



ELSEVIER

Available online at www.sciencedirect.com

ScienceDirect

Computer Speech & Language xxx (2017) xxx-xxx

www.elsevier.com/locate/csl

Characterisation of voice quality of Parkinson's disease using differential phonological posterior features[☆]

Milos Cernak^{*,a}, Juan Rafael Orozco-Aroyave^{b,e}, Frank Rudzicz^c, Heidi Christensen^d,
Juan Camilo Vásquez^{b,e}, Elmar Nöth^e

^a *Idiap Research Institute, Martigny, Switzerland*

^b *Universidad de Antioquia Medellín, Colombia*

^c *University of Toronto, Canada*

^d *University of Sheffield, UK*

^e *Friedrich-Alexander-Universität Erlangen-Nürnberg, Germany*

Received 13 March 2017; received in revised form 12 May 2017; accepted 12 June 2017

Abstract

Change in voice quality (VQ) is one of the first precursors of Parkinson's disease (PD). Specifically, impacted phonation and articulation causes the patient to have a breathy, husky-semiwhisper and hoarse voice. A goal of this paper is to characterise a VQ spectrum – the composition of non-modal phonations – of voice in PD. The paper relates non-modal healthy phonations: breathy, creaky, tense, falsetto and harsh, with disordered phonation in PD. First, statistics are learned to differentiate the modal and non-modal phonations. Statistics are computed using phonological posteriors, the probabilities of phonological features inferred from the speech signal using a deep learning approach. Second, statistics of disordered speech are learned from PD speech data comprising 50 patients and 50 healthy controls. Third, Euclidean distance is used to calculate similarity of non-modal and disordered statistics, and the inverse of the distances is used to obtain the composition of non-modal phonation in PD. Thus, pathological voice quality is characterised using healthy non-modal voice quality “base/eigenspace”. The obtained results are interpreted as the voice of an average patient with PD and can be characterised by the voice quality spectrum composed of 30% breathy voice, 23% creaky voice, 20% tense voice, 15% falsetto voice and 12% harsh voice. In addition, the proposed features were applied for prediction of the dysarthria level according to the Frenchay assessment score related to the larynx, and significant improvement is obtained for reading speech task. The proposed characterisation of VQ might also be applied to other kinds of pathological speech. © 2017 Published by Elsevier Ltd.

Keywords: Phonological features; Non-modal phonation; Parkinson's disease

1. Introduction

Speech of hypokinetic dysarthria in Parkinson's disease (PD) is characterised by hypokinesia (rigid, less motion describing decreased range and frequency of movement) of the vocal folds and articulators. Besides of impacted

[☆] This paper has been recommended for acceptance by Roger Moore.

* Corresponding author.

E-mail address: milos.cernak@idiap.ch (M. Cernak).

4 prosody and articulation, phonation is impacted by incomplete vocal fold adduction. Clinicians, otolaryngologists
5 and speech-language pathologists, consider hoarseness – a rough quality of voice – as a basic symptom of a voice
6 disorder in PD. When hoarse, the voice may sound breathy, raspy, or strained, and if this abnormal/pathological
7 voice quality is accompanied with relatively constant loudness and pitch deviations, it is diagnosed as Parkinsonian
8 dysphonia (Aronson and Bless, 2011).

9 Healthy subjects may also produce speech sounds of different voice quality based on different modes of vibration
10 of the vocal folds. Laver (1980) defines the term of voice quality in a broad sense as the characteristic auditory col-
11 ouring of an individual speaker’s voice, and not just in a narrow sense coming from the laryngeal activity. The neu-
12 tral mode phonation, often used in *modal voice*, is one against which the other modes can be contrastively described,
13 also called non-modal phonations. Ladefoged and Johnson (2014) describe four basic states of the glottis (which is
14 defined as the space between the vocal folds). The position of the vocal folds is adjusted by the arytenoid cartilages
15 placed toward the back. In (i) a voiced sound, the vocal folds are close together (adducted) and vibrating, whereas in
16 (ii) a voiceless sound, they are pulled apart (abducted). If there is considerable airflow, the abducted vocal folds will
17 be set vibrating – flapping in the airstream – producing what is called (iii) *breathy voice*, or murmur. Alternatively,
18 breathy voice is produced with the vocal folds apart only between the arytenoid cartilages in the lower (posterior)
19 part. If the arytenoid cartilages are tightly together, so that the vocal folds can vibrate only at the anterior end, (iv)
20 *creaky voice* is produced. Creaky-voiced sounds may also be called laryngealised. Besides these basic non-modal
21 phonation, Laver (1980) defines *tense*, *harsh* and *falsetto* phonations. Such voice qualities impact the production of
22 the speech sounds, and we hypothesise that these changes might be captured by changes in phonological features.

23 The goal of this paper is to present a study on the production of speech sounds with healthy non-modal phonation,
24 and project its non-modal statistics to analyse disordered production of speech sounds with pathological phonation.
25 This approach might help to alleviate a problem of missing data in research of pathological speech. Voice quality of
26 the speech sounds can be characterised by phonological features (Cernak et al., 2017b), and the current work pro-
27 poses to use differential phonological posterior features (between modal and non-modal, and between healthy and
28 disordered phonations) for characterisation of both healthy non-modal and parkinsonian phonations. Comparing to
29 the work of Cernak et al. (2017b), the novel aspects of this paper is in using pathological speech and characterisation
30 of pathological voice quality using healthy non-modal voice quality “base/eigenspace”. An Euclidean distance
31 between the non-modal and disordered phonation characterisations quantifies the composition of non-modal voice
32 qualities in PD. This characterisation of non-modal phonation in PD is novel, and shows objective quantification of
33 voice quality using phonological features not investigated in previous approaches.

34 For studying speech with non-modal phonation, the read-VQ database of Kane (2012) is used, the recording of
35 which was inspired by prototype voice quality examples produced by Laver (1980). Laver’s recordings are consid-
36 ered as recordings of non-modal phonation with excellent quality, however only one utterance per phonation type is
37 available, and thus they are speaker-specific. The read-VQ database contains recordings from four speakers. The
38 database covers five different non-modal phonations: falsetto, creaky, harshness, tense and breathiness. For studying
39 speech with pathological phonation, the Colombian-Spanish database (Orozco-Arroyave et al., 2014) is used, which
40 contains speech recordings of 50 patients with PD and 50 healthy controls (HC).

41 The structure of the paper is as follows: Section 2.1 gives an overview of the non-modal (healthy) and pathologi-
42 cal (Parkinsonian) phonation types considered in this work. Section 3 introduces differential phonological posterior
43 (DPP) features used in further characterisation of VQ. Section 4 describes experimental setup and evaluation data-
44 bases, and Section 5 presents results and their validation. Finally, Section 6 concludes the paper.

45 2. Voice quality of Parkinson’s disease

46 2.1. Non-modal (healthy) phonation

47 Different modes of vibration of the vocal folds contribute significantly to VQ. The modal (periodic) phonation
48 can be contrastively described against the other modes, also called non-modal (aperiodic) phonations.

