



**SYLLABLE-LEVEL FEATURES FOR SPEECH  
PATHOLOGY DETECTION: A CASE STUDY OF  
PARKINSON'S DISEASE**

Sevada Hovsepyan<sup>a</sup>      Mathew Magimai-Doss

Idiap-RR-02-2026

MAY 2026

---

<sup>a</sup>Idiap



# Syllable-Level Features for Speech Pathology Detection: A Case Study of Parkinson's Disease

Sevada Hovsepyan<sup>a</sup> and Mathew Magimai.-Doss

*Idiap Research Institute, Martigny, CH-1920, Switzerland*

1 In this paper we explore newly proposed syllable level feature extraction for speech  
2 pathology detection: focusing on the Parkinson's disease detection from speech.  
3 The method is inspired by the spectro-temporal representations used in the neuro-  
4 computational models of speech perception, and simply represents a standardized  
5 representation of the frequency content of the syllable-like segments. In this study,  
6 we demonstrate that this type of representation is powerful enough to successfully  
7 detect various speech pathologies from speech samples, independent of speech type,  
8 language, and dataset. Furthermore, our analyses show that these representations  
9 lead to explainable and interpretable results, and that changes in pathologies are  
10 related to both system- and source-related changes in speech production. Overall,  
11 syllable-level features suggest themselves as a simple but robust and explainable ap-  
12 proach to understand physiological changes in speech production due to different  
13 pathologies, as well as to detect those pathologies from speech.

---

<sup>a</sup>[sevada.hovsepyan@idiap.ch](mailto:sevada.hovsepyan@idiap.ch)

## 14 I. INTRODUCTION

15 Parkinson’s disease (PD) stands out as one of the most prevalent neurodegenerative dis-  
16 orders, affecting millions of individuals worldwide. PD impacts both motor ([Mazzoni \*et al.\*,](#)  
17 [2012](#)) and non-motor systems, with symptoms ranging from motor planning difficulties,  
18 tremors, and bradykinesia to cognitive impairments, neural degeneration, gastrointestinal  
19 issues, and neuropsychiatric conditions. Early diagnosis and intervention are crucial for more  
20 effective treatment and management of the disease. In recent years, numerous studies have  
21 focused on developing non-invasive methods for detecting PD, such as gait and handwrit-  
22 ing analysis, which achieve high detection accuracy by assessing motor deficiencies. Blood  
23 analysis and cognitive tasks are also employed to evaluate cognitive impairments. However,  
24 these methods often require special equipment and settings, and most importantly, rarely  
25 capture both cognitive and motor aspects simultaneously, which would be beneficial for a  
26 more comprehensive assessment of the disorder.

27 From this perspective, speech-based PD detection is unique, as it can capture both motor-  
28 related impairments (e.g., imprecise articulation) and cognitive impairments (e.g., language  
29 planning). Moreover, speech based pathology detection are often easy to deploy, scalable,  
30 and not necessarily require special equipment (e.g. phone’s microphone can be sufficient).  
31 Research in this area has evolved from expert-crafted, hand-engineered features, such as  
32 vowel triangles and formant patterns, hypothesized to reflect motor impairments in vocal  
33 cord vibration and articulatory control ([Karlsson \*et al.\*, 2020](#); [Michaelis \*et al.\*, 1998](#)), to  
34 large consortia of handcrafted feature sets like ComPARE ([Schuller \*et al.\*, 2010, 2013](#)) and

35 eGeMAPS (Eyben *et al.*, 2016). The latter, provided more data-driven approach and often  
36 operated under "brute-force" paradigm, by capturing a broad spectrum of paralinguistic  
37 aspects of speech, though not necessarily relevant to PD. More recently, deep learning ap-  
38 proaches, particularly convolutional neural networks (CNNs) and adaptive architectures,  
39 have achieved strong empirical performance (Escobar-Grisales *et al.*, 2023; van Gelderen  
40 and Tejedor-García, 2024). However, these end-to-end models, while powerful, are typically  
41 generic and lack interpretability, a critical factor for clinical adoption (Ananthanarayanan  
42 *et al.*, 2025).

43 Current approaches predominantly focus on acoustic deficits in PD speech, which are  
44 largely attributed to motor impairments. Yet, PD, being a neurodegenerative disease, also  
45 affects neural functioning and cognitive processing, that can manifest in produced speech.  
46 This oversight is notable given that speech production and perception are tightly coupled  
47 processes, shaped not only by motor control and articulation but also by underlying neural  
48 activity and linguistic content.

49 To address this gap, we recently proposed a novel approach for PD speech analysis (Hov-  
50 sepyan and Magimai.-Doss, 2024) based on the acoustic representations rooted in neurophys-  
51 iologically plausible generative models of speech (Hovsepyan *et al.*, 2020; Nabé *et al.*, 2021).  
52 These models posit that speech is parsed into syllable-like units during perception, and that  
53 disruptions in speech production due to PD may leave detectable traces in these perceptual  
54 representations. By focusing on syllable-level spectrotemporal patterns, our method aims  
55 to capture both motor and neural aspects of PD-related speech changes, offering a more  
56 holistic and interpretable approach to pathology detection.

57 In this study, our objective is to extensively explore and refine our proposed SLF for the  
58 detection of Parkinson’s disease by speech. We test our method on various speech types  
59 and datasets, including the PC-GITA dataset (Orozco *et al.*, 2014) and the Dutch Corpus  
60 of Pathological and Normal Speech (COPAS) (Martens *et al.*, 2011). Our objectives are  
61 to determine the robustness and generalizability of our feature construction method and to  
62 provide insights into the physiological mechanisms underlying PD-related speech changes.

63 The remainder of the manuscript is organized as follows. In the Background, we provide a  
64 quick overview of the aforementioned generative models and delve deeply into the proposed  
65 SLF approach. In the Methods section, we detail the datasets used in the study, outline the  
66 construction of syllable-level features, and describe the methods used for classification and  
67 evaluation of our results. In the Results section, we present our classification outcomes and  
68 compare them with baseline feature sets. In the final section, we discuss our findings and  
69 outline future directions.

## 70 II. BACKGROUND

71 In this section, we present the primary motivation for examining syllables as the basis for  
72 feature extraction in Parkinson’s disease speech analysis. We then briefly describe the main  
73 findings from our pilot investigation using syllable-level features and outline the specific  
74 objectives of this manuscript.

## 75 **A. Syllables as speech building blocks**

76 Syllables are linguistic units widely recognized as fundamental building blocks of speech,  
77 occupying a central role at the intersection of acoustic, linguistic, cognitive, and neural  
78 aspects of speech production and perception. From a cognitive perspective, syllables serve  
79 as discrete planning units for speech production and as perceptual chunks during speech  
80 perception. From a neurophysiological standpoint, it is hypothesized that the brain employs  
81 dedicated mechanisms, such as theta oscillations (4–8 Hz), to segment the continuous speech  
82 stream into syllable-like, discrete units during perception. These multifaceted roles make  
83 syllables ideal candidates for investigating how neurodegenerative diseases such as PD alter  
84 speech. Specifically, syllable-level analysis can reveal the sources of deficiencies at multiple  
85 levels, extending beyond motor dysfunction to include disruptions in neural, cognitive, and  
86 linguistic processing.

## 87 **B. Speech Production-Perception Cycle and Neurocomputational Models of** 88 **SPeech Perception**

89 In the neuroscience of speech processing, a seminal framework is the analysis-by-synthesis  
90 approach ([Bever and Poeppel, 2010](#); [Halle and Stevens, 1962](#); [Nair \*et al.\*, 2008](#)), a theory  
91 rooted in the motor theory of speech perception ([Dogonasheva \*et al.\*, 2025](#); [Hovsepian \*et al.\*,](#)  
92 [2023](#); [Hyafil \*et al.\*, 2015](#); [Nabé \*et al.\*, 2021](#)). This framework posits that speech perception  
93 involves the internal generation of motor commands associated with incoming acoustic sig-  
94 nals. Contemporary developments, such as predictive coding and the Free Energy Principle

95 (FEP) (Buckley *et al.*, 2017; Friston, 2010; Friston and Kiebel, 2009), extend this idea,  
96 proposing that the brain functions as a statistical, generative system. It anticipates sensory  
97 input based on accumulated internal knowledge of the external world and continuously syn-  
98 thesises expected speech patterns to minimize prediction error and infer incoming sensory  
99 signal, such as speech.

100 In recent years, neurocomputational models of speech perception have been developed  
101 that integrate neural oscillations, syllable-centric processing, and hierarchical predictive cod-  
102 ing (Dogonasheva *et al.*, 2025; Hovsepian *et al.*, 2023; Hyafil *et al.*, 2015; Nabé *et al.*, 2021).  
103 These generative models use parameterized spectrotemporal patterns of syllables to simu-  
104 late expected acoustic signals and refine perceptions by minimizing prediction errors. While  
105 primarily designed to explore brain function during speech perception, these models offer  
106 a unique perspective for analyzing deficits in speech production. Given that PD disrupts  
107 neural processing and motor execution of speech, these disruptions should also manifest in  
108 the spectrotemporal patterns of syllables. If these models can infer syllable identity from  
109 such patterns, then differences in these patterns should emerge between healthy and PD  
110 conditions.

111 Figure 1 illustrates this conceptual bridge: while traditional approaches analyze the  
112 acoustic waveform directly (a bottom-up process), our method leverages the structure of  
113 perceptual models to interpret production deficits (a top-down/bidirectional approach).

