

Spectral Study of the Vocal Tract in Vowel Synthesis: A Comparison between 1D and 3D Acoustic Analysis

Negar M. Harandi^{*}, Daniel Aalto[°], Antti Hannukainen[†], Jarmo Malinen[†], Sidney Fels^{*}

^{*}University of British Columbia, Canada. [°]University of Alberta, Canada. [†]Aalto University, Finland.

Abstract

A state-of-the-art 1D acoustic synthesizer has been previously developed, and coupled to speaker-specific biomechanical models of oropharynx in ArtiSynth. As expected, the formant frequencies of the synthesized vowel sounds were shown to be different from those of the recorded audio. Such discrepancy was hypothesized to be due to the simplified geometry of the vocal tract model as well as the one dimensional implementation of Navier-Stokes equations. In this paper, we calculate Helmholtz resonances of our vocal tract geometries using 3D finite element method (FEM), and compare them with the formant frequencies obtained from the 1D method and audio. We hope such comparison helps with clarifying the limitations of our current models and/or speech synthesizer.

1 Introduction

Articulatory speech synthesizers generate sound based on the shape of the vocal tract. Vibration of the vocal folds under the expiratory air flow is the source in the system; and the vocal tract, consisting of the larynx, pharynx, oral and nasal cavities, constitutes a filter where sound frequencies are shaped. This creates a number of resonant peaks in the spectrum, known as formants. The first and second formants (F_1 and F_2) are used to distinguish the vowel phonemes, where the value of F_1 and F_2 is controlled by the height and backness-frontness of the tongue body respectively.

Traditionally, the acoustic system is approximated by a one-dimensional wave equation that associates the slow varying cross-sectional area of a rigid tube to the pressure wave for a low-frequency sound. However, complex shape of the vocal tract, with its side branches and asymmetry, has motivated higher dimensional acoustic analysis. The 3D analysis methods were shown to produce a better representation of the sound spectrum at the price of higher computational cost. However, some studies suggested that the spectrum yielded by 1D acoustic analysis matches closely that of the 3D analysis for frequencies less than 7KHz (Takemoto et al., 2014; Arnela and Gausch, 2014). Aalto et al. (2012) suggested that the discrepancy between the resonance frequencies computed by 3D analysis

of the vocal tract and the formant frequencies of the recorded audio is a result of insufficient boundary conditions in the wave equation especially in case of the open lips and/or velar port.

In this paper, we follow Aalto et al. (2014) in calculating the Helmholtz resonances of our vocal tract geometries using 3D FEM analysis. The resonances are then compared to the formant frequencies obtained from the 1D acoustic synthesizer proposed by Doel and Ascher (2008) and those of the recorded audio.

2 Material and Methods

We use static MRI images acquired with a Siemens Magnetom Avanto 1.5 T scanner. A 12-element Head Matrix Coil, and a 4-element Neck Matrix Coil, allow for the Generalize Auto-calibrating Partially Parallel Acquisition (GRAPPA) acceleration technique. One speaker, a 26-year-old male, was imaged while he uttered four sustained Finnish vowels. The MRI data covers the vocal and nasal tracts, from the lips and nostrils to the beginning of the trachea, in 44 sagittal slices, with an in-plane resolution of 1.9mm. Figure 1 shows the VT surface geometries extracted from MRI data using an automatized segmentation method (Aalto et al., 2013).

For our 1D acoustic analysis, we describe the vocal tract by an area function $A(x, t)$ where $0 \leq x \leq L$ is the distance from the glottis on the tube axis and t denotes the time. We take the similar notion of Doel and Ascher (2008) in defining the variables $u(x, t) = A(x, t)\hat{u}/c$ and $p(x, t) = \hat{p}/\rho_0 - 1$ as the scaled versions of volume-velocity \hat{u} and air density \hat{p} respectively. ρ_0 is the mass density of the air and c is the speed of sound. We solve for $u(x, t)$ and $p(x, t)$ in the tube using derivations of the linearised Navier-Stokes equation (1a) and the equation of continuity (1b) subject to the

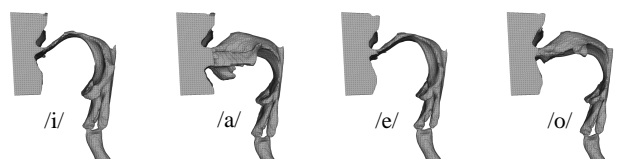


Figure 1: VT geometries extracted from MRI data (Aalto et al., 2013).

