

3D model-based pose invariant face recognition from multiple views

Q. Chen, J. Yao and W.K. Cham

Abstract: A 3D model-based pose invariant face recognition method that can recognise a human face from its multiple views is proposed. First, pose estimation and 3D face model adaptation are achieved by means of a three-layer linear iterative process. Frontal view face images are synthesised using the estimated 3D models and poses. Then the discriminant ‘waveletfaces’ are extracted from these synthesised frontal view images. Finally, corresponding nearest feature space classifier is implemented. Experimental results show that the proposed method can recognise faces under variable poses with good accuracy.

1 Introduction

Face recognition has been an active research area over the past few years and numerous face recognition algorithms have been proposed. Good reviews can be found in [1–5]. Most automatic face recognition (AFR) algorithms are for face recognition under controlled conditions. For example, satisfactory recognition rates on face images which are uncovered, in frontal view, with neutral expression and controlled lighting have been reported in [6–11]. However, when the input face images are not so ideal and have some variations on such conditions, the performance of these AFR algorithms will deteriorate.

The problem related to variations in poses received much attention and many algorithms have been developed to tackle this problem. An early attempt is the 2D appearance based approach that describes faces under varying pose with a set of 2D features and achieves pose analysis (such as pose estimation and pose classifying) and face recognition by comparing these features. Murase and Nayar [12] present a method for pose invariant face recognition in the entire eigenspace. Pentland *et al.* [13] and Huang *et al.* [14] achieved pose invariant face recognition in the viewspace which is a subspace of the eigenspace. Demir’s method [15] is similar to that of Pentland *et al.* [18], but employing a sub-linear discriminant analysis (LDA) space as the view-space. This problem was tackled in the discriminant ‘waveletface’ space [16] and the kernel LDA space [17]. De Vel and Aeberhard [18] describe a line-based algorithm for pose invariant face recognition. When the gallery has an image having a pose similar to that of the test images, these appearance-based methods have good recognition results. Otherwise, their performance deteriorates. Hence, these methods require dense sampling of the continuous pose. This not only increases the gallery size but also makes the recognition process more time-consuming. Therefore the 3D model-based approach was proposed, which has

stronger generalisation to pose invariant face recognition, though its implementation is more complex.

In the 3D model-based approach, a 3D face model is built to represent the 3D geometry of human faces in 2D images. This approach removes the effect of pose variations on face recognition by estimating and aligning poses with a 3D face model and then extracting features under a uniform pose for classification. Generally, pose estimation is the most critical and challenging operation in 3D model based approach. Methods described in [19–24], which may be referred to as single-view-based methods, achieve pose invariant face recognition using only one view image for each candidate in the gallery. In [19, 20], a fixed generic 3D face model was proposed to be used for all candidates. On the basis of this face model, the pose estimation was achieved by using affine transforms. These methods are simple but the pose of a particular face cannot be estimated accurately by using a fixed 3D model. In [23, 24], simple adaptive 3D face models that can adapt to fit a particular person were proposed. The pose estimation and the model adaptation were achieved synchronously using geometrical measurements. As the 3D models used in [27, 28] are simple and rough, these methods cannot achieve accurate pose estimation, especially they cannot estimate seesaw rotation of faces. Recently, Blanz *et al.* [21, 22] built a 3D morphable face model from a large set of real 3D face data for pose invariant face recognition. On the basis of this model, the pose estimation and the model adaptation were achieved by hybrid geometric information and texture information based optimisation. In fact, the 3D geometry and the pose of a particular face can hardly be exactly recovered from only one of its 2D projections. Hence, in [25], rangefinder is used to construct accurate 3D face model for pose invariant face recognition. However, each registered candidate should be captured by rangefinder, the use of which is the limitation of this system as this is expensive equipment. To achieve good performance on pose estimation and model adaptation just from 2D grey-face images, Zhang *et al.* [26] proposed a multiple-view approach which created a 3D face structure and estimated poses from multiple views of a particular face by adapting a generic 3D face model with a cubic explicit polynomial. Pose estimation was achieved by minimising the distance map residual error using the Levenberg–Marquardt optimisation method. Such nonlinear pose

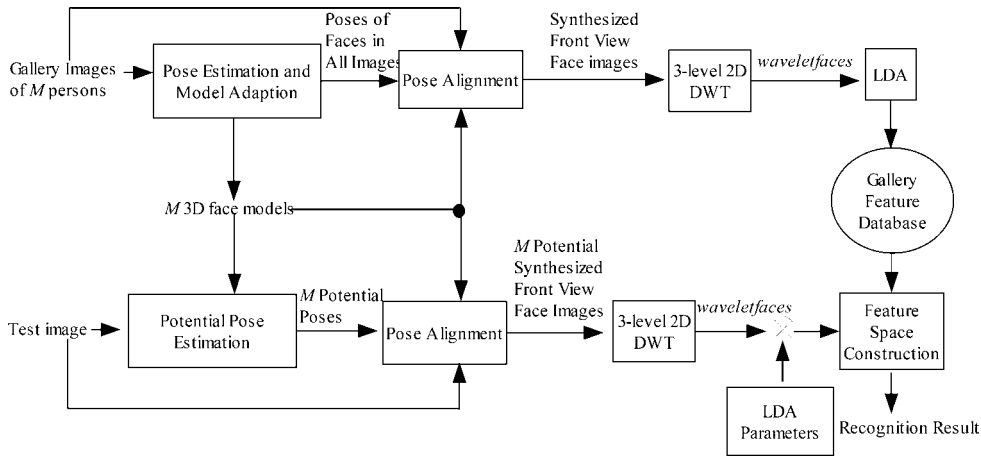


Fig. 1 Overall method architecture of the face recognition

estimation algorithm requires large computation cost and is sensitive to the initial values (local minimal problem).

This paper proposes a 3D model based pose invariant face recognition method using the multiple-view approach. The proposed method retains the advantages of the 3D model based approach and overcomes most limitations of previous methods using this approach. The major contribution of the proposed method is that it uses our linear iterative algorithm [27] to achieve accurate, fast and robust pose estimation from multiple views for pose invariant face recognition. In addition, by performing linear discriminant analysis (LDA) on synthesised frontal view images and employing corresponding nearest feature space (CNFS) classifier, our method makes full use of the statistic information to extract discriminant features and achieves robust classification.

2 Overview of the proposed method

In this paper, we propose a multiple view 3D model based pose invariant face recognition method that can recognise a face from its multiple views. This method is composed of four steps: (1) pose estimation and 3D model adaptation, (2) pose alignment, (3) feature extraction and (4) classification. The major challenge of this approach is step (1) in which we need to estimate the pose of a face and adapt the 3D model to fit the face. In this paper, we solve this problem by using the inverse projection rays based geometric constraint in a three-layer linear iterative process.

The block diagram of the proposed method is given in Fig. 1. First, the poses of the faces in the test image and gallery images are estimated using a three-layer linear iterative process. This algorithm achieves pose estimation in a linear iterative way like [28–30], but does not require the accurate 3D geometry of the face in the images before estimation. The 3D geometry of a face is obtained by iteratively updating the reference face model to fit a particular person. On the basis of the estimated poses and 3D face models, pose alignment is implemented by synthesising frontal view face images from the test face image and gallery face images. We extract the waveletfaces from these synthesised frontal view images. Finally, linear discriminant analysis is performed on these waveletfaces and the CNFS classifier is used for classification.

3 Pose estimation and 3D model adaptation

We adopt our model-based linear pose estimation (MBLPE) algorithm described in [27] for pose estimation and model

adaptation. In the algorithm, we use a 3D reference face model to represent the 3D geometry of a face. The input is a set of N_s face images of the same candidate with the 2D coordinates of needed facial feature points extracted from these face images. The output consists of the pose of a human face for each input face image and an adapted reference face model for the face in the input images.

