

Image-Based Reconstruction for View-Independent Human Motion Recognition

**Robert Bodor
Bennett Jackson
Osama Masoud
Nikolaos Papanikolopoulos**

**Artificial Intelligence, Robotics and Vision Laboratory
Dept. of Computer Science and Engineering
University of Minnesota**

**Technical Report
3/25/2003**

Table of Contents

Introduction.....	2
Description of Work.....	4
Image Capture.....	5
Foreground Segmentation.....	7
Calculation of Motion.....	9
Determination of Virtual Camera Parameters.....	11
Image-Based Rendering of Orthogonal Views	12
Motion Recognition System	13
Results.....	15
Conclusions.....	16
Acknowledgments.....	16
References.....	17
Appendix A: Transformation Matrices.....	19

Introduction

The problem of using computer vision to track and understand the behavior of human beings is a very important one. It has applications in the areas of human-computer interaction, user interface design, robot learning, and surveillance, among others.

Much work has been done in the field of human motion recognition. However, a great deal of this work has been done using a single camera providing input to a training-based learning method ([1], [3], [4], [8], [9], [10], [12], [16], [17], [18], [22], [23]). These methods, which rely on a single camera, are highly view-dependent. They are usually designed to recognize motion when viewing all subjects from a single viewpoint (for example, from the side, with the subject moving in a direction orthogonal to the camera's line of sight (Figure 1)). The training becomes increasingly ineffective as the view angle increases away from orthogonal. As a result, these methods have limited real-world applications, since it is often impossible to limit the direction of motion of people so rigidly.



Figure 1. A single frame of human motion taken from a camera positioned orthogonal to the direction of motion.

Many of the methods track a single person indoors for both gesture recognition and whole body motion recognition ([1], [4], [22], [23]). Bregler details a method for motion analysis by modeling the articulated dynamics of the human body in outdoor scenes taken from a single camera [3].

Analysis of human gaits is also very popular in this domain. Statistical learning methods were developed by Fablet and Black to characterize walking gaits from 12 angles [8]. Additionally, Polana and Nelson studied the periodicity of the legs while walking or running to classify pedestrian motion [17].

Gavrila and Davis [10] used multiple synchronized cameras to reconstruct 3D body pose and study human motion based on 3D features, especially 3D joint angles. Experimental data was taken from four near-orthogonal views of front, left, right, and back.

Oliver *et. al.*, used coupled hidden Markov models along with a Kalman filter in a hierarchical approach to track and categorize motion [16].

Rosales and Sclaroff [18] developed a trajectory-based recognition system that is trained on multiple view angles of pedestrians in outdoor settings. This system can avoid the common problem of being limited in its ability to recognize motion from only a small set of view angles. However, the method requires a tremendous amount of training data.

Difranco *et. al.*, [6] tackle the problem of reconstructing poses from challenging sequences like athletics and dance from a single view point based on 2D correspondences specified either manually or by separate 2D registration algorithms. They also describe an interactive system that makes it easy for users obtain 3D reconstructions very quickly based on a small amount of manual effort.

In [11], Luck *et. al.*, generate a 3D voxelized model of a moving human, extracting the moving form using adaptive background subtraction and thresholding in real-time. Joint angle constraints are used to reconstruct an 11-degree of freedom model of the body.

Additionally, silhouettes have been used widely as tools for recognition. Weik and Liedtke [21] use 16 cameras and a shape-from-silhouette method to build a model template for the body, which is then matched to volume data generated in each successive frame. Pose analysis is performed based on an automatic reconstruction of the subject's skeleton. Silhouettes are also used by [2] and [5] to classify poses.

The image based rendering system we describe in this paper removes the view-dependent constraints of all of the motion recognition systems described above, thus making their application much more possible in a variety of settings, without modification to the recognition systems themselves.

Description of Work

Using image-based reconstruction, orthogonal input views can be created automatically from the renderer to construct the proper view from a combination of non-orthogonal views taken from several cameras (see Figure 2).

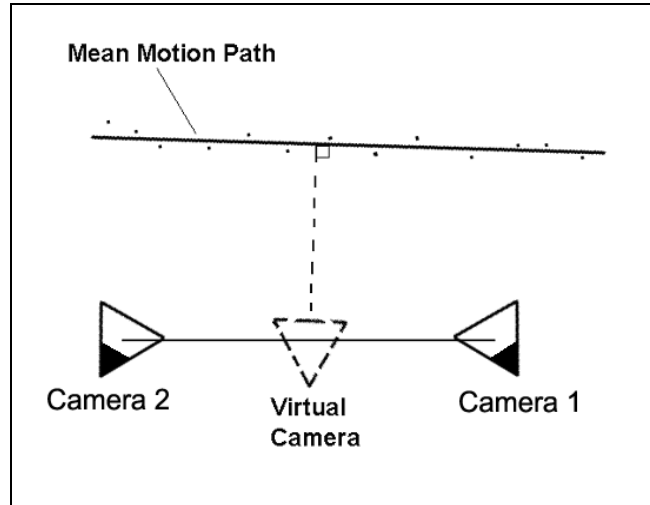


Figure 2. Camera positions relative to motion path.