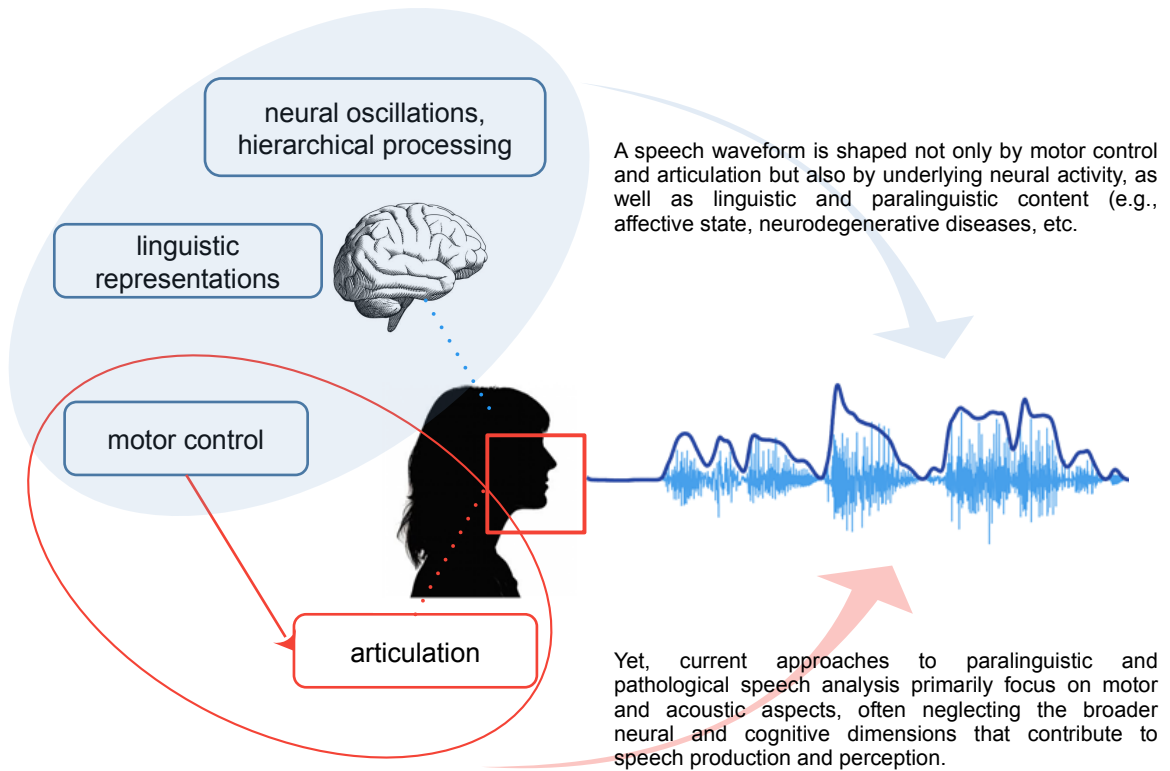


FIG. 1. Speech production and perception are tightly interlinked processes, shaped by motor control, articulatory mechanisms, neural activity (e.g., neural oscillations and hierarchical processing), and linguistic and paralinguistic factors (e.g., emotional tone and cognitive load). This interplay is particularly critical in neurodegenerative disorders, which often disrupt both the motor execution of speech (e.g., articulation) and the neural substrates underlying linguistic and paralinguistic functions (e.g., emotional prosody and cognitive planning). Current approaches to pathological speech analysis, however, predominantly focus on motor and acoustic deficits in the speech waveform. In this manuscript, we argue that valuable insights into pathological speech can also be gained by analyzing components relevant to speech perception.

### 114 C. Syllable Level Features for Parkinson’s disease detection

115 Based on the ideas discussed in the previous two sections, that syllables serve as building  
116 blocks for speech at multiple levels (acoustic, neural, cognitive, linguistic) and that speech  
117 production deficits can be understood from a perception perspective, we recently proposed  
118 a feature engineering pipeline for PD speech analysis. The so-called syllable-level features  
119 (SLFs) are derived from the spectral representations of syllables used in neurocomputational  
120 models of speech perception. The procedure involves calculating a standardized represen-  
121 tation of the spectrotemporal patterns of syllables (with a fixed number of frequency and  
122 temporal bins) and using these patterns (after flattening) as feature vectors for PD speech  
123 classification. In the next section, we detail the construction of SLFs, whereas this section  
124 summarizes the results of our pilot study ([Hovsepyan and Magimai.-Doss, 2024](#)).

125 Neurocomputational models aim to mimic brain activity and be neurophysiologically  
126 plausible; therefore, to calculate the spectrotemporal pattern of syllables or speech, they  
127 often employ a model of the cochlea. Following the same logic, we initially used the same  
128 model of the auditory periphery as the aforementioned models, but we also compared other  
129 approaches for spectrogram calculation, such as Mel-frequency cepstral coefficients (MFCC)  
130 or the Short-Time Fourier Transform (STFT). In addition to comparing different spectro-  
131 gram calculation methods, we aimed to explore the optimal number of frequency channels  
132 used for SLF calculation (specifically, into how many frequency bins the spectrogram should  
133 be collapsed).

134 Our analysis of diadochokinetic (DDK) speech samples from the PC-GITA dataset re-  
135 vealed that the best classification performance was achieved using STFT decomposition,  
136 with performance improving as the number of frequency channels increased (we tested 6,  
137 22, and 46 channels). We also examined whether the method of deriving syllable chunks  
138 affected classification performance, specifically comparing syllable onset-based segmenta-  
139 tion (valleys in the sonority envelope) to syllable nucleus-based segmentation (peaks in the  
140 envelope). The analysis showed that, on average, nucleus-based chunking (peaks in the en-  
141 velope) yielded better performance (83.23% vs. 80.1% AUROC for the 46-channel STFT  
142 case). Furthermore, our pilot analysis demonstrated that segmenting long DDK utterances  
143 into linguistically informed chunks (such as syllable onsets or nuclei) was more advantageous  
144 than using fixed-length chunks of comparable duration (AUROC = 71.85% for 25-50ms seg-  
145 ments, AUROC = 71.48% for 25-450ms segments and AUROC = 71.32% for 300-600ms  
146 segments).

147 The results of the pilot study serve as a starting point for the more in-depth analysis  
148 performed in this manuscript. Based on our previous results, we focus on the STFT spec-  
149 trogram of syllables and use syllable nuclei to chunk syllable-like segments of the analyzed  
150 speech signal. Beyond these parameters, the goal of this manuscript is to conduct a more  
151 intensive and in-depth analysis of syllable-level features. More specifically, we delve into  
152 how SLFs can be applied to speech types other than DDK utterances, what kind of acoustic  
153 information SLFs capture, and whether the method is applicable to other pathologies and/or  
154 languages.

### 155 III. METHODS

#### 156 A. Datasets

157 In this study, we used two well-established datasets for pathological speech research: the  
158 "New Speech Corpus Database for Analysis of People Suffering from Parkinson's Disease"  
159 (PC-GITA) (Orozco *et al.*, 2014) and the Dutch Corpus of Pathological and Normal Speech  
160 (COPAS) (Martens *et al.*, 2011).

161 The PC-GITA dataset consists of recordings from 50 individuals with Parkinson's disease  
162 (PD) and 50 healthy controls (HC) (Orozco *et al.*, 2014). The PD patients were diagnosed  
163 by expert neurologists and were classified as having different levels of severity based on  
164 the Unified Parkinson's Disease Rating Scale (UPDRS) (Goetz *et al.*, 2008) and Hoehn  
165 and Yahr (H&Y) scales (Hoehn and Yahr, 1967). Both groups were matched by age and  
166 gender. The PD group consists of 25 males with a mean age of  $62.2 \pm 11.2$  years and 25  
167 females with a mean age of  $60.1 \pm 7.8$  years. The HC group consists of 25 males with  
168 a mean age of  $61.2 \pm 11.3$  years and 25 females with a mean age of  $60.7 \pm 7.7$  years.  
169 The recordings were performed in a soundproof booth with a professional microphone, and  
170 the resulting audio files were sampled at 44,100 Hz with 16-bit resolution. The PC-GITA  
171 dataset includes various speech types produced by all speakers, such as sustained vowels,  
172 rapid syllable repetitions (diadochokinetic (DDK) utterances), monologues, and read text  
173 passages and sentences. Since this study focuses on syllable-level features, sustained vowels  
174 were excluded from the experiments, and only DDK, read text, monologue, and sentence  
175 utterance types were used.

176 In contrast, the COPAS dataset includes recordings of over 300 speakers performing var-  
177 ious speech tasks, such as repeating sentences and reading texts (Martens *et al.*, 2011). The  
178 dataset includes recordings of healthy speakers and speakers with speech pathologies, such  
179 as hearing impairments, articulation disorders, and dysarthria. Recordings were performed  
180 in a quiet clinical setting and sampled at 16 kHz and 16-bit resolution. However, the num-  
181 ber of recordings varies greatly depending on the type of pathology. Additionally, not all  
182 speakers performed all speech tasks; therefore, the number of recordings per pathology and  
183 task varies. For the current experiment, we only needed long-form speech tasks, such as sen-  
184 tences and reading texts. Therefore, we only used two pathologies: dysarthria and hearing  
185 impairment, as there were enough recordings for analysis.

## 186 B. Control feature sets

187 As mentioned in the Background section, we chose ComPARE (Schuller *et al.*, 2010) and  
188 eGeMAPS (Eyben *et al.*, 2016) feature sets as controls. These feature sets were specifically  
189 designed for paralinguistic challenges and have been utilized in various competitions over the  
190 years. The ComPARE feature set, for instance, was created for speech analysis challenges,  
191 including those at Interspeech (Schuller *et al.*, 2010, 2015, 2016, 2013, 2021), and contains  
192 a large number of features related to various speech characteristics such as pitch variations,  
193 formants, and harmonics.

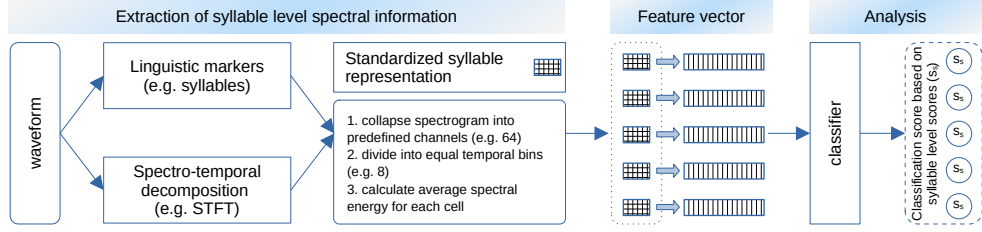
194 eGeMAPS (Eyben *et al.*, 2016), on the other hand, is a derivative of the ComPARE  
195 feature set, where a handful of features are specifically selected for their importance in par-  
196 alinguistic and clinical speech analysis, including PD. After several revisions, the eGeMAPS

197 feature set now includes 88 low-level descriptors, covering spectral, energy, and frequency  
198 aspects of speech. The SMILE toolbox (Eyben *et al.*, 2010), with default parameters, was  
199 used to extract eGeMAPS and ComPARE feature sets from both datasets at the utterance  
200 level, so that for each utterance, we have a fixed number of features suitable for standard  
201 classification procedures.

