# Integrating Face and Gait for Human Recognition at a Distance in Video

Xiaoli Zhou and Bir Bhanu, Fellow, IEEE

Abstract—This paper introduces a new video-based recognition method to recognize noncooperating individuals at a distance in video who expose side views to the camera. Information from two biometrics sources, side face and gait, is utilized and integrated for recognition. For side face, an enhanced side-face image (ESFI), a higher resolution image compared with the image directly obtained from a single video frame, is constructed, which integrates face information from multiple video frames. For gait, the gait energy image (GEI), a spatio-temporal compact representation of gait in video, is used to characterize human-walking properties. The features of face and gait are obtained separately using the principal component analysis and multiple discriminant analysis combined method from ESFI and GEI, respectively. They are then integrated at the match score level by using different fusion strategies. The approach is tested on a database of video sequences, corresponding to 45 people, which are collected over seven months. The different fusion methods are compared and analyzed. The experimental results show that: 1) the idea of constructing ESFI from multiple frames is promising for human recognition in video, and better face features are extracted from ESFI compared to those from the original side-face images (OSFIs); 2) the synchronization of face and gait is not necessary for face template ESFI and gait template GEI; the synthetic match scores combine information from them; and 3) an integrated information from side face and gait is effective for human recognition in video.

*Index Terms*—Biometrics fusion, face recognition, gait recognition, video-based recognition.

## I. INTRODUCTION

**I** T HAS BEEN found to be difficult to recognize a person from arbitrary views when one is walking at a distance. For optimal performance, a system should use as much information as possible from the observations. A fusion system, which combines face and gait cues from video sequences, is a potential approach to accomplish the task of human recognition. The general solution to analyze face and gait video data from arbitrary views is to estimate 3-D models. However, the problem of building reliable 3-D models for nonrigid face, with flexible neck and articulated human body from low-resolution video data, remains a hard one. In recent years, integrated face and gait recognition approaches without resorting to 3-D models have achieved some success [1]–[4].

Most current gait-recognition algorithms rely on the availability of the side view of the subject since human gait or the

Color versions of one or more of the figures in this paper are available online at http://ieeexplore.ieee.org.

Digital Object Identifier 10.1109/TSMCB.2006.889612

style of walking is best exposed when one presents a side view to the camera. For face recognition, on the other hand, it is preferred to have frontal views. These conflicting requirements pose some challenges when one attempts to integrate face and gait biometrics in real-world applications. In previous fusion systems [1]–[3], the side view of gait and the frontal view of a face are used. Kale et al. [1] present a gait-recognition algorithm and a face-recognition algorithm based on sequential importance sampling. The database contains video sequences for 30 subjects walking in a single-camera scenario. For face recognition, only the final segment of the database presents a nearly frontal view of face, and it is used as the probe. The gallery consists of static faces for the 30 subjects. Therefore, they perform still-to-video face recognition. Shakhnarovich et al. [2], [3] compute an image-based visual hull from a set of monocular views of multiple cameras. It is then used to render virtual canonical views for tracking and recognition. They discuss the issues of cross-modal correlation and score transformations for different modalities, and present the cross-modal fusion. In their work, four monocular cameras are used to get both the side view of gait and the frontal view of the face simultaneously. Recently, Zhou et al. [4] propose a system which combines cues of face profile and gait silhouette from the video sequences taken by a single camera. It is based on the fact that a side view of a face is more likely to be seen than a frontal view of a face when one exposes the best side view of the gait to the camera. The data are collected for 14 people with two video sequences per person. Even though the face profile in the work of Zhou *et al.* is used reasonably, it only contains shape information of the side view of a face and misses the intensity distribution on the face. In this paper, an innovative video-based fusion system is proposed, aiming at recognizing noncooperating individuals at a distance in a single-camera scenario. Information from two biometrics sources, side face and gait, from the single-camera video sequence, is combined. We distinguish a side face from a face profile. A face profile refers to the outline of the shape of a face as seen from the side. A side face includes not only the outline of the side view of a face, but also the entire side view of the eye, nose, and mouth, possessing both shape and intensity information. Therefore, a side face has more discriminating power for recognition than a face profile.

Table I presents a summary of a related work and compares it with the work presented in this paper. It is difficult to get reliable information of a side face directly from a video frame for a recognition task because of the limited resolution. To overcome this problem, we construct an enhanced side-face image (ESFI), a higher resolution image compared with the image directly obtained from a single video frame, to fuse

Manuscript received May 16, 2006. This paper was recommended by Associate Editor K. Bowyer.

The authors are with the Center for Research in Intelligent Systems, University of California, Riverside, CA 92521 USA.

	TABLE I				
APPROACH FOR INTEGRATING FACE AND	GAIT FOR HUMAN	RECOGNITION	VERSUS THE	PREVIOUS	Work

Features	Kale et al. [1]	Shakhnarovich et al. [2] [3]	Zhou et al. [4]	This Paper
Biometrics	<ul><li>Frontal face</li><li>Gait</li></ul>	<ul><li>Frontal face</li><li>Gait</li></ul>	<ul><li> Face profile</li><li> Gait</li></ul>	<ul><li>Side face</li><li>Gait</li></ul>
Number of Cameras	1	4	1	1
Face Features and Recognition	<ul> <li>Motion vectors</li> <li>Time series model</li> <li>Posterior distribution</li> <li>Maximum a posteriori</li> </ul>	<ul> <li>PCA features of the detected face.</li> <li>k-NN</li> </ul>	<ul> <li>Curvature based features of face profile from the high-resolution image.</li> <li>Dynamic time warping</li> </ul>	<ul> <li>Face features of Enhanced Side Face Image (ESFI)</li> <li>PCA and MDA combined method</li> <li>k-NN</li> </ul>
Gait Features and Recognition	<ul> <li>Entire canonical view image</li> <li>Template matching based on dynamic time warping.</li> </ul>	<ul> <li>Means and standard deviation of moments and centroid.</li> <li>k-NN</li> </ul>	<ul> <li>Entire Gait Energy image (GEI)</li> <li>Template matching</li> </ul>	<ul> <li>Gait features of Gait Energy Image (GEI)</li> <li>PCA and MDA combined method</li> <li>k-NN</li> </ul>
Data	<ul> <li>30 subjects</li> <li>Number of sequences per person: not specified.</li> <li>Static images as the face gallery</li> </ul>	<ul> <li>26 subjects [2]</li> <li>2 to 14 sequences per person [2]</li> <li>12 subjects [3]</li> <li>2 to 6 sequences per person [3]</li> </ul>	<ul><li>14 subjects</li><li>2 sequences per person</li></ul>	<ul> <li>45 subjects</li> <li>2 to 3 sequences per person</li> </ul>
Fusion Methods	<ul><li>Hierarchical fusion</li><li>Sum/Product rule</li></ul>	<ul> <li>Min, Max, Sum and Product rules [2].</li> <li>Sum rule [3]</li> </ul>	<ul><li>Hierarchical fusion</li><li>Sum and Product rules</li></ul>	• Max, Sum and Product rules
Performance Analysis	• No	• No	• No	• Q statistic

the information of a face from multiple video frames. The idea relies on the fact that the temporally adjacent frames in a video sequence, in which one is walking with a side view to the camera, contain slightly different, but unique information about a side face. Experiments show that better face features can be extracted from a constructed ESFI compared to those from the original side-face images (OSFIs).

The contributions of this paper are as follows.

- 1) We present a system that integrates side-face and gait information from video data. The integration of these two biometrics modalities has not been done before.
- 2) Both face- and gait-recognition systems integrate information over multiple frames in a video sequence for improved performance. High-resolution face images are obtained from video, and features from face profile are used for side-face normalization.
- 3) The fusion of side-face and gait biometrics is done at the match score level by obtaining synthetic match scores and using different fusion schemes. Face features and gait features are obtained separately using principal component analysis (PCA) and multiple discriminant analysis (MDA) combined method from the ESFI and the gait energy image (GEI), respectively. The fusion performance is evaluated using the Q statistic.
- 4) Various experiments are performed on 45 people with data from 100 video sequences collected over seven months. Performance comparisons between different biometrics and different fusion methods are presented.

This paper is organized as follows. Section II presents the overall technical approach. It explains the construction of ESFI and describes the generation of GEI. It presents PCA and MDA combined method for feature extraction using ESFI and GEI templates. It introduces an approach to generate synthetic match scores for fusion and provides a description of the classification method. In Section III, a number of dynamic video sequences are tested in three experiments using the approach presented. Experimental results are compared and discussed. Finally, Section IV concludes this paper.

#### **II. TECHNICAL APPROACH**

The overall technical approach is shown in Fig. 1. We first construct the ESFI as the face template and GEI as the gait template from video sequences. During the training procedure, we perform a component and discriminant analysis separately on face templates and gait templates obtained from all the training videos. As a result, transformation matrices and features that form the feature gallery are obtained. During the recognition procedure, each testing video is processed to generate both face templates and gait templates, which are then transformed by the transformation matrices obtained during the training procedure to extract face features and gait features, respectively. These testing features are compared with the gallery features in the database, and then different fusion strategies are used to combine the results of face classifier and gait classifier to improve recognition performance.

# A. ESFI Construction

Multiframe resolution enhancement seeks to construct a single high-resolution image from multiple low-resolution images. These low-resolution images must be of the same object, taken



Fig. 1. Technical approach for integrating side face and gait in video.

from slightly different angles, but not so much as to change the overall appearance of the object in the image.

