

Single-Shot Deep Volumetric Regression for Mobile Medical Augmented Reality

Florian Karner^{1,2}, Christina Gsaxner^{1,2,3}, Antonio Pepe^{1,2}, Jianning Li^{1,2}, Philipp Fleck¹, Clemens Arth¹, Jürgen Wallner^{2,3}, and Jan Egger^{1,2,3(\boxtimes)}

 1 Institute of Computer Graphics and Vision, Graz University of Technology, Graz, Austria

egger@tugraz.at

² Computer Algorithms for Medicine Laboratory (Café-Lab), Graz, Austria

³ Department of Oral and Maxillofacial Surgery, Medical University of Graz, Graz, Austria

Abstract. Augmented reality for medical applications allows physicians to obtain an inside view into the patient without surgery. In this context, we present an augmented reality application running on a standard smartphone or tablet computer, providing visualizations of medical image data, overlaid with the patient, in a video see-through fashion. Our system is based on the registration of medical imaging data to the patient using a single 2D photograph of the patient. From this image, a 3D model of the patient's face is reconstructed using a convolutional neural network, to which a pre-operative CT scan is automatically registered. For efficient processing, this is performed on a server PC. Finally, anatomical and pathological information is sent back to the mobile device and can be displayed, accurately registered with the live patient, on the screen. Hence, our cost-effective, markerless approach needs only a smartphone and a server PC for image processing. We present a qualitative and quantitative evaluation using real patient photos and CT from the clinical routine in facial surgery, reporting overall processing times and registration errors.

Keywords: Augmented Reality \cdot 3D face reconstruction \cdot 3D registration \cdot Deep learning \cdot Volumetric Regression Network \cdot Handheld devices \cdot Smartphone \cdot Tablet

1 Introduction

In the last years the use of digital imaging tools to support clinical diagnosis and treatment pathways have undergone a remarkable rate of technological revolution in many medical fields. This is especially true for the cranio-maxillofacial and head and neck complex where computer-assisted technologies such as automated segmentation tools, digital software packages for three dimensional (3D) preoperative visualization and operation planning or 3D printed templates of complex facial bone structures used for surgical implant adaption or others

© Springer Nature Switzerland AG 2020

T. Syeda-Mahmood et al. (Eds.): ML-CDS 2020/CLIP 2020, LNCS 12445, pp. 64–74, 2020. https://doi.org/10.1007/978-3-030-60946-7_7

have become the goldstandard in big clinical centers [13,27,28]. These computer assisted technologies mostly work on the basis of routinely performed computed tomography (CT) or on cone beam CT scans which are daily performed in each clinical center [19,20].

In cranio-maxillofacial surgery, digital imaging data acquisition and especially CT scans are one of the standard imaging tools and are increasingly used in today's clinical routine to support clinical diagnosis and treatment plans. The usage of routinely performed digital imaging data sets such as CT scans, has increased in the last years, because of the increased digital data storage in combination with a fast data access in the clinical centers, the improved resolution that comes with new image scanner generations, resulting in more accurate anatomical imaging data, and the simultaneously reduced imaging acquisition time clinically needed for medical image data creation [17]. Therefore, software programs for processing medical image data have become a central tool in clinical medicine and are subject to continuous ongoing technological developments [26].

Augmented Reality (AR) presents an interesting opportunity to display medical imaging data, as they have been shown to provide a more intuitive interface for visualization than the standardized 2D slice view [24]. In medical AR systems, the registration of image data to the patient in the physician's view is the integral part of the technology. A straight forward approach to this problem is to use a fiducial-based method, where markers are rigidly attached to the patient [5]. These markers are tracked, either in an intrinsic fashion by the camera using computer vision methods [11,12], or extrinsically using a tracking system [2,16,18]. However, this requires excessive preparation as well as calibration, and might cause discomfort to the patient. As an alternative, surface registration algorithms present a more natural solution to the image-to-patient registration problem [15]. For obtaining surface information of the patient in this context, several methods, like depth sensors [6,7,21,22,25] or stereo cameras [1,29], are commonly used.