The MBLPE algorithm consists of three linear iterative updating processes as shown in Fig. 2. In the innermost layer, we estimate the rotation and translation parameters of each input face image based on a reference face model and a set of facial feature points. In the middle layer, the facial feature points on the reference model are globally updated by three model scaling factors which are estimated using multiple face images with the recovered poses. Then, we locally update the facial feature points on the reference face model in the outermost layer. Finally, the whole reference model can be obtained by model deformation. Before the three-layer iterative process, the rotation matrix for each face image is initially set equal to an identity matrix ($\text{diag}(1, 1, 1)$), and the default 3D generic wireframe face

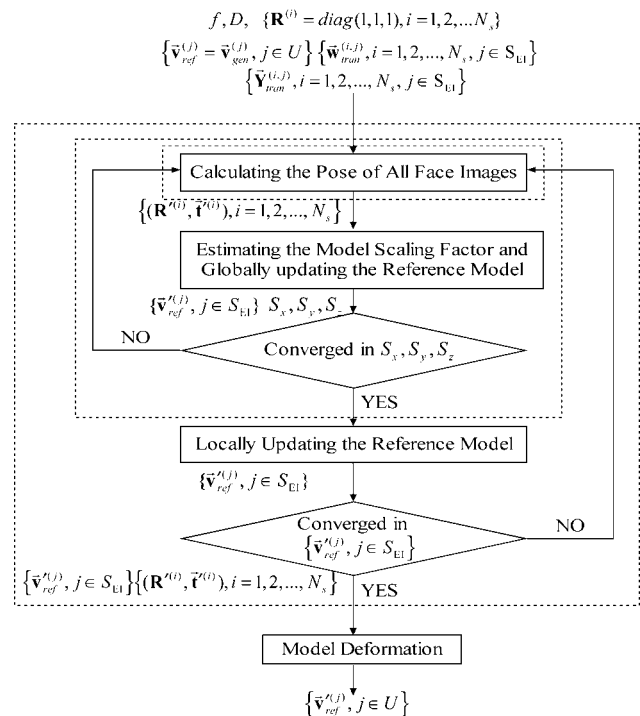


Fig. 2 Flowchart of the MBLPE algorithm

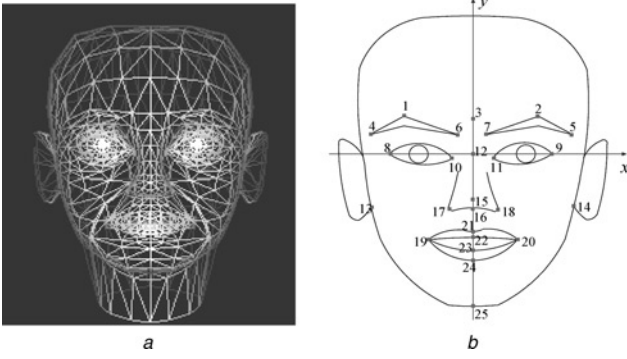


Fig. 3 Generic wireframe face model and facial feature points
 a Generic wireframe face model
 b Facial feature points to be extracted in each face

model, which comprises 1039 control points shown in Fig. 3a, is used as the initial reference face model. This generic wireframe model was developed by Instituto Superior Tecnico [31].

In this paper, the default generic wireframe model is represented by control points $\{\vec{v}_{\text{gen}}^{(j)} = (x_{v_{\text{gen}}}(j), y_{v_{\text{gen}}}(j), z_{v_{\text{gen}}}(j))^T, j \in U\}$, where U is the set containing integers from 1 to 1039. The 25 points, shown in Fig. 3b, are called facial feature points which represent noticeable facial features on a face. We assume that these facial feature points have been extracted, as long as they are not occluded in a face image. If the mouth in an image is closed, facial feature points 1–25 are regarded as forming a rigid body. Otherwise, only facial feature points 1–16 are regarded as forming a rigid body. Only facial feature points on a rigid body S_{EI} are used in the MBLPE algorithm. Actually, S_{EI} is a subset of U . If the distance between points 21 and 24 is shorter than six times the distance between points 22 and 23, the mouth will be considered as open.

Consider an imaging system as shown in Fig. 4 with the origin of the image plane placed at $(0, 0, f)$ from the perspective center O , where f is the focal length of image plane and D is the distance between the origin O and the origin o_v of a 3D reference face model. We assume it is a weak perspective projection system, that is, $D \gg f$. The reference model is represented by control points $\{\vec{P}_{\text{ref}}^{(j)} = (X_{P_{\text{ref}}}(j), Y_{P_{\text{ref}}}(j), Z_{P_{\text{ref}}}(j))^T, j \in U\}$ or $\{\vec{v}_{\text{ref}}^{(j)} = (x_{v_{\text{ref}}}(j), y_{v_{\text{ref}}}(j), z_{v_{\text{ref}}}(j))^T, j \in U\}$ w.r.t. origin O or origin o_v , respectively. $\{\vec{w}_{\text{ref}}^{(j)} = (x_{w_{\text{ref}}}(j), y_{w_{\text{ref}}}(j))^T, j \in U\}$ represent the projected 2D points on the image plane from $\{\vec{P}_{\text{ref}}^{(j)}, j \in U\}$. Now we consider the rigid transformation of the 3D reference face model by rotating and translating the control points $\{\vec{v}_{\text{ref}}^{(j)}, j \in U\}$, yielding $\{\vec{v}_{\text{tran}}^{(j)} = \mathbf{R}\vec{v}_{\text{ref}}^{(j)} + \vec{t}, j \in U\}$, where \mathbf{R} and \vec{t} stand for the rotation matrix and the translation vector, respectively. Therefore

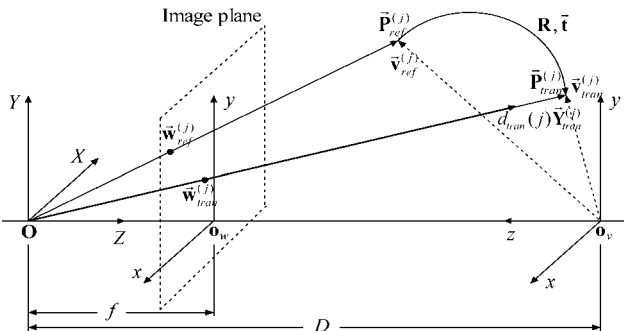


Fig. 4 Perspective projection system

$\{\vec{v}_{\text{tran}}^{(j)}, j \in U\}$ and corresponding $\{\vec{P}_{\text{tran}}^{(j)} = (X_{P_{\text{tran}}}(j), Y_{P_{\text{tran}}}(j), Z_{P_{\text{tran}}}(j))^T, j \in U\}$ represent the control points on the transformed reference model w.r.t. o_v and O , respectively. $\{\vec{w}_{\text{tran}}^{(j)} = (x_{w_{\text{tran}}}(j), y_{w_{\text{tran}}}(j))^T, j \in U\}$ are projection of $\{\vec{P}_{\text{tran}}^{(j)}, j \in U\}$ on the image plane (Fig. 4). In a weak perspective projection system, coordinates of $\vec{P}_{\text{tran}}^{(j)}$ and $\vec{w}_{\text{tran}}^{(j)}$ are related by

$$x_{w_{\text{tran}}}(j) = -f \frac{X_{P_{\text{tran}}}(j)}{Z_{P_{\text{tran}}}(j)} \simeq -f \frac{X_{P_{\text{tran}}}(j)}{D},$$

$$y_{w_{\text{tran}}}(j) = f \frac{Y_{P_{\text{tran}}}(j)}{Z_{P_{\text{tran}}}(j)} \simeq f \frac{Y_{P_{\text{tran}}}(j)}{D} \quad (1)$$

In this case, $\vec{w}_{\text{tran}}^{(j)}$ tracked in an image will give an inverse projection ray with unit vector $\vec{Y}_{\text{tran}}^{(j)}$ specified by $(x_{w_{\text{tran}}}(j)^2 + y_{w_{\text{tran}}}(j)^2 + f^2)^{-1/2} (x_{w_{\text{tran}}}(j), y_{w_{\text{tran}}}(j), f)^T$. So we can calculate the actual coordinates of the control point $\vec{P}_{\text{tran}}^{(j)}$ on the transformed 3D reference face model as follows

$$\vec{P}_{\text{tran}}^{(j)} = d_{\text{tran}}(j) \vec{Y}_{\text{tran}}^{(j)} \quad (2)$$

where $d_{\text{tran}}(j)$ denotes the depth of $\vec{P}_{\text{tran}}^{(j)}$ from the perspective center O . Here, the depth $d_{\text{tran}}(j)$ is unknown and it would be iteratively updated in the innermost-layer iterative process.

3.1 Pose estimation

For convenience, we consider the camera fixed and the head moving. In the innermost-layer iteration, we update the motion of a human head for each input image by computing six absolute motion parameters x, y, z, t_x, t_y and t_z .