We use a simple background subtraction method [5] for human-body segmentation. A human body is divided into two parts according to the proportion of its parts [6]: from the top of the head to the bottom of the chin, and then from the bottom of the chin to the bottom of the foot. A head tall is defined as the length from the top of the head to the bottom of the chin. We regard the adult human body as 7.75 head tall. Another 0.25 of one head length is added when the height of hair and the length of the neck are considered. Therefore, the human head cut from the human body in the image should be 1.25 head tall. The ratio of human head (1.25 head) to human body (7.75 head) is 0.16. Therefore, we assume that the upper 16% of the segmented human body includes the human head. In this paper, original low-resolution side-face images are first localized and then extracted by cutting the upper 16% of the segmented human body obtained from multiple video frames.

1) Side-Face Image Alignment: Before multiple lowresolution face images can be fused to construct a highresolution image, motion estimates must be computed to determine pixel displacements between them. It is very important since the quality of a high-resolution image relies on the correctness of low-resolution image alignment. In this paper, the side-face images are aligned using a two-step procedure. In the first step, an elastic registration algorithm [7] is used for motion estimation in low-resolution side-face images. In the second step, a match statistic is introduced to detect and discard images that are poorly aligned. Hence, the quality of constructed high-resolution images can be improved by rejecting such errors.

*Elastic registration method:* Denote f(x, y, t) and  $f(\hat{x}, \hat{y}, t-1)$  as the reference side-face image and the image to be aligned, respectively. Assuming that the image intensities are conserved at different times, the motion between images is modeled locally by an affine transform

$$f(x, y, t) = f(m_1 x + m_2 y + m_5, m_3 x + m_4 y + m_6, t - 1)$$

where  $m_1$ ,  $m_2$ ,  $m_3$ , and  $m_4$  are the linear affine parameters, and  $m_5$  and  $m_6$  are the translation parameters. To account for intensity variations, an explicit change of local contrast and brightness is incorporated into the affine model. Specifically, the initial model takes the form

$$m_7 f(x, y, t) + m_8 = f(m_1 x + m_2 y + m_5,$$
  
 $m_3 x + m_4 y + m_6, t - 1)$  (1)

where  $m_7$  and  $m_8$  are two new (spatially varying) parameters that embody a change in contrast and brightness, respectively. In order to estimate these parameters, the following quadratic error function is minimized:

$$E(\mathbf{m}) = \sum_{x,y\in\Omega} \left[ m_7 f(x,y,t) + m_8 - f(m_1 x + m_2 y + m_5, m_3 x + m_4 y + m_6, t - 1) \right]^2$$
(2)

where  $\mathbf{m} = (m_1, m_2, \dots, m_8)^{\mathrm{T}}$ , and  $\Omega$  denotes a small spatial neighborhood around (x, y). Since this error function is nonlinear in its unknowns, it cannot be minimized analytically. To simplify the minimization, this error function is approximated by using a first-order truncated Taylor series expansion. It now takes the form below.

$$E(\mathbf{m}) = \sum_{x,y\in\Omega} (k - \mathbf{c}^{\mathrm{T}}\mathbf{m})^2$$
(3)

where the scalar k and vector **c** are given as

$$k = f_t - f + x f_x + y f_y$$
  

$$\mathbf{c} = (x f_x \quad y f_x \quad x f_y \quad y f_y \quad f_x \quad f_y \quad -f \quad -1)^{\mathrm{T}} \quad (4)$$

where  $f_x(\cdot)$ ,  $f_y(\cdot)$ , and  $f_t(\cdot)$  are the spatial/temporal derivatives of  $f(\cdot)$ . Minimization of this error function is accomplished by differentiating  $E(\mathbf{m})$ , setting the result equal to zero and solving for  $\mathbf{m}$ . The solution is

$$\mathbf{m} = \left(\sum_{x,y\in\Omega} \mathbf{c}\mathbf{c}^{\mathrm{T}}\right)^{-1} \left(\sum_{x,y\in\Omega} \mathbf{c}k\right).$$
(5)

Align the low-resolution side face image with the reference side face image
<b>Input:</b> the reference side face image and the side face image to be aligned.
<b>Output:</b> the motion vector $\mathbf{m}$ and the match statistic S of the aligned image.
1. For each pyramid level in global registration
1.1 Estimate ${f m}$ between the newest warped image and the reference image using Equation (5)
1.2 Warp the image to the next level of the pyramid using the newest estimate
2. For each pyramid level in local registration
2.1 Estimate <b>m</b> between the newest warped image and the reference image using Equation (5) with $\Omega = 5 \times 5$
2.2 Warp the image using the newest estimate
2.3 For each iteration
2.3.1 Estimate ${f m}$ between the newest warped image and the reference image using Equation (9)
2.3.2 Warp the image using the newest estimate
2.4 Warp the image to the next level of the pyramid using the newest estimate
3. Compute the match statistic $S$ of the aligned image
4. If $S \ge$ threshold, keep the low-resolution image; otherwise, discard it

#### Fig. 2. Pseudocode for low-resolution image alignment.

Intensity variations are typically a significant source of error in differential motion estimation. The addition of the contrast and brightness terms allows us to accurately register images in the presence of local intensity variations. Another important assumption for the model is that the model parameters **m** vary smoothly across space. A smoothness constraint on the contrast/brightness parameters has the added benefit of avoiding a degenerate solution where a pure brightness modulation is used to describe the mapping between images.

To begin, the error function  $E(\mathbf{m})$  in (3) is augmented as follows:

$$\hat{E}(\mathbf{m}) = E_{\rm b}(\mathbf{m}) + E_{\rm s}(\mathbf{m}) \tag{6}$$

where  $E_{\rm b}({\bf m})$  is defined without the summation:

$$E_{\rm b}(\mathbf{m}) = (k - \mathbf{c}^{\rm T} \mathbf{m})^2 \tag{7}$$

with k and c as in (4). The new quadratic error term  $E_s(\mathbf{m})$  embodies the smoothness constraint

$$E_{\rm s}(\mathbf{m}) = \sum_{i=1}^{8} \lambda_i \left[ \left( \frac{\partial m_i}{\partial x} \right)^2 + \left( \frac{\partial m_i}{\partial y} \right)^2 \right] \tag{8}$$

where  $\lambda_i$  is a positive constant that controls the relative weight given to the smoothness constraint on parameter  $m_i$ . This error function is again minimized by differentiating with respect to the model parameters, setting the result equal to zero and solving  $(d\hat{E}(\mathbf{m})/d\mathbf{m}) = (dE_{\rm b}(\mathbf{m})/d\mathbf{m}) + (dE_{\rm s}(\mathbf{m})/d\mathbf{m}) = 0$ . Since solving for  $\mathbf{m}$  at each pixel location yields an enormous linear system which is intractable to solve, an iterative scheme is used to solve for  $\mathbf{m}$  [8]. Now,  $\mathbf{m}$  is expressed as the following iterative equation:

$$\mathbf{m}^{(j+1)} = (\mathbf{c}\mathbf{c}^{\mathrm{T}} + \mathbf{L})^{-1} \left(\mathbf{c}k + \mathbf{L}\overline{\mathbf{m}}^{(j)}\right)$$
(9)

where  $\overline{\mathbf{m}}$  is the componentwise average of  $\mathbf{m}$  over a small spatial neighborhood, and  $\mathbf{L}$  is an 8 × 8 diagonal matrix with diagonal elements  $\lambda_i$ , and zero off the diagonal. On each iteration  $j, \overline{\mathbf{m}}^{(j)}$  is estimated from the current  $\mathbf{m}^{(j)}$ . The initial estimate  $\mathbf{m}^{(0)}$  is estimated from (5).

In this paper, a two-level Gaussian pyramid is constructed for both the reference side-face image and the side-face image to be aligned. The global parameters **m** are first estimated at each pyramid level as in (5) for the entire image. Then, the local parameters **m** are estimated with  $\Omega = 5 \times 5$  as in (5) using the least square algorithm. This estimate of **m** is used to bootstrap the iterations in (9). At each iteration,  $\lambda_i$ ,  $i = 1, \ldots, 8$ , is constant for all **m** components, and its value is set to  $10^{11}$ .  $\overline{m}_i$  is computed by convolving with the  $3 \times 3$ kernel (1 4 1; 4 0 4; 1 4 1)/20. The number of iterations is 10. This process is repeated at each level of the pyramid. The values of these parameters are chosen empirically and based on the previous motion estimation work. Although the contrast and brightness parameters  $m_7$  and  $m_8$  are estimated, they are not used when the side-face image is aligned to the reference side-face image [7].

*Match statistic:* A match statistic is designed to indicate how well a transformed image aligns with the reference image. It is used to select or reject a low-resolution image during alignment. If the size of the reference image is  $M \times N$ , the mean-square error (mse) between the aligned image and the reference image is

$$E = \sum_{x=1}^{M} \sum_{y=1}^{N} \left[ f(x, y, t) - f(m_1 x + m_2 y + m_5, m_3 x + m_4 y + m_6, t - 1) \right]^2 / MN.$$

The match statistic of the aligned image is defined as

$$S = 1 - \frac{E}{\left[\sum_{x=1}^{M} \sum_{y=1}^{N} f^2(x, y, t)\right] / MN}.$$
 (10)

If the value of S is close to 1, the image at time t - 1 is well aligned with the image at time t. A very low value indicates misalignment. A perfect match is 1. However, even images that are very well aligned typically do not achieve 1 due to error in the transformation and noise. For improving the image quality, the resolution-enhancement method discussed next works most effectively when the match values of aligned images are close to 1. A match threshold is specified, and any aligned image whose match statistic falls below the threshold will not be subsequently used.