### 202 C. Syllable level features

203 Syllable-level features were extracted using a procedure that closely followed the steps  
204 described in a previous study (see (Hovsepyan and Magimai.-Doss, 2024) and Figure 2a).  
205 First, syllable nuclei timesteps were detected using the syllable onset detection algorithm  
206 proposed by (Räsänen *et al.*, 2018). The spectrogram of each utterance was then extracted  
207 using the short-time Fourier transform. Based on the syllable boundaries and STFT in-  
208 formation, a standardized representation of each syllable was constructed by averaging the  
209 spectral energy across neighboring frequency channels and equally distributed temporal bins  
210 (see Figure 2b). Consequently, the resulting spectro-temporal representation maintained  
211 consistent temporal and frequency dimensions (64 frequency bins and 8 temporal bins in  
212 this study), irrespective of the original syllable’s duration. Finally, the standardized repre-  
213 sentation was flattened into a feature vector for use in classification analyses. The details of  
214 each step are outlined below.

**a. Flowchart of proposed methodology**



**b. Syllable level feature extraction**

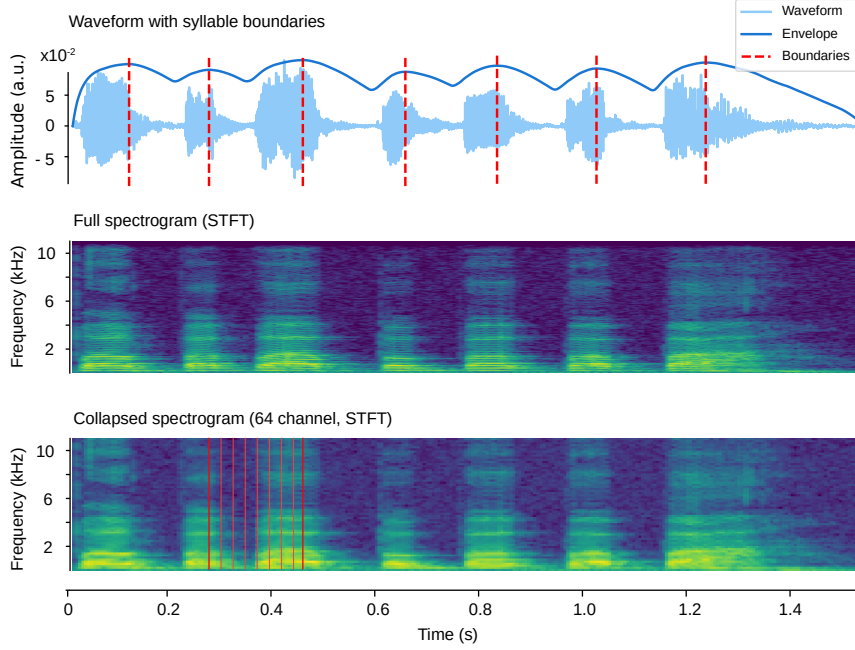


FIG. 2. Section **a** illustrates the pipeline for syllable-level feature extraction. For each speech waveform (top panel, section **b**), linguistic markers, such as syllable boundaries (dashed red lines) and the corresponding spectrogram (middle panel, section **B**), are first extracted. Next, a standardized syllable representation is computed: the spectrogram (bottom panel, section **b**) is collapsed into predefined frequency bins (e.g., 64 bins), and the syllable’s duration is divided into fixed, equal temporal bins (e.g., 8 bins, marked by red vertical lines). The average spectral energy is then calculated for each frequency–temporal bin. As a result, each syllable is represented as a fixed-dimensional matrix ( $8 \times 64$ ), which is flattened into a feature vector for subsequent classification analysis.

216 We used the algorithm proposed by (Räsänen *et al.*, 2018) for syllable boundary detection.  
217 The algorithm is inspired by the oscillatory mechanism that the brain uses to segment con-  
218 tinuous speech into syllabic segments by entraining theta (4–8 Hz) oscillations to the speech  
219 signal. In their model, (Räsänen *et al.*, 2018) and his colleagues first extracted the sonority  
220 envelope of the speech signal by filtering the speech signal into sub-bands using gammatone  
221 filter banks and then combining them to obtain the sonority envelope. The resulting signal  
222 was then fed to a damped theta oscillator, which shifted the phase of the oscillator to align  
223 with the peaks in the sonority envelope, which are typically associated with syllable nuclei.  
224 Syllable boundaries were associated with oscillatory phase resets that occurred at the local  
225 minima (valleys) of the sonority envelope. Additionally, candidate peaks were filtered based  
226 on duration constraints (e.g., peaks closer than 50 ms were merged.). The algorithm’s out-  
227 put is the time steps (in milliseconds) of the syllable boundaries (valleys) and nuclei (peaks).  
228 In this study, we used the information about syllable nuclei as syllable markers.

229 Additionally, for each speech type, we calculated the syllable duration distribution, de-  
230 fined as the duration between two consecutive syllable nuclei or peaks of the sonority enve-  
231 lope. For speech types that had more than one utterance per speaker (e.g., DDK utterances),  
232 the syllable duration distribution was based on syllable segments from all utterances. Fur-  
233 thermore, syllable duration distributions were approximated using Gaussian kernels, such  
234 that for each speech type, there was a Gaussian approximation for the respective syllable

235 duration distribution. This information was further used for control analysis, highlighting  
236 the importance of extracting features based on linguistic markers, such as syllables.

## 237 *2. Spectrogram: Short Term Fourier Transform*

238 During this study, we used the librosa library (McFee *et al.*, 2015) to calculate the power  
239 spectrogram with the short-term Fourier transform of each utterance. The following pa-  
240 rameters were used throughout the simulations: ‘nfft points’ = 511 and ‘hop length’ = 128.  
241 Since the PC-GITA dataset (Orozco *et al.*, 2014) includes audio files sampled at 44.1 kHz,  
242 we loaded those files with the librosa default settings, resulting in a sampling rate of 22.05  
243 kHz. For the COPAS dataset (Martens *et al.*, 2011), which is sampled at 16 kHz, we did  
244 not resample the audio signal during the STFT calculation. This yielded to ‘hop length’ of  
245 23 ms and 30 ms for the PC-GITA and COPAS datasets, respectively.

## 246 *3. Syllable based segmentation*

247 We calculated the syllable-level features for each utterance based on the extracted power  
248 spectrogram and the corresponding syllable boundaries (nuclei timesteps), following the  
249 procedure below.

250 First, we averaged the power in adjacent frequency channels of the power spectrogram to  
251 collapse it into 64 frequency channels. We then divided the resulting collapsed spectrogram  
252 into syllable segments, with each segment consisting of the spectral information between  
253 two consecutive syllable nuclei.

254 Furthermore, each syllable segment was divided into eight equal temporal bins (red ver-  
255 tical lines in the bottom panel of Figure 2b), and the average activity of each frequency  
256 channel within each bin was calculated. Thus, each syllable segment is represented in a  
257 standardized way with sixty-four frequency and eight temporal bins. Finally, we trans-  
258 formed the standardized spectrogram for each syllable into the decibel (db) scale, and then  
259 flattened it to obtain a syllable-level feature vector with dimensions of  $1 \times (8 \times 64 = 512)$ .

#### 260 *4. Control segmentation*

261 To assess whether segmentation based on linguistic markers, such as syllables, is critical  
262 for feature extraction, we extracted SLF-like representations using randomly segmented  
263 utterances. The procedure mirrored that used for SLF extraction, with one key difference:  
264 utterances were not segmented according to syllable markers. Instead, segment durations  
265 were randomly sampled from a Gaussian distribution, derived before, approximating the  
266 natural syllable duration distribution.

267 Next, each random segment, similar to the syllable segment, was divided into 8 equal  
268 parts, and the average power of each bin was calculated. The resulting random segment  
269 spectrogram, with a size of 64 frequency bins and 8 time bins, was transformed into a dB  
270 scale, and the resulting Random Segment Features (RSF) were calculated.

#### 271 *5. Formant based syllable level features*

272 Original SLF construction, the spectrum of a syllable was binned into fixed frequency bins  
273 to get the collapsed spectrogram with reduced number of frequency channels. We further

274 explored the benefits of dynamic frequency binning, e.g. aggregating spectral information  
275 across formant and fundamental frequency patterns. For this analysis, we made the fol-  
276 lowing modifications to the previously described SLF extraction procedure. First, for each  
277 utterance, we extracted the power STFT and the syllable boundaries, along with the patterns  
278 of the fundamental frequency (F0) and the first four formants (F1-F4). Additionally, for-  
279 mant bandwidths were extracted, whereas for the fundamental frequency, a fixed bandwidth  
280 equal to a quarter of the F0 was used. The formant patterns (with their respective band-  
281 widths) and the fundamental frequency patterns were extracted using the PRAAT (Boersma  
282 and Weenink, 2003) implementation in Python via the Parselmouth library (Jadoul *et al.*,  
283 2018). Next, for each syllable-like segment, the average energy (spectrogram amplitude)  
284 around the pitch (F0  $\pm$  fixed bandwidth) and formant (F1-F4) patterns ( $\pm$  respective band-  
285 width) was derived. Each syllable segment was then divided into T=8 equal temporal bins,  
286 and the average energy associated with the F0 and formant patterns and their respective  
287 bandwidths was derived. This procedure results in a 5x8 spectrotemporal representation of  
288 a syllabic segment. Finally, the flattened representation (1x40) is used for the classification  
289 procedure, as in the earlier proposed SLFs.

#### 290 **D. Classification**

291 Throughout the manuscript, we used a Support Vector Machine (SVM) with a cubic poly-  
292 nomial kernel for all classification. We used the standard parameters with a fixed random  
293 state from the *scikit-learn* Python package without hyperparameter tuning. Different clas-  
294 sification protocols were used for the two datasets. For the PC-GITA dataset (Orozco *et al.*,

295 2014), which has balanced samples of healthy and pathological speech, we used the leave-  
296 one-sample-out (LOSO) protocol. For the COPAS dataset (Martens *et al.*, 2011), which has  
297 imbalanced samples of healthy and pathological speech, we used a stratified K-fold protocol  
298 ( $K = 4$  for hearing impairment and  $K = 5$  for dysarthria).

299 In each case, the classification results are reported based on the accumulated classification  
300 scores from each fold. The scores are obtained by either combining the classification scores  
301 for each syllable in the utterance (called syllable-level scores) or calculating the average of  
302 the syllable scores for each utterance, which leads to an utterance-level score.

