

On H1–H2 as an acoustic measure of linguistic phonation type

Yuan Chai^a and Marc Garellek^b

*Department of Linguistics, University of California San Diego, La Jolla, 92093,
USA*

(Dated: 13 October 2022)

1 [This paper is part of a special issue on Reconsidering Classic Ideas in Speech Com-
2 munication.]

3 The measure H1–H2, the difference in amplitude between the first and second
4 harmonic, is frequently used to distinguish phonation types and to characterize dif-
5 ferences across voices and genders. While H1–H2 can differentiate voices and is used
6 by listeners to perceive changes in voice quality, its relation to voice articulation is
7 less straightforward. Its calculation also involves practical issues with error propa-
8 gation. This paper highlights some developments in the use of H1–H2 and proposes
9 a new measure that we call “residual H1.” In residual H1, the amplitude of the
10 first harmonic is normalized against the overall sound energy (as measured by Root-
11 mean-square Energy) instead of against H2. Residual H1 may mitigate some of the
12 issues with using H1–H2. The current study tests the correlation between Residual
13 H1 and electroglottographic contact quotient (CQ) and compares the ability of resid-
14 ual H1 vs H1–H2 to differentiate statistically across phonation types in !Xóõ and
15 utterance-level changes in phonatory quality in Mandarin. The results show that
16 residual H1 has a stronger correlation with CQ and differentiates contrastive and
17 allophonic phonatory quality better than H1–H2, particularly for more constricted
18 phonation types.

^ayuc521@ucsd.edu

^bmgarellek@ucsd.edu

19 I. INTRODUCTION

20 The acoustic measure H1–H2, also known as L_1 – L_2 (Titze *et al.*, 2015), refers to the
21 difference in amplitude between the first and second harmonics. It is probably the most
22 widely-used voice quality measure in linguistic phonetic research, and correlates with changes
23 in phonation type (e.g. breathy vs. modal vowels) in many languages (Esposito and Khan,
24 2020; Gordon and Ladefoged, 2001), as well as non-phonemic changes in phonation (Hanson
25 *et al.*, 2001; Li *et al.*, 2020; Ní Chasaide and Gobl, 1993). In terms of aerodynamics, H1–H2
26 reflects the amount of airflow through the glottis (Sundberg and Gauffin, 1979); in terms of
27 voice articulation, the measure is related to vocal fold (and, perhaps more broadly, laryngeal)
28 constriction vs. spreading. Generally, lower values of H1–H2 are associated with lower glottal
29 Open Quotient (OQ), more constriction, and increased medial vocal fold thickness (Kreiman
30 *et al.*, 2012; Samlan and Story, 2011; Zhang, 2016b).

31 The relationship between H1–H2, aerodynamics, and voice articulation is better studied
32 than for any other acoustic measure of phonatory quality. Nevertheless, researchers occa-
33 sionally find that phonation types are not distinguished by H1–H2, even when non-modal
34 phonation is perceptually strong (Esposito, 2012; Garellek and Esposito, 2021). This sug-
35 gests that H1–H2 may not be ideally suited for indexing changes in vocal fold constriction
36 as generally thought, and/or that the measure can be refined in some way. In this paper,
37 we review the history and use of H1–H2, particularly in linguistic studies of phonation type,
38 and discuss the possible reasons for a lack of effect of H1–H2 when distinguishing modal
39 vs. non-modal phonation types. We then motivate the use of a related measure – residual

40 H1 – to compare modal vs non-modal phonation. Residual H1 is a measure of H1 controlled
41 for overall sound pressure level (SPL). We show that it is as effective as H1–H2 at discrimi-
42 nating modal vs. non-modal phonation in contrastive and allophonic uses of phonation type;
43 further, its use mitigates certain issues inherent to the use of H1–H2.

44 Moving forward, we first highlight some notes on terminology. We will use “H1–H2”
45 instead of “L₁–L₂” because the former is more widely known, particularly in linguistic pho-
46 netic research. “H1–H2” will be used to refer to any measure that compares the amplitude of
47 the fundamental (H1) to that of the second harmonic (H2). But in any given study, H1–H2
48 can be measured in different ways. If estimated from the voice source (e.g. using inverse
49 filtering), we will refer to the measure as “source H1–H2.” If calculated from the audio out-
50 put without any correction, we will call the measure “uncorrected H1–H2.” Finally, if the
51 measure has been obtained from the audio output but corrected for formant frequencies and
52 bandwidths, we will refer to the measure as H1*–H2*, which is in line with current practice.
53 In what follows, we will discuss the reasons for the existence of these different versions of
54 the measure. Our paper also motivates a new measure called “residual H1.” This measure
55 will always be corrected for formant frequencies and bandwidths and so will appear with
56 an asterisk as “residual H1*” unless discussed more abstractly. In early work that included
57 uncorrected H1 from the audio output, we will refer to that measure as “uncorrected H1.”
58 Finally, H1 estimated at the voice source will be called “source H1,” and when referring to
59 both H1 and H1–H2, we will occasionally use “H1(–H2).”

60 **A. The origins of measuring H1(–H2)**

61 We have known since the 1960s that the roll-off or “tilt” of the harmonic spectrum varies
62 as a function of different phonation types. In her pioneering study, [Fischer-Jørgensen \(1967\)](#)
63 described the acoustic differences between modal (“clear”) and breathy (“murmured”) vowels
64 in Gujarati. Through visual inspection of audio spectra, she found that the most important
65 acoustic distinction between these phonation types is the amplitude of the first harmonic (the
66 fundamental, i.e. uncorrected H1). She found that, for Gujarati breathy vowels, uncorrected
67 H1 is generally stronger than for modal vowels.

68 Fischer-Jørgensen was surprised by this finding: “I had expected to find some extra
69 noise in the breathy vowels, but instead I found a reinforcement of the fundamental” (p.
70 71). Further, she knew that to measure H1 on its own would present a confound between
71 phonation differences and differences in sound intensity; for instance, a stronger uncorrected
72 H1 may be due to increased breathiness, but it can also result from an overall higher sound
73 pressure level. The way to disambiguate between these hypotheses is to normalize for SPL
74 in some way. Fischer-Jørgensen did so by subtracting H2 from H1. If the overall signal is
75 relatively strong, this should affect both H1 and H2 equally; thus, H1–H2 can index the
76 strength of the fundamental while normalizing for any differences in SPL across tokens.

77 The choice to normalize for SPL using H2, rather than another spectral landmark, was
78 not motivated *a priori*; indeed, Fischer-Jørgensen also normalized the fundamental relative
79 to other uncorrected harmonics like H3 and the amplitudes of formants 1–4. Interestingly,
80 she found that the spectral tilt differences between breathy and modal vowels were not

81 consistently found across all measures: uncorrected H1 was indeed stronger in breathy vowels
82 compared to modal ones when normalized against uncorrected H2 or the amplitudes of F1,
83 F2, and F4 but not when normalized against uncorrected H3 or the amplitude of F3. At any
84 rate, the implications of her decision to normalize H1 against H2 remain today: even though
85 H1–H2 involves two harmonic amplitudes, the assumption – usually tacit – is that what we
86 wish to compare across phonation types is a difference in the strength of the fundamental,
87 that is to say H1.

88 **B. H1–H2 and its relation to articulation and perception**

89 While [Fischer-Jørgensen \(1967\)](#) established the importance of H1 and spectral tilt mea-
90 sures like H1–H2 as correlates of a breathy-modal contrast, it remained unclear precisely
91 why breathy vowels have a relatively stronger fundamental than modal ones. This question,
92 though still not fully resolved, has been addressed since the 1970s. [Stevens \(1977\)](#) used
93 models of the transglottal area to schematize overall spectral tilt differences as a function
94 of the degree of inter-arytenoid space, predicting that creaky phonation should have the
95 lowest spectral tilt and breathy phonation the highest. Yet H1 and H1–H2 were not explic-
96 itly discussed. In fact, his depictions of differences in spectral tilt suggest that he would
97 not have predicted differences in H1(–H2) for modal vs. creaky phonation; see Figure 5 in
98 that paper, where the increased tilt is schematized in the higher frequencies only. Still, this
99 work is important in highlighting how different vocal fold configurations, and in particular
100 how changes to the cartilaginous glottis, could affect spectral tilt. Around the same time,
101 [Sundberg and Gauffin \(1979\)](#) found that source H1 is related to overall airflow through the

102 glottis. They showed how differences in source H1 are related to changes in overall SPL,
103 and can be regulated by changes in the degree of vocal fold contact during voicing.

104 Another landmark study about H1(-H2) is the MIT Speech Communication working
105 paper by [Bickley \(1982\)](#), who built on Fischer-Jørgensen’s findings for Gujarati, as well as the
106 preliminary analysis by [Ladefoged \(1981\)](#) of phonation types in !Xóõ. Bickley noted (p. 74–
107 76) that the inverse-filtered glottal source in Gujarati breathy vowels had more symmetrical
108 pulses than that of modal vowels, with less abrupt closure and shorter closed intervals [i.e.,
109 with higher glottal open quotient (OQ)]. Further development in voice source models, such
110 as the LF model of [Fant *et al.* \(1985\)](#), also showed how differences in overall pulse shape
111 relate to changes in spectral tilt. For example, [Fant and Lin \(1988\)](#) described how changes
112 in various LF model parameters modulate source H1 relative to H2 and to other harmonics;
113 for a recent overview and reassessment, see [Gobl and Ní Chasaide \(2019\)](#).

114 [Bickley \(1982\)](#) also conducted what is likely the first perceptual assessment of H1–H2,
115 though it is preliminary by today’s standards. In one experiment, she presented two listeners
116 (one a native speaker of English, the other of Gujarati) with ten tokens of breathy vowels
117 from !Xóõ, and asked them to rate the tokens on a four-point scale from “very breathy”-
118 sounding to “not breathy”-sounding. The tokens that sounded very breathy had uncorrected
119 H1–H2 values greater than 10 dB; the two tokens that sounded not breathy had uncorrected
120 H1–H2 values of -4 and 4 dB. She also resynthesized a continuum from modal to breathy
121 vowels [i, a, o] using an earlier version of the Klatt synthesizer ([Klatt and Klatt, 1990](#)), in
122 which uncorrected H1 of the “breathy” vowel was equal to that of the “modal” vowel, or
123 was higher by 9, 12, and 15 dB. The amplitude of spectral noise also varied (orthogonally to

124 uncorrected H1) in 5-dB increments over a range of 20 dB. The resynthesized vowels were
125 then spliced onto natural tokens of CV(C) words. Four native speakers of Gujarati were
126 presented with the stimuli, and were asked to do a two-alternative forced-choice task with
127 minimal pairs; that is, they chose whether they heard a word with a breathy vowel (e.g. [b̤i]
128 “be afraid”) or one with a modal vowel (e.g. [bi] “seed”). The results showed, perhaps
129 surprisingly, that the level of aspiration noise did not affect listeners’ choice. However, an
130 increase in uncorrected H1 was associated with an increase in breathy vowel responses.

131 [Klatt and Klatt \(1990\)](#) outline several studies relating to H1(–H2) and the relationship
132 between articulation, acoustics, and perception. They explicitly argue that source H1 is
133 related to increased OQ, which in turn is related to the size of the posterior glottal opening
134 (i.e., the cartilaginous glottis). [Holmberg et al. \(1995\)](#) found that uncorrected H1–H2 was
135 correlated with airflow and electroglottographic (EGG) in English speakers. Similar rela-
136 tionships between H1–H2 and OQ have been found in studies using natural speech ([Sundberg](#)
137 [et al., 1999](#)) and resynthesized/simulated data ([Stevens and Hanson, 1995](#)). In studies of
138 linguistic phonation, EGG contact quotient (CQ) (sometimes also quoted as “OQ”) is also
139 reasonably well correlated with H1–H2; for example, [DiCanio \(2009\)](#) reports adjusted R^2
140 values ranging from 0.3 to 0.46 between EGG OQ and uncorrected H1–H2 across the reg-
141 isters (phonation types) of Takhian Thong Chong, and [Kuang \(2011\)](#) reports an R^2 of 0.3
142 between EGG CQ and H1*–H2* in Southern Yi.

143 The relationship between H1–H2 and OQ, though robust, is often found to be weak and/or
144 non-linear. Using inverse filtering of oral airflow, [Hanson \(1995\)](#) obtained glottal waveform
145 and measured its OQ. Of four speakers total, three speakers showed a trend whereby larger

146 OQ was related to higher $H1^*-H2^*$ (pp. 81, 85–86). [Kreiman *et al.* \(2012\)](#) found that,
147 although source H1–H2 is closely correlated to OQ, this relationship varies considerably
148 across speakers. H1–H2 is also correlated with other articulatory or aerodynamic parameters,
149 including increased vocal fold process separation ([Samlan *et al.*, 2013](#)), increased medial vocal
150 fold thickness ([Zhang, 2016a](#)), and glottal skew or symmetry ([Doval and d’Alessandro, 1997](#);
151 [Doval *et al.*, 2006](#); [Henrich *et al.*, 2001](#); [Kreiman *et al.*, 2012](#); [Shue *et al.*, 2010](#); [Swerts and](#)
152 [Veldhuis, 2001](#)). The effect of glottal open quotient on source H1–H2 further interacts
153 with glottal skew: [Gobl *et al.* \(2018\)](#) and [Gobl and Ní Chasaide \(2019\)](#) found that, when
154 glottal OQ was high, source H2 was affected by glottal skew, while source H1 was mostly
155 independent of glottal skew. More skewed pulses were related to higher source H2. And
156 computational simulations have shown that the relationship between $H1^*-H2^*$ and vocal
157 fold process separation is non-linear, such that $H1^*-H2^*$ first increases but then decreases
158 with increasingly large separation ([Samlan and Story, 2011](#)).