The pseudocode for the low-resolution image alignment is shown in Fig. 2. Two alignment results with the match statistic S are shown in Fig. 3. The reference images and the images to be aligned are from a video sequence, in which a person is (a)

Fig. 3. Two examples of alignment results with the match statistic S. (a) Well-aligned image with S = 0.95 and (b) badly aligned image with S = 0.86. (Left) Reference image. (Middle) Image to be aligned. (Right) Aligned image.

walking and exposes a side view to the camera. The reference images in Fig. 3(a) and (b) are the same. The time difference between the image to be aligned in Fig. 3(a) and the reference image is about 0.033 s, and the time difference between the image to be aligned in Fig. 3(b) and the reference image is about 0.925 s. The S values are 0.95 and 0.86 for Fig. 3(a) and (b), respectively. Note the differences in the bottom right part of each of the aligned images. We specify the match threshold at 0.9. For 28 out of 100 video sequences used in our experiments, one or two low-resolution images are discarded from each of the sequences during the image-alignment process.

2) *Resolution-Enhancement Algorithm:* An iterative method [9] is used to construct a high-resolution side-face image from the aligned low-resolution side-face images, whose match statistics are above the specified threshold.

*Imaging model:* The imaging process, yielding the observed side-face image sequence  $f_k$ , is modeled by

$$f_k(m,n) = \sigma_k \left( h \left( T_k \left( F(x,y) \right) \right) + \eta_k(x,y) \right)$$
(11)

where

- $f_k$  sensed image of the tracked side face in the kth frame;
- F high-resolution image of the tracked side face in a desired reconstruction view. Finding F is the objective of the superresolution algorithm;
- $T_k$  2-D geometric transformation from F to  $f_k$ , determined by the 2-D motion parameters m of the tracked side face in the image plane, which is obtained in Section II-A1;  $T_k$  is assumed to be invertible and does not include the decrease in the sampling rate between F and  $f_k$ ;
- *h* blurring operator determined by the point spread function (PSF) of the sensor; we use a circular averaging filter with radius 2 as PSF;
- $\eta_k$  additive noise term;
- $\sigma_k$  down sampling operator which digitizes and decimates the image into pixels and quantizes the resulting pixel values.

The receptive field (in F) of a detector whose output is the pixel  $f_k(m, n)$  is uniquely defined by its center (x, y) and its shape. The shape is determined by the region of the blurring operator h, and by the inverse geometric transformation  $T_k^{-1}$ . Similarly, the center (x, y) is obtained by  $T_k^{-1}(m, n)$ . The resolution-enhancement algorithm aims to construct a higher resolution image  $\hat{F}$ , which approximates F as accurately as possible and surpasses the visual quality of the observed images in  $\{f_k\}$ .

Algorithm for resolution enhancement: The algorithm for creating higher resolution images is iterative. Starting with

an initial guess  $F^{(0)}$  for the high-resolution side-face image, the imaging process is simulated to obtain a set of low-resolution side-face images  $\{f_k^{(0)}\}_{k=1}^K$  corresponding to the observed input images  $\{f_k\}_{k=1}^K$ . If  $F^{(0)}$  were the correct high-resolution side-face image, then the simulated images  $\{f_k^{(0)}\}_{k=1}^K$  should be identical to the observed low-resolution side-face image  $\{f_k\}_{k=1}^K$ . The difference images  $\{f_k - f_k^{(0)}\}_{k=1}^K$  are used to improve the initial guess by "back projecting" each value in the difference images onto its receptive field in  $F^{(0)}$ , yielding an improved high-resolution side-face image  $F^{(1)}$ . This process is repeated iteratively to minimize the error function

(b)

$$e^{(n)} = \sqrt{\frac{1}{K} \sum_{k=1}^{K} \left\| f_k - f_k^{(n)} \right\|^2}.$$
 (12)

The imaging process of  $f_k$  at the *n*th iteration is simulated by

$$f_k^{(n)} = \left(T_k(F^{(n)}) * h\right) \downarrow s \tag{13}$$

where  $\downarrow s$  denotes a down sampling operator by a factor *s*, and \* is the convolution operator. The iterative update scheme of the high-resolution image is expressed by

$$F^{(n+1)} = F^{(n)} + \frac{1}{K} \sum_{k=1}^{K} T_k^{-1} \left( \left( \left( f_k - f_k^{(n)} \right) \uparrow s \right) * p \right) \quad (14)$$

where K is the number of low-resolution side-face images.  $\uparrow s$  is an up sampling operator by a factor s, and p is a "back projection" kernel, determined by h.  $T_k$  is 2-D motion parameters. The averaging process reduces the additive noise.

In this paper, we use a sampling factor s = 2. An initial guess  $F^{(0)}$  for the high-resolution image is obtained by up sampling a low-resolution image using bilinear interpolation. Ten low-resolution side-face images contribute to a high-resolution side-face image. The high-resolution image is obtained after ten iterations (N = 10).

The pseudocode for the high-resolution image construction is shown in Fig. 4. Fig. 5 shows four examples of low-resolution face images and reconstructed high-resolution face images. The resolution of the low-resolution side-face images is  $68 \times 68$ , and the resolution of the high-resolution side-face images is  $136 \times 136$ . For comparison, we resize the low-resolution face images using the bilinear interpolation. From this figure, we can see that the quality of the reconstructed high-resolution images is much better than the resized low-resolution images.

Construct the high-resolution side face image from the low-resolution side face images

**Input:** the observed input images  $\{f_k\}_{k=1}^K$  and the corresponding motion vectors  $\{\mathbf{m}_k\}_{k=1}^K$ **Output:** the high-resolution image F.

- 1. Start with iteration n = 0
- 2. Obtain an initial guess  $F^{(0)}$  for the high-resolution image using bilinear interpolation
- 3. Obtain a set of low-resolution images  $\{f_k^{(n)}\}_{k=1}^K$  using Equation (13) 4. Obtain an improved high-resolution image  $F^{(n+1)}$  using Equation (14)
- 5. Let n = n + 1
- 6. If  $n \le N$ , go to step 3; otherwise, stop

Fig. 4. Pseudocode for high-resolution image construction.



Fig. 5. Four examples of (top) resized low-resolution face images and (bottom) constructed high-resolution face images.

3) Side-Face Normalization: Before feature extraction, all high-resolution side-face images are normalized. The normalization is based on the locations of nasion, pronasale, and throat on the face profile. These three fiducial points are identified by using a curvature-based fiducial extraction method [10]. It is explained as follows.

We apply a canny edge detector to the side-face image. After edge linking and thinning, the profile of a side face is extracted as the leftmost points different from background, which contain fiducial points like nasion, pronasale, chin, and throat. The profile consists of a set of points T = (x, y), where x is a row index and y is a column index of a pixel. Then, a Gaussian scale-space filter is applied to this 1-D curve to reduce noise. The convolution between Gaussian kernel  $q(x, \sigma)$  and signal f(x) depends both on x, the signal's independent variable, and on  $\sigma$ , the Gaussian's standard deviation. It is given by

$$F(x,\sigma) = f(x) \oplus g(x,\sigma) = \int_{-\infty}^{\infty} f(u) \frac{1}{\sigma\sqrt{2\pi}} e^{\frac{-(x-u)^2}{2\sigma^2}} du \quad (15)$$

where  $\oplus$  denotes convolution with respect to x. The bigger the  $\sigma$ , the smoother the  $F(x, \sigma)$ . The curve T is parameterized as T(u) = (x(u), y(u)) by the arc length parameter u. An evolved version of T is  $T_{\sigma}(u) = (X(u, \sigma), Y(u, \sigma))$ , where  $X(u, \sigma) =$  $x(u) \oplus g(u, \sigma)$  and  $Y(u, \sigma) = y(u) \oplus g(u, \sigma)$ .

Curvature  $\kappa$  on  $T_{\sigma}$  is computed as

$$\kappa(u,\sigma) = \frac{X_u(u,\sigma)Y_{uu}(u,\sigma) - X_{uu}(u,\sigma)Y_u(u,\sigma)}{(X_u(u,\sigma)^2 + Y_u(u,\sigma)^2)^{1.5}}$$
(16)

where the first and second derivatives of X and Y can be computed as

$$\begin{aligned} X_u(u,\sigma) &= x(u) \oplus g_u(u,\sigma) \quad X_{uu}(u,\sigma) = x(u) \oplus g_{uu}(u,\sigma) \\ Y_u(u,\sigma) &= y(u) \oplus g_u(u,\sigma) \quad Y_{uu}(u,\sigma) = y(u) \oplus g_{uu}(u,\sigma) \end{aligned}$$

where  $g_u(u,\sigma)$  and  $g_{uu}(u,\sigma)$  are the first derivative and the second derivative of the Gaussian kernel.

To localize the fiducial points, the curvature of a profile is first computed at an initial scale, and the locations, where the local maxima of the absolute values occur, are chosen as corner candidates. These locations are tracked down, and the fiducial points are identified at lower scales. The initial scale must be large enough to remove noise and small enough to retain the real corners. Our method has advantages, in that, it does not depend on too many parameters and it does not require any thresholds. It is also fast and simple. The complete process to find the fiducial points is described as follows.

- Step 1) Compute the curvature of a profile at an initial scale, find all points with the large absolute curvature values as corner candidates, and track them down to lower scales.
- Step 2) Regard the rightmost point in the candidate set as the throat.
- Step 3) Regard the pronasale as one of the two leftmost candidate points in the middle part of the profile and then identify it using the curvature value around this point.
- Step 4) Assume that there are no candidate points between pronasale and nasion and identify the first candidate point above the pronasale as nasion.



