

Learning Silhouette Features for Control of Human Motion

LIU REN

Carnegie Mellon University

GREGORY SHAKHNAROVICH

Massachusetts Institute of Technology

JESSICA K. HODGINS

Carnegie Mellon University

HANSPETER PFISTER

Mitsubishi Electric Research Laboratories

and

PAUL VIOLA

Microsoft Research

We present a vision-based performance interface for controlling animated human characters. The system interactively combines information about the user's motion contained in silhouettes from three viewpoints with domain knowledge contained in a motion capture database to produce an animation of high quality. Such an interactive system might be useful for authoring, for teleconferencing, or as a control interface for a character in a game. In our implementation, the user performs in front of three video cameras; the resulting silhouettes are used to estimate his orientation and body configuration based on a set of discriminative local features. Those features are selected by a machine-learning algorithm during a preprocessing step. Sequences of motions that approximate the user's actions are extracted from the motion database and scaled in time to match the speed of the user's motion. We use swing dancing, a complex human motion, to demonstrate the effectiveness of our approach. We compare our results to those obtained with a set of global features, Hu moments, and ground truth measurements from a motion capture system.

Categories and Subject Descriptors: I.3.7 [**Computer Graphics**]: Three-Dimensional Graphics and Realism—*Animation*; H.5.2 [**Information Interfaces and Presentation**]: User Interfaces—*Input devices and strategies, interaction styles (e.g., commands, menus, forms, direct manipulation)*; I.2.6 [**Artificial Intelligence**]: Learning; I.4.8 [**Image Processing and Computer Vision**]: Scene Analysis—*Motion, tracking*

General Terms: Algorithm

Additional Key Words and Phrases: Performance animation, motion control, motion capture, animation interface, machine-learning, computer vision

Funding of the first and third authors was provided in part by NSF-IIS0326322, NSF-IIS0205224, and Mitsubishi Electric Research Laboratories.

The first author was also partially supported by a School of Computer Science Alumni Fellowship from Carnegie Mellon University. Author's addresses: L. Ren and J. K. Hodgins, Computer Science Dept., Carnegie Mellon University, 5000 Forbes Ave., Pittsburgh, PA 15213; email: {liuren,jkh}@cs.cmu.edu; G. Shakhnarovich, MIT CSAIL, 32 Vassar St., Cambridge, MA 02139; email: Gregory@csail.mit.edu; H. Pfister, MERL—Mitsubishi Electric Research Laboratories, 201 Broadway, Cambridge, MA 02139; email: pfister@merl.com; P. Viola, Microsoft, One Microsoft Way, Redmond, WA 98052; email: viola@microsoft.com.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or direct commercial advantage and that copies show this notice on the first page or initial screen of a display along with the full citation. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, to redistribute to lists, or to use any component of this work in other works requires prior specific permission and/or a fee. Permissions may be requested from Publications Dept., ACM, Inc., 1515 Broadway, New York, NY 10036 USA, fax: +1 (212) 869-0481, or permissions@acm.org.
© 2005 ACM 0730-0301/05/1000-1303 \$5.00

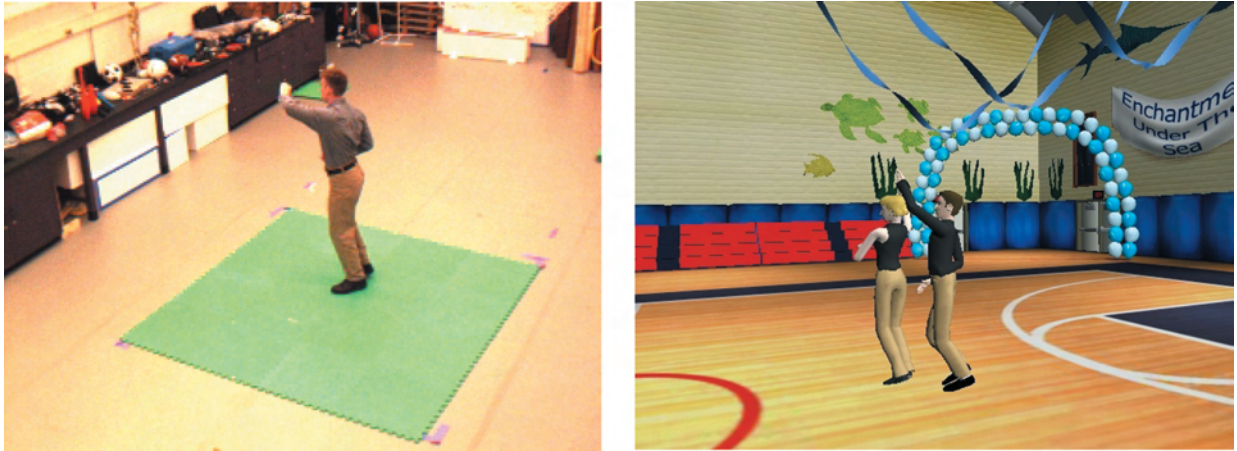


Fig. 1. A dancing user drives the motion of two animated characters.

1. INTRODUCTION

People, even without training, can move in expressive ways to play charades, pantomime a winning touchdown, or act out a story for a child. The goal of our research is to provide untrained users with a vision-based interface for interactively controlling complex human motions. Simple vision processing supplies partial information about the user's movements, and domain knowledge from a previously captured motion database supplements this information to allow plausible movement to be inferred. The process can be performed interactively, and the computation is real time with less than a second of latency between the capture of the video and the rendering of the animated motion.

The experimental setup in this article is based on a single user animating a pair of dancers performing East Coast Swing. We present a low-cost, camera-based system that allows a user to "drive" the dance interactively (Figure 1). The user can change the high-level structure and tempo of the dance, introducing a different sequence of steps, rotations, or movements. The final sequence maintains the quality of the previously captured dance while at the same time reflecting the user's higher-level goals and spontaneity. We envision that such a system might provide the infrastructure in teleconferencing applications, in an interactive game for a dancing child, in an animated karaoke performance for an adult, or, with a sufficiently broad database, in an animation authoring interface for a naive user.

We use an example-based approach that relies on silhouettes extracted from three video streams to match the body configuration of the user against a database of dance motion. This matching process uses a set of discriminative local features that are computed on the silhouette images [Viola and Jones 2001]. The best local features for estimating the yaw orientation and body configuration of the user are selected from a very large set of possible features by training on synthetic data set. The selection is based on a variant of the AdaBoost algorithm [Schapire and Singer 1999; Collins et al. 2000]. The estimation of yaw orientation incorporates a fast search method called *locality-sensitive hashing* (LSH) [Gionis et al. 1999]. The estimation of body configuration relies on the temporal coherence evidenced in a database of human motion. The database is preprocessed into an augmented motion graph that extends the original motion graph [Lee et al. 2002; Kovar et al. 2002] to handle scaling in velocity and two interacting characters. The estimation of the pose of the animated character is a local search in the augmented motion graph around previous pose. A small buffer of input video provides sufficient information for the match and introduces a delay of less than 1 s.

Our experimental results include several users dancing in street clothes and the resulting motions of an animated couple. For evaluation, we capture video and motion capture data simultaneously and compare the motion computed by our system with the ground-truth motion capture data. We also compare the performance of our system with that of a commonly used set of global features, Hu moments [Hu 1962], to demonstrate the effectiveness of the learned local features.

Alternative approaches to creating performance-based animation systems include vision-based tracking algorithms, small marker sets in concert with a database of human motion [Chai and Hodgins 2005], real-time motion capture systems, and such special purpose sensors as accelerometers or the pad used in Footsee [Yin and Pai 2003]. Our system has potential advantages over these approaches because it requires no instrumentation of the user (other than reasonably form-fitting clothing, and a controlled background), produces high-quality motion given an appropriate database, and requires only off-the-shelf cameras rather than special-purpose sensors or cameras. The major limitation of our approach is its reliance on a domain-specific database; however, with the increasing availability of high-quality motion capture data, we believe that this will not prove to be a major impediment for applications in which the behaviors of the user can be reasonably predicted in advance.

This article makes several technical contributions. First, we present a nonintrusive, interactive performance animation system that allows the generation of high-quality and complex human motion from only three video inputs. Unlike other approaches, we use local features that operate only on small, or local, regions of the silhouette image to estimate yaw orientation and body configuration. The key insight in our approach is that the information contained in a motion database can be used to learn the discriminative local features for a particular task. Second, for fast silhouette-based yaw estimation, we present a novel solution that combines AdaBoost and LSH. Finally, we propose an augmented motion graph structure that smoothes the estimated motion, handles variation in the speed of the user's performance and synthesizes the motion of two human characters from the performance of one user.

2. BACKGROUND

We build on research in recovering motion from video, in creating animations from motion capture data, and on performance animation systems. In this section, we review the most relevant research from each of these three areas.

2.1 Motion Recovery from Video

Almost all computer vision approaches rely on tracking to analyze human motion from video input (see Gavrilu [1999] for a survey). Three-dimensional (3D) human motion tracking is most relevant to our work [Yamamoto et al. 1998; Bregler and Malik 1998; Delamarre and Faugeras 1999; Sminchisescu and Triggs 2003]. In general, these systems assume accurate initialization and then track changes in pose based on an articulated 3D human model. For example, Bregler and Malik [1998] used an exponential map formalism for kinematic chains to track simplified human body models. The tracking approach of Delamarre and Faugeras [1999] is closely related to our own in setup: three synchronized cameras were used to compute silhouettes. In their approach, a synthetic articulated model of the human form was first rendered online and the synthetic silhouettes were then compared with the real ones. The difference between the silhouettes was used to generate forces, which were then applied to each part of the 3D articulated model for a better match with the real silhouettes. Most 3D tracking approaches require time-consuming per frame nonlinear optimization, which is not appropriate for such real-time or interactive applications as performance animation systems. Numerical drift is also a significant problem in these approaches, even with manual initialization. As a result, few articles describe systems that can track complex motions or motion with significant yaw rotation for more than a few seconds.

Applications that require tracking only coarse poses can be performed more robustly. Ramanan and Forsyth [2004] described a system that allowed multiple people to be tracked for motion annotation. A robust algorithm for 2D pose tracking was used initially, and the resulting 2D poses were then lifted to a sequence of coarse 3D poses. This approach proved quite effective for classifying such human movements as walking, running, and jumping in an indoor environment.

Alternatively, a robust vision-based interface can directly estimate the 3D body pose without initialization and incremental updates, thus avoiding problems with initialization and drift [Rosales et al. 2001; Mori and Malik 2002; Shakhnarovich et al. 2003]. For example, Rosales et al. [2001] trained a neural network to map each 2D silhouette to 2D positions of body joints and then applied an EM algorithm to reconstruct a 3D body pose based on 2D body configurations from multiple views. Mori and Malik [2002] explored an example-based approach to recover 2D joint positions by matching extracted image features (shape context features) with those of cached examples and then directly estimating the 3D body pose from the 2D configuration considering the foreshortening of each body segment in the image. More recently, Shakhnarovich et al. [2003] combined multiple cues (edges and intensities) and developed an extension of locality sensitive hashing to rapidly estimate the upper body poses of human figures from a large database of example images of figures viewed from the front. Efros et al. [2003] proposed a spatiotemporal motion descriptor based on noisy optical flow measurements for pose estimation from one video camera. Their framework was designed to recognize human actions from distant views. Despite these advances, direct estimation of precise 3D body pose from 2D noisy image features remains a difficult problem. Such systems sometimes yield coarse, imprecise, or incomplete estimations of body pose. We improve on this situation by utilizing the spatial-temporal structure of prerecorded human motion and learning efficient and effective image features from the information in the database.

We extract silhouettes from the video data and use them to drive the character's motion. Silhouettes have been widely used for motion estimation from video data [Brand 1999; Delamarre and Faugeras 1999; Cheung et al. 2000; Rosales et al. 2001; Mikic et al. 2001; Lee et al. 2002; Carranza et al. 2003]. Estimating motions from a single camera is challenging because a single silhouette is not enough to disambiguate the 3D pose. Brand [1999] inferred a 3D motion sequence that was most consistent with a series of input silhouettes using a hidden Markov model that was estimated from 3D motion capture data. Due to the sparsity of input data, their system produced motion that was not an accurate reproduction of human motion and was therefore not suitable for animation. To control an avatar, Lee et al. [2002] explored a single-camera vision interface and retrieved the 3D motion data from the database by extracting silhouettes and comparing global features (Hu moments). Simple global features are easy to extract, but they are not discriminative enough to catch the subtle changes in the silhouette that occur during complex human motions. This limitation can lead to imprecise pose matching as discussed by Rosales et al. [2001]. Lee et al. [2002] introduced a three second delay to reduce the problems caused by both the imprecise matching metric and pose ambiguities from a one camera interface. Although effective for simple motions such as those they tested (walking over and onto a step stool and sitting on it), their approach does not generalize to complex human motions such as dancing.