Image-based rendering can be applied to motion recognition in two ways: i) to generate additional training sets containing a large number of non-orthogonal views and ii) to generate orthogonal views (the views that 2D motion recognition systems are trained for) from a combination of non-orthogonal views taken from several cameras. Figure 3 shows two examples of novel views generated from three different real camera views using image-based rendering. The bulk of this paper focuses on this second approach.

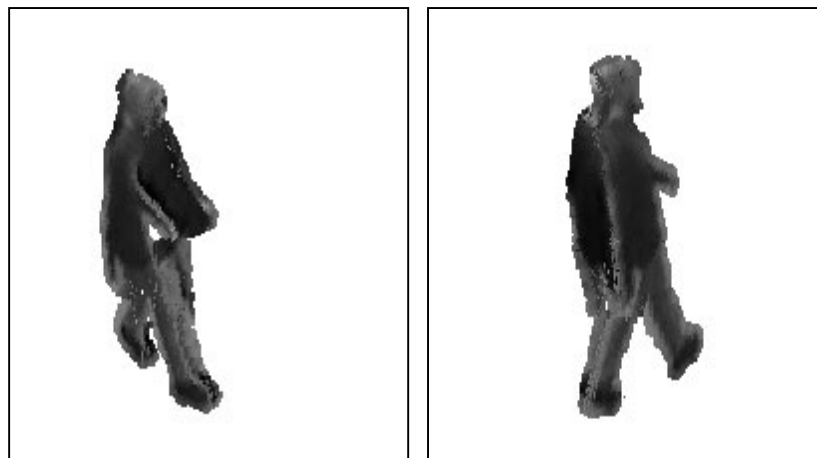


Figure 3. Two novel views generated by the image-based renderer.

To test the system, we fed the images created by the image-based renderer into an existing view-dependent human motion recognition system. This motion recognition system [12] works by first representing an action in terms of a set of points (a manifold) in a low dimensional eigenspace. Training involves computing, for each action, an

“average” manifold, which is used as a representation for that action. The manifold of a new action (whose images are computed orthogonally using image-based rendering) is then compared to all representative manifolds by computing a distance measure. The minimum distance is used to determine the classification result.

The complete process consists of the following steps:

- 1) Capture synchronized video sequences from several cameras spread around the perimeter of an area of interest
- 2) Digitize the video into a sequence of individual frames for processing
- 3) Transform images to compensate for wide angle lens distortions (barrel and pincushion)
- 4) Calibrate the cameras (intrinsic and extrinsic parameters)
- 5) Segment the subject in the images (separate foreground / background) (Ben Jackson helped with this step)
- 6) Reconstruct the motion path of the subject in the 3D world frame
- 7) Calculate the optimal “virtual” camera intrinsic and extrinsic parameters for use in the image-based renderer
- 8) Use an image-based renderer to create orthogonal views as training and test data (Osama developed this)
- 9) Test the views using a view-dependent motion recognition system

Image Capture

The use of image-based rendering requires capturing multiple views of the same scene simultaneously. We implemented this approach for the human motion recognition application with three video cameras positioned around one half of a room (see Figure 4).

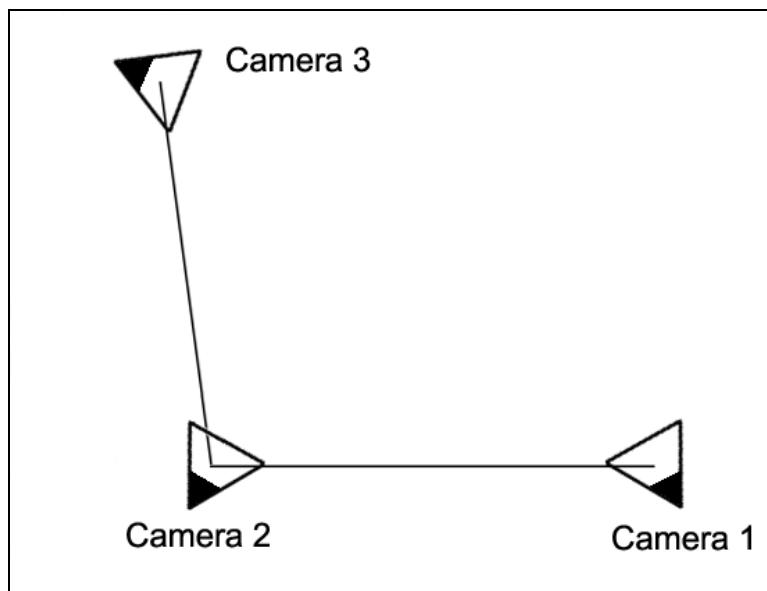


Figure 4. Layout of cameras for synchronized image capture.

