

## USING DETAILED INDEPENDENT 3D SUB-MODELS TO IMPROVE FACIAL FEATURE LOCALISATION AND POSE ESTIMATION

ANGELA CAUNCE\*, CHRIS TAYLOR† and TIM COOTES‡

*Centre for Imaging Science, The University of Manchester  
Manchester, M13 9PT, UK*

*\*angela.caunce@manchester.ac.uk*

*†chris.taylor@manchester.ac.uk*

*‡timothy.f.cootes@manchester.ac.uk*

Received 28 January 2013

Accepted 31 July 2013

Published 20 December 2013

We show that the results from searching 2D images or a video sequence, with a 3D head model can be improved by using detailed sub-models. These parts are initialised with the full model result and are allowed to search independently of that model, and of each other, using the same algorithm. The final results for the sub-models can be reported exactly, or optionally fed back into the full model to be constrained by its parameter space. In the case of a video sequence this can then be used in the initialisation of the next frame. We tested various data sets, constrained and unconstrained, including a variety of lighting conditions, poses, and expressions. Our investigation showed that using the sub-models improved on the original full model result on all but one of the data sets.

*Keywords:* 3D statistical shape models; facial feature tracking; detailed sub-parts; driver monitoring.

### 1. Introduction

Head and facial feature tracking can provide important information in various environments with respect to the activity and attitude of the subject. For example, in a driving scenario, head orientation alone can indicate attentiveness. Feature localisation and subsequent behaviour analysis can give a detailed picture of the driver's state and possible intent. This could identify critical or dangerous situations. Recently, 2D models have been used with great success to localise and analyse features of the face,<sup>1-4</sup> however in some unconstrained scenarios with large pose variation this approach may be limited. Multiple 2D models and detectors may be required for different views.<sup>5,6</sup> There is also the additional problem of view-based occlusion. As a consequence, authors have been experimenting with augmented 2D<sup>7,8</sup> and 3D.<sup>9</sup> Due to pose invariance, a 3D model requires less training data than its 2D counterpart, and is able to report critical pose information directly without additional calculation. In Ref. 10 we showed that our 3D method outperformed an established 2D approach on out of plane rotations and in Ref. 11 we extended this by integrating some limited facial actions for preliminary behaviour

analysis. To do this we used two sparse, largely symmetrical, statistical point models of the whole face. However, the complex interplay of the various parts of the face may not have been fully realised due to global constraints. Also, the ability to deal with some individual quirks, like a crooked smile, is limited by the training set. This is a standard problem with deformable objects containing complex sub-parts. One solution is to define the sub-parts separately and model their inter-relationship. The search is performed by locating the sub-parts and confirming an acceptable configuration. Established methods include: pictorial structures<sup>12</sup>; star models<sup>13</sup>; Hierarchical Deformable Templates<sup>14</sup>; and probabilistic approaches.<sup>15,16</sup> For example, Martinez<sup>17</sup> used a probabilistic approach and weighted abstract sub-parts of the face based on their involvement in the test expression. However working in 3D has advantages in that many configurations of parts (those derived from the object’s pose) are already constrained by the model structure and only the articulation problem remains. Generally, authors working in 3D have taken a bottom up approach, that is finding sub-parts and combining them in some meaningful way. In Ref. 18 Blanz and Vetter show that this produces visually pleasing results. Tena *et al.*<sup>19</sup> illustrated that sub-parts improved performance when reconstructing motion capture data.

### 1.1. Contribution

Our method takes a top down approach, in that we localise the face with a full model using the method in Ref. 11, and then refine the result by allowing the parts to search independently afterwards. By extending the method in this way, we have increased its accuracy by improving its versatility. The system can deal with new feature configurations, without requiring additional training data. This may even reduce the amount of data required for training in future. We demonstrate the improvements by reporting results on 7 large datasets with various challenges (see Figures 6, 10 and 11). These show improvements either overall, or in some localised region of the face. The method has all the advantages of 3D outlined above, and can deal with greater pose variation than demonstrated in other works. We report median average point-to-point errors of less than 15 pixels (inter-ocular distance ~ 100 pixels) on headings up to 70 degrees and on all pitch angles tested (up to 60 degrees).

## 2. Search Method

The search method uses two sparse 3D statistical models<sup>20</sup> of the form:

$$\mathbf{x} = \bar{\mathbf{x}} + \mathbf{P}\mathbf{b} \tag{1}$$

Where each example  $\mathbf{x}$  is represented by a vector of  $n$  3D co-ordinates  $(x_1, y_1, z_1, \dots, x_n, y_n, z_n)$ . Each is expressed in (1) as the mean vector,  $\bar{\mathbf{x}}$ , plus a linear combination of the principal components,  $\mathbf{P}$ , with coefficients,  $\mathbf{b}$ .

One of the two models is built from “identity” training data, i.e. data from 923 individuals, with a close to neutral expression, eyes open, and mouth closed. The other is built from a small set of facial actions created from a neutral base. These are described in

more detail in the next sections. Each 3D point in a model searches for a target independently of the others and then the model is simultaneously fitted to all points subject to its shape constraints. This is described in more detail in Sections 2.4 and 2.5.

## 2.1. Identity model

The identity model was built from 923 head meshes, each having 4923 vertices. Each mesh was created from a manual markup of photographs of an individual. The front and profile shots of each person were marked in detail and the two point sets were combined to produce a 3D representation for that subject (Figure 1, top left).

