

Compact signatures for 3D face recognition under varying expressions

Fahad Daniyal, Prathap Nair and Andrea Cavallaro
Queen Mary University of London
School of Electronic Engineering and Computer Science
Mile End Road, London - E1 4NS

{fahad.daniyal,prathap.nair,andrea.cavallaro}@elec.qmul.ac.uk

Abstract

We present a novel approach to 3D face recognition using compact face signatures based on automatically detected 3D landmarks. We represent the face geometry with inter-landmark distances within selected regions of interest to achieve robustness to expression variations. The inter-landmark distances are compressed through Principal Component Analysis and Linear Discriminant Analysis is then applied on the reduced features to maximize the separation between face classes. The classification of a probe face is based on a nearest mean classifier after transforming the probe onto the subspace. We analyze the performance of different landmark combinations (signatures) to determine a signature that is robust to expressions. The selected signature is then used to train a Point Distribution Model for the automatic localization of the landmarks, without any prior knowledge of scale, pose, orientation or texture. We evaluate the proposed approach on a challenging publicly available facial expression database (BU-3DFE) and achieve 96.5% recognition rate using the automatically localized signature. Moreover, because of its compactness the face signature can be stored on 2D barcodes and used for radio-frequency identification.

1. Introduction

Among the biggest challenges posed by the use of 3D geometrical information for face recognition are its sensitivity to changes in expression and the amount of data required to represent a face. The extraction of an appropriate set of anthropometric landmarks that is robust to variations in expressions, can aid in overcoming these limitations. However, automatic detection of landmarks is usually limited by prior knowledge of orientation and pose of the faces, and also by the availability of a texture map.

In this paper, we propose a compact face signature for 3D

face recognition that is extracted without prior knowledge of scale, pose, orientation or texture. The automatic extraction of the face signature is based on the fitting of a trained Point Distribution Model (PDM) [12]. The recognition algorithm first represents the geometry of the face by a set of Inter-Landmark Distances (ILDs) between the selected landmarks. These distances are then compressed using Principal Component Analysis (PCA) and projected onto the classification space using Linear Discriminant Analysis (LDA). The classification of a probe face is finally achieved by projecting the probe onto the LDA-subspace and using the nearest mean classifier.

The paper is organized as follows: Section 2 discusses prior work in 3D face recognition. Section 3 describes the proposed approach for face recognition and finding the most robust face signature. Section 4 focuses on experimental results and the validation of the algorithm. Finally, in Sec. 5 we draw the conclusions.

2. Prior work

Algorithms for 3D face recognition can be grouped in three main classes of methods, based on: (i) direct comparison of selected regions or of the whole surface [9, 3, 10]; (ii) projecting the faces onto appropriate spaces [2, 1]; and (iii) comparison of features such as landmarks and contours [11, 7]. A common approach to directly compare surfaces is through a rigid registration via the Iterated Closest Point (ICP) algorithm [9, 3, 10]. The main limitation of ICP for face recognition is that the performance of the registration degrades in the presence of deformations due to expressions and outliers in the scans. For the ICP algorithm to converge to the global minimum in terms of mean-square-error (MSE), the surfaces must first be roughly aligned. This requires prior knowledge of face orientation and the localization of specific landmarks on the face [9]. An attempt to overcome the limitation of facial expressions is presented in [3] and [10] through the matching fusion of multiple face

regions. Region detection is obtained with constraints on orientation and thresholds for curvature features.

Statistical classifiers, such as Eigenfaces and Fisherfaces [4], have been extensively used for 2D face recognition due to their efficiency and speed. Chang *et al.* [2] extended the EigenFace approach for use with 3D face meshes. However, this method is highly sensitive to expressions as the whole face surface is projected onto the eigenface space. Bronstien *et al.* [1] modeled faces as isometries of facial surfaces and the face representation is in the form of bending invariant canonicals. This canonization is done on a geodesic mask, and the accuracy of the algorithm highly depends on the accurate detection of the nose tip and other landmarks used in the embedding. Kakadiaris *et al.* [8] performed face recognition with an annotated model that is non-rigidly registered to face meshes through a combination of ICP, simulated annealing and elastically adapted deformable model fitting. A limitation of this approach is the imposed constraints on the initial orientation of the face.

