

# Robust identification and matching of fiducial points for the reconstruction of 3D human faces from raw video sequences

Raffaella Lanzarotti, Alberto Borghese, Paola Campadelli  
Dipartimento di Scienze dell'Informazione  
Università degli Studi di Milano  
Via Comelico, 39/41 20135 Milano, Italy  
{ lanzarotti, borghese, campadelli }@dsi.unimi.it

## Abstract

*We describe an algorithm for the automatic features detection in 2D color images of either frontal or rotated human faces. Such features allow to initialize robustly a bundle-adjustment in order to fit a generic 3D face model to the images. The algorithm first identifies the sub-images containing each feature (eyes, nose and lips), afterwards, it processes them separately to extract fiducial points. The features are looked for in down-sampled images, the fiducial points are identified in the high-resolution ones. The method uses both color and shape information and does not require any manual setting or operator intervention. It has been tested on a database of 130 color images.*

## 1. Introduction

3D Human face digital models are used in a wide range of applications such as medicine [17], video compression/coding, virtual reality [5], and computer graphics [23][21][4][20].

To obtain these models, the common solution is to resort to 3D scanner devices: a set of video cameras recognize either laser spots/stripes [22] or structured light patterns [1] [25] [14] allowing to produce a set of 3D points (range data) sufficiently dense to produce a digital model.

However, apart from their high costs, these systems have two drawbacks. The acquisition time is often in the order of seconds, while it is known that micro head movements require acquisition times of the order of few tenths of a second. Moreover, the subject is forced to be positioned accurately with respect to the scanning structure (the light emitter and the cameras), which makes these system suitable to laboratory environments only.

A more advantageous approach is to resort to multiple natural images picked up from different points of view, by

standard photo cameras, and to reconstruct a 3D model (a 3D mesh) processing the images [25][7][8]. In this framework, the problem is reframed as to find the 3D mesh, whose projections are compatible with all the images. This involves determining the 3D positions of the vertexes of the mesh and the geometrical parameters associated to each view. In [3] these parameters are determined through a global optimization process. The results are quite remarkable, but a fine tuning of the optimization parameters is required and local minima associated to bad solutions do occur. As a result, the method is not suitable for everyday applications.

A different solution is based on estimating first a dense set of corresponding points on the images. A "universal" 3D digital model [21] [16] is then fitted to these points. The identification of the corresponding points is not a trivial task: classical geometrical feature detectors [15] fail on a smooth surface with no texture as the human face is. This difficulty is the reason why in [11] the points were identified by attaching 183 colored markers over the human face. However, the necessity to attach markers before the acquisition makes this approach of little interest for everyday applications. In the last few years, approaches which try to identify corresponding points through adequate algorithms have been developed [7] [8] [6] [19] [9]. They are based on bundle-adjustment techniques [10] regularized by a pre-defined model of the human face. The drawback of these systems lies in their initialization: a few points have to be manually supplied in at least one image.

In [12] and [13] a technique for the automatic fiducial points localization is presented. Although it gives good results, it works only under restricted assumptions (regarding the head position and the illumination conditions) and it is computationally expensive. A solution based on the color is reported in [14]. It is a low-cost algorithm, but, in our experience, color is not sufficient to characterize the features unambiguously.

In this paper, we describe a technique for identifying automatically and with great reliability a set of fiducial points. Our method uses both colour and shape information. Results on a data base of 130 images, taken at different orientations, are reported and discussed.

## 2. Method overview

The method we propose works on images of face' foregrounds. We thus ignore the problem of localizing the faces in more complex scenes [24][26]. We acquire color images with homogeneous and light-colored background, and frontal and diffuse illumination. The images have a wide variety of resolutions. Faces can be either in frontal position or rotated around the head vertical axis of  $30^\circ$  at most, moreover the method is robust to little lateral head tilts of about  $10^\circ$ . In any case, the completeness of the facial images is requested: no occlusions, no wearing glasses, no beard and closed mouth.

