



A perceptually and physiologically motivated voice source model

Gang Chen¹, Marc Garellek^{2,3}, Jody Kreiman³, Bruce R. Gerratt³, Abeer Alwan¹

¹Department of Electrical Engineering, University of California, Los Angeles, USA

²Department of Linguistics, University of California, Los Angeles, USA

³Department of Head and Neck Surgery, School of Medicine, University of California, Los Angeles, USA

{gangchen, alwan}@ee.ucla.edu

{marcgarellek, jkreiman, bgerratt}@ucla.edu

Abstract

Many glottal source models have been proposed, but none has been systematically validated perceptually. Our previous work showed that model fitting of the negative peak of the flow derivative is the most important predictor of perceptual similarity to the target voice. In this study, a new voice source model is proposed to capture perceptually-important source shape aspects. This new model, along with four other source models, was fitted to 40 voice sources (20 male and 20 female) obtained by inverse filtering and analysis-by-synthesis (AbS) of samples of natural speech. We generated synthetic copies of the voices using each modeled source pulse, with all other synthesis parameters held constant, and then conducted a visual sort-and-rate task in which listeners assessed the extent of perceived similarity between the target voice samples and each copy. Results showed that the proposed model provided a more accurate fit and a better perceptual match to the target than did the other models.

Index Terms: voice source model, perceptual validation, analysis-by-synthesis, flow derivative

1. Introduction

According to the linear speech production model [1], speech signals are generated by filtering the voice source by the vocal tract transfer function. Modeling the glottal source has been an important topic for decades and has applications in many areas, such as speech coding and speech synthesis. Many source models have been proposed with varying levels of complexity, such as the Rosenberg [2], Liljencrants-Fant (LF) [3], Fujisaki-Ljungqvist (FL) [4], and Rosenberg++ (R++) [5] models (see [6] for review). With three parameters, the Rosenberg trigonometric model (denoted Ros) has two separate functions for the opening and closing phases to represent the glottal flow volume velocity [2]. The LF and FL models represent the first derivative of the glottal volume velocity pulse, which incorporates lip radiation effects. The four-parameter LF model [3] uses a combination of sinusoidal and exponential functions, and is commonly used in speech synthesis. With six parameters and polynomial functions, the FL model provides greater detail in modeling the glottal pulse shape, but the increased number of parameters also makes it more difficult to use in practice. The R++ model in [5] is computationally more efficient but perceptually equivalent when compared to the LF model. The four-parameter glottal flow model (denoted EE1 [7]) uses a combination of sinusoidal and exponential functions similar to the LF model, but with the ability to adjust the slopes of the opening and closing phases separately. The glottal flow model in [8] (denoted EE2) improves the EE1 model by redefining the model parameters (speed of opening and speed of closing) to allow for

lower computational complexity, faster waveform generation, and more accurate pulse shape manipulation. In that study, the EE2 model was used for automatic glottal flow estimation from acoustic speech signals, and glottal area waveforms extracted from high-speed endoscopic recordings of the laryngeal vibrations were converted to glottal flow in order to evaluate the performance of the glottal flow estimation algorithm.

Research efforts have also been devoted to studying the perceptual importance of changes in source waveform shapes. In [2], listening tests using a variety of glottal excitations showed that simulated excitations with a single slope discontinuity at closure were perceived as more natural-sounding, while very small opening or closing times (or opening times approximately equal to or less than closing times) were not preferred. In [9], the LF model and a turbulent noise generator were used to synthesize four voice quality types (modal, vocal fry, falsetto, and breathy). Perceptual experiments showed that these four voice quality types could be characterized by four parameters: pulse width, pulse skewness, the abruptness of glottal closure, and turbulent noise. In [10], nonmodal phonations were synthesized using a speech synthesizer in which the glottal characteristics were manipulated with quasi-articulatory parameters. In other approaches, voice source waveforms were parameterized to capture variations in voice quality [11, 12, 13, 14, 15, 16, 17], while those characteristics related to vocal intensity were investigated and parameterized in [18, 19, 20, 21]. Data-driven approaches, such as principal component analysis [22, 23] and Gaussian mixture modeling [24], have also been used to model source waveforms. In [25], the LF model was used to modify the glottal pulse shape for synthesis and transformation of singing voice.