Fig. 6. Extracted face profile and the absolute values of curvature.



Fig. 7. Examples of four people. (a) Resized OSFIs and (b) ESFIs.

Fig. 6 shows the extracted face profile and the absolute values of curvature. We amplify the absolute values of curvature 20 times in order to show them more clearly. It is clear that the locations of the fiducial points, including nasion, pronasale, and throat, have large curvature values. Given a set of highresolution images and the three fiducial points of each face image, affine transformations are computed between the first image and all the other images. Subsequently, images are cropped as follows. The highest point is defined as the point six pixels above nasion; the lowest point is defined as the throat; the leftmost point is defined as the point four pixels to the left of pronasion; and the rightmost point is defined as the one, which is half of the height of the cropped image and is to the right of the leftmost point. All cropped images are further normalized to the size of  $64 \times 32$ . We call these images as ESFIs. Similarly, OSFI is a subimage from the normalized version of the lowresolution side-face image. It is obtained by the similar process explained above. The size of OSFI is  $34 \times 18$ . Examples of resized OSFIs and ESFIs for four people are shown for comparison in Fig. 7. Clearly, ESFIs have better quality than OSFIs.

## B. GEI Construction

In recent years, various techniques have been proposed for human recognition by gait. These techniques can be divided as model-based and model-free approaches. Little and Boyd [11] describe the shape of the human motion with scale-independent features from moments of the dense optical flow, and recognize individuals by phase vectors estimated from the feature sequences. Sundaresan et al. [12] propose a hidden-Markovmodels-based framework for individual recognition by gait. Huang et al. [13] extend the template matching method to gait recognition by combining the transformation based on canonical analysis and eigenspace transformation for feature selection. Sarkar et al. [14] directly measure the similarity between the testing and training sequences by computing the correlation of corresponding time-normalized frame pairs. Collins *et al.* [15] first extract key frames from a sequence and then compute the similarity between two sequences using the normalized correlation. Tao et al. [16] introduce a set of Gabor-based human-gait appearance models and propose a general tensor discriminant analysis (GTDA) to solve the carrying status in gait recognition. GTDA incorporates the information about the structure of human gait as a constraint. It shows the results only on human carrying conditions.

In this paper, we focus on a model-free approach that does not recover a structural model of human motion. Regular human walking can be considered as a cyclic motion where human motion repeats at a stable frequency [17]. Therefore, it is possible to divide the entire gait sequence into cycles. Since the humanbody segmentation is performed on the original human-walking sequences, we begin with the extracted binary silhouette image sequences. The silhouette preprocessing includes size normalization (proportionally resizing each silhouette image so that all silhouettes have the same height) and horizontal alignment (centering the upper half silhouette part with respect to its horizontal centroid). In a preprocessed silhouette sequence, the time series signal of lower half silhouette size from each frame indicates the gait frequency and phase information. We estimate the gait frequency and phase by a maximum entropy spectrum estimation [11] from the time series signal.

Given the preprocessed binary gait silhouette image  $B_t(x, y)$  at time t in a sequence, the gray-level GEI is defined as follows [17]:

$$G(x,y) = \frac{1}{N} \sum_{t=1}^{N} B_t(x,y)$$
(17)

where N is the number of frames in the complete cycle(s) of a silhouette sequence, t is the frame number of the sequence (moment of time), and x and y are values in the 2-D image coordinate. Fig. 8 shows the sample silhouette images in a gait cycle from two people, and the rightmost images are the corresponding GEIs. As expected, GEI reflects major shapes of silhouettes and their changes over the gait cycle. It accounts for human walking at different speeds. It is referred as the GEI because: 1) each silhouette image is the space-normalized energy image of human walking at this moment; 2) GEI is the time-normalized accumulative energy image of human walking in the complete cycle(s); and 3) a pixel with a higher intensity value in GEI means that human walking occurs more frequently at this position (i.e., with higher energy). GEI has several advantages over the gait representation of a binary silhouette



Fig. 8. Two examples of normalized and aligned silhouette images in a gait cycle. The rightmost images are the corresponding GEIs.

sequence. GEI is not sensitive to incidental silhouette errors in individual frames. We perform a controlled experiment where GEIs of 16 people are constructed with two or three silhouette images removed. The result demonstrates that with the removal of frames, there is no effect on gait-recognition performance. Note that it is different from ESFI construction where the removal of misaligned images is necessary. Moreover, with such a 2-D template, we do not need to consider the time moment of each frame, and the incurred errors can be therefore minimized.

#### C. Human Recognition Using ESFI and GEI

1) Feature Learning Using the PCA and MDA Combined Method: In this paper, a PCA and MDA combined method [18] is applied to face templates, ESFIs, and gait templates, GEIs, separately to get a low-dimensional feature representation for side face and gait. PCA reduces the dimension of the feature space, and MDA automatically identifies the most discriminating features.

Let  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ ,  $\mathbf{x}_k \in \mathbb{R}^N$ , be *n* random vectors representing *n* ESFIs or *n* GEIs, where *N* is the dimensionality of the image. The covariance matrix is defined as  $\Sigma_{\mathbf{x}} = E([\mathbf{x} - E(\mathbf{x})][\mathbf{x} - E(\mathbf{x})]^{\mathrm{T}})$ , where  $E(\cdot)$  is the expectation operator and *T* denotes the transpose operation. The covariance matrix  $\Sigma_{\mathbf{x}}$  can be factorized into the following form:

$$\Sigma_x = \Phi \Lambda \Phi \tag{18}$$

where  $\Phi = [\Phi_1, \Phi_2, \dots, \Phi_N] \in \mathbb{R}^{N \times N}$  is the orthogonal eigenvector matrix of  $\Sigma_x$ ;  $\Lambda = \{\Lambda_1, \Lambda_2, \dots, \Lambda_N\} \in \mathbb{R}^{N \times N}$  is the diagonal eigenvalue matrix of  $\Sigma_x$  with diagonal elements in descending order. One important property of PCA is its optimal signal reconstruction in the sense of minimum mse when only a subset of principal components is used to represent the original signal. An immediate application of this property is the dimensionality reduction

$$\mathbf{y}_k = \mathbf{P}_{\text{pca}}^{\text{T}} \left[ \mathbf{x}_k - E(\mathbf{x}) \right], \qquad k = 1, \dots, n$$
(19)

where  $\mathbf{P}_{pca} = [\mathbf{\Phi}_1, \mathbf{\Phi}_2, \dots, \mathbf{\Phi}_m], m \leq \min(n, N)$ . The lower dimensional vector  $\mathbf{y}_k \in \mathbb{R}^m$  captures the most expressive features of the original data  $\mathbf{x}_k$ .

MDA seeks a transformation matrix W that maximizes the ratio of the between-class scatter matrix  $\mathbf{S}_{\mathrm{B}}$  to the within-class scatter matrix  $\mathbf{S}_{\mathrm{W}} : J(\mathbf{W}) = (|\mathbf{W}^{\mathrm{T}}\mathbf{S}_{\mathrm{B}}\mathbf{W}|/|\mathbf{W}^{\mathrm{T}}\mathbf{S}_{\mathrm{W}}\mathbf{W}|)$ . Suppose that  $\mathbf{w}_{1}, \mathbf{w}_{2}, \ldots, \mathbf{w}_{c}$  and  $n_{1}, n_{2}, \ldots, n_{c}$  denote the classes and the number of images within each class, respectively, with  $n = n_{1} + n_{2} + \cdots + n_{c}$  and  $\mathbf{w} = \mathbf{w}_{1} \cup \mathbf{w}_{2} \cup \cdots$  $\cup \mathbf{w}_{c}. c$  is the number of classes. The within-class scatter matrix is  $\mathbf{S}_{\mathrm{W}} = \sum_{i=1}^{c} \sum_{\mathbf{y} \in \mathbf{w}_{i}} (\mathbf{y} - \mathbf{M}_{i}) (\mathbf{y} - \mathbf{M}_{i})^{\mathrm{T}}$  and the betweenclass scatter matrix is  $\mathbf{S}_{\mathrm{B}} = \sum_{i=1}^{c} n_{i} (\mathbf{M}_{i} - \mathbf{M}) (\mathbf{M}_{i} - \mathbf{M})^{\mathrm{T}}$ , where  $\mathbf{M}_{i} = (1/n_{i}) \sum_{\mathbf{y} \in \mathbf{w}_{i}} \mathbf{y}$  and  $\mathbf{M} = (1/n) \sum_{\mathbf{y} \in \mathbf{w}} \mathbf{y}$  are the means of the class *i* and the grand mean, respectively.  $J(\mathbf{W})$  is maximized when the columns of  $\mathbf{W}$  are the generalized eigenvectors of  $\mathbf{S}_{\mathrm{B}}$  and  $\mathbf{S}_{\mathrm{W}}$  corresponding to the largest generalized eigenvalues in

$$\mathbf{S}_{\mathrm{B}} \mathbf{\Psi}_{i} = \lambda_{i} \mathbf{S}_{\mathrm{W}} \mathbf{\Psi}_{i}. \tag{20}$$

There are no more than c-1 nonzero eigenvalues  $\lambda_i$  and the corresponding eigenvectors  $\Psi_i$ . The transformed feature vector is obtained as follows:

$$\mathbf{z}_{k} = \mathbf{P}_{\text{mda}}^{\text{T}} \mathbf{y}_{k} = \mathbf{P}_{\text{mda}}^{\text{T}} \mathbf{P}_{\text{pca}}^{\text{T}} [\mathbf{x}_{k} - E(\mathbf{x})]$$
$$= \mathbf{Q} [\mathbf{x}_{k} - E(\mathbf{x})], \qquad k = 1, \dots, n$$
(21)

where  $\mathbf{P}_{mda} = [\Psi_1, \Psi_2, \dots, \Psi_r], r \le c - 1$  and  $\mathbf{Q}$  is the overall transformation matrix. We can choose r to perform feature selection and dimensionality reduction. The choice of the range of PCA and the dimension of MDA reflects the energy requirement. We choose the threshold of 99% in eigenvalue energy for eigenvector selection. The lower dimensional vector  $\mathbf{z}_k \in \mathbb{R}^r$ captures the most expressive and discriminating features of the original data  $\mathbf{x}_k$ .