Face recognition through the use of specific anthropometric landmarks can aid in overcoming the limitations due to expressions mentioned above. Facial features such as landmarks, regions and contours are generally localized based on the surface curvature. Moreno *et al.* [11] proposed an approach based on the extraction of an 86-D feature vector composed of landmarks and regions. The features are evaluated for their discriminatory power and results are demonstrated on different combinations of these features. The main limitation of such approaches lies in the localization of these facial features, which is highly dependent on the prior knowledge of feature map thresholds, face orientation and pose. Gupta *et al.* [7] presented a recognition approach using facial proportions extracted from key landmarks. The selection of the key landmarks is done manually and is based on literature about anthropometric facial proportions. However, automatic localization of these key landmarks (especially around the mouth region) is difficult.

3. Proposed approach

We aim to extract a compact landmark-based *signature* of a 3D face that is robust to changes in facial expressions. To this end, we first select a robust facial representation based on testing extensive sets of manually selected landmarks. Next, we train a Point Distribution Model (PDM) to identify the selected set of landmarks.

3.1. Landmark-based face recognition

Given a set $S = \{\omega_1, \omega_2, \dots, \omega_N\}$ of N 3D landmarks on a face mesh Ψ , where $\omega_i = (x_i, y_i, z_i)$ represents the i^{th} landmark, we extract geometrical information describing the face morphology. To this end, we compute the inter-landmark distances (ILDs), d_{ij} , between pairs of landmarks

and generate a feature vector, Δ , of dimension $N(N-1)/2$, represented as

$$\Delta = (d_{1,2}, d_{1,3}, \dots, d_{1,N}, d_{2,3}, \dots, d_{2,N}, \dots, d_{(N-1),N}), \quad (1)$$

where $d_{i,j} = \|\omega_i - \omega_j\|$. We choose the Euclidean distance for its simplicity of computation and robust representation of face geometry than, for example, geodesic distances. In fact, geodesic distances are highly sensitive to expressions, noise and the resolution of the face meshes. Moreover, the use of the Euclidean distance allows us to obtain a more concise signature as only N landmarks need to be stored, whereas the $N(N-1)/2$ ILDs can be calculated at the recognition stage.

The feature vector Δ is normalized with respect to the size of the face to make it scale invariant, thus generating

$$\tilde{\Delta} = \frac{\Delta}{d_S}, \quad (2)$$

where the scaling factor d_S is the distance between two predefined landmarks.

To reduce the dimensionality of the feature space we apply Subspace Linear Discriminant Analysis (SLDA) [17]. SLDA is the projection of the data onto a LDA space after applying PCA. The use of LDA as a feature space is suited for the task of face recognition, especially when sufficient samples per class are available for training. LDA is a supervised learning algorithm that targets data classification more than feature extraction and finds the classification hyperplane that maximizes the ratio of the between-class variance to the within-class variance, thereby guaranteeing maximal separability. The initial PCA projection allows us to reduce the dimensionality of the data while retaining its discriminative power, which LDA further improves upon by maximizing the class separation.

Let M be the number of faces in the training database and $\tilde{\Delta} = (\tilde{\Delta}^1, \tilde{\Delta}^2, \dots, \tilde{\Delta}^M)$ represent the normalized feature vectors for all the training faces. The initial PCA projection, Λ , is defined as

$$\Lambda = A^T \tilde{\Delta}, \quad (3)$$

where A is the transformation matrix whose columns are the eigenvectors obtained from the covariance matrix, Z_{Δ} , of the data. The LDA projection, Γ , is defined as

$$\Gamma = B^T \Lambda, \quad (4)$$

where the matrix B holds the eigenvectors of $Z_w^{-1} Z_b$. Here Z_w is the *within-class* covariance matrix and Z_b is the *between-class* covariance matrix (see [17] for details).

