# Motion-based Recognition of People in EigenGait Space

Chiraz BenAbdelkader†, Ross Cutler‡,
and Larry Davis†
† University of Maryland, College Park
{chiraz,lsd}@umiacs.umd.edu
‡Microsoft Research
rcutler@microsoft.com

## Abstract

*A motion-based, correspondence-free technique for human gait recognition in monocular video is presented. We contend that the planar dynamics of a walking person are encoded in a 2D plot consisting of the pairwise image similarities of the sequence of images of the person, and that gait recognition can be achieved via standard pattern classification of these plots. We use background modelling to track the person for a number of frames and extract a sequence of segmented images of the person. The self-similarity plot is computed via correlation of each pair of images in this sequence. For recognition, the method applies Principal Component Analysis to reduce the dimensionality of the plots, then uses the k-nearest neighbor rule in this reduced space to classify an unknown person. This method is robust to tracking and segmentation errors, and to variation in clothing and background. It is also invariant to small changes in camera viewpoint and walking speed. The method is tested on outdoor sequences of 44 people with 4 sequences of each taken on two different days, and achieves a classification rate of 77%. It is also tested on indoor sequences of 7 people walking on a treadmill, taken from 8 different viewpoints and on 7 different days. A classification rate of 78% is obtained for near-fronto-parallel views, and 65% on average over all view.*

## 1 Introduction

Recently, gait recognition has received growing interest within the computer vision community, due to its emergent importance as a biometric. The term *gait recognition* is typically used to signify the identification of individuals in image sequences 'by the way they walk'. Gait classification is the recognition of different types of human locomotion, such as running, limping, hopping, etc. Because human ambulation is one form of human movement, gait recognition is closely related to vision methods that detect, track and analyze human movement in general.

Gait recognition research has largely been motivated by Johansson's experiments [19] and the ability of humans to perceive motion from Moving Light Displays (MLDs). In these experiments, human subjects were able to recognize the type of move-

ment of a person solely from observing the 2D motion pattern generated by light bulbs attached to the person. Similar experiments later showed some evidence that the identity of a familiar person ('a friend') [1], as well as the gender of the person [9] might be recognizable from MLDs, though in the latter case a recognition rate of 60% is hardly significantly better than chance (50%).

Despite the agreement that humans can perceive motion from MLDs, there is still no consensus on how humans interpret this MLD-type stimuli (i.e. how it is they use it to achieve motion recognition). Two main theories exist: the first maintains that people use motion information in the MLDs to recover the 3D structure of the moving object (person), and subsequently use the structure for recognition; and the second theory states that motion information is directly used to recognize a motion [7].

The dynamics of gait can be fully characterized via the kinematics of a handful of body landmarks such as limbs and joints [18]. Indeed, one method of motion-based recognition is to first explicitly extract the dynamics of points on a moving object (person). Consider a point $\vec{P}(t) = (x(t), y(t), z(t))$ on a moving object as a function of time $t$. The dynamics of the point can be represented by the phase plot $(\vec{P}(t), d\vec{P}(t)/dt, ...)$. Since we wish to recognize different types of motions (viz. gaits), it is important to know what can be determined from the projection $A$ of $\vec{P}(t)$ onto an image plane, $(u, v) = A(\vec{P})$. Under orthographic projection, and if $\vec{P}(t)$ is constrained to planar motion, the object dynamics are completely preserved up to a scalar factor. That is, the phase space for the point constructed from $(u, v)$ is identical (up to a scalar factor) to the phase space constructed from $\vec{P}(t)$. However, if the motion is not constrained to a plane, then the dynamics are not preserved. Under perspective projection, the dynamics of planar and arbitrary motion are in general not preserved.

Fortunately, planar motion is an important class of motion, and includes "biological motion" [16]. In addition, if the person is sufficiently far from the camera, the camera projection becomes approximately orthographic (with scaling). In this case, and assuming we can accurately track a point $\vec{P}(t)$ in the image plane, then we can completely reconstruct the phase space of the dynamic system (up to a scalar factor). The phase space can then be used directly to classify the object motion (e.g., [6]).

In general, point correspondence is not always possible in realistic image sequences (without the use of special markers), due to

occlusion boundaries, lighting changes, insufficient texture, image noise, etc. However, for classifying motions, we do not necessarily have to extract the complete dynamics of the system; qualitative measures may suffice to distinguish a class of motions from each other. In this paper, we use a correspondence-free image feature for motion-based gait recognition.