The algorithm has two hierarchical processing modules: the first identifies four sub-images, each tightly containing one of the features of interest (the two eyes, the nose and the lips); the latter module is specialized in localizing fiducial points on the found features with high accuracy.

A preliminary solution to this problem, experimented on a low number of images, has been presented in [18]. Now we have modified completely the module for localizing the sub-images making it more general, more robust and efficient. Moreover we have partially changed the method for localizing the fiducial points.

## 3. Identification of sub-images

In this module the images are down-sampled to a resolution in the range between  $150 \times 170$  to  $300 \times 340$  pixels.

At first, we divide the features of interest from the skin and the background, by clustering the grey level image into three clusters through the method presented in [2]. The lightest grey level represents the background, the intermediate the skin and the darkest represents both the features and other dark pixels of the image (for example the hair) [Fig.1].

To localize the features of interest the largest connected component,  $S$ , of the intermediate grey-level cluster (corresponding to the skin) is found. Then all the pixels surrounded by pixels of  $S$  and belonging to the darkest grey cluster are identified and set to 1; all the others are set to 0. What we obtain is the *feature image* [Fig.2].

We observe that, besides the features of interest, also few pixels, such as those corresponding to shadows or to the hair, are set to 1. Further processing is therefore required to isolate the features of interest in separated sub-images.



Figure 1. Clustered images



Figure 2. Features images

We proceed localizing the eyes, then, in order to look for the lips and the nose, the attention can be concentrated upon a restricted image area.

### 3.1 Eyes localization

In order to determine the set of rows (*eyes band*) which contains the eyes, we apply the template matching to the *features image* [Fig.2], searching the two eyes. The difficulty is that we are not looking for an object with a fixed shape: there are inter-personal differences, differences caused by the orientation of the head, and variations introduced by the clustering [Fig.3]. For this reason we adopt a binary template which models the two eyes in a very rough way. It consists of two symmetric blobs placed side by side each being large enough to overlap to the region corresponding to an eye in the *features image*. A single template [Fig.4] has been used for all the images which are of significant different size, thus showing a desirable scale-independence property.

Among the positions with the highest cross-correlation, we maintain the 10 which satisfy also the following symmetry condition: the cross correlation between half the template (one blob) and the sub-images on the left and on the right of the found position are compared. The position is rejected if the results are not similar enough, that is their ratio is lower than 0.7. On all the images of our database most of the points calculated in this way are positioned in



Figure 3. Some example of "clustered" eyes



Figure 4. Eyes template

the *eyes band* (on average 8 out of 10) [Fig.5]. This allows us to select the band easily and with high reliability.

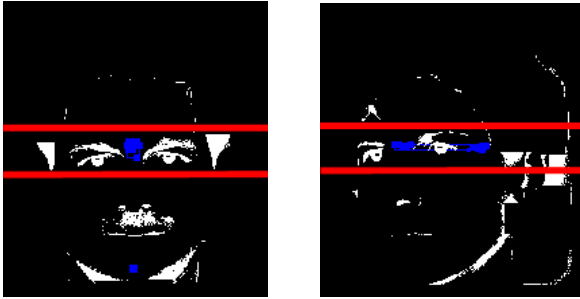


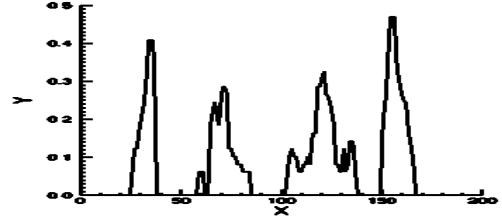
Figure 5. Template matching results and identified eyes band

Furthermore, we apply the vertical projection to the *eyes band* [Fig.6(a)], finding more than two distinct peaks [Fig.6(b)] due to shadows, hair, and ears. For each distinct peak we take the corresponding grey level piece of image and give it as input to a multilayer perceptron which has been trained to distinguish images representing eyes from those representing other features that can be found in the same band. The neural network has one hidden layer with ten units and two outputs neurons; it has been trained by backpropagation on 200 images and tested on 200 images. The classification error is about 9% but only very rarely the network is wrong on both eyes, so we know that an error has been made since we have a face with a number of eyes different from two. This allows us to treat this situation separately and we are studying a way to correct the error.

