

A STATISTICAL SHAPE SPACE MODEL OF THE PALATE SURFACE TRAINED ON 3D MRI SCANS OF THE VOCAL TRACT

Alexander Hewer¹⁻³, Ingmar Steiner^{2,3}, Timo Bolkart^{1,3}, Stefanie Wuhler⁴, Korin Richmond⁵

¹Saarbrücken Graduate School of Computer Science, Germany

²DFKI Language Technology Lab, Saarbrücken, Germany

³Cluster of Excellence Multimodal Computing and Interaction, Saarland University, Germany

⁴INRIA Grenoble Rhône-Alpes, France

⁵Centre for Speech Technology Research, University of Edinburgh, UK

{ahewer|steiner|bolkart}@mmci.uni-saarland.de, stefanie.wuhler@inria.fr, korin@cstr.ed.ac.uk

ABSTRACT

We describe a minimally-supervised method for computing a statistical shape space model of the palate surface. The model is created from a corpus of volumetric magnetic resonance imaging (MRI) scans collected from 12 speakers. We extract a 3D mesh of the palate from each speaker, then train the model using principal component analysis (PCA). The palate model is then tested using 3D MRI from another corpus and evaluated using a high-resolution optical scan. We find that the error is low even when only a handful of measured coordinates are available. In both cases, our approach yields promising results. It can be applied to extract the palate shape from MRI data, and could be useful to other analysis modalities, such as electromagnetic articulography (EMA) and ultrasound tongue imaging (UTI).

Keywords: vocal tract MRI, principal component analysis, palate model

1. INTRODUCTION

The palate plays an important role in articulation; as part of the vocal tract walls it contributes to vowel production, and it is critical for the production of obstruents such as /z/, /ʃ/, or /j/, and for palatalization [11]. Therefore, analyzing its shape and understanding its interaction with other articulators is of great interest in speech science. A shape model of the palate could also contribute to acoustic models of the vocal tract.

Direct measurements of the palate shape are however a challenging task. Nowadays, magnetic resonance imaging (MRI) is the modality of choice for imaging the human vocal tract. This technique is able to provide dense 3D information about the inside of a speaker’s mouth without being hazardous

or invasive. The acquired data, however, has to be further processed to obtain the desired palate shape. In particular, a high-level structured shape representation is desirable, such as a polygonal mesh.

A model of the palate surface can be directly used in various fields of application. For example, for automatic image segmentation of MRI data, it can be used as a prior. It could also provide a persistent landmark for analysis with spatially sparse modalities, such as electromagnetic articulography (EMA) or ultrasound tongue imaging (UTI). Moreover, a palate mesh could be integrated to derive the vocal tract area function for acoustic modeling.

1.1. Related work

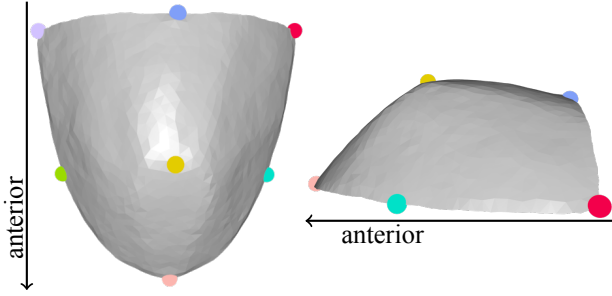
Analyzing the shape of the palate is an active field of research.

Yunusova et al. [16] used a thin plate spline (TPS) technique to estimate the contour of the palate in a palate trace acquired by EMA. TPS is a data-driven method that tries to deform a thin plate such that it passes through a set of control points. Additionally, the resulting plate should have some degree of smoothness. In their work, Yunusova et al. found that the weight for the smoothness constraint had an impact on the result: using values that were too small resulted in an overfitting to the sample points of the palate trace, whereas too large a value prevented the plate from deforming at all. In their experiments, they derived an optimal value for this weight empirically. As the method is purely data-driven, it might produce undesirable results if the data is too sparse.

