

Dense Surface Point Distribution Models of the Human Face

Tim J. Hutton
Biomedical Informatics Unit
Eastman Dental Institute
University College London
London, UK
T.Hutton@eastman.ucl.ac.uk

Bernard F. Buxton
Computer Science Department
University College London
London, UK
B.Buxton@cs.ucl.ac.uk

Peter Hammond
Biomedical Informatics Unit
Eastman Dental Institute
University College London
London, UK
P.Hammond@eastman.ucl.ac.uk

Abstract

In this paper¹ we show how a dense surface model of the human face can be built from a population of examples. A technique that combines active shape models (ASMs) with iterative closest point (ICP) can be used to fit the model to new faces. The model is built by aligning the surfaces using a sparse set of hand-placed landmarks, then using thin-plate spline warping to make a dense correspondence with a base mesh. All of the mesh vertices are then used as landmarks to build a 3D point distribution model. The dense surface point distribution model is more sensitive than the landmark model to correlated facial characteristics such as gender, age and the presence of congenital abnormalities.

1. Introduction

The building of a detailed model of complex objects such as the surface of the human face is problematic because of the large number of correspondences that need to be made to capture the subtle shape changes across the surface. Placing these landmarks manually is implausible not only because of the large number needed (at least 2,000 for the subtleties of the human face) but also because there is insufficient guidance for their location - there are perhaps only 10 or

20 recognisable points, the rest must be interpolated.

Previous approaches to establishing a dense correspondence between 3D surfaces have typically used curvature similarities [15, 2, 17, 10] but such methods are not necessarily ideal for modelling surfaces that exhibit large deformation or for structures that are not initially registered.

The approach suggested in [5] is to first establish a rigid correspondence between pairs of shapes, making use of a decimated version of each, then to fill-in the dense vertices using a brushfire algorithm across each surface. The method requires the surfaces to be sufficiently similar in shape that a rigid-body match will converge to give the correct correspondence. For classes of object that exhibit large shape changes it is not clear that the method will continue to give good results. Also, they demonstrate their method on a set of closed surfaces but do not explain how the method could make use of open surfaces such as the human face where the boundary of the area of interest is poorly defined and the extent of the input data varies from example to example.

Lorenz and Krahnstöver [13] show an improved method for building dense surface models that is similar to ours but which is again only demonstrated on closed surfaces. All the input meshes are warped onto the hand-placed landmarks of a single example, then a coating procedure is used to resample each surface to solve the correspondence problem. A mesh regularization step is included to ensure that folds in the surface introduced by the coating procedure do not appear in the final model. They demonstrate their method on a training set of 31 lumbar vertebrae, using 15 landmarks to produce a point distribution model with ap-

¹to be presented at the Workshop Mathematical Methods in Biomedical Image Analysis (MMBIA) as part of the IEEE International Conference on Computer Vision and Pattern Recognition (CVPR), Kauai, Hawaii, December 2001.

proximately 600 vertices.

As in [13], our approach sacrifices full automation at the model-building stage in exchange for more robustness to large deformation, since it allows the user to annotate whatever shape changes are necessary to bring surfaces into alignment. We find that a very detailed model of the human face can be built using a remarkably sparse set of landmarks.

Our approach is to use thin-plate spline warping [4] to establish a dense set of correspondences between surfaces based on a set of hand-placed landmarks. The dense correspondences are used as the input to build a point distribution model [7] of nearly 2,000 vertices in a surface mesh. The deformable model can then be used to interpret unseen face scans using a combination of iterative closest point [3] and active shape model fitting [7].

2. Method

The training set consisted of 193 scans of the human face. The images were acquired on a DSP400 face-scanner made by TCTi (www.tcti.com) which is based on digital stereo photogrammetry. Each comprised a triangulated surface mesh with between 2,000 and 10,000 vertices (see Fig. 1). The nature of the acquisition is such that structures other than the face area of interest are frequently included. The model-building method deals with these automatically as we shall show.

The selected images covered a broad age range (1 to 70 years) as well as variation in gender, ethnic origin and facial expression, in order to demonstrate the model’s capacity to deal with such shape differences.

