V. Kwatra. Novel Skeletal representation for Articulated Creatures. In *European Conference on Computer Vision*, Prague, Czech Republic, May 11-14, May 2004. Cited on page(s): 6, 75
- [45] J.M. Buades, R. Mas, and F.J. Perales. Matching a Human Walking Sequence with a VRML Synthetic Model. In *Workshop on Articulated Motion and Deformable Objects*, LNCS 1899, Palma de Mallorca, Spain, Sep 2000. Cited on page(s): 71
- [46] J.M. Buades and F. J. Perales. Upper Body Tracking for Interactive Applications. In *Conference on Articulated Motion and Deformable Objects (AMDO)*, Andratx, Mallorca, Spain, 11-14 July 2006. Cited on page(s): 77
- [47] D. Bullock and J. Zelek. Towards Real-Time 3-D Monocular Visual Tracking of Human Limbs in Unconstrained Environments. *Real-Time Imaging*, 11:323–353, 2005. Cited on page(s): 76
- [48] H. Buxton. Learning and Understanding Dynamic Scene Activity: A Review. *Image and Vision Computing*, 21(1):125–136, 2003. Cited on page(s): 4, 74
- [49] Q. Cai, A. Mitiche, and J.K. Aggarwal. Tracking Human Motion in an Indoor Environment. In *International Conference on Image Processing*, Washington DC, USA, Oct 23-26 1995. Cited on page(s): 22
- [50] S. Calinon and A. Billard. Stochastic Gesture Production and Recognition Model for a Humanoid Robot. In *International Conference on Intelligent Robots and Systems*, Alberta, Canada, Aug 2-6, 2005. Cited on page(s): 40, 42, 42, 75
- [51] S. Calinon, F. Guenter, and A. Billard. Goal-Directed Imitation in a Humanoid Robot. In *International Conference on Robotics and Automation*, Barcelona, Spain, April 18-22, 2005. Cited on page(s): 40, 40, 42, 42, 76
- [52] M.B. Capellades, D. Doermann, D. DeMenthon, and R. Chellappa. An Appearance Based Approach for Human and Object Tracking. In *International Conference on Image Processing*, Barcelona, Spain, Sep 14-17 2003. Cited on page(s): 15, 19, 19, 20, 20, 21, 74
- [53] J. Carranza, C. Theobalt, M. Magnor, and H.-P. Seidel. Free-Viewpoint Video of Human Actors. In *ACM SIGGRAPH*, pages 565—577, 2003. Cited on page(s): 7, 7, 8, 9, 26, 30, 43, 74
- [54] C. Cedras and M. Shah. Motion-Based Recognition: A Survey. *Image and Vision Computing*, 13(2):129–155, 1995. Cited on page(s): 4
- [55] T.H. Chalidabhongse, K. Kim, D. Harwood, and L. Davis. A Perturbation Method for Evaluating Background Subtraction Algorithms. In *Workshop on Visual Surveillance and Performance Evaluation of Tracking and Surveillance*, Beijing, China, Oct 15-16 2005. Cited on page(s): 13, 76
- [56] I.C. Chang and C.L. Huang. The Model-Based Human Motion Analysis System. *Image and Vision Computing*, 18:1067–1083, 2000. Cited on page(s): 71

- [57] M. Chen, G. Ma, and S. Kee. Pixels Classification for Moving Object Extraction. In *IEEE Workshop on Motion and Video Computing (MOTION'05)*, Breckenridge, Colorado, Jan 2005. Cited on page(s): 11, 76
- [58] F. Cheng, W.J. Christmas, and J. Kittler. Recognising Human Running Behaviour in Sports Video Sequences. In *International Conference on Pattern Recognition*, Quebec, Canada, August 11-15 2002. Cited on page(s): 37, 73
- [59] G. Cheung, S. Baker, and T. Kanade. Shape-From-Silhouette for Articulated Objects and its use for Human Body Kinematics Estimation and Motion Capture. In *Computer Vision and Pattern Recognition*, Madison, Wisconsin, USA, June 16-22, 2003. Cited on page(s): 6, 9, 9, 25, 25, 42, 74
- [60] K. Cheung, S. Baker, and T. Kanade. Shape-From-Silhouette Across Time Part II: Applications to Human Modeling and Markerless Motion Tracking. *International Journal of Computer Vision*, 63(3):225–245, 2005. Cited on page(s): 76
- [61] S.S. Cheung and C. Kamath. Robust Techniques for Background Subtraction in Urban Traffic Video. In *Video Communications and Image Processing. SPIE Electronic Imaging*, volume 5308, San Jose, California, Januar 2004. Cited on page(s): 13
- [62] K. Choo and D.J. Fleet. People Tracking Using Hybrid Monte Carlo Filtering. In *International Conference on Computer Vision*, Vancouver, Canada, July 9-12 2001. Cited on page(s): 72
- [63] A.R. Chowdhury and R. Chellappa. A Factorization Approach for Event Recognition. In *CVPR Event Mining Workshop*, Madison, Wisconsin, USA, June 16-22, 2003. Cited on page(s): 33
- [64] A.R. Chowdhury and R. Chellappa. A Factorization Approach for Activity Recognition. In *Computer Vision and Pattern Recognition*, Washington DC, USA, June, 2004. Cited on page(s): 33, 74
- [65] C.-W. Chu, O.C. Jenkins, and M.J. Mataric. Markerless Kinematic Model and Motion Capture from Volume Sequences. In *Computer Vision and Pattern Recognition*, Madison, Wisconsin, USA, June 16-22, 2003. Cited on page(s): 6, 74
- [66] Collins, Lipton, Kanade, Fujiyoshi, Duggins, Tsin, Tolliver, Enomoto, and Hasegawa. A System for Video Surveillance and Monitoring: VSAM Final Report. Technical Report Technical report CMU-RI-TR-00-12, Robotics Institute, Carnegie Mellon University, 2000. Cited on page(s): 13, 21
- [67] D. Comaniciu, V. Ramesh, and P. Meer. Kernel-Based Object Tracking. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 25(5):564–575, 2003. Cited on page(s): 14, 15, 74
- [68] T.F. Cootes, C.J. Taylor, D.H. Hopper, and J. Graham. Active Shape Models - Their Training and Application. *Computer Vision and Image Understanding*, 61(1):39–59, 1995. Cited on page(s): 41
- [69] R. Cucchiara, C. Grana, M. Piccardi, and A. Prati. Detecting Moving Objects, Ghosts, and Shadows in Video Streams. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 25(10):1337–1342, 2003. Cited on page(s): 10, 11, 11, 11, 11, 12, 74
- [70] R. Cucchiara, C. Grana, G. Tardini, and R. Vezzani. Probabilistic People Tracking for Occlusion Handling. In *International Conference on Pattern Recognition*, Cambridge, UK, 23-26 August 2004. Cited on page(s): 15, 19, 75
- [71] R. Cucchiara and R. Vezzani. Assessing Temporal Coherence for Posture Classification with Large Occlusions. In *Workshop on Motion and Video Computing (MOTION'05)*, Breckenridge, Colorado, Jan 2005. Cited on page(s): 76

- [72] N. Cuntoor and R. Chellappa. Key Frame-Based Activity Representation Using Antieigenvalues. In *Asian Conference on Computer Vision*, volume 3852 of *LNCS*, Hyderabad, India, Jan, 13-16, 2006. Cited on page(s): 41, 77
- [73] C. Curio and M.A. Giese. Combining View-based and Model-based Tracking of Articulated Human Movements. In *MOTION*, Breckenridge, Colorado, USA, 5-7 Jan 2005. Cited on page(s): 76
- [74] M. Dahmane and J. Meunier. Real-Time Video Surveillance with Self-Organizing Maps. In *Canadian Conference on Computer and Robot Vision*, Victoria, British Columbia, Canada, May 9-11 2005. Cited on page(s): 76
- [75] N. Dalal and B. Triggs. Histograms of Oriented Gradients for Human Detection. In *Computer Vision and Pattern Recognition*, San Diego, CA, June 20-25 2005. Cited on page(s): 16, 21, 76
- [76] N. Dalal, B. Triggs, and C. Schmid. Human Detection Using Oriented Histograms of Flow and Appearance. In *European Conference on Computer Vision*, Graz, Austria, May 7-13, 2006. Cited on page(s): 16, 77
- [77] B. Dariush. Human Motion Analysis for Biomechanics and Biomedicine. *Machine Vision and Applications*, 14:202–205, 2003. Cited on page(s): 39, 42
- [78] N. Date, H. Yoshimoto, D. Arita, and R. Taniguchi. Real-Time Human Motion Sensing based on Vision-Based Inverse Kinematics for Interactive Applications. In *International Conference of Pattern Recognition*, Cambridge, UK, Aug 23-26, 2004. Cited on page(s): 75
- [79] J. Davis and H. Gao. Recognizing Human Action Efforts: An Adaptive Three-Mode PCA Framework. In *International Conference on Computer Vision*, Nice, France, Oct 13-16, 2003. Cited on page(s): 39, 39, 74
- [80] J.W. Davis and A. Bobick. The Representation and Recognition of Action Using Temporal Templates. In *Conference on Computer Vision and Pattern Recognition*, San Juan, Puerto Rico, June 1997. Cited on page(s): 13
- [81] J.W. Davis and H. Gao. Gender Recognition from Walking Movements using Adaptive Three-Mode PCA. In *Computer Vision and Pattern Recognition*, Washington DC, USA, June, 2004. Cited on page(s): 38, 39, 39, 75
- [82] J.W. Davis and V. Sharma. Robust Detection of People in Thermal Imagery. In *International Conference on Pattern Recognition*, Cambridge, UK, 23-26 August 2004. Cited on page(s): 21, 75
- [83] J.W. Davis and S.R. Taylor. Analysis and Recognition of Walking Movements. In *International Conference on Pattern Recognition*, Quebec, Canada, August 11-15 2002. Cited on page(s): 38, 38, 73
- [84] L. Davis, V. Philomin, and R. Duraiswami. Tracking Humans from a Moving Platform. In *International Conference on Pattern Recognition*, Barcelona, Spain, September 3-8 2000. Cited on page(s): 16, 71
- [85] A.J. Davison, J. Deutscher, and I.D. Reid. Markerless motion capture of complex full-body movement for character animation. In *Eurographics Workshop on Computer Animation and Simulation*, Manchester, UK, Sep 2001. Cited on page(s): 26, 72
- [86] Q. Delamarre and O. Faugeras. 3D Articulated Models and Multi-view Tracking with Physical Forces. *Computer Vision and Image Understanding*, 81(3):328–357, 2001. Cited on page(s): 26, 72
- [87] D. Demirdjian. Enforcing Constraints for Human Body Tracking. In *Workshop on Multiple Object Tracking*, Madison, Wisconsin, USA, June 16-22, 2003. Cited on page(s): 74

- [88] D. Demirdjian. Combining Geometric- and View-Based Approaches for Articulated Pose Estimation. In *European Conference on Computer Vision*, Prague, Czech Republic, May 11-14, 2004. Cited on page(s): 75
- [89] D. Demirdjian, T. Ko, and T. Darrell. Constraining Human Body Tracking. In *International Conference on Computer Vision*, Beijing, China, Oct 15-21, 2003. Cited on page(s): 74
- [90] J. Deutscher, A. Blake, and I. Reid. Articulated Body Motion Capture by Annealed Particle Filtering. In *Computer Vision and Pattern Recognition*, Hilton Head Island, South Carolina, June 13-15 2000. Cited on page(s): 26, 26, 30, 43, 71
- [91] J. Deutscher, A. Davison, and I. Reid. Automatic Partitioning of High Dimensional Search Spaces associated with Articulated Body Motion Capture. In *Computer Vision and Pattern Recognition*, Kauai Marriott, Hawaii, December 9-14 2001. Cited on page(s): 72
- [92] J. Deutscher and I. Reid. Articulated Body Motion Capture by Stochastic Search. *International Journal of Computer Vision*, 61(2):185–205, 2005. Cited on page(s): 76
- [93] M. Dimitrijevic, V. Lepetit, and P. Fua. Human Body Pose Recognition Using Spatio-Temporal Templates. In *Workshop on Modeling People and Human Interaction*, Beijing, China, Oct 15-21, 2005. Cited on page(s): 76
- [94] A.A. Efros, A.C. Berg, G. Mori, and J. Malik. Recognizing Action at a Distance. In *International Conference on Computer Vision*, Nice, France, Oct 13-16, 2003. Cited on page(s): 35, 35, 74
- [95] A. Elgammal, R. Duraiswami, and L.S. Davis. Probabilistic Tracking in Joint Feature-Spatial Spaces. In *Computer Vision and Pattern Recognition*, Madison, Wisconsin, June 16-22 2003. Cited on page(s): 74
- [96] A. Elgammal, D. Harwood, and L. Davis. Non-Parametric Model for Background Subtraction. In *European Conference on Computer Vision*, Dublin, Ireland, June 2000. Cited on page(s): 10, 11, 12, 74
- [97] A. Elgammal and C. Lee. Separating Style and Content on a Nonlinear Manifold. In *Computer Vision and Pattern Recognition*, Washington DC, USA, June, 2004. Cited on page(s): 36, 37, 75
- [98] A. Elgammal and C.S. Lee. Inferring 3D Body Pose from Silhouettes using Activity Manifold Learning. In *Computer Vision and Pattern Recognition*, Washington DC, USA, June, 2004. Cited on page(s): 24, 36, 37, 42, 75
- [99] A. Elgammal, V. Shet, Y. Yacoob, and L. Davis. Learning Dynamics for Exemplar-based Gesture Recognition. In *Computer Vision and Pattern Recognition*, Madison, Wisconsin, USA, June 16-22, 2003. Cited on page(s): 37, 74
- [100] A.M. Elgammal and L.S. Davis. Probabilistic Framework for Segmenting People Under Occlusion. In *International Conference on Computer Vision*, Vancouver, Canada, July 9-12 2001. Cited on page(s): 72
- [101] H.L. Eng, K.A. Toh, A.H. Kam, J. Wang, and W.Y. Yau. An Automatic Drowning Detection Surveillance System For Challenging Outdoor Pool Environments. In *International Conference on Computer Vision*, Nice, France, Oct 13-16, 2003. Cited on page(s): 11, 11, 12, 12, 33, 74
- [102] H.L. Eng, J. Wang, A.H.K.S. Wah, and W.Y. Yau. Robust Human Detection Within a Highly Dynamic Aquatic Environment in Real Time. *IEEE Transactions on Image Processing*, 15(6):1583–1600, 2006. Cited on page(s): 11, 77

- [103] C. Fanti, L. Zelnik-Manor, and P. Perona. Hybrid Models for Human Motion Recognition. In *Computer Vision and Pattern Recognition*, San Diego, California, USA, June 20-25, 2005. Cited on page(s): 39, 76
- [104] P.F. Felzenszwalb and D.P. Huttenlocher. Efficient Matching of Pictorial Structures. In *Computer Vision and Pattern Recognition*, Hilton Head Island, South Carolina, USA, June 13-15, 2000. Cited on page(s): 23, 30, 71
- [105] P. Figueroa, N. Leite, and R.M.L. Barros. Tracking Soccer Players using the Graph Representation. In *International Conference on Pattern Recognition*, Cambridge, UK, Aug 2004. Cited on page(s): 75
- [106] P. Figueroa, N. Leite, and R.M.L. Barros. Background Recovering in Outdoor Image Sequences: An Example of Soccer Players Segmentation. *Image and Vision Computing*, 24(4):363–374, 2006. Cited on page(s): 12, 77
- [107] P. Figueroa, N. Leite, and R.M.L. Barros. Tracking Soccer Players aiming their Kinematical Motion Analysis. *Computer Vision and Image Understanding*, 101:122–135, 2006. Cited on page(s): 77
- [108] P. Fihl, M.B. Holte, T.B. Moeslund, and L. Reng. Action Recognition using Motion Primitives and Probabilistic Edit Distance. In *Conference on Articulated Motion and Deformable Objects (AMDO)*, Andratx, Mallorca, Spain, 11-14 July 2006. Cited on page(s): 77
- [109] D.A. Forsythe and M.M. Fleck. Body Plans. In *Computer Vision and Pattern Recognition*, Puerto Rico, June 17-19, 1997. Cited on page(s): 23
- [110] J.P. Foster, M.S. Nixon, and A. Prgel-Bennett. Automatic Gait Recognition Using Area-Based Metrics. *Pattern Recognition Letters*, 24:2489–2497, 2003. Cited on page(s): 34, 74
- [111] P. Fua, A. Gruen, N. D’Apuzzo, and R. Plänkner. Markerless Full Body Shape and Motion Capture from Video Sequences. In *Symposium on Close Range Imaging, International Society for Photogrammetry and Remote Sensing*, Corfu, Greece, September 2002. Cited on page(s): 73
- [112] J. Gao, A.G. Hauptmann, and H.D. Wactlar. Combining Motion Segmentation with Tracking for Activity Analysis. In *International Conference on Automatic Face and Gesture Recognition*, Seoul, Korea, May 17-19 2004. Cited on page(s): 37, 75
- [113] D.M. Gavrila. The Visual Analysis of Human Movement: A Survey. *Computer Vision and Image Understanding*, 73(1):82–98, 1999. Cited on page(s): 4
- [114] P. Gerard and A. Gagalowicz. Human Body Tracking using a 3D Generic Model Applied to Golf Swing Analysis. In *Conference on Model-based Imaging, Rendering, image Analysis and Graphical special Effects*, INRIA Rocquencourt, France, 10-11 March 2003. Cited on page(s): 74
- [115] J. Giebel, D.M. Gavrila, and C. Schnvrr. A Bayesian Framework for Multi-cue 3D Object Tracking. In *European Conference on Computer Vision*, Prague, Czech Republic, May 11-14, 2004. Cited on page(s): 75
- [116] M. Giese and T. Poggio. Neural Mechanisms for the Recognition of Biological Movements. *Nature Reviews*, 4:179–192, 2003. Cited on page(s): 39, 42
- [117] M. Gleicher. Evaluating Video-Based Motion Capture. In *Conference on Computer Animation*, Geneva, Switzerland, June 19-21 2002. Cited on page(s): 73
- [118] B. Gloyer, H.K. Aghajan, K.Y.S. Siu, and T. Kailath. Video-Based Freeway Monitoring System Using Recursive Vehicle Tracking. In *IS&T-SPIE Symposium on Electronic Imaging: Image and Video Processing*, 1995. Cited on page(s): 12

- [119] J. González. *Human Sequence Evaluation: the Key-frame Approach*. PhD thesis, Univeritat Autònoma de Barcelona, Spain, 2004. Cited on page(s): 75
- [120] J. González, J. Varona, F.X. Roca, and J.J. Villanueva. *aSpaces: Action Spaces for Recognition and Synthesis of Human Actions*. In *International Workshop on Articulated Motion and Deformable Objects*, Palma de Mallorca, Spain, Nov 21-23, 2002. Cited on page(s): 32, 41, 73
- [121] J.J. Gonzalez, I.S. Lim, P. Fua, and D. Thalmann. Robust Tracking and Segmentation of Human Motion in an Image Sequence. In *International Conference on Acoustics, Speech, and Signal Processing*, Hong Kong, April 2003. Cited on page(s): 13, 74
- [122] L. Gorelick, M. Galun, E. Sharon, A. Brandt, and R. Basri. Shape Representation and Recognition Using the Poisson Equation. In *Computer Vision and Pattern Recognition*, Washington DC, USA, June, 2003. Cited on page(s): 36, 36
- [123] N. Grammalidis, G. Goussis, G. Troufakos, and M.G. Strintzis. 3-D Human Body Tracking from Depth Images using Analysis by Synthesis. In *International Conference on Image Processing*, Thessaloniki, Greece, October 7-10 2001. Cited on page(s): 72
- [124] K. Grauman, G. Shakhnarovich, and T. Darrell. Inferring 3D Structure with a Statistical Image-based Shape Model. In *International Conference on Computer Vision*, Nice, France, Oct 13-16, 2003. Cited on page(s): 24, 75
- [125] P.J. Green. *Highly Structured Stochastic Systems*, chapter Trans-dimensional Markov chain Monte Carlo. Oxford University Press, 2003. Cited on page(s): 19
- [126] A. Gritai, Y. Sheikh, and M. Shah. On the Use of Anthropometry in the Invariant Analysis of Human Actions. In *International Conference on Pattern Recognition*, Cambridge, UK, Aug 23-26, 2004. Cited on page(s): 38, 39, 75
- [127] K. Grochow, S.L. Martin, A. Hertzmann, and Z. Popovic. Style-based Inverse Kinematics. In *ACM Transactions on Graphics (SIGGRAPH)*, 2004. Cited on page(s): 29, 29
- [128] P. Guha, A. Mukerjee, and K.S. Venkatesh. Efficient Occlusion Handling for Multiple Agent Tracking by Reasoning with Surveillance Event Primitives. In *Workshop on Visual Surveillance and Performance Evaluation of Tracking and Surveillance*, Beijing, China, Oct 15-16 2005. Cited on page(s): 11, 19, 19, 76
- [129] D. Gutchess, M. Trajkovic, E.C. Solal, D. Lyons, and A. Jain. A Background Model Initialization Algorithm for Video Surveillance. In *International Conference on Computer Vision*, Vancouver, Canada, July 9-12 2001. Cited on page(s): 12, 72
- [130] A.K. Halvorsen. *Model-based Methods in Motion Capture*. PhD thesis, Uppsala University, Sweden, 2002. Cited on page(s): 73
- [131] T.X. Han, H. Ning, and T.S. Huang. Efficient Nonparametric Belief Propagation with Application to Articulated Body Tracking. In *Computer Vision and Pattern Recognition*, New York City, New York, USA, June 17-22, 2006. Cited on page(s): 77
- [132] M. Hariadi, A. Harada, T. Aoki, and T. Higuchi. An LVQ-Based Technique for Human Motion Segmentation. In *Asia-Pacific Conference on Circuits and Systems*, Bali, Indonesia, October 28-31 2002. Cited on page(s): 73
- [133] I. Hariatoglu, R. Cutler, D. Harwood, and L.S. Davis. Backpack: Detection of People Carrying Objects using Silhouettes. *Computer Vision and Image Understanding*, 81(3):385–397, 2001. Cited on page(s): 72
- [134] I. Haritaoglu, M. Flickner, and D. Beymer. Ghost3D: Detecting Body Posture and Parts Using Stereo. In *Workshop on Motion and Video Computing*, Orlando, Florida, November 7 2002. Cited on page(s): 17, 18, 73