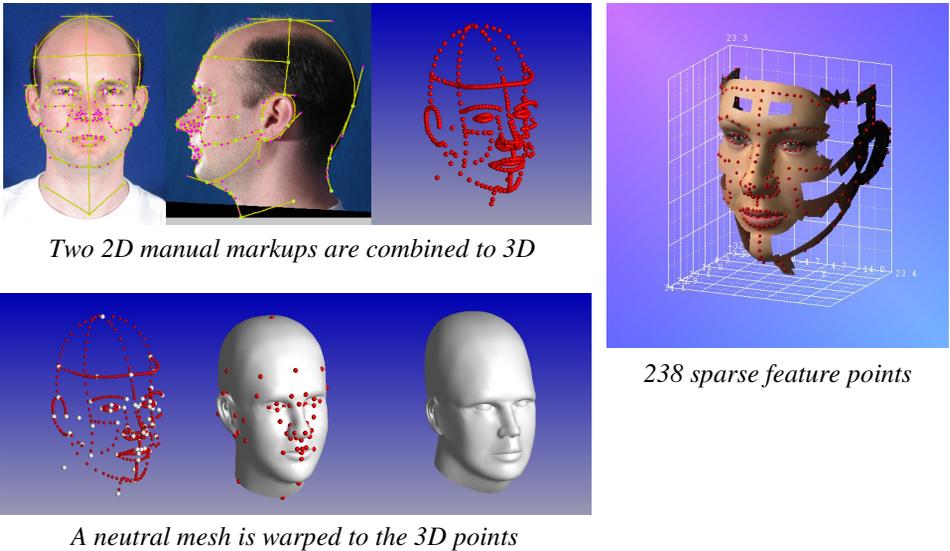


Fig. 1. The front and profile markups are combined to create a 3D point set (*top left*). Using known correspondences between the markup and a generic head mesh, an individual mesh can be created for each subject (*bottom*). Only a subset of the mesh vertices are used to build the statistical shape model (*right*).

A generic mesh, with known correspondence to the 3D points, was warped<sup>21</sup> to fit the markup giving a mesh for each individual (Figure 1, bottom). Since the same mesh was used in each case the vertices are automatically corresponded across the set. Therefore, any subset of vertices from this mesh can be used to build a sparse 3D shape model of the type described above.

Although it is possible to build and use a shape model incorporating all the vertices of the mesh it proved cumbersome in practice. Searching with 4923 vertices was time consuming and computationally expensive, and there is a great deal of redundancy in much of the face, e.g. forehead, cheeks. We therefore chose to use 238 points (Figure 1) which are close to features of interest such as eyes, nose, mouth, etc.

## 2.2. Facial actions model

It is difficult to obtain corresponded 3D head data for a spectrum of emotions, so we chose to build a more versatile model of basic facial actions. Only 8 examples were used to build the model: Mouth open; Eyes closed; Smile; Mouth turned down; Brow raised; Brow lowered; Grimace; and neutral (Figure 2). These were all created by modifying the same neutral head mesh used to generate the examples in the identity model (Figure 1). This meant that they were automatically corresponded with those examples as well as each other. Using a single example of each facial action assumes that they are transferrable between individuals. After construction utilising 99.5% of the variance in the data, the model had 5 modes of variation which, by observation, had the primary functions: mouth open; brow raise/lower; smile; grimace; and eyes close.

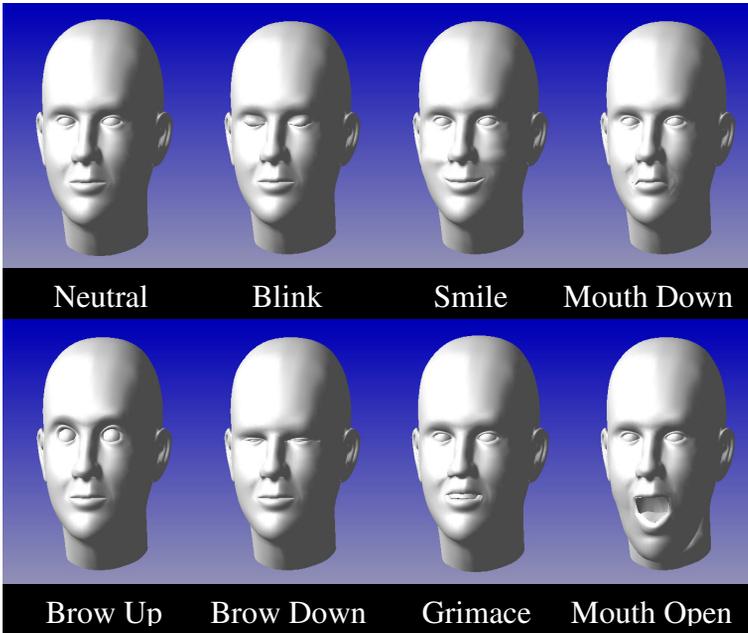


Fig. 2. The neutral mesh (*top left*) was modified to create a set of simple basic facial actions.

## 2.3. Combining the models

Unlike other approaches, which use a combined model strategy,<sup>22,23</sup> these two models are used in an alternating process to localise the features of the face and to provide a basic representation of some simple behaviours.<sup>11</sup> This is done by substituting the results from the ID model into the actions model, and vice-versa, when matching to the target points. Therefore, at each iteration of the algorithm, both models are fitted in sequence to the same target before moving on to the next iteration. This is represented in the following equations:

$$\mathbf{x}^{k(1)} = \bar{\mathbf{x}}^{k(1)} + \mathbf{P}_{\text{ID}} \mathbf{b}_{\text{ID}}^k \quad (2)$$

$$\bar{\mathbf{x}}^{k(1)} = \bar{\mathbf{x}}_{\text{ID}} + \mathbf{P}_{\text{A}} \mathbf{b}_{\text{A}}^{k-1} \quad (3)$$

$$\mathbf{x}^{k(2)} = \bar{\mathbf{x}}^{k(2)} + \mathbf{P}_{\text{A}} \mathbf{b}_{\text{A}}^k \quad (4)$$

$$\bar{\mathbf{x}}^{k(2)} = \bar{\mathbf{x}}_{\text{ID}} + \mathbf{P}_{\text{ID}} \mathbf{b}_{\text{ID}}^k \quad (5)$$

Where  $k(1)$  and  $k(2)$  refer to the 1st and 2nd fit at each iteration  $k$ , (5) is the current identity result and is used as the action model mean, (3) is devised from the current action result and is used as the identity model mean, and  $\mathbf{b}_{\text{A}}^0$  is the zero vector. Notice that the action model mean,  $\bar{\mathbf{x}}_{\text{A}}$ , is not used since this has no meaning in this context.