2.1 Landmarking the training set

The first step is to place landmarks manually on each surface. We found that just 9 landmarks were enough to build a good model with our dataset. Experimentation showed that too few (5) landmarks resulted in anatomical mis-correspondences, while too many (20) introduced noise into the model since there is insufficient visual guidance for their placement.

Figure 1(a) shows an example mesh with the landmarks overlaid. These landmarks are placed manually on the 2D texture and mapped onto the 3D surface by first mapping each mesh vertex to its corresponding 2D texture coordinate. The parametric coordinates of the texture triangle containing the mouse click location determine its location within the corresponding 3D surface triangle. Figure 1(b) shows a detail of the mesh around the eye with the landmarks as spheres.

Placing the landmarks on the 2D texture is intuitive to the user and is often easier than placing the landmarks on

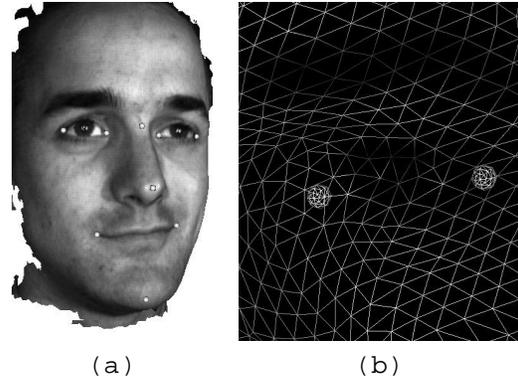


Figure 1. An example scan with our nine landmarks (a) and a detail of the mesh around the eye with two landmarks visible (b).

the 3D surface alone, since some features such as the corners of the eyes and mouth are more easily identifiable by grey-level changes than changes to the surface of the mesh. With only nine landmarks, each example takes only a few seconds to annotate by hand.

A more general method is to place landmarks directly onto the 3D surface using ray-intersection from the mouse-pointer taking account of the camera model being used to render the image. This allows the user to put landmarks on surfaces without texture.

2.2 Forming the dense correspondence

The next step is to establish a dense correspondence between the surface meshes. This could be done using any set of landmarks as a frame of reference but it is desirable that the landmarks are typical of the distribution, so we have used the generalized Procrustes algorithm [9] to compute the mean landmarks. The key step in this process that is used repeatedly is a least-squares alignment of two sets of 3D landmarks for which we use the quaternion method described for example in [11].

Each surface is then warped onto the mean landmarks using the thin-plate spline (TPS) technique [4]. This brings the landmarks into exact alignment and interpolates a smooth transform for the other parts of the mesh, minimizing a bending energy. This is intended to ensure that while all the information implicit in the landmarks is taken advantage of, as little spurious variation as possible is introduced, especially in the vicinity of each landmark.

Having brought all the surfaces into close alignment, the dense correspondence is made by taking the closest point on each surface from each vertex in a base mesh. We used as a base mesh one of the examples in the training set. Experi-

ments showed that as long as the area of interest is covered by an adequate triangulation, the choice of which mesh to use is not critical.

The scans in the training set (including the base mesh) often included significant neck and ear areas that were not present in all the examples. We snip off these areas by using only those vertices where the *maximum* distance from the base mesh to each surface (after alignment) is less than 20mm (found through experiment). While this distance is, of course, application-specific, this technique is very effective in restricting the model to those regions that are well-represented by the training set surfaces. Essentially we are taking an *intersection* of the meshes; keeping only those areas that are well covered in *all* the examples.

Additionally, the scans used as a training set often have occasional small holes or triangulation errors, causing them not to be locally manifold. Because our method uses one mesh to sample the others, it is robust to such errors - any vertex displacements caused by these errors tend not to be correlated with the major shape changes and are thus not represented in the principal components.

When the dense correspondence with a base mesh has been made, the connectivity of the base mesh can be applied to the points of correspondence in each mesh to give a set of new meshes. We can then dispose of both the landmarks and the original meshes. We then apply the inverse of the aligning TPS warp to return each surface to its original location. (While TPS does not have an analytical inverse, an approximation can be computed using Newton’s method, see Acknowledgements section for an available implementation of inverse TPS.)