Finally, the two sub-images chosen in the previous step are refined in order to tailor them around each eye. For this aim, we calculate the vertical derivatives that highlight the eyebrows: the eyebrow position gives the eye upper and lateral limits. As regard the lower eyes limits, we look for them in the corresponding position of the *feature image*



(a)



(b)

Figure 6. (a) Eyes band; (b) Vertical projection

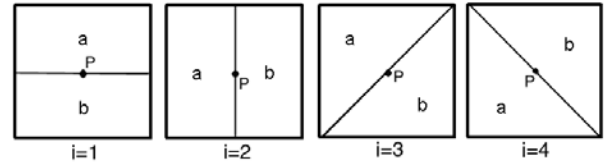


Figure 7. Windows

where the iris are very evident.

### 3.2 Lips and nose localization

Once determined the eyes' bounding boxes, we move to the lips and nose localization. We can limit the research region to the one below the eyes. To build the bounding box around the lips we use both grey level and color information. On the grey level image a non linear edge detector is applied [30]. The detector uses local statistical information: a square window centered around a pixel  $P$  is divided into two sub-windows of equal size in four different ways [Fig.7]

For each sub-division  $i$  the following function  $D_i$  (diversity) is evaluated:

$$D_i = \alpha \Delta m - (1 - \alpha) \Delta \sigma$$

where

$$\Delta m = |m_a - m_b|,$$

$$\Delta \sigma = |\Delta \sigma_a - \Delta \sigma_b|,$$

$m_a, m_b, \sigma_a, \sigma_b$  are respectively the mean and the standard deviation in the regions  $a$  and  $b$ , and  $\alpha$  is a constant ( $0 \leq \alpha \leq 1$ )

The maximum diversity  $D(P) = \max_{i=1..4} \{D_i\}$  is assigned to  $P$ . We use this method with a  $7 \times 7$  window size positioned on a pixel out of 3. The lips and the nose are roughly characterized by pixels with high horizontal and vertical diversity respectively.

To better localize the lips, we combine color information with the output of the edge detector. Since in the Cr color image plane the pixels corresponding to the lips have a very high value [Fig.8], we threshold the Cr image in a rectangular region centered around the line identifying the highest number of horizontal edge pixels. The thresholding operation maintains only the 10% of the pixels with the highest values. The bounding box of the largest region found in this way corresponds to the lips sub-images.

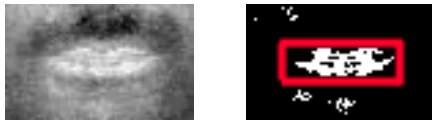


Figure 8. Cr plane and its binarization

The nose is localized in the remaining portion of the image in the region where the horizontal diversity is different from 0.

At this stage, we have defined the four sub-images [Fig.9] that have to be analyzed more accurately in order to locate the fiducial points and, if possible, the contour of the features of interest. To do this we return to the highest-resolution images.

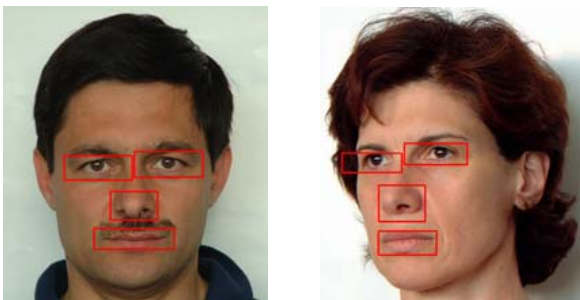


Figure 9. Subdivisions

## 4. Identification of fiducial points

In this module we work on the single sub-images separately and at the highest-resolution.

## 4.1 Eyes

### 4.1.1 Pupil and eye extreme points recognition

The first step in the recognition of the eye fiducial points is the identification of the pupil.