Silhouettes from multiple cameras can be used to reduce the ambiguity problem in pose estimation. Such systems usually estimate motion via polygonal or volumetric modeling of the observed shape, but the additional cameras increase the complexity of modeling and the computational power required [Cheung et al. 2000; Matusik et al. 2001; Mikic et al. 2001]. For accurate 3D reconstruction of body shape, high-quality video cameras are often required. Carranza et al. [2003] incorporated domain knowledge in the form of a human body model; but their approach was computationally expensive and therefore unsuitable for interactive applications.

Domain knowledge, especially that embedded in a human motion database, has also been used by researchers to estimate 3D human motion from video input [Leventon and Freeman 1998; Deutscher et al. 2000; Brand and Hertzmann 2000; Sidenbladh et al. 2002; Lee et al. 2002; Stenger et al. 2003]. Leventon and Freeman [1998] applied dimensionality reduction to 3D motion capture data and then used a Bayesian approach to synthesize a 3D motion sequence from tracked 2D marker positions in a monocular image sequence. Similarly, Sidenbladh et al. [2002] learned a low-dimensional linear model of 3D human motion from a database and presented a probabilistic tree search method for efficient synthesis and tracking. More recently, Stenger et al. [2003] proposed a probabilistic method for tracking. Their method searched a database that was hierarchically partitioned into areas with constant density. These approaches as well as the work reported here demonstrate that the domain knowledge inherent in a human motion database allows efficient motion synthesis from video data.

2.2 Resequencing Motion Capture Data

Resequencing motion capture data is a very successful method for synthesizing new animations and controlling animated characters [Lee et al. 2002; Kovar et al. 2002; Arikan and Forsyth 2002; Arikan et al. 2003]. In this approach, novel motions are created by selecting motion clips from a database in the form of a graph and reassembling them to generate a new motion that approximately follows a path specified by the animator or that matches a set of constraints specified by the animator. This approach produces compelling animation in part because the motion that is selected is not modified and therefore retains the subtle details of the original motion data. We build on this work by using the estimated pose of the user to find similar sequences of motion in a modified motion graph. We extend the motion graph concept to handle two interacting subjects and allow time scaling of the motion sequences to better match the user's performance.

2.3 Performance-Based Animation

Performance-based animation systems have been used quite successfully for characters that respond in real time to the actions of human actors [Granieri et al. 1995; Noser and Thalmann 1997; Shin et al. 2001]. Although problems with retargetting, filtering, and removing collisions remain, all commercially available technologies for motion capture now have the ability to perform real-time capture of a human actor. While these systems are very effective for choreographed performances, the time required to attach markers or sensors prevents employing these systems in applications that will be used by many different people. Like other vision-based systems, our approach requires only a controlled background and reasonably form-fitting clothing and does not require the additional overhead of suiting up the user with markers or sensors.

Another approach to performance-based animation is to use a sensing technology that is less obtrusive than markers or magnetic sensors. Yin and Pai [2003] used a foot pressure sensor pad to control avatars interactively. The pressure patterns for the foot-to-ground contacts of the user were correlated to a motion capture database with matching footpad information to reconstruct full body motion. This interface provides appealing results for a number of kicking and stepping motions but subtle changes in the upper body motion may not result in unique foot pressure patterns.

3. EXAMPLE-BASED APPROACH

As described in Figure 2, our vision-based performance interface allows a user to control an animated human character. The system extracts silhouettes from the input of three cameras. The vision-based performance interface estimates the dancer's motion from these silhouettes. Estimating silhouette-based motion is a difficult problem because the mapping from 2D silhouettes to 3D body configurations is ambiguous. In the work reported here, we used an example-based framework to solve this problem

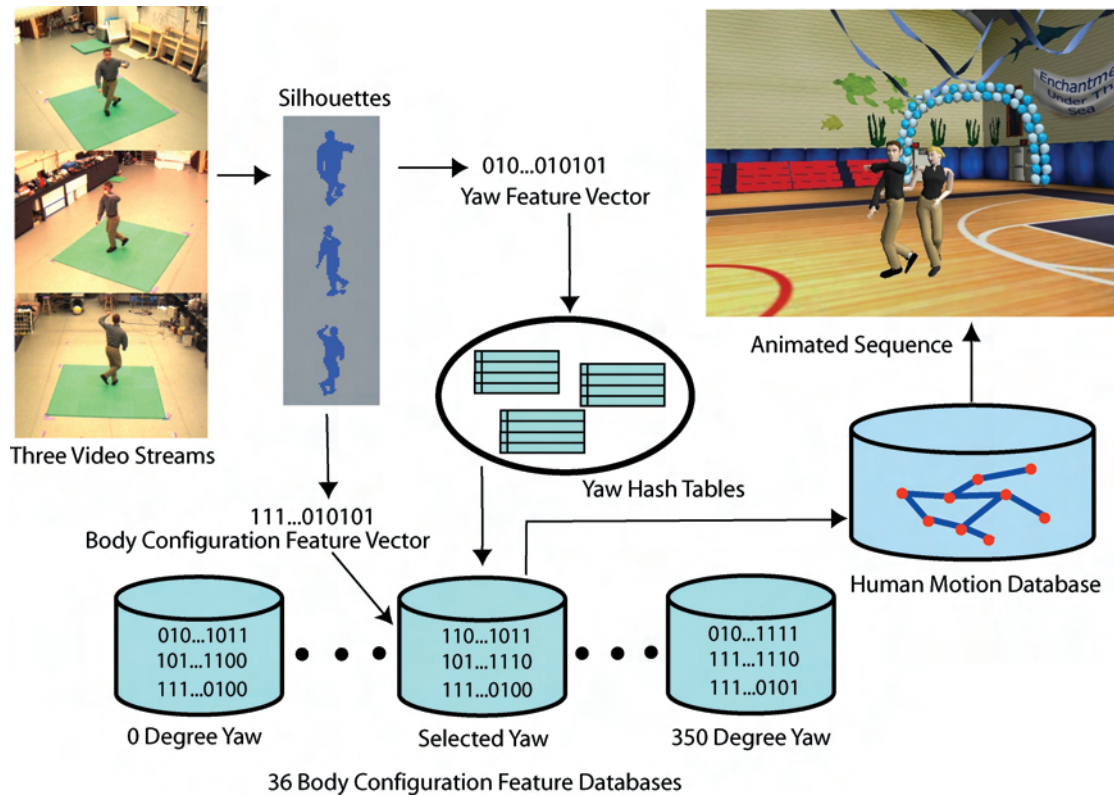


Fig. 2. System diagram. In the preprocessing step, we construct one database for yaw estimation and 36 databases for body configuration estimation. The examples, the binary feature vectors extracted from the synthetic silhouettes and the associated pose parameters are stored in those databases. At runtime, the three video streams are processed to silhouettes. Discriminative local features are extracted from those silhouettes and form binary feature vectors for the estimation of yaw and body configuration. The system first uses the yaw feature vector to search for a similar yaw orientation through a hashing algorithm. The system then searches for a similar body configuration from the database of the estimated yaw angle. The animation of the user and his partner are created from the poses found in the database.

by relying on domain knowledge in the motion capture database. In this section, we first present an overview of our approach and then discuss the major challenges and the proposed solutions.

In our example-based framework, captured motion is stored in a database along with the corresponding synthetic silhouettes. The body configuration for a new silhouette from the cameras is estimated by searching the database for a similar silhouette and retrieving the corresponding body configuration as well as changes in global position and orientation. Despite the simplicity of the example-based approach, there are a few challenges in using this approach for motion estimation from silhouettes:

—*Example storage.* Because the database determines the set of allowable body configurations, it must cover the space of required motions for the application domain. For general human motion, an example-based approach might require a prohibitively large database. The required storage can be reduced by building a behavior-specific motion database of the 3D motion data and 2D silhouette data. We build a database that covers only the basic movements of a specific dance, East Coast Swing. We represent each silhouette as a precomputed bit string based on an efficient binary local feature set [Viola and Jones 2001]. In this format, the silhouettes only require 34 MB of storage.

- False matches.* Example-based motion estimation from silhouettes requires that similar silhouettes correspond to similar poses; however, silhouettes are ambiguous because two quite different poses may create similar silhouettes. Large-scale ambiguity in the silhouettes is caused by the 3D-to-2D projection and mirror symmetry. Small-scale ambiguity also exists because a small change on the boundary of the silhouette (for example, the arm) can imply a large change in the limb position (for example, the hand) that is located inside the silhouette boundary. As a result, two silhouettes may appear similar even if their underlying poses are quite different. Such ambiguities can cause false matches in which the nearest silhouette in the database does *not* correspond to the nearest pose. We reduce the ambiguity by using three cameras and constructing silhouette similarity measures that effectively capture changes in local areas of the silhouette image by selecting good local features from a large set during a training step.
- Robustness.* Silhouettes can be noisy because the separation between background and foreground is imprecise. Our silhouette similarity measures are robust to noisy inputs because each local feature used in the similarity measure is computed based only on the statistics of a local patch and is independent of other parts of the silhouette.
- Computational efficiency.* Example-based estimation requires frequent comparisons of silhouettes with similarity measures and this process must be computationally efficient for online applications. The rectangular local features we chose are efficient [Viola and Jones 2001] and provide a small bit string representation of each silhouette. The silhouette similarity measure is then the Hamming distance between two bit strings, which is also computationally efficient to compute.
- Search efficiency.* Online example-based estimation requires an efficient K-Nearest Neighbor (KNN) search. Real-time performance for a large silhouette database is impractical for current computers because the complexity of naive KNN search is a linear function of the size of the database. We use two methods to improve the search efficiency. We rely on *locality-sensitive hashing* (LSH), an approximate nearest-neighbor search, to speed up the KNN search process for yaw estimation. We then use a local search to retrieve a close body configuration from an *augmented motion graph*, a graph structure whose edges connect body configurations that are nearby in body configuration space to reduce the search space.

We now describe the details of this framework and the results obtained with our system. We first describe the video and motion data processing (Section 4). We then discuss learning good similarity measures for the example-based estimation (Section 5) obtained through an offline training process (Section 6). Next we describe the online motion control in Section 7. In Section 8, we compare the results obtained with our approach to those obtained with a widely used global feature, Hu moments, and to a ground truth measurement from a motion capture system. Finally, we discuss some of the engineering choices we made in the design of our system and its limitations (Section 9).

4. VIDEO AND MOTION DATA PROCESSING

In our real-time video capture system, three Dragonfly cameras [Point Grey Corporation 2001] were connected to a single PC through a IEEE-1394 firewire interface. Synchronized video frames were captured at 15 frames per second (fps) with a resolution of 640×480 . Figure 3 shows the layout of the system. We built a workspace coordinate system with software for multiple camera calibration and determined the position of the real and the corresponding synthetic cameras [California Institute of Technology 2002]. This information was used for rendering synthetic silhouettes for training.

