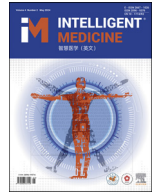


## Detecting somatisation disorder via speech: introducing the Shenzhen Somatisation Speech Corpus

Kun Qian, Ruolan Huang, Zhihao Bao, Yang Tan, Zhonghao Zhao, Mengkai Sun, Bin Hu, Björn W. Schuller, Yoshiharu Yamamoto

### Angaben zur Veröffentlichung / Publication details:

Qian, Kun, Ruolan Huang, Zhihao Bao, Yang Tan, Zhonghao Zhao, Mengkai Sun, Bin Hu, Björn W. Schuller, and Yoshiharu Yamamoto. 2024. "Detecting somatisation disorder via speech: introducing the Shenzhen Somatisation Speech Corpus." *Intelligent Medicine* 4 (2): 96–103. <https://doi.org/10.1016/j.imed.2023.03.001>.



## Research Article

# Detecting somatisation disorder via speech: introducing the Shenzhen Somatisation Speech Corpus

Kun Qian<sup>1,2,#</sup>, Ruolan Huang<sup>3,4,#</sup>, Zhihao Bao<sup>1,2,#</sup>, Yang Tan<sup>1,2</sup>, Zhonghao Zhao<sup>1,2</sup>, Mengkai Sun<sup>1,2</sup>, Bin Hu<sup>1,2,\*</sup>, Björn W. Schuller<sup>5,6</sup>, Yoshiharu Yamamoto<sup>7</sup>

<sup>1</sup> Key Laboratory of Brain Health Intelligent Evaluation and Intervention, Ministry of Education, Beijing Institute of Technology, Beijing 100081, China

<sup>2</sup> School of Medical Technology, Beijing Institute of Technology, Beijing 100081, China

<sup>3</sup> The First School of Clinical Medicine, Southern Medical University, Guangzhou, Guangdong 510515, China

<sup>4</sup> Department of Neurology, Shenzhen University General Hospital, Shenzhen, Guangdong 518055, China

<sup>5</sup> Group on Language, Audio, & Music, Imperial College London, London SW7 2AZ, UK

<sup>6</sup> Embedded Intelligence for Health Care and Wellbeing, University of Augsburg, Augsburg 86159, Germany

<sup>7</sup> Educational Physiology Laboratory, Graduate School of Education, The University of Tokyo, Tokyo 113-0033, Japan



## ARTICLE INFO

## Keywords:

Somatisation disorder  
Machine learning  
Healthcare  
Computer audition

## ABSTRACT

**Objective** Speech recognition technology is widely used as a mature technical approach in many fields. In the study of depression recognition, speech signals are commonly used due to their convenience and ease of acquisition. Though speech recognition is popular in the research field of depression recognition, it has been little studied in somatisation disorder recognition. The reason for this is the lack of a publicly accessible database of relevant speech and benchmark studies. To this end, we introduced our somatisation disorder speech database and gave benchmark results.

**Methods** By collecting speech samples of somatisation disorder patients, in cooperation with the Shenzhen University General Hospital, we introduced our somatisation disorder speech database, the Shenzhen Somatisation Speech Corpus (SSSC). Moreover, a benchmark for SSSC using classic acoustic features and a machine learning model was proposed in our work.

**Results** To obtain a more scientific benchmark, we compared and analysed the performance of different acoustic features, i. e., the full ComPare feature set, or only Mel frequency cepstral coefficients (MFCCs), fundamental frequency (F0), and frequency and bandwidth of the formants (F1–F3). By comparison, the best result of our benchmark was the 76.0% unweighted average recall achieved by a support vector machine with formants F1–F3.

**Conclusion** The proposal of SSSC may bridge a research gap in somatisation disorder, providing researchers with a publicly accessible speech database. In addition, the results of the benchmark could show the scientific validity and feasibility of computer audition for speech recognition in somatization disorders.

## 1. Introduction

According to the World Health Organisation (WHO) report, 1 in 8 people worldwide are suffering from mental disorders in 2019 [1]. Since COVID-19 broke out in 2020, all kinds of mental disorders around the world have become more frequent [2–5]. In particular, the number of major depression and anxiety patients worldwide has increased by more than 25% [6]. Mental disorders have brought severe harm to patients themselves, families, and the community. For the patients themselves, on the one hand, self-harm and attempted suicide are the most harm-

ful behaviours to them, which directly bring physical damage and pain to the patients. As reported in the latest WHO statistics, about 700,000 people worldwide die from suicide every year [7], of which mental disorders account for a large proportion, and many more have attempted suicide. On the other hand, people with psychosis have an increased risk of physical-related illnesses and a shorter life span than the ordinary people [8–10]. For the patients' families, they have to tolerate a series of mental pressures and financial burden brought by the patients. As a result, their life qualities decline fast. For the community, patients with psychotic episodes may cause socially endangering behaviours. They

\* Corresponding author: Bin Hu, Key Laboratory of Brain Health Intelligent Evaluation and Intervention, Ministry of Education, Beijing Institute of Technology, Beijing 100081, China (Email: [bh@bit.edu.cn](mailto:bh@bit.edu.cn)).

# These authors contributed equally to this work.

<https://doi.org/10.1016/j.imed.2023.03.001>

Received 7 October 2022; Received in revised form 23 February 2023; Accepted 7 March 2023

2667-1026/© 2023 Published by Elsevier B.V. on behalf of Chinese Medical Association. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>)

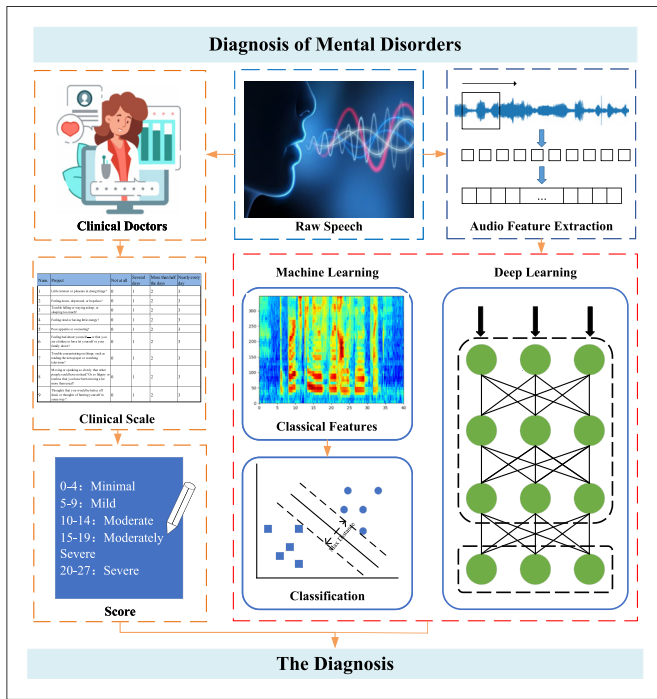


Figure 1. Methods of diagnosing mental disorders.

