

Fitting a 3D Morphable Model to Edges: A Comparison Between Hard and Soft Correspondences

Anil Bas*

William A. P. Smith*

Timo Bolkart[†]Stefanie Wuhrer[‡]

* Department of Computer Science, University of York, UK
{ab1792, william.smith}@york.ac.uk

[†] Multimodal Computing and Interaction, Saarland University, Germany
tbolkart@mmci.uni-saarland.de

[‡] Morpho Team, INRIA Grenoble Rhône-Alpes, France
stefanie.wuhrer@inria.fr

Abstract. In this paper we explore the problem of fitting a 3D morphable model to single face images using only sparse geometric features (edges and landmark points). Previous approaches to this problem are based on nonlinear optimisation of an edge-derived cost that can be viewed as forming soft correspondences between model and image edges. We propose a novel approach, that explicitly computes hard correspondences. The resulting objective function is non-convex but we show that a good initialisation can be obtained efficiently using alternating linear least squares in a manner similar to the iterated closest point algorithm. We present experimental results on both synthetic and real images and show that our approach outperforms methods that use soft correspondence and other recent methods that rely solely on geometric features.

1 Introduction

Estimating 3D face shape from one or more 2D images is a longstanding problem in computer vision. It has a wide range of applications from pose-invariant face recognition [1] to creation of 3D avatars from 2D images [2]. One of the most successful approaches to this problem is to use a statistical model of 3D face shape [3]. This transforms the problem of shape estimation to one of model fitting and provides a strong statistical prior to constrain the problem.

The model fitting objective can be formulated in various ways, the most obvious being an analysis-by-synthesis approach in which appearance error is directly optimised [3]. However, feature-based methods [4,5] are in general more robust and lead to optimisation problems less prone to convergence on local minima. In this paper, we focus on fitting to edge features in images.

Image edges convey important information about a face. The occluding boundary provides direct information about 3D shape, for example a profile view reveals strong information about the shape of the nose. Internal edges, caused by texture changes, high curvature or self occlusion, provide information about the

position and shape of features such as lips, eyebrows and the nose. This information provides a cue for estimating 3D face shape from 2D images or, more generally, for fitting face models to images.

In Section 2 we introduce relevant background. In Section 3 we present a method for fitting to landmarks with known model correspondence. Our key contribution is in Section 4 where we present a novel, fully automatic algorithm for fitting to image edges with hard correspondence. By hard correspondence, we mean that an explicit correspondence is computed between projected model vertex and edge pixel. For comparison, in Section 5 we describe our variant of previous methods [4,6,7] that fit to edges using soft correspondence. By soft correspondence, we mean that an energy term that captures many possible edge correspondences is minimised. Finally, we compare the two approaches experimentally and others from the recent literature in Section 6.

1.1 Related Work

Landmark fitting 2D landmarks have long been used as a way to initialize a morphable model fit [3]. Breuer et al. [8] obtained this initialisation using a landmark detector providing a fully automatic system. More recently, landmarks have been shown to be sufficient for obtaining useful shape estimates in their own right [9]. Furthermore, noisily detected landmarks can be filtered using a model [10] and automatic landmark detection can be integrated into a fitting algorithm [11]. In a similar manner to landmarks, local features can be used to aid the fitting process [5].

Edge fitting An early example of using image edges for face model fitting is the Active Shape Model (ASM) [12] where a 2D boundary model is aligned to image edges. In 3D, contours have been used directly for 3D face shape estimation [13] and indirectly as a feature for fitting a 3DMM. The earliest work in this direction was due to Moghaddam et al. [14] who fitted a 3DMM to silhouettes extracted from multiple views. From a theoretical standpoint, Lüthi et al. [15] explored to what degree face shape is constrained when contours are fixed.

Romdhani et al. [4] include an edge distance cost as part of a hybrid energy function. Texture and outer (silhouette) contours are used in a similar way to LM-ICP [16] where correspondence between image edges and model contours is “soft”. This is achieved by applying a distance transform to an edge image. This provides a smoothly varying cost surface whose value at a pixel indicates the distance (and its gradient, the direction) to the closest edge. This idea was extended by Amberg et al. [6] who use it in a multi-view setting and smooth the edge distance cost by averaging results with different parameters. In this way, the cost surface also encodes the saliency of an edge. Keller et al. [7] showed that such approaches lead to a cost function that is neither continuous nor differentiable. This suggests the optimisation method must be carefully chosen.

Edge features have also been used in other ways. Cashman and Fitzgibbon [17] learn a 3DMM from 2D images by fitting to silhouettes. Zhu et al. [18] present a method that can be seen as a hybrid of landmark and edge fitting. Landmarks that define boundaries are allowed to slide over the 3D face surface

during fitting. A recent alternative to optimisation-based approaches is to learn a regressor from extracted face contours to 3DMM shape parameters [19].

Fitting a 3DMM to a 2D image using only geometric features (i.e. landmarks and edges) is essentially a non-rigid alignment problem. Surprisingly, the idea of employing an iterated closest point [20] approach with hard edge correspondences (in a similar manner to ASM fitting) has been discounted in the literature [4]. In this paper, we pursue this idea and develop an iterative 3DMM fitting algorithm that is fully automatic, simple and efficient (and we make our implementation available¹). Instead of working in a transformed distance-to-edge space and treating correspondences as “soft”, we compute an explicit correspondence between model and image edges. This allows us to treat the model edge vertices as a landmark with known 2D position, for which optimal pose or shape estimates can be easily computed.