159 As mentioned earlier, sometimes H1–H2 does not “behave” as expected (e.g. by not
160 showing a difference across phonation types). The findings from the aforementioned studies
161 imply that the reason for the occasional unexpected behavior of H1–H2 may be that it
162 is affected by other factors that are unrelated to glottal open quotient. Still, support for
163 the continued use of H1–H2 comes from [Kreiman *et al.* \(2007\)](#), who tested the correlation
164 between spectral shape and glottal pulse shape and 78 acoustic measures (e.g. H1–H2, H2–
165 H4, and slope of spectrum at different frequency intervals). Using correlation and principal
166 component analyses, they found that the 78 acoustic measures can be reduced to just four
167 independent ones: source H1–H2, source H2–H4, overall spectral slope, and high-frequency

168 noise. Importantly, they also found that source H1-H2 is related to variability in spectral
169 and glottal pulse shape, and stated that its measured values did not differ appreciably as a
170 function of different glottal source models, implying that H1-H2 is robust to measurement
171 artifacts.

172 Additional support for using H1-H2 comes from perceptual studies. Since [Bickley \(1982\)](#)
173 [confirmed by [Klatt and Klatt \(1990\)](#)], researchers have found that changes in H1-H2 corre-
174 late with perceived changes in breathiness. More recently, [Esposito \(2010\)](#) found that uncor-
175 rected H1-H2 correlated with perceived breathiness, regardless of whether listeners' native
176 language was Gujarati (with contrastive breathiness), English (with allophonic breathiness),
177 or Spanish (with no breathiness), though the language groups relied on uncorrected H1-H2
178 to differing degrees. [Garellek et al. \(2013\)](#) systematically manipulated source H1-H2 (as
179 well as other spectral tilt measures) in White Hmong, a language with contrastive breathy
180 voice on a particular lexical tone. They found that, controlling for all other param-
181 eters, an increase in source H1-H2 or source H2-H4 led to more “breathy tone” responses.
182 Finally, studies of the perceptual sensitivity to the harmonic source spectrum have shown
183 that listeners of various languages are sensitive to changes in source H1-H2, though this
184 varies by language ([Kreiman and Gerratt, 2010](#); [Kreiman et al., 2010](#)); the just-noticeable
185 differences for source H1-H2 were comparable to those for source H2-H4 and source H4-2
186 kHz, suggesting that listeners are particularly sensitive to changes in the lower-frequency
187 harmonic source spectral slope ([Garellek et al., 2016b](#)).

188 C. Filter correction and H1*–H2*

189 In the audio output, harmonic amplitudes from the voice source are significantly affected
190 by the filter function. This leads to a problem with using uncorrected H1–H2, especially when
191 comparing tokens with different vowel qualities: if uncorrected H1–H2 for [i] = H1–H2 for [a],
192 is this because the two vowels' H1–H2 values are the same at the voice source, or could it be
193 because the filtering effects of the vocal tract have resulted in the same output uncorrected
194 H1–H2? Thus, the use of H1–H2 – indeed of all spectral tilt measures – becomes less
195 informative of phonatory quality when compared across different vowel categories, because
196 any effect of phonatory quality could be obscured by the influence of the filter. [Hanson](#)
197 (1995, 1997) proposed the formula $20 * \log_{10}[F1^2/F1^2 - f^2]$, where f refers to the harmonic
198 frequency that need a formant correction. The product of the formula is subtracted from
199 H1 and H2. The value after subtraction reflects the amplitude of H1 and H2 before its being
200 affected by a formant of similar frequency.

201 Until [Hanson \(1995, 1997\)](#), uncorrected H1–H2 was measured from the audio signal gen-
202 erally for tokens with low vowels of the same quality. Alternatively, some studies relied
203 on inverse filtering to subtract the effects of the vocal tract filter, thereby approximating
204 source H1–H2 (e.g. [Bickley, 1982](#); [Huffman, 1987](#)). Other studies (rightfully) avoided using
205 uncorrected H1-based measures because formant correction was not widely used at the time.
206 For example, in their study of tongue root contrasts in Maa, [Guion et al. \(2004\)](#) used A1–A2
207 (the difference in amplitude between F1 and F2) instead of uncorrected H1–H2 because the
208 vowels differed in quality, with some vowels of interest having low F1 (which would interfere

209 with H1 and H2 estimation); see discussion in Section 2.3.2 of that paper. Hanson’s correc-
210 tion and subsequent versions (e.g. [Iseli et al., 2007](#)) have enabled researchers to correct for
211 differences in formant frequency and bandwidths without the need of specialized equipment,
212 such as a Rothenberg mask, for inverse filtering. Today, corrected spectral tilt measures,
213 denoted with asterisks (as in “H1*–H2*”), are the norm in linguistic research on phonation
214 types because they enable researchers to compare tokens with differing formant structures,
215 even within vowel category ([Garellek, 2022](#)).

216 **D. The effectiveness of H1(–H2) in distinguishing phonatory qualities**

217 Studies have found that H1–H2, and sometimes also H1, is an effective measure at distin-
218 guishing differences in phonation between women and men, between contrastive phonation
219 types, and between lexical stress and phrasal accent. In this section, we review the pioneer-
220 ing studies in this area. After the study of [Fischer-Jørgensen \(1967\)](#) on Gujarati modal and
221 breathy vowels, [Ladefoged \(1981\)](#) conducted a preliminary analysis of breathy (“murmured”)
222 vs modal (“clear”) vowels in !Xóõ, in which he found that the amplitude of uncorrected H1
223 was higher for breathy vowels. [Bickley \(1982\)](#) measured uncorrected H1–H2 from spectra
224 of breathy and modal vowels produced by ten speakers of !Xóõ (based on recordings made
225 by Tony Traill and Peter Ladefoged) and by four speakers of Gujarati. She found that, in
226 both languages, breathy vowels consistently had higher H1–H2 than modal vowels. There
227 were also large cross-speaker and between-language variations in these H1–H2 comparisons,
228 leading to the important assumption that phonation differences should be measured not in
229 absolute but instead in relative terms, ignoring the raw values of H1–H2. We note here

230 that it remains an open question whether absolute values are informative for H1–H2, but
231 they likely can be. (For example, while voice onset time (VOT) is also compared relatively
232 – voiceless unaspirated stops have lower VOT than voiceless aspirated ones – it is also the
233 case that we generally don’t expect unaspirated stops to have a VOT greater than about 35
234 ms; see discussion by [Cho and Ladefoged \(1999\)](#) and [Chodroff *et al.* \(2019\)](#).) Finally, Bick-
235 ley called attention to the presence, though variable, of increased spectral noise for breathy
236 vowels. We now know that spectral noise is a very important component to distinguishing
237 phonation types ([Garellek, 2019](#); [Gordon and Ladefoged, 2001](#)).

238 [Maddieson and Ladefoged \(1985\)](#) showed that, in four Tibeto-Burman languages with
239 “tense” (more constricted) vs. “lax” (breathier) vowels (Hani, Jingpho, Yi, and Wa), lax
240 (breathier) phonation had higher uncorrected H1–H2 than “tense” (more constricted) vowels.
241 This was perhaps the first journal article making use of H1–H2 to characterize phonation
242 types that are more constricted than modal voice. Although earlier work did investigate
243 spectral tilt differences between modal and constricted phonation types, we are not aware
244 of a previous study that specifically measured H1–H2; for example, the investigation by
245 [Ladefoged \(1983\)](#) and [Kirk *et al.* \(1984\)](#) of breathy, modal, and laryngealized vowels in !Xóõ
246 and Mazatec measured H1–A1.

247 In what may have been the first study to measure H1–H2 for consonants, [Traill and](#)
248 [Jackson \(1988\)](#) investigated the acoustic differences between breathy and modal nasals in
249 Tsonga, and their effects on following vowels. They measured uncorrected H1–H2, as well
250 as other spectral tilt measures, during the nasal consonant as well as the vowel onset.
251 Generally they found large differences within and across speaker gender on all spectral tilt

252 measures, but reported that breathy vs. modal nasals were more effectively distinguished
253 by two higher spectral tilt measures (the difference in slope between H1 and the harmonic
254 nearest 1400 Hz, and between H1 and the strongest harmonic above 2000 Hz) than by
255 H1–H2 (cf. more recent discussion of the acoustics of breathiness and nasality by [Garellek](#)
256 *et al.* 2016a; [Simpson](#) 2012; [Styler](#) 2017; [Tabain et al.](#) 2022). Few differences in H1–H2 at
257 vowel onset were found. Since the 1980s, H1–H2 has been used to quantify sex/gender-based
258 differences in phonatory quality. For example, in their investigation of male vs. female voice
259 differences among speakers of British English, [Henton and Bladon](#) (1985) used uncorrected
260 H1–H2 to measure whether female speakers of two dialects (Received Pronunciation and
261 Modified Northern British English) differ from male speakers in terms of voice quality.
262 They found that, in both dialects, H1–H2 was higher for women than for men. [Hanson and](#)
263 [Chuang](#) (1999) compared the spectral tilt in the production by male speakers with the data
264 of female speakers collected from [Hanson](#) (1995, 1997), and found that female speakers had
265 higher (by about 3dB) and larger standard deviation of H1*–H2* than male speakers. They
266 suggested that such differences in H1*–H2* indicated that, in terms of voice articulation,
267 female speakers have a larger OQ than male speakers.

268 In the 1990s there was also work investigating spectral tilt as a correlate of lexical stress
269 and phrasal accent and other phonological contrasts. [Sluijter and van Heuven](#) (1996) were
270 perhaps the first to investigate this (for Dutch), but they measured energy in four frequency
271 bands and not H1–H2. In their extension of that work, [Campbell and Beckman](#) (1997)
272 measured uncorrected H1–H2 (labeled there as “H2–H1”) but did not find it to be a reliable
273 correlate of lexical or phrasal prominence. They suggested (p. 70) that this might be due

274 to the fact that two (of four) speakers produced the low intonation tone with creaky voice.
275 Subsequent work on other languages has confirmed that H1–H2 can indeed correlate with
276 lexical/post-lexical prominence (Caballero and Carroll, 2015; Garellek and White, 2015;
277 Guion *et al.*, 2010).

278 E. Issues with H1–H2

279 The association between H1–H2 and glottal OQ, as well as the fact that listeners readily
280 use H1–H2 to perceive changes in voice quality (Esposito, 2010; Garellek *et al.*, 2013, 2016b;
281 Kreiman and Gerratt, 2010; Kreiman *et al.*, 2010), have contributed to the popularity of
282 H1–H2 as a measure of phonatory quality and voice quality more broadly (Garellek, 2022).
283 However, using H2 to normalize for SPL is theoretically arbitrary. Many studies also rely
284 on other landmarks; see Garellek (2019) for an overview. Further, Sundberg (2022) found
285 that H1 and H2 (denoted there as L_1 and L_2) were affected differently by the influence of
286 subglottal pressure: H2 is more sensitive to the pulse amplitude increase than H1, resulting
287 in an inverse relationship between H1–H2 and pulse amplitude. Yet as we mentioned earlier,
288 the use of H2 to normalize the SPL of the signal is based on the assumption that SPL affects
289 H1 and H2 equally. The conclusions from Sundberg (2022) thus provide more support for
290 avoiding the use of H2 as a normalizing landmark.

291 Compared with H1, measuring H1–H2 is also likely to be more prone to error propagation;
292 that is, the transferring of uncertainties in the input variable(s) to the output variable(s)
293 (Arras, 1998). H1–H2 involves the calculation of two measures – H1 and H2 – whereas H1
294 involves only one. The correct estimation of H1–H2 thus requires the correct estimation

295 of both H1 and H2. Error propagation in H1–H2 can be attributed to two sources. First,
296 H2 is estimated based on H1, which requires accurate estimation of f_0 . Yet H1–H2 is often
297 used to measure non-modal voices. Occasionally the aspiration noise in breathy voice, and
298 frequently the decrease in periodicity in creaky voice, can make the correct tracking of f_0
299 difficult or impossible (Garellek, 2019; Keating *et al.*, 2015; Kuang, 2017). Any error in the
300 estimation of f_0 will influence both H1 and H2, leading to more calculation error than if H1
301 alone were estimated.

302 Another issue arises in the application of the filter correction for H1*–H2*. Although
303 corrected H1*–H2* is now the norm when measuring H1–H2 across vowels qualities, errors
304 in estimating formant frequencies and amplitude inevitably arise. For example, when a token
305 has a high f_0 and a low F1 (i.e. for a high vowel), it is possible for the tracking algorithm
306 to mistake f_0 for F1 and F1 for F2. Another common error for formant estimation occurs
307 when F1 and F2 are similar in frequency (e.g. for back vowels); in such cases, F1 and F2 are
308 likely to be mistaken for a single formant, with the real F3 consequently being mistaken for
309 F2. When formant frequencies and amplitudes are thus miscalculated, the corrected H1*
310 and H2* are highly likely to be erroneous as well. We hypothesize that H1* is less likely to
311 have such errors than H1*–H2* because erroneous formant tracking will influence H1* only
312 once, but will affect H1*–H2* twice– when estimating each harmonic level.