The problem of estimating the pose of the face in an input face image w.r.t. a 3D reference face model can be described as follows. Given a set of 2D facial feature points $\{\vec{w}_{\text{tran}}^{(j)}, j \in S_{\text{EI}}\}$ on the input face image and a set of 3D facial feature points $\{\vec{v}_{\text{ref}}^{(j)}, j \in S_{\text{EI}}\}$ on the reference face model, we seek rotation matrix \mathbf{R} , translation vector \vec{t} and $\{d_{\text{tran}}(j), j \in S_{\text{EI}}\}$ by minimising the following error function

$$\varepsilon^2(\mathbf{R}, \vec{t}, \{d_{\text{tran}}(j)\}) = \sum_{j \in S_{\text{EI}}} j \times \left\| (d_{\text{tran}}(j) \vec{Y}_{\text{tran}}^{(j)}) - (\mathbf{R} \vec{v}_{\text{ref}}^{(j)} + \vec{t}) \right\|^2 \quad (3)$$

where $\{\vec{Y}_{\text{tran}}^{(j)}, j \in S_{\text{EI}}\}$ are unit inverse projective ray vectors of $\{\vec{w}_{\text{tran}}^{(j)}, j \in S_{\text{EI}}\}$ from the perspective center O . The coefficient j is determined by the availability of the facial feature point, which corresponds to the j th control point, in the input face image. It is 1 if we can successfully extract this facial feature point in the input face image, otherwise it is 0. The function (\cdot) denotes the coordinate conversion relationship between the coordinates of the point $\vec{P}_{\text{tran}}^{(j)}$ w.r.t. O and the coordinates of the corresponding point $\vec{v}_{\text{tran}}^{(j)}$ w.r.t. o_v , which is given as follows

$$\vec{v}_{\text{tran}}^{(j)} = (\vec{P}_{\text{tran}}^{(j)}) = \text{diag}(-1, 1, -1) (\vec{P}_{\text{tran}}^{(j)} - [0 \ 0 \ D]^T) \quad (4)$$

$$\vec{P}_{\text{tran}}^{(j)} = (\vec{v}_{\text{tran}}^{(j)}) = \text{diag}(-1, 1, -1) (\vec{v}_{\text{tran}}^{(j)} - [0 \ 0 \ D]^T) \quad (5)$$

The computation requirement of the nonlinear pose estimation algorithm increases quickly as the number of used facial feature points increases [32, 33]. Local minima problem is another limitation of the nonlinear pose estimation algorithm. In order to reduce the computation time and to avoid being trapped at a local minimum solution, the pose estimation is divided into three linear iterative

problems that can be solved efficiently [29]. The first stage approximates the global translation vector of the reference model. The second stage updates the depth values $\{d_{\text{tran}}(j), j \in S_{\text{EI}}\}$ of the facial feature points $\{\bar{\mathbf{P}}_{\text{tran}}^{(j)}, j \in S_{\text{EI}}\}$ on the transformed reference face model using recovered motion parameters. The third stage is to determine the rigid motion parameters by using a least-square minimisation algorithm. The above three stages are repeated in turn until convergent values of x, y, z, t_x, t_y and t_z are obtained.

3.1.1 Global translation approximation stage: Given a certain rotation matrix \mathbf{R} , the optimal value for the global translation vector $\vec{\mathbf{t}}$ can be computed in closed form as

$$\vec{\mathbf{t}} = \text{diag}(-1, 1, -1) \sum_{j \in S_{\text{EI}}} j (\mathbf{I} - \bar{\mathbf{Y}}_{\text{tran}}^{(j)} \bar{\mathbf{Y}}_{\text{tran}}^{(j)\text{T}})^{-1} \times \sum_{j \in S_{\text{EI}}} j (\bar{\mathbf{Y}}_{\text{tran}}^{(j)} \bar{\mathbf{Y}}_{\text{tran}}^{(j)\text{T}} - \mathbf{I}) (\mathbf{R} \bar{\mathbf{v}}_{\text{ref}}^{(j)}) \quad (6)$$

where \mathbf{I} is an identity matrix.

3.1.2 Depth approximation stage: After the motion parameters are approximated in the global translation approximation stage or in the previous pose estimation stage, the deviation between the recovered facial feature points and the true facial feature points becomes small. The facial feature points $\{\bar{\mathbf{P}}_{\text{tran}}^{(j)}, j \in S_{\text{EI}}\}$ on the transformed reference model are approximated by the perpendicular intersection of the corresponding facial feature points $\{(\mathbf{R} \bar{\mathbf{v}}_{\text{ref}}^{(j)} + \vec{\mathbf{t}}), j \in S_{\text{EI}}\}$. Hence, the depth values of $\{\bar{\mathbf{P}}_{\text{tran}}^{(j)}, j \in S_{\text{EI}}\}$ can be updated as follows

$$d'_{\text{tran}}(j) = \bar{\mathbf{Y}}_{\text{tran}}^{(j)\text{T}} (\mathbf{R} \bar{\mathbf{v}}_{\text{ref}}^{(j)} + \vec{\mathbf{t}}) \quad (7)$$

In this way, the estimated depth values will be gradually moved toward the true depth values.

3.1.3 Least-square fitting stage: After the depth values $\{d'_{\text{tran}}(j), j \in S_{\text{EI}}\}$ are determined, the pose represented by \mathbf{R} and $\vec{\mathbf{t}}$ can be further updated by minimising the following error function

$$\varepsilon^2(\mathbf{R}, \vec{\mathbf{t}}) = \sum_{j \in S_{\text{EI}}} j \left\| \bar{\mathbf{v}}_{\text{tran}}^{(j)} - (\mathbf{R} \bar{\mathbf{v}}_{\text{ref}}^{(j)} + \vec{\mathbf{t}}) \right\|^2 \quad (8)$$

where $\{\bar{\mathbf{v}}_{\text{tran}}^{(j)}, j \in S_{\text{EI}}\}$ are the updated estimation of facial feature points on the transformed reference model given by

$$\bar{\mathbf{v}}_{\text{tran}}^{(j)} = (\bar{\mathbf{P}}_{\text{tran}}^{(j)}) = (d'_{\text{tran}}(j) \bar{\mathbf{Y}}_{\text{tran}}^{(j)}) \quad (9)$$

The singular value decomposition method described in [32] is used to solve this least-square minimisation problem and we obtain updated pose estimation \mathbf{R}' and $\vec{\mathbf{t}}'$.

3.2 Globally updating the reference model

Only when the reference model represents the 3D geometry of the face in the input images with good accuracy, reliable pose estimation in the innermost-layer iterative process can be obtained. Thus, we need to adapt the reference model to fit the particular face in the input images. The model adaptation is achieved by globally updating the reference model and locally updating the reference model iteratively. Locally updating the reference model will be discussed in Section 3.3. Global updating of the reference model is implemented by scaling the 3D coordinates of the facial feature points on the reference model with scaling matrix

$\mathbf{S} = \text{diag}(S_x, S_y, S_z)$, where S_x, S_y and S_z are three model scaling factors. We compute these model scaling factors using multiple face images with poses recovered in the above pose estimation process by minimising the following error function

$$\varepsilon^2(\mathbf{S}) = \sum_{i=1}^{N_s} \sum_{j \in S_{\text{EI}}} ij \left\| \bar{\mathbf{v}}_{\text{tran}}^{(i,j)} - (\mathbf{R}^{(i)} \mathbf{S} \bar{\mathbf{v}}_{\text{ref}}^{(j)} + \vec{\mathbf{t}}^{(i)}) \right\|^2 \quad (10)$$

where N_s is the total number of input face images. $\mathbf{R}^{(i)}$ and $\vec{\mathbf{t}}^{(i)}$ denote the recovered absolute rotation matrix and absolute translation vector of the i th input face image, respectively. $\{\bar{\mathbf{v}}_{\text{tran}}^{(i,j)}, j \in S_{\text{EI}}\}$ calculated using (9) denote the last updated estimation of facial feature points on the transformed reference model which corresponds to pose in the i th input face image. The coefficient ij is similar to j in (3) but defined just for the i th face image.

By taking the partial derivative of ε^2 in (10) w.r.t. S_x, S_y and S_z and then setting them equal to zeros, three scaling factors can be computed as follows

$$S_x = \frac{\sum_{i=1}^{N_s} \sum_{j \in S_{\text{EI}}} ij (\bar{\mathbf{v}}_{\text{tran}}^{(i,j)} - \vec{\mathbf{t}}^{(i)}) \bar{\mathbf{r}}_1^{(i)} x_{\text{vref}}(j)}{\sum_{i=1}^{N_s} \sum_{j \in S_{\text{EI}}} ij (x_{\text{vref}}(j))^2} \quad (11)$$

$$S_y = \frac{\sum_{i=1}^{N_s} \sum_{j \in S_{\text{EI}}} ij (\bar{\mathbf{v}}_{\text{tran}}^{(i,j)} - \vec{\mathbf{t}}^{(i)}) \bar{\mathbf{r}}_2^{(i)} y_{\text{vref}}(j)}{\sum_{i=1}^{N_s} \sum_{j \in S_{\text{EI}}} ij (y_{\text{vref}}(j))^2} \quad (12)$$

$$S_z = \frac{\sum_{i=1}^{N_s} \sum_{j \in S_{\text{EI}}} ij (\bar{\mathbf{v}}_{\text{tran}}^{(i,j)} - \vec{\mathbf{t}}^{(i)}) \bar{\mathbf{r}}_3^{(i)} z_{\text{vref}}(j)}{\sum_{i=1}^{N_s} \sum_{j \in S_{\text{EI}}} ij (z_{\text{vref}}(j))^2} \quad (13)$$

where $\bar{\mathbf{r}}_1^{(i)}, \bar{\mathbf{r}}_2^{(i)}$ and $\bar{\mathbf{r}}_3^{(i)}$ denote the first, second and third column vectors of the rotation matrix $\mathbf{R}^{(i)}$, respectively, that is, $\mathbf{R}^{(i)} = [\bar{\mathbf{r}}_1^{(i)}, \bar{\mathbf{r}}_2^{(i)}, \bar{\mathbf{r}}_3^{(i)}]$. We obtain the globally updated facial feature points $\{\bar{\mathbf{v}}_{\text{ref}}^{(j)}, j \in S_{\text{EI}}\}$ by $\{\mathbf{S} \cdot \bar{\mathbf{v}}_{\text{ref}}^{(j)}, j \in S_{\text{EI}}\}$.