State of the art The most recent face shape estimation methods are able to obtain considerably higher quality results than the purely model-based approaches above. They do so by using pixel-wise shading or motion information to apply finescale refinement to an initial shape estimate. For example, Suwajanakorn et al. [21] use photo collections to build an average model of an individual which is then fitted to a video and finescale detail added by optical flow and shape-from-shading. Cao et al. [22] take a machine learning approach and train a regressor that predicts high resolution shape detail from local appearance.

Our aim in this paper is not to compete directly with these methods. Rather, we seek to understand what quality of face reconstruction it is possible to obtain using solely sparse, geometric information. The output of our method may provide a better initialisation for state of the art refinement techniques or remove the need to have a person specific model.

2 Preliminaries

Our approach is based on fitting a 3DMM to face images under the assumption of a scaled orthographic projection. Hence, we begin by introducing scaled orthographic projection and 3DMMs.

2.1 Scaled Orthographic Projection

The scaled orthographic, or weak perspective, projection model assumes that variation in depth over the object is small relative to the mean distance from camera to object. Under this assumption, the projected 2D position of a 3D point $\mathbf{v} = [u \ v \ w]^T$ given by $\mathbf{SOP}[\mathbf{v}, \mathbf{R}, \mathbf{t}, s] \in \mathbb{R}^2$ does not depend on the distance of the point from the camera, but only on a uniform scale s given by the ratio of the focal length of the camera and the mean distance from camera to object:

$$\mathbf{SOP}[\mathbf{v}, \mathbf{R}, \mathbf{t}, s] = s \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \mathbf{R}\mathbf{v} + s\mathbf{t} \quad (1)$$

¹ Matlab implementation: github.com/waps101/3DMM.edges

where the pose parameters $\mathbf{R} \in \mathbb{R}^{3 \times 3}$, $\mathbf{t} \in \mathbb{R}^2$ and $s \in \mathbb{R}^+$ are a rotation matrix, 2D translation and scale respectively.

2.2 3D Morphable Model

A 3D morphable model is a deformable mesh whose shape is determined by the shape parameters $\boldsymbol{\alpha} \in \mathbb{R}^S$. Shape is described by a linear model learnt from data using Principal Components Analysis (PCA). So, the shape of any face can be approximated as:

$$\mathbf{f}(\boldsymbol{\alpha}) = \mathbf{P}\boldsymbol{\alpha} + \bar{\mathbf{f}}, \quad (2)$$

where $\mathbf{P} \in \mathbb{R}^{3N \times S}$ contains the S principal components, $\bar{\mathbf{f}} \in \mathbb{R}^{3N}$ is the mean shape and the vector $\mathbf{f}(\boldsymbol{\alpha}) \in \mathbb{R}^{3N}$ contains the coordinates of the N vertices, stacked to form a long vector: $\mathbf{f} = [u_1 \ v_1 \ w_1 \ \dots \ u_N \ v_N \ w_N]^T$. Hence, the i th vertex is given by: $\mathbf{v}_i = [f_{3i-2} \ f_{3i-1} \ f_{3i}]^T$. For convenience, we denote the sub-matrix corresponding to the i th vertex as $\mathbf{P}_i \in \mathbb{R}^{3 \times S}$ and the corresponding vertex in the mean face shape as $\bar{\mathbf{f}}_i \in \mathbb{R}^3$, such that the i th vertex is given by: $\mathbf{v}_i = \mathbf{P}_i\boldsymbol{\alpha} + \bar{\mathbf{f}}_i$. Similarly, we define the row corresponding to the u component of the i th vertex as \mathbf{P}_{iu} (similarly for v and w) and define the u component of the i th mean shape vertex as \bar{f}_{iu} (similarly for v and w).

3 Fitting with Known Correspondence

We begin by showing how to fit a morphable model to L observed 2D positions $\mathbf{x}_i = [x_i \ y_i]^T$ ($i = 1 \dots L$) arising from the projection of corresponding vertices in the morphable model. We discuss in Section 4 how these correspondences are obtained in practice. Without loss of generality, we assume that the i th 2D position corresponds to the i th vertex in the morphable model. The objective of fitting a morphable model to these observations is to obtain the shape and pose parameters that minimise the reprojection error, E_{lmk} , between observed and predicted 2D positions:

$$E_{\text{lmk}}(\boldsymbol{\alpha}, \mathbf{R}, \mathbf{t}, s) = \frac{1}{L} \sum_{i=1}^L \|\mathbf{x}_i - \text{SOP}[\mathbf{P}_i\boldsymbol{\alpha} + \bar{\mathbf{f}}_i, \mathbf{R}, \mathbf{t}, s]\|^2. \quad (3)$$

The scale factor in front of the summation makes the magnitude of the error invariant to the number of landmarks. This problem is multilinear in the shape parameters and the SOP transformation matrix. It is also nonlinearly constrained, since \mathbf{R} must be a valid rotation matrix. Although minimising E_{lmk} is a non-convex optimisation problem, a good initialisation can be obtained using alternating linear least squares and this estimate subsequently refined using nonlinear optimisation. This is the approach that we take.

