Vocal Tract Model Adaptation Using Magnetic Resonance Imaging

Peter Birkholz¹*, Bernd J. Kröger²

¹Institute for Computer Science, University of Rostock, 18051 Rostock, Germany

²Department of Phoniatrics, Pedaudiology, and Communication Disorders University Hospital Aachen (UKA) and Aachen University (RWTH) 52074 Aachen, Germany

piet@informatik.uni-rostock.de, bkroeger@ukaachen.de

Abstract. We present the adaptation of the anatomy and articulation of a 3D vocal tract model to a new speaker using magnetic resonance imaging (MRI). We combined two MRI corpora of the speaker: A corpus of volumetric images of sustained phonemes and a corpus with midsagittal image sequences of dynamic utterances. The volumetric MRI corpus was used for the adaptation of vocalic and (neutral) consonantal target shapes. For each phoneme, the vocal tract parameters were adjusted manually for a close visual match of the MRI-tracings and the model-derived outlines. The resulting acoustic match of the vowels in terms of formant differences was examined and optimized. The dynamic MRI corpus was used to replicate the coarticulation of the speaker. Therefore, we analyzed the MRI tracings of the consonants articulated in the contexts of the vowels /a:/, /i:/, and /u:/. The articulatory differences of the consonants due to the different contexts were translated into a dominance model used to control the simulated vocal tract.

1. Introduction

In the last few years, we have been developing an articulatory speech synthesizer based on a geometric 3D model of the vocal tract (Birkholz, 2005; Birkholz et al., 2006). Our goals are high quality text-to-speech synthesis as well as the application of the synthesizer in a neural model of speech production (Kröger et al., 2006). Till now, the anatomy and articulation of our vocal tract model were based on x-ray tracings of sustained phonemes of a Russian speaker (Koneczna and Zawadowski, 1956). However, these data were not sufficient to reproduce the speakers anatomy and articulation very accurately. They neither provided information about the lateral vocal tract dimensions nor on coarticulation of phonemes. These information had to be estimated and impeded a strict evaluation of the synthesizer.

In this study, we started to close this gap by adapting the anatomy and articulation of our vocal tract model to a new speaker using MRI. Two MRI corpora were available to us: one corpus of volumetric images of sustained vowels and consonants, and one corpus

^{*}Supported by the German Research Foundation.

of dynamic midsagittal MRI sequences. Additionally, we had computer tomography (CT) scans of oral-dental impressions. The CT scans were used to adapt the geometry of the maxilla, the jaw, and the teeth. The articulatory targets for vowels and consonants were determined by means of the volumetric MRI data. The dynamic MRI corpus was used to determine the dominance of the consonants over the individual articulators. This is the basis of the dominance model used for the simulation of coarticulation in our synthesizer.

Section 2 will discuss the analysis and normalization of the images from both corpora, and Sec. 3 introduces the vocal tract model and describes the adaptation of vowels and consonants. Conclusions are drawn in Sec. 4.

2. Magnetic Resonance Image Processing

2.1. Corpora

We analyzed two MRI corpora of the same native German speaker (JD, male, 35 years) that were available to us from other studies (Kröger et al., 2000, 2004). The first corpus contains volumetric images of sustained phonemes including tense and lax vowels, nasals, voiceless fricatives, and the lateral /l/. Each volumetric image consists of 18 sagittal slices with 512 x 512 pixels. The pixel size is 0.59 x 0.59 mm² and the slice thickness is 3.5 mm.

The second corpus contains dynamic MRI sequences of midsagittal slices scanned at a rate of 8 frames/second with a frame resolution of 256 x 256 pixels. The pixel size is $1.18 \text{ x } 1.18 \text{ mm}^2$. The recorded utterances consist of multiple repetitions of the sequences /a:*C*a:/, /i:*C*i:/ and /u:*C*u:/ for nearly all German consonants *C*.

In addition to these two corpora, we had high resolution CT scans of plaster casts of the upper and lower jaws and teeth of the speaker with a voxel size of $0.226 \times 1 \times 0.226 \text{ mm}^3$.