49 Recent work on non-modal phonation focuses on detection (Drugman et al., 2014), analysis (Malyska, 2008;
50 Malyska et al., 2011) and synthesis (Bangayan et al., 1997) of speech with non-modal phonation. Modern computa-
51 tional paralinguistics tries to 1) get rid of non-modal phonation, or 2) model it, for example, for classification

52 purposes (Schuller and Batliner, 2013). Non-modal phonation is also studied in sociolinguistics. For example, creaky
53 and falsetto phonations are used more commonly by women (Anderson et al., 2014; Podesva, 2007).

54 Breathy and creaky voices belong to the most studied non-modal phonation types. In breathy phonation, the vibra-
55 tion of the vocal folds is accompanied by aspiration noise, which causes a higher first formant bandwidth and a miss-
56 ing third formant (Klatt and Klatt, 1990) due to steeper spectral tilt (Hanson, 1997). In creaky phonation (also
57 referred to as vocal fry, laryngealisation), secondary vibrations occur with lower fundamental frequencies.

58 Tense voice is produced with higher degree of overall muscular tension involved in the whole vocal tract. The
59 higher tension of the vocal folds does not result in irregularities that are seen in harsh voice. It is characterised by
60 richer harmonics in higher frequencies due to a less steep spectral tilt. Harsh voice is a result of very high muscular
61 tension at the laryngeal level. Pitch is irregular and low, and the speech spectrum contains more noise.

62 Falsetto voice is the most different with respect to modal voice (Laver, 1980). The voice is produced with thin
63 vocal folds, that results in a higher pitch voice with a steeper spectral slope.

64 2.2. Pathological (Parkinsonian) phonation

65 Auditory-perceptual evaluation of disordered VQ is the most commonly used clinical assessment method, and is
66 considered by clinicians as the “gold standard” for documenting voice impairment severity (Kreiman et al., 1993).
67 Describing a particular voice as breathy and rough, for example, is likely to be more easily interpreted by a wide
68 range of people than a description that specifies the noise-to-harmonic ratio associated with that voice (Oates, 2009).
69 Moreover auditory-perceptual evaluation is cheap and practical. Perceptual analysis is used with the human auditory
70 perceptual system, often in combination with an external rating system, such as the GRBAS protocol (Hirano, 1981)
71 developed by the Japanese Society of Logopedics and Phoniatics. The GRBAS protocol contains 4-point scales for
72 grade (overall severity), roughness, breathiness, asthenia (lack of vocal power), and strain.

73 On the other hand, the perceptual evaluation has been characterised by questionable validity and poor reliability,
74 adding further analysis error via measurement and scaling issues (Aronson and Bless, 2011), and missing consensus
75 on stimulus categories (Barsties and De Bodt, 2015). At present, the Consensus Auditory Perceptual Evaluation of
76 Voice (CAPE-V)¹, containing six primary perceptual parameters (overall severity, roughness, breathiness, strain,
77 pitch, and loudness), is undergoing field testing, and experimental data on its validity and reliability are forthcom-
78 ing.

79 Acoustic analysis is widely employed in clinical and research settings, and focuses on analysis of parkinsonian
80 speech that provides objective measures of vocal function, such as fundamental frequency, signal amplitude, jitter,
81 shimmer, noise-to-harmonic ratios, voice onset time and glottal leakage, and last but not least the spectral features
82 such as spectral tilt (Holmes et al., 2000; Little et al., 2009; Rusz et al., 2011; Bauer et al., 2011). Parkinsonian
83 speech is characterised by higher jitter (more roughness), higher shimmer, decreased pitch range, shorter maximum
84 phonation time and slower diadochokinetic (articulation) rate (Darley et al., 1969). However, acoustic measures can-
85 not be applied to more severe disordered voices due to their nonlinear and non-Gaussian random properties (Little
86 et al., 2007), that limits their clinical usefulness.

87 There is a considerable amount of literature on objective perceptual evaluation based on acoustic and aerody-
88 namic speech production characteristics. For example, Wuyts et al. (2000) propose a Dysphonia Severity Index, con-
89 structed from highest frequency, lowest intensity, maximum phonation time and jitter. Bhuta et al. (2004) and Maryn
90 et al. (2009) provide detailed studies of correlation of acoustic measurements with perceived voice quality. Recent
91 methods include in objective perceptual evaluation also spectral/cepstral features, such as spectrum slope and tilt
92 (Maryn et al., 2010), and cepstral peak prominence (Awan et al., 2009).

93 3. Differential phonological posteriors

94 The probabilities of phonological features inferred from the speech signal – phonological posteriors – can be
95 reliably estimated using a deep learning approach (Cernak et al., 2015). This extraction processes is further called
96 phonological analysis. In this work, the Sound Patterns of English (SPE) feature set of Chomsky and Halle (1968) is

¹ <http://www.asha.org/uploadedFiles/members/divs/D3CAPEVprocedures.pdf>.

used. Motivation to this older phonological system was that (i) it takes the articulatory production mechanism as the underlying principle of phoneme organisation (and thus allows easier interpretation of obtained results), and (ii) SPE assumes that the flat, unstructured binary feature specifications are language independent and characterise the set of possible phonemes in languages of the world (and thus is more suitable for studies with more languages like described in this paper). The mapping from phonemes to SPE phonological classes is taken from Cernak et al. (2017a). The distribution of the phonological labels is non-uniform, driven by mapping different numbers of phonemes to the phonological classes.

Phonological analysis starts with a short-term analysis of speech, which consists of converting the speech signal into a sequence of acoustic feature vectors $X = \{\vec{x}_1, \dots, \vec{x}_n, \dots, \vec{x}_N\}$. Each \vec{x}_n is also known as an acoustic frame or just frame, and can be composed by the conventional Mel frequency cepstral coefficients (MFCC). N is the number of frames and frames are equally spaced in time.

Then, K phonological probabilities z_n^k are estimated for each frame. Each probability is computed independently by using a binary classifier based on deep neural network (DNN) and trained with one class versus the rest. Finally, the acoustic feature observation sequence X into a sequence of phonological vectors $Z = \{\vec{z}_1, \dots, \vec{z}_n, \dots, \vec{z}_N\}$. Each vector $\vec{z}_n = [z_n^1, \dots, z_n^k, \dots, z_n^K]^T$ consists of phonological class posterior probabilities $z_n^k = p(c_k|x_n)$ of K phonological features (classes) c_k . The a posteriori estimates $p(c_k|x_n)$ are $0 \leq p(c_k|x_n) \leq 1$, $\forall k$, and $\max \sum_{k=1}^K p(c_k|x_n) = K$.

The matrix of posteriors Z consists of N rows, indexed by the processed speech frames, and K columns. The following analysis is done on non-silence speech frames of the evaluation data:

$$\mu_k = \frac{1}{N_s} \sum_{n=1}^{N_s} p(c_k|x_n), \forall n \Leftrightarrow p(c_{\text{SIL}}|x_n) < 0.5, \quad (1)$$

where c_{SIL} is a posterior probability of silence class being observed, and N_s is the number of non-silence frames. The probability of c_{SIL} is computed as for the other phonological classes (i.e., the silence versus the rest) but it is only taken into account when computing each μ_k . The statistics μ_k is calculated for different “contrastive” data groups, such as data with modal vs. data with non-modal phonations, and data from healthy speakers vs. data from pathological speakers.

