



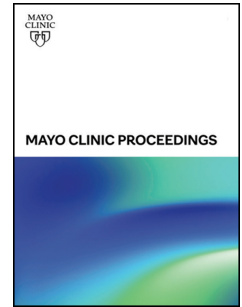
Since January 2020 Elsevier has created a COVID-19 resource centre with free information in English and Mandarin on the novel coronavirus COVID-19. The COVID-19 resource centre is hosted on Elsevier Connect, the company's public news and information website.

Elsevier hereby grants permission to make all its COVID-19-related research that is available on the COVID-19 resource centre - including this research content - immediately available in PubMed Central and other publicly funded repositories, such as the WHO COVID database with rights for unrestricted research re-use and analyses in any form or by any means with acknowledgement of the original source. These permissions are granted for free by Elsevier for as long as the COVID-19 resource centre remains active.

Journal Pre-proof

Guess What We Can Hear – Novel Voice Biomarkers for the Remote Detection of Disease

Jaskanwal Deep Singh Sara, MBChB, Diana Orbelo, PhD, Elad Maor, MD, PhD, Lilach O. Lerman, MD, PhD, Amir Lerman, MD



PII: S0025-6196(23)00130-1

DOI: <https://doi.org/10.1016/j.mayocp.2023.03.007>

Reference: JMCP 4115

To appear in: *Mayo Clinic Proceedings*

Received Date: 25 July 2022

Revised Date: 8 February 2023

Accepted Date: 16 March 2023

Please cite this article as: Singh Sara JD, Orbelo D, Maor E, Lerman LO, Lerman A, Guess What We Can Hear – Novel Voice Biomarkers for the Remote Detection of Disease, *Mayo Clinic Proceedings* (2023), doi: <https://doi.org/10.1016/j.mayocp.2023.03.007>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2023 Published by Elsevier Inc on behalf of Mayo Foundation for Medical Education and Research

Guess What We Can Hear – Novel Voice Biomarkers for the Remote Detection of Disease

Jaskanwal Deep Singh Sara, MBChB¹, Diana Orbelo PhD^{2,3}, Elad Maor MD, PhD⁴, Lilach O. Lerman MD, PhD⁵, Amir Lerman MD^{1*}

¹Department of Cardiovascular Medicine, Mayo Clinic College of Medicine, Rochester, MN, USA

²Division of Laryngology, Mayo Clinic College of Medicine, Rochester, MN, USA

³Chaim Sheba Medical Center, Tel Hashomer, Israel

⁴Sackler School of Medicine, Tel Aviv University, Tel Aviv, Israel

⁵Division of Nephrology and Hypertension, Mayo Clinic, Rochester, MN, USA

Acknowledgment: The authors would like to thank Joseph R Duffy, PhD for his insightful comments during the preparation of this manuscript.

Disclosure: All authors have no relevant conflicts of interest to disclose

Address for Correspondence:

Amir Lerman, MD

Department of Cardiovascular Medicine

Mayo Clinic College of Medicine, 200 First Street SW, Rochester, MN, 55905, USA

Phone: 507-255-4152, Fax: 507-255-2550

E-mail address: lerman.amir@mayo.edu

Keywords

Voice biomarker, Voice signal analysis, Telemedicine, Digital health

Abbreviations

AUC, Area under the curve; CAD, Coronary artery disease; COVID-19, Coronavirus-2019; CVD, Cardiovascular disease; MFCCs, Mel Frequency Cepstral Coefficients; PTSD, Posttraumatic stress disorder

Journal Pre-proof

Abstract

The advancement of digital biomarkers and the provision of remote healthcare has greatly progressed during the coronavirus-19 global pandemic. Combining voice/speech data with artificial intelligence and machine-based learning offers a novel solution to the growing demand for telemedicine. Voice biomarkers, obtained from the extraction of characteristic acoustic and linguistic features, are associated with a variety of diseases and even coronavirus disease-2019. In the current review we i) describe the basis upon which digital voice biomarkers could facilitate “telemedicine,” ii) discuss potential mechanisms that may explain the association between voice biomarkers and disease, iii) offer a novel classification system to conceptualize voice biomarkers depending upon different methods for recording and analyzing voice/speech samples, iv) outline evidence demonstrating an association between voice biomarkers and a number of disease states, and v) describe the process of developing a voice biomarker from recording and storing voice samples and extracting relevant acoustic and linguistic features to training and testing deep and machine based learning algorithms to detect disease. We further explore several important future considerations in this area of research, including the necessity for clinical trials and importance of safeguarding data and individual privacy. To this end, we searched Pubmed and Google Scholar to identify studies evaluating the relationship between voice/speech features and biomarkers and various diseases. Search terms included “digital biomarker,” “telemedicine,” “voice features,” “voice biomarker,” “speech features,” “speech biomarkers,” “acoustics,” “linguistics,” “cardiovascular disease,” “neurologic disease,” “psychiatric disease,” and “infectious disease.” The search was limited to studies published in English in peer-reviewed journals between 1980 and the present day. To identify potential studies not captured by our database search strategy, we also searched studies listed in the bibliography of relevant publications and reviews.

KEY POINTS

- Voice biomarkers offer promise in delivering healthcare remotely by providing adjunctive information that can aid healthcare providers in making decisions about the detection, progression and treatment of disease, and could be used as surrogate endpoints in clinical trials
- The association between voice biomarkers and systemic disease could relate to disease or injury effecting the lungs, vocal folds or articulators, which produce phonation, or the brain and other neurological structures, which influence phonation; or through systemic disease processes effecting the cardiovascular, endocrine or immune systems; or could reflect our psychosocial states
- Voice features and biomarkers are associated with a number of cardiovascular, neurologic and psychiatric diseases and even coronavirus disease-2019
- Classification of voice biomarkers involves outlining the source of vocalization, the method of voice assessment, the method of voice analysis and the location of recordings
- Developing voice biomarkers involves obtaining voice recordings, performing audio pre-processing, undertaking acoustic and linguistic feature extraction and selection, carrying out the training and testing of voice algorithms and evaluating their integration into clinical practice
- For the field of voice biomarkers to evolve there is a need for large prospective clinical trials in demographically and linguistically diverse populations; special attention also needs to be given in defining how to safely integrate voice biomarkers into existing clinical and research work models in a way that prioritizes data protection and individual privacy

Introduction

The use of telemedicine and digital biomarkers has increased over recent years. In this context, the coronavirus-19 (COVID-19) global pandemic has provided a strong impetus for healthcare providers to examine previously accepted practice patterns and to look for alternative ways of delivering healthcare. One consequence of this has been a dramatic shift from in-person care to non-face-to-face healthcare visits,¹ which has laid the foundation for remote clinics forming a more permanent aspect of healthcare delivery.²⁻⁴ Studies have shown the benefits of remotely provided healthcare through event monitors, smart devices and wearables on various disease processes including hypertension and heart failure.^{5,6} Indeed, while traditional healthcare models require in-person evaluations with potentially lengthy visits and costly testing, telemedicine holds the promise of offering simple non-invasive methods of evaluation that can be undertaken remotely from healthcare providers. Non-face-to-face clinical assessments also reduce the risk of transmitting infectious diseases, and are more accessible, time-efficient and cost-effective than face-to-face visits, which are advantageous to healthcare systems that are overloaded or resource-poor.

Voice/Speech Analysis

Voice/speech analysis is an promising non-invasive digital technology that can examine spoken language to provide information on health status and offers promise as a digital biomarker to screen for disease. In this paper, “voice/speech data” refers to the sounds created by vocal fold oscillation, speech and any other sounds created in a human vocal tract. “Language” on the other hand is separate and is created by cognition and linguistic knowledge. Thus, language involves the production of speech sounds into meaningful units (i.e., phonology), the combination of these meaningful units into words (i.e., affixes and base words), and the generation of sentences and paragraphs to communicate thought.

While the vocal tract only produces voice and speech, cognitive and linguistic elements produce language. Generating voice/speech data encompasses a wide range of behaviors, from the production of a simple sound such as a sigh to natural language produced in conversation. Because voice/speech reflects the coordination of a variety of cognitive and motor processes, voice/speech could provide a sensitive snapshot of cognition, tissue integrity, and motor function relevant to many diseases. Indeed, studies show that voice/speech contains acoustic and linguistic features that can be identified by machine learning models to gauge a speaker's behavioral health status.⁷⁻¹¹ Moreover, the collection of data for voice/speech analysis can be relatively simple and engaging for patients and is relatively inexpensive and convenient, requiring just a microphone, a quiet place, and a device to capture audio samples. Voice/speech samples can be collected remotely using a personal device, and, with simple and clear instructions, does not require special training. Furthermore, voice/speech analysis permits repeated and frequent testing reducing measurement error and providing richer and more detailed data across time. Thus, voice/speech-based biomarkers (henceforth referred to as voice biomarkers) could permit sensitive screening for disease, monitoring of progression and response to treatment and, with informed consent, a useful surrogate marker for clinical research.