Figure 2 illustrates the whole process. In 2D the same effect could be accomplished by simply interpolating a fixed number of landmarks along each line segment but this does not extend to surfaces in 3D, unlike the TPS method we propose.

2.3 Building the detailed PDM

Now that we have constructed corresponding vertices in all the surfaces, we can treat them as landmarks. Then, following [7], we first apply the Procrustes algorithm to align all the shapes and produce a mean shape. Because our data is calibrated for size, we do not include scaling in the Procrustes alignment, but instead build a size-and-shape model [8]. Figure 3 shows the mean mesh that was computed for our dataset, a visual check that this mesh is smooth and free of artefacts such as vertex bunching tells us that our landmarks are sufficient and placed correctly. The mesh has been clipped of the poorly represented areas, leaving it with 1688 vertices.

The next step is to apply a principal components analysis (PCA) to the data. Each example can be represented with a

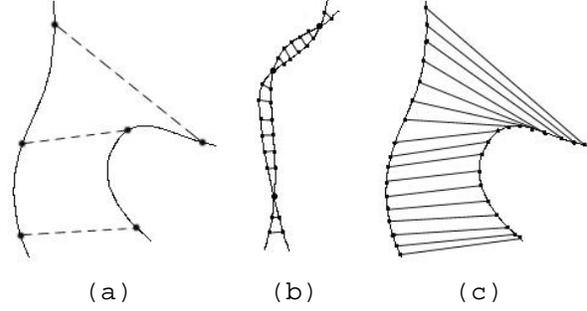


Figure 2. The TPS alignment step illustrated in 2D. The surfaces are landmarked (a), then TPS-warped onto a mean set of landmarks and a dense correspondence is established (b). Finally, the surfaces are unwarped back to their original locations taking the points of correspondence with them (c).

shape vector of concatenated x , y and z coordinates:

$$\mathbf{x}_i = [x_1, y_1, z_1, \dots, x_n, y_n, z_n]^T \quad (1)$$

The mean shape vector is given by:

$$\bar{\mathbf{x}} = \frac{1}{s} \sum_{i=1}^s \mathbf{x}_i \quad (2)$$

where $s = 193$ is the number of examples in the training set.

We compute the $3n \times s$ matrix \mathbf{D} using:

$$\mathbf{D} = [(\mathbf{x}_1 - \bar{\mathbf{x}}) | \dots | (\mathbf{x}_s - \bar{\mathbf{x}})] \quad (3)$$

The covariance matrix \mathbf{S} can then be computed using:

$$\mathbf{S} = \frac{1}{s-1} \mathbf{D} \mathbf{D}^T \quad (4)$$

\mathbf{S} is eigen-decomposed to give a set of eigenvectors ϕ_i and eigenvalues λ_i . However, because there are $n = 1688$ vertices in the mesh, each shape vector \mathbf{x}_i has $3n = 5064$ elements, so \mathbf{S} has 5064^2 elements. Fortunately, we can avoid having to eigen-decompose the very large matrix \mathbf{S} by instead computing the $s \times s$ matrix \mathbf{T} :

$$\mathbf{T} = \frac{1}{s-1} \mathbf{D}^T \mathbf{D} \quad (5)$$

which is somewhat more manageable at 193^2 elements. The first s eigenvalues of \mathbf{S} are the same as those of \mathbf{T} , the remainder are zero. We can compute the first s eigenvectors of \mathbf{S} from the eigenvectors \mathbf{e}_i of \mathbf{T} using:

$$\phi_i = \mathbf{D} \mathbf{e}_i \quad (6)$$

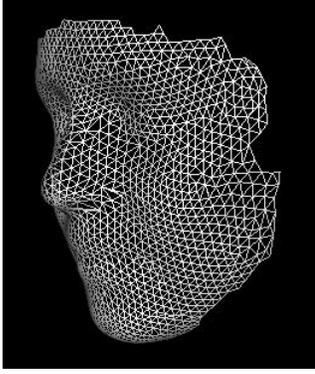


Figure 3. The averaged mesh, restricted to only those vertices that had good correspondences in the training set.