313 A further issue arises with (even slight) nasality: the first nasal pole (P0) can increase the
314 amplitude of either H1 or H2, depending on the f_0 (Dang and Honda, 1996; Simpson, 2012;
315 Styler, 2017). Nasal zeroes will further attenuate the oral resonances (Dang and Honda,
316 1996; Simpson, 2012; Styler, 2017). P0 is usually in the range of 200–450 Hz (Styler, 2017),

317 so for typical adult male speakers with an f_0 of 120 Hz, P0 is more likely to influence H2 (240
318 Hz) and H3 (360 Hz). But for typical adult female speakers with an f_0 above 200 Hz, P0 is
319 more likely to influence H1 (Simpson, 2012). As a result, when a token contains nasalization,
320 H1–H2 will inevitably be influenced by P0 in unpredictable, f_0 -dependent ways. Admittedly,
321 measuring H1 alone does not fully avoid these issues. We advise then that f_0 should be used
322 as a control variable when analyzing either H1 or H1–H2.

323 II. ADDRESSING ISSUES WITH H1(–H2)

324 A. Meta-analysis of H1(–H2) in studies of linguistic phonation type

325 To review the effectiveness of H1–H2 at distinguishing phonation types, we conducted
326 a meta-analysis of studies that compare H1(–H2) between contrastive phonation types in
327 a given language. We focus here on contrastive phonation types, because we expect the
328 H1(–H2) differences to be relatively large and consistent across speakers. In addition to
329 H1–H2, we also include studies that measure H1, a measure that is closely related to the
330 new measure that we elaborate on below.

331 Our survey focused on journal articles, particularly from *Journal of Phonetics*, *Journal*
332 *of the International Phonetic Association*, and *JASA*, though we also included some the-
333 ses that focused on linguistic phonation type. The earliest study was published in 1985.
334 We include data from 39 languages and 76 comparisons of contrastive phonation types
335 (e.g. breathy vs modal vowels). The languages in the survey come from several families
336 (including Otomanguean, Mayan, Indo-European, Kx’a, Taa, Niger-Congo, Austroasiatic,

337 Hmong-Mien, and Sino-Tibetan) spoken in various parts of the world, but especially from
338 Mesoamerica, Southern Africa, and Southeast Asia, where phonation contrasts are more
339 prevalent.

340 The spectral measures used in these comparisons include uncorrected H1–H2, H1*–H2*,
341 and H1*. Most studies published before c. 2010 included uncorrected measures that were es-
342 timated manually, whereas those published after the advent of Praat scripts and VoiceSauce
343 (Shue *et al.*, 2011) generally include corrected measures that were estimated automatically.
344 Of the 76 phonation type comparisons, there are 15 creaky (i.e. more constricted) vs. modal
345 and 35 breathy vs. modal comparisons that included a quantitative analysis of whether
346 the differences between these phonation types were statistically significant. The languages
347 included in the survey for each measure and for the comparisons of breathy and creaky
348 vs. modal phonations are listed in Table I. We summarize the number of comparisons
349 that showed significant differences, partially-significant differences, and non-significant dif-
350 ferences for each contrast and each measure in Table II. The detailed results of the survey,
351 including the language names and the corresponding references, are in supplementary mate-
352 rial S1, available at <https://doi.org/10.17605/OSF.IO/QGBKA>. We define a difference as
353 “significant” when the p value of the comparison is smaller than 0.05. We define a difference
354 as being “partially significant” either when the p value of the comparison is between 0.05
355 and 0.1, or when a significant difference was found only for a subgroup of the speakers or
356 in a subset of the stimuli. We define a difference as “non-significant” when the p value of
357 the comparison was larger than 0.1 for all speakers and all stimuli. Generally, the results in
358 Table II show that, for both breathy–modal and creaky–modal contrasts, the majority show

359 significant differences in either uncorrected H1–H2 or H1*–H2*. This implies that H1–H2 is
360 indeed a robust index of phonation differences, and supports the findings of Kreiman *et al.*
361 (2007) that source H1–H2 is resistant to measurement artifacts. However, there exist cases
362 in which H1–H2 does not distinguish contrastive phonation types: uncorrected H1–H2 did
363 not distinguish breathy from modal phonation in Mon or in Tamang in Esposito (2006);
364 uncorrected H1–H2 did not distinguish creaky from modal phonation in Mpi or in Jalapa
365 Mazatec in Pennington (2005) (cf. Blankenship 2002); and H1*–H2* did not distinguish
366 creaky from modal phonation in White Hmong in Esposito (2012) (cf. Garellek 2012). And
367 while there are fewer studies that use H1–H2 for creaky vs modal comparisons, there are
368 relatively more cases where H1–H2 does not significantly distinguish creaky vs. modal vowels
369 than cases where the measure does not distinguish breathy vs. modal ones (3/15 vs. 2/35).
370 This indicates that H1–H2 may be more effective at capturing breathiness than creakiness,
371 when these phonation types are compared to modal voice. Table II also shows that H1* is
372 rarely used as a measure to distinguish phonation types. To our knowledge, the only existing
373 study that made a quantitative comparison of H1* between contrastive phonation types is
374 Esposito (2012). We therefore need more data to test the effectiveness of H1* as a measure
375 of phonatory quality.

376 B. Error simulations

377 Next we verify our hypothesis that H1–H2 is more prone to error propagation than
378 H1, particularly when these measures are corrected as H1*–H2* and H1*. We created
379 simulations of two circumstances: when f_0 is wrongly estimated and when formants are

TABLE I. Languages in the survey of spectral differences between phonation types

Contrast	Measure	Language
Breathy vs. Modal	H1–H2	Suai; Ju ’hoansi; Jalapa Mazatec; White Hmong; Krathing Chong; Shanghainese; Ningbo Wu; Changyinsha Wu; Wenzhou Wu; Green Mong; Takhian Thong Chong; !Xóõ; Fuzhou Min; Green Mong; SADV Zapotec; SLQ Zapotec; Tlacolula Zapotec; Gujarati; Tsonga; Mon; Tamang
	H1*–H2*	Jalapa Mazatec; Gujarati; Chichimec; White Hmong; !Xóõ; Black Miao; Khmer; Shanghainese; Green Mong; Chrau
	H1*	White Hmong
Creaky vs. Modal	H1–H2	Ju ’hoansi; Coatzospan Mixtec; Jalapa Mazatec; Takhian Thong Chong; Green Mong; Mpi
	H1*–H2*	Jalapa Mazatec; White Hmong; Chichimec; !Xóõ; Black Miao; Green Mong
	H1*	White Hmong

³⁸⁰ wrongly estimated. We then determined how uncorrected H1, H1*, uncorrected H1–H2,

TABLE II. Summary of the survey results of spectral differences between phonation types; The numbers represent the number of studies in each category.

Contrast	Measure	Significant	Partially significant	Non-significant
Breathy–Modal	H1–H2	19	2	2
	H1*–H2*	8	3	0
	H1*	1	0	0
Creaky–Modal	H1–H2	4	1	2
	H1*–H2*	5	1	1
	H1*	1	0	0

381 and H1*–H2* are affected in both circumstances. [See also [Simpson 2012](#) for a simulation
382 of how nasality affects H1–H2.]

383 We synthesized a token of [u] as the base token using the Klatt synthesizer ([Klatt and](#)
384 [Klatt, 1990](#)). The values of f0 and formants are shown in [Table III](#). The segment duration,
385 bandwidth fraction, and formant frequency interval are 0.4s, 0.05, and 1000 Hz. We manually
386 entered the correct values of f0, F1, F2, and F3 for the base token [u] (as in [Table III](#)) in
387 VoiceSauce ([Shue et al., 2011](#)) and used those values to estimate the values of H1, H2, and
388 H1–H2. To simulate cases where f0 is mistracked, we manually changed the f0 between
389 180 Hz to 300 Hz in six 20-Hz increments and then re-estimated the same spectral energy
390 values. The results of the six f0 conditions are shown in [Figure 1](#). We see that, when f0

391 is mistracked, neither H1 nor H1–H2 consistently outperforms the other in terms of being
 392 closer to the true spectral value. However, there is a tendency for H1–H2 to have a larger
 393 deviation than H1 when f_0 is incorrectly estimated. Table IV lists the absolute deviation
 394 of the estimated value from the true value, for different mistracked f_0 values. The mean
 395 deviation for H1–H2 is nearly twice that for H1, for both corrected and uncorrected values
 396 of these measures. Therefore, when f_0 is incorrectly estimated (as is common during creaky
 397 voice), H1 appears more resistant to error than H1–H2.

TABLE III. Formant frequencies and bandwidths (in Hz) for synthesized [u].

f_0	F1	B1	F2	B2	F3	B3	F4
240	453	50	944	18	2899	593	3778

398 Next, we stimulated two common formant tracking errors for [u] using the same stimuli
 399 and method as the f_0 error demonstration. The first error type is when f_0 is mistaken as F1,
 400 and F1 as F2. The second type is when F1 and F2 of [u] are mistracked as just one formant
 401 (F1), and F3 is mistaken as F2. We used the mean of F1 and F2 as “mistracked F1” for the
 402 second type of error. All other parameters were held constant when comparing those two
 403 scenarios with the correctly estimated values. The formant values used to illustrate formant
 404 tracking errors are presented in Table V. As with the f_0 manipulation, we manually entered
 405 the values of F1, F2, and F3 in Table V into VoiceSauce and calculated the corrected and
 406 uncorrected H1 and H1–H2 values for the different conditions. The results are presented in
 407 Figure 2. Uncorrected values of H1, H2, and H1–H2 did not change, as expected. But for

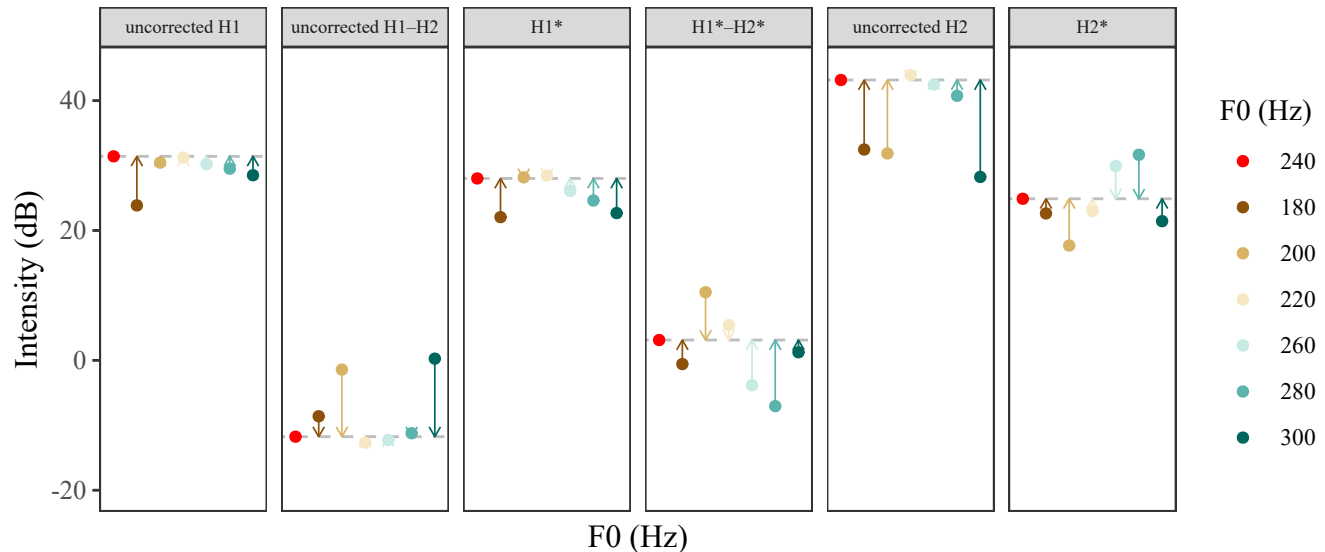


FIG. 1. (color online) Spectral energy with different f_0 estimations. The leftmost red dot of each panel and the dotted line represent the true spectral (slope) value when f_0 is correctly estimated at 240 Hz. The arrows represent the deviation between the true spectral value and the estimated spectral value when f_0 is wrongly estimated, as is the case for all f_0 values not equal to 240 Hz.

408 both conditions with formant tracking errors, $H1^*$ shows a smaller deviation from the true
 409 value than $H1^*-H2^*$. The mean deviation of $H1^*-H2^*$ is nearly three times that of $H1^*$.
 410 Summary statistics appear in Table VI.