The video frames were processed to form real silhouettes. We assumed a static background and fixed lighting to simplify the extraction of the silhouettes. Shadows often cause errors in silhouettes. We

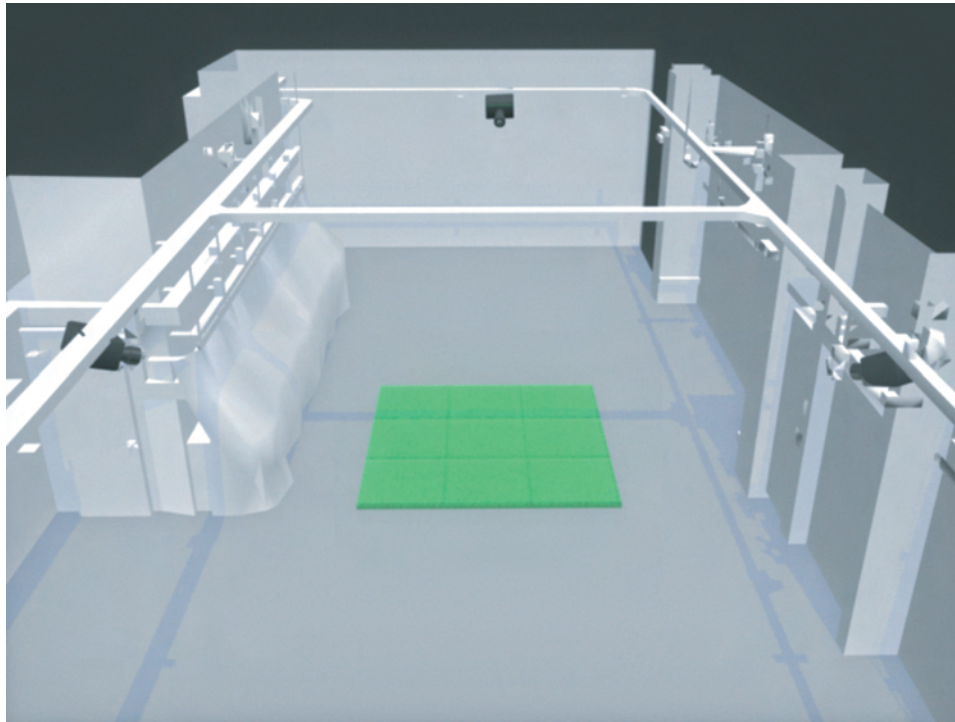


Fig. 3. Our system has three Dragonfly cameras (black) that point downward to capture an $8' \times 8'$ workspace (green). The three cameras form an approximate isosceles triangle with a smallest acute angle about 30 degrees. The distance between each camera and the workspace should be large enough so that an affine camera model provides a reasonable approximation.

avoided this problem with a simple two-step method. We first used a low threshold to extract an initial silhouette, which may include regions of shadow. Then we computed the zero and first order spatial moments to determine the silhouette center and recompute the foreground pixels below the center using a higher threshold. This two-step process eliminated most shadows near the feet but did not mislabel the foreground pixels of the face and arms. An alternative method would be to use an angular difference metric for the hue channel of pixel color, but this approach would be significantly more computationally expensive [Cheung et al. 2000].

As the user danced, he moved around the capture region and his silhouette changed size in the image. We compensated for this by cropping and scaling the silhouettes from the three views (Figure 4). The resulting silhouettes were concatenated and treated as one image in the subsequent analysis. For simplicity, we use the term *silhouette* to refer to this image.

We built a human motion database by capturing the motion of two subjects dancing together. We used a Vicon optical motion capture system [Vicon Motion Systems 2002] with 12 cameras, each of which can record at 120 Hz with images of 1000×1000 resolution. We used a standard 41 marker setup for each subject. The motion data were down-sampled to 15 fps to match the frame rate of our video cameras. The size of the motion capture area was $8' \times 24'$. The motion database contained approximately 6 min of a couple swing dancing (East Coast Swing). The motion data were processed into a hierarchical skeleton (joint angles and pelvis position and orientation for a total of 62 degrees of freedom). Contact constraints were not enforced. As a result, the fixed or sliding contacts made by the hands and arms of the dancing couple were not always realistic.

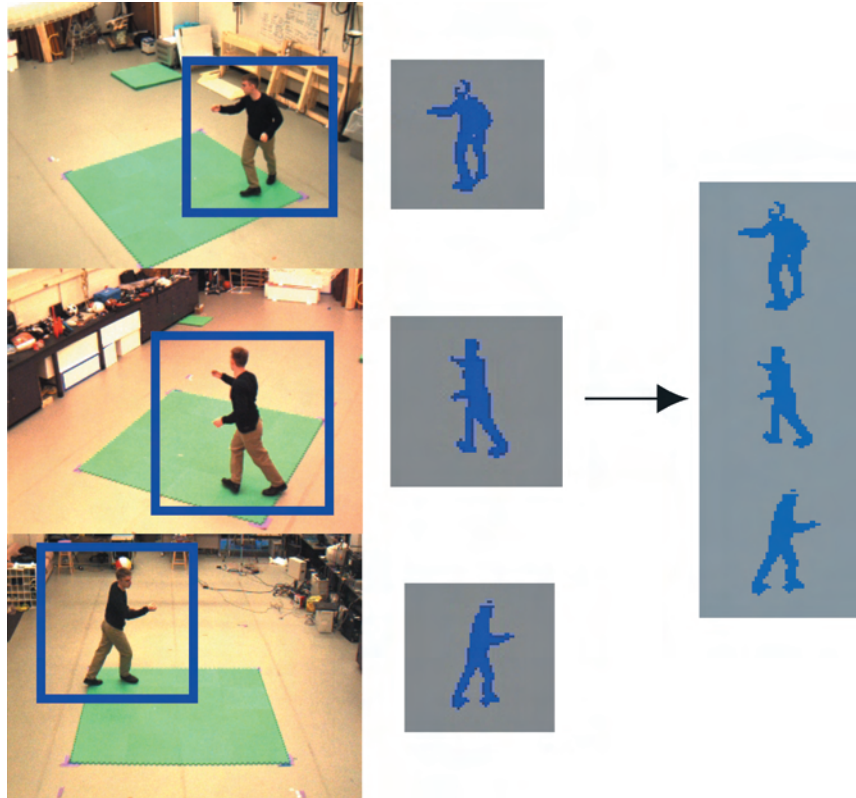


Fig. 4. The silhouette is extracted from each video input and then cropped. To crop a silhouette, the system computes the first-order moment and the second-order central moment from the horizontal and vertical distribution of foreground pixels. The centroid of the foreground pixels and their standard deviation are then calculated from those image moments. We use the center and 2.6 times the largest standard deviation to determine the square bounding box for the cropping. Finally, the cropped silhouette is resampled to a 60×60 pixel window using bilinear interpolation and the three silhouettes are stacked as a combined representation for further analysis. The silhouettes are not perfect, because the background is not controlled. In particular, when the user's head passes in front of a dark object, the silhouette of the head is often incomplete as is shown in the top silhouette of this example. Similar problems occur when the color of the user's clothing is similar to that of the background. Shadows sometimes affect the shape of the feet in the silhouette.

Recorded human motion is employed in several different ways in the system. The motion data are used to create synthetic examples for training by rendering a human model and extracting a silhouette from the rendered image. A motion graph is created from the motion database for synthesis of the motion of the dancing couple. Finally, we performed motion capture and video capture simultaneously and motion capture data that were not included in the database served as ground truth to assess the performance of the system.

5. LEARNING SILHOUETTE SIMILARITY MEASURES

Silhouette similarity measures are used to retrieve examples from the motion database that match the current pose of the user. The matches are found based on silhouettes from the cameras and synthetic silhouettes computed from the motion capture database. The similarity measures must minimize false matches and be computationally efficient. In this section, we describe a machine-learning approach to

constructing such similarity measures. The learning algorithm selects a set of efficient binary features such that the Hamming distance between the feature vectors computed from two silhouettes is small when the underlying poses of the silhouettes are sufficiently close.

5.1 The Learning Task

Similarity relationship in the pose space is captured by the binary neighborhood function:

$$N_\theta(\pi_1, \pi_2) = 1 \text{ if } |\pi_1 - \pi_2| < \theta, \text{ and } -1 \text{ otherwise.} \quad (1)$$

That is, N_θ labels a pair of poses, π_1, π_2 , as positive if the two poses are within θ from each other under a distance measure $|\cdot|$ in pose space, and as negative otherwise. Our motion database consists of examples of poses with corresponding synthetic silhouettes (π_i, S_i) . For pairs of such examples, we can evaluate N_θ on the poses directly, thus producing “ground truth” similarity labels. The same labels are assigned to the silhouette pairs of the corresponding poses: (S_i, S_j) is positive if and only if (π_i, π_j) is also positive.

At runtime our system must infer the pose from silhouettes only with a silhouette similarity measure F . Applied to a pair of silhouettes S_1 and S_2 , $F(S_1, S_2)$ infers the value of N_θ for the corresponding poses π_1 and π_2 : large positive values of $F(S_1, S_2)$ imply $N_\theta(\pi_1, \pi_2) = 1$. The learning task is to construct an F which is as much in agreement with N_θ as possible.

We first consider a similarity measure based on a single binary valued feature h_j : the Hamming distance between the values of $h_j(S_1)$ and $h_j(S_2)$. Zero distance implies that the corresponding poses π_1, π_2 form a positive pair (i.e., are similar); and a distance of one implies that the poses form a negative pair. Thus, this similarity measure is associated with a *silhouette pair classifier* H_j :

$$H_j(S_1, S_2) = \begin{cases} +1 & \text{if } h_j(S_1) = h_j(S_2), \\ -1 & \text{if } h_j(S_1) \neq h_j(S_2). \end{cases} \quad (2)$$

The accuracy of this similarity measure is equivalent to the classification accuracy of H_j with respect to the labels assigned by N_θ . This accuracy may be empirically evaluated for a given h_j using a set of pairs with known labels.

In general, given a restrictive set of binary features, no single H_j can approximate the neighborhood function with satisfactory accuracy. Approximation of N_θ may be achieved by combining many silhouette pair classifiers:

$$F(S_1, S_2) = \sum_{j=1}^M H_j(S_1, S_2). \quad (3)$$

Let h_1, \dots, h_M be M binary features; these features define an encoding of a silhouette S into a binary vector $V(S)$ of length M , in which the j th bit is $V(S)[j] = h_j(S)$. The Hamming distance between encodings of two silhouettes is correlated with the measure in Equation (3).

The sign of $F(S_1, S_2)$ defines a binary classifier on silhouette pairs, which is a combination of M simpler classifiers. In the machine-learning literature, such a classifier is known as an *ensemble classifier*. The accuracy of the Hamming distance in the space of $V(S)$ as a similarity measure is equivalent to the classification accuracy of F :

$$\sum_{i,j} N_\theta(\pi_i, \pi_j) \text{sign}(F(S_i, S_j)), \quad (4)$$

and can be assessed empirically based on a training set of labeled silhouette pairs. As a result, the goal of the learning algorithm is to select a set of about 300 binary features h_j from a very large set (more than one million) so that the resulting ensemble classifier F can maximize the quantity in Equation (4).

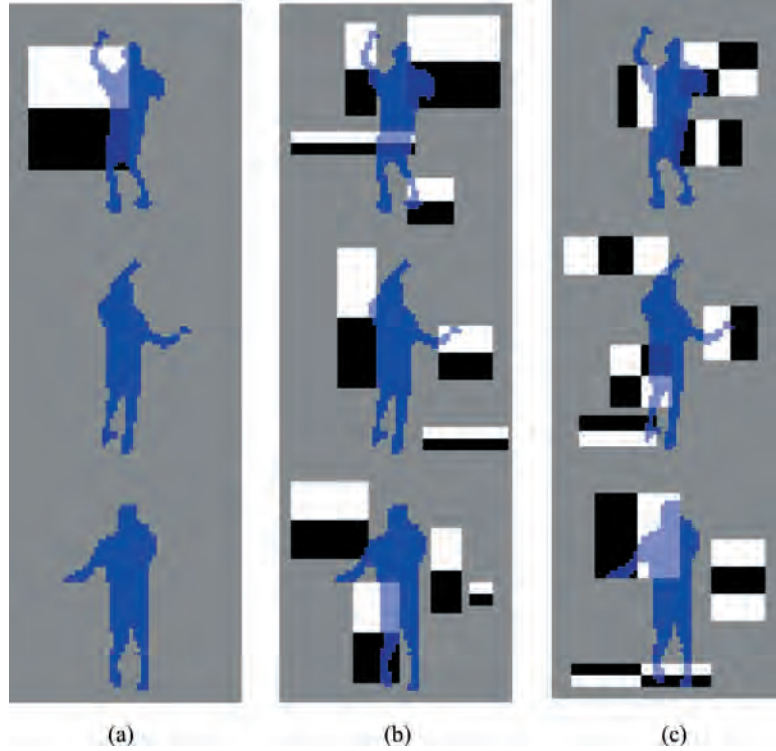


Fig. 5. (a) A rectangular filter on a binary silhouette image. Two rectangles are used in the filter. We compute the value of the filter by summing the intensities of all pixels in the white regions and subtracting the sum of the intensities of all pixels in the black regions. (b) The same type of filter can be defined with a variety of scales, aspect ratios, and locations. (c) Other types of filters used in our system have different orientations and more rectangles. The potential filter set is very large although it can be pruned by considering only those filters that overlap any silhouette in the training set.