3.1 Pose Estimation

We make an initial estimate of \mathbf{R} , \mathbf{t} and s using a simple extension of the POS algorithm [23]. Compared to POS, we additionally enforce that \mathbf{R} is a valid rotation matrix. We begin by solving an unconstrained system in a least squares sense. We stack two copies of the 3D points in homogeneous coordinates, such that $\mathbf{A}_{2i-1} = [u_i \ v_i \ w_i \ 1 \ 0 \ 0 \ 0 \ 0]$ and $\mathbf{A}_{2i} = [0 \ 0 \ 0 \ 0 \ u_i \ v_i \ w_i \ 1]$ and form a long vector of the corresponding 2D points $\mathbf{d} = [x_1 \ y_1 \ \cdots \ x_L \ y_L]^\top$. We then solve for $\mathbf{k} \in \mathbb{R}^8$ in $\mathbf{A}\mathbf{k} = \mathbf{d}$ using linear least squares. We define $\mathbf{r}_1 = [k_1 \ k_2 \ k_3]$ and $\mathbf{r}_2 = [k_5 \ k_6 \ k_7]$. Scale is given by $s = (\|\mathbf{r}_1\| + \|\mathbf{r}_2\|)/2$ and the translation vector by $\mathbf{t} = [k_4/s \ k_8/s]^\top$. We perform singular value decomposition on the matrix formed from \mathbf{r}_1 and \mathbf{r}_2 :

$$\mathbf{USV}^\top = \begin{bmatrix} \mathbf{r}_1 \\ \mathbf{r}_2 \\ \mathbf{r}_1 \times \mathbf{r}_2 \end{bmatrix} \quad (4)$$

The rotation matrix is given by $\mathbf{R} = \mathbf{UV}^\top$. If $\det(\mathbf{R}) = -1$ then we negate the third row of \mathbf{U} and recompute \mathbf{R} . This guarantees that \mathbf{R} is a valid rotation matrix. This approach gives a good initial estimate which we subsequently refine with nonlinear optimization of E_{lmk} with respect to \mathbf{R} , \mathbf{t} and s .

3.2 Shape Estimation

With a fixed pose estimate, shape parameter estimation under scaled orthographic projection is a linear problem. The 2D position of the i th vertex as a function of the shape parameters is given by: $s\mathbf{R}_{1..2}(\mathbf{P}_i\boldsymbol{\alpha} + \bar{\mathbf{f}}_i) + s\mathbf{t}$. Hence, each observed vertex adds two equations to a linear system. Concretely, for each image we form the matrix $\mathbf{C} \in \mathbb{R}^{2L \times S}$ where

$$\mathbf{C}_{2i-1} = s(\mathbf{R}_{11}\mathbf{P}_{iu}^\top + \mathbf{R}_{12}\mathbf{P}_{iv}^\top + \mathbf{R}_{13}\mathbf{P}_{iw}^\top)$$

and

$$\mathbf{C}_{2i} = s(\mathbf{R}_{21}\mathbf{P}_{iu}^\top + \mathbf{R}_{22}\mathbf{P}_{iv}^\top + \mathbf{R}_{23}\mathbf{P}_{iw}^\top)$$

and vector $\mathbf{h} \in \mathbb{R}^{2L}$ where

$$\mathbf{h}_{2i-1} = x_i - s(\mathbf{R}_1\bar{\mathbf{f}}_i + \mathbf{t}_1) \quad \text{and} \quad \mathbf{h}_{2i} = y_i - s(\mathbf{R}_2\bar{\mathbf{f}}_i + \mathbf{t}_2).$$

We solve $\mathbf{C}\boldsymbol{\alpha} = \mathbf{h}$ in a least squares sense subject to an additional constraint to ensure plausibility of the solution. We follow Brunton et al. [24] and use a hyperbox constraint on the shape parameters. This avoids having to choose a regularisation weight but ensures that each parameter lies within k standard deviations of the mean by introducing a linear inequality constraint on the shape parameters (we use $k = 3$ in our experiments). Hence, the problem can be solved in closed form as an inequality constrained linear least squares problem.

3.3 Nonlinear Refinement

Having alternated pose and shape estimation for a fixed number of iterations, finally we perform nonlinear optimisation of E_{lmk} over α , \mathbf{R} , \mathbf{t} and s simultaneously. We represent \mathbf{R} in axis-angle space to ensure that it remains a valid rotation matrix and we retain the hyperbox constraint on α . We minimise E_{lmk} using the trust-region-reflective algorithm [25] as implemented in the Matlab `lsqnonlin` function.

4 Fitting with Hard Edge Correspondence

The method in Section 3 enables a 3DMM to be fitted to 2D landmark positions if the correspondence between landmarks and model vertices is known. Edges, for example caused by occluding boundaries, do not have a fixed correspondence to model vertices. Hence, fitting to edges requires shape and pose estimation to happen in conjunction with establishing correspondence between image and model edges. Our proposed approach establishes these correspondences explicitly by finding the closest image edge to each model boundary vertex (subject to additional filtering to remove unreliable matches). Our method comprises the following steps:

1. Detect facial landmarks
2. Initialise shape and pose estimates by fitting to landmarks only
3. Improve initialisation using iterated closest edge fitting
4. Nonlinear optimisation of hybrid objective function containing landmark, edge and prior terms

We describe each of these steps in more detail in the rest of this section.