Our method maps a sequence of images of a walking person to a *similarity plot* (SP), defined as the matrix of self-similarities between each pair of images of the person in the sequence. We contend that this 2D feature encodes a projection of the planar dynamics of gait, and hence a signature of gait dynamics. The feature vectors used for classification consist of the contiguous square blocks, termed *Units of Self-Similarity* (USS), in the SP of size one gait period each. Our method treats a USS much the same way that the Eigenfaces technique [28] treats a face image; it uses Principal Components Analysis (PCA) to reduce the dimensionality of the feature space, then applies some supervised pattern classification technique (k-nearest neighbor rule in our case) in the reduced feature space for recognition, termed *the Eigengait*.

This method is invariant to background texture and lighting, and clothing, and is robust to segmentation errors. It assumes that people walk on a known plane with constant velocity for about 3-4 seconds (the time it takes to walk 5-8 steps at normal speed), the frame rate is greater than twice the frequency of walking, and the camera is static.

The rest of the paper is organized as follows. Section 2 reviews recent vision literature related to gait recognition. In Section 3 we describe the method in detail. Section 4 describes the experimental methodology and results, and finally in Section 5 we conclude with a brief summary and discussion of future work.

## 2 Related Work

We review vision methods used in detection, tracking and recognition of human movement in general, as they are closely related to gait recognition ([7, 5, 13] are good surveys on this topic). These methods can be divided into two main categories: methods that recover high-level structure of the body and use this structure for motion recognition, and those that directly model how the person moves. We shall describe the latter in more detail as it is more relevant to the gait recognition approach proposed in this paper.

Structure-free methods characterize its motion pattern, without regard to its underlying structure. They can be further divided into two main classes. The first class of methods consider the human action or gait to be comprised of a sequence of poses of the moving person, and recognize it by recognizing a sequence of static configurations of the body in each pose [23, 17, 15]. The second class of methods characterizes the spatiotemporal distribution generated by the motion in its continuum, and hence analyze the spatial and temporal dimensions simultaneously [24, 25, 10, 22, 21, 8].

State-space methods represent human movement as a sequence of static configurations. Each configuration is recognized by learning the appearance of the body (as a function of its color/texture, shape or motion flow) in the corresponding pose. Murase and Sakai [23] describe a template matching method which uses the parametric eigenspace representation as applied in face recognition [28]. Huang et al. [17] use a similar technique, as they apply

PCA to map the binary silhouette of the moving figure to a low dimensional feature space. The gait of an individual person is represented as a cluster in this space, and gait recognition is done by determining if all the input silhouettes belong to this cluster. He and Debrunner [15] recognize individual gaits via an HMM that uses a quantized vector of Hu moments computed from the person's binary silhouette as input.

In spatiotemporal methods the action or motion is characterized via the entire 3D spatiotemporal (XYT) data volume spanned by the moving person in the image. It could for example consist of a sequence of grey-scale images, optical flow images, or binary silhouettes of the person. Of particular interest is the work by Cutler and Davis [8] which is closely related to our method. They use similarity plots to characterize periodicity of human motion, and thereby detect humans in video (and not for gait recognition).

## 3 Method

An overview diagram of the method is shown in Figure 1. An input image sequence is first processed to segment the moving person from the background and track it in each frame. The obtained sequence of blobs are then properly aligned and scaled to a uniform height, to account for detection/tracking errors and any depth changes that occur in non-fronto-parallel walking. A self-similarity plot (SP) of the person is computed by correlating each pair of these blobs, and a set of normalized feature vectors, we call *Units of Self-Similarity*, are then extracted from this SP and used for gait recognition via standard statistical pattern classification technique. In the following two sections, we explain our methods for feature extraction/selection and classification in more detail.
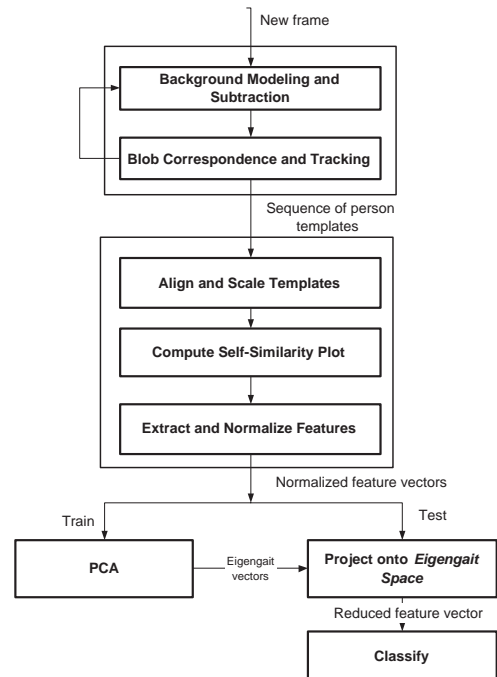