This used three Panasonic GP-KR222 digital video cameras, with Rainbow auto-iris lenses. Two of the lenses were wide angle, with focal lengths of 3.5mm, while the third had a focal length of 6mm. Each camera was connected to a VCR, and the video was recorded onto three VHS tapes. The video from the tapes was then digitized and separated into sequences of single frames, and all three sequences were synchronized.

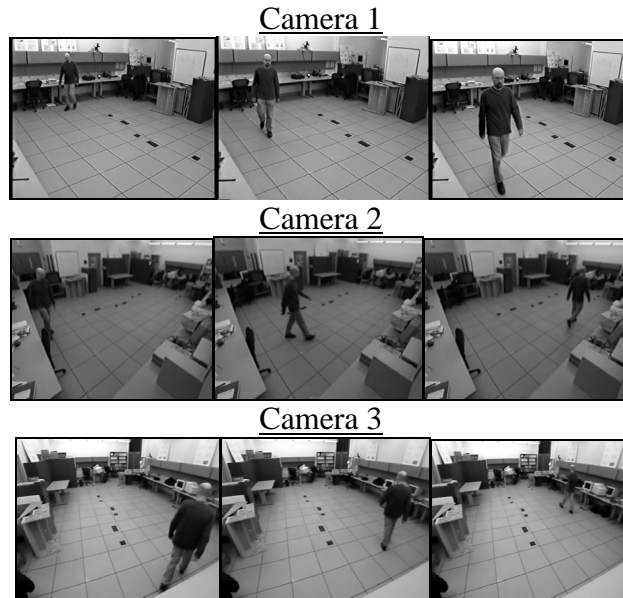


Figure 5. Images of a single walking sequence captured from 3 cameras positioned around a room (undistorted).

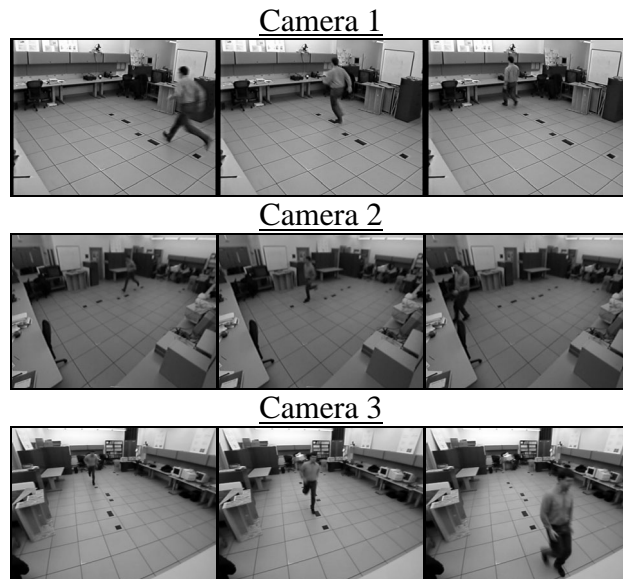


Figure 6. Images of a single running sequence captured from 3 cameras positioned around a room (undistorted).

The cameras were calibrated to determine all intrinsic and extrinsic parameters. We used the method of Masoud *et. al.*, described in [13]. The calibration involved selecting

parallel lines in an image for each camera. The result was series of homography transformation matrices between the cameras, as well as the intrinsic and extrinsic (position and orientation) parameters of each camera (see below, and in Appendix A).

In addition, the images were corrected for lens distortion due to the wide-angle lenses. For part of this procedure, we used the process of Ojanen [15]. We used Matlab software provided by Ojanen to dewarp the images, based on a 2D grid of points. However, we had to modify the Matlab code, as the wide-angle lenses we were using created more distortion than the method accounted for, thus it mis-correlated many of the points. Once this was done, the program output a de-warp correction function that could be applied to each image. We wrote a batch file to run this on all input images, and a second bit of Matlab code to clip the resulting images to their original size and output them in the proper format.

The Ojanen dewarp process accounted for barrel and pincushion wide-angle distortions, as well as translational offset of the lens center. The radial distortion was modeled as a third-order function given by equation 1 below.

$$r' = r + c_1 r^2 + c_2 r^3 \quad (1)$$

Foreground Segmentation

It was necessary to isolate the human for reconstruction in each frame. This required an automated process to separate the human (foreground) from the background of the room. We attempted a variety of approaches to this problem. In fact, this part of the project took about three weeks and was the longest single portion to implement. Ultimately, we found that best results could be achieved by using a combination of simple background subtraction, chromaticity analysis, and morphological operations.

