Speech and Voice Quality as Digital Biomarkers in Depression: A Systematic Review[★]

****Giovanni Briganti, and ***Jerôme R. Lechien, *\$Mons, †Liège, ‡Bruxelles, ¶Baudour, Belgium, and ||**Paris, France

Summary: Objective. To review the current evidence on the use of artificial intelligence-driven speech and voice analysis as a biomarker for depression.

Methods. PubMed, Scopus, and Cochrane databases were reviewed by two independent investigators for studies investigating the use of artificial intelligence-driven speech and voice quality outcomes as biomarkers for depression according to the Preferred Reporting Items for Systematic Reviews and Meta-Analyses statements. The methodological quality and risk of bias of each included study were assessed using the Quality Assessment of Diagnostic Accuracy Studies (QUADAS-2) tool.

Results. Of the 108 identified records, 12 studies met the inclusion criteria. The studies examined 16 872 participants, including patients with major depressive disorder (n = 1535), bipolar disorder (n = 111), schizophrenia spectrum disorders (n = 35), and anxiety disorders (n = 224). Control groups included a total of 1204 healthy individuals. Speech and voice quality outcomes consistently distinguished depression from controls (AUC = 0.71-0.93), with prosodic, spectral, and perturbation measures showing significant correlations with standardized depression scales. Classification accuracies ranged from 78% to 96.5%. Six studies demonstrated high risk of methodological bias, primarily in patient selection and validation techniques. Voice recording contexts varied between clinical settings and mobile technologies.

Conclusion. The findings of this review highlight the potential of voice biomarkers as a novel tool for depression detection and monitoring. While current evidence demonstrates promising classification accuracy, methodological heterogeneity and generalizability concerns must be addressed before widespread clinical adoption.

Key Words: Speech—Voice—Otolaryngology—Otorhinolaryngology—Laryngeal—Larynx—Acoustic—Biomarker—Mood monitoring—Machine learning—Artificial intelligence.

INTRODUCTION

Major depressive disorder (MDD) is a prevalent and debilitating psychiatric condition characterized by persistent low mood, cognitive disturbances, and psychomotor changes. Depression significantly impairs daily functioning and is associated with an increased risk of suicide, cardiovascular disease, and metabolic disorders. Despite advancements in pharmacological and psychotherapeutic interventions, accurate and timely assessment of depressive symptoms remains a major challenge in clinical practice. Current diagnostic and monitoring approaches rely primarily on clinical interviews and self-reported symptom

scales, which, although widely used, are subject to recall bias, variability in patient adherence, and limited accessibility.³ There is a critical need for objective biomarkers that can facilitate real-time, accurate, and scalable monitoring of depression to improve clinical outcomes.⁴

Speech analysis has emerged as a promising noninvasive tool for detecting and monitoring depression. While artificial intelligence applications represent recent innovations, the foundation of speech analysis for detecting emotions and psychological states dates back to the early 1970s. Pioneering research established fundamental connections between vocal acoustics and emotional states, developing theoretical frameworks that linked physiological changes to specific vocal modifications.^{5–7} Alterations in vocal parameters, including pitch, prosody, speech rate, and spectral features, have been hypothesized to reflect the neurophysiological and psychomotor changes associated with depressive states. Depressive speech is often characterized by reduced vocal intensity, slower speech tempo, and increased acoustic perturbations, which may serve as measurable biomarkers for mood disorders. ^{9,10} With the advancement of artificial intelligence-driven voice analysis, particularly through machine learning and deep learning techniques, it has become possible to extract and classify complex speech features associated with depression. The integration of voice biomarkers into digital health platforms, including smartphone-based and telemedicine applications, has expanded the potential for continuous, ecologically valid monitoring of patients in real-world settings.

Address correspondence and reprint requests to: Giovanni Briganti, M.D., M.Sc., Ph.D., Unit of Computational Medicine and Neuropsychiatry, Faculty of Medicine, University of Mons, Avenue du Champ de Mars 6, 7000 Mons, Belgium. E-mail: giovanni.briganti@hotmail.com

Journal of Voice, Vol xx, No xx, pp. xxx-xxx 0892-1997

© 2025 The Voice Foundation. Published by Elsevier Inc. All rights are reserved, including those for text and data mining, AI training, and similar technologies. https://doi.org/10.1016/j.jvoice.2025.05.002

Accepted for publication May 1, 2025.

^{*} This research received no external funding.

From the *Unit of Computational Medicine and Neuropsychiatry, Faculty of Medicine, Pharmacy and Biomedical Sciences, University of Mons (UMONS), Mons, Belgium; †Department of Clinical Sciences, Faculty of Medicine, University of Liège, Liège, Belgium; ‡Faculty of Medicine, Université libre de Bruxelles, Bruxelles, Belgium; §Surgery Department, Research Institute for Language Science and Technology, University of Mons (UMons), Mons, Belgium; ¶Division of Laryngology and Bronchoesophagology, Department of Otolaryngology Head Neck Surgery, EpiCURA Hospital, Baudour, Belgium; ∥Department of Otolaryngology-Head and Neck Surgery, Foch Hospital, School of Medicine, UFR Simone Veil, Université Versailles Saint-Quentin-en-Yvelines (Paris Saclay University), Paris, France; and the **Department of Otolaryngology, Elsan Hospital, Paris, France.

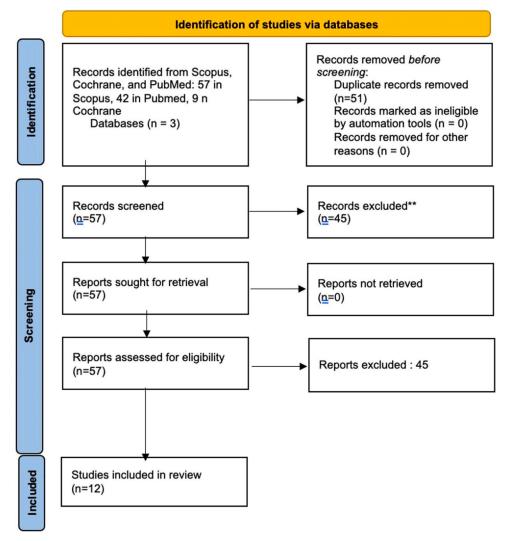


FIGURE 1. PRISMA flowchart.