The computed eigenvectors can be treated as deformations of the whole mesh and can be directly added to the coordinates of the vertices of the mean mesh to synthesise new faces:

$$\mathbf{x}_{\text{new}} = \bar{\mathbf{x}} + \Phi \mathbf{b} \quad (7)$$

where $\Phi = [\phi_1 | \phi_2 | \dots | \phi_t]$ is the matrix of the first t eigenvectors and $\mathbf{b} = [b_1, b_2 \dots b_t]$ is a set of parameters controlling the modes of shape variation. The value of t is typically determined by the number of components that are required to account for 98% of the variation, ie. the lowest value of t such that:

$$\frac{\sum_{i=1}^t \lambda_i}{\sum_{i=1}^s \lambda_j} \geq 98\% \quad (8)$$

We model the distribution using a Gaussian, making the assumption that as long as each parameter is within three standard deviations then the face synthesised will be within the variation seen in the training set and should be a plausible human face. Scatter plots between pairs of modes such as that shown in Fig. 5 can be examined to check that the distribution contains no obvious non-linearities.

2.4 Fitting to unseen face scans

As with active shape models (ASMs) [7], the dense surface PDM can be used to interpret new examples. The search is initiated by placing the mean shape into the scene and iteratively deforming it using the parameters b_i to best match what is found locally around each landmark. With 2D and 3D images, a model of the target grey-level profile must be built in order to drive the search but with surfaces the task is much easier, needing only a simple nearest point

search. The fitting procedure is a hybrid of iterative closest point (ICP) [3] and active shape models.

The template mesh at each stage in the fitting process is specified by a rigid-body transform, \mathbf{T} , and the set of deformation parameters for the mesh, \mathbf{b} . (If we had removed scaling from the model then \mathbf{T} would have to be a similarity transform to allow for this.)

Each vertex in the deformable template mesh finds the nearest point on the target mesh, these points together yield a new shape, \mathbf{x}' . The rigid-body transform, \mathbf{T} , that minimizes the errors between the current mesh and \mathbf{x}' is computed (following [11]) and applied.

The deformation parameters that best model the shape \mathbf{x}' are given by:

$$\mathbf{b} = \Phi^T (\mathbf{x}' - \bar{\mathbf{x}}) \quad (9)$$

The parameters are clamped to the three standard deviation limit given by the model, $-3\sqrt{\lambda_i} \leq b_i \leq 3\sqrt{\lambda_i}$, and then (7) is applied to synthesise the template. This deformed template mesh, together with its transform, forms the input to the next iteration.

As with ASMs, the fitting process works through local improvement and thus will fail to converge if the search is initiated too far from the target. In an automated system, some method for initializing the search in the vicinity of the target would be required.

An important issue with the search procedure is the synchronization of ICP iterations and ASM iterations. Quite often, many ICP iterations are required for the template to converge with the target - it would be inappropriate for the ASM to deform the mesh before at least an initial convergence was attained. On the other hand, after deformations have been introduced some adjustment of position and orientation may still be necessary so it would be equally inappropriate to run ICP to convergence before ASM fitting. The solution we have adopted is to run our fit to a 'schedule' of 100 iterations - in our experience this is sufficient for the ICP to converge (as long as the starting position isn't so far from the target that the search does not converge at all). The ASM modes are introduced successively over the course of the schedule, from zero modes to the full range. Thus, the first few iterations (with zero ASM modes) are effectively a rigid fit, correcting for gross positional and orientation errors. Then the major modes of deformation are introduced to correct for large-scale shape differences, the finer-scale deformations only being introduced towards the end of the fit when the mesh is expected to be close to convergence.

At each iteration, the number of modes to use is given deterministically by:

$$\lfloor t * \frac{\text{CurrentIteration}}{\text{TotalIterations}} \rfloor \quad (10)$$

Such a scheme of introducing the modes gradually was used successfully in previous work [12].

3. Results

3.1 Examining the model

Figure 4 shows the first three modes of our detailed PDM based on 193 scans. Computation of this model takes roughly five minutes with a 400MHz PentiumIII. Synthesis of a new example takes under a tenth of a second. The first mode dominates (at 73.2%) because we left scaling in the model in order to capture the correlation between face size and shape. Doing this gives us a model that is more *specific* - at the fitting stage we don't allow scaling as long as the target surface is calibrated for scale.