### 303 E. Statistical analysis

304 We ran a linear mixed random effects model to test whether eGeMAPS features were  
305 statistically different between conditions. To do this, we constructed the linear model with  
306 fixed effects of eGeMAPS features and PD vs. HC condition, with speakers coded as a  
307 random effect. We then looked specifically at the interaction term between features and  
308 conditions and extracted which features were statistically significant ( $p < 0.05$ ) across con-  
309 ditions (corrected for multiple comparisons using the false discovery rate).

## 310 IV. RESULTS

311 In this section we extensively test the proposed SLF method for speech pathology detec-  
312 tion. We first extend our earlier report on PD detection from DDK speech utterances from  
313 PC-GITA dataset. Consequently we run classification analysis based on SLF on different  
314 speech types, such as text reading, monologues and sentences. Furthermore, we explore

315 what kind of information SLF capture and lastly we generalize the proposed method by  
316 testing it on COPAS dataset on different pathologies/languages.

#### 317 **A. Parkinson’s disease detection from speech with syllable level features**

318 We first expanded our earlier report on SLF’s capacity to detect PD from speech by con-  
319 ducting classification analyses on additional speech samples and types from the PC-GITA  
320 dataset (Orozco *et al.*, 2014). The results are visualized in Figure 3, with detailed perfor-  
321 mance metrics provided in Supplementary Tables 1a and 1b for syllable- and sentence-level  
322 scores, respectively. For comparison, we also include AUC values derived from classification  
323 using established feature sets, such as eGeMAPS and ComPARE (Figure 4), as controls  
324 for each speech type. In Figure 4, each bar plot shows the average AUC for the respective  
325 speech type, for each feature set. Detailed performance metrics for control features sets  
326 can be found in Supplementary Tables 1c and 1d for eGeMAPS and ComPARE features  
327 respectively.

328 We extend our earlier report on SLF’s capacity to detect PD from speech by conducting  
329 classification analyses on additional speech samples and types from the PC-GITA dataset  
330 (Orozco *et al.*, 2014). Results are visualized in Figure 3, with detailed performance met-  
331 rics provided in Supplementary Tables 1a and 1b for syllable- and utterance-level scores,  
332 respectively.

333 For comparison, we also include AUC values obtained using established feature sets,  
334 eGeMAPS and ComPARE, as controls for each speech type (Figure 4). In Figure 4, each  
335 bar represents the mean AUC across participants for the respective speech type (grey dots

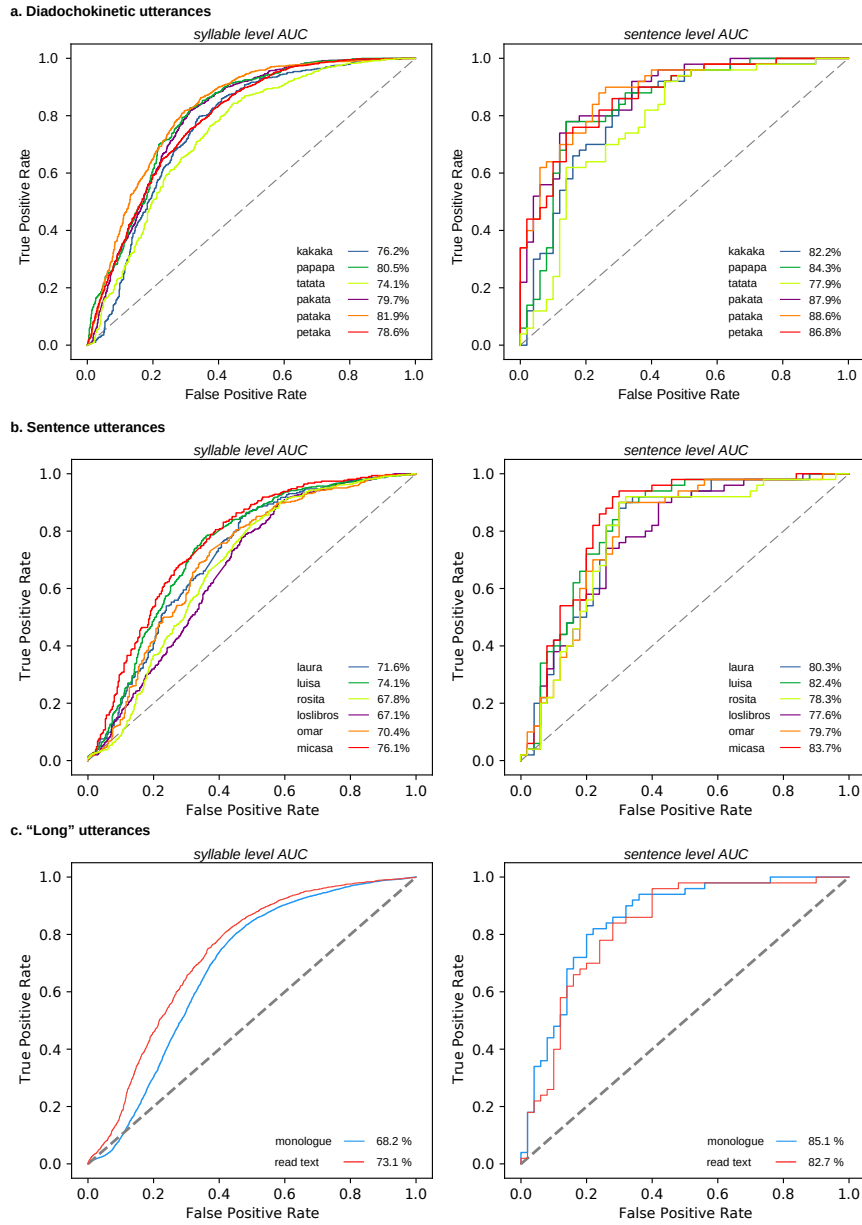


FIG. 3. Classification results on the PC-GITA dataset using SLF (utterance-level scores) and control feature sets are shown in the figure. As illustrated, the proposed SLF approach matches or outperforms the classification performance of all control feature sets across all speech types.

336 represent different utterances within speech type) and feature set. Detailed performance

337 metrics for the control feature sets are reported in Supplementary Tables 1c and 1d for  
338 eGeMAPS and ComPARE, respectively.

339 The results demonstrate consistently strong classification performance across all speech  
340 types, reinforcing the suitability of SLF as a robust candidate for PD detection from speech.  
341 Notably, utterance-level scores systematically outperform syllable-level scores in classifica-  
342 tion accuracy. However, syllable-level scores remain well above chance, indicating that even  
343 at this finer granularity, sufficient discriminative information is captured to detect PD-  
344 related speech alterations. This suggests that while individual syllables provide meaningful  
345 evidence, aggregating information across multiple syllable segments within an utterance  
346 further enhances classification performance by leveraging cumulative evidence.

347 Moreover, when compared to the control feature sets, the SLF approach, focusing on  
348 utterance-level scores (to align with the utterance-based nature of eGeMAPS and Com-  
349 PARE), either outperforms or yields results comparable to those of the established feature  
350 sets (Figure 4). This suggests that the proposed SLF-based classification performance is on  
351 par with that of more traditional feature sets.

## 352 **B. Importance of syllable segments**

353 Having established that SLF provide sufficient information to discriminate between PD  
354 speech and HC speech, we then further focused on the importance of the segments being  
355 based on linguistic markers, such as syllable nuclei. To do this, we conducted an additional  
356 set of analysis comparing the SLF approach with a similar approach where segments are

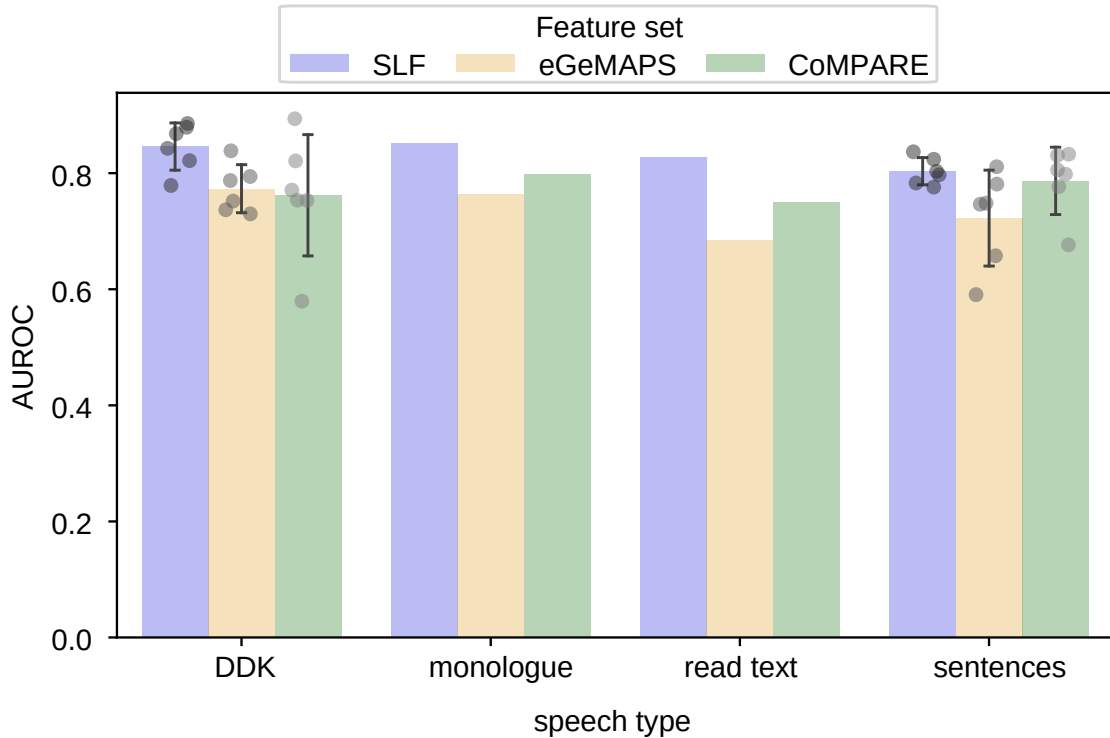


FIG. 4. The bars in the figure represent the classification performance (AUROC values, y-axis) using SLF (blue), eGeMAPS (red), and ComPARE (green) as feature sets, evaluated across all speech types from the PC-GITA dataset (x-axis).