2) Recognition by Integrating ESFI and GEI: We train face templates and gait templates separately for feature extraction. Let {**F**} be the set of all training face templates and  $\mathbf{Q}^{f}$  be the corresponding face transformation matrix. Let {**G**} be the set of all training gait templates and  $\mathbf{Q}^{g}$  be the corresponding gait transformation matrix. Let {**f**<sub>*i*</sub>} be the set of face feature vectors belonging to the *i*th class and {**g**<sub>*i*</sub>} be the set of gait feature vectors belonging to the *i*th class, i = 1, 2, ..., c, where *c* is the number of classes in the gallery. Given a testing video *P*, we follow the procedure explained in Sections II-A and B to generate the set of testing face templates  $\{\mathbf{F}_P\}$  and the set of testing gait templates  $\{\hat{\mathbf{G}}_P\}$ , respectively. The corresponding face and gait feature vector sets are obtained using (21) as follows:

$$\{\hat{\mathbf{f}}_{\mathbf{P}}\}: \hat{\mathbf{f}}_{Pj} = \mathbf{Q}^{\mathrm{f}} \hat{\mathbf{F}}_{Pj}, \qquad j = 1, 2, \dots, n_{\mathrm{f}}$$
$$\{\hat{\mathbf{g}}_{\mathbf{P}}\}: \hat{\mathbf{g}}_{Pj} = \mathbf{Q}^{\mathrm{g}} \hat{\mathbf{G}}_{Pj}, \qquad j = 1, 2, \dots, n_{\mathrm{g}} \qquad (22)$$

where  $n_{\rm f}$  is the number of testing face templates and  $n_{\rm g}$  is the number of testing gait templates.

The Euclidean distance is used as the similarity measure for the face classifier and the gait classifier. From the classifier based on face templates, we obtain

$$D(\hat{\mathbf{f}}_{Pj}, \mathbf{f}_i) = \|\hat{\mathbf{f}}_{Pj} - \mathbf{m}_{fi}\|, \qquad i = 1, 2, \dots, c \quad j = 1, 2, \dots, n_f$$
(23)

where  $\mathbf{m}_{fi} = (1/N_{fi}) \sum_{\mathbf{f} \in \mathbf{f}_i} \mathbf{f}$ , i = 1, 2, ..., c, is the prototype of class *i* for face and  $N_{fi}$  is the number of face feature vectors in  $\{\mathbf{f}_i\}$ . We assign the testing video *P* to class *k* if

$$D(\hat{\mathbf{f}}_{P}, \mathbf{f}_{k}) = \min_{i=1}^{c} \min_{j=1}^{n_{\mathrm{f}}} D(\hat{\mathbf{f}}_{Pj}, \mathbf{f}_{i}).$$
(24)

From the classifier based on gait templates, we obtain

$$D(\hat{\mathbf{g}}_{Pj}, \mathbf{g}_i) = \|\hat{\mathbf{g}}_{Pj} - \mathbf{m}_{gi}\|, \qquad i = 1, 2, \dots, c \quad j = 1, 2, \dots, n_g$$
(25)

where  $\mathbf{m}_{gi} = (1/N_{gi}) \sum_{\mathbf{g} \in \mathbf{g}_i} \mathbf{g}$ , i = 1, 2, ..., c, is the prototype of class *i* for gait and  $N_{gi}$  is the number of gait feature vectors in  $\{\mathbf{g}_i\}$ . We assign the testing video *P* to class *k* if

$$D(\hat{\mathbf{g}}_{P}, \mathbf{g}_{k}) = \min_{i=1}^{c} \min_{j=1}^{n_{g}} D(\hat{\mathbf{g}}_{Pj}, \mathbf{g}_{i}).$$
(26)

Before a combination of the results of face classifier and the results of gait classifier, it is necessary to map the distances obtained from the different classifiers to the same range of values. We use an exponential transformation here. Given that the distance for a probe X are  $D_1, D_2, \ldots, D_c$ , we obtain the normalized match scores as

$$S'_{i} = \frac{\exp(-D_{i})}{\sum_{i=1}^{c} \exp(-D_{i})}, \qquad i = 1, 2, \dots, c.$$
(27)

After normalization, the match scores of face templates and the match scores of gait templates from the same class are fused using different fusion methods. Since face and gait can be regraded as two independent biometrics in our scenario, synchronization is totally unnecessary for them. To take advantage of information for a walking person in video, we use all the possible combinations of face match scores and gait match scores to generate new match scores, which encode information from both face and gait. The new match scores are called synthetic match scores, defined as

$$S_{t}\left(\{\hat{\mathbf{f}}_{P}, \hat{\mathbf{g}}_{P}\}, \{\mathbf{f}_{l}, \mathbf{g}_{l}\}\right) = R\left\{S'(\hat{\mathbf{f}}_{Pi}, \mathbf{f}_{l}), S'(\hat{\mathbf{g}}_{Pj}, \mathbf{g}_{l})\right\}$$
  
$$i = 1, 2, \dots, n_{f} \qquad j = 1, 2, \dots, n_{g}$$
  
$$t = 1, 2, \dots, n_{f} n_{g} \quad l = 1, 2, \dots, c \qquad (28)$$

where S' means the normalized match score of the corresponding distance D and  $R\{,\}$  means a fusion method. In this paper, we use Sum, Product, and Max rules. It is reasonable to generate synthetic match scores using (28), since the ESFI is built from multiple video frames and GEI is a compact spatiotemporal representation of gait in video. In this paper, we use two face match scores and two gait match scores to generate four synthetic match scores for one person from each video.

Distances representing dissimilarity become match scores representing similarity by using (27); therefore, the unknown person would be classified to the class for which the synthetic match score is the largest. We assign the testing video P to class k if

$$S\left(\{\hat{\mathbf{f}}_{P}, \hat{\mathbf{g}}_{P}\}, \{\mathbf{f}_{k}, \mathbf{g}_{k}\}\right) = \max_{l=1}^{c} \max_{t=1}^{n_{f}n_{g}} S_{t}\left(\{\hat{\mathbf{f}}_{P}, \hat{\mathbf{g}}_{P}\}, \{\mathbf{f}_{l}, \mathbf{g}_{l}\}\right).$$
(29)

Since we obtain more than one synthetic match scores after fusion for one testing video sequence, (29) means that the unknown person is classified to the class which gets the maximum synthetic match score out of all the synthetic match scores corresponding to all the classes.

#### **III. EXPERIMENTAL RESULTS**

#### A. Experiments and Parameters

We perform three experiments to test our approach. The data are obtained by a Sony DCR-VX1000 digital video camera recorder operating at 30 frames/s. We collect video sequences of 45 people who are walking in outdoor condition and expose a side view to the camera. The number of sequences per person varies from two to three. The resolution of each frame is  $720 \times 480$ . The distance between people and the video camera is about 10 ft. Each video sequence includes only one person.

In Experiment 1, the data consist of 90 video sequences of 45 people. Each person has two video sequences, one for training and the other one for testing. For the same person, the clothes are the same in the training sequence and the testing sequence. In Experiment 2, the data consist of 90 video sequences of 45 people. Each person has two video sequences, one for training and the other one for testing. For 10 of 45 people, the clothes are different in the training sequences and the testing sequences, and the data are collected on two separate days about one month apart. For the other 35 people, the clothes are the same in the training sequences and the testing sequences. In Experiment 3, we use the same data as in Experiment 2. The difference between them is that we use different number of ESFIs and GEIs in the testing procedure. Table II summarizes the key features of the three experiments.

For gait, we obtain two complete walking cycles from a video sequence according to the gait frequency and gait phase. Each walking cycle includes about 20 frames. We construct two GEIs corresponding to two walking cycles from one video sequence. The resolution of each GEI is  $300 \times 200$ . For face, we also construct two high-resolution side-face images from one video sequence. The match threshold (the match statistic *S*) for aligned low-resolution side-face images is specified at 0.9.

Dut	Experiments		
Data		2	3
Number of subjects	45	45	45
Number of subjects with changed clothes	0	10	10
Number of GEIs for testing per video	2	2	1 or 2
Number of ESFIs for testing per video	2	2	1 or 2

TABLE IISUMMARY OF THREE EXPERIMENTS

Each high-resolution side-face image is built from ten lowresolution side-face images that are extracted from adjacent video frames. The resolution of low-resolution side-face images is  $68 \times 68$ , and the resolution of the reconstructed highresolution side-face images is  $136 \times 136$ . After normalization (see Section II-A3), the resolution of ESFI is  $64 \times 32$ . Recognition performance is used to evaluate our method in the three experiments. For a video sequence, it is defined as the ratio of the number of the correctly recognized people to the number of all the people. To analyze the performance of our method more insightfully, we provide the error index that gives the numbers of misclassified sequences. For comparison, we also show the performance using face features from the OSFIs to demonstrate the performance improvement by using the constructed ESFIs. The resolution of OSFI is  $34 \times 18$ . The procedures of feature extraction, synthetic match score generation, and classification are the same for ESFI and OSFI.

