

1

2

3

4 Identifying bias in models that detect vocal fold paralysis from audio
5 recordings using explainable machine learning and clinician ratings

6 Daniel M. Low^{1,2}, Vishwanatha Rao^{3,4}, Gregory Randolph^{4,5}, Phillip C. Song^{4,5}*,

7 Satrajit S. Ghosh^{1,2,5}*

8

9

10

11 ¹ Program in Speech and Hearing Bioscience and Technology, Harvard Medical School, Boston,
12 MA, USA

13 ² McGovern Institute for Brain Research, MIT, Cambridge, MA, USA

14 ³ Department of Biomedical Engineering, Columbia University, New York, NY, USA

15 ⁴ Department of Otolaryngology–Head and Neck Surgery, Massachusetts Eye and Ear
16 Infirmary, Boston, MA, USA

17 ⁵ Department of Otolaryngology–Head and Neck Surgery, Harvard Medical School, Boston, MA,
18 USA

19 * Equal contribution

20 **Corresponding author**

21 Correspondence can be addressed to Daniel M. Low, Office: 46-4033F, 43 Vassar St,

22 Cambridge, MA 02139, USA. E-mail: dlow@mit.edu.

23

24

25

26

Abstract

27 **Introduction.** Detecting voice disorders from voice recordings could allow for frequent, remote,
28 and low-cost screening before costly clinical visits and a more invasive laryngoscopy
29 examination. Our goals were to detect unilateral vocal fold paralysis (UVFP) from voice
30 recordings using machine learning, to identify which acoustic variables were important for
31 prediction to increase trust, and to determine model performance relative to clinician
32 performance.

33 **Methods.** Patients with confirmed UVFP through endoscopic examination (N=77) and controls
34 with normal voices matched for age and sex (N=77) were included. Voice samples were elicited
35 by reading the Rainbow Passage and sustaining phonation of the vowel "a". Four machine
36 learning models of differing complexity were used. SHapley Additive exPlanations (SHAP) was
37 used to identify important features.

38 **Results.** The highest median bootstrapped ROC AUC score was 0.87 and beat clinician's
39 performance (range: 0.74 – 0.81) based on the recordings. Recording durations were different
40 between UVFP recordings and controls due to how that data was originally processed when
41 storing, which we can show can classify both groups. And counterintuitively, many UVFP
42 recordings had higher intensity than controls, when UVFP patients tend to have weaker voices,
43 revealing a dataset-specific bias which we mitigate in an additional analysis.

44 **Conclusion.** We demonstrate that recording biases in audio duration and intensity created
45 dataset-specific differences between patients and controls, which models used to improve
46 classification. Furthermore, clinician's ratings provide further evidence that patients were over-

47 projecting their voices and being recorded at a higher amplitude signal than controls.
48 Interestingly, after matching audio duration and removing variables associated with intensity in
49 order to mitigate the biases, the models were able to achieve a similar high performance. We
50 provide a set of recommendations to avoid bias when building and evaluating machine learning
51 models for screening in laryngology.

52 **Keywords:** vocal fold paralysis, acoustic analysis, voice, speech, explainability, interpretability,
53 machine learning, bias

54

55

56

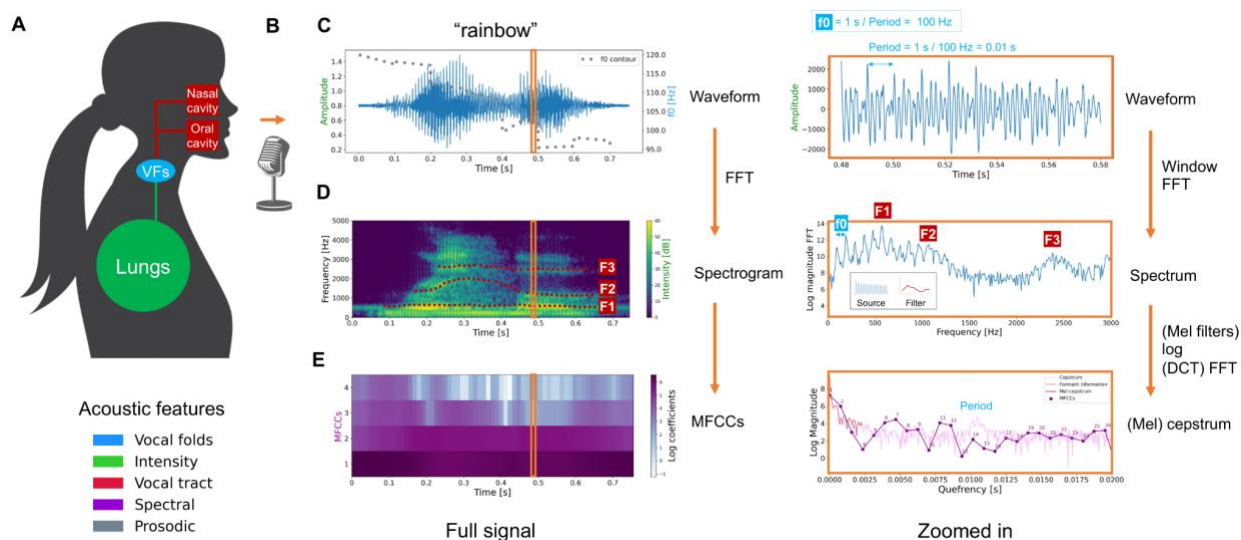
57

58 INTRODUCTION

59 Voice recordings provide a rich source of information related to vocal tract physiology
60 and human physical and mental health. Given advances in smartphones and
61 wearables, these recordings can be made anytime and anywhere. Thus the search for
62 disorder-specific acoustic biomarkers has been gaining momentum. Voice biomarkers
63 have been reported for detecting Parkinson's disease (1) as well as psychiatric
64 disorders including depression, schizophrenia, and bipolar disorder (for a systematic
65 review, see Low et al, 2020 (2)). Given our scientific understanding of the complexity of
66 speech production, multiple acoustic features have been devised for use in machine
67 learning models. In Figure 1, we describe a schematic of speech production and the
68 process of extracting certain acoustic features from an audio signal (see also Quatieri,
69 2008 (3)), which is an important part of explaining how pathophysiology would affect
70 acoustic features that are used in machine learning classifiers. Panel (A) depicts speech
71 as the result of the neural coordination of three subsystems: the respiratory system
72 (lungs), the laryngeal system (vocal folds), and the resonatory system of the vocal tract
73 (pharynx, oral cavity, nasal cavity, articulators, and subglottal effects). Speech
74 production requires air flow from the lungs to generate sound sources that are filtered
75 by the vocal tract. Panel (B) captures the fact that environmental, microphone, and
76 digital sampling characteristics (e.g., background noise, microphone gain, sampling
77 rate) can affect acoustic features. Panel (C) shows the waveform of the audio signal,
78 representing the contraction (positive amplitude) and rarefaction (negative amplitude) of
79 air particles. Higher amplitudes can lead to higher perceived loudness. Prosodic

80 features arise from changes over longer segments of time, which is perceived in the
81 rhythm, stress, and intonation of speech. A segment of the waveform is shown in the
82 right panel, indicating a periodic signal from the vocal folds. Panel (D) shows that for a
83 given time window, a spectrum (right panel) can be obtained through a Fast Fourier
84 Transform (FFT) which represents the magnitude of the frequencies in the signal with
85 peaks (formants F1–F3) due to vocal tract filtering of the source signal produced by the
86 vocal folds. The spectrogram (left panel) is a representation of the spectrum as it varies
87 over time. The approximate location of the F0 and first formants are displayed. Finally,
88 (E) It is possible to separate source and filter components by computing the inverse
89 FFT of the log of the magnitude of the spectrum, called the cepstrum (right panel). The
90 peak in the cepstrum reflects the periodic glottal fold vibration while lower quefrequency
91 components reflect properties of the resonatory subsystem. For speech recognition, Mel
92 filters are applied to the spectrum to better approximate human hearing. A conversion of
93 the Mel-spectrum to a cepstrum using a Discrete Cosine Transform (DCT) generates
94 mel-frequency cepstral coefficients (MFCCs). Similar to the cepstrum, lower MFCCs
95 track vocal-tract filter information.

96



97

98 **Figure 1. Schematic of speech production and the process of extracting certain acoustic features**
 99 **from an audio signal.**

100 (A) Speech production, (B) recording characteristics, (C) waveform of audio signal with fundamental
 101 frequency (f_0), (D) spectrogram with formants F1-F3 and intensity, (E) mel-frequency cepstral coefficients
 102 (MFCCs). Full description in the main text.

103 Furthermore, while machine learning (ML) can be a powerful and successful approach
 104 for diagnostics, they are often treated as "black-boxes". It can be difficult to determine
 105 how the model is making a decision, that is, how it is combining input features from a
 106 given patient to generate a prediction. This is particularly worrisome given ML
 107 algorithms can detect and associate unintended or clinically irrelevant relationships and
 108 introduce bias that may be difficult to anticipate. Explainable ML refers to a series of
 109 methods and quantitative analyses for uncovering and "explaining" the rationale behind
 110 the decision made by complex algorithms, which is particularly critical in the high-stake
 111 decisions of medicine to increase trust among clinicians and patients (4).

112 There are many challenges for applying acoustic analysis to detect specific disorders.
 113 Voice characteristics are highly varied and change over time. Laryngeal pathology, age,

114 gender, size, weight, general state of health, smoking/vaping, and medications can
115 impact vocal acoustic characteristics. Diseases in the larynx and phonatory system (i.e.,
116 larynx, resonating structures, lungs) and/or neurological system, will also affect voice.
117 Compensatory production strategies and environmental conditions can also change the
118 vocal signal. Furthermore, because hoarseness is such a frequent occurrence and
119 specialty voice centers are rare, vocal fold disorders are often undiagnosed, under-
120 reported, or misdiagnosed (5).

121 We chose vocal fold paralysis as the study cohort for several reasons. First, it is
122 clinically important. UVFP can have detrimental effects on voice and quality of life with
123 resultant morbidity related to respiration, swallowing and aspiration. Vocal fold paralysis
124 may occur due to iatrogenic injury, malignancy, idiopathic, and neurological disease.
125 Overall, surgical iatrogenic injury accounts for 46% of all UVFP in adults and thyroid and
126 parathyroid surgeries are responsible for 32% of postsurgical UVFP (6). There is a
127 significant need for a screening tool for the diagnosis and tracking of UVFP because of
128 the high impact of this condition on productivity and quality of life. Screening could be
129 done remotely and frequently, especially when surgical specialists and laryngeal exams
130 are not readily accessible due to geographical, financial, and other barriers (7). Using
131 an explainable ML model as a screening tool for UVFP can provide greater clarity as to
132 who most needs laryngoscopy and provides insight in the key voice characteristics
133 related to the pathophysiology (8–12). The costs associated with UVFP not only relate
134 to patient morbidity and diminished quality of life but also to the economic burden
135 placed on our healthcare system. Greater lengths of hospitalization and increased

136 hospital costs have been associated with postsurgical VFP (13,14). Access to
137 specialists for diagnosis is limited and early detection and management of UVFP appear
138 to improve length of stay and surgical outcomes (15). Special consideration should be
139 given to what the model can actually classify: a model that generalizes well in
140 classifying UVFP from controls may not be able to screen for UVFP out of other voice
141 disorders, but could be used to monitor UVFP patients remotely and affordably during
142 treatment or detect risk for UVFP when it is the most likely cause such as dysphonia
143 after thyroid surgery.

144 Furthermore, UVFP is an ideal model for demonstrating the explainability of ML. UVFP
145 occurs when the mobility of a single vocal fold is impaired as a consequence of
146 neurological injury and diagnosis is consistently verified through routine laryngoscopy;
147 therefore, ground truth labels are available. Second, the clinical signs of UVFP are well-
148 described. These characteristics include a weak, breathy voice quality, early vocal
149 fatigue, reduced cough strength, and aspiration with thin liquids (16,17). Therefore, the
150 acoustic differences between UVFP patients and healthy controls can be interpreted
151 with regards to perceptual symptoms and a well-understood pathophysiology. In
152 contrast, explaining important variables to predict a disorder which is hard to diagnose
153 (e.g., has low inter-rater reliability) and has an unclear pathophysiology would ironically
154 result in a poor explanation, because it would be puzzling how or even if the disorder
155 could modulate the important acoustic variables. Of course, machine learning models
156 can also offer novel explanations into a disorder by characterizing novel characteristics.
157 However, if these models use high-dimensional feature vectors, they are more likely to

158 overfit when using small datasets (18,19), which should lead to more skepticism of
159 these novel explanations.