Differential phonological posterior (DPP) features are obtained by mean normalisation of contrastive data:

$$\begin{aligned} \Delta\mu_{kl}^{NM} &= \mu_{kl}^{\text{non-modal}} - \mu_{kl}^{\text{modal}}, \\ \Delta\mu_k^P &= \mu_k^{\text{PD}} - \mu_k^{\text{HC}}. \end{aligned} \quad (2)$$

Thus, the non-modal mean posteriors are normalised by modal means that yields the normalised statistics $\Delta\mu_l^{NM} = [\Delta\mu_{1l}^{NM}, \dots, \Delta\mu_{kl}^{NM}, \dots, \Delta\mu_{Kl}^{NM}]^T$ for $l \in L$ non-modal phonations, and PD posteriors are normalised by means from healthy speakers that yields pathological (Parkinsonian) statistics $\Delta\mu^P = [\Delta\mu_1^P, \dots, \Delta\mu_k^{NM}, \dots, \Delta\mu_K^{NM}]^T$.

Finally, similarity of non-modal phonation and pathological speech is calculated as the Euclidean distance:

$$q_l = \|\Delta\mu^P - \Delta\mu_l^{NM}\|, \quad (3)$$

for $l \in L$ non-modal phonations, where q_l represents a similarity of the l -th non-modal phonation with VQ in PD. The Euclidean distance was already successfully used as a similarity measure between VQ characterisations in forensic speaker comparison (San Segundo et al., 2017).

The normalisation of the mean posteriors by the posterior features from the modal or healthy speakers is conceptually similar to likelihood ratio test in speaker recognition (Hansen and Hasan, 2015), where likelihoods from the speaker model are subtracted by likelihoods obtained from the background/world model. In the DPP features, the background models represent the modal phonation and healthy speakers.

4. Experimental setup

4.1. Training data

The phonological analyser is trained on the Wall Street Journal (WSJ0 and WSJ1) continuous speech recognition corpora Paul, Baker, (1992). This training database consists primarily of read speech using a close-talking

140 Sennheiser HMD414. The *si_tr_s_284* set of 37,514 utterances was used, split into 90% training and 10% cross-
141 validation sets. Titze (1995) recommends the WSJ database to be used in acoustic analysis research of pathological
142 speech. In addition, Cernak et al. (2015) introduced a deep learning approach using WSJ data to achieve high classi-
143 fication accuracy of phonological features.

144 4.2. Evaluation data

145 Prototype voice quality examples produced by Laver (1980) and the read-VQ database of Kane (2012) were used
146 to obtain characterisation of modal and non-modal phonation. Audio of the read-VQ database was recorded at
147 44.1 kHz using high quality recording equipment: a B&K 4191 free-field microphone and a B&K 7749 pre-
148 amplifier. The microphone was placed at a distance of approximately 30 cm from the speaker and participants were
149 asked to keep this distance as constant as possible throughout the recording session. Recordings were subsequently
150 downsampled to 16 kHz.

151 The read-VQ database contains 4 speakers (2 males and 2 females) who were asked to read 17 sentences in six
152 different phonation types: modal, breathy, tense, harsh, creaky, and falsetto. Participants were given prototype voice
153 quality examples, produced by John Laver and John Kane, and were asked to practise producing them before coming
154 to the recording session. For the recordings, participants were asked to produce the strong versions of each phonation
155 type and to maintain it throughout the utterance. During the recording session, participants were asked to repeat the
156 sentence if it was deemed necessary. The sentences were chosen from the phonetically balanced sentences in the
157 TIMIT corpus (Garofolo et al., 1993), four of which contained all-voiced sounds. 451 sentences were chosen to
158 obtain a wide phonetic coverage, as it is likely that it can be very difficult for speakers to maintain a constant type of
159 phonation over a long utterance. The recordings with modal phonation were 2.2 min long, and the remaining record-
160 ings with non-modal phonation were 2.0 min long each. The read-VQ data was used for estimation of non-modal
161 DPP features $\Delta\mu_i^{NM}$.

162 Speech recordings from the HC and PD patients were obtained from the database provided by Orozco-Arroyave
163 et al. (2014). This database contains speech recordings of 50 patients with PD and 50 HCs sampled at 44.1 kHz with
164 16 resolution-bits. The recordings were captured in noise controlled conditions, in a sound proof booth. All of the
165 speakers are balanced by gender and age. All of the patients were diagnosed and labeled by neurologist experts. The
166 speech samples were recorded with the patients in the ON-state, i.e., no more than 3 h after the morning medication.
167 None of the people in the HC group has a history of symptoms related to PD or any other kind of neurological disorder.
168 The HC and PD data was used for estimation of parkinsonian DPP features $\Delta\mu_i^P$. It is worth to note
169 that the training data of phonological analyser contains English recordings, whereas the HC and PD data contain
170 Columbian-Spanish recordings. It is assumed that phonological features are language independent, and in addition,
171 Cernak et al. (2016) showed effective cross-language usage of phonological posteriors, for English training data and
172 French evaluation data, and vice versa.

173 PD data contains several different speech tasks comprising of isolated and running speech: 24 isolated words, the
174 ‘pataka’ speech task consisting of repeating /pataka.petaka.pakata/ speech production, 10 sentences, one read text
175 with 36 words, and a monologue with an average duration of 44.86s. All data was used for experiments described in
176 Section 5.

177 4.3. Training of phonological analysis

178 The open-source phonological vocoding platform developed by Cernak and Garner (2016) was used to perform
179 phonological analysis and synthesis. Briefly, the platform is based on cascaded speech analysis and synthesis that
180 works internally with the phonological speech representation. In the phonological analysis part, phonological poste-
181 riors are estimated directly from the speech signal by DNNs. Binary (Yu et al., 2012) or multi-valued classification
182 (Stouten and Martens, 2006; Rasipuram and Magimai.-Doss, 2011) may be used. In the latter case, the phonological
183 classes are grouped together based on place or manner of articulation. The binary classification approach is used in
184 this work, and thus each DNN determines the probability of a particular phonological class.

185 To train the DNNs for phonological analysis, a phoneme-based automatic speech recognition system is first
186 trained using Mel frequency cepstral coefficients (MFCC) as acoustic features. The phoneme set comprises 40 pho-
187 nemes (including silence) defined by the CMU pronunciation dictionary. The three-state, cross-word triphone

188 models were trained with the HMM-based speech synthesis system (HTS) variant (Zen et al., 2007) of the Hidden
189 Markov Model Toolkit (HTK) on the 90% subset of the WSJ data. The remaining 10% subset was used for cross-
190 validation. The triphone models are tied with decision tree state clustering based on the minimum description length
191 (MDL) criterion (Shinoda and Watanabe, 1997). The MDL criterion allows an unsupervised determination of
192 the number of states. The trained model had 12685 tied states, and each is modeled with a Gaussian mixture model
193 consisting of 16 Gaussians.