Instead, we propose to use an adapted convolutional neural network, called Volumetric Regression Network (VRN) by Jackson et al. [10], to reconstruct a 3D model of the patient's face from a single 2D patient photograph. We automatically register pre-interventional, volumetric imaging data to this 3D model, which enables a video see-through AR application running on a mobile device, such as a smartphone or tablet computer.

2 Materials and Methods

Our goal was to build an application for image-to-patient registration using only a single 2D facial image and a CT scan from a patient. We built a webbased client/server system by extending the existing StudierFenster framework (http://studierfenster.tugraz.at) [30,31]. The user sends data to the server via a mobile application. After that, the reconstruction and registration takes place on the server. The result of the registration step is sent back to the mobile app,



Fig. 1. Overview of the proposed client/server architecture.

where an augmented view of the patient is displayed accordingly. Figure 1 shows the proposed client/server architecture.

2.1 Data Acquisition and Preprocessing

We collect volumetric imaging data from the head and neck area of patients from clinical routine [8]. The purpose of this data is two-fold: First, it contains medical information to display in the AR environment. Second, it provides an accurate 3D model of the patients skin surface, which we exploit for image-topatient registration. For visualization, structures of interest, such as the skull, are segmented from the imaging data using semi-automatic or manual methods. Then, a Marching Cubes algorithm [14] is applied to extract meshes of these structures. Furthermore, the skin surface is segmented using a simple thresholding approach, and a point cloud representation of the skin surface is created, which is later used for registration to the 3D model of the patient. We transform all 3D models to have their coordinate origin on the tip of the nose, which is simply determined by choosing the point with the largest z coordinate. Finally, the data is deployed to our server.

2.2 Android Application

We developed our mobile application for Android devices using the Unity 3D game engine in combination with the ARCore AR platform. Our application enables two modes, between which the user can switch. One mode is working with the front camera of the device, resulting in something like a "magic mirror" system. The second mode works with the physical back camera of the phone. Depending on the chosen camera, our application follows different pipelines to obtain an AR overlay, as shown in Fig. 2: For the front camera, we use facial detection and tracking for obtaining an augmented view of the patient, for the back camera, we make use of our proposed single-shot pipeline. The reason for

these two modes is, that the front cameras have in general no simultaneous localization and mapping (SLAM) system available, because most SLAM use cases (marker tracking, 3D models, etc.) and applications (like IKEA Place), and games (like Pokémon GO) make no sense for the front camera. The front camera is mostly used for selfies and video chatting (Snapchat, FaceTime, etc.). Hence, AR SDKs do not provide SLAM systems for the front camera at all. However, a SLAM system, like it is available for the back cameras, makes an AR registration much more precise, compared to simple 2D face detection/tracking, like its implemented for the front cameras, because of the precise tracking in three-dimensional environments (continuous tracking of the camera pose and environment extracting features points for calculation of 3D world points).



Fig. 2. Workflow of our augmented reality application. Depending on the chosen camera (front or back), our system follows different pipelines.

Front Camera Pipeline. The front camera uses ARCore's Augmented Faces module. First, a request is sent to the server to deliver the meshes for visualization and registration of the current patient. ARCore detects a person's face and estimates a 3D mesh of the facial surface, characterized by three landmarks: left forehead, right forehead and tip of the nose. We anchor the meshes received from the server at the nose tip detected by ARCore to obtain an AR overlay of the patient with virtual content.

Back Camera Pipeline. For the back camera, our proposed single-shot method is employed. To enable the placement of virtual content in a common world reference frame, we use the environmental understanding and concurrent odometry and mapping capabilities of ARCore. First, our application maps the environment around the patient using simultaneous localization and mapping (SLAM) implemented in ARCore. The user creates a 2D photograph of the patient, which is sent to the server. Simultaneously, we perform a raycast on the created environmental map through the patient's nose tip and save the hit point, which will later act as an anchor for virtual content. On the server side, 2D to 3D reconstruction and registration is performed, which will be explained in greater detail in the following section. The transformed meshes are sent back to our application, and, since they are also centered around the nose tip, we can use the hit point obtained from raycasting to place them in the AR environment accordingly.

2.3 Server Backend

To keep the computational load on the mobile device small, we perform heavy computations on a server PC. The application running on the server requires as input a frontal 2D photo of the patient, which is sent by the client, as well as 3D models from pre-interventional imaging obtained in step Sect. 2.1, which are already stored on the server.