160 There have been several studies detecting unilateral vocal fold paralysis (UVFP) using
161 machine learning (20–28); however, most have included the disorder among a set of
162 voice disorders to be predicted. Limitations of these prior studies could be seen to fall
163 into one of following types: not reporting the performance when classifying the subset of
164 participants with UVFP out of the participants with dysphonia they were trying to detect;
165 small sample sizes given most studies contained 10 participants with UVFP or fewer
166 with one study containing 50 participants (29); a lack of algorithmic explanations: they
167 either do not report on the relative importance of each acoustic variable, use hard-to-
168 interpret input data such as a spectrogram, or use a black-box model such as neural
169 network and do not explain it; using a single type of model or just a few features which
170 may impede high predictive performance and/or obscure a more thorough explanation
171 given a single model or few features may capture only certain aspects of the task; not
172 publicly sharing their trained models to test their generalizability to new data.

173 The objectives of our study were: to detect UVFP using ML; to evaluate the
174 effectiveness of different models in differentiating the acoustic signals between patients
175 with UVFP and patients with normal functioning vocal folds (i.e., controls); to explain
176 which features are most important to the diagnostic models and examine the
177 pathophysiological relevance; and to compare performance to human clinicians
178 evaluating audio recordings. To achieve these objectives, we evaluated four different
179 classes of machine learning algorithms to assess classification performance, obtained

180 the minimal set of features necessary for detection, and identified the most important
181 acoustic features for model construction after removing redundant features. Ultimately,
182 we wanted to see if the most important features identified by the machine learning
183 models matched clinically-known relevant acoustic changes.

184

185 **MATERIALS AND METHODS**

186 This study was approved by the Institutional Review Board at Massachusetts Eye and
187 Ear Infirmary and Partners Healthcare (IRB 2019002711).

188 **Participants and voice samples**

189 Through retrospective chart analysis from 2009 to 2019, a total of 1043 patient charts
190 were reviewed from a tertiary care laryngology practice who underwent endoscopic
191 evaluation and voice testing. Of those, 53 patients with confirmed UVFP were identified.
192 They had documented vocal fold paralysis by endoscopic examination and had
193 undergone acoustic analysis as part of routine clinical care. Each patient had four
194 acoustic recordings. These included three sustained vocalizations of the "a" vowel
195 sound (ɑ in the International Phonetic Alphabet) and a reading of the introductory
196 paragraph of the rainbow passage (30). The acoustic recordings were all taken in an
197 acoustically shielded room. For each of these 53 patients, a board-certified
198 otolaryngologist reviewed their clinical history, video laryngoscopy as well as their audio
199 samples to confirm that they were correctly classified to have UVFP. Voice samples

200 from an additional 24 patients were collected prospectively using a mobile software,
201 OperaVOX™ on an iPad, who were being treated for UVFP. These patients also had
202 the same four acoustic recordings as the patients from retrospective chart review. This
203 combination of data collection yielded a total of 77 UVFP patients for analysis, of which
204 48 had left UVFP and 29 right UVFP.

205 All of the patients were then matched with control samples from a database of patients
206 without UVFP who had also undergone acoustic analysis. Each control was the same
207 sex and had the same smoking status as the UVFP patient and within three years of
208 age, and had documented laryngeal examinations that verified the absence of vocal fold
209 mucosal pathology. The controls were excluded if they had established laryngeal
210 surgery, vocal fold lesions, radiation, head and neck cancer, or neurological disease.
211 The controls had recorded the same four acoustic recordings as the retrospectively
212 gathered UVFP group. A board-certified otolaryngologist confirmed that the voice
213 recordings and video laryngoscopies of these controls matched normal expectancies.
214 The reading samples were divided in thirds to match the amount of vowel production
215 samples, resulting in 6 samples for most participants. Reading recordings were not
216 available for three patients and three patient vowel samples were removed due to
217 containing multiple vowel productions or a cough. The final dataset that was analyzed is
218 described in Table 1. Reading+vowel refers to including all samples (i.e., ~6 samples)
219 from the same participant with the goal of either obtaining higher performance or
220 discovering features that show variation in relation to diagnosis consistently across
221 tasks. Mean (SD) audio lengths were 6.81s (5.47) for reading samples and 3.95s (1.00)

222 for vowel samples. The audio samples were processed using OpenSmile with the
223 eGeMAPS configuration file (article (31) , source code (32)) which applies different
224 summarization statistics to the time series depending on the feature resulting in 88
225 features per sample covering information related to the vocal folds (F0, jitter, shimmer),
226 intensity (loudness, HNR), vocal tract (F1–3 frequency, bandwidth, amplitude), spectral
227 balance (alpha ratio, Hammamberg index, spectral slope, MFCC 1–4, spectral flux), and
228 prosody (voice and unvoiced segments, loudness peaks per second). See section
229 "eGeMAPS features" in Sup. Mat. for full list.

230 **Table 1. Sample sizes and demographic information**

	UVFP	Controls	Total
N	77	77	154
Mean age (SD)	56.4 (18.7)	56.6 (18.8)	56.5 (18.7)
Sex (F/M)	39/38	39/38	78/76
Reading	222	231	453
Vowel	227	231	458
Reading+vowel (total)	449	462	911

231 SD: standard deviation; F: female; M: male.

232 **Machine learning models of increasing complexity**

233 With the goal of classifying voices recording into either UVFP or controls, we used four
234 machine learning algorithms of increasing complexity from the *scikit-learn* (v0.21.3)

235 using the *pydra-ml* (v0.3.1) toolbox (33) (default parameters were used unless
236 otherwise specified):

237 (1) Logistic Regression: a simple linear model that is constrained to use few features
238 due to an L1 penalty making it the simplest model (“liblinear” solver was used which is
239 ideal for smaller datasets).

240 (2) Stochastic Gradient Descent (SGD) Classifier: we used a log loss which implements
241 a logistic regression; therefore, it is also a linear model but tends to use more features
242 due to an elastic net penalty, making it slightly more complex (the `max_iter` parameter
243 was set to 5000 and `early_stopping` was set to True).

244 (3) Random Forest: it is an algorithm that uses simpler decision trees (i.e., weak
245 learners) on feature subsets but then averages the trees’ predictions to create a
246 stronger learner, making it harder to interpret which features are important across trees.

247 (4) Multi-Layer Perceptron: it is a neural network classifier which incorporates, in our
248 case, 100 instances of perceptrons (artificial neurons), which are connected to each
249 input feature through weights with an added ReLU activation function to capture
250 nonlinear relationships in the data. It is not possible to know exactly how the hundreds
251 of internal weights interact to determine feature importance, making the model difficult
252 to interpret directly from its parameters (the `max_iter` parameter was set to 1000; alpha
253 or the L2 penalty parameter was set to 1).

254 To generate independent test and train data splits, a bootstrapped group shuffle split
255 sampling scheme was used. Bootstrapping is more optimal than cross-validation on
256 smaller datasets and provides a measure of uncertainty through a confidence interval
257 (34). For each iteration of bootstrapping, a random selection of 20% of the participants,
258 balanced between the two groups, was used to create a held-out test set. The
259 remaining 80% of participants were used for training. This process was repeated 50
260 times, and the four classifiers were fitted and tested for each test/train split. 50 was
261 chosen empirically to provide a meaningful distribution while reducing the computational
262 time complexity of higher values. The Area Under the Receiver Operating Characteristic
263 Curve (ROC AUC; perfect classification = 1; chance = 0.5) was computed to evaluate
264 the performance of the models on each bootstrapping iteration, resulting in a distribution
265 of 50 ROC AUC scores for each classifier. We did not perform hyperparameter tuning
266 as the models achieve relatively high performance (median ROC AUC ~ 0.85).
267 Additionally, for each iteration, each classifier was trained with randomized
268 patient/control labelings to generate a null distribution of ROC AUC scores (i.e., a
269 permutation test). Each model's performance was statistically compared to other models
270 and to the null distributions using an empirical p-value, a common and effective
271 measure for evaluating classifier performance (see Definition 1 in (35)). The significance
272 level was set to $\alpha = 0.05$.

273 **Assessing feature importance**

274 Kernel SHAP (SHapley Additive exPlanations) was used to determine which acoustic
275 features were most important for each model to detect UVFP. This method is model
276 agnostic in that it can take any trained target model (even “black box” neural networks)
277 and compute feature importance (36). It does so by performing regression with L1
278 penalty between different sets of input features and a single prediction made by the
279 target model. It then uses the coefficients of the additional regression model as a
280 measure of feature importance for a single prediction. We took the average of the
281 absolute SHAP values across all test predictions (positive and negative values are both
282 important for classification). We then weighted the average values by the model’s
283 median performance since an important feature for a bad model could be a less
284 important feature for a good model and vice versa. Since we trained each model 50
285 times (i.e., one for each bootstrapping split), we computed the mean SHAP values
286 across splits for each model. This pipeline (i.e., machine learning models, bootstrapping
287 scheme, SHAP analysis) was done using *pydra-ml*.

288 **Reducing collinearity to do explainability analysis using**

289 **Independence Factor**

290 Highly correlated features (i.e., collinearity) can influence model generation and
291 interpretation. Two models may obtain similar performance while using different
292 features or placing different weights on the same features (i.e., underspecification
293 (18,37)) . This makes it difficult to compare algorithmic explanations across models. For

294 instance, mean F1 frequency may be less important to a given model because the
295 model uses mean F2 frequency which happens to capture very similar information in a
296 particular dataset (i.e., has a high correlation), whereas a different model may use F1
297 instead of F2 or use both but assign less importance to each and still obtain the same
298 performance. To enforce models to use the same features that capture very similar
299 information and be able to compare feature importance across models, we kept a single
300 feature out of the sets of features that share similar information above a given threshold.

301 We used a custom algorithm we call Independence Factor whereby for each
302 feature in alphabetical (i.e., arbitrary) order, we removed features that show strong
303 dependence above a given threshold. The step was repeated for remaining features.
304 We use distance correlation from the Python *dcor* package (v0.4) because, unlike
305 Pearson *r* or Spearman *rho*, it can capture non-monotonic relationships (38,39). We
306 used the following threshold values for the distance correlation [1.0, 0.9, 0.8, 0.7, 0.6,
307 0.5, 0.4, 0.3, 0.2] to compute the Independence Factor, which removed increasingly
308 more features (i.e., 1.0 keeps all features and 0.2 removes features that have a
309 distance correlation above 0.2). We chose the feature size which contains at least one
310 model that scores within three percentage points of the performance using all features,
311 with the goal of obtaining a more parsimonious model for subsequent explanation while
312 maintaining high accuracy. Thus, removing redundant features makes the models
313 easier to interpret for clinical relevance. To visualize the original redundancy across
314 features, we computed clustermaps using *seaborn* package (v0.10.1) performing
315 hierarchical clustering with the average-linkage method and Euclidean distance. This

316 was performed on the pairwise distance correlation, computed separately on data from
317 UVFP, controls, UVFP+controls and on reading, vowel, and reading+vowel.

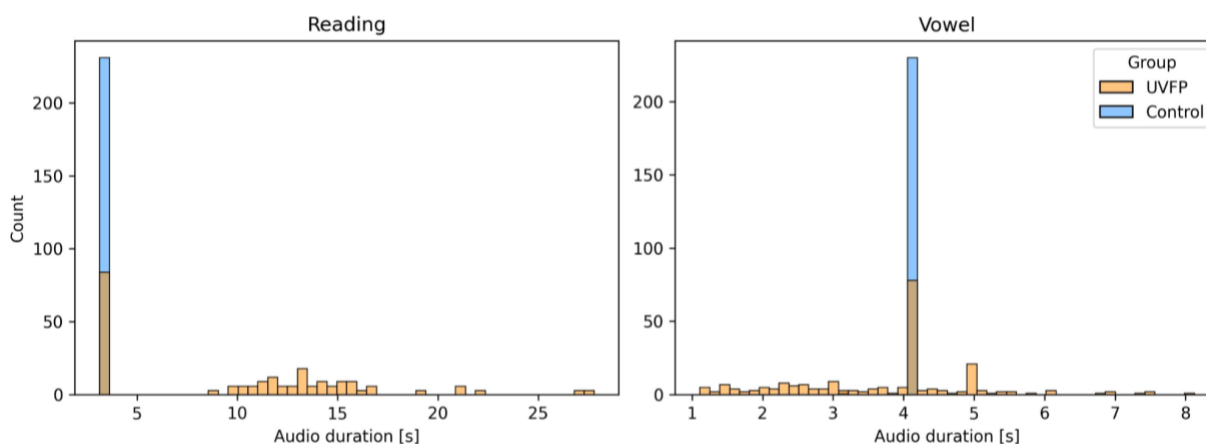
318

319 **Performance using most important and least important features**

320 Studies tend to report and describe the top N features out of M features, but it is not
321 clear what performance the model would obtain when using only those top N features;
322 perhaps it would perform substantially worse than the full model. We will report
323 performance using only top 5 features as well as performance without top 5 features to
324 provide a more realistic evaluation of their importance.