411 C. Current proposal for estimating Residual H1

412 As an alternative to measuring $H1-H2$, we propose factoring out the effect of root-mean-
 413 squared (RMS) energy (henceforth referred as Energy) from $H1$, whether uncorrected $H1$
 414 or $H1^*$ corrected for formant frequencies and bandwidths. We call this “residual $H1$.” Of
 415 course, differences in recording conditions across speakers and studies will affect energy. We

TABLE IV. Deviation (Δ) of estimated uncorrected and corrected H1 and H1–H2 from their true values, for various miscalculated f_0 s

f_0	$\Delta_{\text{uncorrected H1}}$	$\Delta_{\text{uncorrected H1–H2}}$	Δ_{H1^*}	$\Delta_{\text{H1}^*-\text{H2}^*}$
180	7.559	3.146	5.956	3.701
200	0.974	10.331	0.166	7.368
220	0.184	0.933	0.425	2.304
260	1.220	0.515	1.918	6.956
280	1.901	0.535	3.400	10.178
300	2.911	12.018	5.334	1.884
Mean	2.458	4.580	2.867	5.399

TABLE V. Formant values for two different types of formant tracking errors

Condition	F1	F2	F3
True	453	944	2899
f_0 taken as F1	240	453	944
F1 & F2 collapsed into F1	699	2899	3778

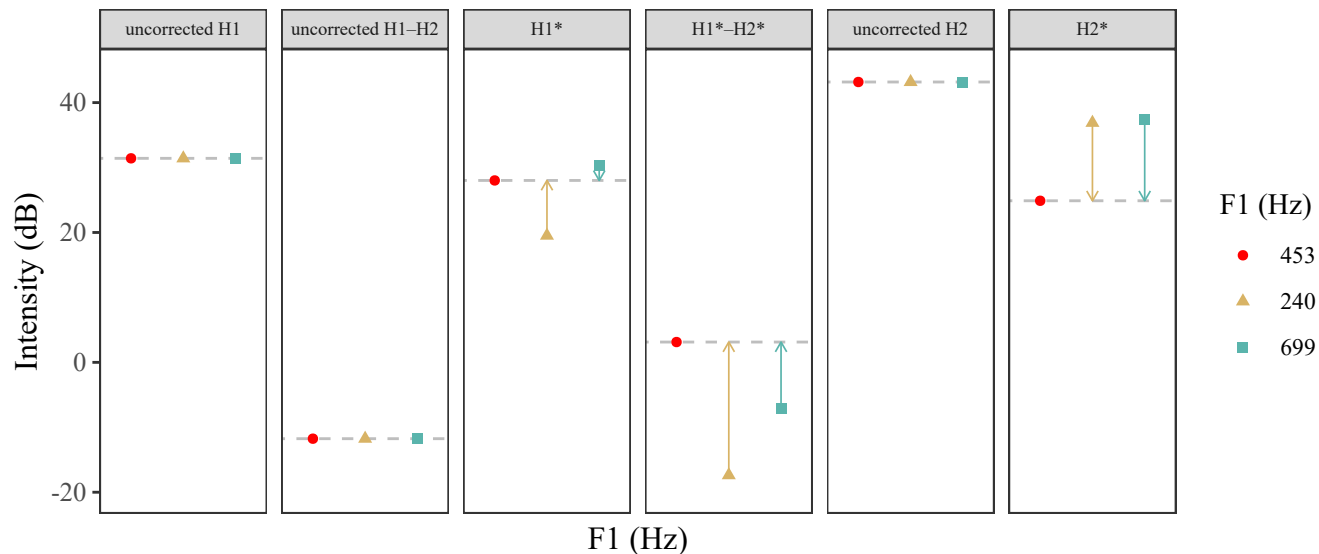


FIG. 2. (color online) Spectral energy with different formant estimations. The leftmost red dot of each panel and the dotted line represent the true spectral value when formants are correctly estimated. The arrows represent the deviation between the true spectral value and the estimated spectral value when formants are wrongly estimated. Only F1 values are listed in the legend. The corresponding F2 and F3 values can be found in Table V.

416 don't view this as problematic for residual H1, because such differences also affect H1: for
 417 example, a quieter signal will result in a lower RMS energy as well as a lower H1. And as
 418 we show below, residual H1 controls for the energy of an individual token on the H1 value
 419 on that same token.

420 Residual H1 avoids some of the issues raised in Section II B regarding the estimation
 421 of H1-H2. In practice, precisely how we control for the effect of energy on H1 varies de-
 422 pending on whether when H1 is a dependent variable or an independent variable. When
 423 H1 is a dependent variable, energy can be added to the model as a covariate: i.e. $H1 \sim$

TABLE VI. Deviation (Δ) of estimated $H1^*$ and $H1^*-H2^*$ from their true values, for two types of mistracking

Condition	F1 (Hz)	$\Delta H1^*$ (dB)	$\Delta H1^*-H2^*$ (dB)
f0 taken as F1	240	8.496	20.503
F1 & F2 collapsed into F1	699	2.261	10.293
Mean		5.379	15.398

424 *main factor(s) + Energy*. By adding Energy as a covariate, the effect of SPL on H1 is con-
 425 trolled for, and the coefficients of the main factors reflect the independent effects of those
 426 factors on H1.

When H1 is an independent variable, the effect of Energy on H1 should be calculated first, and then subtracted from H1, as shown in (1) and (2):

$$\text{Step 1: Get the coefficient of energy } (b_1): H1 \sim b_0 + b_1 * \text{Energy} (\text{logged}) \quad (1)$$

$$\text{Step 2: Calculate Residual H1: } \text{Residual H1} = H1 - b_1 * \text{Energy} (\text{logged}) \quad (2)$$

427 In step 1, the coefficient of energy in Model (1) represents how strongly H1 is correlated with
 428 energy in a given token. Energy is first log-transformed because it is bounded at zero at the
 429 lower end and unbounded at the upper end. In step 2, we multiply the coefficient of energy
 430 with the actual value of energy. We then subtract the product from H1. The residual of H1
 431 after subtraction represents the value of H1 after controlling for the SPL of the recordings.

432 If the data come from multiple speakers, H1 and logged energy can be transformed to z-score
433 to reduce inter-speaker variation.

434 Compared with normalizing against H2, using energy has certain advantages. First, H1
435 requires an accurate estimation of only one spectral value, whereas H1–H2 requires two. H1 is
436 thus less likely to be affected by error propagation than H1–H2. Second, H1 is more resistant
437 to the influence of nasalization (Simpson, 2012). Thus, we hypothesize that H1 normalized
438 for energy (i.e. residual H1) should reflect the degree of constriction or breathiness to the
439 same extent as H1–H2, only with less variability. In Section III, we test the relation between
440 H1 and OQ, to investigate whether Residual H1 has an articulatory basis of vocal fold
441 constriction. We also use two case studies to compare residual H1 with H1–H2 in terms of
442 their effectiveness of representing changes in phonatory quality.

443 III. CONTACT QUOTIENT IN RELATION TO RESIDUAL H1* VS. H1*–H2*

444 Previous work has shown that there is a positive, if sometimes weak and nonlinear,
445 relationship between H1*–H2* and OQ (Kreiman *et al.*, 2012; Samlan *et al.*, 2013). In
446 this section, we test whether residual H1* has a similar or better correlation with OQ, as
447 indexed by electroglottographic CQ, than H1*–H2*. We used data from the “Production and
448 Perception of Linguistic Voice Quality” project at UCLA ([http://www.phonetics.ucla.](http://www.phonetics.ucla.edu/voiceproject/voice.html)
449 [edu/voiceproject/voice.html](http://www.phonetics.ucla.edu/voiceproject/voice.html)). The data and R code for data processing and analysis
450 are available in supplementary material S3 at <https://doi.org/10.17605/OSF.IO/QGBKA>.

451 The corpus includes data from eight languages, 68 speakers, and 9,101 words in total
452 after exclusions, (see summary in Table VII).¹ Each word was measured by nine equal time

TABLE VII. Language data from the UCLA Voice Project used to assess the relationship between CQ and Residual $H1^*$ vs. $H1^*-H2^*$.

Language	Family	Speakers	Phonation types
Bo	Sino-Tibetan	6 (3 F, 3 M)	Tense, Lax
Gujarati	Indo-European	10 (7 F, 3 M)	Breathy, Modal
Luchun Hani	Sino-Tibetan	9 (4 F, 5 M)	Tense, Lax
White Hmong	Hmong-Mien	11 (2 F, 9 M)	Breathy, Modal, Creaky tones
Mandarin	Sino-Tibetan	11 (5 F, 6 M)	Modal, Creaky tones
Black Miao	Hmong-Mien	8 (0 F, 8 M)	Breathy, Modal, Creaky tones
Southern Yi	Sino-Tibetan	7 (4 F, 3 M)	Tense, Lax
Zapotec	Otomanguean	6 (2 F, 4 M)	Breathy, Modal, Creaky

453 intervals, resulting in 81,909 data points in total. This data set included acoustic data of
454 $H1^*-H2^*$ and $H1^*$, calculated using VoiceSauce (Shue *et al.*, 2011). The $H1^*$, $H1^*-H2^*$, and
455 f_0 values were z-scored by speaker to reduce the variation between speakers. Tokens with
456 an absolute z-score value larger than 3 were considered as outliers and were excluded from
457 analyses. Within each vowel category, we calculated the Mahalanobis distance on the F1-F2
458 panel. For tokens with a Mahalanobis distance larger than 6, we regarded their formant
459 values as outliers, similar to what has been done in our previous work (Chai and Ye, 2022;
460 Garellek and Esposito, 2021; Seyfarth and Garellek, 2018). Time points whose f_0 , F1, or F2

461 values were outliers were also excluded from $H1^*$ and $H1^*-H2^*$ analyses, because $H1^*$ and
462 $H1^*-H2^*$ are calculated based on f_0 , $F1$, and $F2$. For energy, we first excluded tokens with
463 a value of zero, then log-transformed to normalize its right-skewed distribution, and then
464 z-scored the logged energy and excluded tokens with a z-score larger than 3. The outlier
465 detection process for the acoustic measures is the same for the following case studies in
466 Sections IV A and IV B.

467 The glottal open quotient was estimated using electroglottographic (EGG) CQ calculated
468 using the hybrid method (Howard, 1995; Orlikoff, 1991).² For CQ, there were 6,943 points
469 with a value of zero; these were first excluded. The remaining CQ values were then z-scored
470 and those with a z-score larger than 3 were considered outliers and therefore were excluded.
471 After outlier exclusion, there were 76,196 valid data points for $H1^*$, 76,570 for $H1^*-H2^*$,
472 81,222 for f_0 , 73,363 for CQ, and 81,473 for energy.

473 To assess the relationship between CQ and $H1^*-H2^*$ and residual $H1^*$, we regressed CQ
474 on both $H1^*-H2^*$ and residual $H1^*$, as in Models (3) and (4). Since $H1^*$ was an independent
475 variable in the model, we factored out the effect of energy from $H1^*$ and calculated residual
476 $H1^*$ using Equations (1) and (2). The coefficient of energy on $H1^*$ was 0.682. The statistics
477 are shown in Table VIII. We use R^2 to represent the effect size of the models. The R^2
478 value is defined as the percentage of variance of the dependent variable that is explained
479 by the independent variables in the model. For linear mixed-effect models, we calculate the
480 **marginal** R^2 of the model, which is defined as the percentage of variance of the dependent
481 variable that is explained by the **fixed** variables in the model (Johnson, 2014). The R^2
482 (for linear models) and marginal R^2 (for linear mixed-effect models) are calculated using

483 the *multilevelTools* package (Wiley, 2020) in *R*. A scatter plot of all data points and the
 484 correlation line between CQ and the two spectral tilt measures are shown in Figure 3.

485 As the results show, residual H1* had a larger absolute coefficient, slightly higher standard
 486 error (0.005 vs. 0.003), and higher t-value than H1*-H2*. Model 4 with H1* as the predictor
 487 had a higher marginal R^2 value than Model 3 with H1*-H2* as the predictor (0.102 vs.
 488 0.060), suggesting that H1* can explain more variance of CQ than H1*-H2*. Figure 3
 489 illustrates that the regression line for residual H1* is steeper than that for H1*-H2*. This
 490 indicates that H1*, after controlling for energy, has a stronger correlation with CQ than
 491 H1*-H2*. By extension, this also confirms the articulatory basis of H1* as an acoustic
 492 correlate of vocal fold approximation.

$$CQ \sim H1^* - H2^* + (1|Speaker) \quad (3)$$

$$CQ \sim Residual\ H1^* + (1|Speaker) \quad (4)$$

TABLE VIII. Correlation between CQ and H1*-H2* and H1*

Model	β	Std. Error	t value	p	R^2
CQ \sim H1*-H2*	-0.230	0.003	-66.090	< .001	0.060
CQ \sim Residual H1*	-0.449	0.005	-88.700	< .001	0.102

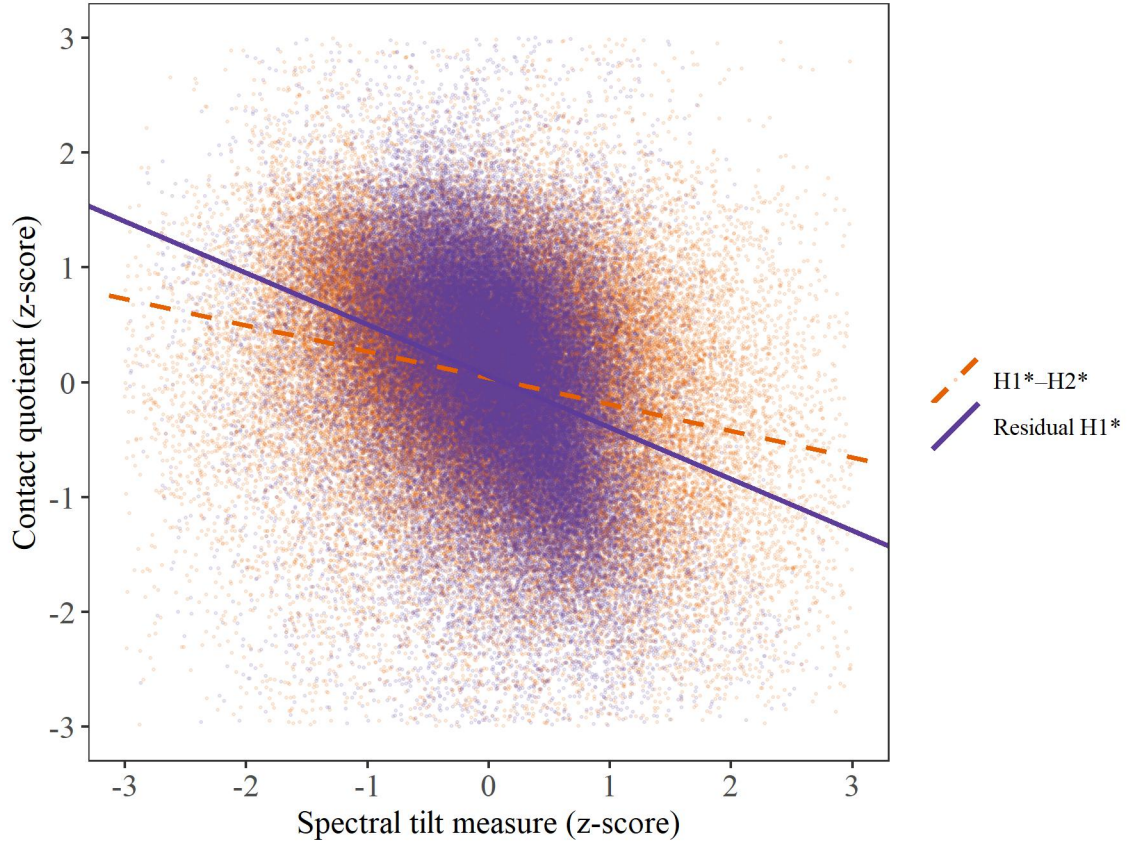


FIG. 3. (color online) Relationship between CQ and $H1^*-H2^*$ and residual $H1^*$. The CQ values have had the random intercept of subjects subtracted. Regression lines were based on results from Model (3) and (4).