- [135] I. Haritaoglu, D. Harwood, and L.S. Davis. Ghost: A Human Body Part Labeling System Using Silhouettes. In *International Conference on Pattern Recognition*, Queensland, Australia, August 17-20 1998. Cited on page(s): 22
- [136] I. Haritaoglu, D. Harwood, and L.S. Davis. W^4 : Who? When? Where? What? - A Real Time System for Detecting and Tracking People. In *International Conference on Automatic Face and Gesture Recognition*, Nara, Japan, 1998. Cited on page(s): 20
- [137] I. Haritaoglu, D. Harwood, and L.S. Davis. W^4 : Real-Time Surveillance of People and Their Activities. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(8):809–830, 2000. Cited on page(s): 10, 12, 13, 16, 17, 19, 38, 71
- [138] K. Hayashi, M. Hashimoto, K. Sumi, and K. Sasakawa. Multiple-Person Tracker with a Fixed Slanting Stereo Camera. In *International Conference on Automatic Face and Gesture Recognition*, Seoul, Korea, May 17-19 2004. Cited on page(s): 17, 18, 18, 75
- [139] M. Heikkila, M. Pietikainen, and J. Heikkila. A Texture-Based Method for Detecting Moving Objects. In *British Machine Vision Conference*, London, UK, September 7-9 2004. Cited on page(s): 11, 12, 13, 75
- [140] M. Heikkila and M. Pietikainen. A Texture-Based Method for Modeling the Background and Detecting Moving Objects. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 28(4):657–662, 2006. Cited on page(s): 10, 12, 13, 77
- [141] L. Herda. *Using Biomechanical Constraints to Improve Video-Based Motion Capture*. PhD thesis, Computer Vision Lab, EPFL, Lausanne, Switzerland, 2003. Cited on page(s): 74
- [142] L. Herda, P. Fua, R. Plänkers, R. Boulic, and D. Thalmann. Using Skeleton-Based Tracking to Increase the Reliability of Optical Motion Capture. *Human Movement Science*, 20(3), 2001. Cited on page(s): 72
- [143] L. Herda, R. Urtasun, and P. Fua. Hierarchical Implicit Surface Joint Limits to Constrain Video-Based Motion Capture. In *European Conference on Computer Vision*, Prague, Czech Republic, May 11-14 2004. Cited on page(s): 6, 75
- [144] L. Herda, R. Urtasun, and P. Fua. Hierarchical Implicit Surface Joint Limits for Human Body Tracking. *Computer Vision and Image Understanding*, 99(2):189–209, 2005. Cited on page(s): 6, 9, 76
- [145] L. Herda, R. Urtasun, P. Fua, and A. Hanson. An Automatic Method for Determining Quaternion Field Boundaries for Ball-and-Socket Joint Limits. In *International Conference on Automatic Face and Gesture Recognition*, Washington DC, USA, May 20-21 2002. Cited on page(s): 73
- [146] A. Hilton, D. Beresford, T. Gentils, R. Smith, and W. Sun. Virtual People: Capturing human models to populate virtual worlds. In *International Conference on Computer Animation*, May 1999. Cited on page(s): 7
- [147] D. Hogg. Model-Based Vision: A Program to See a Walking Person. *Image and Vision Computing*, 1(1):5–20, 1983. Cited on page(s): 22, 25, 43
- [148] T. Horprasert, D. Harwood, and L.S. Davis. A Statistical Approach for Real-time Robust Background Subtraction and Shadow Detection. In *IEEE ICCV'99 FRAME-RATE WORKSHOP*, Corfu, Greece, September 1999. Cited on page(s): 11
- [149] R. Hoshino, S. Yonenmoto, D. Arita, and R.I. Taniguchi. Real-Time Analysis of Human Motion using Multi-View Silhouette Contours. In *The 12th Scandinavian Conference on Image Analysis*, Bergen, Norway, 2001. Cited on page(s): 72

- [150] N Howe. Flow Lookup and Biological Motion Perception. In *International Conference on Image Processing*, Genova, Italy, Sep 11-14 2005. Cited on page(s): 76
- [151] N.R. Howe. Silhouette Lookup for Automatic Pose Tracking. In *Workshop on Articulated and Non-Rigid Motion*, Washington DC, USA, June, 2004. Cited on page(s): 24, 24, 30, 75
- [152] N.R. Howe. Boundary Fragment Matching and Articulated. In *Conference on Articulated Motion and Deformable Objects (AMDO)*, Andratx, Mallorca, Spain, 11-14 July 2006. Cited on page(s): 77
- [153] N.R. Howe, M.E. Leventon, and W.T. Freeman. Bayesian Reconstruction of 3D Human Motion from Single-Camera Video. In *Advances in Neural Information Processing Systems 12*. MIT Press, 2000. Cited on page(s): 29, 71
- [154] M. Hu, W. Hu, and T. Tan. Tracking People through Occlusion. In *International Conference on Pattern Recognition*, Cambridge, UK, Aug 2004. Cited on page(s): 14, 15, 19, 21, 21, 42, 42, 75
- [155] W. Hu, M. Hu, X. Zhou, T. Tan, J. Lou, and S. Maybank. Principal Axis-Based Correspondence between Multiple Cameras for People Tracking. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 28(4):663–671, 2006. Cited on page(s): 17, 77
- [156] W. Hu, T. Tan, L. Wang, and S. Maybank. A Survey on Visual Surveillance of Object Motion and Behaviors. *Transactions on Systems, Man, and Cybernetics - Part C: Applications and Reviews*, 34(3):334–352, 2004. Cited on page(s): 4, 75
- [157] G. Hua, M-H. Yang, and Y. Wu. Learning to Estimate Human Pose with Data Driven Belief Propagation. In *Computer Vision and Pattern Recognition*, San Diego, California, USA, June 20-25, 2005. Cited on page(s): 23, 30, 42, 43
- [158] C.L. Huang and C.C. Lin. Model-based human body motion analysis for MPEG IV video encoding. In *International Conference on Information Technology: Coding and Computing*, Las Vegas, Nevada, April 2-4 2001. Cited on page(s): 72
- [159] C.L. Huang, T.H. Tseng, and H.C. Shih. A Model-Based Articulated Human Motion Tracking System. In *Asian Conference on Computer Vision*, Jeju, Korea, January 27-30 2004. Cited on page(s): 75
- [160] F. Huang, H. Di, and G. Xu. Viewpoint Insensitive Posture Representation for Action Recognition. In *International Conference on Articulated Motion and Deformable Objects, Springer LNCS 4069*, Andratx, Mallorca, Spain, Jul 11-14, 2006. Cited on page(s): 77
- [161] J. Huang, S.R. Kumar, M. Mitra, and W. Zhu. Spatial Color Indexing and Applications. *International Journal of Computer Vision*, 35(3):91–101, 1999. Cited on page(s): 15, 21
- [162] Y. Huang and T.S. Huang. Model-Based Human Body Tracking. In *International Conference on Pattern Recognition*, Quebec, Canada, August 11-15 2002. Cited on page(s): 73
- [163] I. Huerta, D. Rowe, J. González, and J.J. Villanueva. Efficient Incorporation of Motionless Foreground Objects for Adaptive Background Segmentation. In *Conference on Articulated Motion and Deformable Objects (AMDO)*, Andratx, Mallorca, Spain, 11-14 July 2006. Cited on page(s): 77
- [164] A.J. Ijspeert, J. Nakanishi, and S. Schaal. Movement Imitation with Nonlinear Dynamical Systems in Humanoid Robots. In *International Conference on Robotics and Automation*, Washington DC, USA, May, 2002. Cited on page(s): 40, 73

- [165] S.S. Intille and A.F. Bobick. Recognizing Planned, Multiperson Action. *Computer Vision and Image Understanding*, 81(3):414–445, 2001. Cited on page(s): 72
- [166] S. Ioffe and D. Forsyth. Finding people by sampling. In *International Conference on Computer Vision*, Korfu, Greece, 1999. Cited on page(s): 23
- [167] S. Ioffe and D. Forsyth. Human Tracking with Mixtures of Trees. In *International Conference on Computer Vision*, Vancouver, Canada, July 9-12 2001. Cited on page(s): 23, 30, 72
- [168] M. Isard and A. Blake. CONDENSATION - Conditional Density Propagation for Visual Tracking. *International Journal on Computer Vision*, 29(1):5–28, 1998. Cited on page(s): 19
- [169] Y. Ivanov and A. Bobick. Recognition of Visual Activities and Interactions by Stochastic Parsing. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(8):852–872, 2000. Cited on page(s): 41, 41, 42, 42
- [170] Y.A. Ivanov, A.F. Bobick, and J. Liu. Fast Lighting Independent Background Subtraction. *International Journal of Computer Vision*, 37(2):199–207, 2000. Cited on page(s): 17, 17, 71
- [171] S. Iwase and H. Saito. Parallel Tracking of All Soccer Players by Integrating Detected Positions in Multiple View Images. In *International Conference on Pattern Recognition*, Cambridge, UK, Aug 2004. Cited on page(s): 17, 18, 75
- [172] D.S. Jang, S.W. Jang, and H.I. Choi. 2D Human Body Tracking with Structural Kalman Filter. *Pattern Recognition*, 35:2041–2049, 2002. Cited on page(s): 73
- [173] T. Jeaggli, E.K. Meier, and L.V. Gool. Monocular Tracking with a Mixture of View-Dependent Learned Models. In *Conference on Articulated Motion and Deformable Objects (AMDO)*, Andratx, Mallorca, Spain, 11-14 July 2006. Cited on page(s): 77
- [174] O.C. Jenkins and M. Mataric. Automated Modularization of Human Motion into Actions and Behaviors. Technical Report CRES-02-002, Center for Robotics and Embedded Systems, University of S. California, 2002. Cited on page(s): 32, 73
- [175] O.C. Jenkins and M. Mataric. Deriving Action and Behavior Primitives from Human Motion Capture Data. In *International Conference on Robotics and Automation*, Washington DC, USA, May, 2002. Cited on page(s): 40, 73
- [176] O.C. Jenkins and M.J. Mataric. Deriving Action and Behavior Primitives from Human Motion Data. In *International Conference on Intelligent Robots and Systems*, pages 2551–2556, Lausanne, Switzerland, Sept.30 – Oct.4, 2002. Cited on page(s): 40, 73
- [177] A.D. Jepson, D.J. Fleet, and T.F. El-Maraghi. Robust Online Appearance Models for Visual Tracking. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 25(10):1296–1311, 2003. Cited on page(s): 15, 74
- [178] H. Jiang, M. Drew, and Z. Li. Successive convex matching for action detection. In *Computer Vision and Pattern Recognition*, pages 1646–1653, New York City, New York, USA, June 17-22, 2006. Cited on page(s): 36
- [179] H. Jiang, M.S. Drew, and Z.-N. Li. Successive Convex Matching for Action Detection. In *Computer Vision and Pattern Recognition*, New York City, New York, USA, June 17-22, 2006. Cited on page(s): 77
- [180] S. Ju. Human Motion Estimation and Recognition (Depth Oral Report). Technical report, University of Toronto, 1996. Cited on page(s): 4
- [181] I.N. Junejo, O. Javed, and M. Shah. Multi Feature Path Modeling for Video Surveillance. In *International Conference on Pattern Recognition*, Cambridge, UK, Aug 2004. Cited on page(s): 33, 75

- [182] A. Kale, A. Rajagopalan, N. Cuntoor, and V. Krueger. Human Identification Using Gait. In *International Conference on Automatic Face and Gesture Recognition*, Washington DC, USA, May 21-22, 2002. Cited on page(s): 34
- [183] A. Kale, A. Sundaresan, A.N. Rjagopalan, N. Cuntoor, A.R. Chowdhury, V. Krger, and R. Chellappa. Identification of Humans Using Gait. *IEEE Transactions on Image Processing*, 9:1163–1173, 2004. Cited on page(s): 34
- [184] Y. Kameda and M. Minoh. A Human Motion Estimation Method Using 3-Successive Video Frames. In *International Conference on Virtual Systems and Multimedia*, Gifu City, Japan, September 1996. Cited on page(s): 13
- [185] D.W. Kang, Y. Onuma, and J. Ohya. Estimating Complicated and Overlapped Human Body Postures by Wearing a Multiple-Colored Suit Using Color Information Processing. In *International Conference on Automatic Face and Gesture Recognition*, Seoul, Korea, May 17-19 2004. Cited on page(s): 75
- [186] J. Kang, I. Cohen, and G. Medioni. Persistent Objects Tracking Across Multiple Non Overlapping Cameras. In *IEEE Workshop on Motion and Video Computing (MOTION'05)*, Breckenridge, Colorado, Jan 2005. Cited on page(s): 14, 15, 21, 43, 76
- [187] J. Kang, I. Cohen, G. Medioni, and C. Yuan. Detection and Tracking of Moving Objects from a Moving Platform in Presence of Strong Parallax. In *International Conference on Computer Vision*, Beijing, China, Oct 15-21 2005. Cited on page(s): 76
- [188] I.A. Karaulova, P.M. Hall, and A.D. Marshall. A Hierarchical Models of Dynamics for Tracking People with a Single Video Camera. In *British Machine Vision Conference*, Bristol, UK, 11-14 Sep 2000. Cited on page(s): 29, 71
- [189] Y. Ke, R. Sukthankar, and M. Hebert. Efficient Visual Event Detection using Volumetric Features. In *Internatinal Conference on Computer Vision*, Beijing, China, Oct 15-21, 2005. Cited on page(s): 76
- [190] R. Kehl, M. Bray, and L. VanGool. Full Body Tracking from Multiple Views Using Stochastic Sampling. In *Computer Vision and Pattern Recognition*, San Diego, California, USA, June 20-25, 2005. Cited on page(s): 27, 30, 43, 76
- [191] D.G. Kendall, D. Barden, T.K. Carne, and H. Le. *Shape and Shape Theory*. Wiley, 1999. Cited on page(s): 33
- [192] S. Khan, O. Javed, Z. Rasheed, and M. Shah. Human Tracking in Multiple Cameras. In *International Conference on Computer Vision*, Vancouver, Canada, July 9-12 2001. Cited on page(s): 21, 43, 72
- [193] S. Khan and M. Shah. Tracking People in Presence of Occlusion. In *Asian Conference on Computer Vision*, Taipei, Taiwan, January 8-11 2000. Cited on page(s): 14, 15, 19, 20, 20, 21, 21, 42, 42, 71
- [194] S.M. Khan and M. Shah. A Multiview Approach to Tracking People in Crowded Scenes Using a Planar Homography Constraint. In *European Conference on Computer Vision*, Graz, Austria, May 7-13, 2006. Cited on page(s): 18, 77
- [195] K. Kim, T.H. Chalidabhongse, D. Harwood, and L. Davis. Real-Time Foreground-Background Segmentation using Codebook Model. *Real-Time Imaging*, 11(3):172–185, 2005. Cited on page(s): 10, 12, 12, 13, 76
- [196] T. Kohonen. Self-Organized Formation of Topologically Correct Feature Maps. *Biological Cybernetics*, 43:59–69, 1982. Cited on page(s): 41
- [197] A. Koschan, S. Kang, J. Paik, B. Abidi, and M. Abidi. Color Active Shape Models for Tracking Non-Rigid Objects. *Pattern Recognition Letters*, 24:1751–1765, 2003. Cited on page(s): 17, 74

- [198] L. Kovar, M. Gleicher, and F. Pighin. Motion Graphs. In *ACM SIGGRAPH*, pages 473–482, 2002. Cited on page(s): 28, 29
- [199] N. Krahnstoever and R. Sharma. Articulated Models from Video. In *Computer Vision and Pattern Recognition*, Washington DC, USA, June, 2004. Cited on page(s): 6, 75
- [200] N. Krahnstoever, M. Yeasin, and R. Sharma. Automatic Acquisition and Initialization of Articulated Models. *Machine Vision and Applications*, 14(4):218–228, 2003. Cited on page(s): 6, 74
- [201] F. Kristensen, P. Nilsson, and V. wall. Background Segmentation Beyond RGB. In *Asian Conference on Computer Vision*, LNCS 3852, Hyderabad, India, January 13-16 2006. Cited on page(s): 10, 77
- [202] T. Krosshaug and R. Bahr. A Model-Based Image-Matching Technique for Three-Dimensional Reconstruction of Human Motion from Uncalibrated Video Sequences. *Journal of Biomechanics*, 38(4):919–929, 2005. Cited on page(s): 76
- [203] V. Krüger, J. Anderson, and T. Prehn. Probabilistic Model-Based Background Subtraction. In *Scandinavian Conference on Image Analysis*, Joensuu, Finland, Jun 19-22 2005. Cited on page(s): 17, 76
- [204] M.P. Kumar, P.H.S. Torr, and A. Zisserman. Learning Layered Motion Segmentations of Video. In *International Conference on Computer Vision*, Beijing, China, Oct 15-21 2005. Cited on page(s): 76
- [205] C.S. Lee and A. Elgammal. Carrying Object Detection Using Pose Preserving Dynamic Shape Models. In *Conference on Articulated Motion and Deformable Objects (AMDO)*, Andratx, Mallorca, Spain, 11-14 July 2006. Cited on page(s): 77
- [206] C.S. Lee and A. Elgammal. Human Motion Synthesis by Motion Manifold Learning and Motion Primitive Segmentation. In *Conference on Articulated Motion and Deformable Objects (AMDO)*, Andratx, Mallorca, Spain, 11-14 July 2006. Cited on page(s): 77
- [207] D.S. Lee. Effective Gaussian Mixture Learning for Video Background Subtraction. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 27(5):827–832, 2005. Cited on page(s): 11, 76
- [208] J. Lee, J. Chai, P.S.A. Reitsma, J.K. Hodgins, and N.S. Pollard. Interactive Control of Avatars Animated With Human Motion Data. In *ACM SIGGRAPH*, pages 491–500, 2002. Cited on page(s): 28
- [209] M.W. Lee and I. Cohen. Human Upper Body Pose Estimation in Static Images. In *European Conference on Computer Vision*, Prague, Czech Republic, May 11-14 2004. Cited on page(s): 75
- [210] M.W. Lee and I. Cohen. Proposal Maps driven MCMC for Estimating Human Body Pose in Static Images. In *Computer Vision and Pattern Recognition*, Washington DC, USA, June, 2004. Cited on page(s): 28, 43, 75
- [211] M.W. Lee, I. Cohen, and S.K. Jung. Particle Filter with Analytical Inference for Human Body Tracking. In *Workshop on Motion and Video Computing*, Orlando, Florida, November 7 2002. Cited on page(s): 73
- [212] M.W. Lee and R. Navatia. Human Pose Tracking Using Multi-level Structured Models. In *European Conference on Computer Vision*, Graz, Austria, May 7-13, 2006. Cited on page(s): 77
- [213] M.W. Lee and R. Nevatia. Dynamic Human Pose Estimation using Markov chain Monte Carlo Approach. In *Workshop on Motion and Video Computing (MOTION'05)*, Breckenridge, Colorado, Jan 2005. Cited on page(s): 76

- [214] B. Leibe, E. Seemann, and B. Schiele. Pedestrian Detection in Crowded Scenes. In *Computer Vision and Pattern Recognition*, San Diego, CA, June 20-25 2005. Cited on page(s): 16, 76
- [215] I. Leichter, M. Lindenbaum, and E. Rivlin. A General Framework for Combining Visual Trackers The "Black Boxes" Approach. *International Journal of Computer Vision*, 67(3):343–363, 2006. Cited on page(s): 77
- [216] M. Leo, T. D’Orazio, I. Gnoni, P. Spagnolo, and A. Distanto. Complex Human Activity Recognition for Monitoring Wide Outdoor Environments. In *International Conference on Pattern Recognition*, Cambridge, UK, Aug 23-26, 2004. Cited on page(s): 37, 75
- [217] B. Li, R. Chellappa, and H. Moon. Monte Carlo Simulation Techniques for Probabilistic Tracking. In *Thirty-Fifth Asilomar Conference on Signals, Systems and Computers*, Pacific Grove, California, November 4-7 2001. Cited on page(s): 72
- [218] Y. Li, A. Hilton, and J. Illingworth. A Relaxation Algorithm for Real-Time Multiple View 3D-Tracking. *Image and Vision Computing*, 20(12):841–859, 2002. Cited on page(s): 19, 73
- [219] D. Liebowitz and S. Carlsson. Uncalibrated Motion Capture Exploiting Articulated Structure Constraints. *International Journal of Computer Vision*, 51(3):171–187, 2003. Cited on page(s): 74
- [220] H. Lim, V.I. Morariu, O.I. Camps, and M. Sznajder. Dynamic Appearance Modeling for Human Tracking. In *Computer Vision and Pattern Recognition*, New York City, New York, USA, June 17-22, 2006. Cited on page(s): 8, 77
- [221] S.N. Lim, A. Mittal, L.S. Davis, and N. Paragios. Fast Illumination-Invariant Background Subtraction Using Two Views: Error Analysis, Sensor Placement and Applications. In *Computer Vision and Pattern Recognition*, San Diego, CA, June 20-25 2005. Cited on page(s): 17, 17, 18, 76
- [222] X. Liu and C. Chua. Multi-Agent Activity Recognition using Observation Decomposed Hidden Markov Models. *Image and Vision Computing*, 24:166–175, 2006. Cited on page(s): 33, 77
- [223] D.G. Lowe. Distinctive Image Features From Scale-Invariant Keypoints. *International Journal on Computer Vision*, 60(2):91–110, 2004. Cited on page(s): 16
- [224] G. Loy, M. Eriksson, and J. Sullivan. Monocular 3D Reconstruction of Human Motion in Long Action Sequences. In *European Conference on Computer Vision*, Prague, Czech Republic, May 11-14, 2004. Cited on page(s): 27, 75
- [225] C. Lu and N. Ferrier. Repetitive Motion Analysis: Segmentation and Event Classification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 26(2):258–263, 2004. Cited on page(s): 40, 75
- [226] Y. Luo, T.-W. Wu, and J.-N. Hwang. Object-Based Analysis and Interpretation of Human Motion in Sports Video Sequences by Dynamic Bayesian Networks. *Computer Vision and Image Understanding*, 92:196–216, 2003. Cited on page(s): 37
- [227] F. Lv, J. Kang, R. Nevatia, I. Cohen, and G. Medioni. Automatic Tracking and Labeling of Human Activities in a video sequence. In *6th IEEE International Workshop on Performance Evaluation of Tracking and Surveillance*, Prague, May 11-14 2004. Cited on page(s): 75
- [228] F. Lv and R. Navatia. Recognition and Segmentation of 3-D Human Action Using HMM and Multi-class AdaBoost. In *European Conference on Computer Vision*, Graz, Austria, May 7-13, 2006. Cited on page(s): 77