#### 2.4. Obtaining target points

Most of the target points are located using an independent local template matching at each model point. There are 238 points in the model and each can search with a small ( $5 \times 5$ ) view based texture patch using normalised correlation to find the best location in a local neighbourhood ( $9 \times 9$  increasing on each side by 2 at each resolution). This has the advantage of providing some robustness to illumination variation over the face and between images. Each patch is extracted from a mean texture generated from 913 subjects. The individual examples are in the form of faces “unwrapped” from the meshes described in Section 2.1. Because all the vertices of the meshes have the same (UV) coordinates into the texture, all the unwrapped examples correspond directly pixel for pixel and it is easy to obtain the mean (Figure 3).

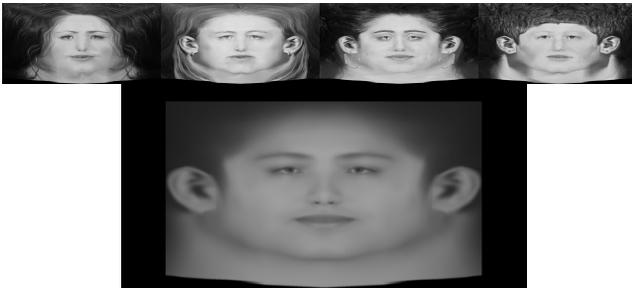


Fig. 3. The texture patches are sampled from a mean texture (*bottom*) averaged over a set of faces “unwrapped” from the head meshes. Some example faces are shown (*top*).

Variation in the texture was not modelled for these experiments. Each patch is always the same size and shape throughout the matching process but changes content at every iteration. It is updated based on the surface normal of the point and the current orientation of the model, and represents the view of the texture at that point (Figure 4). It is assumed for this purpose that the head is a globe and the texture lies tangential to the surface at each point with its major (UV) axes aligned to the lines of latitude and longitude. This is

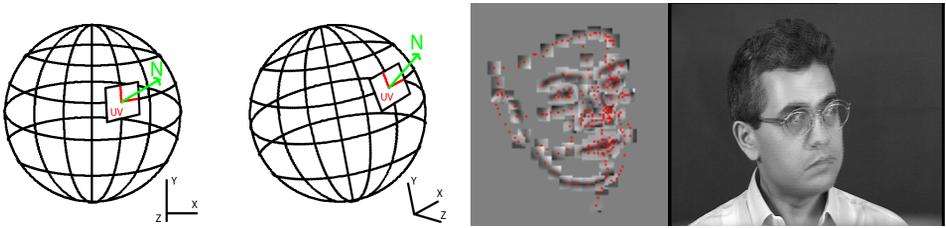


Fig. 4. The local texture patches are generated based on the current pose of the model. The UV axes of the texture are assumed to follow the lines of latitude and longitude of the head at  $0^\circ$  rotation (*left*). During matching the patches are a “window” onto the texture oriented by the estimated pose of the head. Only forward facing points, determined from the surface normal are used to search this way.

different to approaches such as that of Gu & Kanade where reference patches are stored in only a subset of views.<sup>24</sup>

However, searching using only the patches does not take advantage of important image information from the occluding boundaries. Therefore, those points at approximately 90 degrees to the viewpoint do not use the template matching technique but search along the surface normal for the strongest edge. These points are easily identified because the surface normal of the mesh and the current pose are known. The length of the profile is increased as the search resolution is increased. Since the target point is on the strongest edge, the “match value” is the normalised edge strength. Using this information can keep the model from rotating instead of changing shape, and helps to maintain scale, preventing mismatches such as the mouth to the nose.<sup>25</sup>

The search is conducted at multiple resolutions, starting at the lowest, and is completed at each resolution before moving onto the next. The number of iterations at each level is set to a minimum of 10 but can continue, if there are still large matching differences between the model patches and the image, to a maximum of 100.

## 2.5. Model fitting

Once the target points are established the whole model is fitted using the active shape model fitting method in Ref. 26 extended to 3D, assuming an orthogonal projection. An equation of the form in (6) is minimised to find the new model point positions.

$$\min_{\mathbf{b}, \mathbf{t}} \left\{ \left| \mathbf{W} \left( \mathbf{T}_t (\bar{\mathbf{x}} + \mathbf{P}\mathbf{b}) - \mathbf{x}_{\text{obs}} \right) \right|^2 \right\} \quad (6)$$

$\mathbf{T}$  represents the camera transform (in this case orthogonal) and pose with parameters  $\mathbf{t}$ ,  $\mathbf{W}$  is the diagonal weighting matrix, and  $\mathbf{x}_{\text{obs}}$  is the observed, or target, set of points. The weights in  $\mathbf{W}$  are derived from the normalized distances between the current points and the target points of the model, which was found to be a useful weighting scheme.<sup>25</sup> The model fitting is a two-stage process. Firstly, the points are rigidly aligned to minimise the sum of squared distances between matched points, then the shape model parameters,  $\mathbf{b}$ , are updated using a least squares approximation. The global fit has the advantage of minimising the effect of badly matched or occluded points.

Initialisation is achieved using the Viola-Jones (V-J) face detector.<sup>27</sup> For the video sequences this occurs once at the start and again only when the search fails during the sequence. Failure is determined by comparing the average match value of the template patches to some threshold (0.2). If the match does not fail, the search on the next frame is initialised using the latest result. For still images the V-J detector is used on each example.

### 3. Sub-Models

Fitting the model to the target using global constraints, as outlined above, has the advantage of keeping all of the areas of the model in their expected place, as well as having a neutralising effect on rogue matches, such as those found at occluded points. However, once this process is completed, it may be expected that the final result will also suffer, since un-correlated movements of the individual sub-parts may be lost as noise. Plus, individual quirks, such as a crooked smile, cannot be localised unless they are specifically included in the training set.

