MEASUREMENT AND MODELING OF HUMAN FACES FROM MULTI IMAGES

Nicola D'Apuzzo

Institute of Geodesy and Photogrammetry, ETH-Hoenggerberg, 8093 Zurich, Switzerland, nicola@geod.baug.ethz.ch

Commission V, WG V/6

KEYWORDS: Automation, Photogrammetry, Surface, Measurement, Visualization, Photo-Realism

ABSTRACT:

Modeling and measurement of the human face have been increasing by importance for various purposes. Laser scanning, coded light range digitizers, image-based approaches and digital stereo photogrammetry are the used methods currently employed in medical applications, computer animation, video surveillance, teleconferencing and virtual reality to produce three dimensional computer models of the human face. Depending on the application, different are the requirements. Ours are primarily high accuracy of the measurement and automation in the process. The method presented in this paper is based on multi-image photogrammetry. The equipment, the method and results achieved with this technique are here depicted. The process is composed of five steps: acquisition of multi-images, calibration of the system, establishment of corresponding points in the images, computation of their 3-D coordinates and generation of a surface model. The images captured by five CCD cameras arranged in front of the subject are digitized by a frame grabber. The complete system is calibrated using a reference object with coded target points, which can be measured fully automatically. To facilitate the establishment of correspondences in the images, texture in the form of random patterns can be projected from two directions onto the face. The multi-image matching process, based on a geometrical constrained least squares matching algorithm, produces a dense set of corresponding points in the five images. Neighborhood filters are then applied on the matching results to remove the errors. After filtering the data, the three-dimensional coordinates of the matched points are computed by forward intersection using the results of the calibration process; the achieved mean accuracy is about 0.2 mm in the sagittal direction and about 0.1 mm in the lateral direction. The last step of data processing is the generation of a surface model from the point cloud and the application of smooth filters. Moreover, a color texture image can be draped over the model to achieve a photorealistic visualization. The advantage of the presented method over laser scanning and coded light range digitizers is the acquisition of the source data in a fraction of a second, allowing the measurement of human faces with higher accuracy and the possibility to measure dynamic events like the speech of a person.

1. INTRODUCTION

Modeling and measurements of the human face have wide applications ranging from medical purposes (Banda et al., 1992; Koch et al. 1996; Motegi et al., 1996; D'Apuzzo, 1998; Okada, 2001) to computer animation (Pighin et al., 1998; Blanz and Vetter, 1999; Lee and Magnenat-Thalmann, 2000; Liu et al., 2000; Marschner et al., 2000; Sitnik and Kujawinska, 2000), from video surveillance (CNN, 2001) to lip reading systems (Minaku et al., 1995), from video teleconferencing to virtual reality (De Carlo et al., 1998; Borghese and Ferrari, 2000; Fua, 2000; Shan et al., 2001). How realistic and accurate the obtained shape is, how long it takes to get a result, how simple the equipment is and how much the equipment costs are the issues that must be considered to model the face of a real person.

The different approaches to enable the reconstruction of a human face can be classified depending on the requirements. For animation, virtual reality and teleconferencing purposes, the photorealistic aspect is essential. In contrast, high accuracy is required for medical applications. Two major groups can also be distinguished based on their data source: the first using range digitizers and the second using only images.

To date, the most popular measurement technique is laser scanning (Motegi et al., 1996; Hasegawa, 1999; Marschner et al., 2000; Okada, 2001), for example the head scanner of Cyberware (Cyberware, 2002). These scanners are expensive and the data is usually noisy, requiring touchups by hand and sometimes manual registration. Another solution is offered by the structured light range digitizers (Proesmans and Van Gol, 1996; Wolf, 1996; Sitnik and Kujawinska, 2000) which are usually composed of a stripe projector and one or more CCD cameras. These can be used for face reconstruction with relatively inexpensive equipment compared to laser scanners. The accuracy of both systems is satisfactory for static objects, however their acquisition time ranges from a couple of seconds to half of a minute, depending on the size of the surface to measure. Thus, a person must remain stationary during the measurement. Not only does this place a burden on the subject, but it is also difficult to obtain stable measurement results. In fact, even when the acquisition time is short, the person moves slightly unconsciously.

A different approach to face modeling uses images as source data. Various image-based techniques have been developed. They can be distinguished by the type of used image data: a single photograph, two orthogonal photographs, a set of images, video sequences or multi-images acquired simultaneously.