The first mode shows a strong correlation with the age of the subject while modes 2 and 3 show changes in face shape that involve identity and facial expression. Mode 3 is quite strongly correlated with ethnic origin, with Asian faces being found on the left of this mode and Caucasian faces more on the right.

In this model, $t = 33$ modes are required to account for 98% of the total variation. Revealingly, an earlier model built with 72 examples required 25 modes to account for 98% of the variation. This indicates that the model is capturing natural face shape variation quite effectively since only 8 extra modes were required to model the extra 121 faces.

A visual examination of these modes forms part of our evaluation of the efficacy of our method and the correctness of the correspondence. All of the images are plausible human faces and don't appear to have any artefacts such as ridges or distortion.

Figure 5 shows a scatter plot of the first two modes. The distribution shows no obvious correlations between the two components. Plots of other pairs of modes show a roughly similar result; it is hard to tell whether any perceived deviations from a Gaussian pattern are due to sampling or are genuine patterns in the data.

3.2 Fitting to new images

Figure 6 shows the fitting in action on an image that is freely available on the web and that is not part of the training set. The target surface is shown in the first image with its original texture and then as a semi-transparent surface. The deformable surface that is being fitted to it is shown with a texture in order to more clearly illustrate the convergence. Where the image looks mottled, the two surfaces are in close alignment. Figure 6 shows that the model is capable of converging to an unseen face scan which is not limited to the area modelled but extends beyond it. The fit takes under two minutes.

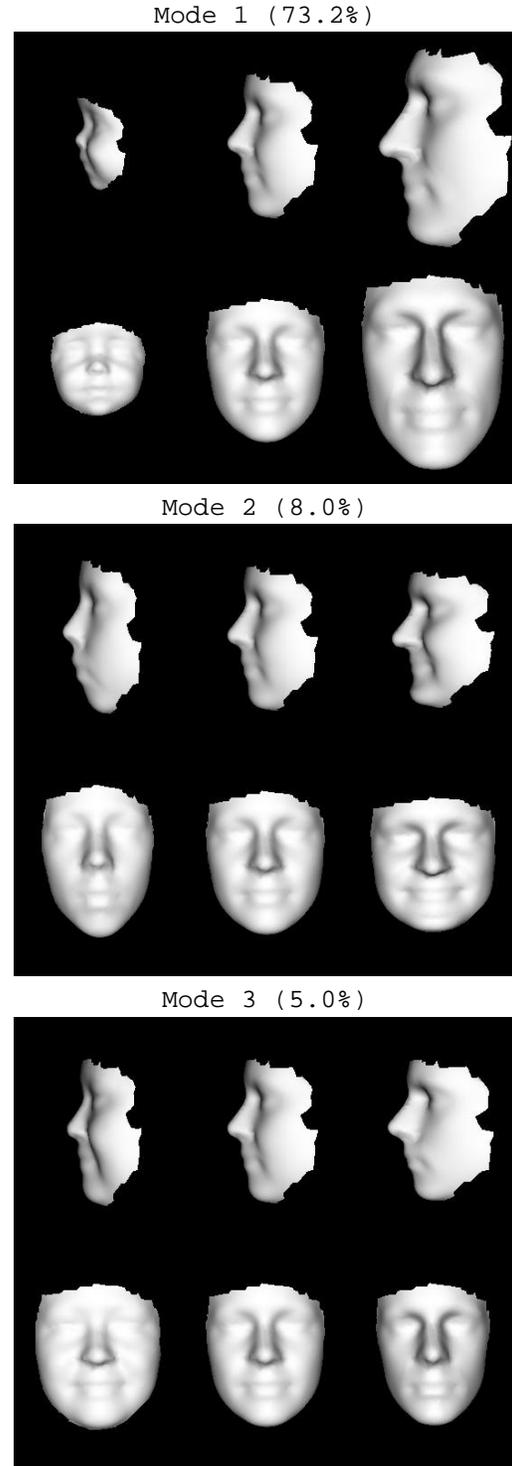


Figure 4. The first three modes, between -3, 0 and +3 standard deviations (columns), side and front views (rows)

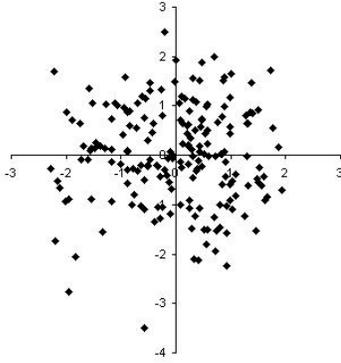


Figure 5. A scatter plot of mode 1 (horizontal) against mode 2 (vertical), showing a reasonably Gaussian distribution. Units are standard deviations.