4.1 Landmarks

We use landmarks both for initialisation and as part of our overall objective function as one cue for shape estimation. We apply a facial landmark detector that is suitable for operating on “in the wild” images. This provides approximate positions of facial landmarks for which we know the corresponding vertices in the morphable model. We use these landmark positions to make an initial estimate of the pose and shape parameters by running the method in Section 3 with only these corresponding landmark positions. Note that any facial landmark detector can be used at this stage. In our experiments, we show results with a recent landmark detection algorithm [26] that achieves state-of-the-art performance and for which code is provided by the authors. In our experimental evaluation, we include the results of fitting to landmarks only.

4.2 Edge Cost

We assume that a subset of pixels have been labelled as edges and stored as the set $\mathcal{E} = \{(x, y) | (x, y) \text{ is an edge}\}$. In practice, we compute edges by applying the Canny edge detector with a fixed threshold to the input image.

Model contours are computed based on the pose and shape parameters as the occluding boundary of the 3D face. The set of occluding boundary vertices, $\mathcal{B}(\boldsymbol{\alpha}, \mathbf{R}, \mathbf{t}, s)$, are defined as those lying on a mesh edge whose adjacent faces have a change of visibility. This definition encompasses both outer (silhouette) and inner (self-occluding) contours. Since the viewing direction is aligned with the z -axis, this is tested simply by checking if the sign of the z -component of the triangle normal changes on either side of the edge. In addition, we check that potential edge vertices are not occluded by another part of the mesh (using z -buffering) and we ignore edges that lie on a mesh boundary since they introduce artificial edges. In this paper, we deal only with occluding contours (both inner and outer). If texture contours were defined on the surface of the morphable model, it would be straightforward to include these in our approach.

We define the objective function for edge fitting with hard correspondence as the sum of squared distances between each projected occluding boundary vertex and the closest edge pixel:

$$E_{\text{edge}}(\boldsymbol{\alpha}, \mathbf{R}, \mathbf{t}, s) = \frac{1}{|\mathcal{B}(\boldsymbol{\alpha}, \mathbf{R}, \mathbf{t}, s)|} \sum_{i \in \mathcal{B}(\boldsymbol{\alpha}, \mathbf{R}, \mathbf{t}, s)} \min_{(x, y) \in \mathcal{E}} \|[x \ y]^T - \text{SOP}[\mathbf{P}_i \boldsymbol{\alpha} + \bar{\mathbf{f}}_i, \mathbf{R}, \mathbf{t}, s]\|^2. \quad (5)$$

Note that the minimum operator is responsible for computing the hard correspondences. This objective is non-convex since the minimum of a set of convex functions is not convex [27]. Hence, we require a good initialisation to ensure convergence to a minimum close to the global optimum. Fitting to landmarks only does not provide a sufficiently good initialisation. So, in the next subsection we describe a method for obtaining a good initial fit to edges, before incorporating the edge cost into a hybrid objective function in Section 4.5.

4.3 Iterated Closest Edge Fitting

We propose to refine the landmark-only fit with an initial fit to edges that works in an iterated closest point manner. That is, for each projected model contour vertex, we find the closest image edge pixel and we treat this as a known correspondence. In conjunction with the landmark correspondences, we again run the method in Section 3. This leads to updated pose and shape parameters, and in turn to updated model edges and correspondences. We iterate this process for a fixed number of iterations. We refer to this process as Iterated Closest Edge Fitting (ICEF) and provide an illustration in Figure 1. On the left we show an input image with the initial landmark detection result. In the middle we show the initial shape and pose obtained by fitting only to landmarks. On the right we show image edge pixels in blue and projected model contours in green (where

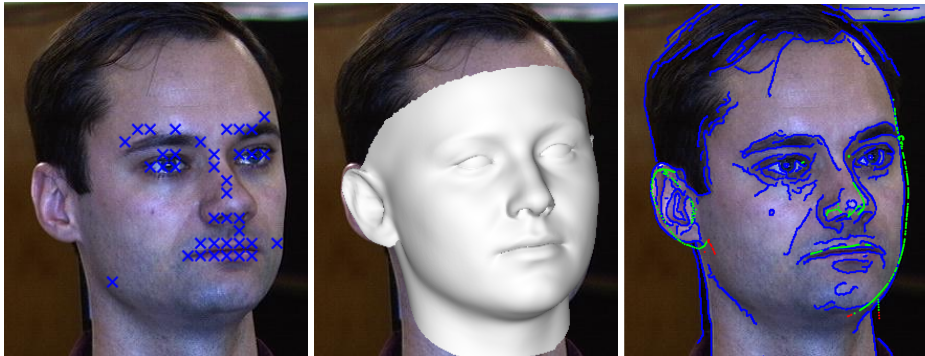


Fig. 1. Iterated closest edge fitting for initialisation of the edge fitting process. Left: input image with automatically detected landmarks. Middle: overlaid shape obtained by fitting only to landmark. Right: image edges in blue, model boundary vertices with image correspondences in green, unreliable correspondences in red.

nearest neighbour edge correspondence is considered reliable) and in red (where correspondence is considered unreliable). The green/blue correspondences are used for the next iteration of fitting.