**Figure 1.** Overview of Method.

## 3.1 Feature Extraction

### 3.1.1 Self-similarity Plots

Given a sequence of images obtained from a static camera, we detect and track the moving person, extract an image template corresponding to the person's motion blob in each frame, then compute the self-similarity plot from the obtained sequence of templates. For this, we use the method described in [8], except that we use background subtraction for foreground detection since the camera is assumed to be stationary [12].

Once a person has been tracked for $N$ consecutive frames, the corresponding sequence of $N$ image templates are scaled to a uniform height, as their sizes may vary due to depth changes and segmentation errors, and the self-similarity plot $S$ is obtained by:

$$S(t_1, t_2) = \min_{|dx, dy| < r} \sum_{(x, y) \in B_{t_1}} |O_{t_1}(x + dx, y + dy) - O_{t_2}(x, y)|.$$

where $1 \leq t_1, t_2 \leq N$, $B_{t_1}$ is the bounding box of the person blob in frame $t_1$, $r$ is a small search radius, and $O_{t_1}, O_{t_2}, .., O_{t_N}$ are the scaled templates. Note that these templates can be either (1) foreground images or (2) binary silhouettes, as shown in Figure 2. There are clearly competing tradeoffs to using either type of template in measuring image similarity: the latter is *more* robust to clothing and lighting variations than the former, but is *less* robust to segmentation errors. We shall later compare these two image similarity measures empirically in the experiments (Section 4).



**Figure 2.** From left-to-right: original image, foreground template and binary template, of a walking person.

### 3.1.2 Properties of Self-similarity Plots

The self-similarity plot $S$ of a walking person has some useful properties[8]. For example, the intersections of its off-diagonals and cross-diagonals, which are also its local minima, encode the frequency and phase of walking. Specifically, each intersection corresponds to a combination of the following four key poses of gait: (i) when the two legs are furthest apart and the left leg is leading, (ii) when the two legs are joined together and the right leg is leading, (iii) when the two legs are furthest apart and the left leg is leading, and (iv) when the two legs are joined together and the left leg is leading, as illustrated by Figure 3. We shall denote these poses as $A$, $B$, $C$, and $D$, respectively. Note that diagonals corresponding to $AC$ and $BD$ only exist in the similarity plot of near fronto-parallel walking. Intuitively, this is because poses $A$ and $C$, and poses $B$ and $D$ are very similar in appearance *only* if the person is walking nearly fronto-parallel to the camera and the

person's gait is almost bilaterally symmetrical (i.e. the right and left leg are functionally indistinguishable which is not the case when the person has a limp).

Thus the frequency and phase of gait can be simply computed by finding the local minima of $S$. However, we can only resolve the phase of gait up to half a period, hence poses $A$ and $C$ and poses $B$ and $D$ are indistinguishable. We currently have no way of determining whether the left or right leg is leading, which is a difficult visual problem in itself (it can be achieved in principle by first determining the direction of walking, then carefully segmenting the two legs to determine which leg occludes the other, i.e. which leg is closer to the camera).
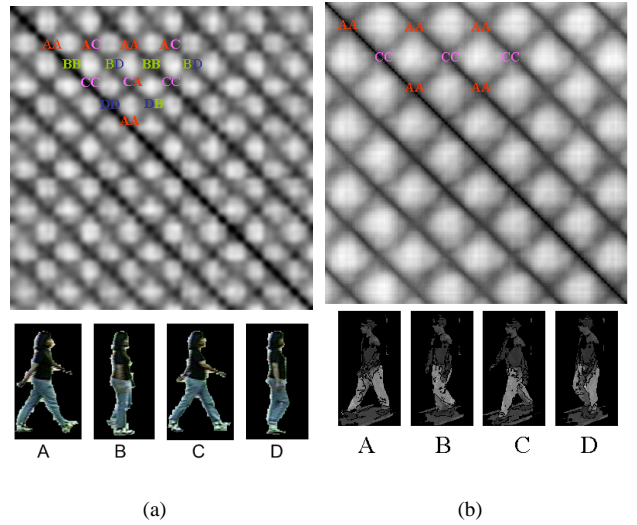


