
This is an electronic reprint of the original article.
This reprint may differ from the original in pagination and typographic detail.

Kadiri, Sudarsana Reddy; Alku, Paavo

Mel-frequency cepstral coefficients derived using the zero-time windowing spectrum for classification of phonation types in singing

Published in:
Journal of the Acoustical Society of America

DOI:
[10.1121/1.5131043](https://doi.org/10.1121/1.5131043)

Published: 08/11/2019

Document Version
Peer-reviewed accepted author manuscript, also known as Final accepted manuscript or Post-print

Please cite the original version:
Kadiri, S. R., & Alku, P. (2019). Mel-frequency cepstral coefficients derived using the zero-time windowing spectrum for classification of phonation types in singing. *Journal of the Acoustical Society of America*, 146(5), EL418-EL423. <https://doi.org/10.1121/1.5131043>

This material is protected by copyright and other intellectual property rights, and duplication or sale of all or part of any of the repository collections is not permitted, except that material may be duplicated by you for your research use or educational purposes in electronic or print form. You must obtain permission for any other use. Electronic or print copies may not be offered, whether for sale or otherwise to anyone who is not an authorised user.

Authors, JASA-EL

Mel-frequency cepstral coefficients derived using the zero-time
windowing spectrum for classification of phonation types in
singing

Sudarsana Reddy Kadiri^{1, a)} and Paavo Alku¹

Department of Signal Processing and Acoustics, Aalto University, Espoo

12200, Finland

sudarsana.kadiri@aalto.fi,

paavo.alku@aalto.fi

(Dated: 4 July 2019)

1 **Abstract:** Voice source characteristics are known to vary in different
2 phonation types due to tension of laryngeal muscles along with respira-
3 tory effort. In the present study, automatic classification of phonation
4 types in singing is studied in four classes: modal/neutral, breathy, flow,
5 and pressed. Existing studies in classification of phonation types in
6 singing use voice source features and conventional mel-frequency cep-
7 stral coefficients (MFCCs) showing poor performance due to high pitch
8 in singing. In this study, high-resolution spectra obtained using the
9 zero-time windowing (ZTW) method is utilized to capture the effect of
10 voice excitation. ZTW does not call for computing the source-filter de-
11 composition which makes it robust to high pitch. For the classification
12 of phonation types in singing, the study proposes extracting MFCCs
13 from the ZTW spectrum. The results show that the proposed features
14 give a clear improvement in classification accuracy compared to the
15 existing voice source features and conventional MFCCs.

© 2019 Acoustical Society of America.

^{a)} Author to whom correspondence should be addressed.

1. Introduction

The human voice production mechanism is capable of affecting the type of phonation in voiced utterances by changing the vibration mode of the vocal folds. For speech signals, this gives rise to the coloring of speech with different voice qualities to signal, for example, vocal emotions (Gobl and Ní Chasaide, 2003). This capability is also manifested in singing affecting the timbre of the voice. Singer identity and the listener's feelings of the singer are expressed through modulations of voice quality. For example, a breathy voice has been reported to express sweetness, a pressed voice stronger expressions, and a flow voice very active singing (Sundberg, 1987, 1999).

According to Sundberg (1987, 1999), phonation types in singing can be categorized into four classes: modal (or neutral), breathy, flow (or resonant) and pressed (or tense). Sundberg expressed phonation types within a two-dimensional space spanned by subglottal pressure and glottal airflow (Sundberg, 1987, 1999). Sundberg's studies indicated that breathy and modal voices involve lower subglottal pressure levels than pressed and flow voices, while modal and pressed voices show reduced glottal airflow than breathy and flow voices.

Different phonation types primarily arise due to the adjustments in the larynx (Gobl and Ní Chasaide, 2003; Sundberg, 1987). In modal phonation, vocal folds vibrate fully along their entire length. In breathy phonation, there is a reduction in vocal fold abduction and minimal vocal fold contact area, which result in a high level of turbulent noise. The perceptual indicator of breathiness is the sensation of excessive laryngeal airflow (Childers

37 and Lee, 1991; Grillo and Verdolini, 2008). Flow phonation corresponds to a vocal technique
38 that is used exclusively in singing (Sundberg, 1995), and it is typically produced using a
39 lowered larynx. Using flow phonation, higher levels of loudness can be achieved with less
40 effort. Pressed phonation is associated with stronger muscular tension with an elevated
41 larynx position, which influences the vocal tract shape.