Possible Mechanisms of Association

Directly attributing specific voice features to disease processes would involve demonstrating the consistent replication of such features in the context of similar pathophysiologic findings, with concurrent evaluation of biologically plausible mechanisms. Treatment or progression of a given disease should correspond with predictable changes in voice features and the identification of reliable changes in the underlying pathophysiology that links the disease to the process of voice generation. Below we outline our current understanding of possible mechanisms for the association between disease and

voice features and biomarkers. For the field of voice/speech analysis-based disease detection to mature large well conducted studies elucidating associations and causal mechanisms will be required to help clinicians and researchers make more sense of these biomarkers.

Organs of Speech

Voice/speech generation is the expression of thoughts and ideas through the transformation of airflow into sound patterns. This process requires the cooperative efforts of the following structures: the lungs, which provide the air to create sound; the larynx, which contains the vocal folds whose vibrations create sound; and vocal tract (the area from the top of the larynx to the lips which contains articulators including tongue, jaw and soft palate) which further shape sounds (**Figure 1**). Thus, there is the **source** (respiratory and laryngeal systems) and the **filter** (resonance and articulation systems). “Voice/speech” therefore refers to the production of sound emanating from air pressure and air flow generated by the respiratory and laryngeal systems, while the vocal tract is responsible for the process of resonance, and the articulators can be positioned to influence resonance and produce speech sounds (i.e., articulation). In other words, sounds come from the larynx, and they are modified by the size and shape of the vocal tract. Abnormalities at any stage in this process may give rise to characteristic voice/speech features that may be detectable by machine and deep learning algorithms. These could include aberrations to the normal structure and function of any of the component organs either from primary disease in organs of phonation, or through secondary organ dysfunction such as central nervous system lesions or degeneration.

Neurocognitive and Neuromotor Function

Figure 1 outlines central contributors to voice/speech aberrations. Ideas form in our mind under the direction of the activity of a number of specific brain networks, such as Broca and Wernicke's areas in the left frontal and temporal cortices, respectively. Further, multiple cranial nerves participate in speech control, including cranial nerves V, VII, IX, X, XI, XII as well as the spinal nerves which control respiration. The vagus nerve in particular plays a critical role in phonation under the direction of the voluntary nervous system (recurrent laryngeal and motor branch of superior laryngeal nerves). Thus, disease or injury directly affecting these or related neurological structures through processes such as ischemia, inflammation, trauma, neoplasm and infection to name a few can provide a plausible basis for changes identified in voice/speech data.

Endocrine and Immune Systems

The link with disease in other systems, such as the endocrinologic and immune systems, is less well understood. Studies have demonstrated a link between systemic immune-mediated diseases, such as rheumatoid arthritis, Sjogren's syndrome and systemic lupus erythematosus and voice disorders with significant correlations between biochemical parameters of inflammation such as sedimentation velocity and anti-DNA antibodies with detectable vocal changes.¹² Endocrinologic disorders such as diabetes mellitus¹³ and hypothyroidism¹⁴ have also been associated with specific voice/speech abnormalities, suggesting that the immunologic and hormonal milieu related to various systemic disorders may be expressed in voice/speech characteristics through secondary effects on organs of phonation. In one study, vocal cord and lung edema were associated with characteristic perturbations in vocal characteristics,¹⁵ which may be underpinned by transudative or exudative causes the latter of which could be caused by inflammatory disorders.

Cardiovascular System

Cardiovascular disease may similarly modify phonation through systemic processes. Atherosclerosis and vascular disease may lead to chronic hypoperfusion of the larynx and other organs of phonation, or to the brain and other neurologic structures that influence voice production. This in turn may lead to structural atrophy, involution, calcification, stiffening and loss of structural integrity and muscle mass.¹⁶ In the larynx this could lead to incomplete closure of the vocal folds, reduced laryngeal mobility, reduced cartilage sliding, decreased mucosal wave vibration and with them, predictable changes in vocal characteristics such as an increase or decrease in fundamental frequency of phonation¹⁶ among others. Moreover, the majority of cardiovascular diseases and their complications are caused by atherosclerotic disease, which is widely considered as a systemic inflammatory process.¹⁷ Thus, a specific cardiovascular inflammatory 'soup' could exert unique secondary effects on the structure and function of organs of phonation resulting in the expression of characteristic voice features. Cardiovascular disease may also play a role in modifying voice features through its intimate relationship with the autonomic nervous system. The vagus nerve innervates the heart through its superior, inferior, and thoracic cardiac branches and its activity is tightly coupled with heart rate control, which are in turn associated with coronary artery disease (CAD)¹⁸ and cardiovascular events.¹⁹ Indeed, both heart failure and CAD are associated with dysregulation in autonomic nervous system activity manifesting with characteristic signs and symptoms of disease. Thus, it follows that similarly stereotypical changes in vocal/speech may also be evident given the key role of the key vagus nerve in phonation.

Mental States

Acoustic and linguistic features may represent an epiphenomenon occurring in parallel to disease rather than being directly caused by disease per se. For example, growing evidence supports the association

between our mental states, and associated bio-psycho-social factors, with disease. Mental stress in particular has been shown to be associated with diabetes mellitus,²⁰ and CVDs including CAD, atrial fibrillation, and stroke.²¹⁻²⁴ Indeed, mental stress is associated with measurable changes in voice/speech characteristics^{25,26} and voice symptoms are common in individuals with high levels of cortisol.²⁷ While these explanations would not account for the high variability of acoustic and linguistic features identified across different diseases mental stress could still be playing a role in the voice/speech analysis of each of these disease states. Mental stress affects our internal dialogue through our cognitions and thoughts, which explains the unique linguistic patterns that can be seen in individuals with depression and other psychiatric diseases.^{28,29} Mental stress is also increasingly being seen to directly affect the neurologic, endocrinologic, immunologic and cardiovascular systems³⁰ through a variety of mechanisms including amygdala and sympathetic nervous system overactivity, upregulation of the hypothalamic-pituitary-adrenal axis, enhanced release of peripherally circulating inflammatory cytokines and greater vascular inflammation, endothelial dysfunction, atherosclerotic disease progression and cardiovascular events.³¹ Thus, mental stress could be a contributing factor to the mechanistic pathways linking a variety of disease states with detectable voice/speech characteristics. More studies will be required to better clarify this and other potential causative mechanisms as we learn more about the potential role of voice biomarkers in clinical practice.

Voice Classification

Currently there is no practical framework for the conceptualization of voice biomarkers that distinguishes the source or characteristics of voice/speech data that can be analyzed and developed into potential biomarkers. These biomarkers may be derived from verbal sounds (speech), vocal sounds, and non-verbal sounds with active or passive tasks (**Figure 2**). On this basis, we propose a novel '4-position'

classification system for the development voice biomarkers. Position I refers to the type of vocalization used for analysis and extraction of a voice biomarker, with S denoting verbal (speech), V denoting vocal, N denoting nonverbal sounds (noise) and D denoting all. Position II refers to the assessment strategy used with A denoting voice recordings during active tasks and processes, P denoting passive voice recordings and D denoting both. Position III refers to the analysis strategy with A denoting analysis of acoustics, L denoting analysis of linguistics and D denoting both. Position IV refers to recordings performed in-person, denoted by I, or remotely, denoted by R (**Table 1**).