$\{\bar{\mathbf{v}}_{\text{ref}}^{(j)}, j \in S_{\text{EI}}\}$ will be used in a new turn of pose estimation process. In the middle-layer iterative process, both the pose estimation process for each face image described in Section 3.1 and the global model updating step are performed in turn until convergent values of S_x, S_y and S_z are obtained.

3.3 Locally updating the reference face model

We have recovered the model scaling factors for globally updating the reference face model to fit the input face images. However, accurate 3D coordinates of the facial feature points are also very important for realistic modelling of a particular person in input face images. On the basis of the recovered face poses in all input face images, we can locally update the locations of facial feature points $\{\bar{\mathbf{v}}_{\text{ref}}^{(j)}, j \in S_{\text{EI}}\}$ on the reference model in order to fit a particular person closely by minimising the following error function

$$\varepsilon^2(\{\bar{\mathbf{v}}_{\text{ref}}^{(j)}\}) = \sum_{i=1}^{N_s} \sum_{j \in S_{\text{EI}}} ij \left\| \bar{\mathbf{v}}_{\text{tran}}^{(i,j)} - (\mathbf{R}^{(i)} \bar{\mathbf{v}}_{\text{ref}}^{(j)} + \vec{\mathbf{t}}^{(i)}) \right\|^2 \quad (14)$$

where $N_s, \mathbf{R}^{(i)}, \vec{\mathbf{t}}^{(i)}, \bar{\mathbf{v}}_{\text{tran}}^{(i,j)}$ and ij have the same meaning as in (10).

By taking the partial derivative of ε^2 in (15) w.r.t. $\bar{\mathbf{v}}_{\text{ref}}^{(j)}$ and then setting it to zero, we locally update $\bar{\mathbf{v}}_{\text{ref}}^{(j)}$ as follows

$$\bar{\mathbf{v}}_{\text{ref}}^{(j)} = \frac{\sum_{i=1}^{N_s} ij (\mathbf{R}^{(i)})^{\text{T}} (\bar{\mathbf{v}}_{\text{tran}}^{(i,j)} - \vec{\mathbf{t}}^{(i)})}{\sum_{i=1}^{N_s} ij} \quad (15)$$

In the system, we assume that the human face is a symmetric body. We keep the symmetry by averaging the coordinate values of symmetric points and keeping x -coordinate values of the points on the symmetric line still at zero, while $\vec{v}_{\text{ref}}^{(j)}$ are being updated. The locally updated $\{\vec{v}_{\text{ref}}^{(j)}, j \in S_{\text{EI}}\}$ will be used in a new turn of pose estimation process. In the outmost-layer iterative process, both the pose estimation process for each face image described in Section 3.1, the global model updating step described in Section 3.2 and the local model updating step are performed in turn until convergent values of $\{\vec{v}_{\text{ref}}^{(j)}, j \in S_{\text{EI}}\}$ are obtained.

On the basis of final set of updated facial feature points $\{\vec{v}_{\text{ref}}^{(j)}, j \in S_{\text{EI}}\}$, the other control points on the reference model $\{\vec{v}_{\text{ref}}^{(j)}, j \in (U - S_{\text{EI}})\}$ can be deformed from the default generic face model by utilising the radial basis function interpolation method [34]. Thus, the final adapted reference model is employed as the estimated face model of the particular person in the input face images, which is represented by control points $\{\vec{v}_{\text{est}}^{(j)}, j \in U\}$. In our method, the MBLPE algorithm will be used for the training images of each candidate in the gallery database. Therefore for each gallery face image, there is an estimated pose, and for each candidate in the gallery, there is an estimated 3D face model. When a test face image is received, the MBLPE algorithm estimates the pose assuming that the test image is an image of the i th candidate for $i = 1, \dots, M$ and M being the number of candidates in the gallery. We estimate the i th potential pose of the face in the test image using the innermost-layer iterative process described in Section 3.1. Thus, we obtain M potential poses for the test face image referring to all candidates in the gallery.

Fig. 5 shows some examples to illuminate the performance of the MBLPE algorithm for pose estimation. Figs 5a and c show four gallery images of a same person with varying poses and the corresponding face model with estimated poses, respectively. The results show that the MBLPE algorithm can successfully estimate the poses and 3D face model from gallery images of different poses. Figs. 5b and d show a set of test images and their poses estimated based on the face model generated from the set of images in Fig. 5a. The result shows that the MBLPE algorithm can obtain a good estimate of poses for these test images. A more detailed evaluation of the MBLPE algorithm will be given in Section 6.

4 Pose alignment

In this section, we propose an algorithm to synthesise a frontal view image from a face image of varying pose. First, we adjust the 3D structure and the orientation of the estimated wireframe model to fit the face in the image. Then the adjusted wireframe model is overlaid on the face image. The frontal view face can be synthesised by rotating the 3D wireframe model and performing texture mapping.

Let $\{\vec{v}_{\text{adj}}^{(j)} = (x_{\text{v}_{\text{adj}}}(j), y_{\text{v}_{\text{adj}}}(j), z_{\text{v}_{\text{adj}}}(j))^T, j \in U\}$ represent the control points of the adjusted face model. Hence, the adjustment should obey the two rules: (i) The 2D projection of $\{\vec{v}_{\text{adj}}^{(j)}, j \in S_{\text{EI}}\}$ on the image plane should coincide with facial feature points on the face image; (ii) the 3D structure of $\{\vec{v}_{\text{adj}}^{(j)}, j \in U\}$ should be as similar as possible to that of $\{\vec{v}_{\text{nor}}^{(j)} = (x_{\text{v}_{\text{nor}}}(j), y_{\text{v}_{\text{nor}}}(j), z_{\text{v}_{\text{nor}}}(j))^T, j \in U\}$ which is obtained by normalising and rotating $\{\vec{v}_{\text{est}}^{(j)}, j \in U\}$ as follows

$$\vec{v}_{\text{nor}}^{(j)} = S_{\text{norm}} \cdot \mathbf{R} \cdot \vec{v}_{\text{est}}^{(j)}, \quad \text{where} \quad (16)$$

$$S_{\text{norm}} = \frac{\|\mathbf{R}_s \cdot \vec{v}_{\text{est}}^{(l)} - \mathbf{R}_s \cdot \vec{v}_{\text{est}}^{(r)}\|}{\|\vec{w}_{\text{tran}}^{(l)} - \vec{w}_{\text{tran}}^{(r)}\|}$$

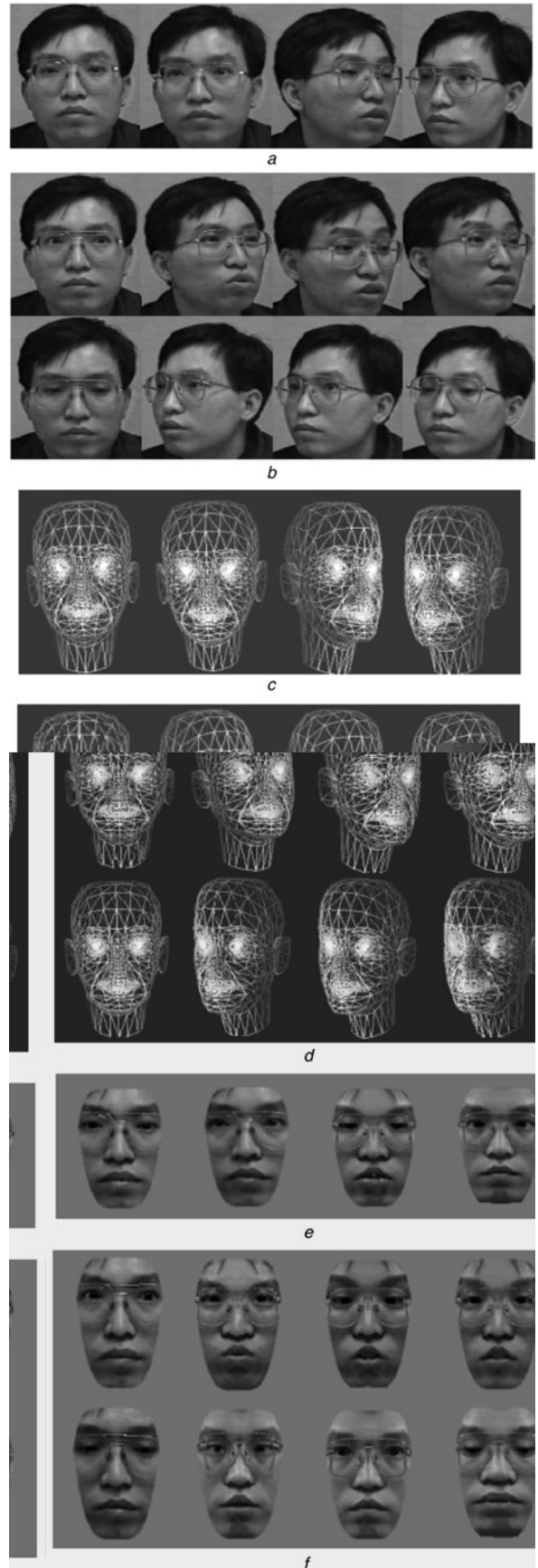


Fig. 5 Performance of the proposed algorithms for pose estimation and pose alignment