- [229] F. Lv and R. Nevatia. Recognition and Segmentation of 3-D Human Action Using HMM and Multi-class AdaBoost. In *European Conference on Computer Vision*, Graz, Austria, May 7-13, 2006. Cited on page(s): 41
- [230] J. MacCormick and M. Isard. Partitioned sampling, articulated objects, and interface-quality hand tracking. In *European Conference on Computer Vision*, Dublin, Ireland, Jun, 2000. Cited on page(s): 26
- [231] O. Masound and N. Papanikolopoulos. A Method for Human Action Recognition. *Image and Vision Computing*, 21:729–743, 2003. Cited on page(s): 35, 74
- [232] S.J. McKenna, S. Jabri, Z. Duric, and H. Wechsler. Tracking Interacting People. In *International Conference on Automatic Face and Gesture Recognition*, Grenoble, France, March 2000. Cited on page(s): 10, 10, 11, 11, 12, 14, 14, 19, 19, 20, 21, 42
- [233] C. Menier, E. Boyer, and B. Raffin. 3D Skeleton-Based Body Pose Recovery. In *International Symposium on 3D Data Processing, Visualisation and Transmission*, 2006. Cited on page(s): 6, 9, 77
- [234] D. Metaxas. From Visual Input to Modeling Humans. In *Conference on Computer Animation*, Geneva, Switzerland, June 19-21 2002. Cited on page(s): 73
- [235] A. Micilotta, E. Ong, and R. Bowden. Detection and tracking of humans by probabilistic body part assembly. In *British Machine Vision Conference*, Oxford, UK, Sep 2005. Cited on page(s): 8, 8, 23, 23, 23, 30, 43, 76
- [236] A.S. Micilotta, E.-J. Ong, and R. Bowden. Real-Time Upper Body Detection and 3D Pose Estimation in Monoscopic Images. In *European Conference on Computer Vision*, Graz, Austria, May 7-13, 2006. Cited on page(s): 77
- [237] I. Mikić, M. Trivedi, E. Hunter, and P. Cosman. Human Body Model Acquisition and Motion Capture Using Voxel Data. In F.J. Perales and E.R. Hancock, editors, *AMDO 2002*, LNCS 2492. Springer-Verlag, 2002. Cited on page(s): 25, 25, 73
- [238] I. Mikić, M. Trivedi, E. Hunter, and P. Cosman. Human Body Model Acquisition and Tracking Using Voxel Data. *International Journal of Computer Vision*, 53(3):199—223, 2003. Cited on page(s): 25, 25, 42, 43, 74
- [239] I. Mikić, M.M. Trivedi, E. Hunter, and P. Cosman. Articulated Body Posture Estimation from Multi-Camera Voxel Data. In *Computer Vision and Pattern Recognition*, Kauai Marriott, Hawaii, December 9-14 2001. Cited on page(s): 72
- [240] K. Mikolajczyk, D. Schmid, and A. Zisserman. Human Detection Based on a Probabilistic Assembly of Robust Part Detectors. In *European Conference on Computer Vision*, Prague, Czech Republic, May 11-14, 2004. Cited on page(s): 23, 43, 75
- [241] J. Mitchelson and A. Hilton. From Visual Tracking to Animation using Hierarchical Sampling. In *Conference on Model-based Imaging, Rendering, image Analysis and Graphical special Effects*, Rocquencourt, France, 10-11 March 2003. Cited on page(s): 74
- [242] J. Mitchelson and A. Hilton. Hierarchical Tracking of Multiple People. In *British Machine Vision Conference*, Norwich, UK, Sep 2003. Cited on page(s): 26, 74
- [243] A. Mittal and L.S. Davis. M2Tracker: A Multi-View Approach to Segmenting and Tracking People in a Cluttered Scene using Region-Based Stereo. In *European Conference on Computer Vision*, Copenhagen, Denmark, 2002. Cited on page(s): 17, 18, 73
- [244] A. Mittal and L.S. Davis. M2Tracker: A Multi-View Approach to Segmenting and Tracking People in a Cluttered Scene. *International Journal of Computer Vision*, 51(3):189–203, 2003. Cited on page(s): 15, 17, 18, 21, 42, 74

- [245] T.B. Moeslund. *Computer Vision-Based Motion Capture of Human Body Language*. PhD thesis, Lab of Computer Vision and Media Technology, Aalborg University, Denmark, 2003. Cited on page(s): 74
- [246] T.B. Moeslund. *Pose Estimating the Human Arm using Kinematics and the Sequential Monte Carlo Framework*, chapter 4 in part IX of *Cutting Edge Robotics*. Pro literatur Verlag. ISBN:3-86611-038-3, 2005. Cited on page(s): 28, 76
- [247] T.B. Moeslund and E. Granum. A Survey of Computer Vision-Based Human Motion Capture. *Computer Vision and Image Understanding*, 81(3):231–268, 2001. Cited on page(s): 1, 4, 4, 4, 5, 5, 5, 5, 7, 22, 25, 72
- [248] T.B. Moeslund and E. Granum. Pose Estimation of a Human Arm using Kinematic Constraints. In *The 12th Scandinavian Conference on Image Analysis*, Bergen, Norway, 2001. Cited on page(s): 6, 72
- [249] T.B. Moeslund and E. Granum. Modelling and Estimating the Pose of a Human Arm. *Machine Vision and Applications*, 14(4):237–247, 2003. Cited on page(s): 74
- [250] T.B. Moeslund and E. Granum. Sequential Monte Carlo Tracking of Body Parameters in a Sub-Space. In *International Workshop on Analysis and Modeling of Faces and Gestures*, Nice, France, October 2003. Cited on page(s): 74
- [251] T.B. Moeslund and E. Granum. Motion Capture of Articulated Chains by Applying Auxiliary Information to the Sequential Monte Carlo Algorithm. In *International Conference on Visualization, Imaging, and Image Processing*, Marbella, Spain, Sep 2004. Cited on page(s): 75
- [252] T.B. Moeslund, C.B. Madsen, and E. Granum. Modelling the 3D Pose of a Human Arm and the Shoulder Complex utilising only Two Parameters. *International Journal on Integrated Computer-Aided Engineering*, 12(2):159–177, 2005. Cited on page(s): 6, 9, 28, 76
- [253] T.B. Moeslund, M. Vittrup, K.S. Pedersen, M.K. Laursen, M.K.D. Sørensen, H. Uhrenfeldt, and E. Granum. Estimating the 3D Shoulder Position using Monocular Vision. In *International Conference on Imaging Science, Systems, and Technology*, Las Vegas, Nevada, June 24-27 2002. Cited on page(s): 6, 73
- [254] A. Mohan, C. Papageorgiou, and T. Poggio. Example-Based Object Detection in Images by Components. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 23(4):349–361, 2001. Cited on page(s): 14, 14, 21, 23, 72
- [255] L. Molina and A. Hilton. Synthesis of Novel Movements from a Database of Motion Capture Data. In *International Conference on Human Motion Analysis*, December 2000. Cited on page(s): 28, 29
- [256] A. Monnet, A. Mittal, N. Paragios, and V. Ramesh. Background Modeling and Subtraction of Dynamic Scenes. In *International Conference on Computer Vision*, Nice, France, 13-16 October 2003. Cited on page(s): 11, 11, 74
- [257] M. Montemerlo, S. Thrun, and W. Whittaker. Conditional Particle Filter for Simultaneous Mobile Robot Localization and People-Tracking. In *International Conference on Robotics & Automation*, Washington DC, USA, May 2002. Cited on page(s): 73
- [258] H. Moon, R. Chellappa, and A. Rosenfeld. Tracking of Human Activities using Shape-Encoded Particle Propagation. In *International Conference on Image Processing*, Thessaloniki, Greece, October 7-10 2001. Cited on page(s): 72
- [259] K. Moon and V. Pavlovic. Impact of Dynamics on Subspace Embedding and Tracking of Sequences. In *Computer Vision and Pattern Recognition*, New York City, New York, USA, June 17-22, 2006. Cited on page(s): 29, 77

- [260] G. Mori and J. Malik. Recovery of 3D Human Body Configurations Using Shape Contexts. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 28(7):1052–1062, 1998. Cited on page(s): 77
- [261] G. Mori, X. Ren, A. Efros, and J. Malik. Recovering Human Body Configurations: Combining Segmentation and Recognition. In *Computer Vision and Pattern Recognition*, Washington DC, USA, June, 2004. Cited on page(s): 75
- [262] J. Mulligan. Upper Body Pose Estimation from Stereo and Hand-Face Tracking. In *Canadian Conference on Computer and Robot Vision*, Victoria, British Columbia, Canada, May 9-11 2005. Cited on page(s): 6, 76
- [263] T. Murakita, T. Ikeda, and H. Ishiguro. Multisensor Human Tracker based on the Markov Chain Monte Carlo Method. In *International Conference on Pattern Recognition*, Cambridge, UK, Aug 2004. Cited on page(s): 75
- [264] H.-H Nagel. From Image Sequences Towards Conceptual Descriptions. *Image and Vision Computing*, 6(2):59–74, 1988. Cited on page(s): 32
- [265] R. Navaratnam, A. Thayananthan, P.H.S. Torr, and R. Cipolla. Hierarchical part-based human body pose estimation. In *British Machine Vision Conference*, Oxford, UK, Sep 2005. Cited on page(s): 28, 43, 76
- [266] P. Nillius, J. Sullivan, and S. Carlsson. Multi-Target Tracking - Linking Identities using Bayesian Network Inference. In *Computer Vision and Pattern Recognition*, New York City, New York, USA, June 17-22, 2006. Cited on page(s): 77
- [267] H. Ning, T. Han, Y. Hu, Z. Zhang, U. Fu, and T. Huang. A Realtime Shrug Detector. In *International Conference on Automatic Face and Gesture Recognition*, Southampton, UK, April 10-12, 2006. Cited on page(s): 39
- [268] K. Ogaki, Y. Iwai, and M. Yachida. Posture Estimation Based on Motion and Structure Models. *Systems and Computers in Japan*, 32(4), 2001. Cited on page(s): 72
- [269] K. Okuma, A. Taleghani, N.D. Freitas, J.J. Little, and David G. Lowe. A Boosted Particle Filter: Multitarget Detection and Tracking. In *European Conference on Computer Vision*, Prague, Czech Republic, May 11-14 2004. Cited on page(s): 14, 14, 15, 21, 21, 42, 42, 75
- [270] N. Oliver, B. Rosario, and A. Pentland. A Bayesian Computer Vision System for Modeling Human Interactions. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(8):831–843, 2000. Cited on page(s): 10, 71
- [271] E.-J. Ong and A. Hilton. Learnt Inverse Kinematics for Animation Synthesis. In *Conference on Vision, Video and Graphics*, Edinburgh, UK, June 2005. Cited on page(s): 29, 29
- [272] E.-J. Ong, A. Hilton, and A.S. Micilotta. Viewpoint Invariant Exemplar-Based 3D Human Tracking. In *First IEEE Workshop on Modeling People and Human Interaction (PHI'05)*, Beijing, China, Oct 15 2005. Cited on page(s): 24, 29, 76
- [273] D. Ormoneit, M.J. Black, T. Hastie, and H. Kjellstrm. Representing Cyclic Human Motion using Functional Analysis. *Image and Vision Computing*, 23(14):1264–1276, 2005. Cited on page(s): 76
- [274] D. Ormoneit, H. Sidenbladh, M.J. Black, and T. Hastie. Learning and Tracking Cyclic Human Motion. In *Workshop on Human Modeling, Analysis and Synthesis at CVPR*, Hilton Head Island, South Carolina, June 13-15 2000. Cited on page(s): 71
- [275] I.B. Ozer and W.H. Wolf. A Hierarchical Human Detection System in (Un)Compressed Domains. *IEEE Transactions on Multimedia*, 4(2):283–300, 2002. Cited on page(s): 14, 14, 36, 73

- [276] C. Pan and S. Ma. Parametric Tracking of Human Contour by Exploiting Intelligent Edge. In *International Conference on Automatic Face and Gesture Recognition*, Seoul, Korea, May 17-19 2004. Cited on page(s): 75
- [277] V. Parameswaran and R. Chellappa. View Invariants for Human Action Recognition. In *Computer Vision and Pattern Recognition*, Madison, Wisconsin, June 16-22 2003. Cited on page(s): 38, 38, 74
- [278] V. Parameswaran and R. Chellappa. View Independent Human Body Pose Estimation from a Single Perspective Image. In *Computer Vision and Pattern Recognition*, Washington DC, USA, June, 2004. Cited on page(s): 6, 75
- [279] V. Parameswaran and R. Chellappa. View Invariance for Human Action Recognition. *International Journal of Computer Vision*, 66(1):83–101, 2006. Cited on page(s): 38, 38, 77
- [280] S. Park and J.K. Aggarwal. Segmentation and Tracking of Interacting Human Body Parts Under Occlusion and Shadowing. In *Workshop on Motion and Video Computing*, Orlando, Florida, November 7 2002. Cited on page(s): 42, 42, 73
- [281] S. Park and J.K. Aggarwal. Semantic-level Understanding of Human Actions and Interactions using Event Hierarchy. In *CVPR workshop on Articulated and non-rigid motion*, Washington DC, USA, June, 2004. Cited on page(s): 41, 42, 42, 75
- [282] S. Park and J.K. Aggarwal. Simultaneous Tracking of Multiple Body Parts of Interacting Persons. *Computer Vision and Image Understanding*, 102(1):1–21, 2006. Cited on page(s): 15, 15, 20, 21, 21, 77
- [283] S. Park and M. Trivedi. A Track-based Human Movement Analysis and Privacy Protection System Adaptive to Environmental Contexts. In *Asian Conference on Computer Vision*, LNCS 3852, Hyderabad, India, January 13-16 2006. Cited on page(s): 77
- [284] A.E.C. Pece. Tracking of Non-Gaussian Clusters in the PETS2001 Image Sequences. In *Workshop on Performance Evaluation of Tracking and Surveillance*, Kauai, Hawaii, December 9 2001. Cited on page(s): 72
- [285] A.E.C. Pece. From Cluster Tracking to People Counting. In *Workshop on Performance Evaluation of Tracking and Surveillance*, Copenhagen, Denmark, June 1 2002. Cited on page(s): 21, 42, 73
- [286] J. Pers, M. Bon, S. Kovacic, M. Sibila, and B. Dezman. Observation and Analysis of Large-Scale Human Motion. *Human Movement Science*, 21:295–311, 2002. Cited on page(s): 73
- [287] R. Plänkers and P. Fua. Tracking and Modeling People in Video Sequences. *Computer Vision and Image Understanding*, 81(3):285–303, 2001. Cited on page(s): 72
- [288] R. Plänkers and P. Fua. Model-Based Silhouette Extraction for Accurate People Tracking. In *European Conference on Computer Vision*, Copenhagen, Denmark, June 2002. Cited on page(s): 73
- [289] R. Plänkers and P. Fua. Articulated Soft Objects for Multiview Shape and Motion Capture. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 25(9):1182–1187, 2003. Cited on page(s): 7, 7, 9, 26, 26, 74
- [290] E. Polat, M. Yeasin, and R. Sharma. Robust Tracking of Human Body Parts for Collaborative Human Computer Interaction. *Computer Vision and Image Understanding*, 89:44–69, 2003. Cited on page(s): 19, 21, 42, 74
- [291] F. Porikli. Trajectory Distance Metric using Hidden Markov Model based Representation. In *6th IEEE International Workshop on Performance Evaluation of Tracking and Surveillance*, Prague, May 11-14 2004. Cited on page(s): 75

- [292] A. Prati, I. Mikić, R. Cucchiara, and M. M. Trivedi. Analysis and Detection of Shadows in Video Streams: A Comparative Evaluation. In *Computer Vision and Pattern Recognition Conference*, Hawaii, USA, December 2001. Cited on page(s): 72
- [293] A. Prati, I. Mikić, M.M. Trivedi, and R. Cucchiara. Detecting Moving Shadows: Algorithms and Evaluation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 25(7):918–923, 2003. Cited on page(s): 10, 11, 74
- [294] M. Siddiqui R. Rosales, J. Alon, and S. Sclaroff. Estimating 3D Body Pose using Uncalibrated Cameras. In *Computer Vision and Pattern Recognition*, Kauai Marriott, Hawaii, December 9-14 2001. Cited on page(s): 24, 72
- [295] M.M. Rahman and A. Robles-Kelly. A Tuned Eigenspace Technique for Articulated Motion Recognition. In *European Conference on Computer Vision*, Graz, Austria, May 7-13, 2006. Cited on page(s): 36, 77
- [296] D. Ramanan, D.A. Forsyth, and A. Zisserman. Strike a Pose: Tracking People by Finding Stylized Poses. In *Computer Vision and Pattern Recognition*, San Diego, California, USA, June 20-25, 2005. Cited on page(s): 8, 8, 9, 23, 23, 23, 30, 43, 76
- [297] D. Ramanan and C. Sminchisescu. Training Deformable Models for Localization. In *Computer Vision and Pattern Recognition*, New York City, New York, USA, June 17-22, 2006. Cited on page(s): 23, 77
- [298] C. Rao, A. Yilmaz, and M. Shah. View-Invariant Representation and Recognition of Actions. *Journal of Computer Vision*, 50(2):203–226, 2002. Cited on page(s): 40, 73
- [299] F. Remondino. 3-D Reconstruction of Static Human Body Shape from Image Sequence. *Computer Vision and Image Understanding*, 93:65–85, 2004. Cited on page(s): 75
- [300] H. Ren and G. Xu. Human Action Recognition with Primitive-based Coupled-HMM. In *International Conference on Pattern Recognition*, Quebec, Canada, August 11-15 2002. Cited on page(s): 38, 73
- [301] H. Ren, G. Xu, and S. Kee. Subject-independent Natural Action Recognition. In *International Conference on Automatic Face and Gesture Recognition*, Seoul, Korea, May 17-19, 2004. Cited on page(s): 38, 75
- [302] X. Ren, A.C. Berg, and J. Malik. Recovering Human Body Configurations using Pairwise Constraints between Parts. In *International Conference on Computer Vision*, Beijing, China, Oct 15-21, 2005. Cited on page(s): 23, 30, 42, 43, 76
- [303] L. Reng, T.B. Moeslund, and E. Granum. Finding Motion Primitives in Human Body Gestures. In S. Gibet, N. Courty, and J.-F. Kamps, editors, *GW 2005*, number 3881 in LNAI, pages 133–144. Springer Berlin Heidelberg, 2006. Cited on page(s): 41, 77
- [304] Y. Ricquebourg and P. Bouthemy. Real-Time Tracking of Moving Persons by Exploiting Spatio-Temporal Image Slices. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(8):797–808, 2000. Cited on page(s): 35, 35, 71
- [305] J. Rittscher, A. Blake, and S.J. Roberts. Towards the Automatic Analysis of Complex Human Body Motions. *Image and Vision Computing*, 20:905–916, 2002. Cited on page(s): 35, 35, 73
- [306] I. Rius, J. Varona, X. Roca, and J. Gonzàlez. Posture Constraints for Bayesian Human Motion Tracking. In *Conference on Articulated Motion and Deformable Objects (AMDO)*, Andratx, Mallorca, Spain, 11-14 July 2006. Cited on page(s): 77
- [307] G. Rizzolatti, L. Fogassi, and V. Gallese. Parietal Cortex: from Sight to Action. *Current Opinion in Neurobiology*, 7:562–567, 1997. Cited on page(s): 39, 42