Table 2 further elaborates on this classification by providing examples of different methods of voice/speech analysis and different locations and modes of obtaining recordings. “Active” recordings require engaging participants in specifically outlined tasks, whereas “passive” recordings involve free and spontaneous speech and sounds. In our previous studies individuals were asked to obtain three separate 30-second voice recordings, i) reading a prespecified text, ii) describing a positive emotional experience, and iii) describing a negative emotional experience.³²⁻³⁴ This partially structured and semi-spontaneous approach provided the opportunity to assess a combination of acoustic (audible and non-audible sound-related) and linguistic (semantic content-related) features; this offers a balance between “controlling” the vocal task and allowing participants to formulate their own responses. We have previously demonstrated that two features derived from machine learning were independently associated with CAD. These features were “i) a statistical operator describing extreme values of the Mel Frequency Cepstral Coefficient (MFCC), which measures craters around specific frequency bands, and quantifies the width and depth of these craters, and ii) a skewness measure that portrays the asymmetry of intensity distribution in a specific frequency band.”³² Both were more strongly associated when participants were asked to describe a negative emotional experience.³² Hence, unique voice/speech characteristics emerge in the acoustics and/or linguistic content produced by patients

when articulating emotively significant ideas. Special attention should be given to the difficulty in obtaining spontaneous speech samples. Patients can often be reluctant to share positive or negative emotional experiences and should be encouraged to expand on what they are saying. Spontaneous speech samples are valuable when performing voice/speech analyses but are more challenging to obtain than reading passages or sentence repetition. The acoustic signal may also vary more in spontaneous conversational speech compared to controlled vocalization tasks (e.g., sustained vowel sounds). With more voice disruptions and quality fluctuations³⁵ spontaneous speech may be better at eliciting changes in the organization and complexity of language. On the other hand, reading scripted text may capture important changes in acoustic characteristics of speech, as speakers tend to be more fluent in read speech. Reading tasks provide more control over the speech task but require instruction and active engagement. Consideration must also be given to individual reading skill. Moreover, sustained vowel phonations can provide information useful for evaluating dysphonia and assess the integrity of the laryngeal and respiratory systems, while eliminating the influence of confounders such as speaking rate, intonation, dialect, accent, and even language.³⁵ Conversely, diadochokinetic tasks (fast repetition of syllables, e.g., /pa/ta/ka/) require rapid movements of the jaw, lips and tongue challenging individuals to maintain the rate and clarity of speech. These tasks provide the greatest level of control over a task, do not require complex instructions, but still require active subject participation and only provide information for a narrow range of conditions. Thus, the selection of approach for voice/speech recordings and analysis may depend on several factors including resources available and the objective of testing, as different types of voice/speech samples assess the mechanism of voice and speech generation in different ways. For example, if the purpose is to screen for a psychiatric condition, where the content and cadence of speech may be characteristic for specific disorders like depression, or when screening for cognitive impairment in an individual at risk of Alzheimer's disease in which language content, flow and semantic errors may provide important diagnostic clues, then linguistic analyses

maybe more relevant. Conversely, if the goals are to monitor treatment response in individuals with Parkinson's disease in whom variations in the quality of voice and articulation are expected, then analyzing acoustic features may be more helpful.

Voice/Speech Analysis Associated with Diseases

Studies have shown that features identified from voice/speech data are associated with several disease states (**Table 3**). An important distinction to note here is that voice 'features' refers to individual acoustic (audible and non-audible sound-related) and linguistic (semantic content-related) characteristics that can be identified from voice/speech data using machine learning. A voice 'biomarker' is typically identifiable when more than one of these unique features is integrated into a single combined measure. Although single feature voice biomarkers are possible, in the current paper we refer to voice biomarkers as reflective of multiple voice/speech features. We searched Pubmed and Google Scholar to identify studies evaluating the relationship between voice/speech features and biomarkers and various diseases. Search terms included "digital biomarker," "telemedicine," "voice features," "voice biomarker," "speech features," "speech biomarkers," "acoustics," "linguistics," "cardiovascular disease," "neurologic disease," "psychiatric disease," and "infectious disease." The search was limited to studies published in English in peer-reviewed journals between 1980 and the present day. To identify potential studies not captured by our database search strategy, we also searched studies listed in the bibliography of relevant publications and reviews.

Cardiovascular Disease

To our knowledge the first study assessing the relationship between acoustic characteristics in voice with cardiovascular health was published in 1990.¹⁶ Eighteen nonsmoking vocally untrained males divided into groups of healthy young and healthy elderly men, and elderly men diagnosed with chronic atherosclerosis without other systemic complaints were asked to enunciate the sound “a” at a manageable pitch sustained at 70-78 decibels. Thus, according to our proposed classification system for voice biomarker analysis (**Table 1**) this study relied on a **Vocal Active Acoustic In-person (VAAI)** biomarker. Significant differences between groups included mean fundamental frequency, the standard deviation of amplitude, percent jitter (cycle-to-cycle variations in vocal fundamental frequency) and shimmer (cycle-to-cycle variations in vocal amplitude). These differences were most apparent when comparing healthy young males to elderly males with atherosclerotic disease, and while only percent jitter significantly differentiated the two elderly groups, those with atherosclerotic disease generally exhibited phonation with greater short- and long-term variability. These findings suggested that the elderly speakers may be more prone to detectable disruptions in vocalization when diseases are present that have the potential to disrupt vocal-fold vibration.¹⁶ This study provided one of the earliest indications that voice/speech analysis could indicate the presence of systemic pathology.

Our group has extended this work by further evaluating the association between voice biomarkers and cardiovascular disease^{32-34, 36} (**Table 4**) using a **Speech Active Acoustic In-person (SAAI)** biomarker. Patients presenting to our institution for a clinically indicated coronary angiography performed 3 separate 30-second speech recordings on an iPad or smartphone prior to their procedure. In the first of our studies, we used the MFCC to extract information from the speech recordings,³⁷ which is a sound processing tool used to classify healthy and abnormal voices.³⁸⁻⁴⁰ Speech recordings were then further analyzed using the Fourier transform mathematical function to convert data into patterns of frequency over time.⁴¹ We identified and extracted 81 prespecified voice features from each participant’s

recordings. In a multivariate binary logistic regression analysis, we identified 2 features that were associated with CAD.³² In our subsequent studies we used a single integrated voice biomarker that was a composite of a total of 223 acoustic features, which included the Mel Cepstrum representation, jitter, shimmer, loudness and pitch. This biomarker was developed and trained in a large cohort of patients with chronic medical conditions including congestive heart failure who were registered to a healthcare call center.³⁶ ‘Low-level’ acoustic features were presented in a temporal resolution of 100 points per second producing a matrix of 2000 columns, with the number of extracted temporal features as rows. ‘High-level’ features were then extracted from these matrices using a moments analysis, and a machine learning model was employed to generate the vocal biomarker, which was optimized to the training cohort using cross-validation techniques. In a cohort of more than 2,000 patients in Israel the voice biomarker was associated with all-cause mortality and heart failure hospitalization at follow-up. When we used the same biomarker in patients who presented to the catheterization laboratory for clinically indicated right heart catheterization procedures, we showed that the biomarker was significantly associated with mean pulmonary arterial pressure and pulmonary vascular resistance. Most recently, we reanalyzed the speech samples of the cohort of patients who underwent coronary angiographies in our first study using the integrated voice biomarker and prospectively followed these individuals for incident CAD events. In a multivariable analysis, patients with a ‘high’ voice biomarker value had a significantly higher risk of incident coronary disease events including being admitted to hospital with chest pain and an acute coronary syndrome (**Table 4**). The studies underscore the association between voice/speech analyses, based on acoustic features, and a range of different cardiovascular diseases.

Neurological Disease

Parkinson's disease currently lacks a definitive diagnostic test, and its detection often requires examination by a movement disorder specialist. Parkinson's disease is associated with well-characterized effects on speech, such as impaired phonation (hypophonia) and articulation (hypokinetic dysarthria), reduced energy in the higher parts of the harmonic spectrum and variations in pitch, and a high prevalence of voice disorders.⁴² Thus, a digital voice biomarker could provide an objective tool to identify and even "quantify" symptoms and assess response to treatment. Studies show that voice dysfunction may sometimes be an early sign of motor impairment in Parkinson's disease with perturbations in the fine motor control of vocalization, and may occur before less subtle changes in the limbs.⁴³ Thus, it follows that biomarkers evaluating verbal and vocal sounds obtained during active or passive tasks using an analysis of acoustic features could provide useful information when evaluating Parkinson's disease. In one study the investigators analyzed a database of speech recordings of Parkinson disease and non-Parkinson disease patients using machine learning models to predict disease with a **Speech Passive Acoustic Remote (SPAR)** biomarker. Their machine learning model could discriminate subjects with Parkinson's disease from healthy controls even amongst those with little to no symptoms. These findings show that voice/speech analysis could facilitate the early and rapid detection of Parkinson's disease,⁴⁴ potentially in the initial stages of disease when changes in voice are otherwise unrecognized. In another study, the investigators obtained voice samples from patients with Parkinson's disease who were asked to prolong "ahhh." They analyzed 132 acoustic features to generate a **Vocal Active Acoustic Remote (VAAR)** biomarker. When using only 10 acoustic features the authors demonstrated almost 99% classification accuracy.⁴⁵ Additional studies have looked at voice/speech analysis as a method to detect the progression of disease and its response to therapy. In one study the investigators extracted features of Parkinson disease progression in 6000 speech recordings from patients with Parkinson's disease to generate a **Speech Active Acoustic Remote (SAAR)** biomarker). The authors demonstrated clinically useful accuracy in estimates of measures of disease

progression when comparing clinician and voice algorithm derived estimates.⁴⁶ In another study using a **Vocal Active Acoustic Remote (VAAR)** biomarker, the authors demonstrated a significant correlation between the degree of improvement in voice quality attributed to therapy with levodopa and the extent of changes in voice quality assessed using their biomarker, suggesting its potential use in helping to target therapy more effectively.⁴⁷