Few studies have attempted to systematically validate glottal source models perceptually, and model development has focused more on replicating observed pulse shapes than on perceptual sufficiency. As a result, it is unclear which (if any) deviations from perfect fit between models and data have perceptual importance. In our previous study [26], the Ros, FL, LF, EE1, and EE2 source models were fitted to 40 natural normal and pathological voice sources (20 male and 20 female) obtained by inverse filtering and analysis-by-synthesis (AbS), subject to mean square error (MSE) criteria for which each point of the waveform was weighted equally. Evaluation of model fit at different parts of the source waveforms showed that the fit to the target pulses was worst at the negative peak of the flow derivative. Synthetic copies of the voices were then created using each modeled source pulse, while holding all other synthesizer parameters constant (including formant frequencies and bandwidths, fundamental frequency (F0) and amplitude contours, and spectral noise levels). These stimuli were compared to the

AbS target in a sort-and-rate listening test (described below). Across models and voices, the perceptual match between the target and synthetic tokens was best predicted by the match between the target and modeled stimuli at the negative peak of the flow derivative ($R^2 = 0.34$). Fit during the opening phase also contributed weakly but significantly ($p < 0.01$) to the perceptual match. In a follow-up experiment, we fitted the models to the AbS sources subject to MSE criteria while constraining the models to fit the negative peak of the flow derivative precisely, which significantly increased the mismatch to the opening phase ($p < 0.01$; see Figure 1). Informal listening tests on several tokens showed that this significant mismatch to the opening phase resulted in a noticeable perceptual difference between the target and modeled stimuli. These results indicate the need for a source model with increased flexibility to provide a close fit to all parts of the voice source signal, especially the opening phase.

In this study, a new voice source model, motivated by data from high-speed laryngeal videoendoscopy, is proposed to capture perceptually-important source shape aspects. This model is then evaluated in comparison to 4 existing source models, with respect to fit in both the MSE and perceptual senses.

2. Data and methods

2.1. Stimuli

Source model comparisons required a target source pulse to which the models could be fitted, and the need for experimental control during perceptual evaluations mandated that this target be synthetic, so that voice stimuli could be created that differed in the source, with all other parameters held constant. To ensure that these synthetic targets were as natural in quality as possible and that they represented a range of naturally-occurring voice qualities, target stimuli were derived via analysis-by-synthesis (AbS [27]) from 40 natural samples (20 male, 20 female) of the vowel /a/. A steady-state vowel was chosen because it is routinely used for evaluating voice quality and carries substantial information about the voice source. Further, the simpler acoustic structure of a steady-state vowel should yield responses from listeners in the perceptual studies reflecting simpler perceptual strategies that are more easily interpreted. Samples were directly digitized at 20 kHz using a Brüel & Kjær microphone (model 4193), and a 1-second-long segment was excerpted for analysis. The synthesizer sampling rate was fixed at 10 kHz. Parameters describing the harmonic part of the voice source were estimated from a representative cycle of phonation for each voice using the inverse filtering method described in [28]. The harmonic and inharmonic components (the noise excitation) were identified using a comb-lifering operation in the cepstrum domain [29]. Spectrally-shaped noise was synthesized by passing white noise through a 100-tap finite impulse response filter fitted to that noise spectrum. F0 was estimated pulse by pulse using the time domain waveform. Formant frequencies and bandwidths were estimated using autocorrelation linear predictive coding analysis with a window of 25.6 ms. The complete synthesized source was then filtered through the vocal tract model, and all parameters were adjusted until the synthetic copy formed an acceptable match to the original natural voice sample. A paired comparison (same/different) task ensured that the AbS tokens were indistinguishable from the natural stimuli: d prime ranged from 0 to 1.32 across voices, with a mean of 0.79 (sd=0.41). Given these results, the AbS tokens were used in place of the natural voice samples as the target stimuli in all subsequent analyses.

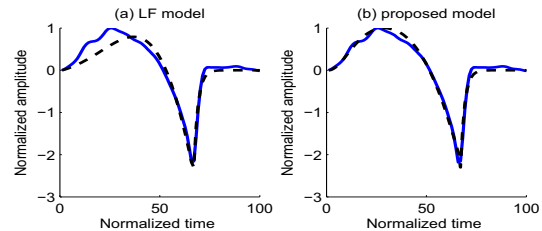


Figure 1: An example of fitting the LF and the proposed models to the same AbS source pulse subject to MSE criteria while constraining the models to fit the negative peak of the flow derivative precisely. Solid line: AbS source. Dashed line: model-fitted source.

