

Chapter 3

A Survey of Computer Vision-Based Human Motion Capture

Synopsis

In this chapter a comprehensive state-of-the-art survey within computer vision-based human MoCap is presented. The chapter consists of a survey published in 2001 [A]. The survey covers more than 130 papers published up until the summer of 2000. 70 new papers have been reviewed in order to bring the survey up-to-date as of December 2002. The additional information gathered from the new papers is presented in appendix A.

Synopsis References

- A. T.B. Moeslund and E. Granum. A Survey of Computer Vision-Based Human Motion Capture. *Computer Vision and Image Understanding*, 81(3), 2001

A Survey of Computer Vision-Based Human Motion Capture

T.B. Moeslund and E. Granum

Abstract

A comprehensive survey of computer vision-based human motion capture literature from the past two decades is presented. The focus is on a general overview based on a taxonomy of system functionalities, broken down into four processes: initialisation, tracking, pose estimation, and recognition. Each process is discussed and divided into subprocesses and/or categories of methods to provide a reference to describe and compare the more than 130 publications covered by the survey. References are included throughout the paper to exemplify important issues and their relations to the various methods. A number of general assumptions used in this research field are identified and the character of these assumptions indicates that the research field is still in an early stage of development. To evaluate the state of the art, the major application areas are identified and performances are analysed in light of the methods presented in the survey. Finally, suggestions for future research directions are offered.

3.1 Introduction

The analysis of human actions by a computer is gaining more and more interest. A significant part of this task is to register the motion, a process known as *human motion capture*. Even though this term covers many aspects, it is mainly used in connection with capturing large scale body movements, which are the movements of the head, arms, torso, and legs. Formally we here define human motion capture as the process of capturing the large scale body movements of a subject at some resolution.

We included *at some resolution* to emphasise that tracking of a subject's limbs, as well as overall tracking of a subject, are considered to fall within the above definition. Hence, human motion capture is used both when the subject is viewed as a single object and when viewed as articulated motion of a high degree of freedom skeleton structure with a number of joints.

What is *not* covered by the above definition is small scale body movements such as facial expressions and hand gestures. A thorough review of hand gestures can be found in the survey by Pavlovic *et al.* [117].

3.1.1 Application Areas

The potential applications of human motion capture are the driving force of system development, and we consider the following three major application areas: surveillance, control, and analysis.

The *surveillance* area covers applications where one or more subjects are being tracked over time and possibly monitored for special actions. A classic example is the surveillance of a parking lot, where a system tracks subjects to evaluate whether they may be about to commit a crime, e.g., steal a car.

The *control* area relates to applications where the captured motion is used to provide controlling functionalities. It could be used as an interface to games, virtual environments, or animation, or to control remotely located implements. For a comprehensive discussion of motion capture in the control application area, see [99].

The third application area is concerned with the detailed *analysis* of the captured motion data. This may be used in clinical studies of, e.g., diagnostics of orthopedic patients or to help athletes understand and improve their performance.

3.1.2 Alternative Technologies for Motion Capture

The systems used to capture human motion consist of subsystems for sensing and processing, respectively. The operational complexity of these subsystems is typically related, so that high complexity of one of them allows for a corresponding simplicity of the other. This trade-off between the complexities also relates to the use of active

versus passive sensing. Active sensing operates by placing devices on the subject and in the surroundings which transmit or receive generated signals, respectively [99].

Active sensing allows for simpler processing and is widely used when the applications are situated in well-controlled environments. That is in particular the case for the third application area, analysis, and in some of the control applications.

Passive sensing is based on "natural" signal sources, e.g., visual light or other electromagnetic wavelengths, and requires no wearable devices. An exception is when markers are attached to the subject to ease the motion capture process. Markers are not as intrusive as the devices used in active sensing. Passive sensing is mainly used in surveillance and some control applications where mounting devices on the subject is not an option.

Computer vision with the passive sensing approach has challenged active sensing within all three application areas. Even though the use of markers may seem a good compromise between passive and active sensing, it is still inconvenient for the subject (sometimes impossible) and computer vision allows in principle for touchfree and more discrete "pure" motion capture systems.

3.1.3 Content of this Paper

This paper is only concerned with computer vision-based approaches, i.e., passive sensing. It provides a compressive survey of publications in computer vision-based human motion capture from 1980 to 2000. The focus is on a general overview in relation to a functionally structured taxonomy rather than extended examples and summaries of individual papers.

Section 3.2 reviews briefly other surveys within the research field and the taxonomy of this survey is presented. It builds on the four primary functionalities of motion capture processing: initialisation, tracking, pose estimation, and recognition. Approaches and techniques are presented in relation to these functionalities in Sections 3.3 to 3.6. The descriptions will focus on general principles and similarities among various systems and methods. Section 3.7 revisits the three application areas and discusses them in light of the survey and in the context of various performance parameters and examples of state-of-the-art systems. Furthermore, suggestions for future research directions are offered. Finally, Section 3.8 concludes the survey.

3.2 Surveys and Taxonomies

Over the last two decades, the number of papers within the field of registering human body motion using computer vision has grown significantly. To structure an overview of the individual papers, their purpose, algorithms, etc. a taxonomy may

be defined to arrange them into various groups having similar characteristics.

Various categories may be used for a taxonomy, e.g. (in random order): kinetic vs kinematic, model-based vs non-model-based, 2D approaches vs 3D approaches, sensor modality (visual light, infrared (IR) light, range data, etc.), number of sensors, mobile vs stationary sensors, tracking vs recognition, pose estimation vs tracking, pose estimation vs recognition, various applications, one person vs multiple persons, number of tracked limbs, distributed vs centralised processing, various motion-type assumptions (rigid, nonrigid, elastic), etc.

Which categories to use depends on the purpose of the survey, and the various published surveys have used different taxonomies.

3.2.1 Previous Surveys

Aggarwal *et al.* [2] give an overview of various methods used prior to 1995, in articulated and elastic nonrigid motion. After a good overview of various motion types the approaches within articulated motion with or without *a priori* shape models are described. Then the elastic motion approaches are described in two categories with and without a shape model.

Cedras and Shah [23] give an overview of methods within motion extraction prior to 1995, which are all classified as belonging to optical flow or motion correspondence. The human motion capture problem is described as action recognition, recognition of the individual body parts, and body configuration estimation.

An overview of the area of human motion estimation and recognition with special focus on optical flow techniques, prior to 1996, is given by Ju [72]. The overall taxonomy is motion estimation and motion recognition, both of which are divided into subclasses.

In the survey by Aggarwal and Cai [1], which describes work prior to 1998, the same taxonomy as in Cedras and Shah [23] is applied even though they use different labels for the three classes. The classes are divided into subclasses yielding a rather comprehensive taxonomy.

A survey by Gavrilu [44] describes work prior to 1998 and gives a good general introduction to the topic with a special focus on applications. The taxonomy covers 2D approaches with and without explicit shape models and 3D approaches. Across these three classes the approaches dealing with recognition are described.

3.2.2 A Taxonomy based on Functionalities

The above taxonomies have different emphasis depending on their purpose. We will focus on more general aspects such as the overall structure of a motion capture system and the various types of information being processed. The functional structure

of a comprehensive motion capture system is shown in Fig. 3.1.

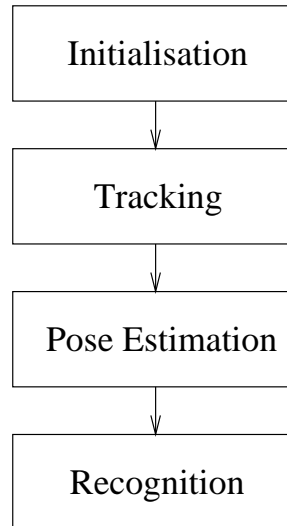


Figure 3.1: A general structure for systems analysing human body motion.

Before a system is ready to process data it needs to be *initialised*; e.g., an appropriate model of the subject must be established. Next the motion of the subject is *tracked*. This implies a way of segmenting the subject from the background and finding correspondences between segments in consecutive frames. The *pose* of the subject's body often needs to be *estimated* as this may be the output of the system, e.g., to control an avatar (the graphical representation of a human) in a virtual environment, or may be processed further by the recognition process. Some higher level knowledge, e.g., a human model, is typically used in *pose estimation*. The final process analyses the pose or other parameters in order to *recognise* the actions performed by the subject.

A system need not include all four processes, especially since many of the systems described in this survey are research, where only a method within one of the processes is investigated. Still, all systems can be described within the structure.

More than 130 human motion capture papers published since 1980 are reviewed for this survey. They are all listed in table 3.1. Within this table the 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 of our taxonomy. 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. The tables show, among other things, that the majority of the work in human motion capture is carried out within tracking and pose estimation. This is reflected in the rest of the paper where these two processes receive more attention than the other two processes.

In earlier surveys detailed summaries of individual papers are used extensively to

Year	First author	Initialisation	Tracking	Pose estimation	Recognition
1980	O'Rourke			[115]	
1983	Hogg		*	[58]	
1984	Akita		*	[3]	
1984	Hogg		*	[59]	
1985	Lee			[83]	
1985	Tsukiyama		[136]		
1987	Bernat		[11]		
1987	Leung		*	[88]	
1987	Leung		[89]		
1989	Attwood			[6]	
1991	Long		*	[92]	
1991	Shio		[130]	*	
1991	Wang		*	[141]	
1991	Yamamoto			[153]	
1992	Kepple			[81]	
1992	Lee			[84]	
1992	Luo		*	[94]	
1992	Wang		*	[142]	
1993	Kameda		*	[79]	
1994	Baumberg	*	[9]		
1994	Bharatkumar		[12]	*	*
1994	Darrell		[32]	*	*
1994	Gu		[48]		
1994	Guo		*	[50]	
1994	Niyogi		[108]	*	*
1994	Perales	*		[118]	
1994	Polana		*		[121]
1994	Rossi	*	[126]		
1994	Schneider		[127]		
1995	Cai		[20]		
1995	Campbell			*	[21]
1995	Campbell			*	[22]
1995	Freeman		[41]		*
1995	Goncalves		*	[47]	
1995	Kakadiaris	[76]		*	
1995	Kameda		*	[80]	
1995	Leung		*	[90]	
1995	Tesei		[135]	*	

Year	First author	Initialisation	Tracking	Pose estimation	Recognition
1996	Azarbayejani	*	[7]	*	
1996	Becker		*		[10]
1996	Bobick		[14]		*
1996	Cai		[19]	*	
1996	Gavrila		*	[45]	
1996	Ju		*	[73]	*
1996	Kahn		[74]		*
1996	Kakadiaris		*	[77]	
1996	Kameda		*	[78]	
1996	Luc				[93]
1996	Moezzi	[103]			
1996	Turk		[137]	*	*
1997	Bregler		*		[17]
1997	Christensen		*	[26]	
1997	Christensen		*	[27]	
1997	Davis		*		[33]
1997	Hunter		*	[62]	
1997	Iwasawa		[66]	*	
1997	Lerasle	*		[86]	
1997	Meyer	*	*	[97]	
1997	Oren		*		[114]
1997	Rohr	*	*	[124]	*
1997	Wachter		*	[139]	
1997	Wren	*	[143]	*	
1998	Bottino		*	[15]	
1998	Bregler	*	*	[18]	
1998	Chomat		*		[25]
1998	Chung		*	[28]	
1998	Corlin		[29]	*	
1998	Cretual		[30]		
1998	Davis		[34]		
1998	Davis		[35]		*
1998	Davis		[36]	*	*
1998	Fua		*	[42]	
1998	Fujiyoshi		[43]	*	*
1998	Goncalves			*	[46]
1998	Gu	[49]		*	
1998	Haritaoglu		[52]	*	*
1998	Haritaoglu		[53]	*	*

Year	First author	Initialisation	Tracking	Pose estimation	Recognition
1998	Heisele		[54]		*
1998	Isard		[64]		
1998	Jojić		*	[71]	
1998	Kakadiaris	*	*	[75]	
1998	Li		*	[91]	
1998	Munkelt		*	[104]	
1998	Nakazawa		[106]		
1998	Narayanan	[107]			
1998	Nordlund		[110]		
1998	Pinhanez		[119]		
1998	Silaghi			[132]	
1998	Sul		*		[133]
1998	Utsumi		[138]	*	
1998	Wren		*	[146]	
1998	Wren		[148]		*
1998	Yacoob		*	*	[150]
1998	Yamada		[151]		
1998	Yamamoto			*	[152]
1998	Yaniz		*	[155]	
1998	Zheng	*		[156]	
1999	Amat		[4]	*	*
1999	Andersen		[5]	*	
1999	Brand		*	[16]	
1999	Cham		[24]	*	
1999	Cutler		[31]	*	
1999	Delamarre		*	[37]	
1999	Douros	[40]			
1999	Haritaoglu		*		[51]
1999	Hilton	[55]	*	*	
1999	Hilton	[56]		*	
1999	Ioffe		*	[63]	*
1999	Iwai		*	[65]	
1999	Iwasawa	*	[67]	*	
1999	Lerasle	*		[87]	
1999	Njåstad	*	*	[109]	
1999	Ohya		[111]	*	
1999	Ong		*	[113]	
1999	Pavlović		*	[116]	*
1999	Plänkers		*	[120]	
1999	Rittscher		*		[123]
1999	Segawa		[129]	*	
1999	Wachter		*	[140]	
1999	Wren		*	[147]	

Year	First author	Initialisation	Tracking	Pose estimation	Recognition
2000	Hilton	[57]		*	
2000	Hu		*	[60]	
2000	Iwasawa	*	[68]		
2000	Jojic		[70]	*	*
2000	McKenna		[96]		*
2000	Moeslund	*	*	[100]	
2000	Moeslund		*	[101]	
2000	Okada	*	*	[112]	
2000	Rigoll	*	[122]		
2000	Rosales		*	[125]	
2000	Segawa		*	[128]	
2000	Sidenbladh	*	*	[131]	
2000	Wren		*	*	[144]
2000	Wren		*	[145]	
2000	Wu		*	[149]	
2000	Yamamoto		*	[154]	
Σ	Total=136	8	48	64	16

Table 3.1: The different publications' relation to the taxonomy. Note that this table is placed in an appendix in the original publication.

exemplify the individual classes in the taxonomies. Generally we will not do this. First because we focus more on overall methods and general characteristics and second to avoid drowning the essence in implementation details. Therefore, only the ideas of individual papers are described, when relevant. Furthermore, when presenting an idea, concept, or method we will generally refer only to one paper where a good description is given. Readers interested in the contents of the papers can, beside earlier surveys, refer to [98] where detailed summaries of many of the papers used for this survey are presented. Actually, [98] could be considered an appendix to this paper.