Some authors (e.g. [27], [29]) suggest finding the darkest pixels. However, this approach is suitable only to particular illumination conditions. Instead, in general, we cannot guarantee this: as can be seen in figure 10, a reflex spot makes several pixels internal to the pupil very bright.



Figure 10. Eye with a reflex and its localization

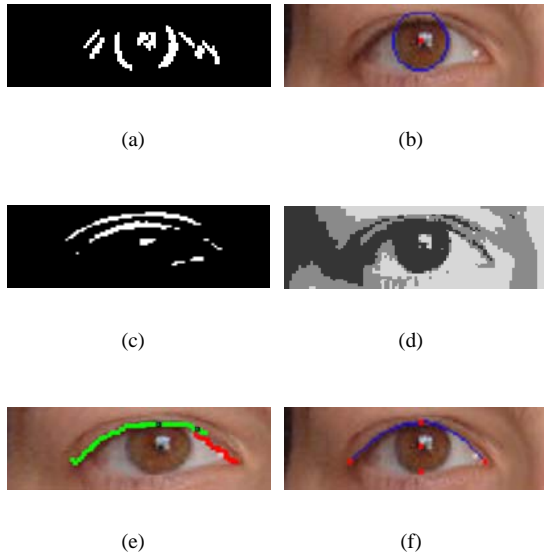
In our previous work [18] we resorted to template matching. As we looked for the iris, the template had the shape of a circle (with ray equal to 8 pixels). We remind here that at this stage we are working on images of different dimensions which means that the use of a fix template is not suitable. We could introduce adaptive templates, but we found more flexible the use of the Hough transform to look for circumferences. The transform is applied to the binarized horizontal derivative of the eye grey-level image [Fig.11(a)] since it highlights the two iris' vertical sides, which are the best references to find the iris' border. Among the found circumferences we choose the one which obtained the most of the votes [Fig.11(b)]. The centre of the pupil  $P$  corresponds to the centre of the found circumference.

Then, the lower extreme of the iris is identified as the intersection of the found circumference and the vertical axis  $A$  passing through  $P$ . It is not possible to look for the upper visible extreme of the iris in an analogue way since the point found would exceed the correct position: in most people the iris' upper half is partially occluded by the eyelid.

What we do is to calculate the absolute value of the vertical derivative,  $\Delta_y$ , of the eye image and to threshold it keeping the 10% of the pixels with the highest values [Fig.11(c)]. The upper extreme of the eye is localized on the axis  $A$  in correspondence to the first pixel, starting from the point  $P$ , that belongs to the region obtained by the thresholding. This is done after the identification and elimination of reflex spot eventually present in the iris [Fig.10].

### 4.1.2 Eye corners and upper arc determination

To define the shape of the eye, we apply an edge following algorithm based on the idea of the hysteresis thresholding [28]. The two thresholds are automatically determined in order to keep the 10% and the 40% of the pixels with the



**Figure 11. Processing of the eye's image:** (a) Input to the Hough transform; (b) Circumference found by the Hough transform ; (c) Binarization of the vertical derivative; (d) Clustering of the equalized image; (e) Edge-following; (f) Detection of the eye upper arc and of the eye fiducial points.

highest values of  $\Delta_y$ . The contour  $C$  obtained in this way [Fig.11(e), *green line*] well-defines part of the shape of the eye, but it does not identify precisely the corners. This is the reason why we have the necessity to adopt another method. We start from the consideration that the white internal part of the eye is quite evident. In order to emphasize the contrast between it and the surrounding regions, we equalize the grey-level image and then cluster it requiring 3 clusters [Fig.11(d)]. The upper border,  $U$ , of such white part [Fig.11(e), *red line*] gives a good indication to detect both the internal corner of the eye and its outline: combining  $U$  and the contour  $C$ , we determine the parabola which better approximates them. Then we localize the eye internal corner on the parabola at the height of the end of  $U$ , and the external corner on the parabola at the height of the end of  $C$ . This allows us to find a good approximation of the upper part of the eye outline [Fig.11(f)].