Parametric face modeling techniques (Blanz and Vetter, 1999) start from a single photograph to generate a complete 3-D model of the face. Exploiting the statistics of a large data set of 3-D face scans, the face model is built by applying pattern classification methods. The results are impressively realistic, however the accuracy of the reconstructed shape is low.

A number of researchers have proposed creating models from two orthogonal views (Ip and Yin, 1996). Manual intervention is required for the modeling process by selecting feature points in the images. It is basically a simplified method to produce realistic models of human faces. The obtained shape does however not reproduce the real face precisely. To solve this problem, some solutions (Lee and Magnenat-Thalmann, 2000) work in combination with range data acquired by laser scanners. Another image-based method consists of automatically extracting the contour of the head from a set of images acquired around the person (Matsumoto et al., 1999; Zengh, 1994). The obtained data are combined to form a volumetric model of the head. The set of images can be generated moving a single camera around the head or having the camera fixed and the face turning. The systems are fast and completely automatic, however the accuracy of the method is low.

Video sequences based methods (Pighin et al., 1998; Fua, 2000; Liu et al., 2000; Shan et al., 2001) uses photogrammetric techniques to recover stereo data from the images. A generic 3-D face model is then deformed to fit the recovered (usually noisy) data. These techniques are full automatic but may perform poorly on face with unusual features or other significant deviations from the normal.

High accuracy measurement of real human faces can be achieved by photogrammetric solutions which combine a thorough calibration process with the use of synchronized CCD cameras to acquire simultaneously multi-images (Banda, 1992; D'Apuzzo, 1998; Minaku et al., 1999; Borghese and Ferrari, 2000; D'Apuzzo, 2001). To increase the reliability and robustness of the results some techniques use the projection of an artificial texture on the face (Banda, 1992; D'Apuzzo, 1998). The high accuracy potential of this approach results however in a time expensive processing.

For our purposes, we are interested in an automatic system to measure the human face relatively fast and with high accuracy. We have therefore chosen a photogrammetric solution. Five synchronized CCD cameras are used to acquire simultaneously multi-images of a human face and artificial random texture is projected onto the face to increase the robustness of the measurement. The processing consists of five steps: acquisition of images of the face from different directions, determination of the camera positions and internal parameters, establishment of dense set of corresponding points in the images, computation of their 3-D coordinates and generation of a surface model. Due to the simultaneous acquisition of all the required data, the proposed method offers the additional opportunity to measure dynamic events.

In this paper, we present the equipment used, the method and the achieved results.

2. METHOD

In this section, are described the system for data acquisition and the method used for its calibration and depicted the methods for the measurement and modeling of the human face from the acquired multi-images.

An advantage of our method is the acquisition of the source data in fractions of a second, allowing the measurement of human faces with high accuracy and the possibility of measuring dynamic events such as speech. Another advantage of our method is that the developed software can be run on a normal home PC reducing the costs of the hardware. We are developing a portable, inexpensive and accurate system for the measurement and modeling of the human face.

2.1 Data acquisition and calibration

Figure 1 shows the setup of the used image acquisition system. It consists of five CCD cameras arranged convergently in front of the subject. The cameras are connected to a frame grabber which digitizes the images acquired by the five cameras at the resolution of 768x576 pixels with 8 bits quantization.



Figure 1. Setup of cameras and projectors

A color image of the face without random pattern projection is acquired by an additional color video camera placed in front of the subject. It is used for the realization of a photorealistic visualization.

Since the natural texture of the human skin is relatively uniform, the projection of an artificial texture onto the face is required to perform robustly the matching process. A random pattern (see figure 2) is preferred to regular patterns to avoid possible mismatches and its resolution has to be fine enough to result in the images in structures the size of few pixels. The use of two projectors enables a focused texture even on the lateral sides of the face; figure 3 shows the five images acquired by the CCD cameras.



Figure 2. Projected random pattern



Figure 3. Multi-images of a face with random pattern projection

The system is calibrated using a 3-D reference frame with coded target points whose coordinates in space are known (see figure 4). These are fully automatically recognized and measured in the images (Niederoest, 1996). The results of the calibration process are the exterior orientation of the cameras (position and rotations: 6 parameters), parameters of the interior orientation of the cameras (camera constant, principle point, sensor size, pixel size: 7 parameters), parameters for the radial and decentering distortion of the lenses and optic systems (5 parameters) and two additional parameters modeling differential scaling and shearing effects (Brown, 1971). A thorough determination of these parameters modeling distortions and other effects is required to achieve high accuracy in the measurement.