Several studies have investigated the relationship between voice characteristics and depression, applying a wide range of methodological approaches and artificial intelligence models: these studies differ in the type of speech data analyzed, whether from controlled reading tasks or spontaneous speech, the machine learning techniques applied, and the validation methods used. While findings suggest that voice-based biomarkers hold potential for depression detection, the variability in study designs and the lack of standardized methodologies limit their generalizability and clinical implementation.

The objective of this systematic review was to summarize the current evidence on the use of artificial intelligence-driven voice analysis as a biomarker for depression. The authors will analyze existing studies on acoustic and linguistic features of speech in depressed individuals, as well as the machine learning models used for classification. This review aims to assess the reliability, validity, and limitations of voice-based biomarkers in the context of depression detection Figure 1.

MATERIALS AND METHODS

Framework for data extraction

Two independent investigators conducted the systematic review and data extraction following the Preferred Reporting Items for Systematic Reviews and Meta-Analyses guidelines for systematic reviews. The inclusion and exclusion criteria were based on the Population, Intervention, Comparison, Outcome, Timing, and Setting framework. Data were extracted systematically using a predefined approach to ensure comprehensive and reproducible findings. Two reviewers performed data extraction independently to maintain accuracy and consistency, with discrepancies resolved through discussion.

Patient population

The systematic review included prospective, retrospective, controlled, and uncontrolled studies investigating the use of artificial intelligence-based voice analysis in patients with

depression. Eligible studies focused on adults or adolescents diagnosed with MDD or related depressive conditions according to standardized diagnostic criteria, including DSM-5 and ICD-10. Studies specifying sample populations, including demographics, depression severity categorized as mild, moderate, or severe, comorbid psychiatric conditions, and control groups, were included. Studies analyzing mixed diagnostic groups without separate analyses for MDD were excluded.

Voice biomarkers and AI models

Studies assessing acoustic, prosodic, spectral, and linguistic voice features in individuals with depression were included. Controlled speech tasks, such as reading standardized texts, and naturalistic speech samples, including spontaneous speech, telemedicine calls, and smartphone-based recordings, were considered. The extracted voice features included fundamental frequency, jitter, shimmer, harmonic-to-noise ratio, MFCCs, spectral tilt, speech rate, intensity variation, and word choice analyzed through natural language processing. Studies implementing machine learning or deep learning models to classify or predict depression status based on these features were reviewed. Machine learning approaches included supervised models such as SVMs, random forests, and deep neural networks, unsupervised learning methods such as clustering and principal component analysis, and hybrid models.

Outcomes

Primary outcomes included classification accuracy, sensitivity, specificity, and area under the curve of artificial intelligence-based models for detecting depression from voice features. Studies that assessed the predictive validity of voice biomarkers for depression severity or treatment response were included. Secondary outcomes focused on correlations between voice features and clinical rating scales, including the Hamilton Depression Rating Scale,² the Patient Health Questionnaire-9,24 and the Montgomery-Åsberg Depression Rating Scale.²⁵ Studies that provided comparisons between depressed and nondepressed individuals or examined longitudinal variations in speech features across different mood states were analyzed. Validation strategies, including k-fold cross-validation, external test datasets, and reproducibility checks, were documented where reported.

Timing and settings

The included studies covered a range of settings, including clinical environments such as hospitals, psychiatric clinics, and structured laboratory assessments, as well as smartphone-based assessments conducted through mobile applications and telemedicine platforms. Some studies collected voice samples in real-world conditions, including conversational artificial intelligence interactions, wearable

devices, and telephone calls. Studies employing single-session assessments without concurrent clinical evaluation or without specification of timing and settings were excluded. Studies with longitudinal monitoring of mood variations over weeks or months were included.

Search strategy

A systematic search was conducted using PubMed, Scopus, and Cochrane databases to identify studies evaluating artificial intelligence-driven voice analysis for depression detection. Studies published between January 2015 and March 2025 were considered. The search strategy combined terms related to depression and voice. The extracted data included study design, country, setting, sample characteristics such as age, gender, and depression severity, voice biomarkers analyzed, artificial intelligence modeling techniques, validation methods, and outcome measures. Only peer-reviewed studies published in English were included.

Bias analysis

The methodological quality and risk of bias of each included study were assessed using the Quality Assessment of Diagnostic Accuracy Studies (QUADAS-2) tool. This framework evaluates the risk of bias across four domains, including patient selection, index test, reference standard, and flow and timing, along with applicability concerns related to patient selection and the generalizability of the artificial intelligence model. Each study was rated based on transparency in data reporting, patient inclusion methodology, model validation, and the presence of potential confounding factors.

RESULTS

A total of 108 records were identified through database searches, including 57 from Scopus, 42 from PubMed, and 9 from Cochrane. Following the removal of duplicates and screening for eligibility, 12 studies met the inclusion criteria for systematic review. These studies included nine cross-sectional studies 11-14,16,17,19,20,22 and three prospective studies. Eight studies incorporated a control group for comparison. L2-17,19,20,22 Study demographics, designs, and outcomes are presented in Table 1, with detailed population characteristics summarized in Table 2.

The included studies examined a total of 16 872 participants, comprising patients with MDD (n=1535), bipolar disorder (n=111), schizophrenia spectrum disorders (n=35), and anxiety disorders (n=224). Control groups included a total of 1204 healthy individuals. The mean age of participants across studies varied between 26.2 and 64 years. Two studies did not report age data, while one study provided only an age range of 23-69 years. Gender data were inconsistently reported, with 754 females and 464 males identified across studies, leaving 15 654 participants without specified gender information in the original