To provide the added flexibility necessary we propose extending the method by allowing sub-models to continue searching independently, after the full model search has completed. In addition, because the sub models will cover a much smaller area of the face we can afford to use all the vertices in the selected area rather than the sparse subset used for the full face, which was chosen to reduce computation time and redundancy. Since large areas of the face have little or no movement, we concentrated on the areas around the eyes and mouth. These are the most expressive, mobile, and informative areas when considering a face tracking and analysis system. We thus built two sub-models from the areas shown in Figure 5. As with the full face models, an identity and an actions model were constructed for each sub-part.

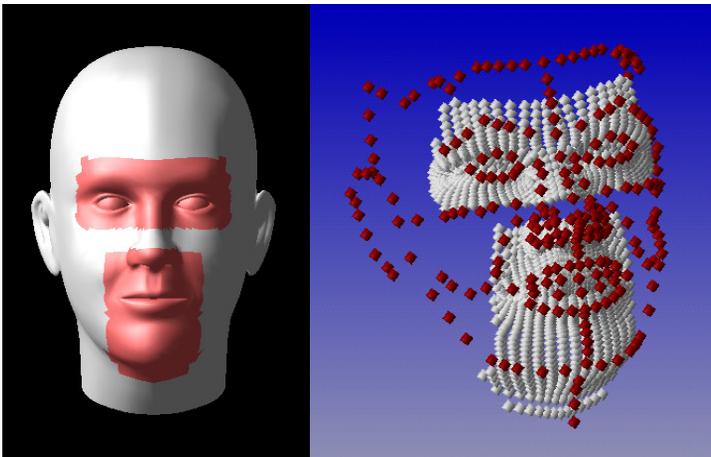


Fig. 5. The two sub-models are built from all the vertices in the eye and mouth areas (*left*). The relationship between the sparse full face model and the sub-models is shown *right*.

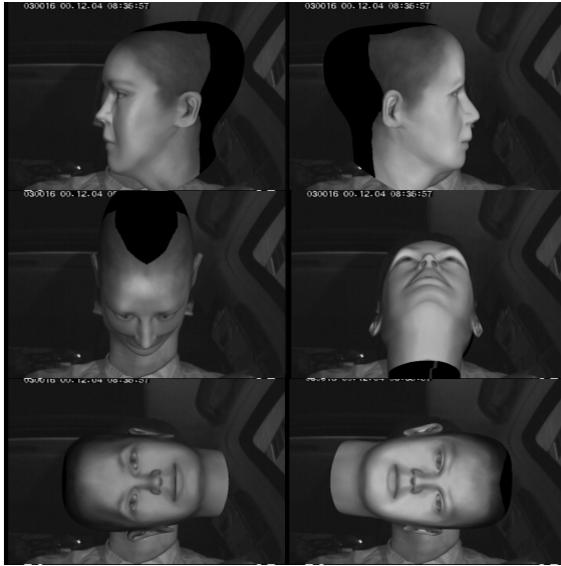


Fig. 6. The Artificial Driver data set contains extreme poses.

#### 4. Sub-Model Search

The full model search is conducted over a series of three increasing image/texture resolutions. The search is completed at each resolution before continuing to the next. For the sub-models we use exactly the same search method except that searching always begins at the penultimate resolution. Experiments indicated that results were improved if two resolutions were used but there appeared to be no advantage to using all three.

When the full model search is complete the sub-models are initialised in their respective positions (Figure 7). Since the sub-models have many more points than the full model an initialisation target is constructed for each sub-part from just those points in the full face model which are common to both. To do this the common points have a weighting of 1 whereas all other points have weight 0. The sub-models are then fitted to the sparse model result using same method used in the image search (section 2.5), i.e. the 3D extension of that in Ref. 26. The identity and action sub-models are fitted to this target alternately using equations (2)–(5). After this initialisation, each sub-part is allowed to search the image independently of the other and of the full model (Figure 7). Algorithm 1 outlines the search process.

Currently the original full model search method<sup>11</sup> has been optimized to run in real-time on multiple processors (up to 24 fps on a quad core PC at ~40% CPU usage), however the sub-model search has not yet been integrated into this online system and the offline search for each frame can take just over 3 seconds on a single processor. Since we used code optimization and multiple processors to improve the original offline search from a time of ~0.4 seconds up to real-time, with CPUs under capacity, we are confident that the time for the parts search can also be reduced to allow real-time performance in the future.



Fig. 7. The sub-models are initialized from the full model search (*left*). Each model searches independently of the other and of the full model (*Middle and Right*).

---

### Algorithm 1

---

1. Use a face detector or the result from the previous frame to initialize the model
  2. For resolutions Low/Medium/High
    - Do SEARCH
      - Obtain patches for current view
      - Search for point targets using patch matching
      - Fit ID model
      - Fit actions model
    - Until fail or end condition
    - If fail then stop
  3. For each part model
    - Initialise the part using the sparse full model result
    - For resolutions Medium/High
      - Do SEARCH (as above) with part model only
  4. Optionally feed back the parts results into the full sparse model by model fitting.
- 

#### 4.1. Feedback and reporting strategies

Once the results from the sub-models are obtained there are several ways that they can be integrated with the full model results for analysis. The most obvious is to use the parts to replace all the points in the full model that are common to the sub-models. We refer to this as “Exact Parts”. The alternative is to feed back (FB) the common points into the full model by fitting it to the Exact Parts result. This fitting follows the same iterative alternating process described above for the sub-model initialisations. Normally, when analysing a video stream, the result of each frame is used to initialise the search on the next. This means that if the sub-model results are fed back into the full model these will influence subsequent frames. For a series of still images this will have no effect.