493 IV. CASE STUDIES

494 In this section, we provide two case studies where we compare residual $H1^*$ to $H1^*-$
 495 $H2^*$ and their ability to track changes in phonation in two languages, !Xóõ and Mandarin.
 496 We use datasets that have previously been analyzed for phonation: Garellek (2020) on !Xóõ
 497 phonation types and Chai (2019, 2021) on Mandarin utterance-level changes in voice quality.
 498 In neither paper did we look specifically at $H1$, so for these case studies we were particularly

499 interested in seeing if the phonation differences are better differentiated acoustically using
500 residual $H1^*$ instead of $H1^*-H2^*$.

501 **A. Phonation types in !Xóõ (Taa)**

502 *1. Corpus*

503 !Xóõ (also known as Taa) is a Tuu language spoken in Botswana, whose phonation types
504 have recently been analyzed acoustically by Garellek (2020). That study only measured
505 $H1^*-H2^*$; here, we compare $H1^*$ (with energy as a covariate) to $H1^*-H2^*$ for three of the
506 phonation types: breathy, modal, and creaky.

507 The recordings are of the East !Xóõ dialect, and were made in the late 1970s by Peter
508 Ladefoged and Tony Traill. They are available for download from the UCLA Phonetics Lab
509 Archive at <http://archive.phonetics.ucla.edu/Language/NMN/nmn.html>. We used the
510 same data as (Garellek, 2020), and thus we refer the reader to that source for details on
511 the segmentation criteria and data segmentation procedures. All the words had /a/ vowels
512 (which varied considerably in phonetic quality due to coarticulation), and were produced
513 by ten speakers of !Xóõ. The corpus had 369 words, containing six phonation types. We
514 only compared the spectral values of breathy, modal, and creaky phonations, resulting in
515 175 tokens for analysis (breathy: 83; modal: 54; creaky: 38). The word list of the stimuli
516 is in supplementary material S2 at <https://doi.org/10.17605/OSF.IO/QGBKA>. The data
517 and R code for data processing and analysis are available in supplementary material S3 at
518 the same URL as S2.

519 The acoustic measures of the recordings were calculated using VoiceSauce every mil-
520 lisecond. Each token was divided into nine equal intervals. The mean of each interval was
521 calculated. We used nine points to represent each token such that the duration of the tokens
522 was normalized. In total, 1,575 data points were measured (175 tokens * 9 time points). The
523 outlier detection method is the same as described in III. After the outlier exclusion, there
524 were 1,351 valid data points for H1*, 1,363 for H1*–H2*, 1,507 for f0, 1,380 for formants,
525 and 1,492 for energy. We calculated the mean acoustic values for each individual words for
526 the statistical analysis in Section IV A 2.

527 2. Results

To compare how effectively H1* and H1*–H2* differentiate the three phonation types,
we regressed both H1* and H1*–H2* on phonation type. The model in which H1* was the
dependent variable also had energy as an independent variable to control for the SPL. The
models for H1*–H2* and H1* are in (5) and (6):

$$H1^* - H2^* \sim \textit{Phonation} \tag{5}$$

$$H1^* \sim \textit{Phonation} + \textit{Energy}(\textit{logged}) + (1|\textit{Speaker}) \tag{6}$$

528 For each model, the modal phonation was set as the baseline for comparison. Random
529 intercepts or slopes by speaker were not included for Model (5), because they resulted in
530 singular fits and did not improve the model. We compared the effectiveness of H1* and
531 H1*–H2* in differentiating creaky and breathy phonation types from modal phonation by
532 looking at the estimate coefficient, standard error, and t-value (estimate/standard error) of

533 the phonation variable. A higher coefficient means that two phonation types have a larger
534 difference in the acoustic measure. A lower standard error indicates that the values of the
535 acoustic measure of each phonation group are less variable. A higher t-value represents a
536 relatively high coefficient and a relatively low standard error, indicating a better separation
537 between two phonation groups.

538 The statistics of Models (5) and (6) are shown in Table IX. For the differentiation between
539 creaky and modal phonation, the model with $H1^*$ as the dependent variable and energy as
540 the covariate had a higher estimate of coefficient, lower standard error, and higher t-value
541 for creaky vs. modal comparison than the model with $H1^*-H2^*$ as the dependent variable.
542 This indicates that $H1^*$ (after controlling for energy) is better at distinguishing creaky from
543 modal phonation than $H1^*-H2^*$.

544 When comparing breathy and modal phonation, $H1^*-H2^*$ behaved similarly to $H1^*$.
545 $H1^*-H2^*$ had a higher coefficient estimate and standard error than did $H1^*$. The t-value of
546 $H1^*-H2^*$ was similar to that of $H1^*$ (10.974 vs. 11.085), whereas both models had p-values
547 smaller than 0.001. Thus, in terms of distinguishing breathy from modal phonation, we
548 consider $H1^*-H2^*$ and $H1^*$ to be equally effective.

549 We also calculated the effect sizes of Models (5) and (6) using R^2 , as shown in Table
550 IX³. The marginal R^2 of Model (6) is higher than the R^2 of Model (5) (0.838 vs. 0.592),
551 indicating that the model with $H1^*$ as the dependent variable has a larger effect size (and
552 thus more variance is explained) than the model with $H1^*-H2^*$ as the dependent variable.

553 The distributions of $H1^*-H2^*$ and $H1^*$ for different phonation types in !Xóõ are shown in
554 Figure 4. For the $H1^*$ data in Figure 4 to show how $H1^*$ distinguished the three phonation

555 types after controlling for energy, the residual $H1^*$ was calculated by subtracting the product
556 of the coefficient of energy in Model (6) ($b = 0.606$) and the z-scored energy from the z-
557 scored $H1^*$ value. Comparing $H1^*-H2^*$ with $H1^*$ in Figure 4, we see that for all the three
558 phonation types, the $H1^*$ values are less variable within group, and there is less overlap
559 between modal and creaky phonation types in $H1^*$ than $H1^*-H2^*$. In sum, after controlling
560 for energy, $H1^*$ distinguished creaky phonation from modal phonation in !Xóõ better than
561 $H1^*-H2^*$, in terms of having a larger effect size and smaller standard errors. However, $H1^*$
562 and $H1^*-H2^*$ do not differ in the effectiveness of distinguishing breathy phonation from
563 modal phonation in !Xóõ.

TABLE IX. Model comparison between $H1^*-H2^*$ and $H1^*$ in distinguishing !Xóõ phonation types

Phonation contrast	Model	β	Std. Error	t value	p
Creaky–Modal	$H1^* - H2^* \sim \mathbf{Phonation}$	-0.462	0.141	-3.278	0.0013
	$H1^* \sim \mathbf{Phonation} + Energy$	-0.554	0.091	-6.069	< .001
Breathy–Modal	$H1^* - H2^* \sim \mathbf{Phonation}$	1.221	0.111	10.974	< .001
	$H1^* \sim \mathbf{Phonation} + Energy$	0.671	0.061	11.085	< .001

TABLE X. R^2 of Model (5) and marginal R^2 of Model (6)

Model		(Marginal) R^2
(5)	$H1^* - H2^* \sim \text{Phonation}$	0.592
(6)	$H1^* \sim \text{Phonation} + \text{Energy} + (1 \text{Speaker})$	0.838

564 B. Phrasing in Mandarin

565 1. Corpus

566 [Chai \(2019\)](#) found that the final position of declarative sentences in Mandarin had more
567 creak than non-final positions, after controlling for f0. They assumed then that vowels in
568 utterance-final position should be more constricted acoustically than non-final positions, but
569 did not find differences in H1*–H2* according to position. In a follow-up study, [Chai \(2021\)](#)
570 increased the sample size and found a correlation between low H1*–H2* and utterance-
571 final position in declaratives. This suggests that the discrepancy in findings between [Chai](#)
572 [\(2019\)](#) and [Chai \(2021\)](#) was due to noisiness in H1*–H2*, requiring a larger data set for
573 effects to emerge. In the present study, we aim to determine whether utterance-final creak
574 in Mandarin is indeed associated with vocal fold constriction, as measured by residual H1*
575 instead of H1*–H2*.

576 We combined the data sets from both [Chai \(2019\)](#) and [Chai \(2021\)](#). There were 823 target
577 declarative sentences produced by 64 Mandarin speakers. Phonologically identical words
578 were placed in the initial, medial, and final position of each sentence. The stimuli include

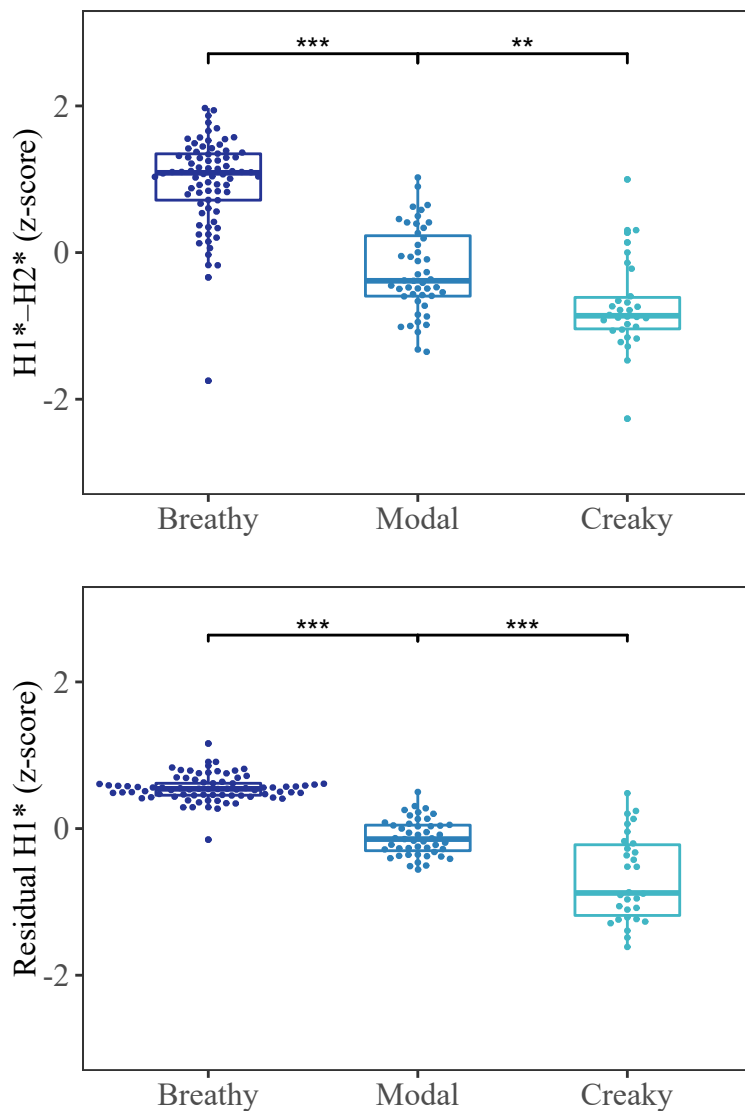


FIG. 4. (color online) $H1^* - H2^*$ (top) and Residual $H1^*$ (bottom) in different phonations in !Xóó

*** $p < .001$; ** $0.001 < p < 0.01$; * $0.01 < p < 0.05$

Residual $H1^* = H1^* - \text{Energy} * \text{Energy coefficient (0.606)} - \text{random intercept in Model (6)}$

579 1,889 target words in total (initial: 631; medial: 628; final: 630). The stimuli sentence list

580 is in supplementary material S2 at <https://doi.org/10.17605/OSF.IO/QGBKA>. The data

581 and R code for data processing and analysis are available in supplementary material S3 at
582 the same URL as S2.

583 The recordings were processed using VoiceSauce, which output a value for $H1^*$, $H1^*-H2^*$,
584 energy, and $f0$ every millisecond. Energy values were first log-transformed. All the acoustic
585 measurements were z-scored by speaker and word. 217,378 data points were generated in
586 total. The outlier detection procedure is the same as described in III. After outlier exclusions,
587 there were 200,505 valid data points for $H1^*$, 204,497 for $H1^*-H2^*$, 210,994 for formants,
588 213,601 for $f0$, and 214,592 for energy. We calculated the mean value for each individual
589 word for the statistical analysis in Section IV B 2.