## 4.2 Lips

Our goal for the lips is the determination of the corners and of their middle point at least. To get a more robust estimate, the entire upper lip outline is determined, and, for frontal face images only, also the lower lip outline is de-

scribed. The technique adopted is based on the following steps: localization of the lips cut and of the lips lowest point, localization of the lips corners and of the lips outline.

### 4.2.1 Recognition of the lips cut and of the lips lowest point

To determine the lips cut we apply the Sobel vertical derivative operator to the mouth sub-image. We then cluster it into three clusters: one associated to light-to-dark vertical transitions, one to dark-to-light and one to no meaningful transitions [Fig.12(a)]. We then determine the largest dark connected region and we extract, for every  $x$ , the upper pixel belonging to it [Fig.12(b)]. The line  $L$  connecting these pixels well represents the lips cut apart from its extremes: in some cases, the line ends before the corners, in other, it exceeds them. At this point we are able to recognize the lips lowest point as the lower extreme of the largest connected white region under the line  $L$  [Fig.12(b)].

### 4.2.2 Recognition of the lips corners and of the lips outline

We go back to the original color image considering the mouth box only and cluster it into four clusters. The darkest one identifies the shadow corresponding to the lips cut and gives a precise information regarding the horizontal position of the corners. Combining this with the information given by the line  $L$ , we obtain the correct corners position [Fig.12(c)].

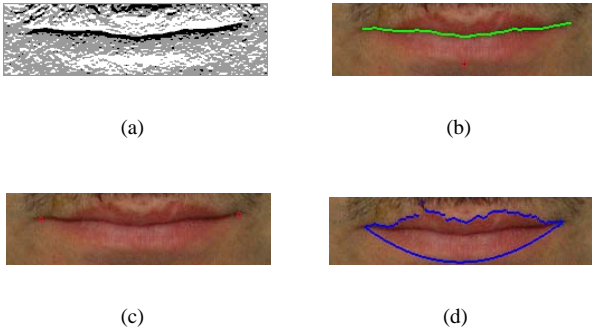
We apply the algorithm in the CIE-Luv color space, being the one which has given experimentally the best results consistently, with respect to RGB, CIE-Lab and HSV color spaces.

In order to identify the lips outline, we apply a second time the same clustering algorithm requiring two clusters. This lead to obtain one cluster,  $M$ , which is roughly associated to the lips and the other to the surrounding skin. The outline of the region  $M$ , is quite precise in its upper part,  $B$ , apart from its corners. We thus obtain the upper lips outline straightening  $B$  to finish in the found corners.

Moreover, for frontal face images, we can identify the lips lower outline as the parabola passing trough the two corners and the lips lowee extreme point. [Fig.12(d)].

## 4.3 Nose

Regarding this feature, we are interested in finding the nose tip. Observing numerous faces' images, we concluded that the nose is characterized by two dark regions, corresponding to the nostrils, and a light region, corresponding to the reflex of the light on the nose tip. In order to identify these regions, we applied the clustering algorithm to the nose grey-level sub-image requiring four clusters. The



**Figure 12. Processing of the lips' image: (a) Clustering of the vertical derivative with 3 clusters; (b) Detection of the lips cut and of the lips lowest point; (c) Corners recognition; (d) Outline.**

lightest grey-level pixels correspond to the reflex, and the darkest grey-level pixels correspond to the nostrils [Fig. 13]. Moreover we observed that the middle point between the nostrils gives a good vertical localization of the nose tip and that the lower extreme of the region corresponding to the reflex gives a good horizontal localization. The intersection of these two axis gives the nose tip.



**Figure 13. Clustering and Tip of the nose**

## 5. Discussion

### 5.1 Results

The method described has been experimented on 130 images acquired in different illumination conditions. They represent either frontal or rotated faces of Caucasian women and men. We asked to the people to have a neutral expression keeping the mouth closed and the eyes opened. We have not dealt with the case of men with beard.