325 **Performance using audio duration**

326 Figure 2 indicates clear differences in the distributions of audio recording duration between
327 UVFP patients and controls. This is due to how recordings were processed and saved and not
328 necessarily due to an intrinsic property of UVFP (e.g., slower speech), which reveals a bias that
329 models can leverage but is not expected to generalize well under different audio processing
330 procedures. Therefore, we examine whether audio duration alone could perform well in
331 classification of UVFP.



332

333 **Figure 2. Distribution of audio duration for reading and vowel tasks split by group**

334 **reveals a dataset bias.** The mode of the audio durations for the controls is 3.5 s for reading

335 samples and 4.11 s for vowel samples.

336

337

338

339 **Performance using cepstral peak prominence**

340 To evaluate whether results are sensitive to choice of features, we use a different set of

341 features derived from cepstral peak prominence (CPP) given it has been shown to be a

342 good measure of breathiness and dysphonia (40,41). We match the summary statistics

343 across the audio recording that eGeMAPS uses: CPP mean, CPP coefficient of

344 variation (standard deviation normalized by the mean), CPP 20th percentile and CPP

345 80th percentile. We use our custom Python implementation which matches MatLab's

346 COVAREP output (42).

347 **Clinician ratings**

348 In order to corroborate whether there are unintended recording differences between
349 UVFP patients and controls that may lead to bias, one otorhinolaryngologist and two
350 speech-language pathologists rated each audio recording of the reading task (one per
351 participant, not split in three) for the following variables (and possible responses), in
352 order: background noise (None, Some, High); UVFP (yes, no), CAPE-V severity (0 to
353 100), CAPE-V roughness (0 to 100), CAPE-V breathiness (0 to 100), CAPE-V strain (0
354 to 100), CAPE-V pitch (0 to 100), CAPE-V loudness (0 to 100; estimated loudness as if
355 the rater were in the recording room), recording loudness (low, medium, high; loudness
356 of the recording). Inter-rater agreement was assessed using intra-class correlation for
357 all numerical variables and Light's k for the binary presence of UVFP (43) using the R
358 package *irr* (v0.84.1) (44). The entire reading task was provided instead of the task split
359 in three to make assignment easier for clinicians. The reading task was chosen over the
360 sustained vowel because we expected it to be easier for clinicians to detect UVFP.

361 **RESULTS**

362 **Performance of models using acoustic features**

363 In Table 2, we report performance for models using all features, models after removing
364 redundant features, models using only top 5 features (to understand their unique role in
365 performance), models using all 88 features without 5 features (to understand whether
366 the top 5 features are necessary for high performance), models using audio duration

367 length, and models using a different feature set based on CPP. Performance was found
368 to be high across most models except CPP-based models. Some of the models just
369 using audio duration length were able to achieve close to the highest performance,
370 which reflects the expected effect of the difference in the dataset. Given dependent
371 features provide similar information (see Supplementary Figures S1, S2, S3, S4, S5,
372 S6, S7, S8, and S9) and distort feature importance analyses, we then tested
373 performance after removing redundant features using the Independence Factor method
374 previously described. Supplementary Figure S10 shows performance for different
375 feature set sizes with increasing amounts of redundant features. From this analysis, we
376 selected the feature-set size that resulted in best performance using the least amount of
377 features for subsequent analyses: 39 features (reading), 13 (vowel), 19
378 (reading+vowel). After removing related features (i.e., reducing collinearity) from the
379 original 88 features, similar performance was obtained (median ROC AUC = 0.84–0.87)
380 using fewer features. Supplementary Materials "Feature selection" section describes an
381 analysis of how this method compares to removing features across each train set (see
382 Sup. Mat. Table S1).

383

384
385
386

387 **Table 2. Model performance**

	Features	LogisticRegression	MLP	RandomForest	SGDClassifier
Reading	88	.87 (.78–.93; .50)	.87 (.80–.93; .50)	.87 (.76–.91; .49)	.83 (.76–.89; .50)
Vowel	88	.84 (.77–.89; .50)	.86 (.79–.91; .50)	.86 (.79–.91; .51)	.80 (.72–.87; .50)
Reading+Vowel	88	.84 (.76–.91; .50)	.86 (.74–.92; .48)	.85 (.77–.92; .49)	.79 (.72–.86; .51)
Reading	39	.84 (.76–.92; .50)	.83 (.76–.91; .50)	.87 (.77–.91; .51)	.78 (.71–.86; .51)
Vowel	13	.80 (.70–.90; .50)	.81 (.74–.91; .50)	.84 (.75–.90; .52)	.74 (.58–.87; .51)
Reading+Vowel	19	.79 (.70–.84; .50)	.82 (.75–.88; .51)	.84 (.77–.91; .51)	.70 (.61–.77; .52)
Reading	Top 5	.81 (.73–.89; .50)	.86 (.78–.92; .47)	.85 (.77–.90; .50)	.75 (.56–.87; .57)
Vowel	Top 5	.78 (.67–.87; .50)	.82 (.74–.92; .53)	.81 (.72–.87; .50)	.72 (.57–.82; .49)
Reading+Vowel	Top 5	.80 (.70–.86; .50)	.82 (.74–.88; .50)	.81 (.74–.89; .53)	.72 (.55–.83; .52)
Reading	88 - Top 5	.85 (.76–.92; .50)	.87 (.77–.92; .49)	.85 (.77–.90; .52)	.82 (.71–.89; .51)
Vowel	88 - Top 5	.84 (.75–.93; .50)	.86 (.72–.93; .51)	.84 (.74–.94; .52)	.80 (.70–.90; .48)
Reading+Vowel	88 - Top 5	.84 (.74–.89; .50)	.85 (.76–.91; .50)	.85 (.76–.91; .50)	.79 (.71–.87; .50)

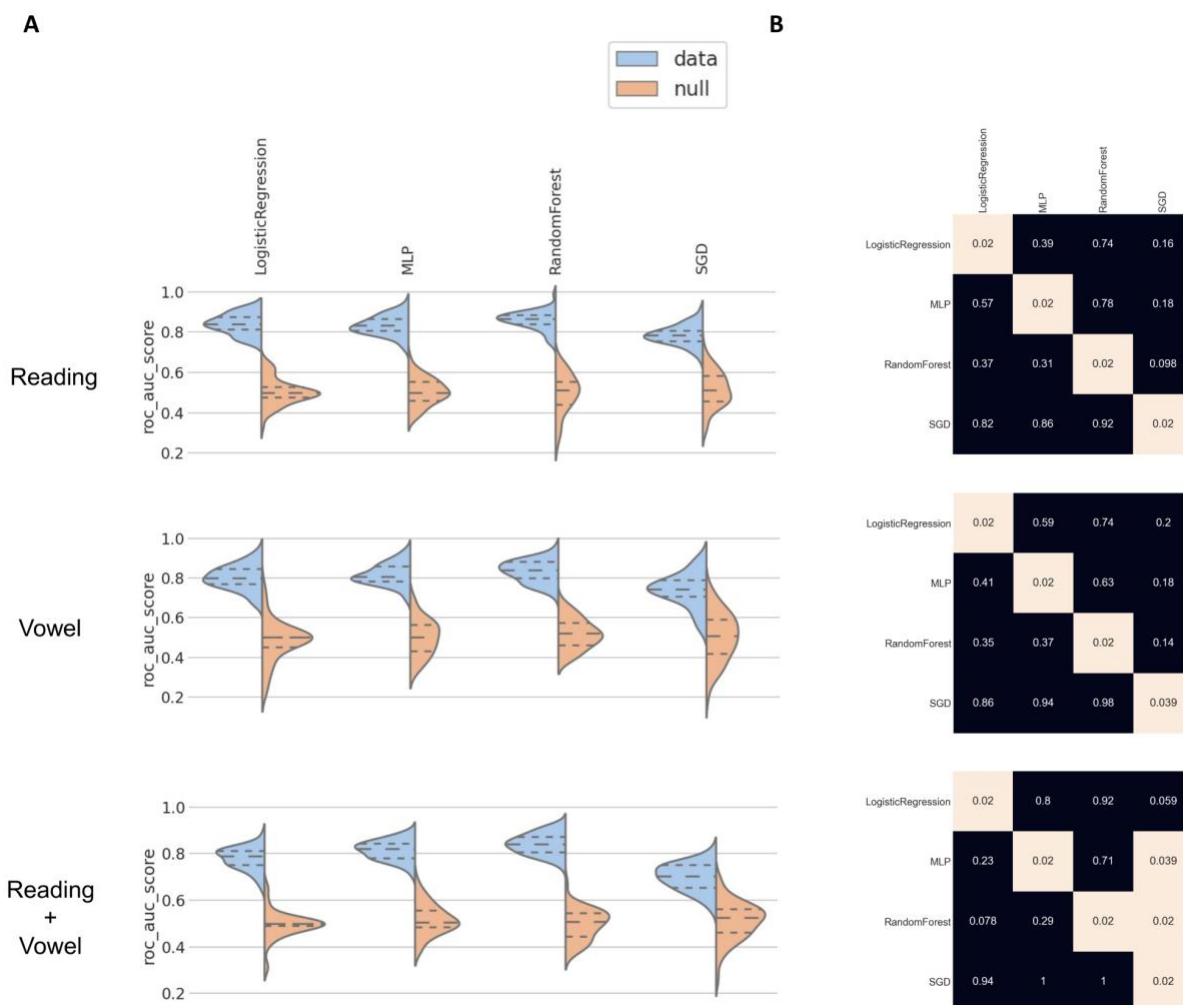
Reading	Duration 1	.81 (.73–.88; .50)	.81 (.73–.88; .50)	.85 (.77–.93; .50)	.76 (.50–.88; .50)
Vowel	Duration 1	.70 (.61–.77; .50)	.80 (.70–.91; .51)	.86 (.76–.94; .52)	.50 (.31–.68; .51)
Reading+Vowel	Duration 1	.70 (.64–.76; .50)	.76 (.67–.84; .50)	.86 (.73–.92; .50)	.64 (.45–.70; .50)
Reading	CPP 4	.76 (.64–.84; .50)	.76 (.64–.84; .46)	.71 (.64–.78; .55)	.74 (.60–.84; .50)
Vowel	CPP 4	.82 (.73–.90; .50)	.82 (.71–.90; .53)	.77 (.65–.85; .50)	.77 (.40–.86; .49)
Reading+Vowel	CPP 4	.72 (.65–.80; .50)	.74 (.68–.84; .53)	.72 (.65–.78; .50)	.68 (.44–.78; .49)

388 Performance of models using either all 88 features, non-redundant features (39, 13, 19), top five most
 389 important features, all 88 features minus top 5 most important features using eGeMAPS features. We
 390 then compared this to using just audio duration as well as a different feature set based on CPP. Median
 391 ROC AUC score from 50 bootstrapping splits (90% confidence interval; median score of null model
 392 trained on permuted labels which should be at .50 if at chance). For full distributions of scores see Figure
 393 S10 in Supplementary Materials. Removing features is a post-hoc analysis because features were
 394 selected based on observing performance on the test sets, and therefore performance might be slightly
 395 overly optimistic and would need to be tested on an independent test set for further validation. MLP: Multi-
 396 Layer Perceptron; SGD: Stochastic Gradient Descent Classifier; CPP: Cepstral Peak Prominence.

397 The bootstrapped ROC AUC distributions and permutation tests for the reduced
 398 (parsimonious) models using the non-redundant feature set are shown in Figure 3. The
 399 figure reports a one tailed statistical comparison (row > column) of models using an
 400 empirical p-value, which represents the fraction of column-model scores where the row-
 401 model classifier had a higher mean performance (e.g., a p-value of 0.02 indicates that
 402 the mean score of a row model is higher than 98% of column-model scores).

403

404



405

406 **Figure 3. Model performance comparison using a permutation test using non-redundant features.**
 407 **(A)** Scores from models trained on true labels (blue) and trained on permuted labels (orange) over
 408 bootstrapping splits. **(B)** Statistical comparison between models (annotation = p-value, highlighted =
 409 significant results).