(a)          (b)

**Figure 3.** Combinations of key poses of human gait correspond with local minima of the self-similarity plot.

### 3.1.3 Units of Self-Similarity

Because gait consists of periodic contiguous steps, the similarity plot can be tiled into contiguous rectangular (square for constant-speed walking) blocks, termed *Units of Self-Similarity* (USS), each of which consists of the person's self-similarity over two periods of gait (Figure 4). Clearly a different such tiling is obtained for each starting phase of the periods.

In this paper, we propose to use the set of all USS's starting at key-pose $A$ or $C$ (respectively the blue and green tiles in Figure 4) as our input feature vectors for gait recognition. We only use the USS's in the top half of the similarity plot (the solid tiles), since they are symmetric to the USS's in the bottom half (the broken-line tiles). Hence, for a sequence containing $k$ gait periods ($k = 4$ in the Figure 4), we can extract $2\frac{k(k+1)}{2} = k(k+1)$ such USS's. Note, however, that since we can only resolve the phase of gait up to half a period (as discussed in Section 3.1.2), the USS's starting at pose $A$ are in fact indistinguishable from those starting at pose $C$, which is why we extract both for recognition.
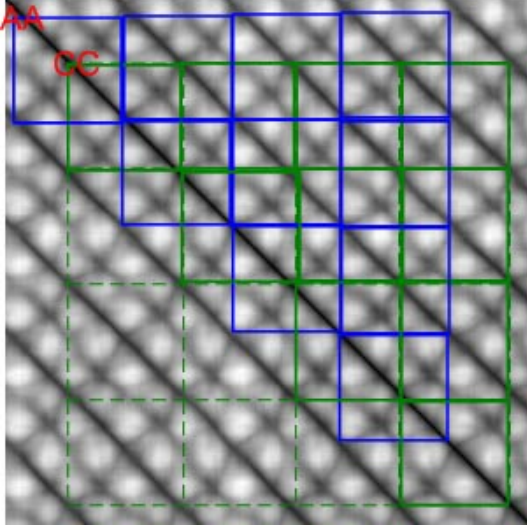
**Figure 4.** Extracting units of self-similarity from the similarity plot. Blue and green USS's start at pose $A$ and $C$, respectively.

### 3.1.4 Normalization

In any pattern classifier, it is important to determine which sources of variation in the input data are irrelevant to classification, and normalize them prior to classification [11]. In our case, a USS of the *same* walking person will vary with: (i) clothing, (ii) the background scene, (iii) number of pixels on target, (iv) camera viewpoint, and (v) walking speed.

By using background subtraction to obtain person templates, we effectively normalize for background variations. A simple way to normalize for variation due to clothing and lighting is by using a color-invariant image similarity measure, such as absolute correlation of binary silhouettes or chamfer matching of the edge maps. However this will not normalize for the style of clothing (for example pants vs. skirts), nor for detection/segmentation errors.

The number of pixels-on-target (POT) is a function of camera depth and image resolution. Assuming the roles of these two are interchangeable, the POT is normalized by scaling down each template such that $h \cdot dpy$ is some fixed constant, where $h$ is the height of the person in the scaled image and $dpy$ is the effective $y$ dimension of a pixel in the frame grabber [27].

Since the size of a USS (equal to one period gait) may vary depending on the actual gait period (speed of walking), we scale them to a uniform $T$x$T$ square tile. This is equivalent to temporal-warping. Note, however, that this does not normalize for the different walking speeds in any *qualitative* way. It only transforms all USS's to feature vectors of equal length ($T^2$-dimensional), to be able to use them as input to the same statistical pattern classifier. Since gait dynamics are at root *not* invariant to speed of walking, there is no (direct) way to normalize for this variation qualitatively.