may attack people around them or even strangers, out of their control. What is more, neurosis is also characterised by a higher prevalence in children and adolescents [11] and in women than in men [12]. Compared to depression and anxiety, somatisation disorder (SD) is a mental disorder that is less regarded. In fact, it also has a higher prevalence and risk. According to a report, the prevalence rate of SD is 10.1% in General Hospital Psychiatric Units Tertiary Care Centres in India [13]. In addition, the prevalence rate in women is higher than that in men. In the study of Babu et al. [14] the prevalence rate of SD in adult women was 40.8%. SD is a neurosis characterised by persistent fear or belief in dominance of various somatic symptoms. Because SD causes patients to shift from emotional and mental distress to the body, they often suffer from unexplained physical discomfort [15]. Such physical discomfort includes stomach pain, back pain, joint pain, and further more. This medically unexplained symptom is characterised by multiple occurrences, persistence, and recurring. Patients are often unaware that they have a SD. Attributing this pain to a physical illness, they repeatedly seek clinical medical advice and treatment, but to no avail [16]. In terms of economic impact, patients frequently seeking medical advice in the wrong direction greatly increase their cost of medical care. The per capita expenditure for health care of patients is up to nine times of the average amount [17]. It also places an indirect burden on the healthcare system [18–20]. According to the report, the annual medical costs of SD in the United States are approximately 2,560 billion dollars [21]. Moreover, unresolved pain and confusion lead to lower quality of life and greater mental stress for patients, which aggravates their condition, or even cause a higher suicide risk [22].

However, the vast majority of patients could not receive effective treatment. This is because the mental health system is severely under-resourced and patients lack of relevant mental health knowledge (especially for some less common mental disorders). Early diagnosis of mental illness is the first step in obtaining beneficial treatment for patients, but the approaches are scarce. As shown in Figure 1, clinical scales and speech analysis are used as the two main methods for the diagnosis of mental disorders introduced in this paper. The scale, combined with the communication and observation of the patient's behaviour, facial expressions, and speech, among others by doctors, is the common form of psychiatric diagnosis by far. But this approach relies on subjective in-

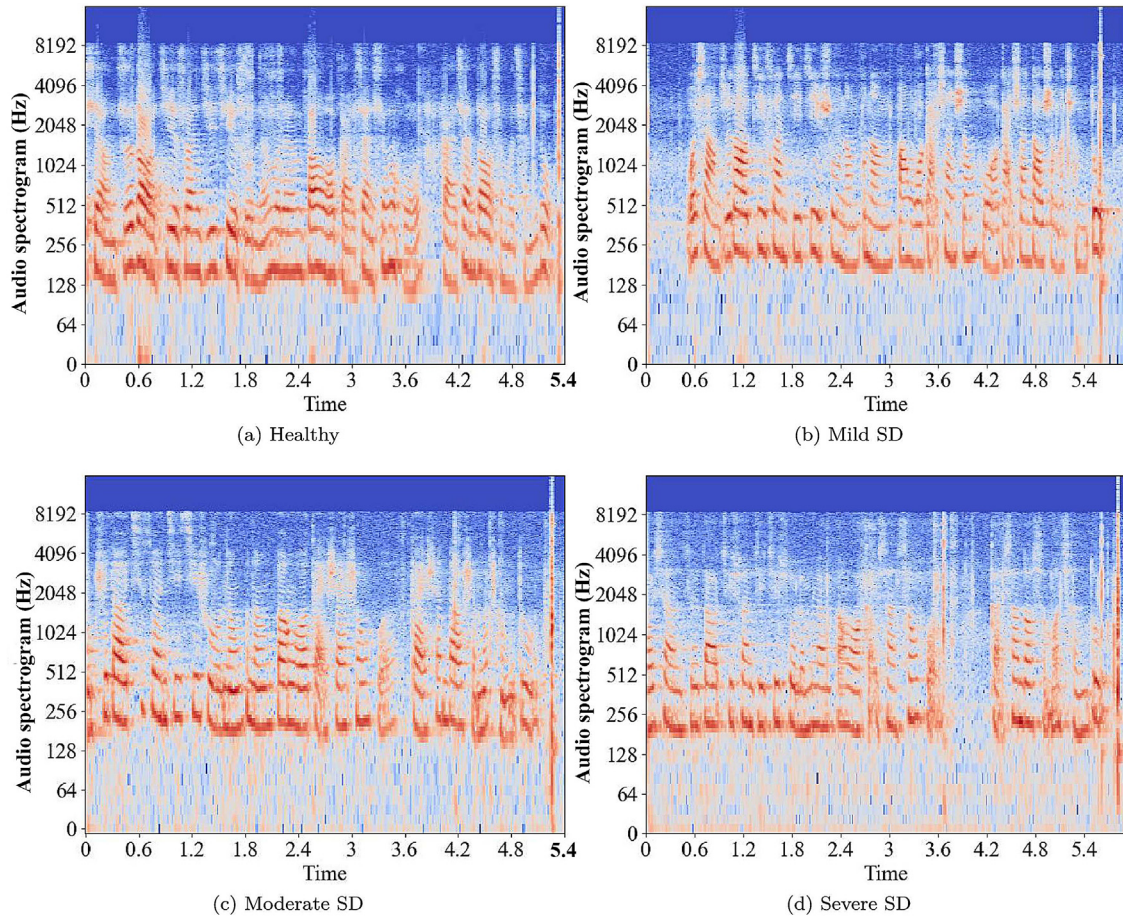
terviews with patients and the clinical experience of the doctors, which results in a partial deviation of the diagnosis from the reality. With the development of computer technology, it is encouraging to see that, artificial intelligence (AI) has been applied to the classification of mental disorders, which avoids the pitfalls of the scale. In particular, audio signals, because of their 'non-invasive' nature, combined with the rapidly developing computer audition (CA) [23] technology are becoming a popular topic of digital medicine research in the search for new digital phenotypes. Speech signals, as a subclass of audio signals, have been demonstrated to be reliable in the diagnosis of certain mental disorders, such as depression and anxiety [24–27]. In particular, related studies have shown that speech features performed better than visual features or text in depression prediction tasks [28–29].