3D Face Reconstruction Using VRN. To reconstruct a 3D model of the patient's face from a single image, we used a Convolutional Neural Network (CNN) developed by Jackson et al. [10], denoted VRN. Contrary to other solutions to the 3D face reconstruction problem, VRN does not require multiple facial images, but instead works on a single 2D image. It is able to reconstruct the entire facial geometry from a large variety of input poses and facial expressions. VRN performs this task by directly regressing a complete 3D volume using a CNN with an hourglass-like architecture, without the need to fit a 3D model to the input. Therefore, it produces fast and reliable outputs. A RGB photograph of the patient, captured with a mobile device running our application, is sent to the server and serves as input to VRN. The image has to be downsampled and rescaled to fit with the input size expected by the VRN. The output of the network is a 3D reconstructed mesh of the patient's face, from which we extract a point cloud for surface registration. Since the mesh obtained from VRN has an arbitrary unit of scale - we scale it to align it with the scale of pre-interventional imaging, which is acquired in millimeters (mm).

Surface Registration. To register the 3D model of the patients face reconstructed by VRN with the pre-interventional imaging data of the patient, we use a global registration method, followed by a refinement stage using iterative closest point (ICP), as proposed by Holz et al. [9]. For global registration, we compute fast point feature histograms [23] in both point clouds and match them iteratively using a random sample consensus algorithm [4]. This results in a coarse alignment of point clouds, which we refine by using point-to-plane ICP [3].

Mesh Transmission. The result of the surface registration is a 4×4 transformation matrix. With this matrix, content in the CT coordinate frame, such as an anatomical structure mesh (e.g., the bones of the skull) or pathological structure meshes (e.g., tumors) are transformed and sent back to the client. The meshes are transmitted by HTTP messages.

3 Results

We evaluated our system with ten medical scans from human subjects, seven CT scans and three magnetic resonance imaging (MRI) scans. The preoperative CT scans are from head and neck cancer patients from the clinical routine. Besides the CT, we collected frontal, high-resolution photographs of the patients, which are routinely taken before the facial operations by our clinical partners. These photos served as input to the VRN for the reconstruction and registration with the corresponding CT. In addition, we used several MRI scans of healthy subjects for testing.

3.1 Quantitative Evaluation

To quantitatively evaluate the registration error between the 3D reconstructed patient photograph and pre-interventional imaging data, we calculate the closest point registration error (CPRE) between the two point clouds. CPRE measures the average distance between points in a reference point cloud \mathbf{P}_u^R with points $u = 1, 2, ..., N_R$ to their nearest points in a model point cloud \mathbf{P}_v^M with v = $1, 2, ..., N_M$, which is registered to the reference using transformation T:

$$CPRE = \frac{1}{N_R} \sum_{u=1}^{N_R} \min_{v \in [1, N_M]} || \mathbf{P}_u^R - T \cdot \mathbf{P}_v^M ||.$$
(1)

Table 1 presents the total CPRE, as well as CPRE in x, y and z directions separately. Additionally, the times our system needs from end-to-end are presented. The overall time measures the following steps: (1) sending a 2D photo to the server backend, (2) performing 2D to 3D reconstruction, (3) registration of the 3D photo to the pre-operative CT and (4) transmission of registration results back to the mobile application. The measurements were repeated five times and the average was determined.

3.2 Qualitative Evaluation

For qualitative analysis, Fig. 3 displays the reconstruction and registration results of three patients. The first row shows the surface extraction from the patient's pre-operative CT. The second row shows the 3D reconstructed patient photos (anonymized with a black bar in the eye area due to patient privacy). The bottom row shows the results of the automatic registration between the surface extracted from CT (first row) and the 3D photo reconstruction (second row).

We tested our application on three healthy humans to simulate the use case with live people. All three test subjects had an MRI scan done beforehand, from which we extracted point cloud information for registration, as well as anatomical information for visualization. Figure 4 shows qualitative results of this experiment. Furthermore, we obtained qualitative feedback from facial- and neurosurgeons, who confirmed an applicability of our application for tumor cases (including biopsies) were accurate navigation is not necessary, but a better understanding of three-dimensional relations might still be beneficial. They stated that a 3D visualization provided on a tablet or smartphone screen can easily, quickly and efficiently support clinical diagnosis and treatment pathways in all patients with maxillofacial or cranial tumors, cysts or any other expansive soft or hard tissue processes.