boundary conditions described in equation 1c:

$$\frac{\partial(u/A)}{\partial t} + c \frac{\partial p}{\partial x} = -d(A)u + D(A) \frac{\partial^2 u}{\partial x^2} \quad (1a)$$

$$\frac{\partial(Ap)}{\partial t} + c \frac{\partial u}{\partial x} = -\frac{\partial A}{\partial t} \quad (1b)$$

$$u(0, t) = u_g(t), \quad p(L, t) = 0 \quad (1c)$$

where $d(A) = d_0 A^{-3/2}$ and $D(A) = D_0 A^{-3/2}$ with the wall loss coefficient $d_0 = 1.6 \text{ m s}^{-1}$ and $D_0 = 0.002 \text{ m}^3 \text{ s}^{-1}$; and $u_g(t)$ is the source volume velocity at the glottis. We couple the vocal tract to a two-mass glottal model (Ishizaka and Flanagan, 1972) and solve equation 1 in the frequency domain using a digital ladder filter defined based on the cross-sectional areas of 20 segments of the vocal tract. We refer to Doel and Ascher (2008) for full details of the implementation.

For our 3D acoustic analysis, we calculate the vowel formants directly from the wave equation by finding the eigenvalues, λ , and their corresponding velocity potential eigenfunction, Φ_λ , from the Helmholtz resonance problem:

$$\lambda^2 \Phi_\lambda - c^2 \Delta \Phi_\lambda = 0 \quad \text{on } \Omega \quad (2a)$$

$$\Phi_\lambda = 0 \quad \text{on } \Gamma_1 \quad (2b)$$

$$\alpha \lambda \Phi_\lambda + \frac{\partial \Phi_\lambda}{\partial \nu} = 0 \quad \text{on } \Gamma_2 \quad (2c)$$

$$\lambda \Phi_\lambda + c \frac{\partial \Phi_\lambda}{\partial \nu} = 0 \quad \text{on } \Gamma_3 \quad (2d)$$

where $\Omega \in \mathbb{R}^3$ is the air column volume and $\partial\Omega$ is its surface including the boundary at mouth opening (Γ_1), at air-tissue interface (Γ_2) and at a virtual plane above glottis (Γ_3); and $\frac{\partial \Phi_\lambda}{\partial \nu}$ denotes the exterior normal derivative. The value of α regulates the energy dissipation through tissue walls, and the case $\alpha = 0$ corresponds with hard, reflecting boundaries. We calculate the numerical solution of equation 2 by Finite Element Method (FEM) using piecewise linear shape functions and approximately 10^5 tetrahedral elements. The imaginary parts of the first two smallest eigenvalues λ_1 and λ_2 give first two Helmholtz resonances of the vocal tract. We refer to Aalto et al. (2014) and Kivelä et al. (2013) for details of implementation.

In order to distinguish the effects of dimensionality (1D vs. 3D) from the effects of different boundary conditions in equations 1 and 2, we also compute the Webster resonances by interpreting equation 2 in one dimension:

$$\left(\frac{\lambda^2}{c^2} \frac{1}{\Sigma^2} + \lambda \frac{2\pi\alpha W}{A} \right) \Phi_\lambda = \frac{1}{A} \frac{\partial}{\partial s} \left(A \frac{\partial \Phi_\lambda}{\partial s} \right) \quad \text{on } [0, L] \quad (3a)$$

$$\lambda \Phi_\lambda - c \Phi'_\lambda = 0 \quad \text{at } s = 0 \quad (3b)$$

$$\Phi_\lambda = 0 \quad \text{at } s = L \quad (3c)$$

where Σ denotes the sound speed correction factor that depends on the curvature of the vocal tract; $A(x)$ is the area function and s is the implicit parameter to Φ_λ , A , W and Σ . We refer to Kivelä (2015) for details of implementation and parameter values.

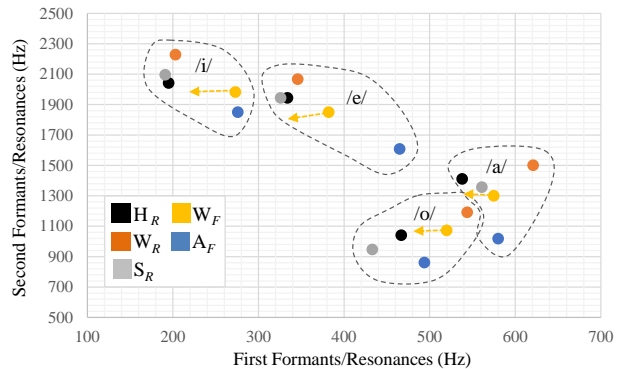


Figure 2: Simulation results for first and second formant/resonance frequencies for different vowels: Helmholtz resonances (H_R), Webster resonances (W_R) and their scaled version (S_R), Webster formants (W_F) and formants from audio signal (A_F).