357 not aligned on linguistic markers (syllable nuclei), but rather were randomly drawn from  
 358 respective syllable duration distribution (see Methods for more details).

359 The classification procedure was identical to that used for SLF (e.g. in Figure 4), and the  
 360 results for each speech type are presented in Figure 5. For consistency, we again aggregated  
 361 scores from all segments into utterance-level scores; classification results are shown in Figure  
 362 5. Our (null) hypothesis is that syllable-based features are more advantageous and lead to  
 363 better discrimination of Parkinsonian speech than random segment-based features. To test  
 364 this hypothesis, we performed a Wilcoxon signed-rank test (restricted to speech types with

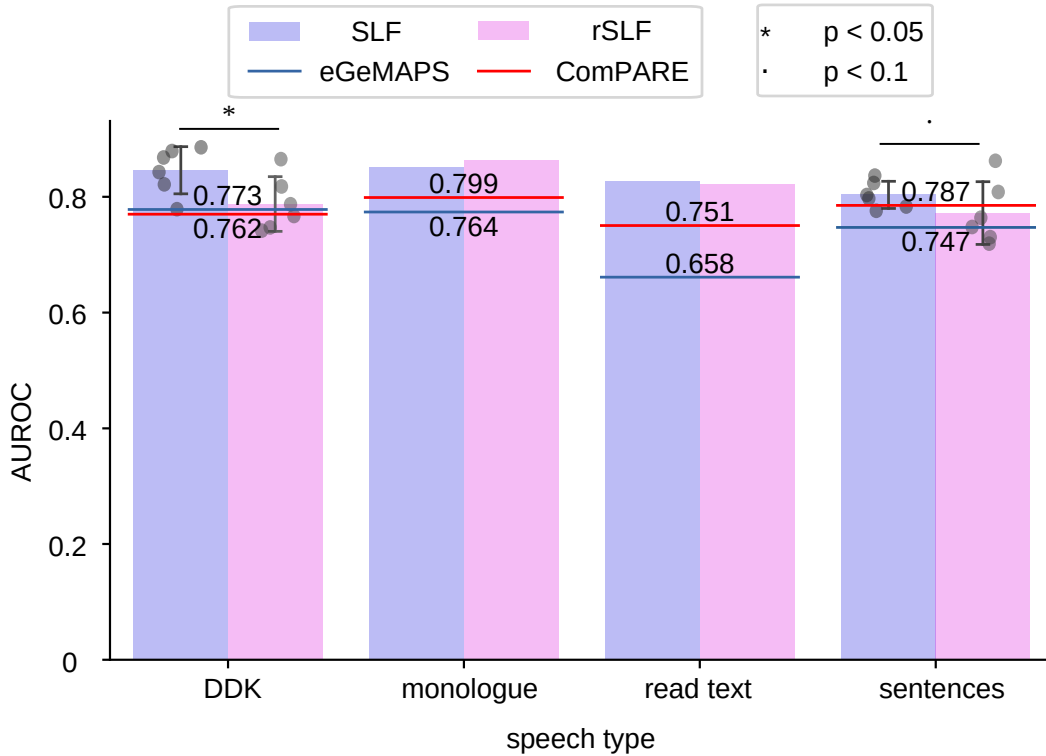


FIG. 5. The figure represents the classification results (measured as AUROC) with syllable based features (blue) versus random syllable-like segment (magenta).

365 multiple utterances; DDK and sentences). The difference was especially prominent for artic-  
 366 ulatorily constrained speech types, specifically, DDK speech (Wilcoxon  $W = 21$ ,  $p = 0.016$ ,  
 367 rank biserial correlation = 1.0), and for shorter “free form” speech, sentences (Wilcoxon  
 368  $W = 18$ ,  $p = 0.078$ , rank biserial correlation = 0.714). Overall, these results indicate that  
 369 aligning segments with linguistic markers, such as syllable nuclei, is an important factor in  
 370 the proposed feature extraction procedure.

### 371 C. Interpretability

372 To elucidate the information captured by SLF that enables effective classification of PD  
373 speech relative to HC speech, we investigated whether standardized syllable spectrograms  
374 encode formant patterns. Previous studies have reported alterations in formant patterns in  
375 the speech of PD patients (Convey *et al.*, 2023; Liu *et al.*, 2021, 2023). We therefore hypothe-  
376 sized that these patterns contribute to classification performance and sought to demonstrate  
377 their presence in SLF representations. To test this, we modified the construction of pre-  
378 viously reported SLFs by replacing constant frequency binning with a targeted analysis of  
379 average energy around formant frequencies and the fundamental frequency. To further val-  
380 idate our findings, we conducted a statistical analysis of the eGeMAPS feature set across  
381 conditions.

382 Figure 6 demonstrates that the formant-based SLF, despite its lower dimensionality ( $5 \times 8$   
383 vs.  $64 \times 8$ ), achieves performance well above chance and its results are comparable to the  
384 classification results coming from established feature sets. More specifically, for DDK utter-  
385 ances, the average AUROC was 84.6% for SLF, 76.2% for ComPARE, 77.2% for eGeMAPS,  
386 and 75.75% for the formant-based SLF. This trend holds across other speech types: formant-  
387 based SLF performance is 5–10% lower than SLF but remains comparable to eGeMAPS and  
388 ComPARE.

389 These findings indicate that the standard SLF encodes formant and F0 information while  
390 also capturing supplementary features that enhance classification performance.

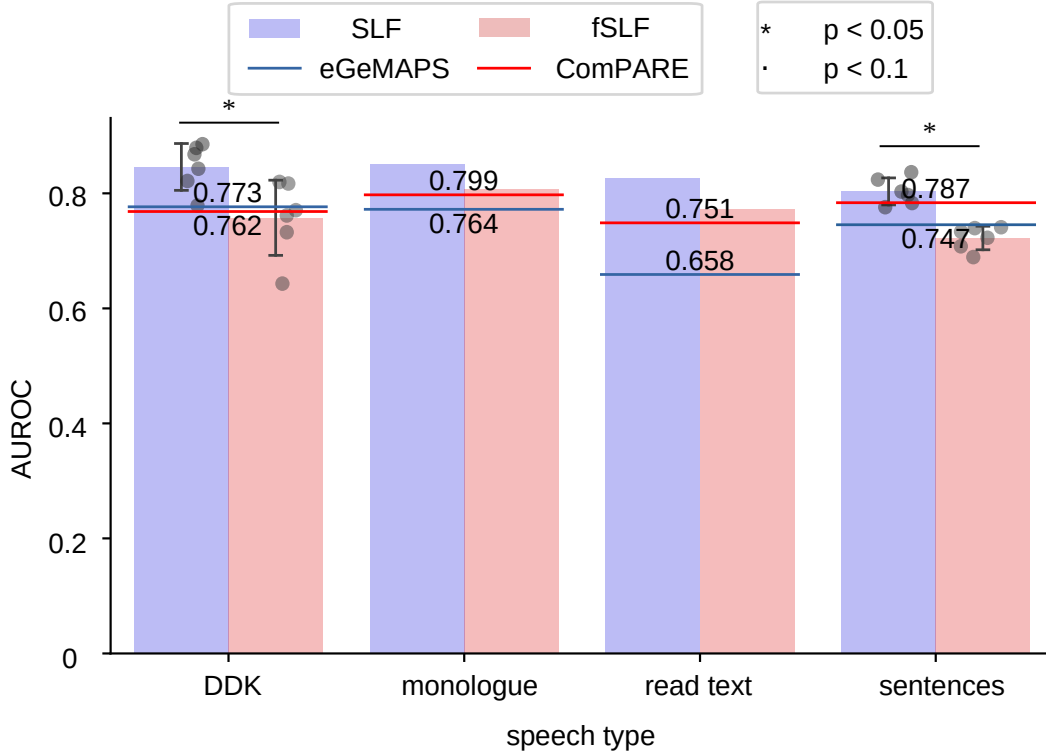


FIG. 6. Classification results for SLF (blue) and fSLF (red). The Wilcoxon signed-rank test indicates that SLF yields significantly better discriminatory power than fSLF for the DDK (Wilcoxon  $W = 21$ ,  $p=0.031$ , rank biserial correlation=1) and sentence (Wilcoxon  $W = 21$ ,  $p=0.031$ , rank biserial correlation=1) speech types.

391 To further validate our hypothesis, we conducted a linear mixed-model analysis of the  
 392 eGeMAPS feature set. We selected eGeMAPS over ComPARE due to its curated set of  
 393 88 features, which directly relate to both the source and system related components of  
 394 speech production. This feature set offers greater interpretability and a clearer link to the  
 395 physiological mechanisms underlying speech (Dubagunta *et al.*, 2022).

396 In our analysis, we modeled fixed effects for eGeMAPS features and PD vs. HC condition,  
 397 with speaker included as a random effect. We then examined the interaction between features

feature group	DDK monologue read text sentences			total
<b>F0</b>	20	3	6	29
<b>F1</b>	6	1	6	13
<b>F2</b>	6	2	6	14
<b>F3</b>	10	1	2	13
<b>mffc</b>			1	3
<b>other</b>	1			1

TABLE I. The number of significant interactions between eGeMAPS features and PD vs HC conditions for each speech type in PC-GITA dataset.

398 and conditions, identifying statistically significant differences ( $p < 0.05$ ) after correcting  
399 for multiple comparisons using the false discovery rate. The results for each speech and  
400 utterance type are presented in Table I.

401 Our findings reveal statistically significant differences across conditions for nearly all  
402 speech types, particularly in features related to formants and/or fundamental frequency. In  
403 some cases, MFCC-related features also showed significance, for more details, such as exact  
404 feature name for each speech type and utterance, as well as respective effect size and exact  
405 p-value, please see the Supplementary table 2a-d. These results align with our formant-based  
406 SLF analyses, reinforcing the conclusion that SLF effectively captures critical information  
407 about source and system of speech production, such as formant patterns and fundamental  
408 frequency at the syllable level.