Therefore for still images we report these methods: “NoParts” (results from the full model search); “NFBEExact” (exact parts substituted); and “FB” (feedback — i.e. full model fitted to parts results). For video data, method “FB” will affect future results and there is an additional combination: “FBExact” (feedback affects future results but exact parts reported).

## 5. Data Sets

For our analysis we used 7 different data sets presenting various challenges (Figures 6, 10 and 11):

- XM2VTS<sup>28</sup>: 2344 still images of 295 individuals at  $720 \times 576$  pixels. The subjects are posed against a fairly uniform backdrop and in general have neutral expressions and near frontal poses.
- BioID<sup>29</sup>: 1520 still images of 23 individuals at  $384 \times 286$  pixels. The data was acquired in an office setting so has cluttered backgrounds. It shows much more natural poses (although still mainly near frontal) and a variety of expressions.
- Expressions: 401 still images of 103 individuals at  $1024 \times 768$  pixels. Each person is making some or all expressions selected from: eyes closed; smile; frown; and surprise. These images are near-frontal and posed against a uniform background so the main challenge is from the facial actions.
- 3 video sequences of 3 different drivers. These were taken in real driving conditions. Each sequence is 2000 frames but only a subset of each was used for evaluation: 150, 156, and 136 frames. There is a great deal of lighting variation, within and between frames, and there is a wider variation in poses than in the still images.
- Artificial driver short sequences.<sup>10</sup> This publicly available dataset was devised to assess the ability of search methods to deal with large poses. It comprises of a series of images of 20 synthetic subjects in known poses. The subjects were arranged against a real in-car background (Figure 6). Each sequence starts at zero rotations and runs in a single direction. The sequences are as follows: Heading  $\pm 90^\circ$ ; Pitch  $\pm 60^\circ$ ; and Roll  $\pm 90^\circ$ .

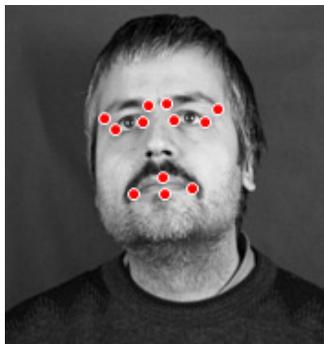


Fig. 8. The points used in the performance assessment. 12 points around the eyebrows, eyes, and mouth.

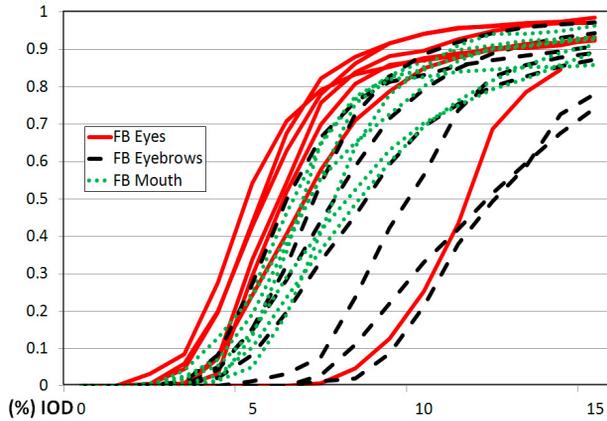


Fig. 9. (Color online) The performance of the Feedback method on different areas of the face over all datasets (cumulative point to point errors).



Fig. 10. The method was presented with a variety of challenges including occlusion; eyes closed; mouth open; strong lighting; low contrast; pose; glasses; facial hair; and expression. Black points are the full model result and white points show the integrated results from the sub-models (common points only). The general method performs well and in these examples only subtle improvements were necessary, but see Figure 11 for more extreme corrections.

## 6. Results

The available annotations for the different data sets do not mark the same features. Due to this, and the differences between the annotations and the model points, only a subset of 12 points was used for point to point error evaluation (Figure 8). The points chosen are located on the better defined features, common to all sets: the ends of the eyebrows; the corners of the eyes; the corners of the mouth; and the top and bottom of the mouth.

On the still image databases the search is initialised using the V-J detector on every image. Since the method relies on initialisation in the area of the face, we assessed the detector's performance by comparing the 12 points of interest to the box it returned. If any manually marked points fell outside the box the detection was considered a failure. It

was found that the detector failed on 8% of the BioID data set. These examples were therefore excluded from the analysis.

### 6.1. Feature location

The generic model can handle a variety of challenges including: occlusion; pose variation; and variable lighting.<sup>11</sup> This includes real-world driving situations where the lighting is changing rapidly between frames and where there are areas of strong shadow on the face (Figure 10, top right). It is able to do this because the points locate their targets independently and use normalised texture patches.

In many cases the sub-models will make only a small correction (Figure 10). However, Figure 11 illustrates the ability of the sub-parts to correct poorer results



Fig. 11. Examples showing how the sub-models can correct a poor full model result. Top row shows real-world driver video. The second image illustrates the independence of the models: only one has recovered. The third image shows that if the result is very poor the parts cannot always recover. The middle row shows BioID images, including glasses. The first image illustrates how the sub-models can allow adaptation to non-symmetric features (crooked mouth). The bottom row shows Expressions Surprise and Expressions Frown. Black points are the full model result and white points show the integrated results from the part-models (common points).

from the full model. If this is very inaccurate, sometimes only one part can recover (top row, column 2) and sometimes recovery is not possible (top row, column 3). The middle left image of Figure 11 illustrates the ability of the method to adapt to non-standard configurations, in this case the mouth is not symmetric with the face. Normally this would need to be included in a training set to be handled by a full face model of this kind.

Figure 12 shows the cumulative error proportion plotted against average point-to-point (P2P) distance for each dataset. In most cases this is a percentage of the interocular distance (IOD). This normalises for the fact that the head size can vary across the data set, which is particularly true of the BioID images. For the artificial data the errors are presented as pixels because of the large pose variation. However, the head size is fairly constant on all images and IOD on the frontal faces is ~100 pixels.