Changes in language ability may be apparent many years before patients or family members recognize any symptoms of cognitive impairment.^{48,49} Indeed, Alzheimer's disease, and cognitive impairment more generally, is known to affect spoken language in several ways including alterations in verbal fluency reflected in hesitant speech and slow speech rate, word finding difficulties that result in the use of filler sounds such as "uh" and "um," circumlocutions, neologisms, grammatical simplification, semantic errors, repetitions etc.⁵⁰ Thus, voice/speech analysis that incorporated an evaluation of speech or voice sounds with a focus on principally linguistic, with some acoustic, features would seem most relevant in this population. In one study, speech samples were obtained from individuals in whom Alzheimer's disease was identified during life (and confirmed at post-mortem examination) and each participant exhibited objective signs of mild cognitive impairment 6 to 18 months prior to developing overt Alzheimer's disease. Speech samples were analyzed to develop a **Speech Active Linguistic In-Person (SALI)** biomarker. Subtle changes in speech patterns were noticeable even during the prodromal stages of Alzheimer's disease, and linear trends were evident in successive stages of disease progression.⁵¹ In an additional study that also developed and used a **Speech Active Linguistic In-Person (SALI)** biomarker the authors showed significant differences in linguistic features, such as "hesitation ratio, speech tempo, length and number of silent and filled pauses and length of utterance," between healthy controls and adults with mild cognitive impairment while performing memory tasks.⁵² Interestingly, acoustic features have also been used to screen for Alzheimer's disease in its pre-clinical

states⁵³ and to distinguish healthy individuals from those with mild cognitive impairment and Alzheimer's disease⁵⁴ using a **Speech Active Acoustic In-Person (SAAI)** biomarker. These findings support the notion that voice/speech analysis could provide an objective means of discriminating pathological cognitive changes from normal age-related changes. Another study used both linguistic features, measured from transcripts, and acoustic features, measured from audio files to develop a **Speech Active Dual (linguistic and acoustic) In-Person (SADI)** biomarker. The biomarker achieved > 80% accuracy in discriminating individuals with and without Alzheimer's disease.⁵⁵

Voice/speech analysis has also provided useful information in children with neurodevelopmental disorders. In autistic spectrum disorders deficits in social interaction are evident, and in those children who are able to speak, abnormalities in speech patterns and voice quality are often exhibited. Yet the diagnostic utility of these findings is unknown, and they remain challenging to assess objectively. In a study evaluating speech samples of children with and without autistic spectrum disorders using a **Speech Active Acoustic In-person (SAAI)** biomarker the investigators demonstrated significant differences between groups. Pitch variability yielded 90% specificity and 80% sensitivity in classifying autism in this sample.⁵⁶

Infectious Diseases

COVID-19 is a systemic illness with primarily upper and lower respiratory tract symptoms. A suitable voice biomarker might evaluate acoustic features of non-speech airway sounds such as breathing and coughing, but potentially also speech production. Ideally these samples could be obtained remotely to limit risk of disease transmission. One study used recordings of a forced cough over a cell phone to develop a **Non-verbal Active Acoustic Remote (NAAR)** biomarker, which was trained on 4,256 subjects and tested on 1,064 subjects. The biomarker achieved sensitivity for COVID-19 based on formal

microbiologic testing of 98.5% and a specificity of 94.2%, with an AUC of 97%.⁵⁷ In another study, the authors developed a **Vocal Active Acoustic Remote** (VAAR) biomarker in a training cohort of 434 voice samples in 272 participants, 160 (59%) of whom had COVID-19 on formal testing.⁵⁸ Transfer learning and adaptation methods were employed for feature extraction, and each 10-second recording of speech was converted to a Mel spectrogram, and was passed through a pre-trained convolutional neural network architecture. Two classification models, random forest and support vector machine, were then evaluated at different regularization levels and the best model was selected and tested in 80 subjects, 40 of whom tested positive for COVID-19. The voice biomarker differed significantly among infected compared to non-infected individuals with an area under the curve (AUC) of 72%. When the voice biomarker was added to a self-reported symptoms the AUC significantly improved to 85%.⁵⁸

Psychiatric Conditions

Psychiatric evaluation often includes judgments of speech loudness, rate, and prosody (fluctuations in rate, pitch and loudness that signal emphasis and emotion) as well as the organization, content and meaning of language which, taken together, convey information about emotional state, mood, clarity of thought, all which can contribute to psychiatric diagnosis. Posttraumatic stress disorder (PTSD) is a psychiatric condition whose diagnosis typically relies on interview and measures of self-report, which can lead to under- and over-diagnosis. Given the nature of PTSD, both acoustic and linguistic features derived from speech are likely to be informative in evaluating this condition. One study analyzed speech samples, obtained during clinical interviews from warzone-exposed veterans and included 52 cases of individuals with PTSD and 77 controls using a **Speech Active Dual** (acoustic and linguistic) **In-person** (SADI) biomarker. From 40,526 speech features that were input to a random forest algorithm 18 speech features were ultimately used, and collectively provided an AUC for a diagnosis of PTSD of 95.4%.⁵⁹

Another area in which clinical interviews and self-report measures are used and often lack accuracy is suicidal ideation, which was evaluated in US veterans again using a **Speech Active Dual** (acoustic and linguistic) **In-person (SADI)** biomarker. Fifteen acoustic and linguistic features derived from speech correctly classified suicidal ideation in veterans with an AUC of 80%.⁶⁰

Voice analysis has also been used to facilitate the detection of depression.⁶¹ In one online study investigators extracted features from anonymous speech samples of depressed subjects for depression severity and suicidality using a **Speech Active Dual** (acoustic and linguistic) **Remote (SADR)** biomarker. Voice features predicted the patient health questionnaire-9 (PHQ-9) scores with an AUC of 82.1%.⁶² Similar results were obtained using a biomarker developed in-person,⁶³ using linguistic assessment only during non-structured interviews with research psychologists,⁶⁴ and using acoustic features only in individuals asked to read a standardized passage in a nearly closed-room environment.⁶⁵ Interestingly, different acoustic features have been shown to be helpful across men and women with spectral and energy-related features being helpful in discriminating depression in males (sensitivity 0.95, specificity 0.88), while prosody-related features being more helpful in females (sensitivity 0.73, specificity 0.86) suggesting the need for stratifying biomarkers across different demographic groups.⁶⁶ **Speech Active Acoustic Remote (SAAR)** biomarkers have also been used to follow therapy in patients with depression, including in a double blinded trial of 105 adults with major depression who were randomized to sertraline 50 mg/day versus placebo. Clinical changes in depressive symptoms among responders to pharmacotherapy were associated with significant differences in speech patterns that were not observed, or were directionally opposite to changes seen in non-responders.⁶⁷ Similar findings were observed in an observational study in which acoustic measures correlated significantly with baseline depression severity and were significantly different between treatment responders and non-

responders.⁶⁸ Lastly, acoustic features extracted from voice recordings made in naturalistic settings obtained using smartphones have also been used to discriminate bipolar disorder using a **Speech Active Dual (acoustic and linguistic) Remote (SADR) biomarker**.⁶⁹ Acoustic features have also been used to distinguish depression and mixed depression, and hypomania and mixed hypomania in individuals with bipolar disorder during an acute episode using an in-person verbal fluency task.⁷⁰ Thus, voice analysis technologies may increase the range of available methods that can be used to explore the neurobiological substrates of psychiatric conditions given their inherent dynamic interplay of systems that direct the motor, cognitive, and emotional basis of speech production.

Developing the Biomarker

Table 5 lists important characteristics of voice biomarkers. There is currently no standard protocol for obtaining and analyzing voice/speech samples for the development of biomarkers. **Figure 3** outlines potential steps for developing and using voice biomarkers in research and clinical practice.

Quality Control

Voice/speech data need to be recorded in conjunction with the collection of relevant clinical information to provide a reference dataset of “labelled gold-standards” of outcomes to facilitate the training and testing of voice biomarkers. Audio recordings then must be pre-processed. This involves steps to improve the performance of feature detection: resampling allows improvement of poor quality recordings and increases the number of high-quality recordings to reduce measurement error; normalization reduces extraneous information collected without altering differences in the ranges of acoustic values; noise detection and reduction generates a “clean voice estimation” through passing the

'noisy voice' through a linear filter and framing involves dividing the voice/speech sample into a number of samples, which are then multiplied by a "window function" to limit signal leakage effects.⁷¹

The Digital Medicine Society has provided a framework for consistent evaluation of voice/speech-based biomarkers,⁷² which comprised three principal components: verification, analytical validation, and clinical validation.