Each frame in the motion sequence was subtracted from the background image in grayscale. This is a simple, standard approach, and is well known to fail in many cases where adaptive backgrounding has been found to be more robust, particularly in changing environments such as outdoor scenes [20]. Since the sequences were taken indoors, background subtraction performed quite well, with the exception of leaving significant shadows below the subject in the foreground image (see Figure 7). This was particularly problematic in this case, since articulation of the legs is critical for motion recognition.

To address this problem, we tried a number of things, but finally settled on an approach similar to that described by Elgammal et. al., [7]. This involved converting the RGB images into chromaticity space (equations 2 below). We then applied a subtraction and threshold operation to each frame (see Figure 7). This approach worked very well at removing shadows, but had the negative effect of also removing much of the valid foreground.

$$\begin{aligned}
 r &= R / (R + G + B), \\
 g &= G / (R + G + B), \\
 b &= B / (R + G + B).
 \end{aligned}
 \tag{2}$$

where $r + g + b = 1$.



Figure 7. Foreground process (left to right): Source image, background subtracted image containing shadow, chromaticity image, weighted chromaticity image, and composite mask with shadow removed.

To retain the best qualities of each method, we generated a composite mask for the final foreground image. This mask was a linear combination of the mask generated by background subtraction and the chromaticity mask. The chromaticity mask was first multiplied by a linear confidence weighting, with full confidence at the bottom of the bounding box of the subject (to remove shadows) and zero confidence at the top (to avoid the chromaticity mask removing valid foreground). Lastly, we applied a series of standard morphological operations such as hole filling, dilation, and erosion to the composite mask. The resultant masks may be seen in Figure 8.



Figure 8. Foreground results for one set of synchronized images.

Foreground segmentation is still an unsolved problem. We are very interested in studying this process further. The current approach does not use all of the available information, as it segments individual images, without using the time component (the fact that this is a sequence where the subject is tracked), or the multiple camera views at each time step. The benefit of the method is that it requires only foregrounding to work, and is quite fast. However, in the future we would like to investigate optical flow and multi-view methods as well to increase the robustness of the system in more challenging environments such as outdoor surveillance.

Calculation of Motion

From each foreground mask, a perimeter outline image was calculated. This perimeter was used to calculate input silhouette contours for the image-based rendering (see below) and it was also used to extract the centroid (based on the perimeter) and the bottom center of the subject in each image (see Figure 9).

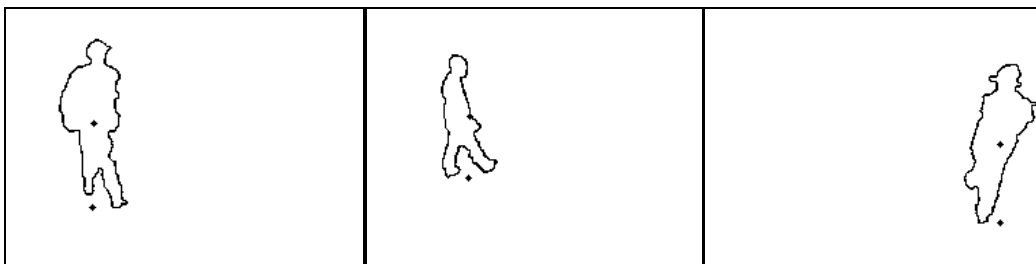


Figure 9. Silhouette perimeter images with centroid and bottom center point displayed.

The bottom center point is simply the point where a vertical line drawn through the centroid of the figure intersects the horizontal line at the bottom of the mask in the perspective of each camera. These points are used to calculate the location of the subject on the plane of the floor in the world frame. We chose the world frame relative to camera 1 to accommodate the assumptions the image-based renderer uses, and constructed the geometry using the transformations from the calibration step.

The bottom center point is calculated relative to the ground in each camera's reference frame. To convert them to a common world frame (the ground plane relative to camera 1), we multiplied each point by the inverse of its camera's homography matrix, and then the transformation matrix between the camera and camera 1. The transformation matrix encodes the change in translation and orientation difference between the cameras. Lastly, we then multiplied the point by the homography matrix of camera 1 to move it to the ground plane of camera 1 (see Appendix A).

