

Using ambulatory voice monitoring to investigate common voice disorders: Research update

Daryush D. Mehta^{1, 2, 3*}, Jarrad H. Van Stan^{1, 3}, Matías Zañartu⁴, Marzyeh Ghassemi⁵, John V. Guttag⁵, Víctor M. Espinoza^{4, 6}, Juan P. Cortés⁴, Harold A. Cheyne⁷, Robert E. Hillman^{1, 2, 3}

¹Center for Laryngeal Surgery and Voice Rehabilitation, Massachusetts General Hospital, USA, ²Department of Surgery, Harvard Medical School, USA, ³MGH Institute of Health Professions, Massachusetts General Hospital, USA, ⁴Department of Electronic Engineering, Universidad Técnica Federico Santa María, Chile, ⁵Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, USA, ⁶Department of Music and Sonology, Faculty of Arts, Universidad de Chile, Chile, ⁷Laboratory of Ornithology, Bioacoustics Research Lab, Cornell University, USA

Submitted to Journal:
Frontiers in Bioengineering and Biotechnology

Specialty Section:
Bioinformatics and Computational Biology

ISSN:
2296-4185

Article type:
Original Research Article

Received on:
17 Jun 2015

Accepted on:
23 Sep 2015

Provisional PDF published on:
23 Sep 2015

Frontiers website link:
www.frontiersin.org

Citation:
Mehta DD, Van_stan JH, Zañartu M, Ghassemi M, Guttag JV, Espinoza VM, Cortés JP, Cheyne HA and Hillman RE(2015) Using ambulatory voice monitoring to investigate common voice disorders: Research update. *Front. Bioeng. Biotechnol.* 3:155. doi:10.3389/fbioe.2015.00155

Copyright statement:
© 2015 Mehta, Van_stan, Zañartu, Ghassemi, Guttag, Espinoza, Cortés, Cheyne and Hillman. This is an open-access article distributed under the terms of the [Creative Commons Attribution License \(CC BY\)](http://creativecommons.org/licenses/by/2.0/). The use, distribution and reproduction in other forums is permitted, provided the original author(s) or licensor are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.

This Provisional PDF corresponds to the article as it appeared upon acceptance, after peer-review. Fully formatted PDF and full text (HTML) versions will be made available soon.

Provisional

Conflict of interest statement

The authors declare a potential conflict of interest and state it below.

Patent application for methodology of subglottal impedance-based inverse filtering:

Zañartu M, Ho JC, Mehta DD, Wodicka GR, Hillman RE. System and methods for evaluating vocal function using an impedance-based inverse filtering of neck surface acceleration. International Patent Publication Number WO 2012/112985. Published August 23, 2012.

Provisional

1 **Using ambulatory voice monitoring**
2 **to investigate common voice disorders:**
3 **Research update**

4
5 **Daryush D. Mehta^{1,2,3*}, Jarrad H. Van Stan^{1,3}, Matías Zañartu⁴, Marzyeh Ghassemi⁵, John V.**
6 **Guttag⁵, Víctor M. Espinoza^{4,6}, Juan P. Cortés⁴, Harold A. Cheyne II⁷, Robert E. Hillman^{1,2,3}**

7 ¹Center for Laryngeal Surgery and Voice Rehabilitation, Massachusetts General Hospital, Boston,
8 Massachusetts, USA

9 ²Department of Surgery, Harvard Medical School, Boston, Massachusetts, USA

10 ³Institute of Health Professions, Massachusetts General Hospital, Boston, Massachusetts, USA

11 ⁴Department of Electronic Engineering, Universidad Técnica Federico Santa María, Valparaíso,
12 Chile

13 ⁵Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology,
14 Cambridge, Massachusetts, USA

15 ⁶Department of Music and Sonology, Faculty of Arts, Universidad de Chile, Santiago, Chile

16 ⁷Bioacoustics Research Lab, Laboratory of Ornithology, Cornell University, Ithaca, New York, USA

17 *** Correspondence:**

18 Daryush D. Mehta
19 Center for Laryngeal Surgery and Voice Rehabilitation
20 Massachusetts General Hospital
21 One Bowdoin Square, 11th Floor
22 Boston, MA, 02114, USA
23 mehta.daryush@mgh.harvard.edu

24
25 **Keywords: voice monitoring, accelerometer, vocal function, voice disorders, vocal**
26 **hyperfunction, glottal inverse filtering, machine learning.**

27 **Abstract (1415/2000 characters)**

28 Many common voice disorders are chronic or recurring conditions that are likely to result from
29 inefficient and/or abusive patterns of vocal behavior, referred to as vocal hyperfunction. The clinical
30 management of hyperfunctional voice disorders would be greatly enhanced by the ability to monitor
31 and quantify detrimental vocal behaviors during an individual's activities of daily life. This paper
32 provides an update on ongoing work that uses a miniature accelerometer on the neck surface below
33 the larynx to collect a large set of ambulatory data on patients with hyperfunctional voice disorders
34 (before and after treatment) and matched control subjects. Three types of analysis approaches are
35 being employed in an effort to identify the best set of measures for differentiating among
36 hyperfunctional and normal patterns of vocal behavior: 1) ambulatory measures of voice use that
37 include vocal dose and voice quality correlates, 2) aerodynamic measures based on glottal airflow
38 estimates extracted from the accelerometer signal using subject-specific vocal system models, and 3)
39 classification based on machine learning and pattern recognition approaches that have been used
40 successfully in analyzing long-term recordings of other physiological signals. Preliminary results
41 demonstrate the potential for ambulatory voice monitoring to improve the diagnosis and treatment of
42 common hyperfunctional voice disorders.

43 **1. Introduction**

44 Voice disorders have been estimated to affect approximately 30 % of the adult population in the
45 United States at some point in their lives, with 6.6 % to 7.6 % of individuals affected at any given
46 point in time (Roy et al., 2005;Bhattacharyya, 2014). While many vocally-healthy speakers take
47 verbal communication for granted, individuals suffering from voice disorders experience significant
48 communication disabilities with far-reaching social, professional, and personal consequences
49 (NIDCD, 2012).

50 Normal voice sounds are produced in the larynx by rapid air pulses that are emitted as the vocal cords
51 (folds) are driven into vibration by exhaled air from the lungs. Disturbances in voice production (i.e.,
52 voice disorders) can be caused by a variety of conditions that affect how the larynx functions to
53 generate sound, including 1) neurological disorders of the central (Parkinson's disease, stroke, etc.)
54 or peripheral (e.g., damage to laryngeal nerves causing vocal fold paresis/paralysis) nervous system;
55 2) congenital (e.g. restrictions in normal development of laryngeal/airway structures) or acquired
56 organic (e.g. laryngeal cancer, trauma, etc.) disorders of the larynx and/or airway; and 3) behavioral
57 disorders involving vocal abuse/misuse that may or may not cause trauma to vocal fold tissue (e.g.
58 nodules). The most frequently occurring subset of voice disorders is associated with *vocal*
59 *hyperfunction*, which refers to chronic "conditions of abuse and/or misuse of the vocal mechanism
60 due to excessive and/or 'imbalanced' [uncoordinated] muscular forces" (p. 373) (Hillman et al.,
61 1989). Over the years, our group has begun to provide evidence for the concept that there are two
62 types of vocal hyperfunction that can be quantitatively described and differentiated from each other
63 and normal voice production using a combination of acoustic and aerodynamic measures (Hillman et
64 al., 1989; 1990).

65 *Phonotraumatic vocal hyperfunction* (previously termed adducted hyperfunction) is associated with
66 the formation of benign vocal fold lesions—such as nodules and polyps. Vocal fold nodules or
67 polyps are believed to develop as a reaction to persistent tissue inflammation, chronic cumulative
68 vocal fold tissue damage, and/or environmental influences (Titze et al., 2003;Czerwonka et al.,

69 2008;Karkos and McCormick, 2009). Once formed, these lesions may prevent adequate vocal fold
70 contact/closure that reduces the efficiency of sound production and can cause individuals to
71 compensate by increasing muscular and aerodynamic forces. This compensatory behavior may result
72 in further tissue damage and become habitual due to the need to constantly maintain functional voice
73 production during daily life in the presence of a vocal fold pathology. In contrast, *non-*
74 *phonotraumatic vocal hyperfunction* (previously termed non-adducted hyperfunction)—often
75 diagnosed as muscle tension dysphonia (MTD) or functional dysphonia—is associated with
76 symptoms such as vocal fatigue, excessive intrinsic/extrinsic neck muscle tension and discomfort,
77 and voice quality degradation in the absence of vocal fold tissue trauma. There can be a wide range
78 of voice quality disturbances (e.g., various degrees of strain or breathiness) whose nature and severity
79 can display significant situational variation, such as variation associated with changes in levels of
80 emotional stress throughout the course of a day (Hillman et al., 1990). MTD can be triggered by a
81 variety of conditions/circumstances, including psychological conditions (traumatizing events,
82 emotional stress, etc.), chronic irritation of the laryngeal and/or pharyngeal mucosa (e.g.,
83 laryngopharyngeal reflux), and habituation of maladaptive behaviors such as persistent dysphonia
84 following resolution of an upper respiratory infection (Roy and Bless, 2000).

85 To assess the prevalence and persistence of hyperfunctional vocal behaviors during diagnosis and
86 management, clinicians currently rely on patient self-report and self-monitoring, which are highly
87 subjective and prone to be unreliable. In addition, investigators have studied clinician-administered
88 perceptual ratings of voice quality and endoscopic imaging and the quantitative analysis of objective
89 measures derived from acoustics, electroglottography, imaging, and aerodynamic voice signals (Roy
90 et al., 2013). Among work that sought to automatically detect voice disorders including vocal
91 hyperfunction, acoustic analysis approaches have employed neural maps (Hadjitodorov et al., 2000),
92 nonlinear measures (Little et al., 2007), and voice source-related properties (Parsa and Jamieson,
93 2000) from snapshots of phonatory recordings obtained during a single laboratory session. Because
94 hyperfunctional voice disorders are associated with daily behavior, the diagnosis and treatment of
95 these disorders may be greatly enhanced by the ability to unobtrusively monitor and quantify vocal
96 behaviors as individuals go about their normal daily activities. Ambulatory voice monitoring may
97 enable clinicians to better assess the role of vocal behaviors in the development of voice disorders,
98 precisely pinpoint the location and duration of abusive and/or maladaptive behaviors, and objectively
99 assess patient compliance with the goals of voice therapy.

100 This paper reports on our ongoing investigation into the use of a miniature accelerometer on the neck
101 surface below the larynx to acquire and analyze a large set of ambulatory data from patients with
102 hyperfunctional voice disorders (before and after treatment stages) as compared to matched control
103 subjects. We have previously reported on our development of a user-friendly and flexible platform
104 for voice health monitoring that employs a smartphone as the data acquisition platform connected to
105 the accelerometer (Mehta et al., 2012b;Mehta et al., 2013). The current report extends on that pilot
106 work and describes data acquisition protocols, as well as initial results from three analysis
107 approaches: 1) existing ambulatory measures of voice use, 2) aerodynamic measures based on glottal
108 airflow estimates extracted from the accelerometer signal, and 3) classification based on machine
109 learning and pattern recognition techniques. Although the methodologies of these analysis
110 approaches largely have been published, the novel contributions of the current paper include
111 ambulatory voice measures from the largest cohort of speakers to date (142 subjects), initial
112 estimation of ambulatory glottal airflow properties, and updated machine learning results for the
113 classification of 51 speakers with phonotraumatic vocal hyperfunction from matched control
114 speakers.

115 2. Materials and Methods

116

117 This section describes subject recruitment, data acquisition protocols, and the three analysis
118 approaches of existing voice use measures, aerodynamic parameter estimation, and machine learning
119 to aid in the classification of hyperfunctional vocal behaviors.

120

121 2.1. Subject Recruitment

122 Informed consent was obtained from all the subjects participating in this study, and all experimental
123 protocols were approved by the institutional review board of Partners HealthCare System at
124 Massachusetts General Hospital.

125 Two groups of individuals with voice disorders are being enrolled in the study: patients with
126 phonotraumatic vocal hyperfunction (vocal fold nodules or polyps) and patients with non-
127 phonotraumatic vocal hyperfunction (muscle tension dysphonia). Diagnoses are based on a complete
128 team evaluation by laryngologists and speech-language pathologists at the Massachusetts General
129 Hospital Voice Center that includes 1) a complete case history, 2) endoscopic imaging of the larynx
130 (Mehta and Hillman, 2012), 3) aerodynamic and acoustic assessment of vocal function (Roy et al.,
131 2013), 4) patient-reported Voice-Related Quality of Life (V-RQOL) questionnaire (Hogikyan and
132 Sethuraman, 1999), and 5) clinician-administered Consensus Auditory-Perceptual Evaluation of
133 Voice (CAPE-V) assessment (Kempster et al., 2009).

134 Matched-control groups are obtained for each of the two patient groups. Each patient typically aids in
135 identifying a work colleague of the same gender and approximate age (± 5 years) who has a normal
136 voice. The normal vocal status of all control subjects is verified via interview and a laryngeal
137 stroboscopic examination. Each control subject is monitored for one full 7-day week.