3.2.3 Assumptions

As for computer vision papers in other fields, various assumptions on the conditions for motion capture are associated with the individual contributions. The actual assumptions made characterise the various systems and provide a useful reference for evaluation.

The typical assumptions may be divided into two classes: *movement assumptions* and *appearance assumptions*. The former concerns restrictions on the movements of the subject and/or the camera(s) involved. The latter concerns aspects of the environment and the subject. In Table 3.2 the relevant assumptions and their association

with the two classes are listed.

Assumptions related to movements	Assumptions related to appearance
1. The subject remains inside the workspace	Environment
2. None or constant camera motion	1. Constant lighting
3. Only one person in the workspace at the time	2. Static background
4. The subject faces the camera at all time	3. Uniform background
5. Movements parallel to the camera-plane	4. Known camera parameters
6. No occlusion	5. Special hardware
7. Slow and continuous movements	
8. Only move one or a few limbs	Subject
9. The motion pattern of the subject is known	1. Known start pose
10. Subject moves on a flat ground plane	2. Known subject
	3. Markers placed on the subject
	4. Special coloured clothes
	5. Tight-fitting clothes

Table 3.2: The typical assumptions made by motion capture systems listed in ranked order according to frequency.

The first three assumptions related to movements are very general and used in every system with a few exceptions; see e.g., [36, 19, 106]. The fourth assumption is mainly used in human computer interaction (HCI) applications and simplifies the calculation of the overall body pose. The next assumption reduces the dimensionality of the problem from 3D to 2D and is often used in applications such as gait analysis. The sixth assumption concerns occlusion and simplifies the task of tracking the subject and limbs since the entire posture of the subject is visible in every frame. The next assumption goes both for the movement of the camera (if it is allowed to move) and the subject. No sudden movements are allowed and the movements follow a simple and continuous trajectory. This assumption simplifies the calculation of the velocity of the subject and of the camera. The eighth assumption allows tracking to focus on only one or a few body parts. The next assumption is used to simplify the tracking and pose estimation problems by reducing the solution space. The final assumption allows a calculation of the distance between the camera and the subject using the camera geometry and the size of the subject.

The first environmental assumption in practice constrains the scene to be indoors. The next assumption requires a static background which makes it possible to segment the subject based on motion information. The third environmental assumption constrains the background further to have a uniform colour and a simple thresholding may be used to segment the subject. The first two assumptions are used in many systems while the third assumption is used in approximately half of the systems. The fourth assumption concerns camera parameters which are necessary to know in order to obtain absolute measures in the registration. The last environmental

assumption concerns the use of special hardware such as multiple cameras or an IR-camera.

The first subject assumption about the known start pose is introduced in many systems to simplify the initialisation problem. The next assumption concerns prior knowledge of the subject, e.g., in terms of specific model parameters such as the subject's height, length, and width of limbs. The last three subject assumptions reduce the segmentation problems by making the subject's structure easier to detect.

The above assumptions are used to make the human motion capture problem tractable and they are applied in varying number and selection in all the reviewed papers. Which assumptions a particular system uses depends on its goals. Generally the complexity of a system is reflected in the number of assumptions introduced; i.e., the fewer assumptions, the higher the complexity.

3.3 Initialisation

As the first of the four major categories of our taxonomy we discuss initialisation. Initialisation covers the actions needed to ensure that a system commences its operation with a correct interpretation of the current scene. Sometimes the term initialisation is also used for preprocessing of data; see, e.g., Meyer *et al.* [97], or Rossi and Bozzoli [126]. We discuss preprocessing in Section 3.4 as a part of the tracking procedure. Some of the initialisation may be performed offline prior to the start of operation, while other parts preferably are included as the first phase of operation. Initialisation may be simplified by relying on some of the assumptions discussed above. Initialisation mainly concerns camera calibration, adaption to scene characteristics, and model initialisation.

As for other computer vision systems the parameters of the camera often need to be known. These can be obtained through offline camera calibration, and for a stationary camera setup occasionally recalibration will suffice. If something in the setup regularly changes, a procedure for online calibration may be preferred as in the work by Azarbayejani *et al.* [7]. However, virtually all other systems are based on offline calibration.

Initialisation to adapt the scene characteristics mainly relates to the appearance assumptions and the segmentation methods (described in Section 3.4) using them. In systems based on these assumptions a typical offline initialisation is carried out to find the thresholds and capture reference images which will be used during processing. In some systems initialised parameters are used in an adaptive procedure to calculate (and update) scene characteristics on the fly [53].

Model initialisation is concerned with two things: the initial pose of a subject and the model representing the subject. Both are closely related to model-based pose estimation described in Section 3.5.

Rohr [124] uses a model-based approach to estimate the pose of a subject and he describes the overall problem as first finding the initial pose of the human and then incrementing it from frame to frame. This structure represents the design of many model-based pose estimation systems. In the majority of these systems the overall problem is reduced either by assuming that the subject's initial pose is known as a special start pose [27] or by having the operator of the system specify it [153]. Perales and Torres [118] take this idea to the extreme by having the operator specify the pose in every single frame. Zheng and Suezaki [156] use a similar approach but they only manually fit the pose at some key frames and have the system interpolate, using correlation, between frames.

Only a few systems actually have a special initialisation phase where the start pose is found automatically [124]. In some systems the same algorithm is used for initialisation and during tracking/pose estimation [109]. This indicates that no temporal information is used and nothing is learned by the system during processing. These systems are usually not considered initialisation since they do not address the initialisation problem in a general sense, but rather in a special situation constrained by the assumption of a known motion pattern.

The result of a model-based approach is usually dependent on how well the subject fits the human model in the system, i.e., the complexity of the model. Some systems use a general model which is an average of many individuals [9]. Others measure the current subject and generate a model based on these data. This may be done offline [49, 86] or online as by Wren *et al.* [143]. In the latter, analysis of the subject's initial pose is carried out to build an initial model which is refined as more information becomes available to the system. Generating a personalised model mainly relies on building up a 3D shape of the subject through multiple cameras, e.g., stereo reconstruction [103, 107], and then mapping texture onto the shape model. The personalised model may also be obtained by fitting a generic model to current data [57, 156].

In computer graphics the concept of using real images of humans to animate personalised human models is becoming more popular. A clear tendency towards the merging of computer vision and computer graphics is apparent [85] and personalised models are being incorporated into computer vision-based motion capture systems [55] to improve performance.

3.4 Tracking

Tracking is a well-established research field which may be addressed from various viewpoints. In this context we define tracking as establishing coherent relations of the subject and/or limbs between frames. What is needed to achieve this depends on the context. Tracking may be seen as a separate process, as a means to prepare data for pose estimation, or as a means to prepare data for recognition.

If considered a separate process the subject is typically tracked as a single object (without any limbs) and no high-level knowledge is used. An example is the work by Tsukiyama and Shirai [136]. They detect moving people in a hallway by first detecting moving objects, then finding the objects corresponding to humans, and finally tracking the moving humans over time.

If the tracking process prepares data for pose estimation its purpose is to extract specific image information, either low level, such as edges, or high level, such as hands and head. An example of low-level tracking is seen in the Walker system by Hogg [59]. Here edges are extracted from the image and matched against the edges of a human model to determine the pose of the model/human. An example of high-level tracking is seen in the Pfinder system by Wren *et al.* [143] where the human body is tracked in 2D using statistical models of the background and foreground to segment the image into blobs. Some of these blobs represent the hands and feet of the subject yielding some sort of limb representation.

If tracking prepares data for recognition, the task is usually to represent data in an appropriate manner. An example of this was published by Polana and Nelson [121] where flow information and down-sampling are used to represent image information in a compact manner which is processed by a classifier to recognise six different classes, e.g., walking and running.

Independent of the context of tracking, three common aspects can be identified. First, nearly every tracking algorithm within human motion capture starts with the figure-ground segmentation problem, i.e., segmenting the human figure from the rest of the image. Second, these segmented images are transformed into another representation to reduce the amount of information or to suit a particular algorithm. Third, how the subject should be tracked from frame to frame is defined.

3.4.1 Figure-Ground Segmentation

Figure-ground segmentation may be based on either temporal or spatial information. Table 3.3 shows how the data may be characterised in further detail.

Temporal data		Spatial data	
Subtraction	Flow	Threshold	Statistics
Two images	Points	Chroma-keying	Pixels
Three images	Features	Special clothes	Blobs
Background	Blobs	IR	Contours

Table 3.3: The various characteristics and subclasses of figure-ground segmentation approaches.

Temporal Data

The use of temporal data is mostly based on the assumption of a static background (and camera). In this case the differences between images from a sequence must originate from the movements of the subject. Two subclasses may be introduced: subtraction and flow.

Subtraction is widely used by simply subtracting the current image from the previous image in a pixel by pixel fashion, using either the intensity values [121] or the gradients [4]. An improved version is to use three consecutive images instead of two [78]. The result reflects movements (and noise) between the images unless the subject has the same intensity or colour as the background. The use of background subtraction is very popular. If the scene is static, a noise-free background image without any subject(s) may be recorded and used as a reference in a subtraction scheme [106]. A more advanced version is to update the background image during processing [53].

Flow is here used as a general term describing coherent motion of points or features between image-frames. An example of such a flow using points is described by Yamamoto and Koshikawa [153]. They find the motion parameters of a human body part by calculating the optical flow of several points within this part and compare them to the movements of a model of the part. Gu *et al.* [48] instead track edge features, defined by their length and contrast, in consecutive images using the optical flow constraint. In the work by Bregler [17] each pixel is represented by its optical flow which is grouped into blobs having coherent motion and represented by a mixture of multivariate Gaussians.

Figure-ground segmentation based on temporal data assumes the subject to be the only moving object in the scene (or rather in the region of interest). In many cases temporal data are a strong alternative to spatial data. Temporal data are usually simpler to extract and focus directly on the target of motion capture. A number of systems are based on mapping the motion data directly to human pose through an inverse kinematic framework, see e.g. [18].

Spatial Data

The use of spatial data falls into two distinct subclasses: *thresholding* and *statistical approaches*. The former is simple processing based on special environmental assumptions. The latter is a rather advanced class where some of the appearance assumptions exploited by the *subtraction* methods are relaxed.

If the subject's colour or intensity appears different from the rest of the scene it may be segmented using simple thresholding. A good example is Chroma-keying where a person appears in front of a one-colour, usually blue, screen wearing nonblue clothes. By thresholding, the person can easily be separated from the background [32, 65]. The opposite approach, where the subject wears one colour, usually dark, clothes

against a different background is also very popular [12]. A special version of this idea is to make the subject wear markers (passive or active) which are easily segmented by thresholding [21, 46]. A related approach is to use an IR camera. Thermal images can be obtained where the subject is easy to segment through thresholding as the only "hot" object in the scene [66].

The *statistical* approaches use the characteristics of individual pixels or groups of pixels to extract the figure from the background. The characteristics are typically colours and edges. Some approaches are inspired by the background subtraction methods described above. A sequence of background images of the scene is recorded and the mean and variance intensity or colour of each pixel are calculated over time. In the current image each pixel is compared to the statistics of the background image and classified as belonging to the background or not [151]. This approach is becoming increasingly popular due to its robustness compared to the subtraction approaches. In the work by McKenna *et al.* [96] the approach is combined with the statistics of pixel gradients to remove the shadows cast by subjects. A more advanced version is used in the blob approaches, where the subject is modelled by a number of blobs with individual colour and spatial statistics. Each pixel in the current image is then classified as belonging to one of the blobs according to its colour and spatial properties [143].

Another statistical approach is to use static or dynamic contours, where dynamic is usually named active contours, and static refers to the use of predefined static structures representing a part of the subject's outline. These structures consist of edge segments and other attributes. Long and Yang [92] uses *Logs*, which is an area defined by two parallel edge segments and a number of attributes. These are found in the image and used to extract human body parts by comparing them with a human model consisting of Logs. The active contours are used in a similar manner except they have the ability to adjust their shape on the fly, hence active. Their deformation is controlled by external and internal energy functions. The former fits the curve to the image features, e.g., edges, while the latter adjusts the smoothness of the curve. This type of active contours is also called *snakes* and works relatively well when the changes in the structure of the object are unknown. If a shape model is used the active contour is known as a deformable template [13]. Active contours may be used to extract the entire outline of the subject [9] or to extract individual body parts [75]. A related approach is to extract the silhouette instead of the contour. Rigoll *et al.* [122] use a stochastic approach to silhouette extraction. They use pseudo-2D Hidden Markov models (HMM) (nested one-dimensional HMMs), to extract the silhouette in a discrete cosine transformed (DCT) representation of the image.