Figure 4. Calibration frame with coded targets

2.2 Matching process

Our approach is based on multi-image photogrammetry using images acquired simultaneously by synchronized cameras. The multi-image matching process is based on the adaptive least squares method (Gruen, 1985) with the additional geometrical constraint of the matched point lying on the epipolar line. Figure 5 shows an example of the result of the least squares matching (LSM) algorithm: the black boxes represent the patches selected in the template image (left) and the affine transformed in the search images (center and right), the epipolar lines are drawn in white.



Figure 5. Geometrical constrained LSM; left: template image, center and right: two search images

The automatic matching process produces a dense and robust set of corresponding points, starting from few seed points. The seed points may be manually defined in each image, generated semi-automatically (defining them only in one image) or fully automatically. The manual mode is used for special cases where the automatic modes could fail; the seed points have to be selected manually with an approximation of at least 2 pixels in each image: LSM is then applied to find the exact position. In the semi-automated mode the seed points have to be selected manually only in the template image; the corresponding points in the other images are established automatically by searching for the best matching results along the epipolar line (see figure 6). This mode is the most convenient for normal cases of static surface measurement: it is fast but leave the operator the choice where to set the seed points. The fully automatic mode is useful in cases with dynamic surface measurement from multi-image video sequences, where the number of multi-image sets to be processed could be very large. In this case, Foerstner interest point operator (Foerstner and Guelch, 1987) is used to

automatically determine in the template image marking points where the matching process may perform robust results; the corresponding points in the other images are then established with the same process as for the semi-automatic mode.



Figure 6. Semi automated seed point definition

After the definition of the seed points, the template image is divided into polygonal regions according to which of the seed points is closest (Voronoi tessellation). Starting from the seed points, the set of corresponding points grows automatically until the entire polygonal region is covered (see figure 7).



Figure 7. Search strategy for the matching process

The matcher uses the following strategy: the process starts from the seed point, shifts horizontally in the template and in the search images and applies the least squares matching algorithm in the shifted location. If the quality of the match is good, the shift process continues horizontally until it reaches the region boundaries; if the quality of the match is not satisfactory, the algorithm computes the matching again, changing some parameters (e.g. smaller shifts from the neighbor, bigger sizes of the patches). The covering of the entire polygonal region of a seed point is achieved by sequential horizontal and vertical shifts. The process is repeated for each polygonal region until the whole image is covered.

To evaluate the quality of the result of the matching, different indicators are used: a posteriori standard deviation of the least squares adjustment, standard deviation in x and y directions, displacement from the start position in x and y directions and distance to the epipolar lines. Thresholds for these values can be defined for different cases, according to the texture and the type of the images.

Before beginning the three dimensional processing, filters can be applied to the 2-D matching data to minimize the number of possible errors. The Voronoi tessellation produces an irregular grid (see figures 8 and 9, left) of points in the template image, therefore, the set of matched points has first to be uniformed to a regular grid before the application of any filters. This is achieved by matching all the points shifted to the regular grid (see figures 8 and 9, right).



Figure 8. Regularization of the matched point grid



Figure 9. Regularization; left: seed points and matched points in the template image; right: after regularization of the grid

For the definition of the filter, the smoothed characteristic of the surface of the human face is taken in account: as shown in figure 10, the transformed image patches of neighboring points belonging to a common smoothed surface have similar shapes. A neighborhood filter is therefore applied to the set of matched points checking for the local uniformity of the shape of the transformed image patches.



Figure 10. Points matched in the neighborhood

The complete matching process (definition of seed points, automatic matching, filtering) is flexible and can also be performed without orientation and calibration information. This functionality can be useful, for example, if the orientation is not accurate enough or unknown. In these special cases, only the image information is used by the least squares matching algorithm. Obviously, the robustness of the result of the process decreases; however the quality of the set of matched points remains satisfactory.

A dedicated software was developed for the face measurement process. Figures 11 and 12 show its user friendly graphical interface.