2.2. The proposed model

The proposed model is based on the models in [7, 8], which were motivated by shapes of glottal area waveforms extracted from laryngeal high-speed videoendoscopy. The model is a combination of sinusoidal and exponential functions shown to be effective in approximating a wide range of glottal flow pulse shapes. The model is then refined using AbS to eventually capture the shapes of the glottal flow derivative, as the LF model does. The model has six parameters: the time of the positive peak (t_i), the shape of the opening (S_1 ; amplitude of the waveform at $t_i/2$), the time of the peak flow (t_p ; zero-crossing of the flow derivative), the time of the negative peak (t_e), the amplitude of the negative peak (E_e), and the slope of the return phase (t_a). The latter four parameters (t_p , t_e , E_e , and t_a) were originally defined in the four-parameter LF model [3]. The first two parameters were added in the proposed model to provide an additional degree of freedom, so that the timing of the positive peak and the shape from the start to the positive peak can be manipulated directly, independent of the negative peak of the flow derivative. The parameters are perceptually-motivated, as mentioned in the Introduction. With these parameters, the glottal opening phase could be modeled more accurately. Recall that our previous studies showed that a significant mismatch to the opening phase could lead to a noticeable perceptual difference between the target and the modeled stimuli. An example of a model waveform is shown in Figure 2. Given the six parameters described above, mathematically the glottal flow derivative $u(t)$ is defined as:

$$u(t) = \begin{cases} f\left(\frac{t}{t_i}, \lambda_1\right) & (0 \leq t \leq t_i) \\ \left[f\left(\frac{2t_e - t_i - t}{2(t_e - t_i)}, \lambda_2\right) - 1 \right] \frac{12(1 + E_e)}{6 + \lambda_2} + 1 & (t_i < t \leq t_e) \\ \frac{-E_e}{\epsilon t_a} \left[e^{-\epsilon(t - t_e)} - e^{-\epsilon(t_c - t_e)} \right] & (t_e < t \leq 1) \end{cases}$$

$$f(t, \lambda) = \frac{1}{\pi(e^\lambda + 1)} \{ e^{\lambda t} [\lambda \sin(\pi t) - \pi \cos(\pi t)] + \pi \}$$

$$\lambda_1 = 12 \cdot (0.5 - S_1)$$

$$\lambda_2 = \arg \min_{\lambda} \left| f\left(\frac{2t_e - t_p - t_i}{2(t_e - t_i)}, \lambda\right) - \frac{12E_e + 6 - \lambda}{12(E_e + 1)} \right|$$

$$\epsilon = \frac{1}{t_a} [1 - e^{-(t_c - t_e)/t_a}]$$

t_c is the time of closure. In practice it is convenient to set $t_c = 1$, i.e., the complete fundamental period [3]. ϵ , λ_1 and λ_2 are intermediate parameters. As illustrated in Figure 2, the proposed parameters can be easily derived from the inverse-filtered differential glottal waveform, and directly control the shape of the glottal waveform in a straightforward way. Unlike the LF model, which describes the open phase ($0 < t < t_e$) using one

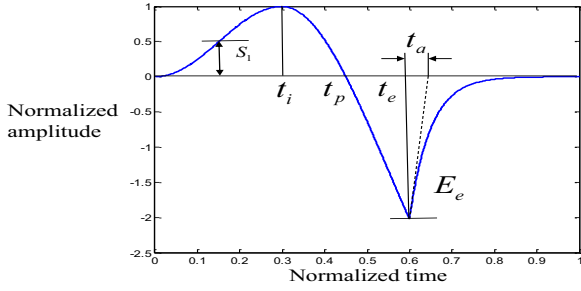


Figure 2: An example of the proposed model with $S_1 = 0.5$, $t_i = 0.3$, $t_p = 0.45$, $t_e = 0.6$, $E_e = 2$, and $t_a = 0.05$.

function, the proposed model uses two functions ($0 < t < t_i$ and $t_i < t < t_e$) to describe the open phase, allowing for more flexibility in modeling. Figure 1 (b) shows an example of constraining the proposed model to fit the negative peak of the flow derivative precisely, while still achieving satisfactory fittings in other parts.