For classification, we project the probe onto the created LDA-subspace and use the nearest mean classifier. Given a

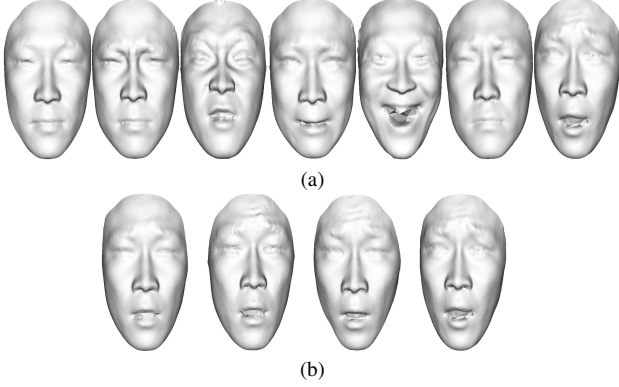


Figure 1. Sample subject scans from the BU-3DFE database: (a) 7 expressions (*Neutral, Anger, Disgust, Fear, Happiness, Sadness, Surprised*), (b) 4 intensity levels of the *Surprise* expression

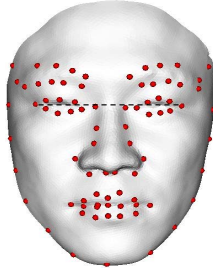


Figure 2. Sample face scan showing the annotated landmarks and the scaling distance d_S (dotted line) used in the tests

probe face Ψ^p and its landmarks $(\omega_1^p, \omega_2^p, \dots, \omega_N^p)$, we compute the feature vector $\tilde{\Delta}^p$ of normalized ILDS (Eq. 2). $\tilde{\Delta}^p$ is then projected onto the LDA-subspace using Eq. 3 and Eq. 4. The identity Δ_p^* is then chosen according to

$$\Delta_p^* = \arg \min_j \left\| \Gamma^p - \bar{\Gamma}^j \right\|, \quad (5)$$

where $\|\cdot\|$ is the Euclidean distance, $\bar{\Gamma}^j$ is the mean for class j and Γ^p is the projected probe face.

We evaluated extensively the proposed recognition algorithm on the BU-3DFE database [16], which includes a challenging range of expressions and 83 manually annotated landmarks for each face. The database contains 25 face meshes per person with four degrees (intensities) of expressions for each of the six available expressions, namely *anger, disgust, fear, happiness, sadness* and *surprise*, in addition to one *neutral* expression. A sample subject showing the range of expressions and intensities is shown in Fig. 1 while the annotated landmarks are shown in Fig. 2. We evaluated various subsets (regions) of the 83 ground-truth landmarks on 100 individuals (56 females and 44 males). The regions included the left and right eyes and eyebrows, the nose, the mouth and the boundary of the face. The scaling distance d_S used for feature normalization is the distance between the two outer eye points, as shown in Fig. 2. An

exhaustive combination of landmarks from the five regions results in 31 different models ($2^5 - 1$), ranging from single regions to all the regions. As expected, the single region that led to the worst recognition results is the mouth region, as it is most affected by variations in expressions. The most compact representation that led to the best result is the combination of the eyes, eyebrows and nose regions (48 landmarks). We refer to this model as "EY2N". This result is in line with recent literature [3, 7, 10, 13] showing robustness of the eyes and nose regions to expressions. This combination has the same recognition results as the full model and was therefore chosen for its compactness.

To automatically detect the EY2N landmarks, we generate a Point Distribution Model (PDM) [12] that includes statistical information of the shape variation of the landmarks over a training set, and then fit it to each probe and training mesh.