Generally, including the parts improves on the original “NoParts” result. However, examining the cumulative error curve for Video 1 up to the 15% threshold, we observed that performance got worse when the parts were introduced. Extending the curve to show errors beyond 15% revealed that, although the errors were increased the proportion of failures was reduced albeit at larger thresholds (shown in Figure 12). Examining the results for different parts of the face revealed that the mouth area was much improved even at the lower thresholds (also shown in the figure). This indicates that the mouth model was able to correct many of the failures thus pushing the curve higher. For the other data sets the eyes showed the best performance. Figure 9 shows the superimposed results from all datasets broken down by face area. The eye corners and mouth generally fall into well-defined areas of the graph but the eyebrow curves show a great deal of variation. This is unsurprising since the corners of the eyes are the most easily located when marking ground truth on the data and the eyebrows are the most difficult to mark. The eyes are therefore likely to produce the lowest errors in a model search and the eyebrows are likely to produce the highest and most variable. The variant giving the best results for the eye corners is therefore shown on the graphs of Figure 12. For the artificial data the results of using a 2D CLM<sup>2,10</sup> are included for comparison.

In a further breakdown, Figure 13 shows the comparative median average point to point errors for the four different categories in the Expressions database versus the 4 different face areas. The width of the bubble represents the comparative error value. Since some of the values were quite similar, a baseline of 5.4% IOD was subtracted to emphasise the relative performance. Three methods are shown: No parts; Exact parts; and the Feedback results. From this figure we can see that the eyes generally perform well and the eyebrows badly, as already noted. There is indication that the eye area, in particular, seems to benefit from the introduction of parts. The best performance under expression seems to be frowning and it also displays the best improvement using the method. Notice that the exact parts method can make the result worse, whereas in every combination but one (closed eyes and eyebrow area) the feedback method improves the result. This would indicate that this is the best method, certainly on this database, however it is less clear with the videos whether this is the case.

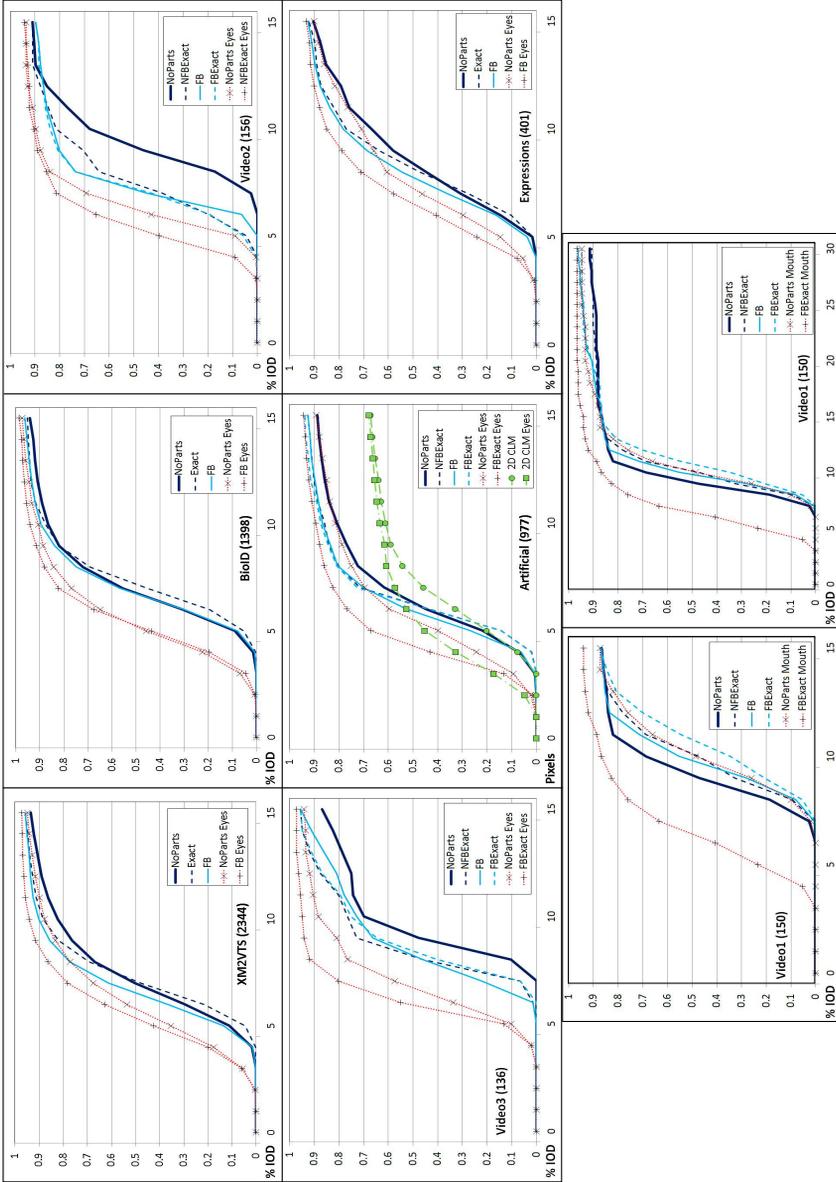


Fig. 12. (Color online) Cumulative Error Distributions. The proportion of data is plotted against the point-to-point error as % inter-ocular distance (IOD), except for the artificial data which is shown in pixels. The number of images in the data set is given in brackets. The best local result is also shown, which is the eyes for all but Video 1. 2D CLM results are also given for the artificial data and an extended error plot is shown for Video 1.