2.3. Model fitting

In this study, each of the 40 target AbS-derived source functions was fitted with 5 source models: the Ros, LF, EE1, EE2, and the proposed model. The FL model, which provided the worst fit to the target sources in our previous experiment, was excluded from further experiments. First-derivative representations were calculated mathematically for the Ros, EE1, and EE2 models, which describe flow pulses in the time domain, so that all models were fitted to the target AbS source functions in the flow derivative domain. One cycle of the AbS source signal for each speaker was normalized to a maximum amplitude of 1. Each derivative-domain model was fitted to all of the AbS source functions using MSE criteria, for which each point of the waveform was weighted equally. Additionally, the proposed model was fitted a second time to the AbS source function with the constraint of exactly matching the first point, the positive peak of the flow derivative, the time of maximum flow (zero-crossing of flow derivative), and the negative peak of the flow derivative. This procedure was included in order to assess the perceptual importance of the landmarks of the voice source signal. Note that it is not always possible to exactly match ALL landmarks for the other models, due to constraints inherent in the models and their parameters. Because of the increased flexibility, especially in modeling the opening phase, the proposed model is able to match all landmarks well. Target AbS source pulses and the corresponding least-MSE-fitted sources using the proposed model for six different speakers are shown in Figure 3. As this figure shows, the proposed model is able to approximate a wide range of pulse widths, pulse skewnesses, and abruptnesses of glottal closure. Because this model fitting is a non-linear optimization problem and suboptimal solutions might be found using standard optimization methods, model fitting was implemented using a codebook search schema (exhaustive search) similar to that in [8] in order to achieve nearly optimal solutions. The codebook of each model has a size of 10000.

2.4. Perceptual experiment

To determine the perceptual importance of these results, we generated synthetic copies of the voices using each modeled source pulse for each voice, with all other synthesizer parameters held constant at the values derived during AbS, as illustrated in Figure 4. For the proposed model, only the model-

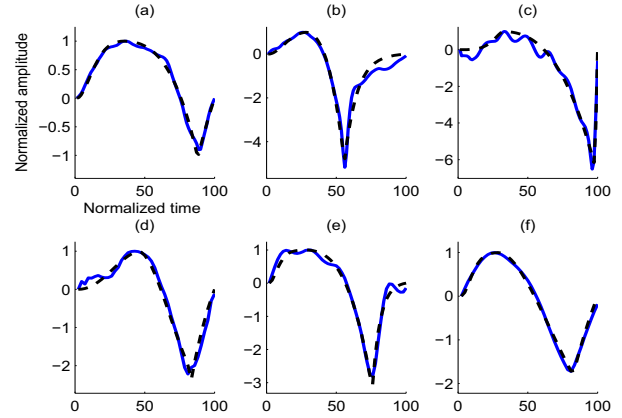


Figure 3: Target AbS source pulses and the corresponding least-MSE-fitted sources using the proposed model for six different speakers. Panels (a), (b), and (c): male speakers. Panels (d), (e), and (f): female speakers. Solid line: AbS source. Dashed line: the proposed model.

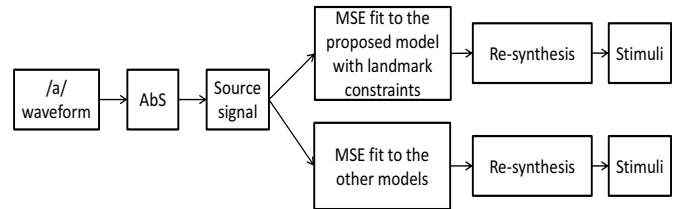


Figure 4: Flowchart showing how stimuli were generated for the perceptual experiment.

fitted sources with exact matching at the landmark points were used in this experiment (denoted “Proposed-LM”). 40 listeners (UCLA students and staff; 18-33 years of age; $M=21.15$ years; $sd=3.03$ years) assessed the similarity of all versions of each voice in a visual sort-and-rate task [30, 31], in which listeners assessed the extent of perceived match between the original voice samples and each copy. Each listener heard 10 voice “families”, where each family included an original natural voice sample, the corresponding target AbS token, and the 5 model-synthesized tokens of the same voice, such that across subjects each family was judged by 10 listeners. The stimuli were presented as distinct icons on the screen. For each family (each trial), listeners were asked to play the stimuli by clicking the icons, and to place perceptually similar sounds close together on a line on the screen, while perceptually dissimilar sounds were to be placed farther away. Listeners were instructed to use as much of the line for sorting the stimuli as they wished. They could listen to the stimuli as often as they like, and the study was not timed.