138 Figure 1 displays the treatment sequences (tracks) and time points at which patients in the study are
139 monitored for a full week. Patients with phonotraumatic vocal hyperfunction may follow one of three
140 usual treatment tracks (Figure 1A). The particular treatment track chosen depends upon clinical
141 management decisions regarding surgery or voice therapy. In Track A, individuals are monitored
142 before and after successful voice therapy and do not need surgical intervention (therapy may involve
143 sessions spanning several weeks or months). In Track B, patients initially attempt voice therapy but
144 subsequently require surgical removal of their vocal fold lesions to attain a more satisfactory vocal
145 outcome; a second round of voice therapy is then typically required to retrain the vocal behavior of
146 these patients to prevent the recurrence of vocal fold lesions. In Track C, patients undergo surgery
147 first followed by voice therapy. Finally, patients with non-phonotraumatic vocal hyperfunction
148 typically follow one treatment track and thus are monitored for one week before and after voice
149 therapy (Figure 1B).

150 Data collection is ongoing, as Figure 1 lists patient enrollment along with the number of vocally
151 healthy speakers who have been able to be recruited to be matched to a patient. For an initial analysis
152 of a complete data set, results are presented for subjects with available data from matched control
153 subjects. In addition, because the prevalence of these types of voice disorders is much higher in
154 females (hence, more data acquired from female subjects) and to eliminate the impact on the analysis
155 of known differences between male and female voice characteristics (such as fundamental
156 frequency), only female subject data were of focus in the current report.

157 Table 1 lists the occupations and diagnoses of the 51 female participants with phonotraumatic vocal
158 hyperfunction in the study who have been paired with matched control subjects (there were only 4
159 male subject pairs). All participants were engaged in occupations considered to be at a higher-than-
160 normal risk for developing a voice disorder. The majority of patients (37) were professional, amateur,
161 or student singers; every effort was made to match singers with control subjects in a similar musical
162 genre (classical or non-classical) to account for any genre-specific vocal behaviors. Forty-four
163 patients were diagnosed with vocal fold nodules, and seven patients had a unilateral vocal fold polyp.
164 The average (standard deviation) age of participants within the group was 24.4 (9.1) years.

165 Table 2 lists the occupations of the 20 female participants with non-phonotraumatic vocal
166 hyperfunction in the study who have been paired with matched control subjects (there were 6 male
167 subject pairs). All patients were diagnosed with muscle tension dysphonia and did not exhibit vocal
168 fold tissue trauma. The average (standard deviation) age of participants within the patient group was
169 41.8 (15.4) years.

170 **2.2. Data Acquisition Protocol**

171 Prior to in-field ambulatory voice monitoring, subjects are assessed in the laboratory to document
172 their vocal status and record signals that enable the calibration of the accelerometer signal for input to
173 the vocal system model that is used to estimate aerodynamic parameters.

174 **2.2.1. In-Laboratory Voice Assessment**

175 Figure 2A illustrates the in-laboratory multisensor setup consisting of the simultaneous acquisition of
176 data from the following devices:

- 177 1) Acoustic microphone placed 10 cm from the lips (MKE104, Sennheiser, Electronic GmbH,
178 Wennebostel, Germany)
- 179 2) Electroglossograph electrodes placed across the thyroid cartilage to measure time-varying
180 laryngeal impedance (EG-2, Glottal Enterprises, Syracuse, NY)
- 181 3) Accelerometer placed on the neck surface at the base of the neck (BU-27135; Knowles Corp.,
182 Itasca, IL)
- 183 4) Airflow sensor collecting high-bandwidth aerodynamic data via a circumferentially-vented
184 pneumotachograph face mask (PT-2E, Glottal Enterprises)
- 185 5) Low-bandwidth air pressure sensor connected to a narrow tube inserted through the lips in the
186 mouth (PT-25, Glottal Enterprises)

187 In particular, the use of the pneumotachograph mask to acquire the high-bandwidth oral airflow
188 signal is a key step in calibrating/adjusting the vocal system model described in Section 2.4 (Zañartu
189 et al., 2013) so that aerodynamic parameters can be extracted from the accelerometer signal. All
190 subjects wore the accelerometer below the level of the larynx (subglottal) on the front of the neck just
191 above the sternal notch. When recorded from this location, the accelerometer signal of an unknown
192 phrase is unintelligible. The accelerometer sensor used is relatively immune to environmental sounds
193 and produces a voice-related signal that is not filtered by the vocal tract, alleviating confidentiality
194 concerns because speech audio is not recorded.

195 The in-laboratory protocol requires subjects to perform the following speech tasks at a comfortable
196 pitch in their typical speaking voice mode:

- 197 1) Three cardinal vowels (“ah”, “ee”, “oo”) sustained at soft, comfortable, and loud levels
198 2) First paragraph of the Rainbow Passage at a comfortable loudness level
199 3) String of consonant-vowel pairs (e.g., “pae pae pae pae pae”)

200 The sustained vowels provide data for computing objective voice quality metrics such as perturbation
201 measures, harmonics-to-noise ratio, and harmonic spectral tilt. The Rainbow Passage is a standard
202 phonetically-balanced text that has been frequently used in voice and speech research (Fairbanks,
203 1960). The string of /pae/ syllables is designed to enable non-invasive, indirect estimates of lung
204 pressure (during lip closure for the /p/ when airway pressure reaches a steady state/equilibrates) and
205 laryngeal airflow (during vowel production when the airway is not constricted) during phonation
206 (Rothenberg, 1973). Figure 2B displays a snapshot of synchronized in-laboratory waveforms from
207 the consonant-vowel task for a 28-year-old female music teacher diagnosed with vocal fold nodules.

208 2.2.2. In-Field Ambulatory Monitoring of Voice Use

209 In the field, an Android smartphone (Nexus S; Samsung, Seoul, South Korea) provides a user-
210 friendly interface for voice monitoring, daily sensor calibration, and periodic collection of subject
211 responses to queries about their vocal status (Mehta et al., 2012b). The smartphone contains a high-
212 fidelity audio codec (WM8994; Wolfson Microelectronics, Edinburgh, Scotland, UK) that records
213 the accelerometer signal using sigma-delta modulation (128x oversampling) at a sampling rate of
214 11,025 Hz. Of critical importance, operating system root access allows for control over audio settings
215 related to highpass filtering and programmable gain arrays prior to analog-to-digital conversion. By
216 default, highpass filter cutoff frequencies are typically set above 100 Hz to optimize cellphone audio
217 quality and remove low-frequency noise due to wind noise and/or mechanical vibration. These cutoff
218 frequencies undesirably affect frequencies of interest through spectral shaping and phase distortion;
219 thus, for the current application, the highpass filter cutoff frequency is modified to a high-fidelity
220 setting of 0.9 Hz. Smartphone rooting also enables setting the analog gain to maximize signal
221 quantization; e.g., the WM8994 audio codec gain values can be set between -16.5 dB and +30.0 dB
222 in increments of 1.5 dB.

223 Figure 3 displays the smartphone-based voice health monitor system. Each morning, subjects affix
224 the accelerometer—encased in epoxy and mounted on a soft silicone pad—to their neck halfway
225 between the thyroid prominence and the suprasternal notch using hypoallergenic double-sided tape
226 (Model 2181, 3M, Maplewood, MN). Smartphone prompts then lead the subject through a brief
227 calibration sequence that maps the accelerometer signal amplitude to acoustic sound pressure level
228 (Švec et al., 2005). Subjects produce three “ah” vowels from a soft to loud (or loud to soft) level that
229 are used to generate a linear regression between acceleration amplitude and microphone signal level
230 (dB-dB plot) so that the uncalibrated acceleration level can be converted to units of dB SPL (dB re
231 20 μ Pa). The acoustic signal is recorded using a handheld audio recorder (H1 Handy Recorder, Zoom
232 Corporation, Tokyo, Japan) at a distance of 15 cm to the subject's lips. The microphone is not needed
233 the rest of the day.

234 With the smartphone placed in the pocket or worn in a belt holster, subjects engage in their typical
235 daily activities at work and home and are able to pause data acquisition during activities that could
236 damage the system, such as exercise, swimming, showering, etc. The smartphone application
237 requires minimal user interaction during the day. Every five hours, users are prompted to respond to
238 three questions related to vocal effort, discomfort, and fatigue (Carroll et al., 2006):

- 239 1) *Effort*: Say “ahhh” softly at a pitch higher than normal. Then say “ha ha ha ha ha” in the same
240 way. Rate how difficult the task was.
241 2) *Discomfort*: What is your current level of discomfort when talking or singing?
242 3) *Fatigue*: What is your current level of voice-related fatigue when talking or singing?

243 The three questions are answered using slider bars on the smartphone ranging from 0 (no presence of
244 effort, discomfort, or fatigue) to 100 (maximum effort, discomfort, or fatigue).

245 At the end of the day, the accelerometer is removed, recording is stopped, and the smartphone is
246 charged as the subject sleeps. A brief daily email survey asks subjects about when their work/school
247 day began and ended and if anything atypical occurred during the day.

248 2.3. Voice Quality and Vocal Dose Measures

249 Voice-related parameters for voice disorder classification fall into the following two categories: (1)
250 time-varying trajectories of features that are computed on a frame-by-frame basis and (2) measures of
251 voice use that accumulate frame-based metrics over a given duration (i.e., vocal dose measures).
252 These measures may be computed offline in a *post hoc* analysis of data or online on the smartphone
253 for real-time display or biofeedback.

254 Table 3 describes the suite of current frame-based parameters computed over 50-ms, non-overlapping
255 frames. These modifiable frame settings currently mimic the default behavior of the Ambulatory
256 Phonation Monitor (KayPENTAX, Montvale, NJ) and strikes a practical balance between the
257 requirement of real-time computation and capturing temporal and spectral voice characteristics
258 during time-varying speech production. The measures quantify signal properties related to amplitude,
259 frequency, periodicity, spectral tilt, and cepstral harmonicity: SPL and f_0 (Mehta et al., 2012b),
260 autocorrelation peak magnitude, harmonic spectral tilt (Mehta et al., 2011), low- to high-frequency
261 spectral power ratio (LH ratio) (Awan et al., 2010), and cepstral peak prominence (CPP) (Mehta et
262 al., 2012c). Figure 4A illustrates the computation of these measures from the time, spectral, and
263 cepstral domains. In the past, we have set *a priori* thresholds on signal amplitude, fundamental
264 frequency, and autocorrelation amplitudes to decide whether a frame contains voice activity or not
265 (Mehta et al., 2012b). Since then, additional signal measures have been implemented to improve
266 voice disorder classification and refine voice activity detection. Table 3 also reports the default
267 ranges for each measure for a frame to be considered voiced.

268 The development of accumulated vocal dose measures (Titze et al., 2003) was motivated by the
269 desire to establish safety thresholds regarding exposure of vocal fold tissue to vibration during
270 phonation, analogous to Occupational Safety and Health Administration guidelines for auditory noise
271 and mechanical vibration exposure. The three most frequently used vocal dose measures to quantify
272 accumulated daily voice use are phonation time, cycle dose, and distance dose. Phonation (voiced)
273 time reflects the cumulative duration of vocal fold vibration, also expressed as a percentage of total
274 monitoring time. The cycle dose is an estimate of the number of vocal fold oscillations during a given
275 period of time. Finally, the distance dose estimates the total distance traveled by the vocal folds,
276 combining cycle dose with vocal fold vibratory amplitude based on the estimates of acoustic sound
277 pressure level.

278 Additionally, attempts were made to characterize vocal load and recovery time by tracking the
279 occurrences and durations of contiguous voiced and non-voiced segments. From these data,
280 occurrence and accumulation histograms provide a summary of voicing and silence characteristics
281 over the course of a monitored period (Titze et al., 2007). To further quantify vocal loading,

282 smoothing was performed over the binary vector of voicing decisions such that contiguous voiced
 283 segments were connected if they were close to each other based on a given duration threshold
 284 (typically less than 0.5 s). The derived contiguous segments approximate speech phrase segments
 285 produced on single breaths to begin to investigate respiratory factors in voice disorders (Sapienza and
 286 Stathopoulos, 1994).

287 Amplitude, frequency, and vocal dose features are traditionally believed to be associated with
 288 phonotraumatic hyperfunctional behaviors (e.g., talking loud, at an inappropriate pitch, or too much
 289 without enough voice rest) (Roy and Hendarto, 2005; Karkos and McCormick, 2009). However, our
 290 previous work demonstrated that overall average signal amplitude, fundamental frequency, and vocal
 291 dose measures were not different between 35 patients with vocal fold nodules or polyps and their
 292 matched-controls (Van Stan et al., 2015b). The results provided in this manuscript replicate our
 293 previous findings with a larger group of 51 matched pairs and extend the analysis approach by (1)
 294 adding novel measures related to voice quality and (2) completing novel comparisons among patients
 295 with non-phonotraumatic vocal hyperfunction versus matched controls and between both sets of
 296 patients with vocal hyperfunction.

297 **2.4. Estimating Aerodynamic Properties from the Accelerometer Signal**

298 Subglottal impedance based inverse filtering (IBIF) is a biologically-inspired acoustic transmission
 299 line model that allows for the estimation of glottal airflow from neck-surface acceleration (Zañartu et
 300 al., 2013). This vocal system model follows a lumped-impedance parameter representation in the
 301 frequency domain using a series of concatenated T-equivalent segments of lumped acoustic elements
 302 that relate acoustic pressure to airflow. Each segment includes terms for representing key
 303 components for the subglottal system such as yielding walls (cartilage and soft tissue components),
 304 viscous losses, elasticity, and inertia. Then, a cascade connection is used to account for the acoustic
 305 transmission associated with the subglottal system based upon symmetric anatomical descriptions for
 306 an average male (Weibel, 1963). In addition, a radiation impedance is used to account for neck skin
 307 properties (Franke, 1951; Ishizaka et al., 1975) and accelerometer loading (Wodicka et al., 1989). The
 308 DC level of the airflow waveform is not modeled by IBIF due to the accelerometer waveform only
 309 being an AC signal. Thus, this overall approach provides an airflow-to-acceleration transfer function
 310 that is inverted when processing the accelerometer signal.