- [308] G. Rizzolatti, L. Fogassi, and V. Gallese. Neurophysiological Mechanisms Underlying the Understanding and Imitation of Action. *Nature Reviews*, 2:661–670, Sept. 2001. Cited on page(s): 39, 42
- [309] T.J. Roberts, S.J. McKenna, and I.W. Ricketts. Adaptive Learning of Statistical Appearance Models for 3D Human Tracking. In *British Machine Vision Conference*, Cardiff, UK, 2002. Cited on page(s): 8, 73
- [310] T.J. Roberts, S.J. McKenna, and I.W. Ricketts. Human Pose Estimation using Learnt Probabilistic Region Similarities and Partial Configurations. In *European Conference on Computer Vision*, Prague, Czech Republic, May 11-14, 2004. Cited on page(s): 8, 8, 9, 23, 23, 30, 75
- [311] N. Robertson and I. Reid. Behaviour Understanding in Video: A Combined Method. In *International Conference on Computer Vision*, Beijing, China, Oct 15-21, 2005. Cited on page(s): 35, 42, 76
- [312] M. Roh, B. Christmas, J. Kittler, and S. Lee. Robust Player Gesture Spotting and Recognition in Low-Resolution Sports Video. In *European Conference on Computer Vision*, Graz, Austria, May 7-13, 2006. Cited on page(s): 35
- [313] M.-C. Roh, B. Christmas, J. Kittler, and S.-W. Lee. Robust Player Gesture Spotting and Recognition in Low-Resolution Sports Video. In *European Conference on Computer Vision*, Graz, Austria, May 7-13, 2006. Cited on page(s): 77
- [314] R. Ronfard, C. Schmid, and B. Triggs. Learning to Parse Pictures of People. In *European Conference on Computer Vision*, Copenhagen, Denmark, June 27-31, 2002. Cited on page(s): 8, 9, 23, 23, 30, 43, 73
- [315] R. Rosales and S. Sclaroff. Learning and Synthesizing Human Body Motion and Posture. In *The fourth International Conference on Automatic Face and Gesture Recognition*, Grenoble, France, March 2000. Cited on page(s): 24, 42, 43
- [316] D. Roth, P. Doubek, and L.V. Gool. Bayesian Pixel Classification for Human Tracking. In *IEEE Workshop on Motion and Video Computing (MOTION'05)*, Breckenridge, Colorado, Jan 2005. Cited on page(s): 14, 15, 19, 20, 20, 20, 21, 42, 76
- [317] S. Roweis and L. Saul. Nonlinear Dimensionality Reduction by Locally Linear Embedding. *SCIENCE*, 290:2323–2327, 2000. Cited on page(s): 37
- [318] M.S. Ryoo and J.K. Aggarwal. Recognition of Composite Human Activities through Context-Free Grammar based Representation. In *Computer Vision and Pattern Recognition*, New York City, New York, USA, June 17-22, 2006. Cited on page(s): 77
- [319] A. Sanfeliu and J.J. Villanueva. An Approach of Visual Motion Analysis. *Pattern Recognition Letters*, 26(3):355–368, 2005. Cited on page(s): 76
- [320] P. Sangi, J. Heikkilä, and O. Silven. Extracting Motion Components from Image Sequences using Particle Filters. In *The 12th Scandinavian Conference on Image Analysis*, Bergen, Norway, 2001. Cited on page(s): 13, 72
- [321] K. Sato and J.K. Aggarwal. Tracking and Recognizing Two-person Interactions in Outdoor Image Sequences. In *Workshop on Multi-Object Tracking*, Vancouver, Canada, July 8 2001. Cited on page(s): 37, 72
- [322] S. Schaal. Is Imitation Learning the Route to Humanoid Robots? *Trends in Cognitive Sciences*, 3(6):233–242, 1999. Cited on page(s): 39, 42
- [323] K. Schindler and H. Wang. Smooth Foreground-Background Segmentation for Video Processing. In *Asian Conference on Computer Vision*, LNCS 3852, Hyderabad, India, January 13-16 2006. Cited on page(s): 11, 77

- [324] G. Schwarz. Estimating the Dimension of a Model. *Annals of Statistics*, 6:461–464, 1978. Cited on page(s): 33
- [325] M. Shah. Understanding Human Behavior from Motion Imagery. *Machine Vision and Applications*, 14(4):210–214, 2003. Cited on page(s): 44, 74
- [326] G. Shakhnarovich, P. Viola, and T. Darrell. Fast Pose Estimation with Parameter-Sensitive Hashing. In *International Conference on Computer Vision*, Nice, France, Oct 13-16, 2003. Cited on page(s): 24, 24, 30, 42, 43, 74
- [327] Y. Sheikh and M. Shah. Bayesian Modelling of Dynamic Scenes for Object Detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 27(11):1778–1792, 2005. Cited on page(s): 11, 11, 76
- [328] Y. Sheikh, M. Sheikh, and M. Shah. Exploring the Space of Human Action. In *International Conference on Computer Vision*, Beijing, China, Oct 15-21, 2005. Cited on page(s): 39, 42, 76
- [329] J. Shi and C. Tomasi. Good Features to Track. In *Computer Vision and Pattern Recognition*, Seattle, Washington, June 1994. Cited on page(s): 15
- [330] Y. Shi, A. Bobick, and I. Essa. Learning Temporal Sequence Model from Partially Labeled Data. In *Computer Vision and Pattern Recognition*, New York City, New York, USA, June 17-22, 2006. Cited on page(s): 37, 77
- [331] H. Sidenbladh. Detecting Human Motion with Support Vector Machines. In *International Conference on Pattern Recognition*, Cambridge, UK, Aug 2004. Cited on page(s): 13, 75
- [332] H. Sidenbladh and M.J. Black. Learning Image Statistics for Bayesian Tracking. In *International Conference on Computer Vision*, Vancouver, Canada, July 9-12 2001. Cited on page(s): 8, 29, 43, 43, 72
- [333] H. Sidenbladh and M.J. Black. Learning the Statistics of People in Images and Video. *International Journal of Computer Vision*, 54(1/2/3):183–209, 2003. Cited on page(s): 8, 29, 74
- [334] H. Sidenbladh, M.J. Black, and D.J. Fleet. Stochastic Tracking of 3D Human Figures Using 2D Image Motion. In *European Conference on Computer Vision*, Dublin, Ireland, Jun, 2000. Cited on page(s): 29, 30, 43
- [335] H. Sidenbladh, M.J. Black, and L. Sigal. Implicit Probabilistic Models of Human Motion for Synthesis and Tracking. In *European Conference on Computer Vision*, Copenhagen, Denmark, 2002. Cited on page(s): 29, 29, 73
- [336] L. Sigal, S. Bhatia, S. Roth, M.J. Black, and M. Isard. Tracking Loose-Limbed People. In *Computer Vision and Pattern Recognition*, Washington DC, USA, June, 2004. Cited on page(s): 29, 30, 75
- [337] L. Sigal and M.J. Black. Measure Locally, Reason Globally: Occlusion-sensitive Articulated Pose Estimation. In *Computer Vision and Pattern Recognition*, New York City, New York, USA, June 17-22, 2006. Cited on page(s): 77
- [338] L. Sigal and M.J. Black. Predicting 3D People from 2D Pictures. In *International Conference on Articulated Motion and Deformable Objects*, Springer LNCS 4069, Andratx, Mallorca, Spain, Jul 11-14, 2006. Cited on page(s): 77
- [339] C. Sminchisescu. Consistency and Coupling in Human Model Likelihoods. In *International Conference on Automatic Face and Gesture Recognition*, Washington DC, USA, May 20-21 2002. Cited on page(s): 73

- [340] C. Sminchisescu, A. Kanaujia, Z. Li, and D. Metaxas. Discriminative Density Propagation for 3D Human Motion Estimation. In *Computer Vision and Pattern Recognition*, San Diego, California, USA, June 20-25, 2005. Cited on page(s): 24, 30, 76
- [341] C. Sminchisescu, A. Kanaujia, and D. Metaxas. Learning Joint Top-Down and Bottom-up Processes for 3D Visual Inference. In *Computer Vision and Pattern Recognition*, New York City, New York, USA, June 17-22, 2006. Cited on page(s): 77
- [342] C. Sminchisescu and B. Triggs. Covariance Scaled Sampling for Monocular 3D Body Tracking. In *Computer Vision and Pattern Recognition*, Kauai Marriott, Hawaii, December 9-14 2001. Cited on page(s): 72
- [343] C. Sminchisescu and B. Triggs. Estimating Articulated Human Motion with Covariance Scaled Sampling. *International Journal of Robotics Research*, 22(6):371–391, 2003. Cited on page(s): 27, 30, 43, 74
- [344] C. Sminchisescu and B. Triggs. Kinematic Jump Processes for Monocular 3D Human Tracking. In *Computer Vision and Pattern Recognition*, Madison, Wisconsin, USA, June 16-22, 2003. Cited on page(s): 27, 74
- [345] K. Smith, D.G. Perez, and J.M. Odobez. Using Particles to Track Varying Numbers of Interacting People. In *Computer Vision and Pattern Recognition*, San Diego, CA, June 20-25 2005. Cited on page(s): 19, 21, 42, 76
- [346] P. Smith, N. da Vitoria Lombo, and M. Shah. Temporal Boost for Event Recognition. In *International Conference on Computer Vision*, Beijing, China, Oct 15-21, 2005. Cited on page(s): 76
- [347] P. Smith, M. Shah, and N. da Vitoria Lobo. Integrating Multiple Levels of Zoom to Enable Activity Analysis. *Computer Vision and Image Understanding*, 103:33–51, 2006. Cited on page(s): 37, 77
- [348] Y. Song, L. Goncalves, and E.D. Bernardo. Monocular Perception of Biological Motion in Johansson Display. *Computer Vision and Image Understanding*, 81(3):303–327, 2001. Cited on page(s): 72
- [349] Y. Song, L. Goncalves, and P. Perona. Learning Probabilistic Structure for Human Motion Detection. In *Computer Vision and Pattern Recognition*, Kauai Marriott, Hawaii, December 9-14 2001. Cited on page(s): 72
- [350] Y. Song, L. Goncalves, and P. Perona. Unsupervised Learning of Human Motion. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 25(7):814–827, 2003. Cited on page(s): 6, 9, 74
- [351] J. Starck, G. Collins, R. Smith, A. Hilton, and J. Illingworth. Animated Statues. *Journal of Machine Vision and Applications*, 14(4):248–259, 2003. Cited on page(s): 7, 73
- [352] J. Starck and A. Hilton. Model-based Multiple View Reconstruction of People. In *International Conference on Computer Vision*, Nice, France, Oct 13-16, 2003. Cited on page(s): 7, 7, 8, 9, 74
- [353] J. Starck and A. Hilton. Spherical Matching for Temporal Correspondence of Non-Rigid Surfaces. In *ICCV*, Beijing, China, Oct 15-21, 2005. Cited on page(s): 25, 76
- [354] C. Stauffer and W.E.L. Grimson. Adaptive Background Mixture Models for Real-Time Tracking. In *Computer Vision and Pattern Recognition*, Santa Barbara, CA, June 1998. Cited on page(s): 10, 10, 12

- [355] C. Stauffer and W.E.L. Grimson. Learning Patterns of Activity Using Real-Time Tracking. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(8):747–757, 2000. Cited on page(s): 32, 71
- [356] A. Stolcke. An Efficient Probabilistic Context-Free Parsing Algorithm That Computes Prefix Probabilities. *Computational Linguistics*, 21(2):165–201, 1995. Cited on page(s): 41
- [357] M. Störring, T. Kocka, H.J. Andersen, and E. Granum. Tracking Regions of Human Skin through Illumination Changes. *Pattern Recognition Letters*, 24:1715–1723, 2003. Cited on page(s): 74
- [358] A. Sundaresan and R. Chellappa. Acquisition of Articulated Human Body Models Using Multiple Cameras. In *International Conference on Articulated Motion and Deformable Objects, Springer LNCS 4069*, Andratx, Mallorca, Spain, Jul 11-14, 2006. Cited on page(s): 77
- [359] K. Takahashi, T. Sakaguchi, and J. Ohya. Remarks on a Real-Time, Noncontact, Nonwear, 3D Human Body Posture Estimation Method. *Systems and Computers in Japan*, 31(14):1–10, 2000. Cited on page(s): 71
- [360] L. Taycher, G. Shakhnarovich, D. Demirdjian, and T. Darrell. Conditional Random People: Tracking Humans with CRFs and Grid Filters. In *Computer Vision and Pattern Recognition*, New York City, New York, USA, June 17-22, 2006. Cited on page(s): 77
- [361] C.J. Taylor. Reconstruction of Articulated Objects from Point Correspondences in a Single Image. *Computer Vision and Image Understanding*, 80(3):349–363, 2000. Cited on page(s): 6, 71
- [362] J. Tenenbaum, V. de Silva, and J. Langford. A Global Geometric Framework for Nonlinear Dimensionality Reduction. *SCIENCE*, 290:2319–2323, 2000. Cited on page(s): 37
- [363] M.N. Thalmann and H. Seo. Data-Driven Approaches to Digital Human Modeling. In *International Symposium on 3D Data Processing, Visualization, and Transmission, Thessalonica, Greece*, Sep 2004. Cited on page(s): 7, 75
- [364] C. Theobalt, M. Magnor, P. Schuler, and H.-P. Seidel. Combining 2D Feature Tracking and Volume Reconstruction for Online Video-Based Human Motion Capture. In *Pacific Conference on Computer Graphics and Applications*, Tsinghua University, Beijing, China, October 8-11 2002. Cited on page(s): 73
- [365] C. Tomasi and T. Kanade. Shape and Motion from Image Streams under Orthography: A Factorization Method. *International Journal of Computer Vision*, 9:137–154, 1992. Cited on page(s): 33
- [366] S. Toyosawa and T Kawai. Crowdedness Estimation in Public Pedestrian Space for Pedestrian Guidance. In *Conference on Intelligent Transportation Systems*, Vienna, Austria, September 13-16 2005. Cited on page(s): 76
- [367] M. Trivedi, K. Huang, and I. Mikić. Intelligent Environments and Active Camera Networks. In *Conference on System, Man and Cybernetics*, Nashville, Tennessee, October 8-11 2000. Cited on page(s): 71
- [368] M.M. Trivedi, I. Mikić, and S.K. Bhonsle. Active Camera Networks and Semantic Event Databases for Intelligent Environments. In *Workshop on Human Modeling, Analysis and Synthesis at CVPR*, Hilton Head Island, South Carolina, June 13-15 2000. Cited on page(s): 71
- [369] N. Ukita and T. Matsuyama. Real-Time Cooperative Multi-Target Tracking by Communicating Active Vision Agents. *Computer Vision and Image Understanding*, 97(2):137–179, 2005. Cited on page(s): 21, 43, 76

- [370] R. Urtasun, D.J. Fleet, and P. Fua. Monocular 3-D Tracking of the Golf Swing. In *Computer Vision and Pattern Recognition*, San Diego, California, USA, June 20-25, 2005. Cited on page(s): 29, 76
- [371] R. Urtasun, D.J. Fleet, and P. Fua. 3D People Tracking with Gaussian Process Dynamical Models. In *Computer Vision and Pattern Recognition*, New York City, New York, USA, June 17-22, 2006. Cited on page(s): 29, 77
- [372] R. Urtasun, D.J. Fleet, A. Hertzmann, and P. Fua. Priors for people tracking from small training sets. In *International Conference on Computer Vision*, Beijing, China, Oct 15-21, 2005. Cited on page(s): 29, 30, 43, 76
- [373] R. Urtasun and P. Fua. 3D Human Body Tracking using Deterministic Temporal Motion Models. In *European Conference on Computer Vision*, Prague, Czech Republic, May 11-14, 2004. Cited on page(s): 29, 29, 43, 75
- [374] A. Utsumi and N. Tetsutani. Human Detection using Geometrical Pixel Value Structures. In *International Conference on Automatic Face and Gesture Recognition*, Washington DC, USA, May 20-21 2002. Cited on page(s): 14, 14, 73
- [375] N. Vasvani, A. Roy Chowdhury, and R. Chellappa. Activity Recognition Using the Dynamics of the Configuration of Interacting Objects. In *Computer Vision and Pattern Recognition*, Madison, Wisconsin, USA, June 16-22, 2003. Cited on page(s): 33, 33, 74
- [376] D.D. Vecchio, R.M. Murray, and P. Perona. Decomposition of Human Motion into Dynamics-based Primitives with Application to Drawing Tasks. *Automatica*, 39(12):2085–2098, 2003. Cited on page(s): 40, 74
- [377] A. Veeraraghavan, R. Chellappa, and A.K. Roy-Chowdhury. The Function Space of an Activity. In *Computer Vision and Pattern Recognition*, New York City, New York, USA, June 17-22, 2006. Cited on page(s): 77
- [378] A. Veeraraghavan, A.K. Roy-Chowdhury, and R. Chellappa. Matching Shape Sequences in Video with Applications in Human Movement Analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 27(12):1896–1909, 2005. Cited on page(s): 76
- [379] A. Veeraraghavan, R. Chellappa, and A. Roy-Chowdhury. The Function Space of an Activity. In *Computer Vision and Pattern Recognition*, New York City, New York, USA, June 17-22, 2006. Cited on page(s): 32, 42
- [380] P. Viola and M. Jones. Rapid Object Detection using a Boosted Cascade of Simple Features. In *Computer Vision and Pattern Recognition*, Kauai Marriott, Hawaii, December 9-14 2001. Cited on page(s): 8, 13
- [381] P. Viola, M.J. Jones, and D. Snow. Detecting Pedestrians Using Patterns of Motion and Appearance. In *International Conference on Computer Vision*, Nice, France, 13-16 October 2003. Cited on page(s): 74
- [382] P. Viola, M.J. Jones, and D. Snow. Detecting Pedestrians Using Patterns of Motion and Appearance. *International Journal of Computer Vision*, 63(2):153–161, 2005. Cited on page(s): 13, 76
- [383] S. Wachter and H.-H. Nagel. Tracking of Persons in Monocular Image Sequences. In *Workshop on Motion of Non-Rigid and Articulated Objects*, Puerto Rico, USA, 1997. Cited on page(s): 22
- [384] S. Wachter and H.-H. Nagel. Tracking Persons in Monocular Image Sequences. *Computer Vision and Image Understanding*, 74(3):174–192, 1999. Cited on page(s): 26, 27, 27

- [385] H. Wang and D. Suter. Background Initialization with a New Robust Statistical Approach. In *Int. Workshop on Visual Surveillance and Performance Evaluation of Tracking and Surveillance*, Beijing, China, Oct 15-16 2005. Cited on page(s): 12, 13, 76
- [386] H. Wang and D. Suter. A Novel Robust Statistical Method for Background Initialization and Visual Surveillance. In *Asian Conference on Computer Vision*, LNCS 3851, Hyderabad, India, January 13-16 2006. Cited on page(s): 12, 13, 77
- [387] J. Wang and B. Bodenheimer. An Evaluation of a Cost Metric for Selecting Transitions between Motion Segments. In *SIGGRAPH Symposium on Computer Animation*, 2003. Cited on page(s): 40, 74
- [388] L. Wang, W. Hu, and T. Tan. Recent Development in Human Motion Analysis. *Pattern Recognition*, 36(3):585–601, 2003. Cited on page(s): 4, 74
- [389] L. Wang, H. Ning, T. Tan, and W. Hu. Fusion of Static and Dynamic Body Biometrics for Gait Recognition. In *International Conference on Computer Vision*, Nice, France, 13-16 October 2003. Cited on page(s): 38, 74
- [390] L. Wang, T. Tan, H. Ning, and W. Hu. Silhouette Analysis-Based Gait Recognition for Human Identification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 25(12):1505–1518, 2003. Cited on page(s): 34, 38, 74
- [391] Y. Wang and G. Baci. Human Motion Estimation from Monocular Image Sequence based on Cross-Entropy Regularization. *Pattern Recognition Letters*, 24(1-3), 2003. Cited on page(s): 74
- [392] Y. Wang, H. Jiang, M. Drew, Z. Li, and G. Mori. Unsupervised Discovery of Action Classes. In *Computer Vision and Pattern Recognition*, New York City, New York, USA, June 17-22, 2006. Cited on page(s): 37, 77
- [393] D. Weinberg, R. Ronfard, and E. Boyer. Motion History Volumes for Free Viewpoint Action Recognition. In *Workshop on Modeling People and Human Interaction (PHI'05)*, Beijing, China, Oct 15 2005. Cited on page(s): 36, 76
- [394] C.R. Wren, A. Azarbayejani, T. Darrell, and A.P. Pentland. Pfunder: Real-Time Tracking of the Human Body. *Transactions on Pattern Analysis and Machine Intelligence*, 19(7):780–785, 1997. Cited on page(s): 10, 22
- [395] B. Wu and R. Navatia. Tracking of Multiple, Partially Occluded Humans based on Static Body Part Detection. In *Computer Vision and Pattern Recognition*, New York City, New York, USA, June 17-22, 2006. Cited on page(s): 17, 19, 77
- [396] B. Wu and R. Nevatia. Detection of Multiple, Partially Occluded Humans in a Single Image by Bayesian Combination of Edgelet Part Detection. In *International Conference on Computer Vision*, Beijing, China, Oct 15-21 2005. Cited on page(s): 17, 23, 76
- [397] Q.Z. Wu, H.Y. Cheng, and B.S. Jeng. Motion Detection via Change-Point Detection for Cumulative Histograms of Ratio Images. *Pattern Recognition Letters*, 26(5):555–563, 2004. Cited on page(s): 76
- [398] Y. Wu, G. Hua, and T. Yu. Tracking Articulated Body by Dynamic Markov Network. In *International Conference on Computer Vision*, Nice, France, 13-16 October 2003. Cited on page(s): 74
- [399] Y. Wu and T. Yu. A Field Model for Human Detection and Tracking. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 28(5):753–765, 2006. Cited on page(s): 16, 77
- [400] T. Xiang and S. Gong. Beyond Tracking: Modelling Action and Understanding Behavior. *International Journal of Computer Vision*, 67(1):21–51, 2006. Cited on page(s): 33, 42, 77