There are a mount of AI researches on speech representation of mental disorders. For anxiety, Wang et al. [24] proposed a new Fourier parameter model using the perceptual content of voice quality. Dan et al. [25] used K-nearest neighbours as classifier and focused on the fundamental frequency for classification. For depression, much work has been done to detect depression through speech. Pan et al. [26] extracted 988 speech features from speech data and established a logistic regression model to achieve a better depression classification rate. Rejaibi et al. [27] proposed a deep learning-based method to assess depression and predict its severity. They aimed to extract Mel frequency cepstral coefficients (MFCCs) from speech and used long short-term memory networks. For insomnia, Espinoza et al. [30] studied obstructive sleep apnea among sleep disorders using a larger speech database containing 426 participants. They used supervector or i-vector techniques to model speech spectral information and predicted by support vector regression. The above researches demonstrated that speech has the ability to represent mental disorders. For SD, unfortunately, due to its mix-up feature of psychologically and physically, it is challenging for clinicians to recognise this psychiatric disorder masquerading as physical pain [31–32]. Not only that, the actual clinical condition may also be diagnosed as a SD, thereby hindering the patient's real need for treatment and longer-term pain [33]. However, there are few reports on the application of AI technology in this field. In particular, we found rare work on speech as raw data for classification. According to our search, Idenfors et al. [34] analysed brain images for global-brain functional connectivity (GFC). The combination of GFC values and support vector machine (SVM) were used to distinguish patients from the controls. The results showed that the patients' GFC was abnormal. Lyu et al. [35] constructed an improved bacterial foraging optimisation-based kernel extreme learning machine (IBFO-KELM) model based on the symbol self-assessment scale (SCL-90) data for the diagnosis of patients with SD. Human communication relies on language communication. From speech, we can feel the others' emotional or mental state. Numerous studies have shown that speech-language pathologists can diagnose mental disorders through speech [36–38]. Therefore, ubiquitous, inexpensive, and easily-acquired speech signals can be used as raw and reliable information for diagnosing or predicting mental disorders.

In this work, we demonstrate the collected data and publish the speech database for the classification of SD. The name we use is the Shenzhen Somatisation Speech Corpus (SSSC). Unlike the automatic speech recognition technology related to the UASpeech databases [39], the related research on the SSSC will belong to the category of mental emotion recognition, although these databases are both composed of speech. A conventional and reproducible benchmark for this publicly accessible speech database is also announced. The main contributions of this work are: (1) We demonstrate the feasibility of using speech to classify SD through the benchmark. (2) We provide a standard speech database for the classification of SD. (3) We compare several typical acoustic features to illustrate the usability of this database.

The remainder of this paper is organised as follows: At first, the materials and methods are presented in Section 2. Subsequently, the results of the benchmark work are given in Section 3. We give an experimental discussion, current limitations, and outlooks in Section 4.





**Figure 2.** Audio spectrogram examples of different PHQ-15 outcomes. SD: somatisation disorder.

## 2. Methods

In this section, we firstly introduce the proposed publicly somatisation speech database, i. e., SSSC. Successively, openSMILE, the toolkit that was used to extract acoustics features in our study is briefly introduced. Then, we give more detailed descriptions of our database work. In the final part, the machine learning method and the optimised strategy we adopted is given.

### 2.1. SSSC database

#### 2.1.1. Data collection

This study was approved by the ethic committee of the Shenzhen University General Hospital. All the participants involved were informed that their voice data will be used only for research purposes. Their agreements for this study were recorded as one of the five following original speech phrases. The data was collected in out-patient-department in the Shenzhen University General Hospital. We asked the participants to speak five sentences (with neutral contextual meaning). At the same time, two self-report questionnaires were answered by the participants regarding their anxiety and SD (physical discomfort). The questionnaire is constructed by two widely used separated questionnaires. The two questionnaires are GAD-7 (Patient anxiety questionnaire 7) and PHQ-15 (Patient health questionnaire 15) [40–42]. All the data were collected from 12 November 2020 to 5 April 2021. We used a Shinco RV-18 recording pen with 32 GB of storage to record all the participants' voices which have a sample rate of 32,000 Hz and a bit rate of 16 bps.

We give examples of the recorded sentences which were spoken in Chinese:

- (1) Today, I am in Shenzhen, Guangdong Province, China.
- (2) I want to know if computer technology can help me improve my life and to what extent.
- (3) She / He is my friend.
- (4) Time is money, efficiency is life.
- (5) I agree to use my voice for emotion recognition.

Figure 2 shows the spectrograms (extracted from *dev0136.wav*, *dev0003.wav*, *dev0010.wav*, and *dev0004.wav*, respectively) corresponding to the second sentence in Chinese.

#### 2.1.2. Data pre-processing

As described, we executed a series of data pre-processing stages before establishing the 'standard' SSSC, which includes data cleansing, voice activity detection, speaker diarisation, and speech transcription [43]. First, we excluded recordings with low quality (e.g., the level of the speech is low compared to the background noise). Then, we removed the non-speech parts (e.g., non-subjects' speech, breathing, and coughing) from each recording, which resulted in maintaining only the segments including voice from the recordings. The segments containing solely the target patient and scripted content (e.g., excluding laughing) were kept. Finally, we obtained 705 audio recordings from 141 participants. To attenuate the effects of the audio recording equipment, the background noise condition and the level of the recording, all files were first high-pass filtered to eliminate low-frequency background noise (cut-off frequency: 120 Hz, 10<sup>th</sup>-order Chebyshev filter) and then their waveforms were normalised individually (peak amplitude set to -3 dB).

#### 2.1.3. Tasks definition

We define two tasks for the benchmark setup: First, the Anxiety Degree should be grouped into: Yes (labelled as "1" when GAD-7  $\geq$  11),

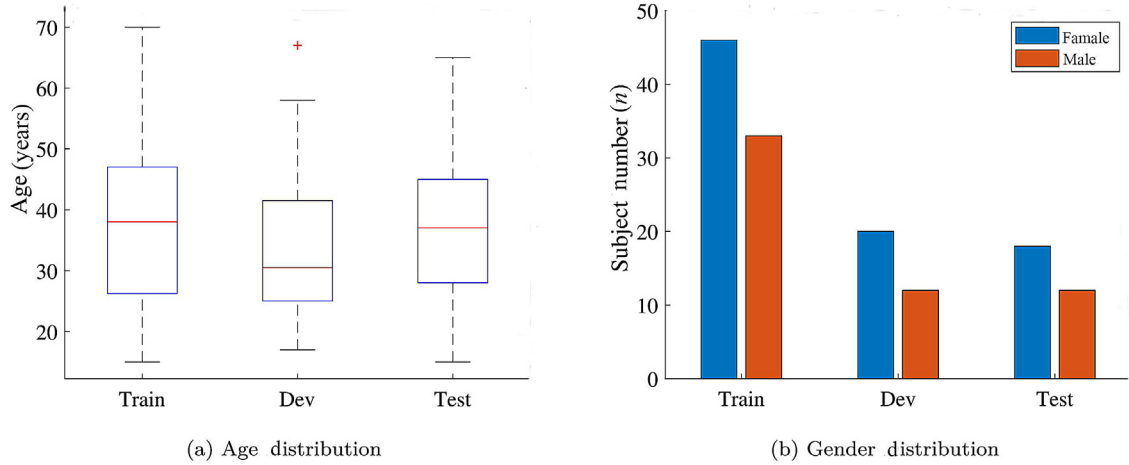


Figure 3. Demographic statistic of the SSSC.