The algorithm works well on images of very different scales and finds fiducial points using both color and shape information.

The decomposition module localizes the sub-images representing the features of interest with high confidence on both frontal and rotated images. The most of the errors are introduced by the classifications done by the neural network. On our database this module fails on the 5% of the images.

The fiducial points are detected with high accuracy (errors of 1 or 2 pixels are negligible). The only critical point is the eye external corner, where we find an error that, normalized with respect to the eye dimension, can be evaluated as about 10% in the worst case. This has to be considered in the reconstruction phase.

Further improvements are required. As regard the sub-images identification, different learning algorithms can be experimented to recognize the eye with higher accuracy (e.g. SVM). As regard fiducial points determination, the algorithm is not completely satisfactory on rotated images, since it does not determine the lips lower outline correctly and it does not localize the nose tip precisely.

A final consideration has to be done about the program running time. The time necessary to process the first module on a Pentium III, 800MHz, 256Mb of RAM is about 5 seconds, independently of the original dimension of the images (on the contrary the second module computational cost depends strictly on the dimension of the images). The running time can be certainly reduced developing the algorithm in a compiled language and optimizing the code.

### 5.2 Application to the human faces 3D reconstruction

In this work we have presented a method for localizing a certain number of fiducial points on 2D images. Applying it to sequences of images, it allows to generate the 3D model of the acquired human face: having at least 6 corresponding points on the different views, we can recover the cameras parameters automatically as well as the 3D positions of the fiducial points in space. Subsequently, a dense set of points (*range data*) can be obtained applying a least square matching algorithm [10]. Then the range data can be used to deform a generic 3D face model [21] to fit the face of the particular human subject.

We observe that, determining the fiducial points automatically, some of them, in some image, could be localized wrongly. What we have to do is to recognize and reposition them. First, we can proceed evaluating the quality of the found fiducial points: after having determined the range data, we compute, for each fiducial point, the distance between its position and the back-projection, on the same image, of the corresponding range data point (*residual*). Erroneously identified fiducial points have high residual, allowing us to recognize and discard them. Then we have to reposition the rejected points: for each of them, pairing its view parameters with the ones of another view where the same fiducial point has been localized correctly, we work out the essential matrix. Using it we can calculate the epipolar line associated to the correctly localized fiducial point and plot it onto the other image, solving the problem.

After having checked and correct the fiducial points,

the parameters have to be computed again and a bundle-adjustment procedure [10] can be run on all the images to refine both the 3D positions of the fiducial points and the parameters of the views. We conclude with the alignment between the corrected range data and the generic 3D face model.

As a net result, a very accurate 3D model can be obtained in a completely automatic fashion, with no intervention from the operator.