- [401] L.Q. Xu and P. Puig. A Hybrid Blob- and Appearance-Based Framework for Multi-Object Tracking through Complex Occlusions. In *Int. Workshop on Visual Surveillance and Performance Evaluation of Tracking and Surveillance*, Beijing, China, Oct 15-16 2005. Cited on page(s): 14, 19, 20, 20, 76
- [402] C. Yam, M. Nixon, and J. Carter. On the Relationship of Human Walking and Running: Automatic Person Identification by Gait. In *International Conference on Pattern Recognition*, Quebec, Canada, August 11-15 2002. Cited on page(s): 34, 73
- [403] M. Yamamoto, H. Mitomi, F. Fujiwara, and T. Sato. Bayesian Classification of Task-Oriented Actions Based on Stochastic Context-Free Grammar. In *International Conference on Automatic Face and Gesture Recognition*, Southampton, UK, April 10-12, 2006. Cited on page(s): 41, 42, 42, 77
- [404] C. Yang, R. Duraiswami, and L. Davis. Fast Multiple Object Tracking via a Hierarchical Particle Filter. In *International Conference on Computer Vision*, Beijing, China, Oct 15-21 2005. Cited on page(s): 14, 15, 15, 21, 21, 42, 42, 76
- [405] D.B. Yang, H.H.G. Banos, and L.J. Guibas. Counting People in Crowds with a Real-Time Network of Simple Image Sensors. In *International Conference on Computer Vision*, Nice, France, 13-16 October 2003. Cited on page(s): 17, 18, 74
- [406] H.D. Yang and S.W. Lee. Multiple Pedestrian Detection and Tracking based on Weighted Temporal Texture Features. In *International Conference on Pattern Recognition*, Cambridge, UK, 23-26 August 2004. Cited on page(s): 17, 20, 75
- [407] M.T. Yang, Y.C. Shih, and S.C. Wang. People Tracking by Integrating Multiple Features. In *International Conference on Pattern Recognition*, Cambridge, UK, 23-26 August 2004. Cited on page(s): 17, 18, 75
- [408] T. Yang, S.Z. Li, Q. Pan, and J. Li. Real-Time Multiple Objects Tracking with Occlusion Handling in Dynamic Scenes. In *Computer Vision and Pattern Recognition*, San Diego, CA, June 20-25 2005. Cited on page(s): 11, 12, 19, 76
- [409] H. Yi, D-Rajan, and L.-T. Chia. A New Motion Histogram to Index Motion Content in Video Segments. *Pattern Recognition Letters*, 26:1221–1231, 2004. Cited on page(s): 36, 75
- [410] A. Yilmaz and M. Shah. Actions Sketch: A Novel Action Representation. In *Computer Vision and Pattern Recognition*, San Diego, California, USA, June 20-25, 2005. Cited on page(s): 36, 76
- [411] A. Yilmaz and M. Shah. Recognizing Human Actions in Videos Acquired by Uncalibrated Moving Cameras. In *International Conference on Computer Vision*, Beijing, China, Oct 15-21, 2005. Cited on page(s): 38, 39, 76
- [412] H. Yu, G.-M. Sun, W.-X. Song, and X. Li. Human Motion Recognition Based on Neural Networks. In *International Conference on Communications, Circuits and Systems*, Hong Kong, China, May 2005, 2005. Cited on page(s): 36, 76
- [413] X. Yu and S.X. Yang. A Study of Motion Recognition from Video Sequences. *Computing and Visualization in Science*, 8:19–25, 2005. Cited on page(s): 41, 75
- [414] J. Zhang, R. Collins, and Y. Liu. Bayesian Body Localization Using Mixture of Nonlinear Shape Models. In *International Conference on Computer Vision*, Beijing, China, Oct 15-21, 2005. Cited on page(s): 76
- [415] J. Zhao, L. Li, and K.C. Keong. 3D Posture Reconstruction and Human Animation from 2D Feature Points. *Computer Graphics Forum*, 24(4):759–771, 2005. Cited on page(s): 76
- [416] L. Zhao and L.S. Davis. Closely Coupled Object Detection and Segmentation. In *International Conference on Computer Vision*, Beijing, China, Oct 15-21 2005. Cited on page(s): 16, 76

- [417] L. Zhao and C.E. Thorpe. Stereo- and Neural Network-Based Pedestrian Detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 1(3):148–154, 2000. Cited on page(s): 16, 71
- [418] T. Zhao and R. Nevatia. Stochastic Human Segmentation from a Static Camera. In *Workshop on Motion and Video Computing*, Orlando, Florida, November 7 2002. Cited on page(s): 73
- [419] T. Zhao and R. Nevatia. Bayesian Human Segmentation in Crowded Situations. In *Computer Vision and Pattern Recognition*, Madison, Wisconsin, June 16-22 2003. Cited on page(s): 17, 74
- [420] T. Zhao and R. Nevatia. Tracking Multiple Humans in Complex Situations. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 26(9):1208–1221, 2004. Cited on page(s): 11, 75
- [421] T. Zhao and R. Nevatia. Tracking Multiple Humans in Crowded Environments. In *Computer Vision and Pattern Recognition*, Washington DC, USA, June 2004. Cited on page(s): 10, 14, 15, 15, 19, 21, 42, 75
- [422] T. Zhao, R. Nevatia, and F. Lv. Segmentation and Tracking of Multiple Humans in Complex Situations. In *Computer Vision and Pattern Recognition*, Kauai Marriott, Hawaii, December 9-14 2001. Cited on page(s): 72
- [423] J. Zhong and S. Sclaroff. Segmenting Foreground Objects from a Dynamic Textured Background Via Robust Kalman Filter. In *International Conference on Computer Vision*, Nice, France, 13-16 October 2003. Cited on page(s): 11, 11, 74
- [424] Z. Zivkovic. Improved Adaptive Gaussian Mixture Model for Background Subtraction. In *International Conference on Pattern Recognition*, Cambridge, UK, Aug 2004. Cited on page(s): 11

Table 2

Publications on human motion capture and analysis from 2000-2006(inclusive). Papers are ordered first by the year of publication and second by the surname of the first author. Four columns allow the clarification of the contributions of the papers within the four processes. The location of the reference number (in brackets) indicates the main topic of the work and an asterisk (*) indicates that the paper also describes work at an interesting level regarding this process.

Publications 2000 - 2006 (inclusive).

Year	First author	Initialisation	Tracking	Pose estimation	Recognition
2000	Barron			[26]	
2000	Buades			[45]	
2000	Chang		*	[56]	*
2000	Davis		[84]		
2000	Deutscher		*	[90]	
2000	Felzenszwalb			[104]	
2000	Haritaoglu	*	[137]	*	*
2000	Howe		*	[153]	
2000	Ivanov		*	[170]	
2000	Karaulova	*	*	[188]	
2000	Khan	*	[193]		
2000	Oliver		[270]		
2000	Ormonoit	[274]	*	*	
2000	Ricquebourg		[304]		*
2000	Stauffer		*		[355]
2000	Takahashi		[359]	*	
2000	Taylor	[361]		*	
2000	Trivedi		[367]		
2000	Trivedi		*		[368]
2000	Zhao		[417]		
Σ	Total=20	2	8	8	2

Year	First author	Initialisation	Tracking	Pose estimation	Recognition
2001	Ambrosio			[16]	
2001	Ambrosio			[17]	
2001	Barron			[27]	
2001	Bobick				[38]
2001	Bradski		*	*	[40]
2001	Choo			[62]	
2001	Davison			[85]	
2001	Delamarre		*	[86]	
2001	Deutscher			[91]	
2001	Elgammal	*	[100]		
2001	Grammalidis	*		[123]	
2001	Gutchess		[129]		
2001	Haritaoglu		*		[133]
2001	Herda	*	*	[142]	
2001	Hoshino		*	[149]	
2001	Huang		*	[158]	
2001	Intille				[165]
2001	Ioffe		*	[167]	
2001	Khan		[192]		
2001	Li			[217]	
2001	Mikić	*	*	[239]	
2001	Moeslund	*	*	[247]	*
2001	Moeslund	*	*	[248]	
2001	Mohan		[254]		
2001	Moon		*	[258]	
2001	Ogaki		*	[268]	
2001	Pece	*	[284]		
2001	Plänkers		*	[287]	
2001	Prati		[292]		
2001	Rosales		*	[294]	
2001	Sangi		[320]		
2001	Sato		[321]		*
2001	Sidenbladh	*	*	[332]	
2001	Sminchisescu		*	[342]	
2001	Song			[348]	
2001	Song	[349]			
2001	Zhao		[422]		*
Σ	Total=37	1	9	23	4

Year	First author	Initialisation	Tracking	Pose estimation	Recognition
2002	Allen	[14]			
2002	Atsushi		[21]		
2002	Ben-Arie		*	*	[31]
2002	BenAbdelkader		*		[32]
2002	Bradski		*		[40]
2002	Cheng		*		[58]
2002	Davis		*	*	[83]
2002	Fua			[111]	
2002	Gleicher			[117]	
2002	González				[120]
2002	Halvorsen		*	[130]	
2002	Hariadi		[132]		
2002	Haritaoglu		[134]		*
2002	Herda	*		[145]	
2002	Huang		*	[162]	
2002	Ijspeert				[164]
2002	Jang		[172]	*	
2002	Jenkins				[174]
2002	Jenkins				[175]
2002	Jenkins				[176]
2002	Lee	*	*	[211]	
2002	Li		*	[218]	
2002	Metaxas	[234]			
2002	Mikić	*	*	[237]	
2002	Mittal		[243]		
2002	Moeslund	[253]	*	*	
2002	Montemerlo		[257]		
2002	Ozer		[275]	*	*
2002	Park		[280]	*	
2002	Pece		[285]		*
2002	Pers		[286]		
2002	Plänkers		*	[288]	
2002	Rao	*	*		[298]
2002	Ren		*	*	[300]
2002	Rittscher	*	*	*	[305]
2002	Roberts		*	[309]	
2002	Ronfard			[314]	
2002	Sidenbladh	*		[335]	
2002	Sminchisescu		*	[339]	
2002	Starck	[351]			
2002	Theobalt	*	*	[364]	
2002	Utsumi	*	[374]		
2002	Yam		*		[402]
2002	Zhao		[418]		
Σ	Total=44	4	12	14	14

Year	First author	Initialisation	Tracking	Pose estimation	Recognition
2003	Allen	[15]			
2003	Azoz		*	[22]	
2003	Babu				[23]
2003	Barron	[28]		*	
2003	Buxton				[48]
2003	Capellades		[52]		*
2003	Carranza	*	*	[53]	
2003	Cheung	*	*	[59]	
2003	Chowdhury				[64]
2003	Chu	*		[65]	
2003	Comaniciu		[67]		
2003	Cucchiara		[69]		
2003	Davis				[79]
2003	Demirdjian	*		[87]	
2003	Demirdjian	*		[89]	
2003	Efros				[94]
2003	Elgammal		[95]		
2003	Elgammal		[96]		
2003	Elgammal				[99]
2003	Eng		[101]		*
2003	Foster		*		[110]
2003	Gerard	*		[114]	
2003	Gonzalez		[121]	*	
2003	Herda		*	[141]	
2003	Jepson		[177]		
2003	Koschan		[197]		
2003	Krahnstoever	[200]	*	*	
2003	Liebowitz	*		[219]	
2003	Masoud				[231]
2003	Mikić	*	*	[238]	
2003	Mitchelson			[241]	
2003	Mitchelson		*	[242]	
2003	Mittal		[244]		
2003	Moeslund	*	*	[245]	
2003	Moeslund		*	[249]	
2003	Moeslund		*	[250]	
2003	Monnet		[256]		
2003	Parameswaran				[277]
2003	Plänkers		*	[289]	
2003	Polat		[290]		
2003	Prati		[293]		
2003	Shah		[325]	*	*
2003	Shakhnarovich			[326]	
2003	Sidenbladh	*	[333]	*	
2003	Sminchisescu		*	[343]	
2003	Sminchisescu		*	[344]	
2003	Song	[350]	*		*
2003	Starck	[352]		*	
2003	Störring		[357]		
2003	Vasvani				[375]
2003	Vecchio				[376]
2003	Viola		[381]		
2003	Wang				[387]
2003	Wang		[388]	*	*
2003	Wang		*	*	[389]
2003	Wang		*		[390]
2003	Wang		[391]		
2003	Wu			[398]	
2003	Yang		[405]		
2003	Zhao		[419]		
2003	Zhong		[423]		
Σ	Total=61	5	22	20	14

Year	First author	Initialisation	Tracking	Pose estimation	Recognition
2004	Agarwal			[6]	
2004	Agarwal		*	[7]	
2004	Agarwal				[12]
2004	Billard				[34]
2004	Bregler			[43]	
2004	Brostow	[44]			
2004	Calinon				[50]
2004	Cucchiara		[70]		
2004	Date			[78]	
2004	Davis				[81]
2004	Davis		[82]		
2004	Demirdjian			[88]	
2004	Elgammal				[97]
2004	Elgammal			[98]	
2004	Figueroa		[105]		
2004	Gao		[112]		*
2004	Giebel			[115]	
2004	González				[119]
2004	Grauman			[124]	
2004	Gritai				[126]
2004	Hayashi		[138]		
2004	Heikkila		[139]		
2004	Herda			[143]	
2004	Howe	*		[151]	
2004	Hu		[154]		
2004	Hu		[156]	*	*
2004	Huang	*	*	[159]	*
2004	Iwase		[171]		
2004	Junejo	*			[181]
2004	Kang		[185]	*	
2004	Krahnstoever	[199]		*	
2004	Lee	*	*	[209]	
2004	Lee	*	*	[210]	
2004	Leo				[216]
2004	Loy			[224]	
2004	Lu				[225]
2004	Lv		*		[227]
2004	Mikolajczyk		*	[240]	
2004	Moeslund		*	[251]	
2004	Mori			[261]	
2004	Murakita		[263]		
2004	Okuma		[269]		
2004	Pan		[276]		
2004	Parameswaran	[278]		*	
2004	Park				[281]
2004	Porikli				[291]
2004	Remondino	[299]			
2004	Ren				[301]
2004	Roberts			[310]	
2004	Sidenbladh	*	[331]		
2004	Sigal			[336]	
2004	Thalmann	[363]			
2004	Urtasun		*	[373]	
2004	Yang		[406]		
2004	Yang		[407]		
2004	Yi				[409]
2004	Yu				[413]
2004	Zhao		[420]		
2004	Zhao	*	[421]		
Σ	Total=59	5	18	20	16

Year	First author	Initialisation	Tracking	Pose estimation	Recognition
2005	Andersen		[18]		
2005	Balan			[25]	
2005	Beleznai		[29]		
2005	Blank				[36]
2005	Boiman				[39]
2005	Bullock	*	*	[47]	*
2005	Calinon				[51]
2005	Chalidabhongse		[55]		
2005	Chen	*	[57]		
2005	Cheung	*		[60]	
2005	Cucchiara		*		[71]
2005	Curio			[73]	
2005	Dahmane		[74]		*
2005	Dalal		[75]		
2005	Deutscher			[92]	
2005	Dimitrijevic	[93]		*	
2005	Fanti		*		[103]
2005	Guha		[128]		
2005	Herda	[144]		*	
2005	Howe			[150]	
2005	Kang		[186]		
2005	Kang		[187]		
2005	Ke				[189]
2005	Kehl			[190]	
2005	Kim		[195]		
2005	Krosshaug			[202]	
2005	Krüger	*	[203]		
2005	Kumar	[204]	*	*	
2005	Lee		[207]		
2005	Lee	*	*	[213]	
2005	Leibe		[214]		
2005	Lim		[221]		
2005	Micilotta			[235]	
2005	Moeslund	*	*	[246]	
2005	Moeslund	[252]	*	*	
2005	Mulligan		*	[262]	
2005	Navaratnam			[265]	
2005	Ong			[272]	
2005	Ormoneit			[273]	
2005	Ramanan	*		[296]	
2005	Ren			[302]	
2005	Robertson				[311]
2005	Roth		[316]		
2005	Sanfeliu		[319]		
2005	Sheikh		[327]		
2005	Sheikh				[328]
2005	Sminchisescu			[340]	
2005	Smith		[345]		
2005	Smith				[346]
2005	Starck	*	*	[353]	
2005	Toyosawa		[366]		*
2005	Ukita		[369]		
2005	Urtasun		*	[370]	
2005	Urtasun		*	[372]	
2005	Veeraraghavan		*		[378]
2005	Viola		[382]		
2005	Wang	*	[385]		
2005	Weinberg				[393]
2005	Wu	*	[396]		
2005	Wu		[397]		
2005	Xu		[401]		
2005	Yang		[404]		
2005	Yang		[408]		
2005	Yilmaz				[410]
2005	Yilmaz				[411]
2005	Yu		*		[412]
2005	Zhang			[414]	
2005	Zhao			[415]	
2005	Zhao		[416]		
Σ	Total=69	4	28	23	14

Year	First author	Initialisation	Tracking	Pose estimation	Recognition
2006	Agarwal		*	[8]	
2006	Ahmad				[13]
2006	Antonini	*	[19]		
2006	Balan			[24]	
2006	Berclaz	[33]			
2006	Bissacco				[35]
2006	Bray			[42]	
2006	Buades	[46]	*	*	
2006	Cuntoor				[72]
2006	Dalal		[76]		
2006	Eng		[102]		
2006	Figueroa		[106]		
2006	Figueroa		[107]		
2006	Fihl		*		[108]
2006	Han			[131]	
2006	Heikkila		[140]		
2006	Howe		*	[152]	
2006	Hu		[155]		
2006	Huang				[160]
2006	Huerta		[163]		
2006	Jaeggli			[173]	
2006	Jiang				[179]
2006	Khan		[194]		
2006	Kristensen		[201]		
2006	Lee		[205]	*	
2006	Lee	[206]			
2006	Lee			[212]	
2006	Leichter		[215]		
2006	Lim	[220]	*		
2006	Liu				[222]
2006	Lv				[228]
2006	Menier	*		[233]	
2006	Micilotta			[236]	
2006	Moon			[259]	
2006	Mori	[260]		*	
2006	Nillius	[266]			
2006	Parameswaran				[279]
2006	Park		[282]	*	
2006	Park		[283]		
2006	Rahman				[295]
2006	Ramanan			[297]	
2006	Reng				[303]
2006	Rius			[306]	
2006	Roh				[313]
2006	Ryoo				[318]
2006	Schindler		[323]		
2006	Shi				[330]
2006	Sigal			[337]	
2006	Sigal			[338]	
2006	Sminchisescu			[341]	
2006	Smith				[347]
2006	Sundaresan	[358]			
2006	Taycher			[360]	
2006	Urtasun			[371]	
2006	Veeraraghavan				[377]
2006	Wang		[386]		
2006	Wang				[392]
2006	Wu	*	[395]		
2006	Wu	*	[399]		
2006	Xiang				[400]
2006	Yamamoto				[403]
Σ	Total=61	7	18	17	19
00-06	Total= 351	28	115	125	83

Using Hidden Markov Models for Recognizing Action Primitives in Complex Actions

V. Krüger and D. Grest

Aalborg Media Lab
Aalborg University, Copenhagen
Lautrupvang 15
2750 Ballerup
Denmark
vok@media.aau.dk

Abstract. There is biological evidence that human actions are composed out of action primitives, similarly to words and sentences being composed out of phonemes. Similarly to language processing, one possibility to model and recognize complex actions is to use *grammars* with action primitives as the alphabet. A major challenge here is that the action primitives need to be recovered first from the noisy input signal before further processing with the action grammar can be done. In this paper we combine a Hidden Markov Model-based approach with a simplified version of a condensation algorithm which allows to recover the action primitives in an observed action. In our approach, the primitives may have different lengths, no clear “divider” between the primitives is necessary. The primitive detection is done online, no storing of past data is necessary. We verify our approach on a large database. Recognition rates are slightly smaller than the rate when recognizing the singular action primitives.

1 Introduction

There is biological evidence that actions and activities are composed out of action primitives similarly to phonemes being concatenated into words [21; 7; 20].

In this sense, one can define a hierarchy of *action primitives* at the coarsest level, and then *actions* and *activities* as the higher abstract levels where actions are composed out of the action primitives while activities are, in turn, a composition of the set of actions [2; 16]¹. If the action primitives are used as an *alphabet* one can use action grammars [12; 23] to model actions and activities.

It is an open problem how to define and detect these action primitives and how to define these grammars. It is reasonable to assume that these things can only be defined in context of the specific application at hand.

If an observed complex action is given and a grammar should be used for parsing and recognition, then the first necessary step is to recover the *letters* in

¹ In the following, we define the term *action* as a sequence of action primitive of arbitrary length.

this observed action, i.e., the action primitives. Once the observed (continuous) sequence has been translated into a discrete set of symbols (letters), parsing based on the grammar description can be done.

In other words, if we have given an alphabet of action primitives P and if we define any *action* O to be a composition $O = a_1 a_2 a_3 \dots a_T$ of these action primitives, then our goal is to recover these primitives and their precise order. The same problem is also found in speech recognition where the goal is to find the right sequences of phonemes (see Sec. 2). Once we have recovered the sequence of action primitives in the observed sequence, we can identify the action through parsing. (In speech recognition, the sequence of detected phonemes is used to identify the corresponding word.)

The recovery of the action primitives is a non-trivial problem. Unlike phonemes (see also discussion in Sec. 2), action primitives can have a “long” durations and the variance of their execution speed may vary greatly. Also, action primitives can be heavily smeared out which complicates the distinction between them.

In this paper we deal with the recovery of the sequence of the action primitives out of an action, when a set (or alphabet) of action primitives is given.