Lammert et al. [12] used realtime MRI to investigate the morphological variation of the palate and the posterior pharyngeal wall. They extracted the shape information from mid-sagittal slices of the vocal tract. Afterwards, they applied a principal component analysis (PCA) to the obtained data to extract the principal modes of variation of both structures.

This study uses data from work supported by EPSRC Healthcare Partnerships Grant number EP/I027696/1 (“Ultrax”).

Figure 1: Palate mesh with landmark vertices shown as colored spheres. **Left:** View from the top. **Right:** Side view.



In their study, they found that the obtained principal modes could actually be related to anatomical variation, such as the degree of concavity of the palate. However, this study was restricted to the 2D case.

1.2. Our contribution

In this work, we present a minimally-supervised method for training a statistical model of the 3D shape of the palate surface. Our approach consists of two steps. We first extract the full shape of the palate surface from static 3D MRI scans of different speakers, where we use a polygonal mesh as the shape representation. Afterwards, we apply a PCA to this data in order to train the model. As the whole process is minimally-supervised, it is relatively easy to include additional MRI scans to improve the coverage of the model.

Such a statistical model can be helpful: for example, it could be useful for investigating the anatomical variation of the 3D shape of the human palate. Moreover, it represents a shape space that is able to generate new palate shapes and evaluate the probability of a specific shape. This property can be used for detecting and reconstructing the shape of the palate from data that is incomplete or very sparse, such as EMA or UTI.

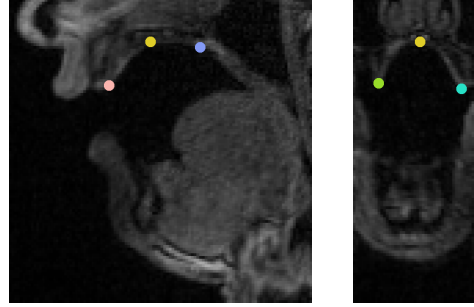
2. METHODS

Before outlining our approach, we first want to give a definition of the polygonal mesh $M := (V, F)$ used as a shape representation. $V := \{\vec{v}_i\}$ with $\vec{v}_i \in \mathbb{R}^3$ is called the vertex set of the mesh and F its face set. A face $f \in F$ is a set of vertices that form a surface patch in the form of a polygon, e.g., a triangle, if linked by edges. Stitching all faces together results in the full surface. An example mesh can be seen in Figure 1.

2.1. Shape extraction

In the first step of our approach, we focus on extracting 3D meshes representing the palate surface from

Figure 2: Landmarks selected on an MRI scan. **Left:** Sagittal slice. **Right:** Coronal slice.



volumetric vocal tract MRI scans collected from a number of different speakers (cf. Section 3, below). Here, we are using the minimally-supervised method of [10] that can be summarized as follows:

2.1.1. Surface point extraction

First, the region belonging to tissue is identified using an automatic image segmentation technique. In particular, each scan is interpreted as a 3D image with gray values in the interval $[0, 255]$. In our case, we chose to use a basic thresholding method to identify the tissue: some tissue like the palate surface appears much brighter than material with lower hydrogen density, such as air or bone. Thus, we automatically classify each point with a brightness higher than some threshold value as tissue.

The method then proceeds by extracting the surface points of the identified tissue regions, which produces a point cloud $P := \{\vec{p}_i\}$ with $\vec{p}_i \in \mathbb{R}^3$.

2.1.2. Template fitting

Then, a template fitting technique is applied to align a provided template mesh to the obtained point cloud. Two manual components were required for this step, viz. (a) a template mesh for the palate surface created beforehand from a single 3D MRI scan, using a medical imaging software [14]; and (b) a set of 7 manually selected vertices used as landmarks in a rigid alignment initialization step (cf. Section 2.4.1). These are required to identify the correct subset of points representing the palate and to deform the template mesh accordingly.

The palate mesh extraction step results in a collection of training meshes $M_i = (V_i, F)$ with $i \in [1, n]$, where n is the number of speakers. We remark that the meshes may differ in the position of their vertices, i.e., $\vec{v}_k \in V_i \neq \vec{v}_k \in V_j$ for $i \neq j$. Their faces, however, still consist of the same vertices which differ only in their position.