Table 1 Subjects information in data splits

Groups	Average age (years)	Male (n)	Female (n)
Train	38.2	46	33
Dev	35.0	20	12
Test	36.7	18	12
$\Sigma$	37.1	84	57

No (labelled as “0” when GAD-7 < 11). Then, an estimation of the Physical Discomfort Disorder Degree should be made as: Yes (labelled as “1” when PHQ-15  $\geq 10$ ), No (labelled as “0” when PHQ-15 < 10).

#### 2.1.4. Data partitioning

Totally, the database contains as mentioned audio samples of 705 speech events from 141 subjects who were checked without organic disease by experts. The number of males in all subjects is 84, and the number of females is 57. The average age of all subjects is ( $37.1 \pm 13.2$ ) years (range 15 to 70 years). The total duration of audio in the database is 3,039.192 s, equalling roughly 50 minutes. The average sample duration is ( $4.311 \pm 2.297$ ) s (range 0.864 to 16.920 s). In order to carry out the experiment, we randomly partitioned the data into a train, a development (Dev) and a test set, which are subject independence. Table 1 gives more information about data partitions in detail. Figure 3 shows data distribution details.

#### 2.2. openSMILE and acoustic features

openSMILE (open-source Speech and Music Interpretation by Large-space Extraction) is an open source toolkit which is widely applied in the field of acoustic representation extractions [44]. openSMILE can provide features commonly used in classic acoustic signal processing methods, such as short-time zero-crossing rate, energy spectrum features, and Mel frequency cepstrum coefficients. To get the statistical information of an audio signal sample, openSMILE firstly extracts the low-level descriptors (LLDs) from the original frame-level audio signals, then performs the statistical information extraction to the frame-based LLDs by functionals. By this method, the limitations of static machine learning models such as SVM are unlocked from the inconsistency of sample duration.

To train the baseline system, we use the ComPare feature set, which includes 65 LLDs (Table 2). The configuration of the 2.3 version openSMILE is ComPare 2016. The more specific information and details are described [45]. Moreover, they also features a detailed introduction of the used functionals (Table 3). The ComPare feature set includes in

Table 2 The low level descriptors for ComPare feature set

Low level descriptors (LLDs)	Groups
4 Energy related LLDs	
RMSE, zero-crossing rate	Prosodic
Sum of auditory spectrum (loudness)	Prosodic
Sum of RASTA-filtered auditory spectrum	Prosodic
6 Voicing related LLDs	
Probability of voicing	Voice quality
F0 (SHS and Viterbi smoothing)	Prosodic
log HNR, jitter (local and $\delta$ ), shimmer (local)	Voice quality
55 Spectral LLDs	
MFCCs 1–14	Cepstral
Spectral energy 250–650 Hz, 1 k–4 kHz	Spectral
Spectral flux, centroid, entropy, slope	Spectral
Spectral roll-off point 0.25, 0.5, 0.75, 0.9	Spectral
Spectral variance, skewness, kurtosis	Spectral
Psychoacoustic sharpness, harmonicity	Spectral
RASTA-filtered auditory spectral bands 1–26 (0–8kHz)	Spectral

Table 3 The functionals applied to LLDs in the ComPare feature set

Functionals
Arithmetic or positive arithmetic mean
Inter-quartile ranges 1–2, 2–3, 1–3,
Linear regression slope, offset
Linear regression quadratic error
Linear Prediction gain and coefficients 1–5
Mean and std. dev. of peak to peak distances
Peak-valley-peak slopes mean and std. dev.
Peak and valley range (absolute and relative)
Peak mean value and distance to arithmetic mean
Quadratic regression coefficients
Quadratic regression quadratic error
Root-quadratic mean, flatness
Rise time, left curvature time
Relative position of max. and min. value
Range (difference between max. and min. values)
Segment length mean, min., max., std. dev.
Standard deviation, skewness, kurtosis, quartiles 1–3
Temporal centroid
Up-level time 25%, 50%, 75%, 90%
99-th and 1-st percentile, range of these

total 6,373 features. The mechanism of functionals is to map the time series based LLDs to a scalar value per each applied functional (e. g., mean, standard deviation, maximum); then, a single, fixed dimension vector which is independent of the audio signal sample’s time duration is generated [45].

**Table 4** Number of per scales and score in the data splits

Scales	Scores	Label	Train	Dev	Test	$\Sigma$
GAD-7	0~4	0	130	40	45	215
	5~9	1	125	55	50	230
	10~14	2	80	25	30	135
	15~21	3	60	40	25	125
	$\Sigma$		395	160	150	705
PHQ-15	0~4	0	90	15	35	140
	5~9	1	135	80	65	280
	10~14	2	110	45	35	190
	$\geq 15$	3	60	20	15	95
	$\Sigma$		395	160	150	705

### 2.3. Database related questionnaires

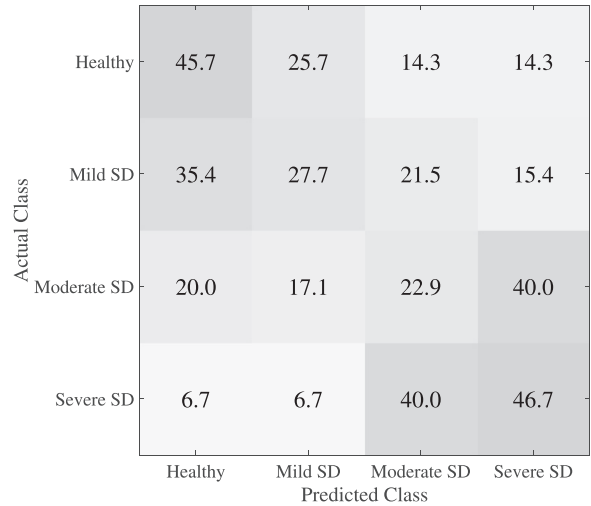
In our study, we carried out measurement work on two scales as mentioned (i.e., GAD-7 and PHQ-15). Based on the scores of these two scales, the Shenzhen University General Hospital's psychologists enrolled in the study classified the conditions into four classes. We give a separate introduction of each scale in this part. Specifically, Table 4 shows the data distribution in details.

Generalised anxiety disorder (GAD) is a common and disabling illness that is often underdiagnosed and undertreated [46]. With the influence of the COVID-19 in China [47], more and more people are suffering from symptoms of anxiety. Therefore, we want to get more information of the GAD for further research. GAD-7 is a brief clinical measure for assessing GAD, which consists of 7 items about anxiety self-report [48]. The score of each item is between 0–3 and the GAD-7's total score is in the range of 0 to 21. Based on the total scores, the Shenzhen University General Hospital's mental health experts classified the participants into 4 types (i.e., no-anxiety for 0–4, mild anxiety for 5–9, moderate anxiety for 10–14, and severe anxiety for 15–21).