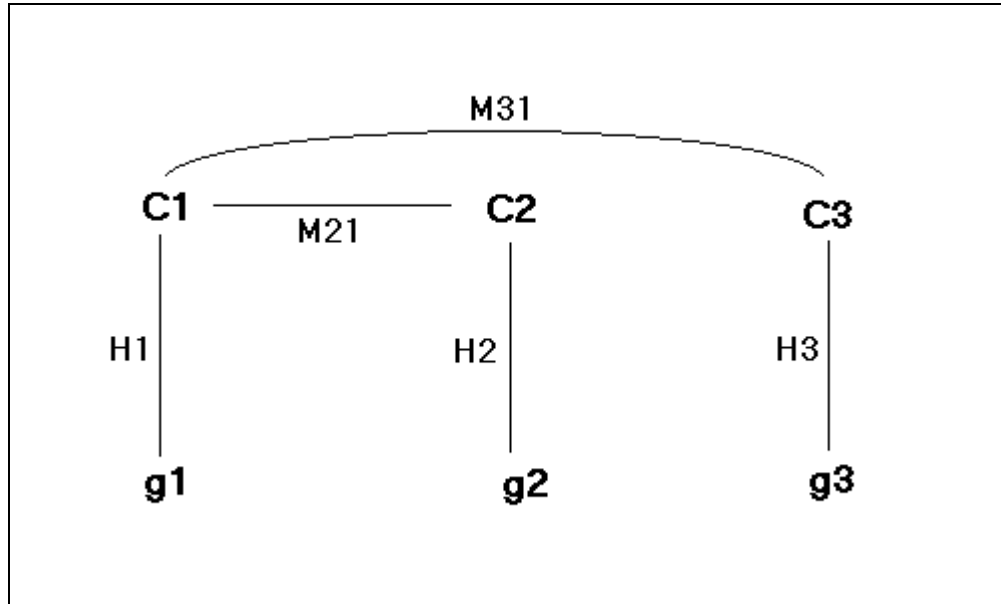


Figure 10. Coordinate system transformations.

For example, for camera 3, the equation would be:

$$g31 = g3 * H3^{-1} * M31 * H1 \quad (3)$$

Once this is done, the subject's position is calculated by projecting each of the three points into the world frame and calculating the Euclidean mean. When this is done at each time step, the subject's entire motion path along the floor can be tracked. Figure 10 shows the three bottom center points projected into the world frame for the first and last positions in the walking sequence, along with the averaged position estimates.

Determination of Virtual Camera Parameters

To create optimal input for the motion recognition system, it is critical to position the virtual camera such that it is orthogonal to the mean motion path of the subject. These calculations were done in the world frame using the results of the calibration procedure.

It is also important to consider the limitations of the image-based renderer when choosing the viewpoint for the virtual camera. The quality of reconstructed views provided by image-based rendering degrade as the viewpoint of the virtual camera moves away from the views of the real cameras in the scene. Therefore, we chose to position the virtual camera center in a plane between the real cameras (see Figures 10 and 11). The virtual camera's orientation was calculated to intersect the midpoint of the motion path at 90 degrees. In addition, the field of view of the virtual camera was automatically calculated to encompass the entire motion sequence without introducing wide-angle distortions. These calculations involved trigonometric rules such as the law of cosines, and finding the intersection of two lines in three-space.

A line can be defined by the slope intercept form:

$$\begin{aligned}y &= mx + b, \\m &= (y_2 - y_1)/(x_2 - x_1), \\b &= y_1 - mx_1\end{aligned}\tag{4}$$

Now, if we consider an alternate line form:

$$ax + by = c\tag{5}$$

where $a_n = -m$, $b_n = 1$, $c_n = b$.

We can solve the system of two equations to find the point where the two lines intersect:

$$\begin{bmatrix} a_1 & b_1 \\ a_2 & b_2 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} c_1 \\ c_2 \end{bmatrix}\tag{6}$$

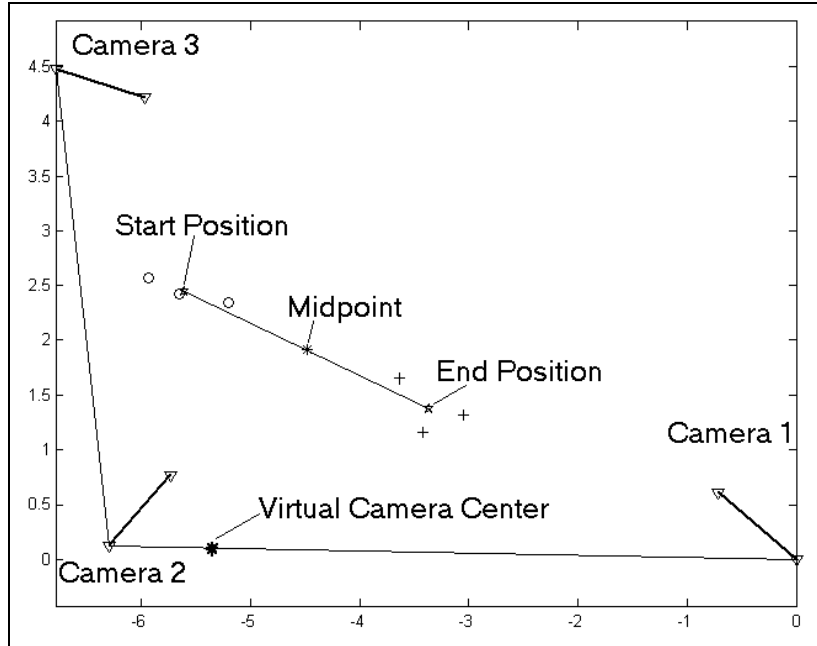


Figure 10. Walking path and virtual camera relative to fixed cameras for subject 1's motion sequence.