194 The acoustic models were used to get phonetic alignment. Each phoneme was mapped to the 13 SPE phonological
195 classes or the one silence class, and thus 14 DNNs were trained as phonological/silence analysers using the frame
196 alignment with a particular phonological/silence label scheme that took two binary values: the phonological/silence
197 class exists for the aligned phoneme or not. In other words, the two DNN outputs correspond to the target class vs.
198 the rest.

199 Some classes might seem to have unbalanced training data; for example, the two labels for the nasal class are
200 associated with the speech samples from just 3 nasal phonemes /m/, /n/, and /ŋ/, and with the remaining 36 (non-
201 nasal) phonemes. However, this split is necessary to appropriately train a discriminative classifier, as all the remain-
202 ing phonemes convey information about all different phonological classes. Each DNN was trained on the whole
203 training set. Several DNN sizes were tested, from 3 to 6 hidden layers with 500–2000 neurons. Finally, the selected
204 size of the DNNs was $351 \times 1024 \times 1024 \times 1024 \times 2$ neurons, is a balance between the training time and the per-
205 formance. Sigmoid activation functions were used in the hidden layers. The input feature vectors consisted of Energy
206 plus 12 MFCC (13 parameters) with the first and second time derivatives. The temporal context from 7 to 11 succes-
207 sive frames was tested with no particular performance increase, so the temporal context of 9 frames was used for the
208 training.

209 The parameters were initialised using deep belief network pre-training following the single-step contrastive diver-
210 gence (CD-1) procedure of Hinton et al. (2006). The DNNs with the softmax output function were then trained using
211 a mini-batch based stochastic gradient descent algorithm with the cross-entropy cost function of the KALDI toolkit
212 (Povey et al., 2011).

213 5. Results

214 5.1. Analysis of non-modal phonation

215 Fig. 1a shows the analysis of the read-VQ evaluation data. Table 2 shows the results of further statistical analysis
216 performed by using the two-sample *t*-test without assuming equal variance, that was carried out to study the differen-
217 ces between speech with modal and non-modal phonations. The significance of the test also allows for the determina-
218 tion of invariant phonological features, listed in Table 2.

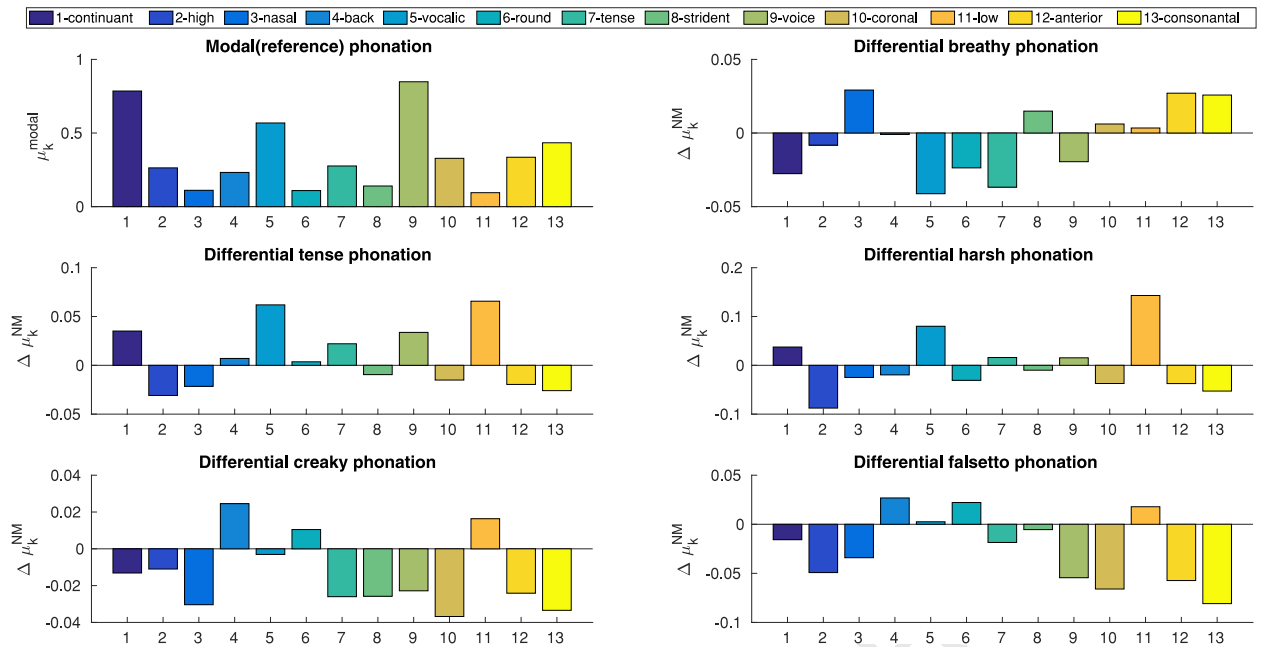
219 According to Table 2, the [Strident] phonological feature is more invariant – “resistant” – to non-modal phona-
220 tion, whereas the [Nasal], [Voice], [Anterior] and [Consonantal] features are heavily impacted (they are not invariant
221 for any phonation type). The [Strident] feature is significantly different only in creaky phonation, which indicates its
222 usefulness, for example, in creaky voice detection. On the contrary, the invariant [Tense] feature might indicate
223 harsh, and the invariant [Low] feature may indicate breathy phonation.

224 The number of invariant features also indicates the impact of non-modal phonation on phonological features.
225 While breathy, creaky and tense phonations keep 4 invariant features, harsh and falsetto phonation keep only 2
226 invariant features.