Similarly, the USS's corresponding to different camera viewpoints are qualitatively different, since a different (planar) projection of gait dynamics is captured in the image plane from any one

camera viewpoint. We currently also have no way of normalizing for this variation.

### 3.2 Gait Classifier

As mentioned in the previous section, the similarity plot is a projection of the dynamics of the walking person that preserves the frequency and phase of the gait. The question then arises as to whether this projection preserves more detailed (higher-dimensional) aspects of gait dynamics, that capture the unique way an individual person walks.

We build a gait pattern classifier that takes USS's as input feature vectors. Our classifier is very much analogous to the 'Eigenface' approach [28], in that we treat a USS much the same way that a face image is treated in that method. Specifically, we apply principal components analysis (PCA) to reduce the dimensionality of the input feature vectors, and use a simple non-parametric pattern classification technique to classify new feature vectors in the subspace spanned by the first few principal components.

#### 3.2.1 Training the Classifier

Let $U_1, U_1, .., U_M$ be a given training set of $M$ labelled (i.e. each corresponding to a known person) USS's, of size $T$x$T$ each, and let $u_i$ be the vector of length $T^2$ corresponding to $U_i$ (obtained by concatenating all its rows). Note, however, that strictly speaking the feature vectors (i.e. the USS's) extracted from the same walking sequence are *not* independent.

We then compute the principal components [20] of the space spanned by $u_1, .., u_M$ by computing the eigenvalue decomposition (also called Karhunen-Loeve expansion) of their covariance matrix $C_u = \frac{1}{M} \sum_{i=1}^{M} (u_i - \bar{u}_i)(u_i - \bar{u}_i)^T$, where $\bar{u}$ is the simple mean of all training vectors $u_1, .., u_M$. This can be efficiently computed in $O(M)$ time (instead of the brute force $O(T^2)$) [28]. The subspace spanned by the $d$ most significant eigenvectors, $v_1, .., v_d$, that account for (say) $95\%$ of the variation in the training vectors, is denoted the *Eigengait*.

#### 3.2.2 Classification

Gait recognition now reduces to standard pattern classification in a $d$-dimensional Eigengait space. Let $x_1, x_2, .., x_l$ be the feature vectors corresponding to the $l$ USS's extracted from a given sequence of an unknown walking person. To classify this sequence, we project each of these vectors $x_i$ to a $d$-dimensional vector in Eigengait space, and determine its class, denoted $c_i$, based on the *k-nearest neighbor rule* [3, 26]. We then decide the class $C$ of the sequence itself as the most frequent $c_i$.

## 4 Experiments

We test our method on four different data sets in order to evaluate its performance across natural variability of individual walking, as well as its sensitivity when other factors are varied: camera viewpoint, walking cadence, image similarity measure, and the KNN parameter $k$. The classifier is trained and tested separately for each combination of the above factors. The person templates

are invariably scaled to a height of 50 pixels before computing the SP, and the USS's are each normalized to a size 32x32, hence spanning a 1024-dimensional feature space. We use the leave-one-out cross-validation to estimate the classification error rate [26, 29].

Since the USS's are not normalized for variation caused by different camera viewpoint and walking cadence, classification is indexed by cadence and camera viewpoint. That is, a different gait classifier is built for each camera viewpoints (for the Keck multiview dataset) and range of cadences (for the CMU MoBo dataset).

## 4.1 Fronto-parallel Datasets

The method is first tested on two different sets of fronto-parallel sequences. The first dataset is the same used by Little and Boyd in [21], and consists of 42 image sequences with six different subjects (4 males and two females) and 7 sequences of each, taken from a static camera. The second datatset contains 108 fronto-parallel sequences taken in an outdoor environment on 2 different days and with 44 different subjects (10 females and 34 males), hence 2 sequences per subject per day. The sequences were captured at 20 fps and a full color resolution of 644x484. Each subject walked a fixed straight path back and forth at their natural pace, as shown in Figure 5.



**Figure 5.** An example of 4 outdoor walking sequences of one person. The top and bottom two sequences were each taken on two different days.

|        | *Little and Boyd* |         | *UMD2 Dataset* |         |
|--------|---------|---------|---------|---------|
| $K_n$  | BC Rate | FC Rate | BC Rate | FC Rate |
| 1      | .93     | .90     | .75     | .77     |
| 3      | .90     | .90     | .72     | .72     |
| 5      | .93     | .87     | .73     | .70     |

**Table 1.** Classification rates for the two fronto-parallel datasets with two different measures of image similarity, Foreground (FC) and Binary (BC), and three different values of $k$ for the KNN classifier, $k = 1, 3, 5$.