410

411 Given 24 UVFP patients were recorded with a different device, an iPad, we trained
 412 models without their samples to make sure these differences in recordings were not
 413 driving performance. There was a small drop in performance, which could be due to a

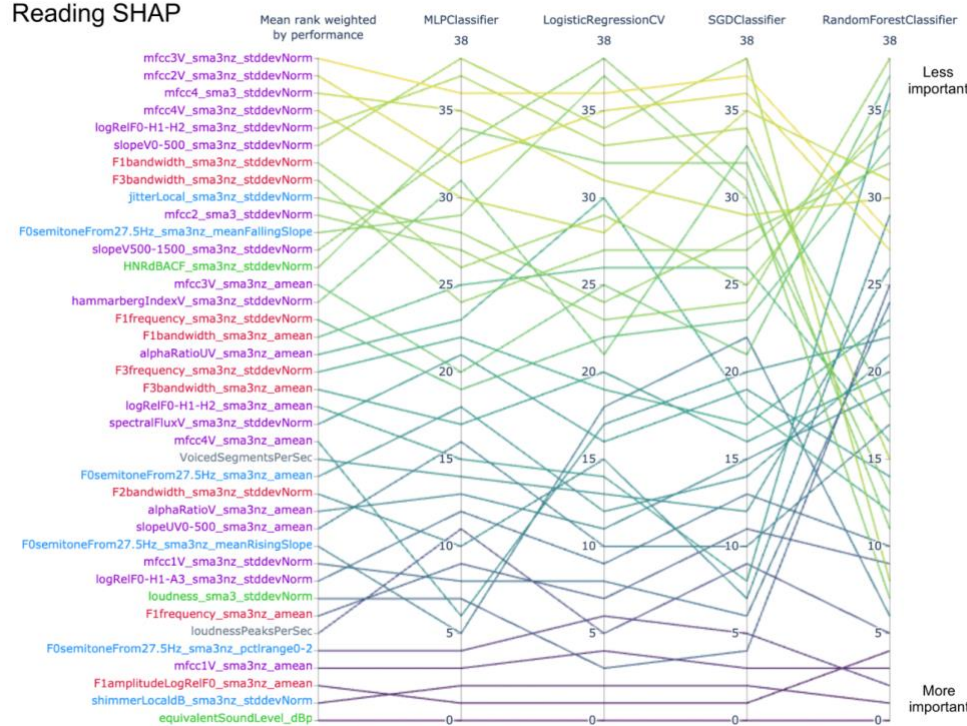
414 bias (the full, original model using information of the recording device), but could also be
415 due to removing training samples. The drop in performance is not large enough to
416 suspect that differences in recording are driving the full original model's performance
417 (see Sup. Mat. Table S2, Table S3, and analysis in Supplementary section
418 "Performance removing participants that used other recording system").

419 **Assessing feature importance**

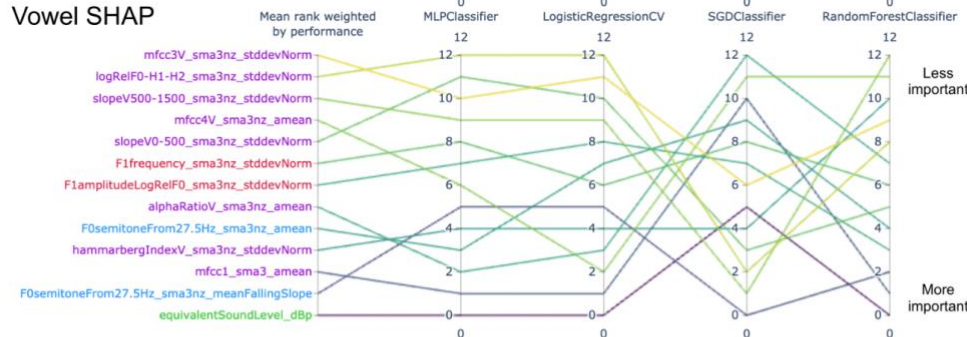
420 Figure 4 reports feature importance using SHAP for all models. For the reading-based
421 models, all models tend to use the same top 5 features except SGD, which also has the
422 lowest performance. For further description of features and the chosen classification of
423 features, see Eyben et al. (2015) (31) and Low et al. (2020) (2). When reviewing
424 important features, it is key to note that any of the features with which it is codependent
425 or associated could be a reasonable important feature (see clusters of redundant
426 features in Supplementary Figures S1-S9). All models except SGDClassifier tend to
427 have high performance, but the variance on feature importance rank is evidence that
428 models can use different feature information and still obtain similar high –although not
429 perfect– performance. We further display the distribution of each top feature and its
430 individual performance in Figure 5, which shows that no single feature is enough to
431 dissociate groups with high performance. This figure also revealed the bias: the
432 intensity-related feature equivalent sound level was counterintuitively higher for UVFP
433 patients than controls. Figure 6 reports similarity between top 5 features and all original
434 88 eGeMAPS features. Features that have a high dcor or distance correlation (i.e.,

435 cluster) with top 5 features were not used in models to avoid redundancy, but still share
436 similar information and can therefore be considered important features as well.
437 Hierarchically-clustered heatmaps for other data types (vowel, reading, both) and
438 groups (UVFP patients, controls, both) are displayed in Supplementary Figures S1, S2,
439 S3, S4, S5, S6, S7, S8, and S9. Clustering tends to reflect pre-defined features types
440 such as those reflecting patterns from vocal folds, intensity, vocal tract, spectral
441 analyses, and prosody.

Reading SHAP



Vowel SHAP



Reading + Vowel SHAP

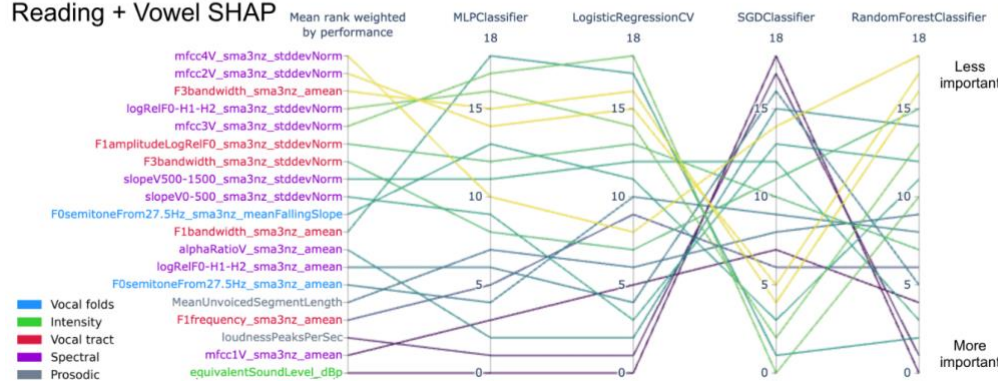
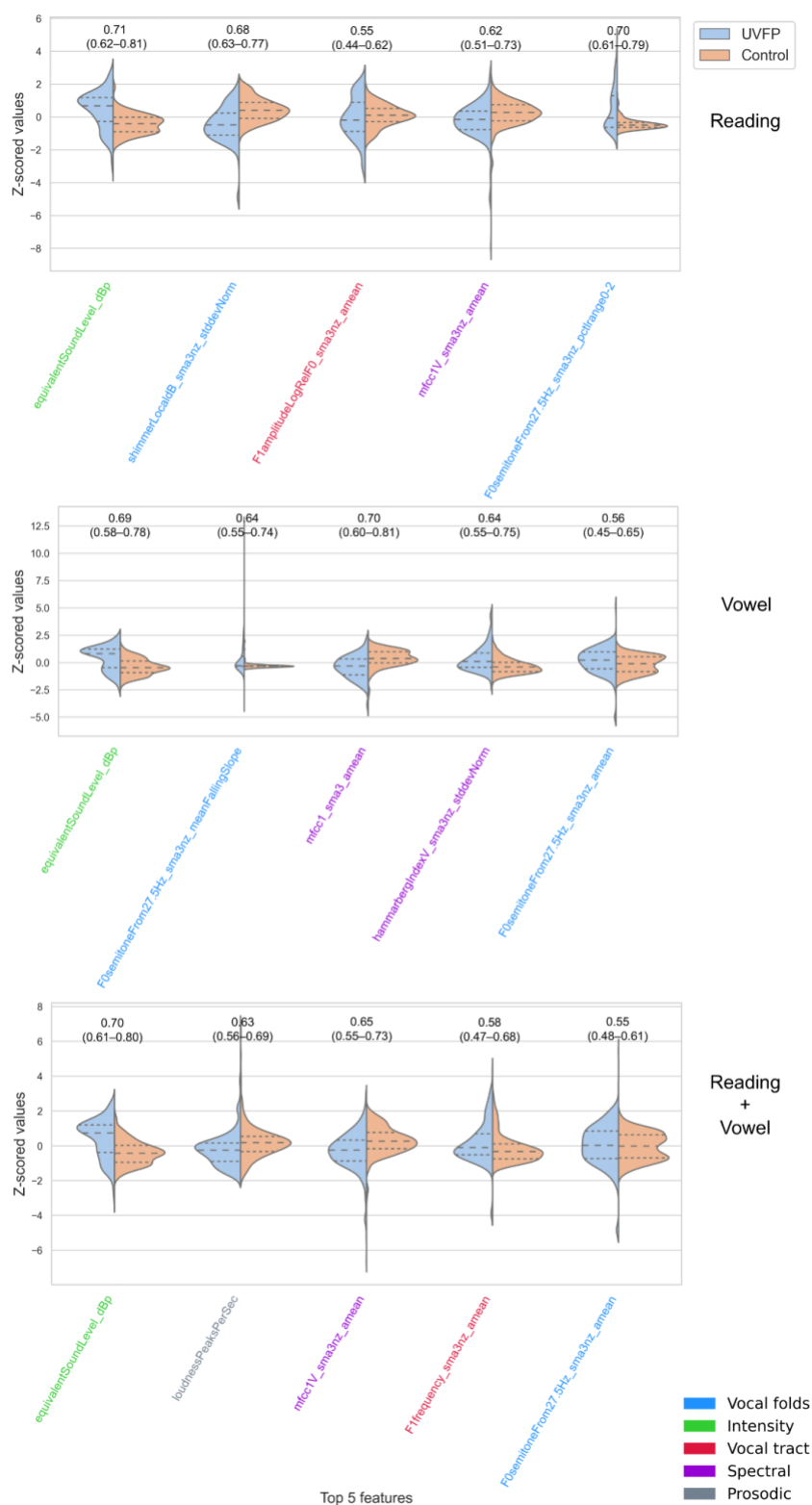


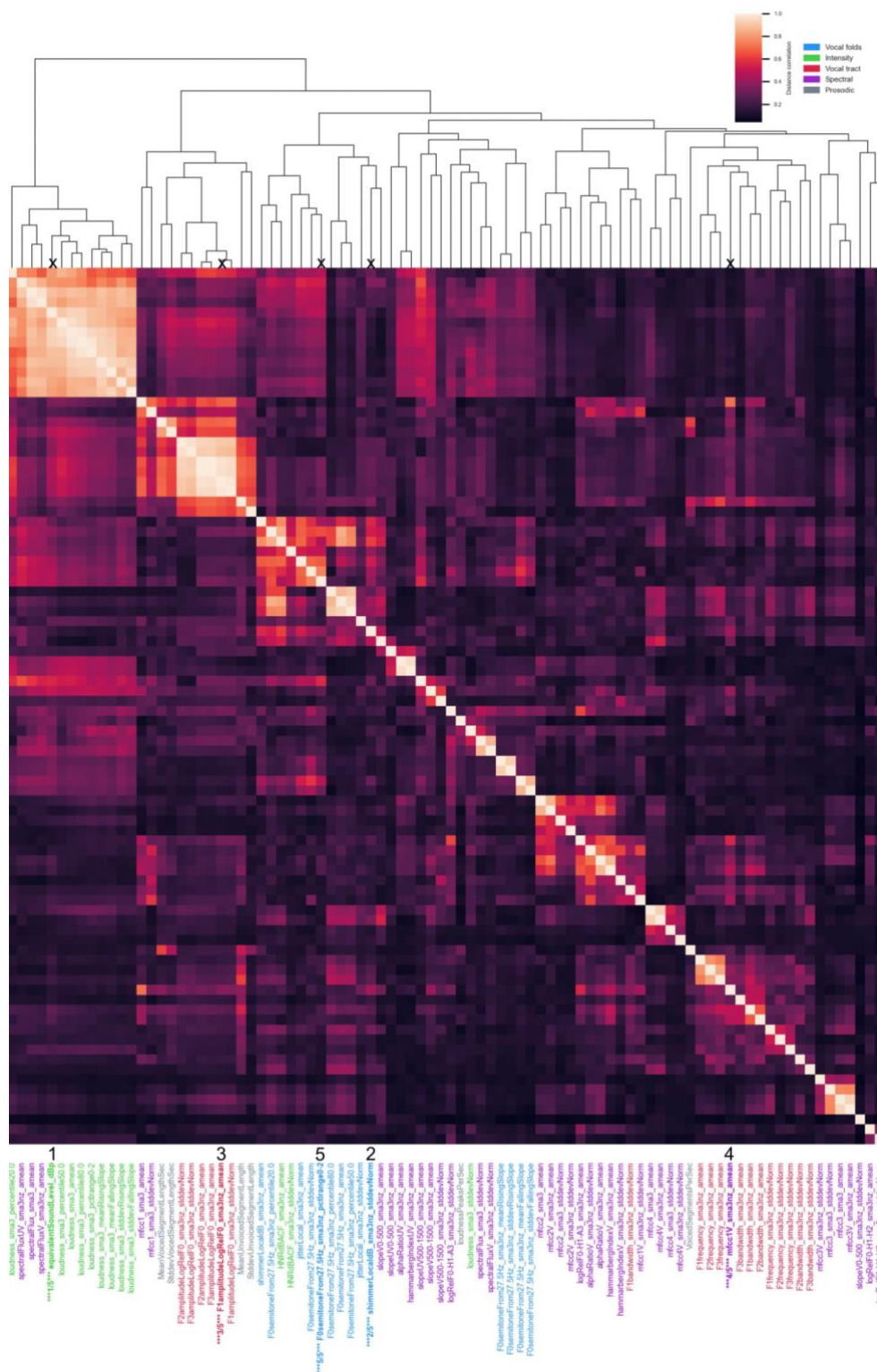
Figure 4. Feature importance parallel coordinate plot. Rank reads from bottom (most important) to top (least important). Mean rank is weighted by performance of each model to avoid a lower performing model biasing the mean rank.