TABLE 1. Study Design and Outcomes	comes						
References	Design	Settings	Ν	F/M	Age (m/md)	Outcomes	Results
Mazur et al 2025 ¹¹	Cross- sectional	Remote recording (personal devices)	14 898	69.4%/ 27.9%*	37.3 (md)	Detection of depression with acoustic (F0, jitter shimmer) and prosodic features	Sensitivity: 71.3% Specificity: 73.5%
Menne et al 2024 ¹²	Cross- sectional Controlled	Laboratory (tablet microphone)	96 (44 MDD 52 HC)	38/58	26.2 (m)	Detection of depression with acoustic temporal and lexical features	AUC = 0.93 (speech model)
Ghosh et al 2024 ¹³	Cross- sectional Controlled	Laboratory (DAIC-WOZ dataset)	189 (59 depressed 130 HC)	SZ	SZ	Depression detection with acoustic textual and visual features	Accuracy: 81.6% (BERT) 68% (BiLSTM)
Huang et al 2024 ¹⁴	Cross- sectional Controlled	Laboratory (DAIC-WOZ dataset)	189 (56 depressed 133 HC)	87/102	S	Depression detection and severity classification with way2vec 2.0 features	Binary accuracy: 96.49% Multiclass: 94.81% RMSE: 0.1875
Ronneberg et al 2024 ¹⁵	s RCT- controlled	Voice-based Al coaching	200 (mild- moderate depression)	SZ	v 18	Depression symptom reduction with voice-based Al intervention	Neural engagement correlated with symptom reduction
Cansel et al 2023 ¹⁶	Cross- sectional Controlled	Smartphone	104 (15 SZ 24 ANX 78/26 25 MDD 15 BD 25 HC)	78/26	28.9- 37.5	Disorder classification with acoustic features (MFCC GTCC)	SVM: 96.48% kNN: 96.94% accuracy
Berardi et al 2023 ¹⁷	Cross- sectional Controlled	Digital voice recorder	60 (20 SZ 20 MDD 20 HC)	23/37	39.3- 41.6	Disorder classification with acoustic features	HC vs MDD: 92.0% SZ vs MDD: 93.2%
Li et al 2023 ¹⁸	Cross- sectional	Clinical interview recordings	329 (233 MDD 96 BD)	217/112	31.4	Depression severity classification with NLP	4-level severity: F1 = 0.719 Binary (mild vs severe): F1 = 0.884
Kim et al 2023 ¹⁹	Cross- sectional Controlled	Smartphone	318 (153 MDD 165 HC)	225/93	37.2- 38.6	Depression detection with acoustic features	CNN: 78.14% RF: 72.80% AUC: 0.86
Wasserzug et al 2023 ²⁰		Mobile phone conversations	144 (40 MDD including 14 in remission 104 HC)	73/71	S	Depression detection and remission with prosodic features	Significant differentiation (P < 0.0001) correlation with HAM-D (r=0.364)

TABLE 1 (Continued)	ted)						
References	Design	Settings	N	F/M	Age (m/md)	Age (m/md) Outcomes	Results
Lin et al 2022 ²¹	Cross- sectional	Smartphone app	263 (240 with	155/105*	63.0	Depression detection with NLP	AUC (PHQ-8) = 0.82 PPV=0.54 NPV = 0.88
	Controlled		depression history 23 HC)				
Hansen et al 2022 ²²	Cross-	Clinical interview	82 (40 MDD including	64/18	30.9-	Depression detection and remission	Depression detection and AUC = 0.71 (MDD vs HC) remission
	Controlled	recordings	25 in remission 42 HC)		2	with acoustic and emotion features	

The gender of some individuals was not specified. Abbreviations: ANX, anxiety disorder; AUC, area under the curve; BD, bipolar disorder; BERT, bidirectional encoder representations from transformers; female/male ratio; GTCC, gamma-tone cepstral coefficients; HAM-D/HAM-D, Hamilton Depression Rating Scale; HC, healthy controls; kNN, k-nearest neighbors; m, mean; N, sample size; NLP, natural language processing; NPV, negative predictive value; NS, not specified; PHQ-8/9, Patient Health Questionnaire (8 or 9 items); PPV, positive predictive value; RCT, randomized controlled trial; RF, random forest; RMSE, root mean square error; SD, standard deviation; SVM, support vector machine; SZ MDD, major depressive disorder; MFCC, mel-frequency cepstral coefficients; mean of precision and recall); F/M, schizophrenia

TABLE 2.
Summary of Demographics, Populations, and Outcomes

Category	N
Gender	
Females	754
Males	464
Unaccounted gender	15 654
Mean age (years)	Varied by study,
	range: 26-64;
	mean 37.5
Populations	
Major depressive disorder	1535
Bipolar disorders	111
Schizophrenia spectrum	35
disorders	
Anxiety disorder	224
Healthy individuals	1204
Unaccounted population	13 763
Setting/Devices	
Clinical evaluations (hospital)	8
Smartphone	4
Outcomes	
F0 (standard deviation,	10
variability, and mean)	
Spectral signal and formants	8
Percent jitter	6
Voice intensity	5
Speech rate	4
Percent shimmer	6
Pause duration (connected	3
speech)	
Harmonic-to-noise ratio	3
Mel-frequency cepstral	5
coefficients	

publications. Our analysis identified significant reporting gaps: 13 763 participants (81.6% of the total sample) lacked clear diagnostic categorization, while 15 654 participants (92.8%) had unspecified gender information. These gaps primarily stem from one study of 14 898 participants, which reported depression screening but provided limited demographic and diagnostic details.

Voice biomarkers and Al models

All included studies assessed voice as a potential biomarker for depression, incorporating various acoustic, spectral, and linguistic analyses. Prosodic features were examined in ten studies, ^{11–14},16,17,19,20,22 primarily focusing on fundamental frequency variations, speech rate, and voice intensity. These studies consistently reported decreased fundamental frequency and reduced speech tempo in depressed individuals. While voice intensity was analyzed in several studies, quantitative findings specific to this parameter were not consistently reported, representing a gap in the current literature that warrants further investigation. Spectral features, including MFCCs and spectral tilt, were analyzed in eight studies. ^{12–14},16,17,19,20,22</sup> Jitter and