3.3 Classification of surfaces

The nature of the dense surface model is such that subtle relationships in the distribution of the vertices can be used for delineation of group characteristics. One such relationship concerns the difference between male and female faces. The space into which the faces have been transformed can be used directly for such comparisons since its *metric* is the similarity in shape between examples (ie. faces with similar shape will be close together).

One simple way to maximise the difference between male and female faces is to construct a line between the average male face in the space and the average female. Figure 7 shows a plot of the examples using this line as the horizontal axis, with the average male located at -1 and the average female at $+1$. The examples show considerable separation, with only 23 out of 193 examples being misclassified on this (simple) basis. Figure 8 shows examples from along this line, corresponding to the marked locations on Fig. 7. Statistical learning techniques such as support vector machines [16] could be used on this distribution in the t -dimensional *face-shape-space* to infer classifications in a more sophisticated way.

Crucially, if we do the same analysis using just the original 9 landmarks, there is virtually no separation of the two groups - the inclusion of the extra surface points has improved the sensitivity of the model. This feature of the technique suggests that it should be of use for delineating characteristics that are known to be shared between subgroups, such as male-female differences as we have shown. An important application of this would be for comparing the surface shape of anatomical components between individuals with and without a congenital abnormality.

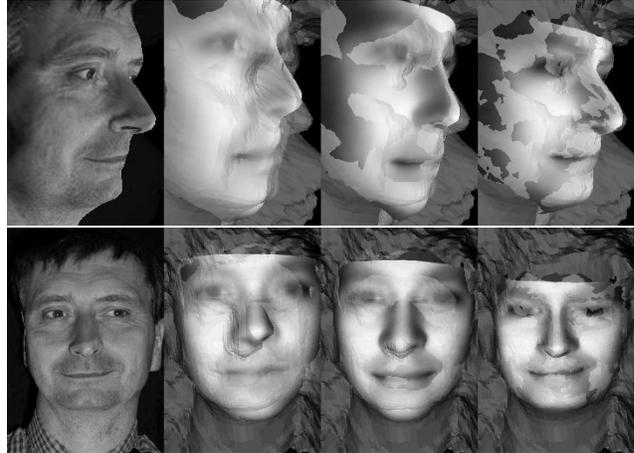


Figure 6. The fitting in action on an unseen example, showing the target face scan, the initial placement of the template, an intermediate stage in the process and the final fit (side and front views).

To illustrate the potential clinical use of our technique we will briefly show some preliminary results from a study of Noonan Syndrome [1]. A collection of 3D scans of 62 children with ages ranging from 4 months to 16 years were landmarked with the same 9 points. 22 of the children had Noonan Syndrome, the other 40 did not. The dataset is approximately balanced in terms of gender and age between the two classes. Figure 9 (a) shows the scatter of points when plotted on the Noonans positive-negative line as a horizontal axis. Figure 9 (b) shows the same plot but taken from the dense surface point distribution model distribution (Fig. 10 shows examples from this plot).

In Fig. 9 (a), 9 of the 62 examples are misclassified, while in (b) only one example is misclassified, showing again that the dense surface point distribution model is more sensitive to correlated features than the landmark model. In Fig. 10 there is good visual agreement between the features identified by the model as being correlated with Noonan Syndrome and those seen in children with the syndrome. While there is insufficient data to demonstrate a clear clinical benefit, we are encouraged that the use of dense surface point distribution models will help delineate and potentially diagnose facial dysmorphologies such as Noonan Syndrome.