2.2. Training the model

In order to train our statistical model, we have to ensure that the palate meshes only differ from each other in their shape. To this end, we first apply a Procrustes alignment [9] to the collection of all extracted meshes. This serves to remove differences in their location, orientation, and scale, which enables us to analyze the features of their shape.

Next, we convert the transformed meshes to feature vectors such that the coordinates belonging to each vertex are located in consecutive rows. Finally, we apply a PCA to these vectors. Such methods are often used in literature to analyze data, cf., e.g., [2, 5, 6, 7]. This provides us with the set of principal directions $\vec{e}_i \in \mathbb{R}^k$ with $k = 3|V|$ of the training data. Interpreting these vectors as a basis gives us access to a space of palate shapes. Afterwards, we project the training data into this shape space and learn its probability distribution by fitting a multivariate Gaussian [8]. Thereby we obtain the variances $\lambda_i \in \mathbb{R}$ and means $m_i \in \mathbb{R}$ along the associated principal directions \vec{e}_i of our training data. Thus, we can also measure the probability of a specific shape in our learned shape space.

2.3. Generating palate shapes

The trained model can be used as follows to generate a new palate mesh $M^* = (V^*, F^*)$: first, we generate a vector \vec{x} representing a palate shape by computing

$$(1) \quad \vec{x} = \sum_i ((m_i + c_i)\vec{e}_i) = \vec{m} + \sum_i (c_i\vec{e}_i)$$

where \vec{m} is the mean of our training data and $c_i \in \mathbb{R}$ is the provided coefficient for the principal direction \vec{e}_i . Then, we convert \vec{x} to a vertex set V^* and assign $F^* = F$, the face set of our template mesh.

2.4. Using the model to register new data

In order to use our model to reconstruct a palate from a point cloud, we perform the following steps:

2.4.1. Rigid alignment

First, we have to find the optimal scale and location in the point cloud for the mesh generated by the mean of our model. This step is necessary because our shape space is not able to produce rigid transformations like translations or rotations. We use the following approach to facilitate this process: on the mesh, we selected 7 vertices as landmarks, as shown in Figure 1. Here, we see that we used three landmarks along the mid-sagittal line of the palate: one at the incisors, one at the hard/soft palate boundary, and another at the point of greatest curvature. In order to

add lateral information, we used the latter two landmarks as the anchor for two additional landmarks at either side of the palate.

Afterwards, we find the points in the data corresponding to these landmarks. If the used cloud originates from an MRI scan, this scan can be used to derive the coordinates like in Figure 2. Here, it is evident that the landmark locations are relatively easy to identify for a user. The scale and position of the mesh are then determined by finding the best rigid transformation that maps the user-provided coordinates to the landmarks on the mesh. Additionally, an iterative closest point (ICP) approach [4] was applied to further improve this rigid alignment.

2.4.2. Fitting the model

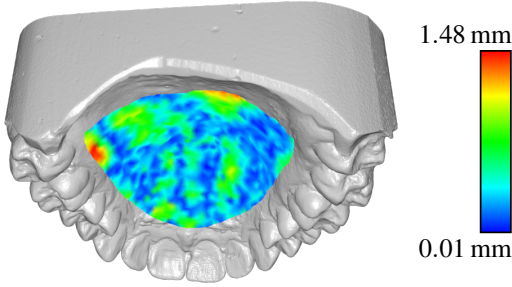
In the final step, we find the coefficients c_i for the principal directions of our model such that the resulting mesh is near the data in the provided point cloud. However, we limit the values for the coefficient c_i to the interval $[-\sqrt{\lambda_i}, \sqrt{\lambda_i}]$. This means we only consider values with a distance of no more than $\sqrt{\lambda_i}$ from the corresponding mean m_i of the coefficient in the training data, which serves to avoid unlikely palate shapes and prevent overfitting. In order to find these coefficients, we minimize an energy where we use a quasi-Newton scheme [13] to find a minimizer.