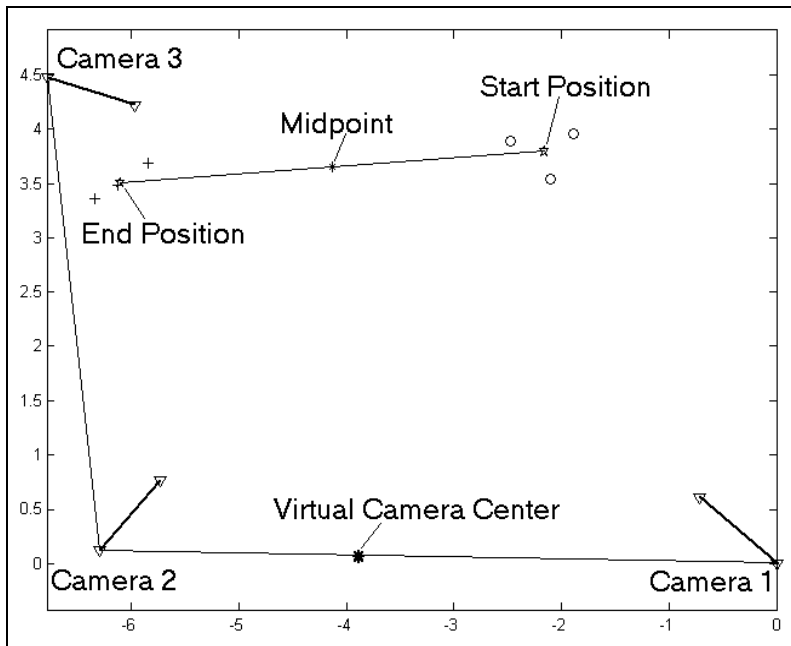


Figure 11. Running path and virtual camera relative to fixed cameras for subject 2's motion sequence.

Image-Based Rendering of Orthogonal Views

To generate novel views, an image-based approach was used. Given the parameters of a novel virtual camera, the renderer uses an approach similar to view morphing [19]. The views used to render the new image are those of the nearest two cameras. This approach differs in that it does not restrict the novel camera center to lie on the line connecting the

centers of the two cameras. In order for view morphing to work, depth information in the form of pixel correspondences needs to be provided. To compute this, we used an efficient epipolar line-clipping algorithm [14] which is also image-based. This method uses the silhouettes of an object to compute the depth map of the object's visual hull. Once the depth map is computed, pixel correspondences are calculated.

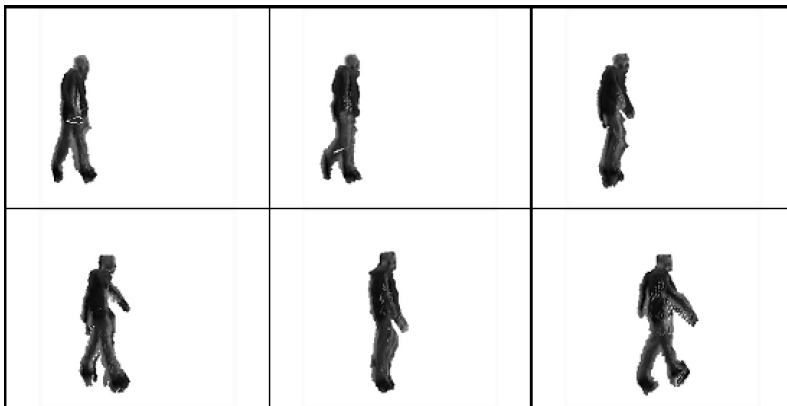


Figure 12. Rendered orthogonal images of the walking sequence.



Figure 13. Rendered orthogonal images of the running sequence.

Motion Recognition System

As mentioned earlier, we used a view-dependent human motion recognition method [12] to study the feasibility of the approach. This recognition method uses principal component analysis to represent a periodic action as a manifold in eigenspace. The images of an action are first processed to produce *feature images*. These images capture instantaneous motion since they are computed by recursively filtering the sequence images.

$$F_i = |M_i - I_i| \quad (7)$$

where:

$$M_i = \alpha \times I_{i-1} + (1 - \alpha) \times M_{i-1} \quad (8)$$

and where I is the image at time i .

Figure 14 shows some examples of feature images computed for the sequence in Figure 12.



Figure 14. Examples of feature images corresponding to the snapshots shown in Figure 12. These images are computed by recursive filtering and are used by the recognition algorithm.