Verification

Verification involves evaluating the quality of recordings and comparing the effects of recording devices, hardware and sensors as well as recording conditions, including the effects of ambient noise, on the integrity of the recordings.⁷² Layout of the user interface must also be considered including the provision of instructions, whether there are empty periods in the recordings to determine ambient noise, and whether feedback is provided using voice prompts or beeps to improve adherence to the task. Ideally, the suitability of conditions can be communicated to the user via intrinsic calibration tasks so to ensure that the quality of recordings are sufficient for accurate analysis.

Analytical Validation

Analytical validation includes assessing the validity and reliability of voice/speech sample processing and includes reproducibility and demographic-dependence (e.g., age, education, diagnosis and accent). Analytical validation also involves comparing voice/speech features against reference standards.⁷² The processing steps involved here can often be complex as the features that can be extracted from voice/speech are highly variable and fluctuate according to the task.

For feature extraction, features are defined as “the most dominating and discriminating characteristics” of a speech/voice sample⁷¹ and ultimately contribute to the training of machine learning algorithms for the detection of disease. The type of features used to generate a particular voice biomarker depends on the disease of interest and the recording and processing methods available. These include acoustic features that reflect mathematical properties of the sound wave (**Table 2**). Many acoustic features are computed via the mathematical transformations of the sound wave and can thus be validated using mathematical models of voice/speech production. MFCCs are the most frequently used acoustic feature extracted and used in speech analysis, followed by perceptual linear prediction coefficients, and linear frequency cepstral coefficients.³² Acoustic features can be identified by applying the temporal, frequency, cepstral, wavelet, and time-frequency domains.⁷³ On the other hand, linguistic features, reflecting the content and organization of language (**Table 2**), can be obtained using natural language-processing tools or manual coding methods. These techniques require prolonged speech recordings to ensure features may be extracted at all linguistic levels. Features must then be “selected” for which various statistical techniques can be employed including Minimum redundancy maximum relevance⁷⁴ or Gram-Schmidt orthogonalization.⁷⁵ These allow a population of features to be selected without alteration by removing highly correlated features and those with missing values or low variance. Selected features can then be used individually or in various combinations as a composite biomarker to train machine and deep learning algorithms to classify clinical outcomes of interest.⁷¹

Clinical Validation

Clinical validation verifies the correspondence of a voice biomarker to a suitable clinical reference standard or several surrogate measures. In clinical validation studies, investigators should report clinically useful performance metrics such as accuracy, specificity, sensitivity, precision, and AUC to

guide clinical application. Where voice biomarkers are used, if they comprise multiple aggregated voice features, it is useful to also determine which individual features contribute most to classification models to help with clinical interpretation and to better understand the changes that accompany a disease. Lastly, if such voice biomarkers are found to be valid, reliable and of potential clinical utility additional steps will be required to ensure their safe and effective use in clinical practice (see below).

Future Directions

Table 6 summarizes key areas in voice/speech analysis that require further investigation. First, there needs to be greater standardization in voice/speech data collection and analysis protocols.

Voice/speech recordings will almost exclusively be done outside of a laboratory setting, in day-to-day situations and environments using an range of devices. Such standardization would permit cross-comparisons and transferability of files and data sharing. This would pave the way for large-scale voice/speech sample repositories, labelled with clinical outcomes, essential for ongoing research.

Moreover, further clinical trials with large, demographically, and linguistically diverse populations from different geographic regions assessing the use of voice biomarkers in evaluating a variety of disease states will be needed to ensure that they are providing universally reliable and valid measurements.

This process of estimating and then improving upon the accuracy of voice biomarkers should be viewed as an iterative process rather than a conclusive outcome, particularly when considering how machine learning performs better with more input. Consideration needs to be given to the broader integration of voice/speech analysis technology into clinical practice and society. Voice biomarkers should be viewed as adjunctive tools, which can aid physicians in making decisions while offering advantages of limiting unnecessary in-person evaluations that may be lengthy and expensive. In the modern era telecommunications are by in large provided digitally, which coupled with the high proportion of

individuals owning a smartphone or tablet, provides the perfect alliance for the remote recording, transmission, and analysis of digital voice/speech recordings. One potential scenario is the adjustment of medications based on serial measurements of voice biomarkers across time such as in patients with heart failure when titrating diuretics or other guideline directed medical therapy. Another scenario could be the weekly or even daily *physical* done from home using voice biomarkers that could provide feedback to patients and healthcare providers on clinical status. Voice biomarker algorithms, embedded within medical devices would raise further questions about agency approval, such as from the Food and Drug Administration; business models for device and biomarker developers; and whether insurance companies would provide coverage. In clinical research, if voice biomarkers prove to be valid and reliable markers of disease progression they could potentially be used as cost and time efficient proxy endpoints in clinical trials. This would help make trials less expensive allowing more participants to be recruited and more endpoints to be studied, and specifically in pharmaceutical trials, could help identify drugs with (and without) promise sooner.

Safety Considerations

These exciting prospects should be viewed with caution, however. Relying on such home monitoring systems to prompt treatment decisions and whether patients should see their providers makes ensuring the reliability and validity of voice biomarkers even more important. The potential adverse consequences of false positives and false negatives must be safeguarded against, further underscoring the need for large prospective studies evaluating the use of voice biomarkers in a real-life setting. Further, the implications of such transitions to remote healthcare on real-life clinical practice patterns and patient care and outcomes need to be determined. This is important given the lack of established guidelines outlining best practice for remote care and the potential for unintended consequences like

those created by the so-called 'digital divide' whereby specific groups of patients, namely those who are elderly, from ethnic minority groups, and those with a higher prevalence of comorbidities are likely less able to use remote care through lack of access to the internet or technology literacy. Conversely, digital health may have the potential for greater inclusivity and helping to break down existing societal divides. Older patients, those from ethnic minority backgrounds, or those with more medical comorbidities may find remote monitoring for healthcare purposes more attractive as they have greater difficulty accessing in-person visits due to barriers to transportation or scheduling.⁷⁶ Thus, the potential effects that voice biomarkers, and digital healthcare biomarkers in general, integrated into care practice models have on patient outcomes need to be studied in large prospective trials.

As part of these studies, a greater understanding of the mechanisms that link detectable changes in voice/speech with various pathological processes should be pursued as well. Healthcare providers will be more likely to use a biomarker to help inform their clinical decisions if they feel they have some appreciation for how it relates to the underlying pathophysiology of a disease, and in the relationship between voice/speech analysis and diseases such as CAD this remains unclear. It may be that voice/speech analysis provides information on an individual's general wellbeing or it could relate more closely to their emotional or psychosocial state, which could be *associated with* the underlying disease of interest. These relationships will need to be disentangled. Another perspective relates to issues with using clinical tools that rely on artificial intelligence, including to what extent we feel comfortable with acting on the output of a data-processing system that we do not intuitively understand. Lastly, privacy concerns should be considered as voice/speech data is sensitive and may be used to uncover a person's identity, ethnicity, and, after all, health status. For example, insurers could use these systems to analyze speech samples without consent to discriminate against customers, while workplaces could do similarly

when selecting prospective employees. Appropriate regulation will be vital, as will measures that encrypt recordings and prevent privacy leakage to address some of these ethical concerns.

Conclusion

Digital biomarkers, aided by artificial intelligence, will be increasingly used in medicine going forward. In the current review we have outlined why voice biomarkers offer promise in remote healthcare monitoring both in clinical practice, by providing adjunctive information to aid healthcare providers make decisions about the detection, progression and treatment of disease, and in research by offering useful surrogate endpoints for clinical trials. However, for voice biomarkers to find an established role in medicine of the future much work remains including the systematic and rigorous evaluation of voice biomarkers in demographically and linguistically diverse populations using large prospective clinical trials, defining how to most appropriately integrate voice biomarkers into existing clinical and research work models, and cooperating with industry, regulatory agencies, insurance companies and, most importantly, our patients to ensure these biomarkers are used safely and effectively while prioritizing data protection and individual privacy.