We note that this approach to fitting the model to the data is more robust than applying a template fitting technique directly. The model contains a whole space of palate shapes, whereas a template only represents a single shape. In contrast to a template mesh, it also allows to evaluate the probability of a generated shape. Furthermore, a template fitting offers many more degrees of freedom to align the template to the data, which means that it would also be possible for a palate mesh to be deformed into an implausible shape.

3. DATASETS

We used scans from two datasets for training our model: the full dataset of the Ultrax project [1] and that of Adam Baker [3]. Both were recorded using a Siemens MAGNETOM Verio at the Clinical Research Imaging Centre in Edinburgh for the purpose of observing the vocal tract configuration for different phones. The Baker dataset consists of static 3D MRI scans of a single male speaker. It was recorded as part of the Ultrax project, but released separately. The Ultrax dataset itself contains static 3D MRI scans of 11 adult speakers where seven are female and four are male. Each considered scan con-

Figure 3: Colored fitted palate mesh of first experiment and maxillar dental cast. Color indicates distance to nearest point on dental cast.



sists of 44 sagittal slices with a thickness of 1.2 mm and size (whole head) of 320×240 pixels with a voxel size of $1.1875 \times 1.1875 \times 1.2 \text{ mm}^3$.

To evaluate our approach, we used the volumetric MRI subset of the mngu0 corpus [15], which contains data from one male speaker, including high-resolution 3D scans of a plaster cast of his teeth and palate. Here, each MRI scan consists of 26 sagittal slices of 4 mm thickness. The size of each slice is given by 256×256 pixels with a corresponding voxel size of $1.1 \times 1.1 \times 4 \text{ mm}^3$. We see that compared to the Ultrax data, these scans offer a lower spatial resolution along the sagittal dimension.

4. EXPERIMENTS

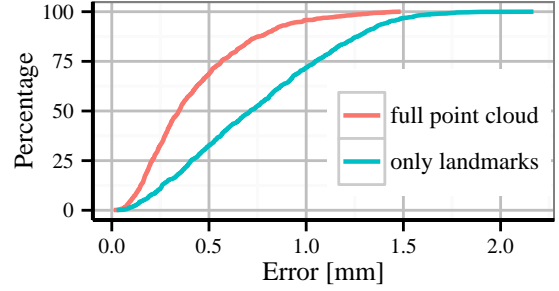
For the training, we used all twelve speakers of the Baker and Ultrax datasets. We selected for each speaker a scan where the palate was clearly visible, with no lingual contact. We then cropped each scan to a region of interest containing only the palate in order to reduce the memory requirements for the point cloud. Afterwards, we applied the methods described in Section 2.1 to extract the palate shapes. Here, we used the value $t = 25$ for the thresholding parameter to perform the image segmentation. All extracted meshes were then used to train the model.

4.1. Experiment setup

In the first experiment, we wanted to investigate if our model could handle data of a speaker it was not trained with. To this end, we selected the /v/ scan of the volumetric data of mngu0. We prepared the data as follows: the scan was once again cropped to a region containing only the palate. We then extracted the surface points of the tissue where we used the threshold $t = 25$ in the image segmentation step. Additionally, we distributed the landmarks needed for the rigid alignment by using the cropped scan. Finally, we fitted our trained model to the obtained point cloud.

Afterwards, we analyzed in a second experiment how our model behaves if *only* the 7 landmarks cho-

Figure 4: Cumulative error functions for the two experiments.



sen in the first experiment are used for the fitting.

4.2. Evaluation

We used the maxillary dental cast of the speaker as the reference solution for the shape of the hard palate. However, the obtained palate meshes differed from the dental cast in their location and orientation. Therefore, we again used landmarks and an ICP technique to perform a rigid alignment to remove these differences. This time, no scaling was applied in order to preserve the original shape. We then measured for each vertex of the palate mesh the distance to the closest point on the dental cast. A heat map visualizing these distances for the mesh of the first experiment can be seen in Figure 3. Afterwards, we interpreted this distance as an error measure and computed the cumulative error function. In Figure 4, we see that in the first experiment nearly 75 % of the error are below 0.5 mm.