446
447

Figure 5. Distributions for top 5 features and corresponding performance for single features. Logistic

448 Regression with L1 penalty was used. No single feature is enough to dissociate groups with high
449 performance. Null models' median performance was 0.5.



450
 451 **Figure 6. Feature redundancy with top 5 features highlighted.** Top 5 features are highlighted in bold and
 452 their rank is displayed. Squares are clusters of redundant features. Computed with all participants on the
 453 reading task.

454 **Clinician ratings**

455 The median ROC AUC for humans was 0.78 (min. = 0.74 to max. = 0.81) meaning the
456 machine learning models performed better than the highest performing clinician.

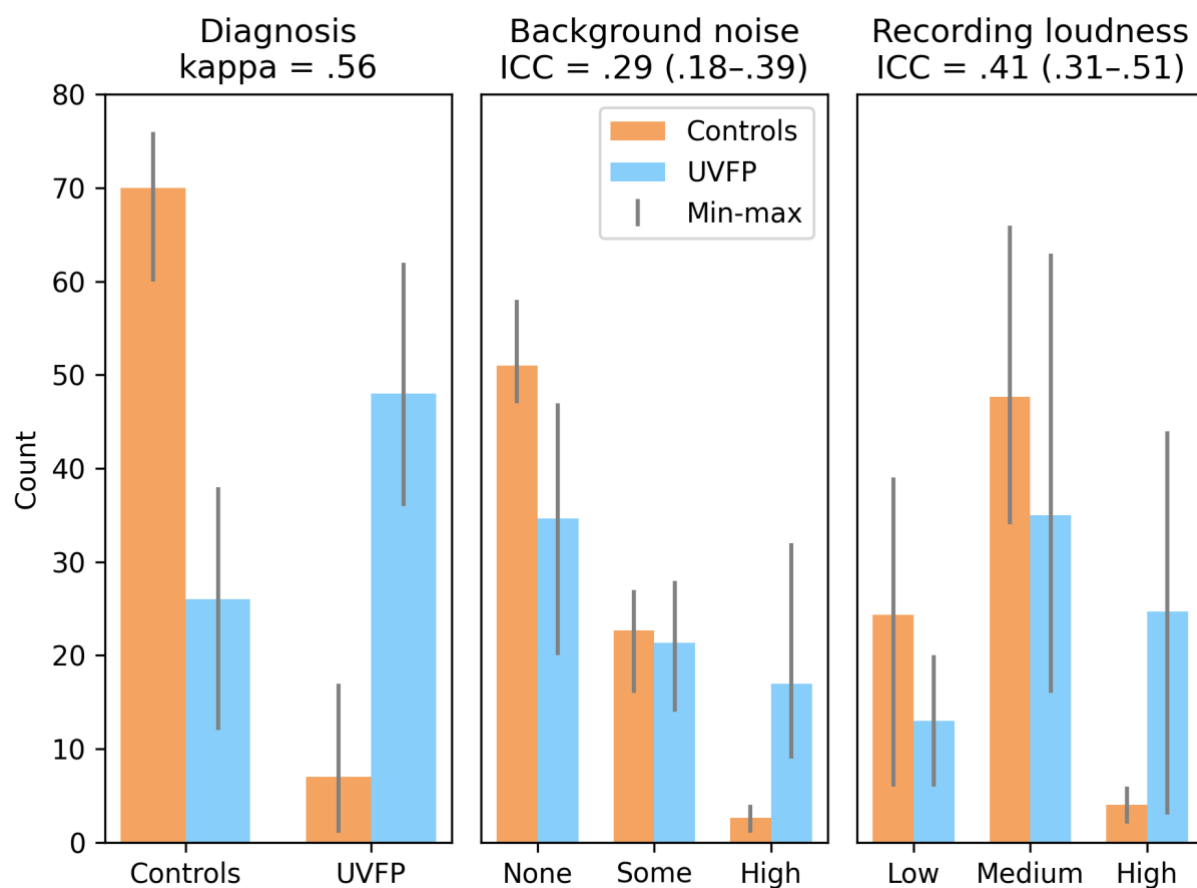
457 Interestingly, using the average clinician's CAPE-V ratings within machine learning models
458 was able to obtain a maximum median ROC AUC of 0.84 (0.72–0.92) with the Random
459 Forest model (Table 3). Using clinicians' perceptual ratings of background noise and
460 recording loudness achieved a maximum median ROC AUC of 0.77 (.63–.87).

461 **Table 3. Performance using clinician ratings as variables for machine learning models**

	Features	LogisticRegression	MLP	RandomForest	SGD
CAPE-V	6	.80 (.69–.88; .50)	.81 (.71–.90; .50)	.84 (.72–.92; .49)	.77 (.45–.92; .51)
Noise+ loudness	2	.76 (.59–.86; .50)	.77 (.63–.87; .50)	.73 (.62–.83; .52)	.64 (.45–.78; .50)

462 Median ROC AUC score from 50 bootstrapping splits (90% confidence interval; median score of null model
463 trained on permuted labels which should be at .50 if at chance).

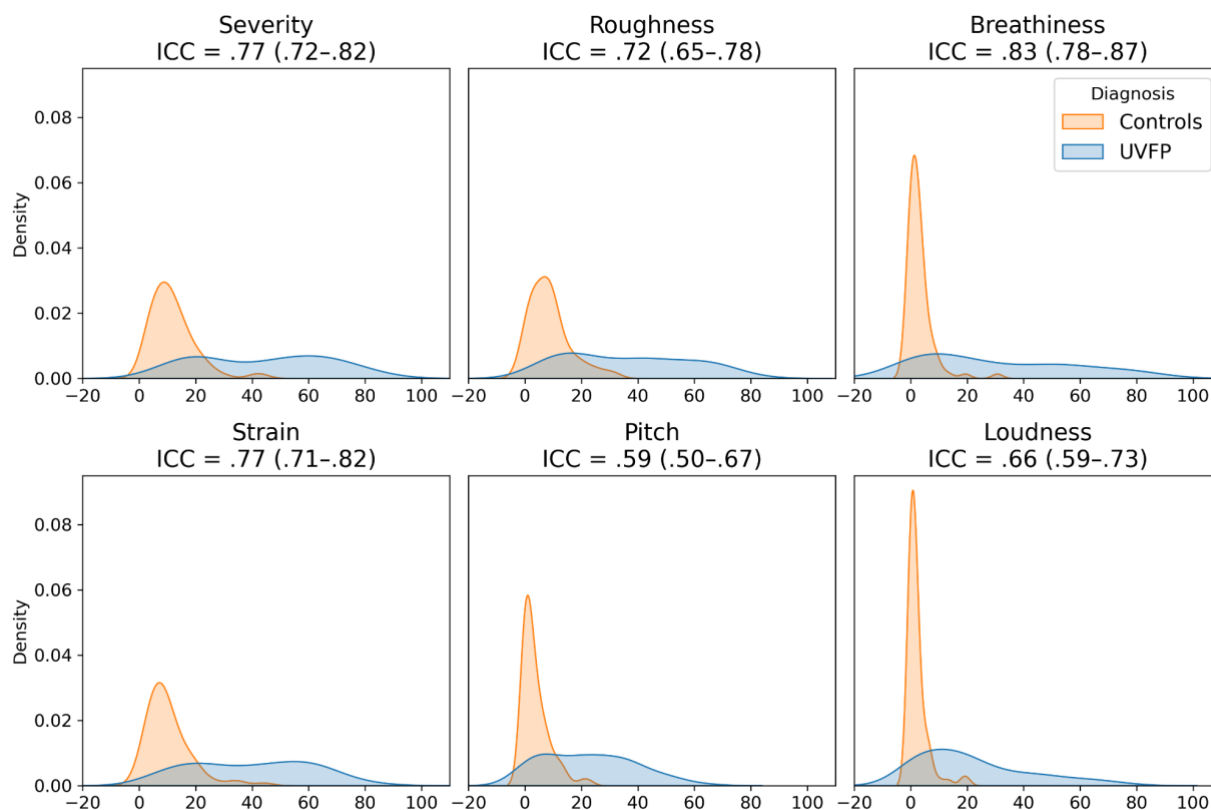
464 In Figures 6 and 7 we report the inter-rater reliability (Flight's kappa and ICC) along with
465 the distribution of the ratings. Common cutoffs for inter-rater agreement are poor for values
466 less than .40, fair for values between .40 and .59, good for values between .60 and .74,
467 and excellent for values between .75 and 1.0 (45). Background noise had poor reliability
468 across rater, UVFP and recording loudness had fair reliability (see Figure 7) and CAPE-V-
469 inspired ratings scored good to excellent except for pitch which was fair (see Figure 8).



470

471 **Figure 7. Descriptive statistics and inter-rater reliability of clinician ratings for unilateral vocal fold**
472 **paralysis (UVFP), background noise, and recording loudness indicating likely bias.** Error bars indicate
473 maximum and minimum count across the three raters. The disproportionate amount of UVFP samples rated
474 as having high background noise and high loudness indicates likely bias, where the gain might have been
475 raised for some UVFP patients and they may have phonated more intensely. kappa: Light's kappa; ICC:
476 intra-class correlation coefficient.

477



478

479 **Figure 8. How clinicians rate the audio recordings of read speech: descriptive statistics and inter-**
480 **rater reliability of average clinician ratings.** The average across raters was taken for each recording. ICC:
481 intra-class correlation coefficient.

482

483 Bias mitigation: matching audio duration and removing features associated to
484 intensity

485 We trimmed the longer UVFP samples so they were matched to control samples, removing the
486 audio duration difference. In Table 4, we show results on these samples after additionally removing
487 all intensity features as well as variables that have a distance correlation (dcor) with any of them
488 ≥ 0.3 and 0.4 based on the reading samples. Models have comparable performance to models
489 with the original duration and intensity-related biases. See section "Biased features" and Table S4

490 in Sup. Mat. for a list of the 44 features associated with audio duration and the 14 intensity related
491 features.

492

493 **Table 4. Performance keeping features least associated to intensity features**

	Features	LogisticRegression	MLP	RandomForest	SGD
dcor<0.4	44	.88 (.80–.92; .50)	.87 (.81–.92; .47)	.87 (.78–.93; .45)	.83 (.76–.90; .48)
dcor<0.3	20	.84 (.78–.89; .50)	.83 (.76–.9; .49)	.85 (.78–.91; .53)	.79 (.66–.87; .51)

494 Median ROC AUC score from 50 bootstrapping splits (90% confidence interval; median score of null model
495 trained on permuted labels which should be at .50 if at chance).