3.2. Model fitting for pose and scale invariant face recognition

To represent the 48 EY2N landmarks, we build a parameterized model $\Omega = \Upsilon(\mathbf{b})$, where $\Omega = \{\omega_1, \omega_2, \dots, \omega_N\}$, with $N = 49$. The extra landmark (nose tip) was included to facilitate the model fitting process that will be described later. The vector \mathbf{b} holds the parameters that can be used to vary the shape and Υ defines the function over the parameters. To obtain the model, a training set of manually localized landmarks from L face meshes is used. Training shapes are aligned and scaled to the same co-ordinate frame to eliminate global transformations using Procrustes analysis [5]. PCA is then applied to capture the variations of the shape cloud formed by the training shapes in the (3×49) -dimensional space, along the principal axes of the point cloud. The principal axes and corresponding variations are represented by the eigenvectors and eigenvalues obtained from the covariance matrix, Z_Ω , of the data.

Let ϕ contain the t eigenvectors corresponding to the largest eigenvalues. Then any shape, Ω , similar to those in the training set can be approximated as

$$\Omega \approx \bar{\Omega} + \phi \mathbf{b}, \quad (6)$$

where $\bar{\Omega}$ is the mean shape, $\phi = (\phi_1 | \phi_2 | \dots | \phi_t)$ and $\mathbf{b} = \phi^T (\Omega - \bar{\Omega})$ is a t dimensional vector. The value of t is chosen such that the model represents 98% of the shape variance, ignoring the rest as noise [5]. The mean shape is obtained when all parameters are set to zero.

The PDM Ω is fitted onto a probe mesh Ψ^p through similarity transformations of the model, estimated using three control points of the mean shape $\bar{\Omega}$. These control points are the inner eye points $(\omega_r$ and $\omega_l)$ and the nose tip (ω_n) , with $\{\omega_r, \omega_l, \omega_n\} \in \bar{\Omega}$ [12]. The inner eye and nose tip areas on a face are normally unique based on local curvature and can be robustly isolated. In order to character-

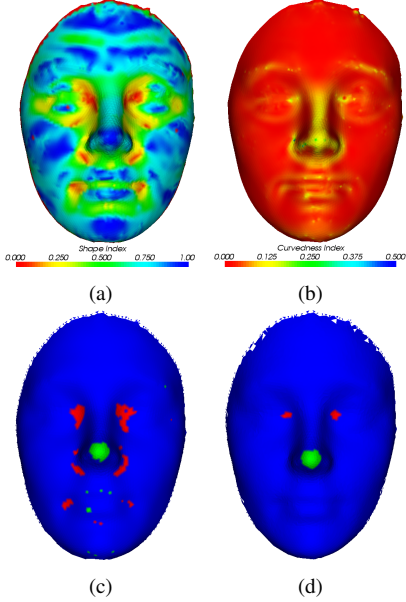


Figure 3. Feature maps used to isolate candidate vertices: (a) shape index, (b) curvedness index. Candidate vertices (regions in green are candidate nose tip vertices and regions in red are candidate eye tip vertices): (c) without decimation and without averaging, (d) with decimation and averaging

ize the curvature property of each vertex (v_i) on the face mesh, two features maps are computed, namely the *shape index*, $\rho(\cdot)$, and the *curvedness index*, $\sigma(\cdot)$ [6]. The shape index describes shape variations from concave to convex, whereas the curvedness index indicates the scale of curvature present at each vertex. These feature maps are computed after Laplacian smoothing to reduce the outliers arising from the scanning process. Figure 3(a-b) shows the two feature maps obtained on a sample face after the smoothing process. Moreover, to reduce the computational overhead, the original mesh is first decimated and then the features are averaged across vertex neighbors according to