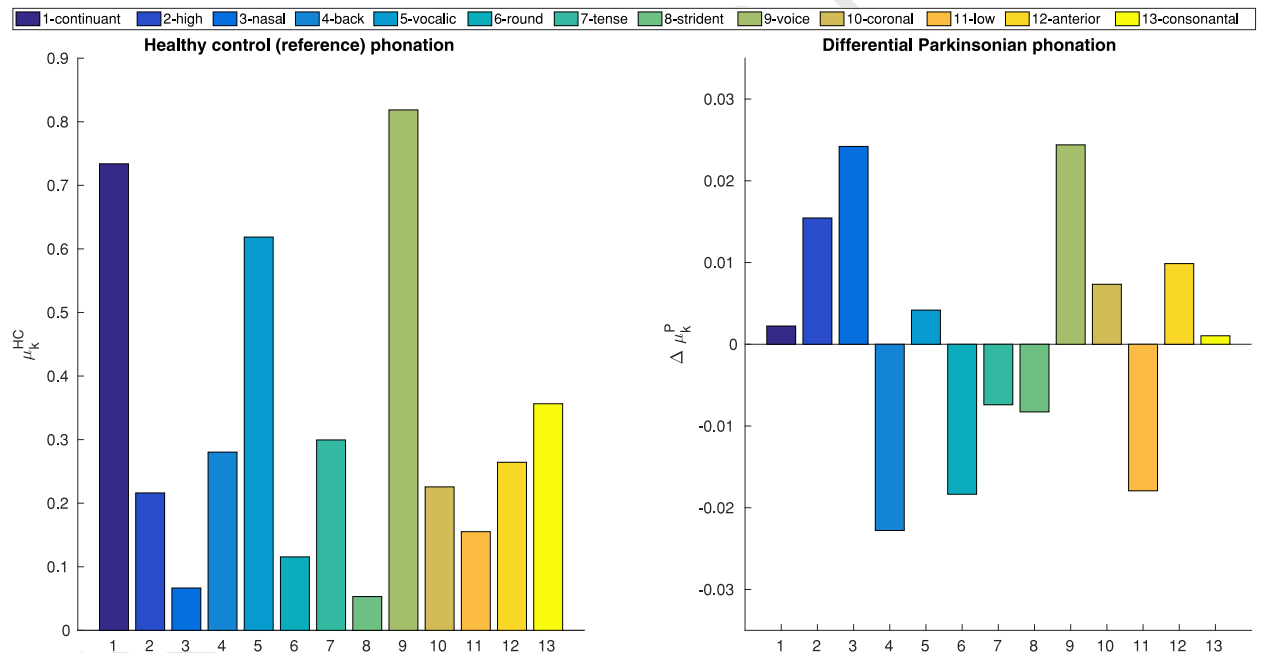
227 5.2. Analysis of Parkinsonian speech

228 Fig. 1b shows the analysis of the HC and PD non-silence speech data: 10 ms framed 805,511 phonological poste-
229 rior vectors of the HC group, and 10 ms framed 784,128 vectors of the PD group.

230 Statistical analysis using the two-sample *t*-test, without assuming equal variances, of the differences between HC
231 and PD speech, resulted into the only invariant [Consonantal] feature with $p = 0.1029$, which is in contradiction
232 with non-modal analysis above, where the [Consonantal] feature was significantly different between modal and non-
233 modal phonations. PD speech exhibited higher values of the [Nasal], [Voice] and [High] features, and lower values
234 of the [Back], [Low], and [Round] features. Validation of these findings is discussed further in Section 5.3.



(a) Analysis of the read-VQ recordings visualizing mean difference of non-modal and modal DPP features.



(b) Analysis of the HC and PD recordings visualizing mean difference of PD and HC DPP features.

Fig. 1. Mean modal/HC SPE posteriors μ_k (top-left figures) and differentials $\Delta\mu_k$ of non-modal/PD phonations with respect to the modal/HC voice.

Table 1

Statistical significance (p -values) of difference between μ_{kl}^{modal} and $\mu_{kl}^{\text{non-modal}}$ for $k \in K = 13$ SPE features and $l \in L = 5$ non-modal phonations. For of the level of significance $\alpha = 0.001$, the bold items represent the invariance of a particular pair of the SPE feature and non-modal phonation, i.e., the SPE features unaffected by non-modal phonations, where statistical significance of differences is $p > \alpha$. The other items shown by ‘-’ represent the SPE features affected by non-modal phonation, with significance $p < \alpha$.

SPE/Phonation	Breathy	Creaky	Tense	Harsh	Falsetto
Continuant	-	0.0042	-	-	-
High	0.0958	0.0267	-	-	-
Nasal	-	-	-	-	-
Back	0.8261	-	0.1308	-	-
Vocalic	-	0.5948	-	-	0.6657
Round	-	0.0031	0.3114	-	-
Tense	-	-	-	0.0012	-
Strident	-	-	0.0251	0.0208	0.2198
Voice	-	-	-	-	-
Coronal	0.2413	-	0.0041	-	-
Low	0.2902	-	-	-	-
Anterior	-	-	-	-	-
Consonantal	-	-	-	-	-

Table 2

The impact of non-modal phonation on phonological features, measured as a positive (+) or negative (-) difference between the mean phonological posteriors of speech with modal phonation, and the mean phonological posteriors with non-modal phonation. The three features with the greatest differences are listed. Invariance is concluded based on statistics in Table 1.

Phonation	Invariant features	Most different features
Breathy	High, Back, Coronal, Low	-Vocalic, -Tense, +Nasal
Creaky	Continuant, High, Vocalic, Round	-Coronal, -Consonantal, -Nasal
Tense	Back, Round, Strident, Coronal	+Low, +Vocalic, +Continuant
Harsh	Strident, Tense	+Low, -High, +Vocalic
Falsetto	Strident, Vocalic	-Consonantal, -Coronal, -Anterior

235 Having the statistics of mean DPP features, we calculated Euclidean distances using Eq. (3) between parkinsonian
 236 DPP $\Delta\mu^P$ (visualised at right of Fig. 1b), and L non-modal DPP $\Delta\mu_i^{NM}$. Table 3 lists obtained Euclidean distances.
 237 As said in Section 1, non-modal phonation modes are contrastive against modal phonation modes, in other words,
 238 they are di-similar. The q_l quantities represent the similarity measures, so to be used for characterisation of Parkinsonian
 239 non-modal phonation, they are turned into di-similarity measures by calculating their inverse, $1/q_l$. Finally, we
 240 assume that each of the non-modal phonation partially (relatively) contributes to the perceived overall non-modal
 241 phonation.

242 Fig. 2 shows composition of voice quality in parkinsonian speech. It might be interpreted as: a voice of an average
 243 patient with Parkinson’s disease contains “a voice quality spectrum” composed of 30% breathy voice, 23% creaky

Table 3

Euclidean distances q_l between non-modal and Parkinsonian DPP features. As Euclidean distance is a similarity measure, whereby smaller is more similar, we calculate an inverse of the Euclidean distance to plot composition of non-modal voice quality in Parkinsonian speech in Fig. 2.

Voice quality	q_l	$1/q_l$
Breathy	0.0935	10.69
Creaky	0.1240	8.06
Tense	0.1417	7.06
Falsetto	0.1904	5.25
Harsh	0.2321	4.31

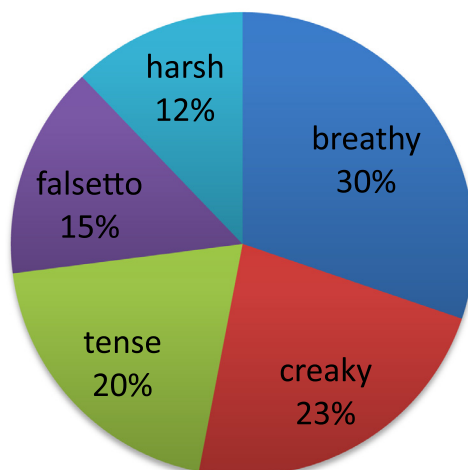


Fig. 2. Composition of voice quality in Parkinson's speech.

244 voice, 20% tense voice, 15% falsetto voice and 12% harsh voice, where about 75% of overall voice quality on aver-
 245 age is composed of breathy, creaky and tense phonations.