The simple binary local feature h_j can be implemented in many different ways. In this work, we use the one originally proposed by Viola and Jones [2001]. These features have been successfully used in such real-time object detection problems as face detection and pedestrian detection [Viola and Jones 2001; Viola et al. 2003]. They are defined as

$$h_j(S) = \begin{cases} 1 & \text{if } \phi_j(S) > T_j, \\ 0 & \text{otherwise,} \end{cases} \quad (5)$$

where T_j is a threshold and a *rectangular filter* $\phi_j(S)$ is a scalar linear function of the pixels in a silhouette image S . Rectangular filters are displayed graphically in Figure 5 as rectangles with white and black regions. By using summed area tables [Crow 1984] or integral images, we can compute the value of each rectangular filter in constant time independent of its scale and location [Viola and Jones 2001]. The number of available features is very large because the available filter set itself is very large and there are many choices for the thresholds. We use a machine-learning process to select a subset of M filters and a best threshold for each.

5.2 Feature Selection Using AdaBoost

An ensemble classifier can be efficiently trained with a variant of the AdaBoost algorithm [Schapire and Singer 1999; Collins et al. 2000] adapted to classify example pairs instead of single examples. In the work of Viola and Jones [2001] a collection of binary features was used to classify each example image into one of two classes (“face” or “not face”). We adopt the form of the pair classifier used by Shakhnarovich and colleagues [2003] to classify example pairs as similar or not in the AdaBoost framework.

AdaBoost proceeds in rounds, adding one feature at a time. Initially each training example pair is assigned a weight of $1/C$, where C is the number of example pairs. In each round, the features are computed for each example pair, and each feature is assigned an error based on the weights of the example pairs that it classified incorrectly. AdaBoost selects the feature with the lowest weighted error. Example pairs that are correctly classified by the new ensemble classifier receive a lower weight while the weights of incorrectly classified example pairs are increased. These new weights are then used to select another feature in the next round. The weight for each selected feature $h_j(S)$ is computed based on the weighted error made by that feature: lower errors receive higher weights.

The best binary feature of each round is determined by selecting not only a filter but also a threshold for that filter in each round. The cost for a threshold T of a filter ϕ has two terms. The first is the number of example pairs that are near in pose space but split by the threshold:

$$N_\theta(\pi_i, \pi_j) = 1 \text{ but } \phi(S_i) < T < \phi(S_j) \text{ or } \phi(S_j) < T < \phi(S_i). \quad (6)$$

The second is the number of examples which are distinct in pose space but on the same side of the threshold:

$$N_\theta(\pi_i, \pi_j) = -1 \text{ but either } \phi(S_i) \leq T \text{ and } \phi(S_j) \leq T \text{ or } \phi(S_j) > T \text{ and } \phi(S_i) > T. \quad (7)$$

For a given filter ϕ , the optimal T can be computed by sorting the examples by $\phi(S_i)$ and then traversing the sorted list while recording pairs that are erroneously split or combined by the threshold. We select the T with the smallest number of errors as the optimal threshold for filter ϕ .

Our feature selection algorithm can be summarized as follows: construct a training set of pairs of silhouettes, labeled by N_θ , and run AdaBoost for M iterations yielding M binary features. In each iteration, we determine the optimal threshold T_i for each selected filter ϕ_i , thus defining a pair classifier H_i . AdaBoost then selects the H_i with the lowest weighted error, reweights the examples (pairs), and proceeds to the next iteration.

AdaBoost produces an ensemble classifier which is a weighted combination of pair classifiers. However, for efficiency we do not weight the selected features during retrieval.

6. TRAINING

We use synthetic silhouettes rendered from the motion capture database to create training examples. The selection of an appropriate set of training examples is crucial for creating estimators that are robust to changes in yaw orientation (rotation of the pelvis about the vertical axis) and translation as the user spins and steps around in the workspace. The potential training set therefore includes each body configuration in the database rendered in five locations in the workspace with 36 different yaw orientations (Figure 6). These transformations create a potential training set of 921,600 silhouettes from 5120 frames of motion data.

Pairs of these silhouettes could, in principle, be used to learn only one similarity measure for pose that includes both body configuration and orientation. However, this approach appears to be infeasible due to the wide variety of silhouettes that this similarity measure would have to handle. We choose instead to separate the estimation of yaw orientation from that of the internal body configuration. We construct one yaw estimator and 36 rotation-specific body configuration estimators, one for each 10 degrees of yaw

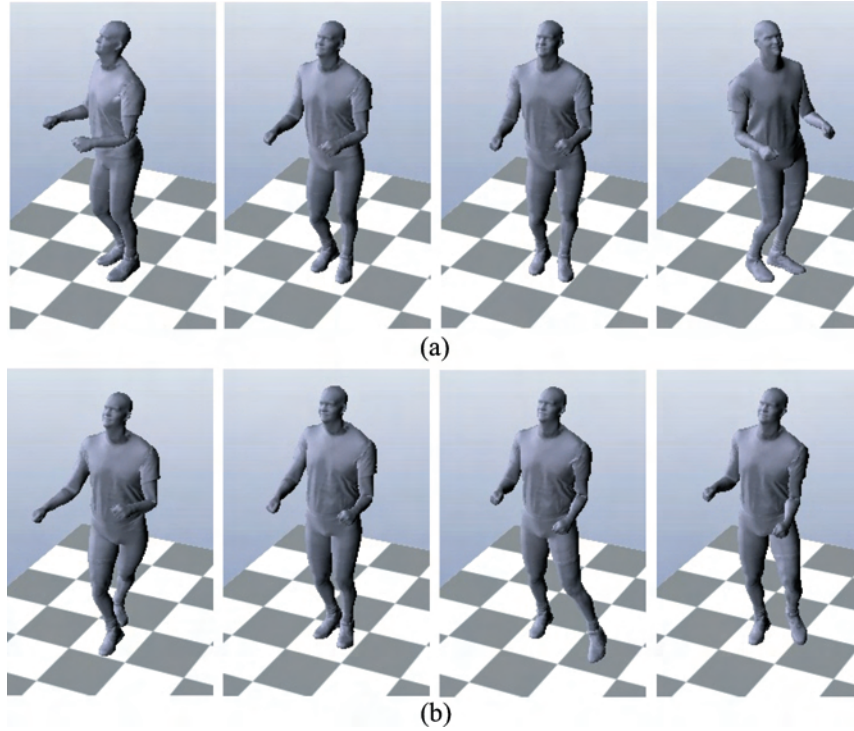


Fig. 6. (a) One-body configuration rendered from different yaw orientations. (b) Different body configurations at the same yaw orientation. In (a) and (b), all body configurations are rendered at the center of the virtual workspace and the checkerboard is the same size as the real workspace. We also render each body configuration with fixed yaw orientation at four corners of the virtual workspace because the affine camera model assumes similar silhouette shapes for a given body configuration and yaw orientation throughout the workspace. We generate the synthetic silhouettes from images rendered from three virtual cameras. Because the synthetic silhouettes must be rendered using virtual cameras at fixed positions, significant movement of the cameras in the real setting would require recalibrating the cameras, regenerating these synthetic silhouettes, and repeating the selection process for the features.

rotation. Each body configuration estimator is then trained to distinguish body configurations within a narrow range of yaw angles (± 5 degrees), which is a much easier problem than determining body configuration independent of yaw angle. To find the complete body pose, the system first estimates yaw orientation and then uses the appropriate body configuration estimator for that yaw angle to determine the rest of the pose. We describe the similarity measure for yaw estimation in the next section and then discuss the construction of body configuration similarity measures in Section 6.2.

6.1 Similarity Measure for Yaw Estimation

Yaw estimation is a multilabel classification task with 36 classes each corresponding to a 10 degree bin centered at 0, 10, 20, \dots , 350 degrees. A pair of silhouettes is labeled positive if their yaw angles fall into the same bin. The body configuration is not taken into consideration. See Figure 7 for examples.

Because we use pairs of silhouettes for training, the set of 921,600 synthetic silhouettes would provide more than 23 billion positive example pairs and three orders of magnitude more negative example pairs. We use only a small, representative set of body configurations as training examples. We find the set by using the k -means algorithm [Duda et al. 2000] to find 50 clusters in the set of body configurations and then selecting the body configuration closest to the center of each cluster. The distance metric for

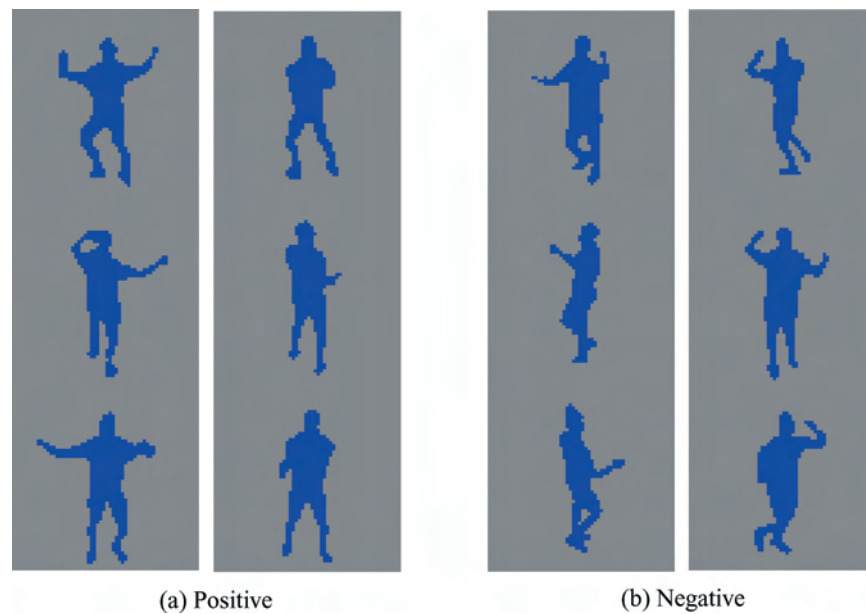


Fig. 7. Examples of training pairs used to select yaw similarity features. (a) Same yaw angle but different body configurations. (b) Different yaw angles but same body configuration.



Fig. 8. The body configuration parameters are a set of joint positions in the root segment coordinate system of the human body. They are concatenated to form a vector in high dimension space. The L2 distance metric is used in the clustering.

clustering is the L2 distance between two body configuration parameters (Figure 8). For each of these 50 body configurations, we generate a multiview silhouette of a synthetic human model rendered at a random orientation within each of 36 yaw bins (yaw of $0, 10, 20, \dots, 350 \pm 5$ degrees) and at each of five locations (the center and the corners of the workspace). Because AdaBoost involves computing