Experiment 1: Figs. 9 and 10 show the data used in Experiment 1. We name 45 people from 1 to 45, and each person has two video sequences. For each of the 45 people, some frames of the training sequence and the testing sequence are shown. Since we construct two GEIs and two ESFIs for each sequence, we totally obtain 90 ESFIs and 90 GEIs as the gallery and another 90 ESFIs and 90 GEIs as the probe. After fusion, as explained in Section II-C2, four synthetic match scores are generated based on two face match scores and two gait match scores for one person from each video. Totally, we have 180 synthetic match scores corresponding to 45 people in the gallery and 180 synthetic match scores corresponding to 45 people in the probe. The dimensionality of PCA features is 72 for GEI, 56 for ESFI, and 65 for OSFI. After MDA, the dimensionality of features is 17 for GEI, 35 for ESFI, and 25 for OSFI. Table III shows the performance of single biometrics. Table IV shows the performance of fusion using different combination rules. In Tables III and IV, the error index gives the number of misclassified sequences.

From Table III, we can see that 73.3% people are correctly recognized by OSFI (12 errors out of 45 people), 91.1% people are correctly recognized by ESFI (four errors out of 45 people), and 93.3% people are correctly recognized by GEI (three errors out of 45 people). Among the three people misclassified by GEI, the person (26) has a backpack in the testing sequence, but not in the training sequence. The difference causes the body shape to change enough to make a recognition error. The

changes of the walking style for the other two people (4, 15) also cause the recognition errors. We show GEIs of the people who are misclassified by the gait classifier in Fig. 11. Among the performances of fusion methods in Table IV, Max rule based on ESFI and GEI performs the best at the recognition rate of 97.8% (one error out of 45 people), followed by Sum rule and Product rule at 95.6% (two errors out of 45 people). For fusion based on OSFI and GEI, the best performance is achieved by Product rule at 95.6%, followed by Sum rule and Max rule at 93.3%. It is clear that the fusion based on ESFI and GEI always has better performance than the fusion based on OSFI and GEI, except using Product rule where they are the same. Fig. 12 shows the people (video sequences) misclassified by integrating ESFI and GEI using different fusion rules. It is clear that the only person (26) who is misclassified by the Max rule has a backpack in the testing sequence that does not occur in the training sequence. This difference makes both the gait classifier and the fused classifier unable to recognize him.

Experiment 2: The data used in Experiment 2 are obtained by substituting ten testing video sequences of Experiment 1 with the other ten testing video sequences shown in Fig. 13. We use the same order as in Experiment 1 to name 45 people. Compared with the data in Experiment 1, the ten replaced testing video sequences are {1, 2, 5, 6, 8, 9, 10, 13, 19, 40}. Therefore, 10 out of 45 people in Experiment 2 wear different clothes in the training sequences and the testing sequences, and for each of the ten people, two video sequences are collected on two separate days about one month apart. We construct two GEIs and two ESFIs from each sequence, so we totally obtain 90 ESFIs and 90 GEIs as the gallery and another 90 ESFIs and 90 GEIs as the probe for 45 people. After fusion, as explained in Section II-C2, we have 180 synthetic match scores corresponding to 45 people in the gallery and 180 synthetic match scores corresponding to 45 people in the probe. The dimensionality of PCA features is 72 for GEI, 56 for ESFI, and 65 for OSFI. After MDA, the dimensionality of features is 17 for GEI, 35 for ESFI, and 25 for OSFI. Table V shows the performance of individual biometrics. Table VI shows the performance of fusion using different combination rules. In Tables V and VI, the error index gives the number of misclassified sequence.

From Table V, we can see that 64.4% people are correctly recognized by OSFI (16 errors out of 45 people), 80% people are correctly recognized by ESFI (nine errors out of 45 people), and 82.2% people are correctly recognized by GEI (eight errors out of 45 people). Compared with the performance of individual biometrics in Experiment 1 in Table III, all the performance of individual biometrics in Experiment 2 decreases to some extent. It is reasonable since gait recognition based on GEI is not only affected by the walking style of a person, but also by the shape of a human body. Changing clothes causes the difference in the shape of the training sequence and the testing sequence for the same person. Also, the lighting condition and the color of clothes cause human-body segmentation to be inaccurate. Fig. 14 shows GEIs of the people who are misclassified by the gait classifier. Meanwhile, since the face is sensitive to noise as well as facial expressions, the different condition in the two video sequences that are taken one month apart brings facerecognition errors. Fig. 15 shows ESFIs of the people who



Fig. 9. Data in Experiment 1. Video sequences from number 1 to 23.

are misclassified by the face classifier. Note the differences in the training and testing GEIs and ESFIs in Figs. 14 and 15. From Table VI, we can see that when ESFI and GEI are fused using appropriate fusion methods, the performance improves. Specifically, Sum rule and Max rule based on ESFI and GEI perform the best at the recognition rate of 88.9% (five errors



Fig. 10. Data in Experiment 1. Video sequences from number 24 to 45.

out of 45 people), and the performance improvement is 6.7% compared with that of the gait classifier. Fig. 16 shows the people (video sequences) misclassified by integrating ESFI and GEI using different fusion rules. For fusion based on OSFI and GEI, there is no improvement compared with the indi-

vidual classifier. These results demonstrate the importance of constructing the ESFI. From ESFI, we can extract face features with more discriminating power. Therefore, the performance improvement is still achieved when ESFI instead of OSFI is used for fusion.

 TABLE III

 Experiment 1: Single-Biometrics Performance and Error Index of Individuals

Derfermen	Biometric				
Performance	Original Face (OSFI)	Enhanced Face (ESFI)	Gait (GEI)		
Recognition Rate	73.3%	91.1%	93.3%		
Error Index	1, 6, 10, 12, 14, 18, 20, 22, 26, 28, 42, 43	13, 16, 21, 35	4, 15, 26		

 TABLE IV

 EXPERIMENT 1: FUSED BIOMETRICS PERFORMANCE AND ERROR INDEX OF INDIVIDUALS

Fusion Method		Sum Rule	Product Rule	Max Rule
OSFI &	Recognition Rate	93.3%	95.6%	93.3%
GEI	Error Index	4, 10, 26	4, 26	4, 10, 26
ESFI &	Recognition Rate	95.6%	95.6%	97.8%
GEI	Error Index	4, 26	4, 26	26



Fig. 11. Experiment 1: GEIs of people misclassified by the gait classifier (see Table III). For each person, two GEIs of the training video sequence and two GEIs of the testing video sequence are shown for comparison.



Fig. 12. Experiment 1: People misclassified by the integrated classifier based on ESFI and GEI using different fusion rules (see Table IV). For each person, one frame of the training video sequence and one frame of the testing video sequence are shown for comparison. (a) Errors by Sum rule. (b) Errors by Product rule. (c) Errors by Max rule.

*Experiment 3:* The data used in Experiment 3 are the same as the data used in Experiment 2. Experiment 3 studies the effect of using the different number of GEIs and ESFIs in the testing procedure. In the gallery, we still use two GEIs and two ESFIs for each of the 45 people. While for the probe, we vary the number of GEIs and ESFIs for each person. Table VII shows the performance of fusion by different combination rules when

the different number of GEIs and ESFIs is used. Except the performance of fusion using two GEIs and two ESFIs, which is obtained from Experiment 2, the other performance is the average value on different combination of GEI and ESFI.

From Table VII, it is clear that if more GEIs and ESFIs are used, i.e., more information in video sequences is used, better performance can be achieved. Meanwhile, this experiment



Fig. 13. Data in Experiment 2: Ten updated video sequences {1, 2, 5, 6, 8, 9, 10, 13, 19, 40}.

 TABLE
 V

 EXPERIMENT 2: SINGLE-BIOMETRICS PERFORMANCE AND ERROR INDEX OF INDIVIDUALS

Daufammanaa	Biometric						
Original Face (OSFI)		Enhanced Face (ESFI)	Gait (GEI)				
Recognition Rate	64.4%	80%	82.2%				
Error Index	1, 2, 5, 6, 8, 9, 13, 18, 19, 20, 26, 28, 34, 40, 42, 43	1, 2, 5, 8, 11, 13, 30, 35, 42	2, 5, 6, 8, 13, 19, 26, 40				

 TABLE
 VI

 EXPERIMENT 2: FUSED BIOMETRICS PERFORMANCE AND ERROR INDEX OF INDIVIDUALS

Fusion Method		Sum Rule	Product Rule	Max Rule
OSFI &	Recognition Rate	82.2%	82.2%	82.2%
GEI	Error Index	2, 5, 6, 8, 13, 19, 26, 40	2, 5, 6, 8, 13, 19, 26, 40	2, 5, 6, 8, 13, 19, 26, 40
ESFI &	Recognition Rate	88.9%	82.2%	88.9%
GEI	Error Index	2, 5, 6, 8, 13	2, 5, 6, 8, 13, 19, 26, 40	2, 5, 6, 8, 13

shows that our method to generate the maximum number of synthetic match scores is a reasonable way to use all the available information.

# B. Performance Analysis

1) Discussion on Experiments: From Experiments 1 and 2, when ESFI and GEI are used, we can see that Max rule

always achieves the best fusion performance, Sum rule has the same fusion performance as the Max rule in Experiment 2, and Product rule does not achieve performance improvement after fusion.