246 5.3. Validity

247 Oates (2009) describes basic pathological phonations as a breathy voice that arises from incomplete glottal clo-
 248 sure and/or the presence of a posterior glottal chink, a rough voice that arises from irregular vocal fold vibration pat-
 249 terns, and a strained or pressed voice that is due to excess laryngeal muscle tension. Barsties and De Bodt (2015)
 250 review three ratings schemes that are the most frequently reported and accepted: (i) the GRBAS scale that includes
 251 *R* for roughness, *B* for breathiness and *S* for strain; (ii) the CAPE-V that includes in the standard analysis the same
 252 parameters as the GRBAS; and (iii) the RBH scale that focus on only three dimensions: roughness, breathiness, and
 253 hoarseness.

254 We objectively estimated that the majority of the VQ spectrum of PD is composed of 30% breathy voice,
 255 23% creaky voice, 20% tense voice; all the three most-important VQs expected/evaluated by perceptual
 256 assessment of hypokinetic dysarthria in PD. Breathly phonation causes breathiness, creaky phonation contrib-
 257 utes significantly to roughness, and tense phonation results into vocal strain (known also as muscle tension
 258 dysphonia).

259 Severity of dysarthria in PD is also rated by the the Frenchay Dysarthria Assessment (FDA-2) score (Enderby,
 260 1983; Enderby and Palmer, 2008). The assessment includes 28 relevant perceptual dimensions of speech, namely
 261 related to the following dimensions:

- 262 • Respiration: noting running out of breath when speaking, and breathy voice.
- 263 • Laryngeal: noting weather the patient has clear phonation with the vocal folds, without huskiness.
- 264 • Tongue: noting accurate tongue movements (positions) with correct articulation.
- 265 • Palate: noting nasal resonance in spontaneous conversation, without hypernasality or nasal emission.
- 266 • Lips: observing the movements of lips in conversation, noting correct shape of lips.

267 While the first dimension is similar to the perceptual assessment of the three rating schemes described above, fur-
 268 ther dimensions are more related to articulation. According to Fig. 1b, PD speech data exhibits:

- 269 1. Greater values of the [Voice] and [Nasal] phonological features that might be related to the Laryngeal and Palate
 270 dimensions. It can be interpreted as more analysed speech frames having higher values of these phonological

- 10 *M. Cernak et al. / Computer Speech & Language xxx (2017) xxx-xxx*
- 271 features, as compared to HC speech data. Thus, patients produced more nasal or voiced sounds compared to
 272 HCs.
- 273 2. Lower values of the [Round] that might be related to the Lips dimension (i.e., patients produced less rounded
 274 sounds).
- 275 3. Lower [Back] and [Low], and greater [High] vales that might be related to Tongue dimension (i.e., the patients
 276 articulated more central speech sounds, that might indicate weaker articulation of PD patients).

277 5.3.1. Prediction of laryngeal FDA scores

278 To validate usefulness of the proposed characterisation of the VQ of Parkinson's disease, we investigated
 279 using the q_l features for the prediction of the dysarthria level according to a modified version of the Frenchay
 280 assessment score. This perceptual evaluations includes the following aspects of speech: respiration, lips move-
 281 ment, palate/velum movement, larynx, tongue, and intelligibility. We hypothesised that the DPP features
 282 should be useful for prediction of the FDA scores related particularly to the larynx, which impacts the VQ
 283 the most.

284 The baseline features include articulation and prosody-based features, which are concatenated to form a
 285 724-dimensional feature vector per utterance (Orozco-Arroyave, 2016; Vasquez-Correa et al., 2017). The
 286 articulation-based features includes 86 descriptors such as the energy content distributed in 22 Bark bands in
 287 the transition from voiced to unvoiced segments (22 descriptors), and from unvoiced to voiced segments
 288 (22 descriptors) (Orozco-Arroyave et al., 2016). The feature set is augmented with the first and second formant
 289 frequencies, and 12 MFCC with their derivatives. The extracted features are grouped and four functionals are
 290 computed (mean, standard deviation, skewness, and kurtosis), forming a 344-dimensional feature vector per
 291 utterance. The second feature set contains prosody-based features computed with the Erlangen prosody module
 292 (Zeiler et al., 2006), using voiced segments as speech unit. The set of features comprises a total of 95 features.
 293 19 of them are based on duration and include among others the number and the length of voiced frames, and
 294 duration of pauses. 36 of the features are based on the F_0 contour, including the mean, standard deviation, jitter,
 295 and others. The energy-based features include measures of the energy within the voiced frames, shimmer, posi-
 296 tion of the maximum energy, and others. The features are grouped into one feature vector and four functionals
 297 are also computed: mean, standard deviation, maximum, and minimum, forming a 380-dimensional feature
 298 vector per utterance.

299 The evaluated features consisted of the concatenated baseline and q_l features calculated per speaker. All 50 PD
 300 speakers were considered in this evaluation. For the prediction task, we used the same Super Vector Regression as
 301 described by Vasquez-Correa et al. (2017), using a leave-one-subject-out (LOSO) cross-validation. The performance
 302 is evaluated using the Spearman's correlation coefficient between the predicted scores and the real scores. The real
 303 scores were obtained by three professional phoniatricians, with the inter-rater reliability of 0.86 measured as the
 304 average Spearman's correlation coefficient obtained between all the evaluators.

305 Table 4 shows the correlation achieved with the baseline and the q_l features. Improvements are obtained for the
 306 monologue and reading speech tasks, of 3% and 16%, respectively, whereas no improvement is obtained with
 307 the pataka speech task. The results imply that the proposed q_l features depend on statistics (μ_k as the mean values of
 308 phonological probabilities), and better results are obtained with more observed (recorded) data. For example, while
 309 the pataka tasks contain speech samples with repeated single word, the read text task includes speech samples of
 310 36 spoken words.

Table 4
 The Spearman's correlation coefficients between the real and predicted modified FDA scores related to the larynx. Median values are calculated for the correlations with the three evaluators. Results obtained for the three sub-sets of the PD data (see Section 4.2) are reported.

Speech task	Baseline features	Proposed features
Pataka	0.56	0.56
Read text	0.39	0.47
Monologue	0.55	0.57

311 6. Conclusions

312 The paper has proposed the characterisation of voice quality (VQ) applied to pathological speech in PD. Often,
313 the analysis of pathological speech is limited by available data, and advanced deep machine learning techniques can-
314 not be fully applied. The lack of proper perceptual labels of pathological speech adds further complication. There-
315 fore, the proposed characterisation learns statistics from healthy speech data that is more widely available, and
316 calculates similarity with disordered speech by using the Euclidean distance.

317 The results obtained by DPP features have been validated by matching the obtained most significant, non-modal
318 phonation types with evaluating parameters of the perceptual assessments. In addition, DPP features of PD have
319 been interpreted by the Frenchay assessment. This interpretation ability can be directly used in clinical assessment.