42 Several studies (Airas and Alku, 2007; Proutskova *et al.*, 2013; Rouas and Ioannidis,
43 2016; Stoller and Dixon, 2016) have investigated phonation types in speech and singing with
44 voice source features that have been derived from glottal flow estimates computed with glottal
45 inverse filtering. The closing quotient of the glottal flow pulse and the difference between the
46 amplitudes of the first two harmonics (H1-H2) in the voice source spectrum were shown to
47 correlate with the amount of pressedness (Millgård *et al.*, 2016). The normalized amplitude
48 quotient (NAQ), which measures the relative length of the glottal closing phase, was shown
49 to be a more robust parameter than the conventional closing quotient for discriminating
50 breathy, neutral, and pressed vowels (Airas and Alku, 2007; Kane and Gobl, 2013; Rouas
51 and Ioannidis, 2016). The harmonics-to-noise ratio (HNR) as well as jitter and shimmer
52 features were investigated by Wakasa *et al.* (2017) to discriminate pressed and neutral voices.
53 The effects of the phonation type on the voice source were studied by Alku and Vilkman
54 (1996) in female and male voices, and it was found that the maximum flow amplitude and
55 the maximum flow declination rate (i.e., the negative peak of the glottal flow derivative)
56 discriminated phonation types.

57 In conventional studies of phonation types, glottal features are typically first extracted
58 from the source waveforms after which statistical tests are conducted to the computed fea-
59 tures to compare their performance to discriminate phonation types. A different and more
60 modern approach is to use machine learning to build an automatic classifier that is trained in
61 a data-driven manner using voice source features. The first study in automatic classification
62 of phonation types in singing was carried out by [Proutskova *et al.* \(2013\)](#) using voice source
63 features derived with inverse filtering. Classification accuracy of phonation types (breathy,
64 modal, flow, and tense) varied from 55% to 70% for various vowels and it was concluded
65 that the voice source features alone are not sufficient for classification. This is mainly due
66 to reduced accuracy of inverse filtering for singing voices, as singing voices are typically of
67 high pitch and they show strong source-filter coupling. Several features including harmonic
68 amplitudes, formant frequencies, formant bandwidths and amplitudes, HNR, and different
69 voice source features were studied by [Rouas and Ioannidis \(2016\)](#). Their study showed that
70 there are confusions between breathy and modal voices, and between voices produced in
71 flow and pressed phonation. A large number of spectral statistics such as spectral centroid,
72 spectral flux, spectral energies in different bands along with various voice source features and
73 MFCCs were investigated by [Stoller and Dixon \(2016\)](#) for classification of phonation types
74 in singing. Recently [Kadiri and Yegnanarayana \(2018a\)](#) proposed using cepstral features de-
75 rived from the single frequency filtering (SFF) method that provides higher spectro-temporal
76 resolution.

77 In the current study, high-resolution spectra obtained using the zero-time windowing
78 (ZTW) method ([Yegnanarayana and Gowda, 2013](#)) are used to capture the effect of the voice
79 excitation. For automatic classification of phonation types in singing, MFCCs computed
80 using the ZTW spectrum are proposed as features. These novel features are used to train
81 a support vector machine (SVM) classifier to conduct the automatic classification of the
82 phonation type in singing.

83 **2. Zero-time windowing (ZTW) and extraction of MFCCs using the ZTW spec-** 84 **trum**

85 This section describes the signal processing methods used for deriving high-resolution spec-
86 trum with the ZTW method ([Yegnanarayana and Gowda, 2013](#)). In addition, the extraction
87 of MFCCs using the ZTW spectrum is described. It is to be noted that the ZTW method
88 does not assume the source-filter model of speech production.