311 Subject-specific parameters need to be obtained to use subglottal IBIF as a signal processing
 312 approach for the accelerometer signal. Five parameters are estimated for each subject—three
 313 parameters for the skin model (skin inertance, resistance, and stiffness) and two parameters for
 314 tracheal geometry (tracheal length and accelerometer position relative to the glottis). The most
 315 relevant parameter values are searched for using an optimization scheme that minimizes the mean-
 316 squared error between oral airflow-derived and neck surface acceleration-derived glottal airflow
 317 waveforms. A default parameter set is fine tuned to a given subject by means of five scaling factors
 318 Q_i , with $i=1, \dots, 5$, which are designed to be estimated from a stable vowel segment. Since the
 319 subglottal system is assumed to remain the same for all other conditions (loudness, vowels, etc.), the
 320 estimated Q parameters may only need to be obtained once for each subject.

321 The subglottal IBIF scheme was initially evaluated for controlled scenarios that represented different
 322 glottal configurations and voice qualities in sustained vowel contexts (Zañartu et al., 2013). Under
 323 these conditions, a mean absolute error of less than 10% was observed for two glottal airflow
 324 measures of interest: maximum flow declination rate (MFDR) and the peak-to-peak glottal flow (AC
 325 Flow). Recently, the method was adapted for a real-time implementation in the context of ambulatory

326 biofeedback (Llico et al., 2015), but again tested and validated only in sustained vowel contexts.
 327 Therefore, an evaluation of the subglottal IBIF method under continuous speech conditions is a
 328 natural next step. Continuous speech is the scenario where subglottal IBIF has the most potential to
 329 contribute to the field of voice assessment, as it can provide aerodynamic measures in the context of
 330 an ambulatory assessment of vocal function.

331 In this paper, we provide an initial assessment of the performance of the subglottal IBIF scheme for
 332 the phonetically-balanced Rainbow Passage obtained in the laboratory, as well as for the data
 333 obtained from a weeklong recording in the field. Multiple measures of vocal function were extracted
 334 from each cycle and averaged over 50-ms frames (50% overlap), including AC Flow, MFDR, open
 335 quotient (OQ), speed quotient (SQ), spectral slope (H1-H2), and normalized amplitude quotient
 336 (NAQ). Figure 4 illustrates the extraction of these measures from the inverse-filtered acceleration
 337 waveform in the time and spectral domains. OQ is defined as $t_o/(t_o + t_c)$, and SQ is defined as
 338 $100(t_{op}/t_{cp})$. NAQ is a measure of the closing phase and is defined as the ratio of AC Flow to MFDR
 339 normalized by the period duration ($t_o + t_c$) (Alku et al., 2002).

340 The in-laboratory voice assessment described in Section 2.2.1 enables a direct comparison of the
 341 subglottal IBIF of neck-surface acceleration with vocal tract inverse-filtering of the oral airflow
 342 waveform. It is noted that inverse filtering of oral airflow for time-varying, continuous speech
 343 segments is a topic of research unto itself, and there are no clear guidelines to best approach the
 344 problem. Thus, we selected a simple but clinically-relevant method of oral airflow processing based
 345 on single formant inverse filtering (Perkell et al., 1991) that has been used for the assessment of
 346 vocal function in speakers with and without a voice disorder (Hillman et al., 1989; Perkell et al.,
 347 1994; Holmberg et al., 1995). Subglottal IBIF with a single set of Q parameters was used to estimate a
 348 continuous glottal airflow signal for each speaker's ambulatory time series.

349 **2.5. Machine Learning and Pattern Recognition Approaches**

350 Machine learning and pattern recognition approaches have become strong tools in the analysis of
 351 time series data. This has been particularly true in wireless health monitoring (Clifford and Clifton,
 352 2012), where multiple levels of analysis are needed to abstract a clinically-relevant diagnosis or state.
 353 Learning problems can be mapped onto a set of four general components: 1) choice of training data
 354 and evaluation method, 2) representation of examples (often called feature engineering), 3) choice of
 355 objective function and constraints, and 4) choice of optimization method. Choosing these
 356 components should be dictated by the goal at hand and the type of data available.

357 We first considered the case of patients with phonotraumatic vocal hyperfunction prior to any
 358 treatment and their matched controls. Each subject (patient or control) had a week of ambulatory
 359 neck-surface acceleration data related to voice use. Previous work suggested that long-term averages
 360 of standard voice measures did not capture differences between patients with vocal fold nodules or
 361 polyps and their matched controls (Mehta et al., 2012a). Thus we hypothesize that the tissue
 362 pathology (nodules or polyps) could create aggregate differences at the extremes of the recorded time
 363 series rather than at the averages. We had some initial success examining whether statistical features
 364 of fundamental frequency (f_0) and SPL, such as skewness, kurtosis, 5th percentile, and 95th percentile,
 365 could capture this more extreme information and lead to an accurate patient classifier in our
 366 population.

367 Briefly, we first extracted SPL, f_0 , and voice quality measures described in Section 2.3 from 50-ms,
 368 non-overlapping frames. From these frames, we built 5-minute, non-overlapping windows (i.e., 6000

369 frames per window) over each day in a subject’s entire weeklong record. We then took univariate
 370 statistics of feature histograms and the cumulative vocal dose measures from windows containing at
 371 least 30 frames labeled as voiced (0.5% phonation time). Normalized versions of the statistics were
 372 obtained by converting each statistic into units of standard deviation based on that feature’s baseline
 373 distribution over an average hour in the first half of the day. Additional methodological details are
 374 available in a previous publication (Ghassemi et al., 2014).

375 Here, a concatenated feature matrix represented each subject’s week. The features from each 5-
 376 minute window were associated with a patient or control label and used to create an L1-regularized
 377 logistic regression using a least absolute shrinkage and selection operator (LASSO) model. The
 378 LASSO model was used to classify 5-minute windows from a held-out set of data from patient and
 379 control subjects. We used leave-one-out-cross-validation (LOOCV) to partition our dataset of 51
 380 paired adult female subjects into 51 training and test sets such that a single patient-control pair was
 381 the held-out test set at each of the 51 iterations. If more than a given proportion of the test subject’s
 382 windows were classified with a patient label, we predicted that subject as being a patient; otherwise,
 383 the subject was classified as a normal control. Classification performance was evaluated across the 51
 384 LASSO models by the proportion of the test set correctly predicted, as well as by the area under the
 385 receiver operating characteristic curve (AUC), F-score, sensitivity (correct labeling of patients), and
 386 specificity (correct labeling of controls).

387 **3. Results**

388 Selected results from applying the three analysis approaches to the current data set of phonotraumatic
 389 and non-phonotraumatic vocal hyperfunction groups are reported as an initial demonstration of the
 390 potential discriminative performance and predictive power of these methods. Patients and their
 391 matched control subjects continue to be enrolled and followed throughout their treatment stages.

392 **3.1. Summary Statistics of Voice Quality Measures and Vocal Dose**

393 Figure 5 illustrates a daylong voice use profile of a 34-year-old adult female psychologist prior to
 394 surgery for a left vocal fold polyp and right vocal fold reactive nodule. Phonation time for her day
 395 reached 20.3% with a mean (SD) SPL of 81.8 (6.4) dB SPL and f0 mode (SD) of 194.5 (51.2) Hz.
 396 Such visualizations (made interactive through navigable graphical user interfaces) of measures such
 397 those described in Section 2.3 may ultimately enable clinicians to identify certain patterns of voice
 398 features related to vocal hyperfunction and subsequently make informed decisions regarding patient
 399 management.

400 As an initial description of the pre-treatment patient data, summary statistics were computed from the
 401 weeklong time series of SPL, f0, voice quality features, and vocal dose measures. The 5th percentile
 402 and 95th percentiles were used to compute minimum, maximum, and range statistics. A four-factor,
 403 one-way analysis of variance was carried out for each summary statistic in the comparison of the two
 404 patient groups and their respective matched-control groups. The between-group comparisons
 405 consisted of the phonotraumatic patients versus their matched controls (51 pairs), the non-
 406 phonotraumatic patients versus their matched controls (20 pairs), and the phonotraumatic group
 407 versus the non-phonotraumatic group.

408 Table 4 reports the group-based mean (SD) for voice use summary statistics of SPL, f0, and vocal
 409 dose measures for weeklong data collected from the phonotraumatic patient and matched-control
 410 groups and the non-phonotraumatic patient and matched-control groups. Based on a *post hoc*
 411 analysis, measures that exhibited statistically significant differences between the two patient groups

412 are highlighted and significant differences between patient and matched-control groups are boxed.
 413 The table also reports voice quality summary statistics of the autocorrelation peak magnitude,
 414 harmonic spectral tilt, LH ratio, and CPP.

415 Individuals with vocal fold nodules and/or polyps exhibited statistically significant differences
 416 compared to individuals with muscle tension dysphonia for all parameters except f_0 . Of note, except
 417 for a few instances, the patient groups and their respective matched-control groups had remarkably
 418 similar accumulated/averaged measurement values (i.e., few statistically significant differences).
 419 These results replicate previously reported findings that, on average, individuals with nodules or
 420 polyps do not speak more often, at a different vocal intensity, or at a different habitual pitch
 421 compared to matched individuals with healthy voices (Van Stan et al., 2015b). Furthermore, the
 422 results provide initial evidence that patients with muscle tension dysphonia also do not differ in these
 423 metrics compared to their matched controls (although CPP trended toward being higher in the
 424 normative group). More sensitive approaches are thus warranted to increase the discriminatory power
 425 among the groups, and the applications of the next two analysis frameworks yield promising,
 426 complementary perspectives.

427 **3.2. Examples of Subglottal Impedance-Based Inverse Filtering**

428 The results of both in-laboratory and in-field assessments are illustrated for a single normal female
 429 subject. The subglottal IBIF yielded estimates of glottal airflow from the neck surface accelerometer
 430 for both assessments. Figure 6 shows a direct contrast of the glottal airflow estimates from oral
 431 airflow and neck-surface acceleration for a portion of the Rainbow Passage. Both waveforms and
 432 derived measures are presented, where it can be seen that, although the fit between signals can be
 433 adequate, the IBIF-based signal is less prone to inverse filtering artifacts than its oral airflow-based
 434 counterpart. This is due to the more stationary underlying dynamic behavior of the subglottal system
 435 relative to that of the time-varying vocal tract, thus constituting a more tractable inverse filtering
 436 problem. As a result, the measures of vocal function derived from the subglottal IBIF processing
 437 appear to be more reliable. Improving upon methods for inverse filtering of oral airflow in running
 438 speech is a current focus of research, which would also allow for testing the assumption that Q
 439 parameters in the IBIF scheme should remain constant in continuous speech conditions.

440 Figure 7 presents histograms of SPL and MFDR derived from the weeklong neck-surface
 441 acceleration recording. The SPL/MFDR relation provides insights on the efficiency in voice
 442 production, which was found to be 9 dB per MFDR doubling in sustained vowels for normal female
 443 subjects (6 dB per MFDR doubling for male subjects) (Holmberg et al., 1988). It is noted in Figure 7
 444 that when a linear scale is used for MFDR, the histogram peak appears skewed to the left. However,
 445 when applying a logarithmic transform to MFDR (Holmberg et al., 1988; Holmberg et al., 1995), both
 446 SPL and MFDR histograms become Gaussian with different means and variances. The ambulatory
 447 relation provides a slope of 1.13 dB/dB, which is similar to the 1.5 dB/dB slope (9 dB per MFDR
 448 doubling) reported for oral airflow-based inverse filtering features under sustained vowel conditions
 449 (Holmberg et al., 1988). This result is encouraging as it provides initial validation for ambulatory
 450 MFDR estimation using subglottal IBIF and also provides an indication that average behaviors in
 451 normal subjects could be related to simple sustained vowel tasks in a clinical assessment. The
 452 relationship warrants further investigation, with challenges foreseen for subjects with voice disorders.

453 3.3. Classification Results Using Machine Learning

454 Figure 8 shows that we were able to correctly classify 74 out of 102 subjects (72.5%) using a
 455 threshold of 0.68. Intuitively, this means that a subject is predicted to be a patient with
 456 phonotraumatic vocal hyperfunction if more than 68% of their windows were classified similarly to
 457 those from the other patients the LASSO model was trained on. The mean (standard deviation) of
 458 performance across the 51 LASSO models was 0.739 (0.274) for AUC, 0.766 (0.204) for F-score,
 459 0.739 (0.296) for sensitivity, and 0.767 (0.288) for specificity.