590 2. Results

As with the !Xóõ case study, here two models were fit to test whether $H1^*-H2^*$ or $H1^*$
best distinguishes voice qualities associated with different utterance positions. The models
for $H1^*-H2^*$ and $H1^*$ were (7) and (8). Since Chai (2019) suggested that utterance-final
position was creakier than non-final positions after controlling for $f0$, $f0$ was added to Model
(7) and (8). The criteria of a better model were the same as the !Xóõ case study: larger
coefficient estimate, smaller standard error, and larger t-value.

$$H1^* - H2^* \sim Position + f0 + (1|Speaker) + (f0 + Position|Speaker) \quad (7)$$

$$H1^* \sim Position + Energy(\text{logged}) + f0 + (1|Speaker) + (f0 + Position|Speaker) \quad (8)$$

591 The statistics of Models (7) and (8) are shown in Table XI. In terms of distinguishing
592 initial position from final position, the coefficient of the position variable in the $H1^*$ model

593 was four times larger than in the H1*-H2* model. The standard errors of the two models
594 were similar. The t-value was higher in the H1* model than the H1*-H2* model. Similarly,
595 when distinguishing medial position from final position, the position variable in the H1*
596 model had larger coefficient, similar standard error, and higher t-value than in the H1*-H2*
597 model.

598 We also calculated the effect sizes of Models (7) and (8) using marginal R^2 of the fixed
599 variables, as shown in Table XII. The marginal R^2 of Model (8) is larger than that of Model
600 (7) (0.805 vs 0.254), indicating that the model with H1* as the dependent variable has a
601 larger effect size than the model with H1*-H2* as the dependent variable; that is, more
602 variance in the dependent variable is explained.

603 Figure 5 shows residual H1*-H2* and residual H1* in utterance-initial, medial, and final
604 position. The H1*-H2* distributions of the three positions are very similar, whereas in H1*
605 we find that final position has overall lower values than the non-final positions.

606 In sum, while H1*-H2* did not distinguish the three utterance positions in Chai 2019,
607 an effect of utterance position on H1*-H2* emerged after we increased the number of data
608 points and subjects by adding on data from a subsequent study, Chai 2021. This suggests
609 that the creakier voice quality of utterance-final position in Chai 2019 was indeed produced
610 with more constriction. The effect likely did not emerge in that original study due to a lack of
611 statistical power. In addition, the comparison between the H1* and H1*-H2* models reflects
612 the fact that H1* captures the difference in vocal fold constriction better than H1*-H2* and
613 requires less statistical power.

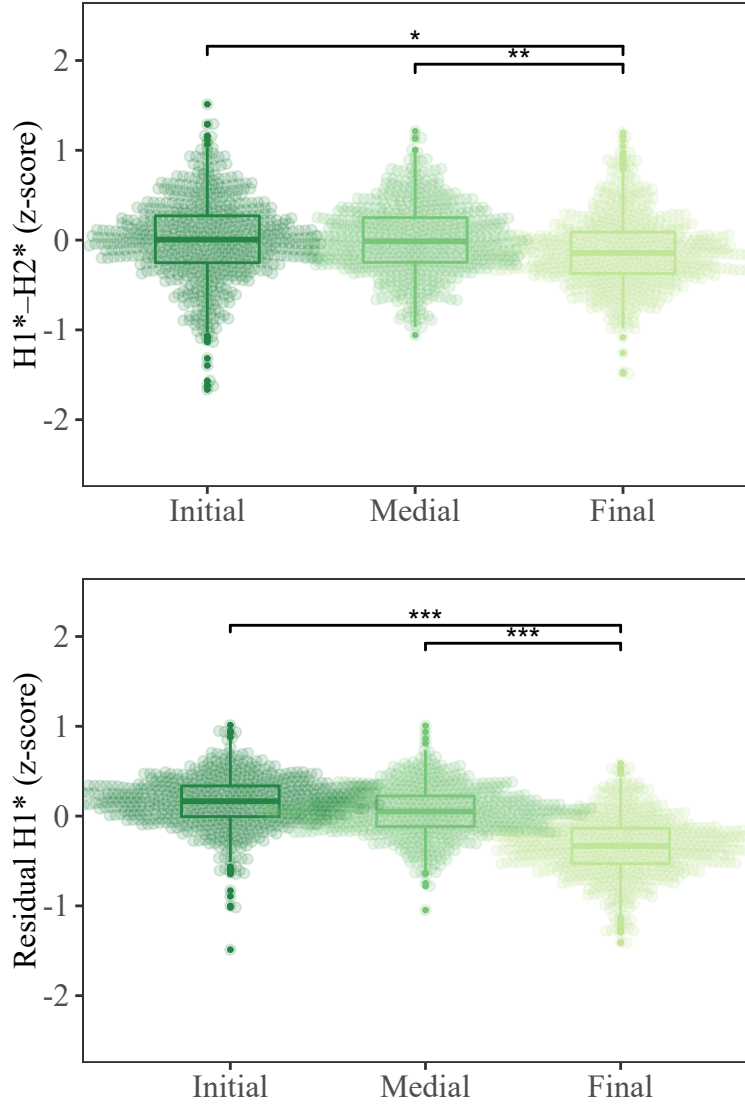


FIG. 5. (color online) Residual $H1^*-H2^*$ (top) and Residual $H1^*$ (bottom) in different positions in Mandarin. *** $p < .001$; ** $0.001 < p < 0.01$; * $0.01 < p < 0.05$

Residual $H1^*-H2^* = H1^*-H2^* - f_0 * f_0$ coefficient (0.279) – random intercept and slopes in Model (7)

Residual $H1^* = H1^* - f_0 * f_0$ coefficient (0.025) – Energy * Energy coefficient (0.541) – random intercept and slopes in Model (8)

TABLE XI. Model comparison between $H1^*$ - $H2^*$ and $H1^*$ in distinguishing utterance positions in Mandarin

Position comparison	Model	β	Std. Error	t value	p
Initial-Final	$H1^* - H2^* \sim \mathbf{Position} + f0$	0.124	0.048	2.565	.013
	$H1^* \sim \mathbf{Position} + \mathbf{Energy} + f0$	0.494	0.047	10.464	<.001
Medial-Final	$H1^* - H2^* \sim \mathbf{Position} + f0$	0.129	0.044	2.914	.005
	$H1^* \sim \mathbf{Position} + \mathbf{Energy} + f0$	0.395	0.037	10.724	<.001

TABLE XII. Marginal R^2 of Model (7) and (8)

Model	Marginal R^2
(7) $H1^* - H2^* \sim \mathbf{Position} + f0 + (1 \mathbf{Speaker}) + (f0 + \mathbf{Position} \mathbf{Speaker})$	0.254
(8) $H1^* \sim \mathbf{Position} + \mathbf{Energy} + f0 + (1 \mathbf{Speaker}) + (f0 + \mathbf{Position} \mathbf{Speaker})$	0.805

614 V. DISCUSSION AND CONCLUSION

615 The goals of this paper were to review the history of $H1(-H2)$ as acoustic measures
616 of phonation type. We trace their origin back to the pioneering work by [Fischer-Jørgensen](#)
617 (1967) on breathy vs. modal vowels in Gujarati, and highlight later studies that advanced our
618 understanding of these measures in terms of their aerodynamic and articulatory correlates,

619 their perceptibility by listeners, and their use in indexing various phonological contrasts and
620 social factors.

621 We then highlighted several issues for using H1–H2 as the indicator for vocal fold constric-
622 tion. First, we reviewed the literature that has made use of H1–H2 in studies of linguistic
623 phonation type. We found that the measure frequently succeeds at distinguishing both
624 breathy and creaky phonation types from modal voice. However, there is a tendency for the
625 measure to be less effective at creaky vs. modal contrasts than breathy vs. modal ones. We
626 attribute this to errors in f_0 estimation during irregular voicing associated with creaky voice.
627 We further argue that H1–H2 (and particularly H1*–H2*) is prone to error propagation, in
628 that it measures two spectral values, both of which are affected by f_0 and vowel formants.
629 H1–H2 is also affected by other glottal features besides OQ, in part because H2 is affected
630 by other factors like glottal skew. Finally, H1–H2 is affected by nasal poles and zeroes. In
631 the current study, we show that using “residual H1*,” for which Energy is used to normalize
632 H1* (instead of H2*), can to some extent mitigate these issues.

633 Limited previous work has already made use of H1 instead of or in addition to H1–H2.
634 For instance, [Esposito \(2012\)](#) found that H1* differentiated the phonation types in White
635 Hmong better than H1*–H2* (in the sense that significant differences between phonation
636 types were found more often), and H1* had a stronger correlation with CQ than H1*–H2*.
637 However, she did not normalize the amplitude of H1*, meaning that there is a potential
638 for a confound between phonation and SPL differences. The current study used energy to
639 normalize H1* either by adding energy as a covariate of the H1* model or by subtracting
640 the effect of energy from H1*, resulting in a new measurement of residual H1*.

641 Using three data sets, we also show how residual $H1^*$ can be used in practice. A corpus
642 analysis of natural speech data (taken from the “Production and Perception of Linguistic
643 Voice Quality” project at UCLA, which includes EGG and audio recordings from eight
644 languages and dozens of speakers), revealed that residual $H1^*$ has a stronger relationship to
645 glottal OQ than $H1^*-H2^*$. Second, we showed that residual $H1^*$ better differentiated the
646 phonation types in !Xóõ than $H1^*-H2^*$, particularly for modal vs. creaky vowels, as expected
647 from our error simulations. Finally, we found that residual $H1^*$ better differentiated the
648 changes in phonatory quality by utterance position in Mandarin than $H1^*-H2^*$. We therefore
649 suggest that researchers consider using RMS energy to normalize for the amplitude of $H1$,
650 treating residual $H1$ as an acoustic correlate of vocal fold constriction— instead of, or in
651 addition to, $H1-H2$. It is worth noting the one context in which $H1-H2$ might be preferred:
652 when directly comparing just two tokens. In such a case, the overall effect of energy on $H1$
653 cannot be estimated. But given the move towards larger data sets in the phonetic sciences,
654 it is exceedingly rare for researchers to describe a contrast using only two tokens, except for
655 the purposes of general illustration. Certainly, in a phonetic analysis of phonation type that
656 makes use of multiple tokens from several speakers, an estimate of the effect of energy on
657 $H1$ can be made, and we argue here that it is desirable for researchers to calculate residual
658 $H1$.

659 A claim can also be made that residual $H1$ is better motivated theoretically than $H1-H2$.
660 As we discussed earlier, early uses of $H1-H2$ were motivated by observable differences in $H1$,
661 rather than by any theoretical import assigned to the slope of $H1-H2$ or to $H2$ in particular.
662 After all, $H1$ is the amplitude of the fundamental, which as the primary correlate of vocal

663 pitch clearly matters for overall voice quality perception. Thus, residual H1 is correlated
664 with how loud the fundamental – and thus pitch – is perceived to be relative to the overall
665 loudness of the signal. We also know that SPL is an important component to voice and
666 signal perception, and it is included in psychoacoustic models of the voice (Kreiman *et al.*,
667 2014). Residual H1 therefore captures information about two important cues: f0 as a cue
668 to pitch, and SPL as a cue to loudness. In contrast, a measure like H1–H2 includes a
669 component of the source spectrum – H2 – that is not known to matter intrinsically for
670 voice quality perception. Clearly, future work is needed to examine how listeners assess
671 H1 as a cue relative to other spectral landmarks and the signal more broadly. This should
672 also include comparisons between H1 and spectral tilt as measured with reference to vowel
673 formants; namely, H1–A1, H1–A2, or other formant-based measures like A1–A2 that make
674 no reference to the fundamental. As Garellek (2019, p. 88) mentioned, the use of formant-
675 based measures carries the assumption that voice quality depends on vowel quality. That
676 assumption may ultimately prove correct, but so far it has gone untested.

677 The conceptualization of H1 relative to overall energy also has implications for how we
678 model the voice source spectrum. Earlier work in this regard (Garellek *et al.*, 2013, 2016a;
679 Kreiman *et al.*, 2012) explicitly models source H1–H2 as a harmonic slope, in addition to
680 H2–H4, H4–H2kHz (the spectral slope from H4 to the harmonic closest to 2000 Hz) and
681 H2kHz–H5kHz (the spectral slope between the harmonics closest to 2000 and 5000 Hz). But
682 if what matters is H1 and not H1–H2, then perhaps a more suitable model of the harmonic
683 source spectrum could include only H1 instead of H1–H2. Practically this would involve
684 only a minor change to the model: instead of H1–H2 and three additional spectral slopes,

685 the updated harmonic source spectrum would include H1 and those same additional slopes.
686 But there are important theoretical implications to this change, because H1 would not be
687 compared directly to another harmonic or to any other segment of the harmonic source
688 spectrum; its raw amplitude is what would matter, just like its raw frequency (that is,
689 the f_0) matters. Of course, to control for overall SPL, H1 (and the other spectral slopes)
690 should be modeled as a sub-component of a larger psychoacoustic model of the voice that
691 includes overall energy, as done already in the psychoacoustic model of [Kreiman *et al.* \(2014\)](#).
692 Ultimately, the choice of whether to include H1 or H1–H2 in a psychoacoustic model of the
693 voice should depend on which of the two measures provides a better link between voice
694 production and voice quality perception. Much more work is therefore needed to determine
695 whether H1 provides a closer link between voice production and perception than H1–H2.