Table 1 gives the recognition rates when using the different image similarity measures and values of the KNN parameter $k$. Note

that correlation of binary silhouettes (denoted BC) gave slightly better results than FC for the first dataset, and almost the reverse is true for the second dataset. Also, the performance slightly degrades for higher values of $k$, which maybe because the training points of any one person form multiple clusters in Eigengait space.

## 4.2 Keck Lab Multiview Dataset

Here we test the method on a database consisting of 7 people (3 females and 4 males) walking on a treadmill, taken on 7 different days and captured simultaneously from 8 different cameras. An average of 56 sequences is provided for each subject. The multiple viewpoints correspond to different pan angles of the camera that are at 15 degree intervals and span a range of about 120 degrees of the camera field of regard. Figure 6 illustrates the eight camera viewpoints used in this experiment. The data sequences were captured in the Keck multi-perspective lab at a frame of 60 fps and using greyscale 644x488 images [4].

The treadmill speed was set to match the natural walking pace for each subject, which typically varied between 2.5 and 3.5 miles per hour. The results are shown in Figure 7. Note that performance is significantly better at nearly fronto-parallel views (95 and 115 deg) and that the best performance of 78% is obtained at a 95 deg with $k = 1$ and binary correlation.
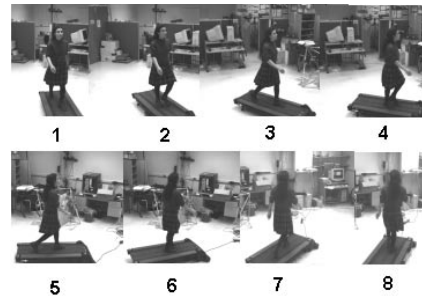


**Figure 6.** Eight camera viewpoints of the sequences in second test data set.
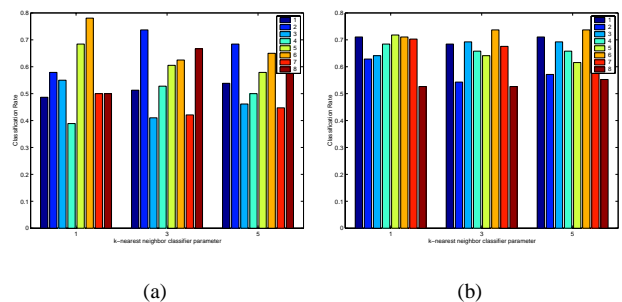


(a)                    (b)

**Figure 7.** Classification rates for Dataset 3 for the 8 viewpoints with $k = 1, 3, 5$ and using (a) Absolute correlation of binary silhouettes (BC). (b) Normalized cross-correlation of foreground images (FC).

### 4.3 CMU MoBo Dataset

To investigate the effect of variation in walking speed, we used an existing dataset [14] of 50 sequences and 25 people walking on a treadmill at a slow (2.06 miles/hr) and moderate (2.82 miles/hr) pace. A recognition rate of 12% was obtained when training on slow-speed sequences and testing on moderate-speed sequences, or vice versa. However, when we both train and test on slow sequences, we obtain a recognition rate of 72%, and 76% for fast sequences. This seems to confirm the expectation that our gait recognition method is sensitive to large changes in walking-speed.

## 5 Conclusions and future work

In this paper, we have used a correspondence-free motion-based method to recognize the gaits of small populations of people. The method is view dependent, and performs best when fronto-parallel images are used. Clothing, lighting, and other variations may degrade the performance of the classifier

When tested with fronto-parallel sequences, the method achieved a recognition rate of 93% on a small dataset of 6 people and 7 sequences each taken on the same day, and 77% on a dataset of 44 subjects and 4 sequences each taken on two different days. The method was also tested on multi-view sequences of 7 people captured from 8 different viewpoints and taken on different days. The best recognition result (78%) was achieved using correlation on binary silhouettes from a near-fronto-parallel viewpoint.