	i
TABLE 3. QUADAS-2 Analysis	

UUADAS-Z Analysis	ılysis					
			Reference			Final QUADAS-2
References	Patient Selection Bias Index Test Bias	Index Test Bias	Standard Bias	Flow and Timing Bias	Applicability Concerns	Judgment
Mazur et al,	High (social media	Low (validated ML	High (self-reported	High (cross-sectional,	Moderate (US and	High risk of bias due to
2025	recruitment,	system,	PHQ-9, no	:	Canada sample,	selection
	selection bias)	rigorous feature	independent clinical	no predictive validity	but limited	method and lack of
		extraction)	Verification)	assessment)	generalizability)	clinical Validation
Menne et al,	High (small	Low (validated MIL	Moderate	High (cross-sectional,	Moderate (Germany	Moderate risk of blas
2024	psycniatric sample,	model,	(psycniatric evaluation,	no rollow-up	and France	aue to
	voluntary	rigorous speech	but no external	to assess predictive	sample, limited	limited sample and lack
	recruitment)	feature extraction)	validation)	validity)	generalizability)	of external validation
Ghosh et al,	High (public dataset,	Low (transformer-	High (clinical	High (cross-sectional,	Moderate (DAIC-WOZ	High risk of bias due to
2024 ¹³		pased	interview dataset,	dataset	dataset,	reliance on
				augmentation		
	potential sample bias)		but no direct	instead of	but no broad	a single dataset and lack
		validation)	clinician validation)	longitudinal validation)	demographic diversity)	of real-world validation
Huang et al,	High (small dataset,	Low (wav2vec 2.0	High (PHQ-8 labels,	High (cross-sectional	Moderate (DAIC-WOZ	High risk of bias due to
2024 ¹⁴	limited	pretrained	no clinician	study, no	dataset,	dataset
	linguistic and cultural	model, automatic	confirmation)	long-term follow-up)	unclear applicability to	constraints and no
	diversity)	leature extraction)			wider populations)	external cillical validation
Ronneberg	Moderate	Low (validated	Low (validated	Moderate (short-term	Moderate (US-based	Moderate risk of bias
et al, 2024 ¹⁵	(randomized trial, limited to	AI-PST	psychological	follow-up	clinical sample,	due to trial
	mild-moderate	system, previous	scales, clinical	but no long-term	potential	limitations but higher
	depression)	pilot RCT)	supervision)	validation)	recruitment bias)	validity than others
Cansel et al,	High (small,	Low (clearly	High (clinical	High (cross-	Moderate (single-site	High risk of bias due to
202316	nonrandom sample,	defined ML	diagnosis, but	sectional, no	study,	sample
	voluntary	model, rigorous	no independent	follow-up for	no linguistic diversity)	selection and lack of
- :	participation)	teature extraction)	validation)	predictive validity)		independent validation
Berardi et al, 2023 ¹⁷	Hign (small,	Low (validated ML	High (clinical	High (cross-sectional	High (single-language stridy	High risk of bias due to small
222			diagnosis commined	3144y,	Study,	singil
	sample, psychlatric in-/outpatients)	rigorous acoustic feature extraction)	by SCID-I, but no external validation)	no predictive validation)	IImited generalizability)	sample size and lack or longitudinal validation
Li et al, 2023 ¹⁸	High (no control	Moderate	High (HAM-D-17	High (cross-sectional	High (linguistic	High risk of bias due to
	group, only	(validated NLP	clinical validation,	study,	based only,	missing speech
	depressive patients)	approach, but no	but no external	no predictive	lacks multimodal	features and lack of
		speech features)	standard)	tracking)	assessment)	observational behavioral data
						555

TABLE 3 (Continued)