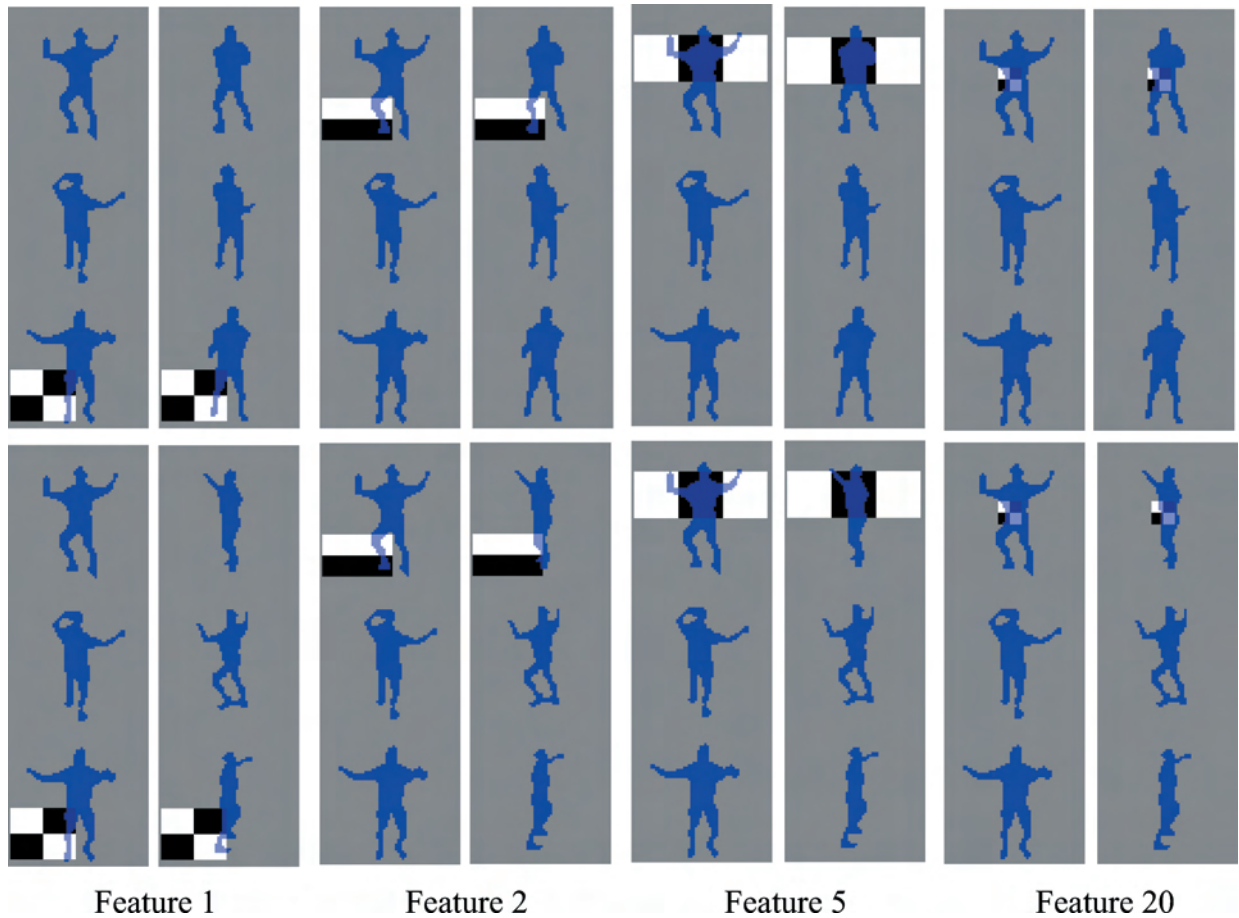


Fig. 9. The figure shows examples of features that were selected by AdaBoost for the yaw similarity measure. Each feature is also associated with a threshold. The top row shows several features for a positive training example. The second row shows features for a negative training example.

the value of each filter for every training example, the resulting set of pairs $(50 \times 36 \times 5)^2$ is still too large. We randomly select 4000 positive and 6000 negative training pairs from this set and use these examples in AdaBoost learning.

Given the richness of the moves in swing dancing, a training set of 10,000 is unlikely to span the space of silhouettes that will be seen when the user dances. We address this problem by *resampling* the training data set [Shipp and Kuncheva 2002; Jones and Viola 2003] after a number of rounds. We compute the AdaBoost weights on a much larger set of examples based on the classification by the set of binary features selected so far, then select a new training set according to these weights, taking care to maintain a balanced representation of positive and negative examples. See Figure 9 for examples of features selected by AdaBoost for the yaw similarity measure.

6.2 Body Configuration Similarity

Thirty-six different silhouette similarity measures for body configuration are learned, one for each of the 36 bins for yaw rotation. Using the synthetic model, we render the 5,120 body configurations at five

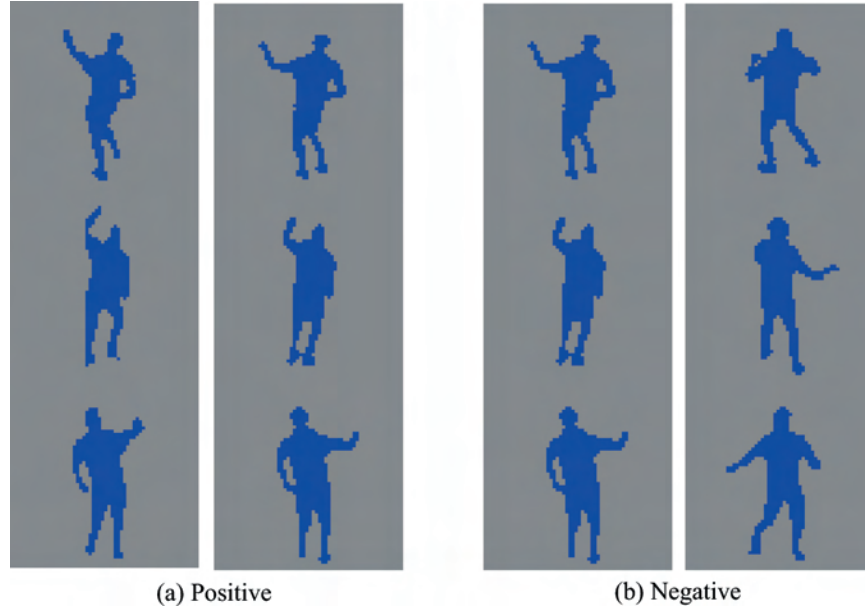


Fig. 10. Examples of training pairs for body configuration similarity (yaw is 90 degrees).

locations, yielding 25,600 examples from which the training pairs for each bin can be constructed. The true label for a pair of poses is determined by a threshold on the L2 distance between the two sets of body configuration parameters (the threshold was set by inspection). The distance metric is the same as that used for k -means clustering. Examples of a positive and a negative training pair are shown in Figure 10.

The number of positive pairs for body configuration is much lower than it is for yaw, and we use clustering to reduce the number of negative pairs. The binary similarity functions are selected by AdaBoost as they were for yaw orientation.

7. ONLINE MOTION CONTROL

For online motion control, the system estimates the yaw orientation of the human body from the silhouette. Then it finds a close match for the body configuration by searching the database using the feature set selected for that yaw orientation. Those two estimation processes must be fast and must result in a smooth motion sequence. We first discuss how we preprocess the data to improve the performance for online motion control (Section 7.1). Then we explain the details of the robust online yaw estimation (Section 7.2) and the body configuration estimation (Section 7.3).

7.1 Preprocessing

Because the evaluation of the similarity function involves comparing only feature vectors, the feature vectors for the silhouette images in the database can be precomputed. Each feature vector contains about 300 bits and requires far less space to store than the full silhouette.

Given an input silhouette I from the video cameras and a yaw similarity measure of M features, the yaw estimation can be performed as follows. We first compute the M binary features from the silhouette image I and concatenate them to create a *bit vector* $[h_1(I) \cdots h_M(I)]$. The K examples most similar to the input are those with the smallest Hamming distance between their bit vector representations and

that of the input. The yaw angle can be estimated as the majority vote among these top K matches. Retrieving K similar examples is a computationally challenging problem that involves the search of more than 900,000 high-dimensional bit vectors. For efficiency, we preprocess the data and use a fast algorithm for approximate similarity search: locality-sensitive hashing (LSH) [Gionis et al. 1999]. LSH is described in the next section.

Given an input silhouette and an estimated yaw angle, the system then uses the selected features in the appropriate yaw bin to compute a body configuration feature vector and to find a close pose in the motion capture database. If performed globally, this search would not necessarily result in a smooth animation sequence. Instead, the system searches locally near the previously selected pose in an *augmented motion graph*. The augmented motion graph is constructed as a preprocessing step and includes two extensions to the original single-character motion graph structure [Lee et al. 2002; Kovar et al. 2002]. One extension allows the animation of two interacting characters and the other permits speed changes in the motion playback to accommodate tempo variations in the music and dancing motion. These extensions are described in detail in Section 7.1.2.

7.1.1 Locality-Sensitive Hashing. LSH proceeds by randomly selecting k functions among those features chosen by AdaBoost, thus defining a k -bit hash function:

$$g(x) = [h^{(1)}(x), h^{(2)}(x), \dots, h^{(k)}(x)]. \quad (8)$$

The entire database is indexed by a hash table with 2^k buckets: a silhouette I is associated with the bucket $g(I)$. In order to increase the probability of success (finding examples similar to the input), LSH constructs l hash tables with independently constructed keys g_1, \dots, g_l . Then, given an input I , the set of candidate matches consisting of the union of the buckets $g_1(I), \dots, g_l(I)$ is exhaustively searched for examples similar to I .

The parameters k (the number of bits per hash key) and l (the number of hash tables) affect the probability of finding a good match for the input and the expected search time. A large k makes spurious matches less likely and thus speeds up the search but also increases the probability of missing good matches. A large l increases the probability of success but also increases the cost of the exhaustive search. We set these parameters based on an empirical evaluation that uses two disjoint sets of silhouettes with known yaw: an example set and a test set. For k in a feasible range, we build a hash table from an example set with a k -bit hash key randomly constructed according to Equation (8). Then we use the test set to evaluate the probability P_k of finding an example similar to the test input and the expected number R_k of candidate examples that need to be searched per input. This procedure is repeated a number of times. For a given P_k and R_k with l independent hash tables, the probability of success is $1 - (1 - P_k)^l$; we choose a fixed high probability (0.96) as our target value and find the smallest number l_k for which the target is achieved. We also need to balance the probability of success with the search time. Therefore, we choose the value of k that minimizes the expected search time given by $l_k R_k$, while achieving the desired probability of success. With 260 features selected by AdaBoost for yaw estimation, the best performance was achieved with $l = 15$ hash tables and with $k = 18$ bits per hash key. These parameters led to a significant acceleration in performance over a brute force search. Instead of searching 900,000 examples for each yaw estimation, the system needed to compare only about 100 examples.

7.1.2 Augmented Motion Graph. A single-character motion graph models the motion data as a first-order Markov process [Lee et al. 2002; Kovar et al. 2002]. The transition from one state to the next depends only on the current state (a single motion frame). The probability of a transition between two frames is computed from a distance function measuring the similarity of the body poses and velocities. The insight that made this approach effective is that the absolute location and orientation of the

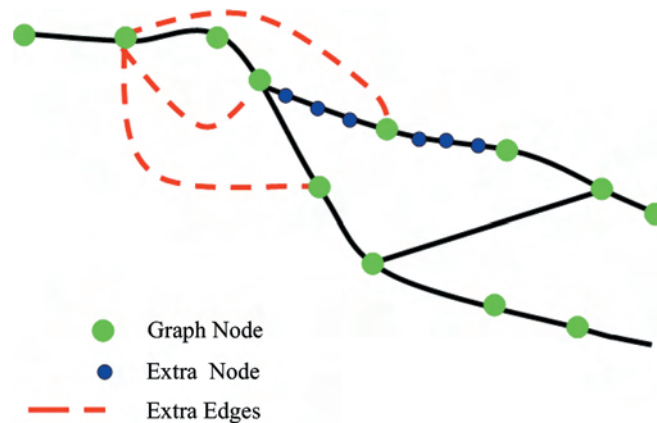


Fig. 11. Extra edges and nodes in the augmented motion graph. Three extra nodes between every two original motion frames allows for the refinement of the pose estimation and the slowdown of the playback speed. The extra edges allow one or two original nodes to be skipped.

character is not important for many recorded motions and the state can be represented in a body local or relative coordinate system. The relative coordinate system then allows transitions to similar motions recorded anywhere in the workspace and requires far less data to control the character's motion [Lee et al. 2002].

Because the motion of the two dancing partners performing East Coast Swing is highly correlated, we can easily build a two-character motion graph. We extend the relative coordinate system to two characters by representing the *partner's* position and orientation as relative to that of the *driven character* (the character that the user controls). The distance function used to compute the transitions is extended to two characters by comparing the body configuration and velocity for each character, the relative external pose of the partner, and the relative velocity between the two characters. In the distance function calculation, the body configuration is a high-dimensional vector representing joint positions in the root segment coordinate system. We use a strategy similar to that used by Lee et al. [2002] for pruning the transition edges. Foot contact states are also used in the pruning process; however, we use only the foot contact states of the driven character, not those of both characters. We blend transitions and maintain foot contact constraints for both characters using a hierarchical motion-fitting algorithm [Lee and Shin 1999]. We represent the partner's foot contact constraints in the root coordinate system of the driven character during the blending process.