$$\tilde{\rho}(v_i) = \frac{1}{P} \sum_{p \in \mathcal{P}(v_i)} \rho(v_p), \quad \tilde{\sigma}(v_i) = \frac{1}{P} \sum_{p \in \mathcal{P}(v_i)} \sigma(v_p), \quad (7)$$

where $\mathcal{P}(v_i)$ is the set of P neighboring vertices of v_i . If $\tilde{\sigma}(\cdot) > \sigma_s$, then v_i is in a salient high-curvature region. The condition $\tilde{\rho}(\cdot) < \rho_e$ identifies concave regions; while $\tilde{\rho}(\cdot) > \rho_n$ identifies convex regions. We can therefore relax thresholds to segregate candidate inner eye vertices from the nose tip ones. The thresholds $\sigma_s = 0.1$, $\rho_e = 0.3$ and $\rho_n = 0.7$ were found to be adequate for the entire database. Figure 3(c-d) shows a comparison of the isolated candidate inner eye vertices (red) and nose tip vertices (green) obtained with and without the mesh decimation and feature averaging steps.

A further reduction in outlier candidate combinations is performed by checking the triangle formed by each com-

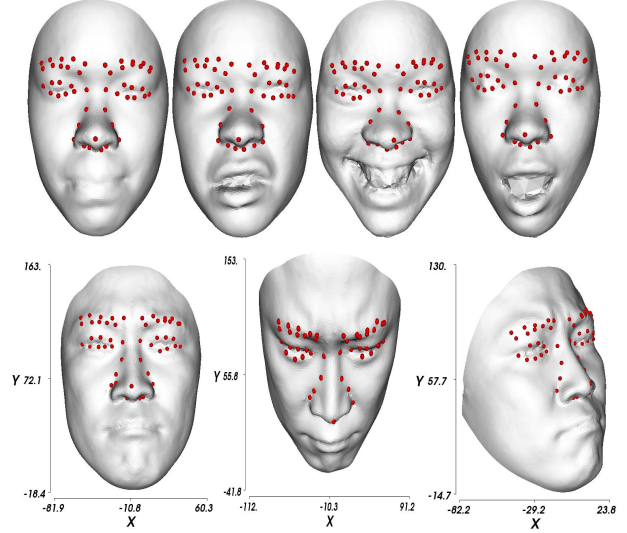


Figure 4. Examples of scale and pose invariant model fitting on faces with different expressions (top), and faces with different pose and scale (bottom)

bination of two candidate inner eye points (α_r, α_l) and a nose tip point (α_n). A plausible eyes-nose triangle should be acute angled with

$$\begin{cases} d_{rl}^2 + d_{rn}^2 > d_{ln}^2 \\ d_{rl}^2 + d_{ln}^2 > d_{rn}^2 \\ d_{rn}^2 + d_{ln}^2 > d_{rl}^2 \end{cases}$$

where d_{rl} , d_{rn} and d_{ln} are the lengths of the sides of the triangle. Plausible combinations of the candidate inner eye vertices and candidate nose tip vertices on Ψ^p are used as target points to transform the model. Next, the remaining points of Ω are moved to the closest vertices on Ψ^p . Ω is then projected back onto the model space and the parameters of the model, \mathbf{b} , are updated. Based on this selective search over the isolated candidate vertices, the transformation exhibiting the minimum deviation from the mean shape is chosen as the fit for the model. Sample face meshes with the fit model are shown in Fig. 4.

4. Experiments and discussions

We evaluate here the performance of the proposed face recognition algorithm and compare it with the 3D eigenface approach. Moreover, we discuss the influence of the expression intensities in the training and the memory requirements of the automatically detected EY2N signature based on 49 landmarks. The proposed PDM discussed in this section is trained with manually annotated landmarks from 100 (out of the total 2500) face meshes. The landmarks of the probe and training faces are detected automatically.