For the second experiment using only the seven landmark points, nearly 75 % of the error are below 1 mm, which indicates that our model can produce acceptable results even with only very sparse information.

5. CONCLUSION

In this work, we described a minimally-supervised method for training a statistical shape model of the palate surface. We saw that a model trained with the palate shapes of twelve speakers was already useful. In particular, we found that it could be used to extract palate information from an MRI scan of a new speaker. Furthermore, even when only a handful of points are used, our model can be fit with acceptable precision, which allows us to use sparse input data, such as from an EMA palate trace.

Further experiments are scheduled to obtain reference data from more speakers, using an intraoral scanner, such as a 3shape TRIOS. Moreover, we plan to investigate how the trained model can be used to reconstruct palate information from existing EMA data of these speakers. Moreover, we plan to acquire more MRI data to increase our training set.

REFERENCES

- [1] 2014. Ultrax: Real-time tongue tracking for speech therapy using ultrasound. <http://www.ultrax-speech.org/>.
- [2] Allen, B., Curless, B., Popović, Z. 2003. The space of human body shapes: Reconstruction and parameterization from range scans. *ACM Transactions on Graphics* 22(3), 587–594.
- [3] Baker, A. 2011. A biomechanical tongue model for speech production based on MRI live speaker data. <http://www.adambaker.org/qmu.php>.
- [4] Besl, P. J., McKay, N. D. 1992. A method for registration of 3-D shapes. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 14(2), 239–256.
- [5] Brunton, A., Salazar, A., Bolkart, T., Wuhler, S. 2014. Review of statistical shape spaces for 3D data with comparative analysis for human faces. *Computer Vision and Image Understanding* 128, 1 – 17.
- [6] Cootes, T. F., Taylor, C. J. 2001. Statistical models of appearance for medical image analysis and computer vision. *Proc. Medical Imaging 2001: Image Processing* 236–248.
- [7] Cootes, T. F., Taylor, C. J., Cooper, D. H., Graham, J. 1995. Active shape models – their training and application. *Computer Vision and Image Understanding* 61(1), 38–59.
- [8] Davies, R., Twining, C., Taylor, C. 2008. *Statistical Models of Shape: Optimisation and Evaluation*. Springer.
- [9] Dryden, I. L., Mardia, K. V. 1998. *Statistical shape analysis*. Wiley.
- [10] Hewer, A., Steiner, I., Wuhler, S. 2014. A hybrid approach to 3D tongue modeling from vocal tract MRI using unsupervised image segmentation and mesh deformation. *Proc. Interspeech* 418–421.
- [11] Hiki, S., Itoh, H. 1986. Influence of palate shape on lingual articulation. *Speech Communication* 5(2), 141–158.
- [12] Lammert, A., Proctor, M., Narayanan, S. 2013. Morphological variation in the adult hard palate and posterior pharyngeal wall. *Journal of Speech, Language, and Hearing Research* 56(2), 521–530.
- [13] Liu, D. C., Nocedal, J. 1989. On the limited memory BFGS method for large scale optimization. *Mathematical Programming* 45(1-3), 503–528.
- [14] Rosset, A., Spadola, L., Ratib, O. 2004. OsiriX: an open-source software for navigating in multidimensional DICOM images. *Journal of Digital Imaging* 17(3), 205–216.
- [15] Steiner, I., Richmond, K., Marshall, I., Gray, C. D. 2012. The magnetic resonance imaging subset of the mngu0 articulatory corpus. *Journal of the Acoustical Society of America* 131(2), EL106–EL111.
- [16] Yunusova, Y., Baljko, M., Pintilie, G., Rudy, K., Faloutsos, P., Daskalogiannakis, J. 2012. Acquisition of the 3D surface of the palate by in-vivo digitization with Wave. *Speech Communication* 54(8), 923–931.