The captured motion reflects the choreographic rules of swing dancing, that is, the motion consists of basic movements and begins and ends with specific poses. As observed by Kim et al. [2003], the transitions learned from the example motions in the motion graph may not perfectly reflect those choreographic rules. We use the beginning and ending poses of basic movements to connect nodes that meet choreographic rules but that were left out in the automatic motion graph construction. In practice, the transitions computed automatically cover most basic movements, and less than 1% of the transitions are added manually.

In the preprocessing step, the two-character motion graph is further augmented automatically by adding extra edges and nodes to improve the efficiency of the online search and to allow the speed of the motion to be synchronized to that of the user's performance (Figure 11). Extra edges allow frames to be skipped for faster playback and extra nodes allow the motion to be slowed. Only the feature vectors are stored for the extra nodes, and the joint angles are interpolated at runtime.

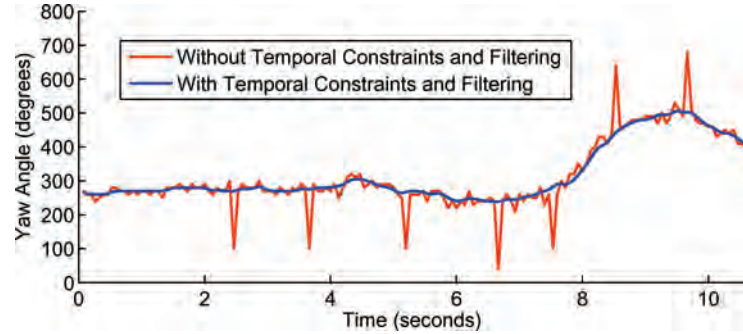


Fig. 12. Yaw estimation comparison. The red line shows the result of yaw estimation using a simple majority voting scheme. Temporal constraints and filtering eliminate errors of 180 degrees and smooth the estimated yaw. The testing video trial was captured at 15 frames per second.

7.2 Robust Online Yaw Estimation

Online yaw estimation improves the estimate by removing outliers caused by the inherent 180 degree ambiguity in silhouettes. First, we use LSH to retrieve a set of nearest examples (20 in the results reported here) and build a histogram of their yaw angles. Then we smooth the histogram using a low-pass filter. We choose the two peak values from the smoothed histogram, θ_1 and θ_2 , and form a candidate set $\theta_1, \theta_2, \theta_1 \pm 180, \theta_2 \pm 180$. We compare those candidates with the yaw value estimated in previous frames and choose the closest one as the new yaw estimation. Finally, we combine the yaw estimations from nearby video frames and apply a Gaussian low-pass filter which spans five frames. In our experience, a simple majority vote from 20 nearest neighbors yields an acceptable yaw estimate most of the time, but the method described above is more robust (Figure 12).

7.3 Online Body Configuration Estimation

Given the yaw estimate, the system searches locally in the augmented motion graph and computes the corresponding 3D motion frame with a full set of parameters (global translation, global orientation, and body configuration). The fast search is performed in two stages: first using multiple video frames to perform a coarse mapping and then refining that match using the extra nodes in the augmented motion graph.

To perform the coarse mapping, the system computes the matching cost for all the child nodes of the previous matched node, W , in the augmented motion graph (Figure 11) and selects the node ($W + 1$) with the smallest cost. This step may skip frames if the best match can be achieved with a faster motion (Figure 13). We use a look-ahead and look-back of r frames (a total of $2r + 1$ video frames) and compute the Hamming distance between the two sets of body configuration feature vectors (bit vectors). The matching cost is a weighted average using a Gaussian function centered on the current node W . There will likely be multiple segments of length $2r + 1$ centered at W because of branching in the augmented motion graph (Figure 13). The smallest of those costs for those segments is assigned to W . The look-ahead video buffer adds an r frame delay to the system, but improves robustness and removes much of the mapping ambiguity between the video input and the motion capture data.

The system refines this coarse match by searching the extra nodes around node ($W + 1$) for a lower cost match. This refinement uses only the current video frame for the computation of matching cost.

The extra nodes and edges in the augmented motion graph speed up the search and allow the original motion to be scaled in time to match the user's motion. To prevent a sudden change in velocity, we average the scaling factors across m frames (thus introducing an additional m frame delay for a total

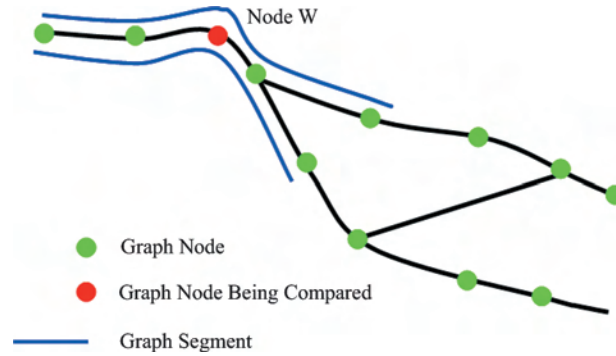


Fig. 13. Two different segments through a node. For simplicity, we illustrate the use of two look-ahead and two look-back video frames for the matching. In the experiments described in this article, ten look-ahead and ten look-back video frames were used.

delay of $r + m$ frames). With $r = 10$ and $m = 3$ and a frame rate of 15 fps, the delay is 0.87 s for the examples presented in this article. The delay introduced by Gaussian smoothing of the yaw estimates is also included in the 0.87 s delay because the yaw orientation is estimated independently of body configuration.

8. RESULTS

Using these algorithms for training and search, we built a vision-based performance interface for interactively controlling the motion of a dancing couple. We tested the performance of our system with users whose motions were not in the database and who were dancing to different music with a different tempo than that used when the database was recorded. Figure 14 shows three dancing sequences from two users.

To learn the yaw and body configuration features, we ran the AdaBoost algorithm until the learned ensemble classifier achieved a small testing error. Figure 15 shows the error of the yaw classifier on the test set as the number of features increases; qualitatively similar behavior was observed while learning body configurations features (Figure 16). The training was done on a Beowulf cluster with 90 nodes, each with a 1-GHz CPU and 512-MB RAM. Table I gives some of the parameters used in training, the training time required, and the performance of the classifiers. The result of the offline learning phase was a set of 260 binary features for the yaw similarity measure and a set of 320 features for each of the 36 body configuration similarity measures.

The yaw estimation uses the 20 top matches retrieved by LSH. We compared the yaw estimate obtained by our system to a ground truth provided by motion capture data (Figure 17). For 73% of the frames, the LSH-aided yaw estimation has an error of less than 10 degrees. The estimation error is less than 20 degrees for 92.5% of the frames and less than 30 degrees for 98% of the frames. Because errors in yaw estimation may affect the body configuration estimation, we tested our body configuration similarity measure in the presence of errors in yaw; the performance as a function of yaw error is plotted in Figure 18 for one yaw bin. Similar performance was seen for other yaw angles. This test demonstrates that the body configuration estimation is robust to the yaw estimation errors seen in our system.

We compare the performance of our system with that of Hu moments, a simple set of global features that have often been used in silhouette-based systems [Davis and Bobick 1997; Rosales and Sclaroff 2000; Lee et al. 2002]. Hu Moments were introduced in the early 1960s for pattern recognition of objects [Hu 1962] and are seven quantities that are invariant under translation, rotation, and scaling of the silhouette. Hu moments are also global features because the real value of each moment invariant



Fig. 14. Three dancing sequences created by two subjects wearing street clothes. The sample frames show the matching between the video, the male dancer, and the couple.

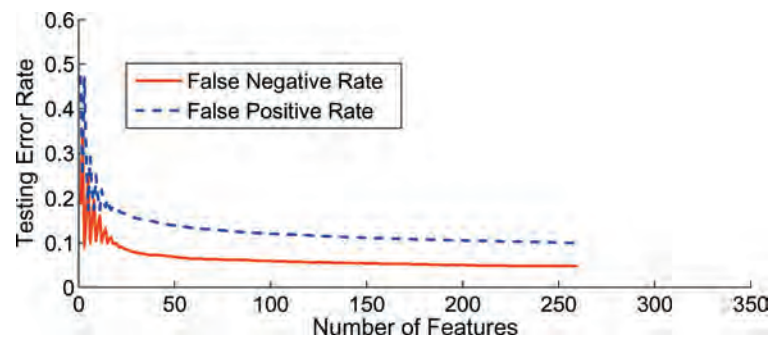


Fig. 15. The error of the yaw classifier on the test set as the number of features increases.

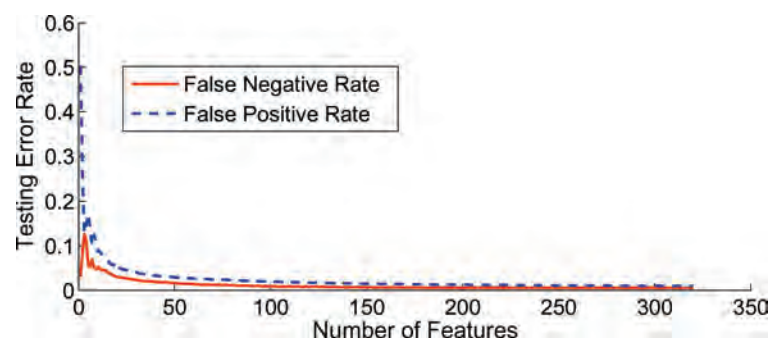


Fig. 16. The error of the body configuration classifier (60 degree yaw angle) on the test set as the number of features increases.

Table I. Parameter and Error Rates for the Learning Algorithms
(The statistics for each of the 36 body configuration similarity functions are approximately the same.)

| | Yaw | Body Configuration (60 degree yaw) |
|---|--------------------|---------------------------------------|
| Features | 260 | 320 |
| Training examples (positive/negative) | 800,000/2,000,000 | 200,000/10,000,000 |
| Training examples per round (positive/negative) | 4,000/6,000 | 3,000/5,000 |
| Resampling interval (rounds) | 40 | 80 |
| Testing examples (positive/negative) | 300,000/37,000,000 | 50,000/2,000,000 |
| Testing error rates (false negative/false positive) | 4.86%/9.92% | 0.5%/0.9% |
| Training Time (hours) | 2.5 | 1.8 |

is computed from the coordinates of all the pixels in the silhouette, not merely from the pixels in a local patch.

The comparison with Hu moments used motion data captured simultaneously with video data as the ground truth. To create a common reference frame, we placed a black and white checkerboard calibration pattern with motion capture markers in the field of view of both video and motion capture cameras. The subject then performed while both video data and motion capture data were recorded. The ground truth motion capture data were not included in the database so the best match from the database will not be an exact match.

We found the best match with respect to the global Hu moment metric and with respect to the local rectangular features. To make the comparison with Hu moments fair, we did not use temporal

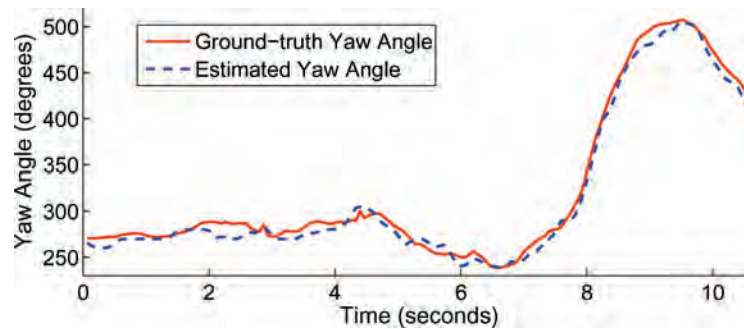


Fig. 17. Comparison of the estimated yaw with that computed from motion capture data that was recorded simultaneously with the input video data.