Figure F							
High (selection Bias Index Test Bias Standard Bias Flow and Timing Bias Applicability Concerns High (clinical Migh (clinical Migh (clinical Migh (cross-sectional High (limited to Korean recording environment) Feature extraction on caterial sample, model, recruitment, rigorous acoustic by DSM-5, but no no predictive controlled speech task, and locks and potential locks and lockers and				Reference			Final QUADAS-2
High (selection bias, CNN model, aganosis confirmed study, and accontrolled recutation) external validation, and controlled speech task, and according environment) High (small sample, norder) Alph (small sample, self-reported history but no accoustic but no model, first-episode MDD that is a patients) High (small sample, norder) High (small controlled speech task, and lacks norder) High (small sample, norder) High (small speech, norder) Hig	References	Patient Selection Bias	Index Test Bias	Standard Bias	Flow and Timing Bias	Applicability Concerns	Judgment
no external CNN model, diagnosis confirmed study, and controlled speech task, and controlled secure extraction) external validation) tracking) recording environment) High (small sample, imited to model, self-replored history but no accustic but no independent of depression) High (small sample, limited to model, environment) High (small sample, but no external validation) tracking) driven, limited to model, environment high (cross-sectional High (imited to model, environment) tracking) High (small sample, but no independent no predictive model, environmential bias) High (HAM-D used, High (cross-sectional High (limited to Danish limited to model, environment) tracking) High (small sample, but no independent no predictive may not generalize to patients) analysis) clinical validation) racking, may not generalize to patients) and no follow-up on languages) High (lamel control model, environmental stracking) analysis) clinical validation) no predictive may not generalize to patients) and no follow-up on languages)	Kim et al,	High (selection bias,	Low (validated	High (clinical	High (cross-sectional	High (limited to Korean	High risk of bias due to
validation, rigorous acoustic by DSM-5, but no no predictive controlled speech task, and controlled returne extraction) external validation) tracking) and controlled model, and potential control work of depression) and potential sample, but no acoustic but no independent nonusers and potential control walidated MLP (AMD-D used, study, and potential control walidated MLP (AMD-D used, study, and potential control walidated MLP (AMD-D used, study, but no acoustic but no independent no predictive but no acoustic but no independent no predictive multimodal assessment) al, High (small sample, but no acoustic clinical validation) tracking) analysis) analysis) trained on clinical validation no predictive may not generalize to patients) and no clinical validation) no predictive may not generalize to and no clinical validation) nonremitters)	2023	no external	CNN model,	diagnosis confirmed	study,	language,	selection bias,
and controlled feature extraction) external validation) tracking) and lacks recording recording environment) High (small sample, normandom rigorous acoustic by HAM-D, but no analysis) and potential sample, but no external strict-episode MDD trained on patients) High (small control of depression) and potential of depression) High (small control of depression) and potential of depression) High (HAM-D used, High (cross-sectional high (smertphone-street) and no clinical validation) tracking) and no clinical validation) tracking) trained on patients) High (small sample, moderate High (HAM-D used, High (cross-sectional High (limited to Danish Imited to model, and no clinical validation) tracking, and no clinical validation) tracking, and no clinical validation no predictive may not generalize to and no clinical validation no predictive may not generalize to and no clinical validation no premitters) High (small sample, moderate High (HAM-D used, High (cross-sectional High (limited to Danish Imited to model, and no clinical validation) non-memitters) High (small sample, moderate High (HAM-D used, High (cross-sectional High (limited to may not generalize to and no clinical validation) non-memitters)		validation,	rigorous acoustic	by DSM-5, but no	no predictive	controlled speech task,	controlled recording,
recording environment) High (small sample, noderate ML High (clinical validation) High (small sample, but no acoustic of depression) High (small sample, commercial bias) Woderate High (small sample, but no acoustic of depression) High (small sample, commercial bias) Woderate High (HAM-D used, High (cross-sectional High (limited to model, first-episode MDD trained on patients) High (small sample, world acoustic but no independent no predictive multimodal assessment) High (small sample, world acoustic but no independent no predictive multimodal assessment) High (small sample, world acoustic but no independent no predictive multimodal assessment) High (small sample, world acoustic but no independent no predictive multimodal assessment) High (small sample, world acoustic but no independent no predictive multimodal assessment) High (small sample, world acoustic but no independent no predictive multimodal assessment) High (small sample, world acoustic but no external study, and no clinical validation) no predictive may not generalize to and no clinical validation nonremitters) High (small sample, model, first-episode MDD trained on clinical validation) no predictive may not generalize to and no clinical validation nonremitters)		and controlled	feature extraction)	external validation)	tracking)	and lacks	and lack of
environment) High (small sample, model, model, moderate diagnosis confirmed and potential control walidated NL Moderate history but no acoustic high (small sample, model, analysis) High (small sample, woderate history but no acoustic high (small sample, walidated ML Moderate high (small sample, walidated ML Moderate high (small sample, walidated ML model, first-episode MDD trained to model, tracking) High (ham-D used, high (cross-sectional high (imited to Danish first-episode MDD trained on emotional speech, and no follow-up on languages) High (small sample, walidated ML but no external study, first-episode MDD trained on emotional speech, and no follow-up on languages) High (small sample, walidated ML but no external study, languages) High (ham-D used, High (cross-sectional high (limited to Danish study, languages) High (ham-D used, High (cross-sectional high (limited to Danish study, languages)) High (ham-D used, High (cross-sectional high (limited to Danish study, languages)) High (ham-D used, High (cross-sectional high (limited to Danish study, languages)) High (small sample, walidation) no predictive may not generalize to tracking, and no follow-up on languages)		recording				real-world	external validation
High (small sample, model, and potential control depression) High (small sample, Low (validated ML High (clinical validation) tracking) Tecruitment, rigorous acoustic by HAM-D, but no no predictive commercial bias) Self-reported history but no acoustic self-reported history but no acoustic depression) High (HAM-D used, High (cross-sectional High (limited to Danish limited to model, first-episode MDD trained on patients) High (clinical validation) tracking) High (HAM-D used, High (cross-sectional High (limited to Danish limited to model, and no clinical validation) tracking, and no clinical validation) tracking, and no clinical validation) no predictive may not generalize to tracking, and no clinical validation) no predictive may not generalize to and no clinical depression) and no clinical validation) noremitters)		environment)				generalizability)	
nonrandom model, diagnosis confirmed study, independent recruitment, rigorous acoustic by HAMD, but no predictive but commercially driven, limited dataset) commercial bias) High (small control Moderate NLP GAD-7 used, study, but no acoustic of depression) High (small sample, (validated MLP Moderate (validated MLP) but no external validation) High (small sample, (validated ML) but no external model, chained to model, and no clinical validation) High (small sample, model) First-episode MDD tracking) First-episode MDD tracking and no clinical validation) First-episode MDD tracking and no clinical validation) First-episode MDD tracking and no clinical validation) First-episode MDD tracking and no clinical validation non-germanic and no clinical depression) First-episode MDD tracking and no clinical validation non-germanic and no clinical depression) First-episode MDD tracking and no clinical validation non-germanic and no clinical validated more tracking and no clinical more tracking and no clinical more tracking and no clinical more tracking non-germanic put tracking non-germanic put tracking datasets.	Wasserzug	High (small sample,	Low (validated ML	High (clinical	High (cross-sectional	High (language-	High risk of bias due to
recruitment, rigorous acoustic by HAM-D, but no predictive but commercially and potential feature extraction) external validation) tracking) driven, limited dataset) commercial bias) High (small control Moderate High (PHQ-8 and High (cross-sectional High (smartphonegroup, model, model, first-episode MDD trained on patients) recruitment, recruitment feature extraction) external wild fracking) tracking, and locational speech, and no clinical validation) recruitment in organization patients) recruitment to model, first-episode MDD trained on patients) and no clinical validation) no predictive may not generalize to tracking, and no follow-up on languages) recruitment tracking driven, limited dataset) driven, limited dataset) and no clinical validation) no predictive may not generalize to tracking, and no clinical depression) nonremitters)	et al, 2023 ²⁰	nonrandom	model,	diagnosis confirmed	study,	independent	small sample,
and potential feature extraction) external validation tracking) driven, limited dataset) commercial bias) High (small control Moderate High (PHQ-8 and High (cross-sectional High (smartphonegroup, model, but no acoustic but no independent no predictive model assessment) High (small sample, Moderate High (HAM-D used, High (cross-sectional High (limited to model, model, model, and no clinical validation) rained on patients) and no clinical validation no predictive may not generalize to tracking and no clinical validation nonremitters) depression) Areaching driven, limited dataset) High (PHQ-8 and High (cross-sectional High (smartphone-based, excludes mountained and no clinical validation) no predictive may not generalize to tracking, and no clinical and no clinical validation) nonremitters)		recruitment,	rigorous acoustic	by HAM-D, but no	no predictive	but commercially	commercial funding,
commercial bias) High (small control Moderate High (PHQ-8 and group, model, self-reported history but no acoustic of depression) High (small sample, (validated ML model, model, model, first-episode MDD trained on patients) High (small sample, model, model, first-episode MDD trained on patients) High (small sample, model, model, first-episode MDD trained on patients) High (small sample, model, model, first-episode MDD trained on patients) High (ham-D but no external study, moderate emotional speech, and no clinical validation) no predictive may not generalize to tracking, and no clinical depression) High (small sample, model, first-episode MDD trained on patients) and no clinical depression) High (small control High (HAM-D used, High (cross-sectional High (limited to Danish study, mon-Germanic and no clinical validation) no predictive may not generalize to tracking, and no clinical validation) no predictive may not generalize to tracking, and no clinical validations no remitters)		and potential	feature extraction)	external validation)	tracking)	driven, limited dataset)	and lack of external
High (small control group, model, self-reported history of depression)Moderate (validated NLP analysis)High (PHQ-8 and GAD-7 used, study, tracking)High (cross-sectional tracking)High (small sample) study, tracking)High (small sample) tracking)High (ham-D used, tracking)High (cross-sectional study, model, first-episode MDDHigh (HAM-D used, trained on emotional speech, and no clinical and no clinicalHigh (HAM-D used, trained on tracking, and no clinical non-Germanic non-GermanicHigh (smartphone- high (smartphone- study, non-Germanic and no clinical non-model,High (small sample, model, first-episode MDD and no clinical depression)High (HAM-D used, high (cross-sectional study, non-Germanic and no clinical non-model, and no clinical hon-model,High (cross-sectional study, non-Germanic and no clinical non-model, and no clinical non-model,							validation
group, (validated NLP GAD-7 used, study, model, model, self-reported history but no acoustic but no independent no predictive nonusers, and lacks of depression) analysis) clinical validation) tracking) tracking) analysis) clinical validation) tracking) multimodal assessment) High (small sample, Moderate High (HAM-D used, High (cross-sectional High (limited to Danish limited to model, model, trained on clinical validation) no predictive may not generalize to patients) and no clinical and no clinical depression) non-remitters)	Lin et al, 2022 ²¹		Moderate	High (PHQ-8 and	High (cross-sectional	High (smartphone-	High risk of bias due to
model, self-reported history but no acoustic but no independent no predictive nonusers, and lacks of depression) analysis) clinical validation) tracking) trained on patients) and no clinical validation) no predictive may not generalize to tracking, and no clinical validation) no predictive non-Germanic and no clinical nonremitters)		group,	(validated NLP	GAD-7 used,	study,	based, excludes	small control
self-reported history but no acoustic but no independent no predictive nonusers, and lacks of depression) analysis) clinical validation) tracking) tracking) multimodal assessment) High (small sample, Moderate High (HAM-D used, High (cross-sectional High (limited to Danish Imited to model, first-episode MDD trained on clinical validation) no predictive may not generalize to patients) emotional speech, and no follow-up on languages) non-Germanic and no clinical nonremitters)			model,				
of depression) analysis) clinical validation) tracking) multimodal assessment) High (small sample, Moderate High (HAM-D used, High (cross-sectional High (limited to Danish Imited to Moderate ML but no external study, model, first-episode MDD trained on clinical validation) no predictive may not generalize to patients) emotional speech, and no follow-up on languages) non-Germanic and no clinical nonremitters)		self-reported history	but no acoustic	but no independent	no predictive	nonusers, and lacks	group, lack of
High (small sample, Moderate High (HAM-D used, High (cross-sectional High (limited to Danish model, model, first-episode MDD trained on clinical validation) no predictive may not generalize to patients) emotional speech, and no clinical and no clinical non-remitters) and no clinical non-remitters)		of depression)	analysis)	clinical validation)	tracking)	multimodal	independent clinical
High (small sample, Moderate High (HAM-D used, High (cross-sectional High (limited to Danish Inmited to Nalidated ML but no external study, and no clinical validation) no predictive may not generalize to tracking, and no clinical and no clinical non-germanic and no clinical non-generalize to tracking, and no follow-up on languages)						assessment)	
High (small sample, Moderate High (HAM-D used, High (cross-sectional High (limited to Danish Innited to (validated ML but no external study, model, first-episode MDD trained on clinical validation) no predictive may not generalize to patients) and no clinical and no clinical non-Germanic and no clinical homemitters)							validation, and
High (small sample, Moderate High (HAM-D used, limited to Danish limited to limited to (validated ML put no external model, model, lirst-episode MDD trained on patients) study, language, language, language, may not generalize to tracking, and no clinical validation)							commercial funding
limited to (validated ML but no external study, language, model, model, first-episode MDD trained on clinical validation) no predictive may not generalize to patients) emotional speech, and no clinical and no follow-up on languages) depression) nonremitters)	Hansen et al,	High (small sample,	Moderate	High (HAM-D used,	High (cross-sectional	High (limited to Danish	High risk of bias due to
trained on clinical validation) no predictive may not generalize to emotional speech, tracking, non-Germanic and no clinical and no follow-up on languages) depression)	2022 ²²	limited to	(validated ML model,	but no external	study,	language,	small
emotional speech, tracking, non-Germanic and no clinical and no follow-up on languages) depression)		first-episode MDD	trained on	clinical validation)	no predictive	may not generalize to	sample size, lack of
l and no follow-up on languages) nonremitters)		patients)	emotional speech,		tracking,	non-Germanic	external
nonremitters)			and no clinical		and no follow-up on	languages)	validation, and limited
			depression)		nonremitters)		generalizability