When we compare Experiment 1 and Experiment 2, it can be seen that the recognition rates in Experiment 2 decrease compared with Experiment 1, since 10 out of 45 people change their clothes in the testing sequences. As explained before, gait



Fig. 14. Experiment 2: GEIs of people misclassified by the gait classifier (see Table V). For each person, two GEIs of the training video sequence and two GEIs of the testing video sequence are shown for comparison.



Fig. 15. Experiment 2: ESFIs of people misclassified by the face classifier (see Table V). For each person, two ESFIs of the training video sequence and two ESFIs of the testing video sequence are shown for comparison.

recognition based on GEI is not only affected by the walking style of a person, but also by the shape of human body. Face is sensitive to noise as well as facial expressions; therefore, the different condition in the training sequence and the testing sequence affects its reliability. All these factors contribute to recognition errors of the individual classifiers. However, the fusion system based on side face and gait overcomes this problem to some extent. In Experiment 2, there are some people who are not correctly recognized by gait, but when side-face information is integrated, the recognition rate is improved. It is because the clothes or the walking style of these people are much different between the training and testing video sequences, so the gait classifier cannot recognize them correctly. However, the side face of these people does not change so



Fig. 16. Experiment 2: People misclassified by the integrated classifier based on ESFI and GEI using different fusion rules (see Table VI). For each person, one frame of the training video sequence and one frame of the testing video sequence are shown for comparison.

Fusion Method	1 GEI & 1 ESFI	1 GEI & 2 ESFI	2 GEI & 1 ESFI	2 GEI & 2 ESFI
Sum Rule	82.8%	84.4%	84.4%	88.9%
Product Rule	77.2%	81.1%	82.2%	82.2%
Max Rule	81.1%	80%	84.4%	88.9%

 TABLE
 VII

 EXPERIMENT 3: FUSED BIOMETRICS PERFORMANCE USING DIFFERENT NUMBER OF GEI AND ESFI

much in the training and testing sequences, and it brings useful information for the fusion system and corrects some errors. Specifically, in Experiment 2, the gait classifier misclassifies eight people {2, 5, 6, 8, 13, 19, 26, 40}, and after fusion with ESFI using Sum rule or Max rule, three errors {19, 26, 40} are corrected. On the other hand, since the face classifier is comparatively sensitive to the variation of facial expressions and noise, it cannot get a good recognition rate by itself. When the gait information is combined, the better performance is achieved.

Our experimental results demonstrate that the fusion system using side face and gait has the potential since face and gait are two complementary biometrics. Compared with gait, face images are readily interpretable by humans, which allows people to confirm whether a computer system is functioning correctly, but the appearance of a face depends on many factors: incident illumination, head pose, facial expressions, moustache/beard, eyeglasses, cosmetics, hair style, weight gain/loss, aging, and so forth. Although gait images can be easily acquired from a distance, the gait recognition is affected by clothes, shoes, carrying status, and specific physical condition of an individual. The fusion system is relatively more robust compared with the system that uses only one biometrics. For example, face recognition is more sensitive to low lighting conditions, whereas gait is more reliable under these conditions. Similarly, when the walker is carrying a heavy baggage or he/she is injured, the captured face information may contribute more than gait.

In Experiment 1, the gait recognition misclassifies three people and achieves the recognition rate of 93.3%. The fusion by using Max rule performs best at 97.8% with one error, followed by Sum rule and Product rule at 95.6% with two errors. It may seem that the improvement is not significant in the number of people because of the size of our database. In Experiment 2, where ten of the subjects wear different clothes in the training data and the testing data, the performance of gait recognition decreases to 82.2% with eight errors. For this more difficult database, there is a larger improvement in fusion performance. The Sum rule and Max rule have an improvement of 6.7% with the fusion performance at 88.9% with five errors. These results demonstrate the effectiveness of integrating ESFI and GEI for human recognition since the proposed fusion system still achieves improvement, even a larger improvement for the more challenging database.

The experimental results in Experiments 1 and 2 clearly demonstrate the importance of constructing ESFI. From ESFI, we can extract face features with more discriminating power. Therefore, better performance is achieved when ESFI instead of OSFI is used for both of the individual classifier and the fused classifier. For example, in Experiment 2, OSFI has bad performance at 64.4%, but ESFI still achieves the recognition rate of 80%. Fusion based on ESFI and GEI achieves the performance improvement of 6.7% (from 82.2% to 88.9%) using the Sum rule and Max rule, while there is no performance improvement by fusion of OSFI and GEI using any combination rule (see Table VI). Furthermore, from Experiment 3, we can see that more information means better performance. This also explains why the ESFI always performs better than the OSFI, since ESFI fuses information from multiple frames.

These results also demonstrate that the match score fusion cannot rectify the misclassification achieved by both of the face classifier and the gait classifier. People misclassified by the individual classifiers are likely to be classified correctly after fusion on the condition that there is at least one of the two classifiers that works correctly. For example, in Table V, there are four misclassified people {2, 5, 8, 13} overlapped between classification using ESFI only and GEI only. There are eight misclassified people  $\{2, 5, 6, 8, 13, 19, 26, 40\}$ overlapped between classification using OSFI only and GEI only. From Table VI, we can see that the set of misclassified people {2, 5, 8, 13} is always a subset of the error indexes when ESFI and GEI are combined by any fusion rule. Similarly, the set of misclassified people  $\{2, 5, 6, 8, 13, 19, 26, 40\}$  is always a subset of the error indexes when OSFI and GEI are combined by any fusion rule. It is also the reason that the fusion performance based on OSFI and GEI can never be better than the performance of the gait classifier.

2) Performance Characterization Statistic Q: For the performance improvement by fusion compared with the individual biometrics, if the different classifiers misclassify features for the same person, we do not expect as much improvement as in the case where they complement each other [19]. We use a statistic to demonstrate this point. There are several methods to assess the interrelationships between the classifiers in a classifier ensemble [20], [21]. Given classifiers *i* and *j* corresponding to feature vectors  $f_i$  and  $f_j$  from the same person, respectively, we compute the Q statistic

$$Q_{i,j} = \frac{N^{11}N^{00} - N^{01}N^{10}}{N^{11}N^{00} + N^{01}N^{10}}$$
(30)

where  $N^{00}$  is the number of misclassification by both i and j;  $N^{11}$  is the number of correct classification by both i and j;  $N^{10}$  and  $N^{01}$  are the numbers of misclassification by i or j, but not by both. It can be easily verified that  $-1 \le Q \le 1$ . The Q-value can be considered as a correlation measure between the classifier decisions. The best combination is the

TABLE VIIIEXPERIMENT 1: Q STATISTICS

Fused Templates	$N^{11}$	$N^{00}$	$N^{01}$	$N^{10}$	Q Statistic
OSFI & GEI	31	1	11	2	0.1698
ESFI & GEI	38	0	4	3	-1

TABLE IXEXPERIMENT 2: Q STATISTICS

Fused Templates	$N^{11}$	$N^{00}$	$N^{01}$	$N^{10}$	Q Statistic
OSFI & GEI	29	8	8	0	1
ESFI & GEI	32	4	5	4	0.7297

one that minimizes the value of Q statistic, which means that the smaller the Q-value is, the greater is the potential for performance improvement by fusion.

Tables VIII and IX show the Q-values in Experiments 1 and 2.  $N^{01}$  is defined as the number of people misclassified by the face classifier but correctly recognized by the gait classifier.  $N^{10}$  is defined as the number of people misclassified by the gait classifier but correctly recognized by the face classifier. The Q-value based on OSFI and GEI in Experiment 2 is 1, which means that the performance improvement by fusion will be zero. The experimental results in Table VI verify it. The Q-value based on OSFI and GEI in Experiment 1 is 0.1698, which explains the fact that the fusion performance increases to 95.6% when Product rule is used (see Table IV). When we compare the Q-values between the fusion of OSFI and GEI and fusion of ESFI and GEI, the results show that the Q-values based on ESFI and GEI are always smaller than the Q-values based on OSFI and GEI in both of the experiments. It indicates that the expected performance improvement using ESFI and GEI is higher than using OSFI and GEI. For example, in Experiment 1, the Q-value based on fusion of ESFI and GEI is -1 and the Q-value based on fusion of OSFI and GEI is 0.1698. The maximum performance increase is 4.5% (from 93.3% to 97.8%) by fusion of ESFI and GEI, while the performance increase by fusion of OSFI and GEI is only 2.3% (from 93.3%) to 95.6%). On the other hand, even though the Q-value of 0.7297 for fusion performance of ESFI and GEI is smaller than the Q-value of 1 for fusion performance of OSFI and GEI in Experiment 2, it is positive and relatively high. This indicates that many times the gait classifier and the face classifier are both performing correct classification or incorrect classification for the same person. In spite of this, our video-based fusion method using ESFI and GEI always achieves better performance than either of the individual classifier when the appropriate fusion strategy is used.

To visualize the correlation of the face classifier and the gait classifier, we plot the normalized match scores of the two classifiers. Fig. 17 shows the correlation of the normalized match scores of the two classifiers in Experiment 1. We can see that the match scores of the gait classifier using GEI and the face classifier using OSFI are more correlated than the match scores of the gait classifier using GEI and the face classifier using GEI



Fig. 17. Experiment 1: (a) Correlation of the normalized match scores of the two classifiers using GEI and OSFI. (b) Correlation of the normalized match scores of the two classifiers using GEI and ESFI.