496

497 **Discussion**

498 This study achieves high performance in detecting UVFP from healthy voices using a few
499 seconds of audio recordings and surpassing clinician evaluations even after mitigating the
500 biases we found in the dataset. As a result of performing the explainability analysis, we
501 discovered a likely bias: intensity features were higher for UVFP patients than controls on
502 average (Figure 5) when UVFP patients should have weaker voices. There are two likely
503 causes. A first cause is that the software that had been used prompted users to speak
504 louder if they had a weak voice in order to achieve an audible recording. A second cause
505 was supported by clinicians' ratings: clinicians rated UVFP patients as having louder
506 recordings and more background noise than controls on average –when they should have
507 similar levels–, which are proxies for microphone gain having been increased. This would

508 have helped models improve performance using characteristics stemming from the
509 recording idiosyncrasies instead of from pathophysiology. However, we removed features
510 correlating with the clearly biased features and still achieved high performance.

511 Our study expands on prior studies which have used pre-existing commercial databases,
512 smaller sample sizes, fewer features, and/or methods for model evaluation that can be
513 biased in small datasets given the test sets may not be representative (for a discussion on
514 bootstrapping for clinical datasets, see Figure 6⁽²⁾). Critically, we provide a roadmap for
515 evaluating models more thoroughly including quantitatively explaining models and
516 checking the robustness of the models to different choices of speech-eliciting tasks,
517 algorithms, and feature sets. All of this should increase trust when no bias is found and
518 when explanations are robust across models and make sense to experts. Such a model
519 could fulfill several clinical needs: (1) postoperative screening for thyroid surgery-related
520 UVFP since after thyroid surgery, UVFP is common, occurring in up to 5 to 10% of
521 cases²⁷. Furthermore, laryngoscopy is not readily available to all postoperative
522 populations and symptomatic changes are notoriously variable. An ML-based screening
523 could help identify patients needing further workup and treatment, and earlier diagnosis is
524 essential to optimize long-term outcomes^{28,29}. (2) Monitoring voice during speech therapy
525 and after surgical treatment for confirmed UVFP to measure when and if the patient's voice
526 is approximating a healthy voice. (3) Preoperative screening prior to surgeries that are at
527 high risk for developing UVFP such as thyroid, head and neck, cardiac, thoracic,
528 esophageal, and cervical spine operations.

529 In Table 5 we summarize several key recommendations to avoid bias when building and
530 explaining machine learning tools for laryngology, although more could be added, and we
531 expand upon how we dealt with some of these steps in the following sections.

532

533 **Table 5. Recommendations to avoid bias for explainable machine learning models that use**
 534 **audio recordings for screening in laryngology**

Recommendations	Description
Before data collection	- Pre-register hypotheses as to which variables should be important for predicting the target group
During recording	- In a <i>controlled</i> setting: models could use any unintended differences between groups to improve classification; therefore, it is important to make sure microphone gain, background noise, instructions are consistent across participants and reflect how recordings will be done once deployed. - In a <i>remote</i> setting: we would want models to work on people's mobile devices outside the clinic. Since we cannot fully control the recording procedure, we should make sure there are no biases affecting one group more than another, test pilot instructions, and collect much more data to weaken the effect of individual recording idiosyncrasies. - Perform pilot studies to do an initial quality control - Collect representative samples so models generalize to different protected groups (e.g., ages, genders, races) or provide appropriate warnings. - Providing instructions so participants do not overproject their voice and control recording procedure so a minimum loudness threshold is not needed
Preprocessing and exploratory data analysis	- Quality control: keep natural outliers but not non-natural outliers due to measurement errors, wrong data collection, or wrong data entry (e.g., fixing mislabeled files, unexpected silent recordings, recordings with extreme much background noise) - Avoid or be cautious with preprocessing steps that might reduce the properties associated with the disorder (e.g., denoising may remove breathiness information which may be useful for prediction). - Observe distribution of variables between groups (e.g., audio duration) to make sure there are no differences that are not intrinsic to the disorder. Extra inspection of the data should be taken with retrospective studies where recording protocols were not controlled.
During training and evaluation	- Train multiple machine learning models of different complexity: two models may perform similarly but use input variables in different ways. If after training a model we only explain one of them, we might have biased conclusions of what variables characterize the disorder. - Avoid overfitting (i.e., finding patterns that do not generalize to new samples). Simple held-out test sets (e.g., of 20%) may not be representative of the population or the dataset, and therefore resampling methods (cross-validation, bootstrapping) are better. If performing hyperparameter tuning, nested resampling is needed to avoid overfitting (2). Avoid feature selection and dimensionality reduction using information from the test set/s. - Report performance on most and remaining important features
During explainability analyses	- Choosing one of the variables that are highly dependent due to collinearity (e.g., that correlate above 0.8 Spearman rho or dcor above a threshold that does not reduce performance as we did in this study) or due to multicollinearity (remove variables if variance inflation factor > 5 or 10); ; grouping correlated variables using leave-one-feature-out (LOFO); obtaining one variable from the correlated variables through dimensionality reduction (without using the test set which could lead to overfitting). - Make conclusions from the features that are robustly important <i>across</i> models; here we take the average importance rank weighted by model performance. - Evaluate potential bias: do important features match hypotheses? Do they dissociate groups in the expected direction? Do certain recording conditions perform better than others and were these done for only one group? Does the

	model work worse for certain races or age groups? Several metrics can evaluate this (e.g., see packages AIF360, fairlearn, and EqualityML). - Use expert ratings to evaluate any potential sources of bias. - Understandability: are the explanations understandable for the engineer, the clinician, and/or the patient?
If bias is detected	- Use bias mitigation strategies either during pre-processing (removing variables generating the bias along with variables correlated with these ones), training (adversarial debiasing, prejudice remover), or evaluation (equalized odds, reject option classification). See packages AIF360, fairlearn, and EqualityML.
After deployment	- Continuous assessment: we need to review predictions and re-assess accuracy once deployed as new environments and populations could change performance (i.e., dataset shift (46)).

535

536 **Explaining acoustic features relevant to detecting vocal fold paralysis**

537 Objective acoustic measurement changes associated with vocal fold paralysis have been
538 described and these changes include reduced loudness and maximum phonation time,
539 higher perturbation measurements such as jitter and shimmer, and increased signal to
540 noise ratio (17,47,48); however these were univariate models, and we have demonstrated
541 that using single variables does not seem to provide high predictive performance. While
542 other multivariate machine learning models have been used, these used few features and
543 small or undefined samples and only report feature importance results for one model;
544 therefore it is not clear whether the important features reported would hold using larger
545 feature sets or how other models would perform. Using a much larger initial set of acoustic
546 features for analysis, we demonstrate that several machine learning algorithms of
547 increasing complexity (using more parameters) identify vocal fold paralysis from healthy
548 voices. We also report that these models can use different features to achieve similar
549 performance. Different models emphasize different features not simply because of its
550 relevance to a disorder, but because of the mathematics associated with the model (e.g.,

551 containing different degrees of interaction effects, regularization, or propensity to
552 underfitting or overfitting) (49). The variability of the ranking of features used by our
553 individual models also illustrates the potential danger of using the single highest
554 performing model, which is commonly seen in published literature.

555 Instead of simply reporting the important features from the highest performing model, we
556 analyzed the models to find common features. Some of the most important features across
557 models were: intensity (especially equivalent sound pressure level which was redundant
558 with multiple loudness features and seems to be due to some patients trying to use more
559 breath for projection or being recorded with a higher microphone gain), Mel Frequency
560 Cepstral Coefficients (especially the first coefficient, which captures spectral envelope or
561 slope), mean F0 semitones (given F0 originates from vocal-fold oscillation, a vocal-fold
562 paralysis is expected to alter F0), mean F1 amplitude and frequency (influenced by how
563 the vocal tract filters F0 and the shape of the glottal pulse which would be affected by
564 UVFP), and voiced and unvoiced segments (prosodic and speech articulation features
565 which may be altered due to changes in the periodicity of F0). Shimmer variability was
566 important just for reading, and it captures variability in glottal pulses and pressure patterns
567 which ultimately affect F0. When we removed these top 5 features from the full feature set,
568 performance is practically equivalent to using 88 features, as expected, since there are
569 features that are redundant with these top 5 features. Therefore, it is not that only these 5
570 specific features drive performance, but rather the information they contain, which in this
571 dataset is also captured by other features as shown in Figure 6.

572 These acoustic features corroborate our clinical understanding of glottal incompetence
573 from UVFP and with common patient complaints of reduced loudness, vocal instability,
574 hoarseness, and rough voice. Uncovering and understanding the basic mechanisms and
575 features that models use to generate predictions and outcomes are important as these
576 tools become part of the clinical decision making process.

577 **Identifying and addressing bias**

578 Equivalent Sound Level was higher in UVFP patients than controls. This is counter-intuitive
579 because UVFP patients are known to have softer voices as already described; however,
580 clinicians rated most UVFP samples as being louder than controls. The bias discovered
581 was likely due to increasing the gain on the microphone for some UVFP patients, which
582 would explain the increased background noise in UVFP patients' recordings. A second
583 source of bias may have occurred from requesting UVFP patients to speak louder in order
584 to meet the minimum intensity threshold on the recording softwares Computerized Speech
585 Lab™ and OperaVOX, or patients could have tried this on their own knowing they were
586 being recorded. This behavioral compensation is likely to occur in biomarker research
587 when the participant has a soft voice, especially in retrospective studies like ours where
588 the study goal is not known at the time of recording or when certain software properties
589 lead individuals with weak voices to speak louder. Even though the current models perform
590 better than the clinicians, a systematic comparison would require more clinician and model
591 assessments across datasets. It is likely a model trained on a single dataset might learn

592 intrinsic characteristics of that dataset that do not generalize as well as clinical expertise
593 might.

594 Having said this, this line of research would help us understand the extent to which UVFP
595 detection is generalizable from acoustic data alone. Finding an objective measure of
596 hoarseness is important given a "normal voice" is a fundamentally subjective classification
597 that is not well defined (50,51) and varies with training (52,53), which may result in low
598 reliability of evaluation of disordered voices among clinical rating scales (54).

599 As a post hoc analysis, we address bias by trying to mitigate its effect: we removed
600 variables associated with intensity variables. After removing these features, the models
601 were able to obtain similar performance using a very different set of features. It is possible
602 that these remaining features better reflect pathophysiology or that the features extracted
603 are still influenced by intensity, but further studies should address their generalizability or
604 their relation to intensity variation.

605 **Evaluating the sensitivity to tasks, model complexity, and features used**

606 In addition to getting a better understanding of features, we explored performance in the
607 context of different vocal tasks. Participants carried out two different tasks to elicit voice,
608 *reading*, which captures more complex speech dynamics, and *sustaining vowels*, which is
609 a simpler measure of vocalization and the respiratory subsystem. Overall, these dynamics
610 from the speech task may have improved model performance as was observed.

611 Comparing simpler and more complex models is important because simpler models such
612 as Logistic Regression could be preferred because they tend to generalize better given

613 they are less at risk for overfitting the training set and they are more interpretable and thus
614 biases can be assessed more directly (55).

615 By removing redundant features, we can concentrate on finding the most useful features
616 for further analysis. Performance decreased only slightly while we made models more
617 parsimonious and explainable. This approach is key given the curse of dimensionality in
618 machine learning that may make models unnecessarily complex and harder to generalize
619 (18).

620 Often studies will report the top N features but not how predictive they are in isolation. In
621 our study we ran models on the top 5 features together (Table 2). The lower performance
622 of these top 5 features relative to a richer feature set helps demonstrate that model
623 performance is dependent on interactions across multiple additional features (with the
624 exception of samples from the reading task which obtained an AUC of 0.86 using just the 5
625 features). We also ran models without top 5 features to demonstrate that leaving features
626 that are redundant with these top features results in almost equivalent high performance to
627 using all 88 features since the redundant features share information. Furthermore, when
628 training models on the individual features from within these top 5 one at a time, the
629 performance was reduced considerably with scores from 0.55 to 0.71. This indicates the
630 need for these models to combine multiple features to achieve high performance and any
631 model evaluation should not focus on only the common or top features without testing their
632 predictive performance.

633 **Limitations and future directions**

634 We cannot determine how the bias will affect the model's performance on future samples,
635 but it will likely underperform in samples where length was not different between groups,
636 where gain cannot be changed, and where participants are instructed to not overproject
637 their voice; however, it is possible the model could underperform for other reasons
638 including dataset shift (e.g., the distribution of voice characteristics or demographics is
639 different in a new sample).