SD is a prevalent condition which is not well treated by many psychiatrists [49]. Patients usually seek care in the medical setting convinced that they suffer from physical discomfort rather than a mental disorder. Based on this situation, the SD patients often encounter ineffective care and even harm. That somatising patients may represent 40% or more of the ambulatory medicine patient population greatly magnifies the problem [50]. To reduce unnecessary expenses, it is very significant for patients with SD to be fully recognised and treated. The classification of SD by the speech of the patients is not only feasible, but also appears efficient and convenient [51]. The PHQ-15 comprises 15 somatic symptoms from the Patient Health Questionnaire (PHQ), each symptom scored from 0 (“not bothered at all”) to 2 (“bothered a lot”) [52]. According to the PHQ-15 scale's total score ranging from 0 to 30, the participants are divided into 5 types (i.e., minimal for 0–4, low for 5–9, medium for 10–14, and high for 15–30).

### 2.4. Machine learning method and optimising strategy

Support vector machine (SVM) is a stable and by now ‘traditional’ classifier. To make this study comparable and reproducible, we use an SVM classifier with linear kernel to conduct all experiments. We train an SVM model with the complexity parameter in the range of  $\{10^{-8}, 10^{-7}, \dots, 10^{-1}, 1.0\}$ . Then, we choose the complexity that performs best on the development set to classify the test set. Moreover, both the training set and the development set train-level set were joined to predict on test data. Upsampling the training set and the train-level set was used for balancing the dataset. At the same time, we processed the feature sets with feature normalisation. In order to obtain more specific results, we evaluate the performance of different acoustic features, i.e., the full ComPare feature set, or selectively only MFCCs, fundamental frequency (F0), and frequency and bandwidth of the formants (F1–F3). All the features are extracted by the openSMILE feature extraction and audio analysis tool.



**Figure 4.** Normalised confusion matrices (values are shown in (%)) of the best-performing (development) model for the PHQ-15 test set. SD: somatisation disorder.

### 3. Results

Although the database includes two scales, our experiment series focuses on the PHQ-15 outcome. This can reflect the order of severity of the subjects with SD. Moreover, we conduct related experiments to analyse if there is any connection between GAD-7 and PHQ-15. Classification results are shown in Tables 5 and 6 with feature type and the final feature number. The best mean unweighted average recall (UAR) per feature on the development set and test set are highlighted.

In order to obtain a better understanding of the data set, we make some changes for the listed labels in Table 4. We set a threshold on the scores. We consider the participants whose scores higher than the threshold as affected by the condition and the healthy when their scores lower than the threshold. According to different thresholds, we set two discrimination modes named “A” and “B” for PHQ-15 (i.e., threshold of each mode, GAD-7: 10, PHQ-15A: 5, PHQ-15B: 10). This modification helps us distinguish the participants with specific condition from the healthy subjects. As described above, PHQ-15 reflects the degree of physical discomfort. We process the labels based on the mentioned threshold above, and then operate ‘AND’ or ‘OR’ on them given their relatedness under the umbrella of being psychological disorders.

According to the UAR indicators, it appears remarkable that the formants F1–F3 perform best on most of the classification models. This means, F1–F3 will provide more information in our tasks. Except for PHQ-15 labelled with four types scores, the classification results of other tasks achieve the UAR higher than 50.0%. Figures 4 and 5 show confusion matrices of the best-performing setup for the test set. From Figure 4, we can find that the predicted labels focus on label 0 and label 3. As could be suspected, the classifiers tend to prefer predictions of healthy subjects or such with severe conditions. Although the classifier performs better on this task, unlike humans, the machine requires more samples to understand the features.

### 4. Discussion

Referring to Tables 5 and 6, we can see that, although F1–F3 has the least number of dimensions, it performs better than the other feature sets. However, the performance is relatively poor when we fed all the features in the ComPare feature set into the model. We have noted that some redundant features led to decrease model performance. In future work, we will analysis the contribution of features. Overall, the formant-based features can represent the phenomenon of interest efficiently.

**Table 5** Classification results of different feature types and tasks (%)

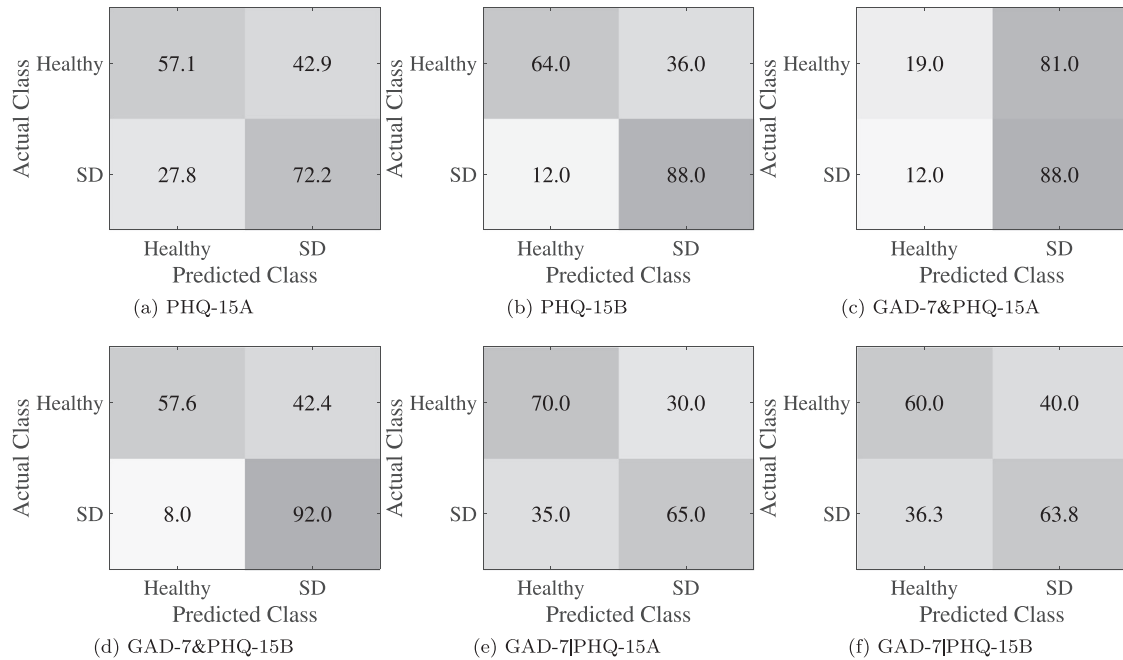
Feature type	Number of features (n)	PHQ-15	PHQ-15A	PHQ-15B	GAD-7&PHQ-15A
ComPare	6,373	24.4/22.1	53.4/45.9	50.8/49.0	49.8/50.0
MFCCs (all)	1,400	28.3/29.5	<b>58.4/64.7</b>	54.8/61.5	51.6/53.5
MFCCs (only coef)	756	26.4/35.7	58.2/62.7	57.4/68.0	50.4/47.0
MFCCs (only delta)	644	25.7/27.1	52.3/52.9	53.2/56.5	<b>53.8/50.5</b>
F0	24	23.9/24.7	47.2/49.1	54.8/69.5	47.8/50.5
F1-F3	14	<b>33.1/28.3</b>	49.3/52.0	<b>57.6/76.0</b>	51.0/54.5

PHQ-15: Labels by way of four scores. PHQ-15A: Labels by way of A; PHQ-15B: Labels by way of B; GAD-7&PHQ-15A: AND operation on GAD-7 and PHQ-15A labels; \*/\*: The first one means the best result of unweighted average recall (UAR) on the development set in all experiments, the last one means the test result by the model that performed best on the development set.