- a Four gallery face images of a same person with different poses
- b Eight test face images of the same person with different poses
- c Estimated face model with the estimated poses for gallery images
- d Face model with the estimated poses for test images
- e Synthesised frontal view images obtained from gallery images in a
- f Synthesised frontal view images obtained from test images in b

Here, $\vec{w}_{\text{tran}}^{(l)}$ denotes the left outside eye corner on the face image if it can be extracted, otherwise $\vec{w}_{\text{tran}}^{(l)}$ denotes the left inside eye corner. $\vec{w}_{\text{tran}}^{(r)}$ has similar definition as $\vec{w}_{\text{tran}}^{(l)}$ but for the right one. $\vec{v}_{\text{est}}^{(l)}$ and $\vec{v}_{\text{est}}^{(r)}$ denote the 3D coordinates of the corresponding eye corners in the estimated 3D face model and \mathbf{R}_s is a 2×3 matrix, which is the first two rows of matrix \mathbf{R} .

Therefore, we achieve model adjustment as follows

$$\forall j \in S_{\text{EI}} \cap j = 1: x_{v_{\text{adj}}}(j) = x_{w_{\text{tran}}}(j),$$

$$y_{v_{\text{adj}}}(j) = y_{w_{\text{tran}}}(j), z_{v_{\text{adj}}}(j) = z_{v_{\text{nor}}}(j) \quad (17)$$

$$\forall j \in U - S_{\text{EI}} \cup j = 0:$$

$$x_{v_{\text{adj}}}(j) = \frac{\sum_{k \in S_{\text{EI}}} (|x_{v_{\text{nor}}}(j) - x_{v_{\text{nor}}}(k)|) \cdot (x_{w_{\text{tran}}}(k) + (x_{v_{\text{nor}}}(j) - x_{v_{\text{nor}}}(k)))}{\sum_{k \in S_{\text{EI}}} (|x_{v_{\text{nor}}}(j) - x_{v_{\text{nor}}}(k)|)}$$

$$y_{v_{\text{adj}}}(j) = \frac{\sum_{k \in S_{\text{EI}}} (|y_{v_{\text{nor}}}(j) - y_{v_{\text{nor}}}(k)|) \cdot (y_{w_{\text{tran}}}(k) + (y_{v_{\text{nor}}}(j) - y_{v_{\text{nor}}}(k)))}{\sum_{k \in S_{\text{EI}}} (|y_{v_{\text{nor}}}(j) - y_{v_{\text{nor}}}(k)|)}$$

$$(r) = e^{-r/\epsilon^2}$$

$$z_{v_{\text{adj}}}(j) = z_{v_{\text{nor}}}(j) \quad (18)$$

where $\{\vec{w}_{\text{tran}}(x_{w_{\text{tran}}}(j), y_{w_{\text{tran}}}(j)), j \in S_{\text{EI}}\}$ represent the facial feature points on the face image and ϵ is a constant. Then we synthesise the frontal view by two steps. The first step is to transform the adjusted 3D model $\{\vec{v}_{\text{adj}}^{(j)}, j \in U\}$ into frontal view w.r.t. the image plane and the second step is the triangulation texture mapping [35].

When the face in an image deviates significantly from the frontal view, part of the face will be occluded. With our assumption that human faces are of bilateral symmetry, whenever there is a rotation around y -axis, we can reconstruct the occluded part according to the bilateral symmetry of human faces. In our method, if $|y| > 15^\circ$, the frontal view face image is synthesised using symmetry information to compensate the missing information, otherwise symmetry information will not be used.

The pose alignment operations for gallery images and test images are the same. Figs. 5e and f give the synthesised frontal view face images for the gallery images shown in Fig. 5a and for the test images shown in Fig. 5b, respectively. The results indicate that good synthesised frontal images can be obtained by using the proposed pose alignment algorithm.

5 Discriminant waveletface analysis and CNFS classifier

For synthesised frontal view images, many existing face recognition algorithms such as those in [7–11, 16, 17] can be applied. In this section, we adopt discriminant waveletface analysis and CNFS classifier for feature extraction and classification similar to Chien's method [16].

First, 2D discrete wavelet transform is used to decompose a synthesised frontal view image into subimages via the high-pass and low-pass filtering. In our method, the three-level lowest frequency sub-image is extracted as the waveletface represented by vector \mathbf{y} . In general, low frequency components are the most informative subimages gearing with the highest discriminating power. Then we apply LDA to convert the vector \mathbf{y} into a new discriminant feature vector $\mathbf{z} = (\mathbf{W}_{\text{lda}} \cdot \mathbf{y})$. The optimal transform matrix

\mathbf{W}_{lda} is estimated by maximising the ratio of the determinants of between-class scatter matrix \mathbf{S}_b and within-class scatter matrix \mathbf{S}_w of the transformed training samples, i.e. $|\mathbf{W}^T \mathbf{S}_b \mathbf{W}| / |\mathbf{W}^T \mathbf{S}_w \mathbf{W}|$. LDA can pull apart the centroid of different classes and reduce the degree of data scattering within the same class. In our method, the discriminant feature vector \mathbf{z} is used for classification.

The NFS classifier is more robust to facial variations than traditional nearest neighbour and nearest feature line classifiers [16]. In our method, we adapt NFS classifier to a CNFS classifier which can achieve classification for a test image with a set of M potential discriminant feature vectors. Let $\{\mathbf{z}^{(c,1)}, \mathbf{z}^{(c,2)}, \dots, \mathbf{z}^{(c, N_s(c))}\}$ denote the independent discriminant feature vectors associated with class c and $N_s(c)$ be the number of gallery images of the c th candidate. $\{\mathbf{z}^{(c,1)}, \mathbf{z}^{(c,2)}, \dots, \mathbf{z}^{(c, N_s(c))}\}$ span a feature space $\mathcal{S}_{\text{gal}}^{(c)} = \text{sp}\{\mathbf{z}^{(c,1)}, \mathbf{z}^{(c,2)}, \dots, \mathbf{z}^{(c, N_s(c))}\}$. A test image has M potential discriminant feature vectors $\mathbf{z}_{\text{test}}^{(i)}$, $i = 1, \dots, M$ corresponding to its M potential synthesised frontal view images. Let $\mathbf{P}_{\text{proj}}^{(i)}$ be the projection of the i th potential discriminant feature vector $\mathbf{z}_{\text{test}}^{(i)}$ on the feature space $\mathcal{S}_{\text{gal}}^{(i)}$. The distance between the potential discriminant feature $\mathbf{z}_{\text{test}}^{(i)}$ and the gallery feature space $\mathcal{S}_{\text{gal}}^{(i)}$ is calculated by $d(\mathbf{z}_{\text{test}}^{(i)}, \mathcal{S}_{\text{gal}}^{(i)}) = \|\mathbf{z}_{\text{test}}^{(i)} - \mathbf{P}_{\text{proj}}^{(i)}\|$. On the basis of Gram-Schmidt process, a set of orthogonal basis $\{\mathbf{u}^{(i,1)}, \mathbf{u}^{(i,2)}, \dots, \mathbf{u}^{(i, N_s(i))}\}$ on the feature space $\mathcal{S}_{\text{gal}}^{(i)}$ can be calculated, and the projection $\mathbf{P}_{\text{proj}}^{(i)}$ is determined by

$$\mathbf{P}_{\text{proj}}^{(i)} = \frac{\mathbf{z}_{\text{test}}^{(i)} \cdot \mathbf{u}^{(i,1)}}{\mathbf{u}^{(i,1)} \cdot \mathbf{u}^{(i,1)}} + \frac{\mathbf{z}_{\text{test}}^{(i)} \cdot \mathbf{u}^{(i,2)}}{\mathbf{u}^{(i,2)} \cdot \mathbf{u}^{(i,2)}} + \dots$$

$$+ \frac{\mathbf{z}_{\text{test}}^{(i)} \cdot \mathbf{u}^{(i, N_s(i))}}{\mathbf{u}^{(i, N_s(i))} \cdot \mathbf{u}^{(i, N_s(i))}} \quad (19)$$

Thus, for a test image, with all distances $\{d(\mathbf{z}_{\text{test}}^{(c)}, \mathcal{S}_{\text{gal}}^{(c)}), c = 1, 2, \dots, M\}$, the classification result \hat{c} is obtained by

$$d(\hat{c}, \mathcal{S}_{\text{gal}}^{(\hat{c})}) = \min_{1 \leq c \leq M} d(\mathbf{z}_{\text{test}}^{(c)}, \mathcal{S}_{\text{gal}}^{(c)})$$

$$= \min_{1 \leq c \leq M} \|\mathbf{z}_{\text{test}}^{(c)} - \mathbf{P}_{\text{proj}}^{(c)}\| \quad (20)$$

6 Experimental results

Experimental results on real data for pose estimation and frontal view image synthesis are shown in Fig. 5. In this section, synthetic data were utilised to test the MBLPE algorithm in the proposed face recognition method to achieve pose estimation and model adaptation from multiple views. Then the overall performance of the proposed face recognition method was evaluated using two real face image databases.