Finding the image edge pixel closest to a projected contour vertex can be done efficiently by storing the image edge pixels in a kd -tree. We filter the resulting correspondences using two commonly used heuristics. First, we remove 5% of the matches for which the distance to the closest image edge pixel is largest. Second, we remove matches for which the image distance divided by s exceeds a threshold (chosen as 10 in our experiments). The division by scale factor s makes this choice invariant to changes in image resolution.

4.4 Prior

Under the assumption that the training data of the 3DMM forms a Gaussian cloud in high dimensional space, then we expect that each of the shape parameters follows a normal distribution with zero mean and variance given by the eigenvalue, λ_i , associated with the corresponding principal component. We find that including a prior term that captures this assumption significantly improves performance over using the hyperbox constraint alone. The prior penalises deviation from the mean shape as follows:

$$E_{\text{prior}}(\boldsymbol{\alpha}) = \sum_{i=1}^S \left(\frac{\alpha_i}{\sqrt{\lambda_i}} \right)^2. \quad (6)$$

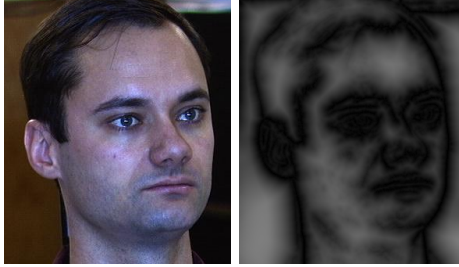


Fig. 2. Edge cost surface with soft correspondence (right) computed from input image (left)

4.5 Nonlinear Refinement

Finally, we perform nonlinear optimisation of a hybrid objective function comprising landmark, edge and prior terms:

$$E(\boldsymbol{\alpha}, \mathbf{R}, \mathbf{t}, s) = w_1 E_{\text{lmk}}(\boldsymbol{\alpha}, \mathbf{R}, \mathbf{t}, s) + w_2 E_{\text{edge}}(\boldsymbol{\alpha}, \mathbf{R}, \mathbf{t}, s) + w_3 E_{\text{prior}}(\boldsymbol{\alpha}), \quad (7)$$

where w_1 , w_2 and w_3 weight the contribution of each term to the overall energy. The landmark and edge terms are invariant to the number of landmarks and edge vertices which means we do not have to tune the weights for each image (for example, for the results in Table ?? we use fixed values of: $w_1 = 0.15$, $w_2 = 0.45$ and $w_3 = 0.4$). We retain the hyperbox constraint and so the hybrid objective is a constrained nonlinear least squares problem and we again optimise using the trust-region-reflective algorithm.

For efficiency and to avoid problems of continuity and differentiability of the edge cost function, we follow [6] and keep occluding boundary vertices, \mathcal{B} , fixed for a number of iterations of the optimiser. After a number of iterations, we recompute the vertices lying on the occluding boundary and restart the optimiser.

5 Fitting with Soft Edge Correspondence

We compare our approach with a method based on optimising an edge cost function, in the same spirit as previous work [6,7,4]. We follow the same approach as Amberg et al. [6] to compute the edge cost function, however we further improve robustness by also integrating over scale. For our edge detector, we use gradient magnitude thresholding with non-maxima suppression. Given a set of edge detector sensitivity thresholds \mathcal{T} and scales \mathcal{S} , we compute $n = |\mathcal{T} \times \mathcal{S}|$ edge images, E^1, \dots, E^n , using each pair of image scale and threshold values. We compute the Euclidean distance transform, D^1, \dots, D^n , for each edge image (i.e. the value of each pixel in D^i is the distance to the closest edge pixel in E^i). Finally, we compute the edge cost surface as:

$$S(x, y) = \frac{1}{n} \sum_{i=1}^n \frac{D^i(x, y)}{D^i(x, y) + \kappa}. \quad (8)$$



Fig. 3. Synthetic input images for one subject

Method	Rotation angle									Mean
	-70°	-50°	-30°	-15°	0°	15°	30°	50°	70°	
Average face	3.35	3.35	3.35	3.35	3.35	3.35	3.35	3.35	3.35	3.35
Proposed (landmarks only)	2.67	2.60	2.58	2.64	2.56	2.49	2.50	2.54	2.63	2.58
Aldrian and Smith [9]	2.64	2.60	2.55	2.54	2.49	2.42	2.43	2.44	2.54	2.52
Romdhani et al. [4] (soft)	2.65	2.59	2.58	2.61	2.59	2.50	2.50	2.46	2.51	2.55
Proposed (ICEF)	2.38	2.40	2.51	2.38	2.52	2.45	2.43	2.38	2.3	2.42
Proposed (hard)	2.35	2.26	2.38	2.40	2.51	2.39	2.40	2.20	2.26	2.35

Table 1. Mean Euclidean vertex distance (mm) with ground truth landmarks

The parameter κ determines the influence range of an edge in an adaptive manner. Amberg et al. [6] suggest a value for κ of 1/20th the expected size of the head in pixels. We compute this parameter automatically from the scale s . An example of an edge cost surface is shown in Figure 2. To evaluate the edge cost, we compute model contour vertices as in Section 4.2, project them into the image and interpolate the edge cost function using bilinear interpolation:

$$E_{\text{softedge}}(\boldsymbol{\alpha}, \mathbf{R}, \mathbf{t}, s) = \frac{1}{|\mathcal{B}(\boldsymbol{\alpha}, \mathbf{R}, \mathbf{t}, s)|} \sum_{i \in \mathcal{B}(\boldsymbol{\alpha}, \mathbf{R}, \mathbf{t}, s)} S(\text{SOP} [\mathbf{P}_i \boldsymbol{\alpha} + \bar{\mathbf{f}}_i, \mathbf{R}, \mathbf{t}, s]). \quad (9)$$