In the training phase, a basis consisting of the principle components is first computed using all feature images of the training data. Then, each feature image is projected using the basis. Thus, an action can be represented as a sequence of points in eigenspace (a manifold). For every action, the average manifold is computed and stored in a database. Testing is performed by computing the manifold of the test action and finding the action in the database whose manifold is most similar.

For two manifolds, A and B , their distance is defined by:

$$d(A, B) = \frac{1}{l} \sum_i \min_{1 \leq j \leq h} \left\| \frac{a_i}{\|a_i\|} - \frac{b_j}{\|b_j\|} \right\| \quad (9)$$

To ensure symmetry, the actual distance measure used was:

$$D(A, B) = d(A, B) + d(B, A) \quad (10)$$

In the experiments, the system was trained to classify an action into one of eight categories: walk, run, skip, march, hop, line-walk, side-walk, and side-skip. Training sequences were provided by eight subjects performing all eight actions (a total of 64 sequences). The subjects used in testing were not among those eight training subjects.

Results

We tested the virtual views produced by the image-based renderer on the walking and running sequences shown above, and using the original training data captured from sequences taken orthogonal to the subject's motion, as in Figure 1. Each image sequence was processed in three ways and each was input to the motion recognition system.

The first set of three input sequences were the source video sequences taken from all three cameras. In all cases, the system failed to identify the motion. This was true even when part of the sequence was orthogonal to the camera and the mean motion path was near orthogonal (as with subject A camera 2 sequence).

<u>Test Sequence</u>	<u>Angle from Orthogonal</u>	<u>Result</u>
<u>Subject A Walking</u>		
Camera 1 source images	74.1°	Failed
Camera 2 source images	15.9°	Failed
Camera 3 source images	81°	Failed
Foreground images	15.9°	Failed
Image-based rendered orthogonal images	0°	Motion correctly recognized
<u>Subject B Running</u>		
Camera 1 source images	44.5°	Failed
Camera 2 source images	45.5°	Failed
Camera 3 source images	69.9°	Failed
Foreground images	44.5°	Failed
Image-based rendered orthogonal images	0°	Motion correctly recognized

The fourth input sequence in each case was the set of foreground images of the camera that had the smallest angle away from orthogonal (subject A camera 2 and subject B camera 1, see Figure 8). We tested the foreground sequence to ensure that any classification result was not due to simply removing the background and thus somehow better isolating the subject's motion for the recognition system. In this case, the system also failed to classify the motion correctly.

The fifth input sequence was the set of virtual images generated by the image-based renderer from a view orthogonal to the motion path (see Figures 10 through 13). This sequence was correctly identified for both motions.

Conclusions

We used image-based rendering to generate additional training sets for view-dependent human motion recognition systems. We believe this process is a novel method for employing image-based rendering to extend the range of use of human motion recognition systems. Input views orthogonal to the direction of motion are created automatically to construct the proper view from a combination of non-orthogonal views taken from several cameras.

The method described here may be used in several ways to reduce the constraints necessary to use 2D recognition in many environments. This may be done by creating a comprehensive set of training data for 2D motion recognition methods with views of motion from all angles, or by converting non-orthogonal views taken from a single camera into orthogonal views for recognition. In addition, the method can be used to provide input to a three-dimensional motion recognition system because the system creates a 3D volume model of the subject over time using only foreground segmentation.

In the future, we intend to test the method comprehensively with a large base of subjects and motion types. In addition, we intend to test the system in outdoor as well as indoor environments, and attempt to determine image capture parameters such as the minimum angle where recognition begins to degrade and the number of cameras necessary to achieve various levels of performance.

Lastly, the current method assumes a relatively linear motion path. We would like to extend it to work for arbitrary motion paths.

Acknowledgements

We would like to thank the National Science Foundation for their generous support of this work.