89 *2.1 ZTW*

90 The objective of the ZTW method ([Yegnanarayana and Gowda, 2013](#)) is to derive the instan-
91 taneous spectrum so that the time-varying voice production characteristics can be captured.
92 In this method, a voice signal to be analyzed is windowed with a heavily decaying window
93 that provides high emphasis on the samples near the starting (zeroth) sampling instant,
94 and hence the name zero-time windowing (ZTW). This happens for every instant of time
95 and hence the technique provides high temporal resolution. Spectral characteristics are es-
96 timated using group delay, which provides good spectral resolution. Hence, the method
97 provides higher temporal resolution while simultaneously maintaining spectral resolution.

98 Previous speech analysis studies have shown that the ZTW spectrum captures various char-
 99 acteristics of the excitation effectively, such as glottal opening and open phase, and also the
 100 vocal tract system, such as formants (Prasad and Yegnanarayana, 2016; Yegnanarayana and
 101 Gowda, 2013). Steps involved in extracting the instantaneous spectral characteristics using
 102 the ZTW method are as follows (Yegnanarayana and Gowda, 2013).

- 103 • The voice signal ($s[n]$) is first pre-emphasized in order to reduce the effects of low-
 104 frequency trend in the signal.
- 105 • Voice segment of L ms (number of samples: $M = Lf_s/1000$) is considered at each
 106 instant, i.e., $s[n]$ is defined for $n = 0, 1, \dots, M - 1$. The segment is multiplied with a
 107 window $w_1^2[n]$, where

$$\begin{aligned}
 w_1[n] &= 0, & n &= 0, \\
 &= \frac{1}{4 \sin^2(\pi n/2N)}, & n &= 1, 2, \dots, N - 1.
 \end{aligned} \tag{1}$$

108 N is the number of samples used in the computation of Discrete Fourier Transform (DFT)
 109 ($N \gg M$). Multiplying the signal with the window $w_1^2[n]$ is approximately equivalent to
 110 integration in the frequency domain (Yegnanarayana and Gowda, 2013). Here, $L=5$ ms and
 111 $N=1024$ are used.

- 112 • Truncation of the signal at the instant $n = M - 1$ may result in a ripple effect in the frequency
 113 domain. The ripple effect is reduced by using another window ($w_2[n]$) for $n = 0, 1, \dots, M - 1$,

114 defined as

$$w_2[n] = 2(1 + \cos(\pi n/M)) = 4 \cos^2(\pi n/2M). \quad (2)$$

- 115 • The spectrum of the windowed signal (i.e., $x[n] = w_1^2[n]w_2[n]s[n]$) is estimated using the
 116 numerator of the group delay (NGD) function ($g_n[k]$) given by

$$g_n[k] = X_R[k]Y_R[k] + X_I[k]Y_I[k], \quad k = 0, 1, 2, \dots, N - 1. \quad (3)$$

117 where $X_R[k]$ and $X_I[k]$ are the real and imaginary parts of the N -point DFT $X[k]$ of $x[n]$.
 118 Likewise, $Y_R[k]$ and $Y_I[k]$ are the real and imaginary parts of the N -point DFT $Y[k]$ of
 119 $y[n] = nx[n]$.

- 120 • In order to highlight the hidden spectral characteristics due to heavily decaying window,
 121 the NGD function is differentiated twice. The spectral features, i.e., peaks in the spectrum
 122 correspond to the resonances of the vocal tract system.
- 123 • The Hilbert envelope ([Yegnanarayana and Gowda, 2013](#)) of the double-differentiated NGD
 124 is computed and is referred to as the ZTW spectrum, denoted by $S_n[k]$.

125 The steps involved in the ZTW method are shown in the schematic block diagram in
 126 Fig. 1. Figure 2 gives an illustration of ZTW spectrograms for soprano voices of different
 127 phonation types (breathy, modal, flow, and tense). It can be clearly seen that there exist
 128 remarkable spectral variations due to the voice excitation effects on the system characteris-
 129 tics.