Although listeners saw no numerical values associated with the endpoints of the line, the left and right endpoints were assigned values of 0 and 1000, respectively. Thus, a numerical value could be assigned to the position of each token. We then calculated the distance of each modeled token from the target AbS voice, and this value was subsequently normalized within family for the range of values used on that given trial by that listener. The absolute values of these normalized distances were used in subsequent analyses, because the orientation of the line was arbitrary and varied from listener to listener.

3. Results

3.1. Overall model fit

Table 1 shows MSE values for fit, of each of the source models under study, to the target AbS sources. (See table caption for the meaning of model labels.) A two-way repeated measures ANOVA (model by speaker sex) showed significant main effects of model [$F(5, 190) = 12.99, p < 0.0001$] and sex [$F(1, 38) = 8.71, p < 0.01$] on mean MSE, as well as a significant model by sex interaction effect [$F(5, 190) = 4.27, p < 0.001$]. Tukey post-hoc t-tests (with Bonferroni adjustment for multiple comparisons) indicated that no cross-model differences were significant for female speakers. For male speakers, a separate t-test showed that the “Proposed” model had lower MSE values than the Ros, LF, EE1, and EE2 models ($p < 0.05$).

Table 1: MSE values (in %) of fitting models to the AbS sources. “Proposed” denotes fitting the proposed model subject to MSE criteria. “Proposed-LM” denotes fitting the proposed model subject to MSE criteria with the constraint of exact landmark matching.

	Ros	LF	EE1	EE2	Proposed	Proposed-LM
Male	27.8	14.1	25.8	21.6	3.9	6.9
Female	11.3	3.6	3.8	3.5	1.2	1.6

3.2. Perceptual experiment

Results of the perceptual experiment are shown in Table 2. Recall that 40 listeners participated in this task, but each only heard 10 of the 40 voices. Thus, every 4 subjects heard the stimuli from all 40 voices. Because a pre-test showed no significant differences in rating, we averaged the results of every 4 subjects, to make 10 “metasubjects”, where each “metasubject” (consisting of 4 listeners) heard all 40 voices. This enabled us to run an ANOVA with “metasubject” as the error term. A two-way (model by sex of voice) repeated-measures ANOVA showed significant main effects of model [$F(4, 36) = 155.77, p < 0.0001$] and sex [$F(1, 9) = 26.49, p < 0.001$] on mean perceptual distance, as well as a significant model by sex interaction effect [$F(4, 36) = 10.62, p < 0.001$]. Tukey post-hoc t-tests (with Bonferroni adjustment for multiple comparisons) indicated that the proposed-LM model formed a significantly better match to the target AbS stimulus (lower mean perceptual distance) than the other models ($p < 0.0001$). The perceptual distance to the target token for the LF model was only lower than that of the Ros model ($p < 0.0001$), but not statistically different from those of the EE1 and EE2 models. The difference between male and female voices in perceptual distances between the modeled and target tokens was significant only for the Ros model, for which male voices were closer perceptual matches to the AbS voice than female voices ($p < 0.0001$). For both sexes, the Ros model had a higher perceptual distance than the other models ($p < 0.0001$).

Table 2: Normalized perceptual distances (range from 0 to 1) between the model-fitted voices and the target AbS voice, for male and female voices. A smaller number indicates a closer perceptual distance (closer match) to the target AbS voice.

	Ros	LF	EE1	EE2	Proposed-LM
Male	0.57	0.46	0.38	0.40	0.26
Female	0.71	0.42	0.46	0.43	0.32

4. Relation to prior work

This paper presented a systematic perceptual evaluation of various source models, and proposed a new model to capture perceptually-relevant information. The study in [9] investigated the factors of vocal quality that might be affected by changes in voice source signals but only 3 listeners were involved. In that study, only the LF model was used to generate the source signal. In [4], 6 models were evaluated but were only used in a task to minimize the linear predictive error from the original voice. In this study, 5 models were evaluated in terms of both physical fits (MSE) to the AbS source and perceptual matches to the target AbS stimuli. Results were based on perceptual experiments with 40 listeners and 40 voice samples.