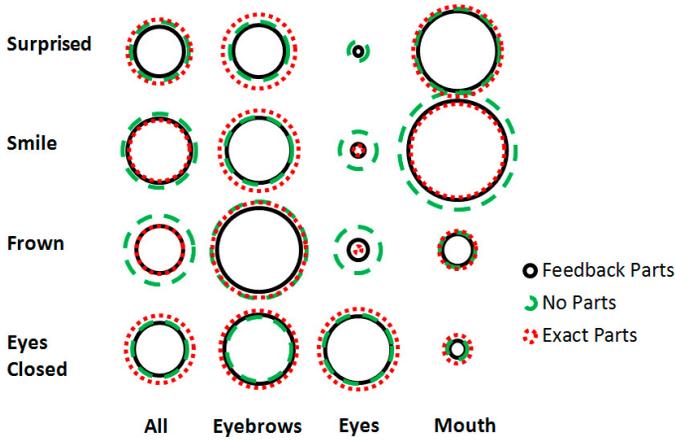


Fig. 13. (Color online) The comparative performance of parts methods on the areas versus expressions of the face (Expressions database only). The diameter of each circle is proportional to the median average point to point error after subtracting a baseline of 5.4%.

## 6.2. Pose estimation

Figure 14 shows the P2P errors broken down by rotation angle for the artificial images. Also shown are the errors on the estimated pose. Here the pose is reported as a quaternion, which represents a rotation about an axis vector. We examined the error on both the angle and the axis. The latter is calculated as the rotation angle between the

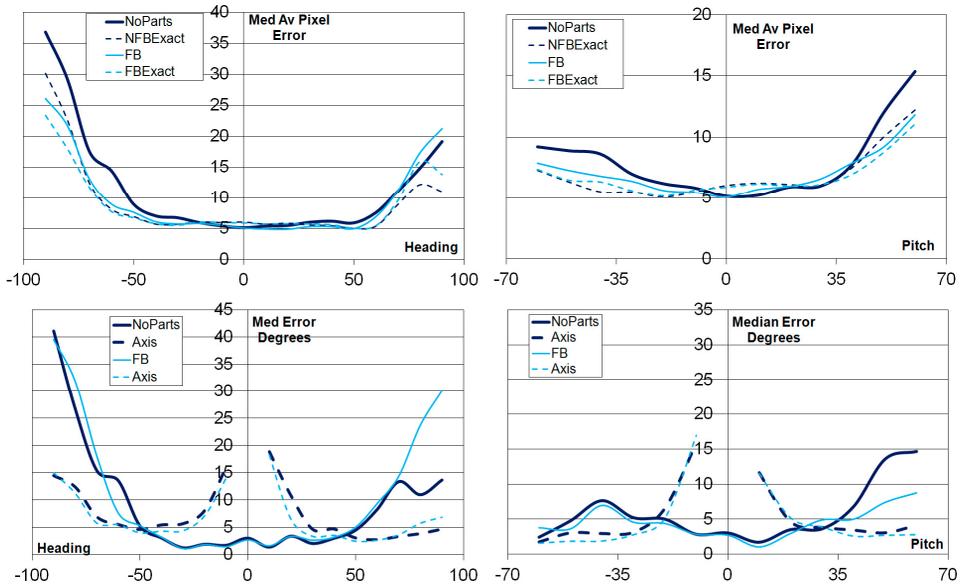


Fig. 14. (Color online) The top row shows the point-to-point errors for the artificial data, broken down by out of plane rotation angle. The bottom row shows pose estimation errors over the same ranges.

known and estimated axes. Since the pose is not affected when using “Exact” methods only the original and FB methods are shown. There are no axis results for 0 degrees rotation because of the high level of uncertainty on the rotation axis at that point. As with Figure 12, including the sub-parts has had a generally beneficial effect. For heading, the improved P2P errors are below 8 pixels for angles between  $\pm 50$  degrees, and the estimated angle error is 5 degrees or less for the same range. For pitch the P2P error is below 13 pixels as far as  $+60$  degrees and below 8 for  $-60$  degrees. The angle estimation error is no more than 9 degrees for the entire range. Since the results for roll (i.e. in plane rotations) were consistent over all angles and methods they are not shown. The median average P2P error varied between 4.75 and 7.7 pixels. The angle estimation error was consistently less than 2.9 degrees.

### 6.3. Discussion

The easiest data sets in this group are XM2VTS and BioID and, as might be expected, the overall performance (all methods) is relatively better on these sets, closely followed by the expressions and the artificial data. The real-world videos obviously present the greatest challenge and the relative performance reflects this.

If we allow that Video 1 was improved in one area only, the mouth, then there is an overall improvement on every other data set from adding in the subparts. The other two videos and the artificial data show the most obvious improvements which implies that using sub-models has the most beneficial effect when used in the more difficult situations, as might be hoped, since it is a corrective technique. For the still images it seems that the best improvement is when the full model is re-fitted to the parts result (FB). This implies that the ability of the parts to correct the full model and deal with non-standard examples has aided the method but a final regularisation step is still needed. However, in the case of the sequence data this is less clear cut. For Videos 2 & 3 and the artificial data, “NoParts” is the worst performer but whether exact or corrected parts should be reported needs further investigation.

When compared to state-of the-art 2D systems such as Refs. 1, 3, 4, our method does not perform as well on the common BioID data set. However, in Refs. 3 and 4 the authors use near frontal faces only and in Ref. 1 state that their experimental data set, acquired from the internet, does not have profile or near profile images and that all faces were detected by an off the shelf face detector. This will tend to exclude not only extreme poses but also unusual lighting and some occlusions. We have shown that our method deals with large rotations, occlusion, a wide variety of data types and conditions and, in addition, can provide an estimate of pose, for those gaze-critical applications.

## 7. Conclusions

We have presented an extension to 3D model search which allows refinement of the results using independent 3D sub-parts. From Figures 12–14 it can be seen that, for all data sets but Video 1, including the sub-parts has had a clearly positive effect. Even in the

case of Video 1, when breaking down the errors between parts of the face, it can be seen that the mouth area is vastly improved. However, there is uncertainty as to whether reporting the exact parts points or refitting the model gives the best individual result. This requires further study and implies that an iterative approach, alternating between the full and parts models, may yield even more improvements.

## Acknowledgments

This project is funded by Toyota Motor Europe who provided the driver videos. We would like to thank Genemation Ltd. for the 3D data markups and head textures.