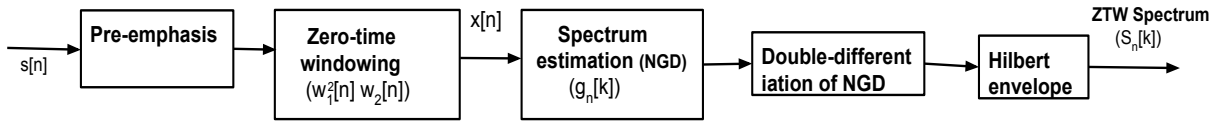


Fig. 1. Schematic block diagram describing the steps in the ZTW method.

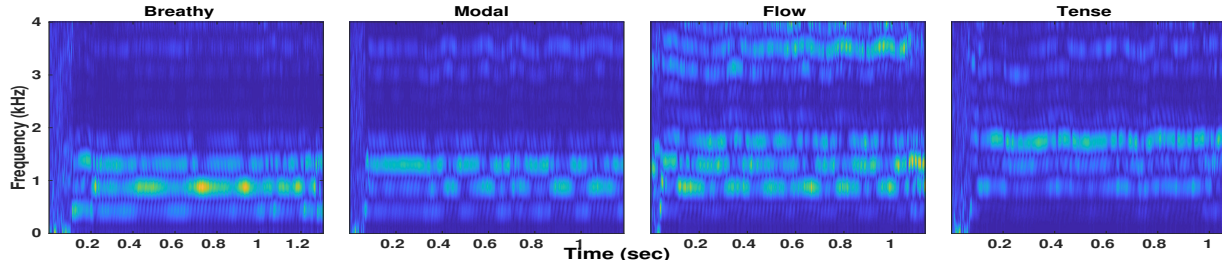


Fig. 2. (color online) An illustration of ZTW spectrograms for different phonation types (breathy, modal, flow, and tense) in soprano voices (vowel /A/).

130 2.2 Extraction of MFCCs using the ZTW spectrum

131 The schematic block diagram describing the steps involved in the extraction of MFCCs using
 132 the ZTW spectrum is shown in Fig. 3. The method performs the mel-filter bank analysis on
 133 the spectrum given by the ZTW method ($S_n[k]$), followed by logarithm and discrete cosine
 134 transform (DCT) operations, and can be expressed as follows

$$C_n[k] = DCT(\log(\text{Mel}(|S_n[k]|^2))), \quad (4)$$



Fig. 3. Extraction of MFCCs using the ZTW spectrum.

135 where $c_n[k]$ denotes the mel-cepstrum. The resulting cepstral coefficients are referred
 136 to as MFCC-ZTW, and they represent compactly the effect of the excitation on vocal tract
 137 system characteristics. The MFCC-ZTW can in principle be obtained at each time instant.
 138 In this study, however, MFCC-ZTW is computed at glottal closure instants (GCIs). From
 139 the mel-cepstrum, the first 13 cepstral coefficients (including the zeroth coefficient) are con-
 140 sidered for each frame. Delta and double-delta coefficients are also computed from the static
 141 coefficients.

142 3. Experimental protocol

143 This section describes the singing voice databases, the reference features, and the classifier.

144 3.1 Singing voice databases

145 Two databases (the soprano database and the baritone database) consisting of singing voices
 146 of different phonation types are used. The soprano database contains sustained vowels sung
 147 by a professional female Russian singer ([Proutskova et al., 2013](#)). The phonation types
 148 correspond to Sundberg’s definitions of breathy, modal, flow, and pressed voice ([Sundberg,](#)
 149 [1987](#)). The database consists of 763 recordings in nine different vowels: /A, AE, I, O, U,
 150 UE, Y, OE, and E/. Pitch ranges from A3 to G5. The baritone database contains sustained
 151 vowels sung by a professional male Greek singer ([Rouas and Ioannidis, 2016](#)). The database
 152 consists of 487 singing voice samples of five vowels /A, O, E, I, U/. Pitch ranges from A2 to
 153 G4. Both of the databases were recorded at a sampling frequency of 44.1 kHz. More details
 154 of the databases are described by [Proutskova et al. \(2013\)](#) and [Rouas and Ioannidis \(2016\)](#).