## References

- [1] Web address: <http://www.minolta-3d.com/>.
- [2] M. Arbib and T. Uchiyama. Color image segmentation using competitive learning. *IEEE Transactions on pattern analysis and machine intelligence*, 16:1197–1206, 1994.
- [3] V. Blanz and T. Vetter. A morphable model for the synthesis of 3d faces. *Proceedings ACM Siggraph '98*, pages 81–89, 1998.
- [4] V. Blanz and T. Vetter. A morphable model for the synthesis of 3d faces. *Proceedings ACM Siggraph '99*, pages 187–194, 1999.
- [5] G. Breton, C. Bouville, and D. Pelè. A 3d facial animation engine for real time applications. *Proceedings WEB3d*, pages 15–22, 2001.
- [6] Q. Chen and G. Medioni. Building human faces models from two images. *Proceedings IEEE second workshop multimedia signal processing*, pages 117–122, 1998.
- [7] N. D'Apuzzo. Automated photogrammetric measurement of human faces. *Int. Archives of Photogrammetry and Remote Sensing, Hakodate, Japan*, B5(32):402–407, 1998.
- [8] P. Fua. Regularized bundle-adjustment to model heads from image sequences without calibration data. *Int. J. Computer Vision*, 2(38):153–171, 2000.
- [9] P. Fua and C. Miccio. Animated heads from ordinary images: A least squares approach. *Computer Vision and Image Understanding*, 3(75):247–259, 1999.
- [10] A. Gruen and H. Beyer. System calibration through self-calibration. *Calibration and Orientation of Cameras in Computer Vision. Washington D.C.*, 1992.
- [11] B. Guenter, C. Grimm, D. Wood, H. Malvar, and F. Pighin. Making faces. *Proceedings Siggraph98*, pages 55–67, 1998.
- [12] R. Herpers, M. Michaelis, K. H. Lichtenauer, and G. Sommer. Edge and keypoint detection in facial regions. *Proceedings International conference on automatic face and gesture recognition*, pages 212–217, 1996.
- [13] R. Herpers and G. Sommer. An attentive processing strategy for the analysis of facial features. In H. Wechsler, P. J. Phillips, V. Bruce, F. F. Soulié, and T. S. Huang, editors, *Face recognition*, pages 457–468. London, 1998.
- [14] R. Hsu and A. Jain. Face modeling for recognition. *Proceedings IEEE International Conference Image Processing*, pages 693–696, 2001.
- [15] J. Shi and C. Tomasi. Good features to track. *IEEE Conference on Computer Vision and Pattern Recognition*, pages 593–600, 1994.
- [16] P. Kalra, A. Mangili, N. M. Thalmann, and D. Thalmann. Simulation of muscle actions based on rational free form deformations. *Proceedings Eurographics*, pages 113–121, 1992.
- [17] R. Koch, M. Gross, F. Carls, D. C. Buren, G. Fankhauser, and Y. Parish. Simulating facial surgery using finite element models. *Proceedings Siggraph '96*, pages 421–428, 1992.
- [18] R. Lanzarotti, N. Borghese, and P. Campadelli. Automatic features detection for overlapping face images on their 3d range models. *Proceedings ICIAP, Palermo*, 2001.
- [19] W. Lee and N. Magnenat-Thalmann. Fast head modeling for animation. *Proceedings Image and Vision Computing*, 4(18):355–364, 2000.
- [20] Y. Lee, D. Terzouloulos, and K. Waters. Realistic facial modeling for animation. *Proceedings Siggraph '95*, pages 55–62, 1995.
- [21] F. Parke and K. Waters. Computing facial animation. *Wellesly MA*, 1998.
- [22] M. Petrov, A. Talapov, T. Robertson, A. Lebedev, A. Zhilayaev, and L. Polonskiy. Optical 3d digitisers: Bringing life to the virtual world. *IEEE Computer Graphics and Applications*, pages 28–37, 1998.
- [23] F. Pighin, J. Hecker, D. Lischinski, R. Szeliski, and D. H. Salesin. Synthesizing realistic facial expressions from photographs. *Proceedings Siggraph98*, pages 75–84, 1998.
- [24] T. Poggio and K. Sung. Example-based learning for view-based human face detection. *Proceedings of the ARPA Image Understanding Workshop*, II:843–850, 1994.
- [25] M. Proesman and L. V. Gool. One-shot 3d shape and texture acquisition of facial data. *Proceedings ICCV98*, pages 215–221, 1998.
- [26] H. Rowley, S. Baluja, and T. Kanade. Neural networkbased face detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20(1):23–38, 1998.
- [27] R. Stiefelwagen, J. Yang, and A. Weibel. Tracking eyes and monitoring eye gaze. *Proceedings of the Workshop on Perceptual User Interfaces (PUI)*, pages 98–100, 1997.
- [28] E. Trucco and A. Verri. *Introductory techniques for 3d computer vision*. Prentice-Hall, Inc., New Jersey, 1998.
- [29] Y. Yan and K. Challapali. A system for the automatic extraction of 3-d facial feature points for face model calibration. *Proceedings ICIP*, 2000.
- [30] P. Zamperoni. Feature extraction. In H. Maitre and J. Zim-Justin, editors, *Progress in feature processing*. 1996.