References

1. Uscher-Pines L, Sousa J, Jones M, et al. Telehealth Use Among Safety-Net Organizations in California During the COVID-19 Pandemic. *JAMA*. 2021;325:1106-1107.
2. Dorsey ER, Topol EJ. State of Telehealth. *N Engl J Med*. 2016;375:154-161.
3. Hunt TL, 2nd, Hooten WM. The Effects of COVID-19 on Telemedicine Could Outlive the Virus. *Mayo Clin Proc Innov Qual Outcomes*. 2020;4:583-585.
4. Contreras CM, Metzger GA, Beane JD, Dedhia PH, Ejaz A, Pawlik TM. Telemedicine: Patient-Provider Clinical Engagement During the COVID-19 Pandemic and Beyond. *J Gastrointest Surg*. 2020;24:1692-1697.
5. Margolis KL, Asche SE, Dehmer SP, et al. Long-term Outcomes of the Effects of Home Blood Pressure Telemonitoring and Pharmacist Management on Blood Pressure Among Adults With Uncontrolled Hypertension: Follow-up of a Cluster Randomized Clinical Trial. *JAMA Netw Open*. 2018;1:e181617.
6. Merriel SW, Andrews V, Salisbury C. Telehealth interventions for primary prevention of cardiovascular disease: a systematic review and meta-analysis. *Prev Med*. 2014;64:88-95.
7. Ringeval F. AVEC 2019 Workshop and Challenge: State-of-mind, Detecting Depression with AI, and Cross-Cultural Affect Recognition," *Proc. of the 9th International on Audio/Visual Emotion Challenge and Workshop, 2019*, pp. 3-12.
8. T. Rutowski AH, Y. Lu, E. Shriberg. *Optimizing Speech-Input Length for Speaker-Independent Depression Classification, Proc. Interspeech, 2019*, pp.3023-3027.
9. N. Cummins SS, J. Krajewski, S. Schnieder, J. Epps, T. F. Quatieri. *A Review of Depression and Suicide Risk Assessment Using Speech Analysis, Speech Communication, vol. 71, pp. 10-49, 2015*.
10. Pampouchidou A. Depression Assessment by Fusing High- and Low-Level Features from Audio, Video, and Text, *Proc. of the 6th International Workshop on Audio/Visual Emotion Challenge, 2016*, pp. 27-34.
11. Williamson JR. Detecting Depression using Vocal, Facial and Semantic Communication Cues," *Proc. of the 6th International Workshop on Audio/Visual Emotion Challenge, 2016*, pp. 11-18.
12. Sanz L, Sistiaga JA, Lara AJ, Cuende E, Garcia-Alcantara F, Rivera T. The prevalence of dysphonia, its association with immunomediated diseases and correlation with biochemical markers. *J Voice*. 2012;26:148-153.
13. Hamdan AL, Jabbour J, Nassar J, Dahouk I, Azar ST. Vocal characteristics in patients with type 2 diabetes mellitus. *Eur Arch Otorhinolaryngol*. 2012;269:1489-1495.
14. Mohammadzadeh A, Heydari E, Azizi F. Speech impairment in primary hypothyroidism. *J Endocrinol Invest*. 2011;34:431-433.
15. Murton OM, Hillman RE, Mehta DD, et al. Acoustic speech analysis of patients with decompensated heart failure: A pilot study. *J Acoust Soc Am*. 2017;142:EL401.
16. Orlikoff RF. The relationship of age and cardiovascular health to certain acoustic characteristics of male voices. *J Speech Hear Res*. 1990;33:450-457.
17. Hansson GK. Inflammation, atherosclerosis, and coronary artery disease. *N Engl J Med*. 2005;352:1685-1695.
18. Liao D, Cai J, Rosamond WD, et al. Cardiac autonomic function and incident coronary heart disease: a population-based case-cohort study. The ARIC Study. Atherosclerosis Risk in Communities Study. *Am J Epidemiol*. 1997;145:696-706.
19. Tsuji H, Larson MG, Venditti FJ, Jr., et al. Impact of reduced heart rate variability on risk for cardiac events. The Framingham Heart Study. *Circulation*. 1996;94:2850-2855.
20. Hackett RA, Steptoe A. Type 2 diabetes mellitus and psychological stress - a modifiable risk factor. *Nat Rev Endocrinol*. 2017;13:547-560.

21. Kivimaki M, Kawachi I. Work Stress as a Risk Factor for Cardiovascular Disease. *Curr Cardiol Rep.* 2015;17:630.
22. Dragano N, Siegrist J, Nyberg ST, et al. Effort-Reward Imbalance at Work and Incident Coronary Heart Disease: A Multicohort Study of 90,164 Individuals. *Epidemiology.* 2017;28:619-626.
23. Kivimaki M, Jokela M, Nyberg ST, et al. Long working hours and risk of coronary heart disease and stroke: a systematic review and meta-analysis of published and unpublished data for 603,838 individuals. *Lancet.* 2015;386:1739-1746.
24. Huang Y, Xu S, Hua J, et al. Association between job strain and risk of incident stroke: A meta-analysis. *Neurology.* 2015;85:1648-1654.
25. Johannes B, Wittels P, Enne R, et al. Non-linear function model of voice pitch dependency on physical and mental load. *Eur J Appl Physiol.* 2007;101:267-276.
26. Protopapas A, Lieberman P. Fundamental frequency of phonation and perceived emotional stress. *J Acoust Soc Am.* 1997;101:2267-2277.
27. Holmqvist-Jansen S, Johansson A, Santtila P, Westberg L, von der Pahlen B, Simberg S. Investigating the Role of Salivary Cortisol on Vocal Symptoms. *J Speech Lang Hear Res.* 2017;60:2781-2791.
28. M. Vermani MM, M. Katzman. Rates of Detection of Mood and Anxiety Disorders in Primary Care: A Descriptive, Cross-Sectional Study,” *Prim. Care Companion CNS Disord.* 2011, PCC.10m01013.
29. M. Carey KJ, G. Meadows, R. Sanson-Fisher, C. D’Este, K. Inder, S. L. Yoong, G. Russell. Accuracy of general practitioner unassisted detection of depression, *Aust. N. Z. J. Psychiatry*, vol. 48, no. 6, pp. 571–578, 2014.
30. Sara JDS, Lerman LO, Lerman A. The endothelium is a key player in the vascular response to acute mental stress. *Eur Heart J.* 2021;42:4089-4091.
31. Kivimaki M, Steptoe A. Effects of stress on the development and progression of cardiovascular disease. *Nat Rev Cardiol.* 2018;15:215-229.
32. Maor E, Sara JD, Orbelo DM, Lerman LO, Levanon Y, Lerman A. Voice Signal Characteristics Are Independently Associated With Coronary Artery Disease. *Mayo Clin Proc.* 2018;93:840-847.
33. Sara JDS, Maor E, Borlaug B, et al. Non-invasive vocal biomarker is associated with pulmonary hypertension. *PLoS One.* 2020;15:e0231441.
34. Sara JDS, Maor E, Orbelo D, Gulati R, Lerman LO, Lerman A. Noninvasive Voice Biomarker Is Associated With Incident Coronary Artery Disease Events at Follow-up. *Mayo Clin Proc.* 2022.
35. Gerratt BR, Kreiman J, Garellek M. Comparing Measures of Voice Quality From Sustained Phonation and Continuous Speech. *J Speech Lang Hear Res.* 2016;59:994-1001.
36. Maor E, Perry D, Mevorach D, et al. Vocal Biomarker Is Associated With Hospitalization and Mortality Among Heart Failure Patients. *J Am Heart Assoc.* 2020;9:e013359.
37. Picone JW. *Signal modeling techniques in speech recognition*1993.
38. O’Shaughnessy D. Invited paper: automatic speech recognition: history, methods and challenges. *Pattern Recogn.* 2008;41:2965-2979.
39. J.I. Godino-Llorente PG-V, M. Blanco-Velasco. Dimensionality reduction of a pathological voice quality assessment system based on Gaussian mixture models and short-term cepstral parameters. *IEEE Trans Biomed Eng.* 2006;53:1943-1953.
40. Ö. Eskidere AG. Voice disorder classification based on multitaper mel frequency cepstral coefficients features. *Comput Math Methods Med.* 2015;2015:956249.
41. Trevisan J, Angelov PP, Carmichael PL, Scott AD, Martin FL. Extracting biological information with computational analysis of Fourier-transform infrared (FTIR) biospectroscopy datasets: current practices to future perspectives. *Analyst.* 2012;137:3202-3215.