640 The classification using just duration itself varied across models and clinicians who
641 listened to the reading passage in its entirety did not achieve as good a classification as
642 the best performing models. Duration itself was not included as a feature in the eGeMaps-
643 based models and has a complex effect on both machines and humans. For example,
644 duration could have affected eGeMaps features (e.g, introduce more variability to the
645 functionals that are computed over sliding time windows) and duration of vowels varied
646 extensively across the UVFP group thus cannot itself be tied to underlying
647 pathophysiology. Therefore, important future work should analyze how duration may affect
648 these features, should address the intrinsic variability in durations of UVFP patients in
649 responding to speech items, and should incorporate models of production that include a
650 consideration of respiratory capabilities, articulation changes, and vocal fold
651 pathophysiology.

652 It is not clear whether these models could detect UVFP from other voice disorders or just
653 healthier voices; however, a model that generalizes well in classifying UVFP from controls

654 could be used to monitor UVFP patients remotely and affordably during treatment or detect
655 risk for UVFP when it is the most likely cause (e.g., dysphonia after thyroid surgery).
656 Larger sample sizes with curated examinations can help increase diverse representation
657 across voice quality and thereby potentially reduce bias in classifier performance. We did
658 not analyze potential racial bias given this data was not extracted from the chart review.
659 Our choice of a standardized feature set worked well in this setting, but may fail to work for
660 differential voice disorder diagnosis or when generalizing to larger datasets, which may
661 bring in additional sources of variance unaccounted for in this dataset. With the availability
662 of more data, additional features could be extracted that better capture changes in
663 coordination (e.g., XCORR (56)) or speech rate (i.e., given UVFP patients may speak
664 slower). Furthermore, while our feature importance evaluation method, SHAP, shows a
665 certain amount of robustness across models, alternative model-agnostic feature-
666 importance methods (e.g., LOFO, permutation importance) as well as model-specific
667 methods (coefficient values for linear models, mean decrease in impurity for Random
668 Forest) could be compared. Model understandability –how easily are the explanations
669 understood by a speech scientist or a clinician– could be assessed rigorously (57). Finally,
670 debiasing the models by removing features correlated with the biased ones was attempted
671 although it is not clear how exactly intensity may influence certain features; we assume if
672 intensity is influencing a variable, it generally should create some considerable association
673 which we discarded using dcor. Therefore, the effect of the bias can be assessed by
674 testing the model's generalizability to new unbiased datasets.

675 **Conclusion**

676 Using one of the largest UVFP datasets to date, our study demonstrates the importance of
677 checking for biases using explainable machine learning and clinician perceptual ratings. In
678 order to first explain models, we tackle collinearity (i.e., redundant or highly correlated
679 independent variables), which biases feature importance, using a custom method called
680 Independence Factor that selects one out of a set of associated features without losing
681 predictive performance. We then compare how results change across different speech-
682 eliciting tasks, training algorithms, features, features set sizes, and highest and lowest
683 performing features to better understand the process that models use to predict vocal
684 changes associated with laryngeal disease, since analyzing a single model will result in a
685 biased view of how predictions are achieved. During this process, we discovered there
686 was a difference in audio duration between groups clearly not related to intrinsic
687 differences in UVFP speech rate, but in cropping all control recordings to a certain length
688 during audio storage. We also discovered that sound equivalent level was
689 counterintuitively higher in UVFP patients, a likely bias resulting from the weak or
690 underprojected voice that characterizes many UVFP patients: patients were prompted by
691 the recording software to speak louder and the microphone gain was likely raised
692 selectively for these patients with weaker voices, possibly generating higher background
693 noise which was detected through clinician's ratings; therefore the models picked up on
694 the acoustic correlates of this increased intensity, which would impede generalization
695 under different recording procedures and natural audio durations. This is more likely to
696 occur in laryngology datasets when patients have a softer voice.

697 Interestingly, we found that matching audio duration between groups and removing all
698 variables that were clearly related to intensity (e.g., bias mitigation) resulted in similar high
699 performance. In this case, the model may be using information more related to
700 pathophysiology, which would need to be further confirmed by future unbiased samples.
701 Machine learning models tended to surpass clinician's evaluation of the same audio
702 recordings. Interestingly, using clinician's voice quality ratings on the recordings in
703 machine learning models performed better than their binary evaluation on whether
704 recordings contained a sample of UVFP voice or not.

705 We hope to promote moving beyond using a single model and only reporting top features
706 to a better explanation of how these models work as well as being able to understand
707 variance across modeling and evaluation choices. We believe these are all aspects of
708 machine learning that clinicians need to understand prior to using such applications.

709 With these considerations along with the recommendations we make, machine learning
710 applications could aid in laryngology screening, allowing for the potential development of
711 in-home screening assessments and continuous pre- and post-treatment monitoring.

712 **Acknowledgments**

713 We would like to thank Cody Sullivan and Carolyn Hsu for their help in rating the audio
714 samples and thank Daryush Mehta, Robert Hillman, and John Guttag for their feedback on
715 an earlier version of this study. DML was supported by a National Institute on Deafness
716 and Other Communication Disorders T32 training grant [5T32DC000038-28], a RallyPoint
717 Fellowship, and an Amelia Peabody Professional Development Award. The work was

718 supported by a gift to the McGovern Institute for Brain Research at MIT. SSG was partially
719 supported by National Institutes of Health grants for the development of pydra-ml [R01
720 EB020740], for reproducible practices [P41 EB019936], and the Bridge2AI voice data
721 generation project [1OT2OD032720-01]. The authors declare that there is no conflict of
722 interest.

723

724 **Data Availability Statement**

725 All data and code are available through Github (<https://github.com/danielmlow/vfp>) and
726 Zenodo (<https://doi.org/10.5281/zenodo.5009208>) including a tutorial to test our models on
727 your own data (https://github.com/danielmlow/vfp/blob/main/vfp_detector.ipynb).

728

729 **Author Contributions**

730 Daniel M. Low: Data curation, Methodology, Formal analysis, Software, Writing - Original
731 Draft; Vishwanatha Rao: Data Curation, Formal analysis, Writing - Original Draft; Gregory
732 Randolph: Writing - Review & Editing; Philip C. Song: Conceptualization, Methodology,
733 Writing - Original Draft, Supervision, Data curation; Satrajit S. Ghosh: Conceptualization,
734 Methodology, Writing - Original Draft, Supervision, Software

735

736 **References**

737 1. Wroge TJ, Özkanca Y, Demiroglu C, Si D. Parkinson's disease diagnosis using machine learning and
738 voice. 2018 IEEE signal [Internet]. 2018; Available from:
739 [https://ieeexplore.ieee.org/abstract/document/8615607?casa_token=qI93B6R4GIYAAAAA:_IQjuRrle_k](https://ieeexplore.ieee.org/abstract/document/8615607?casa_token=qI93B6R4GIYAAAAA:_IQjuRrle_kQ01FDfUiOzPAXp2Gb8sHtO9NeMDjF3yhJqMO7MQoXWgb6jtZbP6SfQADgNzxdk0Kt)
740 [Q01FDfUiOzPAXp2Gb8sHtO9NeMDjF3yhJqMO7MQoXWgb6jtZbP6SfQADgNzxdk0Kt](https://ieeexplore.ieee.org/abstract/document/8615607?casa_token=qI93B6R4GIYAAAAA:_IQjuRrle_kQ01FDfUiOzPAXp2Gb8sHtO9NeMDjF3yhJqMO7MQoXWgb6jtZbP6SfQADgNzxdk0Kt)

- 741 2. Low DM, Bentley KH, Ghosh SS. Automated assessment of psychiatric disorders using speech: A
742 systematic review. *Laryngoscope Investig Otolaryngol*. 2020 Feb;5(1):96–116.
- 743 3. Quatieri TF. *Discrete-Time Speech Signal Processing: Principles and Practice*. Pearson Education;
744 2008. 816 p.
- 745 4. Molnar C. *Interpretable Machine Learning*. Lulu.com; 2019. 319 p.
- 746 5. Stachler RJ, Francis DO, Schwartz SR, Damask CC, Digoy GP, Krouse HJ, et al. Clinical practice
747 guideline: Hoarseness (dysphonia) (update). *Otolaryngol Head Neck Surg*. 2018 Mar;158(1_suppl):S1–
748 42.
- 749 6. Sritharan N, Chase M, Kamani D. The vagus nerve, recurrent laryngeal nerve, and external branch of
750 the superior laryngeal nerve have unique latencies allowing for intraoperative documentation of The
751 [Internet]. 2015; Available from:
752 [https://onlinelibrary.wiley.com/doi/abs/10.1002/lary.24781?casa_token=JxEJ7oj9-](https://onlinelibrary.wiley.com/doi/abs/10.1002/lary.24781?casa_token=JxEJ7oj9-aoAAAAA:L5Y3SkZqJRFDBni7dC0KFfIXPulpDWCYTmNZPfdGtCGVriXPohVZkC2ER_9wgl-DTtvv7xabx2JeBTPa7g)
753 [aoAAAAA:L5Y3SkZqJRFDBni7dC0KFfIXPulpDWCYTmNZPfdGtCGVriXPohVZkC2ER_9wgl-](https://onlinelibrary.wiley.com/doi/abs/10.1002/lary.24781?casa_token=JxEJ7oj9-aoAAAAA:L5Y3SkZqJRFDBni7dC0KFfIXPulpDWCYTmNZPfdGtCGVriXPohVZkC2ER_9wgl-DTtvv7xabx2JeBTPa7g)
754 [DTtvv7xabx2JeBTPa7g](https://onlinelibrary.wiley.com/doi/abs/10.1002/lary.24781?casa_token=JxEJ7oj9-aoAAAAA:L5Y3SkZqJRFDBni7dC0KFfIXPulpDWCYTmNZPfdGtCGVriXPohVZkC2ER_9wgl-DTtvv7xabx2JeBTPa7g)
- 755 7. Randolph GW, Kamani D. The importance of preoperative laryngoscopy in patients undergoing
756 thyroidectomy: voice, vocal cord function, and the preoperative detection of invasive thyroid malignancy.
757 *Surgery*. 2006 Mar;139(3):357–62.
- 758 8. Colton RH, Paseman A, Kelley RT, Stepp D, Casper JK. Spectral moment analysis of unilateral vocal
759 fold paralysis. *J Voice*. 2011 May;25(3):330–6.
- 760 9. Balasubramaniam RK, Bhat JS, Fahim S 3rd, Raju R 3rd. Cepstral analysis of voice in unilateral
761 adductor vocal fold palsy. *J Voice*. 2011 May;25(3):326–9.
- 762 10. Little M, Costello D, Harries M. Objective dysphonia quantification in vocal fold paralysis: comparing
763 nonlinear with classical measures. *Nature Precedings*. 2009 Apr 21;1–1.
- 764 11. Bielowicz S, Stager SV. Diagnosis of unilateral recurrent laryngeal nerve paralysis: laryngeal
765 electromyography, subjective rating scales, acoustic and aerodynamic measures. *Laryngoscope*. 2006
766 Mar;116(3):359–64.
- 767 12. Hartl DAM, Hans S, Vaissière J, Brasnu DAMF. Objective acoustic and aerodynamic measures of
768 breathiness in paralytic dysphonia. *Eur Arch Otorhinolaryngol*. 2003 Apr;260(4):175–82.
- 769 13. Francis DO, Pearce EC, Ni S, Garrett CG, Penson DF. Epidemiology of vocal fold paralyses after total
770 thyroidectomy for well-differentiated thyroid cancer in a Medicare population. *Otolaryngol Head Neck*
771 *Surg*. 2014 Apr;150(4):548–57.
- 772 14. Jeannon JP, Orabi AA, Bruch GA, Abdalsalam HA, Simo R. Diagnosis of recurrent laryngeal nerve palsy
773 after thyroidectomy: a systematic review. *Int J Clin Pract*. 2009 Apr;63(4):624–9.
- 774 15. Bhattacharyya N, Kotz T, Shapiro J. Dysphagia and aspiration with unilateral vocal cord immobility:
775 incidence, characterization, and response to surgical treatment. *Ann Otol Rhinol Laryngol*. 2002
776 Aug;111(8):672–9.
- 777 16. Pinho CMR, Jesus LMT, Barney A. Aerodynamic measures of speech in unilateral vocal fold paralysis
778 (UVFP) patients. *Logoped Phoniatr Vocol*. 2013 Apr;38(1):19–34.
- 779 17. Hartl DM, Crevier-Buchman L, Vaissière J, Brasnu DF. Phonetic effects of paralytic dysphonia. *Ann Otol*
780 *Rhinol Laryngol*. 2005 Oct;114(10):792–8.