The use of thresholding to process spatial data is strongly dependent on a number of appearance assumptions. This might seem unreasonable if the goal is an assumption-free vision system with human capabilities. However, there will always be a range of applications where the environment and the subject can be controlled,

as underlined by the fact that one of the most widely used and robust figure-ground segmentation method is Chroma-keying. In more unconstrained applications, the statistical methods are a far better choice due to their adaptability. The use of pixel statistics is a good concept, but region-based methods, using e.g., blobs, tend to be more reliable. On the other hand it is more difficult to model larger entities of correspondingly higher complexity. The active contours aim directly at extracting the shape of the subject and can be very efficient. However, they require a good initial fit and have difficulties with complex articulated objects such as the human body. The best way to exploit them seems to require an active contour for each body part.

3.4.2 Representation

Segmented entities are described compactly by some convenient representation. There are in principle two types of representations: the *object-based representation* which is based on the figure-ground segmentation and the *image-based representation* which is derived directly from the image. In Table 3.4 the various types of representations are shown.

Object-based	Image-based
Point	Spatial
Box	Spatio-temporal
Silhouette	Edge
Blob	Features

Table 3.4: The various types of data representations.

Object-Based Representation

The object-based representations rely mainly on the output from the figure-ground segmentation. Therefore some of the arguments and descriptions from Section 3.4.1 also apply here.

The point representation is sufficient in systems using passive or active markers. The active markers yield a high contrast in the images and provide a robust representation [132]. If more than one camera is used in a marker-based system a 3D representation may be obtained [104].

The box representation is used in many systems. The philosophy is to represent the subject by a set of boundary boxes containing the pixels or regions found in the figure-ground segmentation process. Some systems track these boxes over time [106], while others use them as intermediate representations [143], which are processed (rerepresented) further in *pose estimation*.

The silhouette representation is popular due to its simplicity. It can be obtained using the thresholding or subtraction methods from figure-ground segmentation. It is used in both 2D and 3D. The 2D representation is usually straightforward [91], but can also be more complex as in the work by Baumberg and Hogg [9], where the silhouette (or rather active contour) is represented even further using closed uniform B-splines with a fixed number of control points equally spaced around the silhouette. This actually makes it a contour representation but equivalent to a silhouette. The 3D (volumetric) silhouette can be obtained using combined 2D silhouettes [15] or directly using stereo approaches [71]. The silhouette representation may, as was the case for the box representation, be tracked directly [53] or more likely processed (rerepresented) further in pose estimation [52].

The blob representation typically follows some of the figure-ground segmentation approaches described under flow and statistics. The subject is represented as a blob or a number of blobs each having some similar characteristics. The similarities can be coherent flow [73], similar colours [54], or both [17]. The main philosophy of grouping information according to similarities is inspired by research into the human visual system by the Gestalt school in the 1930s [82].

Image-Based Representation

This class of representations is based directly on the pixels of the image. The representations are either derived from the image independent of the presence of an object or possibly constrained to the interior of one of the representations described above. These images (or image parts) may be transformed into another space spanned by non-Cartesian basis functions, yielding a more compact representation of the data or image. Transformations used are, e.g., Fourier, principle component analysis (PCA), DCT, and Wavelets.

A more advanced representation is obtained when including the temporal dimension in the representation [25]. This allows for motion-related features to be included in the description.

A third subclass of image-based representations is edges. They may be represented by points [45] or as line segments which are more robust to noise [124].

The last form of representation is features. Features are usually computed from one of the previously mentioned types of representation combined with additional information. Christensen and Corneliussen [27] use the length, area, and colour to represent the individual body parts found by thresholding an image of a subject wearing special coloured clothes.

3.4.3 Tracking over Time

Tracking over time is finding corresponding objects in consecutive frames where the objects may be any of the representations from Table 3.4. The difficulties of this task are related to the complexity of the scene and the complexity of the tracked objects. The latter is again related both to the degrees of freedom of individual objects and to their representation. Tracking more points in an object is equivalent to tracking multiple objects simultaneously. The points or objects (hereafter objects) may split and merge into new objects due to occlusion or image noise, or the appearance of an object may change due to shadows and changes in the lighting.

The correspondence analysis is often supported by prediction. Based on previously detected objects and possibly high level knowledge the state of the objects (appearance, position, etc.) in the next frame is predicted and compared (using some metric) with the states of objects found in the actual image. Prediction introduces a region-of-interest in both image space and state space and hereby reduces the overall need for processing. The prediction of the various state parameters is based on a model of how they evolve over time. A model of velocity and acceleration [101] or more advanced models of movements such as walking [124] may be used. An alternative approach is to learn probabilistic motion models prior to operation [116]. A commonly used method for prediction is the Kalman filter, which is also capable of estimating the uncertainties of the prediction. These uncertainties may be used to determine the regions-of-interest.

The Kalman filter is unfortunately restricted to situations where the probability distribution of the state-parameters is unimodal. In the presence of occlusion, cluttered background resembling the tracked objects, and complex dynamics, the distribution is likely to be multimodal. Alternative tracking algorithms have therefore been developed capable of tracking multiple hypotheses, i.e., support multimodal distributions. Most recognised is perhaps the Condensation algorithm [64]. It is based on sampling the posterior distribution estimated in the previous frame and propagating these samples to form the posterior for the current frame. The method has shown to be a powerful alternative to the Kalman filter [113, 131]. However, since it is nonparametric it requires a relatively large number of samples to ensure a fair maximum likelihood estimate of the current state. In high-dimensional problems a more efficient method might be necessary. In the work by Cham and Rehg [24] only the peaks of the posterior distribution are sampled and propagated to the next frame, resulting in relatively few samples. Furthermore, a distribution is formed by piecewise Gaussians where the means are given by the propagated (predicted) samples and the covariances by the uncertainties of the predictions. Hence the distribution is parametric and allows for a direct maximum likelihood estimate.

Another tracking aspect arises when multiple cameras are used. Some systems have access to more cameras than required and therefore need to decide which camera(s) or image(s) to use at each time instant. Various measurements are introduced to

quantify the individual viewpoints in terms of ambiguity [113], occlusion [75], and in general reliability of data [67, 138].

3.5 Pose Estimation

Pose estimation is the process of identifying how a human body and/or individual limbs are configured in a given scene. Pose estimation can be a postprocessing step in a tracking algorithm or it can be an active part of the tracking process. Various levels of accuracy may be required in pose estimation. At one extreme coarse estimations are carried out yielding only information about the subject's hands and head (or even simpler, the body's centre of mass) which could be used in a surveillance system or in a HCI system. At the other extreme the precise pose in terms of positions, orientation, width, etc. of each limb is estimated. The latter allows for, e.g., direct copy-cat interfaces to virtual environments or as input data to medical studies. Due to the complexity involved in this type of pose estimation, only one subject or a few body parts are considered.

A common aspect of pose estimation is the use of a human model. Usually a geometric model of the human body is applied, but other models, such as motion models, may also be applied. The model may be used for various purposes and at various levels. The general concept of using a human model is to exploit the fact that the system is dedicated to analyse a human and therefore may incorporate knowledge about humans into its processing.

Due to the extensive use of human models it is a good cue for a taxonomy. In former surveys the dimensionality of the applied (geometric) model is used to distinguish between various model-based approaches. This provides a good overview of 2D and 3D approaches, as seen in the surveys by Aggarwal and Cai [1], and Gavrilu [44], respectively. However, in some cases it is not clear whether a system is a 2D or a 3D approach, e.g., when 2D data are combined into 3D data using triangulation. Instead, we suggest investigating how various systems apply a human model. Pose estimation is therefore divided into three classes. The first class, *model-free* covers approaches where no *a priori* model is used. The other two classes of pose estimation methods, *indirect model use* and *direct model use*, are characterised by having a human model beforehand. In the indirect case the model is used as a reference to constrain and guide the interpretation of measured data. In the direct case the model is maintained and updated by the data, and hence at any time it includes an estimate of the pose.

3.5.1 Model-Free

In this class of methods no *a priori* model is used. The methods do, however, build some sort of model to represent the pose. These pose representations are points,

simple shapes, and stick-figures. The two first are similar to the object-based representations described in the previous section, while the stick-figure representation is more advanced.

The pose of the subject may be represented by a set of points and is widely used when markers are attached to the subject. Without the markers, the hands and the head may be estimated from the image represented by just three points [143]. This is a very compact representation which suffices for a variety of applications. These three points are usually found using colour segmentation [101] or blob segmentations [143].

A subject may be represented by simple boundary boxes [32]. The box representation is, however, mainly used as an intermediate representation during processing and not as the final representation. Instead, shapes which are more human-like, such as ellipses [106], may be used as a final representation.

The stick-figure representation includes structure information resembling the human skeleton. It is a popular representation in systems where the gait of a subject is studied. The stick-figure is obtained in various ways, e.g., directly using a medial axis transformation [12] or a distance transformation [66]. A distance transformation is slightly more advanced allowing for suppression of noninteresting parts, e.g., arms when searching for the torso. Both methods give a direct but noisy estimate of the human skeleton.

A rather different approach to pose estimation without the use of an *a priori* model is to learn a mapping directly from image features to pose data. These type of systems rely on extensive training using ground truth data obtained with commercial motion capture systems. Below two examples are given.

Rosales and Sclaroff [125] use 3D ground truth data of a subject's joints to train their system. For a number of viewpoints the 3D data are mapped into 2D joint positions. For the same viewpoints a cylinder model is used to synthesise a silhouette and its Hu-moments are calculated. The 2D joint positions are for all training sets clustered into a number of Gaussians using the expectation maximization (EM) algorithm [39]. For each cluster a neural network is trained which maps Hu-moments into 2D joint positions. During processing the Hu-moments are calculated from an image silhouette and feed to each neural network, resulting in a pose candidate for each cluster. The candidates are each synthesised and their Hu-moments are compared to those found from the image in order to evaluate the different hypotheses. This comparison is similar to the use of silhouettes described in Section 3.5.3.

In the work by Brand [16], the paths through the 3D state space obtained through training are modelled using an HMM where the states are linear paths modelled by multivariate Gaussians. As above the moments (central) are found by synthesising various poses. The moments are associated to the HMM; i.e., a new HMM is obtained with moments as input and pose parameters as hidden states. Altogether a sequence of moments is mapped to the most likely sequence of 3D poses. Obviously

this approach is either offline or includes a significant time-lag, but it has the ability to resolve ambiguities using hindsight and foresight.

3.5.2 Indirect Model Use

The methods in this class use an *a priori* model when estimating the pose of a subject. They use the model as a reference or look-up table from which relevant information may be extracted to guide the interpretation of measured data.

Various types of models and various levels of detail are used. The level of details ranges from just the height of the subject to both structure and dynamic information about subjects. The estimated pose of systems in this class is generally not very detailed. Typically positions of the head, hands, and feet or a rough description of the entire human body is used to represent the pose.

A simple human model is the aspect ratios between the various limbs [20, 53] which may be used to guide the pose estimation. In the work by Leung and Yang [90] the outline of the subject is estimated as edge regions described using 2D ribbons which are U-shaped edge segments. A 2D ribbon model guides the labelling of the image data by searching for structures similar to the model. To select among the various labelled image parts they apply more high-level knowledge such as expected motion. Njåstad *et al.* [109] also assume a known motion pattern which is used to define where to search for the individual body parts with respect to the torso. A related concept is to have key frames beforehand as used by Attwood *et al.* [6] and Akita [3], who exploits them to predict occlusion in the next frame, which again guides the processing. Haritaoglu *et al.* [52] went a step further by first calculating which overall pose is present (standing, crawl-bend, laying down, and sitting) and then using this information to find the individual parts. Wren and Pentland [146] generalised the concept of knowing the motion beforehand by suggesting that various behaviour models may be used to improve the pose estimation. An example of this was also published by Iwai *et al.* [65] where six different motion models are used.

On the borderline to direct model use O'Rourke and Badler [115] describe a rather detailed human model. The model is used to ensure that the predicted human pose is realistic. The same approach is used in Hunter *et al.* [62] where the estimated pose is projected onto a feasible manifold in the solution space, ensuring a realistic pose. Ioffe and Forsyth [63] use the probability of various model configurations to guide the estimation of the pose of one or more people. In the work by Wren and Pentland [145] model information is used to estimate the pose of the elbows after the head and hands have been located.

3.5.3 Direct Model Use

By direct model use we mean that an *a priori* human model is used as the model representing the observed subject. It is then continuously updated by the observations. Hence, this model provides any desired information including pose at any time. Approximately 40% of the surveyed papers use a model in this manner and they are all listed in Table 3.5.

The models used in this section are generally very detailed. They explicitly exist within the computer program and are used intensively in the processing phase. An important benefit gained by introducing a human model is the ability to handle occlusion and the ease by which various kinematic constraints may be incorporated into a system.

A human model is represented by a number of joints and the sticks (bones) connecting them. The sticks and the "flesh" surrounding them may be represented in various ways depending on the level of detail needed by a system. In Table 3.5 the type of model used by each system, the number of segments in the model, and the parts of the subject being estimated, are listed in columns 3, 4, and 5, respectively. The more complex the model, the better results may be expected but on the expense of more processing and training.

The concrete representation of the human model is a state space where each axis represents a degree of freedom of a joint in the model. One pose of the subject may be expressed as one point in the state space as opposed to many points in the 2D image space. The problem is how to use this state space representation and, hence, how to relate image data to pose data. The general approach is known as *analysis-by-synthesis* and used in a *predict-match-update* fashion. The philosophy is to predict the pose of the model corresponding to the next image. The predicted model is then synthesised to a certain abstraction level for comparison with the image data.

When comparing the real and synthesised data a similarity measure is used to evaluate how alike the image and the synthesised model are. Typically this is done for a number of predicted model poses until the correct (best) pose is found and used to update the model in the system.