460 Table 5 summarizes the performance of the statistical measures in classifying phonotraumatic vocal
 461 hyperfunction. As shown, subjects with vocal fold nodules tended to have f_0 and SPL distributions
 462 that were right-shifted from their previous values, i.e., an increased Normalized F_0 95th percentile
 463 and an increased Normalized SPL Skew. We contrast this with the vocally normal group, which had
 464 a right-shifted (non-normalized) SPL distribution, i.e., increased SPL Skew. We could interpret the
 465 right-shifting of Normalized features in subjects with vocal fold nodules to mean that they tended to
 466 deviate from their baseline f_0 and SPL as their days progressed, possibly reflecting increased
 467 difficulty in producing phonation. For the controls, the fact that their absolute SPL Skew was
 468 increased without a corresponding increase to their Normalized distribution suggests that even when
 469 control subjects exhibited higher SPL ranges, they tended to stay within their baseline ranges.

470 While a majority of subjects were correctly classified in this framework, the predicted labels for
 471 some subjects are notably incorrect. One possible reason the classification is more accurate for the
 472 patient versus the control group (19 incorrectly labeled patients versus 9 incorrectly labeled controls)
 473 might stem from our strong labeling assumptions. It is likely that not all frames (and therefore not all
 474 statistical features of 5-minute windows) of a patient exhibit vocal behavior associated with
 475 phonotraumatic hyperfunction. This creates a potentially large set of false-positive labels that can
 476 cause classification bias.

477 4. Discussion

478 An understanding of daily behavior is essential to improving the diagnosis and treatment of
 479 hyperfunctional voice disorders. Our results indicate that supervised machine learning techniques
 480 have the potential to be used to discriminate patients from control subjects with normal voice. It is
 481 important to note, however, that this work did not account for time of day, sequence of window
 482 occurrence, or ordered loading of features. For an example of time-ordered analysis, Figure 9 shows
 483 a three-dimensional distribution showing the occurrence histograms of unvoiced segment durations
 484 that immediately followed successively longer voiced-segment durations over the course of a day.
 485 This analysis approach attempts to reflect a speaker's vocal behavior in terms of how much voice rest
 486 follows bursts of voicing activity. Similarly, ongoing monitoring of phonation time after a particular
 487 vocal load in a preceding window represents additional methods that may lead to complementary
 488 pieces of information that can aid in the successful detection of hyperfunctional vocal behaviors.

489 The subglottal IBIF measures for continuous speech appear more accurate than the oral airflow based
490 due to the additional challenges associated with performing time-varying inverse filtering for the
491 vocal tract. Improving upon methods for inverse filtering of oral airflow in continuous speech is a
492 current focus of research, which would also allow for testing the assumption that Q parameters
493 remain constant during speech production. The evaluation of subglottal IBIF using weeklong
494 ambulatory data acquired with the VHM illustrates that the relation between SPL and MFDR is very
495 well aligned with previous observations for sustained vowels for adult female subjects (Holmberg,
496 Hillman, and Perkell 1988). This result provides initial validation of using IBIF to estimate MFDR
497 from the acceleration signal; however, further analysis using normative speaker populations and
498 individuals with varying voice disorder severity is required.

499 In order to make the most use of our data without re-using any training data in the test set, we trained
500 51 separate L1-regularized logistic regression LASSO models. For a fair comparison of the collective
501 performance of these models on test input, we used a uniform threshold of 0.5 to classify the output
502 of each 5-minute window passed through the LASSO model. This created a set of predicted binary
503 labels (0, 1) for all windows in any subject's entire record. The proportion of each subject's windows
504 that are classified as a 1 in this process is plotted in Figure 8, ranging from 0 to 100%. For example, a
505 subject very near the top of the graph would have had almost all of their 5-minute windows over the
506 course of the week classified as a 1. Using this output, we can perform inter-model comparisons. In
507 the paper, we report the "optimal threshold" (0.68) that created the highest accuracy measure. It is
508 possible to improve the sensitivity or specificity of our results by lowering or raising this threshold
509 appropriately.

510 One of the most challenging aspects of voice treatment is achieving carryover (long-term retention)
511 of newly established vocal behaviors from the clinical setting into the patient's daily environment
512 (Ziegler et al., 2014). Adding biofeedback capabilities to an ambulatory monitor has significant
513 potential to address this carryover challenge by providing individuals with timely information about
514 their vocal behavior throughout their typical activities of daily living. Pilot work has shown that
515 speakers with normal voices exhibit a biofeedback effect by modifying their SPL levels in response
516 to cueing from an ambulatory voice monitoring device (Van Stan et al., 2015a). Long-term retention,
517 however, was not observed and may require the use of alternative biofeedback schedules (e.g.,
518 decreasing the frequency and delaying the presentation of biofeedback) that have been well-studied
519 in the motor learning literature.

520 **5. Conclusion**

521 Wearable voice monitoring systems have the potential to provide more reliable and objective
522 measures of voice use that can enhance the diagnostic and treatment strategies for common voice
523 disorders. This report provided an overview of our group's approach to the multilateral
524 characterization and classification of common types of voice disorders using a smartphone-based
525 ambulatory voice health monitor. Preliminary results illustrate the potential for the three analysis
526 approaches studied to help improve assessment and treatment for hyperfunctional voice disorders.
527 Delineating detrimental vocal behaviors may aid in providing real-time biofeedback to a speaker to
528 facilitate the adoption of healthier voice production into everyday use.

529 **Acknowledgments**

530 The authors acknowledge the contributions of R. Petit for aid in designing and programming the
531 smartphone application; M. Bresnahan, D. Buckley, M. Cooke, and A. Fryd, for data segmentation

532 assistance; J. Kobler and J. Heaton for help with voice monitor system design; C. Andrieu and F.
 533 Simond for Android audio codec advice; and J. Rosowski and M. Ravicz for use of their
 534 accelerometer calibration system. This work was supported by the Voice Health Institute and the
 535 National Institutes of Health (NIH) National Institute on Deafness and Other Communication
 536 Disorders under Grants R33 DC011588 and F31 DC014412. The paper's contents are solely the
 537 responsibility of the authors and do not necessarily represent the official views of the NIH.
 538 Additional support received from MIT-Chile grant 2745333 through the MIT International Science
 539 and Technology Initiatives (MISTI) program, Chilean CONICYT grants FONDECYT 11110147 and
 540 Basal FB0008, and scholarships from CONICYT, Universidad Federico Santa María, and
 541 Universidad de Chile. Further funding provided by the Intel Science and Technology Center for Big
 542 Data and the National Library of Medicine Biomedical Informatics Research Training Grant
 543 (NIH/NLM 2T15 LM007092-22).

544 **References**

- 545 Alku, P., Backstrom, T., and Vilkman, E. (2002). Normalized amplitude quotient for parametrization
 546 of the glottal flow. *J. Acoust. Soc. Am.* 112, 701–710.
- 547 Awan, S.N., Roy, N., Jetté, M.E., Meltzner, G.S., and Hillman, R.E. (2010). Quantifying dysphonia
 548 severity using a spectral/cepstral-based acoustic index: Comparisons with auditory-perceptual
 549 judgements from the CAPE-V. *Clin. Linguist. Phon.* 24, 742–758.
- 550 Bhattacharyya, N. (2014). The prevalence of voice problems among adults in the United States.
 551 *Laryngoscope* 124, 2359–2362.
- 552 Carroll, T., Nix, J., Hunter, E., Emerich, K., Titze, I., and Abaza, M. (2006). Objective measurement
 553 of vocal fatigue in classical singers: A vocal dosimetry pilot study. *Otolaryngol. Head. Neck.*
 554 *Surg.* 135, 595–602.
- 555 Clifford, G.D., and Clifton, D. (2012). Wireless technology in disease management and medicine.
 556 *Annu. Rev. Med.* 63, 479–492.
- 557 Czerwonka, L., Jiang, J.J., and Tao, C. (2008). Vocal nodules and edema may be due to vibration-
 558 induced rises in capillary pressure. *Laryngoscope* 118, 748–752.
- 559 Fairbanks, G. (1960). *Voice and Articulation Drillbook*. New York: Harper and Row.
- 560 Franke, E.K. (1951). Mechanical impedance of the surface of the human body. *J. Appl. Physiol.* 3,
 561 582–590.
- 562 Ghassemi, M., Van Stan, J.H., Mehta, D.D., Zañartu, M., Cheyne Ii, H.A., Hillman, R.E., and Guttag,
 563 J.V. (2014). Learning to detect vocal hyperfunction from ambulatory neck-surface
 564 acceleration features: Initial results for vocal fold nodules. *IEEE Trans. Biomed. Eng.* 61,
 565 1668–1675.
- 566 Hadjitodorov, S., Boyanov, B., and Teston, B. (2000). Laryngeal pathology detection by means of
 567 class-specific neural maps. *IEEE Trans. Inf. Technol. Biomed.* 4, 68–73.

Ambulatory monitoring of voice disorders

- 568 Hillman, R.E., Holmberg, E.B., Perkell, J.S., Walsh, M., and Vaughan, C. (1989). Objective
569 assessment of vocal hyperfunction: An experimental framework and initial results. *J. Speech*
570 *Hear. Res.* 32, 373–392.
- 571 Hillman, R.E., Holmberg, E.B., Perkell, J.S., Walsh, M., and Vaughan, C. (1990). Phonatory function
572 associated with hyperfunctionally related vocal fold lesions. *J. Voice* 4, 52–63.
- 573 Hogikyan, N.D., and Sethuraman, G. (1999). Validation of an instrument to measure voice-related
574 quality of life (V-RQOL). *J. Voice* 13, 557–569.
- 575 Holmberg, E.B., Hillman, R.E., and Perkell, J.S. (1988). Glottal airflow and transglottal air pressure
576 measurements for male and female speakers in soft, normal, and loud voice. *J. Acoust. Soc.*
577 *Am.* 84, 511–529.
- 578 Holmberg, E.B., Hillman, R.E., Perkell, J.S., Guiod, P.C., and Goldman, S.L. (1995). Comparisons
579 among aerodynamic, electroglottographic, and acoustic spectral measures of female voice. *J.*
580 *Speech Hear. Res.* 38, 1212–1223.
- 581 Ishizaka, K., French, J., and Flanagan, J.L. (1975). Direct determination of vocal tract wall
582 impedance. *IEEE Transactions on Acoustics, Speech and Signal Processing* 23, 370–373.
- 583 Karkos, P.D., and McCormick, M. (2009). The etiology of vocal fold nodules in adults. *Current*
584 *Opinion in Otolaryngology & Head & Neck Surgery* 17, 420–423.
- 585 Kempster, G.B., Gerratt, B.R., Verdolini Abbott, K., Barkmeier-Kraemer, J., and Hillman, R.E.
586 (2009). Consensus auditory-perceptual evaluation of voice: Development of a standardized
587 clinical protocol. *Am. J. Speech Lang. Pathol.* 18, 124–132.
- 588 Little, M.A., Mcsharry, P.E., Roberts, S.J., Costello, D.A., and Moroz, I.M. (2007). Exploiting
589 nonlinear recurrence and fractal scaling properties for voice disorder detection. *Biomed. Eng.*
590 *Online* 6, 23.
- 591 Llico, A.F., Zañartu, M., González, A.J., Wodicka, G.R., Mehta, D.D., Van Stan, J.H., and Hillman,
592 R.E. (2015). Real-time estimation of aerodynamic features for ambulatory voice biofeedback.
593 *J. Acoust. Soc. Am.* 138, EL14–EL19.
- 594 Mehta, D.D., and Hillman, R.E. (2012). Current role of stroboscopy in laryngeal imaging. *Curr.*
595 *Opin. Otolaryngol. Head Neck Surg.* 20, 429–436.
- 596 Mehta, D.D., Woodbury Listfield, R., Cheyne Ii, H.A., Heaton, J.T., Feng, S.W., Zañartu, M., and
597 Hillman, R.E. (2012a). Duration of ambulatory monitoring needed to accurately estimate
598 voice use. *Proceedings of InterSpeech: Annual Conference of the International Speech*
599 *Communication Association.*
- 600 Mehta, D.D., Zañartu, M., Feng, S.W., Cheyne Ii, H.A., and Hillman, R.E. (2012b). Mobile voice
601 health monitoring using a wearable accelerometer sensor and a smartphone platform. *IEEE*
602 *Trans. Biomed. Eng.* 59, 3090–3096.