320 A drawback of the presented experimental study was in missing VQ perceptual labels of PD data. To the authors'
321 knowledge, the used PD database is the biggest open-source database available, containing both isolated and con-
322 nected speech, and was selected primarily for its size. By missing perceptual labels, validation of the proposed VQ
323 characterisation thus has been done on all speakers focusing on differentiating HC and PD speech, and its direct
324 application in diagnosis and therapy is limited. In future, we plan to obtain PD data with labeled VQ, and validate
325 the VQ characterisation on individual patients, looking for example for regression of the perceptual scores.

326 Acknowledgements

327 The work reported here was partially carried out during the 2016 Jelinek Memorial Summer Workshop on Speech
328 and Language Technologies, which was supported by Johns Hopkins University via DARPA LORELEI Contract no.
329 [HR0011-15-2-0027](#), and gifts from Microsoft, Amazon, Google, and Facebook. It was also partially supported by
330 CODI from Universidad de Antioquia, project [2015-7683](#), and COLCIENCIAS project # 111556933858.

331 References

- 332 Anderson, R.C., Klofstad, C.A., Mayew, W.J., Venkatachalam, M., 2014. Vocal fry may undermine the success of young women in the labor
333 market. *PLoS One* 9 (5), e97506+. doi: [10.1371/journal.pone.0097506](#).
- 334 Aronson, A.E., Bless, D., 2011. *Clinical Voice Disorders*. Thieme, New York.
- 335 Awan, S.N., Roy, N., Dromey, C., 2009. Estimating dysphonia severity in continuous speech: application of a multi-parameter spectral/cepstral
336 model. *Clin. Linguist. Phon.* 23 (11), 825–841.
- 337 Bangayan, P., Long, C., Alwan, A.A., Kreiman, J., Gerratt, B.R., 1997. Analysis by synthesis of pathological voices using the Klatt synthesizer.
338 *Speech Commun.* 22 (4), 343–368. doi: [10.1016/s0167-6393\(97\)00032-0](#).
- 339 Barsties, B., De Bodt, M., 2015. Assessment of voice quality: current state-of-the-art. *Auris Nasus Larynx* 42 (3), 183–188.
- 340 Bauer, V., Alerić, Z., Jančić, E., Miholović, V., 2011. Voice quality in Parkinson's disease in the Croatian language speakers. *Coll Antropol.*
341 35 (2), 209–212.
- 342 Bhuta, T., Patrick, L., Garnett, J.D., 2004. Perceptual evaluation of voice quality and its correlation with acoustic measurements. *J. Voice* 18 (3),
343 299–304.
- 344 Cernak, M., Asaei, A., Honnet, P.-E., Garner, P.N., Boulard, H., 2016. Sound pattern matching for automatic prosodic event detection. In: Pro-
345 ceedings of Interspeech, pp. 170–174.
- 346 Cernak, M., Benus, S., Lazaridis, A., 2017a. Speech vocoding for laboratory phonology. *Comput. Speech Lang.* 42, 100–121.
- 347 Cernak, M., Garner, P.N., 2016. PhonVoc: a phonetic and phonological vocoding toolkit. In: Proceedings of Interspeech. San Francisco, CA, USA,
348 pp. 988–992.
- 349 Cernak, M., Nöth, E., Rudzicz, F., Christensen, H., Orozco-Arroyave, J.R., Arora, R., Bocklet, T., Chinaei, H., Hannink, J., Nidavolu, P.S.,
350 Vasquez, J.C., Yancheva, M., Vann, A., Vogler, N., 2017. On the impact of non-modal phonation on phonological features. In: Proceedings of
351 International Conference on Acoustics, Speech and Signal Processing, ICASSP. IEEE.
- 352 Cernak, M., Potard, B., Garner, P.N., 2015. Phonological vocoding using artificial neural networks. In: Proceedings of International Conference on
353 Acoustics, Speech and Signal Processing, ICASSP. IEEE, pp. 4844–4848.
- 354 Chomsky, N., Halle, M., 1968. *The Sound Pattern of English*. Harper & Row, New York, NY.
- 355 Darley, F.L., Aronson, A.E., Brown, J.R., 1969. Differential diagnostic patterns of dysarthria. *J. Speech Lang. Hear. Res.* 12 (2), 246–269.
- 356 Drugman, T., Kane, J., Gobl, C., 2014. Data-driven detection and analysis of the patterns of creaky voice. *Comput. Speech Lang.* 28 (5), 1233–
357 1253. doi: [10.1016/j.csl.2014.03.002](#).
- 358 Enderby, P.M., 1983. *Frenchay Dysarthria Assessment*. College Hill Press.
- 359 Enderby, P.M., Palmer, R., 2008. *FDA-2: Frenchay Dysarthria Assessment: Examiner's Manual*. Pearson.
- 360 Garofolo, J.S., Lamel, L.F., Fisher, W.M., Fiscus, J.G., Pallett, D.S., 1993. *Darpa Timit Acoustic-phonetic Continous Speech Corpus CD-ROM*.
361 NASA STI/Recon technical report 93.
- 362 Hansen, J.H., Hasan, T., 2015. Speaker recognition by machines and humans: a tutorial review. *IEEE Signal Process. Mag.* 32 (6), 74–99.