155 3.2 Reference features

156 Five sets of reference features are considered for comparison based on recent studies on dis-
 157 crimination of phonation types (Kadiri and Yegnanarayana, 2018a,b; Kane and Gobl, 2013).
 158 These reference feature sets are: (1) conventional MFCCs, (2) voice quality (VQ) features,
 159 (3) excitation features derived from the modified zero frequency filtering (ZFF) method, (4)
 160 zero time windowing cepstral coefficients (ZTWCCs), and (5) single frequency filtering cep-
 161 stral coefficients (SFFCCs). Conventional MFCCs are computed using Hamming-windowed
 162 frames of 25 ms with a shift of 5 ms. The VQ features are derived from glottal flow waveforms
 163 estimated by the iterative adaptive inverse filtering method (Alku, 2011). The VQ features
 164 consist of the normalized amplitude quotient (NAQ) (Alku *et al.*, 2002), the quasi-open quo-
 165 tient (QOQ) (Airas and Alku, 2007), H1-H2 (Airas and Alku, 2007; Mehta *et al.*, 2019), the
 166 parabolic spectral parameter (PSP) (Airas and Alku, 2007), the harmonic richness factor
 167 (HRF) (Airas and Alku, 2007), and the maximum dispersion quotient (MDQ) (Kane and
 168 Gobl, 2013). The excitation features are derived from approximate source waveforms com-
 169 puted using the modified ZFF method (Murty and Yegnanarayana, 2008). These features
 170 consist of the strength of excitation (SoE), the energy of excitation (EoE), the loudness mea-
 171 sure, and the ZFF signal energy (Kadiri and Yegnanarayana, 2018a). Cepstral features are
 172 derived from the ZTW spectrum (Kadiri and Yegnanarayana, 2018b) and the SFF spectrum
 173 (Kadiri and Yegnanarayana, 2018a), and they are referred to as ZTWCCs and SFFCCs,
 174 respectively. For MFCCs, ZTWCCs, and SFFCCs, first 13 cepstral coefficients (static coeffi-

175 cients), delta, and double-delta coefficients are computed, which results in a 39-dimensional
 176 feature vector.

177 *3.3 Classifier*

178 Support vector machine (SVM) with a radial basis function (RBF) kernel is used as a clas-
 179 sifier (Chang and Lin, 2011). Classification experiments are conducted using 10-fold cross-
 180 validation. One fold is held out to be used for testing, with the remaining nine folds used
 181 for training.

182 **4. Classification experiments and results**

183 Table 1 shows the results of the 10-fold cross validation experiments in terms of mean and
 184 standard deviation of the classification accuracy for the soprano and baritone voices. From
 185 the table, it can be seen that the proposed MFCC-ZTW features provide the highest average
 186 classification accuracy for both databases compared to the reference features. It is to be
 187 noted that the VQ features are not able to discriminate phonation types in singing as well
 188 as in speech (Kadiri and Yegnanarayana, 2018a; Kane and Gobl, 2013). This is because the
 189 VQ features (except MDQ) are derived from glottal flow waveforms estimated by inverse
 190 filtering whose accuracy is known to deteriorate for high-pitched voices (Alku, 2011). The
 191 performance of the existing MFCCs and SFFCCs is similar in both databases.

192 Tables 2 and 3 show the confusion matrices using the proposed MFCC-ZTW features
 193 for the soprano and baritone voices, respectively. From the tables, it can be observed that
 194 there is clear confusion between breathy and modal voices, and between flow and tense voices.
 195 It can also be observed that flow phonation shows confusion with voices of modal phonation.

Table 1. Mean and standard deviation of classification accuracy after 10-fold cross validation with different feature vectors for the soprano and baritone datasets.