References

- [1] Ben-Arie, J., Wang, Z., Pandit, P., and Rajaram, S. "Human Activity Recognition Using Multidimensional Indexing," *IEEE Trans. on Pattern Analysis and Machine Intelligence*, vol. 24, No. 8, August, 2002.
- [2] Bradski G. and Davis, J. "Motion Segmentation and Pose Recognition with Motion History Gradients," *Int'l Journal of Machine Vision and Applications*, Vol. 13, No. 3, 2002, pp. 174-184.
- [3] Bregler, C. "Learning and Recognizing Human Dynamics in Video Sequences," *IEEE Conf. on Computer Vision and Pattern Recognition*, pp 568-574, 1997.
- [4] Cutler, R. and Turk, M. "View-Based Interpretation of Real-Time Optical Flow for Gesture Recognition," *Proc. Third IEEE Conf. on Face and Gesture Recognition*, Nara, Japan, April 1998.
- [5] Davis J. and Bobick, A. "The Representation and Recognition of Human Movement Using Temporal Templates," *Proc. Computer Vision and Pattern Recognition*, pp.928-934, June, 1997.
- [6] DiFranco, D.E., Cham, T-J. and Rehg, J.M., "Reconstruction of 3-D Figure Motion from 2-D Correspondences," *Proc. of the 2001 IEEE Computer Society Conf. on Computer Vision and Pattern Recognition*, vol. 1, pp. 307 –314, 2001.
- [7] Elgammal, A., Duraiswami, R., Harwood, D., and Davis, L.S. "Background and Foreground Modeling Using Nonparametric Kernel Density Estimation for Visual Surveillance," *Proc. of the IEEE*, vol. 90, 2002.
- [8] Fablet, R. and Black, M.J. "Automatic Detection and Tracking of Human Motion with a View-Based Representation," *European Conf. On Computer Vision, ECCV'02*, May 2002.
- [9] Gavrilu, D.M. "The Visual Analysis of Human Movement: A Survey," *Computer Vision and Image Understanding*, vol. 73, No. 1, 1999.
- [10] Gavrilu, D.M. and Davis, L.S., "3-D model-based tracking of humans in action: a multi-view approach," *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition*, pp. 73 –80,1996.
- [11] Luck, J., Small, D., and Little, C.Q. "Real-Time Tracking of Articulated Human Models Using a 3D Shape-from-Silhouette Method," *Intl. Workshop on Robot Vision 2001*, LNCS 1998, pp. 19-26.
- [12] Masoud, O. "Tracking and Analysis of Articulated Motion with an Application to Human Motion," Doctoral Thesis, Univ. of Minnesota, March, 2000.
- [13] Masoud, O., Rogers, S., and Papanikolopoulos, N.P., "Monitoring Weaving Sections," ITS Institute Technical Report CTS 01-06, October 2001.
- [14] Matusik, W., C. Buehler, R. Raskar, S. Gortler, and L. McMillan, "Image-Based Visual Hulls," *Proc. of ACM SIGGRAPH 2000*, pp. 369-374, July 2000.
- [15] Ojanen, H., "Automatic correction of lens distortion by using digital image processing," Rutgers University, Dept. of Mathematics technical report, July 1999.
- [16] Oliver, N., Rosario, B., and Pentland, A. "A Bayesian Computer Vision System for Modeling Human Interactions," *IEEE Trans. on Pattern Analysis and Machine Intelligence*, vol.22, no.8 August 2000.
- [17] Polana, R. and Nelson, R. "Nonparametric Recognition of Nonrigid Motion," Technical Report, University of Rochester, New York, 1994.

- [18] Rosales, R. and Sclaroff, S. "3D Trajectory Recovery for Tracking Multiple Objects and Trajectory Guided Recognition of Actions," *Proc. of IEEE Conf. On Computer Vision and Pattern Recognition*, June 1999.
- [19] Seitz, S.M. and Dyer, C.R., "View Morphing," *Proc. of ACM SIGGRAPH 1996*, pp. 21-30, 1996.
- [20] Stauffer, C. and Grimson, W.E. "Learning Patterns of Activity Using Real-Time Tracking," *IEEE Trans. on Pattern Analysis and Machine Intelligence*, vol. 22, No. 8, August, 2000.
- [21] Weik, S. and Liedtke, C.E. "Hierarchical 3D Pose Estimation for Articulated Human Body Models from a Sequence of Volume Data," *Intl. Workshop on Robot Vision 2001*, LNCS 1998, pp. 27-34.
- [22] Wren, C.R. and Pentland, A.P. "Dynamic Models of Human Motion," *Proc. Third IEEE Intl. Conf. Automatic Face and Gesture Recognition*, April, 1998.
- [23] Wren, C.R., Azarbayejani, A., Darrell, T., and Pentland, A.P. "Pfinder: Real-Time Tracking of the Human Body," *IEEE Trans. on Pattern Analysis and Machine Intelligence*, vol. 19, No.7, July 1997.

Appendix A: Transformation Matrices

Homography matrices

H1			H2			H3		
-0.0023	-0.0012	2.5641	-0.0045	0.0020	-0.2152	-0.0057	0.0011	0.9934
-0.0029	0.0006	-0.5462	0.0040	0.0023	-3.1493	0.0020	0.0029	-3.0388
0.0001	-0.0016	-0.0132	0.0000	-0.0025	0.0857	-0.0000	-0.0025	0.0994

Reverse transformations from camera 1

M12			M13		
0.2846	2.6353	-533.9295	-0.5726	2.6609	-278.9510
-0.0002	0.7789	-18.5578	-0.0361	0.8501	-32.0144
0.0013	0.0036	-0.5155	0.0001	0.0064	-1.0269

Forward transformations to camera 1

M21			M31		
-0.8977	-1.5873	987.0013	-1.8359	2.5822	418.2002
-0.0638	1.4318	14.4890	-0.1110	1.6957	-22.7163
-0.0027	0.0062	0.5975	-0.0009	0.0109	-1.0751