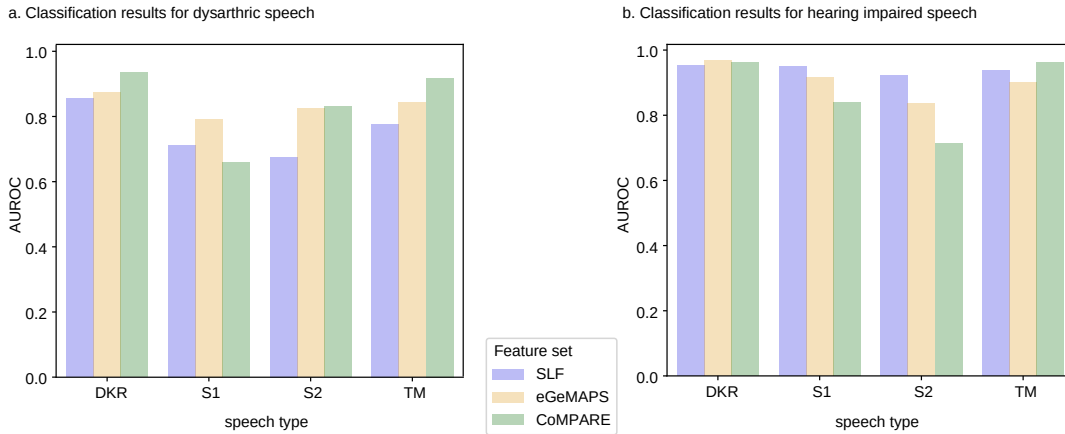


FIG. 7. Classification results on COPAS dataset

#### D. Speech Pathology detection / Generalizability beyond PD

To evaluate the generalizability of the SLF approach, we applied it to the COPAS dataset (Martens *et al.*, 2011), a well-established resource for pathological speech research. As described in the Methods section, this dataset comprises Dutch recordings from participants with various speech pathologies and includes multiple speech types. For this analysis, we focused exclusively on recordings from individuals with hearing impairment and dysarthria, as these were the pathologies with sufficient amount of recordings for robust classification analysis.

Classification results, obtained via K-fold cross-validation, are presented in Figure 7, with more detailed performance metrics presented in Supplementary Table 3a and 3b for dysarthric and hearing impairment cases, respectively. Due to differences in sample size, we used K=4 folds for hearing loss and K=5 folds for dysarthria. The results demonstrate that SLF achieves consistently high performance across all speech types for hearing loss, either

422 outperforming or matching the control feature sets. For dysarthria, however, performance  
423 was more modest, with control feature sets outperforming SLF across all speech types.

## 424 V. DISCUSSION

425 In this manuscript, we build upon our earlier proposal of a novel feature extraction  
426 method for speech pathology detection (Hovsepyan and Magimai.-Doss, 2024), drawing in-  
427 spiration from neurocomputational models of speech perception (Hovsepyan *et al.*, 2020;  
428 Su *et al.*, 2023; Yildiz *et al.*, 2013). We conduct extensive testing of the new feature set  
429 across diverse speech types, expanding our initial analysis, previously limited to DDK tasks,  
430 to include longer speech samples such as sentences, monologues, and read text passages.  
431 Furthermore, we evaluate the proposed SLF method using Dutch dataset for pathological  
432 speech (Martens *et al.*, 2011), which includes recordings from individuals with pathologies  
433 such as hearing impairment and dysarthria. Our findings demonstrate that the SLF method  
434 is not constrained by speech type, language, or dataset, indicating its potential for broader  
435 generalizability.

436 Furthermore, to assess the role of syllable boundary segmentation, we compared linguis-  
437 tically informed segments with randomly segmented intervals of equivalent duration. The  
438 results reveal that linguistically derived segments yield superior performance, underscoring  
439 the critical role of syllable-like units in this approach.

440 Finally, to explore the interpretability of our results, we investigated whether SLF cap-  
441 tures PD-related disruptions in vowel and formant articulation. Accordingly, we adapted  
442 the SLF creation process to capture dynamic frequency patterns, specifically those related

443 to formants and pitch, rather than relying solely on average spectral energy across uniformly  
444 distributed channels. To validate these findings, we employed classifier-independent statis-  
445 tical analyses to identify features with significant differences between PD and HC groups.  
446 These analyses corroborate that formant- and pitch-related information is most prominently  
447 affected in PD compared to HC conditions.

#### 448 **A. Importance of syllable segments**

449 A critical aspect of our methodology is the segmentation of continuous speech into  
450 syllable-like units, a choice motivated by analysis-by-synthesis theories of speech percep-  
451 tion and production (Bever and Poeppel, 2010; Halle and Stevens, 1962; Liberman and  
452 Mattingly, 1985). These theories posit that the brain decodes speech by generating and  
453 comparing internal articulatory gestures to match the incoming speech signal. Impairments  
454 in speech production may therefore manifest as detectable traces in perceptual analysis.  
455 We designed the SLF extraction process to mirror generative models of speech perception,  
456 which recognize syllables by iteratively comparing predicted and actual spectrograms. This  
457 approach has already demonstrated its ability to distinguish between PD patients and HC,  
458 suggesting that syllable-based representations can capture condition-specific differences.

459 A key question is whether syllable-based segmentation, or any segmentation with similar  
460 durations, is inherently advantageous. To test this, we first derived the syllable duration  
461 distribution for each speech type in the PC-GITA dataset (Orozco *et al.*, 2014). We then  
462 recalculated SLF using randomly sampled segment durations from these distributions.

463 Results revealed a significant performance drop for nearly all speech and utterance types.  
464 On average, performance declined by 7% for DDK, 15% for monologues, and 3% for sen-  
465 tences, with minimal change for read text. The latter two cases were somewhat unexpected:  
466 most probably the classification results for sentences were influenced by an outlier (micasa).  
467 By contrast, the results for the read text task can likely be explained by the fact that read  
468 text is more distinctly articulated (compared to monologues) and far longer (compared to  
469 sentences and DDK tasks). As a result, even randomly segmented chunks may closely re-  
470 semble syllabic segments and, due to their longer durations, accumulate sufficient segmental  
471 information to yield performance comparable to syllable-based segmentation at the utter-  
472 ance level. Importantly, all utterance types except the outlier exhibited some performance  
473 drop, confirming the importance of syllable-based segmentation in our method.

474 Interestingly, random segmentation results closely matched those of control feature sets  
475 such as ComPARE and eGeMAPS, sometimes even surpassing them. One possible explana-  
476 tion for this similarity could arise because both control features and random SLF disregard  
477 linguistic boundaries, relying instead on fixed/small sliding windows that blur these bound-  
478 aries. Random segment-based features thus function similarly to sliding windows, with  
479 classifiers identifying patterns across segments. The proposed approach, even with random  
480 segments, occasionally outperforms control features by providing multi-point evidence per  
481 utterance, in contrast to the single-point evidence offered by control sets. This suggests that  
482 pathological speech differences may not be uniformly distributed across utterance duration  
483 but are instead concentrated in specific segments. By dividing speech into multiple parts and  
484 analyzing each segment, we capture the most prominent differences, thereby enhancing clas-

485 sification performance. In contrast, the sliding window technique employed by eGeMAPS  
486 and ComPARE averages frame-level features across the entire utterance, tending to dilute  
487 these differences. This aligns with another key finding: syllable-level classification yielded  
488 lower AUC values than utterance-level classification, which aggregates syllable-level scores.

489 Lastly, it is well established that PD effects on speech is also manifested with disrupted  
490 speech rhythmicity, which is reflected in the syllable duration and rate of pathological speech.  
491 While the SLF method does not directly encode syllable duration, durations are normalized  
492 into fixed temporal bins, it may still indirectly capture temporal irregularities in PD speech.  
493 For example, nucleus-to-nucleus segmentation, as used in SLF extraction in this study,  
494 inherently incorporates the increased intersyllabic or inter-word intervals characteristic of  
495 PD. When spectral energy is averaged across fixed bins, the inclusion of "silent" segments  
496 reduces overall energy levels, particularly in free-form speech such as monologues, potentially  
497 improving classification performance.

498 Though effective, this approach highlights opportunities for refinement. Introducing dy-  
499 namic sampling rates, using unequally spaced temporal bins to target specific syllable seg-  
500 ments, could enhance temporal resolution. Such an approach, which has proven successful  
501 in syllable recognition, may further improve the SLF method's sensitivity and utility in  
502 detecting speech pathologies.

## 503 **B. Interpretability**

504 To better understand the information captured by the SLF method, we pursued two com-  
505plementary approaches: first, we tested the hypothesis that SLF encodes formant and pitch-

506 related information; second, we performed a classifier-independent analysis of the eGeMAPS  
507 feature set to identify systematically differing features across conditions.

508 This hypothesis is grounded in research demonstrating that PD disrupts both vowel  
509 articulation and pitch control. These variations are particularly evident in sustained vowel  
510 tasks, where maintaining a stable vocal tract configuration over time is challenging for  
511 individuals with motor control impairments, such as those with PD.

512 To test this, we transformed the spectrogram representation of syllables by averaging  
513 spectral energy across formant- and pitch-related channels, rather than uniformly distributed  
514 frequency bins. This reduced the SLF feature vector from 64 to just 5 channels—over  
515 a tenfold reduction—while maintaining classification performance comparable to control  
516 feature sets. This result suggests that formant and pitch information is indeed central to  
517 SLF’s effectiveness.

518 We further validated this by statistically comparing eGeMAPS features across conditions.  
519 Formant-related features, fundamental frequency, and select MFCCs consistently emerged  
520 as significant discriminators across all speech types.

521 While these findings confirm that SLF retains formant and pitch information through  
522 average spectral energy, they also reveal that SLF captures additional, richer details. Per-  
523 formance declined when using only formant- and pitch-related channels, indicating that  
524 SLF’s strength lies in its ability to encode consonant articulation—information not directly  
525 measured by eGeMAPS or ComPARE. Since formants are tied to vowels and our segmenta-  
526 tion focuses on syllable nuclei, SLF uniquely captures transitions between vowels, including

527 the execution of intervening consonants. This broader scope explains why the full SLF  
528 method outperforms feature sets limited to vowel- and pitch-related metrics.

### 529 **C. Generalizability of the SLF Method**

530 In this study, we expanded our neurocomputationally inspired method for speech pathol-  
531 ogy detection, which uses syllable spectrogram representations. This approach is both simple  
532 and effective, achieving classification accuracy comparable to established feature sets. We  
533 tested its generalizability by applying it to PD detection across diverse speech types and  
534 datasets, including the COPAS dataset in Dutch, which features hearing impairment and  
535 dysarthria.