In order to take into account possible noise and imperfect data, we base our approach on Hidden Markov Models (HMMs) [9; 17] and represent our action primitives with HMMs.

Thus, given a set of action primitives P where each action primitive is represented by an HMM and given an observed sequence O of these action primitives where

1. the order of the action primitives and
2. the duration of each single action primitive and the position of their boundaries

are unknown, we would like to identify the most likely sequence of action primitives in the observation sequence O .

According to the biological findings, the representation for action recognition is closely related to the representation for action synthesis (i.e. the motor representation of the action) [21; 7; 20]. This motivates us to focus our considerations in this paper to actions represented in joint space. Thus, our actions are given as sequences of joint settings. A further justification for this approach is that this action representation can then be used, in future work, to bias 3D body trackers as it operates directly on the 3D parameters that are to be estimated by the 3D tracker. However, our focus on joint data is clearly without limiting generality and our technique can be applied also to other types of action representations as long as the features live in a metric space. In our on-going research we have applied the same techniques of this paper also action recognition based on silhouettes.

This paper is structured as follows: In Sec. 2 will give an overview of related work. In Sec. 3 we will discuss our approach for the HMM-based recognition of the action primitives. In Sec. 4 we present our extensive experimental results. The paper is concluded then in Sec. 5 with final comments.

2 Related Work

The recovery of phonemes in speech recognition is a closely related to our problem and thus the techniques applied there were worthwhile to be investigated. In speech recognition, acoustic data gets sampled and quantized, followed by using Linear Predictive Coding (LPC) to compute a *cepstral* feature set. Alternatively to LPC, a Perceptual Linear Predictive (PLP) analysis [8] is often applied. In a later step, time slices are analyzed. Gaussians are often used here to compute the likelihoods of the observations of being a particular phoneme [10]. An alternative way to the Gaussians is to analyze time slices with an Artificial Neural Network [3]. Timeslices seem to work well on phonemes as the phonemes have usually a very short duration. In our case, however, the action primitives have usually a much longer duration and one would have to deal with a combinatorial problem when considering time slices.

In the following we will shortly review the most recent publications that consider the action recognition problem based on action primitives.

As mentioned above, the human visual system seems to relate the visual input to a sequence of motor primitives when viewing other agents performing an action [21; 7; 20]. These findings have gained considerable attention from the robotics community [22; 6]. In *imitation learning* the goal is to develop a robot system that is able to relate perceived actions to its own motor control in order to learn and to later recognize and perform the demonstrated actions.

In [14; 13], Jenkins *et al.* suggest applying a spatio-temporal non-linear dimension reduction technique on manually segmented human motion capture data. Similar segments are clustered into primitive units which are generalized into parameterized primitives by interpolating between them. In the same manner, they define action units (“behavior units”) which can be generalized into actions. In [11] the problem of defining motor primitives is approached from the motor side. They define a set of nonlinear differential equations that form a control policy (CP) and quantify how well different trajectories can be fitted with these CPs. The parameters of a CP for a primitive movement are learned in a training phase. These parameters are also used to compute similarities between movements. In [5; 1; 4] a HMM based approach is used to learn characteristic features of repetitively demonstrated movements. They suggest to use the HMM to synthesize joint trajectories of a robot. For each joint, one HMM is used. In [5] an additional HMM is used to model end-effector movement. In these approaches, the HMM structure is heavily constrained to assure convergence to a model that can be used for synthesizing joint trajectories.

In [15], Lu *et al.* also approach the problem from a system theoretic point of view. Their goal is to segment and represent repetitive movements. For this, they model the joint data over time with a second order auto-regressive (AR) model and the segmentation problem is approached by detection significant changes of the dynamical parameters. Then, for each motion segment and for each joint, they model the motion with a damped harmonic model. In order to compare actions, a metric based on the dynamic model parameters is defined. An approach of using meaningful instants in time is proposed by Reng *et al.* [19] where key

poses are found based on the curvature and covariance of the normalized trajectories.

3 Representing and Recognizing Action Primitives using HMMs

In order to approach the action recognition problem, we model each of the action primitives $P = \{a^1, a^2, \dots, a^N\}$ with a continuous mixture-HMM. A Hidden Markov Model (HMM) probabilistic version of a finite state machine. It is generally defined as a triplet $\lambda = (A, B, \pi)$, where A gives the transition likelihoods between states, B the observation likelihoods, conditioned on the present state of the HMM, and the starting state π (see the classics [9; 17] for a further introduction). In case of the continuous mixture HMM, the observation likelihoods are given as Gaussian mixtures with $M \geq 1$ mixtures.

Our HMMs are trained on demonstrations of different individuals and the Gaussian mixtures are able to capture the variability between them. The training results into a set of HMMs $\{\lambda_i | i = 1 \dots N\}$, one for each action primitive.

Once each action primitive is represented with an HMM, the primitives can generally simply be recognized with the classical recognition technique for HMMs, a maximum likelihood or a maximum a-posteriori classifier: Given an observation sequence O_t of an action primitive, and a set of HMMs λ_i , the maximum likelihood (ML)

$$\max_i P(O_t | \lambda_i) \quad (1)$$

identifies the most likely primitive. An alternative to the ML technique is the maximum a-posteriori (MAP) estimate that allows to take into account the likelihood of observing each action primitive:

$$\max_i P(\lambda_i | O_t) = \max_i P(O_t | \lambda_i) P(\lambda_i) \quad , \quad (2)$$

where $P(\lambda_i)$ is the likelihood that the action, represented by the HMM λ_i appears.

Recognition with HMMs

In general, the likelihood of an observation for some HMM λ_i can be computed as

$$P(O | \lambda_i) = \sum_S P(O, S | \lambda_i) \quad (3)$$

$$= \sum_S P(O | S, \lambda_i) P(S | \lambda_i) \quad (4)$$

$$= \sum_S \prod_{t=0}^T P(O_t | S_t, \lambda_i) \prod_{t=0}^T P(S_t | S_{t-1}, \lambda_i) \quad . \quad (5)$$

Here, one needs to marginalizes over all possible state sequences $S = \{S_0, \dots, S_T\}$ the HMM λ_i can pass through.

To apply this technique to our problem directly is difficult in our case: In Eq. 3-5 we evaluate at the end of the observation O and select the HMM which explains this observation best. In case of our problem, we are not able to identify when one primitive ends and where a new one stats. The problem is that we do not know *when* to evaluate, i.e. at what time steps t we should stop and do the maximum-likelihood estimation to find the most likely action primitive that was just now being observed.

Instead of keeping the HMMs distinct, our suggestion is to insert the “action identifier” i of the HMM λ_i as a random variable into Eq. (5) and to rewrite it as

$$P(O|a) = \sum_S \prod_{t=0}^T P(O_t|S_t, i_t)P(S_t, i_t|S_{t-1}, i_{t-1}) . \quad (6)$$

In other words, we would like to estimate at each time step the action i and the state S from the previously seen observations, or, respectively, the probability of λ_i being a model of the observed action:

$$P(S_T, i_T|O_{0:T}) = \prod_{t=0}^T P(O_t|S_t, i_t)P(S_t, i_t|S_{t-1}, i_{t-1}) \quad (7)$$

The difference in the interpretation becomes more clear when we write Eq. (7) in a recursive fashion:

$$P(S_{t+1}, i_{t+1}|O_{0:t+1}) = P(O_{t+1}|S_{t+1}, i_{t+1})P(S_{t+1}, i_{t+1}|S_t, i_t)P(S_t, i_t|O_{0:t}) \quad (8)$$

This is the classical Bayesian propagation over time. It computes at each time step t the likelihood of observing the action i_t while having observed $O_{0:t}$. If we ignore the action identifier i_t , then Eq. (8) explains the usual efficient implementation of the forward algorithm [9]. Using the random variable i_t , Eq. (8) defines a pdf across the set of states (where the state vector S_t is the concatenation of state vectors of each individual HMM) and the set of possible actions. The effect of introducing the action i might not be obvious: using i , we do not any more estimate the likelihood of an observation, given a HMM λ_i . Instead, we compute *at each time step* the probability mass function (pmf) $P(S_t, i_t|O_{0:t})$ of each state and each identity, given the observations. By marginalizing over the states, we can compute the pmf $P(i_t|O_{0:t})$ for the action at each time step. The likelihood $P(i_t|O_{0:t})$ converges to the most likely action primitive as time progresses as more data becomes available (see Fig. 1). From Fig. 1 it is apparent that the pmf $P(i_t|O_{0:t})$ will remain constant after convergence as one action primitive will have the likelihood 1 and all other primitive likelihoods have vanished. To properly evaluate the entire observation sequence, we apply a voting scheme that counts the votes after each convergence and then restarts the

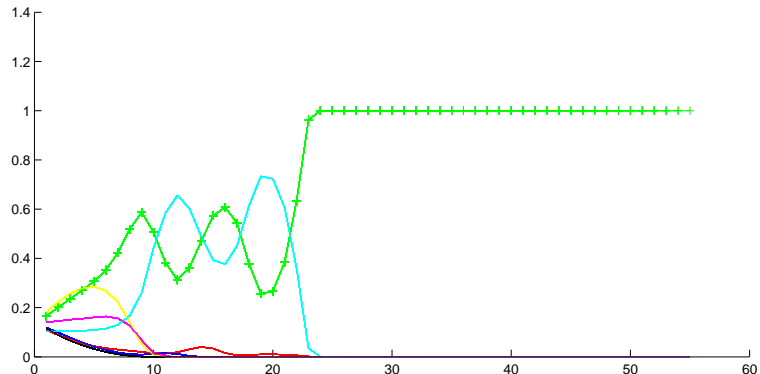


Fig. 1. shows an example for a typical behavior of the pmf $P(i_t|O_{0:t})$ for each of the actions i as time t progresses. One can see that the likelihood for one particular action (the correct one in this example, marked with "+") converges to 1 while the likelihoods for the others vanish.

HMMs. The states are initialized with the present observation likelihoods and then propagated with the transition matrix as usual. Fig. 2 shows the repeated convergence and the restarting of the HMMs. In the example shown in Fig. 2 we have used two concatenated action primitives, denoted by the green curve with the "+" and by the blue curve with the "o", respectively. The first action primitive was in the interval between 0 and 51, while the second action primitive was from sample 52 to the end. One can see that the precise time step when primitive 1 ended and when primitive 2 started cannot be identified. But this does not pose a problem for our recovery of the primitives as for us the order matters but not their precise duration. In Fig. 1 a typical situation can be seen where the observed data did not give enough evidence for a fast recognition of the true action.

4 Experiments

For our experiments, we have used our MoPrim [18] database of human one-arm movements. The data was captured using a **FastTrack** Motion capture device with 4 electromagnetic sensors. The sensors are attached to the torso, shoulder, elbow and hand (see Fig. 3). Each sensor delivers a $6D$ vector, containing $3D$ position and $3D$ orientation thus giving a $24D$ sample vector at each time-step (4 sensors with each $6D$). The MoPrim database consists of 6 individuals, showing 9 different actions, with 20 repetitions for each. The actions in the database are simple actions such as *point forward*, *point up*, *come here*, *stop!*. Each sequence consists of ≈ 60 -70 samples and each one starts with 5 samples of the arm in a resting position where it is simply hanging down.

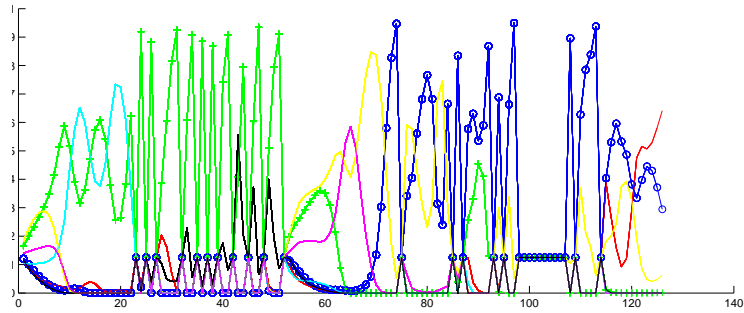


Fig. 2. shows an example for a typical behavior of the pmf $P(i_t|O_{0:t})$ as time t progresses. The input data consisted of two action primitives: first, action primitive “2”, marked with “+”, then, action primitive “3”, marked with “o”. One can see that until \approx sample 52 the system converges to action “2”, after sample 70, the system converges to primitive 3. The length of the first sequence is 51 samples, the length of sequence 2 is 71 samples.

Instead of using the sensor positions directly, we transform the raw $24D$ sensor data into joint angles: one elbow angle, one shoulder angle between elbow, shoulder and torso and a 3D orientation of the normal of this shoulder-elbow-torso-triangle. The orientation of the normal is given with respect to the normal of this triangle when the arm is in resting position. All angles are given in radians. No further processing of the MoPrim data was done.

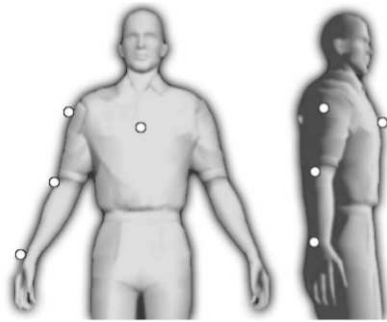


Fig. 3. marks the positions of the magnetic sensor on the human body.

We have carried out several different experiments:

1. In the first test, we tested for invariance with respect to the performing human. We have trained nine HMM for nine action. Each of the HMMs was trained on 6 individuals and all the 20 repetitions of the actions. The

recognition testing was then carried out on the remaining individual (leave-one-out-strategy). The HMMs we use were mixture HMMs with 10 states and 5 mixtures per state.

2. In this test, we tested for invariance with respect to the variations within the repetitions. We have trained nine HMMs for nine actions. Each HMM was trained on all individuals but only on 19 repetitions. The test set consisted of the 20th repetition of the actions.
3. As a base line reference, we have tested how good the HMMs are able to recognize the actions primitives by testing action primitive sequences of length 1. Here, the HMMs were trained as explained under 2 above. This test reflects the recognition performance of the classical maximum-likelihood approach.
4. We have repeated the above three experiments after having added Gaussian noise with zero mean and a standard deviation of $\sigma = 0$, $\sigma = 0.3$ and $\sigma = 1$ to the training and testing data. As all angles are given in radians, thus, this noise is considerable.

To achieve a good statistic we have for each test generated 10.000 test actions of random length ≤ 100 . Also, we have systematically left out each individual (action) once and trained on the remaining ones. The results below are averaged across all leave-one-out tests. In each test action, the action primitives were chosen randomly, identically and independently. Clearly, in reality there is a strong statistical dependency between action primitives so that our recognition results can be seen as a lower bound and results are likely to increase considerably when exploiting the temporal correlation by using an action grammar (e.g. another HMM).

The results are summarized in Table 1. One can see that the recognition rates of the individual action primitives is close to the general base-line of the HMMs. The recognition rates degrade with increasing noise which was to be expected, however, the degradation effect is the same for all three experiments (identities, repetition, baseline).

All actions in the action database start and end in a resting pose. To assure that the resting pose does not effect the recognition results, we have repeated the above experiments on the action primitives where the rest poses were omitted. However, the recognition results did not change notably.

5 Conclusions

In this work we have presented an approach to recover the motion primitives from an action where the motion primitives are represented with a Hidden Markov Model. The approach we have taken is to consider the joint distribution of the state and the action at the same time instead of using the classical maximum likelihood approach. The experiments show that the approach is able to successfully recover the action primitives in the action with a large likelihood. It is worth pointing out that in our experiments the pairwise appearance of action primitives was statistically independent. Thus, for the recovery of the action primitives no temporal constraints between the action primitives were used or

Leave-one-Out experiments		
Test	Noise σ	Recognition Result
Identities (Test 1)	0	0.9177
Repetitions (Test 2)	0	0.9097
Baseline (Test 3)	0	0.9417
Identities (Test 1)	0.5	0.8672
Repetitions (Test 2)	0.5	0.8710
Baseline (Test 3)	0.5	0.8649
Identities (Test 1)	1	0.3572
Repetitions (Test 2)	1	0.3395
Baseline (Test 3)	1	0.3548

Table 1. summarizes the results of our various experiments. In the experiments, the training of the HMMs were done without the test data. We tested for invariance w.r.t. identity and w.r.t. the action. The *baseline* shows the recognition results when the test action was a single action primitives.

exploited. Temporal constraints between the action primitives are later introduced at a higher level through action grammars.

In future work we will use a further HMM to learn sequences of action primitives from training examples to learn such an action grammar.

Acknowledgement This work was partially funded by PACO-PLUS (IST-FP6-IP-027657).

References

1. A. Billard, Y. Epars, S. Calinon, S. Schaal, and G. Cheng. Discovering Optimal Imitation Strategies. *Robotics and Autonomous Systems*, 47:69–77, 2004.
2. A.F. Bobick. Movements, Activity, and Action: The Role of Knowledge in the Perception of Motion. In *Royal Society Workshop on Knowledge-based Vision in Man and Machine*, London, England, February 1997.
3. H. Bourlard and N. Morgan. *Connectionist Speech Recognition: a Hybrid Approach*. Kluwer Press, 1994.
4. S. Calinon and A. Billard. Stochastic Gesture Production and Recognition Model for a Humanoid Robot. In *International Conference on Intelligent Robots and Systems*, Alberta, Canada, Aug 2-6, 2005.
5. S. Calinon, F. Guenter, and A. Billard. Goal-Directed Imitation in a Humanoid Robot. In *International Conference on Robotics and Automation*, Barcelona, Spain, April 18-22, 2005.
6. B. Dariush. Human Motion Analysis for Biomechanics and Biomedicine. *Machine Vision and Applications*, 14:202–205, 2003.
7. M. Giese and T. Poggio. Neural Mechanisms for the Recognition of Biological Movements. *Nature Reviews*, 4:179–192, 2003.
8. H. Hermansky. Perceptual linear predictive (plp) analysis of speech. *Journal of Acoustical Society of America*, 87(4):1738–1725, 1990.

9. X.D. Huang, Y. Ariki, and M.A. Jack. *Hidden Markov Models for Speech Recognition*. Edinburgh University Press, 1990.
10. X.D. Huang and M.A. Jack. Semi-continuous hidden markov models for speech signals. *Computer Speech and Language*, 3:239–252, 1989.
11. A.J. Ijspeert, J. Nakanishi, and S. Schaal. Movement Imitation with Nonlinear Dynamical Systems in Humanoid Robots. In *International Conference on Robotics and Automation*, Washington DC, USA, May, 2002.
12. Y. Ivanov and A. Bobick. Recognition of Visual Activities and Interactions by Stochastic Parsing. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(8):852–872, 2000.
13. O.C. Jenkins and M. Mataric. Deriving Action and Behavior Primitives from Human Motion Capture Data. In *International Conference on Robotics and Automation*, Washington DC, USA, May, 2002.
14. O.C. Jenkins and M.J. Mataric. Deriving Action and Behavior Primitives from Human Motion Data. In *International Conference on Intelligent Robots and Systems*, pages 2551–2556, Lausanne, Switzerland, Sept.30 – Oct.4, 2002.
15. C. Lu and N. Ferrier. Repetitive Motion Analysis: Segmentation and Event Classification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 26(2):258–263, 2004.
16. H.-H. Nagel. From Image Sequences Towards Conceptual Descriptions. *Image and Vision Computing*, 6(2):59–74, 1988.
17. L. R. Rabiner and B. H. Juang. An introduction to hidden Markov models. *IEEE ASSP Magazine*, pages 4–15, January 1986.
18. L. Reng, T. Moeslund, and E. Granum. Finding motion primitives in human body gestures. In S. Gibet, N. Courty, and J.-F. Kamp, editors, *GW 2005*, pages 133–144. Springer, 2006.
19. L. Reng, T.B. Moeslund, and E. Granum. Finding Motion Primitives in Human Body Gestures. In S. Gibet, N. Courty, and J.-F. Kamps, editors, *GW 2005*, number 3881 in LNAI, pages 133–144. Springer Berlin Heidelberg, 2006.
20. G. Rizzolatti, L. Fogassi, and V. Gallese. Parietal Cortex: from Sight to Action. *Current Opinion in Neurobiology*, 7:562–567, 1997.
21. G. Rizzolatti, L. Fogassi, and V. Gallese. Neurophysiological Mechanisms Underlying the Understanding and Imitation of Action. *Nature Reviews*, 2:661–670, Sept. 2001.
22. S. Schaal. Is Imitation Learning the Route to Humanoid Robots? *Trends in Cognitive Sciences*, 3(6):233–242, 1999.
23. A. Stolcke. An Efficient Probabilistic Context-Free Parsing Algorithm That Computes Prefix Probabilities. *Computational Linguistics*, 21(2):165–201, 1995.

Learning and Recognition of Object Manipulation Actions Using Linear and Nonlinear Dimensionality Reduction

Isabel Serrano Vicente and Danica Kragic
Computational Vision and Active Perception Laboratory
Centre for Autonomous Systems
Royal Institute of Technology, Stockholm, Sweden
497600@celes.unizar.es, dani@kth.se

Abstract—Neuroscientific and physiological literature states that the core of developmental learning in humans is by watching another person performing a task. This has also motivated the research in the robotics area of learning by imitation and robot programming through demonstration. There is an extensive amount of work dealing with issues of *what, when and how* to imitate.

In this work, we perform an extensive statistical evaluation for learning and recognition of object manipulation actions. We concentrate on single arm/hand actions but study the problem of modeling and dimensionality reduction for cases where actions are very similar to each other in terms of arm motions. For this purpose, we evaluate a linear and a nonlinear dimensionality reduction techniques: Principal Component Analysis and Spatio-Temporal Isomap. Classification of query sequences is based on different variants of Nearest Neighbor classification. We thoroughly describe and evaluate different parameters that affect the modeling strategies and perform the evaluation with a training set of 20 people.