42. Dashtipour K, Tafreshi A, Lee J, Crawley B. Speech disorders in Parkinson's disease: pathophysiology, medical management and surgical approaches. *Neurodegener Dis Manag.* 2018;8:337-348.
43. Harel B, Cannizzaro M, Snyder PJ. Variability in fundamental frequency during speech in prodromal and incipient Parkinson's disease: a longitudinal case study. *Brain Cogn.* 2004;56:24-29.
44. Tracy JM, Ozkanca Y, Atkins DC, Hosseini Ghomi R. Investigating voice as a biomarker: Deep phenotyping methods for early detection of Parkinson's disease. *J Biomed Inform.* 2020;104:103362.
45. Tsanas A, Little MA, McSharry PE, Spielman J, Ramig LO. Novel speech signal processing algorithms for high-accuracy classification of Parkinson's disease. *IEEE Trans Biomed Eng.* 2012;59:1264-1271.
46. Tsanas A, Little MA, McSharry PE, Ramig LO. Accurate telemonitoring of Parkinson's disease progression by noninvasive speech tests. *IEEE Trans Biomed Eng.* 2010;57:884-893.
47. Cushnie-Sparrow D, Adams S, Abeyesekera A, Pieterman M, Gilmore G, Jog M. Voice quality severity and responsiveness to levodopa in Parkinson's disease. *J Commun Disord.* 2018;76:1-10.
48. Forbes-McKay KE, Venneri A. Detecting subtle spontaneous language decline in early Alzheimer's disease with a picture description task. *Neurol Sci.* 2005;26:243-254.
49. Oulhaj A, Wilcock GK, Smith AD, de Jager CA. Predicting the time of conversion to MCI in the elderly: role of verbal expression and learning. *Neurology.* 2009;73:1436-1442.
50. Reilly J, Peelle JE, Antonucci SM, Grossman M. Anomia as a marker of distinct semantic memory impairments in Alzheimer's disease and semantic dementia. *Neuropsychology.* 2011;25:413-426.
51. Ahmed S, Haigh AM, de Jager CA, Garrard P. Connected speech as a marker of disease progression in autopsy-proven Alzheimer's disease. *Brain.* 2013;136:3727-3737.
52. Toth L, Hoffmann I, Gosztolya G, et al. A Speech Recognition-based Solution for the Automatic Detection of Mild Cognitive Impairment from Spontaneous Speech. *Curr Alzheimer Res.* 2018;15:130-138.
53. Martinez-Sanchez F, Meilan JJG, Carro J, Ivanova O. A Prototype for the Voice Analysis Diagnosis of Alzheimer's Disease. *J Alzheimers Dis.* 2018;64:473-481.
54. Konig A, Satt A, Sorin A, et al. Automatic speech analysis for the assessment of patients with predementia and Alzheimer's disease. *Alzheimers Dement (Amst).* 2015;1:112-124.
55. Fraser KC, Meltzer JA, Rudzicz F. Linguistic Features Identify Alzheimer's Disease in Narrative Speech. *J Alzheimers Dis.* 2016;49:407-422.
56. Bonnef YS, Levanon Y, Dean-Pardo O, Lossos L, Adini Y. Abnormal speech spectrum and increased pitch variability in young autistic children. *Front Hum Neurosci.* 2011;4:237.
57. J. Laguarda FHaBS. COVID-19 Artificial Intelligence Diagnosis Using Only Cough Recordings. *IEEE Open Journal of Engineering in Medicine and Biology.* 2020;1:275-281.
58. Maor E, Tsur N, Barkai G, et al. Noninvasive Vocal Biomarker is Associated With Severe Acute Respiratory Syndrome Coronavirus 2 Infection. *Mayo Clin Proc Innov Qual Outcomes.* 2021;5:654-662.
59. Marmar CR, Brown AD, Qian M, et al. Speech-based markers for posttraumatic stress disorder in US veterans. *Depress Anxiety.* 2019;36:607-616.
60. Belouali A, Gupta S, Sourirajan V, et al. Acoustic and language analysis of speech for suicidal ideation among US veterans. *BioData Min.* 2021;14:11.
61. I. Goodfellow YB, A.Courville. Deep Learning. MIT Press, 2016.
62. Zhang L, Duvvuri R, Chandra KKL, Nguyen T, Ghomi RH. Automated voice biomarkers for depression symptoms using an online cross-sectional data collection initiative. *Depress Anxiety.* 2020;37:657-669.

63. Taguchi T, Tachikawa H, Nemoto K, et al. Major depressive disorder discrimination using vocal acoustic features. *J Affect Disord.* 2018;225:214-220.
64. Yamamoto M, Takamiya A, Sawada K, et al. Using speech recognition technology to investigate the association between timing-related speech features and depression severity. *PLoS One.* 2020;15:e0238726.
65. Hashim NW, Wilkes M, Salomon R, Meggs J, France DJ. Evaluation of Voice Acoustics as Predictors of Clinical Depression Scores. *J Voice.* 2017;31:256 e251-256 e256.
66. Lee S, Suh SW, Kim T, et al. Screening major depressive disorder using vocal acoustic features in the elderly by sex. *J Affect Disord.* 2021;291:15-23.
67. Mundt JC, Vogel AP, Feltner DE, Lenderking WR. Vocal acoustic biomarkers of depression severity and treatment response. *Biol Psychiatry.* 2012;72:580-587.
68. Mundt JC, Snyder PJ, Cannizzaro MS, Chappie K, Geralts DS. Voice acoustic measures of depression severity and treatment response collected via interactive voice response (IVR) technology. *J Neurolinguistics.* 2007;20:50-64.
69. Faurholt-Jepsen M, Busk J, Frost M, et al. Voice analysis as an objective state marker in bipolar disorder. *Transl Psychiatry.* 2016;6:e856.
70. Weiner L, Guidi A, Doignon-Camus N, Giersch A, Bertschy G, Vanello N. Vocal features obtained through automated methods in verbal fluency tasks can aid the identification of mixed episodes in bipolar disorder. *Transl Psychiatry.* 2021;11:415.
71. Fagherazzi G, Fischer A, Ismael M, Despotovic V. Voice for Health: The Use of Vocal Biomarkers from Research to Clinical Practice. *Digit Biomark.* 2021;5:78-88.
72. Goldsack JC, Coravos A, Bakker JP, et al. Verification, analytical validation, and clinical validation (V3): the foundation of determining fit-for-purpose for Biometric Monitoring Technologies (BioMeTs). *NPJ Digit Med.* 2020;3:55.
73. Sharma G, Umapathy K, Krishnan S. Trends in audio signal feature extraction methods. *Applied Acoustics.* 2020;158:107020.
74. Sajal MSR, Ehsan MT, Vaidyanathan R, Wang S, Aziz T, Mamun KAA. Telemonitoring Parkinson's disease using machine learning by combining tremor and voice analysis. *Brain Inform.* 2020;7:12.
75. Arora S, Baghai-Ravary L, Tsanas A. Developing a large scale population screening tool for the assessment of Parkinson's disease using telephone-quality voice. *J Acoust Soc Am.* 2019;145:2871.
76. Syed ST, Gerber BS, Sharp LK. Traveling towards disease: transportation barriers to health care access. *J Community Health.* 2013;38:976-993.

Figure Legends

Figure 1

Classification of voice biomarkers with different active and passive methods for making voice recordings

Figure 2

Outline of the potential underlying mechanisms that could explain the relationship between disease states and voice acoustic and linguistic features.

Figure 3

Outline of the process required for developing a voice biomarker for use in clinical practice

Abbreviation: FDA, food and drug administration

Table 1**Proposed novel classification system for voice biomarkers**

I	Source of Vocalization	S – speech V – vocal sounds N – non-verbal sounds D – all types
II	Method of Voice Assessment	A – active recording P – passive recording D – both types
III	Method of Voice Analysis	A – acoustic L – linguistic D – both types
IV	Location of Recording	I – in-person R – remote

Table 2

Examples of Different Methods of Voice Analysis and Different Location and Modes of Voice Recordings

		Examples	Potential Clinical application
Method of Voice Analysis	Linguistic analysis	<ul style="list-style-type: none"> • Phonetic and phonological (e.g., number of pauses, total pause time, hesitation ratio, speech rate) • Linguistic prosody (e.g., melody of speech conveying attitude and emotion, high and low formants) • Lexico-semantic (e.g., average rate of occurrence for each part of speech, number of repetitions, semantic errors, and closed-class word errors) • Morphosyntactic and syntactic (e.g., number of words per clause, number of dependent and simple clauses, number of clauses per utterance, mean length of utterances) • Discourse-pragmatic (e.g., cohesion, coherence) 	Linguistic features extracted from spontaneous or semi-spontaneous speech for the detection of Alzheimer's disease, mild cognitive impairment or mental health disorders (e.g., depression and post-traumatic stress disorder)
	Acoustic analysis	<ul style="list-style-type: none"> • Attitudinal and affective prosody (e.g., frequency over time) • Fundamental frequency-based measures (e.g., zero-crossing rate, harmonic-to-noise ratio, noise-to-harmonic ratio, jitter and shimmer) • Spectral characteristics (e.g., spectral flux, slope, centroid, entropy, roll-off, and flatness) • Phonation (e.g., pitch period entropy) • Nonlinear dynamic features (e.g., correlation dimension, fractal dimension, recurrence period density entropy, Lempel-Ziv complexity) • Segmental features (e.g., MFCCs, perceptual linear prediction coefficients, linear frequency cepstral coefficients) 	Acoustic features extracted from sustained vowel phonations or diadochokinetic tasks used in the detection of Parkinson's disease or from individuals actively coughing and/or passively breathing used in the detection of COVID-19 pneumonia