Features	Soprano ($\mu \pm \sigma$)[%]	Baritone ($\mu \pm \sigma$)[%]
VQ	47.18 \pm 7.02	57.19 \pm 7.04
MFCCs	66.97 \pm 6.67	76.35 \pm 6.53
Excitation	52.11 \pm 5.92	47.11 \pm 6.12
ZTWCCs	66.06 \pm 6.13	74.62 \pm 4.38
SFFCCs	67.83 \pm 3.99	75.24 \pm 4.44
MFCC-ZTW	80.32 \pm 3.74	85.17 \pm 4.24

Table 2. Confusion matrix with 10-fold cross validation for the soprano database using the proposed MFCC-ZTW features.

	Breathy [%]	Modal [%]	Flow [%]	Tense [%]
Breathy	87.79	10.98	0	1.23
Modal	16.01	72.66	5.33	6.00
Flow	0.65	11.76	76.47	11.11
Tense	2.81	4.28	7.99	84.91

Table 3. Confusion matrix with 10-fold cross validation for the baritone database using the proposed MFCC-ZTW features.

	Breathy [%]	Modal [%]	Flow [%]	Tense [%]
Breathy	86.75	7.64	5.61	0
Modal	6.90	87.25	4.73	1.11
Flow	08.66	6.56	75.40	9.38
Tense	0	1.89	8.16	89.95

¹⁹⁶ These observations are also in line with the results reported by Proutskova *et al.* (2013), by

197 [Stoller and Dixon \(2016\)](#), and by [Rouas and Ioannidis \(2016\)](#). Even though the proposed
198 features show a remarkable improvement in accuracy, there is still confusion between modal
199 and breathy, and between flow and tense voices. Hence, there is a need for exploring features
200 that can capture changes in voice production characteristics, especially for the discrimination
201 between breathy and modal voices, and between flow and tense voices.

202 **5. Summary and conclusions**

203 In this study, MFCCs derived using the ZTW spectrum were proposed for automatic classifi-
204 cation of phonation types in singing. The ZTW method provides high-resolution spectra and
205 captures the effect of excitation on the vocal tract characteristics. From the experimental
206 results, it was shown that the proposed MFCC-ZTW features provide a better discrimina-
207 tion of phonation types in singing voices compared to several known reference features. The
208 voice source features derived using glottal inverse filtering show less accuracy due to poor
209 performance of inverse filtering in the estimation of the glottal flow for high-pitched voices
210 in singing. The proposed features, however, do not suffer from this problem because they
211 do not use source-filter separation for deriving features. This suggests that the proposed
212 features can be useful for analyzing spontaneous, continuous speech in addition to singing
213 voices. Despite the large improvement in classification accuracy achieved by the proposed
214 MFCC-ZTW features in the current study, new voice production feature extraction tech-
215 niques are needed to better discriminate breathy and modal voices, as well as voices in flow
216 and tense phonation.

217 **6. Acknowledgements**

218 This study was partly funded by the Academy of Finland (project no. 312490).

219 **References and links**

220

221 Airas, M., and Alku, P. (2007). “Comparison of multiple voice source parameters in different
222 phonation types,” in *Proc. INTERSPEECH*, pp. 1410–1413.

223 Alku, P. (2011). “Glottal inverse filtering analysis of human voice production—a review of
224 estimation and parameterization methods of the glottal excitation and their applications,”
225 *Sadhana* **36**(5), 623–650.

226 Alku, P., Backstrom, T., and Vilkmán, E. (2002). “Normalized amplitude quotient for
227 parametrization of the glottal flow,” *J. Acoust. Soc. Am.* **112**(2), 701–710.

228 Alku, P., and Vilkmán, E. (1996). “A comparison of glottal voice source quantification
229 parameters in breathy, normal and pressed phonation of female and male speakers,” *Folia*
230 *Phoniatr Logop.* **48**, 240 – 254.

231 Chang, C.-C., and Lin, C.-J. (2011). “Libsvm: A library for support vector machines,”
232 *ACM Trans. Intell. Syst. Technol. (TIST)* **2**(3), 27.