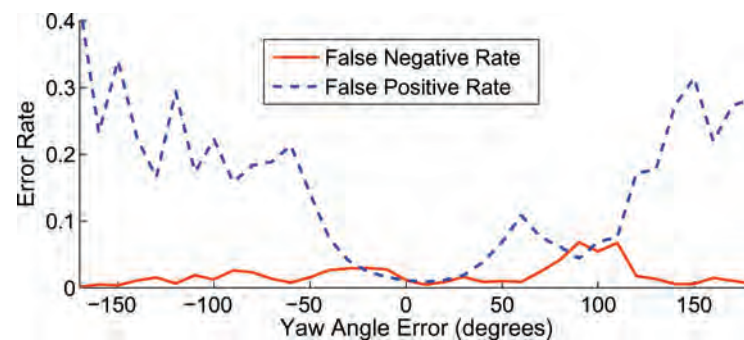


Fig. 18. Performance of body configuration classifiers in the presence of errors in yaw estimation. We tested the performance of each of the 36 body configuration classifiers using test data with a yaw angle of zero. The graph shows the false positive rate and the false negative rate of the body configuration classifier of each yaw bin. Similar performance was seen for other yaw angles. The result shows that body configuration estimation is robust to yaw estimation errors.

constraints or filtering in the yaw estimation, but otherwise the approach using local features was as described earlier. To find the global best match, we computed the joint positions in the ground truth motion capture data and used those 3D locations to retrieve the closest body pose from the database by exhaustive search. For each of these methods, we computed the error as the L2 distance in 3D parameter space between the retrieved match from the database and the recorded motion capture data for each testing frame. The histogram of the difference between the error obtained with Hu moments and the error with the local rectangular features is shown in Figure 19. The performance of local features was better than that of Hu moments on 86% of the frames. Figure 20 shows a sample of results from the three methods.

To evaluate whether the temporal information inherent in the augmented motion graph would correct body configurations that are poorly estimated by Hu moments, we created an animation using Hu moments and compared it to that created using local features. Seven Hu moments were extracted from each of the three silhouettes and stored in the feature database. The process was otherwise identical except that the system computed the Hamming distance between the body configuration bits in the estimated yaw bin in the first case, whereas the system calculated the Euclidean distance between the Hu moments of the input silhouette and those of the matched motion frame in each yaw bin and chose the smallest one as the similarity measure in the second case. A similar method was used in the previous approaches based on Hu moments [Davis and Bobick 1997; Lee et al. 2002]. Figure 21 demonstrates

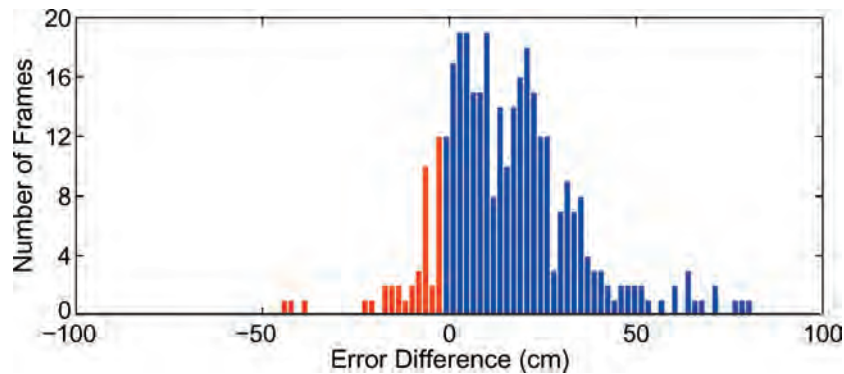


Fig. 19. Performance comparison between local features (rectangular features) and global features (Hu moments) for body configuration. The graph shows a histogram of the L2 estimation error for Hu moments minus that achieved with local features. Blue bars indicate the frames on which the local features performed better; red bars show the frames on which Hu moments performed better.

the performance of the two approaches on a motion sequence of about 7 s (15 fps). The estimation error is quite stable for rectangular local features but shows serious drifting problems for Hu moments.

9. DISCUSSION

We have described a silhouette-based vision interface that allows the user to control animated characters that reproduce his or her motion and that of his virtual dance partner. This interface is based on a set of learned local features. We have demonstrated that this approach produces visually pleasing motion for long sequences of user input. We have assessed this approach both qualitatively through animation and quantitatively through error measures and a comparison with an alternative approach that uses Hu moments.

Instead of using 2D local features from a multiview silhouette, we could have based our approach on visual hulls or 3D volumes. However, the computation cost for online reconstruction and the storage cost for those 3D examples would likely be prohibitive for an interactive vision-based interface. Matching the 3D representations to the motion in the database would be difficult because only a coarse and view-dependent 3D reconstruction can be obtained from three camera views. We chose to work on 2D silhouettes primarily for reasons of efficiency.

Silhouette similarity measures can be implemented in many different ways. They can be based on measures such as Hamming distance between foreground pixel masks, Hausdorff distance between points on contours, or Euclidean distance between global image features such as Hu moments. However, none of these similarity measures directly addresses the problem of “false matches,” where similar metric values do not always imply close poses. Because our feature set was selected from a very large set with the goal of avoiding false matches, our approach is likely to be more robust to this problem than a single, fixed metric.

Efficiency is also a concern in our interactive approach. Direct comparison between foreground pixel masks is not memory efficient while the Hausdorff distance metric is not computationally efficient. Hu moments, which are based on global statistics of silhouette pixels, are efficient but not robust because noisy silhouettes affect the values of all components. The local feature set we chose is particularly appealing for real-time or interactive vision-based applications because it requires only a small amount of storage, is computationally efficient, and the large number of features should make it robust to noise.

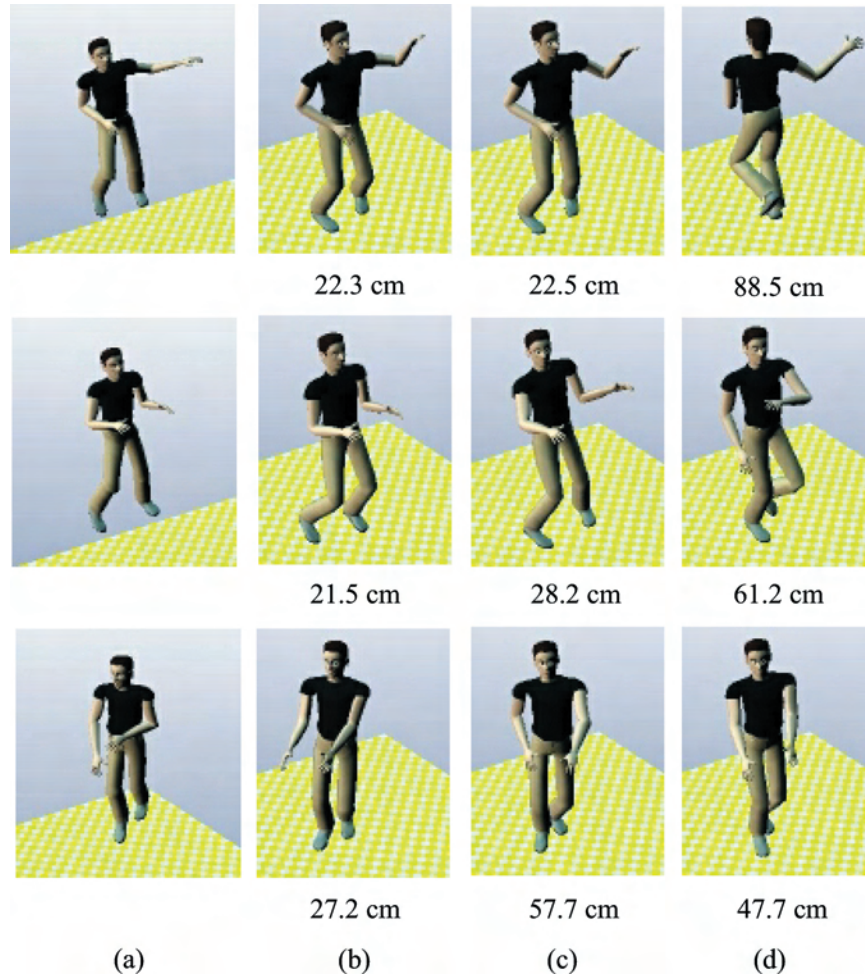


Fig. 20. Comparison of the poses retrieved by three different approaches. (a) The ground truth pose from a commercial motion capture system. The image shows the translation, orientation, and body configuration of the motion capture subject. (b) The best available match in the database found by a brute force search on the 3D locations of the joints in the coordinate system of the root segment. The image uses the retrieved body pose and the ground truth yaw angle. (c) The closest match found with rectangular local features. (d) The closest match found with Hu moments. For a fair comparison, temporal constraints and filtering are not used for yaw estimation in (c). No temporal information is used in any of the methods. Both (c) and (d) show the estimated yaw angle and body configuration. For each body configuration, we compute the L2 distance between the retrieved body configuration and the ground truth as a measure of the error in estimation. Those errors are given below each image. Because the body configuration is defined in the root coordinate system of the human body, the estimation error in yaw is ignored when computing those errors. The first row shows an example where Hu moments failed to correctly estimate the yaw orientation. In the second row, Hu moments had a good estimate for yaw orientation but failed to estimate the body configuration accurately. The third row shows an example where the approach based on Hu moments performs slightly better than the rectangular features. For our testing sequences, the average baseline error (distance to the best match from the true 3D parameters) was 25 cm; the average error with our method was 44 cm; the Hu moments had an average error of 76 cm.

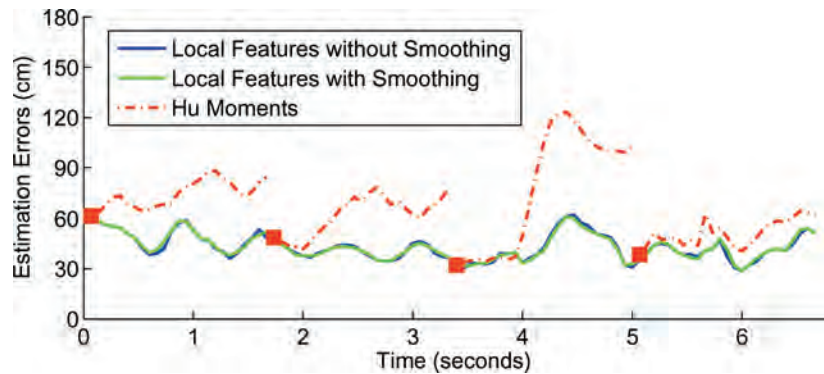


Fig. 21. Performance comparison between rectangular local features and global features (Hu moments). The testing video was captured at 15 frames per second. The estimation error is the L2 distance between the retrieved body configuration and the ground truth from the simultaneous motion capture recording. The local features provide a good match to the user's motion but the Hu moments tend to drift from the captured ground truth. Every 1.7 s, we manually restarted the search with Hu moments from the body configuration found with local features. The graph also shows the estimation error of local features when smoothing is disabled. The smoothing averages the scaling factors across three frames and improves the visual quality of the synthesized motion although this improvement is not visible in the computed error.

In our example-based approach, we use LSH for fast yaw estimation. The procedure can also be used in other fast parameter estimation applications. We compared our fast parameter estimation approach with the work of Shakhnarovich et al. [2003]. Instead of using AdaBoost for feature selection, they directly selected features or hash functions based on the single feature's classification performance on the original, unweighted training data. Their approach did not take into account dependence between those hash functions. If the feature space is highly redundant, as the rectangular feature set is, many variants of the same feature may be selected, creating a dependency between the hash functions. In contrast, we first used AdaBoost to learn a set of independent hash functions and then randomly select the locality-sensitive hash functions among them. We also set parameters for LSH based on an empirical evaluation rather than using a closed-form formula to choose the bounds for those parameters as was done by Shakhnarovich et al. [2003].

The major limitation to our example-based approach is that it requires a database that is appropriate for the problem domain. We believe that this is not a serious limitation for many applications where the likely human behaviors can be predicted in advance. For example, we envision that such a system might allow a little girl to spin and twirl in her family room and see an animated version of her motion applied to a cartoon hippo, a ballerina, or a little girl of the same age. Similarly, adults might be able to see an animated version of themselves gesturing while crooning into the microphone during a karaoke session. Although less structured than dancing, expressive gestures made during meetings might also be contained in a database of sufficient size and used for teleconferencing. We have not yet tested to see how much generality we can achieve with a given database. For applications where precise reproduction of the user's motion is not needed, a basic database may be sufficient.

We demonstrated that our system is robust to users with different body types as the synthetic character used for training was not specifically tuned to our subjects (Figure 22); however, we expect that the current system would not be robust to large changes in body type. We believe that this problem could be solved by training for a small set of different body types and then choosing the appropriate set of classifiers automatically based on the silhouettes extracted from a few calibration poses.