Figure 5 shows a comparison of the recognition rates obtained with different training and probe combinations to analyze the influence of the expression intensities used in the

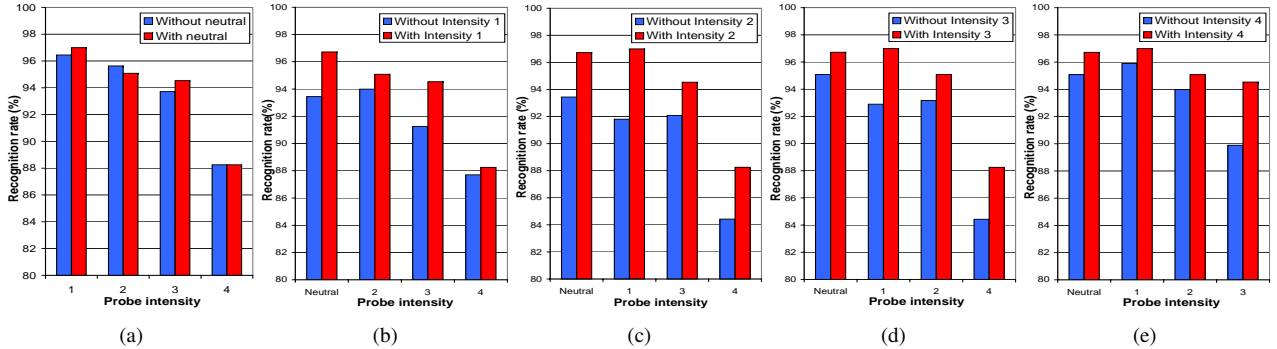


Figure 5. Comparison of recognition results using different combinations of probe and training sets. The training set was varied to include (red) and exclude (blue): (a) neutral, (b) intensity-1 (lowest intensity), (c) intensity-2, (d) intensity-3, and (e) intensity-4 (highest intensity)

SLDA training. Note that the neutral intensity has the least influence on training (Fig. 5(a)), while the inclusion of the remaining intensities have a larger effect (Fig. 5(b-e)), because neutral samples are fewer than the other intensities (1 neutral and 6 each of intensity 1-4 per person). Higher accuracy is achieved when more samples having a wide expression range are used in the training. More accurate recognition is obtained when the neutral and intensity-1 samples are used as probe, while intensity-4 provides the least accurate results. The best result (96.5% recognition accuracy) is achieved using intensity-1 as probe and the remaining samples in the training. The reduced recognition accuracy for the highly expressive samples is to be attributed to the reduced accuracy in the PDM fitting.

The use of PCA before applying LDA allows us to considerably reduce the dimensions of the feature space while retaining the most relevant information. To analyze the influence of the number of dimensions in the identification accuracy, Figure 6 shows the rank-1 recognition results obtained when varying the amount of feature energy retained by the eigenmodes after PCA. The number of dimensions that lead to the highest accuracy was 115, which corresponds to 10.20% of the original size of the feature vector (1128), using the manual landmarks. In the case of automatic landmarks, the maximum accuracy was obtained with 22.53% (265 modes) of the feature vector. This is due to the fact that as automatic landmarks contain a larger amount of noise as compared to the manual landmarks, they require more information to represent a face. These reduced dimensions correspond to 99.97% and 99.57% of the signal energy for automatic and manual landmarks, respectively.

Figure 7 shows the Receiver Operating Characteristics (ROC) curves that compares the proposed approach with a 3D eigenface method replicating [15], where depth-maps of entire face meshes were used in the PCA projection. The 3D eigenface method has lower accuracy results, with 60.48% rank-1 recognition rate, as it cannot properly handle large expression changes. The 3D eigen-face approach is also presented in [14] and [2] with multimodal data and with

3D modality only, with recognition rates of 85% and 88.9% respectively, being reported.