- 603 Mehta, D.D., Zaňartu, M., Quatieri, T.F., Deliyiski, D.D., and Hillman, R.E. (2011). Investigating
604 acoustic correlates of human vocal fold vibratory phase asymmetry through modeling and
605 laryngeal high-speed videoendoscopy. *J. Acoust. Soc. Am.* 130, 3999–4009.
- 606 Mehta, D.D., Zaňartu, M., Van Stan, J.H., Feng, S.W., Cheyne Ii, H.A., and Hillman, R.E. (2013).
607 Smartphone-based detection of voice disorders by long-term monitoring of neck acceleration
608 features. *Proceedings of the 10th Annual Body Sensor Networks Conference*.
- 609 Mehta, D.D., Zeitels, S.M., Burns, J.A., Friedman, A.D., Deliyiski, D.D., and Hillman, R.E. (2012c).
610 High-speed videoendoscopic analysis of relationships between cepstral-based acoustic
611 measures and voice production mechanisms in patients undergoing phonemicsurgery. *Ann.*
612 *Otol. Rhinol. Laryngol.* 121, 341–347.
- 613 Nidcd (2012). *2012-2016 Strategic Plan*. Bethesda, MD: National Institute on Deafness and Other
614 Communication Disorders (NIDCD), U.S. Department of Health and Human Services.
- 615 Parsa, V., and Jamieson, D.G. (2000). Identification of pathological voices using glottal noise
616 measures. *J. Speech. Lang. Hear. Res.* 43, 469–485.
- 617 Perkell, J.S., Hillman, R.E., and Holmberg, E.B. (1994). Group differences in measures of voice
618 production and revised values of maximum airflow declination rate. *J. Acoust. Soc. Am.* 96,
619 695–698.
- 620 Perkell, J.S., Holmberg, E.B., and Hillman, R.E. (1991). A system for signal-processing and data
621 extraction from aerodynamic, acoustic, and electroglottographic signals in the study of voice
622 production. *J. Acoust. Soc. Am.* 89, 1777–1781.
- 623 Rothenberg, M. (1973). A new inverse filtering technique for deriving glottal air flow waveform
624 during voicing. *J. Acoust. Soc. Am.* 53, 1632–1645.
- 625 Roy, N., Barkmeier-Kraemer, J., Eadie, T., Sivasankar, M.P., Mehta, D., Paul, D., and Hillman, R.
626 (2013). Evidence-based clinical voice assessment: A systematic review. *Am. J. Speech Lang.*
627 *Pathol.* 22, 212–226.
- 628 Roy, N., and Bless, D.M. (2000). Personality traits and psychological factors in voice pathology: A
629 foundation for future research. *J. Speech. Lang. Hear. Res.* 43, 737–748.
- 630 Roy, N., and Hendarto, H. (2005). Revisiting the pitch controversy: Changes in speaking
631 fundamental frequency (SFF) after management of functional dysphonia. *J. Voice* 19, 582–
632 591.
- 633 Roy, N., Merrill, R.M., Gray, S.D., and Smith, E.M. (2005). Voice disorders in the general
634 population: Prevalence, risk factors, and occupational impact. *Laryngoscope* 115, 1988–1995.
- 635 Sapienza, C.M., and Stathopoulos, E.T. (1994). Respiratory and laryngeal measures of children and
636 women with bilateral vocal fold nodules. *J. Speech. Lang. Hear. Res.* 37, 1229–1243.
- 637 Švec, J.G., Titze, I.R., and Popolo, P.S. (2005). Estimation of sound pressure levels of voiced speech
638 from skin vibration of the neck. *J. Acoust. Soc. Am.* 117, 1386–1394.

Ambulatory monitoring of voice disorders

- 639 Titze, I.R., Hunter, E.J., and Švec, J.G. (2007). Voicing and silence periods in daily and weekly
640 vocalizations of teachers. *J. Acoust. Soc. Am.* 121, 469–478.
- 641 Titze, I.R., Švec, J.G., and Popolo, P.S. (2003). Vocal dose measures: Quantifying accumulated
642 vibration exposure in vocal fold tissues. *J. Speech. Lang. Hear. Res.* 46, 919–932.
- 643 Van Stan, J.H., Mehta, D.D., and Hillman, R.E. (2015a). The effect of voice ambulatory biofeedback
644 on the daily performance and retention of a modified vocal motor behavior in participants
645 with normal voices. *J. Speech. Lang. Hear. Res.* ePub, 1–9.
- 646 Van Stan, J.H., Mehta, D.D., Zeitels, S.M., Burns, J.A., Barbu, A.M., and Hillman, R.E. (2015b).
647 Average ambulatory measures of sound pressure level, fundamental frequency, and vocal
648 dose do not differ between adult females with phonotraumatic lesions and matched control
649 subjects. *Ann. Otol. Rhinol. Laryngol.* ePub, 1–11.
- 650 Weibel, E.R. (1963). *Morphometry of the Human Lung*, 1st ed. New York: Springer. p. 139.
- 651 Wodicka, G.R., Stevens, K.N., Golub, H.L., Cravalho, E.G., and Shannon, D.C. (1989). A model of
652 acoustic transmission in the respiratory system. *IEEE Trans. Biomed. Eng.* 36, 925–934.
- 653 Zañartu, M., Ho, J.C., Mehta, D.D., Hillman, R.E., and Wodicka, G.R. (2013). Subglottal impedance-
654 based inverse filtering of voiced sounds using neck surface acceleration. *IEEE Trans. Audio
655 Speech Lang. Processing* 21, 1929–1939.
- 656 Ziegler, A., Dastolfo, C., Hersan, R., Rosen, C.A., and Gartner-Schmidt, J. (2014). Perceptions of
657 voice therapy from patients diagnosed with primary muscle tension dysphonia and benign
658 mid-membranous vocal fold lesions. *J. Voice* 28, 742–752.
- 659

660 TABLES

661 Table 1. Occupations of adult females with phonotraumatic vocal hyperfunction and matched-control
 662 participants analyzed to date (51 pairs). Diagnoses for the patient group are also listed for each
 663 occupation.

| Occupation | No. Subject Pairs | Patient Diagnosis |
|----------------------------------|-------------------|---------------------------|
| Singer | 37 | Nodules (32) Polyp (5) |
| Teacher | 5 | Nodules |
| Consultant | 2 | Nodules (1) Polyp (1) |
| Psychotherapist/ Psychologist | 2 | Nodules |
| Recruiter | 2 | Nodules |
| Marketer | 1 | Nodules |
| Media relations | 1 | Nodules |
| Registered nurse | 1 | Polyp |

664

Provisional

665 Table 2. Occupations of adult females with non-phonotraumatic vocal hyperfunction and matched-
666 control participants analyzed (20 pairs). All patients were diagnosed with muscle tension dysphonia.

667

| Occupation | No. Subject Pairs |
|--------------------------|-------------------|
| Registered nurse | 3 |
| Singer | 3 |
| Teacher | 3 |
| Administrator | 2 |
| At-home caregiver | 2 |
| Student | 2 |
| Social worker | 1 |
| Actress | 1 |
| Administrative assistant | 1 |
| Exercise instructor | 1 |
| Systems analyst | 1 |

Provisional

668 Table 3. Description of frame-based signal features computed on in-field ambulatory voice data.

| Feature | Units | Voicing criteria | Description |
|--------------------------------|-----------|------------------|-------------------------------------------------------------------------------------------------------------|
| Sound pressure level at 15 cm | dB SPL | 45–130 | Acceleration amplitude mapped to acoustic sound pressure level (Švec et al., 2005) |
| Fundamental frequency | Hz | 70–1000 | Reciprocal of first non-zero peak location in the normalized autocorrelation function (Mehta et al., 2012b) |
| Autocorrelation peak amplitude | | 0.60–1 | Relative amplitude of first non-zero peak in the normalized autocorrelation function (Mehta et al., 2012b) |
| Subharmonic peak | | 0.25–1 | Relative amplitude of a secondary peak, if it exists, located around half way to the autocorrelation peak |
| Harmonic spectral tilt | dB/octave | –25–0 | Linear regression slope over the first 8 spectral harmonics (Mehta et al., 2011) |
| Low-to-high spectral ratio | dB | 22–50 | Difference between spectral power below and above 2000 Hz (Awan et al., 2010) |
| Cepstral peak prominence | dB | 10–35 | Magnitude of the highest peak in the power cepstrum (Mehta et al., 2012c) |
| Zero crossing rate | | 0–1 | Proportion of frame that signal crosses its mean |

669

Provisional

Ambulatory monitoring of voice disorders

670 Table 4. Group-based mean (SD) of summary statistics of weeklong vocal dose and voice quality
 671 data collected from adult females in the phonotraumatic vocal hyperfunction (n = 51) and non-
 672 phonotraumatic vocal hyperfunction (n = 20) patient groups. Statistically significant differences
 673 between means are highlighted (p < 0.001). Minimum, maximum, and range are trimmed estimators
 674 reporting 5th percentile, 95th percentile, and range of the middle 90% of the data, respectively.

| Summary statistic | Phonotraumatic controls | Phonotraumatic group | Non-phonotraumatic group | Non-phonotraumatic controls |
|----------------------------------------|-------------------------|----------------------|--------------------------|-----------------------------|
| Monitoring duration (hh:mm:ss) | 81:11:49 (13:13:35) | 77:21:43 (15:36:33) | 73:44:37 (10:04:12) | 78:59:16 (13:50:13) |
| <i>SPL (dB SPL re 15 cm)</i> | | | | |
| Mean | 83.9 (4.6) | 85.2 (4.1) | 80.1 (6.0) | 83.0 (5.2) |
| Standard deviation | 12.5 (2.4) | 11.8 (1.9) | 9.9 (3.1) | 11.2 (3.3) |
| Minimum | 62.7 (5.8) | 64.5 (4.9) | 63.3 (7.0) | 64.5 (6.3) |
| Maximum | 104.2 (6.7) | 103.5 (5.9) | 96.3 (8.3) | 101.7 (9.5) |
| Range | 41.4 (8.5) | 39.0 (6.7) | 33.0 (10.6) | 37.2 (11.6) |
| <i>f0 (Hz)</i> | | | | |
| Mode | 201.4 (19.1) | 197.2 (22.3) | 193.8 (31.1) | 192.9 (25.7) |
| Standard deviation | 89.6 (17.5) | 75.3 (17.3) | 73.5 (24.9) | 70.1 (14.3) |
| Minimum | 170.3 (14.9) | 166.7 (17.4) | 160.0 (20.5) | 163.2 (22.2) |
| Maximum | 440.6 (58.9) | 392.4 (65.5) | 382.4 (81.4) | 374.6 (62.3) |
| Range | 270.3 (55.9) | 225.7 (56.7) | 222.4 (81.2) | 211.4 (49.4) |
| <i>Phonation time</i> | | | | |
| Cumulative (hh:mm:ss) | 7:24:08 (2:33:32) | 7:33:45 (2:36:34) | 4:25:14 (2:31:57) | 5:46:13 (2:16:17) |
| Normalized (%) | 9.2 (2.9) | 9.7 (2.6) | 6.0 (3.1) | 7.3 (2.7) |
| <i>Cycle dose</i> | | | | |
| Cumulative (millions of cycles) | 7.121 (2.76) | 6.718 (2.495) | 3.708 (2.202) | 4.814 (1.831) |
| Normalized (cycles/hr) | 87,954 (30,508) | 85,719 (25,633) | 49,892 (26,997) | 61,310 (22,241) |
| <i>Distance dose</i> | | | | |
| Cumulative (m) | 26,769 (11,815) | 26,689 (10,999) | 12,254 (8,284) | 18,084 (8,466) |
| Normalized (m/hr) | 330.0 (129.3) | 340.7 (112.1) | 165.1 (102.4) | 228.0 (98.4) |
| <i>Autocorrelation peak</i> | | | | |
| Mean | 0.851 (0.018) | 0.843 (0.015) | 0.827 (0.022) | 0.837 (0.014) |
| Standard deviation | 0.080 (0.004) | 0.079 (0.004) | 0.082 (0.007) | 0.079 (0.004) |
| Minimum | 0.677 (0.020) | 0.672 (0.016) | 0.657 (0.024) | 0.668 (0.014) |
| Maximum | 0.941 (0.010) | 0.934 (0.011) | 0.926 (0.014) | 0.928 (0.010) |
| Range | 0.263 (0.015) | 0.262 (0.014) | 0.269 (0.021) | 0.260 (0.013) |
| <i>Harmonic spectral tilt (dB/oct)</i> | | | | |
| Mean | -14.1 (0.6) | -14.4 (0.6) | -13.6 (1.1) | -14.1 (0.8) |
| Standard deviation | 2.4 (0.3) | 2.4 (0.2) | 2.5 (0.3) | 2.4 (0.2) |
| Minimum | -17.8 (0.8) | -18.2 (0.8) | -17.5 (1.0) | -17.8 (1.1) |
| Maximum | -9.9 (0.8) | -10.5 (0.6) | -9.3 (1.5) | -9.8 (1.0) |
| Range | 8.0 (1.0) | 7.7 (0.8) | 8.2 (1.2) | 8.0 (0.8) |
| <i>LH ratio (dB)</i> | | | | |
| Mean | 30.5 (1.1) | 30.5 (1.3) | 30.1 (1.3) | 30.7 (1.5) |
| Standard deviation | 4.4 (0.4) | 4.5 (0.4) | 4.1 (0.5) | 4.5 (0.5) |
| Minimum | 24.0 (0.6) | 23.8 (0.7) | 23.8 (0.5) | 24.1 (0.7) |
| Maximum | 38.3 (1.6) | 38.6 (1.8) | 37.3 (2.1) | 38.8 (2.2) |
| Range | 14.3 (1.3) | 14.8 (1.3) | 13.5 (1.7) | 14.7 (1.6) |
| <i>CPP (dB)</i> | | | | |
| Mean | 22.9 (1.0) | 23.2 (1.1) | 21.4 (2.1) | 22.8 (1.1) |
| Standard deviation | 4.5 (0.3) | 4.4 (0.3) | 4.2 (0.5) | 4.4 (0.3) |
| Minimum | 15.1 (0.5) | 15.3 (0.6) | 14.3 (0.8) | 14.9 (0.7) |
| Maximum | 29.6 (1.2) | 29.7 (1.2) | 28.0 (2.3) | 29.3 (1.1) |
| Range | 14.5 (1.0) | 14.4 (0.9) | 13.8 (1.6) | 14.4 (1.0) |