Fig. 3. Reconstruction and registration results of three patients. The first row shows the surface extraction from the patient's CT. The second row shows the 3D reconstructed patient photos (anonymized with a black bar in the eye area due to patient privacy). The bottom row shows the results of the automatic registration between the surface extracted from CT and the 3D photo reconstruction.

4 Discussion

For an objective evaluation of our proposed reconstruction and registration pipeline, we used real patient photos and CT scans from the clinical routine in facial surgery. Since our system requires minimal setup time and expenditure, it can be used for applications where the usage of a commercial navigation



Fig. 4. Examples of several live cases. Left: ARCore is using its SLAM system to generate a point cloud of the environment (white points). Middle: Two examples of our AR visualization on live cases. We display the subject's skulls from the front and from the side. Right: Application usage from the user's perspective, showing the subjects skull and a tumor in the facial area (red). (Color figure online)

Table 1. Total closest point registration error (CPRE), as well as CPRE in x, y and z directions between the 3D reconstructed 2D (patient) photos and the pre-operative CT/MRI scans. For subjects 1–7, CT scans and patient photos were available. Subjects 8–10 are live human subjects with MRI scans. In addition, we present the overall processing time of our system consisting of following steps: (1) sending a 2D photo from smartphone to the server, (2) performing the 2D to 3D reconstruction, (3) registration of the 3D photo reconstruction to the pre-operative CT, and (4) transmission of the registration result back to the mobile application.

Subject		1	2	3	4	5	6	7	8	9	10	Total
CPRE (mm)	Avg	6.1	5.0	5.3	5.7	6.9	5.5	6.0	7.3	6.5	7.7	6.2
	Avg_X	2.5	2.2	2.5	2.7	2.2	2.3	2.2	3.1	3.3	4.0	2.7
	Avg_Y	2.6	3.2	3.2	3.4	5.3	3.4	4.2	4.5	2.7	5.1	3.7
	Avg_Z	3.8	2.1	2.2	2.3	2.1	2.3	2.0	2.8	4.0	3.0	2.4
Time (s)	Mean	10.1	9.5	8.7	10.4	10.4	11.6	10.4	9.5	10.3	11.0	10.2

system, which requires preparation and a setup times of up to 30 min, is not feasible. Our quantitative results show that the VRN is able to reconstruct a patients face sufficiently well for successful registration with volumetric medical imaging data. The demands in the registration accuracy of medical AR systems for image guided interventions are exceptionally high, and conventional mobile hardware is, at this point, usually not able to meet these requirements. Therefore, we focus our system on an application which does not require sub-millimeter precision: Pre-operative visualization, facilitating a rough estimation of the target localization or surgical entry point. Our medical partners attested that for this application, our registration error of around six mm is acceptable. Compared to navigation systems with much higher accuracy (but that might not be available in all medical centers, especially in smaller ones), our application runs on a low-cost smartphone and does not need any preparation or set up time.

In addition, our evaluation grants a deeper insight into VRN. By default, the 3D reconstruction from the 2D photo with the VRN does not provide any real-world size unit, because the distance of the image plane to the face and the person's head size is not known from a single photo. The availability of a millimeter-precise 3D face model from medical scans, such as CT and MRI, gives us the unique opportunity to also estimate the scale of the reconstruction returned by the VRN, which is an aspect that has so far been overlooked by the non-medical computer vision community.

5 Conclusion

In this contribution, we introduced an AR video see-through system for medical applications, which performs image-to-patient registration using only a single patient photo and volumetric imaging data of the patient. Our application runs solely on a mobile device, such as a smartphone or tablet. It allows the visualization of anatomical structures (like bones) and pathological structures (like tumors) for an augmented view of the patient. One main advantage is that our approach, unlike previous work, does not depend on any external devices, such as navigation systems or depth sensors; it only needs a standard smartphone or tablet, which makes it very cost-effective. Furthermore, it works without any markers and does not require complicated calibration.