Clearly the state space, even with a coarse resolution, describes a very large number of possible model poses unreasonable to synthesis for matching. Therefore constraints are introduced to prune the state space. An obvious and substantial reduction of the range of each parameter and its derivatives in the state space is achieved by introducing the kinematic constraints of the human motor system, e.g., the bending of the elbow is between 0 and 145°. This may be used directly to partition the state space into legal and illegal regions, as in [100], or the constraints may be defined as forces acting on an unconstrained state phase, as in [144]. The fact that two human body parts cannot pass through each other also introduces constraints. Another approach to reduce the number of possible model poses is to assume a known

Year	First author	Model type	Parts	Object	Ab. level	Dim.
1983	Hogg [58]	Cylinders	14	Body	Edges	$2\frac{1}{2}$
1984	Hogg [59]	Cylinders	14	Body	Edges	$2\frac{1}{2}$
1985	Lee [83]	Stick-Figure	17	Body	Joints	3
1989	Attwood [6]	Cylinders	16	Body	Joints	3
1991	Wang [141]	Cylinders	2	Leg	Motion	3
1991	Yamamoto [153]	CAD Model	4	U Body	Motion	$2\frac{1}{2}$
1992	Lee [84]	Stick-Figure	17	Body	Joints	3
1992	Luo [94]	Stick-Figure	6	Body	Silhouettes	3
1992	Wang [142]	Cylinders	2	Leg	Motion	3
1993	Kameda [79]	Patches	17	Body	Silhouettes	$2\frac{1}{2}$
1994	Guo [50]	Stick-Figure	10	Body	Sticks	2
1995	Goncalves [47]	R-circular cones	2	Arm	Edges	3
1995	Kameda [80]	Patches	15	Body	Silhouettes	$2\frac{1}{2}$
1996	Gavrila [45]	Super-Quadrics	10	Body	Edges	3
1996	Ju [73]	Planar Patches	2	Leg	Motion	2
1996	Kakadiaris [77]	D Silhouettes	2	Arm	Silhouettes	3
1996	Kameda [78]	Patches	9	U Body	Silhouettes	$2\frac{1}{2}$
1997	Christensen [27]	Cylinders	10	Body	Sticks	2
1997	Lerasle [86]	CAD Model	2	Leg	Texture	3
1997	Meyer [97]	Boxes	6	Body	Contours	$2\frac{1}{2}$
1997	Rohr [124]	Cylinders	14	Body	Edges	$2\frac{1}{2}$
1997	Wachter [139]	R-Elliptical Cones	10	Body	Edges	3
1998	Bregler [18]	Ellipsoids	10	Body	Motion	3
1998	Fua [42]	Ellipsoids	2	Arm	Silhouettes	3
1998	Jojić [71]	D Super-Quadrics	(5)6	(U) Body	Contours	3
1998	Kakadiaris [75]	D Silhouettes	4	Arms	Silhouettes	3
1998	Li [91]	Rectangles	14	Body	Silhouettes	2
1998	Munkelt [104]	Modified Cylinders	10	Body	Joints	3
1998	Silaghi [132]	Stick-Figure	16	Body	Sticks	3
1998	Wren [146]	Stick-Figure	5	U Body	Blobs	3
1998	Yamamoto [152]	CAD Model	9	Body	Motion	3
1998	Yaniz [155]	Stick-Figure	16	Body	Sticks	3
1999	Cham [24]	Scaled Primitives	10	Body	Texture	2
1999	Delamarre [37]	Cones and Spheres	15	Body	Silhouettes	3
1999	Iwai [65]	Stick-Figure	6	U Body	Silhouettes	3
1999	Lerasle [87]	CAD Model	2	Leg	Texture	3
1999	Njåstad [109]	Cylinders	10	Body	Contours	$2\frac{1}{2}$
1999	Ong [113]	Stick-Figure	4	Arms	Contour	3
1999	Pavlović [116]	Scaled Primitives	6	Body	Texture	2
1999	Plänkers [120]	Ellipsoids	6	Arm	Depth	3
1999	Wachter [140]	R-Elliptical Cones	10	Body	Edges	3
1999	Wren [147]	Stick-Figure	5	U Body	Blobs	3

Year	First author	Model type	Parts	Object	Ab. level	Dim.
2000	Hu [60]	Rectangles	10	Body	Silhouette	2
2000	Moeslund [100]	Cylinders	2	Arm	Silhouette	3
2000	Moeslund [101]	Cylinders	2	Arm	Silhouette	3
2000	Okada [112]	Boxes & Cylinders	4	U Body	Motion/Depth	3
2000	Rosales [125]	Cylinders	10	Body	Silhouette	2
2000	Sidenbladh [131]	Cylinders	10	Body	Texture	3
2000	Wren [145]	Stick-Figure	5	U Body	Blobs	3
2000	Yamamoto [154]	CAD Model	11	Body	Motion	3

Table 3.5: The papers where a human model is used directly. The first and second column states the publication year and first author. The third shows how the human is modelled. The fourth shows the number of parts in the model. The fifth shows which part of the subject is analysed. The sixth shows which abstraction level the system works at and the last shows the dimensionality of the pose estimated data. D, deformable; R, right; U, upper.

motion pattern - especially cyclic motion such as walking and running. In the work by Rohr [124] gait parallel to the image plane is considered. Using a cyclic motion model of gait all pose parameters are estimated by just one parameter, which specifies the current phase of the cycle. This is the most efficient pruning encountered in the surveyed literature. Ong and Gong [113] map training data into the state space and use PCA to find a linear subspace where the training data can be compactly represented without losing too much information. Pavlović *et al.* [116] take this idea a step farther by learning the possible or rather likely trajectories in the state space from training data; i.e., dynamic models are learned. Yet another approach is to rerepresent the state space more efficiently (without losing any information as in PCA). In the work by Moeslund and Granum [100] the state space representation of an arm is given in just two parameters instead of the usual four (two for the shoulder and two for the elbow) using geometric reasoning.

The systems also differ in the way they compare the synthetic data and the image data. Even after the constraints have been applied a brute force comparison is seldom realistic. Therefore, the matching problem is formulated as a function which is optimised. Due to the high dimensionality of the problem an analytic solution is not possible and a numeric iterative approach is generally used [140]. An alternative approach is to predict just one state and let the difference between the synthesised data and the measurements be used as an error signal to correct the state of the model [47, 147]. An excellent framework for this type of approach is the Kalman filter. Yet another approach is to predict the most likely states using a multiple hypothesis framework (see Section 3.4.3) [24].

If proper constraints are introduced, the state space is reduced significantly. Com-

binning this with an effective match scheme the analysis-by-synthesis approach can be very successful. The papers in Table 3.5 all use the analysis-by-synthesis approach and they have produced some of the best results within human pose estimation. To give a better understanding of how the analysis-by-synthesis approach is used the various abstraction levels which may be used in a system are described in the following.

Abstraction Levels

The various abstraction levels used for comparing image data and synthesised data are edges, silhouettes, contours, sticks, joints, blobs, depth, texture, and motion. In Table 3.5 (column 6) each system and its level of abstraction can be seen. Below the various levels are exemplified using a number of references.

The *edges* of the model and the subject in the image are relatively easy to find, especially when the appearance-related assumptions are used. Edges are therefore a good representation for the matching process. One of the first analysis-by-synthesis publications was by Hogg [58] in 1983. He used image subtraction to obtain a boundary box of a human. In this box he found edges and compared them with edges projected from a human model. He successfully showed the effect of the analysis-by-synthesis approach. Rohr [124] used the same idea, but in a more sophisticated system where he used Kalman filter, edge segments, and a motion model tuned to walking to obtain a more robust result. Gavrilu and Davis [45] worked with the same problem but without assuming a known motion model. They instead utilised four camera views and tight-fitting coloured clothes to obtain good results. They compare image and model edges using a robust variant of chamfer matching. In the work by Wachter and Nagel [140] both edges and regional information is applied in the matching. They also work with unconstrained motion patterns, but only using monocular vision. Another interesting issue is that they only use edge segments which are expected to be visible in the image.

A *silhouette* and its *contours* are, like edges, relatively easy to extract from both the model and the image. A silhouette has the advantage over edges of being less sensitive to noise, since it is a region-based data type. On the other hand, fine details may be lost in the extraction of the silhouette. In the work by Kameda *et al.* [79] the silhouette extracted from the image is matched against the silhouette of the model and the similarity measure is simply given as the area difference. They use a local match strategy, i.e., one limb at the time. Hu *et al.* [60] also use a local match strategy but they consider both positive and negative matching errors and introduce, weighting of the two in their similarity measure. Furthermore, they apply multiscale morphologic matching to obtain better results. The idea is that the main structural information is reserved in large morphologic scale. So both the input and model silhouettes are dilated in large scale, removing a significant amount of noise and thereby improving the result of the matching.

If the contour, rather than the silhouette, is used the fitting between the image and model data is usually done via active contours, i.e., the forces which should be applied to the model (a deformable rather than a geometric model) to make it fit the image contours are calculated [75, 71]. In the work by Meyer *et al.* [97] a geometric model is used. The parts of the subject are found from optical flow and represented by their contours. These are obtained by active rays which are similar to active contours, except the problem is reduced from 2D to 1D. The contours are then compared to predicted model contours. In some systems a stick-figure model is matched against an image silhouette. This may, for example, be seen in the work by Luo *et al.* [94] where the silhouette of the subject from two different views are compared with a synthesised stick-figure model using kinematic match criteria.

The human is represented by its *joints* or *stick-figure* since it reflects anatomic features of the human. Both may be easily obtained from the model since it is basically defined in these terms. However, it may be hard to obtain them directly from the image and usually various assumptions are introduced to simplify matters. In the work by Lee and Chen [83] the positions of the joints in the image and the 3D length of each segment are known beforehand. Given the 3D position of the neck a partial tree is build. At each node one of two solutions for the next joint is possible, since the joint's projection in the image and the 3D length are known. A path through the tree is equal to one body pose. The tree is pruned using kinematic constraints and the assumption that the subject is walking. In [84] they improve their system by introducing a smooth motion constraint. In the work by Attwood [6] a similar approach is taken, except that he uses three static postures (standing, kneeling, and sitting) instead of a walking assumption to prune the solution space. In Munkelt *et al.* [104] the 3D positions of the joints are estimated using markers and stereo. These are compared with 3D model joints using a graph-based scheme to find the correct pose of the subject.

A stick-figure representation is closely related to a joint representation and also reflects anatomic features which make it tractable. In the work by Guo *et al.* [50] a model stick-figure is compared to the image-skeleton found from the silhouette. To reduce the complexity of the matching, a potential field is introduced. It transforms the problem into one of finding a stick-figure with minimal energy in a potential field. Prediction and angle constraints of the individual joints are introduced to simplify the matching process further.

In the work by Wren and Pentland [145, 146, 147] *blobs* are used as the abstraction level. The Pfinder algorithm [143] (see Section 3.4) is run on two images from two different cameras producing 3D blobs for the hands and the head. Together with a dynamic model, kinematic constraints, and a Kalman filter the blobs provide enough information to estimate the pose parameters of the upper body. Another interesting aspect of this work is, as mentioned earlier, the inclusion of behaviour models into the control loop. These models resemble the effect of an active complex controller (nervous system) as the kinematic and dynamic models resemble the motor system.

Due to the difficulties in explicitly modelling the nervous system they apply learned probabilistic models, each corresponding to a prototypical behaviour. They choose a model according to current observations and thereby allow multiple hypotheses as in the Condensation algorithm.

In the work by Plänklers *et al.* [120] *depth data* of a subject's right arm is estimated using three cameras. The subject wears tight fitting textured clothes to simplify the correspondence problem. The depth data are compared to model ellipsoids in an optimisation framework.

In the work by Lerasle *et al.* [86, 87] a *texture* representation is used; i.e., the image data is used without much preprocessing. Offline they generate a model of a subject's leg, which is texture mapped using the real image data where the subject wears a pair of textured tights. During processing they compare the texture in the image to the texture on the model using correlation. This is done for several cameras, and by merging the results the 3D motion of the leg is estimated. In the work by Sidenbladh *et al.* [131] a texture model of each limb is generated offline. They use a commercial motion capture system to obtain ground truth 3D pose data. Based on these data and a camera model they derive a mapping between an image pixel and its position on a cylinder - modelling a limb. In this way they automatically build a weighted texture model of each cylinder/limb. Using PCA the models are compactly represented in a linear subspace. During processing they synthesise various poses, i.e., texture maps, and compare them to the image data to evaluate the similarity. The similarities of the synthesised poses constitute the posterior distribution in the Condensation algorithm where either a smooth motion model or a gait model is used to propagate the distribution over time.

The previous abstraction levels have been used in a spatial context where the structure of the subject is synthesised and compared with image measurements in each frame. Obviously the temporal context is also considered, but it is mainly to constrain the search space and thereby making the approach suitable for real-time implementation. The systems described later all use the temporal context more intensively by estimating the motion between consecutive frames.

Instead of first finding the pose of the subject in several individual frames and then using this information to calculate the motion, one may measure the *motion* of the subject between images and set up an inverse kinematic framework which makes it possible to calculate the corresponding motion in the 3D model. That is, the model and therefore the pose is updated based on the motion in the images. This was first done by Yamamoto and Koshikawa [153] in 1991. They measured optical flow within various body parts and used that through a Jacobian matrix to update the model. Ju *et al.* [73] use two planar patches to model a leg - a cardboard model. The motion of each patch is defined by eight parameters. For each frame the eight parameters are estimated by applying the optical flow constraint on all pixels in the predicted patches. The distance between the corners of the predicted patches are constrained to reduce the complexity of the estimation. Bregler and Malik [18]

extended the concept by introducing a twist motion model and exponential maps which simplify the relation between image motion and model motion. Their novel formulation also has the advantage of being open to both single and multiple views. In addition to lab video sequences they also tested their system on a number of the famous Muybridge images recorded in 1884.