4. Summary and Conclusions

We have described a method for creating detailed deformable models from a set of surface meshes. The method described will work on surfaces of any topology, including

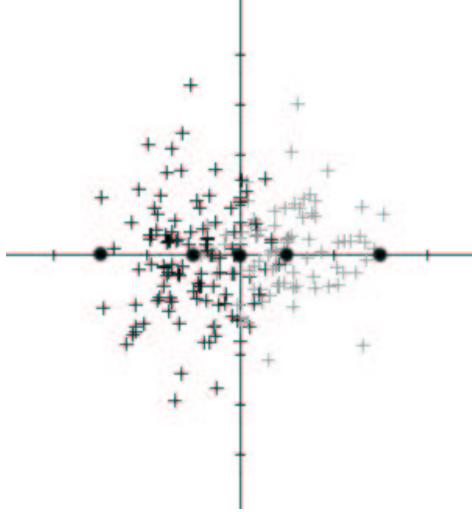


Figure 7. A scatter plot of the male-female mode (see text) against mode 2, with the darker crosses showing the males and the lighter crosses the females. The circular blobs correspond to the five faces in Fig. 8, from left to right.

those with holes and unconnected regions, as long as the topology is consistent across the training set. If the topology is not consistent then the surface used as the base mesh will be used to sample the others and may give incorrect results. The method will work with surfaces of regular or irregular triangulation, as well as with surfaces of varying densities of triangulation - both for building the model and for fitting. Indeed, non-polygonal surface representations could also be used, as long as they are amenable to TPS warping and querying for closest point on the surface. If required, a triangulation could be created to use as a base mesh.

The model-building method will work for shapes displaying considerable deformation, including articulation, as long as enough landmarks are used to bring the surfaces into sufficient correspondence that the nearest-point sampling from the base mesh remains valid. If the shape changes are such that modelling the distribution with a Gaussian becomes invalid a non-linear approach to modelling the shape space such as described in [6, 14] might become appropriate. If the deformation in an example is very large then the fitting algorithm will not always converge correctly, this is true of both ICP and ASM separately, as well as our hybrid fitting method.

The method is robust to occasional errors in the training set data such as noisy or missing vertices. Such errors are not likely to be correlated with the other shape changes and

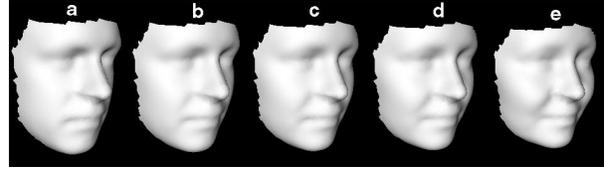


Figure 8. Five faces from along the male-female line (see text), corresponding to the circular blobs in Fig. 7. Faces *b* and *d* are the average male and female faces, while faces *a* and *e* are exaggerations, to show more clearly the changes being captured.

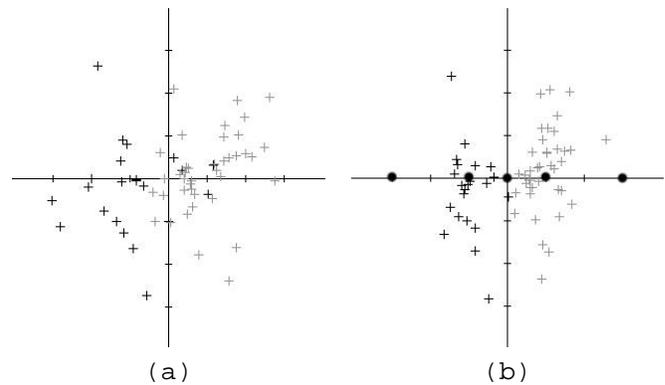


Figure 9. Scatter plot with the Noonans positive-negative line horizontal and mode 1 vertical, using the original 9 landmarks (a) and the same plot using the dense surface point distribution model, showing greater separation (b).

should therefore only appear in the less significant modes.

We have shown how the deformable model can be used to automatically interpret new scans. This scheme could be used to implement a boot-strapping procedure for building the model. One issue to be addressed is the use of a fixed schedule to synchronize the ICP and ASM iterations. The intention is that the amount of deformation allowed in the ASM model corresponds to the degree of convergence of the ICP fit, to prevent the template deforming to fit the surface too closely before a good overall correspondence has been established. A better way to achieve this might be some form of adaptive scheduling to introduce modes based on some measure of convergence of the ICP algorithm.