536 The COPAS dataset allowed us to evaluate SLF on two distinct pathologies. For  
537 dysarthria, performance varied: SLF matched control feature sets in structured tasks like  
538 diadochokinetic utterances (DKR) but lagged in longer speech types. This discrepancy  
539 may stem from the diverse etiologies of dysarthria in COPAS, which can disrupt syllable  
540 boundary detection and reduce SLF’s effectiveness compared to PD. In contrast, hearing  
541 impairment, a pathology rooted in perceptual deficits, aligned well with SLF’s theoretic-  
542 al foundation. Analysis-by-synthesis theories suggest that production artifacts in hearing  
543 impairment arise from impaired efference copy, a feedback mechanism critical for motor  
544 control. SLF achieved over 90% AUROC across all speech types for hearing impairment,  
545 often surpassing control features and demonstrating remarkable consistency.

546 These results highlight SLF’s robustness in capturing production deficits linked to per-  
547 ceptual mechanisms. The method’s ability to generalize across languages, speech types, and  
548 pathologies underscores its potential as a versatile tool for speech pathology detection.

## 549 VI. CONCLUSION AND OUTLOOK

550 In this manuscript, we evaluated the use of syllable-level features (SLF) for detecting  
551 speech pathologies. Inspired by neurophysiologically plausible generative models of speech  
552 perception, SLF is grounded in the analysis-by-synthesis theory, which posits that the brain  
553 internally generates speech during perception. Our hypothesis was that production deficits  
554 due to pathology would manifest in perceptual representations, enabling detection from the  
555 perception side: a departure from traditional motor-focused approaches.

556 Syllables were chosen as the feature extraction unit because they are fundamental to  
557 speech production and perception, with distinct acoustic patterns and cognitive relevance.  
558 By extracting features from linguistically motivated segments, SLF captures linguistic,  
559 acoustic, and cognitive aspects of speech. Our results confirm that this approach reli-  
560 ably detects pathologies across datasets, speech types, and languages, with particularly  
561 strong performance for PD and hearing impairment. While SLF performed respectably for  
562 dysarthria, further refinement may be needed for pathologies with complex or variable effects  
563 on speech.

564 The SLF method is inherently flexible. It is, in principle, agnostic to specific segmen-  
565 tation or spectrogram calculation methods, allowing for optimization tailored to different  
566 pathologies. For example, dynamic frequency binning, focusing on formant and F0 patterns

567 affected by PD, could enhance targeted feature construction. Similarly, dynamic tempo-  
568 ral sampling could improve resolution in regions critical for consonant articulation, further  
569 reducing feature dimensionality.

570 Future directions include exploring dynamic temporal and frequency sampling to focus  
571 on pathology-specific regions, as well as adapting SLF for long-term monitoring in noisy  
572 environments or on-device applications.

573 Overall, SLF demonstrates that linguistically driven feature extraction can advance  
574 speech pathology detection, offering a biologically plausible, low-profile alternative to tradi-  
575 tional methods.

## 576 **ACKNOWLEDGMENTS**

## 577 **VII. REFERENCES**

- 578
- 579 Ananthanarayanan, A., Senivarapu, S., and Murari, A. (2025). “Towards Causal Inter-  
580 pretability in Deep Learning for Parkinson’s Detection from Voice Data” doi: [10.1101/  
581 2025.04.25.25326311](https://doi.org/10.1101/2025.04.25.25326311).
- 582 Bever, T. G., and Poeppel, D. (2010). “Analysis by Synthesis: A (Re-)Emerging Program  
583 of Research for Language and Vision,” *Biolinguistics* 27, doi: [10.5964/bioling.8783](https://doi.org/10.5964/bioling.8783).
- 584 Boersma, P., and Weenink, D. (2003). “Praat: doing phonetics by computer,” [https:  
585 //api.semanticscholar.org/CorpusID:60594797](https://api.semanticscholar.org/CorpusID:60594797).

COPAS				
Condition	Speech Type	Subjects	M (n, age $\pm$ std)	F (n, age $\pm$ std)
control	DKR	83	34 (40.50 $\pm$ 16.37)	49 (33.27 $\pm$ 17.05)
	S1	82	33 (41.09 $\pm$ 16.25)	49 (33.27 $\pm$ 17.05)
	S2	82	33 (41.12 $\pm$ 16.21)	49 (33.27 $\pm$ 17.05)
	TM	82	34 (40.50 $\pm$ 16.37)	48 (32.96 $\pm$ 17.09)
dysarthria	DKR	55	32 (44.56 $\pm$ 21.78)	23 (44.52 $\pm$ 25.05)
	S1*	50	28 (43.46 $\pm$ 22.62)	21 (44.76 $\pm$ 24.94)
	S2*	50	28 (43.57 $\pm$ 23.11)	21 (44.76 $\pm$ 24.94)
	TM*	48	28 (43.46 $\pm$ 22.43)	19 (41.79 $\pm$ 26.09)
hearing loss	DKR	24	9 (31.33 $\pm$ 7.97)	15 (35.93 $\pm$ 15.36)
	S1	24	9 (31.33 $\pm$ 7.97)	15 (35.93 $\pm$ 15.36)
	S2	24	9 (31.33 $\pm$ 7.97)	15 (35.93 $\pm$ 15.36)
	TM	26	9 (31.33 $\pm$ 7.97)	17 (38.29 $\pm$ 15.98)
PC-GITA				
Healthy	DDK, sentences,	50	25 (61.2 $\pm$ 11.3)	25(60.7 $\pm$ 7.7)
Control	read text, monologue			
Parkinson's	DDK, sentences,	50	25 (62.2 $\pm$ 11.2)	25(60.1 $\pm$ 7.8)
Disease	Read text, monologue			

TABLE II. The table shows the number of speakers, their mean age ( $\pm$  standard deviation) by gender for each pathology and speech type in the COPAS and PC-GITA datasets. An asterisk indicates that, for a particular speech type, there was a speaker with no gender information.

586 Buckley, C. L., Kim, C. S., McGregor, S., and Seth, A. K. (2017). “The free energy principle  
587 for action and perception: A mathematical review,” *Journal of Mathematical Psychology*  
588 **81**, 55–79, doi: [10/gcnkg7](https://doi.org/10/gcnkg7).

589 Convey, R. B., Ihalainen, T., Liu, Y., Räsänen, O., Ylinen, S., and Penttilä, N. (2023).  
590 “A comparative study of automatic vowel articulation index and auditory-perceptual as-  
591 sessments of speech intelligibility in Parkinson’s disease,” *International Journal of Speech-*  
592 *Language Pathology* 1–11, doi: [10.1080/17549507.2023.2251725](https://doi.org/10.1080/17549507.2023.2251725).

593 Dogonasheva, O., Doelling, K. B., Zakharov, D., Giraud, A.-L., and Gutkin, B. (2025).  
594 “Rhythm-based hierarchical predictive computations support acoustic-semantic trans-  
595 formation in speech processing,” *Nature Computational Science* 5(10), 915–926, doi:  
596 [10.1038/s43588-025-00876-9](https://doi.org/10.1038/s43588-025-00876-9).

597 Dubagunta, S. P., Magimai.-Doss, M., Eleni Theocharopoulos, and Mathew Magimai Doss  
598 (2022). “Towards Automatic Prediction of Non-Expert Perceived Speech Fluency Rat-  
599 ings,” *ICMI Companion* doi: [10.1145/3536220.3563689](https://doi.org/10.1145/3536220.3563689).

600 Escobar-Grisales, D., Ríos-Urrego, C. D., and Orozco-Arroyave, J. R. (2023). “Deep Learn-  
601 ing and Artificial Intelligence Applied to Model Speech and Language in Parkinson’s Dis-  
602 ease,” *Diagnostics* 13(13), 2163, doi: [10.3390/diagnostics13132163](https://doi.org/10.3390/diagnostics13132163).

603 Eyben, F., Scherer, K. R., Schuller, B. W., Sundberg, J., Andre, E., Busso, C., Dev-  
604 illers, L. Y., Epps, J., Laukka, P., Narayanan, S. S., and Truong, K. P. (2016). “The  
605 Geneva Minimalistic Acoustic Parameter Set (GeMAPS) for Voice Research and Af-  
606 fective Computing,” *IEEE Transactions on Affective Computing* 7(2), 190–202, doi:  
607 [10.1109/TAFFC.2015.2457417](https://doi.org/10.1109/TAFFC.2015.2457417).

608 Eyben, F., Wöllmer, M., and Schuller, B. (2010). “Opensmile: The munich versatile and  
609 fast open-source audio feature extractor,” in *Proceedings of the 18th ACM International*  
610 *Conference on Multimedia, MM ’10*, Association for Computing Machinery, New York,

611 NY, USA, pp. 1459–1462, doi: [10.1145/1873951.1874246](https://doi.org/10.1145/1873951.1874246).

612 Friston, K. (2010). “The free-energy principle: A unified brain theory?,” *Nature Reviews*  
613 *Neuroscience* **11**(2), 127–138, doi: [10/bj3r7f](https://doi.org/10/bj3r7f).

614 Friston, K., and Kiebel, S. (2009). “Predictive coding under the free-energy principle,”  
615 *Philosophical Transactions of the Royal Society B: Biological Sciences* **364**(1521), 1211–  
616 1221, doi: [10/fcswfc](https://doi.org/10/fcswfc).

617 Goetz, C. G., Tilley, B. C., Shaftman, S. R., Stebbins, G. T., Fahn, S., Martinez-Martin, P.,  
618 Poewe, W., Sampaio, C., Stern, M. B., Dodel, R., Dubois, B., Holloway, R., Jankovic, J.,  
619 Kulisevsky, J., Lang, A. E., Lees, A., Leurgans, S., LeWitt, P. A., Nyenhuis, D., Olanow,  
620 C. W., Rascol, O., Schrag, A., Teresi, J. A., Van Hilten, J. J., and LaPelle, N. (2008).  
621 “Movement Disorder Society-sponsored revision of the Unified Parkinson’s Disease Rat-  
622 ing Scale (MDS-UPDRS): Scale presentation and clinimetric testing results,” *Movement*  
623 *Disorders* **23**(15), 2129–2170, doi: [10.1002/mds.22340](https://doi.org/10.1002/mds.22340).

624 Halle, M., and Stevens, K. (1962). “Speech recognition: A model and a program for re-  
625 search,” *IRE Transactions on Information Theory* **8**(2), 155–159, doi: [10.1109/TIT.1962.](https://doi.org/10.1109/TIT.1962.1057686)  
626 [1057686](https://doi.org/1057686).