Our results show that an accurate overlay of virtual content with the real scene can be achieved with our pipeline. Therefore, our system could be used for various medical applications involving the head and face, for example, preoperative visualization or educational purposes. Future work will evaluate our approach in clinical routine and introduce more sophisticated visualizations and interactions.

Acknowledgment. This work received funding from the Austrian Science Fund (FWF) KLI 678-B31 (enFaced - Virtual and Augmented Reality Training and Navigation Module for 3D-Printed Facial Defect Reconstructions). Further, this work sees the support of CAMed - Clinical additive manufacturing for medical applications (COMET K-Project 871132), which is funded by the Austrian Federal Ministry of Transport, Innovation and Technology (BMVIT), and the Austrian Federal Ministry for Digital and Economic Affairs (BMDW), and the Styrian Business Promotion Agency (SFG), and the TU Graz Lead Project (Mechanics, Modeling and Simulation of Aortic Dissection). Moreover, the Summer Bachelor (SB) Program of the Institute of Computer Graphics and Vision (ICG) of the Graz University of Technology (TU Graz). Finally, we want to point out to our medical online framework Studierfenster (www.studierfenster.at), where an automatic single-shot 3D face reconstruction and registration module has been integrated, and a video tutorial is available on YouTube (3D Face Reconstruction and Registration with Studierfenster: https://www.youtube.com/watch?v=DbbFm9XxlGE).

73

References

- Ahn, J., Choi, H., Hong, J., Hong, J.: Tracking accuracy of a stereo-camera-based augmented reality navigation system for orthognathic surgery. J. Oral Maxillofac. Surg. 77(5), 1070.e1–1070.e11 (2019)
- Chen, X., et al.: Development of a surgical navigation system based on augmented reality using an optical see-through head-mounted display. J. Biomed. Inform. 55, 124–131 (2015)
- Chen, Y., Medioni, G.: Object modelling by registration of multiple range images. Image Vis. Comp. 10(3), 145–155 (1992)
- 4. Choi, S., Zhou, Q.Y., Koltun, V.: Robust reconstruction of indoor scenes. In: Conference on Computer Vision and Pattern Recognition (CVPR) (2015)
- Eggers, G., Mühling, J., Marmulla, R.: Image-to-patient registration techniques in head surgery. Int. J. Oral Maxillofac. Surg. 35(12), 1081–1095 (2006)
- Fan, Y., Jiang, D., Wang, M., Song, Z.: A new markerless patient-to-image registration method using a portable 3D scanner. Med. Phys. 41(10), 101910 (2014)
- Gsaxner, C., Pepe, A., Wallner, J., Schmalstieg, D., Egger, J.: Markerless imageto-face registration for untethered augmented reality in head and neck surgery. In: Shen, D., et al. (eds.) MICCAI 2019. LNCS, vol. 11768, pp. 236–244. Springer, Cham (2019). https://doi.org/10.1007/978-3-030-32254-0_27
- Gsaxner, C., Wallner, J., Chen, X., Zemann, W., Egger, J.: Facial model collection for medical augmented reality in oncologic cranio-maxillofacial surgery. Scientific Data 6(1), 310 (2019)
- Holz, D., Ichim, A.E., Tombari, F., Rusu, R.B., Behnke, S.: Registration with the point cloud library: a modular framework for aligning in 3-D. IEEE Robot. Autom. Mag. 22(4), 110–124 (2015)
- Jackson, A.S., Bulat, A., Argyriou, V., Tzimiropoulos, G.: Large pose 3D face reconstruction from a single image via direct volumetric CNN regression. In: International Conference on Computer Vision (ICCV) (2017)
- Jayender, J., Xavier, B., King, F., Hosny, A., Black, D., Pieper, S., Tavakkoli, A.: A novel mixed reality navigation system for laparoscopy surgery. In: Frangi, A.F., Schnabel, J.A., Davatzikos, C., Alberola-López, C., Fichtinger, G. (eds.) MICCAI 2018. LNCS, vol. 11073, pp. 72–80. Springer, Cham (2018). https://doi.org/10. 1007/978-3-030-00937-3_9
- Jiang, T., Zhu, M., Chai, G., Li, Q.: Precision of a novel craniofacial surgical navigation system based on augmented reality using an occlusal splint as a registration strategy. Sci. Rep. 9(1), 501 (2019)
- Lamecker, H., et al.: Automatic segmentation of mandibles in low-dose CT-data. Int. J. Comput. Assisted Radiol. Surg. 1, 393 (2006)
- Lorensen, W.E., Cline, H.E.: Marching cubes: a high resolution 3d surface construction algorithm. ACM Siggraph Comput. Graph. 21(4), 163–169 (1987)
- Markelj, P., Tomaževič, D., Likar, B., Pernuš, F.: A review of 3D/2D registration methods for image-guided interventions. Med. Image Anal. 16(3), 642–661 (2012)
- Maruyama, K., et al.: Smart glasses for neurosurgical navigation by augmented reality. Operative Neurosurgery 15(5), 551–556 (2018)
- McCann, M.T., Nilchian, M., Stampanoni, M., Unser, M.: Fast 3d reconstruction method for differential phase contrast x-ray CT. Optics Express 24(13), 14564– 14581 (2016)
- Meulstee, J.W., et al.: Toward holographic-guided surgery. Surgical Innov. 26(1), 86–94 (2019)