Abbreviations: Al-PST, artificial intelligence-problem-solving therapy; CNN, convolutional neural network; DAIC-WOZ, Distress Analysis Interview Corpus—Wizard of Oz; DSM-5, Diagnostic and Statistical Manual of Mental Disorders, 5th Edition; GAD-7, Generalized Anxiety Disorder 7-item scale; HAM-D/HAM-D, Hamilton Depression Rating Scale; HAM-D-17, Hamilton Depression Rating Scale; HAM-D-17, Hamilton Depression Rating Scale; 17-item version); ML, machine learning; MDD, major depressive disorder; NLP, natural language processing; PHQ-8, Patient Health Questionnaire (8-item version); QUADAS-2, Quality Assessment of Diagnostic Accuracy Studies, 2nd version; RCT, randomized controlled trial; SCID-1, Structured Clinical Interview for DSM Disorders, Axis I; US, United States; wav2vec 2.0, speech representation learning framework.

shimmer, representing measures of voice perturbation, were investigated in six studies, ^{11–14,16,19} with findings suggesting increased perturbation among individuals with depression compared with healthy controls. Pause duration and harmonic-to-noise ratio were assessed in three studies, ^{12,17,20} indicating longer pause times and decreased harmonicity in depressed speech patterns.

Supervised machine learning was implemented in all 12 studies, with SVMs and random forest classifiers being the most used algorithms. Classification accuracies varied across studies, with SVMs achieving performance ranging from 78% to 96.48%. Deep learning models, particularly convolutional neural networks, were utilized in three studies, ^{13,14,21} demonstrating good predictive performance in voice-based depression detection. Cross-validation techniques were applied in seven studies, ^{12–14,16,19,20,22} with five-fold and tenfold cross-validation being the most frequently used approaches.