As with the hard edge cost, we found that the best performance was achieved by also including the landmark and prior terms in a hybrid objective function. Hence, we minimise:

$$E(\boldsymbol{\alpha}, \mathbf{R}, \mathbf{t}, s) = w_1 E_{\text{lmk}}(\boldsymbol{\alpha}, \mathbf{R}, \mathbf{t}, s) + w_2 E_{\text{softedge}}(\boldsymbol{\alpha}, \mathbf{R}, \mathbf{t}, s) + w_3 E_{\text{prior}}(\boldsymbol{\alpha}). \quad (10)$$

We again initialise by fitting to landmarks only using the method in Section 4.1, retain the hyperbox constraint and optimise using the trust-region-reflective algorithm. We use the same weights as for the hard correspondence method in our experiments.

6 Experimental Results

We present two sets of experimental results. First, we use synthetic images with known ground truth 3D shape in order to quantitatively evaluate our method and provide comparison to previous work. Second, we use real images to provide qualitative evidence of the performance of our method in uncontrolled conditions. For the 3DMM in both sets of experiments we use the Basel Face Model [28].

Method	Landmark noise std. dev.					
	$\sigma = 0$	$\sigma = 1$	$\sigma = 2$	$\sigma = 3$	$\sigma = 4$	$\sigma = 5$
Proposed (landmarks only)	2.58	2.60	2.61	2.68	2.76	2.85
Aldrian and Smith [9]	2.52	2.53	2.55	2.62	2.65	2.73
Romdhani et al. [4] (soft)	2.55	2.57	2.57	2.62	2.70	2.76
Proposed (ICEF)	2.42	2.43	2.43	2.50	2.57	2.60
Proposed (hard)	2.35	2.36	2.35	2.39	2.47	2.50

Table 2. Mean Euclidean vertex distance (mm) with noisy landmarks

Method	Rotation angle									Mean
	-70°	-50°	-30°	-15°	0°	15°	30°	50°	70°	
Proposed (landmarks only)	6.79	6.84	5.19	5.74	5.68	6.34	6.48	7.04	7.74	6.43
Zhu et al. [18]	N/A	N/A	4.63	5.09	4.19	5.22	4.92	N/A	N/A	N/A
Romdhani et al. [4] (soft)	4.46	3.42	3.66	3.78	3.77	3.57	4.31	4.19	4.73	3.99
Proposed (ICEF)	3.70	3.32	3.26	3.23	3.37	3.50	3.43	4.07	3.52	3.49
Proposed (hard)	3.43	3.20	3.19	3.09	3.30	3.36	3.36	3.84	3.41	3.35

Table 3. Mean Euclidean vertex distance (mm) with automatically detected landmarks

6.1 Quantitative Evaluation

We begin with a quantitative comparative evaluation on synthetic data. We use the 10 out-of-sample faces supplied with the Basel Face Model and render orthographic images of each face in 9 poses (rotations of 0° , $\pm 15^\circ$, $\pm 30^\circ$, $\pm 50^\circ$ and $\pm 70^\circ$ about the vertical axis). We show sample input images for one subject in Figure 3. In all experiments, we report the mean Euclidean distance between ground truth and estimated face surface in mm after Procrustes alignment.

In the first experiment, we use ground truth landmarks. Specifically, we use the 70 Farkas landmarks, project the visible subset to the image (yielding between 37 and 65 landmarks per image) and round to the nearest pixel. In Table 1 we show results averaged over pose angle and over the whole dataset. As a baseline, we show the error if we simply use the average face shape. We then show the result of fitting only to landmarks, i.e. the method in Section 3. We include two comparison methods. The approach of Aldrian and Smith [9] uses only landmarks but with an affine camera model and a learnt model of landmark variance. The soft edge correspondence method of Romdhani et al. [4] is described in Section 5. The final two rows show two variants of our proposed methods: the fast Iterated Closest Edge Fitting version and the full version with nonlinear optimisation of the hard correspondence cost. Average performance over the whole dataset is best for our method and, in general, using edges over landmarks only and applying nonlinear optimisation improves performance. The performance improvement of our methods over landmark-only methods improves with pose angle. This suggest that edge information becomes more salient for non-frontal poses.