6.1 Testing the MBLPE algorithm using synthetic data

In these experiments, the width and the height of the generic head model are 170 and 270 pixels, respectively. We set the focal length of lens $f = 10$ pixels and the distance between the origin O of the camera coordinate system and the origin o_v of the 3D face models $D = 2000$. The projection system used here is assumed to be a weak perspective projection system. In addition, we assume that 3D face models will be projected on an image plane of size 2×2 and then digitised to a screen resolution of 256×256 . As a result, the width and the height of the observed head models on the screen plane are about 109 and 173, respectively.

In our experiments, a 3D synthetic face model for a particular person was generated by scaling the default generic wireframe model (Fig. 3) with the model scaling factors $\mathbf{S} = (S_x, S_y, S_z)$ and then adding Gaussian random noise on 3D coordinates of facial feature points on this scaled model. The standard variance of this Gaussian random noise is denoted as σ_{3d} . It was introduced in a symmetric way and the 3D facial feature points of a synthetic face model are described as $\{\vec{v}_{\text{syn}}^{(j)} = (x_{v_{\text{syn}}}(j), y_{v_{\text{syn}}}(j), z_{v_{\text{syn}}}(j))^T, j \in S_{\text{EI}}\}$. Then the synthetic face model was rotated and translated with rotation parameter sets $(\alpha_x(i), \alpha_y(i), \alpha_z(i)), i = 1, \dots, N_s$ and translation vectors $(t_x(i), t_y(i), t_z(i))^T, i = 1, \dots, N_s$, where N_s denotes the total number of the face images for a particular person. The facial feature points on the transformed synthetic faces model were projected to the 2D image plane. In order to simulate the facial feature extraction error, Gaussian random noise with fixed standard variance $\sigma_{2d} = 3$ was added on the x - and y -coordinates of the projected facial feature points in the following experiments except in those to test the performance using different 2D noise. We utilised the MBLPE algorithm to recover the rotation parameter sets $(\alpha'_x(i), \alpha'_y(i), \alpha'_z(i)), i = 1, \dots, N_s$, the translation vectors $(t'_x(i), t'_y(i), t'_z(i))^T, i = 1, \dots, N_s$, and the 3D structure of facial feature points on the synthetic face model $\{\vec{v}_{\text{est}}^{(j)}, j \in S_{\text{EI}}\}$. We utilise the absolute pose error defined as $\Delta \theta_{\text{err}} = (|\alpha'_x(i) - \alpha_x(i)| + |\alpha'_y(i) - \alpha_y(i)| + |\alpha'_z(i) - \alpha_z(i)|)/3$ to measure the accuracy of the recovered pose. We did not test the performance of the recovered translation vectors since they have nothing to do with the following steps in our proposed face recognition method.

First, we investigated the performance of the pose estimation using different N_s which corresponds to the number of input face images. We randomly generated 100 3D synthetic faces with $S_x = 1, S_y = 0.9, S_z = 1.1$ and $\sigma_{3d} = 3$. Each synthetic face was rotated and translated with N_s uniformly distributed random rotation parameter sets $(\alpha_x(i), \alpha_y(i), \alpha_z(i) \in [-50^\circ, +50^\circ]), i = 1, \dots, N_s$, and N_s uniformly distributed random translation vectors $(t_x(i), t_y(i), t_z(i) \in [-100, +100])^T, i = 1, \dots, N_s$. From the 2D projected faces of these face models, the MBLPE algorithm recovered the pose parameter sets for each synthetic face. The average absolute pose error for all the projected face images $\Delta \theta_{\text{err}}$ was calculated. We did the

same experiments on other two sets of 100 synthetic faces generated with parameters $(S_x = 1, S_y = 0.85, S_z = 1.15, \sigma_{3d} = 5)$ and $(S_x = 1, S_y = 0.8, S_z = 1.2, \sigma_{3d} = 7)$ to test the pose estimation performance of the synthetic face model with different extent of deviation from the default generic face model. The results are shown in Fig. 6a from which we find that the larger the N_s is, the more accurate the recovered poses are. Especially the performance of $N_s \geq 2$ is much better than that of $N_s = 1$. It means that the MBLPE algorithm is more appropriate for pose estimation from multiple views than from a single view. However, the performance does not improve significantly when the number of input images N_s is larger than 4.

Model adaptation is also important in the MBLPE algorithm. The performance of model adaptation was tested in the following experiment. For given scaling factors S_x, S_y and S_z and a given σ_{3d} , we randomly generated 100 3D synthetic faces and each synthetic face was rotated and translated with four uniformly distributed random rotation parameter sets $(\alpha_x(i), \alpha_y(i), \alpha_z(i) \in [-50^\circ, +50^\circ]), i = 1, \dots, 4$, and four uniformly distributed random translation vectors $(t_x(i), t_y(i), t_z(i) \in [-100, +100]), i = 1, \dots, 4$. Then we projected these transformed synthetic face models on the 2D image plane and utilised these projected face images to recover the poses and the 3D structure of facial feature points on the synthetic face models. In the MBLPE algorithm, the absolute size of the synthesised model cannot be estimated because the size of the synthesised model is related to the translation vector in a weak perspective projection system. Thus, each set of recovered facial feature points $\{\vec{v}_{\text{est}}^{(j)}, j \in S_{\text{EI}}\}$ in 3D space is just a scaled estimation of the facial feature points $\{\vec{v}_{\text{syn}}^{(j)}, j \in S_{\text{EI}}\}$ on corresponding synthetic face model with an unknown scaling factor. We develop a normalised distance

$$D_{\text{norm}}(\mathbf{v}_{\text{obj}}, \mathbf{v}_{\text{src}}) = \min_{s_{\text{opt}}} \frac{\sum_{j \in S_{\text{EI}}} \|\vec{v}_{\text{obj}}^{(j)} - s_{\text{opt}} \cdot \vec{v}_{\text{src}}^{(j)}\|}{\sum_{j \in S_{\text{EI}}} \|\vec{v}_{\text{obj}}^{(j)}\|} \quad (21)$$

to measure the similarity between the objective facial feature points set $\mathbf{v}_{\text{obj}} = \{\vec{v}_{\text{obj}}^{(j)}, j \in S_{\text{EI}}\}$ and the source facial feature points set $\mathbf{v}_{\text{src}} = \{\vec{v}_{\text{src}}^{(j)}, j \in S_{\text{EI}}\}$ by ignoring the scaling difference. For all the 100 synthetic face

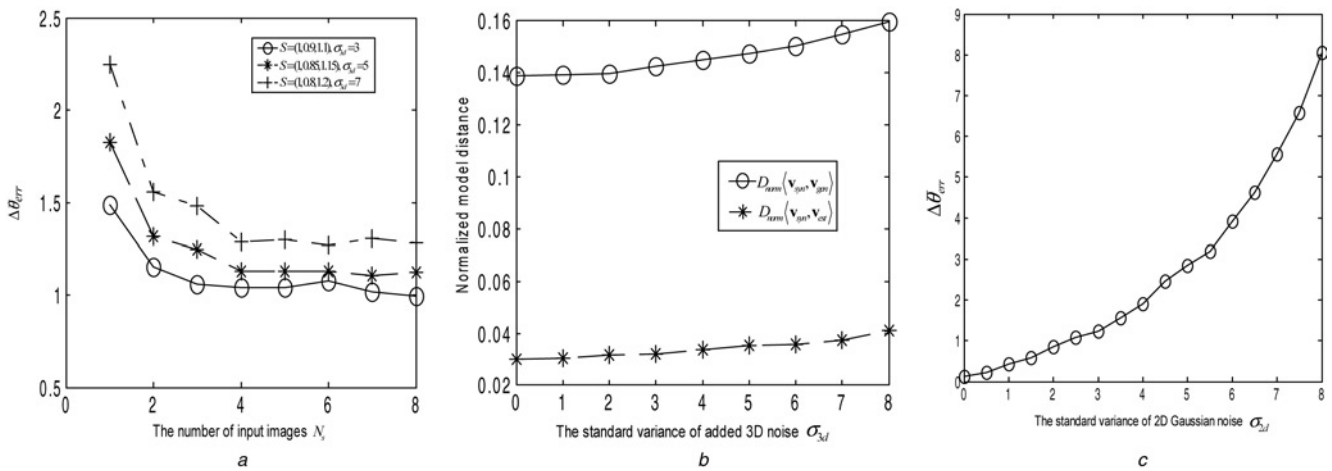


Fig. 6 Results of the pose estimation performance

- a Average absolute errors of the recovered poses using different number of input images with different 3D noise on 3D coordinates of feature points on the scaled face models with different scaling factors
- b Average normalised model distances $D_{\text{norm}}(\mathbf{v}_{\text{syn}}, \mathbf{v}_{\text{est}})$ and $D_{\text{norm}}(\mathbf{v}_{\text{syn}}, \mathbf{v}_{\text{gen}})$ under different 3D noise on 3D coordinates of feature points on scaled face models
- c Average absolute errors of the recovered poses of the test images under different 2D noise on measurement (pixel) of 2D projected feature points