627 Hoehn, M. M., and Yahr, M. D. (1967). “Parkinsonism: Onset, progression and mortality,”  
628 *Neurology* **17**(5), 427–442, doi: [10.1212/wnl.17.5.427](https://doi.org/10.1212/wnl.17.5.427).

629 Hovsepyan, S., and Magimai.-Doss, M. (2024). “Syllable Level Features for Parkinson’s  
630 Disease Detection from Speech,” in *ICASSP 2024 - 2024 IEEE International Conference*  
631 *on Acoustics, Speech and Signal Processing (ICASSP)*, IEEE, Seoul, Korea, Republic of,  
632 pp. 11416–11420, doi: [10.1109/ICASSP48485.2024.10446484](https://doi.org/10.1109/ICASSP48485.2024.10446484).

633 Hovsepyan, S., Olasagasti, I., and Giraud, A.-L. (2020). “Combining predictive coding and  
634 neural oscillations enables online syllable recognition in natural speech,” *Nature Commu-  
635 nications* **11**(1), 3117, doi: [10/gg2n6w](https://doi.org/10/gg2n6w).

636 Hovsepyan, S., Olasagasti, I., and Giraud, A.-L. (2023). “Rhythmic modulation of pre-  
637 diction errors: A top-down gating role for the beta-range in speech processing,” *PLOS  
638 Computational Biology* **19**(11), e1011595, doi: [10.1371/journal.pcbi.1011595](https://doi.org/10.1371/journal.pcbi.1011595).

639 Hyafil, A., Fontolan, L., Kabdebon, C., Gutkin, B., and Giraud, A. L. (2015). “Speech  
640 encoding by coupled cortical theta and gamma oscillations,” *eLife* **4**(MAY), 1–45, doi:  
641 [10/gdh97p](https://doi.org/10/gdh97p).

642 Jadoul, Y., Thompson, B., and De Boer, B. (2018). “Introducing Parselmouth: A Python  
643 interface to Praat,” *Journal of Phonetics* **71**, 1–15, doi: [10.1016/j.wocn.2018.07.001](https://doi.org/10.1016/j.wocn.2018.07.001).

644 Karlsson, F., Schalling, E., Laakso, K., Johansson, K., and Hartelius, L. (2020). “Assess-  
645 ment of speech impairment in patients with Parkinson’s disease from acoustic quantifica-  
646 tions of oral diadochokinetic sequences,” *The Journal of the Acoustical Society of America*  
647 **147**(2), 839–851, doi: [10.1121/10.0000581](https://doi.org/10.1121/10.0000581).

648 Liberman, A. M., and Mattingly, I. G. (1985). “The motor theory of speech perception  
649 revised,” *Cognition* **21**(1), 1–36, doi: [10/d79rrw](https://doi.org/10/d79rrw).

650 Liu, Y., Penttilä, N., Ihalainen, T., Lintula, J., Convey, R., and Räsänen, O. (2021).  
651 “Language-Independent Approach for Automatic Computation of Vowel Articulation Fea-  
652 tures in Dysarthric Speech Assessment,” *IEEE/ACM Transactions on Audio, Speech, and  
653 Language Processing* **29**, 2228–2243, doi: [10.1109/TASLP.2021.3090973](https://doi.org/10.1109/TASLP.2021.3090973).

654 Liu, Y., Reddy, M. K., Penttilä, N., Ihalainen, T., Alku, P., and Räsänen, O. (2023).  
655 “Automatic Assessment of Parkinson’s Disease Using Speech Representations of Phonation  
656 and Articulation,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*  
657 **31**, 242–255, doi: [10.1109/TASLP.2022.3212829](https://doi.org/10.1109/TASLP.2022.3212829).

658 Martens, J., Bodt, M., Nuffelen, G. V., and Middag, C. (2011). “Corpus of Pathological  
659 and Normal Speech (COPAS),” .

660 Mazzoni, P., Shabbott, B., and Cortés, J. C. (2012). “Motor Control Abnormalities in  
661 Parkinson’s Disease,” *Cold Spring Harbor Perspectives in Medicine* **2**(6), a009282, doi:  
662 [10.1101/cshperspect.a009282](https://doi.org/10.1101/cshperspect.a009282).

663 McFee, B., Raffel, C., Liang, D., Ellis, D., McVicar, M., Battenberg, E., and Nieto, O.  
664 (2015). “Librosa: Audio and Music Signal Analysis in Python,” in *Python in Science*  
665 *Conference*, Austin, Texas, pp. 18–24, doi: [10.25080/Majora-7b98e3ed-003](https://doi.org/10.25080/Majora-7b98e3ed-003).

666 Michaelis, D., Fröhlich, M., and Strube, H. W. (1998). “Selection and combination of acous-  
667 tic features for the description of pathologic voices,” *The Journal of the Acoustical Society*  
668 *of America* **103**(3), 1628–1639, doi: [10.1121/1.421305](https://doi.org/10.1121/1.421305).

669 Nabé, M., Schwartz, J.-L., and Diard, J. (2021). “COSMO-Onset: A Neurally-Inspired  
670 Computational Model of Spoken Word Recognition, Combining Top-Down Prediction and  
671 Bottom-Up Detection of Syllabic Onsets,” *Frontiers in Systems Neuroscience* **15**, 75, doi:  
672 [10.3389/fnsys.2021.653975](https://doi.org/10.3389/fnsys.2021.653975).

673 Nair, V., Susskind, J., and Hinton, G. E. (2008). “Analysis-by-Synthesis by Learning to  
674 Invert Generative Black Boxes,” in *Artificial Neural Networks - ICANN 2008*, edited by  
675 V. Kůrková, R. Neruda, and J. Koutník, **5163** (Springer Berlin Heidelberg, Berlin, Hei-

676 delberg), pp. 971–981, doi: [10.1007/978-3-540-87536-9\\_99](https://doi.org/10.1007/978-3-540-87536-9_99).

677 Orozco, J. R., Arias-Londoño, J. D., Vargas-Bonilla, J., González-Rátiva, M., and Noeth,  
678 E. (2014). “New Spanish speech corpus database for the analysis of people suffering from  
679 Parkinson’s disease,” in *International Conference on Language Resources and Evaluation*  
680 (*LREC*).

681 Räsänen, O., Doyle, G., and Frank, M. C. (2018). “Pre-linguistic segmentation of speech into  
682 syllable-like units,” *Cognition* **171**, 130–150, doi: [10.1016/j.cognition.2017.11.003](https://doi.org/10.1016/j.cognition.2017.11.003).

683 Schuller, B., Steidl, S., Batliner, A., Burkhardt, F., Devillers, L., Müller, C., and Narayanan,  
684 S. S. (2010). “The INTERSPEECH 2010 paralinguistic challenge,” in *Interspeech 2010*,  
685 ISCA, pp. 2794–2797, doi: [10.21437/Interspeech.2010-739](https://doi.org/10.21437/Interspeech.2010-739).

686 Schuller, B., Steidl, S., Batliner, A., Hantke, S., Hönl, F., Orozco-Arroyave, J. R., Nöth,  
687 E., Zhang, Y., and Weninger, F. (2015). “The INTERSPEECH 2015 computational par-  
688 alinguistics challenge: Nativeness, Parkinson’s & eating condition,” in *Interspeech 2015*,  
689 ISCA, pp. 478–482, doi: [10.21437/Interspeech.2015-179](https://doi.org/10.21437/Interspeech.2015-179).

690 Schuller, B., Steidl, S., Batliner, A., Hirschberg, J., Burgoon, J. K., Baird, A., Elkins, A.,  
691 Zhang, Y., Coutinho, E., and Evanini, K. (2016). “The INTERSPEECH 2016 Computa-  
692 tional Paralinguistics Challenge: Deception, Sincerity & Native Language,” in *Interspeech*  
693 *2016*, ISCA, pp. 2001–2005, doi: [10.21437/Interspeech.2016-129](https://doi.org/10.21437/Interspeech.2016-129).

694 Schuller, B., Steidl, S., Batliner, A., Vinciarelli, A., Scherer, K., Ringeval, F., Chetouani,  
695 M., Weninger, F., Eyben, F., Marchi, E., Mortillaro, M., Salamin, H., Polychroniou, A.,  
696 Valente, F., and Kim, S. (2013). “The INTERSPEECH 2013 computational paralinguis-  
697 tics challenge: Social signals, conflict, emotion, autism,” in *Proceedings of the Annual*

698 *Conference of the International Speech Communication Association, INTERSPEECH*, pp.  
699 148–152, doi: [10.21437/Interspeech.2013-56](https://doi.org/10.21437/Interspeech.2013-56).

700 Schuller, B. W., Batliner, A., Bergler, C., Mascolo, C., Han, J., Lefter, I., Kaya, H.,  
701 Amiriparian, S., Baird, A., Stappen, L., Ottl, S., Gerczuk, M., Tzirakis, P., Brown,  
702 C., Chauhan, J., Grammenos, A., Hasthanasombat, A., Spathis, D., Xia, T., Cicuta,  
703 P., Rothkrantz, L. J. M., Zwerts, J. A., Treep, J., and Kaandorp, C. S. (2021).  
704 “The INTERSPEECH 2021 Computational Paralinguistics Challenge: COVID-19 Cough,  
705 COVID-19 Speech, Escalation & Primates,” in *Proc. Interspeech 2021*, pp. 431–435, doi:  
706 [10.21437/Interspeech.2021-19](https://doi.org/10.21437/Interspeech.2021-19).

707 Su, Y., MacGregor, L. J., Olasagasti, I., and Giraud, A.-L. (2023). “A deep hierarchy of  
708 predictions enables online meaning extraction in a computational model of human speech  
709 comprehension,” *PLOS Biology* **21**(3), e3002046, doi: [10.1371/journal.pbio.3002046](https://doi.org/10.1371/journal.pbio.3002046).

710 van Gelderen, L., and Tejedor-García, C. (2024). “Innovative Speech-Based Deep Learn-  
711 ing Approaches for Parkinson’s Disease Classification: A Systematic Review,” *Applied*  
712 *Sciences* **14**(17), 7873, doi: [10.3390/app14177873](https://doi.org/10.3390/app14177873).

713 Yildiz, I. B., von Kriegstein, K., and Kiebel, S. J. (2013). “From Birdsong to Human Speech  
714 Recognition: Bayesian Inference on a Hierarchy of Nonlinear Dynamical Systems,” *PLoS*  
715 *Computational Biology* **9**(9), e1003219–e1003219, doi: [10/gf26tp](https://doi.org/10/gf26tp).