**Table 6** Classification results of different feature types and tasks (%)

Feature type	Number of features (n)	GAD-7&PHQ-15B	GAD-7 PHQ-15A	GAD-7 PHQ-15B
ComPare	6,373	51.4/47.6	56.8/50.0	53.9/54.5
MFCCs(all)	1,400	52.1/59.2	59.8/53.8	54.3/57.2
MFCCs(only coef)	756	51.2/60.4	<b>61.4/56.2</b>	54.3/56.6
MFCCs(only delta)	644	50.7/51.2	53.3/48.8	53.2/60.5
F0	24	<b>55.9/56.8</b>	47.9/50.0	<b>55.3/57.4</b>
F1-F3	14	55.0/74.8	51.0/67.5	53.5/61.9

GAD-7&PHQ-15B: AND operation on GAD-7 and PHQ-15B labels; GAD-7|PHQ-15A: OR operation on GAD-7 and PHQ-15A labels; GAD-7|PHQ-15B: OR operation on GAD-7 and PHQ-15B labels; \*/\*: The first one means the best result of unweighted average recall (UAR) on the devolpment set in all experiments, the last one means the test result by the model that performed best on the development set.

**Figure 5.** Normalised confusion matrices (values are shown in (%)) of the best-performing (development) model for the test set. SD: somatisation disorder.

On the other hand, we only use an SVM classifier to test the representational ability of the extracted features. Furthermore, the number of feature types selected are few. In the future, we expect other work on the database to increase the number of features and compare more models to improve the ability to represent this task.

SSSC is the first speech dataset available for SD classification. Of great interest, SD could only be recognised shortly when the doctor is skilful and major physical disorders are excluded. Moreover, some symptoms similar to SD should be excluded carefully. The training set in this study had been reviewed beforehand to exclude severe physical disorders such as stroke or coronary artery diseases. The models in current

benchmarks achieve over 50% of UAR, yet more would be needed to expand the reliance. We provide a public database for the study and use of AI and CA. Recognition and primary estimation of the mental state or mood by CA could be the first step. Adjusting the response or reaction could be the next step follow up, which may generate more importance. When different international standards are used to classify categories, the standard we use generates different results. As the first database, the future of SSSC is undoubted. One of the key features of mental disorders is that they are diagnosed without an obvious objective criterion or examinations. Therefore, independent and skilful doctors with psychological experiences are of great importance. Unfortunately, it is be-



coming more difficult and expensive for doctors likewise to diagnose SD. AI recognition would be the first step for screening before searching for medical help. Similar to AI assistance in the medical imaging or pathological field, CA could surely be of helpful assistance for the primary and unskillful individuals to recognise SD, and even reduce unnecessary anxiety.

SSSC is the first speech dataset available for SD classification, a heretofore untapped resource for such public data. Most of the models in current benchmarks achieve over 50% of UAR on the development set. However, SSSC also inevitably underlies limitations:

(1) Limited dataset. Due to monocentricity, it is difficult for us to recruit a sufficient number of well-represented subjects. This limited amount of data is not conducive to the application of deep learning on it. But in fact, that SSSC can be trusted enough to be adopted for the study of SD was demonstrated through our benchmark experiments. Therefore, we encourage researchers to give more consideration to the application of traditional machine learning which may be more suitable for the small-scale databases. More importantly, in future works, data augmentation is worth looking into and researching, which should be expected to reduce the model's lack in generalisation. An easy way is to use some toolboxes for data enhancement. For instance, Maguolo et al. [53] published a matlab toolkit that provides 15 different augmentation algorithms for raw audio data and 8 for spectrograms.

(2) Overlapping symptoms. A study has shown overlap between somatic symptoms, anxiety, and depression [54]. In other words, the three are co-morbid and triggering each other. Thus, different mental illnesses may show the same features. This may be challenging to model, yet, on the contrary, may facilitate better results. As an emerging trend, multi-label learning [55] requires more discovery and exploration in this field.

(3) Reliability of labels. The labels stem from scoring participants' scale questionnaires. As mentioned above, a small number of labels generated by the scale score may be deviating from the actual situation, as these are greatly influenced by the subjective interviews of participants and the clinical experience of doctors. This is an inevitable mistake caused by the lack of objective factors. In the future, one can involve more professional doctors in the scoring, set more differential versions to recognise SD and use multi-modality to ensure reliability of the label. Nevertheless, manual annotation of speech data is an expensive and time-consuming task. In order to overcome this difficulty, in the AI filed, we think that self-supervised learning [56–57] could be used to reduce the reliance on labels, as well as introducing active learning [58] on this database.

(4) Availability of samples. There are incomplete or redundant samples of participants' speech content in the database. Therefore, those samples affect the SD audio analysis. We will exhaustively screen and remove these samples to provide researchers with a more scientific and credible new version.

## Conflicts of interest statement

The authors declare that they have no conflicts of interest.

## Funding

This work was partially supported by the Ministry of Science and Technology of the People's Republic of China with the STI2030-Major Projects (Grant No. 2021ZD0201900), the National Natural Science Foundation of China (Grant Nos. 62227807 and 62272044), the Teli Young Fellow Program from the Beijing Institute of Technology, the Shenzhen Municipal Scheme for Basic Research (Grant Nos. JCYJ20210324100208022 and JCYJ20190808144005614), China, the JSPS KAKENHI (Grant No. 20H00569), the JST Mirai Program (Grant No. 21473074), the JST MOONSHOT Program (Grant No. JPMJMS229B).

## Author contributions

**Kun Qian:** Visualization, Conceptualization, Funding acquisition, Writing – original draft, Supervision. **Ruolan Huang:** Data curation, Funding acquisition, Writing – original draft. **Zhihao Bao:** Writing – original draft. **Yang Tan:** Investigation, Writing – original draft. **Zhonghao Zhao:** Methodology, Writing – original draft. **Mengkai Sun:** Writing – review & editing. **Bin Hu:** Writing – review & editing, Supervision, Funding acquisition. **Björn W. Schuller:** Writing – review & editing. **Yoshiharu Yamamoto:** Writing – review & editing.