models, we calculated the average normalised distances $D_{\text{norm}}(\mathbf{v}_{\text{syn}}, \mathbf{v}_{\text{est}})$ and $D_{\text{norm}}(\mathbf{v}_{\text{syn}}, \mathbf{v}_{\text{gen}})$, where $\mathbf{v}_{\text{syn}} = \{\bar{\mathbf{v}}_{\text{syn}}^{(j)}, j \in S_{\text{EI}}\}$, $\mathbf{v}_{\text{est}} = \{\bar{\mathbf{v}}_{\text{est}}^{(j)}, j \in S_{\text{EI}}\}$ and $\mathbf{v}_{\text{gen}} = \{\bar{\mathbf{v}}_{\text{gen}}^{(j)}, j \in S_{\text{EI}}\}$ which represent the facial feature points on the default generic model. Two curves of $D_{\text{norm}}(\mathbf{v}_{\text{syn}}, \mathbf{v}_{\text{est}})$ and $D_{\text{norm}}(\mathbf{v}_{\text{syn}}, \mathbf{v}_{\text{gen}})$ with $S_x = 1, S_y = 0.8, S_z = 1.2$ and different $_{3d}$ are shown in Fig. 6b from which we observe that $D_{\text{norm}}(\mathbf{v}_{\text{syn}}, \mathbf{v}_{\text{est}})$ is evidently smaller than $D_{\text{norm}}(\mathbf{v}_{\text{syn}}, \mathbf{v}_{\text{gen}})$ at the same $_{3d}$. It indicates that in average the estimated models are closer to the synthesised faces than the default generic model. Also it can be observed that the larger the $_{3d}$ is, the larger the distance between $D_{\text{norm}}(\mathbf{v}_{\text{syn}}, \mathbf{v}_{\text{est}})$ and $D_{\text{norm}}(\mathbf{v}_{\text{syn}}, \mathbf{v}_{\text{gen}})$ is. It means that the larger the dissimilarity of structure of facial feature points between the synthesised and the default generic model is, the more obvious the effect of model adaptation is.

The preceding experiments were carried out to examine the performance of the pose estimation and model adaptation algorithms for gallery images. We also examined the performance of the pose estimation algorithm for a test image under the condition that the 3D model used in the pose estimation process for the test face image belongs to the same candidate. One hundred synthetic face models were randomly generated with $S_x = 1$ and uniformly distributed $S_y \in [0.8, 1.2]$, $S_z \in [0.8, 1.2]$ and $_{3d} \in [0, 8]$. Each synthetic face model is rotated and translated with five uniformly distributed random rotation parameter sets ($_{x}(i) \in [-30^\circ, 30^\circ]$, $_{y}(i) \in [-40^\circ, 40^\circ]$, $_{z}(i) \in [-20^\circ, 20^\circ]$, $i = 1, \dots, 5$, and five uniform translation vectors ($t_x(i), t_y(i), t_z(i) \in [-100, 100]^T$, $i = 1, \dots, 5$). We use the face images projected from the first four transformed synthetic face model as gallery data, and the face image projected from the fifth transformed synthetic face model is regarded as test data. From the gallery data, each synthetic face model can be estimated using the MBLPE algorithm. Then we recovered the pose of the test images based on these estimated face models. Finally, the average absolute pose error $_{\text{err}}$ for these test images is calculated. The performance of the recovered test images with different $_{2d}$ is shown in Fig. 6(c) from which we observe that $_{\text{err}}$ exponentially increases with the increase of $_{2d}$. Although the MBLPE is affected by 2D noise, fortunately, the absolute pose error is smaller than 4° when $_{2d} \leq 6$. Hence, pose with good accuracy can be obtained when the 2D noise is not too high.



Fig. 7 Examples of face images in ORL face image database

6.2 Overall performance of the face recognition method

IIS face image database (accessible at <http://smart.iis.sinica.edu.tw/html/download.html>) and ORL face image database (accessible at <http://www.cam-orl.co.uk/facedatabase.html>) were used to evaluate the overall performance of the proposed face recognition method. There are 100 persons and 30 pictures for each person in IIS database (Fig. 5). For the 30 pictures of the same candidate, 10 are in nearly frontal view, 10 are in left-side view and 10 are in right-side view. ORL database include 40 persons and each person has 10 pictures with pose variations as shown in Fig. 7. We manually marked the facial feature points as shown in Fig. 3b on each image in these two databases. These marked facial feature points were used in the following experiments. Actually, extracting facial feature points automatically for pose estimation and face recognition should be a possible future work of the proposed approach.

First, we did experiments on IIS database. For each candidate, 4 images were randomly selected to construct the gallery database and the remaining 26 images were used as test images. In the training process, first we performed pose estimation for each image in the gallery database and obtained an estimated 3D model for each candidate by the MBLPE algorithm. On the basis of the estimated pose and the corresponding estimated model, a frontal view face image (size of 112×168 as shown in Fig. 5) was synthesised for each gallery image. Through three-level Harr DWT, we obtained the lowest frequency subimage (size of 14×21) of each synthesised frontal view image. LDA was then performed on these waveletfaces. We obtained the linear transformation matrix and a discriminant feature vector for each gallery sample. The dimension of the

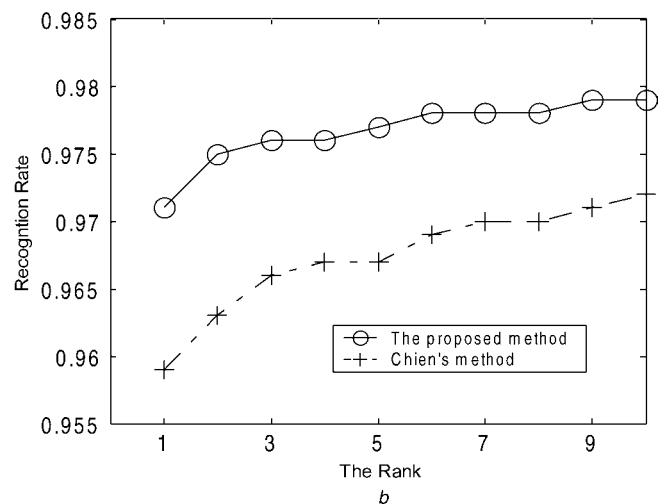
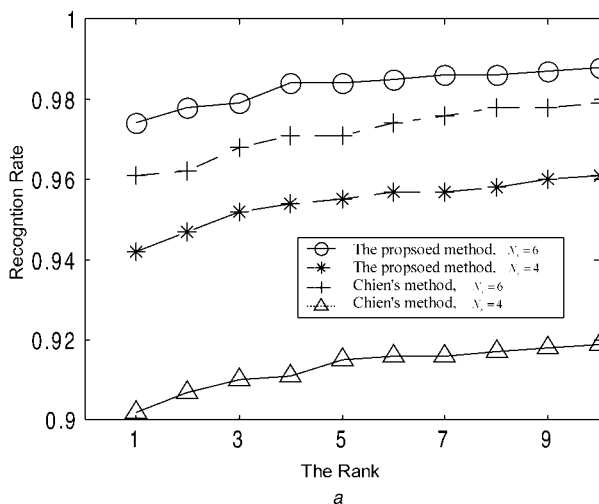


Fig. 8 Face recognition results

a Results on IIS database of the proposed method and Chien's method using four and six training images
b Results on ORL database of the proposed method and Chien's method using five training images

discriminant feature vector in our experiments was set as 60. These feature vectors of each candidate span a feature space. For a test image, we obtained 100 potential pose estimations, 100 corresponding potential synthesised frontal view images, 100 corresponding potential waveletfaces and 100 corresponding potential discriminant feature vectors, since there are 100 candidates in the gallery database. The distance between each potential feature vector and its corresponding feature space, which is spanned by the corresponding feature vectors in the gallery, is calculated. The candidate corresponding to the shortest distance is the first rank recognition result. The candidate corresponding to the second shortest distance is the second-rank recognition result and so on. We developed a program using Visual C++ 6.0 to do these experiments on a P4-1.8 GHz PC. The offline computation of pose estimation and model adaptation for a person with four gallery images generally spends 0.7 s. In the case of 100 candidates in the gallery, the recognition time (not including feature points extraction) is not more than 1 s for a test face.