Predictive validity of voice biomarkers

The diagnostic accuracy of voice-based artificial intelligence models varied depending on the specific classification task and methodological approach. In distinguishing depressed individuals from healthy controls, classification performance ranged from AUC values of 0.71 to 0.93. 12,22 The differentiation between depression and other psychiatric conditions yielded mixed results, performing in distinguishing depression from schizophrenia with 93.2% accuracy¹⁷ and in differentiating MDD from bipolar disorder. 18 Longitudinal studies demonstrated the feasibility of using voice biomarkers for symptom tracking, with individual-specific models showing stronger predictive correlations compared with population-level approaches. 15,21 Home monitoring systems reported promising results, with specific correlations between vocal features and standardized depression scales. 11,19,20 Sensitivity values in real-world applications ranged from 71.3% to 96.05%, 11,14 while specificity values ranged from 73.5% to 99.23%. 11

Correlation with symptom scales and clinical validation

The correlation between voice biomarkers and established depression rating scales was assessed in most included studies. Significant correlations were observed between fundamental frequency variability and Hamilton Depression Rating Scale scores, with vocal parameters associated with symptom severity. Several studies demonstrated moderate correlations between voice-derived features and depression severity as measured by the Patient Health Questionnaire-9, while others reported significant associations between jitter, shimmer, and Clinical Global Impression scores.

Temporal stability of these correlations was evaluated in longitudinal studies, ^{15,21} with follow-up durations varying across studies. Individual-specific models demonstrated

stronger clinical correlations than population-based models, supporting the potential for personalized voice-based mental health monitoring.¹⁵

Timing and settings

Data collection strategies varied among the studies, as four integrated mobile sensing technologies, ^{11,19–21} while eight implemented structured voice recording tasks in clinical settings. ^{12–18,22} The integration of artificial intelligence voice analysis into digital health platforms was evident, with studies using real-world voice samples collected through smartphone applications and telemedicine platforms. Recording frequency ranged from single-session clinical assessments to repeated voice recordings collected over several months. Circadian influences on voice characteristics were addressed in two studies, ^{18,21} with findings suggesting potential variations in voice parameters depending on time-of-day recording.

Bias analysis

The methodological quality of included studies was assessed using the QUADAS-2 tool. High risk of bias was identified in six studies, \$11,14,18-21\$ primarily due to patient selection biases and lack of external validation. Issues related to reference standard bias were present in four studies, \$11,14,18,19\$ where voice-derived features were compared against self-reported symptom scales rather than clinician-rated assessments. Cross-validation and independent test set validation were adequately performed in seven studies, \$12-14,16,19,20,22\$ while three studies lacked robust validation techniques. \$11,18,21\$ The generalizability of findings was limited in studies that focused on a single language or cultural group, with applicability concerns highlighted in four studies. \$17-19,22\$

DISCUSSION

This systematic review synthesized evidence from twelve studies investigating the role of voice biomarkers in the detection and monitoring of depression. The findings suggest that acoustic and prosodic features of speech, when analyzed using machine learning techniques, can provide valuable insights into depression severity, symptom progression, and mood state classification. While the results demonstrate promising classification performance, methodological heterogeneity, limitations in generalizability, and potential biases warrant further investigation.

Across studies, fundamental frequency (F0) variability, speech rate, and voice intensity were consistently identified as key prosodic markers of depression. ^{11–14,17,19} Depressed individuals exhibited lower mean F0, decreased speech tempo, and reduced vocal intensity, supporting prior literature on psychomotor retardation in depression. Spectral features, including mel-frequency cepstral coefficients (MFCC) and spectral tilt, were also implicated. ^{12–14,16,17,19}

Additionally, voice perturbation measures such as jitter and shimmer were frequently elevated in individuals with depression, 11–13,16,19 reflecting increased vocal instability.

Machine learning models demonstrated moderate-to-high accuracy in differentiating depressed individuals from healthy controls. Support vector machines (SVM) and random forest classifiers were the most frequently employed supervised learning techniques, ^{12,13,16,17,19} achieving classification accuracies ranging from 70.9% to 96.5%. Deep learning models, particularly convolutional neural networks (CNNs), were implemented in select studies, ^{13,14,19} showing good performance in voice-based depression detection. However, despite encouraging results, variability in feature selection, classification tasks, and performance metrics complicates direct comparison across studies.

Longitudinal studies assessing the predictive validity of voice biomarkers revealed that individualized models outperformed population-based approaches. ^{15,21} Personalized models demonstrated stronger correlations with symptom severity as measured by standard depression scales, ^{12,18–20,22} reinforcing the potential for real-world deployment of voice-based monitoring tools. Notably, smartphone-based assessments emerged as a growing trend, with studies incorporating passive voice monitoring and structured voice recording tasks. ^{11,19–21}

Despite promising results, several limitations must be addressed. First, methodological heterogeneity across studies limits comparability. Differences in sample sizes, diagnostic criteria, recording environments, and feature extraction methods introduce variability in findings. Some studies 11,14,21 relied on self-reported symptom scales for depression classification rather than clinician-rated assessments, raising concerns about reference standard bias. Additionally, the absence of standardized voice data collection protocols poses challenges for reproducibility and external validation.

Second, demographic disparities and reporting inconsistencies hinder generalizability. A substantial proportion of participants had unspecified demographic information, particularly regarding gender distribution. 11,13–15,18,21 Where reported, female participants were overrepresented in MDD cohorts, potentially affecting model performance and applicability across populations. Furthermore, the age range of participants varied widely, from early adulthood to late-life depression, suggesting that voice biomarkers may exhibit different predictive properties across developmental stages.

Third, cultural and linguistic variability remains a critical issue. Most studies did not account for language differences in voice analysis, which may influence the acoustic properties of speech and the effectiveness of AI-based classification models. As voice biomarkers are inherently influenced by phonetic structure, future research should consider cross-linguistic validation to enhance model robustness and global applicability.

To advance the field, standardized protocols could include recording both structured speech tasks and

naturalistic conversation in the same individuals; collecting samples at multiple time points to account for circadian variations; documenting acoustic parameters using open-source analysis tools with transparent parameters; and implementing machine learning with mandatory cross-validation techniques. Additionally, future datasets should prioritize demographic diversity through balanced gender representation, inclusion of multiple age groups (adolescent, adult, and geriatric), cross-linguistic sampling with standardized translation protocols, and representation of varied depression severities and subtypes verified through standardized clinical assessments.

To advance the field of voice-based depression detection, future studies should prioritize methodological standardization, demographic diversity, and cross-linguistic validation. Establishing large-scale, multisite datasets with standardized voice recording protocols will enhance reproducibility and facilitate model comparison. Additionally, leveraging explainable AI techniques could improve clinical interpretability by identifying which specific vocal features contribute most significantly to depression classification.