Dimensionality of Pose Estimation Approaches

The distinction between 2D and 3D systems concerns the dimensionality of the pose estimation approaches based on a direct use of a human model. Classification of the relevant systems is provided in the last column in Table 3.5 where the dimensionality is indicated as 3D, 2D, and $2\frac{1}{2}$ D according to the definitions below.

3D refers to estimation of 3D movements. 2D refers to estimating either motion carried out in 2D, e.g., swinging an arm fronto-parallel to the camera, or to the motion observed directly in the image, which per definition is 2D.

$2\frac{1}{2}$ D refers to either estimating 3D pose data based on 2D processing or testing a 3D pose estimating framework on pseudo 3D data. The former occurs, for example, when a subject is walking fronto-parallel to the camera. Estimating the pose of the arm and leg closest to the camera is a 2D problem but due to symmetry and *a priori* knowledge the pose of the opposite arm and leg may be estimated producing 3D pose data, or rather $2\frac{1}{2}$ D. The latter refers to situations where a 3D model is used but for some reason the test data is only $2\frac{1}{2}$ D, e.g., when the subject is walking fronto-parallel to the camera. This makes it hard to judge the success of a 3D pose estimating system and therefore these systems are also referred to as being $2\frac{1}{2}$ D.

Evaluating Pose Estimation

Another aspect arising when discussing systems estimating 3D (and pseudo-3D) poses is how to evaluate their results. In the case of 2D pose data a straight forward comparison with the image data may be used in most cases, and pose estimation tailored to a recognition task should of course be tested in this context; see Section 3.6.

The problem of evaluating the estimated poses is that usually no ground truth is available. Therefore alternative test methods are used which may be divided into quantitative tests and qualitative tests.

Quantitative tests rely on estimating the ground truth in some way and comparing it to the estimated data. One way is to move the subject and/or limbs along a well-defined path where the ground truth is known, e.g., a rectangular pattern on a table [47], or a circular groove engraved in a glass plate [75]. Another approach is to use a static object which is measurable, e.g., a doll [71], and yet another is to use synthetic data [77] or hand segmented data [18].

Qualitative tests are widely used and rely on visual inspection. The most straightforward way is to project the estimated 3D pose into the image and inspect, visually, how alike the two are. The projection may be on top of the subject [45] or next to him or her [80]. Another form of visual inspection is to apply the estimated motion to a virtual character and see (from various viewpoints) if the movements seem human-like [120].

In the future when the motion capture devices based on active sensing, see Section 3.1.2, become cheaper, less noisy, and easier to use, this might be a good way to obtain ground truth.

3.6 Recognition

The recognition aspect of human motion capture can be seen as a kind of post processing. It is relevant to include since the recognition guides the development of many motion capture systems as it is their final or long time goal. The recognition is usually carried out by classifying the captured motion as one of several types of actions. The actions are normally simple, such as walking and running, but more advanced actions such as different ballet dance steps have also been studied.

Traditionally two different paradigms exist: recognition by reconstruction and direct recognition. The former is based on the concept of first reconstructing the scene and then recognising it, while the latter recognises directly on the low-level data, e.g., motion, without much preprocessing. Which one is preferable is hard to say and both paradigms may be supported by studies into the human visual system. Johansson [69] showed in his moving lights displays (MLD) experiments that the actions of a human may be recognised solely on motion (of the lights). This totally agrees with the idea of recognition directly through motion. Later Sumi [134] tried to redo Johansson's experiments but turned the data upside down. This resulted in a very poor recognition rate suggesting that a variant model of some kind is used when recognising motion¹. This could, of course, be a motion model but it could also be a geometric model which agrees with the former paradigm. Another fact which supports the latter is that humans may recognise different postures from a single frame, i.e., without any motion cues.

The reconstruction paradigm is most relevant for motion capture due to its strong relation to pose estimation, but both paradigms can be seen in the literature reviewed for this survey. Common for the reviewed recognition systems is that if pose

¹In a panel session at a CVPR'2000 postconference workshop "Human Modeling, Analysis and Synthesis" Pietro Perona performed a similar experiment. He showed a MLD sequence which everybody in the audience recognised as human motion. He then showed another MLD sequence which the majority of the audience believed originated from a huge spider (without considering the inherent difficulties involved in actually estimating the motion of a spider's limbs). In fact, it was the exact same MLD sequence, but visualised from above!

estimation is used, it is simply a tool for generating higher level data. The various systems may be taxonomised with respect to the two paradigms, but some of the systems do not follow one of the paradigms in a strict sense, e.g., when recognition is based directly on image data (not motion). Instead a structure of first representing and then classifying the data is used.

A more relevant distinction is to look at whether the recognition is static or dynamic, i.e., whether the recognition is based on one or more frames.

3.6.1 Static Recognition

Static recognition is concerned with spatial data, one frame at a time. The approaches usually compare prestored information with the current image. The information may be templates [133], transformed templates [114], normalised silhouettes [52], or postures [22]. The goal of static recognition is mainly to recognise various postures, e.g., pointing [74], standing and sitting [6], or specially defined postures used in interfaces [4, 41].

In the interactive Karaoke system build by Sul *et al.* [133] the postures of the subject are used to trigger and control the system. The postures are recognised by comparing the image data to different prerecorded templates. Templates are also used in the work by Oren *et al.* [114]. Offline they segment pedestrians in a number of images and generate a common template based on Haar wavelets. In run-time the template is compared to various parts of an image to find pedestrians. In the work by Freeman *et al.* [41] a computer chip capable of doing on-board image calculations is presented. The chip is used to calculate the orientation histogram of an image in real time. These histograms are matched against prerecorded histograms of various human postures and the current pose may be found. In the work by Jovic *et al.* [70] a dense depth map of the scene wherein a subject is pointing is used. After a depth-background subtraction the data are classified into points belonging to the arm and points belonging to the rest of the body. The index finger and top of the head are found as the two extreme points in the two classes and the line through them defines the pointing direction.

3.6.2 Dynamic Recognition

These approaches use temporal characteristics in the recognition task. Relatively simple activities, such as walking, are typically used as test scenarios. The systems may use low-level or high-level data.

Low-level recognition is typically based on spatio-temporal data without much processing. The data are spatio-temporal templates [25] and motion templates [121]. The goal is usually to recognise whether a human is walking in the scene or not [54]. More high level methods are usually based on pose estimated data. Such methods

vary from correlation [12] and silhouette matching [32] to HMMs [17] and neural networks [54]. The objective is to recognise actions such as walking [124], carrying objects [51], removing and placing objects [96], pointing and waving [32], gestures for control [4], standing vs walking [53], walking vs jogging [116], walking vs running [43], and classifying various aerobic exercises [123] or ballet dance steps [21].

Chomat and Crowley [25] generate motion templates by using a set of temporal-spatial filters computed by PCA. A Bayes classifier is used to perform action selection. In the work by Polana and Nelson [121] motion templates are also used but in a different way. Six subsequent motion images are computed. Each motion image is rerepresented by a subsampling where each new pixel contains the number of motion pixels from the original motion image. By representing the six subsampled images as a vector it can be classified using standard techniques. Niyogi and Adelson [108] also use temporal templates, but in a very special way. They generate an XT-slice by concatenating one of the lowest rows from each image in a sequence. If a walking figure is present an ankle profile of the figure can be seen in the XT-slice. By comparing this to prerecorded templates a walking figure can be recognised, and using a k -nearest neighbours with Euclidean distance measures various persons may be recognised by their walking pattern. The idea of using only a small part of the body to recognise the walking person is also used in the work by Heisele and Wohler [54]. They segment the area containing the legs over time and process it by a neural network. In this way they can recognise whether a pedestrian is present or not. Davis and Bobick [33, 35] also use temporal templates which are created based on motion. They use the information about where and how much motion has been present in a sequence of frames. Both a motion-energy image and a motion-history image are used. By representing the templates by its seven Hu-moments [61] a Mahalanobis distance can be used to classify the action of the subject by comparing it to the Hu-moments of prerecorded actions.

More high-level dynamic recognition is usually based on pose estimated data for the different limbs. In the work by Ju *et al.* [73] the recognition of movements is based on the motion parameters of the individual body parts (legs). The problem is viewed as matching the curves of the motion parameters against a set of known curves. Yacoob and Black [150] use the results from [73] and build a classifier on top. After translating and scaling the motion parameters a PCA is used to obtain a more compact and discriminative representation. Four different activities are recognised (four variants of walking) by comparing the PCA-transformed data to offline generated (and PCA-transformed) data sets. In the work by Bharatkumar *et al.* [12] a stick-figure representation of the legs of a walking human are matched against a kinematic human walking model. This is done by correlating the two data sets with each other. Fujiyoshi and Lipton [43] also use a stick-figure representation of the legs of a person walking. They transform the time signal of the various parameters into the frequency domain, enabling them to separate walking from running. In the work by Bregler [17] the idea of representing motion data by *movemes* (similar to

phonemes in speech recognition) is suggested. This makes it possible to compose a complex activity (word) out of simple movemes. An HMM is used to classify three different gait categories: running, walking, and skipping. This type of high level symbolic representation is also used in the work by Wren *et al.* [144]. They automatically build a behaviour alphabet (a behaviour is similar to a moveme) and model each behaviour using an HMM. The alphabet is used to classify different types of actions in a simple virtual reality game and to distinguish between the playing style of different subjects.

3.7 Discussion

It is difficult to compare and grade the different systems and methods described in the four previous sections. The systems are based on different test data and different assumptions. However, some general comments can be made concerning the application areas in relation to some performance parameters.

3.7.1 Performance Characterisation

In the Introduction, three main application areas was be identified. In Table 3.6 these are related to three concepts which may be considered the main performance parameters in any motion capture system: robustness, accuracy, and speed. The symbols in the table indicate the performance requirements for each of the applications.

	Surveillance	Control	Analysis
Robustness	+	+/-	-
Accuracy	-	+	+
Speed	+	+	-

Table 3.6: The three application areas and their requirements of the three main performance parameters.

The *robustness* of a system is here related to the various assumptions shown in Table 3.2. The fewer assumptions a system imposes on its operational conditions, the more robust it is considered to be. Surveillance systems are aiming at very robust performance since they will often be working continuously at remote and uncontrolled locations, e.g., a parking lot. They should not be sensitive to changes in lighting, weather, number of people, clothes, etc. Furthermore, they are required to work autonomously and for long periods of time. Some control applications, e.g., gestures for signalling, are for the same reasons also subjected to this high level of robustness. Other systems in the control applications, e.g., direct avatar

control, do not necessarily need this level of robustness, since they may operate in well-defined environments and for shorter periods of time where a number of assumptions conveniently may be applied. In the analysis applications the situation is similar. It is often possible to use a highly controlled setup and therefore the robustness is not the most important issue.

The *accuracy* of a system refers to how close the captured motion corresponds to the actual motion performed by the subject. Generally speaking the accuracy is directly proportional to the size of the human in the image. In surveillance systems the accuracy is rarely a key issue. It may not matter whether a surveyed subject is recognised to be walking around a car in a 3- or 4-m radius. The important thing is that the behaviour is recognised. For the control purpose the situation is different. In many applications a one-to-one mapping between the movement carried out by the subject and the action he or she controls is critical, e.g., in tele-surgery and other futuristic telepresence applications. In the last area the need for accuracy is even more pronounced, since the applications rely directly on the pose estimation output.

The processing *speed* of a system is usually divided into real-time processing and offline processing. In this context the definition of real-time is not clear even though many researchers use it about their systems. One definition is that each frame is processed before a new frame is recorded. Another definition relate to the motion which is being captured. A simpler way to view speed is to divide it into online and offline processing. In that sense surveillance systems require high speed since the images need to be processed before the car is stolen! In the control applications an even higher requirement for speed is needed. Actually the second definition of real-time would apply in this case. In the last application area the processing may be performed offline, and if so no special requirements for high speed are needed.

3.7.2 State of the Art

In this section we will discuss state of the art within the three application areas. Furthermore, we will relate the three areas to the major issues which have emerged from the survey.

Surveillance applications are mainly carried out in uncontrolled environments. Therefore the figure-ground segmentation relies mostly on motion data, since these are less dependent on various assumptions such as a known subject, known lighting, and different markers. For the same reason, object-based representation (Table 3.4) is the natural way of representing the images at a higher level. Surveillance applications are generally more focused on tracking than on any of the three other processes in Fig. 3.1. Therefore, the surveillance applications mainly use one of the two simple forms of pose estimation where no model or only an indirect use of a model is used. Due to the nature of this application area the dynamic recognition approach is widely used.

An example of state of the art is the W4-system by Haritaoglu *et al.* [53] where the aim is to survey and recognise interactions between people and people or objects in an outdoor setting². They detect and track multiple people and their body parts. The system works with monocular gray-scale images and infrared images. It uses a standard predict-match-update scheme, where it matches predicted objects or persons with measured objects or persons (in the image). The objects are obtained by detecting movements using an adaptive background subtraction, yielding a motion boundary box. The position (and motion) parameters of a person are estimated in two steps, first median matching for a coarse matching and then silhouette correlation between two consecutive frames for the fine matching. The individual body parts, head, torso, hands, legs, and feet, are found using a cardboard model of a walking human as reference. Online, the system is able to track multiple people and their limbs and coping with occlusions. Furthermore, it can detect and track objects carried and exchanged by people [51].