676 Table 5. Association of summary statistics features of sound pressure level (SPL) and fundamental
 677 frequency (f0) with group label across the 51 LASSO models. The maximum number that the
 678 “association count” field can have is 51. This occurs when that particular variable (row) has a
 679 statistically significant effect ($p < 0.001$, absolute average odds ratios ≥ 1.10) in each model. Many
 680 associations persisted across all models and also tended to agree well on the magnitude of the
 681 association. The 95% confidence interval (CI) is from the lowest bound across subsets to the highest
 682 bound across subsets.

| Summary statistic | Association Count | | Multivariate LASSO Association | |
|--------------------------------------------|-------------------|---------|--------------------------------|-----------------------------|
| | Patient | Control | Beta Mean (SD) | Odds Ratio Mean (95% CI) |
| Normalized SPL Skew | 51 | 0 | 1.11 (0.04) | 3.03 (2.72–3.69) |
| Normalized f0 95 th percentile | 51 | 0 | 0.86 (0.03) | 2.36 (2.16–2.70) |
| f0 Skew | 51 | 0 | 0.53 (0.09) | 1.69 (1.42–2.35) |
| Normalized SPL Kurtosis | 51 | 0 | 0.28 (0.02) | 1.32 (1.22–1.44) |
| Normalized SPL 5 th percentile | 51 | 0 | 0.14 (0.03) | 1.16 (1.05–1.30) |
| Normalized Percent Phonation | 51 | 0 | 0.12 (0.02) | 1.13 (1.07–1.20) |
| Normalized F0 5 th percentile | 0 | 50 | -0.10 (0.02) | 0.91 (0.85–1.00) |
| Normalized SPL 95 th percentile | 0 | 51 | -0.17 (0.03) | 0.84 (0.77–0.91) |
| SPL Kurtosis | 0 | 51 | -0.28 (0.02) | 0.76 (0.69–0.82) |
| Normalized f0 Skew | 0 | 51 | -0.41 (0.07) | 0.66 (0.51–0.77) |
| SPL Skew | 0 | 51 | -2.84 (0.12) | 0.06 (0.03–0.08) |

683

Provisional

684 **FIGURES**

685 Figure 1: Treatment tracks for patients exhibiting phonotraumatic and non-phonotraumatic
 686 hyperfunctional vocal behaviors. Week numbers (W1, W2, W3, and W4) refer to time points during
 687 which ambulatory monitoring of voice use is being acquired using the smartphone-based voice health
 688 monitor. The current enrollment of each patient and matched-control pairing is listed above each
 689 week number.

690 Figure 2. In-laboratory data acquisition setup. (A) Synchronized recordings are made of signals from
 691 an acoustic microphone (MIC), electroglottography electrodes (EGG), accelerometer sensor (ACC),
 692 high-bandwidth oral airflow (FLO), and intraoral pressure (PRE). (B) Signal snapshot of a string of
 693 “pae” tokens required for the estimation of subglottal pressure and airflow during phonation.

694 Figure 3: Ambulatory voice health monitor: (A) Smartphone, accelerometer sensor, and interface
 695 cable with circuit encased in epoxy; (B) the wired accelerometer mounted on a silicone pad affixed to
 696 the neck midway between the Adam’s apple and V-shaped notch of the collarbone.

697 Figure 4: Parameterization of the (A) original and (B) inverse-filtered waveforms from the oral
 698 airflow (black) and neck-surface acceleration (ACC, red-dashed) waveform processed with subglottal
 699 impedance-based inverse filtering. Shown are the time waveform, frequency spectrum, and cepstrum,
 700 along with the parameterization of each domain to yield clinically salient measures of voice
 701 production.

702 Figure 5: Illustration of a daily voice use profile for an adult female diagnosed with bilateral vocal
 703 fold nodules. Shown are five-minute moving averages of the median and 95th percentile of frame-
 704 based voice quality measures, along with self-reported ratings of effort, discomfort, and fatigue at the
 705 beginning and end of day. The daylong histograms of each measure are shown to the right of each
 706 time series. The plots below display the occurrence histograms of contiguous voiced segments (left)
 707 and estimates of speech phrases between breaths (right).

708 Figure 6: Time-varying estimation of measures derived from the airflow-derived (black) and
 709 accelerometer-derived (red-dashed) glottal airflow signal using subglottal impedance-based inverse
 710 filtering. Trajectories are shown for an adult female with no vocal pathology for the difference
 711 between the first two harmonic amplitudes (H1-H2), peak-to-peak flow (AC Flow), maximum flow
 712 declination rate (MFDR), open quotient (OQ), speed quotient (SQ), and normalized amplitude
 713 quotient (NAQ).

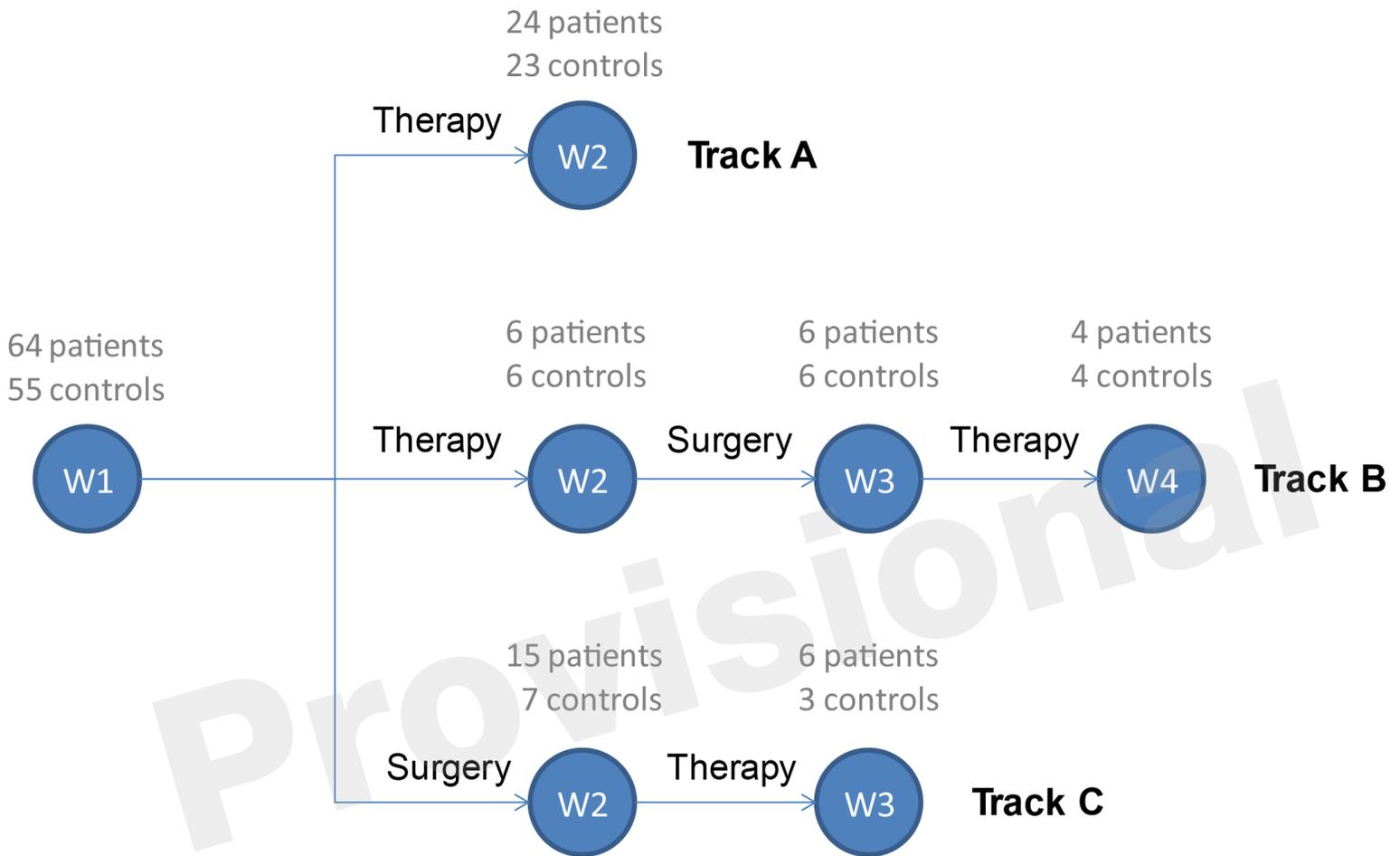
714 Figure 7: Exemplary results using subglottal impedance-based inverse filtering of a weeklong neck-
 715 surface acceleration signal from an adult female with a normal voice. Histograms of the maximum
 716 flow declination rate (MFDR) measure are displayed in physical and logarithmic units. The logarithm
 717 of MFDR is plotted against sound pressure level (SPL) to confirm the expected linear correlation
 718 ($r = 0.94$) and slope (1.13 dB/dB).

719 Figure 8: Classification results on 102 adult female subjects, 51 with vocal fold nodules and 51
 720 matched-control subjects with normal voices. Per-patient unbiased model performance using
 721 summary statistics of sound pressure level and fundamental frequency from non-overlapping, five-
 722 minute windows.

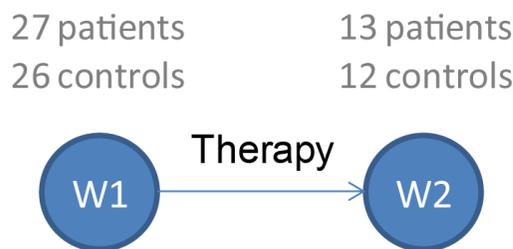
723 Figure 9: Occurrence histogram of voiced/unvoiced contiguous segment pairs. The figure includes
724 the number of times (per hour) that a voiced segment of a given duration is followed by an unvoiced
725 segment of a given duration.