I. INTRODUCTION

Human-computer interaction, surveillance, video retrieval are just a few example areas that require human activity recognition, [1]. In robotics, recognition of human activity has been used extensively for robot task learning through imitation and demonstration, [2], [3], [4], [5], [6], [7], [8], [9], [10]. For humans, one of the fundamentals of social behaviors is the understanding of each others intentions through perception and recognition of performed actions. However, the neural and functional mechanisms underlying this ability in human are still poorly understood, [11] which makes it difficult to develop the necessary models needed for designing a robot system that can learn just by observing a human or another robot performing an action. The recent discovery of *mirror neurons* in monkey's brain [12], [13] has nevertheless introduced new hypotheses and ideas about the process of imitation and its role in the evolution.

It has been shown in [14] that an action perceived by a human can be represented as a sequence of clearly segmented *action units*. This motivates the idea that the action recognition process may be considered as an interpretation of the continuous human behaviors which, in its turn, consists of a sequence of action primitives [8] such as *reaching, picking up, putting down*. In relation, learning *what* and *how* to imitate has been recognized as an important problem, [10]. It has been argued that the data used for imitation has statistical

dependencies between the activities one wishes to model and that each activity has a rich set of features that can aid both the modeling and recognition process. While in the computer vision community, most work on modelling human motion has concentrated on cyclic motions such as walking or running, [1], examples in robotics consider mainly non-cyclic actions. In [5], [6], a framework for acquiring hand-action models by integrating multiple observations based on gesture spotting is proposed. [7] present a gesture imitation system where the focus is put on the coordinate system transformation (*View-Point Transformation*) so that the teacher induced gesture is transformed into the robot's egocentric system. This way the robot *observes* the gesture as it was generated by the observer himself. [8] approaches the task learning problem by proposing a system for deriving behavior vocabularies or simple action models that can be used for more complex task extraction and learning. [10] presents a learning system for one and two-hand motions where the robot's body constraints are considered as a part of the optimal trajectory generation process. An interesting trend to note here is that most of the studies are based on a **single** user generated motion. A natural question to pose here is how the underlying modeling methods scale and apply for cases when the robot is supposed to learn from multiple teachers. The experimental evaluation conducted in our work is based on 20 people.

In robotics, many of the systems used for imitation are based on generative models such as Hidden Markov Models, [5], [10]. Generative models define a joint probability distribution over observations and state variables. For modeling of the observation process and enumerating all possible sequence of observations, it is commonly assumed that these are atomic and independent. This affects the inference problem which makes generative models intractable for multiple overlapping features of the observation or complex dependencies of observations at multiple time steps. One of the solutions to this problem may be the use of discriminative models such as Conditional Random Fields, [15].

In this work, we perform an extensive statistical evaluation for learning and recognition of object manipulation actions. Single arm/hand actions are considered with a specific focus on the problem of modeling and dimensionality reduction for cases where actions are very similar to each other in terms

of arm motions. For this purpose, we evaluate a linear and a nonlinear dimensionality reduction techniques: Principal Component Analysis and Spatio-Temporal Isomap. Classification of query sequences is based on a combination of clustering and different variants of Nearest Neighbor classifiers. For both methods, we thoroughly describe and evaluate different parameters that affect the modeling strategies and perform the evaluation with a training set of 20 people. To our knowledge, there are no examples in the field of robotics where such a large set of people was considered. Similar to [16], the results can be used to enable a more sophisticated probabilistic modeling and recognition of actions and provide a modeling basis for methods such as those presented in [8], [10].

Thus, the questions we wanted to answer with the current study were:

- What modeling strategies are suitable for action representation and recognition purposes?
- Is it possible to learn action when we do not have the knowledge of the task or the embodiment (kinematic structure) of the teacher?
- Is it possible to distinguish between very similar actions such as *pick up* and *push* an object?
- Is it enough to only observe the motion of the arm/hand or does the motion of the object have to be included in the modeling process?

This paper is organized as follows. In Section II we describe the experimental setting and collection of training data. In Section III we give a short overview of dimensionality reduction techniques and present details of their implementation in Section IV. Experimental evaluation is summarized in Section V and paper concluded in Section VI.

II. DATA COLLECTION AND PREPROCESSING

We follow the classical approach to activity recognition through training and testing steps. Training step refers to a learning step where the data is collected, labelled and used to find an appropriate representation space for the data. The system learns a model for each activity which is then used for the classification of new actions in the testing step. The four activities considered in this work are:

- 1) Push forward an object placed on a table (P);
- 2) Rotate an object placed on table (R);
- 3) Pick up the object placed on the table (PU) and
- 4) Put down an object on a table (PD).

Notations P, R, PU, PD are used to denote different actions in the experimental evaluation in Section V.

Fig. 1 shows two example images stored during a push activity training - the activity is performed with the object being placed at two different heights. To motivate the choice of these activities, let us consider a robot being a part of a coffee drinking scenario. A *pick up* activity could be representing the fact of picking up the cup to take a swig of coffee; *put down* an object could represent leaving the cup of coffee after taking a swig, *rotate* an object would be similar to fold a napkin placed on the table, and finally, let

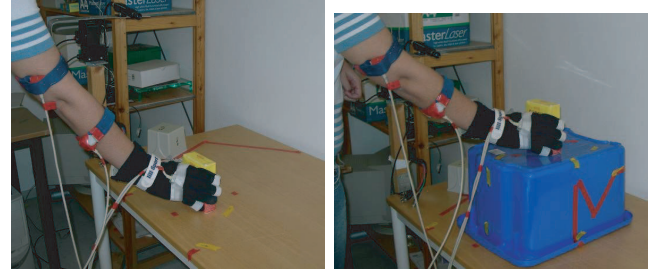


Fig. 1. Left) An example of pushing forward an object on the table and Right) An example of pushing forward an object on the box

us suppose that the person who sat down in front of you taking a coffee asks for the sugar bowl close to you and you *push* the bowl sliding over the table to bring it closer to him/her. The activities considered in this work are major building blocks of any similar task.

To generate the measurements for the training data, a Nest of Birds sensor was used, see Fig. 2 (right). The Nest of Birds simultaneously tracks the position and orientation of four sensors, referred to transmitter emitting pulsed DC magnetic field. The placement of the sensors is shown Fig. 1: thumb, hand, lower arm and upper arm. The persons involved in the study were not trained in any special way - each action started by having an arm in a relaxed, vertical position. Apart from the variation in their height and velocity with which an action was performed, the following variations were introduced to the training data:

- The objects were put on two different heights
- The person was standing at three different angles with respect to the table: 0, 30 and 60 degrees

Each action was performed three times for all combinations of the above heights and orientations resulting in total 18 training sequences per person and action thus 360 training sequences for each action.

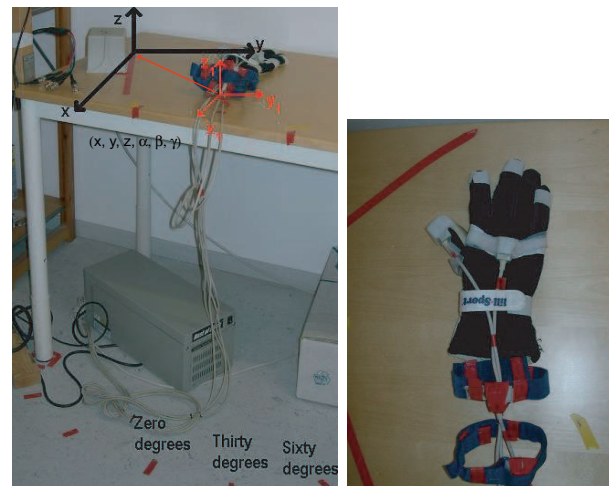


Fig. 2. Left) Nest of Birds sensor, and Right) Glove with the four sensors.

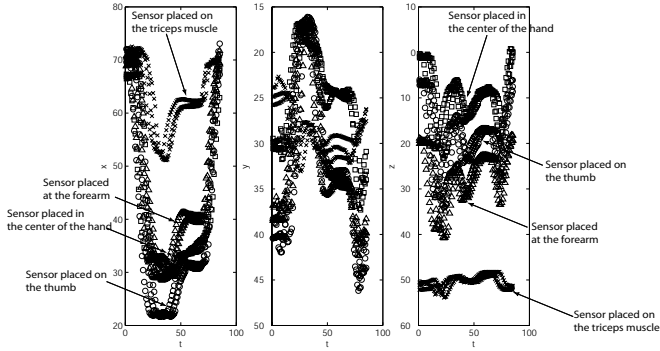


Fig. 3. Sensor measurements retrieved for three trials of a "rotate" activity.

III. DIMENSIONALITY REDUCTION

Finding low-dimensional data model hidden in the high-dimensional observations is one of the key problems in the area of activity modeling and recognition. In the current study, we have evaluated two dimensionality reduction methods. The first is the classical PCA which finds a low-dimensional embedding of the input data where the principal components are chosen such that they maximally explain the variance in the data. Since each data point is reconstructed by a suitable linear combination of the principal components, this method is applicable for cases where the assumption of linearity holds. However, for cases where the data represents essential nonlinear structures, PCA and similar techniques fail to detect the intrinsic dimensionality and model for the data. Therefore, we also evaluate a nonlinear dimensionality reduction approach proposed in [8] which is based on the isometric feature mapping or Isomap, [17].

A. Principal Component Analysis - PCA

PCA is commonly used for data dimensionality reduction, [18]. This method retains those characteristics of the data set that contribute most to its variance, by keeping lower-order principal components and ignoring higher-order ones. The idea is that such low-order components often contain the "most important" aspects of the data and the high-order components often introduce more redundant information than new one. Therefore, the error introduced by ignoring the higher-order components is not significant if the assumption of linearity holds.

In relation to human motion modeling, the use of PCA for representation of temporal curves is common. It provides a statistical model of the variation present in the training set and can thus be used to construct a probabilistic prior for motion tracking based on Bayesian methods, [16].

B. Isometric Feature Mapping - Isomap

The main idea of Isomap, [17] is to find the intrinsic geometry of the data by computing the geodesic manifold distances between all pairs of data points. Once the geodesic distances are estimated, multidimensional scaling is applied which removes nonlinearities in the data and produces a coordinate space intrinsic to the underlying manifold.

Since the training data in our system are represented in a global coordinate system (robot centered), the system should be able to perform disambiguation of spatially proximal data that are structurally different (*pick up* and *put down*) as well as model the correspondence of spatially distal data points that share common structure (actions performed at different heights). An extension of the classical Isomap, the ST-Isomap, proposed in [8] is a method that satisfies the above requirements. Implementation details are presented in Section IV-C.

C. Clustering Methods

We have evaluated two clustering techniques in connection to PCA based action classification: *k*-means clustering and Gustafson-Kessel clustering. *k*-means clustering [18] is a partitioning method in which clusters are mutually exclusive (hard partitioning method). Clustering algorithms group sample points, \mathbf{m}_j into c clusters. The set of cluster prototypes or centers is defined as $\mathbf{C} = [\mathbf{c}^{(1)}, \dots, \mathbf{c}^{(c)}]$ where

$$\mathbf{c}^{(i)} = \frac{\sum_{j=1}^d u_{ij} \mathbf{m}_j}{\sum_{j=1}^d u_{ij}} \quad i = 1, 2, \dots, c \quad (1)$$

where $u_{ij} \in \mathbf{U}$ denotes the membership of \mathbf{m}_j in the i th cluster and \mathbf{U} is known as the *partition matrix*.

For the classical *k*-means clustering, the hard partitioning space is defined as:

$$M_h = \{\mathbf{U} \in V_{cd} : u_{ij} \in \{0, 1\}, \forall (i, j); \sum_{i=1}^c u_{ij} = 1; 0 < \sum_{i=1}^d u_{ij} < d, \forall i\} \quad (2)$$

The objective function we have to minimize is:

$$J_h(\mathbf{M}; \mathbf{U}, \mathbf{C}) = \sum_{i=1}^c \sum_{j=1}^d u_{ij} d_A^2(\mathbf{m}_j, \mathbf{c}^{(i)}) \quad (3)$$

where A is a norm-inducing matrix and d_A^2 represents the distance measure

$$d_A^2 = (\mathbf{m}_j, \mathbf{c}^{(i)}) = \|\mathbf{m}_j - \mathbf{c}^{(i)}\|_A^2 = (\mathbf{m}_j - \mathbf{c}^{(i)})^T \mathbf{A} (\mathbf{m}_j - \mathbf{c}^{(i)}) \quad (4)$$

The above condition of hard membership can be relaxed so that each sample point has some graded or "fuzzy" membership in a cluster. The incorporation of probabilities (or graded memberships) may improve the convergence of the clustering method compared to the classical *k*-means method. In addition, we do not have to assume anymore that the samples belong to spherical clusters.

We shortly describe the method used in our work also known as Gustafson-Kessel (GK) clustering. First, we define a fuzzy partition space as:

$$M_f = \{\mathbf{U} \in V_{cd} : u_{ij} \in [0, 1], \forall (i, j); \sum_{i=1}^c u_{ij} = 1; 0 < \sum_{i=1}^d u_{ij} < d, \forall i\} \quad (5)$$

Here, fuzzy objective function is a least-squares functional:

$$J_f(\mathbf{M}; \mathbf{U}, \mathbf{C}) = \sum_{i=1}^c \sum_{j=1}^d (u_{ij})^w d_A^2(\mathbf{m}_j, \mathbf{c}^{(i)}) \quad (6)$$

where w is a weighting factor $w = [1, \infty)$. Gustafson-Kessel method is a variation of fuzzy clustering algorithms which allows the samples to belong to several clusters simultaneously, with different degrees of membership. It employs an adaptive distance norm in order to detect clusters of different geometrical shapes in the data set. Specifically, each cluster has its own norm-inducing matrix $\mathbf{A}^{(i)}$:

$$d_{\mathbf{A}^{(i)}}^2 = (\mathbf{c}_l^{(i)} - \mathbf{m}_j)^T \mathbf{A}^{(i)} (\mathbf{c}_l^{(i)} - \mathbf{m}_j) \quad (7)$$

where

$$\mathbf{A}^{(i)} = (|\mathbf{F}^{(i)}|)^{1/(r+1)} (\mathbf{F}^{(i)})^{-1} \quad (8)$$

and

$$\mathbf{F}^{(i)} = \frac{\sum_{j=1}^d (u_{ij})^w (\mathbf{m}_j - \mathbf{c}^{(i)}) (\mathbf{m}_j - \mathbf{c}^{(i)})^T}{\sum_{j=1}^d (u_{ij})^w} \quad (9)$$

IV. IMPLEMENTATION

We give a short overview and implementation details for the methods used in this study.

A. PCA without temporal dependencies

The basic idea investigated here was that each action consists of a set of discrete poses that are represented in some high-dimensional space since. These action are gathered by 1,2,3 or 4 sensors where each sensor provides a full pose estimate (3 translations and 3 rotations). Since the sensor used to capture the training data provides Euler angles in the reference coordinate system, we represent each angle by its sine and cosine value resulting in 9 measurements in total per sensor. This then defines the dimension of the covariance matrix, estimated in the PCA process, [18].

Our reasoning here was that different actions will vary differently along different directions. If we are able to find this directions, each action may be represented only with those ones along which the data varies the most, precisely what PCA gives us. The implementation follows the classical PCA approach: we first estimate the mean of the data, subtract it from all the samples, estimate the covariance matrix and estimate its SVD, [18]. Finally, we keep only the eigenvectors that for which eigenvalues $\lambda_n > 0.005\lambda_{max}$. In our evaluation, dependant of the number of sensors used to measure an action, the dimensionality reduction was following: single sensor (from 9 to 3), two sensors (18 to 5), three sensors (27 to 6) and four sensors (36 to 7). These values are easy to understand due to the constraints posed by the kinematic structure of the arm. Once the basic set of eigenvectors is chosen, the training data is projected to this reduced action representation space. This is done for each action separately. To ease the classification, we cluster each action representation space. For this purpose, we have used k -means and GK clustering presented in Section III-C.

In the classification stage, each testing sequence is first projected to the reduced action representation space. For each sample point in an action, the distance to the closest cluster center is estimated and the classification is based on the minimum Euclidean distance sum.

B. PCA with Temporal Dependencies

We have also evaluated a PCA approach where, similar to the studies performed on cyclic motions, [1], we took into account the temporal dependencies of the activities. To be able to estimate the covariance matrix using whole sequences, we normalized them to equal length - 85 sample points per sequence. According to the procedure described in the previous section, the dimensionality reduction was following: single sensor (from 765 to 17), two sensors (1530 to 22), three sensors (2225 to 24) and four sensors (3060 to 26). Training sequences are then projected to separate decreased spaces where each represents one of the actions. Classification of a new sequence is performed based on the minimum Euclidean distance sum.

C. ST-Isomap

For the implementation of Isomap, we adopted the approach proposed in [8]. As in the case of temporal PCA, the sequences are first normalized to equal length of 85 sample points. We shortly explain the basic idea behind the method.

- Calculate a distance matrix D^l between N local neighbors using Euclidean distances. In the current implementation, $N = 10$. For each data sample, identify common temporal neighbors (CTN) and adjacent temporal neighbors (ATN). We refer to [8] and [19] for a more detailed definition of these.
- Reduce the distances in the original matrix taking into account spatio-temporal correspondences

$$D_{S_i, S_j}^0 = \begin{cases} D_{S_i, S_j}^l / (c_{CTN} c_{ATN}) & \text{if } S_j \in CTN(S_i) \text{ and } j = i + 1, \\ D_{S_i, S_j}^l / c_{CTN} & \text{if } S_j \in CTN(S_i), \\ D_{S_i, S_j}^l / c_{ATN} & \text{if } j = i + 1, \\ \text{penalty}(S_i, S_j) & \text{otherwise.} \end{cases} \quad (10)$$

where c_{ATN} and c_{CTN} are scalar parameters and $CTN()$ denotes common temporal neighbors. In the current implementation, we set $c_{ATN} = 1$ and varied value for $c_{CTN} = [2 \ 5 \ 10 \ 100]$. Fig. 4 shows the effect of c_{CTN} parameter to the resulting embedding of the activities.

- Use D_0 to compute shortest path distance matrix D_g using Dijkstra's algorithm, [20]
- Use Multidimensional Scaling [21] to embed D_g to a lower dimensional space. We have evaluated the system for [3 4 5 6] dimensions.

V. EXPERIMENTAL EVALUATION

We present the results for i) PCA without temporal dependencies, ii) PCA with temporal dependencies and iii) ST-Isomap.

A. PCA without temporal dependencies

We have trained the system with 1, 5, 10 or 20 people. In case of a single person, we split the data in three possible combinations of two trials for training and one trial for evaluation. Similarly, this was done for the case of five and

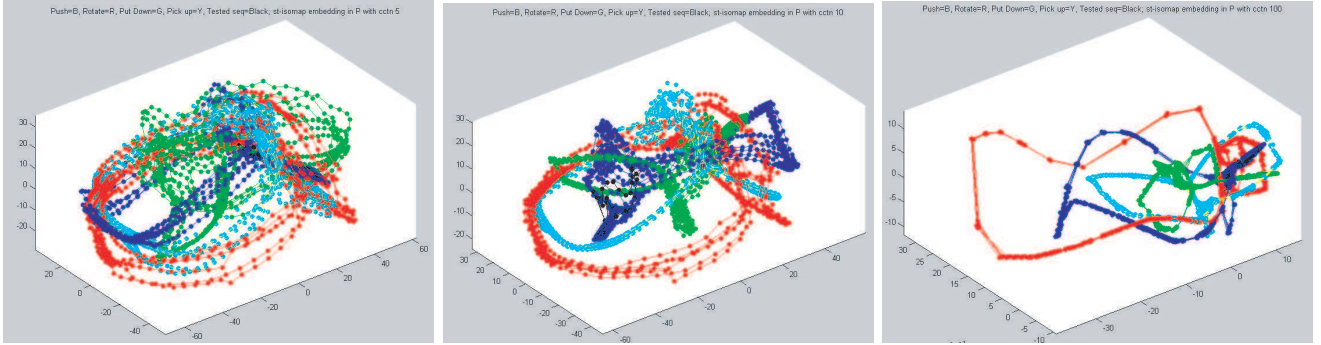


Fig. 4. Training data after estimating ST-Isomap and MDS embedding in 3 dimensions. The figures show the influence of the $cCTN$ parameter to the embedding: higher $cCTN$ brings sequences closer to each other.

ten people. For the case of twenty people, we split the trials in the three possible combinations of two for training the system and one for testing it, so we test the system three times with the demonstrations of all people. As we evaluated the system for each person three times and twenty people demonstrated the actions, it made a total of sixty tests.

Furthermore, in all cases we clustered the data using two both k -means and GK-clustering algorithm. We clustered the data to three, five and eight clusters in both cases. Here, we show only the resulting average of all the experiments and refer to [19] for a more detailed evaluation. In the forthcoming tables, the actions in the upper row are the tested sequences and the actions in the left column are the result of the classification. The results are expressed in percentage.

As explained in Section II, for each action, we have varied the position of the object (two heights) and the relative orientation of the person with respect to the table. The first experimental evaluation considered only two actions (push and rotate) where training and testing was performed on sequences captured under the same conditions (same orientation and height of the object). The average results considering different number of people in the training set as well as different numbers of sensors are summarized in Table I. We note here that we present the results for 5 clusters in more detail since it gave the highest classification rate on average. It can be seen that for only two actions, a classification rate of close to 90% is achieved. The presented results use are based on k -means clustering. GK-clustering gave approximately the same classification rate.