The ultimate surveillance system tracks multiple humans and all their (inter)actions in real-time. Haritaoglu *et al.* [53] may be heading in the right direction, but improved performance is required for handling a dynamic background and nontrivial movement patterns. Furthermore, a more intelligent handling of multiple objects and their occlusion is needed. Perhaps a further step could be to use 3D data in a direct model-based pose estimation scheme.

Some control applications are concerned with the recognition of gestures. In this case the methods used are generally similar to the ones used in the surveillance applications. However, if the application is more in the form of direct animation, e.g., avatar control, different methods are used. This type of application is carried out in an indoor setting where a number of assumptions may be introduced, e.g., known subject, known background, and known start pose. Then the appearance-based figure-ground segmentation methods are applied. To obtain good accuracy, direct use of a human model is usually used.

As an example of state of the art we consider the work by Wren [144]. First of all they use the Pfinder algorithm [143] as the underlying tracking methods. It is a probabilistic methods which segments the subject into a number of blobs and track those over time. This method has proven to be fast, robust, and able to directly estimate the positions of the head and hands, which are of great importance in control applications. They apply two Pfinder algorithms to obtain 3D estimates of the hands and head. Using a human model and kinematic constraints they estimate the 3D pose of the upper body. In the framework of a Kalman filter the model is predicted into the next frame to support the blob segmentation and tracking. The innovation of the Kalman filter is used to learn the various motion patterns (behaviours) of the subject. These can then be incorporated into the filter to improve the state estimates and predictions, i.e., a better pose estimation result.

²The W4-system has also been used in an indoor setting [111] where it was extended with Kalman filters and kinematic constraints from [145].

The last application area is concerned with analysis of the captured motion and is typically used for clinical studies. These applications are carried out in well-controlled environments meaning that a number of assumptions may be introduced. In commercial systems markers are used which allow point representation of the data. A model of the human is necessary for interpreting the data. Usually it is not used directly in the pose estimation, but rather indirectly. The use of markers yields stable tracking, but the obtained points are not placed directly on the skeleton (for obvious reasons). Therefore an offset distance is introduced between the markers and the physical skeleton. This is, besides initialisation and calibration, where the main problems are in state of the art commercial systems. In future systems one goal is to move away from the marker approach and aim at the more pure computer vision solution without the use of markers [38]. This will make systems more flexible and less cumbersome. The solution may be based on detailed human models used directly in the processing.

An example of direct use of a model is the work by Gavrilu and Davis [45]. They use a model-based approach to track a subject in 3D. A recognition cycle goes as follows. Based on the current and previous states, the allowed intervals for each body-parameter (e.g., joint angles) are predicted. For each combination of the 22 body parameters the human model is synthesised from the cameras' point of view. They compare edges between the synthesised model and the images and thereby (re)formulate the problem as a search problem - how to compare two edge images (a real image with a synthesised image). The search problem is solved using a robust variant of chamfer matching. When they find the best fit (highest similarity measure) the model is updated using these parameters. They use four synchronised sequences from four different cameras and run the algorithm for each view. In order to obtain stable edges they wear tight-fitting coloured clothes. The high number of joints in their relatively detailed model, the four cameras, and relatively few assumptions make it a rather complex system which, to some extent, is able to estimate the pose of an entire subject.

Future work in marker-free systems includes improved initialisation to obtain a good model and the pose of the current subject fast and reliably. Although the analysis-by-synthesis approach seems to be the right one, it is still rather slow and computational demanding. Methods to prune the state space and faster optimisation schemes are required. Generally we may expect to see workable systems in the analysis area before we see them in the control area because the requirements for speed are relaxed in the analysis applications.

3.7.3 Future Directions

Although assumptions might be acceptable (e.g., Chroma-keying) for some application, it is evident from the number of assumptions applied in the papers reviewed for this survey that the research field is still in a phase of development. Perhaps

inspiration may be found in related research fields, e.g., speech recognition. First of all a tremendous amount of time is spent on recording and labelling training data in this field. These data are of a general nature, i.e., they suit a number of speech recognition tasks, and represented in a well-defined modelling language which are the atoms (e.g., phonemes) of the spoken language. One reason for not spending the same amount of time on the training phase in computer vision might be the lack of a general underlying modelling language, i.e., how to map the images into symbols. An alphabet consisting of motion-entities would make computer vision-based human motion capture much easier, since it will transform the pose estimation problem into a recognition problem, i.e., recognise a sequence of symbols. This has already happened in Bregler [17] where the letters are called movemes, and Wren and Pentland [144] where the letters are called behaviours. Even though their alphabets are rather limited it is still a step in a very interesting direction. Furthermore, by having such an alphabet a vocabulary may be introduced to constrain the task at hand, as is the case in speech recognition.

Besides the current lack of a general alphabet, another reason for not using extensive training data is the amount of time required to actually capture and label human motion data. One solution is to use commercial motion capture systems (e.g., magnetic sensors) [16] which, when calibrated, easily produce thousands of labelled data sets. Another solution is to apply computer graphics to synthesise the appearance of a human model from various viewpoints, as in [125].

Another aspect of speech recognition, which is actually being seen more and more in computer vision, is the use of probabilistic models for aspects other than recognition, e.g., modelling the position of the head using a Gaussian density. Some of these models are learned automatically using unsupervised methods, such as the EM algorithm. The entire tracking framework is also widely based on probabilistic methods such as the Kalman filter and the Condensation algorithm. Also, HMMs [122] and neural networks [125] have found their way into tracking and pose estimation. In future systems more of this may be expected due to the methods' ability to handle uncertainties and to suppress noise.

Even though interesting results such as [16] have recently arisen from methods not using a human model, the direct use of a human model seems to be the preferred trend. From Table 3.5 it can be seen that the choice of model type differs while silhouettes seem to be the preferred abstraction levels.

The use of silhouettes is motivated by the presence of simple algorithms for their estimation. They are easier to estimate than joints and a stick-figure, and their region-based nature makes them more robust to noise than local information such as edges. Furthermore, useful silhouettes might be extracted from relatively low resolution images. Due to the global nature of silhouettes details are likely to be missing. This results in additional complexity when trying to estimate a 3D pose from 2D silhouettes. Future work should therefore consider combining a silhouette representation with data capable of representing the interior of the silhouette, i.e.,

its relation to the human skeleton structure. An example is to combine silhouettes with the positions of the hands and head as seen in [100] and [113].

The use of motion as an abstraction level in pose estimation is also rather popular due to its inherent relation to the application. The motion in the images may be linked directly to the motion of the various limbs. Furthermore, many image points might be used to estimate the motion parameters. We expect the use of motion as a cue to be used more extensively in the future. However, to achieve success a number of issues still need to be addressed. First, the methods are based on incremental updates which rely on local (both spatial and temporal) smoothness. Therefore they often rely on a number of assumptions such as no occlusion and the subject being the only moving object in the image. Moreover, due to the incremental update, the initial pose is required and the systems have no way to recover after a total loses of track, lacking a mechanism for globally searching the entire image. Another problem is the risk of accumulating errors due to the incremental procedure. One solution is to use key frames as suggested in [154]. Given the initial and final pose parameters, both forward and backward iteration may ensure a consistent pose sequence. Alternatively, one might combine a motion-based method with a method based on spatial data. For example, it could be interesting to see image measurements and a human model linked by both the motion framework of Bregler and Malik [18] and the edge comparisons of Gavrilu and Davis [45].

In addition to the problems related to incremental updates, another issue also has to be considered. Many movements become ambiguous when projected into the image plane; e.g., rotation about an axis parallel to the image plane will produce the same optical flow field as a translation in a certain direction will. Furthermore, movements along the optical axis are difficult to register robustly. To solve these problems multiple cameras are required or multiple data types as in the work by Okada *et al.* [112]. They combine motion data and depth data to resolve the ambiguities, thereby making the pose estimation more robust. In [110] a more general discussion on combining motion and depth data is given.

Generally speaking it seems to be a good and necessary approach to combine various data types to broaden invariance and robustness to all possible situations. Another promising approach is to use future measurements when processing the current data, i.e., allow a lag in the output. This helps to resolve ambiguities [128][16].

3.8 Conclusion

Human motion capture goes back to at least the 1870s when Marey [95] and Muybridge [105] started their work. But recently, new technologies have made the motion capture problem popular as more convenient and affordable devices such as cameras, magnetic trackers, and computer power have become available.

Advances in active sensors, e.g., magnetic trackers, are making them cheaper, smaller,

more precise, and generally easier to use. They will, however, still be cumbersome and limited in their use due to the need for special hardware. Therefore, computer vision could provide an attractive touch-free alternative.

The solutions developed to date are all based on a number of assumptions to make the problem tractable. This, together with the relatively simple methods being used, can be seen as an indication of the current state of the field - as being in its early development. The latest systems, however, use more advanced methods based on comprehensive probabilistic models and advanced training. Nevertheless, some assumptions are still required and we are far from a general solution to the human motion capture problem. Some of the general key issues needing to be addressed are *initialisation*, *recover from failure*, and *robustness*.

Many systems are based on knowing the initial state of their system and/or a well-defined model fitted (offline) to the current subject. In a real life scenario we may expect a system to run on its own, i.e., adapt to the current situation. This might seem a minor problem, but what if none of the current research directions result in a system capable of autonomy? Should two parallel direction be followed: one for initialisation and one for processing or should we aim at a common solution?

Related is the problem of how to recover from failure. A number of systems are based on incremental updates or searching around a predicted value. Many of these fail due to occlusion, bad predictions, and a change in the framerate/camera focus/image resolution and are not able to recover. This is an important problem since real life applications are likely to challenge a system by new situations not included in the design and training and hereby making it fail from time to time.

The robustness relates to the number of assumptions applied in systems, but also to the fact that most systems are tested on less than 1000 frames. How can one justify to evaluate the robustness of a system within such a short lifespan? Long test sets available for everybody need to be generated (as in the face recognition community) to evaluate the robustness of individual systems and compare various systems.

For future systems to be more successful and less dependent of various assumptions new methods and a combination of current methods should be developed, i.e., the combination of various image cues, such as motion and silhouettes, and more extensive and adaptive use of human models. Furthermore, new sensors or combinations of sensors might also be an interesting path into the future.

The rapid developments in computer graphics may benefit human motion capture. Until recently the computer graphics field has been mostly interested in visual realism³ and personalised human models, while the motion capture community has been more interested in spatial accuracy of the human models. We expect that the commercial interest in both fields will accelerate the development in human modelling and make the two fields approach and benefit from each other.

³A more comprehensive discussion on human animation can be found in [8].

The applications of human motion capture are numerous and it is expected that we will see a continuous growth in the resources devoted to this topic and hence that interesting new results in spite of everything will appear in the not too distant future.

Acknowledgement

We would like to thank Moritz Störring and Hanne E. Andreassen for help in editing this document, and the Danish National Research Councils, who through the project: "The Staging of Virtual Inhabited 3D Spaces" funded this work.

References

- [1] J.K. Aggarwal and Q. Cai. Human Motion Analysis: A Review. *Computer Vision and Image Understanding*, 73(3), 1999.
- [2] 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*, pages 2–14, Austin, Texas, USA, 1994.
- [3] K. Akita. Image Sequence Analysis of Real World Human Motion. *Pattern Recognition*, 17(1):73–83, 1984.
- [4] J. Amat, M. Casals, and M. Frigola. Stereoscopic System for Human Body Tracking in Natural Scenes. In *International Workshop on Modeling People at ICCV'99*, Corfu, Greece, September 1999.
- [5] B. Andersen, T. Dahl, M. Iversen, M. Pedersen, and T. Søndergaard. Human Motion Capture. Technical report, Laboratory of Image Analysis, Aalborg University, Denmark, January 1999.
- [6] C.I. Attwood, G.D. Sullivan, and K.D. Baker. Model-based Recognition of Human Posture Using Single Synthetic Images. In *Fifth Alvey Vision Conference*, University of Reading, UK, 1989.
- [7] A. Azarbayejani, C.R. Wren, and A.P. Pentland. Real-Time 3-D Tracking of the Human Body. In *IMAGE'COM 96*, Bordeaux, France, May 1996.
- [8] N. Badler. Virtual Humans for Animation, Ergonomics, and Simulation. In *Workshop on Motion of Non-Rigid and Articulated Objects*, Puerto Rico, USA, 1997.
- [9] A.M. Baumberg and D.C. Hogg. An Efficient Method for Contour Tracking using Active Shape Models. In *Workshop on Motion of Non-Rigid and Articulated Objects*, pages 2–14, Austin, Texas, USA, 1994.