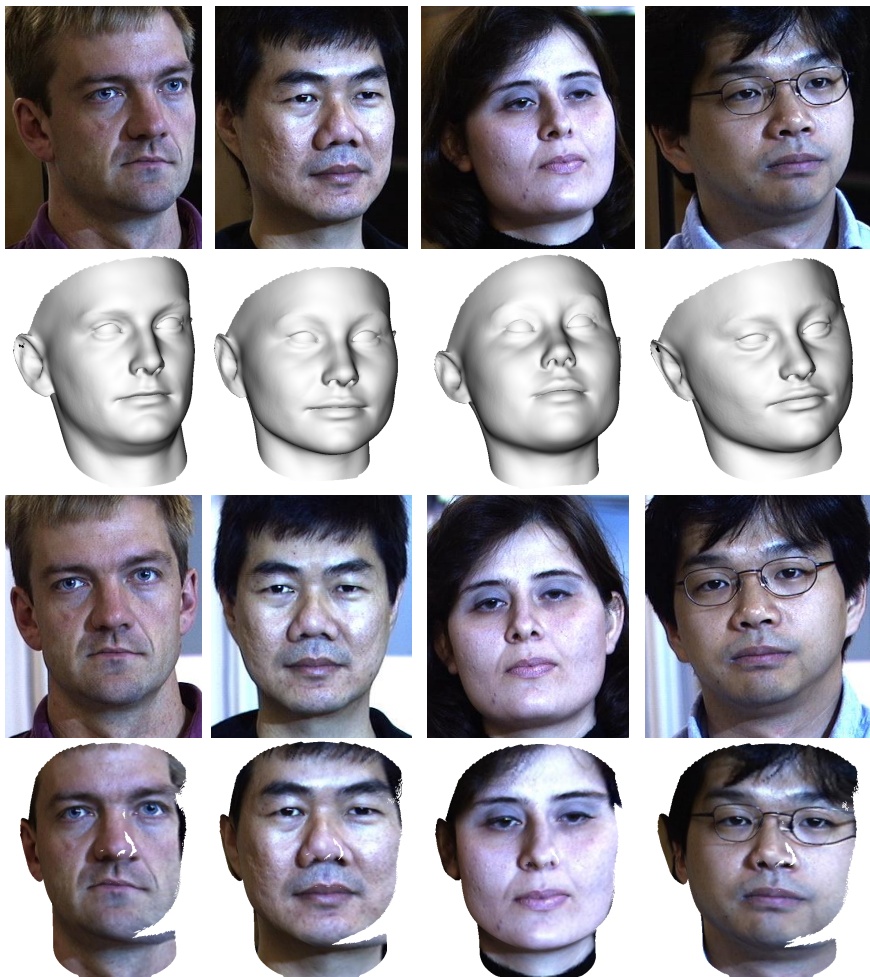


Fig. 4. Qualitative frontalisation results

The second experiment is identical to the first except that we add Gaussian noise of varying standard deviation to the ground truth landmark positions. In Table 2 we show results averaged over all poses and subjects.

In the final experiment we use landmarks that are automatically detected using the method of Zhu and Ramanan [26]. This enables us to include comparison with the recent fitting algorithm of Zhu et al. [18]. We use the author's own implementation which only works with a fixed set of 68 landmarks. This means that the method cannot be applied to the more extreme pose angles where fewer landmarks are detected. In this more challenging scenario, our method again gives the best overall performance and is superior for all pose angles.



Fig. 5. Qualitative pose editing results

6.2 Qualitative Evaluation

In Figure 4 we show qualitative examples from the CMU PIE [29] dataset. Here, we fit to images (first row) in a non-frontal pose using automatically detected landmarks [26] and show the reconstruction in the second row. We texture map the image onto the mesh, rotate to frontal pose (bottom row) and compare to an actual frontal view (third row). Finally, we show qualitative examples from the Labelled Faces in the Wild dataset [30] in Figure 5. Again, we texture map the image to the mesh and show a range of poses. These results show that our method is capable of robustly and fully automatically fitting to unconstrained images.

7 Conclusions

We have presented a fully automatic algorithm for fitting a 3DMM to single images using hard edge correspondence and compared it to existing methods using soft correspondence. In 3D-3D alignment, the soft correspondence of LM-ICP [16] is demonstrably more robust than hard ICP [20]. However, in the context of 3D-2D nonrigid alignment, a soft edge cost function is neither continuous nor differentiable since contours appear, disappear, split and merge under parameter changes [7]. This makes its optimisation challenging, unstable and highly dependent on careful choice of optimisation parameters. Although our proposed algorithm relies on potentially brittle hard correspondences, solving for shape and pose separately requires only solution of a linear problem and, together, optimisation of a multilinear problem. This makes iterated closest edge fitting very fast and it provides an initialisation that allows the subsequent nonlinear optimisation to converge to a better optimum. We believe that this explains the improved performance over edge fitting with soft correspondence.

There are many ways this work can be extended. First, we could explore other ways in which the notion of soft correspondence is formulated. For example,

we could borrow from SoftPOSIT [31] or Blind PnP [32] which both estimate pose with unknown 3D-2D correspondence. Second, we could incorporate any of the refinements to standard ICP [33]. Third, we currently use only geometric information and do not fit texture. Finally, we would like to extend the method to video using a model that captures expression variation and incorporating temporal smoothness constraints.