To quantify the decrease in recognition accuracy when reducing the precision of the proposed facial signature, Figure 7 compares results obtained with automatic landmarks of the probe using 32-bit floating point and 16-bit integer representation, where the landmark coordinates were rounded to the nearest integer. While the storage of the signature with floating point representation requires 588 bytes only, with an integer representation we achieve a further 50% reduction in the storage requirements. This would allow the 3D face signature to be stored not only on devices such as RFIDs, but also in 2D barcodes. The rank-1 identification rate is 96.5% with the floating point representation and 92.64% with integer representation. In summary, with a significant decrease in the signature size (50%) using the integer representation, there was only a 3.86% decrease in rank-1 recognition.

5. Conclusions

We proposed a novel approach to scale and pose invariant 3D face recognition with the use of a facial signature and Subspace Linear Discriminant Analysis (SLDA). The signature is a concise representation of a face that is robust to facial expressions. The approach first extracts geometric features of the face in the form of inter-landmark distances (ILDs) within a set of regions of interest. Dimensionality reduction is applied through Principal Component Analysis (PCA) to compress the data, and Linear Discriminant Analysis (LDA) is used on the reduced features to maximize the separation between the classes. The classification of a probe face is based on the nearest mean classifier after transforming its signature onto the SLDA space. We determined the most robust model to expressions using a large set of candidate landmarks and demonstrated the improved accuracy of using the sub-space LDA transformation compared to PCA or LDA alone. The automatic extraction of landmarks is based on the training and fitting of a point dis-

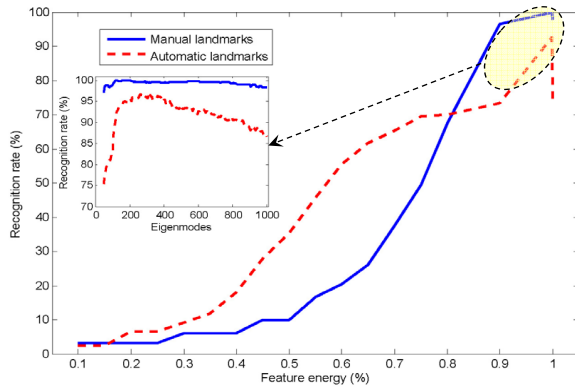


Figure 6. Comparison of rank-1 recognition accuracies of manually and automatically localized signatures on varying the amount of feature energy retained by the PCA eigenmodes. The inset is a comparison of the signature accuracy against the number of eigenmodes for the highlighted region

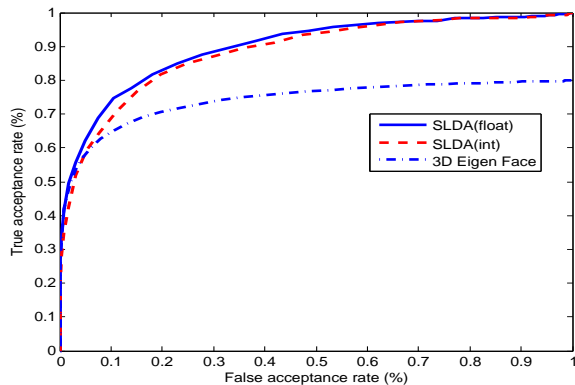


Figure 7. Comparison of results using signatures with 32-bit float and 16-bit integer representations, and the baseline 3D eigen-face approach

tribution model that eliminates the need for prior knowledge of orientation, pose or texture information.

Current work includes the validation of the proposed approach on additional datasets (such as FRGC and 3D_RMA) and on improving the fitting with local neighborhood constraints and global optimization strategies.

Acknowledgment

We acknowledge Dr. Lijun Yin, Department of Computer Science, The State University of New York at Binghamton, and The Research Foundation of State University of New York, for providing the BU-3DFE database.