We have shown that the extension of a landmark model of the face to a dense surface point distribution model improves sensitivity to subtle correlated features such as the

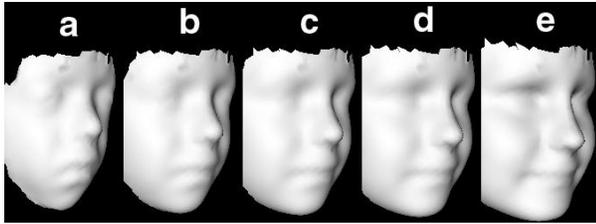


Figure 10. Five faces from along the Noonans positive-negative line (see text), corresponding to the circular blobs in Fig. 9 (b). Faces b and d are the average Noonans and non-Noonans faces, while faces a and e are exaggerations, to show more clearly the changes being captured.

shape differences between male and female faces. It is hoped that the technique can be used to gain a quantitative understanding of the changes across the entire face seen in syndromic facial dysmorphologies.

Acknowledgements

The TCTi face scanner was acquired with funding from the Birth Defects Foundation.

Many of the techniques used in this work are implemented in the Visualization Toolkit (VTK) which is an open-source C++ library available from:

<http://public.kitware.com/>

In particular, the inverse TPS warping was contributed by David G. Gobbi.

Our thanks are also due to the many volunteers who have allowed their scans to be used in our work.

References

- [1] J. Allanson, J. Hall, H. Hughes, M. Preus, and D. Witt. Noonan syndrome: An evolving phenotype. *Am. J. Genet.*, 21:507–514, 1985.
- [2] A. Benayoun, N. Ayache, and I. Cohen. Adaptive meshes and nonrigid motion computation. In *International Conference on Pattern Recognition*, pages 730–732, 1994.
- [3] P. Besl and N. McKay. A method for registration of 3-D shapes. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 14:239–256, 1992.
- [4] F. Bookstein. Shape and the information in medical images: A decade of the morphometric synthesis. *Computer Vision and Image Understanding*, 66:99–118, 1997.
- [5] A. Brett and C. Taylor. A method of automated landmark generation for automated 3D PDM construction. In *British Machine Vision Conference*, pages 914–923. BMVA, 1998.
- [6] T. Cootes and C. Taylor. A mixture model for representing shape variation. In A. Clark, editor, *British Machine Vision Conference*, pages 110–119, 1997.
- [7] T. Cootes, C. Taylor, D. Cooper, and J. Graham. Active shape models - their training and application. *Computer Vision and Image Understanding*, 61(1):38–59, 1995.
- [8] I. Dryden and K. Mardi. *Statistical Shape Analysis*. Wiley, 1998.
- [9] J. Gower. Generalized procrustes analysis. *Psychometrika*, 40:33–51, 1975.
- [10] E. Guest, E. Berry, and D. Morris. Using the CSM correspondence calculation algorithm to quantify differences between surfaces. In *British Machine Vision Conference*. BMVA, 2000.
- [11] B. Horn. Closed-form solution of absolute orientation using unit quaternions. *Journal of the Optical Society of America A*, 4:629–642, 1987.
- [12] T. Hutton, S. Cunningham, and P. Hammond. An evaluation of active shape models for the automated landmarking of cephalograms. *European Journal of Orthodontics*, 22:499–508, 2000.
- [13] C. Lorenz and N. Krahnstöver. Generation of point-based 3d statistical shape models for anatomical objects. *Computer Vision and Image Understanding*, 77:175–191, 2000.
- [14] B. Scholkopf, A. Smola, and K.-R. Muller. Nonlinear component analysis as a kernel eigenvalue problem. Technical report, Max-Planck-Institut für Biologische Kybernetik Arbeitsgruppe Bulthoff, Tübingen, Germany, December 1996. TR 44.
- [15] L. VanGool, P. Kempenaers, and A. Oosterlinck. Recognition and semi-differential invariants. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 454–460, 1991.
- [16] V. Vapnik. *The nature of statistical learning theory*. Springer, New York, 1995.
- [17] S. Yamany and A. Farag. Free-form surface registration using surface signatures. In *International Conference on Computer Vision*, 1998.