Fig. 18. Experiment 2: (a) Correlation of the normalized match scores of the two classifiers using GEI and OSFI. (b) Correlation of the normalized match scores of the two classifiers using GEI and ESFI.

using ESFI. It is consistent with the Q statistics in Table VIII. Fig. 18 shows the correlation of the normalized match scores of the two classifiers in Experiment 2. We can see that the match scores of the gait classifier using GEI and the face classifier using OSFI are more correlated than the match scores of the gait classifier using GEI and the face classifier using ESFI. It is also consistent with the Q statistics in Table IX. When we compare Figs. 17 and 18, it is clear that the correlation of the match scores of the two classifiers in Experiment 2 is higher than in Experiment 1.

# IV. CONCLUSION

This paper proposes an innovative video-based fusion system, which aims at recognizing noncooperating individuals at a distance in a single-camera scenario. Information from two biometrics sources, side face, and gait, is combined using different fusion methods. Side face includes the entire side views of eye, nose, and mouth, possessing both shape information and intensity information. Therefore, it has a more discriminating power for recognition than face profile. To overcome the problem of a limited resolution of a side face at a distance in video, we use ESFI, a higher resolution image constructed from multiple video frames instead of OSFI directly obtained from a single video frame, as the face template for an individual. ESFI serves as a better face template than OSFI since it generates more discriminating face features. Synthetic match scores are generated for fusion based on the characteristics of face and gait. The experimental results show that the integration of information from side face and gait is effective for individual recognition in video. The performance improvement is always archived when appropriate fusion rules, such as the Max rule

and the Sum rule, are used to integrate information from ESFI and GEI. Consequently, our fusion system is relatively robust compared with the system using only one biometrics in the same scenario.

However, our system has some limitations: 1) gait recognition based on GEI is affected by the shape of human body to some extent; 2) the side face contains less information compared with the frontal face, and it is sensitive to noise as well as facial expressions; and 3) the system has been tested on limited video sequences. In spite of these limitations, we demonstrate that the integration of face and gait can achieve better recognition performance at a distance in video. Although our database is not very big, it is of reasonable size (45 people with 100 video sequences) and shows how the proposed ideas work. In the future, we will collect more data to evaluate the performance of our system. We will also focus on problems that are not addressed in this paper. We will use multiple cameras to capture different views of a person. To get face images with high quality, we will actively track the whole human body first and then zoom in to get better face images. We will speed up the process of ESFI and GEI constructions so that our system can operate in real time.

## REFERENCES

- A. Kale, A. Roy-chowdhury, and R. Chellappa, "Fusion of gait and face for human identification," in *Proc. Acoust., Speech, and Signal Process.*, 2004, vol. 5, pp. 901–904.
- [2] G. Shakhnarovich and T. Darrell, "On probabilistic combination of face and gait cues for identification," in *Proc. Automat. Face Gesture Recog.*, 2002, pp. 169–174.
- [3] G. Shakhnarovich, L. Lee, and T. Darrell, "Integrated face and gait recognition from multiple views," in *Proc. Comput. Vis. and Pattern Recog.*, 2001, vol. 1, pp. 439–446.
- [4] X. Zhou, B. Bhanu, and J. Han, "Human recognition at a distance in video by integrating face profile and gait," in *Proc. Audio and Video-Based Biometric Person Authentication*, 2005, pp. 533–543.
- [5] J. Han and B. Bhanu, "Performance prediction for individual recognition by gait," *Pattern Recognit.*, vol. 26, no. 5, pp. 615–624, Apr. 2005.
- [6] P. A. Hewitt and D. Dobberfuhl, "The science and art of proportionality," *Sci. Scope*, vol. 27, no. 4, pp. 30–31, Jan. 2004.
- [7] S. Periaswamy and H. Farid, "Elastic registration in the presence of intensity variations," *IEEE Trans. Med. Imag.*, vol. 22, no. 7, pp. 865–874, Jul. 2003.
- [8] B. K. P. Horn, Robot Vision. Cambridge, MA: MIT Press, 1986.
- [9] M. Irani and S. Peleg, "Motion analysis for image enhancement: Resolution, occlusion and transparency," J. Vis. Commun. Image Represent., vol. 4, no. 4, pp. 324–335, Dec. 1993.
- [10] B. Bhanu and X. Zhou, "Face recognition from face profile using dynamic time warping," in *Proc. Int. Conf. Pattern Recog.*, 2004, vol. 4, pp. 499–502.
- [11] J. J. Little and J. E. Boyd, "Recognizing people by their gait: The shape of motion," *Videre: J. Comput. Vis. Res.*, vol. 1, no. 2, pp. 1–32, 1998.
- [12] A. Sundaresan, A. Roy-Chowdhury, and R. Chellappa, "A hidden Markov model based framework for recognition of humans from gait sequences," in *Proc. Int. Conf. Image Process.*, 2003, vol. 2, pp. 93–96.
- [13] P. S. Huang, C. J. Harris, and M. S. Nixon, "Recognizing humans by gait via parametric canonical space," *Artif. Intell. Eng.*, vol. 13, no. 4, pp. 359–366, Oct. 1999.
- [14] S. Sarkar, P. J. Phillips, Z. Liu, I. R. Vega, P. Grother, and K. W. Bowyer, "The human ID gait challenge problem: Data sets, performance, and analysis," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 27, no. 2, pp. 162–177, Feb. 2005.
- [15] R. T. Collins, R. Gross, and J. Shi, "Silhouette-based human identification from body shape and gait," in *Proc. IEEE Conf. Autom. Face Gesture Recog.*, 2002, pp. 351–356.
- [16] D. Tao, X. Li, X. Wu, and S. Maybank, "Human carrying status in visual surveillance," in *Proc. IEEE Conf. Comput. Vis. and Pattern Recog.*, 2006, vol. 2, pp. 1670–1677.

- [17] J. Han and B. Bhanu, "Individual recognition using gait energy image," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 28, no. 2, pp. 316–322, Feb. 2006.
- [18] P. N. Belhumeur, J. P. Hespanha, and D. J. Kriegman, "Eigenfaces vs. fisherfaces: Recognition using class specific linear projection," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 19, no. 7, pp. 711–720, Jul. 1997.
- [19] T. Kinnune, V. Hautamaki, and P. Franti, "Fusion of spectral feature sets for accurate speaker identification," in *Proc. Int. Conf. Speech and Comput.*, Sep. 2004, pp. 361–365.
- [20] C. A. Shipp and L. I. Kuncheva, "Relationships between combination methods and measures of diversity in combining classifiers," *Inf. Fusion*, vol. 3, no. 2, pp. 135–148, Jun. 2002.
- [21] G. Brown, J. Wyatt, R. Harris, and X. Yao, "Diversity creation methods: A survey and categorisation," *Inf. Fusion*, vol. 6, no. 1, pp. 5–20, Mar. 2005.



Xiaoli Zhou received the B.S. and M.S. degrees in electrical engineering from the Beijing University of Posts and Telecommunications, Beijing, China, in 1998 and 2001, respectively. She is currently working toward the Ph.D. degree at the University of California, Riverside (UCR).

She is working on her research with the Center for Research in Intelligent Systems, UCR. Her research interests are in computer vision, pattern recognition, and image processing. Her recent research has been concerned with fusion of biometrics for human recognition at a distance in video.



**Bir Bhanu** (S'72–M'82–SM'87–F'96) received the S.M. and E.E. degrees in electrical engineering and computer science from the Massachusetts Institute of Technology, Cambridge, the Ph.D. degree in electrical engineering from the University of Southern California, Los Angeles, and the M.B.A. degree from the University of California, Irvine.

He was the Founding Professor of electrical engineering and served its first Chair with the University of California, Riverside (UCR). He has been the Cooperative Professor of computer science and

engineering and the Director with the Visualization and Intelligent Systems Laboratory (VISLab) since 1991. Currently, he also serves as the Founding Director of an interdisciplinary Center for Research in Intelligent Systems (CRIS) with UCR. Previously, he was a Senior Honeywell Fellow with Honeywell Inc., Minneapolis, MN. He has been with the faculty of the Department of Computer Science, University of Utah, Salt Lake City, and has worked with Ford Aerospace and Communications Corporation, CA, INRIA-France, and IBM San Jose Research Laboratory, CA. He has been a Principal Investigator of various programs for NSF, DARPA, NASA, AFOSR, ARO, and other agencies and industries in the areas of video networks, video understanding, learning and vision, image understanding, pattern recognition, target recognition, biometrics, navigation, image databases, and machine vision applications. He has coauthored Computational Learning for Adaptive Computer Vision (Springer-Verlag, 2007), Evolutionary Synthesis of Pattern Recognition Systems (Springer-Verlag, 2005), Computational Algorithms for Fingerprint Recognition (Kluwer, 2004), Genetic Learning for Adaptive Image Segmentation (Kluwer, 1994), and Qualitative Motion Understanding (Kluwer, 1992), and coedited a book on Computer Vision Beyond the Visible Spectrum (Springer-Verlag, 2004). He holds 11 U.S. and international patents and over 250 reviewed technical publications in the areas of his interest.

Dr. Bhanu has received two Outstanding Paper Awards from the Pattern Recognition Society and has received industrial and university awards for research excellence, outstanding contributions, and team efforts. He has been on the editorial board of various journals and has edited special issues of several IEEE Transactions (PAMI, IP, SMC-B, R&A, IFS) and other journals. He was the General Chair for the IEEE Conference on Computer Vision and Pattern Recognition, IEEE Workshops on Applications of Computer Vision, IEEE Workshops on Learning in Computer Vision and Pattern Recognition; Chair for the DARPA Image Understanding Workshop and Program Chair for the IEEE Workshops on Computer Vision Beyond the Visible Spectrum and Multi-modal Biometrics. He is a Fellow of the American Association for the Advancement of Science, International Association of Pattern Recognition, and the International Society for Optical Engineering (SPIE).