696 To conclude, we have shown that residual H1 has fewer error propagation issues than
697 H1–H2; using residual H1 can therefore lead to more accurate measurements, and thus
698 better description, of the acoustic correlates of vocal fold constriction. Future studies should
699 investigate what the specific articulatory and aerodynamic correlates of H1 are: does the
700 measure more closely reflect changes in vocal fold constriction, medial fold thickness, or
701 glottal skew? Additionally, future work could investigate the extent to which H1 outperforms
702 H1–H2 when comparing the voice quality across nasal vowels and whether listener judgments
703 of voice quality are better predicted by H1 than by H1–H2.

704 ¹Words that have nasal vowels; are marked as “do not use”; or have unmatched annotations between the
705 acoustic and EGG results files were excluded.

706 ${}^2CQ = 1 - OQ$.

707 ³The effect size of Model (6) is represented by the marginal R^2 of the fixed variables.

708

709 Arras, K. O. (1998). "An Introduction To Error Propagation: Derivation, Meaning and
710 Examples of Equation $Cy = Fx Cx FxT$," Technical Report, [http://hdl.handle.net/](http://hdl.handle.net/20.500.11850/82620)
711 [20.500.11850/82620](http://hdl.handle.net/20.500.11850/82620), doi: [10.3929/ETHZ-A-010113668](https://doi.org/10.3929/ETHZ-A-010113668), artwork Size: 22 p. Medium:
712 application/pdf.

713 Bickley, C. (1982). "Acoustic analysis and perception of breathy vowels," MIT Speech Com-
714 munication Working Papers **1**, 71–81.

715 Blankenship, B. (2002). "The timing of nonmodal phonation in vowels," Jour-
716 nal of Phonetics **30**(2), 163–191, [https://linkinghub.elsevier.com/retrieve/pii/](https://linkinghub.elsevier.com/retrieve/pii/S009544700190155X)
717 [S009544700190155X](https://linkinghub.elsevier.com/retrieve/pii/S009544700190155X), doi: [10.1006/jpho.2001.0155](https://doi.org/10.1006/jpho.2001.0155).

718 Caballero, G., and Carroll, L. (2015). "Tone and stress in Choguita Rarámuri (Tarahumara)
719 word prosody," International Journal of American Linguistics **81**, 457–493.

720 Campbell, N., and Beckman, M. (1997). "Stress, prominence, and spectral tilt," in *Intona-*
721 *tion: Theory, Models, and Applications*, International Speech Communication Association,
722 Athens, Greece, pp. 67–70.

723 Chai, Y. (2019). "The source of creak in Mandarin utterances," in *Proceedings of the 19th*
724 *International Congress of Phonetic Sciences, Melbourne, Australia 2019*, edited by S. Cal-
725 houn, P. Escudero, M. Tabain, and P. Warren, Australasian Speech Science and Technology
726 Association Inc, Canberra, Australia, pp. 1858–1862.

727 Chai, Y. (2021). “The source of creak in Mandarin utterances,” UC San Diego: San Diego
728 Linguistic Papers 8, 1–32, <https://escholarship.org/uc/item/8mg0x5pb>.

729 Chai, Y., and Ye, S. (2022). “Checked Syllables, Checked Tones, and Tone Sandhi in Xiapu
730 Min,” Languages 7(1), 47, <https://www.mdpi.com/2226-471X/7/1/47>, doi: 10.3390/
731 languages7010047.

732 Cho, T., and Ladefoged, P. (1999). “Variation and universals in VOT: Evidence from 18
733 languages,” Journal of Phonetics 27(207-229).

734 Chodroff, E., Golden, A., and Wilson, C. (2019). “Covariation of stop voice onset time across
735 languages: Evidence for a universal constraint on phonetic realization,” The Journal of the
736 Acoustical Society of America 145(1), EL109–EL115, [http://asa.scitation.org/doi/](http://asa.scitation.org/doi/10.1121/1.5088035)
737 [10.1121/1.5088035](http://asa.scitation.org/doi/10.1121/1.5088035), doi: 10.1121/1.5088035.

738 Dang, J., and Honda, K. (1996). “Acoustic characteristics of the human paranasal si-
739 nuses derived from transmission characteristic measurement and morphological obser-
740 vation,” The Journal of the Acoustical Society of America 100(5), 3374–3383, [http:](http://asa.scitation.org/doi/10.1121/1.416978)
741 [//asa.scitation.org/doi/10.1121/1.416978](http://asa.scitation.org/doi/10.1121/1.416978), doi: 10.1121/1.416978.

742 DiCanio, C. T. (2009). “The phonetics of register in Takhian Thong Chong,” Jour-
743 nal of the International Phonetic Association 39(2), 162–188, [https://www.cambridge.](https://www.cambridge.org/core/product/identifier/S0025100309003879/type/journal_article)
744 [org/core/product/identifier/S0025100309003879/type/journal_article](https://www.cambridge.org/core/product/identifier/S0025100309003879/type/journal_article), doi: 10.
745 [1017/S0025100309003879](https://www.cambridge.org/core/product/identifier/S0025100309003879/type/journal_article).

746 Doval, B., and d’Alessandro, C. (1997). “Spectral correlates of glottal waveform models:
747 an analytic study,” in *1997 IEEE International Conference on Acoustics, Speech, and*
748 *Signal Processing*, IEEE Comput. Soc. Press, Munich, Germany, Vol. 2, pp. 1295–1298,

749 <http://ieeexplore.ieee.org/document/596183/>, doi: 10.1109/ICASSP.1997.596183.

750 Doval, B., d'Alessandro, C., and Henrich, N. (2006). "The spectrum of glottal flow models,"
751 Acta Acustica united with Acustica **92**(6), 1026–1046, <https://www.ingentaconnect.com/content/dav/aaua/2006/00000092/00000006/art00021>.

752

753 Esposito, C. M. (2006). "The Effects of Linguistic Experience on the Perception of Phona-
754 tion," Ph.D. Dissertation, University of California, Los Angeles, Los Angeles, CA, USA,
755 http://phonetics.linguistics.ucla.edu/research/Esposito_diss.pdf.

756 Esposito, C. M. (2010). "The effects of linguistic experience on the perception of phonation,"
757 Journal of Phonetics **38**(2), 306–316, <https://linkinghub.elsevier.com/retrieve/pii/S0095447010000203>, doi: 10.1016/j.wocn.2010.02.002.

758

759 Esposito, C. M. (2012). "An acoustic and electroglottographic study of White Hmong tone
760 and phonation," Journal of Phonetics **40**(3), 466–476, <https://linkinghub.elsevier.com/retrieve/pii/S0095447012000174>, doi: 10.1016/j.wocn.2012.02.007.

761

762 Esposito, C. M., and Khan, S. u. D. (2020). "The cross-linguistic patterns of phona-
763 tion types," Language and Linguistics Compass **14**(12), 1–25, <https://onlinelibrary.wiley.com/doi/10.1111/lnc3.12392>, doi: 10.1111/lnc3.12392.

764

765 Fant, G., Liljencrants, J., and Lin, Q.-g. (1985). "A four-parameter model of glottal flow,"
766 STL-QPSR **26**, 1–13.

767 Fant, G., and Lin, Q.-g. (1988). "Frequency domain interpretation and derivation of glottal
768 flow parameters," STL-QPSR **4**, 1–21.

769 Fischer-Jørgensen, E. (1967). "Phonetic analysis of breathy (murmured) vowels in Gu-
770 jarati," Indian Linguistics **28**, 71–139.

771 Garellek, M. (2012). “The timing and sequencing of coarticulated non-modal phonation in
772 English and White Hmong,” *Journal of Phonetics* **40**(1), 152–161, [https://linkinghub.
773 elsevier.com/retrieve/pii/S0095447011000969](https://linkinghub.elsevier.com/retrieve/pii/S0095447011000969), doi: 10.1016/j.wocn.2011.10.003.

774 Garellek, M. (2019). “The phonetics of voice,” in *Routledge Handbook of Phonetics*, edited
775 by W. Katz and P. Assmann (Routledge, Oxford), pp. 75–106.

776 Garellek, M. (2020). “Acoustic Discriminability of the Complex Phonation System
777 in !Xóõ,” *Phonetica* **77**(2), 131–160, [https://www.degruyter.com/document/doi/10.
778 1159/000494301/html](https://www.degruyter.com/document/doi/10.1159/000494301/html), doi: 10.1159/000494301.

779 Garellek, M. (2022). “Theoretical achievements of phonetics in the 21st century: Phonetics
780 of voice quality,” *Journal of Phonetics* **94**, 101155, [https://linkinghub.elsevier.com/
781 retrieve/pii/S0095447022000304](https://linkinghub.elsevier.com/retrieve/pii/S0095447022000304), doi: 10.1016/j.wocn.2022.101155.

782 Garellek, M., and Esposito, C. M. (2021). “Phonetics of White Hmong vowel and tonal con-
783 trasts,” *Journal of the International Phonetic Association* 1–20, [https://www.cambridge.
784 org/core/product/identifier/S0025100321000104/type/journal_article](https://www.cambridge.org/core/product/identifier/S0025100321000104/type/journal_article), doi: 10.
785 1017/S0025100321000104.

786 Garellek, M., Keating, P., Esposito, C. M., and Kreiman, J. (2013). “Voice quality
787 and tone identification in White Hmong,” *The Journal of the Acoustical Society of*
788 *America* **133**(2), 1078–1089, <http://asa.scitation.org/doi/10.1121/1.4773259>, doi:
789 10.1121/1.4773259.

790 Garellek, M., Ritchart, A., and Kuang, J. (2016a). “Breathy voice during nasality: a cross-
791 linguistic study,” *Journal of Phonetics* **59**, 110–121.

792 Garellek, M., Samlan, R., Gerratt, B. R., and Kreiman, J. (2016b). “Modeling the
793 voice source in terms of spectral slopes,” *The Journal of the Acoustical Society of*
794 *America* **139**(3), 1404–1410, <http://asa.scitation.org/doi/10.1121/1.4944474>, doi:
795 [10.1121/1.4944474](https://doi.org/10.1121/1.4944474).

796 Garellek, M., and White, J. (2015). “Phonetics of Tongan stress,” *Journal of the Interna-*
797 *tional Phonetic Association* **45**, 13–34.

798 Gobl, C., Murphy, A., Yanushevskaya, I., and Chasaide, A. N. (2018). “On the relation-
799 ship between glottal pulse shape and its spectrum: correlations of open quotient, pulse
800 skew and peak flow with source harmonic amplitudes,” in *Proceedings of Interspeech 2018*,
801 Hyderabad, India, pp. 222–226.

802 Gobl, C., and Ní Chasaide, A. (2019). “Time to Frequency Domain Mapping of the Voice
803 Source: The Influence of Open Quotient and Glottal Skew on the Low End of the Source
804 Spectrum,” in *Interspeech 2019*, ISCA, pp. 1961–1965, [https://www.isca-speech.org/](https://www.isca-speech.org/archive/interspeech_2019/gobl19_interspeech.html)
805 [archive/interspeech_2019/gobl19_interspeech.html](https://www.isca-speech.org/archive/interspeech_2019/gobl19_interspeech.html), doi: [10.21437/Interspeech.](https://doi.org/10.21437/Interspeech.2019-2888)
806 [2019-2888](https://doi.org/10.21437/Interspeech.2019-2888).

807 Gordon, M., and Ladefoged, P. (2001). “Phonation types: a cross-linguistic overview,”
808 *Journal of Phonetics* **29**, 383–406.

809 Guion, S. G., Amith, J. D., Doty, C. S., and Shport, I. A. (2010). “Word-level prosody in
810 Balsas Nahuatl: The origin, development, and acoustic correlates of tone in a stress accent
811 language,” *Journal of Phonetics* **38**, 137–166.

812 Guion, S. G., Post, M. W., and Payne, D. L. (2004). “Phonetic correlates of tongue root
813 vowel contrasts in Maa,” *Journal of Phonetics* **32**, 517–542.

814 Hanson, H. M. (1995). “Glottal characteristics of female speakers,” PhD Thesis, Harvard
815 University.

816 Hanson, H. M. (1997). “Glottal characteristics of female speakers: Acoustic correlates,”
817 Journal of the Acoustical Society of America **101**, 466–481.

818 Hanson, H. M., and Chuang, E. S. (1999). “Glottal characteristics of male speakers: Acous-
819 tic correlates and comparison with female data,” The Journal of the Acoustical Society of
820 America **106**(2), 1064–1077, <http://asa.scitation.org/doi/10.1121/1.427116>, doi:
821 [10.1121/1.427116](http://dx.doi.org/10.1121/1.427116).

822 Hanson, H. M., Stevens, K. N., Kuo, H.-K. J., Chen, M. Y., and Slifka, J. (2001). “Towards
823 models of phonation,” Journal of Phonetics **29**, 451–480.

824 Henrich, N., d’Alessandro, C., and Doval, B. (2001). “Spectral Correlates of Voice Open
825 Quotient and Glottal Flow Asymmetry : Theory, Limits and Experimental Data,”
826 in *Proceedings of EUROSPEECH-2001*, Aalborg, Denmark, pp. 47–50, [https://www.](https://www.isca-speech.org/archive_v0/eurospeech_2001/e01_0047.html)
827 [isca-speech.org/archive_v0/eurospeech_2001/e01_0047.html](https://www.isca-speech.org/archive_v0/eurospeech_2001/e01_0047.html).

828 Henton, C. G., and Bladon, R. A. W. (1985). “Breathiness in normal female speech: Inef-
829 ficiency versus desirability,” Language and Communication **5**, 221–227.