-
- [10] D.A. Becker and A. Pentland. Staying Alive: A Virtual Reality Visualization Tool for Cancer Patients. In *AAAI'96 Workshop on Entertainment and Alife/AI*, Portland, Oregon, USA, August 1996.
 - [11] A.P. Bernat, J. Nelan, S. Riter, and H. Frankel. Security Applications of Computer Motion Detection. *Applications of Artificial Intelligence V*, 786, 1987.
 - [12] A.G. Bharatkumar, K.E. Daigle, M.G. Pandey, Q. Cai, and J.K. Aggarwal. Lower Limb Kinematics of Human Walking with the Medial Axis Transformation. In *Workshop on Motion of Non-Rigid and Articulated Objects*, Austin, Texas, USA, 1994.
 - [13] A. Blake and M. Isard. *Active Contours*. Springer, 1998.
 - [14] A.F. Bobick and J.W. Davis. An Appearance-based Representation of Action. In *International Conference on Pattern Recognition*, 1996.
 - [15] A. Bottino, A. Laurentini, and P. Zuccone. Toward Non-intrusive Motion Capture. In *Asian Conference on Computer Vision*, 1998.
 - [16] M. Brand. Shadow Puppetry. In *International Conference on Computer Vision*, Corfu, Greece, September 1999.
 - [17] C. Bregler. Learning and Recognizing Human Dynamics in Video Sequences. In *Conference on Computer Vision and Pattern Recognition*, San Juan, Puerto Rico, 1997.
 - [18] C. Bregler and J. Malik. Tracking People with Twists and Exponential Maps. In *International Conference on Computer Vision and Pattern Recognition*, Santa Barbara, California, June 1998.
 - [19] Q. Cai and J.K. Aggarwal. Tracking Human Motion Using Multiple Cameras. In *International Conference on Pattern Recognition*, 1996.
 - [20] Q. Cai, A. Mitiche, and J.K. Aggarwal. Tracking Human Motion in an Indoor Environment. In *International Conference on Image Processing*, 1995.
 - [21] L. Campbell and A. Bobick. Recognition of Human Body Motion Using Phase Space Constraints. In *International Conference on Computer Vision*, Cambridge, Massachusetts, 1995.
 - [22] L. Campbell and A. Bobick. Using Phase Space Constraints to Represent Human Body Motion. In *International Workshop on Automatic Face- and Gesture-Recognition*, Zurich, Switzerland, 1995.
 - [23] C. Cedras and M. Shah. Motion-Based Recognition: A Survey. *Image and Vision Computing*, 13(2):129–155, 1995.

-
- [24] T.J. Cham and J.M. Rehg. A Multiple Hypothesis Approach to Figure Tracking. In *Conference on Computer Vision and Pattern Recognition*, Fort Collins, Colorado, USA, June 23-25 1999.
- [25] O. Chomat and J.L. Crowley. Recognizing Motion Using Local Appearance. In *International Symposium on Intelligent Robotic Systems*, University of Edinburgh, 1998.
- [26] C. Christensen and S. Corneliussen. Tracking of Articulated Objects using Model-Based Computer Vision. Technical report, Laboratory of Image Analysis, Aalborg University, Denmark, June 1997.
- [27] C. Christensen and S. Corneliussen. Visualization of Human Motion using Model-based Vision. Technical report, Laboratory of Image Analysis, Aalborg University, Denmark, January 1997.
- [28] J.M. Chung and N. Ohnishi. Cue Circles: Image Feature for Measuring 3-D Motion of Articulated Objects Using Sequential Image Pair. In *International Conference on Automatic Face and Gesture Recognition*, Nara, Japan, 1998.
- [29] C.R. Corlin and J. Ellesgaard. Real Time Tracking of a Human Arm. Technical report, Laboratory of Image Analysis, Aalborg University, Denmark, January 1998.
- [30] A. Cretual, F. Chaumette, and P. Bouthemy. Complex Object Tracking by Visual Servoing Based on 2D Image Motion. In *International Conference on Pattern Recognition*, 1998.
- [31] R. Cutler and L. Davis. Real-Time Periodic Motion Detection, Analysis, and Applications. In *Conference on Computer Vision and Pattern Recognition*, Fort Collins, Colorado, USA, June 23-25 1999.
- [32] T. Darrell, P. Maes, B. Blumberg, and A.P. Pentland. A Novel Environment for Situated Vision and Behavior. In *Workshop for Visual Behaviors at CVPR-94*, 1994.
- [33] J.W. Davis and A. Bobick. The Representation and Recognition of Action Using Temporal Templates. In *Conference on Computer Vision and Pattern Recognition*, 1997.
- [34] J.W. Davis and A. Bobick. SIDeshow: A Silhouette-based Interactive Dual-screen Environment. Technical Report 457, MIT Media Lab, 1998.
- [35] J.W. Davis and A. Bobick. Virtual PAT: A Virtual Personal Aerobics Trainer. In *Workshop on Perceptual User Interfaces*, San Francisco, November 1998.

-
- [36] L. Davis, S. Fejes, D. Harwood, Y. Yacoob, I. Hariatoglu, and M.J. Black. Visual Surveillance of Human Activity. In *Asian Conference on Computer Vision*, Mumbai, India, 1998.
- [37] Q. Delamarre and O. Faugeras. 3D Articulated Models and Multi-view Tracking with Silhouettes. In *International Conference on Computer Vision*, Corfu, Greece, September 1999.
- [38] B. Delaney. On the Trail of the Shadow Women: The Mystery of Motion Capture. *Computer Graphics and Applications*, 18(5):14–19, 1998.
- [39] A. Dempster, N. Laird, and D. Rubin. Maximum Likelihood Estimation from Incomplete Data via the *EM* Algorithm. *Journal of the Royal Statistical Society (B)*, 39(1), 1977.
- [40] I. Douros, L. Dekker, and B.F. Buxton. An Improved Algorithm for Reconstruction of the Surface of the Human Body from 3D Scanner Data Using Local B-Spline Patches. In *International Workshop on Modeling People at ICCV'99*, Corfu, Greece, September 1999.
- [41] W.T. Freeman, K. Tanaka, J. Ohta, and K. Kyuma. Computer vision for computer games. In *International Workshop on Automatic Face- and Gesture-Recognition*, Zurich, Switzerland, 1995.
- [42] P. Fua, A. Gruen, R. Plänkers, N. D'Apuzzo, and D. Thalmann. Human Body Modeling and Motion Analysis From Video Sequences. In *International Symposium on Real Time Imaging and Dynamic Analysis*, Hakodate, Japan, June 1998.
- [43] H. Fujiyoshi and A.J. Lipton. Real-Time Human Motion Analysis by Image Skeletonization. In *Workshop on Applications of Computer Vision*, 1998.
- [44] D.M. Gavrila. The Visual Analysis of Human Movement: A Survey. *Computer Vision and Image Understanding*, 73(1):82–98, 1999.
- [45] D.M. Gavrila and L.S. Davis. 3-D Model-Based Tracking of Humans in Action: A Multi-View Approach. In *Conference on Computer Vision and Pattern Recognition*, San Francisco, USA, 1996.
- [46] L. Goncalves, E.D. Bernardo, and P. Perona. Reach Out and Touch Space (Motion Learning). In *International Conference on Automatic Face and Gesture Recognition*, Nara, Japan, 1998.
- [47] L. Goncalves, E.D. Bernardo, E. Ursella, and P. Perona. Monocular Tracking of the Human Arm in 3D. In *International Conference on Computer Vision*, Cambridge, Massachusetts, 1995.

-
- [48] H. Gu, Y. Shirai, and M. Asada. MDL-Based Spatiotemporal Segmentation from Motion in a Long Image Sequence. In *Conference on Computer Vision and Pattern Recognition*, 1994.
- [49] J. Gu, T. Chang, I. Mak, S. Gopalsamy, H.C. Shen, and M.M.F. Yuen. A 3D Reconstruction System for Human Body Modeling. In *Lecture Notes in Artificial Intelligence 1537. Modeling and Motion Capture Techniques for Virtual Environments*, 1998.
- [50] Y. Guo, G. Xu, and S. Tsuji. Tracking Human Body Motion Based on a Stick Figure Model. *Journal of Visual Communication and Image Representation*, 5(1):1–9, 1994.
- [51] I. Haritaoglu, R. Cutler, D. Harwood, and L.S. Davis. *Backpack*: Detection of People Carrying Objects Using Silhouettes. In *International Conference on Computer Vision*, Corfu, Greece, September 1999.
- [52] I. Haritaoglu, D. Harwood, and L.S. Davis. Ghost: A Human Body Part Labeling System Using Silhouettes. In *International Conference on Pattern Recognition*, 1998.
- [53] 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.
- [54] B. Heisele and C. Wohler. Motion-Based Recognition of Pedestrians. In *International Conference on Pattern Recognition*, 1998.
- [55] A. Hilton. Towards Model-Based Capture of a Persons Shape, Appearance and Motion. In *International Workshop on Modeling People at ICCV'99*, Corfu, Greece, September 1999.
- [56] 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*, pages 174–185, May 1999.
- [57] A. Hilton, D. Beresford, T. Gentils, R. Smith, and W. Sun. Whole-body modelling of people from multi-view images to populate virtual worlds. *The Visual Computer*, pages 411–436, 2000.
- [58] D. Hogg. Model-Based Vision: A Program to See a Walking Person. *Image and Vision Computing*, 1(1), February 1983.
- [59] D.C. Hogg. *Interpreting Images of a Known Moving Object*. PhD thesis, University of Sussex, UK, 1984.

- [60] C. Hu, Q. Tu, Y. Li, and S. Ma. Extraction of Parametric Human Model for Posture Recognition Using Genetic Algorithm. In *The fourth International Conference on Automatic Face and Gesture Recognition*, Grenoble, France, March 2000.
- [61] M. Hu. Visual Pattern Recognition by Moment Invariants. *IRE Trans. Information Theory*, IT-8(2):179–187, 1962.
- [62] E.A. Hunter, P.H. Kelly, and R.C. Jain. Estimation of Articulated Motion Using Kinematically Constrained Mixture Densities. In *Workshop on Motion of Non-Rigid and Articulated Objects*, Puerto Rico, USA, 1997.
- [63] S. Ioffe and D. Forsyth. Finding people by sampling. In *International Conference on Computer Vision*, Corfu, Greece, September 1999.
- [64] M. Isard and A. Blake. CONDENSATION - conditional density propagation for visual tracking. *International Journal on Computer Vision*, pages 5–28, 1998.
- [65] Y. Iwai, K. Ogaki, and M. Yachida. Posture Estimation using Structure and Motion Models. In *International Conference on Computer Vision*, Corfu, Greece, September 1999.
- [66] S. Iwasawa, K. Ebihara, J. Ohya, and S. Morishima. Real-Time Estimation of Human Body Posture from Monocular Thermal Images. In *Conference on Computer Vision and Pattern Recognition*, 1997.
- [67] S. Iwasawa, J. Ohya, K. Takahashi, T. Sakaguchi, S. Kawato, K. Ebihara, and S. Morishima. Real-Time Estimation of Human Body Posture from Trinocular Images. In *International Workshop on Modeling People at ICCV'99*, Corfu, Greece, September 1999.
- [68] S. Iwasawa, J. Ohya, K. Takahashi, T. Sakaguchi, K. Ebihara, and S. Morishima. Human Body Postures from Trinocular Camera Images. In *The fourth International Conference on Automatic Face and Gesture Recognition*, Grenoble, France, March 2000.
- [69] G. Johansson. Visual Motion Perception. In *Scientific American*, pages 76–88. June 1975.
- [70] N. Jovic, B. Brumitt, B. Meyers, S. Harris, and T. Huang. Detection and Estimation of Pointing Gestures in Dense Disparity Maps. In *The fourth International Conference on Automatic Face and Gesture Recognition*, Grenoble, France, March 2000.
- [71] N. Jovic, J. Gu, H.C. Shen, and T. Huang. 3-D Reconstruction of Multipart Self-Occluding Objects. In *Asian Conference on Computer Vision*, 1998.

-
- [72] S. Ju. Human Motion Estimation and Recognition (Depth Oral Report). Technical report, University of Toronto, 1996.
- [73] S.X. Ju, M.J. Black, and Y. Yacoob. Cardboard People: A parameterized Model of Articulated Image Motion. In *International Conference on Automatic Face and Gesture Recognition*, Killington, Vermont, USA, 1996.
- [74] R.E. Kahn and M.J. Swain. Gesture Recognition Using the Perseus Architecture. Technical Report TR-96-04, Department of Computer Science, University of Chicago, 1996.
- [75] I. Kakadiaris and D. Metaxas. Vision-Based Animation of Digital Humans. In *Conference on Computer Animation*, pages 144–152, 1998.
- [76] I.A. Kakadiaris and D. Metaxas. 3D Human Body Model Acquisition from Multiple Views. In *International Conference on Computer Vision*, pages 618–623, Cambridge, Massachusetts, June 20-23 1995.
- [77] I.A. Kakadiaris and D. Metaxas. Model-Based Estimation of 3D Human Motion with Occlusion Based on Active Multi-Viewpoint Selection. In *Conference on Computer Vision and Pattern Recognition*, 1996.
- [78] Y. Kameda and M. Minoh. A Human Motion Estimation Method Using 3-Successive Video Frames. In *International Conference on Virtual Systems and Multimedia*, 1996.
- [79] Y. Kameda, M. Minoh, and K. Ikeda. Three Dimensional Pose Estimation of an Articulated Object from its Silhouette Image. In *Asian Conference on Computer Vision*, Osaka, Japan, November 1993.
- [80] Y. Kameda, M. Minoh, and K. Ikeda. Three Dimensional Motion Estimation of a Human Body Using a Difference Image Sequence. In *Asian Conference on Computer Vision*, 1995.
- [81] T.M. Kepple. MOVE3D - Software for Analyzing Human Motion. In *Proc. of Johns Hopkins National Search for Computing Applications to Assist Persons with Disabilities*, Laurel, Maryland, February 1992.
- [82] K. Koffka. *Principle of Gestalt Psychology*. New York, Harcourt Brace, 1935.
- [83] H.J. Lee and Z. Chen. Determination of 3D Human Body Posture from a single View. *Computer Vision, Graphics, and Image processing*, 30:148–168, 1985.
- [84] H.J. Lee and Z. Chen. Knowledge-Guided Visual Perception of 3-D Human Gait from a Single Image Sequence. *Transactions on Systems, Man, and Cybernetics*, 22(2):336–342, 1992.