Figure 11. Graphical user interface of the face measurement software; seed points definition



Figure 12. Graphical user interface of the face measurement software; matching results and visualization of the computed 3-D point cloud

The required intervention of the operator for the matching process is reduced to the semi-automatic definition of about ten seed points and the selection of a contour of the region to measure. The operation can be performed in a couple of minutes, then the process will continue completely automatically. On a Pentium III 600 MHz machine, about 20,000 points are matched on half of the face in approximately 10 minutes.

2.3 Modeling and visualization

Since the human face is a steep surface and both sides of the face are not visible to the same camera, the five acquired images are used as two separate set of triplets, one for each side of the face. They are processed separately and at the end, the results are merged into a single data set.

The 3-D coordinates of the matched points are computed by forward ray intersection using the orientation and calibration data of the cameras. The achieved accuracy of the 3-D points is about 0.2 mm in the sagittal direction and about 0.1 mm in the lateral direction.



Figure 13. Top: measured 3-D point cloud (45,000 points), bottom: after filtering and thinning (10,000 points)

As shown in figure 13 (top), the point cloud is very dense (45,000 points) and the region of overlap of the two joined data set can be observed in the center line of the face. To overcome the redundant data and remove eventual outliers, Gaussian filters (Borghese and Ferrari, 2000) are applied to the 3-D point cloud and the data is afterwards thinned (see figure 13 bottom). For surface measurement purposes, the computed 3-D point cloud is satisfactory. In case of visualization, a complete model of the face with texture has to be produced. A meshed surface is therefore generated from the 3-D point cloud by 2.5-D Delauney triangulation and to achieve photorealistic visualization, the natural texture acquired by the color video camera is draped over the model of the face. Figure 14 shows the surface model, the texture image and two views of the resulted face model with texture, figure 15 shows two other examples of face models.



Figure 14. Photorealistic visualization; top: shaded surface model and texture image; bottom: face model with texture

3. CONCLUSIONS

A process for an automated measurement of the human face from multi-images acquired by five synchronized CCD cameras has been presented. The main advantages of this method are its flexibility, the reduced costs of the hardware and the possibility to perform surface measurement of dynamic events.

ACKNOWLEDGEMENT

The work reported here was funded in part by the Swiss National Science Foundation.

REFERENCES

Banda, F.A.S. et al., 1992. Automatic Generation of Facial DEMs. In: *Int. Archives of Photogrammetry and Remote Sensing*, Vol. XXIX, Part B5, pp. 893-896.

Blanz, V. and Vetter, T., 1999. A Morphable Model for the Synthesis of 3D Faces. In: *SIGGRAPH'99 Conf. Proc.*, pp. 187-194.

Borghese, A. and Ferrari, S., 2000. A Portable Modular System for Automatic Acquisition of 3-D Objects. *IEEE Trans. on Instrumentation and Measurement*, 49(5), pp. 1128-1136.

Brown, D.C., 1971. Close-Range Camera Calibration. *Photogrammetric Engineering and Remote Sensing*, 37(8), pp. 855-866.

CNN, 2001. Facing Up to Airport Security Fear. http://europe.cnn.com/2001/US/09/28/facial.recognition.QandA (accessed 11 June 2002)

Cyberware, 2002. Head and Face Color 3D Scanner Model 3030. http://www.cyberware.com/products/psInfo.html (accessed 11 June 2002)

D'Apuzzo, N., 1998. Automated Photogrammetric Measurement of Human Faces. In: *Int. Archives of Photogrammetry and Remote Sensing*, Hakodate, Japan, Vol. XXXII, Part B5, pp. 402-407.

D'Apuzzo, N., 2001. Photogrammetric Measurement and Visualisation of Blood Vessel Branching Casting: A Tool for Quantitative Accuracy Tests of MR-, CT- and DS-Angiography. In: *Videometrics and Optical Methods for 3D Shape Measurement,* San Jose, USA, Proc. of SPIE, Vol. 4309, pp. 204-211.

D'Apuzzo, N., 2002. Surface measurement and surface tracking of human body parts from multi image video sequences. *ISPRS Journal of Photogrammetry and Remote Sensing*, 56(4).

De Carlo, D. et al., 1998. An Anthropometric Face Model Using Variational Techniques. In: *SIGGGRAPH'98 Conf. Proc.*, pp. 67-74.