References

1. Blanz, V., Vetter, T.: Face recognition based on fitting a 3D morphable model. *IEEE Trans. Pattern Anal. Mach. Intell.* **25** (2003) 1063–1074
2. Ichim, A.E., Bouaziz, S., Pauly, M.: Dynamic 3D avatar creation from hand-held video input. *ACM Trans. Graph.* **34** (2015) 45
3. Blanz, V., Vetter, T.: A morphable model for the synthesis of 3D faces. In: *SIGGRAPH*. (1999)
4. Romdhani, S., Vetter, T.: Estimating 3D shape and texture using pixel intensity, edges, specular highlights, texture constraints and a prior. In: *CVPR*. (2005)
5. Huber, P., Feng, Z., Christmas, W., Kittler, J., Räscher, M.: Fitting 3D morphable models using local features. In: *ICIP*. (2015)
6. Amberg, B., Blake, A., Fitzgibbon, A., Romdhani, S., Vetter, T.: Reconstructing high quality face-surfaces using model based stereo. In: *ICCV*. (2007)
7. Keller, M., Knothe, R., Vetter, T.: 3D reconstruction of human faces from occluding contours. In: *MIRAGE*. (2007)
8. Breuer, P., Kim, K., Kienzle, W., Schölkopf, B., Blanz, V.: Automatic 3D face reconstruction from single images or video. In: *Proc. FG*. (2008) 1–8
9. Aldrian, O., Smith, W.A.P.: Inverse rendering of faces with a 3D morphable model. *IEEE Trans. Pattern Anal. Mach. Intell.* **35** (2013) 1080–1093
10. Amberg, B., Vetter, T.: Optimal landmark detection using shape models and branch and bound. In: *Proc. ICCV*. (2011)
11. Schönborn, S., Forster, A., Egger, B., Vetter, T.: A monte carlo strategy to integrate detection and model-based face analysis. *Patt. Rec.* (2013) 101–110
12. Cootes, T.F., Taylor, C.J., Cooper, D., Graham, J.: Active shape models – their training and application. *Comput. Vis. Image Underst.* **61** (1995) 38–59
13. Atkinson, G.A., Smith, M.L., Smith, L.N., Farooq, A.R.: Facial geometry estimation using photometric stereo and profile views. In: *Proc. ICB*. (2009)
14. Moghaddam, B., Lee, J., Pfister, H., Machiraju, R.: Model-based 3D face capture with shape-from-silhouettes. In: *Proc. FG*. (2003)
15. Lüthi, M., Albrecht, T., Vetter, T.: Probabilistic modeling and visualization of the flexibility in morphable models. In: *Math. of Surf. XIII*. (2009)
16. Fitzgibbon, A.W.: Robust registration of 2D and 3D point sets. *Image Vis. Comput.* **21** (2003) 1145–1153
17. Cashman, T.J., Fitzgibbon, A.W.: What shape are dolphins? Building 3D morphable models from 2d images. *IEEE Trans. Pattern Anal. Mach. Intell.* **35** (2013) 232–244
18. Zhu, X., Lei, Z., Yan, J., Yi, D., Li, S.Z.: High-fidelity pose and expression normalization for face recognition in the wild. In: *Proc. CVPR*. (2015) 787–796
19. Sánchez-Escobedo, D., Castelán, M., Smith, W.: Statistical 3D face shape estimation from occluding contours. *Comput. Vis. Image Underst.* **142** (2016)

20. Besl, P.J., McKay, N.D.: A method for registration of 3-D shapes. *IEEE Trans. Pattern Anal. Mach. Intell.* **14** (1992) 239–256
21. Suwajanakorn, S., Kemelmacher-Shlizerman, I., Seitz, S.M.: Total moving face reconstruction. In: *ECCV*. (2014)
22. Cao, C., Bradley, D., Zhou, K., Beeler, T.: Real-time high-fidelity facial performance capture. *ACM Trans. Graph.* **34** (2015) 46
23. Dementhon, D.F., Davis, L.S.: Model-based object pose in 25 lines of code. *Int. J. Comput. Vis.* **15** (1995) 123–141
24. Brunton, A., Salazar, A., Bolkart, T., Wuhrer, S.: Review of statistical shape spaces for 3D data with comparative analysis for human faces. *Comput. Vis. Image Underst.* **128** (2014) 1–17
25. Coleman, T., Li, Y.: An interior, trust region approach for nonlinear minimization subject to bounds. *SIAM J. Optimiz.* **6** (1996) 418–445
26. Zhu, X., Ramanan, D.: Face detection, pose estimation, and landmark localization in the wild. In: *Proc. CVPR*. (2012)
27. Grant, M., Boyd, S., Ye, Y.: Disciplined convex programming. In: *Global Optimization: From Theory to Implementation*. Springer (2006) 155–210
28. Paysan, P., Knothe, R., Amberg, B., Romdhani, S., Vetter, T.: A 3D face model for pose and illumination invariant face recognition. In: *Proc. AVSS*. (2009)
29. Sim, T., Baker, S., Bsat, M.: The CMU pose, illumination, and expression database. *IEEE Trans. Pattern Anal. Mach. Intell.* **25** (2003) 1615–1618
30. Huang, G.B., Ramesh, M., Berg, T., Learned-Miller, E.: Labeled faces in the wild: A database for studying face recognition in unconstrained environments. Technical Report 07-49, University of Massachusetts, Amherst (2007)
31. David, P., DeMenthon, D., Duraiswami, R., Samet, H.: SoftPOSIT: Simultaneous pose and correspondence determination. In: *Proc. ECCV*. (2002) 698–714
32. Moreno-Noguer, F., Lepetit, V., Fua, P.: Pose priors for simultaneously solving alignment and correspondence. In: *Proc. ECCV*. (2008) 405–418
33. Rusinkiewicz, S., Levoy, M.: Efficient variants of the ICP algorithm. In: *Proc. 3DIM*. (2001)