Further research is also needed to explore the ecological validity of voice biomarkers in real-world settings. The transition toward passive voice monitoring via smartphones and wearable devices presents an opportunity for continuous, nonintrusive mental health assessment. 24,25 However, ensuring data privacy and ethical considerations remains paramount. Finally, integrating voice biomarkers with multimodal data—such as facial expression analysis and physiological signals—could enhance predictive accuracy and enable a more comprehensive approach to digital mental health monitoring.

CONCLUSION

The findings of this review highlight the potential use of AI to analyze voice biomarkers as a novel tool for depression detection and monitoring. While current evidence demonstrates promising classification accuracy, methodological heterogeneity and generalizability concerns must be addressed before widespread clinical adoption. Standardized protocols, diverse datasets, and real-world validation efforts will be crucial for translating voice-based AI models into practical psychiatric applications. With continued advancements in AI and digital health, voice biomarkers could play a pivotal role in the future of remote mental health assessment and personalized treatment strategies.

Data Availability Statement

Not applicable.

Declaration of Competing Interest

The authors have no financial interest in the subject under discussion. All authors have read and approved the paper.

Acknowledgments

None.

Author Contributions

Giovanni Briganti: Design, acquisition of data, data analysis and interpretation, drafting, final approval, and accountability for the work; final approval of the version to be published; agreement to be accountable for all aspects of the work in ensuring that questions related to the accuracy or integrity of any part of the work are appropriately investigated and resolved. Jerome R. Lechien: Design, acquisition of data, data analysis and interpretation, drafting, final approval, and accountability for the work; final approval of the version to be published; agreement to be accountable for all aspects of the work in ensuring that questions related to the accuracy or integrity of any part of the work are appropriately investigated and resolved.

Institutional Review Board Statement Not required.

Informed Consent Statement

Not applicable.

Supplementary Materials

None.

References

- Goldberg D. The heterogeneity of "major depression". World Psychiatry. 2011;10:226–228.
- Belmaker RH, Agam G. Major depressive disorder. N Engl J Med. 2008;358:55–68
- 3. Marx W, Penninx BW, Solmi M, et al. Major depressive disorder. *Nat Rev Dis Primer*. 2023;9:44.
- 4. Strawbridge R, Young AH, Cleare AJ. Biomarkers for depression: recent insights, current challenges and future prospects. *Neuropsychiatr Dis Treat.* 2017;13:1245–1262.
- Alpert M, Kurtzberg RL, Friedhoff AJ. Transient voice changes associated with emotional stimuli. Arch Gen Psychiatry. 1963;8:362–365.
- Scherer KR. Vocal affect expression: a review and a model for future research. *Psychol Bull.* 1986;99:143.
- Scherer KR, Banse R, Wallbott HG, Goldbeck T. Vocal cues in emotion encoding and decoding. *Motiv Emot.* 1991;15:123–148.

- Arevian AC, Bone D, Malandrakis N, et al. Clinical state tracking in serious mental illness through computational analysis of speech. *PloS One*. 2020;15:e0225695.
- Koops S, Brederoo SG, De Boer JN, Nadema FG, Voppel AE, Sommer IE. Speech as a biomarker for depression. CNS Neurol Disord - Drug Targets. 2023;22:152–160.
- Talavera JA, Saiz-Ruiz J, Garcia-Toro M. Quantitative measurement of depression through speech analysis. Eur Psychiatry. 1994;9:185–193.
- Mazur A, Costantino H, Tom P, Wilson MP, Thompson RG. Evaluation of an AI-based voice biomarker tool to detect signals consistent with moderate to severe depression. *Ann Fam Med.* 2025;23:60-65.
- Menne F, Dörr F, Schräder J, et al. The voice of depression: speech features as biomarkers for major depressive disorder. BMC Psychiatry. 2024;24:794.
- Ghosh D, Karande H, Gite S, Pradhan B. Psychological disorder detection: a multimodal approach using a transformer-based hybrid model. *MethodsX*. 2024;13:102976.
- 14. Huang X, Wang F, Gao Y, et al. Depression recognition using voice-based pre-training model. *Sci Rep.* 2024;14:12734.
- Ronneberg CR, Lv N, Ajilore OA, et al. Study of a PST-trained voice-enabled artificial intelligence counselor for adults with emotional distress (SPEAC-2): design and methods. *Contemp Clin Trials*. 2024;142:107574.
- 16. Cansel N, Faruk Alcin Ö, Furkan Yılmaz Ö, Ari A, Akan M, Ucuz İ.A. New artificial intelligence-based clinical decision support system for diagnosis of major psychiatric diseases based on voice analysis. *Psychiatr Danub*. 2023;35:489–499.
- Berardi M, Brosch K, Pfarr JK, et al. Relative importance of speech and voice features in the classification of schizophrenia and depression. *Transl Psychiatry*. 2023;13:298.
- Li N, Feng L, Hu J, et al. Using deeply time-series semantics to assess depressive symptoms based on clinical interview speech. Front Psychiatry. 2023;14:1104190.
- Kim AY, Jang EH, Lee SH, Choi KY, Park JG, Shin HC. Automatic depression detection using smartphone-based text-dependent speech signals: deep convolutional neural network approach. *J Med Internet* Res. 2023;25:e34474.
- Wasserzug Y, Degani Y, Bar-Shaked M, et al. Development and validation of a machine learning-based vocal predictive model for major depressive disorder. J Affect Disord. 2023;325:627–632.
- Lin D, Nazreen T, Rutowski T, et al. Feasibility of a machine learningbased smartphone application in detecting depression and anxiety in a generally senior population. *Front Psychol.* 2022;13:811517.
- 22. Hansen L, Zhang YP, Wolf D, Sechidis K, Ladegaard N, Fusaroli R. A generalizable speech emotion recognition model reveals depression and remission. *Acta Psychiatr Scand.* 2022;145:186–199.
- 23. Hamilton M. A rating scale for depression. *J Neurol Neurosurg Psychiatry*. 1960;23:56–62.
- Kroenke K, Spitzer RL, Williams JB. The PHQ-9: validity of a brief depression severity measure. J Gen Intern Med. 2001;16:606–613.
- Davidson J, Turnbull CD, Strickland R, Miller R, Graves K. The Montgomery-Asberg depression scale: reliability and validity. *Acta Psychiatr Scand.* 1986;73:544–548.