2.2. Outline Tracing

The midsagittal airway boundaries of all MR images were hand-traced on the computer for further processing. The manual tracing was facilitated by applying an edge detector (Sobel operator) to the images. Examples of MR images from corpora 1 and 2 are shown in Fig. 1 (a) and (d), respectively. Pictures (b) and (e) show the corresponding results of the Sobel edge detector, and the tracings are depicted in (c) and (f). For corpus 1 phonemes, we additionally traced the tongue outlines approximately 1 cm left from the midsagittal plane (dashed curve in Fig. 1 (c)).

In corpus 2, we were interested in the articulation of the consonants in the context of the vowels /a:/, /i:/, and /u:/. The analysis of the dynamic MRI sequences revealed, that the sampling rate of 8 frames/second was to low to capture a clear picture of each spoken phoneme. But among the multiple repetitions that we had of each spoken /VCV/sequence, we identified for each consonant+context at least 2 (usually 4-5) candidate frames, where the consonantal targets were met with sufficient precision. One of these candidates was chosen as template for tracing the outlines. The chosen candidate frame was supposed to be the one that best represented the mean of the candidate set. Therefore, we chose the frame that had the smallest sum of "distances" to all other frames in that set. The distance between two images was defined as the signal energy of the difference



Figure 1. (a) Original image of corpus 1. (b) Edges detected by the Sobel operator for (a). (c) Tracing result for (b). (d)-(f) Same as (a)-(c) for an image of corpus 2.

image. The volumetric CT images of the plaster casts of the upper and lower jaw were exactly measured both in the lateral and coronal planes to allow a precise reconstruction of these rigid parts in the vocal tract model.

2.3. Contour Normalization

The comparison of Fig. 1 (c) and (f) shows that the head was not held in exactly the same way in both corpora. In corpus 1, the neck is usually more "stretched" than in corpus 2, resulting in a greater angle between the rear pharyngeal wall and the horizontal¹. Small variations of this angle also exist within each corpus. For the vocal tract adaptation it was essential to equalize/normalize these differences in head postures.

Our basic assumption for the normalization was, that there is a fixed point R (with respect to the maxilla) in the region of the soft palate, around which the rear pharyngeal outline rotates when the head is raised or lowered. Given this assumption, the straight lines approximating the rear pharyngeal outlines of all tracings should intersect in R. Therefore, R was determined solving the minimization problem

$$\sum_{i=1}^{N} d^2(R, l_i) \to min,$$

where N is the total number of traced images from both corpora, and $d(R, l_i)$ denotes the shortest distance from R to the straight line l_i that approximates the rear pharyngeal wall of the *i*th image. Each MRI-tracing was then warped such that its rear pharyngeal outline was oriented at a predefined constant angle. Warping was implemented using the

¹All tracings were rotated such that the upper row of teeth was oriented horizontally.



Figure 2. Warping of the MRI-tracing of the consonant /b/ in /ubu/.



Figure 3. (a) 3D-rendering of the vocal tract model. (b) Vocal tract parameters.

method by Beier and Neely (1992) with 3 corresponding pairs of vectors as exemplified in Fig. 2. The horizontal vectors on top of the palate and the vertical vectors at the chin are identical for the original and the warped image, keeping these parts of the vocal tract equal during warping. Only the vectors pointing down the pharyngeal outline make the vocal tract geometry change in the posterior part of the vocal tract. Both of these vectors only differ in the degree of rotation around R. Figure 2 (b) shows the MRI-tracing in (a) before warping (dotted curve) and after warping (solid curve). This method proofed to be very effective and was applied to all MRI-tracings.

3. Adaptation

3.1. Vocal Tract Model

Our vocal tract model consists of different triangle meshes that define the surfaces of the tongue, the lips and the vocal tract walls. A 3D rendering of the model is shown in Fig. 3 (a) for the vowel /a:/. The shape of the surfaces depends on a number of predefined parameters. Most of them are shown in the midsagittal section of the model in Fig. 3 (b). The model has 2 parameters for the position of the hyoid (HX,HY), 1 for the velic aperture (VA), 2 for the protrusion and opening of the lips (LP, LH), 3 for the position and rotation of the jaw (JX, JY, JA) and 7 for the midsagittal tongue outline (TRE, TCX, TCY, TBX,



Figure 4. MRI outlines (dotted curves) and the matched model-derived outlines (solid curves) for the vowels /a:/, /i:/, and /u:/.