3 Results and Discussion

Figure 2 shows the first two formant/resonance frequencies, computed for the four Finnish vowels. Webster formants (W_F) are calculated by solving Equation 1, as suggested by Doel and Ascher (2008). Helmholtz (H_R) and Webster resonances (W_R) are obtained from equations 2 and 3, respectively (Aalto et al., 2014). S_R denotes the scaled version of W_R . The figure also includes the formant frequencies (A_F) computed from audio signals recorded in an anechoic chamber (Aalto et al., 2014). The values are averaged over 10 repetitions of each vowel utterance.

As we can see in Figure 2, the resonance values (H_R , W_R and S_R) lie close together for vowels /i/ and /e/, with S_R being closer to H_R , as expected. For vowels /o/ and /a/ there is more difference in the first resonances of H_R and W_R ; For /o/, although S_R lies closer to H_R , its first resonance is surprisingly low. For all of the vowels in Figure 2, the second formant of the audio is less than the computed results. The vowel /i/ is expected to be very sensitive to glottal end position, which, in turn, suggests the significance of adequate MRI resolution and accurate geometry processing for its spectral analysis.

Interestingly, the Webster formants (W_F) remain closer to the audio formants (A_F) than any of the resonances in the case of /i/, /e/, and /a/. For /o/ the distance to the A_F is almost equal for W_F and H_R , with both having similar values for the second formant/resonance; however, the first H_R is lower, and the first W_F is higher, than the first A_F .

The time-domain Webster analysis (Doel and Ascher, 2008) accounts for the VT wall-vibration phenomenon that is missing in the resonance analysis. This is done by substituting $A(x, t)$, from equation 5.3, with $A(x, t) + C(x, t)y(x, t)$: where $C(x, t)$ is the slowly-varying circumference and $y(x, t)$ is the wall displacement governed by a damped mass-spring system. Setting $y(x, t)$ to zero, the Webster formants move along the arrows in Figure 2, reducing in their first formants. This moves the W_F closer to the H_R as both acousti-

cal models now ignore the wall vibration. Meanwhile, W_F moves away from the audio formants in the case of /i/, /e/, and /a/. The distance between W_R and W_F remains large, despite the fact that both acoustical models solve the Webster equation. The results imply that 3D Helmholtz analysis is more realistic than its 1D Webster version, as expected.

Overall, our experiments suggest that the time-domain interpretation of acoustic equations provides more realistic results – even if it requires reducing from 3D to 1D. This may be partially due to the fact that time-domain analysis allows for more complexity in the acoustical model such as inclusion of lip radiation and wall loss. Certainly unknown parameters always remain (such as those involved in glottal flow, coupling between fluid mechanics and acoustical analysis, etc.), which are estimated indirectly, based on observed behaviour in simulations.

It should be noted that our experiments are solely based on data from a single speaker. A larger database – inclusive of more speakers from different genders and languages – is needed in order to confirm the validity/generalizability of our findings.

References

- Aalto D, et al. 2014. Large scale data acquisition of simultaneous MRI and speech. *J Appl Acoust.* 83:64–75.
- Aalto D, et al. 2013. Algorithmic Surface Extraction from MRI Data-Modelling the Human Vocal Tract. *Proceeding of 6th International Joint Conference on Biomedical Engineering Systems and Technologies*; Barcelona, Spain.
- Aalto D, et al. 2012. How far are vowel formants from computed vocal tract resonances? *arXiv:1208.5963*.
- Arnela M, Guasch O. 2014. Three-dimensional behavior in the numerical generation of vowels using tuned two-dimensional vocal tracts. *Proceeding of 7th Forum Acousticum*; Kraków, Poland.
- Doel K van den, Ascher UM. 2008. Real-time numerical solution of Webster's equation on a non-uniform grid. *IEEE Trans Audio Speech Lang Processing* 16:1163–1172.
- Ishizaka K, Flanagan JL. 1972. Synthesis of voiced sounds from a two-mass model of the vocal cords. *J Bell Syst Tech.* 51: 1233–1268.
- Kivelä A. 2015. Acoustics of the vocal tract: MR image segmentation for modelling, Master's thesis, Aalto University School of Science.
- Kivelä A, Kuortti J, Malinen J. 2013. Resonances and mode shapes of the human vocal tract during vowel production. *Proceedings of 26th nordic seminar on computational mechanics*; Oslo, Norway.
- Takemoto H, Mokhtari P, Kitamura T. 2014. Comparison of vocal tract transfer functions calculated using one-dimensional and three-dimensional acoustic simulation methods. *Proceeding of 15th Annual Conference of the International Speech Communication Association*; Singapore, Singapore.