## References

1. P. N. Belhumeur, D. W. Jacobs, D. J. Kriegman and N. Kumar, Localizing parts of faces using a consensus of exemplars, in *Proc. Computer Vision and Pattern Recognition* (2011), pp. 545–552.
2. D. Cristinacce and T. Cootes, Automatic feature localisation with constrained local models, *Pattern Recognition* **41**(10) (2007) 3054–3067.
3. M. Valstar, B. Martinez, X. Binefa and M. Pantic, Facial point detection using boosted regression and graph models, in *Proc. Computer Vision and Pattern Recognition* (2010), pp. 2729–2736.
4. T. Cootes, M. Ionita, C. Lindner and P. Sauer, Robust and accurate shape model fitting using random forest regression voting, in *Proc. European Conf. on Computer Vision* (Florence, 2012).
5. A. Pentland, B. Moghaddam and T. Starner, View-based and modular eigenspaces for face recognition, in *Proc. Computer Vision and Pattern Recognition* (1994), pp. 1–7.
6. A. Athana, T. Marks, M. Jones, K. Tieu and R. MV, Fully automatic pose-invariant face recognition via 3D pose normalization in *Proc. Int. Conf. on Computer Vision* (2011).
7. C. Vogler, Z. Li, A. Kanaujia, S. Goldenstein and D. Metaxas, The best of both worlds: Combining 3D deformable models with active shape models, in *Proc. Int. Conf. on Computer Vision* (2007), pp. 1–7.
8. J. Xiao, S. Baker, I. Matthews and T. Kanade, Real-time combined 2D+3D active appearance models, in *Proc. Conf. on Computer Vision and Pattern Recognition* (2004), pp. 535–542.
9. S. Romdhani, J. Ho, T. Vetter and D. J. Kriegman, Face recognition using 3-D models: Pose and illumination, *Proc. of the IEEE* **94**(11) (2006) 1977–1999.
10. A. Cauce, D. Cristinacce, C. Taylor and T. Cootes, Locating facial features and pose estimation using a 3D shape model, in *Proc. Int. Symp. on Visual Computing* (Las Vegas, 2009), pp. 750–761.
11. A. Cauce, C. Taylor and T. Cootes, Adding facial actions into 3D model search to analyse behaviour in an unconstrained environment, in *Proc. Int. Symp. on Visual Computing* (Las Vegas, 2010), pp. 132–142.
12. M. A. Fischler and R. A. Elschlager, The representation and matching of pictorial structures, *IEEE Transactions on Computers* **C-22** (1973) 67–92.
13. P. F. Felzenswalb, R. B. Girshick and D. McAllester, Cascade object detection with deformable part models, in *Proc. Computer Vision and Pattern Recognition* (2010), pp. 1–8.
14. L. Zhu, Y. Chen and A. Yuille, Learning a hierarchical deformable template for rapid deformable object parsing, *IEEE Trans. Pattern Analysis and Machine Intelligence* **32**(6) (2010) 1029–1043.

15. M. C. Burl, M. Weber and P. Perona, A probabilistic approach to object recognition using local photometry and global geometry, in *Proc. European Conf. on Computer Vision* (1998), pp. 628–641.
16. G. Hua and Y. Wu, Sequential mean field variational analysis of structures deformable shapes, *Computer Vision and Image Understanding* **101** (2006) 87–99.
17. A. M. Martinez, Recognizing imprecisely localized, partially occluded, and expression variant faces from a single sample per class, *Pattern Analysis and Machine Intelligence* **24**(6) (2002) 748–763.
18. V. Blanz and T. Vetter, A morphable model for the synthesis of 3D faces, in *Proc. SIGGRAPH* (1999), pp. 187–194.
19. J. R. Tena, F. D. I. Torre and I. Matthews, Interactive region-based linear 3D face models, in *Proc. SIGGRAPH* (2011).
20. T. F. Cootes and C. J. Taylor, Active shape models — ‘Smart Snakes’, in *Proc. British Machine Vision Conference* (1992), pp. 266–275.
21. F. L. Bookstein, Principal warps: Thin-plate splines and the decomposition of deformations, *IEEE Transactions on Pattern Analysis and Machine Intelligence* **11**(6) (1989) 567–585.
22. B. Amberg, R. Knothe and T. Vetter, Expression invariant 3D face recognition with a morphable model, in *Proc. Int. Conf. on Automatic Face Gesture Recognition* (Amsterdam, 2008), pp. 1–6.
23. C. Basso and T. Vetter, Registration of expressions data using a 3D morphable model, *Journal of Multimedia* **1**(4) (2006) 37–45.
24. L. Gu and T. Kanade, 3D Alignment of face in a single image, in *Proc. Int. Conf. on Computer Vision and Pattern Recognition* (New York, 2006), pp. 1305–1312.
25. A. Cauce, C. Taylor and T. Cootes, Improved 3D model search for facial feature location and pose estimation in 2D images, in *Proc. British Machine Vision Conference* (Aberystwyth, 2010), pp. 81.1–81.10.
26. T. F. Cootes, D. H. Cooper, C. J. Taylor and J. Graham, Active shape models — Their training and application, *Computer Vision and Image Understanding* **61**(1) (1995) 38–59.
27. P. Viola and M. J. Jones, Robust real-time face detection, *Int. Journal of Computer Vision* **57**(2) (2004) 137–154.
28. K. Messer, J. Matas, J. Kittler and K. Jonsson, XM2VTSDB: The extended M2VTS database, in *Proc. Int. Conf. on Audio and Video-based Biometric Person Authentication* (Washington DC, USA, 1999).
29. O. Jesorsky, K. J. Kirchberg and R. W. Frischholz, Robust face detection using the hausdorff distance, in *Proc. Int. Conf. on Audio and Video-Based Person Authentication* (Halmstaad, Sweden, 2001), pp. 90–95.