Provisional

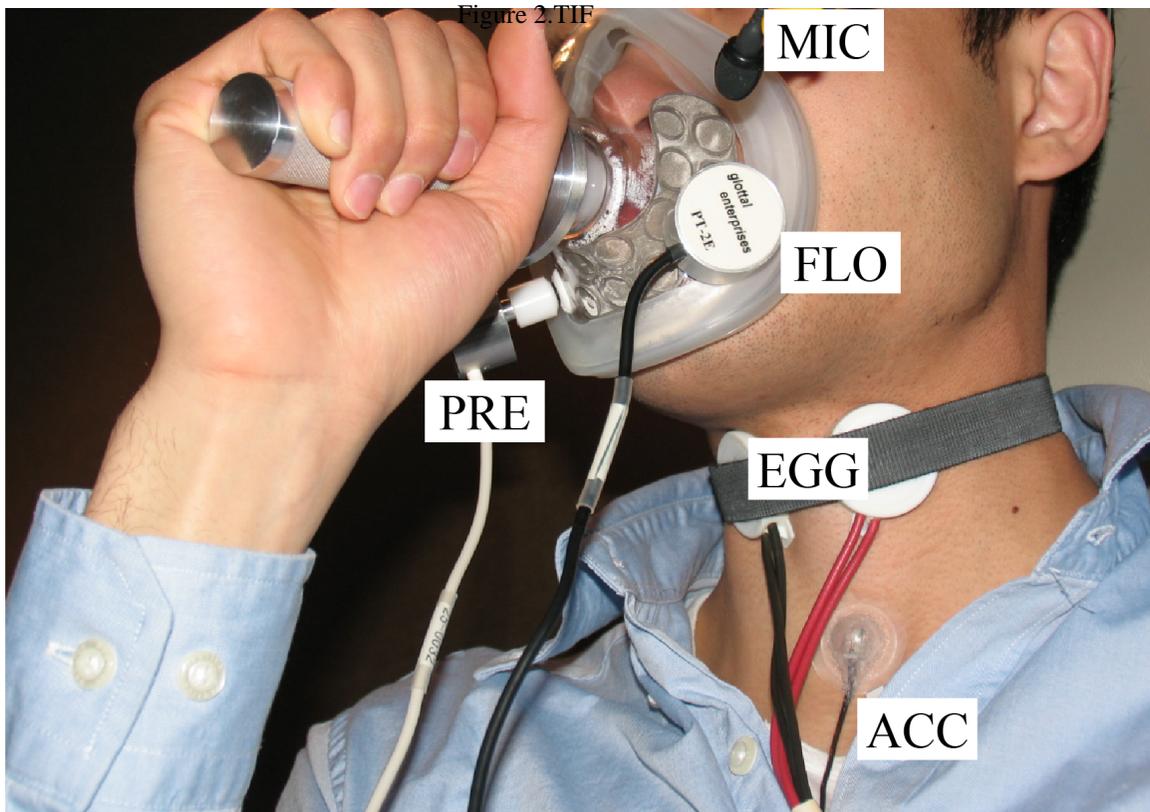
Phonotraumatic vocal hyperfunction



Non-phonotraumatic vocal hyperfunction



A



B

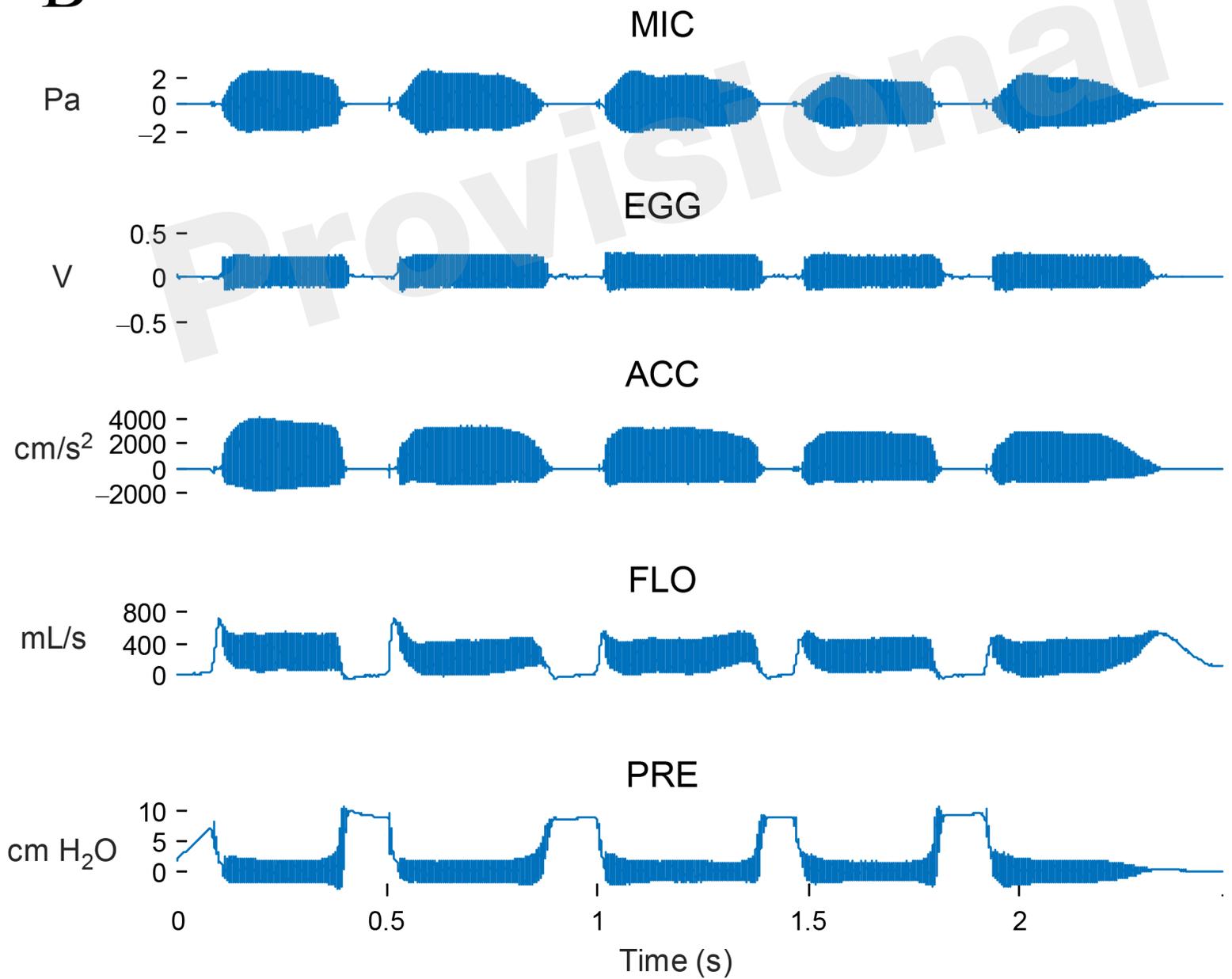


Figure 3.TIF

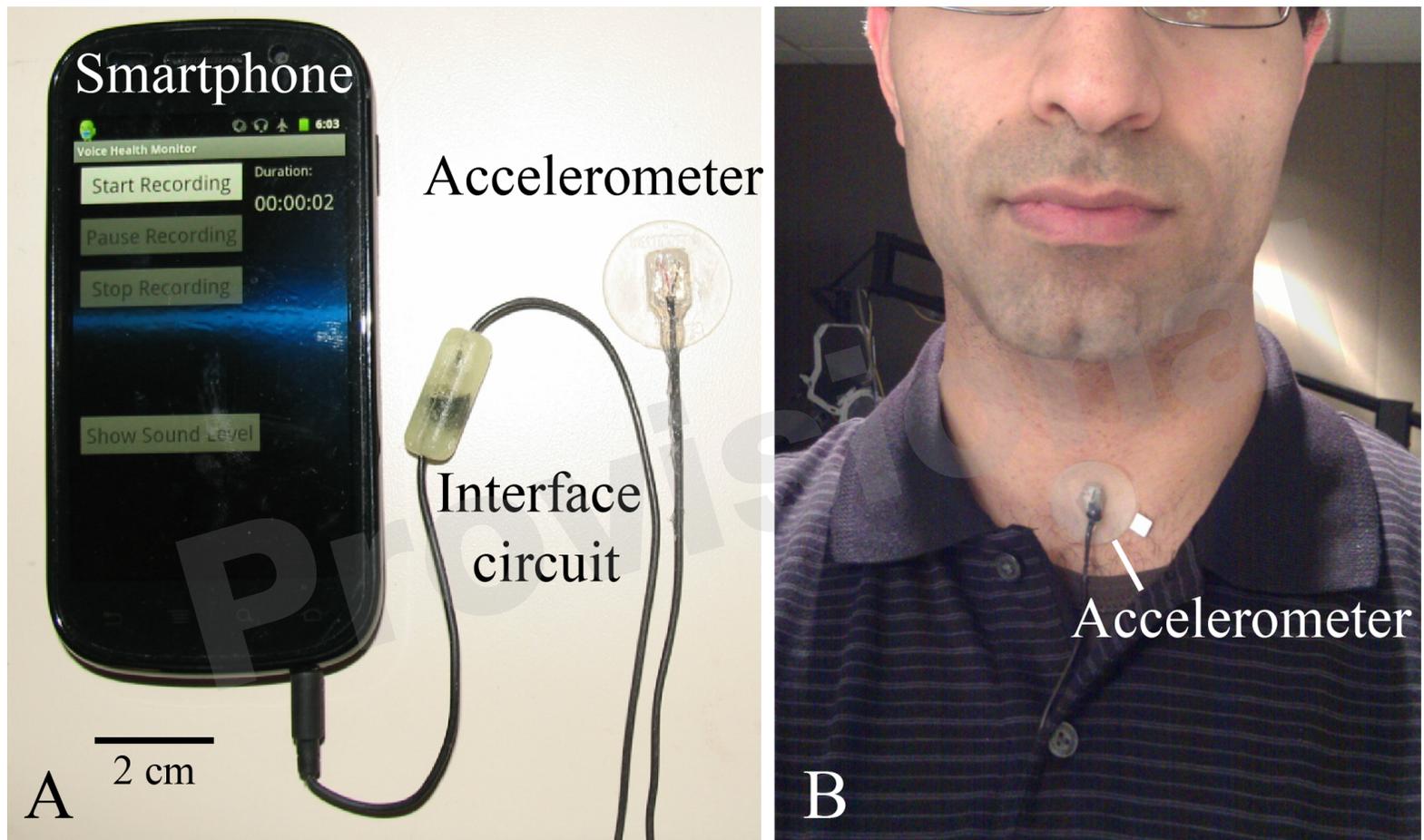


Figure 4.TIF

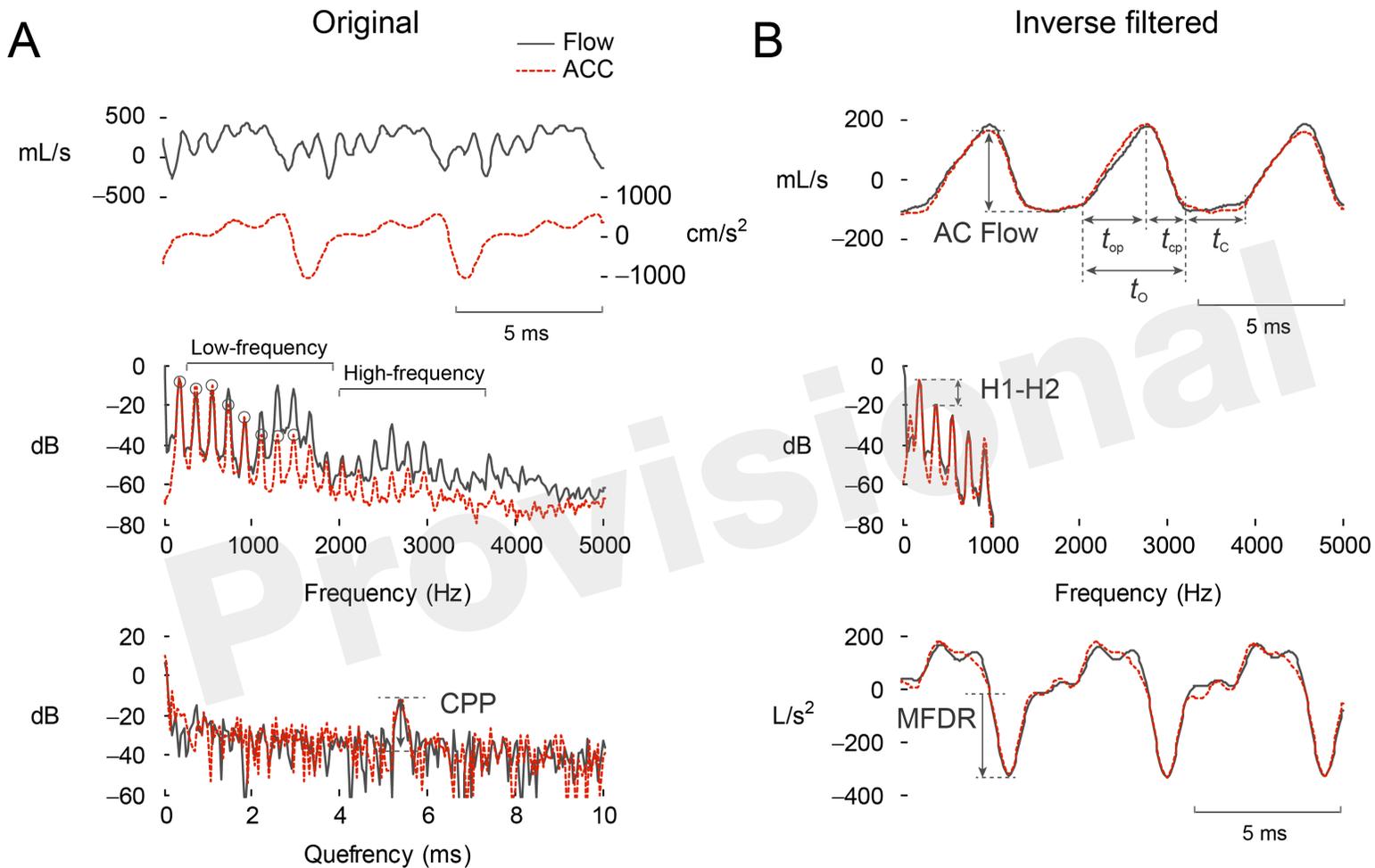


Figure 5.TIF

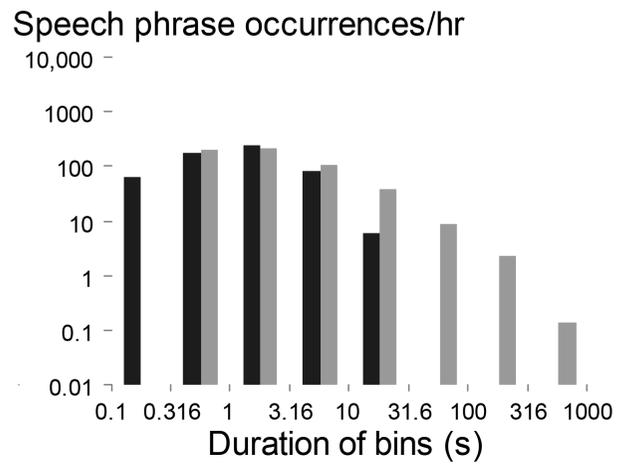
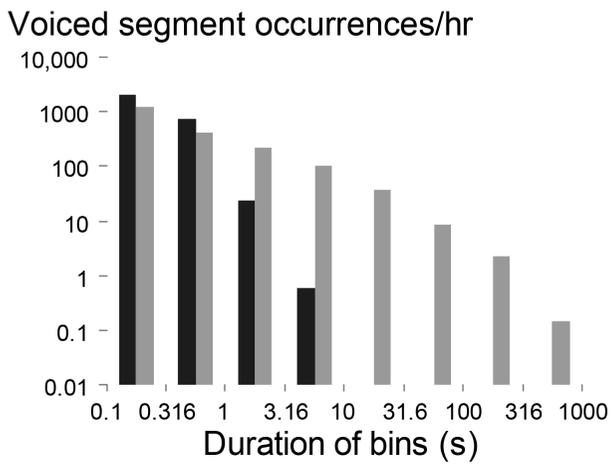