- 363 Hanson, H.M., 1997. Glottal characteristics of female speakers: acoustic correlates. *J. Acoust. Soc. Am.* 101 (1), 466–481.
- 364 Hinton, G.E., Osindero, S., Teh, Y.W., 2006. A fast learning algorithm for deep belief nets. *Neural Comput.* 18 (7), 1527–1554. doi: [10.1162/neco.2006.18.7.1527](https://doi.org/10.1162/neco.2006.18.7.1527).
- 365
- 366 Hirano, M., 1981. *Clinical Examination of Voice. Disorders of human communication*, 5. Springer.
- 367 Holmes, R.J., Oates, J.M., Phylaud, D.J., Hughes, A.J., 2000. Voice characteristics in the progression of Parkinson's disease. *Int. J. Lang. Commun. Disord.* 35 (3), 407–418.
- 368
- 369 Kane, J., 2012. *Tools for Analysing the Voice*. Ph.D. thesis. Trinity College Dublin, Dublin.
- 370 Klatt, D.H., Klatt, L.C., 1990. Analysis, synthesis, and perception of voice quality variations among female and male talkers. *J. Acoust. Soc. Am.* 87 (2), 820–857.
- 371
- 372 Kreiman, J., Gerratt, B.R., Kempster, G.B., Erman, A., Berke, G.S., 1993. Perceptual evaluation of voice quality review, tutorial, and a framework for future research. *J. Speech Lang. Hear. Res.* 36 (1), 21–40.
- 373
- 374 Ladefoged, P., Johnson, K., 2014. *A Course in Phonetics*. 7th ed. Cengage Learning.
- 375 Laver, J., 1980. *The Phonetic Description of Voice Quality*. Cambridge Studies in Linguistics. Cambridge University Press
- 376 Little, M.A., McSharry, P.E., Hunter, E.J., Spielman, J., Ramig, L.O., et al., 2009. Suitability of dysphonia measurements for telemonitoring of Parkinson's disease. *IEEE Trans. Biomed. Eng.* 56 (4), 1015–1022.
- 377
- 378 Little, M.A., McSharry, P.E., Roberts, S.J., Costello, D.A., Moroz, I.M., 2007. Exploiting nonlinear recurrence and fractal scaling properties for voice disorder detection. *BioMed. Eng. OnLine* 6 (1), 23.
- 379
- 380 Malyska, N., 2008. *Analysis of Nonmodal Glottal Event Patterns with Application to Automatic Speaker Recognition*. Ph.D. thesis. Harvard University – MIT Division of Health Sciences and Technology, USA.
- 381
- 382 Malyska, N., Quatieri, T.F., Dunn, R.B., 2011. Sinewave representations of nonmodality. In: *Proceedings of Interspeech*, pp. 69–72.
- 383 Maryn, Y., De Bodt, M., Roy, N., 2010. The acoustic voice quality index: toward improved treatment outcomes assessment in voice disorders. *J. Commun. Disord.* 43 (3), 161–174.
- 384
- 385 Maryn, Y., Roy, N., De Bodt, M., Van Cauwenberge, P., Corthals, P., 2009. Acoustic measurement of overall voice quality: a meta-analysis. *J. Acoust. Soc. Am.* 126 (5), 2619–2634.
- 386
- 387 Oates, J., 2009. Auditory-perceptual evaluation of disordered voice quality. *Folia Phoniatr. Logop.* 61 (1), 49–56.
- 388 Orozco-Arroyave, J.R., 2016. *Analysis of Speech of People with Parkinson's Disease*. 41. Logos Verlag Berlin GmbH.
- 389 Orozco-Arroyave, J.R., Arias-Londoño, J.D., Bonilla, J.F.V., Gonzalez-Rátiva, M.C., Nöth, E., 2014. New Spanish speech corpus database for the analysis of people suffering from Parkinson's disease. In: *Proceedings of Conference on Language Resources and Evaluation, LREC*, pp. 342–347.
- 390
- 391 Orozco-Arroyave, J.R., Vásquez-Correa, J.C., Hönl, F., Arias-Londoño, J.D., Vargas-Bonilla, J.F., Skodda, S., Rusz, J., Nöth, E., 2016. Towards an automatic monitoring of the neurological state of the Parkinson's patients from speech. In: *Proceedings of the 41st International Conference on Acoustic, Speech, and Signal Processing, ICASSP*.
- 392
- 393
- 394 Paul, D.B., Baker, J.M., 1992. The design for the Wall Street Journal-based CSR corpus. In: *Proceedings of the Workshop on Speech and Natural Language*. Association for Computational Linguistics, Stroudsburg, PA, USA, pp. 357–362. doi: [10.3115/1075527.1075614](https://doi.org/10.3115/1075527.1075614).
- 395
- 396 Podesva, R.J., 2007. Phonation type as a stylistic variable: the use of falsetto in constructing a persona. *J. Socioling.* 11 (4), 478–504.
- 397 Povey, D., Ghoshal, A., Boulianne, G., Burget, L., Glembek, O., Goel, N., Hannemann, M., Motlicek, P., Qian, Y., Schwarz, P., Silovsky, J., Stemmer, G., Vesely, K., 2011. The Kaldi speech recognition toolkit. In: *Proceedings of Automatic Speech Recognition & Understanding, ASRU*. IEEE SPS. IEEE Catalog No.: CFP11SRW-USB.
- 398
- 399
- 400 Rasipuram, R., Magimai-Doss, M., 2011. Integrating articulatory features using Kullback–Leibler divergence based acoustic model for phoneme recognition. In: *Proceedings of International Conference on Acoustics, Speech and Signal Processing, ICASSP*. IEEE, pp. 5192–5195. doi: [10.1109/icassp.2011.5947527](https://doi.org/10.1109/icassp.2011.5947527).
- 401
- 402
- 403 Rusz, J., Cmejla, R., Ruzickova, H., Ruzicka, E., 2011. Quantitative acoustic measurements for characterization of speech and voice disorders in early untreated Parkinson's disease. *J. Acoust. Soc. Am.* 129 (1), 350–367.
- 404
- 405 San Segundo, E., Tsanas, A., Gómez-Vilda, P., 2017. Euclidean distances as measures of speaker similarity including identical twin pairs: a forensic investigation using source and filter voice characteristics. *Forensic Sci. Int.* 270, 25–38.
- 406
- 407 Schuller, B., Batliner, A., 2013. *Computational Paralinguistics: Emotion, Affect and Personality in Speech and Language Processing*. Wiley.
- 408 Shinoda, K., Watanabe, T., 1997. Acoustic modeling based on the MDL principle for speech recognition. In: *Proceedings of Eurospeech*, pp. 1–99–102.
- 409
- 410 Stouten, F., Martens, J.-P., 2006. On the use of phonological features for pronunciation scoring. In: *Proceedings of International Conference on Acoustics, Speech and Signal Processing, ICASSP*. IEEE, p. I. doi: [10.1109/icassp.2006.1660024](https://doi.org/10.1109/icassp.2006.1660024).
- 411
- 412 Titze, I.R., 1995. *Workshop on Acoustic Voice Analysis: Summary Statement*. National Center for Voice and Speech.
- 413 Vasquez-Correa, J.C., Orozco-Arroyave, J.R., Arora, R., Nöth, E., Dehak, N., Christensen, H., Rudzicz, F., Bocklet, T., Cernak, M., Chinaei, H., Hannink, J., Nidadavolu, P.S., Yancheva, M., Vann, A., Vogler, N., 2017. Multi-view representation learning via GCCA for multimodal analysis of Parkinson's disease. In: *Proceedings of International Conference on Acoustics, Speech and Signal Processing, ICASSP*.
- 414
- 415
- 416 Wuyts, F.L., De Bodt, M.S., Molenberghs, G., Remacle, M., Heylen, L., Millet, B., Van Lierde, K., Raes, J., Van de Heyning, P.H., 2000. The dysphonia severity index: an objective measure of vocal quality based on a multiparameter approach. *J. Speech Lang. Hear. Res.* 43 (3), 796–809.
- 417
- 418 Yu, D., Siniscalchi, S., Deng, L., Lee, C.-H., 2012. Boosting attribute and phone estimation accuracies with deep neural networks for detection-based speech recognition. In: *Proceedings of International Conference on Acoustics, Speech and Signal Processing, ICASSP*. IEEE SPS.
- 419
- 420 Zeiler, V., Adelhardt, J., Batliner, A., Frank, C., Nöth, E., Shi, R.P., Niemann, H., 2006. The prosody module. *SmartKom: Foundations of Multimodal Dialogue Systems*. Springer, pp. 139–152.
- 421
- 422 Zen, H., Nose, T., Yamagishi, J., Sako, S., Masuko, T., Black, A., Tokuda, K., 2007. The HMM-based speech synthesis system version 2.0. In: *Proceedings of ISCA Speech Synthesis Workshop, SSW6*, pp. 131–136.
- 423