Foerstner, W. and Guelch, E., 1987. A Fast Operator for Detection and Precise Location of Distinct Points, Corners and Centres of Circular Features. In: *Proc. of the Intercommission Conference on Fast Processing of Photogrammetric Data*, Interlaken, Switzerland, pp. 281-305.

Fua, P., 2000. Regularized Bundle-Adjustment to Model Heads from Image Sequences without Calibration Data. *Int. Journal of Computer Vision*, 38(2), pp. 153-171.



Figure 15. Photorealistic visualization; two other examples of face models

Gruen, A., 1985. Adaptive Least Squares Correlation: A Powerful Image Matching Technique. *South African Journal of Photogrammetry, Remote Sensing and Cartography*, 14(3), pp. 175-187.

Hasegawa, K. et al., 1999. A High Speed Face Measurement System. In: *Proc. of Vision Interface '99*, Trois-Rivières, Canada, pp. 196-202.

Ip, H.H.S. and Yin, L., 1996. Constructing a 3D Individualized Head Model from Two Orthogonal Views. *The Visual Computer*, 12, pp. 254-266.

Koch, R.M. et al., 1996. Simulating Facial Surgery Using Finite Element Models. In: *SIGGRAPH96 Conference Proceeding*, New Orleans, USA.

Lee, W.-S. and Magnenat-Thalmann, N., 2000. Fast Head Modeling for Animation. *Image and Vision Computing Journal*, 18(4), pp. 355-364.

Liu, Z. et al., 2000. Rapid Modeling of Animated Faces from Video. In: *Proc. of the 3rd Int. Conf. on Visual Computing (Visual2000)*, Mexico City, pp. 58-67.

Marschner, S.R. et al., 2000. Modeling and Rendering for Realistic Facial Animation. In: *Proc. of the 11th Eurographics Workshop on Rendering*, Brno, Czech Replublic

Matsumoto, Y. et al., 1999. CyberModeler: A Compact 3D Scanner Based on Monoscopic Camera. In: *Three-Dimensional Image Capture and Applications II*, San Jose, USA, Proc. of SPIE, Vol. 3640, pp. 3-10.

Minaku, S. et al, 1995. Three-Dimensional Analysis of Lip Movement by 3-D Auto Tracking System. In: *Int. Archives of Photogrammetry and Remote Sensing*, Zurich, Switzerland, Vol. XXX, Part 5W1.

Motegi, N. et al., 1996. A Facial Growth Analysis Based on FEM Employing Three Dimensional Surface Measurement by a Rapid Laser Device, *Okajimas Folia Anatomica Japonica*, 72(6), pp. 323-328.

Niederoest, M., 1996. *Codierte Zielmarken in der digitalen Nahbereichsphotogrammetrie*, Diplomarbeit, Institut für Geodaesie und Photogrammetrie, ETHZ, Zurich, (in German).

Okada, E., 2001. Three-Dimensional Facial Simulations and Measurements: Changes of Facial Contour and Units Associated with Facial Expression. *Journal of Craniofacial Surgery*, 12(2), pp. 167-74.

Pighin, F. et al., 1998. Synthesizing Realistic Facial Expressions from Photographs. In: *SIGGRAPH'98 Conf. Proc.*, Orlando, USA, pp. 75-84.

Proesmans, M. and Van Gol, L., 1996. Reading Between the Lines. In: *SIGGPRAPH'96 Conf. Proc.*, pp. 55-62.

Shan, Y. et al., 2001. Model-Based Bundle Adjustment with Application to Face Modeling. In: *Proc. of the* 8th *Int. Conf. on Computer Vision (ICCV01) Vol. II*, Vancouver, Canada, pp. 624-651.

Sitnik, R. and Kujawinska, M., 2000. Opto-Numerical Methods of Data Acquisition for Computer Graphics and Animation Systems. In: *Three-Dimensional Image Capture and Applications III*, San Jose, USA, Proc. of SPIE Vol. 3958, pp. 36-43.

Wolf, H.G.E., 1996. Structured Lighting for Upgrading 2D-Vision system to 3D. In: *Proc. of Int. Symposium on Laser, Optics and Vision for Productivity and Manufacturing I*, Besancon, France, pp. 10-14.

Zheng, J.Y., 1994. Acquiring 3-D Models from Sequences of Contours. *IEEE Trans. Patt. Anal. Machine Intell.*, 16(2), pp. 163-178.