TBY, TTX, TTY). Four additional parameters define the height of the tongue sides with respect to the midsagittal outline at the tongue root, the tongue tip, and two intermediate positions. A detailed description of the parameters is given in (Birkholz, 2005; Birkholz et al., 2006). The current version of the model is an extension of the model in the cited references. On one hand, we added the epiglottis and the uvula to the model, which were previously omitted. Furthermore, the 3D-shape of the palate, the mandible, the teeth, the pharynx and the larynx was adapted to the (normalized) MR images.

3.2. Vowels

To reproduce the vowels in corpus 1, the vocal tract parameters were manually adjusted aiming for a close match between the normalized MRI tracings and the model-derived outlines. Furthermore, the tongue side parameters were adjusted for a close match of the tongue side outlines. Figure 4 shows our results for the vowels /a:/, /i:/, and /u:/. The midsagittal model outlines are drawn with solid lines and the tongue sides with dashed lines. The corresponding MRI tracings are drawn with dotted lines. In the case of all examined vowels, we achieved a fairly good *visual* match.

The *acoustic* match between the original and synthetic vowels was tested by comparison of the first 3 formant frequencies. The formants of the natural vowels were determined by standard LPC analysis. The audio corpus was recorded independently from the MRI scans with the speaker in a supine position repeating all vowels embedded in a carrier sentence four times. For each formant frequency of each vowel, the mean value was calculated from the 4 repetitions.

The formant frequencies of the synthetic vowels were determined by means of a frequency-domain simulation of the vocal tract system based on the transmission-line circuit analogy (Birkholz, 2005). The area functions for these simulations were calculated from the 3D vocal tract model. The nasal port was assumed to be closed for all vowels. In all acoustic simulations, we considered losses due to yielding walls, viscous friction, and radiation. The *piriform fossa* side cavity was included in the simulations and modeled after Dang and Honda (1997).

The test results are summarized in Fig. 5 for the first two formants of the tense German vowels. The error between the natural and synthetic formant frequencies averaged over the first three formants of all tense vowels was 12.21%. This error must be mainly attributed to the resolution-limited accuracy of the MRI tracings as well as to the



Figure 5. Formant frequencies for the German tense vowels.

imperfect matching of the outlines. It is well known that in certain regions of the vocal tract, the formant frequencies are quite sensitive to small variations of articulatory parameters (Stevens, 1989). Therefore, the acoustic differences could be caused by only small articulatory deviations due to the above sources of errors. To test how far small corrective variations of the vocal tract parameters can improve the acoustic match, we implemented an algorithm searching the parameter space to minimize the formant errors. Each vocal tract parameter was allowed to deviate *maximally* 5% of its whole range from the value that was determined visually. Figure 5 shows that the formants were much closer to their "targets" after this optimization, while the articulation changed only slightly. The average formant error reduced to 3.41%.

3.3. Consonants

To a certain extend, the articulatory realization of a consonant depends on the vocalic context due to coarticulation. In our synthesizer, we use a dominance model to simulate this effect (Birkholz et al., 2006). The basic idea is that each consonant has a "neutral" target shape (just like the vowels), but in addition, each parameter has a weight between 0 and 1, expressing its "importance" for the realization of the consonantal constriction. For /d/, for example, the tongue tip parameters have a high weight, because the alveolar closure with the tongue tip is essential for /d/. Most of the other parameters/articulators are less important for /d/ and have a lower weight. The other way round, the weight expresses how strong a consonantal parameter is influenced by the context vowels (low weight = strong influencing). Formally, this concept is expressed by

$$x_{c|v}[i] = x_v[i] + w_c[i] \cdot (x_c[i] - x_v[i]),$$
(1)



Figure 6. Articulatory realization of the voiced plosives in the context of the vowels /a:/, /i:/, and /u:/. MRI tracings are drawn as dotted curves and model-derived outlines as solid curves.

where *i* is the parameter index, $x_{c|v}[i]$ is the value of parameter *i* at the moment of the maximal closure/constriction of the consonant *c* in the context of the vowel *v*, $w_c[i]$ is the weight for parameter *i*, and $x_c[i]$ and $x_v[i]$ are the parameter values of the targets for the consonant and vowel.