References

[1] A. M. Bronstein, M. M. Bronstein, and R. Kimmel. Three-dimensional face recognition. *Intl. Journal of Computer Vision*, 64(1):5–30, Aug 2005. 1, 2

[2] K. Chang, K. Bowyer, and P. Flynn. An evaluation of multimodal 2D+3D face biometrics. *IEEE Trans. Pattern Anal. Machine Intell.*, 27(4):619–624, Apr 2005. 1, 2, 5

[3] K. Chang, K. Bowyer, and P. Flynn. Multiple nose region matching for 3D face recognition under varying facial expression. *IEEE Trans. Pattern Anal. Machine Intell.*, 28(10):1695–1700, Oct 2006. 1, 3

[4] R. Chellappa, C. L. Wilson, and S. Sirohey. Human and machine recognition of faces: A survey. *Proceeding of the IEEE*, 83(5):705–740, May 1995. 2

[5] T. F. Cootes, C. J. Taylor, D. H. Cooper, and J. Graham. Active shape models: their training and application. *Computer Vision and Image Understanding*, 61(1):38–59, Jan 1995. 3

[6] C. Dorai and A. K. Jain. Cosmos - a representation scheme for 3D free-form objects. *IEEE Trans. Pattern Anal. Machine Intell.*, 19(10):1115–1130, Oct 1997. 4

[7] S. Gupta, J. Aggarwal, M.K.Markey, and A. Bovik. 3D face recognition founded on the structural diversity of human faces. In *IEEE Conf. on Computer Vision and Pattern Recognition*, pages 1–7, Minneapolis, MN, Jun 2007. 1, 2, 3

[8] I. Kakadiaris, G. Passalis, G. Toderici, N. Murtuza, Y. Lu, N. Karampatziakis, and T. Theoharis. Three-dimensional face recognition in the presence of facial expressions: An annotated deformable model approach. *IEEE Trans. Pattern Anal. Machine Intell.*, 29(4):640–648, Apr 2007. 2

[9] X. Lu, A. K. Jain, and D. Colbry. Matching 2.5D face scans to 3D models. *IEEE Trans. Pattern Anal. Machine Intell.*, 28(1):31–43, Jan 2006. 1

[10] A. Mian, M. Bennamoun, and R. Owens. An efficient multimodal 2D-3D hybrid approach to automatic face recognition. *IEEE Trans. Pattern Anal. Machine Intell.*, 29(11):1927–1943, Nov 2007. 1, 3

[11] A. Moreno, A. Sanchez, J. Velez, and F. Diaz. Face recognition using 3D surface-extracted descriptors. In *Irish Machine Vision and Image Processing Conf.*, Ireland, Sept 2003. 1, 2

[12] P. Nair and A. Cavallaro. Region segmentation and feature point extraction on 3D faces using a point distribution model. In *IEEE Intl. Conf. on Image Processing*, volume 3, pages 85–88, Texas, USA, Sept 2007. 1, 3

[13] P. Sinha, B. Balas, Y. Ostrovsky, and R. Russell. Face recognition by humans: nineteen results all computer vision researchers should know about. *Proc. of the IEEE*, 94(11):1948–1962, Nov 2006. 3

[14] F. Tsakanidou, D. Tzocaras, and M. Strintzis. Use of depth and colour eigenfaces for face recognition. In *Pattern Recognition Letters*, volume 24, pages 1427 – 1435, Jun 2003. 5

[15] C. Xu, Y. Wang, T. Tan, and L. Quan. A new attempt to face recognition using 3D eigenfaces. In *Asian Conf. on Computer Vision*, pages 884–889, Jeju, Korea, Jan 2004. 5

[16] L. Yin, X. Wei, Y. Sun, J. Wang, and M. Rosato. A 3D facial expression database for facial behavior research. In *7th Intl. Conf. on Automatic Face and Gesture Recognition*, pages 211–216, Southampton, UK, Apr 2006. 3

[17] W. Zhao, R. Chellappa, and N. Nandhakumar. Empirical performance analysis of linear discriminant classifiers. In *IEEE Conf. on Computer Vision and Pattern Recognition*, pages 164–169, Washington, DC, 1998. 2