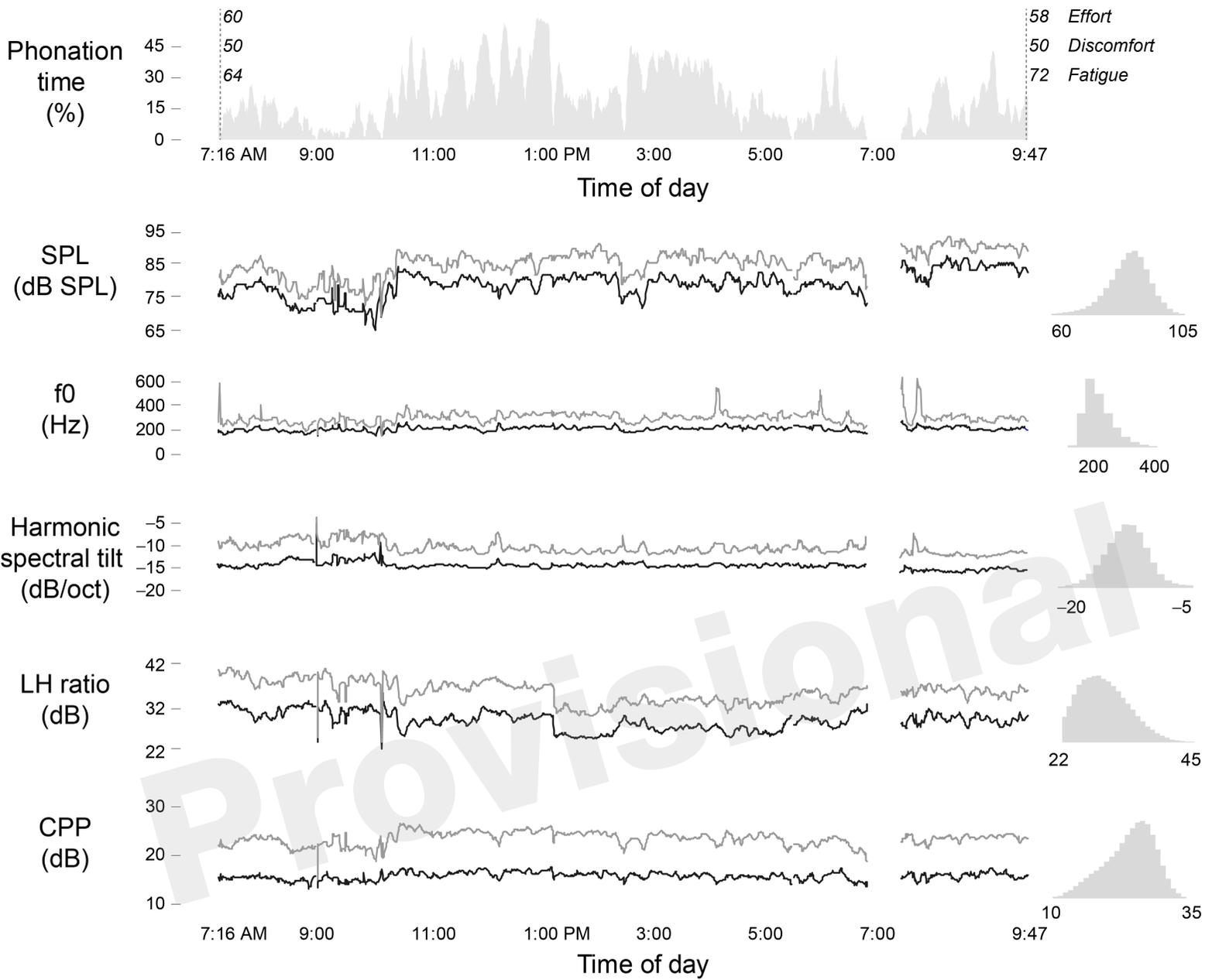


Figure 6.TIF

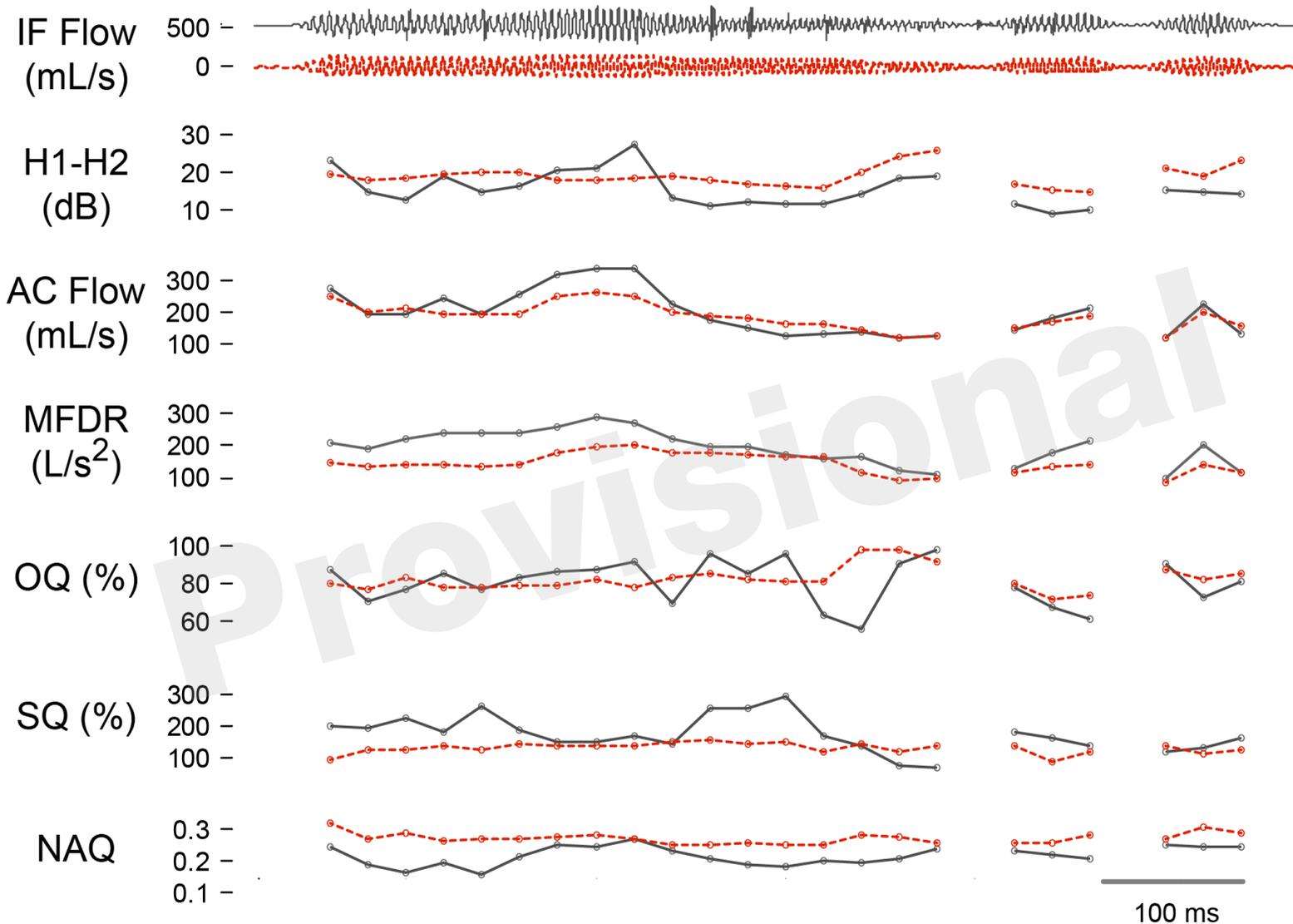
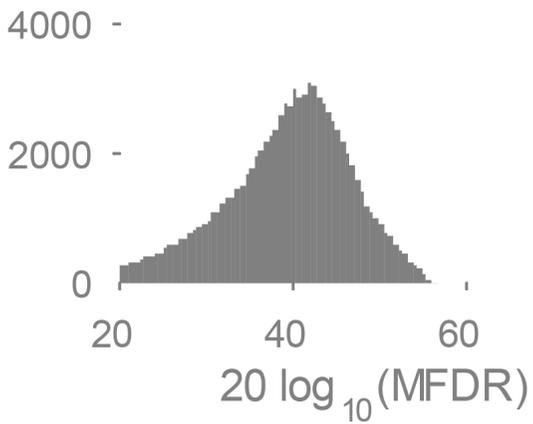
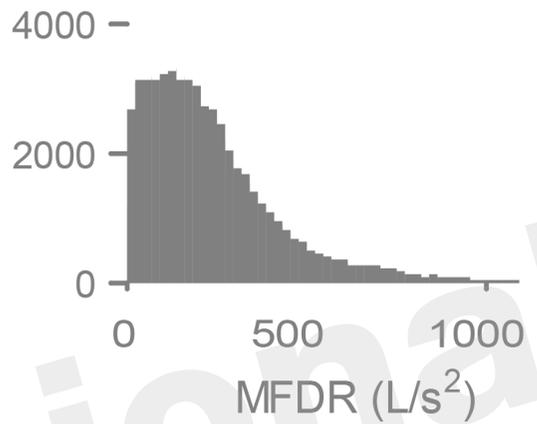


Figure 7.TIF

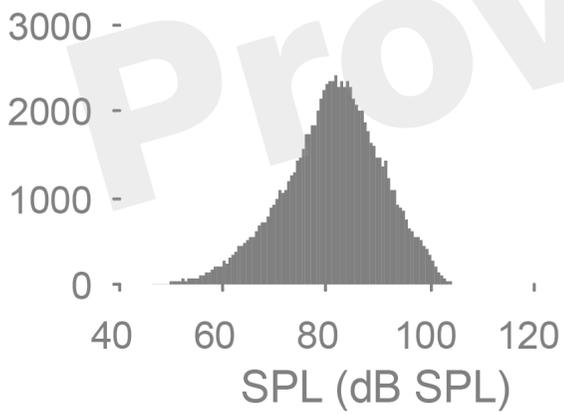
voiced frames



voiced frames



voiced frames



SPL (dB SPL)

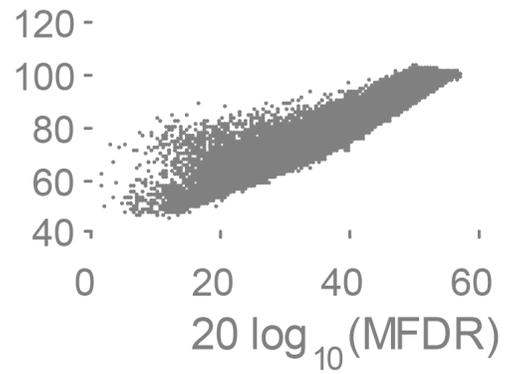


Figure 8.TIF

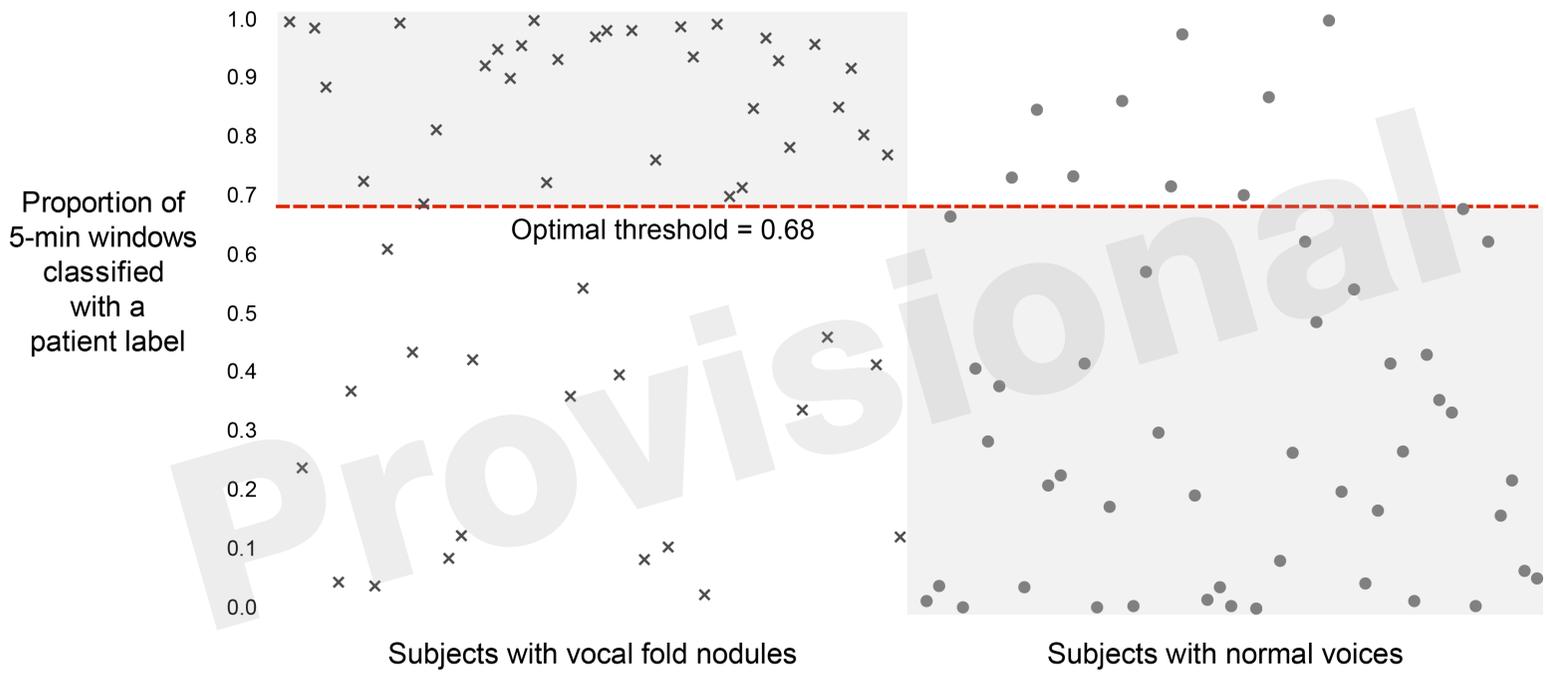


Figure 9.TIF