In order to ensure the statistical robustness of the recognition performance, we performed ten random rounds of face recognition for each situation and all recognition rates were determined by averaging the ten rounds of face recognition. Furthermore, we repeated the preceding experiment using 6 gallery images (the remaining 24 images for test) per candidate. For comparison, we also conducted the experiment for Chien's face recognition method [16] using the same gallery databases and test sets. Fig. 8a shows that our method achieves first-rank recognition rate as high as 94.2% for randomly formed gallery that consists of four images per candidate, whereas Chien's method can just reach 89.3% under the same condition.

We also evaluated the recognition rates of the proposed method using ORL face database. For each person, five images were randomly selected to construct the gallery database and the remaining five images were used as test images. We also performed ten random rounds of face recognition and the recognition rates were determined by averaging the ten rounds of face recognition. For comparison, the same gallery databases and test sets were used to evaluate the performance of Chien's face recognition method. The recognition results shown in Fig. 8b verifies the better performance of the proposed method compared to Chien's method.

7 Conclusions

In this paper, we propose a 3D model based face recognition method which can recognise a human face under variable pose from its multiple views. Using the MBLPE algorithm, we estimate the poses of faces of each person in the gallery and adapt a 3D face model from all faces of this person. The experimental results show that the MBLPE algorithm can achieve efficient 3D model adaptation from multiple faces of a person and estimate poses of these faces with good accuracy. For each face image in the gallery, its frontal view face image is synthesised using the adapted model and estimated pose. Given a test face image, a pose is estimated and a frontal view face image is synthesised using the 3D model for each person in the gallery. Then LDA is performed on these waveletfaces of these synthesised frontal view face images and classification is achieved using CNFS classifier. Experimental results show that the proposed method has a good performance for pose invariant face recognition.

8 References

- Chellappa, R., Willson, C., and Sirohey, S.: 'Human and machine recognition of faces: a survey', *Proc. IEEE*, 1995, **83**, (5), pp. 705–740
- Fromherz, T.: 'Face recognition: a summary of 1995–1997'. Int. Computer Science Inst. ICSI TR-98-027, University of California, Berkeley, 1998
- Zhang, J., Yan, Y., and Lades, M.: 'Face recognition: eigenface, elastic matching, and neural nets', *Proc. IEEE*, 1997, **85**, (9), pp. 1423–1435
- Pentland, A.: 'Looking at people: sensing for ubiquitous and wearable computing', *IEEE Trans. Pattern Anal. Mach. Intell.*, 2000, **22**, (1), pp. 107–119
- Grudin, A.M.: 'On internal representations in face recognition systems', *Pattern Recognit.*, 2000, **33**, (7), pp. 1161–1177
- Chen, C.W., and Huang, J.S.: 'Human face profile recognition from a single front view', *Int. J. Pattern Recognit. Artif. Intell.*, 1992, **6**, (4), pp. 571–593
- Brunelli, R., and Poggio, T.: 'Face recognition: features vs. templates', *IEEE Trans. Pattern Anal. Mach. Intell.*, 1993, **15**, (10), pp. 1042–1053
- Belhumeur, P.N., Hespanha, J.P., and Kriegman, D.J.: 'Eigenfaces vs. fisherfaces: recognition using class specific linear projection', *IEEE Trans. Pattern Anal. Mach. Intell.*, 1997, **19**, (7), pp. 711–720
- Lin, S.H., Kung, S.Y., and Lin, L.J.: 'Face recognition/detection by probabilistic decision-based neural network', *IEEE Trans. Neural Netw.*, 1997, **8**, (1), pp. 114–132
- Li, S., and Lu, J.: 'Face recognition using nearest feature line', *IEEE Trans. Neural Netw.*, 1999, **10**, (2), pp. 439–443
- Bartlett, M.S., Movellan, J.R., and Sejnowski, T.J.: 'Face recognition by independent component analysis', *IEEE Trans. Neural Netw.*, 2002, **13**, (6), pp. 1450–1464
- Murase, H., and Nayar, S.K.: 'Learning and recognition of 3D objects from appearance'. IEEE 2nd Qualitative Vision Workshop, New York, NY, June 1993
- Pentland, A., Moghaddam, B., and Starner, T.: 'View-based modular eigenspaces for face recognition'. Proc. CVPR 1994, June 1994, pp. 84–91
- Huang, F.J., Zhou, Z.H., Zhang, H.J., and Chen, T.: 'Pose invariant face recognition'. Proc. 4th IEEE Int. Conf. Automatic Face and Gesture Recognition, March 2000, pp. 245–250
- Demir, E., Akarun, L., and Alpaydin, E.: 'Two-stage approach for pose invariant face recognition'. Proc. ICASSP2000, June 2000, vol. 6, pp. 5–9
- Chien, J.-T., and Wu, C.-C.: 'Discriminant waveletfaces and nearest feature classifiers for face recognition', *IEEE Trans. Pattern Anal. Mach. Intell.*, 2002, **24**, (12), pp. 1644–1649
- Lu, J.W., Plataniotis, K.N., and Venetsanopoulos, A.N.: 'Face recognition using kernel direct discriminant analysis algorithms', *IEEE Trans. Neural Netw.*, 2003, **14**, (1), pp. 117–126
- De Vel, O., and Aeberhard, S.: 'Line-based face recognition under varying pose', *IEEE Trans. Pattern Anal. and Mach. Intell.*, 1999, **21**, (10), pp. 1081–1088
- Feng, G.C., and Yuen, P.C.: 'Recognition of head-and-shoulder face image using virtual frontal-view image', *IEEE Trans. Syst. Man Cybern. A, Syst. Humans*, 2000, **30**, (6), pp. 871–882
- Chen, Q., Wu, H.Y., Shioyama, S., and Shimada, T.: 'Head pose estimation using both color and feature information'. Proc. 15th ICPR, September 2000, vol. 2, pp. 842–845
- Blanz, V., Romdhani, S., and Vetter, T.: 'Face identification across different poses and illuminations with a 3D morphable model'. Proc. IEEE 5th Int. Conf. Automatic Face and Gesture Recognition, 20–21 May 2002, pp. 100–105
- Blanz, V., and Vetter, T.: 'Face recognition based on fitting a 3D morphable model', *IEEE Trans. Pattern Anal. Mach. Intell.*, 2003, **25**, (19), pp. 1063–1074
- Lam, K.M., and Yan, H.: 'An analytic-to-holistic approach for face recognition based on a single frontal view', *IEEE Trans. Pattern Anal. Mach. Intell.*, 1998, **20**, (7), pp. 673–686
- Gao, Y., Leung, M.K.H., Wang, W., and Hui, S.C.: 'Fast face identification under varying pose from a single 2-D model view', *IEE Proc., Vis. Image Signal Process.*, 2001, **148**, (4), pp. 248–253
- Ishiyama, R., Hamanaka, M., and Sakamoto, S.: 'An appearance model constructed on 3-D surface for robust face recognition against pose and illumination variations', *IEEE Trans. Syst. Man Cybern. C, Appl. Rev.*, 2005, **35**, (3), pp. 326–334
- Zhang, C.Z., and Cohen, F.S.: '3-D face structure extraction and recognition from images using 3-D morphing and distance mapping', *IEEE Trans. Image Process.*, 2002, **11**, (11), pp. 1249–1259

- 27 Yao, J., and Cham, W.K.: 'Efficient model-based linear head motion recovery from movies'. Proc. CVPR2004, July 2004, vol. 2, pp. 414–421
- 28 Or, S.H., Luk, W.S., Wong, K.H., and King, I.: 'An efficient iterative pose estimation algorithm', *Image Vis. Comput.*, 1998, **16**, pp. 353–362
- 29 Lu, C.P., Hager, G.D., and Mjolsness, E.: 'Fast and globally convergent pose estimation from video images', *IEEE Trans. Pattern Anal. Mach. Intell.*, 2000, **22**, (6), pp. 610–622
- 30 Ansar, A., and Daniilidis, K.: 'Linear pose estimation from points or lines', *IEEE Trans. Pattern Anal. Mach. Intell.*, 2003, **25**, (5), pp. 578–589
- 31 Parke, F.I.: 'Parameterized models for facial animation', *IEEE Comput. Graph.*, 1982, **2**, (9), pp. 61–68
- 32 Haralick, R.M., Joo, H., Lee, C.N., Zhuang, X.H., Vaidya, V.G., and Kim, M.B.: 'Pose estimation from corresponding point data', *IEEE Trans. Syst. Man Cybern.*, 1989, **19**, (6), pp. 1426–1446
- 33 Lowe, D.G.: 'Fitting parameterized three-dimensional models to images', *IEEE Trans. Pattern Anal. Mach. Intell.*, 1991, **13**, (5), pp. 441–450
- 34 Schaback, R.: 'Creating surfaces from scattered data using radial basis functions', in Daelhen, M., Lyche, T., and Shumaker, L.L. (Eds.): 'Mathematical methods in computer-aided geometric design III' (Vanderbilt University Press, 1995), pp. 477–496
- 35 Wolberg, G.: 'Digital image warping' (IEEE Computer Society Press, Los Alamos, CA, 1992)