We also showed the robustness of our system to the clothing of the user (Figure 22). In the current system, the user may wear street clothes but they should be fairly form-fitting. This restriction is an

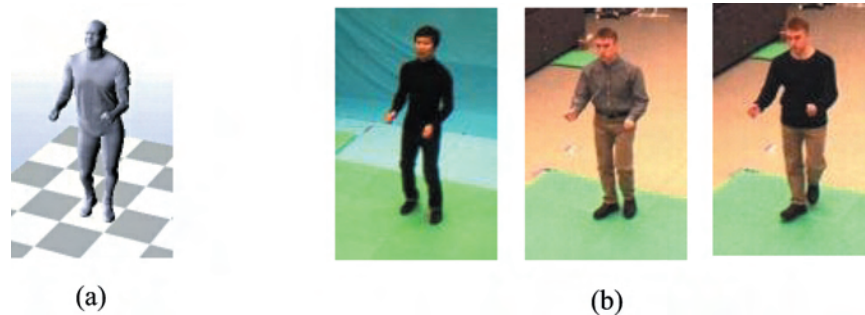


Fig. 22. (a) The human model used to generate synthetic silhouettes for training. (b) Human subjects wearing different street clothes.

inherent weakness of silhouette-based methods because the system must see the same silhouette each time that the user assumes a particular pose. Clothing that has “state” independent of the user’s current pose will introduce errors. We have observed that the system appears robust to a loose shirt but not to loose pants, which cause the silhouette to lose the separation between the legs.

Some restrictions in the setup prevent the current system from being used in an average family room. We assume a static background and reasonable color separation between the user and the background. We believe that more sophisticated silhouette extraction algorithms might make the approach perform better in real-world settings. The current implementation requires three cameras, and the selection of the local features is specific to the positions of the cameras. The framework might work reliably with two cameras, although we have not extensively tested this setup. We have also not assessed how robust the system is to changes in camera position and orientation. It should be robust to translation and rotation of the cameras about the vertical axis centered in the performance area because the performer makes movements that correspond to those changes of view point in the current system without introducing errors.

The implementation of the current system required that we select a number of parameters: the number of local features chosen for yaw and body configuration classification, the number of positive and negative training examples, the number of rounds of training, and the frequency with which the training set is changed for resampling. All of these parameters affect the learning time of the various classifiers. That time is now between 2–3 h per classifier on a large computing cluster (90 processors, 37 classifiers total). Although AdaBoost weights the training examples to ensure that those currently misclassified receive greater emphasis, training could likely be made more rapid by fine-tuning these parameters or by relying on more sophisticated algorithms in selecting the training set. Alternatively, the time required for training could be reduced at the cost of a small impact on performance by further reducing the available feature set. Wu et al. [2004] implemented such a system and used a simple and direct feature selection method to replace AdaBoost.

Although the yaw angle is estimated in the search process, it is not used in the animation of the character. The position and yaw orientation of the character is computed from the trajectory selected in the motion graph. Therefore, the character’s position and yaw orientation will diverge from that of the human dancer as small differences accumulate. We have not found divergence to be a serious drawback in the dancing application as the incremental translations and rotations appear natural. For an application where the absolute facing direction is important, such as karaoke, this problem could be at least partially addressed by including yaw orientation as part of the matching function when the path through the motion graph is selected.

Both the yaw and the body configuration search take recent values for those parameters into account, thus ensuring temporal consistency in the motion. However, if a wrong path is selected in the motion graph, this property could cause the system to stray quite far from the human dancer's motion. We have not seen this problem in practice, perhaps in part because swing dancing is structured with well-defined beginning and ending poses for each move. This structure may allow the search to recover from a less than ideal choice of paths in the motion graph with the start of the next dance move. If increasing error becomes a problem in other applications, we believe that locality-sensitive hashing could also be used to provide a global search for body configuration and to correct errors at the cost of some discontinuity in the motion.

ACKNOWLEDGMENTS

The authors would like to thank Michael Stevens for his assistance collecting and cleaning the motion capture data and Moshe Mahler for modeling the characters and environment.

REFERENCES

- ARIKAN, O. AND FORSYTH, D. A. 2002. Interactive motion generation from examples. *ACM Trans. Graph.* 21, 3, 483–490.
- ARIKAN, O., FORSYTH, D. A., AND O'BRIEN, J. F. 2003. Motion synthesis from annotations. *ACM Trans. Graph.* 22, 3, 402–408.
- BRAND, M. 1999. Shadow puppetry. In *Proceedings of the International Conference on Computer Vision*. 1237–1244.
- BRAND, M. AND HERTZMANN, A. 2000. Style machines. In *Proceedings of SIGGRAPH 2000*. 183–192.
- BREGLER, C. AND MALIK, J. 1998. Tracking people with twists and exponential maps. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 8–15.
- CALIFORNIA INSTITUTE OF TECHNOLOGY. 2002. Camera calibration toolbox for matlab. Available online at <http://www.vision.caltech.edu/bouguetj/calib.doc/>.
- CARRANZA, J., THEOBALT, C., MAGNOR, M. A., AND SEIDEL, H.-P. 2003. Free-viewpoint video of human actors. *ACM Trans. Graph.* 22, 3, 569–577.
- CHAI, J. AND HODGINS, J. 2005. Performance animation from low-dimensional control signals. *ACM Trans. Graph.* 24, 3, 686–696.
- CHEUNG, K. M., KANADE, T., BOUGUET, J.-Y., AND HOLLER, M. 2000. A real time system for robust 3D voxel reconstruction of human motions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 714–720.
- COLLINS, M., SCHAPIRE, R., AND SINGER, Y. 2000. Logistic regression, AdaBoost and Bregman distances. In *Proceedings of Computational Learning Theory*. 158–169.
- CROW, F. C. 1984. Summed-area tables for texture mapping. In *Proceedings of SIGGRAPH 1984*. 207–212.
- DAVIS, J. W. AND BOBICK, A. F. 1997. The representation and recognition of human movement using temporal templates. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 928–934.
- DELAMARRE, Q. AND FAUGERAS, O. D. 1999. 3D articulated models and multi-view tracking with silhouettes. In *Proceedings of the International Conference on Computer Vision*. 716–721.
- DEUTSCHER, J., BLAKE, A., AND REID, I. 2000. Articulated body motion capture by annealed particle filtering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 126–133.
- DUDA, R., HART, P. E., AND STORK, D. G. 2000. *Pattern Classification*. John Wiley & Sons, Inc., New York, NY.
- EFROS, A., BERG, A. C., MORI, G., AND MALIK, J. 2003. Recognizing action at a distance. In *Proceedings of the IEEE International Conference on Computer Vision*. 726–733.
- GAVRILA, D. M. 1999. The visual analysis of human movement: A survey. *Comput. Vis. Image Understand.* 73, 82–98.
- GIONIS, A., INDYK, P., AND MOTWANI, R. 1999. Similarity search in high dimensions via hashing. In *Proceedings of the 25th International Conference on Very Large Data Bases*. 518–529.
- GRANIERI, J. P., CRABTREE, J., AND BADLER, N. I. 1995. Production and playback of human figure motion for visual simulation. *ACM Trans. Model. Comput. Sim.* 5, 3, 222–241.
- HU, M. K. 1962. Visual pattern recognition by moment invariants. *IEEE Trans. Inform. Theor.* 8, 179–187.
- JONES, M. AND VIOLA, P. 2003. Face recognition using boosted local features. MERL Tech. rep. TR2003-25. Mitsubishi Electric Research Laboratories, Cambridge, MA.
- KIM, T. H., PARK, S., AND SHIN, S. Y. 2003. Rhythmic-motion synthesis based on motion-beat analysis. In *ACM Trans. Graph.* 22, 3, 392–401.

- KOVAR, L., GLEICHER, M., AND PIGHIN, F. 2002. Motion graphs. *ACM Trans. Graph.* 21, 3, 473–482.
- LEE, J., CHAI, J., REITSMA, P., HODGINS, J., AND POLLARD, N. 2002. Interactive control of avatars animated with human motion data. *ACM Trans. Graph.* 21, 3, 491–500.
- LEE, J. AND SHIN, S. Y. 1999. A hierarchical approach to interactive motion editing. In *Proceedings of SIGGRAPH 1999*. Computer Graphics Proceedings, Annual Conference Series. ACM Press, New York, NY, 39–48.
- LEVENTON, M. E. AND FREEMAN, W. T. 1998. Bayesian estimation of a 3D human motion from an image sequence. MERL Tech. rep. TR1998-06. Mitsubishi Electric Research Laboratories, Cambridge, MA.
- MATUSIK, W., BUEHLER, C., AND McMILLAN, L. 2001. Polyhedral visual hulls for real-time rendering. In *Proceedings of the 12th Eurographics Workshop on Rendering Techniques*. Springer-Verlag, Berlin, Germany, 115–126.
- MIKIC, I., TRIVERDI, M., HUNTER, E., AND COSMAN, P. 2001. Articulated body posture estimation from multicamera voxel data. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 455–461.
- MORI, G. AND MALIK, J. 2002. Estimating human body configurations using shape context matching. In *Proceedings of the European Conference on Computer Vision*. Vol. 3. 666–680.
- NOSER, H. AND THALMANN, D. 1997. Sensor based synthetic actors in a tennis game simulation. In *Proceedings of the 1997 Conference on Computer Graphics International*. 189.
- POINT GREY CORPORATION. 2001. Dragonfly cameras. Available at <http://www.ptgrey.com>.
- RAMANAN, D. AND FORSYTH, D. A. 2004. Automatic annotation of everyday movements. In *Advances in Neural Information Processing Systems 16*. MIT Press, Cambridge, MA, 1547–1554.
- ROSALES, R. AND SCLAROFF, S. 2000. Specialized mappings and the estimation of body pose from a single image. In *Proceedings of the IEEE Human Motion Workshop*. 19–24.
- ROSALES, R., SIDDIQUI, M., ALON, J., AND SCLAROFF, S. 2001. Estimating 3D body pose using uncalibrated cameras. In Boston University Computer Science Department Tech. rep. 2001-008. Boston University, Boston, MA.
- SCHAPIRE, R. E. AND SINGER, Y. 1999. Improved boosting algorithms using confidence-rated predictions. *Mach. Learn.* 37, 3, 297–336.
- SHAKHNAROVICH, G., VIOLA, P., AND DARRELL, T. 2003. Fast pose estimation with parameter-sensitive hashing. In *Proceedings of the IEEE International Conference on Computer Vision*. 750–757.
- SHIN, H. J., LEE, J., SHIN, S. Y., AND GLEICHER, M. 2001. Computer puppetry: An importance-based approach. *ACM Trans. Graph.* 20, 2, 67–94.
- SHIPP, C. A. AND KUNCHEVA, L. I. 2002. An investigation into how Adaboost affects classifier diversity. In *Proceedings of 9th International Conference on Information Processing and Management of Uncertainty in Knowledge-Based Systems*. 203–208.
- SIDENBLADH, H., BLACK, M. J., AND SIGNAL, L. 2002. Implicit probabilistic models of human motion for synthesis and tracking. In *Proceedings of the European Conference on Computer Vision*. 784–800.
- SMINCHISESCU, C. AND TRIGGS, B. 2003. Estimating articulated human motion with covariance scaled sampling. In *Int. J. Robotics Res.* 22, 6, 371–393.
- STENGER, B., THAYANANTHAN, A., TORR, P., AND CIPOLLA, R. 2003. Filtering using a tree-based estimator. In *Proceedings of the International Conference on Computer Vision*. 1063–1070.
- VICON MOTION SYSTEMS. 2002. Available at <http://www.vicon.com/>.
- VIOLA, P. AND JONES, M. 2001. Rapid object detection using a boosted cascade of simple features. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 511–518.
- VIOLA, P., JONES, M. J., AND SNOW, D. 2003. Detecting pedestrians using patterns of motion and appearance. In *Proceedings of the IEEE International Conference on Computer Vision*. 734–741.
- WU, J., REHG, J. M., AND MULLIN, M. D. 2004. Learning a rare event detection cascade by direct feature selection. In *Advances in Neural Information Processing Systems 16*. MIT Press, Cambridge, MA. 1523–1530.
- YAMAMOTO, M., SATO, A., KAWADA, S., KONDO, T., AND OSAKI, Y. 1998. Incremental tracking of human actions from multiple views. In *Proceedings of the IEEE International Conference on Computer Vision*. 2–7.
- YIN, K. AND PAI, D. K. 2003. Footsee: An interactive animation system. In *Proceedings of the 2003 ACM SIGGRAPH/Eurographics Symposium on Computer Animation*. 329–338.

Received March 2004; revised December 2004; accepted May 2005