## References

- [1] World Health Organisation. Mental disorders. Available from <https://www.who.int/news-room/fact-sheets/detail/mental-disorders> (Accessed on 1 July 2022).
- [2] Vindegaard N, Benros ME. Covid-19 pandemic and mental health consequences: systematic review of the current evidence. *Brain Behav Immun* 2020;89:531–42. doi:10.1016/j.bbi.2020.05.048.
- [3] Racine N, Cooke JE, Eirich R, et al. Child and adolescent mental illness during covid-19: a rapid review. *Psychiatry Res* 2020;292:113307. doi:10.1016/j.psychres.2020.113307.
- [4] Moreira WC, Sousa Ard, Nóbrega MdPSdS. Mental illness in the general population and health professionals during covid-19: a scoping review. *Texto e Contexto Enferm* 2020;29. doi:10.1590/1980-265x-tce-2020-0215.
- [5] Carr MJ, Steeg S, Webb RT, et al. Effects of the covid-19 pandemic on primary care-recorded mental illness and self-harm episodes in the uk: a population-based cohort study. *Lancet Public Health* 2021;6(2):e124–35. doi:10.1016/S2468-2667(20)30288-7.
- [6] World Health Organisation. Mental health and covid-19: early evidence of the pandemic's impact: scientific brief, 2 march 2022. Available from [https://www.who.int/publications/i/item/WHO-2019-nCoV-Sci\\_Brief-Mental\\_health-2022.1](https://www.who.int/publications/i/item/WHO-2019-nCoV-Sci_Brief-Mental_health-2022.1), 2022 (Accessed on 1 July 2022).
- [7] World Health Organisation. World health statistics 2022: monitoring health for the sdgs, sustainable development goals. Available from <https://www.who.int/publications/i/item/9789240051157>, 2022 (Accessed on 1 July 2022).
- [8] Launders N, Kirsh L, Osborn DP, et al. The temporal relationship between severe mental illness diagnosis and chronic physical comorbidity: a uk primary care cohort study of disease burden over 10 years. *Lancet Psychiatry* 2022;9(9):725–35. doi:10.1016/S2215-0366(22)00225-5.
- [9] Firth J, Siddiqi N, Koyanagi A, et al. The lancet psychiatry commission: a blueprint for protecting physical health in people with mental illness. *Lancet Psychiatry* 2019;6(8):675–712. doi:10.1016/S2215-0366(19)30132-4.
- [10] Rodgers M, Dalton JE, Harden M, et al. Integrated care to address the physical health needs of people with severe mental illness: a mapping review of the recent evidence on barriers, facilitators and evaluations. *Int J Integr Care* 2018;18(1):9. doi:10.5334/ijic.2605.
- [11] Merikangas KR, Nakamura EF, Kessler RC. Epidemiology of mental disorders in children and adolescents. *Dialogues Clin Neurosci* 2009;11(1):7–20. doi:10.3188/DCNS.2009.11.1/kmerikangas.
- [12] Gove WR. The relationship between sex roles, marital status, and mental illness. *Soc Forces* 1972;51(1):34–44. doi:10.1093/sf/51.1.34.
- [13] Chander KR, Manjunatha N, Binukumar B, et al. The prevalence and its correlates of somatization disorder at a quaternary mental health centre. *Asian J Psychiatr* 2019;42:24–7. doi:10.1016/j.ajp.2019.03.015.
- [14] Babu AR, Aswathy Sreedevi AJ, Krishnapillai V. Prevalence and determinants of somatization and anxiety among adult women in an urban population in kerala. *Indian J Community Med* 2019;44(Suppl 1):S66–9. doi:10.4103/ijcm.IJCM\_55\_19.
- [15] Witthöft M, Gerlach AL, Bailer J. Selective attention, memory bias, and symptom perception in idiopathic environmental intolerance and somatoform disorders. *J Abnorm Psychol* 2006;115(3):397–407. doi:10.1037/0021-843X.115.3.397.
- [16] Barsky AJ, Orav EJ, Bates DW. Distinctive patterns of medical care utilization in patients who somatize. *Med Care* 2006;44(9):803–11. doi:10.1097/01.mlr.0000228028.07069.59.
- [17] Smith Jr GR, Monson RA, Ray DC. Psychiatric consultation in somatization disorder. *N Engl J Med* 1986;314(22):1407–13. doi:10.1056/NEJM198605293142203.
- [18] Weiss FD, Rief W, Kleinstäuber M. Health care utilization in outpatients with somatoform disorders: descriptives, interdiagnostic differences, and potential mediating factors. *Gen Hosp Psychiatry* 2017;44:22–9. doi:10.1016/j.genhosppsych.2016.10.0.
- [19] Lipowski JJ. Somatization: the concept and its clinical application. *Am J Psychiatry* 1988;145(11):1358–68. doi:10.1176/ajp.145.11.1358.
- [20] Fink P. The use of hospitalizations by persistent somatizing patients. *Psychol Med* 1992;22(1):173–80. doi:10.1017/S0033291700032827.
- [21] Barsky AJ, Orav EJ, Bates DW. Somatization increases medical utilization and costs independent of psychiatric and medical comorbidity. *Arch Gen Psychiatry* 2005;62(8):903–10. doi:10.1001/archpsyc.62.8.903.
- [22] Chioqueta AP, Stiles TC. Suicide risk in patients with somatization disorder. *Crisis* 2004;25(1):3. doi:10.1027/0227-5910.25.1.3.
- [23] Qian K, Li X, Li H, et al. Computer audition for healthcare: opportunities and challenges. *Front Digit Health* 2020;2:5. doi:10.3389/fdgh.2020.00005.