- Olabarriaga, S.D., Smeulders, A.W.: Interaction in the segmentation of medical images: a survey. Med. Image Anal. 5(2), 127–142 (2001)
- Orentlicher, G., Goldsmith, D., Horowitz, A.: Applications of 3-dimensional virtual computerized tomography technology in oral and maxillofacial surgery: current therapy. J. Oral Maxillofacial Surgery 68(8), 1933–1959 (2010)
- Pepe, A., et al.: Pattern recognition and mixed reality for computer-aided maxillofacial surgery and oncological assessment. In: Proceedings Biomedical Engineering International Conference (BMEiCON), pp. 1–5. IEEE, January 2019
- Pepe, A., et al.: A marker-less registration approach for mixed reality-aided maxillofacial surgery: a pilot evaluation. J. Dig. Imag. 32(6), 1008–1018 (2019). https:// doi.org/10.1007/s10278-019-00272-6
- Rusu, R.B., Blodow, N., Beetz, M.: Fast point feature histograms (FPFH) for 3D registration. In: ICRA (2009)
- 24. Sielhorst, T., Feuerstein, M., Navab, N.: Advanced medical displays: a literature review of augmented reality. J. Disp. Technol. 4(4), 26 (2008)
- Sylos Labini, M., Gsaxner, C., Pepe, A., Wallner, J., Egger, J., Bevilacqua, V.: Depth-awareness in a system for mixed-reality aided surgical procedures. In: Huang, D.-S., Huang, Z.-K., Hussain, A. (eds.) ICIC 2019. LNCS (LNAI), vol. 11645, pp. 716–726. Springer, Cham (2019). https://doi.org/10.1007/978-3-030-26766-7_65
- Tucker, S., et al.: Comparison of actual surgical outcomes and 3-dimensional surgical simulations. J. Oral Maxillofacial Surg. 68(10), 2412–2421 (2010)
- 27. Wallner, J., et al.: Clinical evaluation of semi-automatic open-source algorithmic software segmentation of the mandibular bone: practical feasibility and assessment of a new course of action. PLoS ONE **13**(5), 156–165 (2018)
- Wallner, J., Schwaiger, M., Hochegger, K., Gsaxner, C., Zemann, W., Egger, J.: A review on multiplatform evaluations of semi-automatic open-source based image segmentation for cranio-maxillofacial surgery. In: Computer Methods and Programs in Biomedicine, p. 105102 (2019)
- Wang, J., Shen, Yu., Yang, S.: A practical marker-less image registration method for augmented reality oral and maxillofacial surgery. Int. J. Comput. Assisted Radiol. Surg. 14(5), 763–773 (2019). https://doi.org/10.1007/s11548-019-01921-5
- Weber, M., Wild, D., Wallner, J., Egger, J.: A client/server based online environment for the calculation of medical segmentation scores. In: EMBC, pp. 3463–3467 (2019). https://doi.org/10.1109/EMBC.2019.8856481
- Wild, D., Weber, M., Wallner, J., Egger, J.: Client/server based online environment for manual segmentation of medical images. CoRR abs/1904.08610 (2019). http:// arxiv.org/abs/1904.08610