830 Holmberg, E. B., Hillman, R. E., Perkell, J. S., Guiod, P., and Goldman, S. L. (1995).
831 “Comparisons among aerodynamic, electroglottographic, and acoustic spectral measures
832 of female voice,” Journal of Speech and Hearing Research **38**, 1212–1223.

833 Howard, D. M. (1995). “Variation of electrolaryngographically derived closed quotient for
834 trained and untrained adult female singers,” Journal of Voice **9**, 163–172.

835 Huffman, M. K. (1987). “Measures of phonation type in Hmong,” The Journal of the Acous-
836 tical Society of America **81**(2), 495–504, [http://asa.scitation.org/doi/10.1121/1.](http://asa.scitation.org/doi/10.1121/1.394915)
837 [394915](http://asa.scitation.org/doi/10.1121/1.394915), doi: [10.1121/1.394915](https://doi.org/10.1121/1.394915).

838 Iseli, M., Shue, Y.-L., and Alwan, A. (2007). “Age, sex, and vowel dependencies of acoustic
839 measures related to the voice source,” Journal of the Acoustical Society of America **121**,
840 2283–2295.

841 Johnson, P. C. (2014). “Extension of Nakagawa & Schielzeth’s R^2_{glmm} to random slopes
842 models,” Methods in Ecology and Evolution **5**(9), 944–946, [https://onlinelibrary.](https://onlinelibrary.wiley.com/doi/10.1111/2041-210X.12225)
843 [wiley.com/doi/10.1111/2041-210X.12225](https://onlinelibrary.wiley.com/doi/10.1111/2041-210X.12225), doi: [10.1111/2041-210X.12225](https://doi.org/10.1111/2041-210X.12225).

844 Keating, P., Garellek, M., and Kreiman, J. (2015). “Acoustic properties of different kinds of
845 creaky voice,” in *Proceedings of the 18th International Congress of Phonetic Sciences*, Glas-
846 gow, https://idiom.ucsd.edu/~mgarellek/files/Keating_etal_2015_ICPhS.pdf.

847 Kirk, P., Ladefoged, P., and Ladefoged, J. (1984). “Using a spectrograph for measures of
848 phonation types in natural language,” UCLA Working Papers in Phonetics **59**, 102–113.

849 Klatt, D. H., and Klatt, L. C. (1990). “Analysis, synthesis, and perception of voice qual-
850 ity variations among female and male talkers,” The Journal of the Acoustical Society
851 of America **87**(2), 820–857, <http://asa.scitation.org/doi/10.1121/1.398894>, doi:
852 [10.1121/1.398894](https://doi.org/10.1121/1.398894).

853 Kreiman, J., Gerratt, B., and Antoñanzas-Barroso, N. (2007). “Measures of the glottal
854 source spectrum,” Journal of Speech, Language, and Hearing Research **50**, 595–610.

855 Kreiman, J., and Gerratt, B. R. (2010). “Perceptual Assessment of Voice Quality: Past,
856 Present, and Future,” Perspectives on Voice and Voice Disorders **20**, 62–67.

857 Kreiman, J., Gerratt, B. R., Garellek, M., Samlan, R., and Zhang, Z. (2014). “Toward a
858 unified theory of voice production and perception,” *Loquens* **e009**.

859 Kreiman, J., Gerratt, B. R., and Khan, S. u. D. (2010). “Effects of native language on
860 perception of voice quality,” *Journal of Phonetics* **38**(4), 588–593, [https://linkinghub.
861 elsevier.com/retrieve/pii/S0095447010000641](https://linkinghub.elsevier.com/retrieve/pii/S0095447010000641), doi: 10.1016/j.wocn.2010.08.004.

862 Kreiman, J., Shue, Y.-L., Chen, G., Iseli, M., Gerratt, B. R., Neubauer, J., and Alwan,
863 A. (2012). “Variability in the relationships among voice quality, harmonic amplitudes,
864 open quotient, and glottal area waveform shape in sustained phonation,” *Journal of the
865 Acoustical Society of America* **132**, 2625–2632.

866 Kuang, J. (2011). “Production and Perception of the Phonation Contrast in Yi,” Master’s
867 thesis, University of California, Los Angeles, Los Angeles, CA, USA.

868 Kuang, J. (2017). “Covariation between voice quality and pitch: Revisiting the case of
869 Mandarin creaky voice,” *The Journal of the Acoustical Society of America* **142**(3), 1693–
870 1706, <http://asa.scitation.org/doi/10.1121/1.5003649>, doi: 10.1121/1.5003649.

871 Ladefoged, P. (1981). “The relative nature of voice quality,” *The Journal of the Acous-
872 tical Society of America* **69**(S1), S67–S67, [http://asa.scitation.org/doi/10.1121/1.
873 386168](http://asa.scitation.org/doi/10.1121/1.386168), doi: 10.1121/1.386168.

874 Ladefoged, P. (1983). “Cross-linguistic studies of speech production,” in *The Production of
875 Speech*, edited by P. F. MacNeilage (Springer, New York), pp. 177–188.

876 Li, S., Gu, W., Liu, L., and Tang, P. (2020). “The Role of Voice Quality in Mandarin
877 Sarcastic Speech: An Acoustic and Electrolottographic Study,” *Journal of Speech, Lan-
878 guage, and Hearing Research* **63**(8), 2578–2588, <http://pubs.asha.org/doi/10.1044/>

879 [2020_JSLHR-19-00166](#), doi: [10.1044/2020_JSLHR-19-00166](#).

880 Maddieson, I., and Ladefoged, P. (1985). ““Tense” and “lax” in four minority languages
881 of China,” *Journal of Phonetics* **13**(4), 433–454, [https://linkinghub.elsevier.com/
882 retrieve/pii/S0095447019307880](https://linkinghub.elsevier.com/retrieve/pii/S0095447019307880), doi: [10.1016/S0095-4470\(19\)30788-0](#).

883 Ní Chasaide, A., and Gobl, C. (1993). “Contextual Variation of the Vowel Voice
884 Source as a Function of Adjacent Consonants,” *Language and Speech* **36**(2-3),
885 303–330, <http://journals.sagepub.com/doi/10.1177/002383099303600310>, doi: [10.
886 1177/002383099303600310](#).

887 Orlikoff, R. F. (1991). “Assessment of the dynamics of vocal fold contact from the elec-
888 troglottogram,” *Journal of Speech and Hearing Research* **34**, 1066–1072.

889 Pennington, M. (2005). “The phonetics and phonology of glottal manner features,”
890 Dissertation, Indiana University, Bloomington, [https://www.proquest.com/docview/
891 304986742?pq-origsite=gscholar&fromopenview=true](https://www.proquest.com/docview/304986742?pq-origsite=gscholar&fromopenview=true).

892 Samlan, R. A., and Story, B. H. (2011). “Relation of structural and vibratory kinemat-
893 ics of the vocal folds to two acoustic measures of breathy voice based on computational
894 modeling,” *Journal of Speech, Language, and Hearing Research* **54**, 1267–1283.

895 Samlan, R. A., Story, B. H., and Bunton, K. (2013). “Relation of perceived breathiness to
896 laryngeal kinematics and acoustic measures based on computational modeling,” *Journal of
897 Speech, Language, and Hearing Research* **56**, 1209–1223.

898 Seyfarth, S., and Garellek, M. (2018). “Plosive voicing acoustics and voice quality in Yere-
899 van Armenian,” *Journal of Phonetics* **71**, 425–450, [https://linkinghub.elsevier.com/
900 retrieve/pii/S0095447017302711](https://linkinghub.elsevier.com/retrieve/pii/S0095447017302711), doi: [10.1016/j.wocn.2018.09.001](#).

901 Shue, Y.-L., Chen, G., and Alwan, A. (2010). “On the interdependencies between voice
902 quality, glottal gaps, and voice-source related acoustic measures,” in *Proceedings of Inter-*
903 *speech 2010*, pp. 34–37.

904 Shue, Y.-L., Keating, P., Vicenik, C., and Yu, K. M. (2011). “VoiceSauce: A program for
905 voice analysis,” in *Proceedings of the 17th International Congress of Phonetic Science*,
906 Hong Kong, pp. 1846–1849.

907 Simpson, A. (2012). “The first and second harmonics should not be used to measure breath-
908 iness in male and female voices,” *Journal of Phonetics* **40**, 477–490.

909 Sluijter, A., and van Heuven, V. (1996). “Acoustic correlates of linguistic stress and accent in
910 Dutch and American English,” in *Proceeding of Fourth International Conference on Spoken*
911 *Language Processing. ICSLP '96*, IEEE, Philadelphia, PA, USA, Vol. 2, pp. 630–633,
912 <http://ieeexplore.ieee.org/document/607440/>, doi: [10.1109/ICSLP.1996.607440](https://doi.org/10.1109/ICSLP.1996.607440).

913 Stevens, K. N. (1977). “Physics of laryngeal behavior and larynx modes,” *Phonetica* **34**,
914 264–279.

915 Stevens, K. N., and Hanson, H. M. (1995). “Classification of glottal vibration from acoustic
916 measurements,” in *Vocal fold physiology: Voice quality control*, edited by O. Fujimura and
917 M. Hirano (Singular Publishing Group, San Diego, CA), pp. 147–170.

918 Styler, W. (2017). “On the acoustical features of vowel nasality in English and French,”
919 *Journal of the Acoustical Society of America* **142**, 2469–2482.

920 Sundberg, J. (2022). “Objective Characterization of Phonation Type Using Amplitude of
921 Flow Glottogram Pulse and of Voice Source Fundamental,” *Journal of Voice* **36**(1), 4–
922 14, <https://linkinghub.elsevier.com/retrieve/pii/S0892199720301107>, doi: [10.1016/j.jvoice.2021.10.001](https://doi.org/10.1016/j.jvoice.2021.10.001).

923 [1016/j.jvoice.2020.03.018](https://doi.org/10.1016/j.jvoice.2020.03.018).

924 Sundberg, J., Andersson, M., and Hultqvist, C. (1999). “Effects of subglottal pressure vari-
925 ation on professional baritone singers’ voice sources,” *The Journal of the Acoustical Soci-*
926 *ety of America* **105**(3), 1965–1971, <http://asa.scitation.org/doi/10.1121/1.426731>,
927 doi: [10.1121/1.426731](https://doi.org/10.1121/1.426731).

928 Sundberg, J., and Gauffin, J. (1979). “Waveform and spectrum of the glottal voice source,”
929 in *Frontiers of Frontiers of speech communication research, Festschrift for Gunnar Fant*,
930 edited by B. Lindblom and S. E. J. Öhman (Academic Press, London), pp. 301–320.

931 Swerts, M., and Veldhuis, R. (2001). “The effect of speech melody on voice quality,” *Speech*
932 *Communication* **33**, 297–303.

933 Tabain, M., Garellek, M., Hellwig, B., Gregory, A., and Beare, R. (2022). “Voicing in
934 Qaqet: Prenasalization and language contact,” *Journal of Phonetics* **91**, 101138, [https://](https://linkinghub.elsevier.com/retrieve/pii/S0095447022000134)
935 linkinghub.elsevier.com/retrieve/pii/S0095447022000134, doi: [10.1016/j.wocn.](https://doi.org/10.1016/j.wocn.2022.101138)
936 [2022.101138](https://doi.org/10.1016/j.wocn.2022.101138).

937 Titze, I. R., Baken, R. J., Bozeman, K. W., Granqvist, S., Henrich, N., Herbst, C. T.,
938 Howard, D. M., Hunter, E. J., Kaelin, D., Kent, R. D., Kreiman, J., Kob, M., Löfqvist,
939 A., McCoy, S., Miller, D. G., Noé, H., Scherer, R. C., Smith, J. R., Story, B. H., Švec,
940 J. G., Ternström, S., and Wolfe, J. (2015). “Toward a consensus on symbolic notation
941 of harmonics, resonances, and formants in vocalization,” *The Journal of the Acoustical*
942 *Society of America* **137**(5), 3005–3007, [https://asa.scitation.org/doi/citedby/10.](https://asa.scitation.org/doi/citedby/10.1121/1.4919349)
943 [1121/1.4919349](https://doi.org/10.1121/1.4919349), doi: [10.1121/1.4919349](https://doi.org/10.1121/1.4919349) publisher: Acoustical Society of America.

944 Traill, A., and Jackson, M. (1988). “Speaker variation and phonation type in Tsonga nasals,”
945 Journal of Phonetics **16**(4), 385–400, [https://linkinghub.elsevier.com/retrieve/](https://linkinghub.elsevier.com/retrieve/pii/S0095447019305170)
946 [pii/S0095447019305170](https://linkinghub.elsevier.com/retrieve/pii/S0095447019305170), doi: 10.1016/S0095-4470(19)30517-0.

947 Wiley, J. F. (2020). “Multilevel and mixed effects model diagnostics and effect sizes” [https:](https://github.com/JWiley/multilevelTools)
948 [//github.com/JWiley/multilevelTools](https://github.com/JWiley/multilevelTools).

949 Zhang, Z. (2016a). “Cause-effect relationship between vocal fold physiology and voice pro-
950 duction in a three-dimensional phonation model,” The Journal of the Acoustical Society of
951 America **139**(4), 1493–1507, <http://asa.scitation.org/doi/10.1121/1.4944754>, doi:
952 [10.1121/1.4944754](http://asa.scitation.org/doi/10.1121/1.4944754).

953 Zhang, Z. (2016b). “Mechanics of human voice production and control,” Journal of the
954 Acoustical Society of America **140**, 2614–2635.