-
- [85] J. Lengyel. The Convergence of Graphics and Vision. *Computer*, 31(7):46–53, 1998.
- [86] F. Lerasle, G. Rives, and M. Dhome. Human Body Limbs Tracking by Multiocular Vision. In *Scandinavian Conference on Image Analysis*, Lappeenranta, Finland, 1997.
- [87] F. Lerasle, G. Rives, and M. Dhome. Tracking of Human Limbs by Multiocular Vision. *Computer Vision and Image Understanding*, 75(3):229–246, 1999.
- [88] M.K. Leung and Y.H. Yang. A Region Based Approach for Human Body Motion Analysis. *Pattern Recognition*, 20(3):321–339, 1987.
- [89] M.K. Leung and Y.H. Yang. Human Body Motion Segmentation in a Complex Scene. *Pattern Recognition*, 20(1):55–64, 1987.
- [90] M.K. Leung and Y.H. Yang. First Sight: A Human Body Outline Labeling System. *Transactions on Pattern Analysis and Machine Intelligence*, 17(4):359–377, 1995.
- [91] Y. Li, S. Ma, and H. Lu. Human Posture Recognition Using Multi-Scale Morphological Method and Kalman Motion Estimation. In *International Conference on Pattern Recognition*, 1998.
- [92] W. Long and Y.H. Yang. Log-Tracker: An Attribute-Based Approach to Tracking Human Body Motion. *Pattern Recognition and Artificial Intelligence*, 5(3):439–458, 1991.
- [93] E. Luc. Real Time Human Action Recognition for Virtual Environments. In *Computer Science Postgraduate Course*. Computer Graphics Lab, Swiss Federal Institute of Technology, Lausanne, Switzerland, September 1996.
- [94] Y. Luo, F.J. Perales, and J.J. Villanueva. An Automatic Rotoscopy System for Human Motion based on a Biomechanic Graphical Model. *Computers & Graphics*, 16(4), 1992.
- [95] E.J. Marey. *Animal Mechanism: A Treatise on Terrestrial and Aerial Locomotion*. New York: Appleton, 1873. Republished as Vol. XI of the International Scientific Series.
- [96] S.J. McKenna, S. Jabri, Z. Duric, and H. Wechsler. Tracking Interacting People. In *The fourth International Conference on Automatic Face and Gesture Recognition*, Grenoble, France, March 2000.
- [97] D. Meyer, J. Denzler, and H. Niemann. Model Based Extraction of Articulated Objects in Image Sequences. In *Fourth int. conf. on Image Processing*, 1997.

-
- [98] T.B. Moeslund. Summaries of 107 Computer Vision-Based Human Motion Capture Papers. Technical Report LIA 99-01, Laboratory of Image Analysis, Aalborg University, Denmark, 1999.
- [99] T.B. Moeslund. Interacting with a Virtual World through Motion Capture. In Lars Qvortrup, editor, *Interaction in Virtual Inhabited 3D Worlds*, chapter 11. Springer-Verlag, 2000.
- [100] T.B. Moeslund and E. Granum. 3D Human Pose Estimation using 2D-Data and an Alternative Phase Space Representation. In *Workshop on Human Modeling, Analysis and Synthesis at CVPR*, Hilton Head Island, South Carolina, June 16 2000.
- [101] T.B. Moeslund and E. Granum. Multiple Cues used in Model-Based Human Motion Capture. In *The fourth International Conference on Automatic Face and Gesture Recognition*, Grenoble, France, March 2000.
- [102] T.B. Moeslund and E. Granum. A Survey of Computer Vision-Based Human Motion Capture. *Computer Vision and Image Understanding*, 81(3), 2001.
- [103] S. Moezzi, A. Katkere, D.Y. Kuramura, and R. Jain. Reality Modeling and Visualization from Multiple Video Sequences. *Computer Graphics and Applications*, 16(6):58–63, 1996.
- [104] O. Munkelt, C. Ridder, D. Hansel, and W. Hafner. A Model Driven 3D Image Interpretation System Applied to Person Detection in Video Images. In *International Conference on Pattern Recognition*, 1998.
- [105] E. Muybridge. Animal locomotion, 1957. Reprinted in Brown, L.S.(Ed.)(1957). *Animal in motion*. New York: Dover.
- [106] A. Nakazawa, H. Kato, and S. Inokuchi. Human Tracking Using Distributed Video Systems. In *International Conference on Pattern Recognition*, 1998.
- [107] P.J. Narayanan, P.W. Rander, and T. Kanade. Constructing Virtual Worlds Using Dense Stereo. In *International Conference on Computer Vision*, Bombay, India, Januar 1998.
- [108] S.A. Niyogi and E.H. Adelson. Analyzing and Recognizing Walking Figures in XYT. In *Conference on Computer Vision and Pattern Recognition*, 1994.
- [109] J. Njåstad, S. Grinaker, and G.A. Storhaug. Estimating Parameters in a $2\frac{1}{2}$ D Human Model. In *11th Scandinavian Conference on Image Analysis*, Kangerlussuaq, Greenland, 1999.
- [110] P. Nordlund. *Figure-Ground Segmentation Using Multiple Cues*. PhD thesis, Kungl Tekniska Hogskolan, Sweden, 1998.

-
- [111] J. Ohya, J. Kurumisawa, R. Nakatsu, K. Ebihara, S. Iwasawa, D. Harwood, and T. Horprasert. Virtual Metamorphosis. *MultiMedia*, 6(2):29–39, 1999.
- [112] R. Okada, Y. Shirai, and J. Miura. Tracking a Person with 3-D motion by Integrating Optical Flow and Depth. In *The fourth International Conference on Automatic Face and Gesture Recognition*, Grenoble, France, March 2000.
- [113] E.J. Ong and S. Gong. Tracking Hybrid 2D-3D Human Models from Multiple Views. In *International Workshop on Modeling People at ICCV'99*, Corfu, Greece, 1999.
- [114] M. Oren, C. Papageorgiou, P. Sinha, E. Osuna, and T. Poggio. Pedestrian Detection Using Wavelet Templates. In *Conference on Computer Vision and Pattern Recognition*, 1997.
- [115] J. O'Rourke and N.I. Badler. Model-Based Image Analysis of Human Motion Using Constraint Propagation. *Transactions on Pattern Analysis and Machine Intelligence*, 2(6):522–536, 1980.
- [116] V. Pavlović, J.M. Rehg, T.J. Cham, and K.P. Murphy. A Dynamic Bayesian Network Approach to Figure Tracking Using Learned Dynamic Models. In *International Conference on Computer Vision*, Corfu, Greece, September 1999.
- [117] V.I. Pavlovic, R. Sharma, and T.S. Huang. Visual Interpretation of Hand Gestures for Human-Computer Interaction: A Review. *Transactions on Pattern Analysis and Machine Intelligence*, 19(7):677–695, 1997.
- [118] F.J. Perales and J. Torres. A System for Human Motion Matching between Synthetic and Real Images Based on a Biomechanic Graphical Model. In *Workshop on Motion of Non-Rigid and Articulated Objects*, pages 83–88, Austin, Texas, USA, 1994.
- [119] C. Pinhanez and A. Bobick. Using Computer Vision to Control a Reactive Computer Graphics Character in a Theater Play. In *International Conference on Vision Systems*, 1998.
- [120] R. Plänkers, P. Fua, and N. D'Apuzzo. Automated Body Modeling from Video Sequences. In *International Workshop on Modeling People at ICCV'99*, Corfu, Greece, 1999.
- [121] R. Polana and R. Nelson. Low Level Recognition of Human Motion. In *Workshop on Motion of Non-Rigid and Articulated Objects*, Austin, Tx, USA, October 1994.
- [122] G. Rigoll, S. Eickeler, and S. Müller. Person Tracking in Real-World Scenarios Using Statistical Methods. In *The fourth International Conference on Automatic Face and Gesture Recognition*, Grenoble, France, March 2000.

-
- [123] J. Rittscher and A. Blake. Classification of human body motion. In *International Conference on Computer Vision*, Corfu, Greece, September 1999.
- [124] K. Rohr. *Human Movement Analysis Based on Explicit Motion Models*, chapter 8, pages 171–198. Kluwer Academic Publishers, Dordrecht Boston, 1997.
- [125] 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.
- [126] M. Rossi and A. Bozzoli. Tracking and Counting Moving People. Technical Report 9404-03, IRST, Trento, Italy, April 1994.
- [127] M. Schneider and M. Bekker. Tracking of Human Motion. Master’s thesis, LIFIA, Grenoble, France and LIA, AAU, Denmark, 1994.
- [128] H. Segawa, H. Shioya, N. Hiraki, and T. Totsuka. Constraint-Conscious Smoothing Framework for the Recovery of 3D Articulated Motion from Image Sequences. In *The fourth International Conference on Automatic Face and Gesture Recognition*, Grenoble, France, March 2000.
- [129] H. Segawa and T. Totsuka. Torque-based Recursive Filtering Approach to the Recovery of 3D Articulated Motion from Image Sequences. In *International Conference on Computer Vision*, Corfu, Greece, September 1999.
- [130] A. Shio and J. Sklansky. Segmentation of People in Motion. In *Workshop on Visual Motion*, pages 325–332, October 1991.
- [131] H. Sidenbladh, F. De la Torre, and M.J. Black. A Framework for Modeling the Appearance of 3D Articulated Figures. In *The fourth International Conference on Automatic Face and Gesture Recognition*, Grenoble, France, March 2000.
- [132] M.C. Silaghi, R. Plänkers, R. Boulic, P. Fua, and D. Thalmann. Local and Global Skeleton Fitting Techniques for Optical Motion Capture. In *International Workshop on Modelling and Motion Capture Techniques for Virtual Environments*, Geneva, Switzerland, November 1998.
- [133] C. Sul, K. Lee, and K. Wohn. Virtual Stage: A Location-Based Karaoke System. *Multimedia*, 5(2), 1998.
- [134] S. Sumi. Upside-Down Presentation of the Johansson moving light-Spot Pattern. *Perception*, 13:283–286, 1984.
- [135] A. Tesei, G.L. Foresti, and C.S. Regazzoni. Human Body Modeling for People Localization and Tracking from Real Image Sequences. In *Image Processing and Its Applications*, July 1995.

-
- [136] T. Tsukiyama and Y. Shirai. Detection of the Movements of Persons from a Sparse Sequence of TV Images. *Pattern Recognition*, 18(3/4):207–213, 1985.
- [137] M. Turk. Visual Interaction with Lifelike Characters. In *International Conference on Automatic Face and Gesture Recognition*, Killington, VT, USA, 1996.
- [138] A. Utsumi, H. Mori, J. Ohya, and M. Yachida. Multiple-View-Based Tracking of Multiple Humans. In *International Conference on Pattern Recognition*, 1998.
- [139] 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.
- [140] S. Wachter and H.-H. Nagel. Tracking Persons in Monocular Image Sequences. *Computer Vision and Image Understanding*, 74(3):174–192, 1999.
- [141] J. Wang, G. Lorette, and P. Bouthemy. Analysis of Human Motion: A Model-Based Approach. In *Scandinavian Conference on Image Analysis*, 1991.
- [142] J. Wang, G. Lorette, and P. Bouthemy. Human Motion Analysis with Detection of Sub-Part Deformations. *SPIE - Biomedical Image Processing and Three-Dimensional Microscopy*, 1660:329–335, 1992.
- [143] C.R. Wren, A. Azarbajani, 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.
- [144] C.R. Wren, B.P. Clarkson, and A.P. Pentland. Understanding Purposeful Human Motion. In *The fourth International Conference on Automatic Face and Gesture Recognition*, Grenoble, France, March 2000.
- [145] C.R. Wren and A.P. Pentland. Dynaman: Recursive Modeling of Human Motion. *To appear in: Image and Vision Computing*.
- [146] C.R. Wren and A.P. Pentland. Dynamic Models of Human Motion. In *International Conference on Automatic Face and Gesture Recognition*, Nara, Japan, 1998.
- [147] C.R. Wren and A.P. Pentland. Understanding Purposeful Human Motion. In *International Workshop on Modeling People at ICCV'99*, Corfu, Greece, September 1999.
- [148] C.R. Wren *et al.* Perceptive Spaces for Performance and Entertainment. In *ATR Workshop on Virtual Communication Environments: Bridges over Art/Kansei and VR Technologies*, Kyoto, Japan, April 13 1998.

-
- [149] A. Wu, M. Shah, and N. Lobo. A Virtual 3D Blackboard: 3D Finger Tracking using a Single Camera. In *The fourth International Conference on Automatic Face and Gesture Recognition*, Grenoble, France, March 2000.
- [150] Y. Yacoob and M.J. Black. Parameterized Modeling and Recognition of Activities. In *International Conference on Computer Vision*, Bombay, India, 1998.
- [151] M. Yamada, K. Ebihara, and J. Ohya. A New Robust Real-time Method for Extracting Human Silhouettes from Color Images. In *International Conference on Automatic Face and Gesture Recognition*, 1998.
- [152] M. Yamamoto, T. Kondo, T. Yamagiwa, and K. Yamanaka. Skill Recognition. In *International Conference on Automatic Face and Gesture Recognition*, Nara, Japan, 1998.
- [153] M. Yamamoto and K. Koshikawa. Human Motion Analysis Based on A Robot Arm Model. In *Conference on Computer Vision and Pattern Recognition*, 1991.
- [154] M. Yamamoto, Y. Ohta, T. Yamagiwa, and K. Yagishita. Human Action Tracking Guided by Key-Frames. In *The fourth International Conference on Automatic Face and Gesture Recognition*, Grenoble, France, March 2000.
- [155] C. Yaniz, J. Rocha, and F. Perales. 3D Region Graph for Reconstruction of Human Motion. In *Workshop on Perception of Human Motion at ECCV*, 1998.
- [156] J.Y. Zheng and S. Suezaki. A Model Based Approach in Extracting and Generating Human Motion. In *International Conference on Pattern Recognition*, 1998.