Hence, the needed data for the complete articulatory description of a consonant c are the $x_c[i]$ and $w_c[i]$. The parameters for the "neutral" consonantal targets were adjusted analogous to the vowel parameters in Sec. 3.2 using the high resolution MRI data from corpus 1. The weights were determined using the selected MRI tracings from corpus 2, that show the realization of the consonants in symmetric context of the vowels /a:/, /i:/, and /u:/. The vocal tract parameters for these coarticulated consonants were manually adjusted, too. Let us denote these parameters by $x_{c|v_j}$, where $v_j \in \{/a:/, /i:/, /u:/\}$. The optimal weights $w_c[i]$ were determined solving the minimization problem

$$\sum_{j=1}^{N} \left[x_{c|v_j}[i] - x_{v_j}[i] - w_c[i] \cdot (x_c[i] - x_{v_j}[i]) \right]^2 \to min,$$

where N = 3 is the number of context vowels.

Figure 6 contrasts the model-derived outlines of coarticulated consonants using Eq. (1) (solid curves) and the corresponding MRI tracings (dotted curves). Despite some obvious differences in the outlines (especially in the laryngeal region), the basic coarticulatory effects are well reproduced in all examples and are expected to be sufficient for high-quality articulatory speech synthesis.

4. Conclusions

We have presented the anatomic and articulatory adaptation of a 3D vocal tract model to a specific speaker combining data from higher resolution volumetric MRI data and lower resolution dynamic MRI data. We achieved a satisfying visual and acoustic match between the original speaker and the model. The methods proposed in this study can be considered as simple but powerful means for future adaptations to other speakers, provided that the corresponding MRI data are available.

PS: During the presentation of our work, Jean-Luc Boë (ICP, Grenoble) pointed out that the glottis in our vocal tract model is situated at a to high position, which we could attribute to tracing errors in the laryngeal region. Retracing of these regions revealed that the larynx tube is approximately 1 cm longer than shown in this paper. The formant frequencies of the corrected model are actually slightly closer to the measured formants of our subject presented in Sec. 3.2.

References

- Beier, T. and Neely, S. Feature-based image metamorphosis. *Computer Graphics*, 26(5):35–42, 1992.
- Birkholz, P. 3D-Artikulatorische Sprachsynthese. Logos Verlag Berlin, 2005.
- Birkholz, P., Jackèl, D., and Kröger, B. J. Construction and control of a three-dimensional vocal tract model. In *International Conference on Acoustics, Speech, and Signal Processing* (*ICASSP'06*), pages 873–876, Toulouse, France, 2006.
- Dang, J. and Honda, K. Acoustic characteristics of the piriform fossa in models and humans. *Journal of the Acoustical Society of America*, 101(1):456–465, 1997.
- Koneczna, H. and Zawadowski, W. *Obrazy Rentgenograficzne Glosek Rosyjskich*. Panstwowe Wydawnictwo Naukowe, 1956.
- Kröger, B. J., Birkholz, P., Kannampuzha, J., and Neuschaefer-Rube, C. Spatial-to-joint coordinate mapping in a neural model of speech production. In 32. Deutsche Jahrestagung für Akustik (DAGA '06), pages 561–562, Braunschweig, Germany, 2006.
- Kröger, B. J., Hoole, P., Sader, R., Geng, C., Pompino-Marschall, B., and Neuschaefer-Rube, C. MRT-Sequenzen als Datenbasis eines visuellen Artikulationsmodells. *HNO*, 52:837–843, 2004.
- Kröger, B. J., Winkler, R., Mooshammer, C., and Pompino-Marschall, B. Estimation of vocal tract area function from magnetic resonance imaging: Preliminary results. In 5th Seminar on Speech Production: Models and Data, pages 333–336, Kloster Seeon, Bavaria, 2000.

Stevens, K. N. On the quantal nature of speech. Journal of Phonetics, 17:3-45, 1989.