The classification rate has significantly improved compared to the previous version of this method [2], in which we used a (slightly) different feature vector, that consisted of the similarity plot computed over some fixed number of gait periods (typically 3) and starting at some fixed phase. This method achieved a classification rate of at most 28% for the fronto-parallel outdoor dataset (compared to 77% with the new method), and at most 65% for the Keck dataset (compared to 78% with the new method).

We plan to study the sensitivity of this method to changes in camera viewpoint and walking cadence, as well as investigate ways to compensate for this dependence via interpolation techniques in Eigengait space. We are also working to combine Eigengait features obtained from this method with other parametric gait features that can be robustly computed from video, such as cadence, stride length and stature.

### Acknowledgment

## References

[1] C. Barclay, J. Cutting, and L. Kozlowski. Temporal and spatial factors in gait perception that influence gender recognition. *Perception and Psychophysics*, 23(2):145–152, 1978.

[2] C. BenAbdelkader. Gait as a biometric for person identification in video sequences. Technical Report 4289, University of Maryland College Park, 2001.

[3] C. Bishop. *Neural Networks for Pattern Recognition*. Oxford: Clarendon Press, 1995.

[4] E. Borovikov, R. Cutler, T. Horprasert, and L. Davis. Multiperspective analysis of human actions. 1999.

[5] Q. Cai and J. K. Aggarwal. Human motion analysis: a review. In *Proc. of IEEE Computer Society Workshop on Motion of Non-Rigid and Articulated Objects*, 1997.

[6] L. W. Campbell and A. Bobick. Recognition of human body motion using phase space constraints. In *ICCV*, 1995.

[7] C. Cedras and M. Shah. A survey of motion analysis from moving light displays. In *CVPR*, 1994.

[8] R. Cutler and L. Davis. Robust real-time periodic motion detection, analysis and applications. *PAMI*, 13(2), 2000.

[9] J. Cutting and L. Kozlowski. Recognizing friends by their walk: Gait perception without familiarity cues. *Bulletin Psychonomic Soc.*, 9(5):353–356, 1977.

[10] J. W. Davis and A. F. Bobick. The representation and recognition of action using temporal templates. In *CVPR*, 1997.

[11] R. Duda, P. Hart, and D. Stork. *Pattern Classification*. John Wiley and Sons, 2001.

[12] A. Elgammal, D. Harwood, and L. Davis. Non-parametric model for background subtraction. In *ICCV*, 2000.

[13] D. Gavrila. The visual analysis of human movement: a survey. *CVIU*, 73(1), 1999.

[14] R. Gross and J. Shi. The cmu motion of body (mobo) database. Technical report, Robotics Institute, Carnegie Mellon University, 2001.

[15] Q. He and C. Debrunner. Individual recognition from periodic activity using hidden markov models. In *IEEE Workshop on Human Motion*, 2000.

[16] D. Hoffman and B. Flinchbaugh. The interpretation of biological motion. *Biological Cybernetics*, 1982.

[17] P. S. Huang, C. J. Harris, and M. S. Nixon. Comparing different template features for recognizing people by their gait. In *BMVC*, 1998.

[18] V. Inman, H. J. Ralston, and F. Todd. *Human Walking*. Williams and Wilkins, 1981.

[19] G. Johansson. Visual motion perception. *Scientific American*, (232), 1975.

[20] I. T. Joliffe. *Principal Component Analysis*. Springer-Verlag, 1986.

[21] J. Little and J. Boyd. Recognizing people by their gait: the shape of motion. *Videre*, 1(2), 1998.

[22] F. Liu and R. Picard. Finding periodicity in space and time. *ICCV*, 1998.

[23] H. Murase and R. Sakai. Moving object recognition in eigenspace representation: gait analysis and lip reading. *PRL*, 17, 1996.

[24] S. Niyogi and E. Adelson. Analyzing and recognizing walking figures in XYT. In *CVPR*, 1994.

[25] R. Polana and R. Nelson. Detection and recognition of periodic, non-rigid motion. *IJCV*, 23(3), 1997.

[26] B. Ripley. *Pattern Recognition and Neural Networks*. Cambridge University Press, 1996.

[27] R. Tsai. An efficient and accurate camera calibration technique for 3d machine vision. 1986.

[28] M. Turk and A. Pentland. Face recognition using eigenfaces. In *CVPR*, 1991.

[29] S. Weiss and C. Kulikowski. *Computer Systems that Learn*. Morgan Kaufman, 1991.