233 Childers, D. G., and Lee, C. (1991). “Vocal quality factors: Analysis, synthesis, and per-
234 ception,” *J. Acoust. Soc. Am.* **90**(5), 2394–2410.

235 Gobl, C., and Ní Chasaide, A. (2003). “The role of voice quality in communicating emotion,
236 mood and attitude,” *Speech Commun.* **40**(1-2), 189–212.

- 237 Grillo, E. U., and Verdolini, K. (2008). “Evidence for distinguishing pressed, normal, res-
238 onant, and breathy voice qualities by laryngeal resistance and vocal efficiency in vocally
239 trained subjects,” *J. Voice* **22**(5), 546–552.
- 240 Kadiri, S. R., and Yegnanarayana, B. (2018a). “Analysis and detection of phonation modes
241 in singing voice using excitation source features and single frequency filtering cepstral
242 coefficients (SFFCC),” in *Proc. INTERSPEECH*, pp. 441–445.
- 243 Kadiri, S. R., and Yegnanarayana, B. (2018b). “Breathy to tense voice discrimination using
244 zero-time windowing cepstral coefficients (ztwccs),” in *Proc. INTERSPEECH*, pp. 232–236.
- 245 Kane, J., and Gobl, C. (2013). “Wavelet maxima dispersion for breathy to tense voice
246 discrimination,” *IEEE Trans. Audio, Speech & Lang. Process.* **21**(6), 1170–1179.
- 247 Mehta, D. D., Espinoza, V. M., Van Stan, J. H., Zañartu, M., and Hillman, R. E. (2019).
248 “The difference between first and second harmonic amplitudes correlates between glottal
249 airflow and neck-surface accelerometer signals during phonation,” *J. Acoust. Soc. Am.*
250 **145**(5), EL386–EL392.
- 251 Millgård, M., Fors, T., and Sundberg, J. (2016). “Flow glottogram characteristics and per-
252 ceived degree of phonatory pressedness,” *J. Voice* **30**(3), 287–292.
- 253 Murty, K. S. R., and Yegnanarayana, B. (2008). “Epoch extraction from speech signals,”
254 *IEEE Trans. Audio, Speech, & Lang. Process.* **16**(8), 1602–1613.
- 255 Prasad, R. S., and Yegnanarayana, B. (2016). “Determination of glottal open regions by
256 exploiting changes in the vocal tract system characteristics,” *J. Acoust. Soc. Am.* **140**(1),
257 666–677.

- 258 Proutskova, P., Rhodes, C., Crawford, T., and Wiggins, G. (2013). “Breathy, resonant,
259 pressed–automatic detection of phonation mode from audio recordings of singing,” *J. New*
260 *Music Res.* **42**(2), 171–186.
- 261 Rouas, J.-L., and Ioannidis, L. (2016). “Automatic classification of phonation modes in
262 singing voice: Towards singing style characterisation and application to ethnomusicological
263 recordings,” in *Proc. INTERSPEECH*, pp. 150 – 154.
- 264 Stoller, D., and Dixon, S. (2016). “Analysis and classification of phonation modes in
265 singing,” in *Proc. International Society for Music Information Retrieval Conference*.
- 266 Sundberg, J. (1987). *The Science of the Singing Voice* (Illinois University Press).
- 267 Sundberg, J. (1995). “Vocal fold vibration patterns and modes of phonation,” *Folia Phoni-*
268 *atrica et Logopaedica* **47**(4), 218–228.
- 269 Sundberg, J. (1999). “The perception of singing,” in *The Psychology of Music (Second*
270 *Edition)* (Elsevier), pp. 171–214.
- 271 Wakasa, K., Matsubara, M., Hiraga, Y., and Terasawa, H. (2017). “Acoustic characteristics
272 of pressed and normal phonations in choir singing by male singers,” in *Proc. International*
273 *Symposium on Musical Acoustics*, pp. 136–139.
- 274 Yegnanarayana, B., and Gowda, D. (2013). “Spectro-temporal analysis of speech signals
275 using zero-time windowing and group delay function,” *Speech Commun.* **55**(6), 782–795.