The second experiment to conduct was to consider all four actions, again considering the same conditions for training and testing. Due to the limited space, we show only the average classification rates for all four actions. In Table II we show how the size of the training set affects the rate given that the number of clusters is kept constant. In Table III we show how the number of clusters affect the classification rate given that the training set consist of all 20 people. Compared to the previous experiment, we can see that by adding two additional actions, the recognition rate is 30% lower on average. Again, similar results are obtained for both clustering methods.

Finally, we have evaluated the method considering all the

5 clusters								
	push	rot	push	rot	push	rot	push	rot
1pers	1s		2s		3s		4s	
push	91.8	1.6	90.5	1.3	91.5	2	92.1	2
rot	8.2	98.4	9.5	98.7	8.5	98	7.8	98
5pers	1s		2s		3s		4s	
push	80	34.4	83.3	27.8	83.3	19	87.8	30
rot	20	65.6	16.7	72.2	16.7	81	12.2	70
12pers	1s		2s		3s		4s	
push	79.6	18.5	74.5	14.8	82.4	18	82.4	14.8
rot	20.4	81.5	25.5	85.2	17.6	82	17.6	85.2
20pers	1s		2s		3s		4s	
push	83	14.7	91.4	16.7	92.5	11.9	93.1	10.8
rot	17	85.3	8.6	83.3	7.5	88.1	6.9	89.2
3 clusters								
20pers	1s		2s		3s		4s	
push	89.7	28.9	88.6	21.7	93.1	26.4	91.7	21.1
rot	10.3	71.1	11.4	78.3	6.9	73.6	8.3	78.9
8 clusters								
20pers	1s		2s		3s		4s	
push	88.1	15.6	86.9	12.5	90.6	10.6	91.1	8.9
rot	11.9	84.5	13.1	87.5	9.4	89.4	8.9	91.1

TABLE I
CLASSIFICATION RATES FOR TWO ACTIONS (PUSH, ROTATE) WHEN THE TRAINING AND TESTING WAS DONE UNDER SAME CONDITIONS (OBJECT HEIGHT, PERSONS ORIENTATION) USING k -MEANS CLUSTERING.

1 pers	1s	2s	3s	4s
average	91.4	91.1	90.2	90
5 pers	1s	2s	3s	4s
average	61.9	65	68.6	61.1
12 pers	1s	2s	3s	4s
average	60.8	60.8	63.1	61.7

TABLE II
CLASSIFICATION RATES FOR FOUR ACTIONS TRAINED AND TESTED IN SAME CONDITIONS (HEIGHT AND ORIENTATION), WITH VARYING SIZE OF THE TRAINING SET. THE NUMBER OF CLUSTERS IN k -MEANS IS 5.

variance in the data, namely that each action was performed on two different heights and in three orientations. The results are summarized in Table IV. It is obvious that, with the the recognition rates of about 40%, the simple approach considered here is not able to scale accordingly with the variation in the data. The next section presents the results of

3 clusters	1s	2s	3s	4s
average	59.4	61.4	62.2	64.1
5 clusters	1s	2s	3s	4s
average	64.7	68.4	70.6	69.8
8 clusters	1s	2s	3s	4s
average	66.5	68	68.9	70

TABLE III

CLASSIFICATION RATES FOR FOUR ACTIONS AND 20 PEOPLE TRAINED AND TESTED IN THE SAME CONDITIONS (HEIGHT AND ORIENTATION), WITH VARYING NUMBER OF CLUSTERS.

1 pers	1s	2s	3s	4s
average	37.5	30.6	37.5	37.5
5 pers	1s	2s	3s	4s
average	34.7	33.9	38.1	38.9
12 pers	1s	2s	3s	4s
average	34.3	33.7	37.5	35.6
20 pers	1s	2s	3s	4s
average	35.4	37.2	37.3	37.4

TABLE IV

CLASSIFICATION RATES FOR FOUR ACTIONS TRAINED AND TESTED IN DIFFERENT CONDITIONS, WITH VARYING SIZE OF THE TRAINING SET. THE NUMBER OF CLUSTERS USED IN k -MEANS IS FIXED TO FIVE.

the method where temporal dependencies between the data points are taken into account.

1 pers	1s	2s	3s	4s
average	41.7	36.1	38.9	27.8
5 pers	1s	2s	3s	4s
average	35.8	35.8	33.6	36.9
12 pers	1s	2s	3s	4s
average	35.2	38.1	40	40.1
20 pers	1s	2s	3s	4s
average	41	34.3	36.7	36.3

TABLE V

CLASSIFICATION RATES FOR FOUR ACTIONS TRAINED AND TESTED IN DIFFERENT CONDITIONS, WITH VARYING SIZE OF THE TRAINING SET. THE NUMBER OF CLUSTERS USED IN GK CLUSTERING IS FIXED TO FIVE.

B. Temporal PCA

We present here only the results with all four actions, where the training and testing was performed given all 20 people and actions performed in all combinations of orientations and heights. As above, as each action sequence was performed three times in all conditions, we evaluated the system taken all combinations of two testing and one training action sets.

Table VI summarizes the results for one (1s, hand), two (2s, thumb and hand), three (3s, thumb, hand, forearm) and all four (4s) sensors considered. The important thing to note is that the recognition rate is somewhat higher compared to the results in the previous section but the system still has the difficulty of discriminating between some of the actions. We believe that this is an important result. Implementing PCA with temporal dependencies requires aligned and equal length sequences which may be difficult to obtain in an

	1s				2s			
	P	R	PD	PU	P	R	PD	PU
P	50.1	42.5	12.5	29.2	50	43.3	12.5	32.5
R	8.3	33.3	3.3	10	9.2	35	3.3	15
PD	15	3.3	69.2	22.5	15	5.8	69.2	20
PU	25.8	20.8	15	38.3	25.8	15.8	15	32.5
	3s				4s			
	P	R	PD	PU	P	R	PD	PU
P	51.7	42.5	12.5	30	51.7	42.5	12.5	29.2
R	7.5	35	3.3	1.5	8.3	30	3.3	11.7
PD	14.2	5	69.2	21.7	14.2	4.2	66.7	25
PU	26.7	17.5	15	35.8	25.8	23.3	17.5	34.2

TABLE VI

CLASSIFICATION RATES FOR PCA WITH TEMPORAL DEPENDENCIES FOR FOUR ACTIONS AND 20 PEOPLE IN THE TRAINING SET.

online process where we would like to perform recognition during and not after an action has been executed. A simple “voting” approach presented in the previous section may be as suitable. Another issue that we have investigated was if the number of sensors affects the classification rate. The results are summarized in Fig. 5. We note that the difference is only marginal and that almost equal results are obtained with a single or all four sensors. This means that for actions which are very similar in arm motion placing only a single sensor on the hand or tracking only the position and orientation of the hand may be sufficient.

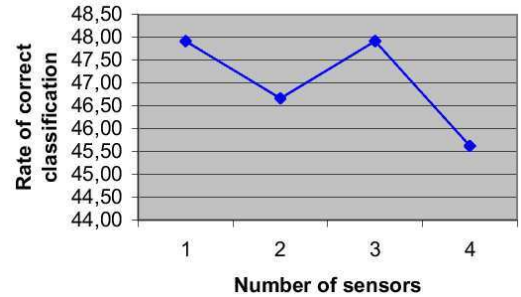


Fig. 5. The effect of number of sensors used to the classification rate.

C. ISOMAP

A non-linear dimension reduction, ST-Isomap was applied to extract a low dimensional representation for the activities. Shepard interpolation [22] was used map a query sequence to the estimated embedding. For classification, minimum Euclidean distance sum between the query samples and samples representing each activity in the embedding was used. From the training set of 20 people, we formed subsets of one, two and three persons. For each person, all four activities were considered using three trials for all combinations of three orientations and two heights. The classification was performed with the queries not included in the training set.

We have evaluated to the system with different numbers and sensors placements. In the forthcoming tables, this is denoted as: sensors placed on the i) hand (s1), ii) hand and thumb (s14), iii) hand, thumb and forearm (s142).

		s1				s14				s142			
		p	r	pd	pu	p	r	pd	pu	p	r	pd	pu
ct=2	p	55.6	0.0	11.1	11.1	61.1	11.1	5.6	22.2	50.0	50.0	33.3	50.0
3dimensions	r	5.6	77.8	0.0	0.0	5.6	61.1	61.1	38.9	44.4	50.0	50.0	16.7
3dimensions	pd	16.7	11.1	38.9	44.4	16.7	27.8	5.6	38.9	5.6	0.0	16.7	33.3
3dimensions	pu	22.2	11.1	50.0	44.4	16.7	0.0	27.8	0.0	0.0	0.0	0.0	0.0
ct=5	p	77.8	0.0	0.0	0.0	33.3	0.0	0.0	0.0	50.0	55.6	5.6	44.4
3dimensions	r	0.0	88.9	5.6	5.6	5.6	94.4	22.2	16.7	33.3	44.4	66.7	38.9
3dimensions	pd	22.2	0.0	66.7	11.1	61.1	5.6	33.3	33.3	0.0	0.0	0.0	0.0
3dimensions	pu	0.0	11.1	27.8	83.3	0.0	0.0	44.4	50.0	16.7	0.0	27.8	16.7
ct=10	p	83.3	0.0	11.1	27.8	38.9	33.3	11.1	5.6	77.8	33.3	55.6	83.3
3dimensions	r	16.7	61.1	5.6	0.0	0.0	50.0	5.6	16.7	5.6	44.4	22.2	16.7
3dimensions	pd	0.0	0.0	0.0	0.0	61.1	0.0	38.9	33.3	5.6	22.2	5.6	0.0
3dimensions	pu	0.0	38.9	83.3	72.2	0.0	16.7	44.4	44.4	11.1	0.0	16.7	0.0
ct=100	p	100.0	0.0	0.0	0.0	27.8	11.1	0.0	22.2	61.1	11.1	50.0	50.0
3dimensions	r	0.0	88.9	0.0	11.1	0.0	50.0	5.6	0.0	33.3	50.0	0.0	0.0
3dimensions	pd	0.0	5.6	61.1	16.7	33.3	0.0	5.6	0.0	0.0	0.0	11.1	0.0
3dimensions	pu	0.0	5.6	38.9	72.2	38.9	38.9	88.9	77.8	5.6	38.9	38.9	50.0

Fig. 6. ST-Isomap results with training based on one person and testing it with another one. The results show how the dimension of the embedding and sensor number affect the classification result.

		s1				s14				s142			
		p	r	pd	pu	p	r	pd	pu	p	r	pd	pu
ct=2	p	72.2	33.3	33.3	16.7	16.7	33.3	27.8	11.1	66.7	22.2	16.7	11.1
3dimensions	r	27.8	55.6	11.1	44.4	83.3	55.6	44.4	72.2	0.0	16.7	61.1	5.6
3dimensions	pd	0.0	0.0	50.0	11.1	0.0	5.6	11.1	0.0	5.6	16.7	55.6	22.2
3dimensions	pu	0.0	11.1	5.6	27.8	0.0	5.6	16.7	16.7	11.1	0.0	22.2	38.9
ct=5	p	77.8	5.6	22.2	33.3	61.1	50.0	27.8	50.0	100.0	5.6	11.1	22.2
3dimensions	r	0.0	77.8	5.6	16.7	0.0	11.1	5.6	5.6	0.0	88.9	5.6	5.6
3dimensions	pd	5.6	0.0	50.0	33.3	38.9	16.7	50.0	11.1	0.0	61.1	38.9	27.8
3dimensions	pu	16.7	16.7	22.2	16.7	0.0	22.2	16.7	33.3	0.0	0.0	44.4	11.1
ct=10	p	94.4	50.0	50.0	44.4	77.8	22.2	22.2	16.7	88.9	11.1	16.7	22.2
3dimensions	r	0.0	38.9	11.1	11.1	0.0	50.0	0.0	0.0	0.0	77.8	5.6	11.1
3dimensions	pd	0.0	0.0	33.3	27.8	11.1	22.2	55.6	44.4	11.1	5.6	44.4	0.0
3dimensions	pu	5.6	11.1	5.6	16.7	11.1	5.6	22.2	38.9	0.0	27.8	55.6	94.4
ct=100	p	77.8	16.7	16.7	38.9	77.8	11.1	61.1	50.0	88.9	5.6	0.0	11.1
3dimensions	r	16.7	66.7	11.1	11.1	5.6	55.6	33.3	16.7	0.0	83.3	5.6	5.6
3dimensions	pd	5.6	11.1	50.0	27.8	16.7	22.2	5.6	22.2	16.7	16.7	5.6	22.2
3dimensions	pu	0.0	5.6	22.2	22.2	0.0	11.1	0.0	11.1	11.1	5.6	38.9	55.6

Fig. 7. ST-Isomap results with training based on 3 persons and testing it with another one. The results show how the dimension of the embedding and sensor number affect the classification result.

Thorough experimental evaluation with different values for c_{CTN} parameter and dimensionality of the embedding space was conducted.

Fig.6 shows the results obtained by ST-Isomap with training based on a single person. The results show how the dimension of the embedding and sensor number affect the classification result. Here, parameter $c_{CTN} = 2$. Fig.7 shows a similar experiment, but here the size of the training set was three. It is interesting to notice that best results are obtained based on the sensor placed on the hand. For the future, this would motivate that only the position of the user's hand and not the complete arm joint motion is needed to recognize object manipulation sequences when ST-Isomap is used. The effect of changing the values of parameter c_{CTN} is shown in Table V-C. On average, the best results are obtained with $c_{CTN} = 5$ and the average values per action are shown in Fig. 8. From the above results, it can be seen that, compared to the PCA, ST-Isomap gives better classification results.

VI. CONCLUSION

In this work, we have performed an initial study on recognition of four object manipulation actions: pick up, put down, rotate and push. Training and testing was performed with 20 people where the manipulated object was placed on two different heights and people performing the actions multiple times at three different orientations. We believe that this study is important and shows how the variation in the training data affects the recognition rate. Most of the current systems that utilize robot imitation learning use a single person to train or teach tasks to the robot. Since the

		push	rot	pd	pu
ct = 2	push	88.9	22.2	29.2	36.1
ct = 2	rot	11.1	70.9	8.3	16.7
ct = 2	pd	0	0	51.4	16.7
ct = 2	pu	0	6.9	11.1	30.5
ct = 5	push	88.9	6.9	19.4	25
ct = 5	rot	0	79.2	4.2	9.7
ct = 5	pd	1.3	0	62.5	27.8
ct = 5	pu	9.7	13.9	13.9	37.5
ct = 10	push	90.3	18.1	25	29.2
ct = 10	rot	1.4	72.2	8.3	13.9
ct = 10	pd	2.8	2.8	50	30.5
ct = 10	pu	5.5	6.9	16.7	26.4
ct = 100	push	84.7	8.3	6.9	23.6
ct = 100	rot	5.6	80.6	5.6	4.2
ct = 100	pd	2.8	6.9	65.3	36.1
ct = 100	pu	6.9	4.2	22.2	36.1

TABLE VII

CLASSIFICATION RESULTS USING A SINGLE SENSOR PLACED ON THE HAND. TRAINING WAS PERFORMED WITH 3 PERSONS. THE RECOGNITION RATES SHOW THE DEPENDENCY ON THE PARAMETER c_{CTN} .

intention for the future is that robots will be able to learn from observing *different* and *multiple* people that perform same actions, we believe that it is important to study how different methods scale with respect to this.

In this work, we have concentrated on evaluation of dimensionality reduction using linear and nonlinear techniques. We have shown how the number of sensors and different parameters affect the classification rate. We are aware of the

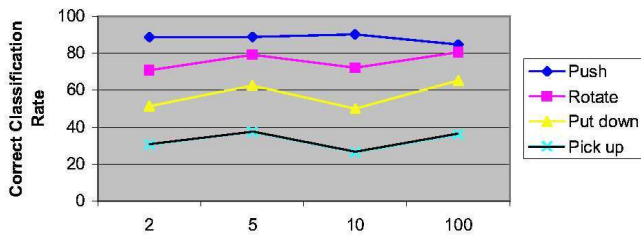


Fig. 8. Analysis of the recognition results when changing $c_{CTN} = 2, 5, 10, 100$ with a single sensor placed on the hand and training with three persons.

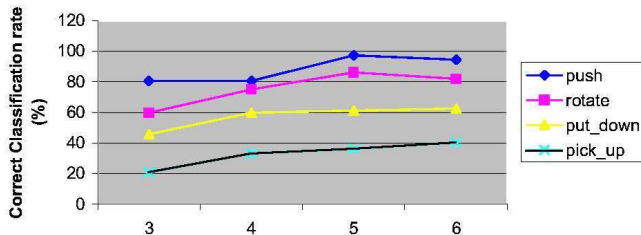


Fig. 9. Analysis of the recognition results when changing the dimension of the embedding space (3,4,5,6) with a single sensor placed on the hand and training with three persons.

fact that PCA and nearest neighbor classification are very simple techniques but we hope that our future work and work of other we evaluate more advanced techniques on the same data (which will be soon available for public access) and compare it to the results obtained in this work. We also believe that this data and evaluation follows the current trend of designing different benchmarking criteria in robotics.

Regarding the four questions posed in Section I we believe that for recognition of actions that are very similar, dimensionality reduction has to be performed with significant care in order to preserve the true variance in the data. We also believe that using the explicit knowledge of kinematic chains (arm model) may not be necessary in order to achieve satisfactory recognition rates. Finally, for some actions it is enough to provide only the measurements of the hand motions while distinguishing between *pick-up* and *put-down* would gain from including the motion of the object as well.

ACKNOWLEDGMENT

This work has been supported by EU through the project PACO-PLUS, FP6-2004-IST-4-27657. We also thank Oddest Chadwicke Jenkins for the valuable input on the implementation of ST-Isomap.

REFERENCES

- [1] J. K. Aggarwal and Q. Cai, "Human motion analysis: A review," *Computer Vision and Image Understanding: CVIU*, vol. 73, no. 3, pp. 428–440, 1999.
- [2] Y. Kuniyoshi, M. Inaba, and H. Inoue, "Learning by watching," in *IEEE Transactiond on Robotics and Automation*, vol. 10(6), pp. 799–822, 1994.
- [3] S. Schaal, "Is imitation learning the route to humanoid robots?," *Trends in Cognitive Sciences*, vol. 3, no. 6, pp. 233–242, 1999.
- [4] A. Billard, "Imitation: A review," *Handbook of brain thory and neural network*, M. Arbib (ed.), pp. 566–569, 2002.

- [5] K. Ogawara, S. Iba, H. Kimura, and K. Ikeuchi, "Recognition of human task by attention point analysis," in *IEEE International Conference on Intelligent Robot and Systems IROS'00*, pp. 2121–2126, 2000.
- [6] K. Ogawara, S. Iba, H. Kimura, and K. Ikeuchi, "Acquiring hand-action models by attention point analysis," in *IEEE International Conference on Robotics and Automation*, pp. 465–470, 2001.
- [7] M. C. Lopes and J. santos Victor, "Visual transformations in gesture imitation: What you see is what you do," in *IEEE International Conference on Robotics and Automation, ICRA04*, pp. 2375– 2381, 2003.
- [8] O. C. Jenkins and M. J. Mataric, "Performance-derived behavior vocabularies: Data-driven acquisition of skills from motion," *International Journal of Humanoid Robotics*, vol. 1, pp. 237–288, Jun 2004.
- [9] S. Ekvall and D. Kragic, "Grasp recognition for programming by demonstration tasks," in *IEEE International Conference on Robotics and Automation, ICRA'05*, pp. 748 – 753, 2005.
- [10] S. Calinon, A. Billard, and F. Guenter, "Discriminative and adaptive imitation in uni-manual and bi-manual tasks," in *Robotics and Autonomous Systems*, vol. 54, 2005.
- [11] M. Iacoboni, I. Molnar-Szakacs, V. Galles, G. Buccino, J. Mazziotta, and G. Rizzolatti, "Grasping the intentions of others with one's own mirror neuron system," *PLoS Biology*, vol. 3, no. 3, 2005.
- [12] V. Ramachandran, "Mirror neurons and imitation learning as the driving force behind the gerat leap forward in human evolution," *Edge*, vol. 69, 2000.
- [13] L. Fadiga, L. Fogassi, V. Gallese, and G. Rizzolatti, "Visuomotor neurons: Ambiguity of the discharge or 'motor perception'?", volume = 35, year = 2000," *International Journal of Psychophysiology*, no. 2-3, pp. 165–177.
- [14] D. N. et al, "The objective basis of behavior unit," *Journal of Personality and Social Pshychology*, vol. 35, no. 12, pp. 847–862, 1977.
- [15] C. Sminchiescu, A. Kanaujia, Z. Li, and D. Metaxas, "Conditional models for contextual human motion recognition," in *International Conference on Computer Vision, ICCV'05*, pp. 1808–1815, 2005.
- [16] J. Tenenbaum, V. de Silva, and J. Langford, "A global geometric framework for nonlinear dimensionality reduction," *Science*, vol. 290, pp. 2319–2323, 2000.
- [17] R. Duda, P. Hart, and D. Stork, *Pattern classification*. New York: Wiley-Interscience, 2001.
- [18] I. S. Vicente, *Human action recognition based on linear and nonlinear dimensionality reduction using PCA and Isomap*. KTH, Stockholm, Sweden: Master thesis, 2006, <http://cogvis.nada.kth.se/~danik/Isabel.pdf>.
- [19] T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein, *Introduction to Algorithms*. Second Edition. MIT Press and McGraw-Hill, 2001.
- [20] M. Cox and M. Cox, *Multidimensional Scaling*. Chapman and Hall, 2001.
- [21] D. Shepard, "A two-dimensional interpolation function for irregularly-spaced data," in *Proc. 23rd National Conference ACM*, pp. 517–524, 1968.