- 781 18. Berisha, V., Krantsevich, C., Hahn, P. R., Hahn, S., Dasarathy, G., Turaga, P., & Liss, J. Digital
782 medicine and the curse of dimensionality. *NPJ Digital Medicine*. 2021 Dec;4(1):s41746–021.
- 783 19. Ruz J, Švihlík J, Krýže P, Novotný M, Tykalová T. Reproducibility of Voice Analysis with Machine
784 Learning. *Mov Disord*. 2021 May;36(5):1282–3.
- 785 20. Schönweiler R, Hess M, Wübbelt P, Ptok M. Novel approach to acoustical voice analysis using artificial
786 neural networks. *J Assoc Res Otolaryngol*. 2000 Dec;1(4):270–82.
- 787 21. Godino-Llorente JI, Gómez-Vilda P. Automatic detection of voice impairments by means of short-term
788 cepstral parameters and neural network based detectors. *IEEE Trans Biomed Eng*. 2004
789 Feb;51(2):380–4.
- 790 22. Fraile R, Saenz-Lechon N, Godino-Llorente JI, Osma-Ruiz V, Fredouille C. Automatic detection of
791 laryngeal pathologies in records of sustained vowels by means of mel-frequency cepstral coefficient
792 parameters and differentiation of patients by sex. *Folia Phoniatr Logop*. 2009;61(3):146–52.
- 793 23. Voigt D, Döllinger M, Yang A, Eysholdt U, Lohscheller J. Automatic diagnosis of vocal fold paresis by
794 employing phonovibrograph features and machine learning methods. *Comput Methods Programs
795 Biomed*. 2010 Sep;99(3):275–88.
- 796 24. Lopes LW, Batista Simões L, Delfino da Silva J, da Silva Evangelista D, da Nóbrega E Ugulino AC,
797 Oliveira Costa Silva P, et al. Accuracy of Acoustic Analysis Measurements in the Evaluation of Patients
798 With Different Laryngeal Diagnoses. *J Voice*. 2017 May;31(3):382.e15–382.e26.
- 799 25. Powell ME, Rodriguez Cancio M, Young D, Nock W, Abdelmessih B, Zeller A, et al. Decoding phonation
800 with artificial intelligence (DeP AI): Proof of concept. *Laryngoscope Investig Otolaryngol*. 2019
801 Jun;4(3):328–34.
- 802 26. Dibazar AA, Narayanan S, Berger TW. Feature analysis for automatic detection of pathological speech.
803 In: *Proceedings of the Second Joint 24th Annual Conference and the Annual Fall Meeting of the
804 Biomedical Engineering Society* [Engineering in Medicine and Biology. 2002. p. 182–3 vol.1.
- 805 27. Seedat N, Aharonson V, Hamzany Y. Automated and interpretable m-health discrimination of vocal cord
806 pathology enabled by machine learning. In: *2020 IEEE Asia-Pacific Conference on Computer Science
807 and Data Engineering (CSDE)*. 2020. p. 1–6.
- 808 28. Mittal V, Sharma RK. Deep Learning Approach for Voice Pathology Detection and Classification. *IJHISI*.
809 2021 Oct 1;16(4):1–30.
- 810 29. Hu HC, Chang SY, Wang CH, Li KJ, Cho HY, Chen YT, et al. Deep Learning Application for Vocal Fold
811 Disease Prediction Through Voice Recognition: Preliminary Development Study. *J Med Internet Res*.
812 2021 Jun 8;23(6):e25247.
- 813 30. Fairbanks G. *Voice and Articulation Drillbook*. Harper; 1960. 196 p.
- 814 31. Eyben F, Scherer KR, Schuller BW, Sundberg J, André E, Busso C, et al. The Geneva Minimalistic
815 Acoustic Parameter Set (GeMAPS) for Voice Research and Affective Computing. *IEEE Transactions on
816 Affective Computing*. 2016 Apr;7(2):190–202.
- 817 32. audEERING GmbH. openSMILE (Version 2.3) [Internet]. 2017. Available from:
818 [https://github.com/naxingyu/opensmile/blob/3a0968e7b36c1b730a4ffd2977031091ee9abf](https://github.com/naxingyu/opensmile/blob/3a0968e7b36c1b730a4ffd2977031091ee9abf7f/config/gemaps/eGeMAPSv01a.conf)
819 [7f/config/gemaps/eGeMAPSv01a.conf](https://github.com/naxingyu/opensmile/blob/3a0968e7b36c1b730a4ffd2977031091ee9abf7f/config/gemaps/eGeMAPSv01a.conf)
- 820 33. S S Ghosh, D M Low, H Rajaei et al. Pydra-ML doi:10.5281/ZENODO.4170850 [Internet]. Available

- 821 from: <https://github.com/nipytype/pydra-ml>
- 822 34. Raschka S. Model Evaluation, Model Selection, and Algorithm Selection in Machine Learning [Internet].
823 arXiv [cs.LG]. 2018. Available from: <http://arxiv.org/abs/1811.12808>
- 824 35. Ojala M, Garriga GC. Permutation Tests for Studying Classifier Performance. In: 2009 Ninth IEEE
825 International Conference on Data Mining. IEEE; 2009. p. 1833–63.
- 826 36. Lundberg S, Lee SI. A Unified Approach to Interpreting Model Predictions [Internet]. arXiv [cs.AI]. 2017.
827 Available from: <http://arxiv.org/abs/1705.07874>
- 828 37. D'Amour A, Heller K, Moldovan D, Adlam B, Alipanahi B, Beutel A, et al. Underspecification presents
829 challenges for credibility in modern machine learning. *J Mach Learn Res.* 2022 Jan 1;23(1):10237–97.
- 830 38. de Siqueira Santos S, Takahashi DY, Nakata A, Fujita A. A comparative study of statistical methods
831 used to identify dependencies between gene expression signals. *Brief Bioinform.* 2014 Nov;15(6):906–
832 18.
- 833 39. Székely GJ, Rizzo ML, Bakirov NK. Measuring and testing dependence by correlation of distances.
834 2007; Available from: <https://projecteuclid.org/journals/annals-of-statistics/volume-35/issue-6/Measuring-and-testing-dependence-by-correlation-of-distances/10.1214/009053607000000505.full>
835
- 836 40. Hillenbrand J, Houde RA. Acoustic correlates of breathy vocal quality: dysphonic voices and continuous
837 speech. *J Speech Hear Res.* 1996 Apr;39(2):311–21.
- 838 41. Murton O, Hillman R, Mehta D. Cepstral Peak Prominence Values for Clinical Voice Evaluation. *Am J*
839 *Speech Lang Pathol.* 2020 Aug 4;29(3):1596–607.
- 840 42. G. Degottex, J. Kane, T. Drugman, T. Raitio and S. Scherer. COVAREP—A collaborative voice analysis
841 repository for speech technologies. *Proc IEEE Int Conf Acoust Speech Signal Process.* 2014.
- 842 43. Hallgren KA. Computing Inter-Rater Reliability for Observational Data: An Overview and Tutorial. *Tutor*
843 *Quant Methods Psychol.* 2012;8(1):23–34.
- 844 44. Gamer M, Lemon J, Gamer MM, Robinson A, Kendall's W. Package "irr." Various coefficients of
845 interrater reliability and agreement. 2012;22:1–32.
- 846 45. Cicchetti DV. Guidelines, criteria, and rules of thumb for evaluating normed and standardized
847 assessment instruments in psychology. *Psychol Assess.* 1994 Dec;6(4):284–90.
- 848 46. Dockès J, Varoquaux G, Poline JB. Preventing dataset shift from breaking machine-learning biomarkers.
849 *Gigascience* [Internet]. 2021 Sep 28;10(9). Available from: <http://dx.doi.org/10.1093/gigascience/giab055>
- 850 47. Ramig LA, Scherer RC, Titze IR, Ringel SP. Acoustic analysis of voices of patients with neurologic
851 disease: rationale and preliminary data. *Ann Otol Rhinol Laryngol.* 1988 Mar-Apr;97(2 Pt 1):164–72.
- 852 48. Morsomme D, Jamart J, Wéry C, Giovanni A, Remacle M. Comparison between the GIRBAS Scale and
853 the Acoustic and Aerodynamic Measures Provided by EVA for the Assessment of Dysphonia following
854 Unilateral Vocal Fold Paralysis. *Folia Phoniatri Logop.* 2001 Nov-Dec;53(6):317–25.
- 855 49. Kriegeskorte N, Douglas PK. Interpreting encoding and decoding models. *Curr Opin Neurobiol.* 2019
856 Apr;55:167–79.
- 857 50. Misono S. The Voice and the Larynx in Older Adults: What's Normal, and Who Decides? *JAMA*
858 *Otolaryngol Head Neck Surg.* 2018 Jul 1;144(7):572–3.

- 859 51. Eadie T, Sroka A, Wright DR, Merati A. Does knowledge of medical diagnosis bias auditory-perceptual
860 judgments of dysphonia? *J Voice*. 2011 Jul;25(4):420–9.
- 861 52. Helou LB, Solomon NP, Henry LR, Coppit GL, Howard RS, Stojadinovic A. The role of listener
862 experience on Consensus Auditory-perceptual Evaluation of Voice (CAPE-V) ratings of
863 postthyroidectomy voice. *Am J Speech Lang Pathol*. 2010 Aug;19(3):248–58.
- 864 53. Eadie TL, Baylor CR. The effect of perceptual training on inexperienced listeners' judgments of
865 dysphonic voice. *J Voice*. 2006 Dec;20(4):527–44.
- 866 54. Karnell MP, Melton SD, Childes JM, Coleman TC, Dailey SA, Hoffman HT. Reliability of clinician-based
867 (GRBAS and CAPE-V) and patient-based (V-RQOL and IPVI) documentation of voice disorders. *J*
868 *Voice*. 2007 Sep;21(5):576–90.
- 869 55. Rudin C. Stop Explaining Black Box Machine Learning Models for High Stakes Decisions and Use
870 Interpretable Models Instead. *Nat Mach Intell*. 2019 May;1(5):206–15.
- 871 56. Williamson JR, Quatieri TF, Helfer BS, Ciccarelli G, Mehta DD. Vocal and Facial Biomarkers of
872 Depression based on Motor Incoordination and Timing. In: *Proceedings of the 4th International*
873 *Workshop on Audio/Visual Emotion Challenge*. New York, NY, USA: Association for Computing
874 Machinery; 2014. p. 65–72. (AVEC '14).
- 875 57. Zhou Y, Ribeiro MT, Shah J. ExSum: From Local Explanations to Model Understanding [Internet]. arXiv
876 [cs.CL]. 2022. Available from: <http://arxiv.org/abs/2205.00130>
- 877
878
879

880 **Figure 1. Schematic of speech production and the process of extracting certain acoustic features**
881 **from an audio signal.** (A) Speech production, (B) recording characteristics, (C) waveform of audio signal
882 with fundamental frequency (f_0), (D) spectrogram with formants F1-F3 and intensity, (E) mel-frequency
883 cepstral coefficients (MFCCs). Full description in the main text.
884

885 **Figure 2. Distribution of audio duration for reading and vowel tasks split by group reveals a dataset**
886 **bias.** The mode of the audio durations for the controls is 3.5 s for reading samples and 4.11 s for vowel
887 samples.

888

889 **Figure 3. Model performance comparison using a permutation test using non-redundant features. (A)**
890 Scores from models trained on true labels (blue) and trained on permuted labels (orange) over bootstrapping
891 splits. **(B)** Statistical comparison between models (annotation = p-value, highlighted = significant results).

892

893 **Figure 4. Feature importance parallel coordinate plot.** Rank reads from bottom (most important) to top
894 (least important). Mean rank is weighted by performance of each model to avoid a lower performing model
895 biasing the mean rank.

896

897 **Figure 4. Distributions for top 5 features and corresponding performance for single features.** Logistic
898 Regression with L1 penalty was used. No single feature is enough to dissociate groups with high
899 performance. Null models' median performance was 0.5.

900 **Figure 6. Feature redundancy with top 5 features highlighted.** Top 5 features are highlighted in bold and
901 their rank is displayed. Squares are clusters of redundant features. Computed with all participants on the
902 reading task.

903 **Figure 7. Descriptive statistics and inter-rater reliability of clinician ratings for unilateral vocal fold**
904 **paralysis (UVFP), background noise, and recording loudness indicating likely bias.** Error bars indicate
905 maximum and minimum count across the three raters. The disproportionate amount of UVFP samples rated
906 as having high background noise and high loudness indicates likely bias, where the gain might have been
907 raised for some UVFP patients and they may have phonated more intensely. kappa: Light's kappa; ICC:
908 intra-class correlation coefficient.

909 **Figure 8. How clinicians rate the audio recordings of read speech: descriptive statistics and inter-**
910 **rater reliability of average clinician ratings.** The average across raters was taken for each recording. ICC:
911 intra-class correlation coefficient.