Location of Recording	In-person	<ul style="list-style-type: none"> • Studio or clinic-based 	Controlled environment where unwanted noise is reduced, requires calibrated equipment, prevents remote testing, requires professional staff e.g., sound engineer
	Remote	<ul style="list-style-type: none"> • Telephone-based 	Allows remote testing, involves variable equipment, restricted recording bandwidth, and does not permit environment control
		<ul style="list-style-type: none"> • Online computer-based 	Accessible and widely available, permits large scale-data collection
		<ul style="list-style-type: none"> • Smartphone-based 	Widely available, low cost, potential for greater bandwidth than telephone recordings, permits remote and frequent monitoring “on-the-go”, does not permit environment control

^a abbreviations: COVID-19, coronavirus disease 2019; MFCC, mel frequency cepstral coefficient

Table 3**Examples of diseases shown to be associated with features and biomarkers derived from voice signal analysis**

Cardiovascular Diseases

- Coronary artery disease at baseline
- Coronary artery disease events at follow-up
- Pulmonary hypertension
- Heart failure hospitalization and mortality

Neurologic Diseases

- Parkinson's disease
- Alzheimer's disease
- Mild cognitive impairment

Infectious Diseases

- Coronavirus disease-19

Psychiatric Conditions

- Autistic spectrum disorder
- Post-traumatic stress disorder
- Major depression
- Bipolar disorder

Table 4

Summary of findings of studies evaluating the association between voice biomarkers and cardiovascular outcomes

Ref.	Year	Sample size and description	Age, yrs	Male (%)	Follow-up	Method of voice sample collection	Voice biomarker used	Cardiovascular outcome	Principal findings
(33)	2018	101 patients undergoing planned coronary angiograms; 37 controls in Rochester, USA Control group: N=15 undergoing noncardiac procedures (i.e., hernia surgery) and N=22 undergoing no procedures	Median 61, IQR 51 - 67	54 (54)	NA (cross-sectional study)	In person, using a smartphone or iPad, 3 separate 30-second voice recordings (R1—reading prespecified text; R2—describing a positive emotional experience; R3—describing a negative emotional experience)	MFCCs used to extract information from the recordings with further analysis using the Fourier transform mathematical function; this allowed the extraction of 81 different prespecified voice features from each recording	Presence of CAD on invasive coronary angiography	Univariate binary logistic regression analysis identified 5 voice features that were associated with CAD ($P < .05$ for all) Multivariate binary logistic regression with adjustment for ASCVD risk score identified 2 voice features associated with CAD (OR, 95% CI: 0.37, 0.18-0.79; and 4.01, 1.25-12.84; $P = 0.009$ and $P = .02$, respectively) Both features were more strongly associated with CAD when patients described an emotionally significant experience

(37)	2020	2,267 patients with CHF	Median 77, IQR 68 - 83	1434 (63%)	Median 20 months, IQR 9 - 34	Analysis of anonymized audio phone recordings of patients	<p>A total of 223 acoustic features were extracted from 20 seconds of speech for each patient</p> <p>The voice biomarker was developed on a training cohort of patients who did not have CHF, N=8316</p>	<p>Primary outcome: all-cause mortality</p> <p>Secondary outcome: hospitalization during follow-up</p>	<p>824 (36%) patients died; 1403 (62%) were admitted for at least 24 hours at follow-up</p> <p>Kaplan-Meier survival analysis showed higher cumulative probability of death with increasing voice biomarker quartiles (23%, 29%, 38%, and 54%; $P<0.001$); There were 313 (54%), 335 (58%), 357 (63%), and 398 (70%) hospitalizations in the four biomarker quartile groups ($P<0.001$)</p> <p>In multivariable analysis, compared with the lowest quartile, patients in the highest quartile (Q4) were 96% more likely to die, $P<0.001$, and 50% more likely to require hospitalization ($P<0.01$)</p>
(34)	2020	83 patients undergoing clinically indicated invasive cardiac catheterization studies	Mean 61.6, SD 15.1	37 (45%)	NA (cross-sectional study)	In person, using a smartphone or iPad, 3 separate 30-second recordings (as above)	<p>The vocal biomarker used was the same biomarker developed in the abovementioned cohort of patients with CHF</p> <p>Intra-class correlation coefficient between the separate voice recording biomarker</p>	Diagnosis of moderate or greater PH (defined as a mean PAP ≥ 35 mmHg)	<p>Patients with a high mean PAP had significantly higher values of the mean voice biomarker compared to a lower mean PAP (0.74 ± 0.85 vs. 0.40 ± 0.88 $p=0.046$)</p> <p>Multivariate logistic regression showed that an increase in the mean voice biomarker by 1 unit was associated with a high PAP, OR 2.31 (95% CI) 1.05–5.07, $p=0.038$</p>

							values was 0.829 (95% CI 0.740–0.889)		
(35)	2022	108 patients undergoing planned coronary angiograms in Rochester, USA	Mean age 59.5, SD 11.4	59 (55%)	Median 24 months, IQR 1 – 60	In person, using a smartphone or iPad, 3 separate 30-second recordings (as above)	The vocal biomarker used was the same biomarker developed in the abovementioned cohort of patients with CHF	<p>Primary outcome: composite of presenting to the emergency department with chest pain and/or being admitted to hospital with chest pain and/or having an acute coronary syndrome</p> <p>Secondary outcome: composite of a positive stress test at follow-up and/or the presence of CAD at follow-up coronary angiography</p>	In multivariable Cox proportional hazards models adjusting for CAD grade on baseline angiography, a high baseline mean voice biomarker was associated with both the primary composite outcome (HR 2.61, 95% CI 1.42 – 4.80, $P=.002$) and the secondary composite outcomes (HR 3.13, 95% CI 1.13 – 8.68, $P=0.03$)

^a abbreviations: ASCVD, atherosclerotic cardiovascular disease; CAD, coronary artery disease; CHF, chronic heart failure; CI, confidence interval; HR, hazard ratio; IQR, interquartile range; MFCC, mel frequency cepstral coefficient; NA, not applicable; OR, odds ratio; PAP, pulmonary artery pressure; PCWP, pulmonary capillary wedge pressure; PH, pulmonary hypertension; PVR, pulmonary vascular resistance; Q, quartile; R, recording; SD, standard deviation; USA, united states of america

Table 5**Important characteristics for Voice Features and Biomarkers**

- Can be done remotely (over the phone or other digital device)
- Non-invasive and non-intrusive
- Acquired in a short period of time
- Independent of language and geographic location
- Independent of age
- Specific for disease groups
- Can be combined with other digital technologies
- Allow patient identification in a secure and encrypted way
- Can be integrated with established clinical protocols and practice patterns
- Can be integrated into the electronic medical record

Table 6

Key areas of further investigation and for future consideration in developing vocal biomarkers in healthcare delivery

- Greater standardization in voice data collection and analysis protocols
- Generating large-scale voice sample repositories labelled with clinical outcomes
- Systematic and rigorous evaluation of voice biomarkers in demographically and linguistically diverse populations from different geographic regions
- Large prospective clinical trials evaluating the use of voice biomarkers as surrogate measures of disease
- Voice biomarkers integrated into clinical practice – at home monitoring using personal device, with feedback etc.
- Individuals using voice biomarkers, amongst other digital biomarkers, at home in smart devices such as smart mirrors etc.
- Agency approval (e.g., the Food and Drug Administration)
- Medical device and biomarker companies' long-term business models
- Insurance companies paying for the use of voice technology in clinical practice
- Voice biomarker as a proxy secondary endpoint in research studies
- Caution about false positives and false negatives and complacency of relying excessively/exclusively on voice biomarkers
- Lack of established guidelines outlining best practice for remote care
- Potential for unintended consequences from the 'digital divide' – both positive and negative
- Understanding the mechanism linking changes in voice signal to disease states
- Using artificial intelligence derived output in clinical practice without understanding 'what's under the car's bonnet'
- Data protection and privacy – need for regulation and sophisticated encryption methods