5. Discussion

Compared to the 4-parameter LF model [3], 2 perceptually-motivated parameters were added in the proposed model to provide more flexibility in matching the glottal opening phase. With the increased number of parameters, it is not surprising that the proposed model provided a better model fit. Nevertheless, the significant improvement achieved by the proposed model over the LF model in perceptual experiments indicated that the source variability at the opening phase (captured by the two additional parameters) is perceptually salient. Recall that the characteristics of the glottal closing phase (e.g., the negative peak of the flow derivative) have usually been assumed to be perceptually important, because of their association with the main acoustic excitation of the vocal tract [32]. However, this study demonstrated the perceptual importance of the glottal source shape at the opening phase, providing insights to modeling studies and synthesis applications. In addition, the parameters in the proposed model are based on the landmarks of the glottal pulse and can be measured directly from the glottal waveform, allowing more efficient source parameterizations in applications such as speech coding.

6. Conclusion and future work

This study presented a new voice source model with increased flexibility to capture the perceptually-important source shape aspects. Five voice source models were fitted to 40 natural voices obtained by inverse filtering and analysis-by-synthesis (AbS). Synthetic copies of the voices were generated using each modeled source pulse. Models were perceptually evaluated using a visual sort-and-rate task in which listeners assessed the extent of perceived match between the AbS copies and stimuli created with model-fitted sources. Compared to the other models, on average, the proposed model provided more accurate fittings (in terms of MSE) to the AbS-derived source. In addition, perceptual experiments showed that the proposed model provided closer perceptual matches to the target AbS voice than the other models. In order to demonstrate the potential applicability of the proposed model for improving the quality of speech synthesis, a preliminary experiment was conducted in which source models were fitted to source signals representing different voice qualities (breathy, modal, and pressed) and F0 levels. Pilot results showed that, on average, the proposed model provided a more accurate fit than did the other models. Future work will examine the effect of using this model in synthesizing continuous speech.

7. Acknowledgements

This work was supported in part by NSF Grant No. IIS-1018863 and by NIH/NIDCD Grant Nos. DC01797 and DC011300.