- [24] Wang K, An N, Li BN, et al. Speech emotion recognition using fourier parameters. *IEEE Trans Affect Comput* 2015;6(1):69–75. doi:10.1109/TAFFC.2015.2392101.
- [25] Zbancioc MD, Feraru SM. A study about the automatic recognition of the anxiety emotional state using emo-db. 2015 E-Health and Bioengineering Conference (EHB). Romania: IEEE; 2015. doi:10.1109/EHB.2015.7391506.
- [26] Pan W, Wang J, Liu T, et al. Depression recognition based on speech analysis. *Chin Sci Bull* 2018;63(20):2081–92. doi:10.1360/N972017-01250.
- [27] Rejaibi E, Komaty A, Meriaudeau F, et al. Mfcc-based recurrent neural network for automatic clinical depression recognition and assessment from speech. *Biomed Signal Process Control* 2022;71:103–7. doi:10.1016/j.bspc.2021.103107.
- [28] Williamson JR, Godoy E, Cha M, et al. Detecting depression using vocal, facial and semantic communication cues. Amsterdam: The Netherlands; 2016. doi:10.1145/2988257.2988263.
- [29] Yang L, Sahli H, Xia X, et al. Hybrid depression classification and estimation from audio video and text information. New York, the: United States; 2017. doi:10.1145/3133944.3133950.
- [30] Espinoza-Cuadros F, Fernández-Pozo R, Toledano DT, et al. Reviewing the connection between speech and obstructive sleep apnea. *Biomed Eng Online* 2016;15(1):1–20. doi:10.1186/s12938-016-0138-5.
- [31] Smith RC, Dwamena FC. Classification and diagnosis of patients with medically unexplained symptoms. *J Gen Intern Med* 2007;22(5):685–91. doi:10.1007/s11606-006-0067-2.
- [32] Raffagnato A, Angelico C, Valentini P, et al. Using the body when there are no words for feelings: alexithymia and somatization in self-harming adolescents. *Front Psychiatry* 2020;11. doi:10.3389/fpsy.2020.00262. 262–261
- [33] De Jonge L, Petrykiv S, Fennema J, et al. Misdiagnosis of loin pain hematuria syndrome as a somatization disorder. *Eur Psychiatry* 2017;41(S1). doi:10.1016/j.eurpsy.2017.01.597. S491–S491
- [34] Pan P, Ou Y, Su Q, et al. Voxel-based global-brain functional connectivity alterations in first-episode drug-naïve patients with somatization disorder. *J Affect Disord* 2019;254:82–9. doi:10.1016/j.jad.2019.04.099.
- [35] Lv X, Chen H, Zhang Q, et al. An improved bacterial-foraging optimization-based machine learning framework for predicting the severity of somatization disorder. *Algorithms* 2018;11(2):17–34. doi:10.3390/a11020017.
- [36] Solot CB, Sell D, Mayne A, et al. Speech-language disorders in 22q11. 2 deletion syndrome: best practices for diagnosis and management. *Am J Speech Lang Pathol* 2019;28(3):984–99. doi:10.1044/2019\_AJSLP-16-0147.
- [37] Brito-Marcelino A, Oliva-Costa EF, Sarmento SCP, et al. Burnout syndrome in speech-language pathologists and audiologists: a review. *Rev Bras Med Trab* 2020;18(2):217. doi:10.47626/1679-4435-2020-480.
- [38] Lanzi AM, Ellison JM, Cohen ML. The “counseling +” roles of the speech-language pathologist serving older adults with mild cognitive impairment and dementia from alzheimer's disease. *Perspect ASHA Spec Interest Groups* 2021;6(5):987–1002. doi:10.1044/2021\_PERSP-20-00295.
- [39] Kim H, Hasegawa-Johnson M, Perlman A, et al. Dysarthric speech database for universal access research. *Proceedings of Ninth Annual Conference of the International Speech Communication Association*; 2008.
- [40] Shih YC, Chou CC, Lu YJ, et al. Reliability and validity of the traditional chinese version of the gad-7 in taiwanese patients with epilepsy. *J Formos Med Assoc* 2022;1–7. doi:10.1016/j.jfma.2022.04.018.
- [41] Cano-García Javier F, Muñoz-Navarro R, Abad Sesé A, et al. Latent structure and factor invariance of somatic symptoms in the patient health questionnaire (phq-15). *J Affect Disord* 2020;261:21–9. doi:10.1016/j.jad.2019.09.077.
- [42] Huang XJ, Ma HY, Wang XM, et al. Equating the phq-9 and gad-7 to the hads depression and anxiety subscales in patients with major depressive disorder. *J Affect Disord* 2022;311:327–35. doi:10.1016/j.jad.2022.05.079.
- [43] Han J, Qian K, Song M, et al. An early study on intelligent analysis of speech under covid-19: severity, sleep quality, fatigue, and anxiety. In: *Proceedings of INTERSPEECH*, Shanghai, China; 2020. doi:10.1109/EHB.2015.7391506.
- [44] Eyben F, Wöllmer M, Schuller B. *Proceedings of the 18th ACM International Conference on Multimedia. Opensmile: the munich versatile and fast open-source audio feature extractor*. Italy: Firenze; 2010. <https://doi.org/10.1145/1873951.1874246>
- [45] Eyben F. *Real-time speech and music classification by large audio feature space extraction*. Cham, Switzerland: Springer International Publishing; 2015. doi:10.1007/978-3-319-27299-3.
- [46] DeMartini J, Patel G, Fancher TL. Generalized anxiety disorder. *Ann Intern Med* 2019;170(7):ITC49–64. doi:10.7326/ATTC201904020.
- [47] Huang Y, Zhao N. Generalized anxiety disorder, depressive symptoms and sleep quality during covid-19 outbreak in china: a web-based cross-sectional survey. *Psychiatry Res* 2020;288:112954. doi:10.1016/j.psychres.2020.112954.
- [48] Spitzer RL, Kroenke K, Williams JB, et al. A brief measure for assessing generalized anxiety disorder: the gad-7. *Arch Intern Med* 2006;166(10):1092–7. doi:10.1001/archinte.166.10.1092.
- [49] Mai F. Somatization disorder: a practical review. *Can J Psychiatry* 2004;49(10):652–62. doi:10.1177/070674370404901002.
- [50] Smith RC. Somatization disorder. *J Gen Intern Med* 1991;6(2):168–75. doi:10.1007/BF02598318.
- [51] Oxman TE, Rosenberg SD, Schnurr PP, et al. Linguistic dimensions of affect and thought in somatization disorder. *Am J Psychiatry* 1985. doi:10.1176/ajp.142.10.1150.
- [52] Kroenke K, Spitzer RL, Williams JB. The phq-15: validity of a new measure for evaluating the severity of somatic symptoms. *Psychosom Med* 2002; 64(2):258–66. doi:10.1097/00006842-200203000-00008.
- [53] Maguolo G, Paci M, Nanni L, et al. Audiogmenter: a matlab toolbox for audio data augmentation. *Appl Comput Informat* 2021. doi:10.1108/ACI-03-2021-0064.
- [54] Fu X, Zhang F, Liu F, et al. Brain and somatization symptoms in psychiatric disorders. *Front Psychiatry* 2019;10:146. doi:10.3389/fpsy.2019.00146.
- [55] Liu W, Wang H, Shen X, et al. The emerging trends of multi-label learning. *IEEE Trans Pattern Anal Mach Intell* 2021;44(11):7955–74. doi:10.1109/TPAMI.2021.3119334.
- [56] Chen T, Kornblith S, Swersky K, et al. Big self-supervised models are strong semi-supervised learners. *Adv Neural Inf Process Syst* 2020;33:22243–55.
- [57] Lee H, Hwang SJ, Shin J. Self-supervised label augmentation via input transformations. *PMLR*; 2020. doi:10.1109/TPAMI.2021.3119334.
- [58] Qian K, Zhang Z, Baird A et al. Active learning for bird sound classification via a kernel-based extreme learning machine. *J Acoust Soc Am* 2017;142(4):1796–804. doi:10.1121/1.5004570.