8. References

- [1] G. Fant, *Acoustic theory of speech production*, 2nd ed. The Hague, Paris: Mouton, 1970, pp. 15-20.
- [2] A. Rosenberg, "Effects of the glottal pulse shape on the quality of natural vowels," *J. Acoust. Soc. Am.*, vol. 49, pp. 583–590, 1971.
- [3] G. Fant, J. Liljencrants, and Q. Lin, "A four-parameter model of glottal flow," *STL-QPSR*, vol. 4, pp. 1–13, 1985.
- [4] H. Fujisaki and M. Ljungqvist, "Proposal and evaluation of models for the glottal source waveform," in *ICASSP*, 1986, pp. 1605–1608.
- [5] R. Veldhuis, "A computationally efficient alternative for the liljencrants–fant model and its perceptual evaluation," *J. Acoust. Soc. Am.*, vol. 103, pp. 566–571, 1998.
- [6] K. Cummings and M. Clements, "Glottal models for digital speech processing: A historical survey and new results," *Digital Signal Processing*, vol. 5, pp. 21–42, 1995.
- [7] Y.-L. Shue and A. Alwan, "A new voice source model based on high-speed imaging and its application to voice source estimation," in *ICASSP*, 2010, pp. 5134–5137.
- [8] G. Chen, Y.-L. Shue, J. Kreiman, and A. Alwan, "Estimating the voice source in noise," in *Interspeech*, 2012, pp. 1600–1603.
- [9] D. Childers and C. Lee, "Vocal quality factors: Analysis, synthesis, and perception," *J. Acoust. Soc. Am.*, vol. 90, pp. 2394–2410, 1991.
- [10] H. M. Hanson, K. N. Stevens, H.-K. J. Kuo, M. Y. Chen, and J. Slifka, "Towards models of phonation," *J. Phonetics*, vol. 29, no. 4, pp. 451–480, 2001.
- [11] P. Alku, T. Bäckström, and E. Vilkmán, "Normalized amplitude quotient for parametrization of the glottal flow," *J. Acoust. Soc. Am.*, vol. 112, pp. 701–710, 2002.
- [12] G. Chen, J. Kreiman, B. R. Gerratt, J. Neubauer, Y.-L. Shue, and A. Alwan, "Development of a glottal area index that integrates glottal gap size and open quotient," *J. Acoust. Soc. Am.*, vol. 133, pp. 1656–1666, 2013.
- [13] J. Kane and C. Gobl, "Wavelet maxima dispersion for breathy to tense voice discrimination," *Audio, Speech, and Language Processing, IEEE Transactions on*, vol. 21, pp. 1170–1179, 2013.
- [14] M. Airas and P. Alku, "Comparison of multiple voice source parameters in different phonation types," in *Interspeech*, 2007, pp. 1410–1413.
- [15] C. T. Ishi, H. Ishiguro, and N. Hagita, "Improved acoustic characterization of breathy and whispery voices," in *Interspeech*, 2011, pp. 2965–2968.
- [16] Y.-L. Shue, G. Chen, and A. Alwan, "On the interdependencies between voice quality, glottal gaps, and voice-source related acoustic measures," in *Interspeech*, 2010, pp. 34–37.
- [17] G. Chen, J. Kreiman, Y.-L. Shue, and A. Alwan, "Acoustic correlates of glottal gaps," in *Interspeech*, 2011, pp. 2673–2676.
- [18] T. Bäckström, P. Alku, and E. Vilkmán, "Time-domain parameterization of the closing phase of glottal airflow waveform from voices over a large intensity range," *IEEE transactions on speech and audio processing*, vol. 10, no. 3, pp. 186–192, 2002.
- [19] G. Seshadri and B. Yegnanarayana, "Perceived loudness of speech based on the characteristics of glottal excitation source," *J. Acoust. Soc. Am.*, vol. 126, pp. 2061–2071, 2009.
- [20] J. Sundberg, E. Fahlstedt, and A. Morell, "Effects on the glottal voice source of vocal loudness variation in untrained female and male voices," *J. Acoust. Soc. Am.*, vol. 117, pp. 879–885, 2005.
- [21] P. Alku, M. Airas, E. Björkner, and J. Sundberg, "An amplitude quotient based method to analyze changes in the shape of the glottal pulse in the regulation of vocal intensity," *J. Acoust. Soc. Am.*, vol. 120, pp. 1052–1062, 2006.
- [22] J. Gudnason, M. R. Thomas, D. P. Ellis, and P. A. Naylor, "Data-driven voice source waveform analysis and synthesis," *Speech Commun.*, vol. 54, no. 2, pp. 199–211, 2012.
- [23] T. Drugman, A. Moinet, T. Dutoit, and G. Wilfart, "Using a pitch-synchronous residual codebook for hybrid hmm/frame selection speech synthesis," in *ICASSP*, 2009, pp. 3793–3796.
- [24] M. R. Thomas, J. Gudnason, and P. A. Naylor, "Data-driven voice source waveform modelling," in *ICASSP*, 2009, pp. 3965–3968.
- [25] A. Roebel, S. Huber, X. Rodet, and G. Degottex, "Analysis and modification of excitation source characteristics for singing voice synthesis," in *ICASSP*, 2012, pp. 5381–5384.
- [26] J. Kreiman, B. Gerratt, G. Chen, M. Garellek, and A. Alwan, "Perceptual evaluation of source models," *J. Acoust. Soc. Am.*, vol. 132, p. 2088, 2012.
- [27] J. Kreiman, N. Antoñanzas-Barroso, and B. Gerratt, "Integrated software for analysis and synthesis of voice quality," *Behavior Research Methods*, vol. 42, pp. 1030–1041, 2010.
- [28] H. Javkin, N. Antoñanzas Barroso, and I. Maddieson, "Digital inverse filtering for linguistic research," *J. Speech Hear. Res.*, vol. 30, pp. 122–129, 1987.
- [29] G. de Krom, "A cepstrum-based technique for determining a harmonics-to-noise ratio in speech signals," *J. Speech Hear. Res.*, vol. 36, pp. 254–266, 1993.
- [30] S. Granqvist, "The visual sort and rate method for perceptual evaluation in listening tests," *Logopedics Phoniatrics Vocology*, vol. 28, pp. 109–116, 2003.
- [31] C. Esposito, "The effects of linguistic experience on the perception of phonation," *J. Phonetics*, vol. 38, pp. 306–316, 2010.
- [32] G. Fant, "Some problems in voice source analysis," *Speech Comm.*, vol. 13, pp. 7–22, 1993.